

# Volltext via OCR

## Möglichkeiten und Grenzen

<charParams l="1120" t="210" r="1235" b="282" suspicious="true" characterHeight="42" hasUncertainHeight="f...  
</charParams></formatting></Line></par>  
<par leftIndent="5" rightIndent="35" lineSpacing="53">  
<line baseline="347" l="491" t="292" r="1801" b="365"><formatting lang="OldGerman" ff="Arial" fs="13." spa...  
<charParams l="491" t="307" r="507" b="339" suspicious="true" characterHeight="32" hasUncertainHeight="fal...  
</charParams>  
<charParams l="509" t="307" r="525" b="339" characterHeight="32" hasUncertainHeight="false" baseLine="0" w...  
</charParams>  
<charParams l="525" t="306" r="541" b="340" suspicious="true" characterHeight="32" hasUncertainHeight="fal...  
</charParams></formatting><formatting lang...  
<charParams l="543" t="302" r="554" b="343" characterHeight="32" hasUncertainHeight="false" baseLine="0" w...  
</charParams>  
<charParams l="558" t="296" r="576" b="348" suspicious="true" characterHeight="32" hasUncertainHeight="fal...  
</charParams>  
<charParams l="573" t="307" r="592" b="338" suspicious="true" characterHeight="32" hasUncertainHeight="fal...  
</charParams>  
<charParams l="592" t="296" r="600" b="347" suspicious="true" characterHeight="32" hasUncertainHeight="fal...  
</charParams>  
<charParams l="600" t="303" r="624" b="336" characterHeight="32" hasUncertainHeight="false" baseLine="0" w...  
</charParams>  
<charParams l="627" t="304" r="646" b="338" characterHeight="32" hasUncertainHeight="false" baseLine="0" w...  
</charParams>  
<charParams l="647" t="303" r="669" b="338" characterHeight="32" hasUncertainHeight="false" baseLine="0" w...  
</charParams>  
<charParams l="669" t="301" r="678" b="347" suspicious="true" characterHeight="32" hasUncertainHeight="fal...  
</charParams>  
<charParams l="678" t="302" r="689" b="338" suspicious="true" characterHeight="32" hasUncertainHeight="fal...  
</charParams>  
<charParams l="689" t="301" r="711" b="339" suspicious="true" characterHeight="32" hasUncertainHeight="fal...  
</charParams>  
<charParams l="713" t="302" r="727" b="338" characterHeight="32" hasUncertainHeight="false" baseLine="0" w...  
</charParams>  
<charParams l="728" t="303" r="752" b="348" characterHeight="32" hasUncertainHeight="false" baseLine="0" w...  
</charParams>  
<charParams l="753" t="303" r="768" b="339" characterHeight="32" hasUncertainHeight="false" baseLine="0" w...  
</charParams>  
<charParams l="768" t="304" r="790" b="338" characterHeight="32" hasUncertainHeight="false" baseLine="0" w...

**Maria Federbusch**  
**Christian Polzin**

# **Volltext via OCR – Möglichkeiten und Grenzen**

Testszzenarien zu den Funeralschriften der  
Staatsbibliothek zu Berlin – Preußischer Kulturbesitz

*Mit einem Erfahrungsbericht von Thomas Stäcker  
aus dem Projekt „Helmstedter Drucke Online“  
der Herzog August Bibliothek Wolfenbüttel*

## Impressum

Beiträge aus der Staatsbibliothek zu Berlin – Preußischer Kulturbesitz  
Band 43

Herausgegeben von Barbara Schneider-Kempf  
Generaldirektorin der Staatsbibliothek zu Berlin

Layout, Druck: Gallery Print, Berlin (Boris Brumnjak, Leo Lamprecht)  
Buchbinderische Verarbeitung: Stein + Lehmann, Berlin  
ISBN 978-3-88053-185-7

© 2013 Staatsbibliothek zu Berlin – Preußischer Kulturbesitz



# Inhalt

Zum Geleit

## 1 Einführung

1.1	Projekthintergrund	7
1.2	Hilfsthese	8
1.3	Kurzfassung der Ergebnisse	9

## 2 Materialspezifik der Funeralschriften

2.1	Gattungsspezifische Aspekte	11
2.2	Sammlungsspezifische Aspekte	12
2.3	Typographische und buchgestalterische Aspekte	14
2.4	Orthographische Besonderheiten	19
2.5	Verwendung gebrochener Schriften	21
2.6	Diversität vorkommender Schrifttypen	23

## 3 Typische OCR-Teilschritte und Einflussgrößen

3.1	Typischer Workflow	26
3.1.1	Bereitstellung des Bildmaterials	26
3.1.2	Bildvorverarbeitung und Binarisierung	28
3.1.3	Segmentierung (Layoutanalyse, Blockerkennung, Segmentierung auf Wort- und Zeichenebene)	29
3.1.4	OCR (inklusive unmittelbar verbundener Fehlerkorrekturmechanismen)	30
3.1.5	Lexikalische und linguistische Nachkorrektur	32
3.1.6	Semantische Erschließung: Strukturdaten und Named Entities	36
3.1.7	Einfluss der Verarbeitung ganzer Seiten	36
3.2	Zusammenfassung	38

## 4 Praktischer Softwaretest und -vergleich I: Die Programme

4.1	Übersicht über Parameter- und Hilfsdateien jedes Produkts	39
4.2	Start der OCR-Verarbeitungskette in der Benutzeroberfläche	41
4.3	Einzelne OCR-Verarbeitungsstufen (ohne Trainings-Phase)	42
4.3.1	Bildvorverarbeitung und Binarisierung	42
4.3.2	Segmentierung	44
4.3.3	OCR	46
4.3.4	Lexikalische Nachkorrektur	48
4.3.5	Semantische Erschließung: Strukturdaten und Named Entities	50
4.3.6	Export	51
4.4	Bereitstellung der OCR-Muster: Anlegen (Training) und Pflege von Musterbibliotheken	52
4.5	Batch-spezifische Aspekte	59
4.6	Voraussetzungen und Randbedingungen des Softwareeinsatzes	61

## 5 Praktischer Softwaretest und -vergleich II:

	<b>Optimierungsversuche durch Sortierung, Binarisierungsparameter, Training und Lexika</b>	<b>62</b>
--	--	-----------

<b>6</b>	<b>Zusammenfassung der Ergebnisse</b>	
6.1	Einflussfaktoren auf die Erkennungsgüte	72
6.1.1	Nachweis des Einflusses von Vorsortierung	72
6.1.2	Nachweis des Einflusses von Binarisierungsparametern	72
6.1.3	Nachweis des Trainingseffekts	73
6.1.4	Nachweis des Wörterbucheffekts	74
6.2	Im Projekt erreichte Erkennungsgüte	75
6.3	Hervorzuhebende Vor- und Nachteile der getesteten Software für die OCR deutschsprachiger Alter Drucke	78
6.3.1	<i>Good practice</i>	78
6.3.2	Gegenwärtige Nachteile	78
6.4	Anforderungen an Software-Weiterentwicklung	79
6.5	Schlussfolgerungen und Tipps; Ausblick	82
6.5.1	Tipps für Anwender der getesteten Software: Wann wäre welches der getesteten Produkte geeignet?	83
6.5.2	Tipps für generelle Planungen von OCR Alter Drucke	85
<b>7</b>	<b>Anhänge</b>	
7.1	Beispiele Exportformat	87
7.2	Kurzanleitungen für einfache Fälle	98
7.2.1	OCR (Muster- und Wortbibliothek müssen vorliegen)	98
7.2.2	Training (Anlegen bzw. Erweitern von Musterbibliotheken)	100
7.3	Details zu einzelnen Parametern	102
7.4	Einige Mindestanforderungen an Bedieneroberflächen	116
7.5	Literatur- und Linkliste	118
<b>8</b>	<b>Thomas Stäcker: Erfahrungsbericht Helmstedter Drucke Online an der Herzog August Bibliothek Wolfenbüttel</b>	
8.1	Zusammenfassung	123
8.2	Bedeutung der Helmstedter Drucke	123
8.3	Umfang und Art der Helmstedter Druckproduktion	124
8.4	Projektziele, Rahmenbedingungen	125
8.5	Ergebnisse, Erfahrungen	126
8.6	Überlegungen zu Textgenauigkeiten und deren Messung	129
8.7	Messung der Textgenauigkeit	131
8.7.1	Testen der Güte einer Software	132
8.7.2	Überprüfen der Güte eines Texts	133
8.8	Ökonomische Betrachtungen	135
<b>9</b>	<b>Register</b>	138

## Zum Geleit

Texterkennung oder auch Optische Zeichenerkennung (OCR) ist ein zentrales Thema der Wissenschaftsgesellschaft am Beginn des 21. Jahrhunderts. Die Image-Digitalisierung macht heute schon Bücher jederzeit, überall und in kürzester Zeit verfügbar, doch erst über durchsuchbare Volltexte potenzieren sich die Möglichkeiten wissenschaftlichen Erkenntnisgewinns.

Bibliotheken mit bedeutenden historischen Beständen können hier entscheidend zum Entstehen entsprechender virtueller Forschungsumgebungen beitragen. Die Staatsbibliothek zu Berlin – Preußischer Kulturbesitz spielt bereits seit Jahren eine wesentliche Rolle beim Auf- und Ausbau der nationalbibliographischen Verzeichnisse der deutschen Drucke des 16. bis 18. Jahrhunderts (VD16, VD17 und VD18). Die Digitalisierung dieser Werke wird über zahlreiche, meist von der Abteilung Historische Drucke betreute Projekte entscheidend vorangetrieben. Die Staatsbibliothek setzt den Fokus in besonderem Maße auf die Erschließung der Images etwa durch eine tiefgehende Strukturdatenerfassung. Seit einigen Jahren engagiert sich die Staatsbibliothek für Einrichtung und Aufbau von Virtuellen Bibliotheken bzw. Portalen wie Europeana.eu, das Zentrale Verzeichnis Digitalisierter Drucke oder die World Digital Library. Der Präsident der Stiftung Preußischer Kulturbesitz ist Vorstandssprecher der Deutschen Digitalen Bibliothek, die 30.000 deutsche Kultur- und Wissenschaftseinrichtungen vernetzen und das nationale Kulturerbe über eine gemeinsame Plattform öffentlich zugänglich machen soll.

Die Staatsbibliothek sieht eine wesentliche Zukunftsaufgabe in der Weiterentwicklung des Angebots digitalisierter Sammlungen hin zu Volltextbibliotheken. Die vorliegende Studie präsentiert die Ergebnisse eines zwölfmonatigen Pilotprojekts zur Evaluierung unterschiedlicher Software-Lösungen für Fraktur-OCR. Ich freue mich sehr, dass die Deutsche Forschungsgemeinschaft mit ihrer Förderung dieses innovative Projekt ermöglicht hat.

So konnte eine für Software-Anbieter eher periphere, für die deutschsprachigen Drucke der Frühen Neuzeit aber entscheidende Problemstellung in einem intensiven Testprozess analysiert und Möglichkeiten und Grenzen der OCR-Anwendung bei diesem Material ausgelotet werden. Als Testbestand wurde eine vom 16. bis 18. Jahrhundert im protestantischen Bereich massenhaft verbreitete Quellengattung ausgewählt, die von den unterschiedlichsten historisch orientierten Disziplinen intensiv beforscht wird und als serielle Quelle besonders hohes Potenzial für eine Volltexterschließung aufweist: Die Staatsbibliothek besitzt eine der umfangreichsten Sammlungen von Funeralschriften (Leichenpredigten und andere Trauerschriften), aus der im Projekt ein kleines Testsample mit OCR erschlossen wurde.

Die Ergebnisse der von der Abteilung Historische Drucke in Zusammenarbeit mit zwei Software-Anbietern durchgeführten komplexen Testverfahren können nun präsentiert werden. Chancen, Herausforderungen und Aufwände einer gattungsspezifischen OCR-Konversion des anspruchsvollen Materials werden umrissen. Erfreulicherweise wird die Publikation durch einen Beitrag von Thomas Stäcker ergänzt, der die einschlägigen Projekterfahrungen aus der Herzog August Bibliothek Wolfenbüttel in die Studie einbringt. Ich hoffe, dass der so entstandene umfassende Praxisbericht eine wesentliche Hilfestellung bei der Etablierung und Weiterentwicklung automatischer Texterkennungsverfahren für Frakturschriften geben kann.

Karl Werner Finger

*Ständiger Stellvertreter der Generaldirektorin der Staatsbibliothek zu Berlin PK*

# 1 Einführung

## 1.1 Projekthintergrund

Seit mehreren Jahren ist die Staatsbibliothek zu Berlin – Preußischer Kulturbesitz (SBB PK) bestrebt, ihre herausragende und überregional bedeutende Funeralschriftensammlung<sup>1</sup> zu digitalisieren. Über eine Imagedigitalisierung hinausgehend sollen auch Volltexte entstehen. Daher wurde 2008 das *Pilotprojekt zum OCR-Einsatz bei der Digitalisierung der Funeralschriften der Staatsbibliothek zu Berlin*<sup>2</sup> konzipiert. Im folgenden Jahr sagte die Deutsche Forschungsgemeinschaft eine Förderung für die Dauer von 12 Monaten zu. Von Oktober 2010 bis September 2011 setzte eine Projektkraft in der Abteilung Historische Drucke dieses Vorhaben um. Ziel war es, die Volltexterfassung mittels optischer Zeichenerkennung (OCR) an bis zu 2.000 Drucken (maximal 50.000 Seiten) aus dem Bestand der Staatsbibliothek zu testen. Diese Anzahl Funeralschriften berücksichtigt einen Ausschnitt aus der mit ca. 15.000 Drucken weltweit zweitgrößten Sammlung ihrer Art. Projektinhalt bildeten Testszenarien für zwei grundlegend verschiedene Softwarelösungen. Die Produkte wurden auf OCR-Erkennungsqualität sowie Konfigurations- und Optimierungsmöglichkeiten hin untersucht und verglichen. Die Umsetzung einzelner Aufgaben erfolgte dabei arbeitsteilig mit den Dienstleistern (zugleich auch Software-Entwicklern). Zur Handhabung der Software bestand mit diesen ein kontinuierlicher Informationsaustausch.

Im Pilotprojekt wurden folgende zwei OCR-Produkte getestet:

- B.I.T. Bureau Ingénieur Tomasi SARL Toulouse<sup>3</sup> – Software: BIT-Alpha
- Herrmann & Kraemer GmbH und Co-KG Garmisch-Partenkirchen<sup>4</sup> – Software: HK-OCR auf Basis der ABBYY FineReader Engine 9.

Durch die inhaltliche Ausrichtung des Materials auf Funeralschriften sowie die daraus resultierenden regionalen und zeitlichen Schwerpunkte ergaben sich zudem gattungsspezifische Optimierungsmöglichkeiten, handelt es sich doch überwiegend um Drucke des 17. und beginnenden 18. Jahrhunderts aus dem mitteldeutschen Raum. Beispielhaft entstand eine Wortliste, die Begriffe versammelt, welche häufig in Funeralschriften benutzt werden. Darüber hinaus integriert diese Wortbibliothek auch die Thesauri der Forschungsstelle für Personalschriften Marburg (THELO<sup>5</sup> und THEPRO<sup>6</sup>) sowie eine Liste historischer Krankheitsbezeichnungen des Vereins für Computergenealogie<sup>7</sup>.

Im Projekt wurden Möglichkeiten zur automatisierten Volltexterfassung größerer Funeralschriftenbestände getestet. Speziell für die weitere Digitalisierung der Funeralschriften der Staatsbibliothek konnten Erkenntnisse für ein Vorgehen in der Bibliothek selbst gewonnen werden.<sup>8</sup> Weitere Etappen auf dem Weg zu einer vollständigen digitalen Funeralschriftensammlung werden planbar. Eine Übertragbarkeit der Ergebnisse auf zukünftige Massendigitalisierungsprojekte aus dem Bereich der frühneuzeitlichen Drucke wird angestrebt.

Der vorliegende Werkstattbericht schildert sowohl die Ausgangslage, den methodischen Ansatz und die Ergebnisse des Softwarevergleichs als auch die erreichte OCR-Optimierung anhand des Beispielmaterials. Im Einzelnen werden neben Aussagen zur Erkennungsqualität Einschätzungen zur Binarisierung (Bildvorbereitung und Wandlung in Schwarz/Weiß) und zur Bild-, Wort- und Zeichensegmentierung getroffen. Zusätzlich bewertet werden die Nutzung von Zeichen- und Wortbibliotheken sowie vorhandene Trainings- und Validierungsmöglichkeiten. Die Dokumentation der gewählten Parameter und Einstellungen eröffnet Bibliotheken die Möglichkeit der späteren Weiterverwendung von Projektergebnissen.

Die Beschreibung von Nutzungsszenarien der untersuchten Softwareprodukte zeigt potentiellen Anwendern sowohl die Möglichkeiten als auch die generellen Probleme der OCR-Prozesse und des Ma-

<sup>1</sup> <http://staatsbibliothek-berlin.de/die-staatsbibliothek/abteilungen/historische-drucke/sammlungen/bestaende/personale-gelegenheitsschriften/> [Stand: 02/2013]

<sup>2</sup> <http://staatsbibliothek-berlin.de/die-staatsbibliothek/abteilungen/historische-drucke/aufgaben-profil/projekte/funeralschriften/> [Stand: 02/2013]

<sup>3</sup> <http://bit.dyndns.biz/> [Stand: 02/2013]

<sup>4</sup> <http://www.hk-gap.de/> [Stand: 02/2013]

<sup>5</sup> <http://www.personalschriften.de/datenbanken/thelo.html> [Stand: 02/2013]

<sup>6</sup> <http://www.personalschriften.de/datenbanken/thepro.html> [Stand: 02/2013]

<sup>7</sup> <http://wiki-de.genealogy.net/Kategorie:Krankheitsbezeichnung> [Stand: 02/2013]

<sup>8</sup> <http://www.slideshare.net/mdz-bsb/digitalisierungspraxis-federbusch-ocrpraxistest> [Stand: 02/2013] und <http://mdzblog.wordpress.com/> [Stand: 02/2013]

terials selbst. Als wesentliches Projektergebnis ist die Abhängigkeit der Softwarebeurteilung von den jeweiligen Anwendungsszenarien zu nennen. Je nach Zielstellung wird eines der beiden Softwareprodukte günstiger erscheinen. Die vorliegende vergleichende Darstellung einzelner Funktionen und Einsatzmöglichkeiten soll Interessenten dabei unterstützen, sachbezogen den Einsatz von OCR-Software zu planen. Auch wenn die Defizite in der OCR-Erkennungsqualität gegenüber Drucken des 19. Jahrhunderts unbestritten sind, ist es möglich, auch die OCR-Konversion ausgewählter frühneuzeitlicher Texte sinnvoll durchzuführen und zu optimieren.

Auch mittels OCR-Konversion entstandene Volltexte eröffnen zusätzliche Suchmöglichkeiten in digitalen Daten. Damit geht die Funktionalität von Digitalisaten über das bloße Lesen der Texte hinaus. Stäcker erläutert in seinem Beitrag in Kap. 8.6 ff die verschiedenen Vorstellungen über ausreichende Textgenauigkeiten und die daraus resultierenden Nutzungsmöglichkeiten. Neben der direkten Einbindung der entstandenen OCR-Daten bietet sich sowohl die manuelle nutzerseitige Korrektur als auch Verbesserung durch nachgelagerte halbautomatische Prozesse (z. B. hinsichtlich der Layouterkennung) sowie semantische Auszeichnung an.

## 1.2 Hilfsthesen

Um die Softwareprodukte nicht nur aufeinander zu beziehen, sondern sie auch ins Umfeld weitergehender Anwender-Erwartungen einzuordnen, wurde während der Projektarbeit von folgenden Annahmen ausgegangen:

- Wesentliche Verarbeitungsschritte der optischen Zeichenerkennung sind auch im bestmöglichen Fall regelmäßig von Informationsverlust betroffen. Ein praxisorientierter Ansatz wird sich zweckmäßigerweise auf die Betrachtung derjenigen Fehlerquellen konzentrieren, die gegenwärtig praktikabel behandelbar oder wenigstens aussichtsreich problematisierbar sind. Dennoch kann es während des Softwarevergleichs sinnvoll sein, auch auf Selbstverständlichkeiten wie vorläufig hinzunehmende Informationsverluste hinzuweisen und die jeweiligen Lösungsangebote hiernach einzuordnen.
- Erkennungsfehlern (systematischen wie unsystematischen) wird auf zwei prinzipiell unterschiedliche Arten begegnet: (a) durch Vorkehrungen zu ihrer Vermeidung und (b) durch nachträgliche Korrektur, d. h. Akzeptanz des Fehlers bei Aussicht auf Kompensation. Unabhängig davon, wann dieser Unterschied praxisrelevant ist oder nicht, ist es sinnvoll, die betrachteten Softwaremerkmale und Lösungsstrategien jeweils einem dieser Typen zuzuordnen.
- In manchen Anwendungsfällen sind Informationsverluste einiger Verarbeitungsstufen unschädlich. Auch wenn sich Anwender in der Praxis dessen bewusst sind und gerade hierauf setzen, kann eine systematische Nennung der Bedingungen, unter denen Informationsverluste harmlos werden können, die Planung von OCR-Vorhaben unterstützen.
- Auch wenn beide Hersteller ihre Software als integrierte Gesamt-OCR-Lösung konzipiert haben und sie primär sogar in eigener Regie als Dienstleister nutzen, könnten Teilschritte identifiziert werden, deren Ein- oder Ausgaben eine modularisierte Verwendung nahelegen, z. B. um ein Produkt speziell nur im Bereich seiner Stärken einzusetzen. Hierauf hinzuweisen lohnt sich auch dann, wenn die Import- bzw. Exportschnittstellen in der betrachteten Software noch nicht realisiert sind.

- Wenn aus Ressourcengründen für Massen-OCR meist allein eine vollautomatische Volltextgewinnung ohne manuelle Eingriffe als akzeptabel gesehen wird, wird dabei ein Verlust an menschlicher, bei manueller Erfassung intuitiv stets eingehender Objektwahrnehmung in Kauf genommen bzw. im besten Fall durch Korrekturverfahren kompensiert. Wenn es gelänge, durch Echtzeit-Evaluation von Teilschritt-Ergebnissen wenige effektive Eingriffspunkte für menschliche Justierungen zu erkennen und dafür eine Bedienoberfläche komfortabel auf sehr schnelle Eingriffe zu optimieren, dann könnte der bisher vorherrschende Einwand eines ökonomisch nicht tragbaren Zeitaufwands hierfür zumindest relativiert werden. Neben der Erkennungsgüte könnte dies die Baustelle sein, an der OCR-Arbeitsumgebungen sich in der nächsten Zeit vorrangig bewähren müssen.

### 1.3 Kurzfassung der Ergebnisse

Beide Programme bedienen einen ungeteilten Gesamt-OCR-Prozess von Bild-Binarisierung bis zur Wortkorrektur und sind im Prinzip für eine auf alte deutschsprachige Drucke abgestimmte Stapelverarbeitung durch Anwender geeignet. Unabhängig davon bieten beide Hersteller die OCR auch als Dienstleistung an.

Markant unterscheiden sich beide Softwareprodukte dadurch, dass BIT-Alpha eine Vielfalt detaillierter Konfigurationsmöglichkeiten bereitstellt, HK-OCR sich dagegen wegen der eingekapselten FineReader-Engine, die viele feste Voreinstellungen mitbringt, auf wenige Konfigurationsmöglichkeiten beschränkt. Es wurde gezeigt, dass folgende Möglichkeiten zur anwenderseitigen Optimierung der OCR-Erkennungsgüte bestehen:

- **Effekt von Binarisierungsoptionen:** Eine vorlagenspezifische Optimierung der Binarisierungs- und Segmentierungsparameter kann die OCR-Ergebnisse deutlich verbessern (nur BIT-Alpha; HK-OCR ist hier nicht konfigurierbar).
- **Effekt von Vorsortierung des Materials und Training eigener Muster:** Bei nach Schriftähnlichkeit vorsortierten Teilbeständen von Leichenpredigten weniger Druckorte und Jahrzehnte des 17. Jahrhunderts war ein Training speziell hierzu passender Muster sinnvoll und erbrachte deutlich bessere Ergebnisse als eine OCR mit allgemeinen Fraktur-Mustern.
- **Effekt geeigneter Wörterbücher:** Auf die Vorlage abgestimmte Lexika (gattungs- und zeitspezifische Wortformen, angereichert durch passende Titel- und Personendaten aus dem Verbundkatalog) erhöhen die Erkennungsgüte deutlich, wenn auch aufgrund verschiedener Mechanismen (in BIT-Alpha durch Nachkorrektur, in HK-OCR integriert in die Zeichenerkennung).

Für auf Seitenebene mehr oder weniger homogene Schriftbilder sind diese deutlichen Steigerungseffekte unmittelbar nutzbar (beide Programme). Für den häufigen Fall von Seiten ohne einheitliche Schriftart (Wechsel von Fraktur zu Antiqua, Wechsel von Sprachen, starke Größenunterschiede) ist gegenwärtig keines der beiden Produkte in der Lage, die Stärken einer selektiven Musterzuweisung auch auf Teilregionen von Seiten anzuwenden. Eine Kombination separat trainierter Musterbestände ist in BIT-Alpha (bzw. dem Museditor BIT-Knowledge) begrenzt, in HK-OCR/FineReader bisher nicht möglich; HK-OCR bietet stattdessen die Wahl, eine Anwendermustermenge mit den vorgegebenen FineReader-Mustern zu kombinieren. Analog muss zur lexikalischen Korrektur für jeden Batchlauf ein Wörterbuch bzw. eine Sprachengruppe (so der Begriff in der HK-OCR-Oberfläche) festgelegt werden; sprachlich gemischte Seiten oder Seitenstapel zwingen daher zur Verwendung eines entsprechend kombinierten Wörterbuchs mit den damit verbundenen Kompromissen in der Korrektursicherheit.

Der Export erfolgt jeweils mindestens in mit Textkoordinaten angereicherten XML-Formaten (BIT-Alpha: ALTO, HK-OCR: FineReader-XML) sowie einigen anderen Formaten (BIT-Alpha: PDF, HTML, Text; HK-OCR: PDF, RTF und ein gegenüber dem einzelzeichenorientierten FineReader-XML deutlich verschlanktes, wortorientiertes XML). Bisher von keinem der beiden Programme geleistet wird eine echte Strukturerkennung (angefragt waren formal gliedernde Elemente z. B. Kolumnentitel, Kustoden, Bogensignaturen u. a.) oder die Auszeichnung von *Named Entities* (Personen, Orte, Berufe, Krankheiten) aufgrund vorgegebener Listen oder von Datumsangaben, Bibelstellen etc.<sup>9</sup> Entsprechend ist eine Ausgabe in semantisch orientierten Formaten wie TEI allenfalls oberflächlich ohne die genannten textspezifischen Zusatzinformationen aus den Exporten ableitbar.

Jede Software hat Stärken und Schwächen bei der Unterstützung eines zügigen und effektiven Vorgehens besonders bei Training und Korrektur/Validierung. Diese Befunde wurden den Herstellern jeweils mitgeteilt, zumal manche Schwächen leicht behebbar erscheinen und möglicherweise einfach auf bisher unzureichenden eigenen Tests der Hersteller beruhen. Beide Firmen haben zugesagt, die Anmerkungen der SBB PK bei der Entwicklung der nächsten Versionen zu berücksichtigen.

Innerhalb dieser Grenzen können beide Programme von selbst trainierenden oder vorhandene Trainingsergebnisse nachnutzenden Anwendern sinnvoll eingesetzt werden. Im Kapitel 6.5 wird versucht, die Vor- und Nachteile jeder Software bezogen auf verschiedene Einsatzszenarien zu nennen. Die dabei beschriebenen Kriterien können zugleich als Hilfe zur Prüfung anderer OCR-Lösungen mit herangezogen werden. Als Nebenergebnis der Software-Tests konnten somit auch Anforderungen an Softwarekomponenten zur OCR Alter Drucke formuliert werden, die hiermit auch potentiellen Anwendern mitgeteilt werden, die nicht zuletzt selbst mit entsprechender Nachfrage die Entwicklungsrichtung der nächsten Versionen mitbestimmen werden.

Nicht eingehend untersucht wurden Open-Source-Lösungen<sup>10</sup>. Ein Vergleich mit diesen könnte weitere Optionen erschließen.

<sup>9</sup> Eine kleine Teilmenge aus dem „Metadatenmodell Funeralschriften (2008)“:

[http://staatsbibliothek-berlin.de/fileadmin/user\\_upload/zentrale\\_Seiten/historische\\_drucke/pdf/Metadaten\\_Funeralschriften.pdf](http://staatsbibliothek-berlin.de/fileadmin/user_upload/zentrale_Seiten/historische_drucke/pdf/Metadaten_Funeralschriften.pdf) [Stand: 02/2013]

<sup>10</sup> Ohne weitere Konfigurationen durchgeführte Tests an acht Seiten einer im Projektmaterial befindlichen Funeralschrift mit der Software OCRopus wurden im Mai 2012 durchgeführt. Das Ergebnis ist durchwachsen, fällt jedoch gegenüber den im Projekt mit optimierten Einstellungen der getesteten Programme erkannten Seiten klar ab.

## 2 Materialspezifik der Funeralschriften

### 2.1 Gattungsspezifische Aspekte

Funeralschriften sind personale Gelegenheitsschriften, deren Anfänge in die Mitte des 16. Jahrhunderts hineinreichen. Starke Verbreitung fand der Brauch, Leichenpredigten zu drucken, vor allem im protestantischen Deutschland des 17. und beginnenden 18. Jahrhunderts. In dieser Zeit, die stark durch die Theologie des Pietismus geprägt wurde, war es üblich, besonders Adelige und Akademiker, aber auch Vertreter des wohlhabenden Bürgertums nach ihrem Tode mit diesen Schriften ausführlich zu würdigen. Funeralschriften bestehen in der Regel neben der Leichenpredigt aus dem Lebenslauf der Verstorbenen. Darüber hinaus enthalten sie oft eine Abdankungsrede (Leichabdankung), ein Programm Academicum, ein Epitaph und verschiedene Epicedien (Trostbriefe, Gedichte, Lieder, Musikbeigaben Verwandter und Bekannter). Als weitere druckgraphische Ausgestaltung findet man Porträts, Wappen, Ahnentafeln, Embleme, aber auch Darstellungen des *Castrum doloris* (hölzerner geschmückter Aufbau zur Ehre des Toten) und des Leichenzugs.<sup>11</sup>

Diese Gebrauchsschriften, deren Ziel es war, neben der Würdigung der Verstorbenen vor allem den Hinterbliebenen Trost zu spenden, sind seit mehreren Jahrzehnten Gegenstand multidisziplinärer Forschung. Mit ihren vielfältigen Informationen bieten sie Einblicke in die damalige Alltagskultur und stellen wichtige Quellen für weitergehende Forschungen dar. Scheibe benennt Forschungsinteressen: „Neben ihrer Bedeutung für Literaturgeschichte und Predigtforschung bieten die Leichenpredigten dem Sozial-, Kultur- und Universitätshistoriker eine Fülle statistisch auswertbarer Daten, dem Kunsthistoriker und Heraldiker reiches Material zur Ikonographie und Emblemik, dem Musikwissenschaftler eine ganze Reihe von nur hier überlieferten Trauerkompositionen; die teilweise sogar mit Abbildungen versehenen Krankheitsberichte schließlich liefern dem Medizinhistoriker wertvolle Informationen. Darüber hinaus aber stellen die Personalschriften ein einzigartiges prosopographisches Nachweisinst-



1 Komposition als Epicedium. Ee 705-1592 S. [85]

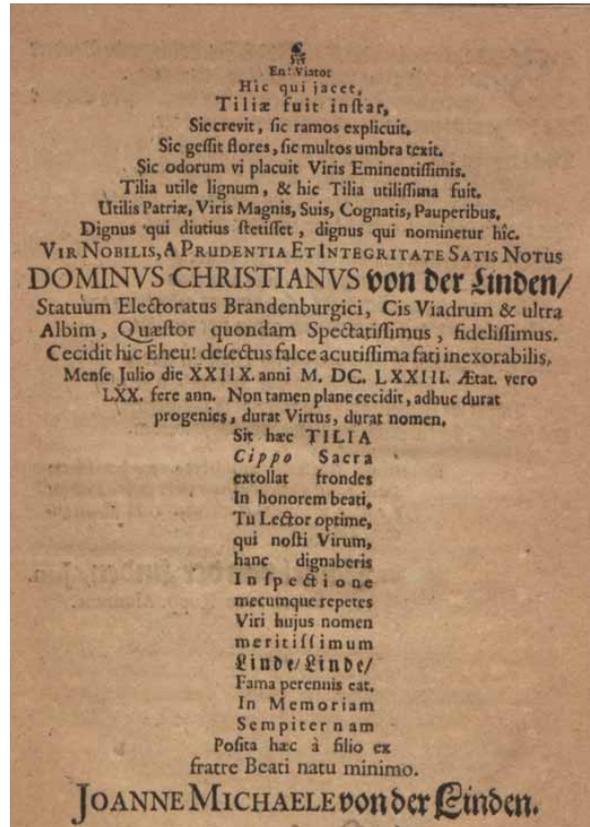


2 Wappen. Ee 708-85 S. 2

<sup>11</sup> Eine ausführliche Beschreibung der Geschichte und des Aufbaus von Leichenpredigten ist zu finden auf der WWW-Seite der Forschungsstelle für Personalschriften Marburg: <http://www.personalschriften.de/leichenpredigten/aufbau.html> [Stand: 02/2013]



3 Frontispiz. Ee 705-28 S. [4]



4 Figuredicht. 3/4 in: Ee 700-1993 S. [16]

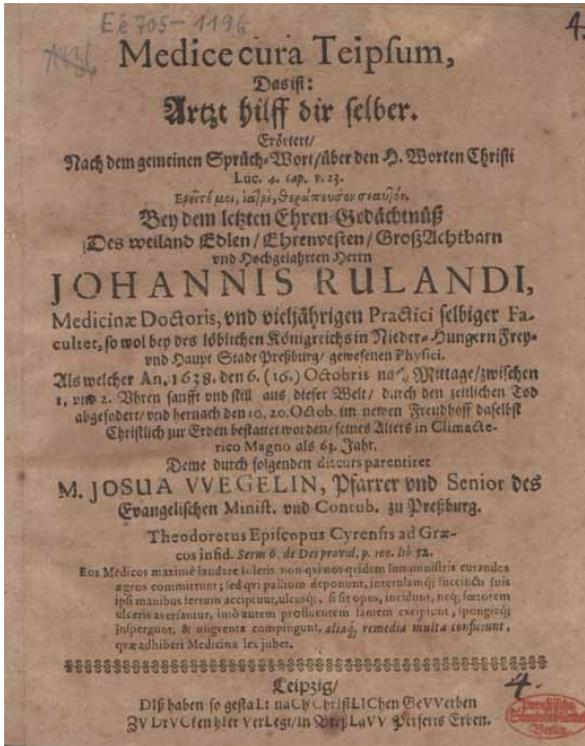
rumment für Genealogie und Biographik dar – eine Qualität, die bereits im frühen 20. Jahrhundert ihre besondere Wertschätzung begründete.<sup>12</sup> Derzeit stellen Volltexte von Funeralschriften ein Desiderat dar. Ein semantisch ausgezeichneter Volltext könnte die oben genannten Forschungsgebiete sinnvoll befördern; statistische Analysen sowie semantische Netze zu extrahierender Personen rücken in greifbare Nähe.

## 2.2 Sammlungsspezifische Aspekte

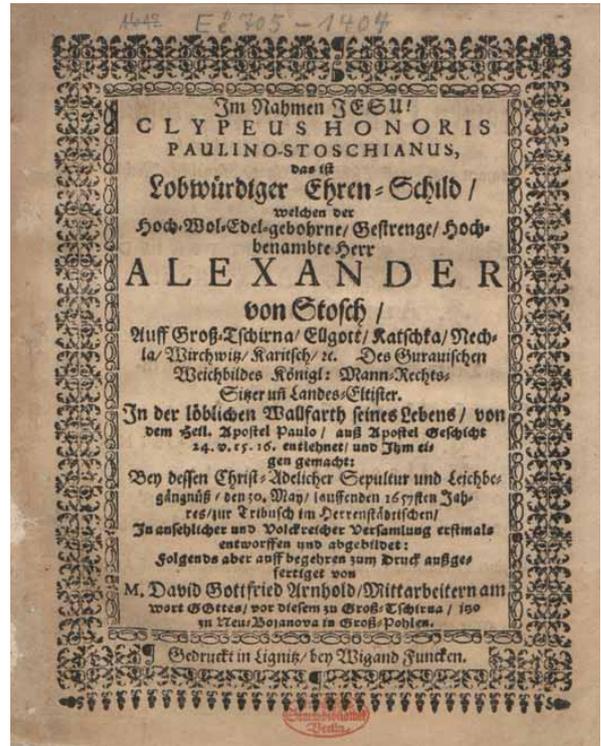
Im historischen Druckschriftenbestand der Staatsbibliothek zu Berlin befinden sich von jeher eine große Anzahl Funeralschriften. In den 30er Jahren des 20. Jahrhunderts gelang es der Bibliothek, vier geschlossene Leichenpredigtensammlungen zu erwerben, so dass heute in der Summe von 15.000 Schriften gesprochen werden kann. Der Bestand wurde um die Leichenpredigten aus der Bibliothek des Bories Freiherrn von Münchhausen auf Windischleuba (bei Altenburg/Thüringen) sowie um die Wernigeroder und die Bückeburger Sammlung vermehrt. Für das durchgeführte OCR-Projekt war vor allem eine vierte vom Antiquariat J. A. Stargardt (Berlin) 1936 vermittelte Leichenpredigtensammlung (Signatur Ee 705) von Interesse. Alle weiteren Aussagen beziehen sich vor allem auf diese Sammlung. Den Schwerpunkt bilden Funeralschriften (1.700 Nummern), die im mitteldeutschen Raum erschienen sind. Dies belegen Erscheinungsorte wie Halle, Magdeburg, Erfurt, Altenburg, Weimar, Wittenberg, Zeitz, Merseburg, Zerbst, Jena, Gotha, Leipzig, Dresden, Zwickau, Hannover, Braunschweig, Helmstedt, Coburg, Hof, Marburg, Kassel, Hanau, Giessen, Schweinfurt und Frankfurt/Main. Hinzu treten ostdeutsche Orte wie Elbing, Küstrin, Königsberg, Stettin, Breslau, Lignitz, Schweidnitz, Groß Glogau

<sup>12</sup> Scheibe, Michaela zur Sammlung Personaler Gelegenheitschriften auf: <http://staatsbibliothek-berlin.de/die-staatsbibliothek/abteilungen/historische-drucke/sammlungen/bestaende/personale-gelegenheitschriften/> [Stand: 02/2013]

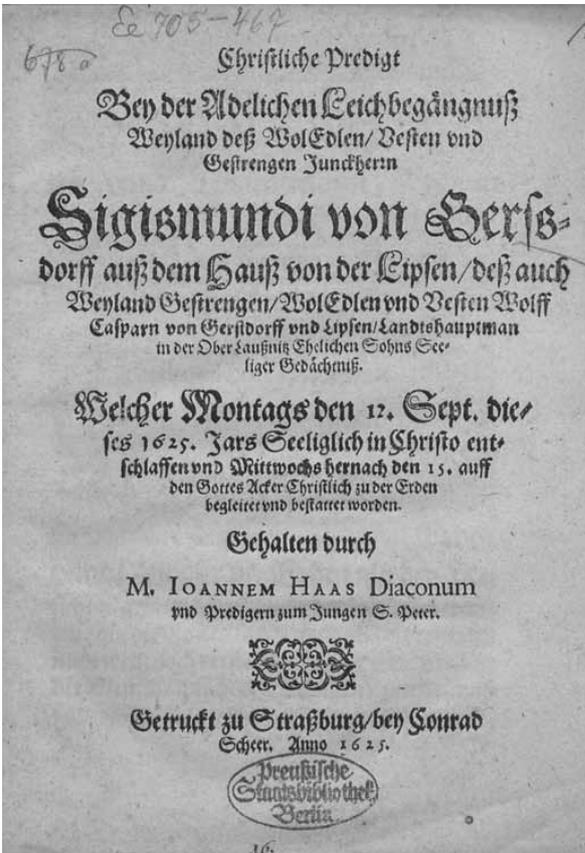
und Brieg. Durch die Verbreitung des Pietismus im süddeutschen Raum (z. B. Württemberger Pietismus) gibt es auch eine nicht geringe Anzahl von Predigten aus diesem Gebiet. Straßburg, Heidelberg, Ulm, Stuttgart und Tübingen sind hier zu nennen. Das breite regionale Spektrum dieser Sammlung belegen ferner Erscheinungsorte wie Berlin, Hamburg, Rostock und Basel.



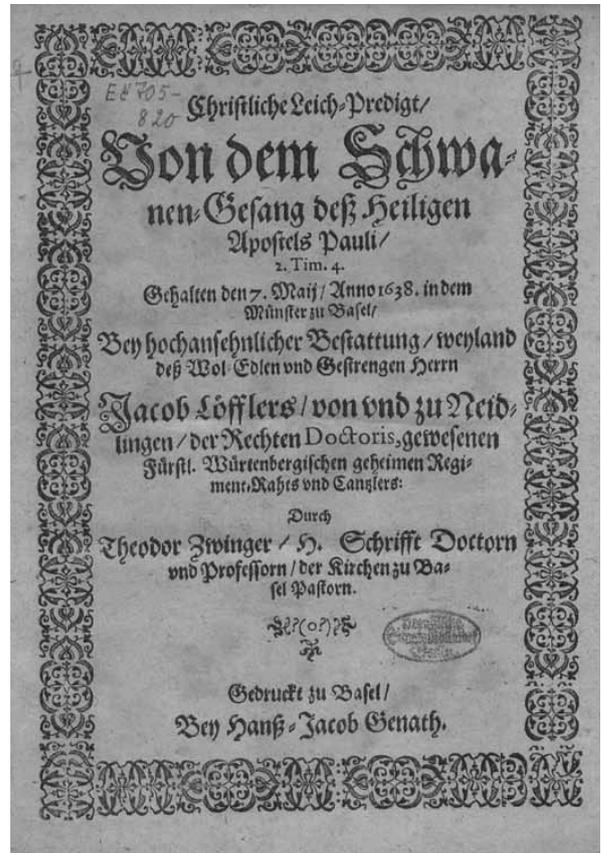
5 Leichenpredigt auf einen Arzt. Leipzig. Ee 705-1196 Titelblatt



6 Leichenpredigt. Lignitz. Ee 705-1404 Titelblatt



7 Leichenpredigt. Straßburg. Ee 705-467 Titelblatt



8 Leichenpredigt. Basel. Ee 705-820 Titelblatt

Die zeitliche Abdeckung der genannten Sammlung spiegelt die Blütezeit der Leichenpredigt<sup>13</sup> wider. So stehen zwei Drittel Schriften des 17. Jahrhunderts einem Drittel des 18. Jahrhunderts gegenüber. Dreißig Funeralschriften datieren ins 16. Jahrhundert. Aufgrund einfacher Auswertungen der Volltexte unter Hinzuziehung des Thesaurus Professionum der Forschungsstelle für Personalschriften Marburg<sup>14</sup> konnten folgende Personenkreise (in originaler Schreibung) ermittelt werden, die in den Schriften Erwähnung finden:

- Prediger / Pfarrer / Pastor / Seelsorger / Bischoff / Beichtvater / Theologus / Diaconus
- Lehrer / Rector / Praeceptor
- General / Hauptmann / Obrister / Leutenant / Soldat / Oberster / Wachmeister / Kriegsmann
- Arzt / Medicus / Medico
- Poet / Mahler
- Rath / Cantzler / Bürgermeister / Amptman / Senator
- Buchhalter / Weber / Zöllner / Pfleger / Wirth

Diese Berufsbezeichnungen beziehen sich nicht nur auf die Verstorbenen, sondern auch auf die Angehörigen und Freunde, sofern sie in der Funeralschrift erwähnt wurden. Ein Abgleich der Liste historischer Krankheitsbezeichnungen des Vereins für Computergenealogie<sup>15</sup> mit den Texten hinsichtlich der genannten Leiden und Todesursachen vermittelt einen ersten Eindruck in die Lebensumstände der Zeit.

*Häufigste Krankheiten:*

- Hitze
- Fieber
- Schlag
- Biß
- Brand
- Flecken
- Pest
- Husten

## 2.3 Typographische und buchgestalterische Aspekte

Neben inhaltlichen Besonderheiten weisen die Funeralschriften einige typographische und buchhistorische Gattungsspezifika auf. Wie andere Druckschriften der jeweiligen Zeit auch zeichnet sich die Mehrzahl der Drucke aus dem 17. Jahrhundert durch ein sehr gedrängtes Druckbild mit übervoll gesetzten Seiten aus. Um Papier zu sparen, wurden zudem anfangs kleine Schriften (12 Pkt) verwendet; der Satzspiegel weist fast keine Absätze im Text auf (vgl. Abb. 13). Üblicherweise gibt es Titelblätter, die in bis zu drei Schriftarten und zusammen zehn verschiedenen Schriftgraden in barocker und ausschweifender Weise den Titel nennen (vgl. Abb. 5-8 u. 10-12). Üblich sind hier zunächst stereotype Titelanfänge wie „Christliche Leichpredigt/ Bey ...“. Wesentlich direkter und variantenreicher erscheinen sie später (z. B. „Tröstende Uranie bey einem zu früh verwelckten Orangen-Zweige, Die ...“ Weimar 1718).

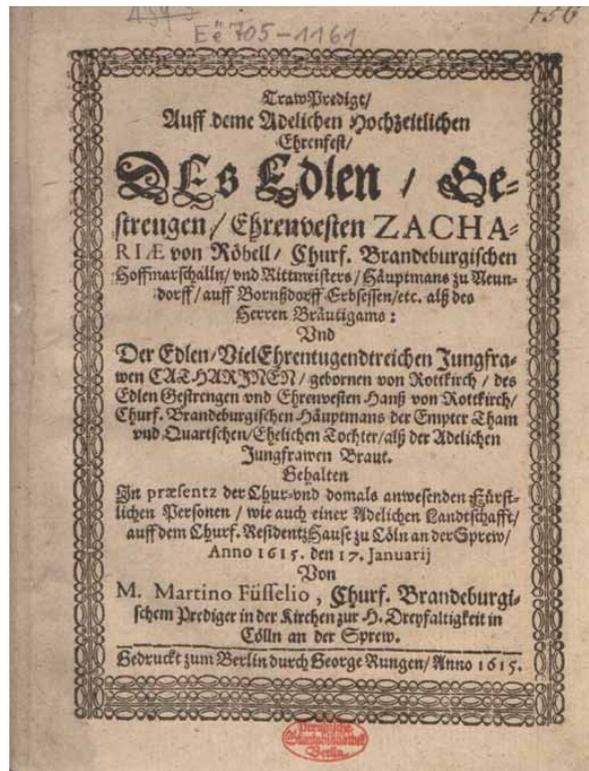
<sup>13</sup> Als Blütezeit gilt die Zeit vor dem 30jährigen Krieg und die zweite Hälfte des 17. Jahrhunderts.

<sup>14</sup> <http://www.personalschriften.de/datenbanken/thepro.html> [Stand: 02/2013]

<sup>15</sup> <http://wiki-de.genealogy.net/Kategorie:Krankheitsbezeichnung> [Stand: 02/2013]



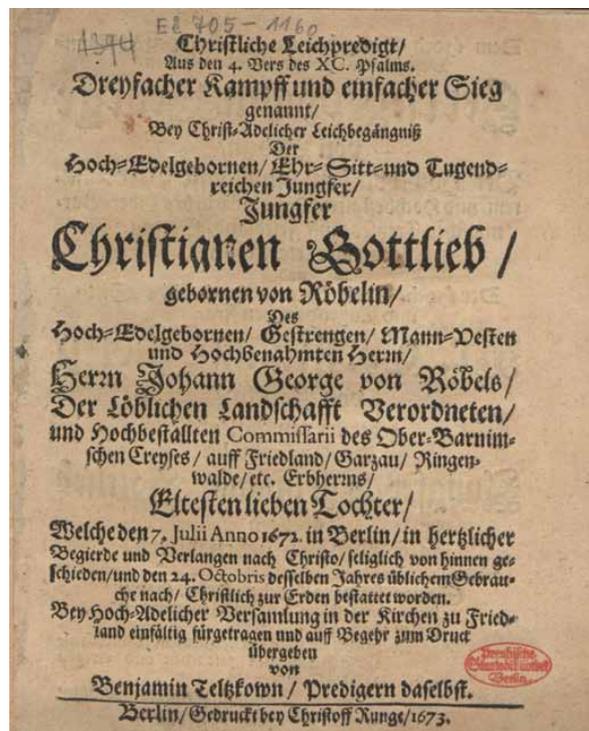
9 Kupfertitel - Hamburg, 1687. Ee 705-1472



10 Titelblatt mit Titelfassung - Berlin, 1615. Ee 705-1161

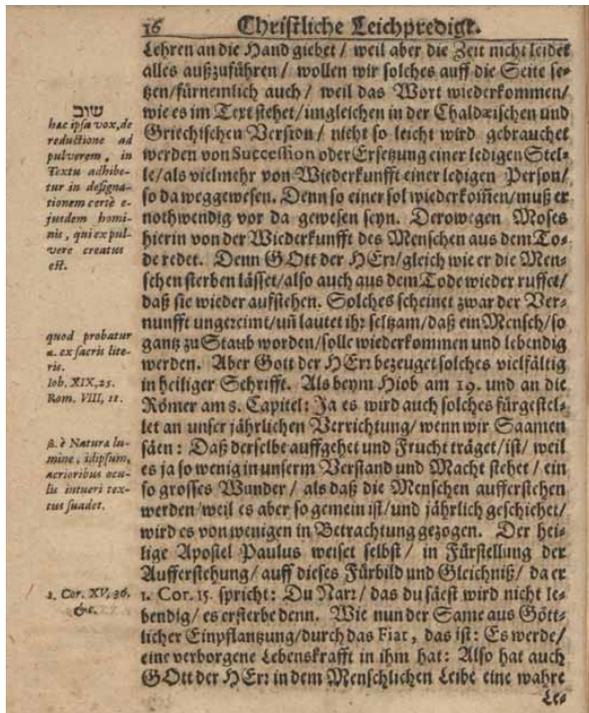


11 Illustriertes Titelblatt 1. Hälfte 17. Jh. Hamburg, 1610. Ee 705-189

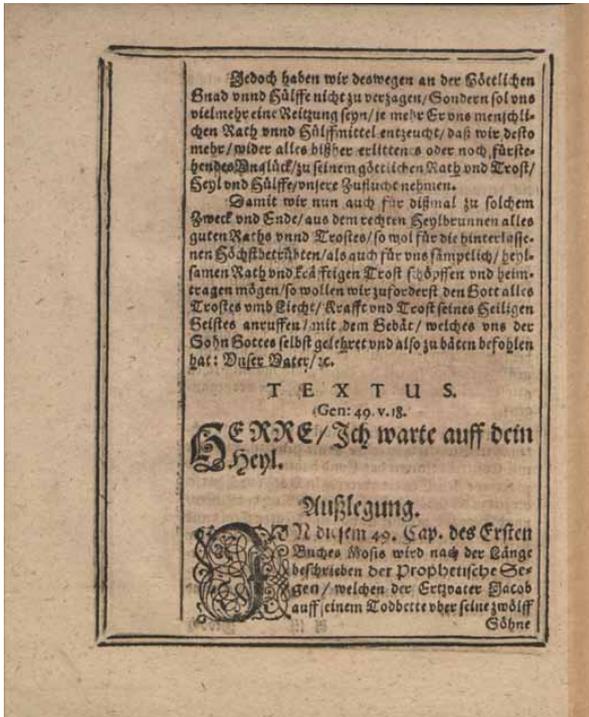


12 Titelblatt 2. Hälfte 17. Jh. Berlin, 1673. Ee 705-1160

Wie bis ins 18. Jahrhundert hinein üblich verwendete man auf den Titelblättern oft Titeleinfassungen (vgl. Abb. 6, 8 u. 10), die zunächst aus ornamentalen Motiven über Holzschnitte, später vermehrt über Kupferstiche realisiert wurden. Auch Kupfertitel und Frontispize sind zu finden (vgl. Abb. 3 u. 9). Die Rückseite des Titelblatts enthält in der Regel eine Widmung, zumeist den Angehörigen und Hauptleidtragenden des Verstorbenen zugeeignet; Vorreden können sich anschließen und leiten zur Predigt über.



13 Text mit Marginalien und Kolummentitel. Berlin, 1673. Ee 705-1160 S. 18



14 Text erscheint in Umrahmung ohne Kolummentitel. Berlin, 1645. Ee 705-669 S. 181

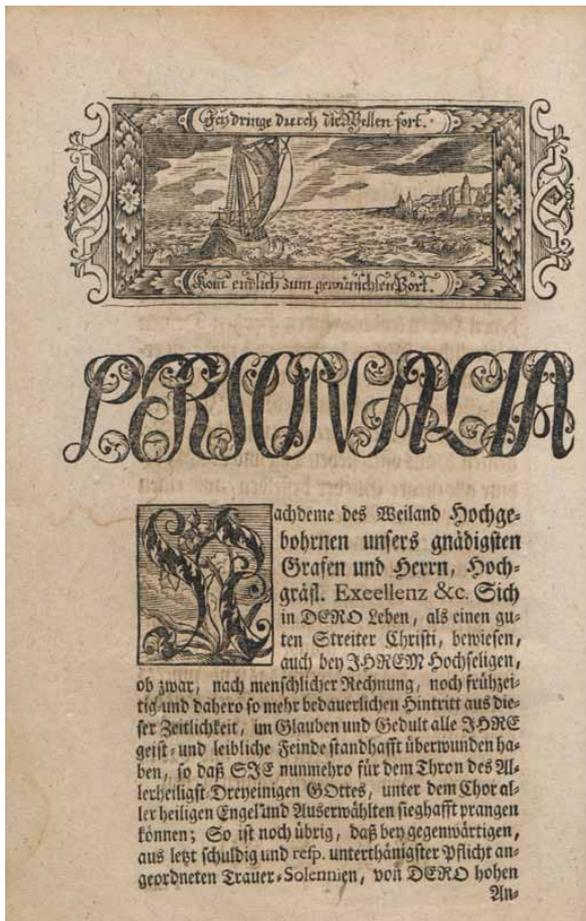
Den Hauptteil einer Funeralschrift bildet die Leichpredigt, die zunächst den Leichtext vorstellt und in den folgenden Teilen eine Vielzahl Bibelzitate enthält, welche in Form von Marginalien am Rand wiederholt werden (vgl. Abb. 13). Partiiell direkt bezeichnet unterscheidet die Predigt Exordium (Vorrede, Einleitung in Redesituation), Propositio (Einleitung des Redegegenstands) und Consolatio (Trostworte). Das Ehrengedächtnis stellt üblicherweise die Schilderung des Lebenslaufs dar. Häufig wird an dessen Ende eine Schlussvignette als Kupferstich gesetzt, die Symbole der Vanitas zeigt (Schädel, Urnen; vgl. Abb. 20). Zumeist finden sich Epicedien gesammelt am Ende der Funeralschrift. Sind diese Beiträge bibliographisch selbständig erschienen, weisen sie ein eigenes Titelblatt auf. Oftmals sind diese Seiten ungezählt und fallen durch eine typographisch abweichende Gestaltung auf. Vor allem verschiedene Gedichte werden nach Nennung des Beitragere hintereinander, zuweilen durch horizontale Linien voneinander getrennt, aufgeführt. An deren Ende erscheint eine weitere Schlussvignette.

Bis zur Mitte des 17. Jahrhunderts wurde häufig der gesamte Text doppelt gerahmt (Randleisten im Inneren des Drucks; vgl. Abb. 14). Anfangs waren Blätter und Seiten oft ungezählt, später erscheinen die Drucke des 17. Jahrhunderts mit Kopfzeilen (Kolummentitel und Paginierung) in diesem Rahmen. Ab der zweiten Jahrhunderthälfte findet man durch horizontale Linien abgeteilte Kopfzeilen.

Sind die meisten Funeralschriften anfangs eher selten und gleichförmig geschmückt (Abbildung von Sarg und Schädel, wiederholte Verwendung gleicher Schlussvignetten mit Motiven<sup>16</sup>, die in der Ikonographie vergangener Jahrhunderte in der Heraldik eine Rolle spielten; vgl. Abb. 19 u. 21), so zeigen sie ab 1650 deutlich mehr Buchschmuck mit vielfältigen Symbolen der Vanitas (Putti pusten Seifenblasen, qualmende Töpfe, ausgeblasene Kerzen; vgl. Abb. 15-18). Dieser über

Kupferstiche realisierte Buchschmuck der Darstellung von Vergänglichkeit nennt des öfteren auch die Initialen des Künstlers. Tragen die Darstellungen im allgemeinen bis zur Mitte des 17. Jahrhunderts

<sup>16</sup> Beispielfhaft soll der Wilde Mann genannt werden, dessen Darstellung auch als Fabelwesen (Halb Mensch, halb Tier) gedeutet werden kann und häufig in den Schlussvignetten des Druckers Runge vorkommt.



15 Gleitendes Schiff als Vanitassymbol. Oehringen, 1737. 1 in: Sh 6986



16 Putto mit qualmender Urne als Vanitassymbol. Dresden, 1713. Ee 705-6 S. 3

eher den Charakter des Buchschmucks, so muss man in der Folgezeit (beispielsweise bei Porträts) zunehmend von Illustrationen sprechen, die einen eigenständigen künstlerischen Beitrag darstellen. Sie wurden individuell für spezielle Drucke angefertigt.

Erscheinen Anfang des 17. Jahrhunderts Funeralschriften geringen Umfangs vorwiegend im Oktavformat, so ist ab 1650 und für das 18. Jahrhundert das Quartformat kennzeichnend. Zudem nehmen die Schriften an Umfang zu, eine vermehrte Spationierung (Vergrößerung horizontaler Zeichenabstände) und Gliederung der Texte ist zu beobachten. Auch größere Schriften kommen zum Einsatz.<sup>17</sup>

Die Entwicklung des Druckbilds, die Verwendung von Buchschmuck und andere typographische Besonderheiten im 17. Jahrhundert wurden vor allem anhand der Funeralschriften der Familie Runge (Georg, Christoph und ihre Witwen) untersucht. Ihre Drucke sind zahlenmäßig stark und zeitlich kontinuierlich im Bestand der SBB vorhanden.

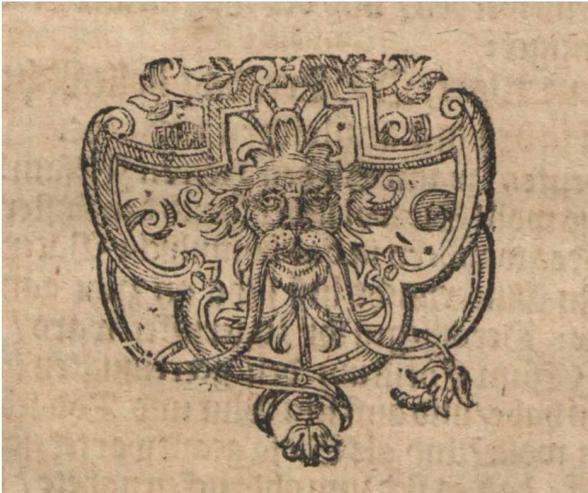


17 Vielfältige Vanitassymbole. Gotha, 1701. Ee 705-44 S. [36]



18 Putto mit Schädel u. Seifenblasen als Vanitassymbol. Leipzig, 1705. Ee 705-21 S. 3

<sup>17</sup> Aussagen zu Buchgewerbe und Druckqualität im 17. Jahrhundert finden sich unter vielfältiger Quellenangabe in KILLIUS 1999, S. 90 ff.



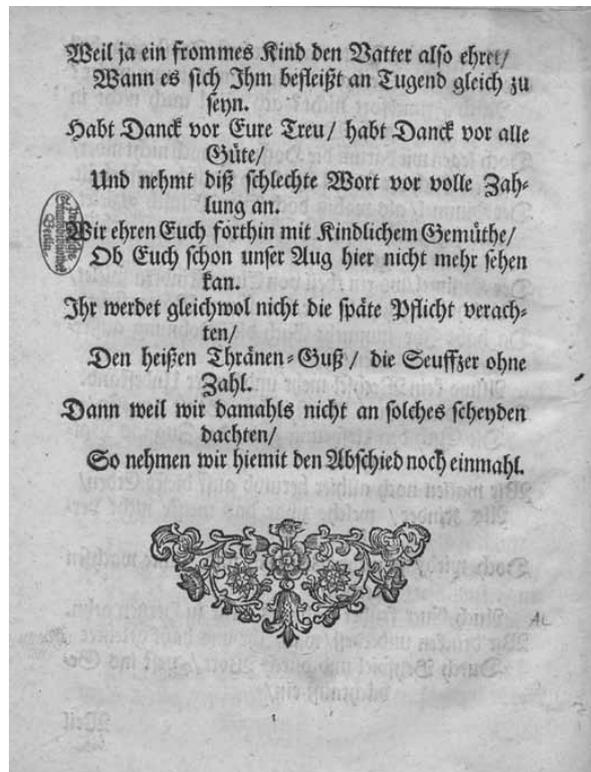
19 Schlussvignette mit Darst. „wilder Mann“. Berlin, 1664. Ee 705-1592 S. 78



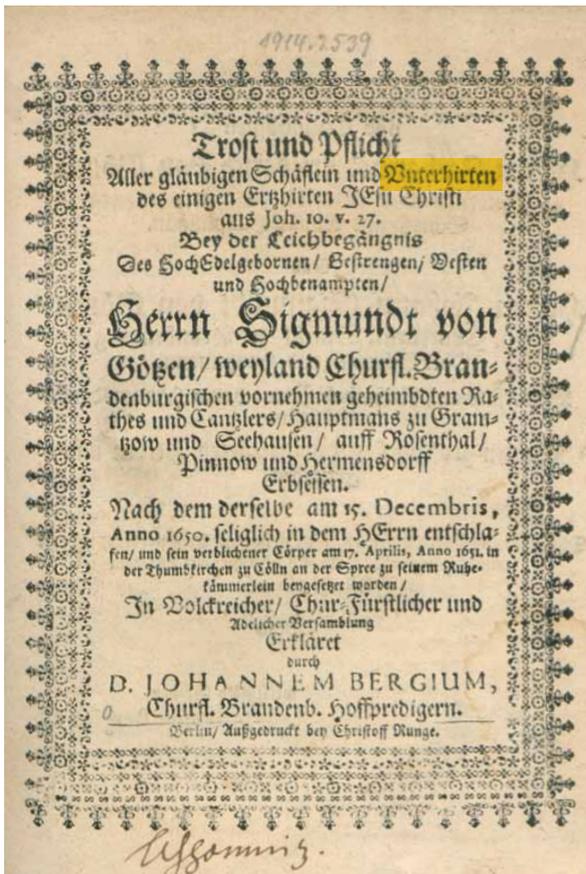
20 Schädel als Vanitassymbol. Berlin, 1674. Ee 705-1159 S. 40



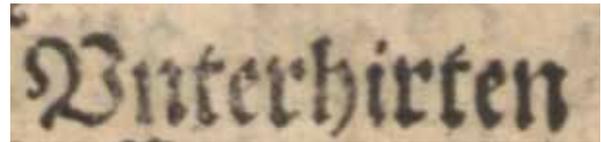
21 Sarg mit Schädel als Vanitassymbol. Berlin, 1666. Ee 705-1590 S. [150]



22 Schlussvignette mit Blumendarstellung. Straßburg, 1699. Ee 705-159 S. [32]



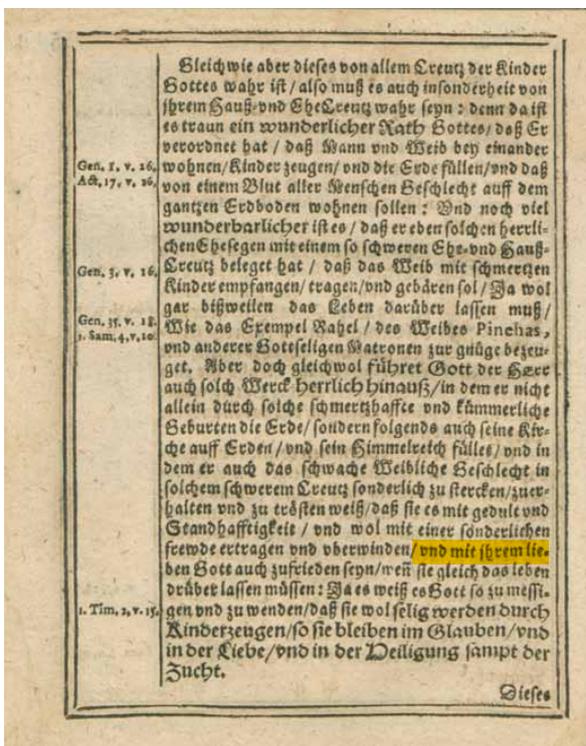
23 Verwendung von „V“ als „U“. Berlin, 1651. 11 in: Ec 1485 Titelblatt



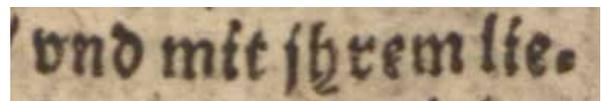
23 Markierter Ausschnitt

## 2.4 Orthographische Besonderheiten

Besonderheiten in der deutschen Rechtschreibung, die ohnehin für das 17. und 18. Jahrhundert nicht einheitlich war und starken Abweichungen unterlag, müssen hervorgehoben werden, auch wenn sie ein allgemeines Phänomen darstellen und nicht auf Funeralschriften beschränkt sind.<sup>18</sup> Bekanntermaßen findet man in den Drucken der ersten Hälfte des 17. Jahrhunderts überwiegend noch keine Unterscheidung von u und v. So lässt sich im Satz feststellen, dass ein gesprochenes „u“ am Wortanfang als „v“ gesetzt wird, im Wortinnern jedoch als „u“. Am häufigsten fällt diese Verwendung beim Wörtchen „vnd“ auf. Auch bei Großbuchstaben wird diese Schreibregel angewandt. Zu nennen wären hier Wörter wie „Vnser“ oder „Vnterthan“. Differenziert man zunächst in



24 Verwendung von „j“ in „jhr““. Berlin, 1633. an: 50 MA 42058 S. 511



24 Markierter Ausschnitt

der zweiten Jahrhunderthälfte die Kleinbuchstaben „u“ und „v“ konsequent lautbezogen, so dauert es bei den Großbuchstaben einige Jahrzehnte länger. Es finden sich reichlich Vorlagen mit dem Gebrauch der alten Regel für Großbuchstaben; parallel werden im Text die Kleinbuchstaben jedoch schon eindeutig verwendet (vgl. Abb. 23).

In diesem Zusammenhang ist auf die ähnliche Problematik bei i/j zu verweisen, wobei hier das große „J“ durchgängig für I oder J in Frakturtexten steht. Hingegen kommt es bei den Kleinbuchstaben längere Zeit zu einem gemeinsamen Gebrauch von „i“ und „j“ am Wortanfang. Auf einer Seite von 1625 beispielsweise findet man in derselben Schrift Worte wie „jhrer“, „jhn“ und „in“, „im“, „inbrünstiglich“, „ist“ gemeinsam, was auf eine Unterscheidung im Gebrauch des langen und kur-

<sup>18</sup> Allen Drucken gleich ist die Darstellung der Umlaute durch Vokale mit übersetztem kleinen „e“. Typisch für lange Zeit ist ebenso die Gliederung der Texte durch Virgel; Kommata kommen erst später in der Frakturschrift in Gebrauch. Weitere Merkmale historischer Drucke sind Kustoden und Bogensignaturen, Hilfsmittel für den Buchbinder. Vgl. weiterführend: SCHNEIDER und ZIPPEL 2009. Zur Varianz der deutschen Orthografie siehe VOESTE 2008.



27 Titelblatt mit verschörkelten Versalien, Regensburg, 1688. Ee 705-1081



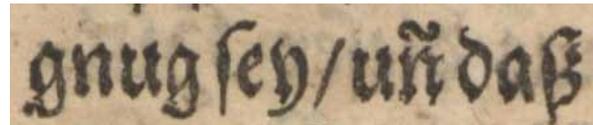
28 Titelblatt mit illustrierten Versalien, Delitzsch, 1716. Ee 705-167

zen i in der Schriftsprache hindeutet und so über mehrere Jahrzehnte zu beobachten ist bis auch hier immer „i“ geschrieben wird, wenn es sich um ein solches handelt (vgl. Abb. 24)<sup>19</sup>.

Weitere Auffälligkeiten entstehen durch die Verwendung des runden „r“<sup>20</sup>. Zunächst verwendet für die Abkürzung „etc.“, erscheint ein rundes „r“ gefolgt von einem „c“ und einem Punkt. Später verbreitet sich der Satz des runden „r“ bei Buchstabenverdopplung (Herr, herrlich, Herrn) oder auch im Zusammentreffen mit „t“ (aufgefordert; vgl. auch Abb. 25). Zusätzlich uneindeutig können Abbrüviaturen (zumeist in Gestalt von Nasalstrichen) über dem Buchstaben „n“ oder „m“ sein. Überwiegend zeigen diese eine Verdopplung an; es gibt aber auch andere Anwendungen wie bei



25 Rundes „R“ bei Buchstabenverdopplung u. in „etc.“. 7 in: Ee 1593-2 Titelblatt



26 Nasalstrich über „n“ zeigt Verdopplung an. 5 in: Ee 1593-2 S. 72

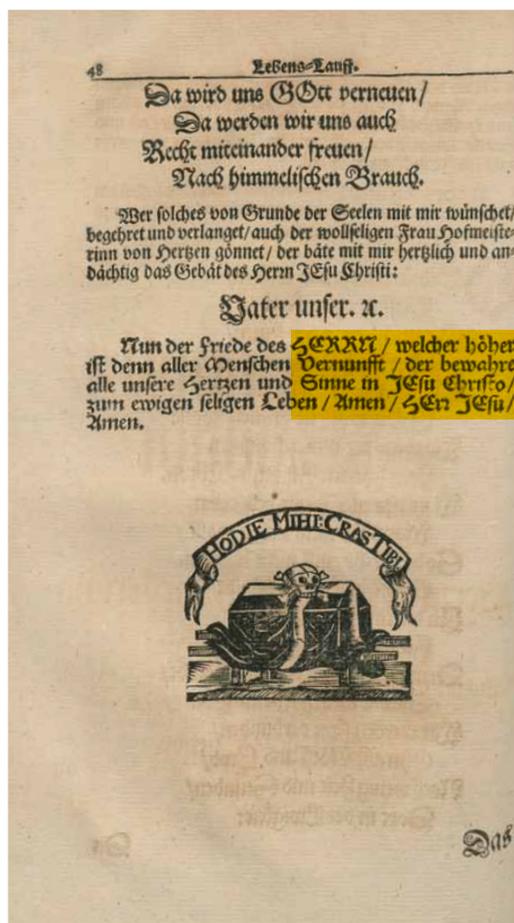
„un“; welches zu „und“ aufzulösen ist (vgl. Abb. 26). In lateinischen Zitaten sind natürlich weitere Abbrüviaturen anzutreffen, deren Auflösung wort- oder wortstammbezogen vorzunehmen ist. Für den geübten Leser ist das meist kein Problem; für die Anwendung automatischer Verfahren der Texterkennung tritt es erschwerend hinzu.

Auch die Groß- und Kleinschreibung soll Erwähnung finden, da sie für die Schriftenvielfalt eine bedeutende Rolle spielt. So schreibt Kapr sehr treffend: „Leider blieb das schnörkelhafte Wesen der barocken Schriftformen nicht ohne Einfluss auf die deutsche Rechtschreibung. Der Wunsch zur Auszeichnung drängte die Schreiber dazu, möglichst viele Wörter mit Großbuchstaben zu schreiben, die mehr Gelegenheit zu schwungvollen Variationen gaben. ... Erst begann man GOTT und HERR groß zu schreiben, dann folgten Papst, Kaiser, Bischof, Kirche, Fürst, und die Ergebnheitsfloskeln in den Vorreden der Buchausgaben des Barocks taten das übrige zur Überhäufung des deutschen Schriftbildes mit Versalien.“<sup>21</sup>

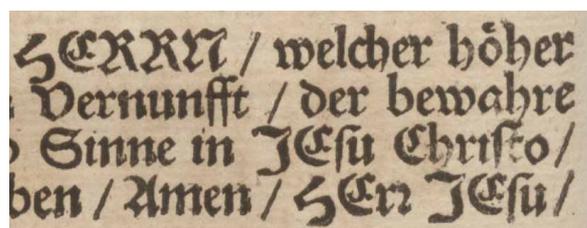
<sup>19</sup> Vgl. weiterführend HORN 1894

<sup>20</sup> „Rundes r; entstanden aus der rechten Hälfte des Unzial-R. Es wurde hinter runde Buchstaben oder als zweites r gesetzt; verlor aber im 19. Jh. an Bedeutung und verschwand.“ REHSE 1998 S. 78

<sup>21</sup> KAPR 1959 S. 61



29 Häufung von Versalien. Berlin, 1683. 3 in 4° Ee 660 S. 48



29 Markierter Ausschnitt

heraus extrahieren

30 Ein in Antiqua gesetztes Fremdwort hat Fraktur-Suffix

So erkennt man an der Entwicklung der Titelblätter der Funeralschriften neben der Neigung zur Verwendung möglichst vieler Schriften eine Zunahme verschnörkelter Versalien bis hin zur Verwendung von Barockinitialen (vgl. Abb. 27 u. 28). Offenbar war es auch im Text nicht ausreichend, Wörter wie GOTT und HERR oder JESU durch Großschreibung hervorzuheben; meist waren sie auch noch in einer größeren Schrift gesetzt (vgl. Abb. 29). Auch der Satz von Kapitalchen wird für Hervorhebungen genutzt. Üblich ist durchaus auch die Hervorhebung nur einer Silbe oder gar das Setzen eines Fremdwords in zwei Schriften (Wortstamm Antiqua, Endung Fraktur; vgl. Abb. 30).

## 2.5 Verwendung gebrochener Schriften

In den hier betrachteten Funeralschriften werden überwiegend gebrochene Schriften für die laufenden Texte verwendet. Es handelt sich vornehmlich um verschiedene Arten von Frakturschriften ergänzt um die Schwabacher Schrift, die sichtbar für Hervorhebungen (v. a. von Namen) eingesetzt wurde. Ferner finden Antiquaschriften für lateinischsprachige Zitate in verschiedenen Graden und Schnitten Verwendung.

Die heutige Definition dieser Schriften (oft verkürzt auch nur als Frakturschriften<sup>22</sup> bezeichnet) macht sich vor allem an den Brechungen der einzelnen Buchstaben fest, hat sie doch ihren Ursprung in der gotischen Textura – einer Schrift ohne jegliche Rundungen. Nachdem in der Renaissance vor allem die aus ihr entstandene offenere und leichter lesbare Schwabacher Schrift Verbreitung fand, ging man in der Mitte des 16. Jahrhunderts dazu über, schlankere Buchstaben mit klarerer Trennung und weniger Verzierungen zu verwenden: die Fraktur war entstanden. Unterschiede zwischen den vorstehend genannten Schriften findet man beim Vergleich der Punzen (Innenraum von Buchstaben). Zeigen gotische Schriften durchweg eckige und schmale Punzen in Form von verschobenen Rechtecken bei den Kleinbuchstaben, so haben Frakturschriften abgerundete Punzen, z. B. beim b, d, o und h. Bei der Schwabacher bilden die Punzen der Kleinbuchstaben Oval und Halboval.<sup>23</sup> Frakturschriften laufen recht eng, d. h. sie haben kleine Wortzwischenräume.

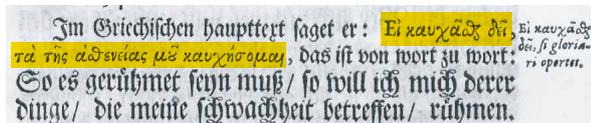
Ein harmonisches und freundliches Schriftbild entsteht aus der Grauwirkung der Fraktur. Kapr beschreibt sie wie folgt: „Das o, das a und alle daraus entwickelten Kleinbuchstaben bestehen aus einem geraden und einem geschwungenen Abstrich mit relativ spitzem Fuß. ... geschnäbelte Oberlängen der Kleinbuchstaben, und die sogenannten Elefantenrüssel, die den Buchstaben einleitenden Schwünge, geben den Versalien A, B, M, N, P, V und W ihren besonderen Ausdruck. [Vgl. Abb. 33 u. 34] Die Versalien sind im allgemeinen breit, die Gemeinen jedoch schmal gehalten. Der Gesamteindruck ist

<sup>22</sup> Die umfassendere Definition der Fraktur ist belegt in Adelung FRACTUR 1808 bzw. im großen Brockhaus FRAKTUR 1954.

<sup>23</sup> Vgl. REHSE 1998 S. 15-21

schwungvoll und reich bewegt, er entspricht unserer konsonantenreichen deutschen Sprache und gestattet das schnelle Erfassen ihrer verhältnismäßig langen Wörter.“<sup>24</sup> „Obwohl die Fraktur viele Ausschmückungen erfahren hat, sind die ursprünglichen Federzüge stets erkennbar geblieben. Vier Jahrhunderte lang war die Fraktur die am meisten verbreitete Schrift in Deutschland. Nahezu alle Bücher, sogar wissenschaftliche, fast alle Tageszeitungen und Zeitschriften, in der Regel aber auch Akzidenzen (Geschäfts- und Kleindrucksachen) wurden in Frakturschrift gedruckt“<sup>25</sup>, wie Rehse weiter ausführt. Als besondere Frakturschrift des ausgehenden 17. Jahrhunderts ist die Luthersche Fraktur zu nennen, die Kapr charakterisiert: „Die Großbuchstaben sind teilweise mit phantasievollen doch überflüssigen Schnörkeln kalligrafisch moduliert. Aber das Hervorstechende bei der Luther-Fraktur ist der Einfluß des Stichels, der den geschriebenen Duktus maßvoll korrigiert.“<sup>26</sup>

Von den oben genannten Schriften finden sich vor allem die verschiedensten Frakturschriften in den betrachteten Leichenpredigten; Zahlen erscheinen immer als Mediävalzahlen, also mit Ober- und Unterlängen. Je nach Entstehungszeit stehen am Anfang eher schmucklose Schriften, die in verschiedenen Schriftgraden kombiniert werden. Auszeichnungen finden zunächst fast immer in der Schwabacher Schrift statt. Im laufenden Text springen derlei Zitate sofort ins Auge. Im Barock hingegen erscheinen viele Frakturschriften, die zunehmend verschnörkelte Versalien, aber schließlich auch verzierte Kleinbuchstaben aufweisen. Auszeichnungen werden nun durch größere und fettere Schriftgrade oder abweichende Frakturschriften erzielt; die Schwabacher Schrift tritt in ihrer Verwendung in den Hintergrund. Wie eingangs bereits erwähnt, kommen zunächst kleinere Schriften (12 pt) zum Einsatz. Im ausgehenden 17. und beginnenden 18. Jahrhundert hingegen findet man oft 16-pt- oder gar 18-pt-Grundschriften auch für längere Texte. Unterbrochen wird das Druckbild im Stile der Zeit durch in Antiqua gesetzte Versalien für Namen oder in Antiqua (auch kursiver Antiqua) gesetzte lateinische Zitate. Bibelzitate sind auch direkt in hebräischer oder griechischer Schrift wiedergegeben (vgl. Abb. 31). Hervorzuheben sind hier ebenfalls die mannigfaltigen Formen und Größen von Initialen, deren Inhalt oft nur im Kontext ge-



31 Zitate in griechischer Sprache



32 Initiale, Leipzig, 1691. Ec 705-32 S. 3



33 Vorkommende Muster für Versalie A

lesen und verstanden werden kann (vgl. Abb. 32). Das Zusammentreffen all dieser Schriften führt zu einer großen Vielfalt, die eine Gruppierung von Drucken außer nach ihrem Entstehungszeitraum sehr erschwert. Möglichkeiten der Gruppierung eröffnen vornehmlich die Großbuchstaben M, I, C und U, deren Formen und Elefantenrüssel sich oft von Schrift zu Schrift unterscheiden. Allerdings ist ihr Auftreten im Verhältnis zu den Kleinbuchstaben eher selten, so dass längere Textabschnitte gelesen werden müssen. Es bleibt abzuwarten, ob sich ein solches Vorgehen als praktikabel erweist.



34 Vorkommende Muster für Versalie M

<sup>24</sup> KAPR 1959 S. 49

<sup>25</sup> REHSE 1998 S. 21

<sup>26</sup> KAPR 1959 S. 59

## 2.6 Diversität vorkommender Schrifttypen

Eine der Schwierigkeiten bei der OCR Alter Drucke besteht in den verwendeten Typen, die von der speziell in den FineReader-Modulen bereits gut erschlossenen Frakturschrift des 19. Jahrhunderts deutlich abweichen. Es sollte daher versucht werden, zum Material passende Muster zu trainieren. Um die Tests unter Laborbedingungen möglichst gut interpretieren zu können, wurde dabei mit einer sehr homogenen Schriftart begonnen. Diese Homogenität zeigte sich in einem Teil der in Königsberg gedruckten Funeralschriften des bekannten Dichters Simon Dach. Wurde diese erste aus 112 Drucken bestehende Marge fast ausschließlich aus Trauergedichten gebildet, so war eine zweite Schriftgruppe etwas offener angelegt. Sie konzentrierte sich aber auch auf denselben Zeitraum, wenige Drucker und zwei dominante Schriftgrößen. Jene Marge bestand aus 141 Leichenpredigten, die Schriften mit hoher Ähnlichkeit aufwiesen, allerdings auch im deutschsprachigen Fließtext mit etlichen Einsprengseln davon abweichender Typen. Um eine schnelle Zuordnung zu erreichen, wurden automatisch und stichprobenartig Teilbildserien erzeugt und visuell in Gruppen aufgeteilt (vgl. Abb. 35 u. 36).

<p>Wittenberg : Hake  <b>Lucius, Johann Andreas</b>          Helms des Heilß der geistlichen          Ritter Jesu Christi: oder die selige          Hoffnung Auß denen Worten S.          Pauli 1. Corinth. 15. vers. 19.          Hoffen wir allein in diesem Leben          auff C...          [3-4] Bl          (Marge 2)</p>								
<p>PPH: 630914079          1667          Zerbst : Palm  <b>Durr, Johann</b>          Scutum Fidei Nobilitissimum:          Aller-Edelstes Glaubens-Schild:          Und in demselben Der Lebendige          Erlöser: Auß den schönen Worten          Hiobs: Ich weiß daß mein Erlöser          lebet [et]c. ...          [36] Bl          (Marge 2)</p>								
<p>PPH: 630907964          1658          Königsberg : Reusner  <b>Dach, Simon</b>          Nil agit in mentem sors mala          morsque bonam. Das ist          Christliches Klag- und Trost-          Geticht Welches bey sehr          frühzeitigen und          hochbetrauerlichem wiewol          seligem Hintritt aus ...          [4] Bl          (Marge 3)</p>								
<p>PPH: 63090121X          1658          Königsberg : Reusner  <b>Dach, Simon</b>          Einfüllige Klag- und Trost-Reime          Welche bey tödlichem und sehr          betrauerlichem Hintritt aus dieser          Welt Der ... Frauen Barbaren          gebornen Korschinn: Des ... Herrn          Reinhol...          [4] Bl          (Marge 3)</p>								
<p>PPH: 630916047          1665          Cöthen : R6el  <b>Lezius, Heinrich</b>          Trau- und Trost-Predigt/ Bey          dem Christlichen Leichbegängnis/          Der ... Frauen Dorotheen Sabinen/          gebornen von Kragen: Des ...          Herrn Ludwigi von Wutenau auff          Trümm Erbh...          [1] Bl, 54 S          (Marge 2)</p>								
<p>PPH: 630918260          1667          Cöthen : Roelen  <b>Raumer, Georg</b>          Hiobs 9. Vörmne: Streit und  </p>								

35 Überblicksdarstellung zur Schriftgruppenanalyse der Drucke

Wer pflegt der Raben junge Zucht  
In Nöthen zu ernehren?  
Gott/welchen alles Fleisch ersucht  
Ihm Speise zu gewehren.

Seyd Ihr in einer frembden Stabt/  
Lasset diß/fals ab zu weinen:  
Wer seinen Gott nur bey sich hat  
Hat über all die Seimen.

So haltet seinem Rahtschluß still.  
Er hat in seinen Händen  
Der Menschen Herr/wohin Er wil  
Dahin kan er es wenden.

Vnd wer ist Ihm an Kräfften gleich?  
Springt Er zu Ewren Sacken/  
So wird Er Freund aus Steinen Euch  
Vnd Ewren Söhnen machen.



und Trostpredigt. 9

de in der Welt mehr haben können: Aber der Herr  
sagete ihnen ein anders. Es ist euch gut / daß ich  
hingehe: Denn so ich nicht hingehe / kömmet  
der Tröster nicht zu euch / so ich aber hingehe /  
wil ich ihn zu euch senden. Der Hingang Christi  
zum Vater / das ist / seine siegreiche Himmelfahrt/  
beraubet uns so gar des Trostes nicht / daß sie uns  
vielmehr zu desto kräftigern Trost dienet. Die Jün-  
ger Christi sind nie freudiger gewesen / als nach der  
Himmelfahrt Christi / da ihnen der Herr Jesus/  
wie Tertullianus redet / vicariam vim Spiritus Sancti  
gesandt hatte / da sie vorher oftters traurig wur-  
den / freueten Sie sich hernach in dem Herrn alle-  
zeit / und war ihre Freude so groß / daß auch ganze  
Ströme der Verfolgungen / ja der Tod selbst / sie  
nicht aufleschen konte.

Und solches Trostes genießen die Kinder Got-  
tes in ihren Trübsalen noch / Sehen Sie schon den  
Herrn Jesum mit leblichen Augen nicht / oder kön-  
nen aus seinem holdseligen Munde die Tröstungen  
unmittelbar nicht anhören. So ist der Herr  
dennoch bey ihnen in der Noth / und wil bey  
Ihnen bleiben bis ans Ende der Welt / Er sie-  
het ihre Thränen / Er zehlet ihre Flücht: Er  
erhält sie durch die rechte Hand seiner Gerech-  
tigkeit: Er tröstet Sie in Trübniß / und ver-  
kehret ihr Trauren in Freude. und zu dem Ende  
hat Er uns auch gelassen sein Wort / ein Wort des  
Trostes und Gedult / das uns tröstet in allerley  
Trübsalen: Er hat versprochen ( nicht nur seinen  
Jüngern /

v. 7.  
de prescript:  
Haret.  
Psal. 90: 15.  
Matt. 28: 20.  
Psal. 56: 9.  
Esa. 41: 10.  
Jer. 31: 14.

Ee 705-1469

LUX IN TENEBRIS:  
Freud im Leid:  
Das ist  
Trost-Predigt/  
Aus dem Siebenden Capitel Micha v. 7. 8. 9.  
Bey Christlicher und Eandes-mäßiger Beerdigung  
Des weiland Wolwürdigen Hoch-Wolgebornen  
Herrn  
**Herrn Hebbhard /**  
Des Heiligen Röm. Reichs Erb-Truchsen/  
und Freyherrn zu Waldenburg / auff Wildenhoffen und Lande-  
burg Erbherrn / Sr. Churfürstl. Durchl. zu Brandenburg  
Cammer-herrn und Obristen-Vicutenants / auch  
des Johanner-Ordens Ritters:  
Nach dem derselbe am 9. Octobris Anno 1664. unfern  
Wien / sein Leben frühzeitig geendet / und dessen Körper am 30.  
Martii 1665. in dem Freyherrl. Schwerinschen Erb-  
Begräbniß zu Alen Landsberg beygesetzt  
worden:  
Daselbst gehalten  
Durch  
JOHANNEM Buntebart/  
Dienern am Worte Gottes bey der Kirchen zu  
St. Peter in Cöln an der Spree.  
Zu Berlin/  
Bedruckt bey Christoff Runge / Anno clb lbc LXV.



36 Beispielseiten der ersten Margen

## Versuch der Gruppierung

### – nach Merkmalen

Zur Gruppierung, also Clustering, wurden unter typenkundlicher Beratung bestimmte individuell verwendete markante Großbuchstaben als Unterscheidungsmerkmal erwogen, sowie einige weniger markante, im Text aber durchgängiger auftretende Kleinbuchstaben. Um für eine effektive OCR noch hinreichend große Mengen von Bildern gemeinsam verarbeiten zu können, konnte und brauchte die Aufteilung andererseits nicht beliebig weit getrieben zu werden. Eine erste Sortierung innerhalb der Drucke aus der Mitte des 17. Jahrhunderts trennte (außer nach Schriftgröße) nach deutlich unterscheidbaren Formen des kleinen runden „s“ (mit Rüssel wie in Fraktur und Textur oder ohne wie in Schwabacher und Rundgotisch, gleich oder ungleich gewichtete vertikale Hälften usw.; vgl. Abb. 37). Für in weiteren Versuchen hinzukommende spätere Drucke eigneten sich etwa der Übergang von Virgeln zu Kommata im Satz oder – bei in das späte 18. Jahrhundert hineinreichenden Teilbeständen – eine bereits merkliche Annäherung an die schlankere Fraktur des 19. Jahrhunderts als initiales Unterscheidungsmerkmal für separat zu trainierende Teilbestände. Im letzten Fall (zunehmende Nähe zu Schriften des 19. Jahrhunderts) spricht viel für den Verzicht auf eigenes Training bei zunehmender Brauchbarkeit der in Softwareprodukten (z. B. in der FineReader-Engine) bereits integrierten Muster (Vgl. Kap. 5).



37 Verschiedene kleine „s“ als Gruppierungsmerkmal

### – alternativ: Gruppierung empirisch

Sobald erste trainierte Muster schon vorliegen, bietet es sich in der Praxis auch an, die Zuordnung von Bildern zu einem Musterbestand nicht vorab nach visuellen Merkmalen, sondern nach einem OCR-Gang empirisch anhand der Ergebnisse zu beurteilen. Im Idealfall liegen für Teile des zu prozessierenden Bestands Transkriptionen oder andere Ground-Truth-Daten (manuell erfasste, auch semantisch ausgezeichnete fehlerfreie Daten) vor, so dass für diese eine echte Erkennungsrate ermittelt und verglichen werden kann. Meist wird solches Referenzmaterial aber fehlen, so dass – mit vorsichtiger Interpretation und je nach Güte der vorliegenden Wörterbücher – hilfsweise z. B. der Anteil der sowohl im Lexikon als auch in den OCR-Texten enthaltenen Wörter als Indikator dienen kann. Wo dieser Anteil deutlich sinkt, sind oft Schriftartwechsel, Sprachwechsel und dergleichen die Ursache; im getesteten Material waren z. B. sporadische Seiten mit deutlichem Antiqua-Anteil auf diese Art schnell nachträglich zu erkennen. Praktisch sinnvoll wird diese nachträgliche Zuordnung immer dann sein, wenn Rechenzeit für Mehrfach-OCR zur Verfügung steht und die Bearbeiter die nachträgliche Auswahl im Vergleich zu einer vorherigen Sichtung jeder Seite schneller vornehmen können (Vgl. Kap. 4.5).

### 3 Typische OCR-Teilschritte und Einflussgrößen

#### 3.1 Typischer Workflow

Für den Zweck der Darstellung des Softwaretests soll die Verarbeitungsschrittfolge wie folgt vergrößert (vgl. Abb. 38) behandelt werden:<sup>27</sup>

- Sichtung und Vorsortierung des Bildmaterials;
- Binarisierung (inkl. vorheriger Bearbeitung des Farb- bzw. Graustufenbilds und Nachbearbeitung des bitonalen Bilds);
- Segmentierung (Layoutanalyse, Blockerkennung, Segmentierung auf Wort- und Zeichenebene);
- OCR (inkl. unmittelbar verbundener Fehlerkorrekturmechanismen);
- Lexikalische Nachkorrektur<sup>28</sup>.

Neben den Einzelschritten begrenzt nicht zuletzt die stapelweise Verarbeitung an sich die mögliche Erkennungsqualität.

Dieser Workflow setzt das Vorhandensein einsatzbereiter OCR-Musterbibliotheken voraus. Wenn diese nicht oder nicht passend bereitstehen, ist deren Erstellung bzw. Bearbeitung voranzuschicken. Im praktischen Softwarevergleich (Vgl. Kap. 4) wird dieser Arbeitsgang separat behandelt.

#### 3.1.1 Bereitstellung des Bildmaterials

##### Stufeninput und -output

*Eingang:* Bildmaterial

*Ausgang:* Bildmaterial selektiert nach Eignung, z.B. gruppiert in Margen, die gemeinsame OCR-Läufe zulassen

##### Typische Informationsverluste und Fehlerquellen

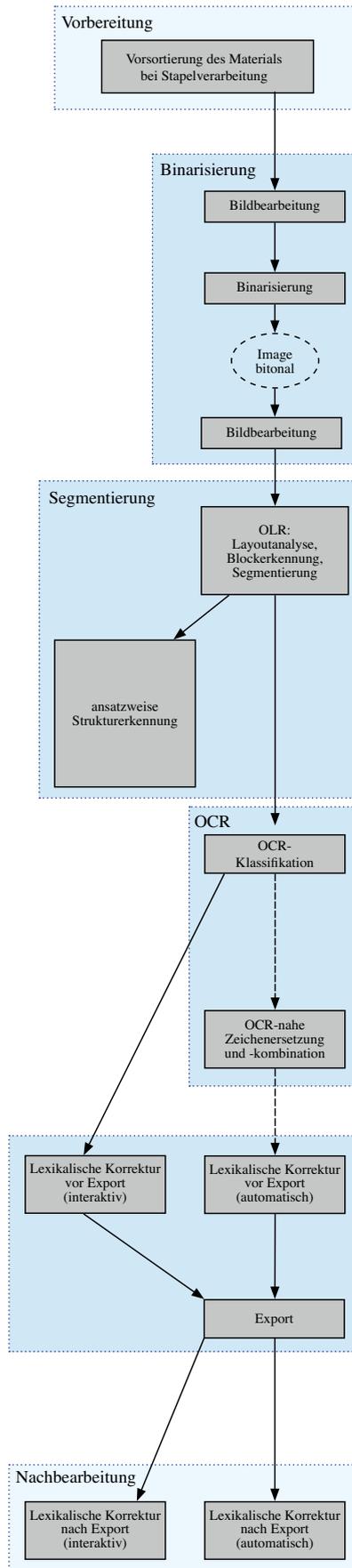
- Unvermeidlich:*
- a) Das menschliche Vorwissen über Textsorten geht i. d. R. nicht in den OLR-/OCR-Prozess mit ein.
  - b) Schlechte Vorlagen aufgrund von Druckqualität, Durchdruck, Papierverfärbung, Verschmutzungen und Annotationen müssen in gewissem Maße hingenommen werden.
- Potentiell beeinflussbar:*
- c) Die Scans schwanken in Bildhelligkeit und -kontrast sowie Randmaßen.
  - d) Innerhalb des OCR-Laufs (Seite, Stapel) kommt es zum Wechsel von Schriftart, Sprache oder typographischen Konventionen.
  - e) Ein Finetuning anhand von bestimmten Beispielseiten kann für andere Seiten desselben Stapels nachteilig sein.

Diese Verluste können oft durch Sichtung und Vorsortierung des Materials, durch eine Vorverarbeitung der Bilder (Randentfernung oder einheitliche Randdarstellung, einheitliche Satzspiegel, Teilung von Margen in rechte/linke Buchseiten usw.) begrenzt werden.

<sup>27</sup> Detaillierte Flussdiagramme gibt es in der Literatur zu konkreten OCR-Untersuchungen: einige Abläufe sind vergleichend dargestellt in KÄMMERER 2009. Nicht dargestellt sind interne Rückkopplungsmechanismen, wie sie ansatzweise auch in den getesteten Programmen integriert sind (je nach Ergebnis adaptiv wiederholte Anwendung von Segmentierungsalgorithmen, Zeichenzuweisungen oder Wortkorrekturschritten).

<sup>28</sup> Linguistische Nachkorrektur wird von beiden Programmen nicht angeboten und daher aus der Betrachtung ausgeschlossen.

Typische OCR-Anwendung



38 Typischer Workflow

### 3.1.2 Bildvorverarbeitung und Binarisierung

Da die zur Zeichenerkennung gespeicherten Muster-Repräsentationen sich auf Schwarzweiß-Zeichen beziehen, muss die Vorlage in ein bitonales Bild gewandelt werden. Vor und/oder nach diesem Schritt werden i. d. R. verschiedene Filterungen angewandt wie Helligkeits- und Kontrastnormalisierungen, Ausgleich von Drehwinkeln, Wölbungen und Verzerrungen („Skew correction“), Ausfilterung von Verunreinigungen, die anhand ihrer (kleinen) Größe oder ihres Abstands zu offensichtlichen Zeichen- und Textregionen erkennbar sind (typisch für altes Papier wie für Mikrofilme), Abschneiden von erkennbaren Rändern und Nicht-Text-Bereichen usw.

#### Stufeninput und -output

**Eingang:** Farb- bzw. Graustufenbild  
+ **Steuerdaten:**  
potentiell: Bildvorverarbeitungs- und Binarisierungsparameter

**Ausgang:** Bitonales, unter Umständen beschnittenes oder gefiltertes Bild

#### Typische Informationsverluste und Fehlerquellen

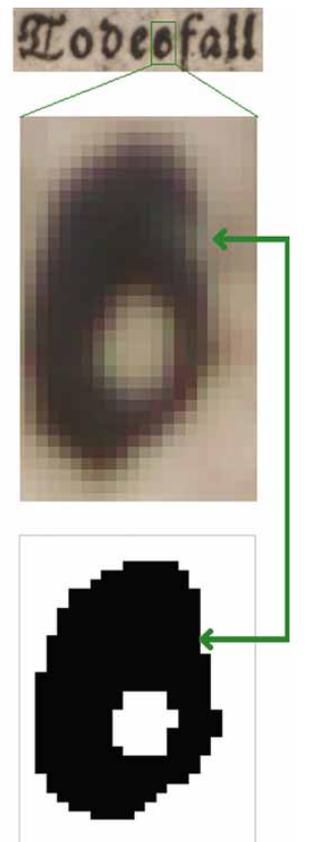
**Unvermeidlich:** a) Die Vergrößerung der Vorlage durch Reduktion der Farb- oder Graustufeninformation auf ein bit je Pixel führt dazu, dass relative Helligkeitsunterschiede innerhalb des hellen oder dunklen Teils des Zeichenbilds der Erkennung nicht mehr zur Verfügung stehen. Wegen der schon bei wenigen Graustufen exponentiell ansteigenden Anforderungen an Musterrepräsentation, Speicherplatz, Rechenzeit und Muster-Trainingszeit wird dieser erste entscheidende Informationsverlust der OCR noch auf längere Sicht nicht überbrückbar sein.

**Potentiell beeinflussbar:** b) Scan-bedingte Verluste entstehen durch:

- ungünstige Farbintensität und Kontrast;
- Verzerrungen (Wölbung, Knicke);
- Verunreinigungen;
- Ränder und graphische Elemente; usw.

*Als wichtigste Abhilfemöglichkeiten sind zu nennen:*

- eine Normalisierung im Voraus oder durch die OCR-Software
- eine Gruppierung nur der einander vergleichbaren Bilder für eine gemeinsame Verarbeitung (s. o.);
- eine Entzerrung im Voraus oder durch die OCR-Software;
- eine softwareseitige Layouterkennung;
- ein manuelles Schneiden der Bilder;
- softwareseitige Filter;
- manuelle Filter.



39 Informationsverlust durch Binarisierung

- c) Die Unterschiede zwischen „gleich hellen“ unterschiedlichen Farben gehen verloren. Abhilfe kann durch eine passende Farbfilterung geschaffen werden.
- d) Die Anwender verzichten auf ein Finetuning wegen fehlender oder zu unkomfortabler und zeitaufwendiger Konfigurierbarkeit der Software.  
*Abhilfe ist allenfalls möglich durch:*
  - den Einsatz von gut eingearbeitetem Personal (anwenderseitig),
  - den Einkauf des Finetunings als Dienstleistung (organisatorisch),
  - die Optimierung der Software auf leichte und verständliche Konfigurations-Interfaces mit übersichtlicher Vorschau der Ergebnisse jeder Parameteränderung (entwicklerseitig).

**Wann unschädlich:**

Unter günstigen Umständen kann dieser Informationsverlust unbedeutend werden (und im wesentlichen darauf sind die in bestimmten Office-Umgebungen heute gelegentlich erreichbaren hundertprozentigen Erkennungsleistungen zurückzuführen): nämlich, wenn z. B. durch hohen und gleichmäßigen Kontrast, durch unbeschädigtes Druckbild (keine „brüchigen“ Lettern, keine bedeutungsunterscheidenden Serifen ...) usw. das relevante Drucktypenbild entweder bereits ohnehin in einer nahezu bitonalen Ausprägung vorliegt oder Farb- und Helligkeitsstufen von vornherein nur redundante Information tragen, und wenn ein geschlossener Typenvorrat eine trennscharfe Abbildung auf die Mustermenge erlaubt. In Alten Drucken ist dieser Idealzustand nur in Ausnahmefällen gegeben.

Fast immer unschädlich bzw. beherrschbar ist der bloße Verlust der Farbinformation gegenüber Graustufen; speziell in Alten Drucken kommen zwar Wechsel zwischen rotem und schwarzem Textdruck vor, aber kaum Situationen, in denen allein die Farbqualität zwischen Vorder- und Hintergrund unterscheidet, die Farbintensität aber nicht.

### 3.1.3 Segmentierung (Layoutanalyse, Blockerkennung, Segmentierung auf Wort- und Zeichenebene)

#### Stufeninput und -output

**Eingang:** Binärbild  
+ *Steuerdaten:*  
potentiell: Segmentierungsparameter wie Abstandsmaße, Mindest- und Höchstgrößen von Zeichen usw.

**Ausgang:** Satz einzelner Binärbilder bzw. Koordinaten der Bildregionen

#### Typische Informationsverluste und Fehlerquellen

**Unvermeidlich:** a) Wo typographische Merkmale (Zwischenraum zwischen Spalten/Absätzen/Wörtern/Zeichen, konsistente Linienführung, Zeilenausrichtung usw.) fehlen, können diese Segmente maschinell nicht richtig erkannt werden.

**Potentiell beeinflussbar:** b) Für menschliche Gestaltwahrnehmung intuitiv bedeutungstragende Layoutmerkmale können nicht in für alle Materialien gleichermaßen gültige Operationalisierungen umgesetzt werden.

- c) Lexikalische und semantische Indikatoren bleiben für die Strukturierung ungenutzt.<sup>29</sup>
- d) Etwaige in Katalogen oder Datenbanken bereits erfasste Strukturinformationen bleiben ebenfalls ungenutzt.

Weil aus verlorener Strukturinformation regelmäßig eine in der Vorlage nicht gegebene Mehrdeutigkeit der Zeicheninterpretation resultieren kann, bildet diese Stufe den zweiten entscheidenden Informationsverlust der maschinellen OCR. *Abhilfe* wird hier vor allem entwicklerseitig durch eine konsequente Unterstützung der Hinzuziehung solcher Strukturinformationen möglich sein. Bis dahin kann eine maschinelle Layoutanalyse stets nur Teile der Anforderungen erfüllen.

**Wann unschädlich:** Am wenigsten problematisch ist dieser Verlust in Fällen sehr einfacher Layouts und strukturell möglichst gleichwertiger Zeilen.  
Manche Strukturinformation wird sich nach der OCR noch aus Merkmalen des gelesenen Texts in Nachkorrekturschritten rekonstruieren lassen. Für eine schon bei der Zeichenerkennung als Musterrestriktion wirksame Vorab-Segmentierung bleibt sie aber auch dann ungenutzt.

### 3.1.4 OCR (inklusive unmittelbar verbundener Fehlerkorrekturmechanismen)

Für die Klassifikation von in bitonaler Form vorliegenden Zeichenbildern gibt es sehr unterschiedliche leistungsfähige Verfahren und einen entwickelten Forschungsstand, so dass bei geeignetem Input an Bildern und Musterklassen potentiell eine nahezu verlustfreie Zuordnung erreichbar ist. Dabei ist es gleich, ob konventionelle Ähnlichkeitsmaße berechnet oder neuronale Netze trainiert, ob eine topologische Abstraktion des Zeichens oder nur aus dem Bild abgeleitete Kennzahlen für die Klassifikation herangezogen werden.<sup>30</sup>

Die Abgrenzung des OCR-Schritts von der vorangehenden Segmentierung einerseits und der nachfolgenden Anwendung erster Zeichenersetzungsregeln andererseits ist unscharf. Ursache dafür ist die Klassifikation, die in guter OCR-Software zum Teil adaptiv auf die Segmentierung zurückwirkt (wenn etwa ein möglicherweise teilbares Zeichensegment mehrfach gelesen wird und je nach Möglichkeit in zwei Segmente geteilt wird oder nicht, oder wenn bestimmte Teilmuster sofort zu einem kombinierten Zeichen zusammengesetzt werden).

#### Stufeninput und -output

**Eingang:** Bildregionen auf Textblock-, Zeilen- und Einzelzeichenebene  
+ *Steuerdaten:*  
Mustersammlung  
potentiell: Wörterbuch  
potentiell: unmittelbar auf Zeichen(-teil)-Ebene ansetzende Ersetzungsregeln

**Ausgang:** OCR-Text  
Koordinaten der gelesenen Zeichen oder „Wörter“  
+ *Evaluationsdaten:*  
potentiell: statistische Auswertungen wie Anteil erkannter Buchstaben und Wörter, Konfidenzwerte zur vermuteten Erkennungssicherheit auf Zeichen- und Wortebene usw.;  
potentiell: beliebige computerlinguistische Auswertungen

<sup>29</sup> Zu Möglichkeiten der regelbasierten Dokumentstrukturanalyse *nach* der OCR vgl. Arbeiten an der Universitäts- und Landesbibliothek Tirol im Rahmen des IMPACT-Projekts, s. MÜHLBERGER 2011; der Nutzen der Strukturinformation wird aber auch dort überwiegend für die Nach-OCR-Phase gesehen.

<sup>30</sup> Zu Klassifikationsverfahren s. u. a. SCHULZ 2006

## Typische Informationsverluste und Fehlerquellen

- Unvermeidlich:*
- a) Es kommt zu falscher Musterzuordnung von beschädigten bzw. verunreinigten Drucktypen, wenn eine hohe Ähnlichkeit zu anderen Zeichen besteht.
  - b) Einige Zeichenwerte auf Muster überlappen sich. Das ist der Fall, wenn im gleichen Musterbestand eine Schrift für ein Zeichen X Druckbilder zeigt, die in einer anderen Schrift desselben Musterbestands für ein Zeichen Y stehen. Fehlerhaft ist hier bereits der Zustand der Musterdatei, besonders wenn die Software eine derartige Vermengung schon beim Anlegen der Musterdatei erzwingt und keine separaten Musterbestände ermöglicht, die später fallweise kombiniert werden können.
  - c) Bei Klassifikation eines Zeichens – auch wenn sie durch Vermerk von Konfidenzwerten gewissermaßen „unter Vorbehalt“ geschieht – geht die Information über die Alternativkandidaten und ihre Rangfolge verloren. Dies muss als der dritte wesentliche Informationsverlust heutiger maschineller OCR beachtet werden: er beeinträchtigt von vornherein die Treffsicherheit jeder nachfolgenden lexikalischen Korrektur, die nun unter allen *generell möglichen* Verwechslungen (im besten Fall sprachbezogen gewichtet) auswählen muss und sich nicht mehr gezielt auf die wenigen am Ort *für das konkrete optische Muster* allein möglichen Verwechslungen beschränken kann.

- Potentiell beeinflussbar:*
- d) Zeichen, die nicht zur aktuell geladenen Musterbibliothek passen, werden nicht erkannt. Eine universelle Mustersammlung (wie z. B. die von FineReader mitgelieferte Fraktur-Bibliothek) bietet zwar einen optimalen Durchschnitt, verschenkt aber potentiell mögliche Erkennungsquoten auf Drucken, in denen eben nicht die ganze Breite der Muster vorkommt und daher zu viele irrelevante optische Muster angewendet werden. Abhilfe ist abhängig davon, ob die Software ein Vorhalten separater „reiner“ Musterbestände zum Zweck fallweiser Kombination überhaupt zulässt.
  - e) Zeichen, die in Alphabet- und Sprachdefinitionen nicht vorhergesehen wurden, werden nicht erkannt.

*Wann unschädlich:* Echte OCR-Fehler sind im Prinzip nur dann unschädlich, wenn sie durch lexikalische oder linguistische Nachkorrektur rückgängig gemacht werden können. Hierzu sind sehr genaue Annahmen über Regelmäßigkeiten des Ausgangstexts nötig, was zumindest in formalisierten Textsorten oder Texten mit inhaltlich „geschlossenen Welten“ oft gegeben sein kann. Wenn im industriellen Umfeld eine praktisch fehlerfreie OCR durchaus anzutreffen ist, dann i. d. R. unter solchen Umständen. Andererseits ist eine Vermeidung von OCR-Fehlern um so aussichtsreicher, je überschaubarer und trennschärfer der verwendete Drucktypenvorrat ist und je vollständiger und exklusiver er durch Training in die Musterbibliothek eingearbeitet werden konnte.

– Zwischenanmerkung –

Will man eher auf Fehlervermeidung als auf nachträgliche Fehlerkorrektur setzen, muss bis zu dieser Stufe das Maximum an Lesetreue erreicht werden.

### 3.1.5 Lexikalische und linguistische Nachkorrektur

Jede auf die OCR folgende Fehlerkorrektur versucht, einen bereits eingetretenen Wissensverlust nachträglich zu kompensieren. Das ist naturgemäß nur unter Heranziehung von Regularitätshypothesen über die Vorlage möglich, die oft heuristisch bleiben müssen und im Einzelfall versagen können. Man wird auf diese Weise auch in Fällen „richtiger“ Korrektur i. d. R. nicht wissen können, ob die Korrektur richtig war. Hinzu kommt die Anforderung, dass reale Druckvorlagen den angenommenen Regularitäten ihrer Textsorten nicht immer folgen (im einfachsten Fall: Druckfehler oder Schreibvarianten enthalten) und auch dann von der OCR treu wiedergegeben werden sollen.<sup>31</sup>

Diesen Nachteilen nachträglicher lexikalischer und „linguistischer“ (morphologischer, syntaktischer, kontextbezogener usw.) Fehlerkorrektur steht der große Vorteil gegenüber, dass ab dieser Stufe auf Textformaten gearbeitet werden kann, d. h. mit der Computerlinguistik die Ressourcen einer ganzen Wissenschaft und eine große Vielfalt an Werkzeugen zur Verfügung stehen.<sup>32</sup>

#### Es gibt ab dieser Stufe daher zusätzliche Alternativen:

- Zwischen in OCR-Software integrierten Korrekturverarbeitungsschritten und einer externen Bearbeitung durch andere Tools kann gewählt werden.
- Ein einmal vorliegendes Roh-OCR-Ergebnis kann (nacheinander oder parallel-verzweigend) mehrfachen Korrekturschritten zugeführt werden.

Erschwerend kommt bei alten Schriften hinzu, dass die zur Nachkorrektur verwendeten Wörterbücher die historische und uneinheitliche Rechtschreibung berücksichtigen müssen; eine gewisse Aufblähung mit Schreibvarianten ist daher nicht zu vermeiden.

Auf echte linguistische Korrekturmöglichkeiten wird hier nicht eingegangen, weil beide untersuchten Softwareprodukte diese nicht anbieten.<sup>33</sup> Der Übergang ist aber fließend, da auch in die lexikalische Korrektur linguistische Komponenten eingehen können, z. B. wenn das Lexikon mit morphologischen Ableitungen, abgeleiteten historischen Schreibvarianten oder Kollokationen angereichert werden konnte.

#### Gesichtspunkte zur Beurteilung der lexikalischen Korrektur können sein:

- ob sie vollautomatisch oder (teil-)manuell erfolgt;
- ob erfasste Wortkoordinaten angepasst werden oder nicht;
- welche sprachbezogene Zusatzinformation verwendet wird (Worthäufigkeiten? Konfusionswahrscheinlichkeiten? Zeichenübergangswahrscheinlichkeiten? usw.);
- ob diese Zusatzinformation implementiert ist oder vom Nutzer fallbezogen parametrisiert werden kann; usw.<sup>34</sup>

<sup>31</sup> Hinweis u. a. von Marco Büchler, der als „big challenge“ der OCR nennt: „How can we correct OCR without correcting spelling errors?“ BÜCHLER 2011

<sup>32</sup> Detaillierte sprach- und dokumentenspezifische Fehlerprofile wurden - mit dem Ziel der Entwicklung gut bedienbarer Korrekturwerkzeuge - u. a. im IMPACT-Projekt untersucht, vgl. REFFLE 2011. - Ein am CIS München vorangehendes DFG-Projekt „Domänen- und dokumentenadaptive Verfahren zur Nachkorrektur von OCR-Ergebnissen“ versuchte die Leistungsfähigkeit von Korrekturverfahren dadurch zu erhöhen, dass explizit „Domäne, sprachliches Bild des Gesamttexts und logischer Kontext von Dokumententeilen bei der Korrektur mitberücksichtigt werden“ (Zu Projektbeschreibung mit Publikationsangaben s. NACHKORREKTUR VON OCR-ERGEBNISSEN). Es ist davon auszugehen, dass das im Anschluss an IMPACT gegründete Kompetenzzentrum einen umfassenden Informationszugang zu solchen Möglichkeiten bereitstellt.

<sup>33</sup> Diese Thematik wird u. a. behandelt in HAUSER 2007, STROHMAIER 2004.

<sup>34</sup> Vgl. SCHULZ 2006, RINGELSTETTER 2003 u. a.

## Stufeninput und -output

- Eingang:* OCR-Textergebnis  
 + *Steuerdaten:*  
 Wortliste  
 potentiell: Gewichte für bestimmte Wörter oder Lexikon-Teilbestände (Vorrang bestimmter Kandidaten vor anderen)  
 Parameter für erlaubte Distanz, bis zu der ersetzt werden darf  
 potentiell: Gewichte für bestimmte Zeichenersetzungsschritte<sup>35</sup>
- Ausgang:* Textexport  
 + *Evaluationsdaten:*  
 potentiell: Vermerke der vorgenommenen Ersetzungen (Original-Lesart, Kandidaten) im Textexport  
 potentiell: beliebige computerlinguistische Auswertungen

## Typische Informationsverluste und Fehlerquellen

### Eine mögliche Fehlerklassifikation<sup>36</sup>

Da sich die lexikalische Nachkorrektur weitgehend materialunabhängig und in computerlinguistisch zugänglichen Textformaten abspielt, gibt es eine reichhaltige, theoretisch gut durchgearbeitete Literatur zu Einflussfaktoren, Fehlersituationen und Korrekturstrategien unter Verwendung verschieden aufbereiteter Korpora, was hier nicht wissenschaftlich referiert werden kann.

Exemplarisch soll hier eine Fehlerklassifikation von Christian M. Strohmaier kurz vorgestellt werden, die zu praktischen Vorüberlegungen bei der Anlage der verwendeten Wörterbücher beitragen kann, an die teils gegensätzliche Anforderungen bestehen. Denn einerseits sollen alle im Original tatsächlich vorkommenden Wörter im Lexikon auch enthalten sein (d. h. akzeptiert und nicht fälschlich korrigiert werden), andererseits soll das Wörterbuch nicht unnötig viele Wörter enthalten, was zur unnötig vermehrten Akzeptanz auch von OCR-Fehlern führt, die eigentlich korrigiert werden müssten.<sup>37</sup>

<sup>35</sup> Die Korrektur von als „häufig“ bekannten Verwechslungen „verbraucht“ dann weniger der zulässigen Distanz, bis zu der die Korrektur stattfinden darf.

<sup>36</sup> STROHMAIER 2004

<sup>37</sup> Im Idealfall würde auch erkannt, was für Wörter in einem Textkontext überhaupt vorkommen dürfen, und es würden auch nur solche als Korrekturkandidaten erscheinen.

### Strohmaier gliedert die Korrekturfehler in sieben Fehlerklassen auf:

- „1. **Falscher Freund (false friend)**. Falsche Freunde sind eine nur schwer aufzuspürende Fehlerklasse. Die Nachkorrektur hat keinen Anlass, an dem vorgefundenen, falschen OCR-Token zu zweifeln, da es lexikalisch ist. Entdeckungsstrategien sind die Verwendung mehrerer OCR-Engines, die Berücksichtigung von Konfidenzwerten der Wort- bzw. Charakter-Erkennung der OCR-Engine selbst und stärkere Gewichtung von Kontextinformationen. Eine Vermeidungsstrategie ist die Verkleinerung des Lexikons. Es gilt, je kleiner das Lexikon  $D$  ist, desto geringer ist die zu erwartende Anzahl falscher Freunde.
2. **Zu vorsichtig (too cautious)**. Das korrekte Wort ist Spitzenkandidat der Vorschlagsliste. Allerdings wird die Korrektur nicht ausgeführt, da der Konfidenzwert zu gering ist  $\text{conf}(w^{\text{ocr}}) < 1$ . Vermeidungsstrategie ist eine mutigere Nachkorrektur, die alle Konfidenzwerte erhöht. Allerdings verhält sich die Fehlerklasse der unglücklichen Korrekturen (5) antagonistisch zu dieser Klasse. Die Einstellung des Konfidenzwertes  $\text{conf}$  muss zwischen diesen beiden Klassen austariert werden.
3. **Falscher Kandidat und falsche Korrekturgrenze (wrong candidate and threshold)**. Das richtige Wort ist im Lexikon enthalten. Allerdings ist es nicht auf die Position des Spitzenkandidaten der Liste gerankt. Ausserdem sind die Konfidenzwerte aller Kandidaten des OCR-Tokens zu niedrig, so dass keine Korrektur vorgenommen wird. Fehlt das richtige Wort ganz in der Kandidatenliste, wurde es entweder zu stark von der OCR-Engine verunstaltet oder die Grobfilterung für die Kandidatenliste ist mangelhaft. Ansonsten ist die Vermeidungsstrategie eine verbesserte Konfidenzwertberechnung, d. h. geeignetere und/oder mehr Faktoren zur Kandidatenbewertung sowie eine bessere Kombination dieser Faktoren.
4. **Falscher Kandidat (wrong candidate)**. Die Nachkorrektur erkennt richtig, dass ein Korrekturvorschlag ausgeführt werden soll und auch das richtige Wort ist in Lexikon enthalten, jedoch nicht an erster Stelle der Kandidatenliste. Auch in diesem Fall sollte die Grobfilterung überprüft werden, wenn das richtige Wort in der Kandidatenliste fehlt. Ansonsten ist die Vermeidungsstrategie wie für die vorherige Fehlerklasse (3) eine verbesserte Konfidenzwertberechnung.
5. **Unglückliche Korrektur (infelicitous correction)**. Korrekturvorschläge, die ausgeführt werden, obwohl das Groundtruth-Token nicht im Lexikon war, führen immer zu einem Fehler ( $w^{\text{orig}} \notin D$ ). Der drastischere Fall tritt ein, wenn die OCR-Engine das Originaltoken bereits richtig gelesen hat  $w^{\text{ocr}} = w^{\text{orig}}$ . Dadurch verschlechtert die Nachkorrektur das Endergebnis. Vermeidungsstrategien sind die Vergrößerung des Korrekturlexikons und eine vorsichtigeren Berechnung des Konfidenzwertes. Allerdings steigen mit einer Lexikonvergrößerung die potentiellen, false friends (1) und mit einer vorsichtigeren Konfidenzwertberechnung die potentiellen Fehler der Klasse too cautious (2).
6. **Unvermeidbar I (no chance I)**. Die Nachkorrektur wird [wegen zu niedrigem Konfidenzwert des Ersetzungskandidaten] nicht aktiv, obwohl eine Fehlererkennung vorliegt. Da aber das richtige Wort nicht im Lexikon und daher auch nicht in der Kandidatenliste ist, ist die Nachkorrektur ohnehin chancenlos. Die einzige Vermeidungsstrategie ist eine Vergrößerung des Korrekturlexikons  $D$ .
7. **Unvermeidbar II (no chance II)**. Die Nachkorrektur erkennt richtig, dass ein Korrekturvorschlag ausgeführt werden soll und wird [aufgrund des ausreichenden Konfidenzwerts des Ersetzungskandidaten] aktiv. Es besteht jedoch keine Chance richtig zu korrigieren, da das originale Wort nicht im Lexikon enthalten ist  $w^{\text{orig}} \notin D$ . Auch hier ist die einzige Vermeidungsstrategie eine Vergrößerung des Korrekturlexikons  $D$ .“

#### und er schließt:

„Je größer das Korrekturlexikon desto geringer die Fehleranzahl des gesamten rechten Astes (Klasse (5), (6) und (7)), aber desto größer ist auch die Anzahl falscher Freunde, Klasse (1). Dies ist allerdings kein streng funktionaler Zusammenhang, sondern eine statistische Gesetzmäßigkeit. Die Fehlerverteilung innerhalb der restlichen Klasse (2), (3) und (4) bleibt im statistischen Sinne konstant.“<sup>38</sup>

<sup>38</sup> STROHMAIER 2004, S. 96f.

**Diese Fehler einbeziehend, wäre folgende Gruppierung möglich:**

*Unvermeidlich:* Dass jede der Fehlerklassen 1-7 in automatischer Korrektur vorkommt, dürfte nicht zu verhindern sein, zumal dieselben Fehler auch in einer manuellen Korrektur mit unvollständiger Sachkenntnis auftreten. Hinzu kommt, dass zwar die Lagekoordinaten der gelesenen Zeichen(ketten) anhand der Bildvorlage ermittelt wurden, für korrigierte Formen aber nur indirekt, z. B. durch Interpolation, auf Koordinaten geschlossen werden kann. Es wird hier daher geringe Abweichungen geben.

*Potentiell beeinflussbar:*

Die Wahrscheinlichkeit des Auftretens aller dieser Fehler ist weitgehend steuerbar, wenn auch zum Teil auf gegenseitige Kosten.

*Aus dem Gesagten ergibt sich als Abhilfe:*

- die Wahl eines Wörterbuchs optimaler Größe;
- gut eingestellte Gewichtungen der Ersetzungskandidaten auf Zeichen- und Wortebene.

*Außerdem sollte beachtet werden:*

- die Wahl geeigneter Distanzmaße<sup>39</sup> zwischen dem als falsch erkannten OCR-Wort und dem Ersetzungskandidaten, anhand derer eine Korrektur vorgenommen oder unterlassen wird. Sprach- und schrifttypisch gewichtete Maße sind hier eindeutig vorzuziehen, da nur sie es ermöglichen, eine im tatsächlichen Schriftbild wahrscheinliche Verwechslung leicht zu machen (z. B. „ni“ > „m“), ohne zugleich auch alle möglichen anderen Ersetzungen mit rechnerisch vergleichbar „niedriger“ Distanz anzustoßen. Die Gewichtungen sollten asymmetrisch (für jede Ersetzungsrichtung separat) zugeteilt werden können.
- dass verschiedene Teile der eingelesenen Vorlage unterschiedliche Wortbestände erwarten lassen, d. h. man Situationen einschränken sollte, in denen Textteile nicht gut zum aktuell eingestellten Wörterbuch passen.
- dass das Lexikon in zwei Rollen auftritt, die möglicherweise unterschiedlich von Unvollständigkeit oder Aufblähung des Wörterbuchs betroffen sein können: (a) als Akzeptor einer OCR-Lesart, (b) als Repertoire für Ersetzungskandidaten.

*Wann unschädlich:*

Per definitionem können Text-Fehler der maschinellen Nachkorrektur nie unschädlich sein, außer im seltenen Fall anschließender manueller Korrektur mit hoher Aufmerksamkeit und perfekter Kenntnis des Textgegenstands. Dagegen geht bei jedem nachgeschalteten weiteren automatischen Korrekturschritt die Information über die Nähe zur optischen Vorlage verloren und hat keine Restriktionswirkung mehr; die in der Nachkorrektur erst entstandenen Fehler können daher in weiteren Korrekturschritten besonders weit vom Original wegführen.

Wenig problematisch ist dagegen der Informationsverlust bezüglich der genauen geometrischen Zeichen- und Wortkoordinaten, solange diese Information nicht zum akkuraten Ausschneiden der Bildteile benötigt wird, sondern z. B. zur Hervorhebung einer Fundstelle auf dem Bild dient. Ob hier z. B. bei einer in der Nachkorrektur erfolgten Zerlegung eines falsch zusammen gelesenen Worts in zwei Teile die interpolierte Lücke auf dem Bildschirm um wenige Pixel neben der Lücke im Original angezeigt wird, dürfte an der Brauchbarkeit der Ortsinformation für die meisten Anwendungsfälle wenig ändern.

<sup>39</sup> Einige Distanzmaße werden vorgestellt in der generell als technische Einführung in OCR lesenswerten Arbeit SCHULZ 2006. Üblich sind zum Beispiel die Levenshtein-Distanz, die die Anzahl notwendiger Zeichenersetzungen zählt bzw. gewichtet aufsummiert.

### 3.1.6 Semantische Erschließung: Strukturdaten und Named Entities

Auf die im OCR-Zusammenhang besonders interessante semantische Auswertung nicht-textueller Merkmale wurde in Abschnitt 3.1.3 bereits eingegangen.

Für die semantische Interpretation von bereits auf Textebene vorliegendem Input (Post-OCR) hat der Anwender prinzipiell die Wahl, auch außerhalb des jeweiligen OCR-Programms nach geeigneten Werkzeugen und Verfahren zu suchen. Im Rahmen des beschriebenen Projekts konnten semantische Auszeichnungen innerhalb der Softwareprodukte praktisch nicht realisiert werden; möglich waren einige automatische statistische Auswertungen durch Anwendung semantischer Entitäten auf die OCR-Exporte. Möglicherweise ist die Erwartung semantischer Auszeichnung aber generell erst dann berechtigt an OCR-Software zu stellen, wenn aus dem Bibliotheksbereich selbst abgesicherte und operationalisierbare Vorgaben zur semantischen Interpretation konkreter Layout- oder textueller Merkmale angegeben werden. Softwarehäuser müssten sich hierzu sehr tief und fundiert mit buchkundlichen Fragen auseinandersetzen und scheinen sich aus Aufwands- und vielleicht auch Gewährleistungsgründen zu scheuen, selbst weitgehende Annahmen darüber in operationalisierte Definitionen zu gießen, die absehbar immer wieder in Einzelfällen oder im nächsten Bilderstapel versagen müssen. Hier „hartcodierte“ Programmierungsleistungen zu bestellen, scheint daher vorläufig keine sinnvolle Option zu sein. Vielleicht ist eine arbeitsteilige Lösung denkbar, indem Softwarehäuser vorkonfigurierte Schnittstellen für typische Strukturelemente und Named Entities<sup>40</sup> programmieren, deren konkrete Erkennungsmerkmale vom Nutzer verfeinert werden müssen.<sup>41</sup> Wenn das gelänge, könnten Anwender, die ihr Material gut kennen, Information gewissermaßen wieder „in die Texte hineintragen“.

### 3.1.7 Einfluss der Verarbeitung ganzer Seiten

Die passgenaue Anwendung der richtigen Spracheinstellungen und Mustersammlungen wird umso anspruchsvoller, je kleinteiliger die Regionen werden, die eine verlustfreie Gleichbehandlung noch vertragen. In Alten Drucken stehen lateinische, griechische oder hebräische Wörter in entsprechender Schrift inmitten von Fraktur-Kontext. Im deutschsprachigen Satz zieht sich das bis in das einzelne Wort hinein, wenn lateinische Wortstämme in Antiqua, deren der deutschen Grammatik entstammenden Flexionssuffixe aber – wie damals häufig üblich – in Fraktur gesetzt sind. Griechische und hebräische Buchstaben dienen gelegentlich als bloße Aufzählungszeichen, stehen damit einzeln und sind der jeweiligen Sprache oft nur durch semantisches Verständnis ihrer Funktion bzw. nur durch Hinzuziehung distalen Kontexts ihrer Schriftfamilie zuzuordnen. Frakturtexte, die keine Zahlen enthalten, könnten aussichtsreicher mit anderen („reineren“) Mustersammlungen gelesen werden als Texte, in denen arabische wie römische Ziffern an unvorhersehbaren Orten stehen können; umgekehrt könnte oft eine textsorten- und layoutgesteuerte Erwartung römischer Zahlen deren richtige Erkennung überhaupt erst ermöglichen usw.

Der vierte wesentliche Informationsverlust der maschinellen OCR ergibt sich daher nicht aus Problemen in einem konkreten der oben genannten OCR-Teilschritte, sondern aus der bei Massen-OCR kaum vermeidbaren Anwendung pauschalisierender Einstellungen, wenn sie im Verarbeitungsstapel weder von Seite zu Seite noch von Textblock zu Textblock oder gar innerhalb einer Zeile bzw. eines Worts flexibel genug gewechselt werden können, um die Potentiale vorhandener Muster- und Sprachinformation tatsächlich wirksam werden zu lassen.

<sup>40</sup> Zumindest für *named entities* können aus der Computerlinguistik formale Kriterien bezogen werden; eine als Anregung lesbare Arbeit ist z. B. GEIERHOS 2007, wo über den Titel „Grammatik der Menschenbezeichner in biographischen Kontexten“ hinaus auch auf grammatische Merkmale weiterer Entitäten Bezug genommen wird.

<sup>41</sup> Beispielsweise als „Kolummentitel“ gelte wahlweise „die erste Zeile“ / „der Text oberhalb der ersten durchgezogenen Linie“ / „die erste Zeile, wenn ihr Abstand zur zweiten größer als NN mm ist“; oder: Als „Bibelstellenangabe“ gelte „jede Zeichenkette, die in Antiqua gesetzt ist, zum Regulären Ausdruck X passt und sich innerhalb eines Strukturelements Y befindet“ usw.

## Batchinput und –output

<b>Eingang:</b>	Menge an Images Musterbibliothek(en) Wortliste(n)
+ <i>Steuerdaten:</i>	Verarbeitungsparameter, Angaben zu Schriftfamilien, Alphabeten und Sprachen, Zielformate, Zielspeicherorte
<b>Ausgang:</b>	Menge an Exportdateien, ggf. in verschiedenen Formaten
+ <i>Evaluationsdaten:</i>	Protokolle beliebige statistische und computerlinguistische Auswertungen

## Typische Informationsverluste und Fehlerquellen

**Unvermeidlich:** Da diese Fehlerquellen allein durch die Verarbeitungsorganisation entstehen, sind sie nicht aus sich heraus unvermeidlich, sondern „nur“ in dem Maße, wie der nötige Aufwand zur Zusammenstellung homogener Stapel nicht geleistet werden kann.

**Potentiell beeinflussbar:** Verluste auf den jeweils schlechter zur aktuellen Konfiguration passenden Regionen/Seiten  
*Abhilfe ist z. B. zu erreichen durch:*

- die Zerlegung von Batches in kleinere;
- die Teilung der Seiten vorab in separat zu verarbeitende Images bzw. Image-Regionen;
- eine Mehrfachverarbeitung unter verschiedenen Bedingungen und den anschließenden Versuch, pro Textblock die im Nachhinein günstigste Ausgabe zu selektieren.

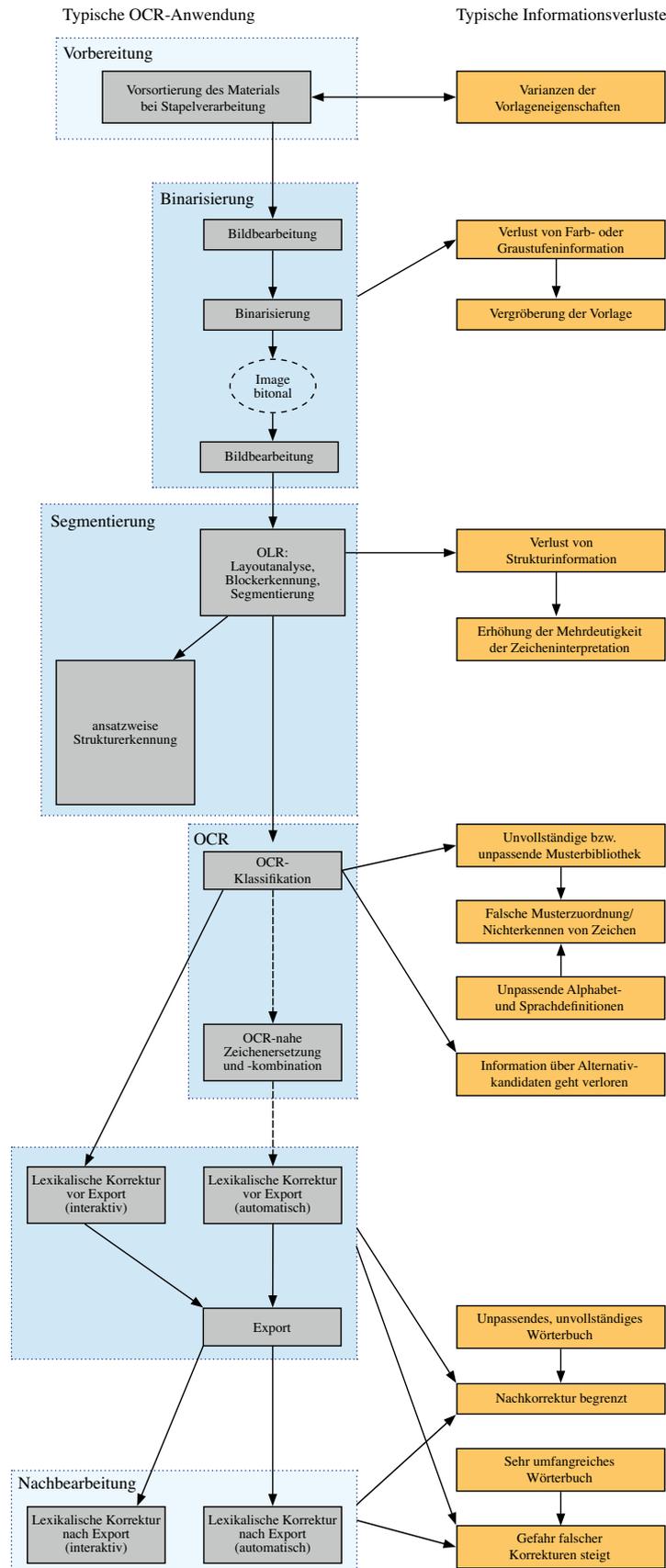
Diese Aufteilung kann idealerweise vor dem Batchlauf vorgenommen werden. Wenn das nicht praktikabel ist, könnten Seiten mit unnötig schlechter Erkennung auch nach dem Batchlauf einer Neubearbeitung mit anderen Parametern zugeführt werden. Software könnte dies unterstützen, wenn sie Evaluationskennwerte zur Ermittlung von zweckmäßiger Stapelaufteilung bereitstellt und Stapelprozesse mit wechselnden Parametern zulässt.

**Wann unschädlich:** Im einfachsten Fall verursacht die Gesamtseiten- und Seitenstapel-Verarbeitung dann keine eigenen Fehler, wenn innerhalb der Seiten bzw. des Stapels gar keine inhomogenen Situationen auftreten, die eine Mischung separater Muster- oder Lexikonbestände erfordern.

Im Fall von inhomogenen Seiten oder Seitenstapeln gilt:  
Bezogen auf die reine Zeichenerkennung können Misch-Musterbibliotheken dann unschädlich sein, wenn die verwendeten Drucktypen der beteiligten Schriftarten ausreichend disjunkt sind, d. h. sich immer nur da überlappen, wo dasselbe Zeichen gelesen werden soll.  
Analog können gemischte Wörterbücher dann unschädlich sein, wenn aufgrund konservativ gewählter Distanzmaße und Konfidenzwerte keine Gefahr besteht, fälschlich „in die falsche Sprache bzw. den falschen Wortbestand hinein“ zu korrigieren.

## 3.2 Zusammenfassung

Wenn man im Vorgriff auf die im nächsten Kapitel folgende Beschreibung der Software auch die Hilfsobjekte (Parameter-, Muster- und Wörterbuchdateien) heranzieht, mit denen die einzelnen Schritte determiniert werden, dann könnten die typischen Fehlerquellen und die wichtigsten vom Anwender zu beachtenden qualitätsbestimmenden Faktoren wie folgt in Abb. 40 dargestellt werden:



40 Workflow und beeinflussende Faktoren sowie resultierende Fehlerquellen

## 4 Praktischer Softwaretest und -vergleich I: Die Programme

### 4.1 Übersicht über Parameter- und Hilfsdateien jedes Produkts

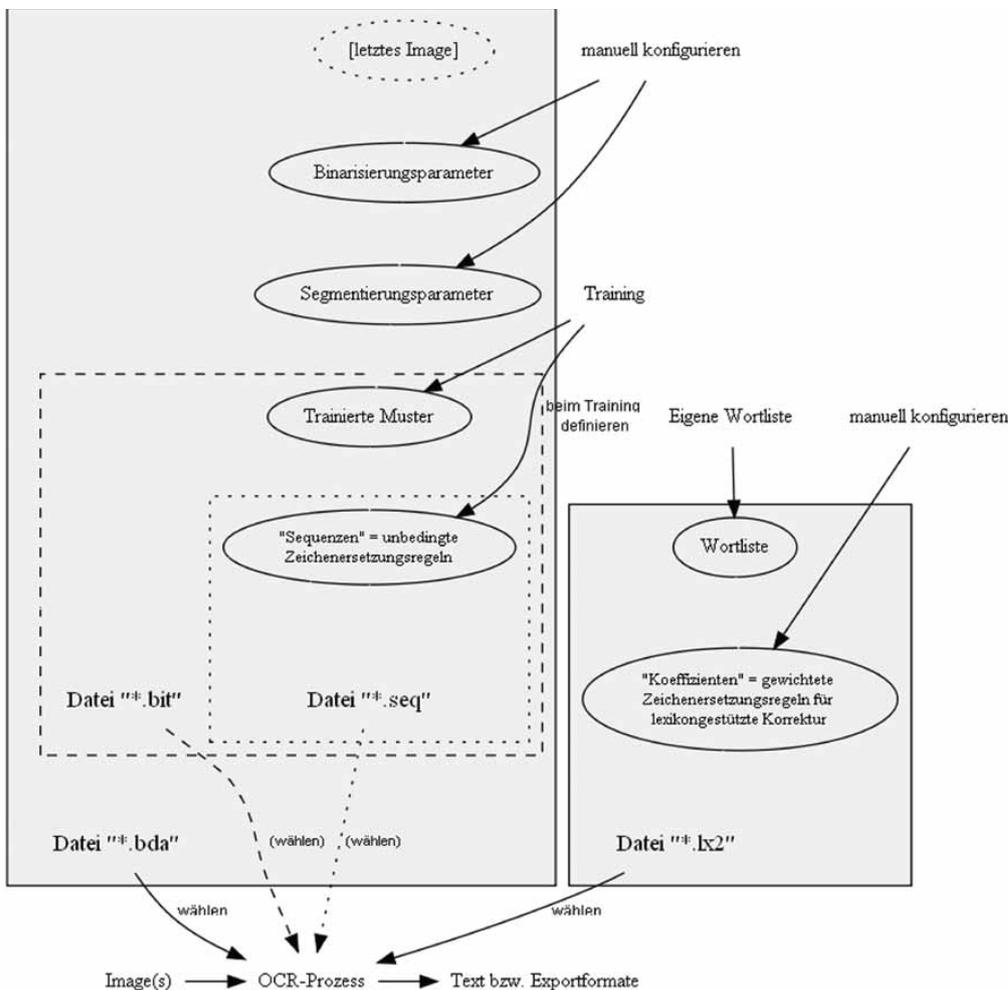
Der wesentliche Teil des Pilotprojekts bestand in der Untersuchung, wie ein Anwender mit Hilfe der jeweiligen Software zu den bestmöglichen Bedingungen eines OCR-Laufs gelangt. Dabei war der Frage nachzugehen, wie nachnutzbar diese sind, z. B. durch Fixierung erreichter Optimierungen in Parameter- bzw. Konfigurationsdateien bzw. in als Dateien speicherbaren Muster- und Wortbibliotheken. Entsprechend nimmt die Erstellung bzw. Beschaffung dieser Parameter- und Hilfsdateien einen großen Teil der (menschlichen) Bearbeitungszeit ein, ehe die eigentliche OCR-Produktion dann automatisch im Stapelbetrieb läuft und vorwiegend Rechenzeit beansprucht.

Zum besseren Verständnis der folgenden einzelschrittbezogenen Details folgt hier zunächst ein Überblick über die Speichermodelle der jeweiligen Software.

#### (a) BIT-Alpha<sup>42</sup>

Beteiligt sind pro OCR-Lauf zwei bis vier Parameterdateien, die in BIT-Alpha explizit als Datei geladen werden müssen und teilweise überlappende Funktionen haben:

- eine Konfigurationsdatei „bda“ für Binarisierungs- und Segmentierungsparameter, optional auch für Musterbibliothek und Sequenzen;<sup>43</sup>
- soweit Muster nicht in „bda“ bereitgestellt werden: ein separates Speicherformat „bit“ für die Musterbibliothek, optional auch für Sequenzen (s. u.);
- soweit Sequenzen nicht in „bda“ oder „bit“ bereitgestellt werden: ein separates Speicherformat „seq“ für Ersetzungsregeln (Sequenzen) zur sofortigen und lexikonunabhängigen Kombination bestimmter Zeichen(teile);
- eine Datei „lx2“ für Wortbibliothek und Parameter der lexikalischen Nachkorrektur.



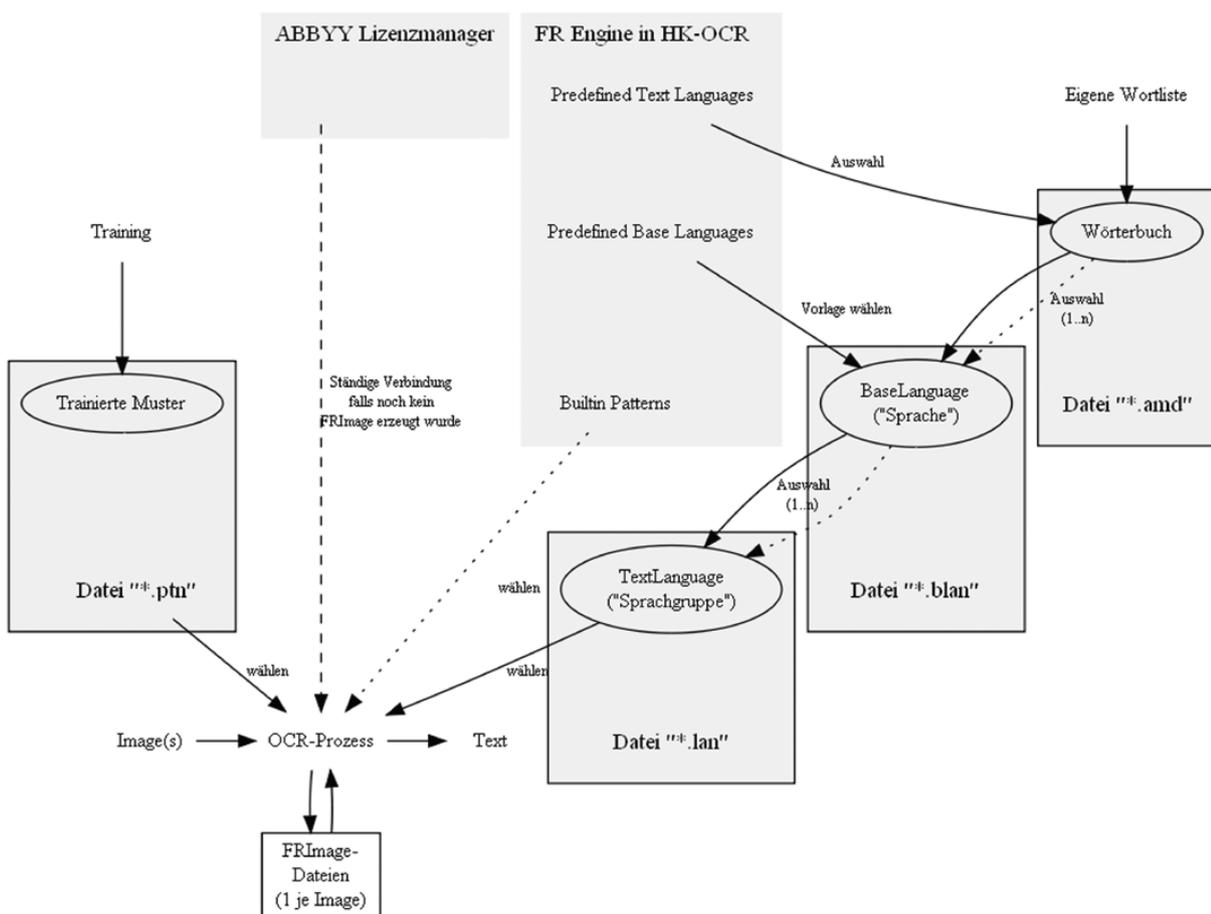
Für die Muster und Sequenzen kann während jeder BIT-Alpha-Sitzung neu entschieden werden, in welchem der alternativen Formate sie abgelegt werden; so würde auch eine Übertragung von Mustern zwischen bda- und bit- bzw. von Sequenzen zwischen bda-, bit- und seq-Dateien stattfinden.

## (b) HK-OCR / FREngine 9

Abgesehen von impliziten Programmkonfigurationsdateien<sup>44</sup> nutzt HK-OCR pro OCR-Lauf vier eigentliche Parameterdateien:

- drei Dateitypen zur Wahl einer eigenen „Sprachgruppe“: eine Wörterbuchdatei „.amd“, eine *Base Language* Datei „.blan“ und eine *Text Language* Datei „.lan“; hier können Wortlisten importiert und Listen erlaubter Zeichen bearbeitet werden. Zu diesen Dateitypen muss nur das Verzeichnis eingegeben werden (Reiter „Weitere“ im Einstellungsbereich von HK-OCR), dann kann im Reiter „Einstellungen“ die Auswahl einer „Sprachgruppe“ erfolgen.
- eine Musterdatei „.ptn“, die explizit als Datei ausgewählt und geladen wird.

Hinzu kann je Einzelbild eine sogenannte (während einer früheren OCR angelegte) FR-Image-Datei kommen, die Vollbild- und Layout-Informationen enthält und die Einzelseitenlizenz verkörpert. Liegt diese FR-Image-Datei für das zu lesende Bild nicht vor, dann wird sie bei der OCR automatisch angelegt; wenn sie bereits im Bildstapelverzeichnis vorliegt, dann ist eine erneute OCR auch ohne weiteren ABBYY-Lizenzverbrauch möglich.<sup>45</sup>



42 Parameter- und Hilfsdateien HK-OCR/FREngine 9

<sup>42</sup> Alle folgenden Anmerkungen beziehen sich auf die im Projekt getesteten Versionen; zuletzt BIT-Alpha 2.0.38.594.

<sup>43</sup> Zur Funktionsweise bei Zeichenkombinationen s. 4.3.3

<sup>44</sup> HK-OCR verwaltet ohne Benutzereingriff im Programmverzeichnis und im Nutzer-Anwendungsverzeichnis je eine XML-Konfigurationsdatei, wo aktuelle Programm-(Session-)Einstellungen gespeichert werden, die aber im gegenwärtigen Design nicht zur Speicherung oder gar Änderung verschiedener später wählbarer Profileinstellungen vorgesehen sind, sondern lediglich – ähnlich wie bei Registry-Einträgen – beim Neustart einige Pfade und Optionen der letzten Sitzung wieder bereitstellen.

<sup>45</sup> Das Lizenzmodell basiert auf der seitenweisen Abrechnung, wobei auch Images Lizenzen verbrauchen, die keine Texte enthalten (z. B. Leerseiten). Die Prozessierung übergroßer Seiten (Folianten) mit viel Text benötigt unter Umständen zwei Lizenzen.

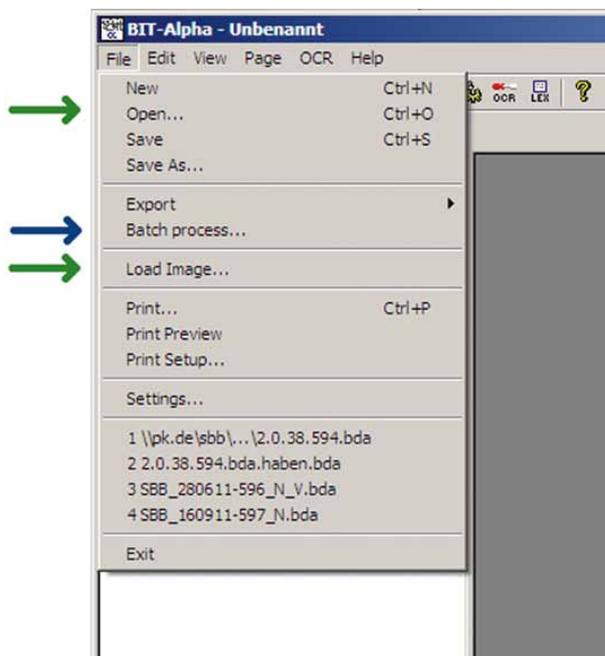
**Desiderata**

Für alle diese Hilfs- und Parameterdateien gilt, dass sie bisher in proprietären Binärformaten vorliegen. Für eine effektive Protokollierung, den Vergleich und die auszugsweise Übernahme von einzelnen Parametern, Mustermengen, Ersetzungsregeln oder Wortbeständen wäre es wünschenswert, dass diese Dateien jeweils

- auch außerhalb eines OCR-Prozesses mindestens passiv auslesbar sind;
- insbesondere miteinander verglichen werden können („In welchen Parametern haben sich OCR-Lauf A und OCR-Lauf B unterschieden?“);
- eine leichte Übernahme von (Teil-)Inhalten einer Konfigurationsdatei in eine andere ermöglichen (Ersatz, Kombination usw.).

Hierzu eignen sich naturgemäß am besten konventionelle offene Textformate wie Config- bzw. INI-Dateien, auch XML, die sowohl von Menschen als auch maschinell gelesen und bearbeitet werden können, die jederzeit mit Standardtextfiltern auswertbar sind und die somit auch von einem beliebigen einbettenden Workflow aus verwaltet werden können.

**4.2 Start der OCR-Verarbeitungskette in der Benutzeroberfläche**



43 Arbeitsbeginn mit BIT-Alpha

**(a) BIT-Alpha**

Die OCR einer Einzelseite startet selbsttätig,

- wenn die bda-Konfigurationsdatei ein gespeichertes Image enthält, und zwar bereits beim Öffnen der bda-Datei<sup>46</sup> oder
- beim Laden eines Bilds („File > Load Image“)<sup>47</sup>.

In der Ansicht (Layer) „OCR“ wird angezeigt, ob und (durch Farben unterschieden) mit welcher Sicherheit BIT-Alpha die Zeichen erkannt hat.

Der Batch-Lauf wird (nach Öffnen/Laden der Konfigurationsdateien und/oder Einstellung aller gewünschten Parameter) angestoßen im Menü „File > Batch process ...“, wo Quell- und Zielverzeichnisse sowie die Exportformate eingestellt werden können, per Button „OK“.

**(b) HK-OCR / FREngine 9**

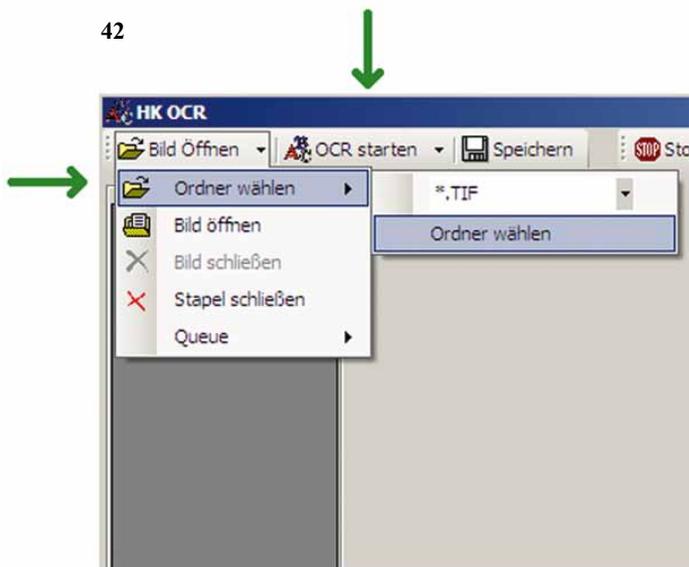
Nachdem die Grundeinstellungen eingerichtet sind (einmalige Angabe von Pfaden im Reiter „Weiter“ im Einstellungsbereich von HK-OCR<sup>48</sup>), wird für den aktuellen Lauf der Bild-Quellordner gewählt („Bild öffnen > Ordner wählen > Ordner wählen ...“) und werden in den Reitern „Einstellungen“ und „Export“ die aktuell gewünschten Einstellungen vorgenommen.

Danach kann die OCR einer Einzelseite gestartet werden durch „OCR starten > aktuelles Bild lesen“. Ein Batch-Lauf wird analog ausgelöst durch „OCR starten > alle Bilder lesen“.

<sup>46</sup> In der getesteten Version von BIT-Alpha war das automatische Abspeichern des letzten Images Standardverhalten und nicht ganz glücklich, da jedes neue Öffnen einer bda-Datei dann aufgrund des damit automatisch beginnenden OCR-Laufs viel Zeit verbraucht. Verhindert werden konnte die Speicherung des je aktuellen Bilds durch explizites Entfernen („Page > Remove“) vor dem Speichern der bda-Datei.

<sup>47</sup> Auch dies hat Nachteile: ein Bearbeiter könnte z.B. nach dem „Laden“ des Bilds noch Einstellungen vornehmen wollen, ehe die OCR beginnt.

<sup>48</sup> Wird über verschiedene Profile (ohne Administrationsrechte) auf die Software zugegriffen, ist darauf zu achten, dass in den profilbezogenen Anwendungsdaten eine vollständige alle Pfade enthaltende Konfigurationsdatei (user.config) liegt. Andernfalls öffnet sich das Programmfenster nicht und es können keinerlei Einstellungen vorgenommen werden.



44 Arbeitsbeginn mit HK-OCR/FREngine 9

## Desiderata

Aus dieser Benutzerführung ergeben sich vor allem einige unnötige Beschränkungen des Batch-Betriebs. So müssen gegenwärtig alle verarbeiteten Bilder in demselben Verzeichnis liegen (das bei HK-OCR zugleich die Exporte aufnehmen muss), was dazu zwingt, große Mengen Bilddaten unnötig zu bewegen. Bildverzeichnisse werden komplett verarbeitet. Und für den gesamten Batchlauf bleibt die Konfiguration starr.

Im Abschnitt 6.4 wird darauf eingegangen, wie deutlich flexiblere Stapelverarbeitungen möglich wären, wenn die Programme die Parameter, die ihnen ohnehin übergeben werden müssen (Namen von Verzeichnissen und Konfigurationsdateien usw., die derzeit einmalig per GUI eingestellt werden und daher starr für den gesamten Batch-Lauf gelten), auch als konventionelle Kommandozeilenparameter akzeptieren würden.

## 4.3 Einzelne OCR-Verarbeitungsstufen (ohne Trainings-Phase)

Jede OCR setzt gespeicherte Repräsentationen optischer Muster voraus, anhand derer die Drucktypen klassifiziert und einem Zeichen bzw. einer Zeichenfolge zugeordnet werden. Die Beschaffung (Training, Bearbeitung, Kombination) der Musterbibliothek geht dem produktiven Einsatz i. d. R. voraus, wird aber in einem späteren Abschnitt dieses Kapitels (s. Kap. 4.4) beschrieben.

### 4.3.1 Bildvorverarbeitung und Binarisierung

Beide Produkte kommen am besten mit Bildern zurecht, die bereits mit normalisierten Helligkeits- und Kontrastwerten, ohne Ränder und Verzerrungen, in Stapeln mit vergleichbaren Weißraum-Maßen vorliegen usw. Dennoch nehmen beide Programme auch selbst etliche Bildbearbeitungen vor, was bei HK-OCR weitgehend automatisch in der eingebetteten FineReader-Engine geschieht und bei BIT-Alpha detailliert konfiguriert werden kann und muss.

#### (a) BIT-Alpha

BIT-Alpha bietet eine Vielzahl von die Binarisierung beeinflussenden Einstellungen, die ein großes Potential zur Anpassung an spezielle Vorlageneigenschaften bieten, aber auf eine schwer durchschaubare Weise komplex zusammenwirken.

Die Einstellungen werden in einer Konfigurationsdatei im proprietären Format „bda“ gespeichert (s. Abb. 41 am Anfang des Kapitels). Die Arbeit in BIT-Alpha beginnt folglich mit dem Anlegen oder der Auswahl solch einer Konfigurationsdatei. Das initiale Programm-Menü „File > Open“<sup>49</sup> bezieht sich hierauf (nicht etwa auf ein zu lesendes Image); entsprechend bedeutet „File > Save“ das Abspeichern der aktuell im Programm BIT-Alpha geltenden Parameter in einer solchen bda-Konfigurationsdatei. Images werden dann zur Verarbeitung hinzugeladen: bei Einzelverarbeitung per „File > Load Image“, zur Stapelverarbeitung über „File > Batch process ...“ und die Angabe des Quellverzeichnisses.

Neben Einstellungen zur Bildvorverarbeitung wie Helligkeits-, Kontrast- und Gammakorrekturreglern, der Behandlung von Umrandungen usw. kann insbesondere sowohl auf Seitenebene (erste Binarisierung für Segmentierungszwecke) als auch auf der Ebene von „Regionen“ (zweite Binarisierung für Zeichenerkennung) zwischen verschiedenen Binarisierungsalgorithmen gewählt werden. Zur Auswahl stehen ein auf Farbanalyse, ein auf Farbintensitätsschwellen beruhender, ein modifizierter Niblack- und ein von B.I.T. selbst entwickelter Algorithmus. B.I.T. empfiehlt für Antiqua den eigenen Algorithmus, für Frakturschrift dagegen den Niblack-Algorithmus, der besonders die lokalen Helligkeitsunterschiede um die Zeichenkonturen herum auswertet.

Mit den von BIT-Alpha angebotenen Vorschau-Ansichten („Layern“) kann die Auswirkung aktueller Einstellungen auf die Binärbildqualität einseitig unmittelbar visuell beurteilt werden. Für den produktiven Einsatz fehlt hier eine Mehr-Seiten-Vorschau, um vorab zu prüfen, wie die anhand einer aktuellen Seite vorgenommenen Parameteränderungen sich auf andere Seiten desselben Stapels auswirken werden. Diese oft erheblichen Nebenwirkungen stellen sich daher häufig erst nach dem OCR-Lauf heraus.

Die binarisierten Seiten-Images können exportiert werden, aber erst beim OCR-Lauf. Nach Bedarf wurden in verschiedenen Projektphasen veränderte Binarisierungsoptionen angewandt. Die insgesamt sehr komplexe Justierung auf bestimmte Vorlageneigenschaften hin wurde meist von B.I.T. übernommen. Zu einigen Parametern sind typische Wertebereiche im Anhang aufgeführt (s. Anh. 7.3). Ein den Nutzer zur eigenständigen Konfiguration befähigendes Handbuch stand bis zum Projektende nicht zur Verfügung.

## **(b) HK-OCR / FREngine 9**

Eingriffsmöglichkeiten für den Anwender fehlen fast völlig; die Bildnormalisierung, Entzerrung und Binarisierung laufen eingekapselt in der FineReader-Engine. Man hat lediglich die Möglichkeit, die automatische Dreh-Ausrichtung der Seite ein- oder auszuschalten. Tests der Auswirkung verschiedener Binarisierungsoptionen konnten daher nicht gemacht werden. Entsprechend gibt es keine Konfigurationsdatei für etwaige Vorverarbeitungs- oder Binarisierungsparameter. Spezielle Probleme wurden mit dem vorliegenden Material dennoch nicht beobachtet, offenbar sind die werkseitigen Einstellungen sehr gut für die meisten Anwendungsfälle geeignet. Ein Export der binarisierten Bilder war in HK-OCR nicht möglich.

## **Desiderata, Kriterien für Alternativen**

### ***Modularität***

Bisher bietet keines der beiden Programme eine Ein- und Ausgabeschnittstelle nach der Binarisierung an.<sup>50</sup> Für eine modulare Verwendbarkeit der Produkte, d. h. entweder das in einem Programm X binarisierte Bild in anderen Workflows weiterzuverarbeiten oder ein bereits binarisiertes Bild direkt in die Weiterverarbeitung mit dem Programm X einzuspeisen, wäre eine Schnittstelle wünschenswert.<sup>51</sup>

<sup>49</sup> Dem entspricht ein Aufruf einer „bda“-Datei im Windows-Explorer, wenn bei der Installation unter Windows ein Dateityp „BIT-Alpha.Document“ mit der Extension „bda“ und BIT-Alpha als Standardanwendung registriert wurde.

<sup>50</sup> Die in BIT-Alpha gegebene Möglichkeit eines Exports des binarisierten Bilds mit der OCR bedeutet zumindest immer einen zeitlichen Mehraufwand.

<sup>51</sup> Im bis 2011 laufenden IMPACT-Projekt wurden u. a. Ein- und Ausgabeanforderungen verschiedener OCR-Werkzeuge beschrieben, um deren Kombinierbarkeit systematisch darzustellen; das 2012 als Nachfolgeinstitution entstandene Kompetenzzentrum kann als Kommunikationsvermittler in Anspruch genommen werden.

## Evaluation

Eine schnelle Vorschau der Binarisierungsqualität mehrerer Seiten (z. B. durch Vorab-Export einer Bildmenge ohne zeitaufwendigen OCR-Lauf) könnte den Zyklus der Parameterbeurteilung vor dem OCR-Lauf wesentlich verkürzen. (s. Kap. 6.4)

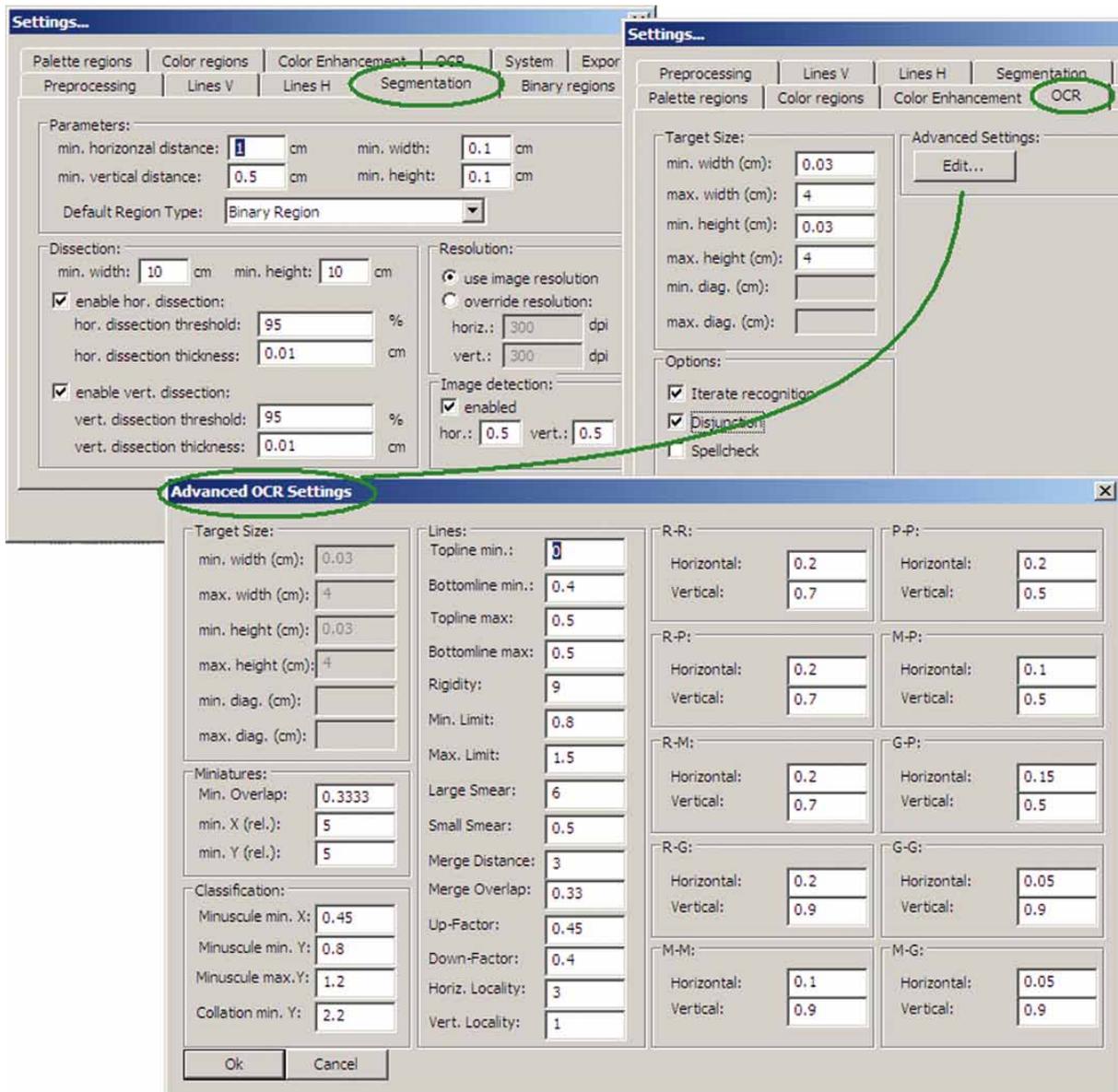
## Benutzerschnittstelle

Die in BIT-Alpha fein konfigurierbaren Möglichkeiten werden in Anwenderhand erst dann wirksam eingesetzt werden können, wenn herstellerseitig eine entsprechende Dokumentation und eine Art Leitfaden zur Parametrierung bereitgestellt wird.

Von HK-OCR (bzw. der eingebauten FineReader-Engine) wäre dagegen zu fordern, überhaupt mehr Konfigurationsmöglichkeiten für den Anwender zu öffnen, die intern durch Abbyy – wie die im Durchschnitt guten Ergebnisse zeigen – durchaus genutzt und beherrscht werden, und die in FineReader-Consumer-Produkten gelegentlich auch nutzerkonfigurierbar sind. Alle in der FineReader-Engine potentiell einstellbaren Optionen<sup>52</sup> sollten in der Programmoberfläche zugänglich gemacht werden.

## 4.3.2 Segmentierung

### (a) BIT-Alpha



Es gilt die gleiche Anmerkung zur Parametervielfalt wie bei der Binarisierung.

Eingestellt werden können u. a. minimale und maximale Distanzen zwischen Textblöcken und Einzelzeichen, die Behandlung von Linien, Parameter der „elastischen“ Zeilenerkennung bei gewölbten Aufnahmen, Erkennung von Nicht-Text-Regionen (Abbildungen, Ornamentik) und etliches mehr.

Die aktuell erreichte Segmentierung kann wieder mit verschiedenen Vorschau-Ansichten beurteilt werden; allerdings ist die Schlussfolgerung, wann welche Parameter geändert werden sollten, sehr kompliziert und bisher nicht in einem Handbuch dokumentiert.

Ebenfalls eingestellt werden können Parameter, die später die OCR steuern wie z. B. Größe und Abstand als irrelevant zu entfernender Punkte, der Größenbereich der zu erkennenden „Targets“, sogar relative Größenverhältnisse der Drucktypen wie Ober-, Mittel- und Untertlänge der Zeichen u. a. m. Alle diese Einstellungen werden mit in der bda-Konfigurationsdatei, die auch die Binarisierungsparameter enthält, gespeichert bzw. daraus gelesen.

### **(b) HK-OCR / FREngine 9**

Wie die Binarisierung verläuft auch die Segmentierung in der getesteten HK-OCR-Version vollautomatisch ohne nutzerseitige Eingriffsmöglichkeiten.

HK-OCR sieht vor dem OCR-Start offenbar keine Ansicht und Beurteilung der entstandenen Bildblöcke vor; sie können dennoch über einen Button der Korrektur-Einstellungen (Reiter „Korrektur“ > „Anzeige“ > „Blöcke“) angezeigt werden.

Layouterkennungs-Ergebnisse dürften mit in die bei der OCR pro Bild entstehenden binären Layout- und FR-Image-Dateien eingehen. Für den Anwender sind diese Informationen nicht auswertbar.<sup>53</sup>

### **Desiderata, Kriterien für Alternativen**

#### ***Modularität***

In beiden Produkten ist weder die Eingabe fremder binarisierter Bilder zur Segmentierung noch die Ausgabe von Segmentierungsergebnissen vor der OCR vorgesehen.

#### ***Evaluation***

Für den Anwender auf dieser Stufe zugängliche Layout-Ergebnisse (Blockkoordinaten o.ä.) gibt es in keinem der beiden Programme.

#### ***Benutzerschnittstelle***

Für BIT-Alpha gilt die bereits im Binarisierungsschritt gemachte Anmerkung zum Erfordernis eines den Nutzer unterstützenden Handbuchs. Für etliche Parameter, die geometrische Größen wie Abstandsmaße, Höchst- und Mindestgrößen beschreiben, könnte alternativ zur numerischen Eingabe eine graphische Eingabe vor dem Hintergrund einer binarisierten Originalseite angeboten werden.

Für HK-OCR gilt hier noch mehr als im vorangehenden Binarisierungsschritt, dass alle in der Fine-Reader-Engine potentiell einstellbaren Optionen in der Programmoberfläche auch zugänglich gemacht werden sollten. Erkennungsrelevante Eigenschaften wie Satzspiegel, Spaltenerkennung, Ausblenden von Ornamentregionen, zulässige Zeichenabstände und Wortzwischenräume usw. sollten vor der OCR eingestellt werden können. Blockgrenzen können in der derzeitigen HK-OCR-Version nicht verschoben werden.<sup>54</sup>

<sup>52</sup> Im Projekt war keine Gelegenheit, sich näher mit ABBYY-Entwicklerhandbüchern zu beschäftigen, dort scheinen aber durchaus Möglichkeiten angelegt zu sein.

<sup>53</sup> Probleme ergaben sich bei der Bearbeitung von Vorlagen, die über die eigentlich zu betrachtende Seite hinaus auch noch „Textreste“ der gegenüberliegenden Seite enthielten, die nach Möglichkeit ausgeblendet werden sollten.

<sup>54</sup> Während im späteren Nachkorrekturschritt („Erweiterte Validierung“) ein Verschieben der Blockgrenzen durchaus angeboten wird.

### 4.3.3 OCR

Optische Zeichenerkennung klassifiziert bildliche Darstellungen und ordnet sie einem Zeichen bzw. einer Zeichenfolge zu. Im Softwarevergleich spielte keine Rolle, wie diese Musterklassifikation im einzelnen mathematisch realisiert wurde (ob als Neuronale Netze oder mittels vorgegebener Ähnlichkeitsfunktionen über abstrakte optische Grundmuster oder Merkmalsvektoren usw.), da herstellerseitig kaum Informationen hierzu gegeben wurden und Qualitätsunterschiede absehbar nicht aus dem je verwendeten Paradigma resultieren würden. BIT-Alpha speichert offenbar verschiedene aus dem optischen Bild extrahierte Kennwerte, während FineReader vermutlich die binären Zeichen-Images selbst in der Musterdatei speichert. Für den Anwender relevant ist, dass BIT-Alpha infolgedessen die gelernten Muster-Images aus der gespeicherten Musterbibliothek nicht wiederherstellen kann.<sup>55</sup>

Hier wird der OCR-Schritt unter der Voraussetzung vorhandener Musterbibliotheken beschrieben; ihre Beschaffung beschreibt ein späterer Abschnitt (s. Kap. 4.4).

#### (a) BIT-Alpha

Vorgegebene Muster sind in BIT-Alpha nicht enthalten. Um überhaupt Texterkennung zu beginnen, muss man daher eine bereits trainierte Musterbibliothek verwenden. Diese Muster können entweder mit in einer bda-Konfigurationsdatei gespeichert sein oder in einer separaten Musterdatei im Format „bit“ vorliegen. Wenn in BIT-Alpha sowohl (per „File > Open“) eine bda-Konfigurationsdatei als auch (per „OCR > Library > Load“) eine bit-Musterdatei geladen sind, dann gelten nur die Muster der bit-Musterdatei. Aus der bda-Konfigurationsdatei sind dann nur die Binarisierungs- und Segmentierungsparameter wirksam.

Wenn die Mustersammlung wie von B.I.T. empfohlen trainiert wurde, sind i. d. R. so genannte „Sequenzen“ für regelmäßige einfache wortunabhängige Zeichenersetzungen angelegt worden. Diese können wahlweise sogar an drei verschiedenen Orten gespeichert werden: wie die Muster entweder in der bda-Konfigurationsdatei oder in der bit-Musterdatei, oder separat in einer eigenen seq-Sequenzdatei (dann zu laden über „OCR > Sequences > Load“)(s. Abb. 41 am Anfang des Kapitels). In der Praxis eignen sich „Sequenzen“ besonders für die Zusammensetzung von Buchstabenteilen<sup>56</sup> oder auch zur Auflösung von Abkürzungen, die aus mehreren gültigen Zeichen bzw. Special Values (Sonderzeichen) bestehen.<sup>57</sup>

Weitere die OCR-Stufe bestimmende Parameter können in den Menüs „File > Settings > OCR“ und „OCR > Library > Settings“ eingestellt werden.<sup>58</sup> Hierzu gehört die Möglichkeit, den Erkennungslauf zu wiederholen sowie ein automatischer Versuch, abgetrennte Teile von Zeichenbildern alternativ als eigenes Muster zu lesen. Des Weiteren sind etliche Größenverhältnisse für die Zeichenklassifikation konfigurierbar.<sup>59</sup>

#### (b) HK-OCR / FREngine 9

In der HK-OCR-Oberfläche wird die OCR gesteuert durch die

- Wahl einer Musterbibliothek und/oder des eingebauten FineReader-Musterbestands;
- Wahl von ein bis drei Schriftfamilien (Fraktur/Antiqua/Typewriter);
- Wahl einer oder mehrerer vorgegebener „Sprachen“ oder selbstdefinierter „Sprachgruppen“ mit Wort- und Zeichenbeständen.

<sup>55</sup> Es ist aber beim OCR-Lauf von Einzelseiten ein (manueller) Export der bitonalen Zeichen-Images möglich; die einzelnen Image-Dateien werden dabei getrennt in Verzeichnissen gespeichert, die nach den zugeordneten Strings benannt sind, so dass außerhalb der Musterbibliothek die Bild-Zeichen-Zuordnung nachvollzogen werden kann.

<sup>56</sup> Wenn z. B. die Zeichensegmentierung dazu neigt, in blasserem Druck Beine der Kleinbuchstaben „m“, „n“, „u“ zu separieren oder „w“ in „v“+„i“ aufzuteilen und dies in der Parameterkonfiguration nur um den Preis einer generell schlechteren Zeichenseparierung zu ändern wäre, dann können beispielsweise die einzelnen Beine als eigene Muster („special values“) trainiert werden und durch Sequenzen zusammengefügt werden. Sequenzen sind daher oft eng auf die trainierten Muster bezogen, d. h. andere Musterbibliotheken verlangen oft andere Sequenzen.

<sup>57</sup> Es bot sich z. B. an, das „runde r“ als *special value* zu codieren, um es im *sequencer* mit folgendem „c.“ zu „etc.“ aufzulösen, sonst aber zu „r“. Historische Abkürzungen können gut als eigene Muster trainiert werden.

<sup>58</sup> Erstere werden wie Binarisierungsoptionen in bda-Dateien gespeichert; letztere wie Muster in bda- oder bit-Dateien.

<sup>59</sup> Vergleiche in der Parameterübersicht in 7.3 folgende Kategorien: Target Size (minimale und maximale akzeptierte Zeichengröße); Miniatures (zur Erkennung von Initialen); Classification (je nach Variabilität der Schriftform relativ zur Minuskelhöhe einstellbare Größenverhältnisse von: [P]latte (senkrecht Bein bzw. kleiner senkrechter Strich, z. B. i ohne Punkt), [M]inuscule (Kleinbuchstabe), [G]rand (Großbuchstabe bzw. großer Strich), [R]este (kleineres Element als Patte, Minuscule oder Grand, z. B. „“ oder „“)). B.I.T. weist darauf hin, dass die Parametrierung dieses Menüs grundsätzlich von B.I.T. selbst und nicht vom Benutzer vorgenommen wird. Die Einstellungen sind für Fraktur, Kursivschrift der Renaissance oder gerade Antiqua unterschiedlich.

Die FineReader-Engine bringt einen Grundbestand an Mustern (bezeichnet als „Builtin-Patterns“) auch für Fraktur mit, man kann also bereits Texterkennung realisieren, bevor eine Musterdatei geladen wird. Für Mainstream-Frakturschriften besonders des 19. Jahrhunderts sind diese *built-in patterns* oft völlig ausreichend und es ist durch eigene Muster keine Verbesserung der Erkennung mehr zu erwarten. Anders gesagt: Es muss nur dann eine eigene Mustersammlung eingesetzt werden, wenn – wie bei Alten Drucken oder bei speziell gestalteten Lettern – die *built-in patterns* den Drucktypen der Vorlagen nicht gut genug entsprechen.

Trainierte Muster sind in Musterdateien im Finereader-Format „.ptn“ gespeichert (s. Abb. 42 am Anfang des Kapitels). Genau eine solche Musterdatei kann in HK-OCR im Reiter „Einstellungen“ ausgewählt und in der Checkbox „Musterdatei verwenden“ aktiviert werden – wahlweise allein ohne die FineReader-Muster (dann Checkbox „interne Muster verwenden“ deaktivieren) oder mit ihnen zusammen (Checkbox „interne Muster verwenden“ aktivieren).<sup>60</sup> Ein Mischen bzw. Zusammenfügen mehrerer in .ptn-Dateien gespeicherter Musterbestände ist in HK-OCR nicht möglich.

Die Auswahl der Schriftfamilie(n) erfolgt über Checkboxes „Normal“, „Fraktur“ und „Schreibmaschine“. Es ist nicht erkennbar, ob dahinter nur eine Kombination von Muster-Teilbeständen steckt<sup>61</sup> oder ob weitere interne Parameter betroffen sind.

Die Wahl einer „Sprache“ bzw. „Sprachgruppe“ bestimmt vor allem die Lexik und den Satz zulässiger Zeichen. Auch wenn eine „Sprache“ „OldGerman“ zur Ausstattung gehört, war die Erkennungsleistung bei Einsatz einer speziell auf den Erscheinungszeitraum (17. Jahrhundert) und die Gattung bezogenen Wortliste besser. Hierzu musste aufbauend auf der vorgegebenen Sprache „OldGerman“ eine eigene „Sprache“ (base language) und „Sprachgruppe“ (text language) angelegt werden, in diese konnte dann der passende Wortbestand importiert werden.

Beobachtungen ließen vermuten, dass der Wortbestand recht früh bei der OCR in die Zeichenerkennung eingeht; mangels exportierbarer Zwischenstufen konnte das nicht genauer untersucht werden. Auch bei vermeintlicher Benutzung ausschließlich des eigenen Wortbestands gab es gelegentlich Erkennungsleistungen, die am ehesten mit weiterhin aktiven FineReader-Wortlisten erklärbar wären. Auch das Verhalten hinsichtlich der prinzipiell in der Konfiguration angebbaren erlaubten und unerlaubten Zeichen einer „Sprache“ ist nicht transparent genug bzw. wurde im Pilotprojekt nicht vollständig verstanden.<sup>62</sup>

Im OCR-Lauf werden von der FineReader-Engine proprietäre sogenannte Layout-Dateien angelegt, die offenbar variierenden Parametern des vorzunehmenden OCR-Laufs entsprechen. Neben diesem binären Format entstehen kurze, als Text lesbare Protokolldateien „layout.txt“ mit einigen statistischen Kenndaten zum OCR-Lauf des jeweiligen *layouts*. Es können pro Image mehrfache *layouts* angelegt werden, auf die dann wahlweise zurückgegriffen werden kann. HK-OCR lässt auch eine automatische Wahl des – nach einer internen Statistik vermeintlich erkannter Zeichen – „besten“ *layouts* zu (Reiter „Einstellungen“, Feld „Qualitätsoptimierung“). Man kann dadurch offenbar einige OCR-Parameter variieren, ohne allerdings zu wissen, welche und in welcher Richtung.

Eine Besonderheit der FineReader-Engine ist die bei der OCR ebenfalls entstehende FR-Image-Datei, die offenbar die volle (Farb-)Bildinformation und einige Zusatzinformationen enthält.<sup>63</sup> Im seitenbezogenen ABBYY-Lizenzmodell verbraucht die Erstellung jeder FR-Image-Datei eine Lizenzeinheit. Diese FR-Image-Datei kann später für erneute OCR-Läufe derselben Seite wieder verwendet werden, ohne dass dann eine neue Lizenzeinheit benötigt wird.<sup>64</sup>

<sup>60</sup> Achtung: Ohne Aktivierung der eigenen Muster (Checkbox „Musterdatei verwenden“) erscheint zwar die Checkbox „interne Muster verwenden“ inaktiv (ausgegraut), die internen Muster sind dann aber automatisch aktiv!

<sup>61</sup> Dazu passt nicht, dass auch bei alleiniger Nutzung einer eigenen (zwangsläufig ungeteilten) Musterbibliothek ein Unterschied der Fraktur-Erkennungsrate bei Zuwahl der Antiqua („Normal“)-Schrift beobachtet wurde. Allerdings schien die „Deaktivierung“ der *built-in patterns* ohnehin nicht vollständig.

<sup>62</sup> Über die Zeichen-Listen in den Menüs „Extras > Sprachen editieren > eigene Sprachen > unzulässige Zeichen“ und „Extras > Sprachen editieren > eigene Base.Languages > Erkennungssatz“ gelang es nicht immer, die zulässigen und nicht zulässigen Zeichen wie erwartet zu steuern. Bei einem offenen Textformat der Sprachdefinitionen würde dieses Problem vermutlich gar nicht erst auftreten oder wenigstens leicht identifizierbar sein.

<sup>63</sup> Entsprechend liegt deren Dateigröße ein wenig über der Dateigröße der originalen, nicht binarisierten TIFF-Bilder, so dass im Dateisystem des Bilddatenverzeichnisses entsprechend Platz vorgehalten werden muss.

<sup>64</sup> Zu beachten ist allerdings, dass spätere Wiederholungen des OCR-Gangs zum Teil nur bei identischen Pfaden (absoluten Speicherorten) der Muster- und Sprach-Dateien möglich sind, und dass die Wiederverwendbarkeit nur für dieselbe Version der FR-Engine zugesichert wird.

## Desiderata, Kriterien für Alternativen

Insbesondere sollte die jederzeitige Kombinierbarkeit separater Muster-Teilbestände ermöglicht bzw. konsequenter angeboten werden. Eine gute Anleitung und Dokumentation der Wirkung jedes einstellbaren Parameters wäre die Voraussetzung, die Möglichkeiten auszuschöpfen. Hinsichtlich der Benutzerschnittstelle konzentrieren sich die Anforderungen an die Programmoberfläche auf die Trainingssituation (s. Kap. 4.4).

### 4.3.4 Lexikalische Nachkorrektur

Idealerweise wird eine Wortliste aller in der Vorlage vorkommenden Wörter und Wortformen benötigt, die alle in der Vorlage nicht vorkommenden Formen nicht enthält. Strikt zu vermeiden sind falsche Wortformen, da sie nicht nur die Akzeptanz von Lesefehlern erhöhen, sondern auch Fehlkorrekturen provozieren.<sup>65</sup> Ähnlich problematisch sind Wortfragmente, wie sie aus Trennungen am Zeilenende resultieren können – es bedarf gründlicher Überlegung, ob und wann solche Wortteile ins Korrekturlexikon aufgenommen werden sollten.

Neben den in Kapitel 3 genannten grundsätzlichen Kriterien kann die Eignung von Nachkorrekturschritten davon abhängen, inwieweit Wort- und Zeichenkoordinaten an den von der Korrektur betroffenen Stellen mitgepflegt werden oder nicht, bzw. ob die Korrektur überhaupt erst auf abgeleiteten Formaten erfolgt. Allerdings sind korrekturbedingte Koordinaten-Ungenauigkeiten im praktischen Gebrauch (Suchanfragen, Finden der Stelle im Faksimile) meist unproblematisch, da auch die korrigierten Zeichen bzw. Wörter sich weiter in unmittelbarer Nähe der nicht mehr ganz zutreffenden Koordinaten befinden.

Eine linguistische Korrektur ist in beiden Produkten nicht enthalten, bzw. es müssten linguistische Kriterien vorab bei der Erstellung der Wortformenliste berücksichtigt werden.

#### (a) BIT-Alpha

##### ***Vor endgültiger Koordinatenspeicherung: automatische lexikalische Korrektur***

Der OCR-Lauf erfolgt mit automatischer lexikalischer Korrektur, wenn in „File > Settings > OCR“ die Checkbox „Spellcheck“ aktiviert und per „OCR > Lexical Correction > Load“ eine lx2-Datei (s. Abb. 41 am Anfang des Kapitels) mit Wortliste und Ersetzungsparametern geladen wird. Das Bearbeitungsfenster für diese Einstellungen erreicht man im Menü „OCR > Lexical Correction > Edit“.

Sowohl in die Pflege der Wortliste als auch besonders in die Justierung der Ersetzungskoeffizienten lohnt es sich zu investieren, da hierdurch die blinde Ersetzung unbekannter Wörter durch *irgendein* bekanntes Wort auf die wahrscheinlicheren Alternativen hin gelenkt werden kann.

Die Wortliste sollte hinsichtlich Wortschatz und Schreibvarianten zwar so eng wie möglich am Textmaterial erstellt werden, aber auch Vollständigkeit anstreben – was besonders dann vorrangig ist, wenn ein „automatisches Mustertraining“ (beschrieben im Abschnitt „OCR-Musterbereitstellung“) genutzt wird.

Die Ersetzungskoeffizienten (Gewichte für Zeichenersetzungen, die zur Überbrückung der eingestellten maximalen Wortdistanz kombiniert werden dürfen)<sup>66</sup> sollten empirisch anhand der OCR-Fehler eingestellt werden, die in einem OCR-Stichprobenlauf des konkreten aktuellen Materials tatsächlich häufig festzustellen waren.

##### ***Nach endgültiger Koordinatenspeicherung: manuelle Korrektur für Export in „Leseformate“***

In der Textansicht („View > Layer > Text“) bietet BIT-Alpha eine einfache Editiermöglichkeit, indem beim Anklicken eines Worts eine Eingabemaske erscheint. Die Arbeitsweise ist dem Vorgehen beim Training analog und intuitiv verständlich. Das Ergebnis der Korrektur kann aber nur in den für Einzelexport zur Verfügung stehenden Exportformaten gespeichert werden („Page > Export“), derzeit sind das insbesondere PDF, Text und HTML; das XML-Exportformat (hier: ALTO) gehört nicht dazu.

<sup>65</sup> Eine unbesehene Übernahme von historischen Transkriptionen bzw. E-Texten ins OCR-Lexikon ist daher auch bei einer „sehr geringen“ Fehlerquote problematisch. Ein Schreibfehler pro Seite kann in 1.000 Seiten Text zu 1.000 falschen Worteinträgen führen. In der OCR-Korrektur ebenfalls nur mit Vorsicht einzusetzen wären – grundsätzlich interessante – halbautomatisch generierte historische Wörterbücher, die von sprachgeschichtlich bekannten Lautwechseln und morphologischen Regelmäßigkeiten ausgehend historischen Wortbestand aufzubauen bzw. Korpora zu ergänzen versuchen, vgl. NEUMANN 2008.]

***Erreichbares Ergebnis:***

Man erhält ein koordinatenrichtiges XML-Ergebnis also nur in der Güte, die per automatischer Korrektur erreichbar war.

**(b) HK-OCR / FREngine 9*****Vor endgültiger Koordinatenspeicherung: manuelle Korrektur („Validierung“)***

Beim Laden eines OCR-erkannten Bilds bietet der Reiter „Korrektur“ im Einstellungsbereich von HK-OCR die Anzeige einer Liste der von der OCR erkannten Wörter nebst einem Korrektoreingabefeld an. In der Wortliste kann mit Cursortasten intuitiv auf- und abwärts navigiert werden, zugleich wird die Fundstelle im Faksimile farbig unterlegt und ein vergrößerter Ausschnitt der Trefferumgebung unterhalb des Seitenfaksimiles dargestellt. Der Korrekturgang („Validierung“) legt eine eigene Layout-Datei an; es kann demnach, solange der Zustand des OCR-Arbeitsverzeichnisses konsistent bleibt, wahlweise auf die OCR-Versionen vor und nach der Korrektur zurückgegriffen werden.

Die im Ergebnis der Korrekturen bzw. Worttrennungen und -zusammenlegungen angepassten Koordinaten werden in die Finereader-XML-Ausgabe eingepflegt.

Über das Menü „Erweiterte Validierung“ ist hier in der Nach-OCR-Phase sogar eine manuelle polygonale Blockgrenzen-Veränderung möglich, um z. B. aus dem Textblock herausgefallene Zeichen mit einzubeziehen oder fälschlich einbezogene Ornamentik auszuschließen.<sup>67</sup>

***Automatische lexikalische Korrektur***

Einige Korrektur-Leistungen werden offenbar OCR-nah bereits in Verbindung mit der Zeichenerkennung erbracht, wie der Einfluss der gewählten Wortliste auf die Erkennung zeigt. Ein separater automatischer Korrekturschritt ist sonst nicht vorgesehen; automatische lexikalische oder linguistische Nachbearbeitung kann also nur mit externen Werkzeugen erfolgen, folglich ohne – oder höchstens mit interpolierend schätzender – Koordinatenpflege.

***Erreichbares Ergebnis:***

Wenn der manuelle Korrekturaufwand geleistet werden kann, erhält man also ein koordinatenrichtiges Fertig-Ergebnis im FineReader-XML-Format und dessen Ableitungen, d. h. allen wählbaren Exportformaten.

**Desiderata, Kriterien für Alternativen**

Gebraucht wird auch hier zuallererst eine Dokumentation der Konfigurationsmöglichkeiten und ihrer Auswirkungen.

Für eine sachgerechte Pflege der Wortbestände und Ersetzungsregeln scheint das ungehinderte Vergleichen, Verändern, Übertragen und Kombinieren von beliebig großen Teilen solcher Wort- und Regellisten unverzichtbar. Beides wird im Regelfall vollständig vom Anwender bereitgestellt und ist hinreichend in Textzeilen beschreibbar. Proprietäre Binärformate, die nach einmaligem Einlesen einer Wortliste vom Anwender selbst weder gelesen noch frei bearbeitet werden können, entziehen dem Anwender dagegen ohne Notwendigkeit die Kontrolle über den ihm gehörenden Bestand. Wortbestände und Ersetzungsregeln könnten dabei unabhängig voneinander gespeichert sein. Gewichte für die Ersetzung eines Zeichens durch ein anderes sollten richtungsabhängig auch unterschiedlich eingestellt werden können (etwa wenn die Ersetzung von „l“ durch „i“ erleichtert werden soll, nicht aber die Ersetzung von „i“ durch „l“).

***Modularität***

Prinzipiell sind lexikalische Korrekturschritte Operationen auf Textdaten und daher zumindest dann nicht zwingend an das OCR-Werkzeug gebunden, wenn keine rückkoppelnden Abgleiche mit den optischen Vorlagen mehr erfolgen. Eingabeseitig läge daher nahe, durch eine Schnittstelle zum Korrektur-

<sup>66</sup> Hoher Koeffizient bedeutet: Ersetzung bevorzugen, niedriger Koeffizient hingegen: Ersetzung „erschweren“.

<sup>67</sup> Es ist nicht recht einzusehen, warum das nicht auch schon für die Segmentierung vor der OCR möglich sein sollte. Offenbar ist die Funktionalität der erweiterten Validierung von der gewählten Lizenzierung abhängig; sie stand später im Projekt nicht mehr zur Verfügung.

turschritt die Einspeisung von bereits OCR-gelesenem Material zu ermöglichen, ohne dass hierzu ein erneuter OCR-Lauf und der nochmalige Zugriff auf das Image Voraussetzung sind.<sup>68</sup> Und warum sollen Korrekturmodule mit speziellen Stärken (wie die durch explizite Gewichte sehr fein steuerbare automatische Korrektur von BIT-Alpha oder die praktische Bedienung der manuellen Korrektur in HK-OCR) sich auf die Verarbeitung nur der in der eigenen OCR-Software erzeugten Texte beschränken?

### **Evaluation**

Aufgrund der in Ansätzen gegebenen statistischen Auswertung<sup>69</sup> sowie der vielfältigen theoretisch einsetzbaren Möglichkeiten zu quantitativen Vergleichen der Versionen vor und nach einer Korrektur wäre es wünschenswert, wenn die Korrekturmodule selbst eine konfigurierbare Auswahl leicht handhabbarer Kennwerte verfügbar machen und vor allem eine Hilfe zur Interpretation anbieten würden.<sup>70</sup> Andererseits lässt das Text-(XML-)Format der Ergebnisse auch Fremdwerkzeuge zur Evaluation zu. Soweit entsprechende Tools zur Verfügung stehen, wäre das für Anwender eine potentiell unabhängige und objektivere Alternative.

### **Benutzerschnittstelle**

Auch wenn hier keine konkreten Anforderungen an eine ideale Korrekturbenutzerführung skizziert werden können: Die Bedeutung der Ergonomie eben des Korrektur-Arbeitsgangs kann nicht hoch genug angesetzt werden, wenn aus kompromissbehafteten Massen-OCR-Produktionen präsentierbare Volltexte hergestellt werden sollen. Das bezieht sich nicht nur auf die prinzipiell erreichbare Arbeitsgeschwindigkeit, sondern auch auf einen einarbeitungsarmen Einstieg externer, womöglich nur sporadisch mitarbeitender Korrektoren.

## **4.3.5 Semantische Erschließung: Strukturdaten und Named Entities**

Weder Strukturerkennung (angefragt waren textsortenspezifisch z. B. Kolummentitel, Kustoden, Bogen-signaturen u. a.) noch die Auszeichnung von „Named Entities“ aufgrund vorgegebener Listen (Orte, Berufe, Krankheiten, Bibelstellenangaben) wird bisher durch die getestete Software unterstützt.

Für die Projektphase war zunächst vorgesehen, den Herstellern ein für 2011 innerhalb der SBB PK und der DBV-AG Alte Drucke erreichbar scheinendes gemeinsam abgestimmtes Application Profile, das die umfangreichen Möglichkeiten der TEI hinsichtlich Elementen und Attributen sinnvoll eingrenzt, vorzugeben. Zu dieser Abstimmung, sicher eine anspruchsvolle Aufgabe,<sup>71</sup> war es während der Projektzeit nicht mehr gekommen, so dass diese Entwicklungsleistung nicht mehr termingerecht verlangt werden konnte.

Für die getestete Software konnte eine Ausgabe in semantisch orientierten Formaten wie TEI daher nur sehr oberflächlich ohne die genannten bibliothekarischen Zusatzinformationen aus den Exporten abgeleitet werden. Die Dienstleister boten entweder eine manuelle Erfassung oder die Programmierung von entsprechenden Zusatzmodulen an, was aber mit weiteren Kosten verbunden gewesen wäre und mangels Demonstration der Leistungsfähigkeit solcher Zusatztools für den aktuellen Test auch nicht angezeigt schien.

Insoweit zur Erkennung von Strukturinformation und Named Entities reine Textmerkmale ausreichen, kann eine diesbezügliche Anreicherung der OCR-Texte selbstverständlich auch nachträglich außerhalb der ursprünglichen OCR-Software erfolgen.

<sup>68</sup> BIT-Alpha führt immer eine Neu-OCR-Analyse durch; HK-OCR kann bei Vorliegen des Images, der Layout-Datei und der Status-XML-Datei im OCR-Arbeitsverzeichnis auf den letzten OCR-Lauf zurückgreifen.

<sup>69</sup> So möglich in statistischen Attributen der XML-Formate und in den layout.txt-Dateien von HK-OCR/FineReader.

<sup>70</sup> Gefragt wären kurze Rekapitulationen des Für und Wider von Anteilsquoten „lexikalischer“ Tokens am Text, von Abdeckungs -Zahlen zum Wörterbuch, der Signifikanz von Zeichen- vs. Wortfehlern usw. - Trotz aller Mängel jedes automatisch gewonnenen Evaluierungskennwerts sind diese nutzbar, wenn sie richtig verstanden werden: „... it will give an overview of where the OCR engine thinks it's succeeding and where it thinks it's failing. If there are wide discrepancies between one batch of files and another, the software log will allow the institution to prioritise those files where OCR accuracy is low, and to manage (and hopefully mitigate) those discrepancies.“ (ANDERSON 2010, S. 2)

### 4.3.6 Export

Von OCR-Software wird man grob zwei Hauptexportmöglichkeiten erwarten:

- ein druckbares, im Idealfall sowohl das Faksimile als auch den Volltext enthaltendes Dokument in einem verbreiteten Format wie PDF, RTF oder HTML oder
- ein die volle Text-, Layout- und Strukturinformation enthaltendes maschinell auswertbares Dokument z. B. in XML, aus dem der Anwender alle weiteren Formate und Funktionen (Volltextpräsentation? Druckpublikation? Hervorheben der Trefferstelle im Faksimile? Datenbank? Suchindex für bloße Volltextsuche?) generieren kann.

#### (a) BIT-Alpha

Für den OCR-Export können die Formate (ALTO-)XML, PDF, HTML, UTF8-Text und ein proprietäres Format „B.I.T. portable“ eingestellt werden; dazu können die Farb- und binarisierten Bilder exportiert werden. Das PDF-Dokument enthält das farbige Faksimile der Seite mit dem für eine Volltextsuche mit Highlighting der Fundstelle hinterlegten erkannten Text, sowie die bitonalen Regionen, die ebenfalls durchsuchbar als einzelne PDFs aufgerufen werden können.

Der Export der Einzelseiten-OCR wird entweder über das Menü „Page > Export“ oder per Kontextmenü (Rechtsklick auf die Seite) veranlasst. Für Batch-OCR müssen die gewünschten Formate im Menü „File > Batch process ...“ ausgewählt werden. Ein automatischer Export der den erkannten Zeichen entsprechenden bitonalen Bildsegmente als einzelne, nach dem Zeichenwert benannte Bilddateien in eine Art Typenbibliothek ist lt. Angaben von B.I.T. gegenwärtig in Vorbereitung. Für eine kurze Beispielseite mit sechs Zeilen Text ist ein Quelltext-Screenshot des ALTO-Outputs im Anhang (s. 7.1) abgebildet.

#### (b) HK-OCR / FReEngine 9

Sowohl für Einzelseiten- als auch für den Batch-Export gelten die im Reiter „Speichern“ im Einstellungsbereich aktivierten Formate; neben FineReader-XML und einem kompakteren, wortorientierten XML wird RTF (mit bildlicher Wiedergabe von nicht als Text erkannten Regionen) und PDF angeboten.<sup>72</sup>

Der Einzelseiten-Export für die gewählten Formate wird bei einer geladenen OCR-gelesenen Seite mit dem Menü „Speichern“ ausgelöst; im Batch-Lauf erfolgt der Export automatisch. Für eine kurze Beispielseite mit sechs Zeilen Text sind Quelltext-Screenshots der genannten XML-Formate im Anhang (s. 7.1) abgebildet.

#### Desiderata, Kriterien für Alternativen

Aufgrund der vielseitigen Verwendungsszenarien für OCR-Volltext wird jeder Satz von Exportformaten gelegentlich Wünsche offen lassen. Den Dienstleistern wurde deshalb als Anregung übermittelt, zusätzlich zu den ausprogrammierten Ausgabeformaten eine templateartige Schnittstelle anzubieten, für die der Benutzer selbst die Umwandlung des primären XMLs bereitzustellen hätte (z. B. als XSLT-Stylesheet, das einfach an ein im Programm vermutlich ohnehin existierendes XSL-Transformationsmodul zu übergeben wäre, oder als Datei von Ersetzungs-Regelsätzen aus Regulären Ausdrücken). Das könnte gerade in Batchläufen eine flexible Nutzung erlauben, ohne dass jeder Anwenderwunsch vom Hersteller ausprogrammiert werden müsste, und vor allem: ohne dass der Anwender sich schon vor Bestellung der Software auf jede Einzelheit seiner späteren Nutzungsanforderungen festlegen muss. Zwar sind Umformungen von XML-Formaten in beliebige Zielformate jederzeit auch extern möglich, wenn man sich

<sup>71</sup> Benötigt würde einerseits eine abschließende Definition einer als Minimalkonsens tauglichen, überschaubaren Untermenge eines aktuellen TEI-Formats (P5 oder Lite), s. a. einschlägige Vorarbeiten wie die an der HAB Wolfenbüttel für frühere TEI-Versionen prototypisch entwickelten „Barock-DTDs“, s. STÄCKER 2002, OPITZ 2002, andererseits hatten verschiedene Bibliotheken unterschiedliche Vorstellungen darüber, wie viele „technische“ OCR-Daten (Verarbeitungsparameter, Protokollwerte usw.) in das TEI-Format übernommen werden sollen und dürfen. Aus Sicht der SBB PK (Anforderungen der internen OCR-AG an ein notwendiges Standardformat wurden von Oliver Duntze formuliert) sollten folgende Daten (auch langfristig) in den Volltexten mitgeführt werden: Bildkoordinaten (buchstaben- oder wortweise), Schriftgröße und -schnitt, „Strukturdaten“ (Textblöcke auf der Seite), Alternative Lesarten, Linguistische Informationen (im Lexikon der OCR vorhanden? Zahlwort? Eigennamen? etc.), Erkennungssicherheit. Hinzu kommen administrative Informationen wie Grundlage der OCR (Dateiname der Bilddatei), Datum des OCR-Laufs, verwendetes OCR-Programm, eingesetzte Wortliste und Musterdateien.

<sup>72</sup> Die Erzeugung des PDF-Formats ist von der Lizenz abhängig, die HK-OCR erfordert.

intensiv genug mit den erhaltenen Exportformaten auseinandersetzt. Eine im Programm bereits angelegte Schnittstelle könnte aber den Aufwand verringern, die Eingriffsstellen etwas vorstrukturieren und die Schwelle zur Gestaltung eigener Exportformate somit senken. Vermieden werden sollte jedenfalls, dass Anwender ihrerseits ihre Volltext-Anwendungen auf zufällige Eigenheiten der Ausgabeformate bestimmter Softwareprodukte abstimmen und damit den zukünftigen Spielraum eigener Entwicklungen unnötig beschränken.

#### 4.4 Bereitstellung der OCR-Muster: Anlegen (Training) und Pflege von Musterbibliotheken

##### (a) BIT-Alpha

###### *Training*

Die Muster können entweder mit in einer bda-Konfigurationsdatei gespeichert sein oder in einer separaten Musterdatei im Format „.bit“ (s. Abb. 41 am Anfang des Kapitels). Nachdem die OCR einer Einzelseite gestartet wurde,<sup>73</sup> zeigt die OCR-Ansicht („View > Layer > OCR“), ob und – durch Farben unterschieden – mit welcher „Sicherheit“ BIT-Alpha die Zeichen erkannt hat. Das Anklicken eines beliebigen Zeichens öffnet das Trainingsfenster (Titelleiste „Verify Recognition“), wo der Zeichenwert des optischen Musters festgelegt oder geändert werden kann.

Mustern, denen kein Zeichenwert aus den beteiligten Alphabeten zugeordnet werden kann, kann ein eigener Code als „Special Value“ zugewiesen werden. Wegen der Möglichkeit, solche Sondercodes OCR-nah durch einen im „Sequencer“ (s. u.) festgelegten Zeichenwert zu ersetzen oder zu kombinieren, kann damit u. a. ein Training von Drucktypenteilen erfolgen, die in der Segmentierung gelegentlich anfallen (z. B. einzelne „Beine“ eines „n“ oder „m“).

Eine Kurzanleitung für einen einfachen Trainingsgang wird im Anhang 7.2.2 gegeben.

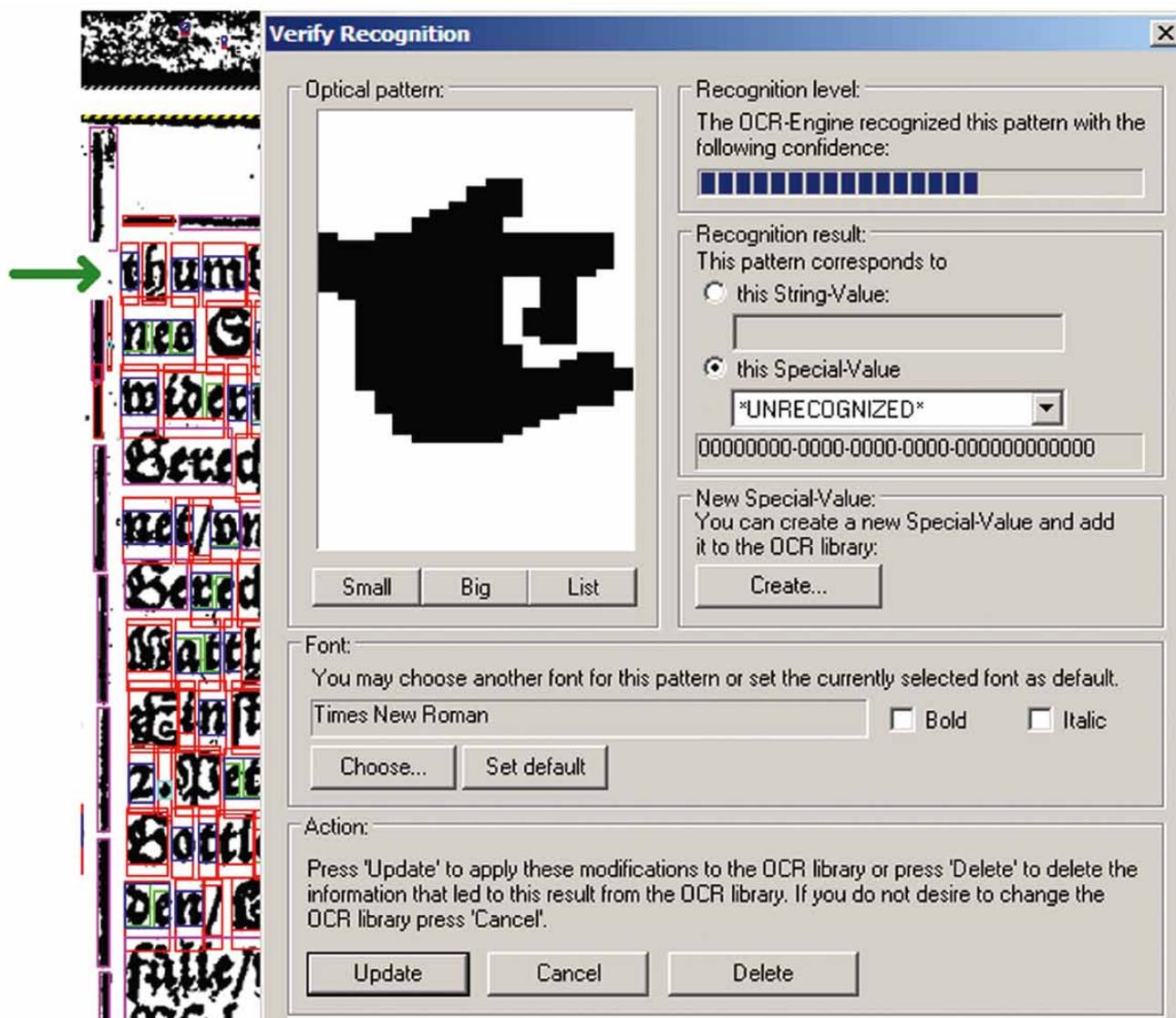
###### *Maschinelle Beschleunigung bzw. Verstärkung des Trainings*

BIT-Alpha versucht mit zwei Verfahren, in manuellem Training begonnene Musterbestände zur automatischen Vergrößerung bzw. zur Verstärkung des Musterbestands zu nutzen. In beiden Fällen werden den manuell trainierten Mustern „ähnliche“ Muster automatisch hinzugefügt, was die richtige Erkennung von Drucktypen befördern *kann*, die dem manuell gelernten Zeichen nicht mehr ganz so ähnlich sind.

- „Learn Page“: Bei einer geöffneten Einzelseite können die als „fraglich“ erkannten Muster einer Seite (in der Ansicht „OCR“ durch nicht-blaue Farben gekennzeichnet), wenn sie vom Bearbeiter als korrekt eingeschätzt wurden, mit dem Befehl „OCR > Library > Learn Page“ der Musterbibliothek insgesamt auf einmal hinzugefügt werden. Dabei ist ein Konfidenzintervall einstellbar, für welche Erkennungswahrscheinlichkeiten das gelten soll.
- „Autolearn“: Etwas unmotiviert in das Menü der exportierenden Batchverarbeitung integriert findet sich die Möglichkeit, eine Anzahl Seiten ohne Export laufen zu lassen und dabei – wieder mit einem (unter „OCR > Library > Learning parameters“) einstellbaren Konfidenzintervall – die auftretenden optischen Muster des ganzen Stapels mit dem erkannten Zeichenwert automatisch der Musterbibliothek hinzuzufügen.<sup>74</sup> Damit können schnell große Mustermengen erreicht werden, deren Erkennungsleistung allerdings vor dem Einsatz kontrolliert werden sollte. In jedem Fall sollte der Zustand der Musterbibliothek vor dem „Autolearn“ gesichert werden, und es empfiehlt sich, durch vergleichende Tests das materialbezogen günstigste Konfidenzintervall zu ermitteln.

<sup>73</sup> Das Starten der OCR geschieht selbsttätig beim Laden eines Bildes („File > Load Image“) oder wenn die bda-Konfigurationsdatei ein gespeichertes Image enthält, dann schon beim Öffnen der bda-Datei („File > Open“).

<sup>74</sup> Achtung, im normalen exportierenden Batch-Lauf sollte „Autolearn“ dagegen deaktiviert werden, um die Musterbibliothek nicht ungewollt zu verändern!



46 Trainingsfenster in BIT-Alpha

**Sequencer**

Für Zeichenwerte, *special values* und Kombinationen aus diesen können im sogenannten Sequencer („OCR > Library > Sequencer“) feste unmittelbar nach der OCR vorzunehmende Ersetzungen definiert werden, die es beispielsweise ermöglichen, als *special value* gelernte Zeichenteile miteinander oder mit einem benachbarten Zeichen zusammzusetzen oder Abkürzungen zu interpretieren. Diese Ersetzungen dienen weniger der Korrektur echter OCR-Lesefehler als vielmehr der Behandlung erwartungsgemäß anfallender punktueller OCR-Ergebnisse, die lediglich einer regulären Umwandlung in einen Ziel-Zeichenwert bedürfen.

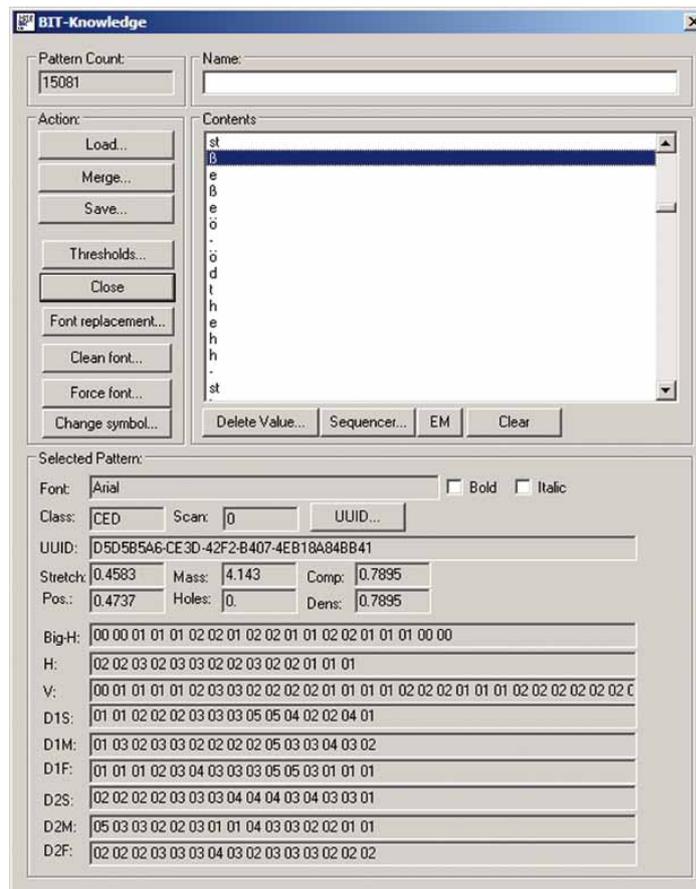
**Speicherung, Bearbeitung, Zusammenführung**

Muster, die im bit-Musterdateiformat vorliegen, können mit dem Mustereditor „BIT-Knowledge“ bearbeitet werden; möglich ist u. a. die Bearbeitung einzelner gespeicherter Muster (Entfernen, Zuweisung eines anderen Zeichens usw.). Hier steht allerdings keine bildliche Ansicht der Muster-Repräsentation mehr zur Verfügung, was den Umgang mit den von BIT-Knowledge angezeigten Muster-Einträgen sehr einschränkt.<sup>75</sup>

<sup>75</sup> Der häufigste Anwendungsfall war das Löschen von einzelnen falsch klassifizierenden Mustern, wie sie z. B. durch Tippfehler beim Training oder durch ein zu großzügiges automatisches Lernen angelegt werden können. Da man nur den Zeichenwert des zu löschenden Musters sieht, ist nicht ganz auszuschließen, dass hierbei vielleicht eine „gute“ Musterrepräsentation gelöscht wird; man löscht in BIT-Knowledge gewissermaßen nur noch „eines der Muster, welche für die falsche Klassifikation zu einem Zeichen x verantwortlich sind“. Für schon in BIT-Alpha bemerkte (oder dort wiederfindbare) Fehler ist die Situation etwas besser: hier kann man im Layer „OCR“ für ein konkretes Zeichensegment (anklicken) dasjenige Muster löschen, „das für diese falsche Klassifikation verantwortlich ist“. Es scheint in jedem Fall ratsam, nach dem Löschen eines Mustereintrags zu dessen Zeichenwert sicherheitshalber einige „gute“ (typische) neue Muster zu lernen.

In BIT-Alpha kann nicht mehr als eine bit-Musterdatei geladen werden, der Mustereditor „BIT-Knowledge“ dient deshalb auch zum Zusammenfügen verschiedener bit-Dateien in eine einzige, damit die Muster aus beiden genutzt werden können. Hierzu müssen die Musterbestände je selbst in einer bit-Musterdatei vorliegen;<sup>76</sup> dann kann zu einer ersten geladenen Musterdatei<sup>77</sup> eine weitere hinzugeladen<sup>78</sup> und die Gesamtmustermenge gemeinsam in einer bit-Musterdatei gespeichert werden.<sup>79</sup> Dieses Hinzumischen einer weiteren Musterdatei kann mit weiteren bit-Dateien wiederholt werden.

Achtung: Wenn außer Mustern auch Sequenzen in den beteiligten bit-Dateien gespeichert sind, werden beim Mischen nur die Sequenzen der zuerst geöffneten bit-Datei übernommen.



47 BIT-Know Mustereditor

<sup>76</sup> D. h. wenn sie bis dahin nur in einer bda-Konfigurationsdatei gespeichert waren, muss daraus erst eine bit-Musterdatei angelegt werden (in BIT-Alpha: „OCR > Library > Save“)

<sup>77</sup> In BIT-Knowledge: „Load ...“

<sup>78</sup> In BIT-Knowledge: „Merge ...“

<sup>79</sup> In BIT-Knowledge: „Save ...“

### **Grafische Benutzeroberfläche (GUI)**

Jede Stelle auf der gelesenen Seite, z. B. in der Ansicht (Layer) „OCR“, ist per Maus frei erreichbar; man kann also wählen, ob man nur schwerpunktmäßig markante Zeichen trainieren oder den ganzen Text durchgehen möchte. Nach Anklicken eines Zeichenbereichs wird das Trainingsfenster „Verify Recognition“ aktiv. Darin wird das optische Muster binarisiert und in der der OCR-Engine vorliegenden Form noch einmal gezeigt, so dass beurteilt werden kann, ob diese Vorlage typisch genug ist und ihr ein Zeichenwert zugewiesen werden soll oder nicht. Statt eines Zeichenwerts kann einem Muster auch ein *special value* (selbstdefinierter Code) zugewiesen werden, auf den im Sequencer kontextabhängig zurückgegriffen werden kann. Sinnvoll ist das z. B. für separat segmentierte Buchstabenteile oder für mehrdeutige Zeichen wie das runde „r“, das in der etc.-Abbeviatur ausnahmsweise als „et“ aufzulösen ist.

Jedem Muster wird ein „Font“ zugewiesen, der zunächst eine Information zur Ausgabeformatierung ist. B.I.T. empfiehlt die Font-Zuweisung als Mittel zur unterschiedlichen Behandlung von Mustern verschiedener Schriftfamilien, d. h. Fraktur einen anderen Font zuzuweisen als Antiqua. Eine andere Art, beim Training die Schriftfamilie zu vermerken, gibt es nicht; insofern auch keine Möglichkeit, aus einer Musterdatei die Muster der Schriftfamilien separat zu extrahieren.

Die Anordnung der Bedienelemente im Fenster „Verify Recognition“ erfordert viel Aufmerksamkeit zur Vermeidung von Fehleingaben und teilweise umständliche Tastatur- und Mausektionen; Hinweise zu einer möglicherweise günstigeren Ergonomie wurden dem Hersteller gegeben.

Per Standardeinstellung („OCR > Autorefresh“ aktiviert) wird nach jeder neuen Musterstringeingabe eine OCR-Analyse mit der neu kompilierten Musterbibliothek gestartet, wodurch die OCR-Ansicht sofort die nach dem Lernschritt neu bzw. anders erkannten Zeichen – wieder nach Erkennungssicherheit farbig abgestuft – anzeigt. Diese Neuanalyse benötigt je etwas Zeit; für ein zügiges Training kann „Autorefresh“ auch deaktiviert werden – dann folgt der Musterstringeingabe keine Verzögerung – und eine Neu-OCR lässt sich jederzeit per „OCR > Refresh“ manuell erzwingen.

Der Sequencer („OCR > Library > Sequencer“) wird in einer Eingabemaske „OCR-Sequence Manager“ bedient; auch hier wurden Hinweise zur Ergonomie an den Hersteller gegeben. In den getesteten Versionen von BIT-Alpha musste speziell beachtet werden, dass entgegen der Intuition eine neue Sequenz offenbar nicht erst bei Betätigung des OK-Buttons dieses Fensters angelegt und gespeichert wird, sondern bereits beim Editieren der Zeichenwerte und ohne Hinweis.

Die Oberfläche zur Musterbearbeitung in BIT-Knowledge eignet sich gut zum Bearbeiten bzw. Löschen einzelner Muster und zum Zusammenfügen zweier bit-Musterdateien, solange keine Konflikte wegen dort mitgespeicherter Sequenzen eintreten können, auf deren implizite Übernahmepriorität nicht hingewiesen wird. Eine effektive Bearbeitung von vielen Mustern auf einmal scheitert an der unübersichtlichen Auflistung (fehlende Sortierbarkeit, keine Mehrfachselektion).

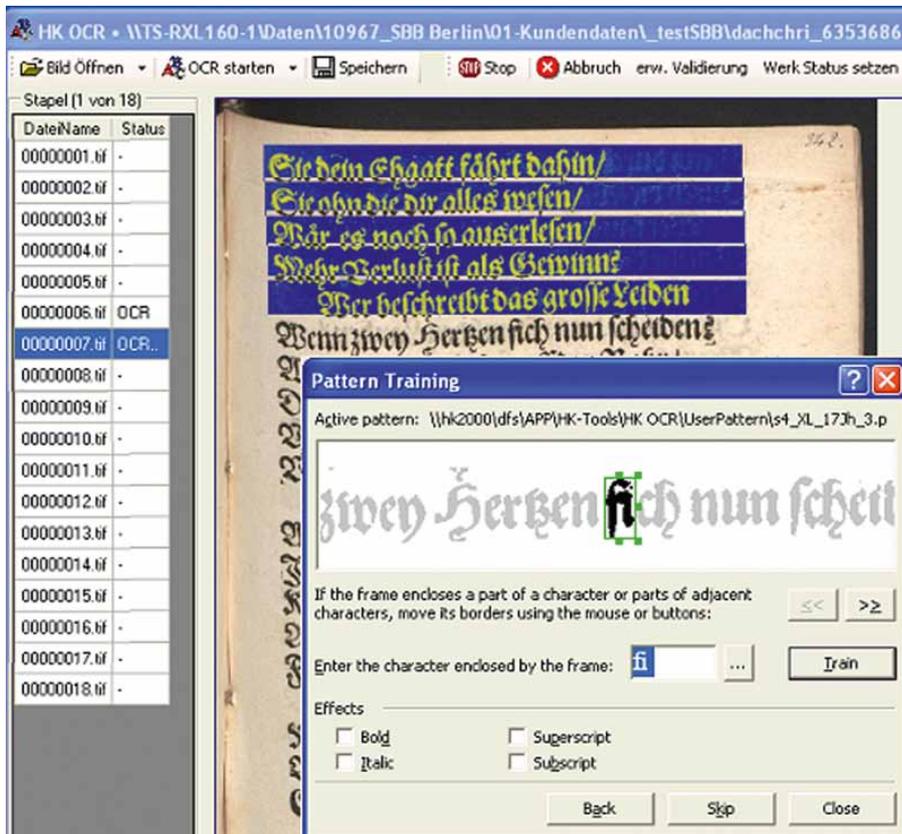
### **(b) HK-OCR / FREngine 9**

Die FineReader-Engine bringt auch für Fraktur einen Grundbestand an Mustern (*built-in patterns*) mit; es muss also nur dann trainiert werden, wenn – wie bei Alten Drucken oder bei speziell gestalteten Lettern – die *built-in patterns* den Drucktypen der Vorlagen nicht gut genug entsprechen.

#### **Training**

Zum Training wird eine Seite wie zur OCR bereitgestellt („Bild öffnen > Ordner wählen > Ordner wählen“); vor dem Trainingslauf muss zunächst im Reiter „Einstellungen“, Feld „Erkennungsmuster“ (bei deaktivierter Checkbox „Musterdatei verwenden“) die Checkbox „Training“ aktiviert werden, und es muss eine ptn-Musterdatei erstellt bzw. – falls schon existent – ausgewählt werden (Button „...“ neben dem Textfeld mit dem Dateinamen der Musterdatei). Wie zur OCR können Sprachen bzw. Sprachengruppen und Schriftfamilien eingestellt werden. Dann kann der Trainingslauf wie ein OCR-Lauf gestartet werden („OCR starten > Aktuelles Bild lesen“).

Die recht problematische, von den FineReader-Bausteinen bestimmte Benutzerführung des Trainings erzwingt im Prinzip ein komplettes lineares Durchgehen der Seite vom ersten Zeichen zum letzten. In einem Fenster wird der Bildausschnitt um das aktuelle Segment herum jeweils binarisiert und vergrößert dargestellt, und dem markierten Muster kann ein Zeichenwert zugeordnet werden. Möglich sind nur Zeichenwerte innerhalb der für die Sprachengruppe eingestellten Zeichen – der Anwender muss also bereits vor dem Training alle im Bestand möglicherweise vorkommenden Sonderzeichen (Diakritika, Zeichen anderer Sprachen) kennen und bereitstellen. Nach Durchlaufen der Seite wird die Musterdatei gespeichert, und es kann eine weitere Seite aufgerufen werden.



48 Trainingsfenster in HK-OCR<sup>80</sup>

### **Speicherung, Bearbeitung, Zusammenführung**

In einem Museditor-Fenster (Fenster „User Pattern“) können die in einer ptn-Datei gespeicherten Muster in der bitonalen optischen Form mit dem Zeichenwert angezeigt, gelöscht und wohl undefiniert werden. In begrenztem Rahmen ist eine Mehrfachselektion zur Löschung mehrerer Muster auf einmal möglich. Nicht vorgesehen ist eine Zusammenführung verschiedener Musterdateien.

### **Grafische Benutzeroberfläche (GUI):**

Speziell die von FineReader vorgegebenen Bildschirminteraktionen beschränken unnötig den Spielraum und die Effizienz beim Training und bei der Bearbeitung der Musterdateien. Entsprechende Hinweise wurden an den HK-OCR-Hersteller gegeben, der trotz Konsultationen mit ABBYY offenbar bisher keine Möglichkeit zur Bearbeitung dieser Bausteine bekommen hat. Auch eigene Nachfragen bei einer ABBYY-Regionalvertretung zeigten, dass die Benutzerfreundlichkeit der Trainingswerkzeuge nicht als aktuelle Aufgabe gesehen wird; verschiedenen während der Projektlaufzeit in Deutschland erreichbaren ABBYY-Vertretern schien die Möglichkeit eines (sinnvollen) anwenderseitigen Muster-Trainings nicht einmal bekannt zu sein.

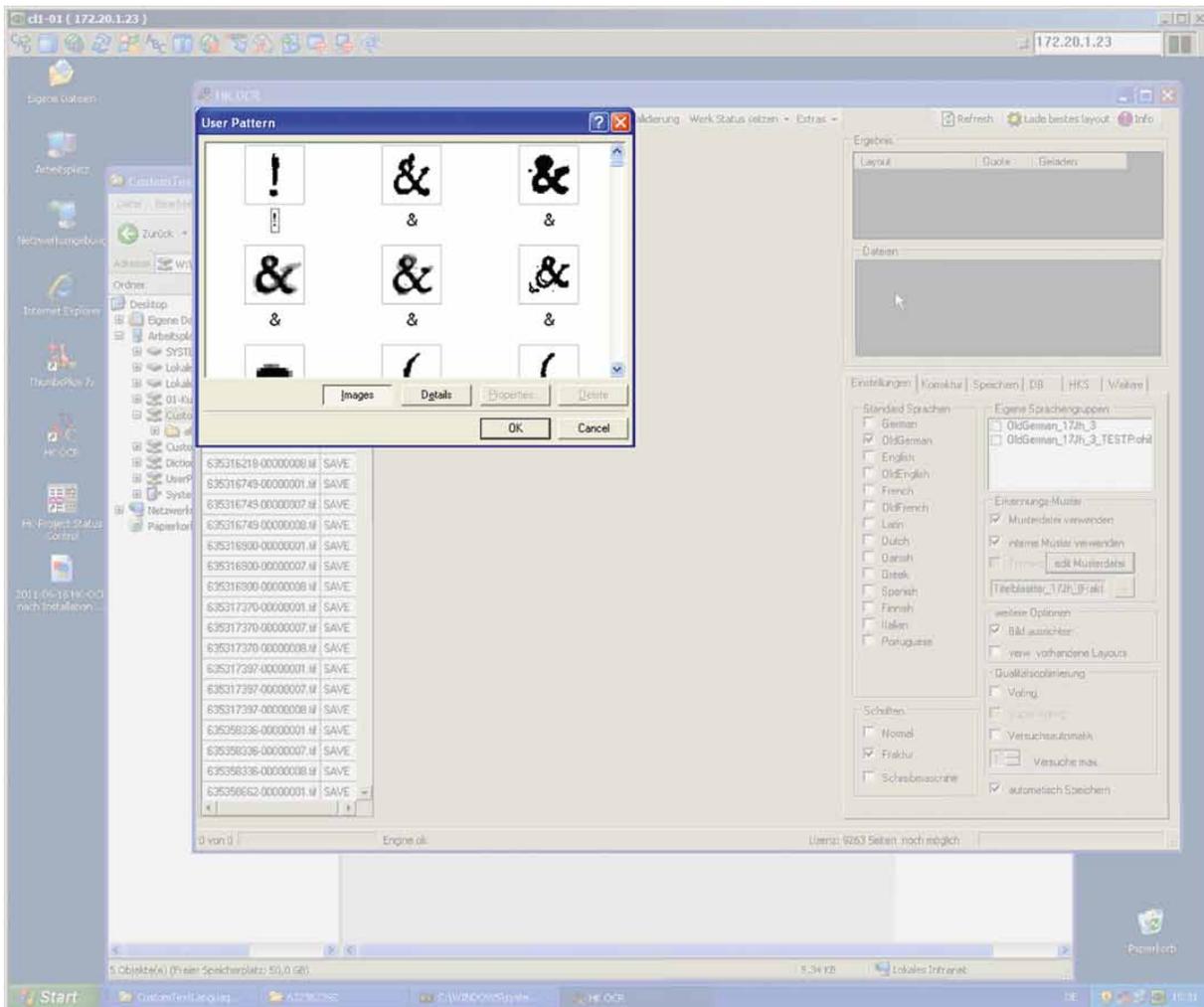
<sup>80</sup> Im Beispiel: Die OCR „erkannte“ die Ligatur „fi“; im Training ist Gelegenheit zur Eingabe der richtigen Zeichenfolge „si“. - Bei diesem eher schlechten Muster wird man auf ein Training aber vermutlich verzichten, um die Erkennung der „fi“-Ligatur nicht zu beschädigen.

In der Trainingssituation (Fenster „Pattern Training“ der FineReader-Engine) wird der Benutzer-Zeit-aufwand durch eine ineffektive Führung durch die Bildseite immens in die Höhe getrieben. Die wichtigsten Nachteile sind:

- Es ist nicht möglich, eine Trainingssequenz an beliebigen Orten innerhalb der Seite zu beginnen und zu beenden. Die Navigation verläuft in Einzelschritten vorwärts. Man kann das Training nur am Anfang der Seite beginnen und muss daher, auch wenn man etwa nur am Ende der Seite ein paar Zeichen trainieren will, jedes Zeichen der ganzen Seite noch einmal bestätigen oder explizit überspringen (mit großem Zeitverlust und der Gefahr, beim „Durchklicken“ bereits gelernte Muster-Zeichen-Paare zu beschädigen).
- Ein Rückschreiten (zur Korrektur von bemerkten Fehleingaben) ist nur in Einzelschritten und nur innerhalb desselben Worts möglich.
- Oft werden dagegen automatisch Zeichen(gruppen) ohne Vorwarnung übersprungen, offenbar dann, wenn sie mit hohem Konfidenzwert gelesen wurden. Abgesehen davon, dass auch dort in Einzelfällen ein Korrektur- oder Trainingsbedarf bestehen kann, zwingt die ständige Unberechenbarkeit, ob nach einer Eingabe wirklich das nächstfolgende oder aber ein viel späteres Muster ins Trainingsfeld eingetragen wird, zu einer unnötigen Verlangsamung des Trainingsgangs. Bei zügiger Arbeitsweise bestätigt man per Enter-Taste oder Button leicht einen String, der auf das nach Schrittfolge erwartete Muster gepasst hätte, versehentlich für ein ganz anderes Muster. Wenn der Sprung dabei die Wortgrenze überschritten hat, ist dadurch ein Zurückschreiten zum eben eingegebenen Zeichen nicht mehr möglich, und die Fehleingabe wird unweigerlich in die Musterdatei übernommen. Eine Korrektur im Mustereditor ist zeitaufwendig und hängt von der Wiedererkennung des richtigen Zeichenbilds ab.
- Die im Fenster „Pattern Training“ angebotenen Tasten sind zumindest im deutschen Tastaturlayout nicht immer günstig, es kommt zu ungünstigen Maus-Tastatur-Wechseln und zu einer Erleichterung unerwünschter Aktionen gegenüber erwünschten.

Im Mustereditor („Edit Musterdatei“), ebenfalls einer FineReader-Komponente, fiel auf:

- Das angezeigte Fenster „User Pattern“ ist viel zu klein, es zeigt je nur 2 x 3 Muster in einer nicht änderbaren Reihenfolge und kann nicht vergrößert werden. Für einen Überblick fehlt die Möglichkeit, das Fenster so groß zu ziehen, dass wesentlich größere Mustergruppen miteinander optisch verglichen werden können.
- Die in „Edit Musterdatei“ vorgesehene Möglichkeit, gespeicherte Muster-Zeichen-Zuordnungen zu ändern, scheint noch nicht zu funktionieren. Während des Tests ist zwar das Löschen einzelner Muster gelungen, nicht aber das Ändern des zugewiesenen Strings.



49 HK-OCR Mustereditor

## Desiderata, Kriterien für Alternativen

### *Modularität*

Um die Stärken selbst trainierter auf konkrete Drucke bezogener Muster ausnutzen zu können, wird eine freie Kombinierbarkeit und Teilbarkeit von Musterbibliotheken benötigt. Hierzu wäre ein dokumentiertes, für den Nutzer verständliches Speicherformat für Muster gefragt, das u. a. Attribute der einzelnen Muster verwalten kann, nach denen Teilbestände selektiert werden können usw. Ideal wäre natürlich eine anwendungsübergreifende Austauschbarkeit trainierter Muster (s. Kap. 6.4).

### *Evaluation*

Ehe überhaupt sinnvolle Evaluationswerkzeuge konzipiert werden können, müssen Musterbibliotheken in einer jederzeit auslesbaren und miteinander vergleichbaren Form vorliegen.

### *Benutzerschnittstelle*

Hinsichtlich der graphischen Benutzerführung könnten beide Programme voneinander lernen. Unerlässlich sind die freie Navigation an eine beliebige Stelle sowie die freie Unterbrechung und unverzögerte Wiederaufnahme des Trainingsgangs.

Restriktionen im Alphabet der zulässigen Zeichenwerte sollte es nicht geben, da prinzipiell jede Druckvorlage, von der noch kein elektronisch durchsuchbarer Volltext vorliegt, etwas Unvorhergesehenes enthalten kann.

Der Mustereditor muss eine übersichtliche Ansicht der gespeicherten Muster bieten, d. h. sortierbar, in größenverstellbarem Fenster, mit Mehrfachauswahl und verständlich bezeichneten Aktionen.

## 4.5 Batch-spezifische Aspekte

Wie bereits angedeutet (s. Kap. 3.1.7), hängt die Qualität von Massen-OCR davon ab, ob die verwendeten Parameter, Musterbestände und Lexika so gut auf die Gruppe der gemeinsam zu verarbeitenden Bilder abgestimmt sind, wie es in der Einzelseitenbearbeitung möglich ist.

Die Software hat hier nur insofern Einfluss, als sie

- diese Gruppierung in irgendeiner Form unterstützen könnte, z. B. durch Kennziffern zur nachträglichen Selektion von Seiten, die gut oder schlecht auf einen OCR-Lauf angesprochen haben oder
- Konfigurationsparameter innerhalb eines Stapels zu wechseln in der Lage ist.

Außerdem wird die Flexibilität der Gruppenbildung davon beeinflusst, ob Quell-(Bild-)Dateien zur Stapelverarbeitung erst in ein gemeinsames Arbeitsverzeichnis kopiert werden müssen, oder auch wie leicht und sicher ggf. der Abbruch und die Weiterführung einer Stapel-Verarbeitung möglich sind usw.

### (a) BIT-Alpha

Eine Vorab-Gruppierung wird nicht unterstützt, sie muss daher vom Anwender geleistet werden.

BIT-Alpha bietet keine für den Anwender erreichbare Auswertung von Kennziffern, die eine nachträgliche Gruppierung der Bilder in „fertige“ und „nachzuprozessierende“ erleichtern würden. Auch der XML-Export enthält keine diesbezüglich auswertbaren Attribute. Für den Ergebnistext könnte der Anwender selbst z. B. den Bezug zum Wörterbuch herstellen und eine Art Worterkennungsrate bestimmen.

Die in einem Batch-Lauf zu bearbeitenden Bilder müssen in ein und demselben Verzeichnis liegen; dann kann nach Laden der gewünschten bda-Konfigurationsdatei und den üblichen Einstellungen der Batch-Prozess per „File > Batch process ...“ konfiguriert (Exportformate festlegen usw.) und angestoßen werden. Zu beachten ist, dass das ebenfalls in diesem Dialog untergebrachte „Autolearn“ (s. Kap. 4.4) im normalen (Export-)Batchlauf deaktiviert sein muss, um die Musterbibliothek nicht ungewollt zu verändern.

Die gewählten Exportformate werden in Unterverzeichnisse eines wählbaren Exportverzeichnisses geschrieben. Wenn außer den sparsamen Text-, HTML- und XML-(ALTO)-Formaten auch das Bilder enthaltende PDF exportiert wird, ist auf ausreichend Speicherplatz im Dateisystem des Exportverzeichnisses zu achten.

Zur Fortsetzung einer abgebrochenen Stapelverarbeitung war es am praktischsten, je die bereits verarbeiteten Bilder aus dem Quellverzeichnis zu entfernen oder überhaupt neue Verzeichnisse zu benutzen.

### (b) HK-OCR / FReEngine 9

Eine Vorab-Gruppierung wird nicht unterstützt, sie muss daher vom Anwender geleistet werden.

HK-OCR legt bei jedem OCR-Lauf zu jedem *layout* auch Protokolldateien an, in denen statistische Kennziffern der aktuellen OCR festgehalten werden und die bei eingeschaltetem *voting* schon von HK-OCR selbst zur automatischen Auswahl des nach diesen Kennziffern „besten“ OCR-Versuchs genutzt werden. Eine genaue Dokumentation der Bedeutung der einzelnen Werte lag der Software nicht bei, sollte aber erhältlich sein. Sinnvolle Interpretierbarkeit vorausgesetzt, könnten sich diese Werte wegen des Textformats der Protokolle auch nach Stapelverarbeitungen maschinell auswerten lassen, ebenso die im Finereader-XML-Export (s. Screenshots einer FineReader-XML-Ausgabe im Anhang 7.1) enthaltenen Attribute zur Zeichenerkennungskonfidenz („charConfidence“) und zum Vorkommen des gelesenen Worts im Wörterbuch („wordFromDictionary“).<sup>81</sup>

Die in einem Batch-Lauf zu bearbeitenden Bilder müssen in ein und demselben Verzeichnis liegen, das zugleich das Ausgabeverzeichnis ist. Wegen der pro Originalbild entstehenden FRImage-Dateien, die stets etwas größer sind als die Originalbilddatei, muss der freie Speicherplatz im Dateisystem des Verzeichnisses noch einmal etwas größer sein als die Gesamtdateigröße der Bilder.

<sup>81</sup> Vgl. die Hinweise zur kritischen Anwendbarkeit solcher Kennzahlen in den Best-Practice-Guides des IMPACT-Projekts (ANDERSON 2010)

Der Status der Verarbeitung jedes Bilds aus dem Batch-Lauf wird in einer Liste in HK-OCR angezeigt und in einer XML-Datei im Ausgabeverzeichnis vermerkt. Solange diese Steuerdatei („\_BildDateien.xml“) erhalten ist, kann eine abgebrochene Stapelverarbeitung leicht wieder aufgenommen werden und überspringt die bereits als verarbeitet und gespeichert vermerkten Bilder. Abbrüche traten vereinzelt spontan auf, möglicherweise wegen geringer Zugriffsverzögerungen auf ein Netzlaufwerk. Sicherheits halber sollte nach jedem unterbrochenen Batch-Prozess die Vollzähligkeit der XML-Exporte überprüft werden; eventuell fehlende einzelne Seiten können dann leicht manuell nachprozessiert werden. Hierzu kann es notwendig werden, per Kontextmenü den in der Stapel-Übersicht von HK-OCR vermerkten „Status“ der jeweiligen Bilddatei zurückzusetzen.

### **Desiderata, Kriterien für Alternativen**

Zu den im Folgenden angedeuteten Anforderungen werden im Abschnitt 6.4 konkrete Abhilfemöglichkeiten beschrieben.

#### ***Granularität***

Batch-Verarbeitung ist aus sich heraus nicht davon abhängig, dass die zu verarbeitenden Bilder in ein und demselben Verzeichnis liegen. Außerdem müssen lange Batch-Läufe gelegentlich unterbrochen und dann nur für den Rest fortgesetzt werden, oder man will von vornherein nur bestimmte Bilder verarbeiten. Es gibt keinen nachvollziehbaren Grund, dass im Dateisystem erreichbare Quellbilder erst in ein Arbeitsverzeichnis gebracht werden müssen, und dass der Stapel vom Inhalt dieses Verzeichnisses bestimmt wird.

Es sollte zudem auch innerhalb eines Stapels möglich sein, für die Bilder wechselnde Parameter, Muster- und Wortbibliotheken vorzusehen. Die innerhalb üblicher Stapel zu beobachtenden erheblichen Qualitätsunterschiede der Fraktur-OCR rechtfertigen den Aufwand allemal, und mit etwas Glück geben bereits die Umfangsangaben im Katalog oder die Metadaten der Digitalisate für Seitenbereiche deren Schriftart und Sprache an.

#### ***Evaluation***

Würde die Software geeignete statistische Indikatoren zur Abschätzung der Erkennungsgüte der einzelnen Seiten bereitstellen, dann könnten so ermittelte ungenügende Seiten nachträglich einer Neu-OCR mit anderen Mustern – oder anderen Wörterbüchern – zugeführt werden. Für Szenarien, in denen genügend Rechenzeit gegeben ist, sollte es möglich sein, heterogenes Material von vornherein parallel mit verschiedenen Musterdateien bzw. Wortlisten zu verarbeiten und aus den Ergebnissen das je geeignetste zu selektieren.

#### ***Benutzerschnittstelle***

In beiden Programmen müssen Konfigurationen bisher auch dann per Bildschirminteraktion ausgewählt werden, wenn sie fertig in Parameterdateien vorliegen. Für eine Stapelverarbeitung mit vorab feststehenden Konfigurationen sollte aber neben der Bildschirminteraktion stets auch ein nicht-interaktives Anstoßen auf Kommandozeilen-, Batch-/Skript- oder Ticket-Ebene möglich sein (konventionelle Übergabe der Arbeitsverzeichnisse und der jeweiligen Parameter-, Muster- und Wörterbuchdateien).

## 4.6 Voraussetzungen und Randbedingungen des Softwareeinsatzes

Für beide Produkte wurden übliche Windows-PCs (Prozessoren Intel Core 2 Duo E8400 mit 3,0 GHz, Arbeitsspeicher 2GB, Betriebssystem Windows XP Professional SP 3) genutzt. Solange nur mit einer Instanz des jeweiligen Programms gearbeitet wird, haben diese Hardwarevoraussetzungen genügt; bestimmte Verarbeitungsphasen der OCR lasten den PC dabei zeitweise sehr aus, so dass von gleichzeitiger Arbeit in anderen Anwendungen abzuraten ist.

### (a) BIT-Alpha

#### *Lizenzmodell*

Die Standalone-Anwendung BIT-Alpha (ebenso der Mustereditor BIT-Knowledge) wird als Anwendung lizenziert, hierzu ist eine Aktivierung unter Administratorrechten nötig.

#### *Installations- und Betriebsvoraussetzungen*

Zum Betrieb von BIT-Alpha genügt ein Einzelplatz-PC mit einem aktuellen Windows-Betriebssystem (ab XP). Eine möglichst gute Prozessorleistung (Mehrkernarchitektur) ist angeraten. Für jede Seite ist größenabhängig mit mehreren Minuten OCR-Analyse zu rechnen. Die Priorität des Programms für das Betriebssystem kann in der Benutzeroberfläche eingestellt werden, z. B. um sie herabzusetzen, wenn der PC gleichzeitig für anderes genutzt werden muss und längere OCR-Zeiten hingenommen werden können („File > Settings > System > Thread Priority“).

#### *Angebot als Dienstleistung*

B.I.T. bietet wahlweise die Komplett-OCR (vom Bild zum Exportformat) oder die Unterstützung anwenderseitiger OCR (Konfiguration der Binarisierungsparameter, Übernahme von Training usw.) auch als Dienstleistung an.

### (b) HK-OCR / FREngine 9

#### *Lizenzmodell*

HK-OCR wurde als Anwendung von Herrmann & Kraemer zeitlich befristet lizenziert und ermöglicht gesondert die Neu-OCR bzw. die Korrektur („Validierung“) von Seiten, für die bereits eine FRImage-Datei im Arbeitsverzeichnis vorliegt. Zur erstmaligen OCR jedes Bilds (zugleich: Erstellung der FRImage-Datei) wird außerdem eine ABBYY-Volumenlizenz benötigt. Hierzu muss entweder auf dem Arbeitsrechner oder auf einem im Netzwerk erreichbaren Rechner der ABBYY-Lizenzmanager installiert sein und erworbene Seitenlizenzen enthalten.

#### *Installations- und Betriebsvoraussetzungen*

Sowohl zum Betrieb von HK-OCR als auch für den ABBYY-Lizenzmanager genügen Einzelplatz-PCs mit einem aktuellen Windows-Betriebssystem (ab XP). Die Installation und Aktivierung erfolgt jeweils unter Administratorrechten. Der ABBYY-Lizenzmanager kann entweder ebenfalls auf dem OCR-Rechner oder auf einem anderen per Netzwerk erreichbaren Rechner installiert werden und betreibt offenbar einen ständig laufenden Dienst. Der Lizenzmanager muss während der OCR ständig und nicht nur „einmal je Image“ erreichbar sein; auch bei nur sehr kurzzeitigem Aussetzen der Verbindung zum Lizenzmanager-Rechner wurden HK-OCR-Stapel ohne direkte Fortsetzungsmöglichkeit abgebrochen. Für jede Seite ist größenabhängig mit mehreren Minuten OCR-Analyse bzw. Exporterstellung zu rechnen.

#### *Angebot als Dienstleistung*

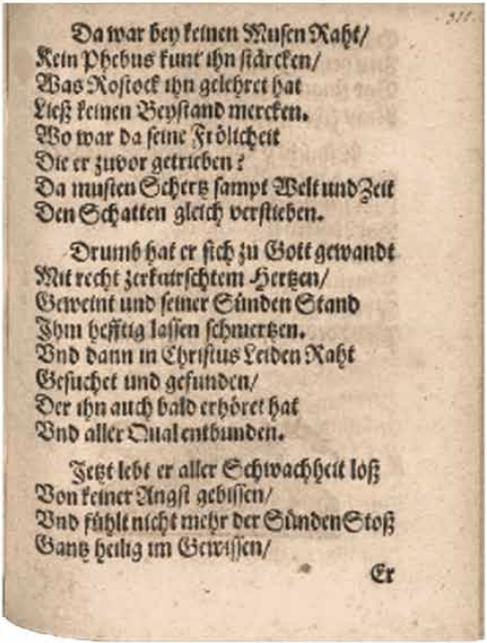
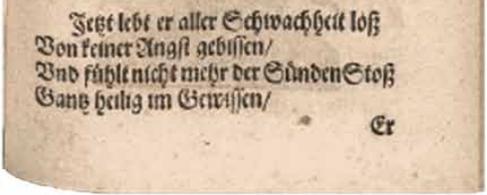
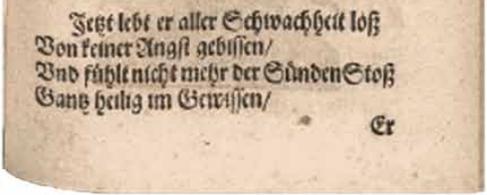
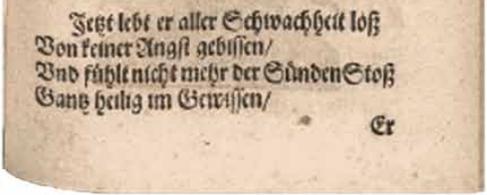
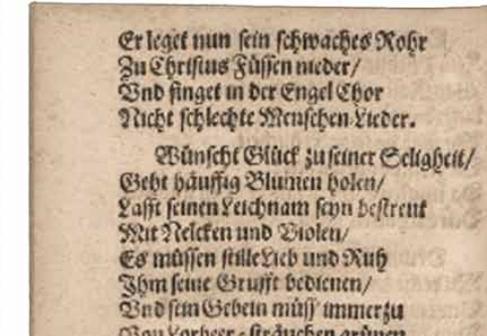
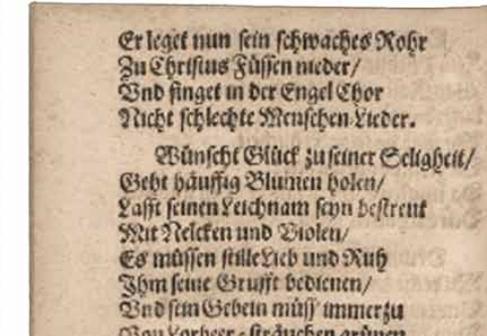
Herrmann & Kraemer bietet die OCR auch als Dienstleistung an, ggf. auch in einem integrierten Workflow mit weiteren Komponenten.

## 5 Praktischer Softwaretest und -vergleich II: Optimierungsversuche durch Sortierung, Binarisierungsparameter, Training und Lexika

### Methodik

Zur Beurteilung der in jeder Software variierbaren OCR-Faktoren waren viele parallele OCR-Läufe erforderlich, die sich jeweils in möglichst wenigen Bedingungen unterscheiden durften.

Als Material verwendet wurden insbesondere eine Marge aus sehr homogenen Gedichtzeilen (Druckort: Königsberg) und eine Marge, in der die zwei nächst kleineren Schriftgrade bei verwandtem Typenbild überwogen. Für bestimmte Musterbibliotheken wurden darüber hinaus bewusst die mit Schriftvarianten überladenen Titelblätter herangezogen, um möglichst schnell eine große Vielfalt zu erreichen (Vgl. Kap. 2.6, Abb. 36).

	HKOCR 110 HUK OG FN s4XL < ff > ss	HKOCR 112 HUK OG17 FN s4XL < 2fQf > Q2ss	PPN 63603108X image 00000007 < 2ff > sss
Raht/ cken/ hat cken. t	Da war bey keinen Musen Raht/ KeinPhebuß kuntihn stärken/ Was Rostock ihn gelehret hat Ließ keinen Beystand mercken. WowardaseineFrölichkeit Die er zuvor getrieben ?	Da war bey keinen Musen Raht/ Kein Phebus kunt ihn stärken/ Was Rostock ihn gelehret hat Ließ keinen Beystand mercken. Wo war da seine Frölichkeit Die er zuvor getrieben ?	
Welt leben.	Da mustenSchertz samptWeltundZeit Den Schatten gleich verstieben. Drumb hat er sich zu Gott gewandt Mit recht zerknirschem Herten/ Geweint und seiner Sünden Stand Ihm hefftig lassen schmerzen. Vnd dann in Christus Leiden Raht Gesuchet und gesunden/ Der ihn auch bald erhöret hat Vnd aller Oval entbunden.	Da musten Schertz sampt Welt undZeit Den Schatten gleich verstieben. Drumb hat er sich zu Gott gewandt Mit recht zerknirschem Herten/ Geweint und seiner Sünden Stand Ihm hefftig lassen schmerzen. Vnd dann in Christus Leiden Raht Gesuchet und gesunden/ Der ihn auch bald erhöret hat Vnd aller Qual entbunden.	
gewandt Hertzen/ Stand tzen. en Raht hat	Drumb hat er sich zu Gott gewandt Mit recht zerknirschem Herten/ Geweint und seiner Sünden Stand Ihm hefftig lassen schmerzen. Vnd dann in Christus Leiden Raht Gesuchet und gesunden/ Der ihn auch bald erhöret hat Vnd aller Oval entbunden.	Drumb hat er sich zu Gott gewandt Mit recht zerknirschem Herten/ Geweint und seiner Sünden Stand Ihm hefftig lassen schmerzen. Vnd dann in Christus Leiden Raht Gesuchet und gesunden/ Der ihn auch bald erhöret hat Vnd aller Qual entbunden.	
heit loß SündenStoß	Jetzt lebt er aller Schwachheit loß Von keinerAngst gebissen/ Vnd süht nicht mehr derSündenStoß Gantz heilig im Gewissen/ Er	Jetzt lebt er aller Schwachheit loß Von keinerAngst gebissen/ Vnd süht nicht mehr derSündenStoß Gantz heilig im Gewissen/ Er	
s Rohr lor Lieder.	Erleget nun sein schwaches Rohr Zu Christus Füßen nieder/ Vnd singet in der Engel Chor Nicht schlechte Menschen Lieder.	Erleget nun sein schwaches Rohr Zu Christus Füßen nieder/ Vnd singet in der Engel Chor Nicht schlechte Menschen Lieder.	
seligkeit/ n/ n bestreut Luh	Wünschet Glück zu seiner Seligkeit/ Geht häufig Blumen holen/ Lässt seinen Leichnam seyn bestreut Mit Nelcken und olen/ Es müssen stille Lieb und Ruh Thun seine Gruff bedienen/	Wünschet Glück zu seiner Seligkeit/ Geht häufig Blumen holen/ Lässt seinen Leichnam seyn bestreut Mit Nelcken und olen/ Es müssen stille Lieb und Ruh Thun seine Gruff bedienen/	

Für beide Schriftgruppen wurden Stichproben festgelegt, die manuell transkribiert wurden und als Referenz zur skriptgesteuerten Ermittlung von OCR-Fehlern der verschiedenen Testläufe dienten. Das Training erfolgte konsequent auf Seiten, die nicht den Stichproben angehörten.<sup>82</sup> Verglichen wurden die unter verschiedenen Bedingungen entstandenen OCR-Ergebnistexte einerseits durch tabellarische Gegenüberstellung der Volltext-Resultate, andererseits durch die Liste der „fehlenden“ und die Liste der „fälschlich gelesenen“ Zeichen einer Seite.

Die Längen dieser Listen dienten als einfache Maßzahlen zum relativen Vergleich verschiedener OCR-Läufe, d. h. um festzustellen, wann eine Konfigurationsänderung einen Fortschritt erbrachte und wann nicht. Als realistische „Fehlerrate“ für die Zeichen-Ebene sind diese Zahlen nicht zu interpretieren.<sup>83</sup>

Wegen der vielen Einflussdimensionen konnten bei weitem nicht alle interessanten Parameterkombinationen getestet werden; oft musste in einer Art Wege-Auswahl frühzeitig entschieden werden, welche günstige Parameteränderung weiterverfolgt wird und welche ungünstige nicht.

(handkorrigiert)	BIT Alpha 159 wie 129 aber Batch10it + BIT12	< 2 > 'z	BIT Alpha 160 wie 129 aber Batch10it + Wortliste3c	< 2z > 'bz	BIT Alpha 157 wie 59 aber Bat
Da war bey keinen Musen Raht/ Kein Phebus kunt ihn stärken/ Was Rostock ihn gelehret hat Ließ keinen Beystand merken. Wo war da seine Frölicheit Die er zuvor getrieben ? Da musten Schertz sampt Welt und Zeit Den Schatten gleich verstieben.	Da war bey keinen Musen Raht / Kein Phebus kunt ' ihn stärken / Was Rostock ihn gelehret hat Ließ keinen Beystand merken . Wo war da seine Frölicheit Die er zuvor getrieben ? Da musten Schertz sampt Welt und Zeit Den Schatten gleich verstieben .		Da war bey keinen Musen Raht/ Kein Phebus kunt ' ihn stärken / Was Rostock ihn gelehret hat Ließ keinen Beystand merken . Wo war da seine Frölicheit Die er zuvor getrieben ? Da musten Schertz sampt Welt und Zeit Den Schatten gleich verstieben .		Da war bey ke Kein Phebus k Was Rostock Ließ keinen B Wo war da sei Die er zuvor g Da musten Sch Zeit Den Schatten Drumb hat er Mit recht zerk Geweint und s Ihm hefftig la Vnd dann in C Gesuchet und o Der ihn auch f Vnd aller Qual
Drumb hat er sich Zu Gott gewandt Mit recht zerknirschem Hertzen/ Geweint und seiner Sünden Stand Ihm hefftig lassen schmerzen. Vnd dann in Christus Leiden Raht Gesuchet und gefunden/ Der ihn auch bald erhöret hat Vnd aller Qual entbunden.	Drumb hat er sich zu Gott gewandt Mit recht zerknirschem Hertzen / Geweint und seiner Sünden Stand Ihm hefftig lassen schmerzen . Vnd dann in Christus Leiden Raht Gesuchet und gefunden / Der ihn auch bald erhöret hat Vnd aller Qual entbunden .		Drumb hat er sich zu Gott gewandt Mit recht zerknirschem Hertzen / Geweint und seiner Sünden Stand Ihm hefftig lassen schmerzen . Vnd dann in Christus Leiden Raht Gesuchet und gefunden / Der ihn auch bald erhöret hat Vnd aller Qual entbunden .		Jetzt lebt er al Von keinerAn Vnd fühlt nich Vnd fühlt nich Gantz heilig ir
Jetzt lebt er aller Schwachheit loß Von keinerAngst gebissen/ Vnd fühlt nicht mehr derSünden Stoß Gantz heilig im Gewissen/ Er	Jetzt lebt er aller Schwachheit loß Von keinerAngst gebissen / Vnd fühlt nicht mehr der Sünden Stoß Gantz heilig im Gewissen / Er		Jetzt lebt er aller Schwachheit loß Von keinerAngst gebissen / Vnd fühlt nicht mehr der Sünden Stoß Gantz heilig im Gewissen / Er		Er
(handkorrigiert)	BIT Alpha 159 wie 129 aber Batch10it + BIT12	< Ein > stu.	BIT Alpha 160 wie 129 aber Batch10it + Wortliste3c	< Sbt > dz.	BIT Alpha 157 wie 59 aber Bat
Er leget nun sein schwaches Rohr Zu Christus Füßen nieder/ Vnd singet in der Engel Chor Nicht schlechte Menschen Lieder.	Erleget nun sein schwaches Rohr Zu Christus Füßen nieder / Vnd singet in der Engel Chor Nicht schlechte Menschen . Lieder .		Erleget nun sein schwaches Rohr Zu Christus Füßen nieder / Vnd singet in der Engel Chor Nicht schlechte Menschen . Lieder .		Erleget nun se Zu Christus Fü Vnd singet in Nicht schlech Wünscht Glüc
Wünscht Glück zu seiner Seligkeit/ Geht häufig Blumen holen/ Lasst seinen Leichnam seyn bestreut Mit Nelcken und Violon/ ES müssen stille Lieb und Ruh Thun seine Geyff bedienen/	Wünscht Glück zu seiner Seligkeit / Geht häufig Blumen holen / Lasst seinen Leichnam seyn bestreut Mit Nelcken und Violon / Es müssen stille Lieb und Ruh Thun seine Geyff bedienen /		Wünscht Glück zu seinrr Seligkeit / Geht häufig Blumen holen / Lasst seinen Leichnam seyn des treue Mit Nelcken und Violon / Es müssen stille Lieb und Ruh Thun seine Geyff bedienen /		Geht häufig E Lasst seinen L Mit Nelcken u

<sup>82</sup> Zu beachten bleibt, dass in OCR-Prozessen immer wieder eine so große Schwankung der Erkennungsrate beobachtet wird, dass jeder Schluss von der Stichprobe auf die Gesamtheit auch hinterfragt werden kann. Die Statistik fordert, dass für solche Inferenzen entweder eine echte Zufallsauswahl oder eine „repräsentative“ (Quoten-)Stichprobe vorliegen muss; beides ist für Images gescannter Buchseiten schwer herzustellen.

Zur Visualisierung der unübersichtlichen Ergebnismengen wurden die verschiedenen konfigurierten Batchläufe in einem Ableitungsdiagramm dargestellt, in dem die gezählten Fehleranteile nicht nur genannt, sondern auch als Farbwert zwischen Rot und Grün dargestellt wurden. So konnten die aussichtsreicheren Zweige der Versuchsanordnung schnell gefunden und auf weniger aussichtsreiche Konfigurationsgruppen früh verzichtet werden (vgl. Abb. 54 u. 58).

### (a) BIT-Alpha

#### **Gruppierung**

Wegen der Möglichkeit, verschiedene Musterdateien zu mischen, ist ein zunächst separates Training relativ reiner Zeichenbestände in BIT-Alpha unproblematisch, wenn die Muster später auch kombiniert genutzt werden sollen. Etwas einschränkend wirkt das Zusammenspiel der Muster mit den eingegeben Sequenzen: verschiedene zu kombinierende Musterdateien sollten auf ähnliche Sequenzen zurückgreifen, um nach dem „Merging“ ihre Erkennungsleistung zu behalten.

Die trainierten Muster einer Schriftgröße haben bis zu einem gewissen Grad auch Zeichen anderer Schriftgrößen richtig klassifiziert, d. h. es gibt eine gewisse Toleranz gegenüber Größenabweichungen.

#### **Binarisierungs- und Segmentierungsparameter**

Wegen der Komplexität der vorzunehmenden Einstellungen wurden die bda-Dateien überwiegend vom Dienstleister konfiguriert und bereitgestellt. Zusätzlich wurden auch im Projekt einige bda-Dateien experimentell abgewandelt oder von Grund auf selbst eingerichtet, um die Wirkung bestimmter Parameter kennenzulernen. Einige Hinweise hierzu finden sich im Anhang 7.3; ein erläuterndes Handbuch bleibt einzufordern.

Durch vergleichende OCR-Läufe verschiedener bda-Dateien auf denselben Seiten mit gleichem Trainingsstand konnten anhand der Fehlerzählung die günstigeren bda-Konfigurationen ausgewählt werden, z. B. wurde so der zur Marge passende Binarisierungsalgorithmus gewählt (hier: „Niblack“).

Gelegentliche sprunghafte Unterschiede in der „Fehlerrate“ derselben Stichprobe konnten auch auf den Ausfall ganzer Absätze hindeuten, wenn bestimmte Segmentierungsparameter auf einzelnen Seiten einen Textblock fälschlich als Bildregion klassifizierten (genauere Hinweise hierzu im Anhang 7.3). Erkennbar war das am ehesten visuell in einer tabellarischen Textübersicht (auffällige Lücken, lange Fehlerlisten) oder durch einen maschinellen Vergleich der Textlänge mit anderen OCR-Läufen derselben Seite (beides wurde von der Software selbst leider nicht bereitgestellt; vgl. Abb. 51).

#### **Training**

Solange man der Benutzerführung im Training folgt und schwerpunktmäßig die nicht erkannten sowie die farbig hervorgehobenen besonders unsicher erkannten Zeichen trainiert, ergibt sich eine automatische Verbesserung der Erkennungsleistung mit fortschreitendem Training von selbst. Im Zweifelsfall kann per manuellem oder automatischem *refresh* der Anzeige nach jedem Trainingsschritt zumindest für den Bereich der aktuell geöffneten Seite sofort kontrolliert werden, ob mit dem Lernfortschritt für das trainierte Zeichen unerwünschte Verschlechterungen der Erkennung anderer Zeichen eingetreten sind (farbige Markierungen ändern sich; vgl. Abb 52).

<sup>83</sup> Auf standardisierte, extern vergleichbare statistische Maße wurde verzichtet, da sie softwaremäßig nicht unmittelbar bereitstanden und ihre aufwendige Einbeziehung in den Software-Test der Fragestellung des Projekts nicht entsprochen hätte. Der wie beschrieben gezählte Anteil „fehlender“ Zeichen kann zumindest als echte Untergrenze der Fehlerrate auf Zeichen-Ebene gelten. Andere Fehlerarten blieben ausgeklammert: z. B. erscheinen Buchstabenvertauschungen („un“ statt „nu“) in diesen Listen nicht; ebenso wurden Weißraum-Fehler wie fehlende oder falsche Spatien bewusst aus der Zählung ausgenommen – sie hätten wegen der sehr unterschiedlichen Verhältnisse im Drucksatz und entsprechenden Abgrenzungsproblemen (ab wann sollen nicht erkannte Wortzwischenräume als OCR-Fehler gelten?) einen unverhältnismäßig hohen und unberechenbar schwankenden Anteil eingenommen und die wichtigeren Differenzen in der echten Zeichenerkennungsleistung vermutlich überdeckt. Die Liste „überschüssiger“ Zeichen kann hingegen Unterschiedliches bedeuten: echte Lesefehler sind es dann, wenn sie falsche Lesarten tatsächlich vorhandener Druckzeichen darstellen; daneben erscheinen hier aber auch Artefakte von fälschlich als Zeichen interpretierten Bildelementen, Durchdruckstellen der Rückseite usw., die eine Nutzbarkeit der so entstandenen Volltexte nicht im selben Maße beeinträchtigen und daher als Fehler schwächer gewichtet werden können.

BIT Alpha 12 SBB_230211-557_V plus Muster von EEE-09	BIT Alpha 13 SBB_230211-557_V Muster von EEE-09 plus bda-Muster	BIT Alpha 14 EEE-09 mit Settings von SBB_230211-557_V
< Maaata > j''''''g'' -5!-...'' '''.isidii si..'''''. ....'rxrlg ..'''.s0	< fumtf > ''''j''-' ...'''.b.. ....iisdii si..'''''' 'rxlnn...' '''''.c.ie( 0)	< upfVkhHes pffggDgenD ielereKrsf fesowirhni htteriteri twenchtss senwzrfei ndetotandT rutsucher
Was . wir hie zuwüfchen pflegejr / Frewde die ohnEnde wehrt / Vnd in keinesHertz hiekehrt / s'	Was wir hie zuwüfchen pflegen / Frewde die ohnEnde wehrt / Vnd in keinesHertz hie kehrt / s'	Gläubt / Gott wird Euch nicht ver- ( lassen / Er ist aller Witwen Schutz / Ewrer Feinde Stoltz vnd Trutz Sucht er endlich heim ohnmasstn /''
DurchdieTochter hohergetzt/' ' /' r	Hört vnd schöpfft erfelbs zugegen '	Weh demselben / deffen^uth
Die Ihm einen Krantz aufftitz .	/' i Durchdie Tochter hoch ergetzt /''	Euch verwegenSchaden thut . .. r^
Gläubt / Gott wirdEuch nichtver - ( lassen / Er ist aller Witwen Schutz / Ewrer Feinde Stoltz vndTrutz Sucht er endlich heim ohnmasstn /''	Die Ihm einen Krantz aufftitz . b .	
r - /5'	Gläubt / Gott wird Euch nicht ver- ( lassen / Er ist aller Witwen Schutz / Ewrer Feinde Stoltz vnd Trutz	
Weh demselben / dessenuth	Sucht er endlich heim ohn jnafft /''	
s -	r -	
Euch verwegenSchaden thut .	Weh demselben / deffen Muth	
.. r	.. ..	

51 Auswirkung falscher Segmentierung auf Fehlerrate



52 Farbliche Darstellung unterschiedlich „sicher“ erkannter Zeichen

Spezieller Evaluierungsbedarf besteht für das in BIT-Alpha mögliche automatische Lernen, bei dem Muster, die (noch) nicht im Trainingsbestand sind, aber einen Zeichenwert mit gewisser „Wahrscheinlichkeit“ klassifizieren, den trainierten Mustern in einem Batch-Lauf (zweckmäßigerweise ohne Export) automatisch hinzugefügt werden. Ob diese Erweiterung der Musterbibliothek real zu einer besseren oder schlechteren OCR-Leistung führt, hängt von der Passung der automatisch hinzugekommenen Muster zur konkreten Vorlage ab und muss an Originalstichproben überprüft werden. Durch die Einstellung des Konfidenzintervalls („OCR > Library > Learning parameters“) kann gesteuert werden, ob das automatische Lernen die Musterbibliothek eher „vorsichtig“ oder eher „großzügig“ erweitert.

Auf diese Weise konnten einige gegenüber der nur manuell trainierten Mustermenge „besser“ erkennende bda-Dateien erstellt werden.

## Wörterbuch

Nachdem die Möglichkeiten der Mustererkennung ausgeschöpft sind, kann eine sorgfältig angelegte lexikalische Korrektur noch einen merklichen Fortschritt bringen. Wieder ist das Zusammenspiel mehrerer Faktoren in BIT-Alpha sehr komplex und erfordert eine gründliche Vorbereitung mit mehrfachem Vergleichstest anhand einer Originalstichprobe. Aufeinander abgestimmt werden müssen:

- eine gut zur Lexik passende Wortliste, die einerseits möglichst vollständig inklusive aller flektierten Formen und vorkommenden Schreibvarianten sein, andererseits möglichst wenig irrelevanten Wortschatz enthalten soll. Es kann versucht werden, einmalige bzw. sehr selten vorkommende Wortformen zu sichten und ggf. auszuschließen. Ob getrennte Wortteile einzeln eingelesen werden können oder zusammengesetzt werden müssen, ist abhängig von den Trennungseinstellungen in der Konfiguration der Lexikalischen Korrektur zu entscheiden. Jedes Vorgehen hat Vor- und Nachteile.
- die Distanz, bis zu der eine Ersetzung überhaupt stattfinden soll (Entscheidung, ob generell eher viel oder wenig ersetzt werden soll);
- die einzelnen Ersetzungskoeffizienten (Erleichterung oder Erschwernis der konkreten Zeichenpaare, die ausgetauscht werden dürfen).

Hier liegt ein hohes Potential an lexikalischer und linguistischer Adaption an das Textmaterial. Die entstehenden Konfigurationen (in lx2-Dateien) können gespeichert werden und wären dann für Material der jeweiligen Sprache und Epoche wiederverwendbar.

<p>BIT Alpha 47 wie 19 ohne Wortliste aber Niblack &lt; fuml̃a &gt; iiiiinn</p> <p>Was weinen wir danvmb jhn viel ? Er wird nicht wieder zu vnskehren . Deslammersist ohn das keinZiel / Wem fehlt eswol anLeidvnd Zehren ? Es wird doch keinem was gespart Ob mmi esgleich nicht offenbahrt .</p> <p>GOttsind invnserm Hertzen stat / Vnd helff' vns allesCreutzbestehen / Wir Armen wissen keinen Raht / Daß wir nicht müstenvntter gehen / Der aber ist der / ehr woldaran Der Todt vnd Welt veriachen kan .</p>	<p>BIT Alpha 58 wie 19 aber Niblack + Wortliste 3b (Split fragments) &lt; fuz &gt; iin</p> <p>Was weinen wir danvmb jhn viel ? Er wird nicht wieder zu vns kehren . Des lammers ist ohn das kein Ziel / Wem fehlt es wol an Leid vnd Zehren ? Es wird doch keinem was gespart Ob mmi es gleich nicht offenbahrt .</p> <p>Gott sind in vnserm Hertzen stat / Vnd helff' vns alles Creutz bestehen / Wir Armen wissen keinen Raht / Daß wir nicht müsten vntter gehen / Der aber ist der / ehr wol daran Der Todt vnd Welt veriachen kan .</p>	<p>FPN 635359987 Image 00000008</p> 
---	--	--

53 Einfluss des Lexikons. Rot über dem Text die Listen der fehlenden (<) und überzähligen (>) Zeichen

Der Aufwand ist allerdings hoch. Im Testprojekt konnten die Möglichkeiten nicht vollständig ausgeschöpft werden; bewährt hat sich die Erstellung von Wörterbüchern aus vergleichbarer Literatur (hier besonders: Personalschriften, Gelegenheitsgedichte und biographisches Material des 17. Jahrhunderts) sowie die Einpflege spezieller Thesauri (hier: zeitspezifische Orts-, Berufs- und Krankheitsbezeichnungen) und bereits vorliegender Katalogdaten (Titel- und Personenfelder) der zu lesenden Schriften selbst, wodurch auch sehr spezifische passende Wortformen ins Wörterbuch gelangten. Die Pflege und Nachbesserung des Wortbestands wird stark behindert durch das binäre, nicht frei editierbare Format der lx2-Dateien; dringend wäre eine freie Bearbeitung der Wortliste in einem offenen Textformat erforderlich, in dem jederzeit Such- und Ersetzungsoperationen, Umsortierungen, Löschung und Austausch größerer Blöcke usw. möglich sind.

Zur Beurteilung des lexikalischen Korrekturschritts sollte neben der Fehleranzahl auch die Art der Fehler betrachtet werden. Wenn etwa eine ursprünglich optisch nahe liegende Verwechslung unter Wörterbucheinfluss zu einem optisch unähnlichen neuen Fehler ersetzt wird, dann wäre trotz gleich gebliebener

Fehlerquote die ursprüngliche Version besser nutzbar als die neue. Insofern können quantitative Maße die Fehler lexikalischer Korrektur nur zum Teil abbilden und gelegentlich sogar in der Vergleichsaussage falsch liegen. Fällt eine falsche Korrektur gar in einen plausiblen Kontext (gesehen z. B.: Vorlage „bei sich“ > OCR „bei stch“ > Korrektur „bei steh“), dann wird auch dem menschlichen Leser, der die fehlerbelassene Version wohl im Sinne der Vorlage richtig interpretiert hätte, eine falsche Lesart aufgedrängt.

## (b) HK-OCR

### *Gruppierung*

Eine erste, auf einer sehr homogenen Schriftgruppe trainierte Mustermenge erkannte die speziell zu dieser Schriftgröße passenden Seiten sehr gut und bereits Seiten des nächst niedrigeren Schriftgrads deutlich schlechter. Die Muster in den ptn-Dateien dürfen demnach als hoch selektiv gelten, was in passenden Einsatzszenarien ein Vorteil sein kann.

Eine Anlage separater später kombinierbarer Musterbibliotheken für bestimmte Schrifttypen und -größen läge daher für HK-OCR bzw. FineReader besonders nahe. Wegen der fehlenden Möglichkeit zum Mischen verschiedener Musterbestände aber erfordert die Mehrfachverwendung eine sorgfältige Planung bzw. erhöht den Aufwand. Ohne weiteres lässt sich z. B. ein Musterbestand für Schrift X trainieren, zunächst in einer X.ptn-Datei speichern, und anschließend könnte für den Einsatz an gemischten Vorlagen der Schriftgruppen X+Y in einer Kopie „XY.ptn“ dieser Musterdatei das Training mit Schriftart Y fortgesetzt werden, um eine kombinierte Bibliothek aus beidem zu erhalten. Eine reine Y-Mustersammlung ist daraus dann aber nicht extrahierbar, wie umgekehrt ein ggf. rein vorliegender Y-Musterbestand nicht der X-Bibliothek hinzugefügt werden kann. Wenn Muster in verschiedenen Mischungen benötigt werden, kann durch diese Beschränkung ein mehrfaches Training bestimmter Vorlagen nötig werden.

Leicht möglich ist andererseits die Kombination der gewählten Anwender-Muster mit den eingebauten FineReader-Mustern, so dass es sich stets lohnt, die drei Alternativen (a) Anwender-Muster, (b) FineReader-Muster, (c) beides zusammen an Originalstichproben miteinander zu vergleichen.

### *Binarisierungsparameter, Training, Wörterbuch*

Wegen der in HK-OCR überschaubaren Konfigurationsmöglichkeiten konnten einige Parametervariationen auf verschiedenen Trainingsstufen systematisch paarweise unter Konstanthaltung der anderen Bedingungen verglichen werden.

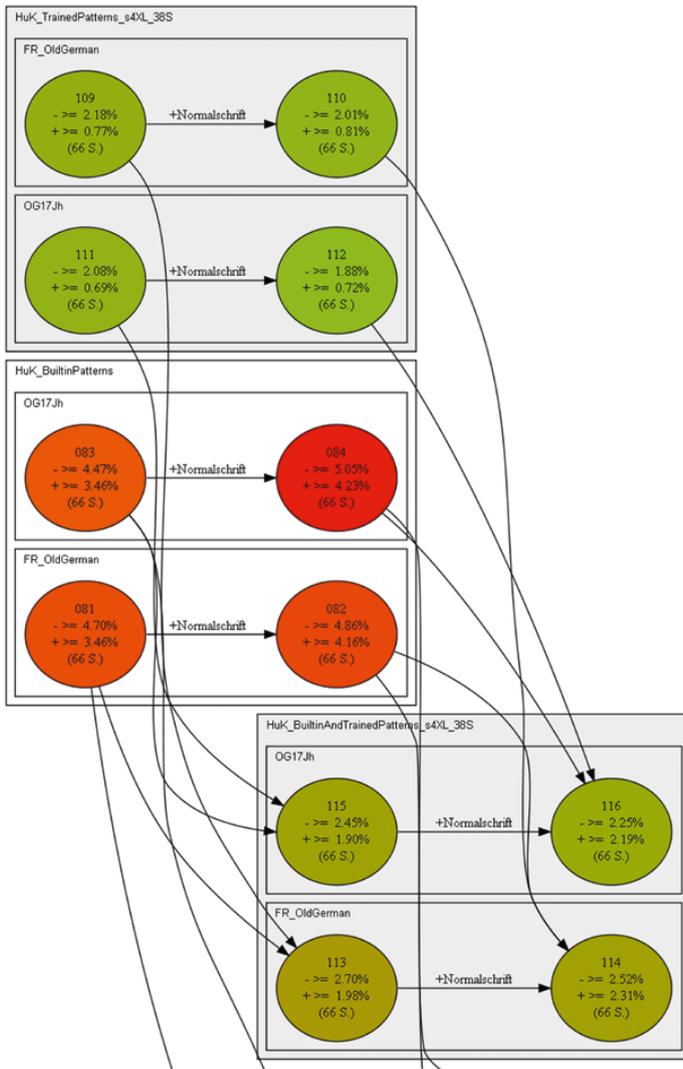
– *Muster (Lohnt sich Training und/oder die Kombination von eingebauten mit trainierten Mustern?)*

Für die Laborsituation der homogenen ersten Marge konnte sehr schnell die Wirksamkeit des Trainings nachgewiesen werden.

Bereits nach einem Training von 37 Seiten aus Gelegenheitsgedichten Simon Dachs war die Erkennung der Stichprobe mit den selbst trainierten Mustern wesentlich besser als die Erkennung mit den FineReader-Mustern allein (die Anzahl fehlender Zeichen fiel auf weniger als die Hälfte, die Anzahl überschüssiger Zeichen ca. auf ein Fünftel)<sup>84</sup> und auch merklich besser als die Erkennung bei Kombination der selbst trainierten mit den FineReader-Mustern (der Anteil fehlender Zeichen sank um ein Fünftel, die Anzahl überschüssiger auf unter die Hälfte). Dieser Effekt war unabhängig davon, ob die FineReader-Sprache „OldGerman“ oder das Anwenderwörterbuch, und ob in den OCR-Optionen allein Fraktur oder die Kombination von Fraktur und Antiqua aktiviert war (vgl. Abb. 54 u. 55).

Für heterogene Muster ist der Fortschritt erwartungsgemäß langsamer, hier wurde nach Training von sechzig Seiten Titelblättern ein Zustand erreicht, in dem der Anteil fehlender Zeichen auch für die OCR der obigen homogenen Gedicht-Seiten bei Verwendung nur selbst trainierter Muster fast an die FineReader-Rate herankam und der Anteil überschüssiger Zeichen bereits deutlich geringer war als bei den FineReader-Mustern. Schon weit vorher, ab ca. zehn Seiten Training, war die Erkennung mit den selbst trainierten Mustern zwar noch deutlich schlechter als die FineReader-Rate, aber bereits hier brachte die

<sup>84</sup> Auch wenn die Fehlerzählung nicht für sich als Gütemaß gelten kann, dürfen die Unterschiede der so gezählten Fehlerfälle wohl als Hinweis auf analoge Bewegungen anderer Fehlermaße interpretiert werden.



Hinzuwahl der selbst trainierten Muster zu den FineReader-Mustern eine stabile Verbesserung.

Für spätere, im 18. Jahrhundert erschienene Schriften des SBB-Funeral-schriftenbestands ändert sich das Bild. Erste Tests deuten an, dass für das späte 17. und beginnende 18. Jahrhundert eine Kombination aus den trainierten und den FineReader-Mustern häufig das bessere Ergebnis erzielt und dass im späteren 18. Jahrhundert häufig die FineReader-Muster auch allein zumindest besser abschneiden können als die trainierten Muster des 17. Jahrhunderts aus dem Projekt. Da speziell zu den Schriften des 18. Jahrhunderts noch keine Proben trainiert wurden, konnte kein echter Vergleich mit solchen stattfinden. Eingehendere Tests wären sicher nicht aussichtslos. Solange allerdings die additive Verwendung speziell trainierter Musterbestände ungeklärt ist, wird der Aufwand ihrer Erstellung wegen der eingeschränkten Nachnutzbarkeit selten vertretbar sein.

54 Graphischer Vergleich aufeinander folgender OCR-Läufe mit verschiedenen Musterdateien

Muster									
"Lohnt sich Training und/oder die Kombination von eingebauten mit trainierten Mustern?"									
Builtin-Patterns		s4XL-Trained Patte		s4XL-Trained + Bui		TitlBI-Trained Patte		TitlBI-Trained + Builtin	
081 HuK OG	4,70	109 HuK OG	2,18	113 HuK OG	2,70	085 HuK OG	12,35	089 HuK OG	5,40
082 HuK OG	4,86	110 HuK OG	2,01	114 HuK OG	2,52	086 HuK OG	14,69	090 HuK OG	5,37
083 HuK OG	4,47	111 HuK OG	2,08	115 HuK OG	2,45	087 HuK OG	12,10	091 HuK OG	4,53
084 HuK OG	5,05	112 HuK OG	1,88	116 HuK OG	2,25	088 HuK OG	14,45	092 HuK OG	4,39
						093 HuK OG	7,19	097 HuK OG	4,24
						094 HuK OG	7,06	098 HuK OG	4,13
						095 HuK OG	7,04	099 HuK OG	4,05
						096 HuK OG	6,95	100 HuK OG	3,91
						117 HuK OG	4,14	145 HuK OG	3,71
						118 HuK OG	4,06	146 HuK OG	3,55
						119 HuK OG	4,14	147 HuK OG	3,55
						120 HuK OG	4,02	148 HuK OG	3,31

55 Vergleich aufeinander bezogener OCR-Läufe mit verschiedenen Musterdateien; die Zahlen zeigen die durchschnittliche Rate fehlender Zeichen pro Seite an

- *Nur-Fraktur- vs. Fraktur-plus-Antiqua-Training* (*„Lohnt es sich, in Fraktur-Musterdateien Antiqua-Muster mitzutrainieren?“*)

Bei der Gelegenheit des Trainings von Titelblättern mit ihrer Vielfalt an Schrifttypen wurde getestet, ob ein Mittrainieren der enthaltenen Antiqua-Zeichen die Erkennung reiner Frakturseiten beeinträchtigt. Wenn ja, würde man darauf zumindest in einer Musterbibliothek verzichten; wenn nein, könnte von Anfang an eine Reserve von Antiqua-Zeichen unbedenklich mitgeführt werden, unabhängig davon, wie viel Antiqua-Einsprengsel der Text tatsächlich enthält.

Wider Erwarten zeigte sich auch auf den Nur-Fraktur-Seiten stets eine leichte Verbesserung der Frakturerkennung, wenn ein gewisser Anteil (im Versuch: ca. fünf Prozent) Antiqua-Zeichen mittrainiert wurde. Eine Erklärung fällt schwer; immerhin bleibt die pragmatische Aussage, hier nicht streng separieren zu müssen (vgl. Abb. 56).

Nur-Fraktur- vs. Fraktur-plus-Antiqua-Training			
"Lohnt es sich, in Fraktur-Musterdateien Antiqua-Muster mitzutrainieren?"			
Nur Fraktur		Fraktur + Antiqua	
093 HuK OG	7,19	101 HuK OG	7,13
094 HuK OG	7,06	102 HuK OG	7,05
095 HuK OG	7,04	103 HuK OG	6,96
096 HuK OG	6,95	104 HuK OG	6,92
097 HuK OG	4,24	105 HuK OG	4,21
098 HuK OG	4,13	106 HuK OG	4,09
099 HuK OG	4,05	107 HuK OG	4,02
100 HuK OG	3,91	108 HuK OG	3,82

56 Zusätzlich trainierte Antiquamuster und durchschnittliche Rate fehlender Zeichen

- *Schriftfamilien-Checkboxen* (*„Lohnt sich die Zuwahl von (eingebauter) Antiqua-Erkennung?“*)
- Ob die in der HK-OCR-Konfiguration mögliche Zuwahl von „Normalschrift“ einen zusätzlichen Musterbestand in die Erkennung einbezieht oder andere Parameter ändert, war nicht dokumentiert. Der Vergleich der Ergebnisse zwischen Nur-Fraktur- und Fraktur-plus-Antiqua-Läufen ergab, dass nur für die FineReader-Builtin-Patterns allein und für noch fast untrainierte eigene Muster durch Zuwahl der Antiqua eine Verschlechterung eintrat, in allen Fällen fortgeschrittenen eigenen Trainings und der Kombination selbst trainierter mit den FineReader-Mustern dagegen eine Verbesserung der Erkennung auch in reinen Fraktur-Seiten. Ohne nähere Kenntnis der Funktion der Normalschrift-Auswahl fällt eine Interpretation dieses Befunds schwer. Zu vermerken ist die Empfehlung, auch für reinen Frakturtext die Normalschrift-Checkbox möglicherweise nutzen zu können (vgl. Abb. 57).

- *Lexikon* (*„Lohnt sich die eigene Wortliste?“*):

In allen ausgewerteten Situationen der ersten Schriftgruppe war eine deutliche Überlegenheit des Einsatzes von altdeutschen Wortlisten gegenüber modernem Wortschatz generell, und fast immer eine leichte Überlegenheit des auf die Funeralschriften angepasst erstellten Lexikons gegenüber dem FineReader-OldGerman festzustellen – obwohl das im Projekt verwendete Wörterbuch bei weitem noch nicht zu Ende optimiert worden war.

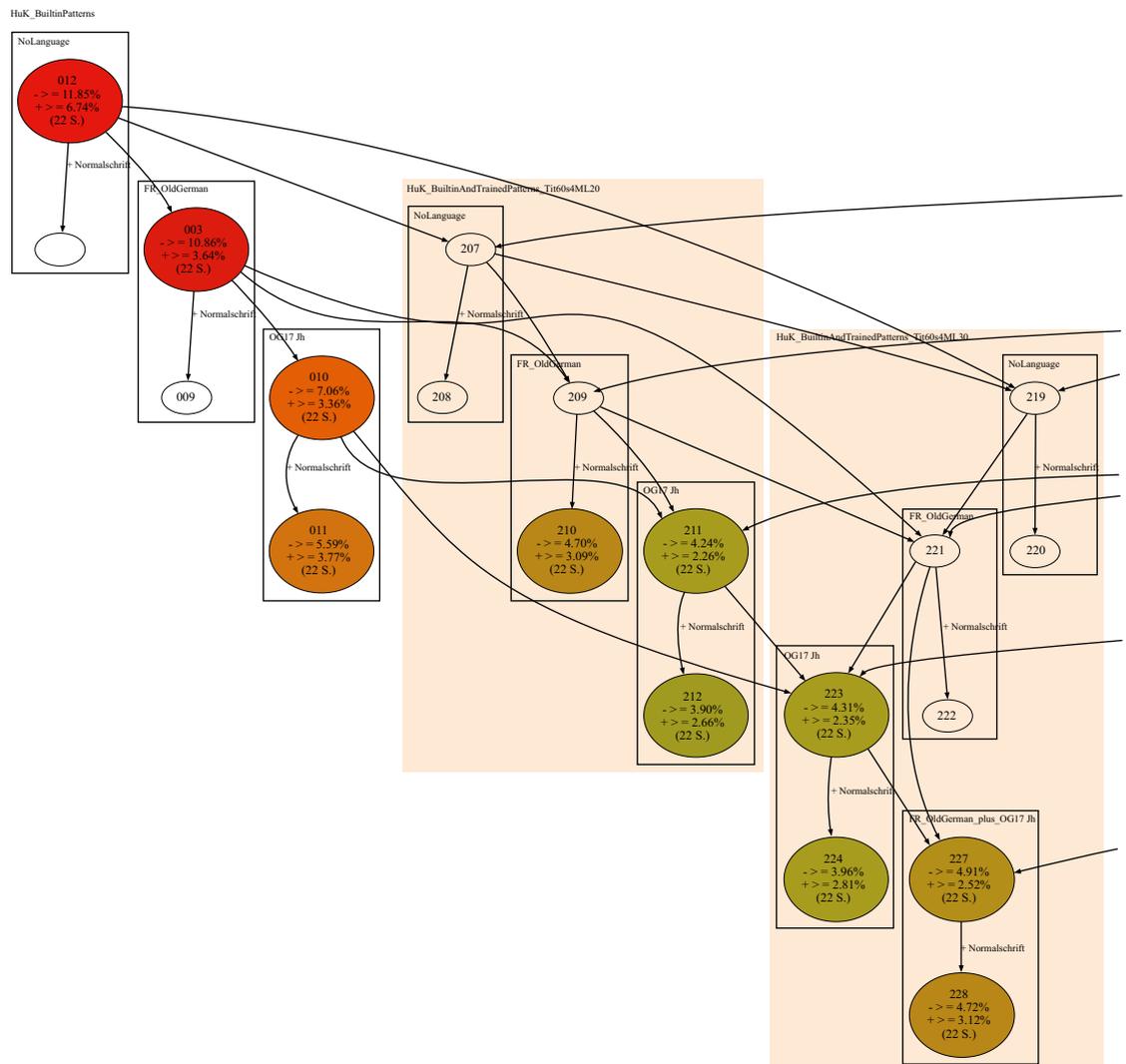
Für die zweite Schriftgruppe (die zugleich aus über Lyrik hinausgehendem Material bestand) fiel

Schriftfamilien-Checkboxen			
"Lohnt sich die Zuwahl von (eingebauter) Antiqua-Erkennung?"			
Fraktur		Fraktur + "Normal"	
081 HuK OG	4,70	082 HuK OG	4,86
083 HuK OG	4,47	084 HuK OG	5,05
085 HuK OG	12,35	086 HuK OG	14,69
087 HuK OG	12,10	088 HuK OG	14,45
089 HuK OG	5,40	090 HuK OG	5,37
091 HuK OG	4,53	092 HuK OG	4,39
093 HuK OG	7,19	094 HuK OG	7,06
095 HuK OG	7,04	096 HuK OG	6,95
097 HuK OG	4,24	098 HuK OG	4,13
099 HuK OG	4,05	100 HuK OG	3,91
109 HuK OG	2,18	110 HuK OG	2,01
111 HuK OG	2,08	112 HuK OG	1,88
113 HuK OG	2,70	114 HuK OG	2,52
115 HuK OG	2,45	116 HuK OG	2,25
117 HuK OG	4,14	118 HuK OG	4,06
119 HuK OG	4,14	120 HuK OG	4,02
145 HuK OG	3,71	146 HuK OG	3,55
147 HuK OG	3,55	148 HuK OG	3,31

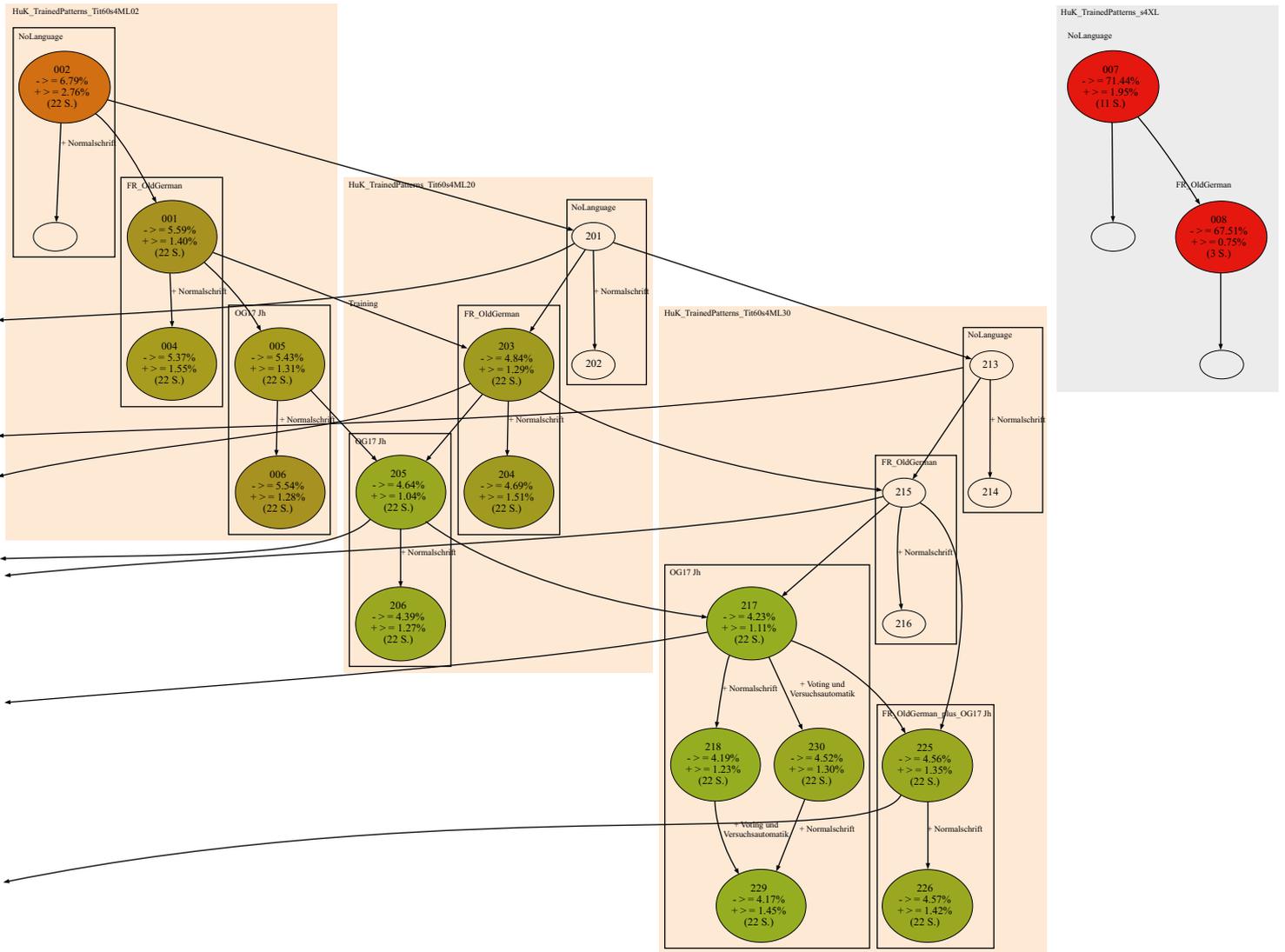
57 Zuwahl mitgelieferter Antiquamuster u. durchschnittliche Rate fehlender Zeichen

der Unterschied noch deutlicher aus. Hier wurde zusätzlich zum Vergleich der FineReader-Sprache „OldGerman“ mit der Anwender-Sprachengruppe auch die in HK-OCR mögliche Kombination beider einbezogen. Sie erreichte die Werte der Anwenderwortbibliothek ebenfalls nicht und blieb meist auf derselben Erkennungsgüte wie die FineReader-Sprache allein. Zu interpretieren ist das wohl so, dass im Fall der Kombination die FineReader-Wortliste klar dominiert: zum Vorteil der Erkennung, wenn sie die passenderen Einträge enthält, sonst zu ihrem Nachteil.

Sowohl die in die Sprach-Definitionen von FineReader eingehenden Wortlisten als auch die Anwenderlexika haben einen starken Einfluss auf die OCR-Lesergebnisse, ohne dass der Anwender hier viel steuern (insbesondere: dämpfen) kann.<sup>85</sup> Auch hier lohnt es sich also, viel Sorgfalt auf die Erstellung des Wörterbuchs zu verwenden, selbst wenn HK-OCR bzw. FineReader bisher keine Feinstuerung der Ersetzungen erlaubt.



<sup>85</sup> Gesehen wurde z.B. eine offenbar lexikongetriebene Hyperkorrektur von zunächst (OCR mit FineReader-Mustern ohne gewählter Sprache:) „dieNen“ zu (FineReader-Muster mit gewählter Sprache:) „vierten“, die nach eigenem Mustertraining bei Verwendung desselben Lexikons richtig als „dienen“ gelesen wurde.



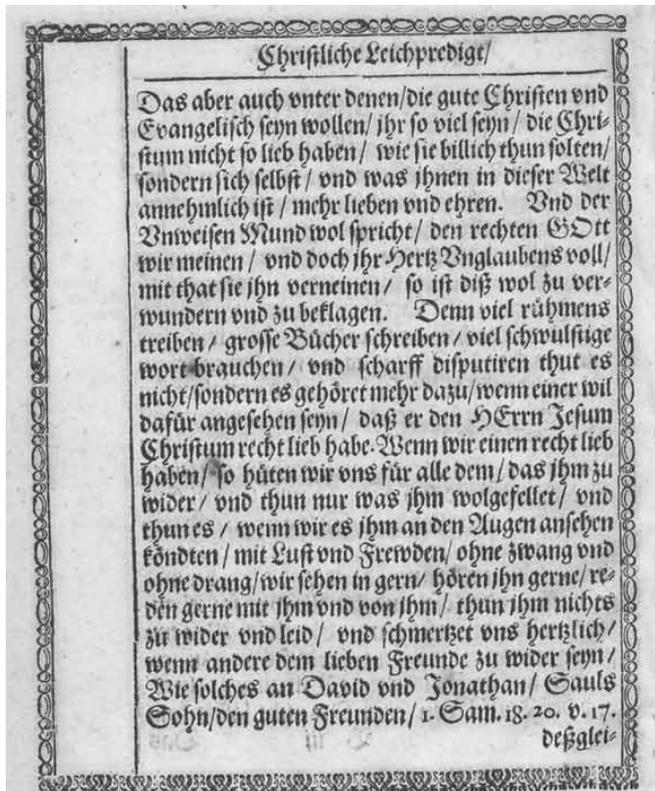
58 Vergleiche in Schriftgruppe 2: Trainingsfortschritt, Musterwahl, Wortbestandwahl und -kombination usw.

## 6 Zusammenfassung der Ergebnisse

### 6.1 Einflussfaktoren auf die Erkennungsgüte

#### 6.1.1 Nachweis des Einflusses von Vorsortierung

Für beide getesteten Softwareprodukte gilt: Seiten, deren Typenvorrat einer passenden trainierten Musterbibliothek zugeordnet werden kann, können auch schon nach mäßigem Training und bei durchschnittlichen sonstigen Randbedingungen eine hohe Erkennungsgüte erreichen:



59 Gute Zuordnung von Drucktypen zu passender Musterbibliothek - bearbeitet mit HK-OCR

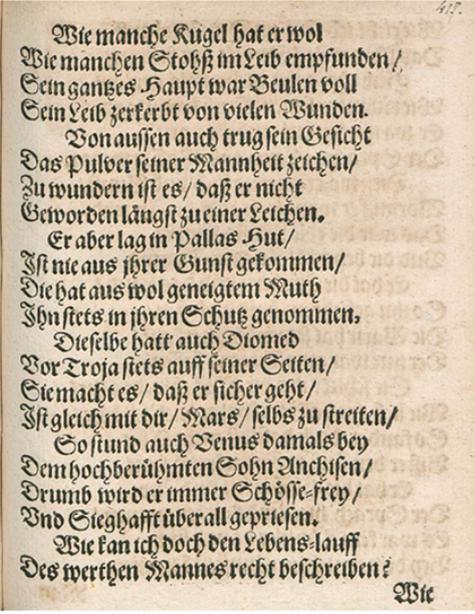
#### Christliche Leichpredigt/

Das aber auch vnter denen/die gute Christen vnd Evangelisch seyn wollen/ jhr so viel seyn / die Christum nicht so lieb haben / wie sie billich thun solten/ sondern sich selbst/ vnd was jhnen in dieser Welt annehmlich ist / mehr lieben vnd ehren. Vnd der Vnweisen Mund wol spricht/ den rechten GOTT wir meinen ^ vnd doch jhr Hertz VnglaubeNS voll/ mit that sie ihn verneinen/ so ist diß wol zu verwundern vnd zu beklagen. Denn viel rühmens treiben / grosse Bücher schreiben / viel schwulstige wort-brauchen/ vnd scharff disputiren thut es nicht/sondern es gehöret mehr dazu/wenn einer wil dafür angesehen seyn/ daß er den HERRN Jesum Christum recht lieb habe. Wenn wir einen recht lieb haben/ so hüten wir ons für alle dem/ das ihm zu wider/ vnd thun nur was ihm wolgefellet/ vnd thun es / wenn wir es ihm an den Augen ansehen köndten / mit Lustvud Frewden/ ohne zwang vnd ohne drang/wir sehen in gern/ hören ihn gerne/ reden gerne mit ihm vnd von ihm/ thun ihm nichts zu wider vnd leid/ vnd schmerzt- vns hertzlich/ wenn andere dem lieben Freunde zu wider seyn/ Wie solches an David vnd Ionathan/ Sauls Sohn/den guten Freunden/ i. Sam. i8. v. i<sup>^</sup> deßglei-

Durch eine entsprechende Gruppierung des Ausgangsmaterials können Stapelverarbeitungen aussichtsreich vorbereitet werden. Ohne Vorsortierung kann diese Erkennungsrate bei weitem nicht erwartet werden, und entsteht die Notwendigkeit, die Spreu erst nach der OCR vom Weizen zu trennen. Dies trifft auch für große Teile der im Projekt verarbeiteten Bände zu, eben weil die Vorsortierung nur in Versuchsmargen (insgesamt ca. 5.000 Seiten) konsequent genug war.

#### 6.1.2 Nachweis des Einflusses von Binarisierungsparametern

Dass Binarisierungs- und Segmentierungsparameter entscheidend für die Möglichkeit einer sinnvollen Musteridentifizierung sind, ist trivial. Zwar ist bemerkenswert, mit wie wenig nutzerseitiger Parametrierung die FineReader-Engine eine *im Durchschnitt* optimale Zeichensegmentierung erreicht, dennoch gehen hier Möglichkeiten zur Nachjustierung verloren. Zu den im Gegenbeispiel BIT-Alpha reichhaltig verfügbaren Optionen und Parametern wird im Anhang versucht, eine vorsichtige Orientierung über sinnvolle Einstellungen und Wertebereiche zu geben (s. Anh. 7.3); verbindliche Information wäre allerdings vom Hersteller anzufordern. An folgendem Beispiel kann verfolgt werden, wie sich die Anzahl fehlender richtiger Zeichen unter dem Gebrauch verschiedener Binarisierungsparameter verändert.

<p>Beispiel B.I.T. Alpha – halbtrainiert, ohne Wörterbuch, Binarisierungsparameter I  <u>BIT Alpha17 SBB 230211-557 V bereinigt Muster von EEE-09 plus bda-Muster</u>  <b>Fehlende Zeichen: 8 mmSoltsr</b>  Überflüssige Zeichen: 11 'ffnnnieni e</p> <p>Wie manche Kugel hat erwoi  Wie manchen Stoß iniLeib empfunden /  Sein ganzes Haupt war Beulen voll  Sein Leib zerkerbt von vielen Wunden .  Vonaußen auch trug fein Gesicht  n^  Das Pulver feiner Mannheit zeichen /  Zuwundern ist es / daß er nicht Geworden längst zu einer Leichen .  Er aber lag in Pallas Hut /  Ist nieaus jhrer Gunst gekennen /  Die hat auswol geneigtem Muth Ihn stets in jhren Schutz genommen .  Dieselbe hatt ' auch Diomed  VorTroja stets auffseiner Seiten /  Istgleich mit dir / Mars / selbszu streiten /  So stund auch Venus damals bey Dem hochberühmten Sohn Anchisen /  Vnd Sieghafft überall gepriesen .  Wie kan ich doch den Lebens - lauff beschreiben ?  Wie</p>	<p>Leichenpredigt 1649</p> 	<p>Beispiel B.I.T. Alpha – halbtrainiert, ohne Wörterbuch, Binarisierungsparameter II  <u>BIT Alpha59 wie 27 ohne Wortliste aber Niblack</u>  <b>Fehlende Zeichen: 3 ene</b>  Überflüssige Zeichen: 10 '..oliccii</p> <p>Wie manche Kugel hat erwol  Wie manchen Stoß imLeib empfuliden /  Sein gantzes Haupt warBeulen voll  Sein Leib zerkerbt von vielen Wunden .  Vonaußen auch trug sein Gesicht  n^  Das Pulver seiner Mannheit zeichen /  Zuwundern ist es / daß er nicht Geworden längst zu einer Leichen .  Er aber lag in Pallas Hut /  Ist nieaus jhrer Gunst gekommen /  Die hat aus wol geneigtem Muth Ihn stets in jhren Schutz genommen .  Dieselbe hatt ' auch Diomed  VorTroja stets auffseiner Seiten /  Istgleich mit dir / Mars / selbszu streiten /  So stund auch Venus damals bey Dem hochberühmten Sohn Anchisen /  Vnd Sieghafft überall gepriesen .  Wie kan ich doch den Lebens - lauff Des werthen Malinesrecht</p>
--	---	---

60 Texterkennung nach verschiedenen Binarisierungsverfahren in BIT-Alpha

### 6.1.3 Nachweis des Trainingseffekts

BIT-Alpha, das ohne Muster ausgeliefert wird, beruht auf der Wirksamkeit des Trainings passender Muster, so dass der hier während der Versuche mit zunehmendem Training beobachtete Erkennungsfortschritt – bzw. die Differenzen bei Verwendung verschiedener Musterbibliotheken – selbstverständlich ist.

Interessanter war die Feststellung, wie sehr sich ein materialbezogenes Training auch bei der mit guter Frakturmuster-Grundversorgung ausgestatteten FineReader-Engine auswirkte. Das zeigen nicht nur die im Text beschriebenen Fehlerzählungen auf den vorliegenden Ground-Truth-Stichproben, sondern auch die in den Layout-Protokollen der FineReader-Engine vermerkten Erkennungsraten (offenbar Anteile der im Wörterbuch enthaltenen Wörter an den insgesamt erkannten). In den Bänden des 17. Jahrhunderts, die sowohl mit Finereader-Mustern als auch mit selbst trainierten Mustern vergleichend prozessiert wurden, wurden für die eigene Musterdatei – trotz der begrenzten Trainingsintensität von 90 Seiten – oft als um zwei bis fünf Prozentpunkte höhere Raten vermerkt als im Lauf mit den FineReader-Frakturmustern.

? L k 8 v n ä l. l. z. "

Weil aber die Furcht Gottes ohne die Liebe des Nächsten z anders nichts als ein thönend Ertzt und klingende Schelle z so hat Sie auch gegen Zhre Neben-Khrsten alle Khrstliche Tugenden erblicken lassen. Gegen die Vor nehmen und höhern ist Sie gewesen ehrerbietig z gegen Zhres gleichen freundlich z gegen die Niedrigen demüthig z gegen die Armen milde z und hat nach dem lobwürdigen Exempel Ihres sel. herm Vaters keinen Nothleydenden/ derumb eine Gabe gebeten/ leer von sich weggehen lassen. Gegen Ihre Freunde aufrichtig Z gegen die Geschwister treuhertzg Z ZHr Gemüth war zur Weißheit derWeltkin- derzdie honig im Munde/und Gall in Ihren hertzen füh ren/ so einfaltig/ wie eine Taubez Sie Ueß weder in Ihren Worten noch Gebärden einige Falschheit henschen Z und in Summa es werden schon Ihres Lebens und Wandels le bendige Zeugen genug verhanden seyn. Zhre Kranckheit betreffende Z so hat man schon etliche lahr her beylhr einen bösen ü ffeK u m verspmet/m nächst verwichenen kwrw aber Z ist Sie auch an einem Fieber krancworden/deßwegen der verordnete Landes p<sup>^y</sup>^cuz» hen l<sup>^ic</sup>enc. Samuel Sturm cvnlulizet wordenz der auch alk hierzu dienliche ^ecjjcamenrgverschleben/unddurch Verleihung göttlicher hülfte es endlich so weit gebracht/ daßSiedasFiebereine<sup>^t</sup>langverlassen/jedochbißweilen wiederkommen. Nichts destoweniger aberz hat Sie gleich- h iij wohl

Gattungsspezifische Wortliste + FineReader-Muster

61 Fortschritt in OCR-Erkennung (hier mit HK-OCR) durch Training

PERSONALIA. ^

Weil aber die Furcht Gottes ohne die Liebe desNäch<sup>^</sup>sten /<sup>^</sup> anders nichts als ein thönend Ertzt und klingende S chelle z so hat Sie auch gegen Ihre Neben-Christen alle Christliche Tugeden erblicken lassen. Gegen die Vor nehmen und Höhern ist Sie gewesen ehrerbietig /<sup>^</sup> gegen Ihres gleichen freundlich /<sup>^</sup> gegen die Niedrigen demüthig <sup>^</sup> gegen die Armen milde /<sup>^</sup> und hat nach dem lobwürdigen Exempel Ihres sel. Herrn Vaters keinen Nothleydenden/ derumb eine Gabe gebeten/ leer von steh weggehen lassen., Gegen Ihre Freunde aufrichtig /<sup>^</sup> gegen die Gefchwister treuhertzg/ Ihr Gemüth war zur Weißheit derWellkm-der/die Honig im Munde/und Gall in Ihren Hertzen füh- ren/ fo einfaltig/ wie eine Taube/ Sie ließ weder in Ihren Worten noch Gebärden einige Falschheit herrschen /<sup>^</sup> und in Summa es werden schon Ihres Lebens und Wandels le bendige Zeugen genug verhanden seyn. Ihre Kranckheit betreffende /<sup>^</sup> so hat man schon etliche lahr her beylhr einen bösen A ssect u m verspüret/in nächst verwichenen Martio aber/ ist Sie auch au einem Fieber krancworden/deßwegenderverordnere Landes Physicus, Herr Licent. Samuel Sturm consuliret worden/ der auch al<sup>^e</sup> hierzu dienliche ^edica<sup>^</sup>eneaverschrieben/ unddurch Verleihung göttlicher Hülffe es endlich so weit gebracht/ daßSiedasFiebereinezeitlangverlassen/jedochbißweilen wiederkommen. Nichts destoweniger aber/ hat Sie gleich<sup>^</sup> H ii<sup>^</sup> wohl

Gattungsspezifische Wortliste + anwendertrainierte Muster

## 6.1.4 Nachweis des Wörterbucheffekts

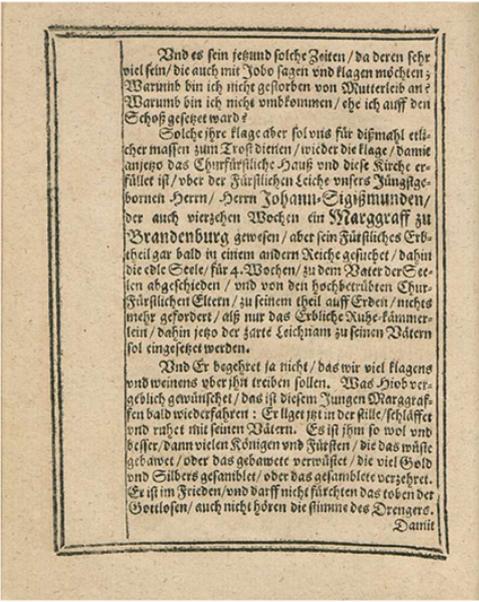
Durch gezielte Gegenüberstellung von OCR-Ergebnissen, die mit bzw. ohne Hinzunahme von Wortlisten entstanden, lässt sich bei beiden Softwareprodukten zeigen, dass sich der Einsatz von Wortbibliotheken positiv auf die Ergebnisse auswirkt. Dieser Effekt beruht in BIT-Alpha auf einer vom Anwender detailliert steuerbaren lexikalischen Korrektur nach der OCR. Bei HK-OCR dagegen beeinflusst die Wortbibliothek durch Definition so genannter „Sprachen“ und „Sprachengruppen“ der FineReader-Engine offenbar schon die Texterkennung. Der maßgebliche Einfluss der Wortliste zeigt sich vor allem in der Erkennung der Wortgrenzen. Vergleicht man das Beispiel in den verschiedenen Bearbeitungsstufen wird vor allem durch die korrigierten Wortgrenzen der Text flüssiger lesbar. Es fällt ebenso auf, dass durch Einsatz der Wortliste<sup>86</sup> neue Fehler entstehen können: „dieNen“ wird fälschlich zu „vierten“.

Die im Projekt erstellte gattungsspezifische Wortliste bestand in der OCR-Optimierungsphase aus ca. 70.000 Einträgen. Wie eingangs beschrieben, enthält sie vor allem Wortschatz aus Leichenpredigten. Dieser Wortschatz entstand durch Transkription von Referenzseiten und speist sich zusätzlich aus Titelmateriale aller über den GBV nachgewiesenen Leichenpredigten. Zudem ergänzt um Gedichte von Simon Dach und einige wenige zeitgenössische Texte wurden vor allem die Thesauri der Forschungsstelle für Personalschriften Marburg THEPRO und THELO sowie eine Liste historischer Krankheitsbezeichnungen des Vereins für Computergenealogie<sup>87</sup> herangezogen. Eine weitere Optimierung in der Testphase selbst konnte nicht vorgenommen werden.

Wörterbücher können deutlich genauer erstellt werden, wenn das Ausgangsmaterial regional, zeitlich und inhaltlich auf den Verwendungszweck abgestimmt im Zuge einer fehlerfreien Erfassung entstanden ist. Ein weiterer Weg erschließt sich über die Hinzunahme aller möglichen Flexionsformen bzw. Rechtschreibvarianten der vorkommenden Wörter (s. Kap. 3.1.5 und 4.3.4).

<sup>86</sup> Die Ersetzung von Wörtern über das Wörterbuch in HK-OCR ist nicht steuerbar.

<sup>87</sup> VEREIN FÜR COMPUTERGEALOGIE

Beispiel HK-OCR/FR9 – ohne Training, ohne „Sprache“	Leichenpredigt 1625	Beispiel HK-OCR/FR9 – ohne Training mit selbst def. „Sprache“
<p><b>HKOCR012 NoLang F Builtins</b>  <b>Fehlende Zeichen: 108</b>                  ö?ß?üßüüü Lülßüüüßä                  Läjüü-jzä äjönüüne üeweeewee                  ewüeeelnle eeleeblee zererFreen                  nrnürendbe nderlenhnh hörendeede  <b>DrenerDa</b>  <b>Überflüssige Zeichen: 108</b>                  !JOTN!ZNNc MccccN,ccu                  cukccukckck cWWHcucuMc                  ffccZffcuM Vkvfcvzv Svmtmkvttm                  ZvvffctMkc tcsktsftu cctSistimm                  sSsZ)mik</p> <p>Vndes fein                  feftundsolcheZsiten/daderen schr                  vie! fein/ die auch mit Jobs sagen                  vnd klagm mochten ; Warumb bin                  ich nichkgsstorben von Mutterletb                  ani Warumb bin ich nicht                  vmbkommen/ ehe ich auff den                  Schoßgesetzet wards iOolchejhr                  klagcaber folvns fur distmahl tli-                  cher massen zum Trost dieNen /                  tvicdr die kiagc / vamt anjctzo das                  ChnrfursilichHaust vnd diese                  Kirche er-fullet ist/vber der                  Fursilichenieiche vnserslungstgc-                  bornett Herrn/ HerrnZohanN-                  StgHMtNden/ der such vierzchen                  Wochen ein Marggraff zu                  Brandenburg</p>		<p><b>HKOCR010 OG17 F Builtins</b>  <b>Fehlende Zeichen: 32 ??-----L</b>                  --äSFeeloe ueeedeener Dt  <b>Überflüssige Zeichen: 43</b>                  ZOTZZZZZZ Z,ZZZMkcZc                  cMZZZZZZv cfcfZccssZ jlk</p> <p>Vnd es fein jetzund solche Zsiten Z                  da deren schr viel sein/ die auch                  mit lobo sagen vnd klagen möchten                  ; Warumb bin ich nicht gestorben                  von Mutterletb ane Warumb bin ich                  nicht vmbkommen / ehe ich auff                  den Schoß gesetzet ward s                  iOolchejhrclage aber sol vns für                  dißmahl etli cher massen zum Trost                  vierten Z wieder die klage z damit                  anjctzo das Churfürstliche Hauß                  vnd diese Kirche er fullet ist/vber                  der Fürstlichenieiche vnsers                  lüngstge bornen HerrnZ Herrn                  lohann-Sigißmunden/ der auch                  vierzehen Wochen ein Marggraff zu                  Brandenburg</p> <p>...</p>

62 Fortschritt in OCR-Erkennung (hier mit HK-OCR) durch Wortliste

## 6.2 Im Projekt erreichte Erkennungsgüte

Eine wie im Beitrag von Thomas Stäcker (vgl. 8.7) dargestellte statistisch abgesicherte Fehleranalyse wurde bisher nicht durchgeführt. Gleichwohl wurden Stichproben aus den verschiedenen Projektphasen untersucht. Diese zum Teil erst aus der Zeit nach der Optimierung entstandenen Daten sollen hier kurz vorgestellt werden, um Eindrücke anhand von kompletten Beispielen zu erlauben. Die Fehlerzählung erfolgte analog zu der in der Herzog August Bibliothek Wolfenbüttel durchgeführten Methode (Worttrennungen, fehlende bzw. überschüssige Spatien fanden keine Berücksichtigung).

Wie bereits geschildert (s. Kap. 5) wurde in der Optimierungsphase nur in sehr bescheidenem Maße trainiert, so dass grundsätzlich starke Probleme bei der Erkennung von in Marginalien gesetzten Texten auftraten, was zum einen an der sehr kleinen und oft schwer lesbaren Schrift, zum anderen an der mangelnden Segmentierung lag. Auch die Erkennung von Ziffern und Zeichen in fremden (als den im Trainingsmaterial vorkommenden) Alphabeten muss daher als unbefriedigend bezeichnet werden. Die folgenden Fälle sind zu unterscheiden:

### 1. Marge:

Homogene Texte von Simon Dach in einer Schrift. Hierzu wurden vorab 50 Seiten Text trainiert. *BIT-Alpha* (vgl. Abb. 63):

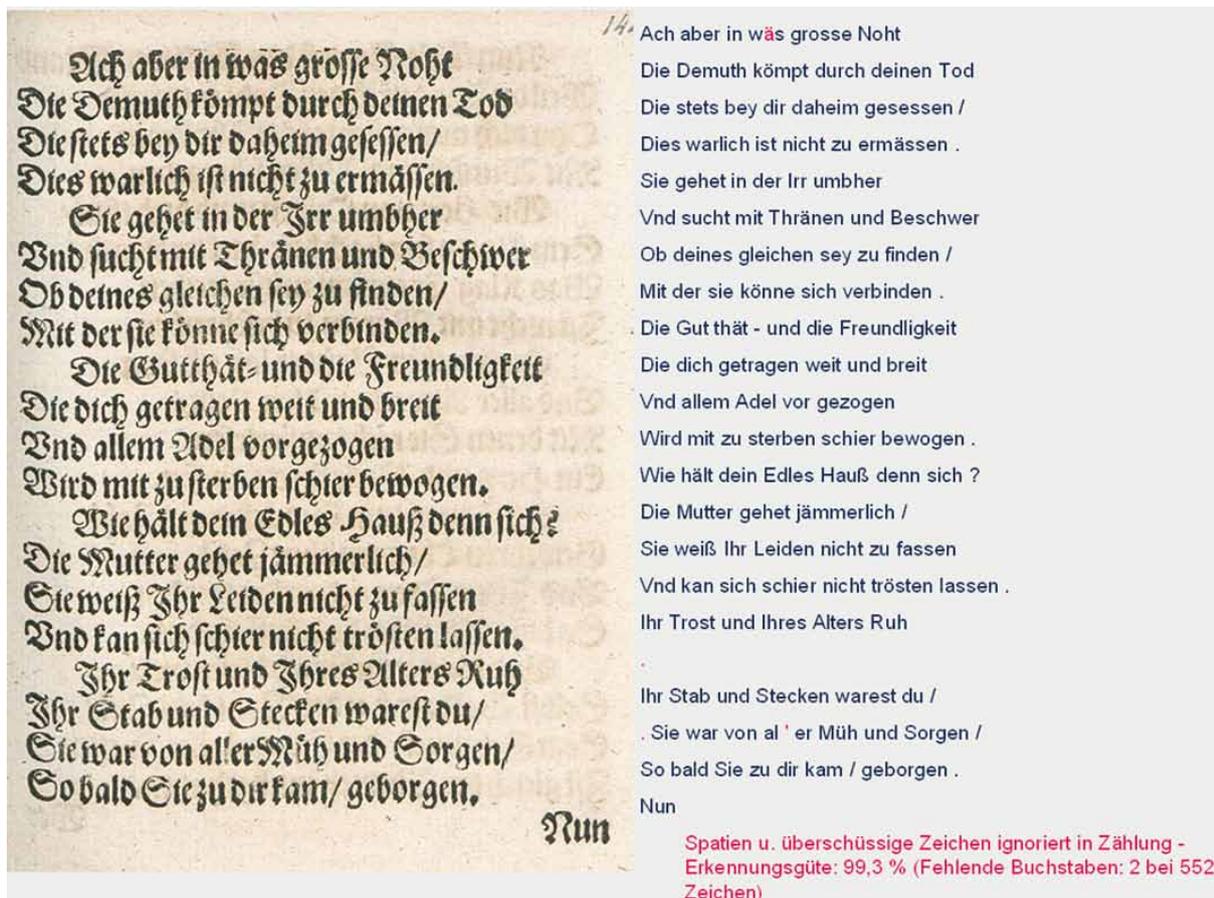
Von sieben wahllos herausgesuchten Seiten wiesen fünf eine Erkennungsgüte von 97,7 % und höher auf; in zwei Fällen wurden über 99 % erreicht.

### 2. Marge:

Weniger homogene Texte von Simon Dach mit wenigen unterschiedlichen Schriften. Hierzu wurden vorab 60 Titelblätter und 30 Seiten Text trainiert.

*HK-OCR:*

Hier wurden nur drei Seiten stichprobenartig untersucht. Die Ergebnisse lagen bei 93,6 % bis 96,4 %.



63 Beispiel BIT-Alpha; 1. Optimierungslauf: 40 in: Yi 851-3, S. 2

### 3. Marge (Massenlauf):

Nicht vorselektierte Texte Funeralschriften vornehmlich aus dem 17. Jh. Kein weiteres Training hat stattgefunden; die vorab erstellten Musterdateien der 2. Marge wurden für das gesamte Material angewendet.

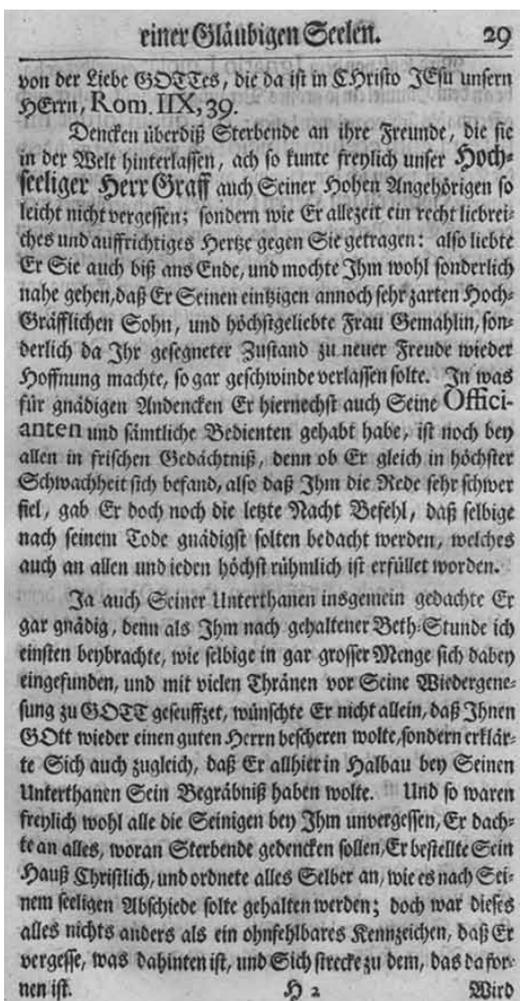
*BIT-Alpha* (vgl. Abb. 64):

Von acht einzelnen Seiten wiesen nur vier eine Erkennungsgüte über 90 % auf (91,8 % bis 95,8 %).

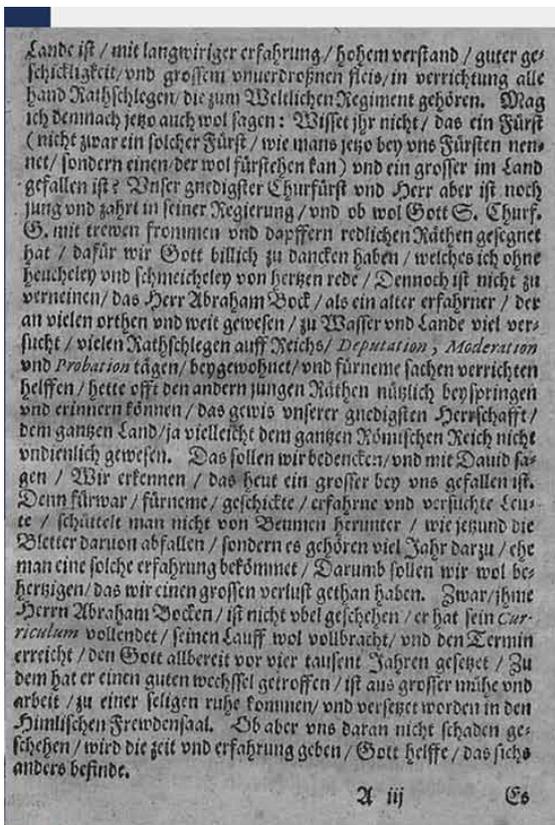
*HK-OCR* (vgl. Abb. 65):

Hier wurden nur sechs Seiten stichprobenartig untersucht. Die Ergebnisse lagen bei fünf Seiten deutlich über 94 % (88 % bis 97,9 %).

Eine Hochrechnung erlauben diese willkürlich erhobenen Daten kaum. Es gilt außerdem zu berücksichtigen, dass im so genannten Massenlauf keine weiteren Selektionen mehr vorgenommen werden konnten, was sicherlich bei weiterem Training eigens getrennter Gruppen zu Qualitätssteigerungen geführt hätte. Speziell für BIT-Alpha ist diese Art der Prozessierung nicht die angezeigte, da wie bereits oben geschildert, die Software in besonderem Maße ein Training erfordert und nicht auf einen großen Fundus vorhandener Muster zurückgegriffen werden kann.



64 Beispiel BIT-Alpha; Massenlauf: Ee 705-1071, S. 29



65 Beispiel HK-OCR; Massenlauf: Ee 705-106, S. 7

ey  
einer Glänbigen See Jen .

oi , n aer giebc q ; OeTes , die da ist in icHrisio : 3Esii unsern HEern , Rom . IIX , g9 .

Dencken überdiß Sterbende an ihre Freunde , die sie in der Welt hinterlassen , ach so kunte freylich unser Hochseeliger HerrGrtff auch Sciner Hohen Angehörigen so leicht nichtvergessen ; sondern wie Er allezeit ein recht liebliches und auffrichtiges Herl ; e gegen Sie getragen : also liebte Er Sie auch biß ansEnde , und mochtelhm wohl sonderlich nahe gehen , daß Et Seinen einzigen annoch sehr zarten HochGräfflichen Sohn , und höchstgeliebte Frau Gemahlin , sonderlich da Ihr gesegneter Zustand zu neuer Freude wieder Hoffnung machte , so gar geschivinde verlassen sollte . Inittias für gnädigen Aindencken Er hiernächst auch Seine Officianten und sämtliche Bedicnten gehabt habe , ist noch bey allen in frischen E ) eddchniß , denn ob Er gleich in höchster Schwachhcitsich befand , also daß lhm die Redc sehrschwefiel , gab Er doch noch die letzte Nacht Befehl , daß selbige nach seinem Todc gnädigst solten bedacht werden , welches auch an allen und iedele höchstrühmlich ist erfüllet worden . Ja auch Seiner Linterthanen insgemeiAt gedachte Et gar gnädig , denn als lhm nach gehaltenerBeth - Stunde ich einsten beybrachte , wie selbige in gar grofferMenge sich dabey eingefunden , und mit vielen Thränen vor Seine Wiedergenesung zuGOTT gesenfftet , wunschte Er nichtallein , daß Ihnen GOTT wieder einen guten Herrn bescheren wolte , sondern erklärte Sich auch zugleich , daß Er allhierin Halbau bey Seinen Unterthanen Sein Begräbniß haben wolte . Und so waren freylich wohl alle die Seinigen bey lhm undergessen , Er dachtean alles , woran Sterbende gedencken sollen , Er bestesiteSein Hauß Christlich , und ordnete alles Selber an , wie es nach Sei -

Spatien u. überschüssige Zeichen ignoriert in Zählung -  
Erkennungsgüte: 97,2 % (Fehlende Buchstaben: 47 bei 1664 Zeichen)

Lande ist / mit langwüiriger erfahrung / hohem verstand / guter geschickligkeit / vnd grossem vnuerdroßnen fleis / in verrichtung alle hand Rathschlegen / die zum Weltlichen Regiment gehören. Mag ich demnach jetzo auch wol sagen : Wisset ihr nicht / das ein Fürst ( nicht zwar ein solcher Fürst i wie mans jetzo bey vns Fürsten nennet / sondern einen / der wol fürstehen kan ) vnd ein grosser im Land gefallen ist? Vnser gnedigster Churfürst vnd Herr aber ist noch jung vnd zahrt in seiner Regierung / vnd ob wol Gott S. Churf. G. mit trewen frommen vnd dapffern redlichen Rätthen gesegnet hat / dafür wir Gott billich zu dancken haben / welches ich ohne heuchelei vnd sthmeicheley von hertzen rede / Dennoch ist nicht zii verneinen / das Herr Abraham Bock / als ein alter erfahrner / der an vielen orthen vnd weit gewesen / zu Wasser vnd Lande viel versuecht / vielen Rathschlegen auff Reichs / Deputation , Moderation vnd Probation lägen / beygewohnet / vnd fürneme sachen verrichten helffen / heue oft den andern jungen Rätthen nützlich beyspringen vnd ei innern können / das gewis vnserer gnedigsten Herrschafft / dem gantzen Land / ja vielleicht dem gantzen Römischen Reich nicht vndienlich gewesen. Das sollen wir bedencken / vnd mit Dauid sagen / Wir erkennen / das heut ein grosser bey vns gefallen ist. Denn fürwar / fürneme / geschickte / erfahrne vnd versuchte Leute / schüttelt man i nicht voit Beunien herunter / wie jetzund die Bletter damon abfallen / sondern es gehören viel Jahr darzu / ehe man eine solche erfahrung bekömmet / Darumb sol. en wir wol behertzen / das wir einen grossen verlust gethan haben. Zwar / ihme Herrn Abraham Bocken / ist nicht vbel geschehen / er hat sein Curriculum vollendet / seinen Lauff wol vollbracht / vnd den Termin erreicht / den Gott allbereit vor vier tausent Jahren gesctzet / Zu dem hat er einen guten wchssel getroffen / ist aus grosser mühe vnd arbeit / zu einer seligen ruhe kommen / vnd ver setzet worden inden Himlischen Frewdensaal. Ob aber vns daran nicht schaden geschehen / wird die zeit vnd erfahrung geben / Gott helffe / das sich anders befinde.

A iij Es

^ softwareseitig nicht erkanntes Zeichen; Spatien u. überschüssige Zeichen ignoriert in Zählung - Erkennungsgüte: 98,3 % (Fehlende Buchstaben: 27 bei 1791 Zeichen)

## 6.3 Hervorzuhebende Vor- und Nachteile der getesteten Software für die OCR deutschsprachiger Alter Drucke

### 6.3.1 *Good practice*

#### (a) BIT-Alpha

BIT-Alpha zeichnet sich aus durch:

- eine reiche Parametervielfalt in Binarisierung und Segmentierung, wodurch auf Layout- und Schrifttypeneigenschaften detailliert eingegangen werden kann;
- die Möglichkeit, auch Schrifttypen-Teile (Schaft, Fuß usw.) als eigene Muster auf „Special characters“ zu trainieren, um die Teile später mittels Sequencer zusammensetzen (so können bestimmte häufige OCR-Fehler wie brüchige „n“, „m“, „u“ repariert werden);
- die mit dem Hilfsprogramm BIT-Knowledge gegebene Möglichkeit zur Mischung mehrerer Musterbibliotheken;
- eine fein steuerbare Wortkorrektur durch gewichtete Zeichenersetzungskoeffizienten.

#### (b) HK-OCR/FREngine 9

HK-OCR besitzt:

- eine aus dem Stand (ohne langwierige Justierungen) gegebene sehr gute Segmentierung auf Zeichenniveau und
- eine gut bedienbare Benutzeroberfläche für Wortkorrekturen mit Berücksichtigung der Korrekturen in der Koordinatenberechnung des XML-Exports.

### 6.3.2 Gegenwärtige Nachteile

#### (a) BIT-Alpha

Die wichtigsten Nachteile von BIT-Alpha sind darin zu sehen, dass

- unnötige „Compound“-Formate der Parameter- und Hilfsdateien (s. Kap. 4.1) zur manchmal unerwünschten Zusammenspeicherung von nicht zusammengehörigen Informationen führen und
- zu den zahlreichen Optionen bzgl. Parameteränderungen eine gründliche Dokumentation und eine übersichtliche Vorschau der Auswirkungen solcher Justierungen fehlte.

#### (b) HK-OCR/FREngine 9

Als für die getesteten OCR-Aufgaben ungünstig erwiesen sich folgende Eigenschaften von HK-OCR, die offenbar auf Komponenten der FineReader-Engine beruhen:

- die stark beschränkte Navigation im Trainingsteil (s. Kap. 4.4);
- die unpraktische Bedienung des Mustereditor-Fensters;
- die fehlende Kombinierbarkeit trainierter Musterbestände sowie
- die Erschwernis durch die Doppel-Lizenzierung einerseits der Software beim Hersteller, andererseits der von ABBYY zu erwerbenden und zu installierenden Einzelseitenlizenzen.

## 6.4 Anforderungen an Software-Weiterentwicklung

Aus Sicht eines an realistischen Produktionsbedingungen für OCR interessierten Anwenders haben einige Produkteigenschaften gefehlt, die für eine Stapel verarbeitende Software eigentlich bewährt sind und die auch mit dem Programmdesign beider Produkte ohne weiteres vereinbar wären.

Ein Anliegen dieser Studie ist es daher, zur Diskussion über nahe liegende Anforderungen bei der Weiterentwicklung von für historische Drucke geeigneten OCR-Tools beizutragen und eine aktive Nachfrage nach wesentlichen Produktmerkmalen anzuregen, die den Herstellern eine adressatengerechte Produktentwicklung ermöglicht.

Die folgenden, für die getesteten Programme BIT-Alpha und HK-OCR den Herstellern sinngemäß mitgeteilten (und in deren Software meist ohne weiteres integrierbar scheinenden) Anregungen können zugleich als allgemeiner Kriterienkatalog für die Auswahl bzw. Entwicklung von OCR-Produkten mit herangezogen werden.

### Aufrufsituation im Betriebssystem

Könnten dem jeweiligen Programm als konventionelle Aufrufparameter schon beim Start alle Angaben (Namen der Konfigurations-, Muster- und Wörterbuch-/Sprachdateien, ebenso Bildverzeichnisse bzw. Bilddateinamen) übergeben werden, sinngemäß etwa:

```
OCRAPP.exe -conffile xx -patternfile yy -lexiconfile zz-outputdir xx/xx image1.tiff
OCRAPP.exe -conffile yy -patternfile zz -lexiconfile xx-outputdir xx/yy image2.tiff
```

...

Dann könnten Batchprozesse viel flexibler gestaltet werden, indem z. B. im selben Stapel zu bestimmten Bilder(gruppe)n bestimmte Konfigurationen oder Muster-/Wortbibliotheken übergeben werden. Auch für Beschränkungen, die aus der gegenwärtig „nur aus einem gemeinsamen Bildverzeichnis heraus“ möglichen Batch-Verarbeitung resultieren, wäre die explizite Angabe der zu prozessierenden Images hilfreich. Dadurch könnten die Images an ihrem Speicherort verbleiben, ohne dass der Bearbeiter ein eigentlich unmotiviertes (und gelegentlich wegen fehlender Rechte oder aus Speicherplatzgründen gar nicht mögliches) temporäres Hin-und-Her-Schieben von Quelldateien organisieren muss. Demzufolge wäre ein Anstoßen beliebig kleinteiliger Test- und Produktionsläufe mit sehr unterschiedlichen Bedingungen kein zeitlicher Faktor mehr.

Zudem wären die OCR-Anwendungen mit dieser konventionellen Aufrufsyntax so auch in beliebige Workflow-Systeme integrierbar, an denen Anwenderinstitutionen aus unterschiedlichen Gründen festhalten wollen. Die Entscheidung zwischen bibliothekseigenem Workflow und der Nutzung von Produkt XY wäre dann nicht mehr nötig.

### Konfigurationsspeichermodell; Formate; Zusammenspiel von Parameter-, Muster- und weiteren Hilfsdateien

Da für eine an einzelne Vorlagen angepasste OCR die eingesetzten Binarisierungsparameter-, Muster-, Ersetzungsregel-, Sprach- und Lexikon-Dateien optimal ausgewählt und kombiniert werden müssen, sollten sie am besten:

- unabhängig voneinander einzeln abgespeichert werden;
- soweit möglich jederzeit auch außerhalb des OCR-Prozesses passiv lesbar sein (so für Protokoll- und Vergleichszwecke, z.B. wenn Nutzer ihre längerfristig Muster-, Parameter- und Regelsammlungen anlegen und pflegen);
- und je möglichst komfortabel editierbar sein.

Im Ergebnis sollte die gezielte Änderung oder Löschung bestimmter Parameter/Muster/Ersetzungsregeln/Wortlisteneinträge jederzeit möglich sein, ohne eine OCR-Analyse anzustoßen und ohne Auswirkungen auf nicht-betroffene Parameter/Muster/Regeln. Im Sinne „atomarer“ Datenablage sollte für jede Steuerinformation nur ein eindeutiger Speicherort vorgesehen sein.

Für alle diese Hilfsdateien (außer den Mustern selbst) gilt, dass eine Speicherung in einem textartigen, les- und editierbaren Format (einfache Listen, Schlüssel-Wert-Paare wie in \*.ini- oder \*.config-Dateien, notfalls XML) möglich sein sollte und folgende Vorteile hätte:

- Es wäre erstmals ein durchgehendes *monitoring* und *logging* der im konkreten OCR-Lauf wirksamen Parameter, Sequenzen, Wortlisten oder Ersetzungskoeffizienten möglich – eine Voraussetzung für sinnvolle langfristige Weiterpflege durch die Anwender.
- Parameterlisten, Sequenzen, Wortlisten, lexikalische Ersetzungen usw. könnten bequem und wesentlich flexibler ergänzt, korrigiert oder in ausgewählten Teilen gelöscht werden.
- Zur Übertragung bewährter Parameter oder Regeln in eine andere Konfiguration ist bisher der manuelle Neueintrag jedes einzelnen Werts in zeitaufwendigen Dialogen der graphischen Benutzerschnittstelle nötig. Aus textartigen Hilfsdateien aber wäre eine selektive und dennoch kompakt „blockweise“ Übertragung von Parametersätzen von einer Konfigurationsdatei in eine andere möglich.<sup>88</sup>

Insbesondere ist kein Grund erkennbar, ausgerechnet auch die vom Anwender selbst in Textform oder durch Tastatureingaben bereitgestellten Wortlisten, Listen erlaubter/verbotener Buchstaben, Ersetzungsregeln usw. in proprietären, vom Anwender nicht mehr lesbaren Binärformaten zu verwalten. Die sachgerechte Pflege dieser Bestände gestaltet sich dadurch unbegründet schwierig bis unmöglich.

Aber auch für die Speicherung der Muster selbst wäre ein dokumentiertes und anwendungsübergreifendes Format anzustreben, das die freie Bearbeitung und nach Möglichkeit sogar den Austausch von Mustern zwischen verschiedenen OCR-Anwendungen erlaubt. Die Formate von Open-Source-OCR-Anwendungen wurden nicht untersucht, könnten aber vermutlich zur Formulierung entsprechender Anforderungen herangezogen werden. Bei geeigneter Serialisierung (Aufeinanderfolge von Datensätzen) wären Kombination und Durchsuchbarkeit solcher Musterdateien trivial.

### **Benutzeroberfläche und Ergonomie**

Während des Pilotprojekts ist viel Zeit mit ineffektiven Bedienungsabläufen verbraucht worden. Das betrifft nicht nur die OCR-bezogenen Softwareeigenschaften, zu denen im folgenden Anregungen genannt werden, sondern offenbar geraten längst bewährte und etabliert geglaubte Merkmale ergonomischen Programm-Designs in Neuentwicklungen gelegentlich in Vergessenheit. Umso wichtiger ist es, dass Anwender, die durch Aufträge an Softwarehäuser Einfluss auf die Programmgestaltung nehmen können, sich rechtzeitig über die Mindestanforderungen klar sind, die für eine Arbeitsweise mit vertretbarem Personalaufwand unverzichtbar sind.

Da die Wirtschaftlichkeit von maschineller OCR eng mit der für Training und Konfiguration aufzuwendenden Arbeitszeit zusammenhängt, sind im Anhang (s. Anh. 7.4) einige solcher allgemeiner Mindestanforderungen genannt, auf die geachtet werden sollte.

Die Zusammenarbeit mit beiden Dienstleistern hat gezeigt, dass Hinweise durchaus willkommen sind – schließlich sind Produkte, die einer schnellen technischen Entwicklung folgen müssen, praktisch immer „in Entwicklung“. Bei konsequenter Nachfrage und rechtzeitiger Übernahme gewünschter Eigenschaften z. B. in ein Pflichtenheft würden diese Merkmale sicher nicht fehlen.

<sup>88</sup> Speziell die Stärke der B.I.T.-Programme, unterschiedliche Kombinationen von Binarisierungsparametern, Mustern, Sequenzen, Wortlisten, Lexikalischen Ersetzungsregeln usw. überhaupt zu ermöglichen, könnte so noch viel konsequenter eingesetzt und produktiv gemacht werden.

OCR-spezifisch wären folgende Anliegen zu nennen:<sup>89</sup>

***Zeitsparende Orientierung über aktuell geltende Eigenschaften, Parameter usw. und ihre Auswirkungen***

Wenn eine reiche Konfiguration von Binarisierungs- und Segmentierungsparametern angeboten wird, dann wird wegen ihrer Vieldimensionalität eine Beurteilung der Wirkung einzelner Änderungen fast unüberschaubar und zwingt bei der Konfiguration eines OCR-Stapels zur Beschränkung auf erfahrungsgestützte Vermutungen, ohne diese auch nur stichprobenartig mit vertretbarem Aufwand verifizieren zu können.

Eine Vorschau auf zu erwartende Ergebnisse könnte leicht gewonnen werden, wenn die Programme

- wie oben beschrieben alle nötigen Parameter auch in der Kommandozeile akzeptieren würden und
- die bisher nur in der Benutzeroberfläche erreichbaren Vorschau-Ansichten als Bilddateien exportieren könnten.

Denn dann könnte durch Batch-Testläufe wie

```
OCRAPP.exe -conffile xx image1.tiff --printlayer=BlockLayout
> image1-confxx-blocklayout.png
OCRAPP.exe -conffile yy image2.tiff --printlayer=CharLayout
> image2-confyy-blocklayout.png
OCRAPP.exe -conffile yy image2.tiff --printlayer=OCR
> image2-confyy-charlayout.png
USW.
```

sehr schnell (jedenfalls ohne von Bearbeitern aufzuwendende Zeit) ein repräsentativer Überblick über die konfigurationsbezogene Segmentierungsgüte und Layouterkennung einer größeren Bildmenge gewonnen und z. B. in einer automatisch erstellten Tabelle visuell vergleichend dargestellt werden.<sup>90</sup> Vermutlich würde es durch solch eine „Viele-Screenshots-Ansicht“ überhaupt erst möglich, die Auswirkung einer konkreten Parameteränderung auf einen vorliegenden Gesamt-Stapel halbwegs kontrolliert abzuschätzen. Bisher entscheidet man über die Anpassung eines Parameters anhand einer eher sehr kleinen Menge gesehener und einzeln analysierter Bilder (dabei mit Minuten Wartezeit je Bild) und muss sich für die Prognose auf die Intuition verlassen – mit der Gefahr, bei der „Rettung“ einer Seite die Erkennung vieler anderer zu beschädigen.

***Interaktivität, frühe Evaluation der Einzelverarbeitungsschritte***

Für alle Stufen gilt: Wo adaptive Selbstkonfiguration und automatische Korrektur nicht stattfinden oder eher heuristisch, d. h. nur möglicherweise optimal arbeiten, könnte Software mit unmittelbaren Warnmeldungen und Hinweisen auf sinnvolle interaktive Eingriffe die schnelle Findung der geeigneten Einstellungen und Konfigurationen unterstützen.

Speziell die Segmentierung sollte auch in der Stapelverarbeitung justierbar präsentiert werden, um auf druckspezifische Schwierigkeiten bei der Layouterkennung reagieren zu können. Gebraucht würden einerseits „halbfertige“ Konfigurationsvorlagen, wo der Anwender bestimmte Abstandsmaße, Eingangstexte und ähnliche Kriterien festlegen kann, anhand derer bestimmten Blöcken *für den aktuellen Stapel* eine bestimmte Textstruktureinheit zugewiesen wird,<sup>91</sup> andererseits eine am Bildschirm gut geführte Möglichkeit zur schnellen graphischen Verschiebung von Blockgrenzen sowohl für Einzelseiten als auch auf Beispielseiten mit Geltung für einen Stapel. Voraussetzung wäre hierfür u. a. eine Möglichkeit, Images eines Stapels auf dem Bildschirm in Schichten übereinander darzustellen und so Gruppen mit Layoutgemeinsamkeiten (einfachster Fall: linke Seiten vs. rechte Seiten) bilden zu können.

<sup>89</sup> Vgl. außerdem die einzelschrittbezogen genannten Desiderata in den Abschnitten des Kapitels 4.

<sup>90</sup> Hingewiesen sei auch auf Vorarbeiten im IMPACT-Projekt zu Evaluierungen bereits jeder Teilstufe des OCR-Prozesses. Das 2012 als Nachfolgeeinrichtung von IMPACT entstandene Kompetenzzentrum sollte als Ansprechpartner in Anspruch genommen werden.

<sup>91</sup> Unter den Konfigurationsparametern von BIT-Alpha gibt es etliche, die hierfür geeignet wären.

### ***Alternative: Vorsorgliche Parallelverarbeitung mit späterer Evaluation***

Eine in OCR-Algorithmen ansatzweise bereits anzutreffende Parallelverarbeitung<sup>92</sup> könnte als Vorbild dienen, bei der OCR inhomogener Vorlagen generell unterschiedliche Lesegänge durchzuführen und deren Ergebnisse für einen anschließenden Vergleich bereitzustellen.

Insoweit in solchen parallelen OCR-Outputs dieselben Teilblock- (Absatz- usw.) Segmentierungen erreicht und ausgezeichnet wurden, sollte hier sogar ein über die Selektion des besseren Gesamtseiten-Ergebnisses hinausgehender Mehrwert erreichbar sein, indem trotz seitenweiser OCR die Evaluation auf Blockebene ansetzen könnte und sich z. B. für einen Absatz der Output aus dem OCR-Lauf mit Konfiguration A, für einen anderen Absatz der Output aus dem OCR-Lauf mit Konfiguration B als passender erweist.

## **Öffnung der Blackboxes**

### ***Transparenz und Konfigurierbarkeit der wirkenden Faktoren***

Alle sinnvoll variierbaren Faktoren sollten für den Anwender konfigurierbar sein und müssen in jedem Fall – auch nachträglich und auch dann, wenn die Einstellung automatisch erfolgte – visuell dargestellt und protokolliert werden können. Erfolge an einem OCR-Stapel sind nur nachnutzbar, wenn sich die Umstände reproduzieren lassen.

### ***Modularität***

Wo immer ein Anwenderinteresse an der Ausgabe von Teilergebnissen oder an der Einspeisung anderen Materials vorliegt, sollte die geschlossene Programmarchitektur geöffnet und eine voraussetzungslose Nutzung auch von Einzelkomponenten ermöglicht werden.

Es gibt beispielsweise keine Notwendigkeit, die gute Binarisierungsleistung einer Software nur im Zusammenhang mit dem OCR-Lauf in derselben Software zu erlauben.

Die Funktion der Korrekturkomponente sollte nur dann an das Vorhandensein der Original-Bilddateien geknüpft sein, wenn tatsächlich neue OCR-Analysen notwendig sind. Soweit hingegen nur die textartigen OCR-Leseergebnisse verarbeitet werden, sollten diese – ohne erneute OCR-Analyse – eingespeist werden können; idealerweise auch Material, das aus anderen Quellen als der OCR des Programms selbst kommt. Gerade wenn ein Korrekturmodul spezielle Stärken hat (HK-OCR seine gut bedienbare „Validierungs“-Oberfläche, BIT-Alpha seine gut justierbare automatische Nachkorrektur anhand gewichteter Zeichenersetzungsregeln), wäre sein Einsatz auf beliebigem Input wünschenswert.

Die jederzeitige Ausgabe von Zwischenergebnissen wird auch schon zur Evaluation von Teilschritt-Konfigurationen benötigt.

## **6.5 Schlussfolgerungen und Tipps; Ausblick**

### **Informationsverluste als „Leitfaden“**

Der vermutlich ergiebigste Hinweis für die OCR Alter Drucke dürfte sein, dass nur eine adäquat vorbereitete OCR Aussicht auf zufrieden stellende Ergebnisse bietet und dass hierzu eine gründliche Kenntnis des Materials notwendig ist.

Je nach Verwendungszweck<sup>93</sup> können unterschiedliche Qualitätsanforderungen bestehen, zu deren Überprüfung zunächst von der gewählten OCR-Software unabhängige Kriterien und Werkzeuge gefunden werden sollten. Zur Mindestausstattung von OCR-Projekten gehören daher Standardwerkzeuge für Textextraktion und einfache statistische Auswertungen; hierfür können die in allen Betriebssystemen gängigen Kommandozeilenfilter und eine Skriptsprache mit komfortabler Textverarbeitung wie Perl bereits genügen.

<sup>92</sup> Zu nennen ist hier z. B. in HK-OCR das mehrfache Lesen unter unterschiedlichen Bedingungen („Layouts“) mit anschließender manueller oder automatischer Auswahl des „besten“ Ergebnisses, s. Kap.4.3.3“

<sup>93</sup> Eine Übersicht der wichtigsten Verwendungsziele gibt z. B. Ralf Stockmann (SUB Göttingen) in einem Vortrag „Was tun mit den Ergebnissen der OCR?“ im Rahmen des IMPACT-Projekts; darunter u. a. eine nach steigender „Sichtbarkeit“ geordnete Aufzählung der Bereitstellungsarten: „Versteckt in Suchindex“ - „Versteckt, aber Image-Highlighting der Fundstelle“ - „Volltext als Layer hinter dem Image (etwa in PDF gebunden)“ - „Volltext sichtbar über / neben dem Image“ - „Nur Volltext sichtbar“ - „Volltext als Download“ - „Volltext für Harvester verfügbar (TEI Datei in OAI)“, s. STOCKMANN 2010

Neben der „absoluten“ Qualitätsbeurteilung (zur Entscheidung, ob und wann sich ein OCR-Projekt lohnt) sind innerhalb der Arbeit häufig relative Entscheidungen zu treffen, z. B. in welche Richtung eine Konfiguration geändert werden soll, wofür es nicht auf die Perfektion von Fehlerkriterien ankommt, sondern nur darauf, ob die Tendenz einer Verbesserung oder Verschlechterung der Erkennungsrate abgelesen werden kann. Vergleiche müssen beliebig oft (und jeweils schon für Veränderungen nur eines Faktors) vorgenommen werden können, um aus den Ergebnissen die richtigen Schlüsse zu ziehen.

Für die Suche nach geeigneter Software oder nach Dienstleistern können dann die im Kapitel 3 genannten typischen Fehlerquellen der OCR als Entscheidungshilfe dienen. Interessenten sollten sich klar werden, welche der absehbaren Informationsverluste auch im eigenen konkreten Material auftreten werden und welche möglicherweise unproblematisch sind. Daraus folgernd könnten sie dann die Software oder den Dienstleister danach auswählen, wie sie mit den betreffenden Fehlerquellen umgeht.

Außerhalb der eigentlichen OCR-Leistungsfähigkeit werden Fragen der Nachhaltigkeit und der Bindung an Lizenzen und Formate eine Rolle spielen wie:

- Korrespondieren von Import- und Exportformaten mit den eigenen Geschäftsgängen;
- Wiederholbarkeit der OCR- und Korrekturläufe;
- Wiederverwendbarkeit von Konfigurationsdateien, Musterbibliotheken und Korrekturlexika;
- Auswirkungen einer eventuellen Softwareinstallation auf betriebliche Abläufe;
- usw.

Nicht zuletzt lohnt es sich, stets die denkbaren Auftragsmodi anhand der konkreten OCR-Aufgabe gegeneinander abzuwägen:

- Arbeit in eigener Regie durch Erwerb bzw. Lizenzierung der Software;
- Beauftragung der Gesamt-OCR als Dienstleistung (hier sollte eine Zielerkennungsgüte vertraglich festgelegt werden; s. auch Kap. 8.8);
- Arbeitsteilige Zusammenarbeit: z. B. mit eigener Produktion, eigenem Training, aber Beauftragung der Konfiguration schwieriger Parameter<sup>94</sup> – oder umgekehrt: externe Beauftragung von Stapelverarbeitungen, die der Anwender vorher selbst konfiguriert hat.<sup>95</sup>

### **6.5.1 Tipps für Anwender der getesteten Software: Wann wäre welches der getesteten Produkte geeignet?**

Um sich *von der Software ausgehend* ein Bild zu machen, werden die Kapitel 4 und 5 empfohlen, wo in den einzelnen Abschnitten auf die Merkmale von HK-OCR und BIT-Alpha eingegangen wird.

Zur *Nachnutzbarkeit* von nutzereigenen Muster- und Wortbibliotheken gilt für beide Produkte ähnlich, dass die erstellten Bibliotheken prinzipiell aufbewahrt und wiederverwendet werden können. Selbstverständlich muss langfristig auf die Kompatibilität etwaiger Nachfolgeversionen der jeweiligen Programme geachtet werden. Eine wirklich flexible Nachnutzbarkeit für nicht vorhergeplante Fälle wird wie beschrieben durch die proprietären, vom Nutzer nur bedingt editierbaren Speicherformate eingeschränkt.

Wenn man *vom Einsatzbedarf* ausgehen möchte, könnten aus den Versuchen im Pilotprojekt folgende Empfehlungen abgeleitet werden:

#### **Wenn eigenes Training erfolgen soll:**

- BIT-Alpha eher für Szenarien, in denen
  - der zu erkennende Zeichenvorrat (insbesondere Diakritika, griechische, hebräische Zeichen u. ä.) im voraus nicht vollständig bekannt ist oder
  - es im Training darauf ankommt, frei und mit häufigen Sprüngen navigieren zu können.

<sup>94</sup> Nachteil dieser Variante wäre, dass unerwünschte Softwareeigenschaften wie eine kompliziert dargebotene oder schlecht dokumentierte Konfiguration honoriert würden.

<sup>95</sup> Ein möglicher Vorteil: die professionelle Materialkenntnis kann optimal eingesetzt werden.

- HK-OCR eher für Szenarien, in denen
  - der zu erkennende Zeichenvorrat (insbesondere Diakritika, griechische, hebräische Zeichen u. ä.) im voraus vollständig bekannt ist und in den Spracheinstellungen angegeben werden kann oder
  - ein Training überwiegend ohne Sprünge, d. h. fortlaufend über die ganze Seite erfolgen kann.

#### **Für eine OCR-Prozessierung von Anfang bis Ende:**

- BIT-Alpha eher für Szenarien, in denen
  - der Anwender vorlagenbedingt Einfluss auf Binarisierung und Segmentierung nehmen möchte oder muss und Zeit für die Konfiguration hat (Alternative: Konfiguration in Form einer bda-Datei beim Softwarehersteller bestellen; hier sollten Qualitätskriterien vereinbart werden);
  - eine ausreichend große ungefähr passende Mustermenge trainiert werden kann oder vorliegt;
  - ein die Vorlage weitgehend abdeckendes Lexikon zur Verfügung steht und eine automatische, gewichtet einstellbare lexikalische Korrektur bevorzugt wird.  
(Als Alternative können Konfiguration und/oder Training und/oder Stapelverarbeitung (mit/ohne Wortkorrektur) beim Softwarehersteller als Dienstleistung bei B.I.T. in Auftrag gegeben werden; hier empfiehlt sich, Zielmaße für die zu liefernde Erkennungsgüte zu vereinbaren.)
- HK-OCR eher für Szenarien, in denen
  - das Bildmaterial keine speziellen Schwierigkeiten für Binarisierung oder Segmentierung aufweist, die eine individuelle Änderung von Parametern erfordern würden;
  - eine sehr gut zur Vorlage passende Mustermenge trainiert werden kann oder vorliegt und die Stapelverarbeitungen sehr gut nach homogener Sprache, Schriftfamilie und -größe portioniert werden können;
  - eine manuelle lexikalische Korrektur bevorzugt wird bzw. auf eine automatische lexikalische Korrektur entweder verzichtet werden kann oder diese nachgelagert – dann ohne Mitpflege der Wortkoordinaten – vorgenommen wird.  
(Die Übernahme von Trainingsleistungen wurde zum Projektzeitpunkt nicht von der Herstellerfirma angeboten. Wenn der Anwender aber Musterdateien und Lexika bereitstellen kann, könnte auch hier die Verarbeitung als Dienstleistung beim Softwarehersteller Herrmann & Kraemer vereinbart werden.)

Nicht zum Untersuchungsgegenstand gehörte die Eignung für jüngere Frakturschriften, es soll aber wenigstens erwähnt werden, dass für zunehmende Nähe zum 19. Jahrhundert die in der FineReader-Engine bereits enthaltenen Fraktur-Muster immer besser passen und ein eigenes Training hierfür meist unnötig wäre.

Bei beiden Produkten bleibt es Aufgabe des Anwenders, die jeweiligen XML-Exporte in ein Zielformat seiner Volltextpräsentation oder in einen Suchindex zur Volltextsuche zu überführen.

#### **Für eine schrittweise, modulare OCR-Verarbeitung**

Wenn Eingriffsmöglichkeiten bzw. Einspeisung und Auswertung vor und nach jedem Zwischenschritt erwünscht sind, sind beide Produkte derzeit nicht darauf vorbereitet. Theoretisch mögliche Kombinations- und Anschlussmöglichkeiten bestehen daher nur am Anfang und Ende der jeweiligen Verarbeitungskette.

Eine wirklich offene und an jedem Verarbeitungsschritt ihre Teilergebnisse liefernde Verarbeitungskette scheint allerdings auch in anderen Produkten nicht Standard zu sein; möglicherweise sollte das 2012 gegründete IMPACT-Kompetenzzentrum hier in Anspruch genommen werden, um einschlägige Anforderungen zu bündeln.

## 6.5.2 Tipps für generelle Planungen von OCR Alter Drucke

### *Modularität des Workflows*

Jedes OCR-Projekt ist bestimmten Eigenheiten und Rahmenbedingungen unterworfen, die wesentlich in zu erschließenden Material selbst begründet sind. Gleichwohl soll versucht werden, einige allgemeine Hinweise zu formulieren. Wie bereits eingangs erwähnt, lässt sich der Workflow auch dank internationaler Projekte und ihrer Ergebnisse zunehmend modularisieren und standardisieren. Wesentliches Ziel ist dabei, sich zukünftig für das jeweilige Projekt einen geeigneten Workflow aus verschiedenen Komponenten zusammenzustellen und somit die für die Aufgabenstellung besten Resultate zu erhalten. Im hier behandelten Projekt bestand diese Möglichkeit nicht bzw. konnte nur sehr begrenzt durch Parametrierung der jeweiligen Software Einfluss genommen werden. Auf die dadurch entstandenen Nachteile wurde wiederholt hingewiesen. Die Softwaretests zeigen, dass sich beide Softwareprodukte durch Stärken und Schwächen in verschiedenen Bereichen auszeichnen. Bestünde die Möglichkeit, die einzelnen starken Module für eine spezielle Anwendung zu kombinieren und auch weitere freie Module in den eigenen Workflow zu integrieren, ergäben sich neue Chancen zur Ergebnisoptimierung. Die Einbeziehung weiterer Software war jedoch nicht Gegenstand dieses Pilotprojekts. In der Fachdiskussion spielen das Workflowmanagement und seine technische Unterstützung eine wachsende Rolle.<sup>96</sup>

### *Materialkenntnis*

Wichtige Kriterien für die Planung von OCR leiten sich zuerst aus der Sichtung und Kenntnis der Drucke selbst ab. Sowohl ihr Erhaltungszustand, typographische Besonderheiten als auch inhaltliche Bezüge bilden Themen, mit denen eine Beschäftigung im Vorfeld lohnt. Der Erhaltungszustand des Drucks sowie die Qualität seiner Digitalisierung erfordern geeignete Aktivitäten für eine optimale Binarisierung. Wie auch im Projekt gezeigt werden konnte, beeinflussen verschiedene Binarisierungsmethoden direkt die OCR-Erkennung. Das Wissen über Inhalt und Darstellung des zu erkennenden Wortmaterials ist wesentlich für die Einbeziehung möglichst passender Musterdateien und Wörterbücher. Neben der regionalen und zeitlichen Einordnung der Drucke ist daher die Kenntnis der überwiegend verwendeten Sprache (auch der Fachsprache) sowie der vorkommenden Schriftenvielfalt von Vorteil. Besteht die Möglichkeit, in einem Projekt geringfügige Anpassungen der Konfigurationsparameter vorzunehmen, können bereits wenige Vorsortierungen das Ergebnis verbessern. Eine grobe Einteilung nach Druckorten (-gebieten) bzw. Erscheinungszeiträumen ist hierbei anzuraten; gegebenenfalls lohnt auch eine Systematisierung nach äußeren Merkmalen oder Gattungen.

### *Einbeziehung der OCR-Ergebnisse*

Eine zusätzliche Optimierungsvariante ist es, die Erkennungsgüte der erstellten Dateien zu analysieren. Hierfür bieten sich verschiedene Methoden an – wichtig ist letztlich die Vergleichbarkeit der Ergebnisse innerhalb eines Projekts. Heranzuziehen sind durchaus auch Werte, die durch die OCR selbst erstellt und den Wörtern oder Zeichen als Information (z. B. Konfidenzwerte) mitgegeben werden. All diese Informationen können Indikatoren für eine nachträgliche Selektion sein, um Material zurückzustellen oder einer gesonderten Behandlung zu unterziehen. Beispielsweise ist folgendes Szenario vorstellbar: die Erkennungsgüte einer Seite ist extrem schlecht. Nach Blick auf das Image wird festgestellt, dass die Seite im Buch kopfüber eingebunden war, wodurch die Zeichenerkennung fast unmöglich wurde. Diesen Fehler könnte man folglich unkompliziert beheben. Ein anderes Beispiel ist die gezielte Analyse der Wörterbuchnutzung: Eine signifikant geringe Wörterbuchnutzung (so die Software derartige Auswertungen zulässt) kann anzeigen, dass sehr wahrscheinlich ein Werk in einer fremden Sprache vorliegt. „Fremd“ bedeutet in diesem Fall: die Sprache gehört nicht zu denen, die im zeitgleich bearbeiteten Material überwiegend vertreten sind und mit deren Wörterbuch aktuell gearbeitet wird. Weitere Einsatzfelder von Evaluationsergebnissen verschiedenster Qualität sind aufwandsabhängig denkbar.

<sup>96</sup> Ein interessanter Vortrag zum Thema Workflowdesign: SCHLARB 2011.

### ***Sekundäre Materialien***

Hinsichtlich der Sortierung wird offenbar, wie wichtig es bei Alten Drucken schon in der vorbereitenden Phase sein kann, die vorhandenen bibliographischen Beschreibungen zu analysieren. Existieren in den oftmals konvertierten Titelaufnahmen gar Informationen über die Schriftart („In Fraktur gedruckt“ oder dergleichen), gibt es eine Sprachbezeichnung? Kann diese Sprachbezeichnung eventuell auf einfache Weise generiert oder übernommen werden? Lässt sich mit diesem Titelmaterialein möglicherweise bereits ein Wortschatz für ein Wörterbuch finden? Gibt es vergleichbare Texte, die als Referenzmaterialien verwendet werden können?

### ***Arbeitsteilung***

Eine andere wesentliche Frage ist die der Nutzung von Dienstleistungen als Alternative zum Inhouse-Einsatz von OCR-Software. Die Beantwortung hängt von vielen Faktoren ab. Nicht zuletzt wird jedoch die technische und personelle Ausstattung in der jeweiligen Einrichtung maßgeblich sein. Nicht zu unterschätzen (und das belegen die Erfahrungen im Projekt) sind Rechenzeiten und Speicherbedarf.<sup>97</sup>

### ***Bibliotheksspezifische Schlussfolgerungen***

Es eröffnen sich vielfältige Möglichkeiten zur Gestaltung eigener Arbeitsabläufe. Mit jedem weiteren Workflowschritt entstehen neue Fragen und Ideen. Mithin ergeben sich für die Funeralschriften/Alten Drucke der Staatsbibliothek zu Berlin weitere Testszenarien, von denen nur einige angerissen werden sollen:

Ist es realisierbar, den OCR-Service der VZG<sup>98</sup> unter Nutzung des Abbyy Recognition Servers über das Ticket-System mit den via HK-OCR erstellten Muster- und Wörterbuchdateien zu „füttern“? Bei erfolgreichem Einsatz könnten auf diese Weise weitere Funeralschriften einer Massen-OCR-Erkennung unterzogen werden. Gleichzeitig wäre zu testen, ob diese Muster und Parameter auch für inhaltlich vielfältigeres Material (jedoch möglichst zeitlich vergleichbares) zu besseren Ergebnissen führen als die Standardeinstellungen.

Grundsätzlich ist für die Staatsbibliothek weiterhin zu untersuchen, ob es für einzelne Workflowschritte andere Softwarelösungen gibt, mit denen bessere Ergebnisse zu erzielen sind. In diesem Zusammenhang könnte auch die Opensource-Software Tesseract<sup>99</sup> Gegenstand vergleichender Betrachtungen sein. Weitere Ansatzpunkte bilden die vielfältigen Bemühungen im von der EU geförderten IMPACT-Projekt, Phasen der OCR-Vorbereitung und OCR-Nachbearbeitung zur Optimierung zu nutzen.<sup>100</sup> Beispielhaft wurden inzwischen Evaluationswerkzeuge (auch einzelschrittbezogen) entwickelt, deren Anwendung vergleichbare Ergebnisse verschiedener Herkunft ermöglichen soll. Der intensive Kontakt und die Zusammenarbeit mit verwandten Projekten ist für den Erfahrungsaustausch von unschätzbarem Wert (Vgl. eine Auflistung von Projekten im Anhang 7.5).

<sup>97</sup> Die Prozessierung eines Images (Quartformat) mit einem Standardbürocomputer dauerte durchschnittlich bis zu 10 Minuten; eine Alto-Datei ist ca. 70 KB, eine FR-XML ca. 500 KB, eine FR-Image-Datei ca. 5-7.000 KB groß.

<sup>98</sup> [http://www.gbv.de/wikis/cls/OCR-Service\\_der\\_VZG](http://www.gbv.de/wikis/cls/OCR-Service_der_VZG) [Stand: 02/2013]

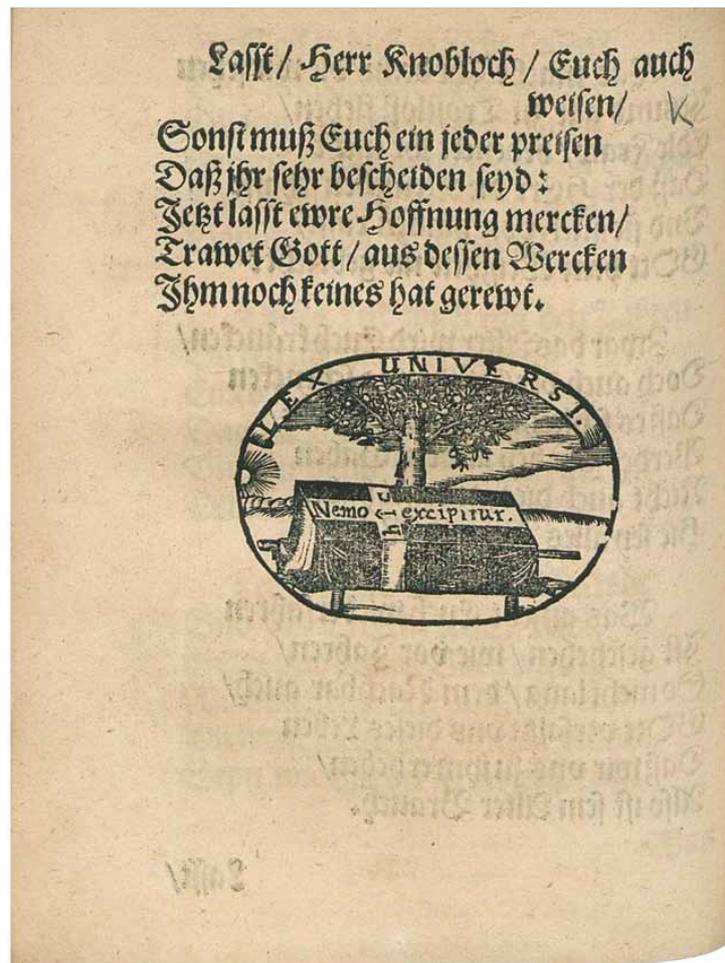
<sup>99</sup> <http://code.google.com/p/tesseract-ocr/> [Stand: 02/2013]

<sup>100</sup> Zu nennen ist hier die Entwicklung des IMPACT Interoperability Framework.

## 7 Anhänge

### 7.1 Beispiele Exportformat

Zur Demonstration der XML-Ausgabeformate wurde eine sehr kurze, nur aus sechs Zeilen bestehende Seite gewählt:



66 Faksimile – Beispielseite

## (a) BIT-Alpha: ALTO

```

1 <?xml version="1.0" encoding="utf-8"?>
2 <alto xmlns="http://schemas.ccs-gmbh.com/ALTO" xmlns:a="http://schemas.ccs-gmbh.com/ALTO" xmlns:xsd="http://www.w3.org/2001/XMLSchema"
3   xmlns:xlink="http://www.w3.org/1999/xlink">
4   <a:Description>
5     <a:MeasurementUnit>mm10</a:MeasurementUnit>
6     <a:sourceImageInformation>
7       <a:fileName>F:\Batch_SBB\dachklag_635359391_orig\00000003.tif</a:fileName>
8     </a:sourceImageInformation>
9     <a:OCRProcessing ID="ID_BIT_1609ec32-7185-4b66-ac85-b6949a62bf81">
10      <a:preProcessingStep>
11        <a:processingDateTime>2011-06-29T19:42:47</a:processingDateTime>
12        <a:processingAgency>Staatsbibliothek zu Berlin - PK</a:processingAgency>
13        <a:processingStepDescription>Color Enhancement</a:processingStepDescription>
14        <a:processingStepSettings>ContrastR=0.5 ContrastG=0.5 ContrastB=1.0 GammaG=1.0 GammaB=1.0 LuminanceR=0.5
15          LuminanceG=0.5 LuminanceB=0.5</a:processingStepSettings>
16        <a:processingSoftware>
17          <a:softwareCreator>B.I.T. Bureau Ingénieur Tomasi</a:softwareCreator>
18          <a:softwareName>BIT-Alpha</a:softwareName>
19          <a:softwareVersion>2.0.38.595 (Rel. 38)</a:softwareVersion>
20        </a:processingSoftware>
21      </a:preProcessingStep>
22      <a:preProcessingStep>
23        <a:processingDateTime>2011-06-29T19:42:47</a:processingDateTime>
24        <a:processingAgency>Staatsbibliothek zu Berlin - PK</a:processingAgency>
25        <a:processingStepDescription>Rotation</a:processingStepDescription>
26        <a:processingStepSettings>Type=None Margins=(Left=0.5 Right=0.5 Top=0.800000011920929 Bottom=0.5)</a:processingStepSettings>
27        <a:processingSoftware>
28          <a:softwareCreator>B.I.T. Bureau Ingénieur Tomasi</a:softwareCreator>
29          <a:softwareName>BIT-Alpha</a:softwareName>
30          <a:softwareVersion>2.0.38.595 (Rel. 38)</a:softwareVersion>
31        </a:processingSoftware>
32      </a:preProcessingStep>
33      <a:preProcessingStep>
34        <a:processingDateTime>2011-06-29T19:42:47</a:processingDateTime>
35        <a:processingAgency>Staatsbibliothek zu Berlin - PK</a:processingAgency>
36        <a:processingStepDescription>Binarisation</a:processingStepDescription>
37        <a:processingStepSettings>SourceBPP=24 Algorithm=Intensity based algorithm</a:processingStepSettings>
38        <a:processingSoftware>
39          <a:softwareCreator>B.I.T. Bureau Ingénieur Tomasi</a:softwareCreator>
40          <a:softwareName>BIT-Alpha</a:softwareName>
41          <a:softwareVersion>2.0.38.595 (Rel. 38)</a:softwareVersion>
42        </a:processingSoftware>
43      </a:preProcessingStep>
44      <a:preProcessingStep>
45        <a:processingDateTime>2011-06-29T19:42:50</a:processingDateTime>
46        <a:processingAgency>Staatsbibliothek zu Berlin - PK</a:processingAgency>
47        <a:processingStepDescription>Cleaning</a:processingStepDescription>
48        <a:processingStepSettings>ContrastR=0.5 ContrastG=0.5 ContrastB=0.5 GammaR=1.0 GammaG=1.0 GammaB=1.0 LuminanceR=0.5
49          LuminanceG=0.5 LuminanceB=0.5</a:processingStepSettings>
50        <a:processingSoftware>
51          <a:softwareCreator>B.I.T. Bureau Ingénieur Tomasi</a:softwareCreator>
52          <a:softwareName>BIT-Alpha</a:softwareName>
53          <a:softwareVersion>2.0.38.595 (Rel. 38)</a:softwareVersion>
54        </a:processingSoftware>
55      </a:preProcessingStep>
56      <a:preProcessingStep>
57        <a:processingDateTime>2011-06-29T19:42:52</a:processingDateTime>
58        <a:processingAgency>Staatsbibliothek zu Berlin - PK</a:processingAgency>
59        <a:processingStepDescription>Remove Dots</a:processingStepDescription>
60        <a:processingStepSettings>ContrastR=0.5 ContrastG=0.5 ContrastB=0.5 GammaR=1.0 GammaG=1.0 GammaB=1.0 LuminanceR=0.5 LuminanceG=0
61          LuminanceB=0.5</a:processingStepSettings>
62        <a:processingSoftware>
63          <a:softwareCreator>B.I.T. Bureau Ingénieur Tomasi</a:softwareCreator>
64          <a:softwareName>BIT-Alpha</a:softwareName>
65          <a:softwareVersion>2.0.38.595 (Rel. 38)</a:softwareVersion>
66        </a:processingSoftware>
67      </a:preProcessingStep>
68      <a:preProcessingStep>
69        <a:processingDateTime>2011-06-29T19:42:57</a:processingDateTime>
70        <a:processingAgency>Staatsbibliothek zu Berlin - PK</a:processingAgency>
71        <a:processingStepDescription>Blackborder elimination</a:processingStepDescription>
72        <a:processingStepSettings>ContrastR=0.5 ContrastG=0.5 ContrastB=0.5 GammaR=1.0 GammaG=1.0 GammaB=1.0 LuminanceR=0.5 LuminanceG=0
73          LuminanceB=0.5</a:processingStepSettings>
74        <a:processingSoftware>
75          <a:softwareCreator>B.I.T. Bureau Ingénieur Tomasi</a:softwareCreator>
76          <a:softwareName>BIT-Alpha</a:softwareName>
77          <a:softwareVersion>2.0.38.595 (Rel. 38)</a:softwareVersion>
78        </a:processingSoftware>
79      </a:preProcessingStep>
80      <a:preProcessingStep>
81        <a:processingDateTime>2011-06-29T19:43:03</a:processingDateTime>
82        <a:processingAgency>Staatsbibliothek zu Berlin - PK</a:processingAgency>
83        <a:processingStepDescription>Detection of horizontal lines</a:processingStepDescription>
84        <a:processingStepSettings>ForceBitmap=true MaxThickness=0.5 MinThickness=0.0 MaxWhiteLength=0.1 MinBlackLength=0.9
85          MinTotalLength=2.0 Margins=(Left=0.5 Right=0.5 Top=0.800000011920929 Bottom=0.5)</a:processingStepSettings>
86        <a:processingSoftware>
87          <a:softwareCreator>B.I.T. Bureau Ingénieur Tomasi</a:softwareCreator>
88          <a:softwareName>BIT-Alpha</a:softwareName>
89          <a:softwareVersion>2.0.38.595 (Rel. 38)</a:softwareVersion>
90        </a:processingSoftware>
91      </a:preProcessingStep>
92      <a:preProcessingStep>
93        <a:processingDateTime>2011-06-29T19:43:04</a:processingDateTime>
94        <a:processingAgency>Staatsbibliothek zu Berlin - PK</a:processingAgency>
95        <a:processingStepDescription>Detection of vertical lines</a:processingStepDescription>
96        <a:processingStepSettings>ForceBitmap=true MaxThickness=0.1 MinThickness=0.0 MaxWhiteLength=0.0 MinBlackLength=0.9
97          MinTotalLength=5.0 Margins=(Left=0.5 Right=0.5 Top=0.800000011920929 Bottom=0.5)</a:processingStepSettings>
98        <a:processingSoftware>
99          <a:softwareCreator>B.I.T. Bureau Ingénieur Tomasi</a:softwareCreator>
100         <a:softwareName>BIT-Alpha</a:softwareName>
101         <a:softwareVersion>2.0.38.595 (Rel. 38)</a:softwareVersion>
102       </a:processingSoftware>
103     </a:preProcessingStep>
104   </a:OCRProcessing>
105 </a:Description>
106 </alto>

```

```

92 <a:softwareCreator>B.I.T. Bureau Ingénieur Tomasi</a:softwareCreator>
93 <a:softwareName>BIT-Alpha</a:softwareName>
94 <a:softwareVersion>2.0.38.595 (Rel. 38)</a:softwareVersion>
95 </a:processingSoftware>
96 </a:preProcessingStep>
97 <a:preProcessingStep>
98 <a:processingDateTime>2011-06-29T19:43:05</a:processingDateTime>
99 <a:processingAgency>Staatsbibliothek zu Berlin - PK</a:processingAgency>
100 <a:processingStepDescription>Segmentation</a:processingStepDescription>
101 <a:processingStepSettings>MinHeight=0.1 MinWidth=0.1 MinVertDist=2.0 MinHorDist=2.0 Margins=(Left=0.5 Right=0.5
Top=0.80000011920929 Bottom=0.5)</a:processingStepSettings>
102 <a:processingSoftware>
103 <a:softwareCreator>B.I.T. Bureau Ingénieur Tomasi</a:softwareCreator>
104 <a:softwareName>BIT-Alpha</a:softwareName>
105 <a:softwareVersion>2.0.38.595 (Rel. 38)</a:softwareVersion>
106 </a:processingSoftware>
107 </a:preProcessingStep>
108 <a:preProcessingStep>
109 <a:processingDateTime>2011-06-29T19:43:06</a:processingDateTime>
110 <a:processingAgency>Staatsbibliothek zu Berlin - PK</a:processingAgency>
111 <a:processingStepDescription>Region identification</a:processingStepDescription>
112 <a:processingStepSettings>Default=Binary ColorRegions=(Detect=false) PaletteRegions=(Detect=false) BinaryRegions=(Detect=false)
</a:processingStepSettings>
113 <a:processingSoftware>
114 <a:softwareCreator>B.I.T. Bureau Ingénieur Tomasi</a:softwareCreator>
115 <a:softwareName>BIT-Alpha</a:softwareName>
116 <a:softwareVersion>2.0.38.595 (Rel. 38)</a:softwareVersion>
117 </a:processingSoftware>
118 </a:preProcessingStep>
119 <a:ocrProcessingStep>
120 <a:processingDateTime>2011-06-29T19:43:16</a:processingDateTime>
121 <a:processingAgency>Staatsbibliothek zu Berlin - PK</a:processingAgency>
122 <a:processingStepDescription>Optical Character Recognition</a:processingStepDescription>
123 <a:processingStepSettings>not implemented.</a:processingStepSettings>
124 <a:processingSoftware>
125 <a:softwareCreator>B.I.T. Bureau Ingénieur Tomasi</a:softwareCreator>
126 <a:softwareName>BIT-Alpha</a:softwareName>
127 <a:softwareVersion>2.0.38.595 (Rel. 38)</a:softwareVersion>
128 </a:processingSoftware>
129 </a:ocrProcessingStep>
130 </a:OCRProcessing>
131 </a:Description>
132 <a:Layout>
133 <a:Page ID="ID_BIT_41462251-f383-4bc1-a3f6-f2301a5327d1" HEIGHT="1939" WIDTH="1424" PHYSICAL_IMG_NR="0" QUALITY="DK" PROCESSING=
"ID_BIT_1609ec32-7185-4b66-ac85-b6949a62bf81">
134 <a:PrintSpace ID="ID_BIT_70278dba-ee8-44ea-b094-d85119a7733b" HEIGHT="1808.87" WIDTH="1323.25" HPOS="50.0" VPOS="80.0">
135 <a:TextBlock ID="ID_BIT_59d0ach9-lbbc-4ef6-9c2d-471d558705a5" HEIGHT="505" WIDTH="991" HPOS="350" VPOS="149">
136 <a:TextLine ID="ID_BIT_c6c8b85f-920d-4b89-blad-30b4e24e6770" HEIGHT="83.82" WIDTH="892.39" HPOS="99.06" VPOS="6.77">
137 <a:String HEIGHT="79.59" WIDTH="117.69" VPOS="6.77" HPOS="99.06" CONTENT="Lassee"/>
138 <a:SP WIDTH="123.61" VPOS="6.77" HPOS="99.06"/>
139 <a:String HEIGHT="79.59" WIDTH="30.48" VPOS="6.77" HPOS="222.67" CONTENT="/">
140 <a:SP WIDTH="38.1" VPOS="6.77" HPOS="222.67"/>
141 <a:String HEIGHT="79.59" WIDTH="137.16" VPOS="6.77" HPOS="260.77" CONTENT="Her"/>
142 <a:SP WIDTH="157.48" VPOS="6.77" HPOS="260.77"/>
143 <a:String HEIGHT="72.81" WIDTH="245.53" VPOS="4.23" HPOS="418.25" CONTENT="Knobloch"/>
144 <a:SP WIDTH="261.62" VPOS="5.5" HPOS="418.25"/>
145 <a:String HEIGHT="70.27" WIDTH="28.79" VPOS="6.77" HPOS="679.87" CONTENT="/">
146 <a:SP WIDTH="46.57" VPOS="6.77" HPOS="679.87"/>
147 <a:String HEIGHT="70.27" WIDTH="130.39" VPOS="6.77" HPOS="726.44" CONTENT="Euch"/>
148 <a:SP WIDTH="157.48" VPOS="5.93" HPOS="726.44"/>
149 <a:String HEIGHT="71.12" WIDTH="107.53" VPOS="5.08" HPOS="883.92" CONTENT="auch"/>
150 </a:TextLine>
151 <a:TextLine ID="ID_BIT_7e62def2-46be-48f6-98ee-d1146calefc2" HEIGHT="2.54" WIDTH="3.39" HPOS="159.17" VPOS="73.66">
152 <a:String HEIGHT="2.54" WIDTH="3.39" VPOS="73.66" HPOS="159.17" CONTENT="."/>
153 </a:TextLine>
154 <a:TextLine ID="ID_BIT_49fadae8-6170-4e08-a73b-95f294a1fb91" HEIGHT="77.05" WIDTH="303.95" HPOS="687.49" VPOS="80.43">
155 <a:String HEIGHT="71.97" WIDTH="158.33" VPOS="80.43" HPOS="687.49" CONTENT="weisen"/>
156 <a:SP WIDTH="159.17" VPOS="80.43" HPOS="687.49"/>
157 <a:String HEIGHT="71.97" WIDTH="24.55" VPOS="80.43" HPOS="846.67" CONTENT="/">
158 <a:SP WIDTH="103.29" VPOS="86.78" HPOS="846.67"/>
159 <a:String HEIGHT="64.35" WIDTH="41.49" VPOS="93.13" HPOS="949.96" CONTENT="."/>
160 </a:TextLine>
161 <a:TextLine ID="ID_BIT_126ea2a4-ba77-44ec-b9b2-ff9421760f70" HEIGHT="80.43" WIDTH="820.42" HPOS="3.39" VPOS="137.16">
162 <a:String HEIGHT="73.66" WIDTH="160.87" VPOS="137.16" HPOS="3.39" CONTENT="Sonst"/>
163 <a:SP WIDTH="193.04" VPOS="137.16" HPOS="3.39"/>
164 <a:String HEIGHT="73.66" WIDTH="96.52" VPOS="137.16" HPOS="196.43" CONTENT="müs"/>
165 <a:SP WIDTH="108.37" VPOS="141.82" HPOS="196.43"/>
166 <a:String HEIGHT="67.73" WIDTH="128.69" VPOS="146.47" HPOS="304.8" CONTENT="Euch"/>
167 <a:SP WIDTH="139.7" VPOS="148.17" HPOS="304.8"/>
168 <a:String HEIGHT="65.19" WIDTH="69.43" VPOS="149.86" HPOS="444.5" CONTENT="ein"/>
169 <a:SP WIDTH="82.97" VPOS="150.28" HPOS="444.5"/>
170 <a:String HEIGHT="66.04" WIDTH="115.99" VPOS="150.71" HPOS="527.47" CONTENT="jeder"/>
171 <a:SP WIDTH="128.69" VPOS="151.55" HPOS="527.47"/>
172 <a:String HEIGHT="65.19" WIDTH="167.64" VPOS="152.4" HPOS="656.17" CONTENT="peissen"/>
173 </a:TextLine>
174 <a:TextLine ID="ID_BIT_7fad1c16-25cb-4ae5-b4d9-fed34033a0ee" HEIGHT="71.97" WIDTH="724.75" HPOS="0.0" VPOS="214.21">
175 <a:String HEIGHT="69.43" WIDTH="86.36" VPOS="214.21" HPOS="0.0" CONTENT="DaS"/>
176 <a:SP WIDTH="115.99" VPOS="214.21" HPOS="0.0"/>
177 <a:String HEIGHT="69.43" WIDTH="86.36" VPOS="214.21" HPOS="115.99" CONTENT="jhr"/>
178 <a:SP WIDTH="101.6" VPOS="216.75" HPOS="115.99"/>
179 <a:String HEIGHT="66.04" WIDTH="92.29" VPOS="219.29" HPOS="217.59" CONTENT="sehr"/>
180 <a:SP WIDTH="105.83" VPOS="219.71" HPOS="217.59"/>
181 <a:String HEIGHT="66.04" WIDTH="248.92" VPOS="220.13" HPOS="323.43" CONTENT="bescheiden"/>
182 <a:SP WIDTH="268.39" VPOS="220.13" HPOS="323.43"/>
183 <a:String HEIGHT="66.04" WIDTH="100.75" VPOS="220.13" HPOS="591.82" CONTENT="seyd"/>
184 <a:SP WIDTH="114.3" VPOS="219.71" HPOS="591.82"/>
185 <a:String HEIGHT="65.19" WIDTH="18.63" VPOS="219.29" HPOS="706.12" CONTENT="."/>
186 </a:TextLine>

```

```

187 <a:TextLine ID="ID_BIT_4abe08d2-e8f9-433a-bbba-9a0880e2e487" HEIGHT="77.89" WIDTH="874.61" HPOS="0.0" VPOS="284.48">
188 <a:String HEIGHT="66.89" WIDTH="124.46" VPOS="284.48" HPOS="0.0" CONTENT="Jetztz"/>
189 <a:SP WIDTH="135.47" VPOS="285.75" HPOS="0.0"/>
190 <a:String HEIGHT="66.04" WIDTH="99.06" VPOS="287.02" HPOS="135.47" CONTENT="lasse"/>
191 <a:SP WIDTH="112.61" VPOS="287.87" HPOS="135.47"/>
192 <a:String HEIGHT="71.12" WIDTH="117.69" VPOS="288.71" HPOS="248.07" CONTENT="evre"/>
193 <a:SP WIDTH="147.32" VPOS="288.71" HPOS="248.07"/>
194 <a:String HEIGHT="71.12" WIDTH="236.22" VPOS="288.71" HPOS="395.39" CONTENT="Hoffnung"/>
195 <a:SP WIDTH="251.46" VPOS="292.1" HPOS="395.39"/>
196 <a:String HEIGHT="66.89" WIDTH="199.81" VPOS="295.49" HPOS="646.85" CONTENT="mercken"/>
197 <a:SP WIDTH="201.51" VPOS="295.49" HPOS="646.85"/>
198 <a:String HEIGHT="66.89" WIDTH="26.25" VPOS="295.49" HPOS="848.36" CONTENT=""/>
199 </a:TextLine>
200 <a:TextLine ID="ID_BIT_fb094242-77e3-4695-aa8e-917cbfb9440f" HEIGHT="79.59" WIDTH="868.68" HPOS="0.0" VPOS="354.75">
201 <a:String HEIGHT="69.43" WIDTH="198.97" VPOS="354.75" HPOS="0.0" CONTENT="Travet"/>
202 <a:SP WIDTH="212.51" VPOS="357.29" HPOS="0.0"/>
203 <a:String HEIGHT="67.73" WIDTH="132.08" VPOS="359.83" HPOS="212.51" CONTENT="Gott"/>
204 <a:SP WIDTH="140.55" VPOS="359.83" HPOS="212.51"/>
205 <a:String HEIGHT="67.73" WIDTH="23.71" VPOS="359.83" HPOS="353.06" CONTENT=""/>
206 <a:SP WIDTH="33.02" VPOS="361.53" HPOS="353.06"/>
207 <a:String HEIGHT="66.04" WIDTH="99.06" VPOS="363.22" HPOS="386.08" CONTENT="aus"/>
208 <a:SP WIDTH="110.07" VPOS="364.07" HPOS="386.08"/>
209 <a:String HEIGHT="66.89" WIDTH="140.55" VPOS="364.91" HPOS="496.15" CONTENT="dessen"/>
210 <a:SP WIDTH="151.55" VPOS="366.18" HPOS="496.15"/>
211 <a:String HEIGHT="66.89" WIDTH="220.98" VPOS="367.45" HPOS="647.7" CONTENT="Wercken"/>
212 </a:TextLine>
213 <a:TextLine ID="ID_BIT_72020365-be74-4556-8710-51632c34990d" HEIGHT="77.05" WIDTH="717.13" HPOS="0.0" VPOS="428.41">
214 <a:String HEIGHT="71.12" WIDTH="132.08" VPOS="428.41" HPOS="0.0" CONTENT="Ihm"/>
215 <a:SP WIDTH="141.39" VPOS="429.68" HPOS="0.0"/>
216 <a:String HEIGHT="71.12" WIDTH="117.69" VPOS="430.95" HPOS="141.39" CONTENT="noch"/>
217 <a:SP WIDTH="127.85" VPOS="432.22" HPOS="141.39"/>
218 <a:String HEIGHT="71.12" WIDTH="151.55" VPOS="433.49" HPOS="269.24" CONTENT="keines"/>
219 <a:SP WIDTH="162.56" VPOS="434.34" HPOS="269.24"/>
220 <a:String HEIGHT="69.43" WIDTH="83.82" VPOS="435.19" HPOS="431.8" CONTENT="hat"/>
221 <a:SP WIDTH="93.13" VPOS="435.19" HPOS="431.8"/>
222 <a:String HEIGHT="70.27" WIDTH="171.87" VPOS="435.19" HPOS="524.93" CONTENT="gerewt"/>
223 <a:SP WIDTH="175.26" VPOS="435.19" HPOS="524.93"/>
224 <a:String HEIGHT="70.27" WIDTH="16.93" VPOS="435.19" HPOS="700.19" CONTENT="."/>
225 </a:TextLine>
226 </a:TextBlock>
227 <a:TextBlock ID="ID_BIT_31e5d647-63a5-4880-9b45-ead04f5b660e" HEIGHT="502" WIDTH="717" HPOS="488" VPOS="723"/>
228 <a:TextBlock ID="ID_BIT_37402e39-2b3d-4b14-b757-c80891a5f62b" HEIGHT="292" WIDTH="168" HPOS="498" VPOS="1435">
229 <a:TextLine ID="ID_BIT_2f4716bb-6e35-4e20-9c1d-262e1e93926d" HEIGHT="25.4" WIDTH="13.55" HPOS="0.0" VPOS="233.68">
230 <a:String HEIGHT="25.4" WIDTH="13.55" VPOS="233.68" HPOS="0.0" CONTENT="."/>
231 </a:TextLine>
232 <a:TextLine ID="ID_BIT_ca23d635-5c2f-4ced-9e6b-0aa38d6c4d71" HEIGHT="2.54" WIDTH="3.39" HPOS="8.47" VPOS="277.71">
233 <a:String HEIGHT="2.54" WIDTH="3.39" VPOS="277.71" HPOS="8.47" CONTENT="."/>
234 </a:TextLine>
235 <a:TextLine ID="ID_BIT_5836981f-c2f8-43ec-bbdd-589b6433557d" HEIGHT="7.62" WIDTH="9.31" HPOS="5.93" VPOS="284.48">
236 <a:String HEIGHT="66.89" WIDTH="199.81" VPOS="295.49" HPOS="646.85" CONTENT="mercken"/>
237 <a:SP WIDTH="201.51" VPOS="295.49" HPOS="646.85"/>
238 <a:String HEIGHT="66.89" WIDTH="26.25" VPOS="295.49" HPOS="848.36" CONTENT=""/>
239 </a:TextLine>
240 <a:TextLine ID="ID_BIT_fb094242-77e3-4695-aa8e-917cbfb9440f" HEIGHT="79.59" WIDTH="868.68" HPOS="0.0" VPOS="354.75">
241 <a:String HEIGHT="69.43" WIDTH="198.97" VPOS="354.75" HPOS="0.0" CONTENT="Travet"/>
242 <a:SP WIDTH="212.51" VPOS="357.29" HPOS="0.0"/>
243 <a:String HEIGHT="67.73" WIDTH="132.08" VPOS="359.83" HPOS="212.51" CONTENT="Gott"/>
244 <a:SP WIDTH="140.55" VPOS="359.83" HPOS="212.51"/>
245 <a:String HEIGHT="67.73" WIDTH="23.71" VPOS="359.83" HPOS="353.06" CONTENT=""/>
246 <a:SP WIDTH="33.02" VPOS="361.53" HPOS="353.06"/>
247 <a:String HEIGHT="66.04" WIDTH="99.06" VPOS="363.22" HPOS="386.08" CONTENT="aus"/>
248 <a:SP WIDTH="110.07" VPOS="364.07" HPOS="386.08"/>
249 <a:String HEIGHT="66.89" WIDTH="140.55" VPOS="364.91" HPOS="496.15" CONTENT="dessen"/>
250 <a:SP WIDTH="151.55" VPOS="366.18" HPOS="496.15"/>
251 <a:String HEIGHT="66.89" WIDTH="220.98" VPOS="367.45" HPOS="647.7" CONTENT="Wercken"/>
252 </a:TextLine>
253 <a:TextLine ID="ID_BIT_72020365-be74-4556-8710-51632c34990d" HEIGHT="77.05" WIDTH="717.13" HPOS="0.0" VPOS="428.41">
254 <a:String HEIGHT="71.12" WIDTH="132.08" VPOS="428.41" HPOS="0.0" CONTENT="Ihm"/>
255 <a:SP WIDTH="141.39" VPOS="429.68" HPOS="0.0"/>
256 <a:String HEIGHT="71.12" WIDTH="117.69" VPOS="430.95" HPOS="141.39" CONTENT="noch"/>
257 <a:SP WIDTH="127.85" VPOS="432.22" HPOS="141.39"/>
258 <a:String HEIGHT="71.12" WIDTH="151.55" VPOS="433.49" HPOS="269.24" CONTENT="keines"/>
259 <a:SP WIDTH="162.56" VPOS="434.34" HPOS="269.24"/>
260 <a:String HEIGHT="69.43" WIDTH="83.82" VPOS="435.19" HPOS="431.8" CONTENT="hat"/>
261 <a:SP WIDTH="93.13" VPOS="435.19" HPOS="431.8"/>
262 <a:String HEIGHT="70.27" WIDTH="171.87" VPOS="435.19" HPOS="524.93" CONTENT="gerewt"/>
263 <a:SP WIDTH="175.26" VPOS="435.19" HPOS="524.93"/>
264 <a:String HEIGHT="70.27" WIDTH="16.93" VPOS="435.19" HPOS="700.19" CONTENT="."/>
265 </a:TextLine>
266 </a:TextBlock>
267 <a:TextBlock ID="ID_BIT_31e5d647-63a5-4880-9b45-ead04f5b660e" HEIGHT="502" WIDTH="717" HPOS="488" VPOS="723"/>
268 <a:TextBlock ID="ID_BIT_37402e39-2b3d-4b14-b757-c80891a5f62b" HEIGHT="292" WIDTH="168" HPOS="498" VPOS="1435">
269 <a:TextLine ID="ID_BIT_2f4716bb-6e35-4e20-9c1d-262e1e93926d" HEIGHT="25.4" WIDTH="13.55" HPOS="0.0" VPOS="233.68">
270 <a:String HEIGHT="25.4" WIDTH="13.55" VPOS="233.68" HPOS="0.0" CONTENT="."/>
271 </a:TextLine>
272 <a:TextLine ID="ID_BIT_ca23d635-5c2f-4ced-9e6b-0aa38d6c4d71" HEIGHT="2.54" WIDTH="3.39" HPOS="8.47" VPOS="277.71">
273 <a:String HEIGHT="2.54" WIDTH="3.39" VPOS="277.71" HPOS="8.47" CONTENT="."/>
274 </a:TextLine>
275 <a:TextLine ID="ID_BIT_5836981f-c2f8-43ec-bbdd-589b6433557d" HEIGHT="7.62" WIDTH="9.31" HPOS="5.93" VPOS="284.48">
276 <a:String HEIGHT="7.62" WIDTH="9.31" VPOS="284.48" HPOS="5.93" CONTENT="."/>
277 </a:TextLine>
278 </a:TextBlock>
279 </a:PrintSpace>
280 </a:Page>
281 </a:Layout>
282 </alto>
283

```

(b) HK-OCR:

FineReader-XML

```

1 <?xml version="1.0" encoding="UTF-8"?>
2 <document xmlns="http://www.abbyy.com/FineReader_xml/FineReader8-schema-v2.xml" xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
3   version="1.0" producer="FineReader 8.0" xsi:schemaLocation="http://www.abbyy.com/FineReader_xml/FineReader8-schema-v2.xml
4   http://www.abbyy.com/FineReader_xml/FineReader8-schema-v2.xml" mainLanguage="OldGerman" languages="OldGerman">
5   <page width="1710" height="2311" resolution="300">
6     <block blockType="Text" blockName="" isHidden="true" l="413" t="170" r="1596" b="795">
7       <region>
8         <rect l="413" t="170" r="1596" b="795"/>
9       </region>
10      <text>
11        <par align="Right" rightIndent="1" lineSpacing="85">
12          <line baseline="260" l="533" t="177" r="1589" b="279">
13            <formatting lang="OldGerman" ff="Arial" fs="17" spacing="-20" style="0">
14              <charParams l="533" t="189" r="574" b="263" characterHeight="52" hasUncertainHeight="false" baseLine="0" wordStart="true"
15                wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="1" meanStrokeWidth="97"
16                charConfidence="81" serifProbability="255">L</charParams>
17              <charParams l="577" t="207" r="610" b="259" characterHeight="52" hasUncertainHeight="false" baseLine="0" wordStart="false"
18                wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="1" meanStrokeWidth="97"
19                charConfidence="98" serifProbability="255">x</charParams>
20              <charParams l="614" t="190" r="641" b="275" characterHeight="52" hasUncertainHeight="false" baseLine="0" wordStart="false"
21                wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="1" meanStrokeWidth="97"
22                charConfidence="100" serifProbability="255">e</charParams>
23              <charParams l="630" t="189" r="659" b="274" characterHeight="52" hasUncertainHeight="false" baseLine="0" wordStart="false"
24                wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="1" meanStrokeWidth="97"
25                charConfidence="100" serifProbability="255">e</charParams>
26              <charParams l="650" t="197" r="673" b="259" characterHeight="52" hasUncertainHeight="false" baseLine="0" wordStart="false"
27                wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="1" meanStrokeWidth="97"
28                charConfidence="100" serifProbability="255">t</charParams>
29              <charParams l="681" t="193" r="714" b="270" characterHeight="52" hasUncertainHeight="false" baseLine="0" wordStart="false"
30                wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="1" meanStrokeWidth="97"
31                charConfidence="100" serifProbability="255">K</charParams>
32              <charParams l="714" t="189" r="733" b="275" characterHeight="52" hasUncertainHeight="false" baseLine="0"> </charParams>
33            </formatting>
34            <formatting lang="OldGerman" ff="Arial" fs="19" spacing="-20" style="1">
35              <charParams l="733" t="190" r="799" b="275" characterHeight="52" hasUncertainHeight="false" baseLine="0" wordStart="true"
36                wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="97"
37                charConfidence="100" serifProbability="255">H</charParams>
38              <charParams l="805" t="205" r="828" b="260" characterHeight="52" hasUncertainHeight="false" baseLine="0" wordStart="false"
39                wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="97"
40                charConfidence="100" serifProbability="255">e</charParams>
41              <charParams l="830" t="206" r="856" b="259" characterHeight="52" hasUncertainHeight="false" baseLine="0" wordStart="false"
42                wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="97"
43                charConfidence="100" serifProbability="255">r</charParams>
44              <charParams l="860" t="206" r="887" b="259" characterHeight="52" hasUncertainHeight="false" baseLine="0" wordStart="false"
45                wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="97"
46                charConfidence="100" serifProbability="255">r</charParams>
47              <charParams l="887" t="190" r="911" b="275" characterHeight="52" hasUncertainHeight="false" baseLine="0"> </charParams>
48              <charParams l="911" t="190" r="966" b="263" characterHeight="52" hasUncertainHeight="false" baseLine="0" wordStart="true"
49                wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="97"
50                charConfidence="100" serifProbability="255">K</charParams>
51              <charParams l="970" t="207" r="1010" b="261" characterHeight="52" hasUncertainHeight="false" baseLine="0" wordStart="false"
52                wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="97"
53                charConfidence="98" serifProbability="255">n</charParams>
54              <charParams l="1013" t="207" r="1043" b="260" characterHeight="52" hasUncertainHeight="false" baseLine="0" wordStart="false"
55                wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="97"
56                charConfidence="100" serifProbability="255">o</charParams>
57              <charParams l="1048" t="192" r="1080" b="262" characterHeight="52" hasUncertainHeight="false" baseLine="0" wordStart="false"
58                wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="97"
59                charConfidence="100" serifProbability="255">b</charParams>
60              <charParams l="1082" t="192" r="1104" b="261" characterHeight="52" hasUncertainHeight="false" baseLine="0" wordStart="false"
61                wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="97"
62                charConfidence="100" serifProbability="255">l</charParams>
63              <charParams l="1104" t="208" r="1149" b="261" characterHeight="52" hasUncertainHeight="false" baseLine="0" wordStart="false"
64                wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="97"
65                charConfidence="81" serifProbability="255">o</charParams>
66              <charParams l="1141" t="192" r="1200" b="279" characterHeight="52" hasUncertainHeight="false" baseLine="0" wordStart="false"
67                wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="97"
68                charConfidence="100" serifProbability="255">o</charParams>
69              <charParams l="1141" t="192" r="1200" b="279" characterHeight="52" hasUncertainHeight="false" baseLine="0" wordStart="false"
70                wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="97"
71                charConfidence="100" serifProbability="255">h</charParams>
72              <charParams l="1200" t="190" r="1219" b="271" characterHeight="52" hasUncertainHeight="false" baseLine="0"> </charParams>
73              <charParams l="1219" t="197" r="1253" b="271" suspicious="true" characterHeight="52" hasUncertainHeight="false" baseLine="0"
74                wordStart="true" wordFromDictionary="false" wordNormal="false" wordNumeric="false" wordIdentifier="false" wordPenalty="0"
75                meanStrokeWidth="97" charConfidence="100" serifProbability="255">K</charParams>
76              <charParams l="1253" t="192" r="1275" b="271" characterHeight="52" hasUncertainHeight="false" baseLine="0"> </charParams>
77              <charParams l="1275" t="194" r="1323" b="263" characterHeight="52" hasUncertainHeight="false" baseLine="0" wordStart="true"
78                wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="97"
79                charConfidence="100" serifProbability="255">E</charParams>
80              <charParams l="1326" t="207" r="1364" b="261" characterHeight="52" hasUncertainHeight="false" baseLine="0" wordStart="false"
81                wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="97"
82                charConfidence="100" serifProbability="255">u</charParams>
83              <charParams l="1370" t="192" r="1429" b="277" characterHeight="52" hasUncertainHeight="false" baseLine="0" wordStart="false"
84                wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="97"
85                charConfidence="98" serifProbability="255">o</charParams>
86              <charParams l="1370" t="192" r="1429" b="277" characterHeight="52" hasUncertainHeight="false" baseLine="0" wordStart="false"
87                wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="97"
88                charConfidence="98" serifProbability="255">h</charParams>
89              <charParams l="1429" t="177" r="1460" b="264" characterHeight="52" hasUncertainHeight="false" baseLine="0"> </charParams>
90              <charParams l="1460" t="205" r="1491" b="257" characterHeight="52" hasUncertainHeight="false" baseLine="0" wordStart="true"
91                wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="97"
92                charConfidence="100" serifProbability="255">a</charParams>
93              <charParams l="1496" t="199" r="1532" b="255" characterHeight="52" hasUncertainHeight="false" baseLine="0" wordStart="false"
94                wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="97"
95                charConfidence="100" serifProbability="255">u</charParams>
96              <charParams l="1536" t="177" r="1589" b="264" characterHeight="52" hasUncertainHeight="false" baseLine="0" wordStart="false"
97                wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="97"
98                charConfidence="80" serifProbability="255">o</charParams>

```

```

45 <charParams l="1536" t="177" r="1589" b="264" characterHeight="52" hasUncertainHeight="false" baseLine="0" wordStart="false"
wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="97"
charConfidence="80" serifProbability="255">h</charParams>
46 </formatting>
47 </line>
48 </par>
49 <par startIndent="806" lineSpacing="85">
50 <line baseline="345" l="1229" t="281" r="1590" b="364">
51 <formatting lang="OldGerman" ff="Arial" fs="19." spacing="-20" style="1">
52 <charParams l="1229" t="290" r="1286" b="347" characterHeight="48" hasUncertainHeight="false" baseLine="0" wordStart="true"
wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="1" meanStrokeWidth="55"
charConfidence="100" serifProbability="255">w</charParams>
53 <charParams l="1290" t="296" r="1313" b="346" characterHeight="48" hasUncertainHeight="false" baseLine="0" wordStart="false"
wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="1" meanStrokeWidth="55"
charConfidence="98" serifProbability="255">s</charParams>
54 <charParams l="1318" t="290" r="1333" b="345" characterHeight="48" hasUncertainHeight="false" baseLine="0" wordStart="false"
wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="1" meanStrokeWidth="55"
charConfidence="96" serifProbability="255">i</charParams>
55 <charParams l="1337" t="281" r="1367" b="359" characterHeight="48" hasUncertainHeight="false" baseLine="0" wordStart="false"
wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="1" meanStrokeWidth="55"
charConfidence="100" serifProbability="255">s</charParams>
56 <charParams l="1355" t="298" r="1376" b="346" characterHeight="48" hasUncertainHeight="false" baseLine="0" wordStart="false"
wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="1" meanStrokeWidth="55"
charConfidence="100" serifProbability="255">e</charParams>
57 <charParams l="1381" t="295" r="1417" b="345" characterHeight="48" hasUncertainHeight="false" baseLine="0" wordStart="false"
wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="1" meanStrokeWidth="55"
charConfidence="98" serifProbability="255">n</charParams>
58 <charParams l="1418" t="285" r="1446" b="358" characterHeight="48" hasUncertainHeight="false" baseLine="0" wordStart="false"
wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="1" meanStrokeWidth="55"
charConfidence="100" serifProbability="255">/</charParams>
59 <charParams l="1446" t="281" r="1533" b="359" characterHeight="48" hasUncertainHeight="false" baseLine="0"> </charParams>
60 <charParams l="1533" t="286" r="1590" b="364" suspicious="true" characterHeight="48" hasUncertainHeight="false" baseLine="0"
wordStart="true" wordFromDictionary="false" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="2"
meanStrokeWidth="55" charConfidence="18" serifProbability="255">^</charParams>
61 </formatting>
62 </line>
63 <line baseline="429" l="423" t="361" r="1392" b="446">
64 <formatting lang="OldGerman" ff="Arial" fs="19." spacing="-20" style="1">
65 <charParams l="423" t="361" r="495" b="435" characterHeight="50" hasUncertainHeight="false" baseLine="0" wordStart="true"
wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="10"
charConfidence="100" serifProbability="255">S</charParams>
66 <charParams l="499" t="379" r="531" b="430" characterHeight="50" hasUncertainHeight="false" baseLine="0" wordStart="false"
wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="10"
charConfidence="100" serifProbability="255">O</charParams>
67 <charParams l="534" t="377" r="571" b="429" characterHeight="50" hasUncertainHeight="false" baseLine="0" wordStart="false"
wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="10"
charConfidence="100" serifProbability="255">n</charParams>
68 <charParams l="576" t="363" r="616" b="445" characterHeight="50" hasUncertainHeight="false" baseLine="0" wordStart="false"
wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="10"
charConfidence="100" serifProbability="255">S</charParams>
69 <charParams l="576" t="363" r="616" b="445" characterHeight="50" hasUncertainHeight="false" baseLine="0" wordStart="false"
wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="10"
charConfidence="100" serifProbability="255">C</charParams>
70 <charParams l="625" t="375" r="684" b="430" characterHeight="50" hasUncertainHeight="false" baseLine="0" wordStart="false"
wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="10"
charConfidence="100" serifProbability="255">M</charParams>
71 <charParams l="688" t="377" r="727" b="429" characterHeight="50" hasUncertainHeight="false" baseLine="0" wordStart="false"
wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="10"
charConfidence="100" serifProbability="255">U</charParams>
72 <charParams l="730" t="361" r="766" b="444" characterHeight="50" hasUncertainHeight="false" baseLine="0" wordStart="false"
wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="10"
charConfidence="100" serifProbability="255">S</charParams>
73 <charParams l="766" t="361" r="778" b="445" characterHeight="50" hasUncertainHeight="false" baseLine="0"> </charParams>
74 <charParams l="778" t="363" r="826" b="429" characterHeight="50" hasUncertainHeight="false" baseLine="0" wordStart="true"
wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="10"
charConfidence="100" serifProbability="255">E</charParams>
75 <charParams l="830" t="377" r="868" b="427" characterHeight="50" hasUncertainHeight="false" baseLine="0" wordStart="false"
wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="10"
charConfidence="100" serifProbability="255">U</charParams>
76 <charParams l="874" t="361" r="931" b="445" characterHeight="50" hasUncertainHeight="false" baseLine="0" wordStart="false"
wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="10"
charConfidence="100" serifProbability="255">C</charParams>
77 <charParams l="874" t="361" r="931" b="445" characterHeight="50" hasUncertainHeight="false" baseLine="0" wordStart="false"
wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="10"
charConfidence="100" serifProbability="255">h</charParams>
78 <charParams l="931" t="361" r="944" b="429" characterHeight="50" hasUncertainHeight="false" baseLine="0"> </charParams>
79 <charParams l="944" t="376" r="966" b="427" characterHeight="50" hasUncertainHeight="false" baseLine="0" wordStart="true"
wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="3" meanStrokeWidth="10"
charConfidence="100" serifProbability="255">E</charParams>
80 <charParams l="970" t="368" r="986" b="424" characterHeight="50" hasUncertainHeight="false" baseLine="0" wordStart="false"
wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="3" meanStrokeWidth="10"
charConfidence="100" serifProbability="255">I</charParams>
81 <charParams l="989" t="376" r="1026" b="425" characterHeight="50" hasUncertainHeight="false" baseLine="0" wordStart="false"
wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="3" meanStrokeWidth="10"
charConfidence="100" serifProbability="255">n</charParams>
82 <charParams l="1026" t="364" r="1042" b="429" characterHeight="50" hasUncertainHeight="false" baseLine="0"> </charParams>
83 <charParams l="1042" t="364" r="1060" b="442" characterHeight="50" hasUncertainHeight="false" baseLine="0" wordStart="true"
wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="10"
charConfidence="100" serifProbability="255">j</charParams>
84 <charParams l="1063" t="378" r="1085" b="427" characterHeight="50" hasUncertainHeight="false" baseLine="0" wordStart="false"
wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="10"
charConfidence="100" serifProbability="255">e</charParams>
85 <charParams l="1089" t="364" r="1122" b="429" characterHeight="50" hasUncertainHeight="false" baseLine="0" wordStart="false"
wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="10"
charConfidence="100" serifProbability="255">d</charParams>
86 <charParams l="1126" t="378" r="1149" b="429" characterHeight="50" hasUncertainHeight="false" baseLine="0" wordStart="false"
wordFromDictionary="true" wordNormal="true" wordNumeric="false" wordIdentifier="false" wordPenalty="0" meanStrokeWidth="10"
charConfidence="100" serifProbability="255">e</charParams>

```









**HK-XML**

Mit Koordinaten auf Wort- statt auf Zeichenebene ist auch ein deutlich kompakteres Format möglich:

```

1 <?xml version="1.0" encoding="utf-8"?>
2 <IMAGE Name="00000008" ImageX="1710" ImageY="2311">
3   <WORD Top="0,07529" Left="0,31053" Right="0,41573" Bottom="0,11164">Lasst</WORD>
4   <WORD Top="0,07702" Left="0,42749" Right="0,51696" Bottom="0,10775">Herr</WORD>
5   <WORD Top="0,07789" Left="0,53158" Right="0,7" Bottom="0,11813">Knobloch</WORD>
6   <WORD Top="0,08265" Left="0,7117" Right="0,73099" Bottom="0,1151">/</WORD>
7   <WORD Top="0,08178" Left="0,74444" Right="0,83392" Bottom="0,11856">Euch</WORD>
8   <WORD Top="0,08741" Left="0,85263" Right="0,92749" Bottom="0,1138" LB="True">auch</WORD>
9   <WORD Top="0,12289" Left="0,71696" Right="0,84327" Bottom="0,15361">weisen</WORD>
10  <WORD Top="0,12289" Left="0,89474" Right="0,92749" Bottom="0,15707" LB="True">^</WORD>
11  <WORD Top="0,14929" Left="0,24503" Right="0,44503" Bottom="0,18693">Sonstmuß</WORD>
12  <WORD Top="0,15232" Left="0,45263" Right="0,54152" Bottom="0,18866">Euch</WORD>
13  <WORD Top="0,15881" Left="0,54912" Right="0,59708" Bottom="0,18044">ein</WORD>
14  <WORD Top="0,15405" Left="0,60702" Right="0,68596" Bottom="0,18347">jeder</WORD>
15  <WORD Top="0,1614" Left="0,69532" Right="0,81111" Bottom="0,18563" LB="True">preisen</WORD>
16  <WORD Top="0,18823" Left="0,24269" Right="0,3269" Bottom="0,22328">Daß</WORD>
17  <WORD Top="0,18953" Left="0,33333" Right="0,38246" Bottom="0,21549">jhr</WORD>
18  <WORD Top="0,18866" Left="0,3924" Right="0,45614" Bottom="0,21679">sehr</WORD>
19  <WORD Top="0,1891" Left="0,46608" Right="0,63743" Bottom="0,21852">bescheiden</WORD>
20  <WORD Top="0,19039" Left="0,65146" Right="0,72047" Bottom="0,22068">seyd</WORD>
21  <WORD Top="0,19732" Left="0,72982" Right="0,74211" Bottom="0,22112" LB="True">:</WORD>
22  <WORD Top="0,22328" Left="0,24152" Right="0,32807" Bottom="0,2527">Jetzt</WORD>
23  <WORD Top="0,22415" Left="0,33626" Right="0,40409" Bottom="0,25314">lasst</WORD>
24  <WORD Top="0,23107" Left="0,41345" Right="0,67895" Bottom="0,26222">evreHoffnung</WORD>
25  <WORD Top="0,23453" Left="0,68889" Right="0,84386" Bottom="0,26222" LB="True">mercken</WORD>
26  <WORD Top="0,26049" Left="0,24211" Right="0,38012" Bottom="0,28949">Travet</WORD>
27  <WORD Top="0,26093" Left="0,3883" Right="0,48012" Bottom="0,28949">Gott</WORD>
28  <WORD Top="0,27174" Left="0,48596" Right="0,50175" Bottom="0,27564">/^</WORD>
29  <WORD Top="0,26785" Left="0,50819" Right="0,57719" Bottom="0,29078">aus</WORD>
30  <WORD Top="0,26395" Left="0,5848" Right="0,68187" Bottom="0,29338">dessen</WORD>
31  <WORD Top="0,26525" Left="0,68947" Right="0,84211" Bottom="0,29641" LB="True">Wercken</WORD>
32  <WORD Top="0,29598" Left="0,24094" Right="0,33333" Bottom="0,32843">Ihm</WORD>
33  <WORD Top="0,3042" Left="0,33977" Right="0,42105" Bottom="0,33535">noch</WORD>
34  <WORD Top="0,29771" Left="0,42807" Right="0,53216" Bottom="0,32843">keines</WORD>
35  <WORD Top="0,29944" Left="0,54035" Right="0,59766" Bottom="0,32843">hat</WORD>
36  <WORD Top="0,30723" Left="0,60468" Right="0,73743" Bottom="0,33276" LB="True">gerewt.</WORD>
37 </IMAGE>

```

## 7.2 Kurzanleitungen für einfache Fälle

### 7.2.1 OCR (Muster- und Wortbibliothek müssen vorliegen):

#### (a) BIT-Alpha (Beispiel mit Musterbibliothek in bda-Datei und lexikalischer Korrektur):

##### Einzelseite

###### *Vorbereitung:*

bda-Datei bereitstellen (für Binarisierungseinstellungen und Muster) (s. Kap. 4.1)

lx2-Datei bereitstellen (für Wortkorrektur) (s. Kap. 4.1)

###### *BIT-Alpha-Sitzung:*

File > Open

    bda-Datei wählen

(Nur falls in einer bit-Datei gespeicherte Muster verwendet werden sollen:)

OCR > Library > Load

    bit-Datei wählen

(Nur falls in einer seq-Datei gespeicherte Sequenzen verwendet werden sollen:)

OCR > Sequences > Load

    seq-Datei wählen

OCR > Lexical Correction > Load

    lx2-Datei wählen

File > Settings > OCR

    aktivieren:

        Iterate recognition

        Disjunction

        Spellcheck

File > Load image

    TIFF-Bild wählen

        => OCR beginnt

File > Export

    gewünschtes Ausgabeformat wählen

##### Stapelverarbeitung

###### *Vorbereitung:*

Zu bearbeitende Bilder in einem Verzeichnis versammeln

bda-Datei bereitstellen (für Binarisierungseinstellungen und Muster) (s. Kap. 4.1)

lx2-Datei bereitstellen (für Wortkorrektur) (s. Kap. 4.1)

###### *BIT-Alpha-Sitzung:*

File > Open

    bda-Datei wählen

OCR > Lexical Correction > Load

    lx2-Datei wählen

File > Settings > OCR

    aktivieren:

        Iterate recognition

        Disjunction

        Spellcheck

File > Batch process ...

Bild-Verzeichnis und Export-Verzeichnis und Export-Formate einstellen (Achtung,  
„Autolearn“ hier im Regelfall deaktivieren!)

> OK

## (b) HK-OCR (Beispiel mit Musterbibliothek und Nutzerwörterbuch):

### Einzelseite und Stapelverarbeitung

#### **Vorbereitung:**

Zu bearbeitende Bilder in einem Verzeichnis versammeln

Bei der Installation werden die Pfade der Muster- und Sprachdateienverzeichnisse eingestellt,  
z. B.:

<VERZEICHNIS>\UserPattern

<VERZEICHNIS>\CustomLanguageFolder

<VERZEICHNIS>\CustomBaseLanguageFolder

<VERZEICHNIS>\DictionaryFolder

In diesen Verzeichnissen Muster- und Sprachdateien bereitstellen (s. Kap. 4.1)

#### **HK-OCR-Sitzung:**

Bild öffnen > Ordner wählen

Bildverzeichnis einstellen (ist zugleich Exportverzeichnis)

Reiter „Einstellungen“

Standardsprachen

alle deaktivieren, wenn eigene Sprachengruppe genommen wird

Eigene Sprachengruppen

aktivieren der gewünschten Sprachengruppe (Haken setzen kann 2 Klicks erfordern)

Erkennungsmuster

aktivieren: „Musterdatei verwenden“

deaktivieren: „interne Muster verwenden“

Button „...“: Musterdatei auswählen

Schriften

aktivieren: „Normal“

aktivieren: „Fraktur“

deaktivieren: „Schreibmaschine“

weitere Optionen

aktivieren: „Bild ausrichten“

aktivieren: „verw. vorhandene Layouts“

Qualitätsoptimierung

aktivieren: „Voting“

aktivieren: „Versuchsautomatik“

„4 Versuche max.“

aktivieren: „Automatisch Speichern“

Reiter „Speichern“

> XML, PDF, RTF

Ausgabeformate einstellen, mindestens „FR-XML“

> Nach Speichern

alles deaktivieren

OCR starten

[für Einzelseite:] > aktuelles Bild lesen

[für Stapelverarbeitung:] > Alle Bilder lesen

## 7.2.2 Training (Anlegen bzw. Erweitern von Musterbibliotheken)

### (a) BIT-Alpha

#### **Musterdatei wählen oder erstellen:**

(Falls bda-Datei schon existiert:)

File > Open  
bda-Datei wählen

(Nur falls in einer bit-Datei gespeicherte Muster verwendet werden sollen:)

OCR > Library > Load  
bit-Datei wählen

(Nur falls in einer seq-Datei gespeicherte Sequenzen verwendet werden sollen:)

OCR > Sequences > Load  
seq-Datei wählen

#### **Bild laden:**

File > Load image  
TIFF-Bild wählen  
=> OCR-Analyse der Seite beginnt, abwarten

#### **Training**

View > Layer > Analyzed  
sollte Zeichensegmentierung erkennen lassen (wenn nicht: Binarisierungseinstellungen verbessern)

View > Layer > OCR  
zeigt in jedem Zeichensegment entweder das binarisierte Bild der Zeichenregion oder, wenn ein mehr oder weniger passendes Muster gefunden wurde, das „erkannte“ Zeichen in einem farbigen Feld von rot („sehr unsicher“) über orange und gelb zu grün („sichere“ Zuordnung zu einem sehr ähnlichen Muster) bis zu blau (für genau dieses Binärbild ist schon beim Training ein Zeichen gelernt worden).

OCR > Autorefresh

- entweder aktivieren (nach jeder Trainingseingabe wird die Musterbibliothek neu kompiliert, das führt jeweils zu mehreren Minuten Wartezeit, zeigt dann aber im Layer „OCR“ sofort den neuen Erkennungsstand der Seite an)
- oder deaktivieren (Trainingseingaben werden ohne Einkompilierung und ohne Neuanzeige der Seite zwischengespeichert), dafür kann ohne Zeitverzug weitertrainiert werden. Eine Neukompilierung der Musterbibliothek und damit Anzeige-Aktualisierung erreicht man jederzeit durch einen Aufruf:

OCR > Refresh

Im Layer „OCR“ kann nun frei das nächste zu trainierende Zeichen<sup>101</sup> gewählt werden (mit Maus anklicken). Es erscheint das Trainingsfenster „Verify Recognition“, wo das binäre Zeichenmuster in verschiedenen Ansichten gezeigt wird und im Feld „Recognition result > this String-Value“ der Zeichenwert für dieses Muster eingegeben werden kann. Im Rahmen „Action“ wird durch „Update“ der Zeichenwert dem Muster zugeordnet.

<sup>101</sup> Zweckmäßigerweise wird man zunächst die sehr schlecht erkannten Zeichen (rot und orange unterlegte) oder die noch gar nicht erkannten Zeichen (binäres Image im Segment) wählen.

**Speichern der Muster nach Ende des Trainings:**

- entweder zusammen mit den Binarisierungs- und Segmentierungsparametern in die bda-Datei:  
File > Save
- oder separat in eine bit-Datei:  
OCR > Library > Save

**(b) HK-OCR****Musterdatei wählen oder erstellen:**

Reiter „Einstellungen“ > Rahmen „Erkennungs-Muster“ > Button „,..“

**Bild laden:**

Bild öffnen > Ordner wählen > Ordner wählen  
In Rahmen „Stapel“ Bilddatei auswählen (z. B. Doppelklick)

**Training:**

Reiter „Einstellungen“

Rahmen „Erkennungs-Muster“:

Checkbox „Musterdatei verwenden“ deaktivieren

Checkbox „Training“ aktivieren

Rahmen „Schriften“

mindestens eine Checkbox für Schriftfamilie aktivieren (z. B. „Fraktur“)

Rahmen „Standardsprachen“ und/oder „Eigene Sprachengruppen“

passende Sprache wählen (hiervon hängt der mögliche Zeichenvorrat ab)

OCR starten > Aktuelles Bild lesen

Nach Einlesen des Bilds erscheint das Trainingsfenster „Pattern Training“ und springt zum ersten bzw. nächsten Zeichensegment.<sup>102</sup> Im Eingabefeld kann ein Zeichen oder eine Zeichenfolge („ligature“) neu eingetragen bzw. geändert und mit „Train“ bestätigt werden. Wenn bereits das richtige Zeichen vorgeschlagen wurde, kann es mit „Train“ bestätigt oder mit „Skip“ übergangen werden. Muster, die nicht trainiert werden sollen, müssen mit „Skip“ übergangen werden. In begrenztem Maße können die Segmentgrenzen vergrößert oder verkleinert werden (Ziehen mit der Maus oder Buttons „<<“, bzw. „>>“, soweit aktiv).

**Speichern der Muster nach Ende des Trainings:**

Nach Schließen des Trainingsfensters können die trainierten Muster in die ptn-Datei gespeichert werden, nach Erreichen des Seitenendes geschieht das automatisch (zur Kontrolle ist empfohlen, Datum und aktuelle Dateigröße der ptn-Datei im Musterdateiverzeichnis zu überprüfen).

<sup>102</sup> Das nächste zu trainierende Zeichen kann nicht frei ausgewählt werden; man wird von Anfang bis Ende durch die Seite geführt.

## 7.3 Details zu einzelnen Parametern

### (a) Parametervielfalt BIT-Alpha

Hinweis: Die angegebenen Informationen stammen teils aus dem knappen Handbuch, teils aus der Schulung des Herstellers und teils aus eigenen Tests.

#### Systemeinstellungen:

*Menü File > Settings > System > Thread Priority:*

<i>Zweck:</i>	Priorität einstellen, mit der BIT-Alpha den Prozessor belasten darf
<i>Werte:</i>	„Idle“ bis „Time critical“ (sehr zu Lasten aller anderen Anwendungen)
<i>Verwendet:</i>	„Normal“

#### Einstellung einer bda-Datei:

Empfehlung: ein Image mit mittlerer Helligkeit und mittlerem Kontrast laden, um Parameteränderungen testen zu können.

#### File > Settings [in empfohlener Reihenfolge der Einstellungen bei Beginn:]

##### – ganze Seite und erste Segmentierung in Textblöcke betreffend:

*Menü File > Settings > Color Enhancement:*

#### **Settings**

<i>Zweck:</i>	Helligkeit, Kontrast und Gammakorrektur für das Image so einstellen, dass eine optimale Binarisierung erreicht wird.
<i>Orientierung:</i>	Ansicht im Layer „Binary“. Änderungen eines der drei Parameter können Nachjustierungen der anderen erforderlich machen.
– Contrast	
verwendete Werte:	0.5–0.6 (Durchschnitt: 0.52, am häufigsten: 0.5)
– Intensity	
typischer Wertebereich:	0.4–0.6; verwendet: 0.4–0.51 (Durchschnitt: 0.49, am häufigsten: 0.5)
Wann erhöhen?	z. B. wenn Buchstaben zusammenkleben
Wann erniedrigen?	z. B. wenn das Bild zu hell ist
– Gamma	verwendet: 1

*Menü File > Settings > Preprocessing:*

#### **Margins**

<i>Zweck:</i>	Randbereiche des Images von der OCR ausschließen
<i>Orientierung:</i>	Ansicht in Layer „Binary“: dunkelblauen Rahmen passend einstellen
– left (cm)	
verwendete Werte:	0.5–0.7 (Durchschnitt: 0.52, am häufigsten: 0.5)
– right (cm)	
verwendete Werte:	0.5–0.7 (Durchschnitt: 0.52, am häufigsten: 0.5)

- top (cm)  
verwendete Werte: 0.5–1.2 (Durchschnitt: 0.71, am häufigsten: 0.7)
- bottom (cm)  
verwendete Werte: 0.5–1.2 (Durchschnitt: 0.72, am häufigsten: 0.7)

### ***Blackborder-Elimination***

*Orientierung:* Ansicht in Layer „Binary“: hellblauen Rahmen passend einstellen

- left (cm)  
verwendete Werte: 0.5–0.7 (am häufigsten: 0.5)
- right (cm)  
verwendete Werte: 0.5–0.7 (am häufigsten: 0.5)
- top (cm)  
verwendete Werte: 0.7–1.2 (am häufigsten: 0.7)
- bottom (cm)  
verwendete Werte: 0.5–1.2 (am häufigsten: 0.7)
- Enabled  
verwendet: in der Regel aktiviert

### ***Binarisation***

*Zweck:* Einstellung der „ersten“ (für die Segmentierung verwendeten) Binarisierung

*Algorithmen:* „Color based algorithm“, „BIT“, „Intensity based algorithm“, „Modified Niblack algorithm“

*Verwendete Auswahl:* „Modified Niblack algorithm“ (für Fraktur empfohlen)

- Configure
  - > Intensity based algorithm
  - > Relative Intensity threshold %  
verwendete Werte: 35–60 (Durchschnitt: 46.94, am häufigsten: 50)  
Wann erhöhen: bei hellem Bild  
Wann erniedrigen: bei dunklem Bild
- „Clean binarised image“  
*Zweck:* entfernt „Wolkenstaub“;  
Wann deaktivieren: bei sehr „dünnen“ Vorlagen mit schwachem Kontrast
- „Remove dots“  
*Zweck:* entfernt „größere“ Punkte  
Wann aktivieren: z. B. bei Mikrofilm; rechnerintensiv
- min. dist. cm  
verwendete Werte: 0.02–0.09 (Durchschnitt: 0.05, am häufigsten: 0.04)
- max. size cm  
verwendete Werte: 0.05–0.15 (Durchschnitt: 0.10, am häufigsten: 0.09)

### ***Rotation***

*Zweck:* Seitenausrichtung anhand erkannter Linien

*Algorithmen:* „Linear Algorithm“ (versucht, Rechtwinkligkeit wiederherzustellen), „Circular Algorithm“ (Rotation, Winkel bleiben erhalten)

*Verwendet:* „Circular Algorithm“ und ohne Rotation

- min. angle (degrees)  
verwendete Werte: 0.01

*Menü File > Settings > Lines V***Limits (im Image gespeicherte Maßeinheiten und Auflösung müssen korrekt sein)**

- Zweck:* Entfernung von vertikalen Linien
- max. thickness cm  
*verwendete Werte:* 0.1–0.8 (Durchschnitt: 0.54, am häufigsten: 0.8)
- min. thickness cm  
*verwendete Werte:* 0.01
- min. black length cm  
*Zweck:* minimale Länge eines gedruckten „Zeichens“ in eventuell unterbrochenen Linien  
*verwendete Werte:* 0.4–0.9 (Durchschnitt: 0.54, am häufigsten: 0.4)
- max. white length cm  
*Zweck:* maximaler weißer Zwischenraum in eventuell gepunkteten oder unterbrochenen Linien  
*verwendete Werte:* 0
- min. total length cm  
*Zweck:* minimale totale Länge, um zu definieren, ab wann eine Linie als solche erkannt werden soll, die eventuell gepunktet oder unterbrochen ist  
*typischer Wertebereich:* -3  
*verwendete Werte:* 0.75–10 (Durchschnitt: 5.24, am häufigsten: 10)  
*Wann erhöhen:* + wenn „falsche Linien“ (z. B. aus übereinander stehenden m-Füßen) erkannt werden
- Use bitmap  
*Zweck:* Gescanntes Bild der Linie wiedergeben, statt sie ideal nachzuzeichnen  
*verwendete Werte:* Aktiviert

*Menü File > Settings > Lines H***Limits (im Image gespeicherte Maßeinheiten und Auflösung müssen korrekt sein)**

- Zweck:* Entfernung von horizontalen Linien zur besseren Absatztrennung (aber nicht erwünscht, wenn keine Linien da sind)
- max. thickness cm  
*Zweck:* maximale Breite des Strichs der Linie  
*verwendete Werte:* 0.1–0.5 (Durchschnitt: 0.36, am häufigsten: 0.4)
- min. thickness cm  
*Zweck:* minimale Breite des Strichs der Linie  
*verwendete Werte:* 0
- min. black length cm  
*Zweck:* minimale Länge eines gedruckten „Zeichens“ in eventuell unterbrochenen Linien  
*verwendete Werte:* 0.4–0.9 (Durchschnitt: 0.59, am häufigsten: 0.5)

- max. white length cm  
*Zweck:* maximaler weißer Zwischenraum in eventuell gepunkteten oder unterbrochenen Linien  
 verwendete Werte: 0.1
- min. total length cm  
*Zweck:* minimale totale Länge, um zu definieren, ab wann eine Linie als solche erkannt werden soll, die eventuell gepunktet oder unterbrochen ist  
 verwendete Werte: 0.75–10 (Durchschnitt: 5.13, am häufigsten: 10)
- Use bitmap  
*Zweck:* Gescanntes Bild der Linie wiedergeben, statt sie ideal nachzuzeichnen  
 verwendete Werte: Aktiviert

*Menü File > Settings > Segmentation*

***Parameters (im Image gespeicherte Maßeinheiten und Auflösung müssen korrekt sein)***

- Zweck:* Einstellungen für die Segmentierung der Seite in verschiedene Regionen, d. h. Paragraphen, Titel, Untertitel und graphische Elemente
- min. horiz. distance cm  
*Zweck:* die minimal erlaubte horizontale Distanz zwischen benachbarten Regionen. Näher beieinander liegende Regionen werden zusammengefügt.  
 verwendete Werte: 1–2 (Durchschnitt: 1.83, am häufigsten: 2)  
 Praxiserfahrungen: Bei zunächst deaktivierter „Dissection“ justieren, bis in der „Binary“-Ansicht der prozessierten Seite die grauen Teilrahmen korrekt sind (z. B. Satzspiegel)
- min. vert. distance cm  
*Zweck:* die minimal erlaubte vertikale Distanz zwischen benachbarten Regionen. Näher beieinander liegende Regionen werden zusammengefügt.  
 verwendete Werte: 0.5–2 (Durchschnitt: 1.69, am häufigsten: 2)
- min. width cm  
*Zweck:* die minimal erlaubte horizontale Breite für die Erstellung einer Region.  
 verwendete Werte: 0.1
- min. height cm  
*Zweck:* die minimal erlaubte vertikale Höhe für die Erstellung einer Region.  
 verwendete Werte: 0.1
- Default Region Type  
*Zweck:* BIT-Alpha teilt jeder bei der Segmentierung kreierten Region Eigenschaften über den Typ des Bilds in der Region zu.  
 Werte: „Binary Region“: monochromes Bild (Schwarz-weiß), „Palette Region“: Graustufenbild, jeder Pixel mit 8 bit kodiert, „color region“: Farbbild, jeder Pixel mit 8 oder 16 bit kodiert  
 verwendete Werte: „Binary Region“

**File > Settings****– betreffend Segmentierung in Unterblöcke:***Menü File > Settings > Segmentation****Dissection***

*Zweck:* betrifft zusätzlich Unterteilung nach Abständen zwischen Textblöcken; Teilblöcke werden dann automatisch separat eingestellt

*Praxiserfahrungen:* Justieren, bis in der Binary-Ansicht der prozessierten Seite die grauen Teilrahmen korrekt sind (z. B. Absätze, Strophen ...). Nachjustieren kann Abhilfe schaffen, wenn in manchen Textblöcken keine Zeichen-Segmentierung erfolgt ist.

## – min. width cm

*Zweck:* Versuch, ab dieser Breite den Textblock zu teilen  
 verwendete Werte: 5-10 (Durchschnitt: 7.94, am häufigsten: 9)

## – min. height cm

*Zweck:* Versuch, ab dieser Höhe den Textblock zu teilen  
 verwendete Werte: 4–10 (Durchschnitt: 5.00, am häufigsten: 4)

## – enable hor. dissection

*Zweck:* versucht, entlang „waagerechter weißer Streifen“ zu teilen  
 verwendete Werte: Aktiviert

## – hor. dissection threshold % (Prozent-Anteil weiß im Zwischenraum)

verwendete Werte: 95–99 (Durchschnitt: 97.83, am häufigsten: 98)

## – hor. dissection thickness cm (Mindestbreite weiß)

verwendete Werte: 0.01-0.3 (Durchschnitt: 0.13, am häufigsten: 0.2)

## – enable vert. dissection

*Zweck:* versucht, entlang „senkrechter weißer Streifen“ zu teilen  
 verwendete Werte: Aktiviert

## – vert. dissection threshold %

verwendete Werte: 95-99 (Durchschnitt: 98.11, am häufigsten: 99)

## – vert. dissection thickness cm

verwendete Werte: 0.01-0.5 (Durchschnitt: 0.24, am häufigsten: 0.2)

Wann verändern: Wenn Textblöcke nicht als solche erkannt werden und „Image detection“ nicht deaktiviert werden soll

***Resolution***

*Zweck:* Angabe der korrekten Auflösung, falls in den Image-Metadaten nicht korrekt enthalten

verwendete Werte: Aktiviert

***Image detection***

*Zweck:* Erkennung von Bildregionen (Illustrationen usw.), die von der OCR ausgenommen werden sollen

Wann deaktivieren: wenn Textblöcke fälschlich als Bild interpretiert werden (keine Zeichensegmentierung in Ansicht „Analyzed“) und eine Änderung der „vert. dissection thickness cm“ keine Abhilfe bringt

- enabled  
verwendete Werte: Aktiviert, Deaktiviert (am häufigsten: Aktiviert)
- hor.  
verwendete Werte: 0.45-1 (Durchschnitt: 0.57, am häufigsten: 0.45)
- vert.  
verwendete Werte: 0.45-1 (Durchschnitt: 0.57, am häufigsten: 0.45)

## File > Settings

### – betreffend Zeichenregionen:

*Menü File > Settings > Binary regions*

#### **Detection**

- Detect binary regions  
*Zweck:* Unterscheidung von Zeichen (für OCR) und Bildregionen (für bildliche Wiedergabe)  
*Hinweis:* Parametrierung unnötig, wenn in Preprocessing > Segmentation > Default Region Type „Binary Region“ eingestellt wurde  
verwendete Werte: Aktiviert, Deaktiviert (am häufigsten: Aktiviert)

#### **Binarisation**

- Praxiserfahrungen:* In Ansicht „Analysed“ und „OCR“ auf gute Zeichensegmentierung achten und darauf, dass keine Textblöcke aus der Zeichensegmentierung herausfallen
- Algorithmen:* „Color based algorithm“  
„BIT-Algorithmus“  
„Intensity based algorithm“  
„Modified Niblack algorithm“  
– *Konfigurationsparameter siehe unten* –
- Clean binarised image  
*Empfehlung:* Deaktivieren bei sehr dünnen (aus „Punktwolken“ bestehenden) Schriftzeichen  
verwendete Werte: Aktiviert, Deaktiviert (am häufigsten: Aktiviert)
- Remove dots:  
verwendete Werte: Aktiviert, Deaktiviert (am häufigsten: Aktiviert)
- min. dist. cm  
verwendete Werte: 0.02-0.13 (Durchschnitt: 0.07, am häufigsten: 0.05)
- max. size cm  
verwendete Werte: 0.05-0.25 (Durchschnitt: 0.12, am häufigsten: 0.08)

#### **Binarisation > Color based algorithm**

- verwendete Werte: Deaktiviert

Menü *File* > *Settings* > *Binary regions* [Forts.]

***Binarisation > BIT-Algorithmus > Configure***

- Hinweis:* vom Hersteller für typische Drucke empfohlen; weniger für Zeitungen und Fraktur
- verwendete Werte: Aktiviert, Deaktiviert (am häufigsten: Deaktiviert)
- Minimum > horiz.  
verwendete Werte: 0.01-0.025 (am häufigsten: 0.02)  
Wann erhöhen: bei blassem Druck
  - Minimum > vert.  
verwendete Werte: 0.01-0.025 (am häufigsten: 0.02)
  - Maximum > horiz.  
verwendete Werte: 2-4 (am häufigsten: 2)  
Wann erhöhen: bei blassem Druck
  - Maximum > vert.  
verwendete Werte: 2-4 (am häufigsten: 2)
  - Binarisation > cut-off  
verwendete Werte: 0.4-0.53 (am häufigsten: 0.5)  
Wann erhöhen: bei blassem Druck (wenn Buchstaben dunkler werden sollen)  
Wann erniedrigen: gegen Durchdruck

***Binarisation > Intensity based algorithm > Configure***

- verwendete Werte: Aktiviert, Deaktiviert (am häufigsten: Deaktiviert)
- Relative Intensity threshold %  
typischer Wertebereich: 50-70  
verwendete Werte: 50, 61 (am häufigsten: 50)  
Wann erhöhen: wenn schwacher Kontrast  
Wann erniedrigen: wenn Zeichen zusammenkleben (zu dunkel)

***Binarisation > Modified Niblack algorithm > Configure***

- Hinweis:* nutzt dynamischen Kontrast (Varianz der Helligkeit) um jeden Buchstaben herum, erkennt eher Kontur als „schwere“ Innenbereiche; für Fraktur empfohlen
- verwendete Werte: Aktiviert, Deaktiviert (am häufigsten: Aktiviert)
- Noise reduction:
    - Enable noise reduction  
verwendete Werte: Aktiviert
    - Vertical noise frequency (mm):  
verwendete Werte: 0.25-0.5 (am häufigsten: 0.5)
    - Horizontal noise frequency (mm):  
verwendete Werte: 0.25-0.5 (am häufigsten: 0.5)
  - Configuration:
    - Vertical adjustment radius (cm):  
verwendete Werte: 0.2-0.25 (am häufigsten: 0.25)
    - Horizontal adjustment radius (cm):  
verwendete Werte: 0.2-0.25 (am häufigsten: 0.25)
    - Noise threshold (%):  
verwendete Werte: 9-10 (am häufigsten: 10)  
Wann erhöhen: wenn zu dunkel  
Wann erniedrigen: wenn Zeichen brüchig sind
    - Noise factor:  
verwendete Werte: 1

Menü File > Settings > OCR

### **Target Size**

- |                           |   |
|---------------------------|---|
| <i>Zweck:</i>             | minimale und maximale akzeptierte Zeichengröße              |
| – min. width (cm)         |   |
| verwendete Werte:         | 0.03–0.05 (Durchschnitt: 0.04, am häufigsten: 0.03)         |
| Wann erniedrigen:         | wenn Text fälschlich als Miniatur erkannt wird              |
| <i>Praxiserfahrungen:</i> | in Ansicht „OCR“ so groß wie ohne Lücken möglich einstellen |
| – max. width (cm)         |   |
| verwendete Werte:         | 0.035–6 (Durchschnitt: 3.84, am häufigsten: 4)              |
| <i>Hinweis:</i>           | ggf. an Initialen denken                                    |
| – min. height (cm)        |   |
| verwendete Werte:         | 0.03–0.05 (Durchschnitt: 0.04, am häufigsten: 0.03)         |
| – max. height (cm)        |   |
| verwendete Werte:         | 3-4 (Durchschnitt: 3.83, am häufigsten: 4)                  |
| <i>Hinweis:</i>           | ggf. an Initialen denken                                    |

### **Options**

- |                                 |  |
|---------------------------------|--|
| <i>Hinweis:</i>                 | Optionen primär für OCR-Lauf sinnvoll; im Training dagegen diese zunächst deaktivieren! Nach dem ersten Training kann erst mit „Iterate“, dann mit „Iterate+Disjunction“ das Training wiederholt werden. |
| – Iterate recognition           |  |
| <i>Zweck:</i>                   | entfernt z. B. Punkte bei nicht erkannten Zeichen.   |
| verwendete Werte (im OCR-Lauf): | Aktiviert  |
| – Disjunction                   |  |
| <i>Zweck:</i>                   | Trennung von zusammengeklebten Charakteren für die OCR-Bearbeitung   |
| verwendete Werte (im OCR-Lauf): | Aktiviert, Deaktiviert (am häufigsten: Aktiviert)  |
| – Spellcheck                    |  |
| <i>Hinweis:</i>                 | Sprachkorrektur per Wortliste für die OCR-Bearbeitung aktivieren (s. „Lexika“)   |
| verwendete Werte (im OCR-Lauf): | Aktiviert, Deaktiviert (am häufigsten: Aktiviert)  |

Menü File > Settings > OCR > Advanced Settings > Edit  
(muss für jede Schriftfamilie separat eingestellt werden)

### **Target Size**

- |                    |                                 |
|--------------------|---------------------------------|
| – min. width (cm)  |                                 |
| verwendete Werte:  | 0.03-0.05 (am häufigsten: 0.04) |
| – max. width (cm)  |                                 |
| verwendete Werte:  | 3-6 (am häufigsten: 4)          |
| – min. height (cm) |                                 |
| verwendete Werte:  | 0.03-0.05 (am häufigsten: 0.04) |
| – max. height (cm) |                                 |
| verwendete Werte:  | 3-4 (am häufigsten: 4)          |

Menü *File* > *Settings* > *OCR* > *Advanced Settings* > *Edit [Forts.]*

### **Miniatures**

- Zweck:* Erkennung von Initialen
- Min. Overlap  
verwendete Werte: 0.02–0.3333 (Durchschnitt: 0.12, am häufigsten: 0.05)
  - min. X (rel.)  
verwendete Werte: 3–9 (Durchschnitt: 5.78, am häufigsten: 6)
  - min. Y (rel.)  
verwendete Werte: 3–9 (Durchschnitt: 5.78, am häufigsten: 6)

### **Classification**

- Zweck:* je nach Variabilität der Schriftform<sup>103</sup>, relativ zur Minuskelhöhe einstellbare Größenverhältnisse von: [P]atte (senkrecht Bein bzw. kleiner senkrechter Strich, z. B. *i* ohne Punkt), [M]inuscule (Kleinbuchstabe), [G]rand (Großbuchstabe bzw. großer Strich, [R]este (kleineres Element als Patte, Minuscule oder Grand, z. B. „*,-*“ oder „*-*“)
- Minuscule min. X  
verwendete Werte: 0.3-0.5 (Durchschnitt: 0.4, am häufigsten: 0.4)
  - Minuscule min. Y  
verwendete Werte: 0.5-0.8 (Durchschnitt: 0.71, am häufigsten: 0.7)
  - Minuscule max. Y  
verwendete Werte: 1.05-1.2 (Durchschnitt: 1.11, am häufigsten: 1.1)
  - Collation min. Y  
*Zweck:* Abschneiden von über die Gesamtgröße hinausgehenden Zeichen  
verwendete Werte: 1.8-2.5 (Durchschnitt: 2.03, am häufigsten: 2)  
*Praxiserfahrungen:* Heraufsetzen kann Abhilfe schaffen, wenn in Bereichen keine Zeichensegmentierung erfolgte.

### **Lines**

- Hinweis:* muss für jede Schriftfamilie separat eingestellt werden
- Zweck:* Zeileneigenschaften
- Hinweis:* In Menü und Handbuch wird „Line“ teils für „Zeile“, teils für „Linie“ verwendet.
- Topline min.  
verwendete Werte: 0
  - Bottomline min.  
verwendete Werte: 0.25-0.5 (Durchschnitt: 0.31, am häufigsten: 0.27)
  - Topline max.  
verwendete Werte: 0.25-0.6 (Durchschnitt: 0.37, am häufigsten: 0.35)
  - Bottomline max.  
verwendete Werte: 0.3-0.6 (Durchschnitt: 0.36, am häufigsten: 0.3)
  - Rigidity  
*Zweck:* Starrheit der (Zeilen-)Linienführung  
verwendete Werte: 5-9 (Durchschnitt: 6.44, am häufigsten: 7)  
Wann erhöhen: + um weniger Krümmung zu erlauben  
Wann erniedrigen: – um mehr Krümmung zu erlauben

<sup>103</sup> Die Parametrierung dieses Menüs wird grundsätzlich von B.I.T. selbst und nicht vom Benutzer vorgenommen. Die Einstellungen sind für Fraktur, Kursivschrift der Renaissance oder gerade Antiqua unterschiedlich.

- Min. Limit  
verwendete Werte: 0.3-0.9 (Durchschnitt: 0.63, am häufigsten: 0.8)
- Max. Limit  
verwendete Werte: 1-1.5 (Durchschnitt: 1.33, am häufigsten: 1.5)
- Large Smear  
Zweck: Berechnung der Zeilenerkennung mittels „Verwischen“ der Großbuchstaben  
verwendete Werte: 3-7 (Durchschnitt: 4.67, am häufigsten: 5)  
Praxiserfahrungen: Erniedrigen kann Abhilfe schaffen bei „verschluckten“ Zeilen
- Small Smear  
Zweck: Berechnung der Zeilenerkennung mittels „Verwischen“ der Kleinbuchstaben  
verwendete Werte: 0.3-0.5 (Durchschnitt: 0.41, am häufigsten: 0.5)
- Merge Distance  
Zweck: Einstellung, ab wann Zeilen als separat gesehen werden sollen  
verwendete Werte: 0.75-3 (Durchschnitt: 1.53, am häufigsten: 1)
- Merge Overlap  
verwendete Werte: 0.3-0.9 (Durchschnitt: 0.56, am häufigsten: 0.9)
- Up-Factor  
verwendete Werte: 0.13-0.5 (Durchschnitt: 0.24, am häufigsten: 0.2)
- Down-Factor  
verwendete Werte: 0.15-0.5 (Durchschnitt: 0.22, am häufigsten: 0.17)
- Horiz. Locality  
verwendete Werte: 2-4 (Durchschnitt: 3.61, am häufigsten: 4)
- Vert. Locality  
verwendete Werte: 1-2 (Durchschnitt: 1.83, am häufigsten: 2)

***R-R (Abstand bei Kombination der Klassen „Reste“-“Reste“)***

- Horizontal  
verwendete Werte: 0.05-0.3 (Durchschnitt: 0.12, am häufigsten: 0.1)
- Vertical  
verwendete Werte: 0.2-0.9 (Durchschnitt: 0.72, am häufigsten: 0.9)

***R-P (Abstand bei Kombination der Klassen „Reste“-“Patte“)***

- Horizontal  
verwendete Werte: 0.1-0.2 (Durchschnitt: 0.13, am häufigsten: 0.1)
- Vertical  
verwendete Werte: 0.7-1.5 (Durchschnitt: 0.98, am häufigsten: 0.9)

***R-M (Abstand bei Kombination der Klassen „Reste“-“Minuscule“)***

- Horizontal  
verwendete Werte: 0.05-0.2 (Durchschnitt: 0.12, am häufigsten: 0.1)
- Vertical  
verwendete Werte: 0.7-1.5 (Durchschnitt: 0.98, am häufigsten: 0.9)

***R-G (Abstand bei Kombination der Klassen „Reste“-“Grand“)***

- Horizontal  
verwendete Werte: 0.1-0.2 (Durchschnitt: 0.13, am häufigsten: 0.1)
- Vertical  
verwendete Werte: 0.4-1.5 (Durchschnitt: 1.02, am häufigsten: 0.9)

Menü File > Settings > OCR > Advanced Settings > Edit [Forts.]

**M-M (Abstand bei Kombination der Klassen „Minuscule“-“Minuscule“)**

- Horizontal  
verwendete Werte: 0.02-0.1 (Durchschnitt: 0.04, am häufigsten: 0.03)
- Vertical  
verwendete Werte: 0.9-1.5 (Durchschnitt: 0.98, am häufigsten: 0.9)

**P-P (Abstand bei Kombination der Klassen „Patte“-“Patte“)**

- Horizontal  
verwendete Werte: 0.05-0.2 (Durchschnitt: 0.12, am häufigsten: 0.1)
- Vertical  
verwendete Werte: 0.5-0.9 (Durchschnitt: 0.85, am häufigsten: 0.9)

**M-P (Abstand bei Kombination der Klassen „Minuscule“-“Patte“)**

- Horizontal  
verwendete Werte: 0.05-0.15 (Durchschnitt: 0.09, am häufigsten: 0.1)
- Vertical  
verwendete Werte: 0.5-0.9 (Durchschnitt: 0.88, am häufigsten: 0.9)

**G-P (Abstand bei Kombination der Klassen „Grand“-“Minuscule“)**

- Horizontal  
verwendete Werte: 0.05-0.2 (Durchschnitt: 0.08, am häufigsten: 0.05)
- Vertical  
verwendete Werte: 0.5-5 (Durchschnitt: 1.29, am häufigsten: 0.9)

**G-G (Abstand bei Kombination der Klassen „Grand“-“Grand“)**

- Horizontal  
verwendete Werte: 0.02-0.07 (am häufigsten: 0.05)
- Vertical  
verwendete Werte: 0.9-2 (Durchschnitt: 1.13, am häufigsten: 0.9)

**M-G (Abstand bei Kombination der Klassen „Minuscule“-“Grand“)**

- Horizontal  
verwendete Werte: 0.02-0.05 (Durchschnitt: 0.04, am häufigsten: 0.05)
- Vertical  
verwendete Werte: 0.9-2 (Durchschnitt: 1.13, am häufigsten: 0.9)

Menü File > Settings > Export

**PDF**

- only text  
Zweck: Bild nicht mitexportieren  
verwendete Werte: Deaktiviert
- one single image ebedded per page  
Zweck: Nur ein Bild pro Seite (Kopie des Original-Faksimiles)  
einbetten. Der per OCR gelesene Text ist für Textsuche und  
Highlighting als unsichtbares Hintergrundbild hinterlegt.  
verwendete Werte: Aktiviert, Deaktiviert (am häufigsten: Deaktiviert)

- one image embedded for each segmented region  
*Zweck:* Jede bei der Segmentierung der Seite erstellte Region auch als Bild einbetten (kann dann separat aus dem PDF geöffnet werden).  
 verwendete Werte: Aktiviert, Deaktiviert (am häufigsten: Aktiviert)
- Create high-resolution links for text regions  
 verwendete Werte: Aktiviert, Deaktiviert (am häufigsten: Deaktiviert)

**METS/ALTO**

- relative to the original image  
*Zweck:* Koordinaten auf das originale Bild beziehen (sinnvoll, wenn dieses und nicht auf das zur OCR ggf. gedrehte Bild präsentiert wird)  
 verwendete Werte: Aktiviert

**OCR**

*OCR > Library > Learning parameters ...*

**Learning configuration**

- min.:  
 typischer Wertebereich: 0.55 bei „use Spellcheck“, 0.8 ohne Spellcheck
- max.:  
 typischer Wertebereich: 0.75 bei „use Spellcheck“, 0.9 ohne Spellcheck
- use Spellcheck  
 verwendete Werte: Deaktiviert

*OCR > Library > Settings...*

**Recognition**

- mass-threshold  
 verwendete Werte: 0.12-0.2 (Durchschnitt: 0.15, am häufigsten: 0.14)
- position-threshold  
 verwendete Werte: 0.12-0.16 (Durchschnitt: 0.14, am häufigsten: 0.14)
- stretching-threshold  
 verwendete Werte: 0.12-0.3 (am häufigsten: 0.14)
- density-threshold  
 verwendete Werte: 0.12-0.2 (Durchschnitt: 0.15, am häufigsten: 0.14)
- confidence-threshold (nur in bestimmten BIT-Alpha-Versionen)  
 verwendete Werte: 0.7-0.8 (Durchschnitt: 0.78, am häufigsten: 0.8)
- blanks low-limit  
 verwendete Werte: 0.05-0.25 (Durchschnitt: 0.16, am häufigsten: 0.2)
- blanks high-limit  
 verwendete Werte: 0.6-1.5 (Durchschnitt: 0.78, am häufigsten: 0.6)
- blanks coefficient  
 verwendete Werte: 0.6-1.2 (Durchschnitt: 0.97, am häufigsten: 1)  
 Wann erhöhen: wenn mehr zusammengehalten werden soll  
 Wann erniedrigen: wenn mehr getrennt werden soll

OCR > Library > Settings... [Forts.]

### **Separation**

- min. confidence  
verwendete Werte: 0.7-0.75 (Durchschnitt: 0.73, am häufigsten: 0.75)
- min. confidence  
verwendete Werte: 0.6-0.75 (Durchschnitt: 0.67, am häufigsten: 0.65)
- min.width  
verwendete Werte: 0.4-0.5 (Durchschnitt: 0.49, am häufigsten: 0.5)
- limit (left)  
verwendete Werte: 0.2-0.3 (Durchschnitt: 0.22, am häufigsten: 0.2)
- limit (right)  
verwendete Werte: 1.2-1.7 (Durchschnitt: 1.44, am häufigsten: 1.5)
- min. (right)  
verwendete Werte: 0.15-0.2 (Durchschnitt: 0.19, am häufigsten: 0.2)
- min. thickness  
verwendete Werte: 0.1-0.5 (am häufigsten: 0.4)
- min. complexity (left)  
verwendete Werte: 1
- min. complexity (right)  
verwendete Werte: 1

### **Algorithm**

- by form  
verwendete Werte: Aktiviert
- by mass  
Zweck: empfohlen z. B. bei chinesischen Zeichen

OCR > Lexical Correction > Edit

### **Coefficients:**

- Default:  
verwendete Werte: 1-1.51 (am häufigsten: 1)
- Split:  
verwendete Werte: 0.70-1 (am häufigsten: 1)  
Wann erhöhen: Trennung erschweren (auf 2 oder 3 um Trennungen zu verhindern)  
Wann erniedrigen: Trennung erleichtern
- Symbol added/deleted:  
verwendete Werte: 1-2.1 (am häufigsten: 1)  
(Text-area)  
Zweck: Für einzelne Zeichenpaare gewichtete „Ersetzungskosten“ einstellen (0 = Ersetzung ungehindert; 1 = keine Ersetzung)

### **Signature:**

- Zweck: wie oft ein Wort korrigiert werden kann
- [Entry]  
verwendete Werte: 2 oder 3 (am häufigsten: 2)
- Disable correction for words starting with capital letters  
verwendete Werte: Deaktiviert

***Correction:***

- max. distance:  
verwendete Werte: 0.2-0.5 (Durchschnitt: 0.35, am häufigsten: 0.20)
- Wann erhöhen: wenn großzügiger ersetzt werden soll
- Wann erniedrigen: wenn zurückhaltender ersetzt werden soll

***Split mode:***

- Recurse fragments  
verwendete Werte: Aktiviert, Deaktiviert (am häufigsten: Deaktiviert)
- Split fragments  
verwendete Werte: Aktiviert, Deaktiviert (am häufigsten: Aktiviert)
- No fragmentation  
verwendete Werte: Deaktiviert

**(b) HK-OCR / FineReader: Gibt es Potentiale durch die Erschließung von ABBYY-seitig angelegten Konfigurationsmöglichkeiten?**

Über die Oberfläche von HK-OCR sind im wesentlichen nur die in den Kapiteln 4 und 5 beschriebenen Parameter konfigurierbar.

Für ein Interesse an eventuell weiterführenden, in HK-OCR derzeit noch nicht realisierten Konfigurationsmöglichkeiten der FineReader-Engine muss auf einschlägige ABBYY-Entwickler-Dokumentationen, speziell die darin enthaltenen Objekt-Spezifikationen verwiesen werden.

## 7.4 Einige Mindestanforderungen an Bedieneroberflächen

Man muss die inzwischen kaum noch überschaubare Menge von Normen, Styleguides und Empfehlungen zu Softwareergonomie<sup>104</sup> keineswegs systematisch durcharbeiten, um zumindest folgende (systematisch meist unter Steuerbarkeit und Erwartungskonformität einzuordnende) Anforderungen für plausibel zu halten. Sie spielten schon im Projektverlauf eine Rolle und können im OCR-Produktionseinsatz erst recht zum kritischen Faktor werden:

– **Jede Arbeitssituation muss jederzeit konsistent beendet werden können.**

Konfigurations-, Muster-, Lexikon- und alle anderen Dateien sollten jederzeit entladen werden können (Fortsetzung der Sitzung erkennbar mit den jeweiligen Default-Werten), auch ohne dass dazu eine neue entsprechende Datei geladen werden muss. (Die Anzeige der aktuellen Parameterdateinamen in Titel- oder Statuszeile muss dem natürlich folgen).

Besonders während Verarbeitungsschritten mit Zeitbedarf (häufig während Layoutanalyse und OCR) muss eine Möglichkeit zum sofortigen Abbruch angeboten werden.

Bei Abbrüchen von Stapelläufen muss es eine einfache Möglichkeit zur Fortsetzung am Abbruchpunkt geben, ohne dass Verzeichnisinhalte umorganisiert werden müssen.

– **Jede Arbeitssituation muss auf den ersten Blick erkennbar sein.**

Die Namen aller aktuell geladenen bzw. wirksamen Parameter-, Image- und sonstigen Dateien sollten dem Bearbeiter jederzeit ohne weitere Handgriffe, Mausektionen, Suchschritte vor Augen stehen, z. B. in der Titelzeile und/oder der Statuszeile.

Es muss erkennbar sein, ob und wann manuelle Parameteränderungen in einer Konfigurationsdatei gespeichert werden, d. h. diese verändern (was man in der aktuellen Sitzung nicht immer will).

Benutzeraktionen sollen stets (auch) über beschriftete Menüs erreichbar sein, nicht allein durch einen ikonographischen Button, dessen Bedeutung erst erlernt werden muss oder erst nach Mausektionen angezeigt wird.

Aktionen, die der Anwender wahlweise anstoßen möchte oder nicht, sollten je explizit durch einen passenden Menübefehl ausgelöst werden und nicht implizit in einem unerwarteten Zusammenhang beginnen.

In Dateiauswahl-Dialogen muss statt der Verwendung des mehrdeutigen Labels „Öffnen“ jeweils erkennbar sein, ob eine gewählte Datei (a) zum Lesen bzw. Betrachten, (b) zum Bearbeiten oder (c) zum Ausführen darin enthaltenen Codes geladen wird.

– **Safety first: Bedienelemente müssen eine sichere Arbeitsweise erleichtern und unsichere Aktionen erschweren.**

In keiner Situation darf die riskantere Aktion leichter fallen als die sicherere Aktion. In Editier-Umgebungen muss das Übergehen bzw. Verwerfen einer Eingabe (per Tasten oder Button-Fokus) mindestens genauso schnell erreichbar sein wie die Bestätigung einer Eingabe, die etwas verändert.

Für Situationen, in denen aus Dateien, Datensätzen usw. nur passiv gelesen werden soll, dürfen zum „Öffnen“ der Datei bzw. des Datensatzes keine Komponenten mit Schreibmöglichkeit verwendet werden.

– **Bedienelemente sollen eine schnelle Arbeitsweise ermöglichen.**

Um sofort in einer bewährten Arbeitsumgebung beginnen zu können, sollte es möglich sein zu wählen, ob Fenstergröße und -position sowie die aktuellen Verzeichnisse der letzten Sitzung wiederverwendet werden oder nicht.

Parameterdateinamen, Bilddateiliste, Ausgabeverzeichnis usw. sollten auch bereits als Programmparameter beim Aufruf übergeben werden können, um (in Batchdateien bzw. Shell-Skripten oder Desktop-Icons usw.) passende Aufrufe für Sitzungen zu ermöglichen.

<sup>104</sup> Ein kompakter Überblick mit weiterführenden Quellenangaben findet man z. B. in RUDLOF 2006; Normen u. a.: EN ISO 9241

Für häufige Aktionen sollte außer einem Maus-Button wenn möglich immer auch eine Taste angeboten werden, möglichst in Anlehnung an allgemein übliche, z. B. für „Refresh“ so etwas wie [F5] oder [Ctrl]+[R], für „weiter“ so etwas wie [Tab] oder Cursorstasten usw.

Solange mit beiden Tastaturhälften gearbeitet werden muss, ist der Wechsel zur Maus und zurück immer ein Zeitverlust. Und umgekehrt: In Situationen, wo die Mausbedienung effektiver oder notwendig ist, sollte nach Möglichkeit nur eine Hand ergänzend auf der Tastatur nötig sein.

Bedienelemente sollen sowohl örtlich als auch hinsichtlich des Fokuswechsels in einer der Bearbeitung entsprechenden Reihenfolge angeordnet sein. Für absehbare Schrittfolgen sollte die [Enter]- oder [Tab]-Taste zum nächsten Schritt führen.

Listen von Einträgen müssen zum schnellen Auffinden eines Eintrags jederzeit in nachvollziehbaren Ordnungen sortierbar sein. Lange Listen sollten eine Mehrfachselektion ermöglichen, um Gruppen von Einträgen zusammen behandeln zu können.

Das erwünschte Speichern von kleinen Parameteränderungen soll nicht dadurch erschwert werden, dass unnötig viel anderes jeweils mitgespeichert werden muss – eine störend lange Speicherzeit verleitet zum Verzicht auf rechtzeitiges Sichern von Einstellungen und führt so leicht zum Verlust bereits erarbeiteter Konfigurationen.

- **Navigation einzuschränken hat nur Sinn, wenn dadurch nicht zeitaufwendige Umwege für absehbar notwendige Arbeitsgänge entstehen.**

Die Geometrie der Oberfläche sollte niemals dazu führen, sinnvolle Aktionen für unmöglich oder zu aufwendig zu halten.

- **Eine schnelle und fehlerfreie Navigation erfordert Anlehnung an die vom Betriebssystem gewohnten Standard-Abläufe.**

Dateiauswahldialoge sollten an einem nachvollziehbaren, aber bequemen Ort beginnen, z. B. im Verzeichnis der aktuell an dieser Stelle geltenden Auswahl – d. h. nicht stets in der Wurzel des Dateisystems und nicht in einem irgendwann früher benutzten, aber aktuell nicht geltenden Ordner.

In Speicherdialogen zu einer bereits geöffneten Datei wird ein editierbarer Vorschlag der Speicherung am bisherigen Ort unter bisherigem Dateinamen, so wie man das als Standardverhalten von Anwendungen kennt, erwartet.

Dialoge zur Dateiauswahl und -speicherung sollten stets in einer echten Combobox bereitgestellt werden, in der alternativ zum „Durchklicken“ der Pfad auch als kompletter String ins Textfeld kopiert werden kann.

Eine routinemäßige Speicherabfrage („Soll .. gespeichert werden?“) sollte nur dann angeboten werden, wenn tatsächlich Einstellungen geändert wurden (nur dann hätte die Abfrage auch eine sinnvolle Hinweisfunktion darauf, dass man in der Sitzung etwas geändert hat, und man hätte Anlass zu überlegen, ob versehentlich oder nicht).

## 7.5 Literatur- und Linkliste

Im Folgenden wird sowohl zitierte als auch weiterführende Literatur angegeben. Für systematischen bzw. weiterführenden Informationsbedarf wird auf die Literaturangaben in den genannten Publikationen verwiesen.

### Materialspezifik der Funeralschriften; Frakturschrift

#### **FRACTUR 1808**

Die Fraktur

in: Adelung, Johann Christoph: Grammatisch-kritisches Wörterbuch der hochdeutschen Mundart.

Wien: Pichler. Bd. 2. 1808. S. 261

#### **FRAKTUR 1954**

Fraktur

in: Der große Brockhaus. 16., völlig Neubearb. Aufl. Wiesbaden : Brockhaus. Bd. 4. 1954. S. 194

#### **HORN 1894**

Horn, E.: Zur Orthographie von U und V, I und J. Eine historisch typographische Erörterung

in: Zentralblatt für Bibliothekswesen 11(1894), S. 385-400

#### **JENSEN 1969**

Jensen, Hans: Die Schrift in Vergangenheit und Gegenwart. 3., Neubearb. u. erw. Aufl. Berlin : Dt. Verl. der Wiss., 1969. 607 S.

#### **KAPR 1959**

Kapr, Albert: Deutsche Schriftkunst. Versuch einer neuen historischen Darstellung. 2., verb. Aufl. - Dresden : Verl. der Kunst, 1959. 287 S.

#### **KAPR 1993**

Kapr, Albert: Fraktur. Form und Geschichte der gebrochenen Schriften. Mainz : Schmidt, 1993. 248 S.

#### **LENZ 1989**

Lenz, Rudolf[Hg.]: Leichenpredigten. Quellen zur Erforschung der Frühen Neuzeit. Marburg/Lahn: Forschungsstelle für Personalschriften, 1989. 20 S.

#### **KILLIUS 1999**

Killius, Christina: Die Antiqua-Fraktur Debatte um 1800 und ihre historische Herleitung. Wiesbaden : Harrassowitz, 1999. 488 S. (=Mainzer Studien zur Buchwissenschaft Bd. 7)

#### **REHSE 1998**

Rehse, Ernst-Günther: Gebrochene Schriften. Alphabete, Abkürzungen, Zeichen, Druckbeispiele, Fonts. Gotisch, Rundgotisch, Schwabacher, Fraktur, Fraktur-Varianten, Schreibschrift. Schaubuch, Nachschlagewerk und Hilfsbuch für den Umgang mit gebrochenen Schriften. Itzehoe : Verl. Beruf + Schule, 1998. 240 S.

#### **SCHNEIDER**

Schneider, Michael: Geschichte der deutschen Orthographie. Marburg. 30 S.

<http://decemsys.de/sonstig/gesch-rs.pdf> [Gesehen am 12.02.2013]

#### **TSCHICHOLD 1952**

Tschichold, Jan: Meisterbuch der Schrift. Ein Lehrbuch mit vorbildlichen Schriften aus Vergangenheit und Gegenwart für Schriftmaler, Graphiker, Bildhauer, Graveure, Lithographen, Verlagshersteller, Buchdrucker, Architekten und Kunstschulen. Ravensburg : Maier, 1952. 238 S.

#### **VEREIN FÜR COMPUTERGEALOGIE**

Alte Krankheitsbezeichnungen. Verein für Computergenealogie

<http://wiki-de.genealogy.net/Kategorie:Krankheitsbezeichnung> [Gesehen am 12.02.2013]

#### **VOESTE 2008**

Voeste, Anja: Orthographie und Innovation. Die Segmentierung des Wortes im 16. Jahrhundert. Hildesheim [u. a.] : Olms, 2008. VIII, 250 S.

[http://ebooks.ciando.com/book/index.cfm/bok\\_id/335655](http://ebooks.ciando.com/book/index.cfm/bok_id/335655) (Online-Publikation 2012)

#### **WALTHER 2006, 1**

Walther, Karl Klaus: Buchschmuck

in: Lexikon der Buchkunst und der Bibliophilie. Hamburg : Nikol, 2006. S. 122-130

#### **WALTHER 2006, 2**

Walther, Karl Klaus: Druckschrift

in: Lexikon der Buchkunst und der Bibliophilie. Hamburg : Nikol, 2006. S. 146-150

#### **WALTHER 2006, 3**

Walther, Karl Klaus: Fraktur

in: Lexikon der Buchkunst und der Bibliophilie. Hamburg : Nikol, 2006. S. 186-187

**WALTHER 2006, 4**

Walther, Karl Klaus: Personalschriften

in: Lexikon der Buchkunst und der Bibliophilie. Hamburg : Nikol, 2006. S. 287-293

**WITZEL 2003**

Witzel, Jörg: Der Tod kam im Examen. Zur Auswertung von Leichenpredigten in der Marburger Forschungsstelle für Personalschriften

in: Marburger Uni-Journal 15(2003), S. 40-43

<http://www.uni-marburg.de/aktuelles/unijournal/15/Personalschriften> [Gesehen am 12.02.2013]

**WITZEL 2010**

Witzel, Jörg: Texte für die Ewigkeit

in: Epoc 3(2010), S. 64-69

**ZIPPEL 2009**

Typographie der Pausen. 2009. 4 S.

[http://www.logo-sz.de/site/ESSAYS/Puenktchen\\_Puenktchen\\_Komma\\_Strich.pdf](http://www.logo-sz.de/site/ESSAYS/Puenktchen_Puenktchen_Komma_Strich.pdf) [Gesehen am 12.02.2013]

**OCR und Korrekturverfahren****ANDERSON 2010**

Anderson, Niall: Best Practice Guide: Optical Character Recognition. 2010

<http://www.impact-project.eu/uploads/media/IMPACT-ocr-bpg-pilot-s3.pdf> [Gesehen am 12.02.2013]

**BÜCHLER 2011**

Büchler, Marco: Textvervollständigung, OCR- und Rechtschreibkorrektur. Drei Sichten auf gleiche Methoden. Vortrag 12.10.2012, München

<http://www.slideshare.net/mdz-bsb/digitalisierungspraxis-bchler-textvervollstndigung> [Gesehen am 12.02.2013]

**FURRER 2011**

Furrer, L.; Volk, M.: Reducing OCR errors in Gothic-script documents

in: 8th International Conference on Recent Advances in Natural Language Processing (RANLP 2011), Hisar, 16 September 2011 – 16 September 2011, 97-103.

[http://www.zora.uzh.ch/49812/4/Furrer\\_Volk\\_2011\\_V.pdf](http://www.zora.uzh.ch/49812/4/Furrer_Volk_2011_V.pdf) [Gesehen am 12.02.2013]

**GEIERHOS 2007**

Geierhos, Michaela: Grammatik der Menschenbezeichner in biographischen Kontexten. Magisterarbeit, CIS München, gekürzte Fassung. 2007. 160 S.

[http://www.cis.uni-muenchen.de/download/cis\\_ais/ais-002.pdf](http://www.cis.uni-muenchen.de/download/cis_ais/ais-002.pdf) [Gesehen am 12.02.2013]

**HAUSER 2007**

Hauser, Andreas W.: OCR Postcorrection of Historical Texts. Schriftliche Hausarbeit im Fach Computerlinguistik, CIS München. 2007. 90 S.

<http://www.cip.ifi.lmu.de/~hauser/papers/histOCRNachkorrektur.pdf> [Gesehen am 12.02.2013]

**HÖHN 2006**

Höhn, Winfried: Mustererkennung in Frühdrucken. Diplomarbeit, Institut für Informatik, Lehrstuhl für Informatik II Würzburg. 2006. 56 S.

[http://opus.bibliothek.uni-wuerzburg.de/volltexte/2008/3042/pdf/Diplomarb\\_Hoehn.pdf](http://opus.bibliothek.uni-wuerzburg.de/volltexte/2008/3042/pdf/Diplomarb_Hoehn.pdf) [Gesehen am 12.02.2013]

**KÄMMERER 2009**

Kämmerer, Carmen: Vom Image zum Volltext – Möglichkeiten und Grenzen des Einsatzes von OCR beim alten Buch

in: Bibliotheksdienst 43(2009) S. 626-659

[http://www.zlb.de/aktivitaeten/bd\\_neu/heftinhalte2009/Technik010609BD.pdf](http://www.zlb.de/aktivitaeten/bd_neu/heftinhalte2009/Technik010609BD.pdf) [Gesehen am 12.02.2013]

**MÜHLBERGER 2011**

Mühlberger, Günter: Strukturanalyse auf der Basis von OCR Ergebnissen. Vortrag 11.10.2012, München.

<http://www.slideshare.net/impactproject/bsb-demo-day-mhlberger-dokumentstrukturanalyse> [Gesehen am 12.02.2013]

**NEUMANN 2008**

Neumann, Andreas: Konstruktion historischer Wörterbücher. München. 2008. 116 S.

[http://www.neumann.biz/data/computerlinguistik/Konstruktion\\_historischer\\_W%C3%B6rterb%C3%BCher.pdf](http://www.neumann.biz/data/computerlinguistik/Konstruktion_historischer_W%C3%B6rterb%C3%BCher.pdf)

[Gesehen am 12.02.2013]

**OCR-WORKSHOP 2011**

Historische Dokumente auf dem Weg zum digitalen Volltext. Workshop 11.-12. Oktober 2011, BSB München: Präsentationen

<http://www.slideshare.net/mdz-bsb> [Gesehen am 12.02.2013]

**OPITZ 2002**

Opitz, Andrea: Document Type Definitions zur Erschließung von Gattungen des Barock im Internet. Ein Projekt an der Herzog August Bibliothek Wolfenbüttel

In: Jahrbuch für Computerphilologie 5(2003) S. 55-64

<http://computerphilologie.uni-muenchen.de/jg03/opitz.html> [Gesehen am 12.02.2013]

**REFFLE 2011**

Reffle, Ulrich: Analyse und Nachkorrektur historischer und OCR- erfasster Ergebnisse. Vortrag 11.10.2012, München  
<http://www.slideshare.net/impactproject/bsb-demo-day-reffle-analyse-und-nachkorrektur> [Gesehen am 12.02.2013]

**RINGLSTETTER 2003**

Ringlstetter, Christoph: OCR-Korrektur und Bestimmung von Lebensstein-Gewichten. Magisterarbeit, CIS München. 2003. 130 S.

<http://www.cis.uni-muenchen.de/people/kristof/Publications/ringlstetter03.pdf> [Gesehen am 12.02.2013]

**RUDLOF 2006**

Rudlof, Christiane: Handbuch Software-Ergonomie. 2., überarb. Aufl. Tübingen : Unfallkasse Post und Telekom, 2006. 144 S.

<http://www.ukpt.de/pages/dateien/software-ergonomie.pdf> [Gesehen am 12.02.2013]

**SCHLARB 2011**

Schlarb, Sven: Entscheidungsfindung in der Digitalisierung durch experimentelle Workflow-Entwicklung. Vortrag, 11.10.2012, München

<http://www.slideshare.net/impactproject/bsb-demo-day-schlarb-workflowdesign> [Gesehen am 12.02.2013]

**SCHULZ 2006**

Schulz, Klaus U.: Nachkorrektur von Ergebnissen einer optischen Charaktererkennung. 2006. 108 S.

[http://www.cis.uni-muenchen.de/~uli/kurse/ws0708/hist\\_ocr/material/schulz\\_ocr.pdf](http://www.cis.uni-muenchen.de/~uli/kurse/ws0708/hist_ocr/material/schulz_ocr.pdf) [Gesehen am 12.02.2013]

**STÄCKER 2002**

Stäcker, Thomas: XML für Alte Drucke – welche Erschließungspotentiale bietet die neue Auszeichnungssprache? 2002. 7 S.

<http://www.hab.de/bibliothek/wdb/barocktdt/staecker.pdf> [Gesehen am 12.02.2013]

**STOCKMANN 2010**

Stockmann, Ralf: Was tun mit den Ergebnissen der OCR? Vortrag 04.03.2010

<http://www.slideshare.net/impactproject/stockmann-endnutzer-impact-workshop-muc-3344329> [Gesehen am 12.02.2013]

**STROHMAIER 2004**

Strohmaier, Christian M.: Methoden der lexikalischen Nachkorrektur OCR-erfasster Dokumente. Dissertation, LMU München: Fakultät für Sprach- und Literaturwissenschaften. 2004. 158 S.

[http://edoc.ub.uni-muenchen.de/3674/1/Strohmaier\\_Christian.pdf](http://edoc.ub.uni-muenchen.de/3674/1/Strohmaier_Christian.pdf) [Gesehen am 12.02.2013]

**WIENERS 2008**

Wieners, Jan Gerrit: Zur Erweiterungsfähigkeit bestehender OCR Verfahren auf den Bereich extrem früher Drucke. Magisterarbeit, Historisch-Kulturwiss. Informationsverarbeitung. Köln, 2008. 81 S.

[http://www.hki.uni-koeln.de/files/MA\\_wieners.pdf](http://www.hki.uni-koeln.de/files/MA_wieners.pdf) [Gesehen am 12.02.2013]

**Projekte und Web-Adressen****BERLINER FUNERALSCHRIFTEN**

Personale Gelegenheitsschriften. Historische Drucke

<http://staatsbibliothek-berlin.de/die-staatsbibliothek/abteilungen/historische-drucke/sammlungen/bestaende/personale-gelegenheitsschriften/> [Gesehen am 12.02.2013]

Pilotprojekt zum OCR-Einsatz bei der Digitalisierung der Funeralschriften der Staatsbibliothek zu Berlin

<http://staatsbibliothek-berlin.de/die-staatsbibliothek/abteilungen/historische-drucke/aufgaben-profil/projekte/funeralschriften/> [Gesehen am 12.02.2013]

Federbusch, Maria: Praxistest zweier OCR-Softwareprodukte am Beispiel ausgewählter Funeralschriftenbestände der SBB.

Vortrag 12.10.2011, München

<http://www.slideshare.net/mdz-bsb/digitalisierungspraxis-federbusch-ocrpraxistest> [Gesehen am 12.02.2013]

Metadaten zu den Berliner Funeralschriften. 2008

[http://staatsbibliothek-berlin.de/fileadmin/user\\_upload/zentrale\\_Seiten/historische\\_drucke/pdf/Metadaten\\_Funeralschriften.pdf](http://staatsbibliothek-berlin.de/fileadmin/user_upload/zentrale_Seiten/historische_drucke/pdf/Metadaten_Funeralschriften.pdf) [Gesehen am 12.02.2013]

**DEUTSCHES TEXTARCHIV / Volltextfassung schlesischer Leichenpredigten**

<http://www.deutschestextarchiv.de/> [Gesehen am 12.02.2013]

**HELMSTEDTER DRUCKE ONLINE**

Helmstedter Drucke Online

Projekt der HAB Wolfenbüttel. Volltextgenerierung durch OCR mit BIT-Alpha

<http://www.hab.de/de/home/wissenschaft/projekte/helmstedter-drucke-online.html> [Gesehen am 12.02.2013]

**IMPACT**

Impact. Centre of competence

<http://www.digitisation.eu/> [Gesehen am 12.02.2013]

Improving Access to Text

<http://www.impact-project.eu/> [Gesehen am 12.02.2013]

Tools and applications

neue Seite: <http://www.digitisation.eu/tools/> [Gesehen am 12.02.2013]

alte Seite: <http://www.impact-project.eu/taa/tech/> [Gesehen am 12.02.2013]

Anderson, Niall: Best Practice Guide: Optical Character Recognition

<http://www.impact-project.eu/uploads/media/IMPACT-ocr-bpg-pilot-s1.pdf> bis ...s3 [Gesehen am 12.02.2013]

Impact Day / Impact Conference: Presentations. Dec. 2009-Dec. 2011

<http://www.slideshare.net/impactproject> [Gesehen am 12.02.2013]

**MARBURGER PERSONALSCHRIFTENSTELLE**

Forschungsstelle für Personalschriften. Akademie der Wissenschaften und der Literatur Mainz

<http://www.personalschriften.de/> [Gesehen am 12.02.2013]

THELO

Thesaurus Locorum

<http://www.personalschriften.de/datenbanken/thelo.html> [Gesehen am 12.02.2013]

THEPRO

Thesaurus Professionum

<http://www.personalschriften.de/datenbanken/thepro.html> [Gesehen am 12.02.2013]

**NACHKORREKTUR VON OCR-ERGEBNISSEN**

Projekt: Domänen- und dokumentenadaptive Verfahren zur Nachkorrektur von OCR-Ergebnissen. CIS München. 2003-2010

Eintrag 5419670 in der DFG-Datenbank GEPRIS

<http://gepris.dfg.de/gepris/OCTOPUS/?module=gepris&task=showDetail&context=projekt&id=5419670> [Gesehen am 12.02.2013]

**8. Erfahrungsbericht Helmstedter Drucke  
Online an der Herzog August Bibliothek  
Wolfenbüttel**

*Thomas Stäcker*

## 8.1 Zusammenfassung

Ziel des Projekts ist die Image- und Volltextdigitalisierung der in Helmstedt gedruckten Werke, die sich in der Herzog August Bibliothek befinden. Die Herzog August Bibliothek besitzt etwa 10.000 Drucke mit dem Druckort Helmstedt, die bis zur Schließung der Universität Helmstedt 1810 in Helmstedt gedruckt worden sind. Die auf das 16. und 17. Jh. entfallenden Drucke sollen als Image und im Volltext per OCR digitalisiert, mit Strukturdaten zur Navigation versehen und im Verbundkatalog sowie im VD 16 und VD 17 nachgewiesen werden. Der aus dem 18. Jh. stammende Anteil wird digitalisiert, auf nationalbibliographischem Niveau katalogisiert und zur Nachnutzung einem möglichen VD 18 zur Verfügung gestellt. Mit der nahezu kompletten Digitalisierung der Produktion eines der wichtigsten norddeutschen Universitätsdruckorte – Helmstedt findet sich unter den 10 am häufigsten nachgewiesenen Druckorten im VD 17 – wird nicht nur ein substantieller Beitrag zur Komplettdigitalisierung des deutschen gedruckten Kulturerbes geleistet, sondern erstmals ein Überblick über die Druckproduktion einer bedeutenden frühneuzeitlichen Universität geschaffen. Das Projekt flankiert den programmatischen und auf mehrere Jahre angelegten Forschungsschwerpunkt der HAB zur Erforschung der Universitätsgeschichte Helmstedts und kommt unmittelbar mehreren Forschungsprojekten in diesem Feld zugute.

## 8.2 Bedeutung der Helmstedter Drucke

Das Vorhaben der Digitalisierung der Helmstedter Drucke steht im Kontext der Erforschung der Universitäts- und Wissenschaftsgeschichte der Frühen Neuzeit.<sup>105</sup> In ihrer Selbstbeschreibung formulierte jede frühneuzeitliche Universität einerseits ihre Ziele und Ansprüche, andererseits deren Einlösung. Da der Status universitären Wissens von spezifischen Geltungsgründen und der Einhaltung bestimmter Verfahren abhing, war eine Hochschule zu intensiver Selbstbeobachtung gezwungen. Deren Ergebnisse bilden einen wichtigen zeitgenössischen Beurteilungsmaßstab für die Qualität und Leistung einer Universität und insofern einen geeigneten Ausgangspunkt für ein retrospektives *reviewing*, zumal sie aufgrund der kompetitiven Grundstruktur des Wissenschafts- und Universitätssystems im Rahmen der fürstenstaatlichen Ordnung und der gelehrten Welt kommuniziert werden mussten. Diese als universitäre Repräsentation und Selbstvergewisserung zu bezeichnende Aufgabe bündelt sich wie in einem Brennglas in der gesamten Druckproduktion der privilegierten Universitätsbuchdruckerei. Das Spezifikum der Druckproduktion einer Universitätsdruckerei<sup>106</sup> ist, dass sie an einer fürstenstaatlich initiierten, finanzierten und kontrollierten Einrichtung stattfand und insofern politisch wie konfessionell autorisiert, legitimiert und funktionalisiert wurde. Die Druckproduktion der Universität kann damit in besonderer Weise über universitäre und staatliche Autorisierungsmechanismen Aufschluss geben. Sie manifestiert die Paradigmen akademischen Wissens in den verschiedenen Fakultäten ebenso wie die verfolgten Lehrpraktiken, denn frühneuzeitliche Universitäten waren noch weit bis ins 18. Jh. hinein in noch stärkerem Maße als heute Unterrichts- und Ausbildungsanstalten.

Wurden die Hochschulschriften, also die akademischen Dissertationen, Reden und Programme, lange Zeit als eine eher uninteressante Quellengruppe eingestuft, hat sich diese Wertung seit den 1980er Jahren geradezu umgekehrt, wobei insbesondere die Dissertationen des 16. bis 18. Jahrhunderts mittlerweile als „ungehobene Schätze“ betrachtet werden.<sup>107</sup> Gerade an diesen universitären Gelegenheitsschriften lässt sich ablesen, wie die eigentliche „Wissensproduktion“ an den Universitäten aussah, womit ein wichtiger Beitrag zur Alltagsgeschichte der Bildung und Ausbildung in der Frühen Neuzeit geleistet

<sup>105</sup> Zur Universitäts- und Bildungsgeschichte der Frühen Neuzeit im Allgemeinen vgl. insb.: Paulsen, Friedrich: Geschichte des gelehrten Unterrichts auf den deutschen Schulen und Universitäten vom Ausgang des Mittelalters bis zur Gegenwart. Mit besonderer Rücksicht auf den klassischen Unterricht, 2 Bde., Leipzig / Berlin 1919 / 1921 (ND 1965); Hammerstein, Notker (Hg.): Handbuch der deutschen Bildungsgeschichte, Bd. I: 15. bis 17. Jahrhundert. Von der Renaissance bis zum Ende der Glaubenskämpfe, München 1996; Rüegg, Walter (Hg.): Geschichte der Universität in Europa, Bd. II: Von der Reformation bis zur Französischen Revolution (1500 – 1800), München 1996; Schindling, Anton: Bildung und Wissenschaft in der Frühen Neuzeit 1650 – 1800, 2. Aufl., München 1999; Hammerstein, Notker: Bildung und Wissenschaft vom 15. bis zum 17. Jahrhundert, München 2003; Bruning, Jens: Innovation in Forschung und Lehre : die Philosophische Fakultät der Universität Helmstedt in der Frühaufklärung 1680 – 1740. Wiesbaden 2012 (Wolfenbütteler Forschungen, 132). – Eine Auswahlbibliographie zum Standort Helmstedt steht über das Helmstedt-Portal zur Verfügung: [http://uni-helmstedt.hab.de/docs/bibliographie\\_helmstedt.pdf](http://uni-helmstedt.hab.de/docs/bibliographie_helmstedt.pdf) [Stand: 02/2013]

<sup>106</sup> Zu dieser: Eule, Wilhelm: Helmstedter Universitätsbuchdrucker. Helmstedt 1921. S.a. Benzing, Josef: Die Buchdrucker des 16. und 17. Jahrhunderts im deutschen Sprachgebiet. 2. verb. Aufl. Wiesbaden 1982, S. 200-203 mit weiterer Literatur zu den einzelnen Druckern.

<sup>107</sup> Eule (wie Anm. 2), S. 29ff.

werden kann. Zu betonen ist dabei, dass den akademischen Druckerzeugnissen nicht nur als Einzelpublikationen, sondern auch als „Massenquelle“ bzw. „serielle“ Quelle besondere Bedeutung zukommt und sie gerade in diesem Charakterzug einen einzigartigen Blick in die Entwicklung der verschiedenen Wissenschaftsdisziplinen ermöglichen.

### 8.3 Umfang und Art der Helmstedter Druckproduktion

Die wichtigsten Drucker des 16. und 17. Jh. in Helmstedt waren Jakob Lucius I (1579-1597), Jakob Lucius II (1598-1616, Erben bis 1633 vermutlich unter der Leitung von Henning Müller d. Ä.), Georg Calixt (1629-1676 mit Pächtern, insb. Henning Müller d. J.), Jacob Lucius III (1634-1639), Henning Müller d. J. (1640-1674), Johann Heitmüller (1656-1677), Jakob Müller (1661-1681), Georg Wolfgang Hamm (1681-1714), Heinrich Hesse (1681-1715/16) und Salomon Schnorr (1690-1723). Für das 18. Jh. sind insbesondere Johann Drimborn und Johann Heinrich Kühnlin zu nennen. Bemerkenswert ist die Betätigung des berühmten Gelehrten und Theologen Georg Calixt als Druckherr. Hier zeigt sich die enge Verschränkung von Universität und Druckerei in besonders augenfälliger Weise.<sup>108</sup> Besonders wichtig war für die Helmstedter Druckgeschichte Henning Müller d. J. Seine Produktion zeigt denn auch das Spektrum der Produktion auf exemplarische Weise. Von den 2.422 im VD 17 nachgewiesenen Drucken aus seiner Offizin entfallen mehr als die Hälfte, nämlich 1.342 auf Dissertationen und 276 auf Hochschulschriften eine bis dato vernachlässigte Quellengruppe.<sup>109</sup> Dazu kommen 91 Leichenpredigten und 590 sonstige Gelegenheitsschriften, die in besonderer Weise dazu geeignet sind, biographische Hintergründe des akademischen Personals aufzuhellen. Rund 500 Werke entfallen auf sonstige Werke und wissenschaftliche Abhandlungen, insbesondere von Professoren der Universität Georg Calixt, Johannes Caselius, Hermann Conring, Hermann v. d. Hardt, Heinrich Meibom d. Ä. und d. J., Johann Lorenz von Mosheim, Heinrich Rixner, Heinrich Julius Scheuerl, Valentin Heinrich Vogler und Giordano Bruno (ein Jahr in Helmstedt, nicht als Professor, hielt aber vermutlich Privatvorlesungen), deren Arbeiten von der hohen Qualität der Forschung in Helmstedt in der Frühen Neuzeit Zeugnis ablegen.

Das im engeren Sinne akademische Schrifttum (Dissertationen, Abhandlungen) wurde, wie üblich in dieser Zeit, fast ausnahmslos in Latein verfasst. Allerdings finden sich in den Drucken immer wieder deutsche Einsprengsel und Zitate. Daneben druckten die Universitätsdrucker vor allem Predigten oder Leichenpredigten in Deutsch. Lateinische Texte wurden in Antiqua, Deutsche in Fraktur gesetzt, wobei die Grenze mitunter im Wort selbst verlief. So finden sich Fälle, wo die Endung „e“ des Worts „Materie“ in Fraktur, der Rest „Materi“ in Antiqua gesetzt wurde, weil es sich um ein lateinisches Wort handelt, das eigentlich die Endung „a“ trägt. Typischerweise enthalten die Drucke Mischformen von Antiqua, Fraktur und Kursive als Brotschriften. Hinzu kommen barocke Auszeichnungsschriften, die die Drucker aber für das in der Masse wenig aufwendig gestaltete Universitätsschrifttum vergleichsweise zurückhaltend einsetzten. Die verwendeten Fonts weisen wenig Varianten auf und scheinen über die Jahre weitgehend stabil gewesen zu sein, da Fonts von Universitätsdrucker zu Universitätsdrucker vererbt und der Erwerb neuer Fonts aus ökonomischen Gründen vermieden wurde.

Obwohl es zur Frage des deutschen bzw. europäischen Schriftenhandels beklagenswert wenige wissenschaftliche Untersuchungen gibt, deutet doch das wenige, was wir haben, darauf hin, dass man auch unabhängig von diesem innerhelmstedtischen Befund von weitgehender Einheitlichkeit bei den Schriften ausgehen kann und dass gerade im 17. Jh. ein starker Zentralisierungsprozess im Handel mit Schriften eingetreten ist.<sup>110</sup> Hervorzuheben ist hier vor allem die erfolgreiche Tätigkeit der Egenolffschen/Lutherschen Schriftgießerei in Frankfurt. Insofern kann man von einem hohen Homogenitätsgrad ausgehen, was für die OCR-Konversion der Barockschriften nicht ohne Belang ist.

<sup>109</sup> Vgl. Komorowski, Manfred: Die Hochschulschriften des 17. Jahrhunderts und ihre bibliographische Erfassung. In *Wolfenbütteler Barock-Nachrichten* 24 (1997) 1, S. 19-42.

<sup>110</sup> Vgl. hierzu Carter, Harry: *A View of Early Typography up to about 1600*. Reprinted. London 2002.

## 8.4 Projektziele, Rahmenbedingungen

Ziel des Projekts ist die Image- und Volltextdigitalisierung sowie Erschließung der in Helmstedt gedruckten Ausgaben, die vom Beginn des Universitätsdrucks bis zur Schließung der Universität 1810 in Helmstedt gedruckt wurden und sich heute in der Herzog August Bibliothek befinden.

Die Wolfenbütteler Sammlung bietet optimale Ausgangsbedingungen für ein Digitalisierungsprojekt der Helmstedter Druckproduktion. Die Helmstedter Universitätsdrucker waren verpflichtet, Exemplare ihrer gesamten Produktion an die eigene Universitätsbibliothek und die fürstliche Bibliothek in Wolfenbüttel zu überweisen.<sup>111</sup> Da beide Sammlungen, die Wolfenbütteler und die Helmstedter nach Schließung der Universität Helmstedt wieder in Wolfenbüttel vereint wurden, kann von einer sehr dichten Überlieferungslage vor Ort ausgegangen werden. Da diese mit Ausnahme der in Helmstedt verbliebenen Drucke<sup>112</sup> und der Drucke des 18. Jh. in den bereits abgeschlossenen nationalbibliographischen Verzeichnissen des VD 16<sup>113</sup> und VD 17<sup>114</sup> durch die HAB Wolfenbüttel katalogisiert wurden, liegen weitgehend alle erforderlichen Metadaten für die Digitalisierung vor. Im VD 16 sind 1.216, im VD 17 sind 7.485 (Stand: 3.10.2009) Drucke für den Druckort Helmstedt nachgewiesen. Von diesen besitzt die HAB gerundet aus dem 16. 789, aus dem 17. 5.366 und aus dem 18./19. Jh.(-1810) 4.390 Ausgaben, also insgesamt rund 10.000 Ausgaben mit geschätzt 1 Mio. Seiten.

Ein Novum für ein Digitalisierungsprojekt zur Frühen Neuzeit war, neben der eigentlichen Image-digitalisierung auch eine OCR-Konversion des Schriftguts in Angriff zu nehmen. Allerdings wurde die ursprünglich intendierte vollständige Bearbeitung des Korpus auf Wunsch der DFG auf ein kleineres Sample von 120.000 Seiten reduziert, um angesichts der erheblichen Kosten zunächst Erfahrungen mit dem Verfahren zu sammeln.

Bei der anvisierten Komplettdigitalisierung der Druckerzeugnisse mit dem Druckort Helmstedt mussten bereits laufende Digitalisierungsprojekte, insbesondere an der Bayerischen Staatsbibliothek München, aber auch andersorts berücksichtigt werden, um Doppeldigitalisierungen so weit wie möglich zu vermeiden. Gerade mit Blick auf die OCR-Bearbeitung schien das Google-Projekt der Bayerischen Staatsbibliothek zunächst von Interesse. Eine genauere Vorprüfung zeigte aber, dass die Google Digitalisate nicht für das Projekt nutzbar waren. Zum einen ist die von Google im Internet bereitgestellte Qualität zu schlecht für eine nachgeordnete OCR-Bearbeitung (Gerüchten zufolge besitzt Google intern bessere Scans, was aber für aktuell durchgeführte Projekte leider wenig nützt), zum anderen ist es nach den Geschäftsbedingungen von Google<sup>115</sup> nicht gestattet, Volltexte außerhalb von Google zu benutzen, was eklatant Forschungsinteressen bezüglich der freien Nutzung der Volltexte widerspricht, z. B. in dem Wunsch nach unabhängiger Eigenindexierung und Bildung von reproduzierbaren Rankingmechanismen.<sup>116</sup> Insofern scheint nach dem gegenwärtigen Stand ein auf Volltextgenerierung angelegtes wissenschaftliches Digitalisierungsprojekt grundsätzlich nur mit Digitalisaten aus öffentlicher Hand möglich zu sein, zumindest sofern entsprechende CC-Rechte am Volltext eingeräumt werden, was vor dem Hintergrund einer DFG-Förderung typischerweise der Fall sein dürfte.

<sup>111</sup> Vgl. Schneider, Heinrich: Beiträge zur Geschichte der Universitätsbibliothek Helmstedt. Helmstedt 1924, S. 75, Anm. 194.

<sup>112</sup> Diese sind in einem eigenständigen Projekt „Katalogisierung der im Juleum Helmstedt verbliebenen Drucke der ehemaligen Universitätsbibliothek Helmstedt“ derzeit in Bearbeitung: <http://www.hab.de/de/home/wissenschaft/projekte/katalogisierung-der-im-juleum-helmstedt-verbliebenen-drucke-der-ehemaligen-universitaetsbibliothek-helmstedt.html> [Stand: 02/2013]

<sup>113</sup> <http://www.vd16.de> [Stand: 02/2013]

<sup>114</sup> <http://www.vd17.de> [Stand: 02/2013]

<sup>115</sup> <https://www.google.de/intl/de/policies/terms/regional.html>; vgl. a. <https://play.google.com/intl/de/about/books-terms.html> [Stand: 02/2013]

<sup>116</sup> Zum Problem der mangelnden Wissenschaftstauglichkeit von Googles Suchmaschine vgl. Elli Pariser: The filter bubble : what the Internet is hiding from you. London 2011.

Die Volltexte sollten nach vorausgegangenen Tests mit der Software BIT-Alpha der Fa. B.I.T. Tomasi generiert werden. Wie oben geschildert, besteht das Charakteristikum der Software darin, dass sie eine effiziente Oberfläche zum schnellen Erlernen von Schriften bietet. Hierzu wurde eine Hilfskraft eingestellt, die Schriften von Helmstedter Druckern trainierte. Dabei handelte es sich überwiegend um Antiquaschriften inklusive Kursive, zu einem kleineren Teil um Fraktur. Anschließend sollten die Drucke bearbeitet und die Ergebnisse in ALTO<sup>117</sup> exportiert und zur Nachnutzung in der Wolfenbütteler Digitalen Bibliothek in ein TEI P 5 kompatibles Format konvertiert werden. Die Suche in den Texten erfolgte mit eXist, einer nativen XML-Datenbank.

## 8.5 Ergebnisse, Erfahrungen

Nach der Digitalisierung von 4.800 Drucken (Stand 01.06.2012) hat sich ein Seitendurchschnitt von 87 Seiten pro Ausgabe ergeben. Die gegenüber sonstigen Projekten<sup>118</sup> eher geringe durchschnittliche Seitenzahl erklärt sich daraus, dass es sich bei Universitätschriften meist um weniger umfangreiche Druckerzeugnisse wie Dissertationen und Gelegenheitschriften handelt. Für die Archivkopie fallen derzeit pro Image durchschnittlich 28,5 MB pro Seite (uncompressed TIFF) an.

Ziel der ersten Phase des Projekts war es, 120.000 Seiten per OCR zu bearbeiten. Die OCR-Komponente des Antrags war ein wichtiger neuer und innovativer Schritt im Prozess der Massendigitalisierung, weil bislang bei Materialien des 16.-18. Jh. fast ausschließlich die Imagedigitalisierung zum Einsatz kam. Wegen der dargelegten Schwierigkeiten der Vorlagen kam nach einigen Tests mit anderen Anbietern die OCR-Software der Firma B.I.T. Bureau Ingénieur Tomasi zum Einsatz. Das Spezifikum dieser Software liegt darin, dass sie nicht, wie z. B. Abbyy, auf weitgehend werkseitig vortrainierte Fonts angewiesen ist, so dass die typischen Mischtexte dieser Zeit (Antiqua, Fraktur, Kursive, Griechisch und andere Typen) nicht korrekt erkannt werden, sondern z. B. entweder nur der Antiqua- oder nur der Frakturanteil. Einen typischen Fall solcher Mischung zeigt sich in Abb. 77.

Die jeweiligen Typen mussten zunächst trainiert werden; dies konnte schriftübergreifend erfolgen. Effekte des „Übertrainierens“, wie man sie wohl bei anderen Produkten wie Abbyy beobachten kann, traten nicht auf. Das Training konzentrierte sich zunächst auf das 17. Jh. und den Helmstedter Großdrucker Henning Müller d. J., der allein für rund 3.000 Druckerzeugnisse verantwortlich ist. Damit verband sich die Absicht, im weiteren Verlauf auf der Basis eines weitgehend homogenen Schriftmaterials arbeiten zu können. Trainiert werden mussten Antiqua und Fraktur, ebenso Kursivschrift und altgriechische Schrift. Letztere wurde allerdings aus Kapazitätsgründen bisher nur in Ansätzen trainiert, weil dafür spezialisierte Kompetenzen erforderlich sind. Sie soll im Laufe eines Fortsetzungsprojekts durch Mitarbeiter der Bibliothek weiter trainiert und die Erkennungsgenauigkeiten für frühneuzeitliche griechische Fonts verbessert werden. Die Trainingsphase verlief weitgehend ohne Probleme; die Software erwies sich als sehr komfortabel in der Handhabung. Besonders effizient war der so genannte „Sequencer“, der es erlaubte typische Buchstabenkombinationen zu identifizieren (s. Kap. 4.4), typische Fälle sind etwa „n+i“ Verbindungen, die als „m“ zu lesen sind. Trainiert wurden von der studentischen Hilfskraft im Laufe der ersten Projektphase (24 Monate) 132.000 Zeichen mit einem Aufwand von durchschnittlich fünf Wochenstunden.

<sup>117</sup> <http://www.loc.gov/standards/alto/> [Stand: 02/2013]

<sup>118</sup> Durchschnittszahlen auf dem Server der HAB: 140 Seiten bei 14.000 Werken. Bei Dünnhaupt Digital sind es im Schnitt 235 Seiten.

putationi includere, & eruditorum placido examini publici exercitij gratia subjicere est animus: & quidem res præcipuas tam olim quam hodiernum in Comitiiis tractari consuetas, in quantum à me ex historiis peti potuerunt, thesibus hisce inserere conabor. Quod ut feliciter succedat, tu Deus largire vires, & fac ut omnia in nominis tui gloriam cedant.

VIII. Res vero quarum causa Comitiiia haberi solent, pro varietate negotiorum Imperio incumbentium, uti olim ita & hodie diversæ sunt. Omnia nimirum illa quæ ad salutem vel Statum Regis & regni pertinebant, in Comitiiis Carolinorum Regum ævo tractabantur. *Hincm. Epist. 3. c. 33.* Hodieque omnia negotia quæ ad Imperium vel conservandum vel defendendum pertinent, in Comitiiis tractari solent *R. A. de anno 1512.* Et ut paucis me absolvam, quicquid est, quod ad gloriam Dei, conservacionem justitiæ, honestatis & publici status ædificationem conferre possit, de eo in conventibus publicis agitur. *Reformation guter Policey zu Augßburg anno 1530. §. 1.* *Sampß ihnen alles das fürzunehmen zu rahschlagen / zu handeln und zu schliessen / daß zu fürderst Gotte dem Allmächtigen zu Ehr und Lob / gemeiner Christenheit und Teutscher Nation zu Wolsfahre Fried und Einigkeit / auch dem H. Röm. Reich zu nutz auffnehmen und gedeyen gereichen möcht. R. A. des Reichstags zu Regensßburg de anno 1541. §. 1.* Alles anders zu handeln und zu schliessen / daß dem H. Reich / und desselben Ständen zu wolsfahre nutz und gutem reichen mög. Ac in novissima Pacis constitutione ea de re ita cavetur. *Gaudeant sine contradictione jure suffragii in omnibus deliberationibus super negotiis Imperii, præsertim ubi leges ferendæ, vel interpretandæ, bellum decernendum, tributa indicenda, delectus aut hospitaciones militum instituendæ.*

A 3

tuendæ

Die Konversionsphase gestaltete sich anders als die Trainingsphase steiniger, was im Wesentlichen auf zwei Gründe zurückzuführen war. Zum einen erwiesen sich Parametrierung und Einstellungen der Software als zu komplex, zum anderen gab es zahlreiche Probleme mit den Vorlagen. Was Ersteres anlangte, so konnte in der fast zweijährigen Pilotphase trotz intensiver Bemühungen nicht nachvollzogen werden, wie sich welche Einstellungen auf das Ergebnis auswirkten (vgl. Kap. 5 u. 6.3 sowie Anh. 7.3). Grund war einerseits das Fehlen einer Dokumentation, die im Detail verdeutlicht hätte, wie sich die vielfältigen Parametrierungsmöglichkeiten auswirken. Andererseits entwickelte die Firma die Software laufend fort, was zwar zu Verbesserungen in der Erkennungsqualität führte, jedoch auch Änderungen bei den Parametereinstellungen zur Folge hatte. Insgesamt war es so weitgehend unmöglich, empirisch reproduzierbare Effekte zu erzielen, da Änderungen an einem Parameter dazu führten, dass andere Einstellungen unwirksam wurden, sich Einstellungen gegenseitig aufhoben oder Einstellungen in einer älteren Version sich nicht mehr so verhielten wie Einstellungen einer neueren Version. Bemühungen, die Software zu „verstehen“ wurden daher nach einer gewissen Zeit abgebrochen und es weitgehend der Firma überlassen, nach Maßgabe der jeweiligen Vorlagen die richtigen Einstellungen festzulegen. Die Zusammenarbeit in dieser Hinsicht verlief jedoch reibungslos und die Firma zeigte sich sehr kooperativ und entgegenkommend. Dass materialabhängig Anpassungen überhaupt sinnvoll sind, hängt mit den zugrundeliegenden Softwarealgorithmen zusammen. Marktübliche Software wie Abbyy oder Omnipage gehen von „durchschnittlichen“ Einstellungen aus, was zwar dazu führt, dass nicht das Optimum aus der Software herausgeholt wird, was es aber dem Nutzer erlaubt, ohne intime Kenntnisse der Einstellungsmöglichkeiten doch noch befriedigende Ergebnisse zu erzielen. Letztlich führten diese Umstände im Projekt zu der Einsicht, dass die Software von B.I.T. für den durchschnittlichen bzw. selbst für den erfahrenen Anwender mit Ausnahme des hervorragenden Trainingsinterfaces eigentlich nicht geeignet ist und legte für das geplante Folgeprojekt einen Wechsel im Vorgehen nahe. Abgesehen vom Training einiger Sonderschriften (vor allem Griechisch), das mangels entsprechender Kompetenzen der Firma in der Bibliothek erfolgen muss, soll die OCR-Wandlung künftig vollständig durch die Firma in einem Dienstleistungsverhältnis durchgeführt werden.

Trotz dieser Einschränkungen konnten aber Stärken und Schwächen der Software zumindest im Allgemeinen bewertet werden. Gemäß den Prozessschritten Binarisierung, Segmentierung bzw. Layoutanalyse, *pattern matching* und *intelligent character recognition* (Wörterbuch) kann man hervorheben, dass die Stärken der Software gerade im Bereich der Binarisierung und des *pattern matching* hervorstechen; einzelne Probleme, wie stark gekrümmte Seiten und die damit verbundene problematische Grundlinienfindung, konnten mittels neuprogrammierter Algorithmen gelöst werden. Bei der Segmentierung gab es verschiedentlich noch Probleme, die teils im Laufe des Projekts behoben werden konnten (z. B. bei der Identifizierung von Initialen), teils noch bestehen, wie bei der Erkennung von Marginalien oder Zeilenzählungen, die zu dicht am Text gedruckt sind. Hier scheinen weitere Verbesserungen möglich, zumal Abbyy bessere Segmentierungsergebnisse zu erzielen vermag. Auch bei der Wörterbuchunterstützung sind noch Verbesserungen möglich. Dies liegt allerdings vor allem daran, dass geeignete Wörterbücher für das Frühneuhochdeutsche oder auch Lateinische fehlen, die das typische universitäre Vokabular der frühen Neuzeit abzubilden vermöchten. Die Ergebnisse der durch Wörterbuch unterstützten Konversion wirkten sich in der ersten Phase des Projekts noch nicht aus bzw. es können dazu noch keine verallgemeinerbaren Aussagen getroffen werden. Die noch nicht repräsentativen Ergebnisse waren einerseits besser, führten andererseits aber auch zu neuen Fehlern durch falsche Korrekturen. Andererseits gibt es eine Reihe von Phänomenen im Text, die nur mit Hilfe von Wörterbüchern gelöst werden können und bei denen reines *pattern matching* nicht weiterführt. Dies ist z. B. der Fall, wenn kleines „c“ und „e“, da unsauber gedruckt, nicht visuell unterschieden oder wo innerhalb eines Worts „i+n“ im Sequencer nur durch ein Wörterbuch disambiguiert werden können, weil beides „in“ und „m“ möglich ist. Der Weg, ein effizientes *pattern matching* mit einem Wörterbuch zu kombinieren, ist daher zwingend, weil eine rein „visuelle“ Erkennung bestimmte Grenzen nicht zu überschreiten vermag.

Neben den Schwierigkeiten mit der Parametrierung der Software bestanden aber auch Probleme, die aus der Qualität der Vorlagen resultierten. Die ursprüngliche Ansicht, man müsse die Software nur einen bestimmten Schrifttyp lernen lassen, um zu guten Ergebnissen zu gelangen, ist zwar im Grundsatz richtig, doch zeigte sich, dass die Schriftvarianten, auf die man bei verschiedenen Druckern trifft, sich weniger negativ auswirkten als mindere Papier- und Druckqualität, Widerdruck oder sonstige Verschmutzungen. Unterstreichungen, die oft als Durchstreichungen ausgeführt sind, machten eine Texterkennung so gut wie unmöglich.

Trotz dieser Probleme sind aber die bisherigen Ergebnisse so gut, dass eine Fortsetzung mit größeren Margen und Komplettkonversion dieses Korpus und anderer Korpora sinnvoll scheint, zumal im Fall der Helmstedter Drucke, bei denen es sich überwiegend um serielles Material handelt und der eigentliche Nutzen sich unstreitig erst aus dem Volltext ergibt.

Um aber hierzu nachvollziehbare Aussagen machen zu können, müssen zunächst zwei Gesichtspunkte näher erörtert werden. Zum einen die Frage, welche Genauigkeit der Texterkennung erforderlich ist, zum anderen wie man diese Genauigkeiten misst (8.7). Erst aus der Beantwortung dieser Fragen lässt sich auch eine Vorstellung über den nötigen und angemessenen ökonomischen Aufwand für ein OCR-Projekt gewinnen (8.8).

## 8.6 Überlegungen zu Textgenauigkeiten und deren Messung

Im allgemeinen herrschen durchaus disparate Vorstellungen davon, was ein qualitativ ausreichender Text ist. An der Bayerischen Staatsbibliothek (BSB) z. B. teilt man Texte in vier Klassen ein: sehr gut (99,6-99,95% Genauigkeit), gut (97-99,5%), durchschnittlich (90-96%) und schlecht (unter 90%).<sup>119</sup> Gewöhnlich wird argumentiert – so auch in den bisherigen DFG-Praxisregeln (gültig bis Ende 2012) –, dass eine in diesem Sinne sehr gute Qualität Voraussetzung für eine wissenschaftliche Nutzung sei. Dem ist nach traditionellem Verständnis durchaus zuzustimmen. Gleichwohl stehen dem Überlegungen entgegen, wie sie von Thaller<sup>120</sup> und anderen angestellt werden, dass auch eine weniger exakte, nach den bisherigen Vorstellungen eher unzulängliche OCR-Konversion noch wissenschaftlichen Nutzen bringt, selbst wenn *argumenta ex silentio* nicht mehr möglich sein sollten. Die heutige Technik erlaubt Mengen von Texten automatisch zu bearbeiten, die im Laufe eines Gelehrtenlebens nicht hätten gelesen werden können, so dass man mit geeigneten Recherchewerkzeugen ein neues, wenn auch nicht ganz zuverlässiges Bild, über Inhalte gewinnt, zu denen man früher ohne dieses Hilfsmittel keinen Zugang gehabt hätte. Angesichts dessen relativiert sich die Frage des zuverlässigen Texts zumindest für Fragestellungen, bei denen ein anderer, meint traditionell „lesender“ Zugang<sup>121</sup> ohnehin verwehrt wäre. Die intensive Nutzung von Google durch die Wissenschaft scheint diese Behauptung zu stützen, ebenso wie der verstärkte Rückgriff auf quantitative Methoden in den *digital humanities*. Dennoch darf man sich nicht täuschen, dass diese Nutzungsszenarien von der Forschung als „provisorisch“ empfunden werden und dass z. B. strategische Großvorhaben wie das EU-Projekt CLARIN<sup>122</sup>, in denen vor allem linguistische Verfahren zum Tragen kommen, großen Wert auf „gute“ Texte legen, weil „schmutziges OCR“ wissenschaftliche Fragestellungen verfälschen und vor allem dann kritisch gesehen werden müssen, wenn die Herstellungsbedingungen solcher Texte, also die eingesetzten Algorithmen der OCR-Konversion und Qualitätsstufe unbekannt sind. Was also für eine wissenschaftliche Nutzung hinreichende Qualität ist und wie sie bewertet oder gemessen wird, ist daher von mehreren Faktoren abhängig, von denen die innere Textqualität, also die höchstmögliche Übereinstimmung einer Kopie mit dem Original, nur einen Faktor darstellt. Wichtige weitere Faktoren sind u. a. die Kosten (Machbarkeit) oder die Forschungsfrage, die sich mit der Textkonversion verbindet. Gerade die Kostenfrage dürfte ein Schlüsselement in der Bewertung von OCR-Verfahren bilden, denn offenbar gibt es in diesem Bereich einen *Pareto-Effekt*,

<sup>119</sup> [http://www.digitale-sammlungen.de/~mdz/mdz/content/service/docs/2011-10-19\\_Brantl\\_Digitization\\_Lifecycle.pdf](http://www.digitale-sammlungen.de/~mdz/mdz/content/service/docs/2011-10-19_Brantl_Digitization_Lifecycle.pdf) [Stand: 02/2013]

<sup>120</sup> <http://www.slideshare.net/mdz-bsb/digitalisierungspraxis-thaller-volltextdigitalisierung?ref=http://mdzblog.wordpress.com/2011/10/page/2/>

<sup>121</sup> Gemeint ist das Intensivlesen. Neuere Konzepte wie *hyperreading*, *machine reading* oder *distant reading* versuchen genau diesem Problem zu begegnen. Vgl. Stäcker, Thomas: Wie schreibt man Digital Humanities richtig? - Überlegungen zum wissenschaftlichen Publizieren im digitalen Zeitalter. In: Bibliotheksdienst 47 (2013) 1, S. 24-50.

<sup>122</sup> <http://www.clarin.eu> [Stand: 02/2013]

der dazu führt, dass der Wunsch nach hoher Textgenauigkeit die Kosten gerade für den Schritt zu sehr guten Texten sprunghaft ansteigen lässt. Um die Kosten und Qualität in eine angemessene Relation zu bringen, bedürfte es jedoch einerseits klarerer Vorstellungen darüber, was unter Qualität zu verstehen, andererseits Verfahren, wie diese Qualität zu messen ist. Beides ist derzeit nicht oder nur mit Einschränkung gegeben. So legen die genannten Zahlen aus der BSB zwar Genauigkeitsgrenzen fest, sagen aber weder, wie diese zu begründen sind, noch, wie sie gemessen werden. Obwohl es nicht Aufgabe des Wolfenbütteler Pilotprojekts Helmstedter Drucke Online I war, diese Aspekte systematisch zu untersuchen, sondern nur ein Praktikabilitätsszenario für die massenhafte OCR-Konversionen frühneuzeitlicher Drucke zu entwickeln, ergaben sich einige Beobachtungen, die es erlauben, sich der Frage anzunähern.

Die Bezifferung von Textgenauigkeiten, also der Übereinstimmung von Original und Kopie, ist der entscheidende Faktor in der Bewertung von OCR bzw. Konversionsverfahren. So simpel diese Aussage ist, so schwierig ist ihre Umsetzung und damit Ermittlung des „Werts“ in der Praxis möglich. Hauptprobleme sind a) die Möglichkeit der Überprüfung (exakter Referenztext, so genannte „ground truth“), b) die Festlegung der Art der Genauigkeitsfeststellung (Buchstaben, Worte, Layout) und c) Normierung des Gegenstands (Digitalisierungsvorlage, Digitalisat) der Messung.

Zu a) Bei der hier bearbeiteten Materie, also universitäres Schrifttum vor allem des 17. und 18. Jahrhunderts liegen keine Referenztexte vor, von denen eine Differenzmenge gebildet und die absolute Genauigkeit des Konversionsergebnisses festgestellt werden könnte. Gleichwohl wurde versucht, die Genauigkeit der OCR zumindest näherungsweise zu ermitteln, indem im Laufe der ersten zwei Jahre des Projekts immer wieder Stichproben genommen wurden. Die Projektmitarbeiterin überprüfte jeweils einzelne Zeilen an definierten Stellen des Buchs, die bei kleineren Drucken den Seiten 4, 5 und 6, bei umfangreicheren den Seiten 10, 15 und 20 entnommen wurden. Auf diese Weise wertete sie 35 Drucke aus. Gezählt wurde, wie gesagt, die korrekte Erkennung einzelner Buchstaben, nicht die von Wörtern. Die Genauigkeit lag im Schnitt bei diesen Werken bei 96%. Die Schwankungsbreite war vorlagenbedingt erheblich, die Genauigkeit fiel jedoch nie unter 90%. Bei optimaler Erkennung wurden 99%, in Einzelfällen sogar mehr erreicht. Angesichts des Umstands, dass das Projekt es mit ausserordentlich anspruchsvollen Materialien zu tun hatte (schlechte Druckvorlagen der Handpressenzeit, Mischung von Antiqua, Kursive und Fraktur), kann man das erzielte Ergebnis als durchaus bemerkenswert betrachten, da andere Softwarelösungen ein solches Ergebnis kaum erzielt hätten (vgl. hierzu die obigen Ergebnisse des Berliner Projekts). Allerdings ist ein Vergleich von BIT-Alpha mit Abbyy oder Omnipage, also einem lernenden mit einem statischen System, nur bedingt möglich, weil ein Versuchsaufbau „unter gleichen Bedingungen“ implizit das statische System privilegiert, denn ein lernendes System setzt voraus, dass die Erkennungssoftware individuell für das jeweilige Objekt optimiert wird. Der eigentliche Vorteil des lernenden Systems wird damit relativiert, ist aber in einem solchen Vergleich unvermeidbar. Umgekehrt kann in Massenprojekten nur ein begrenzter Aufwand getrieben werden, um Software anzupassen, so dass die Frage nach der besseren Software auch ökonomisch entschieden werden muss (vgl. unten).

Ad b) ist festzulegen, was man als Fehler interpretiert. Sind z. B. nicht erkannte Layoutinformationen wie Marginalien, wenn sie zwar richtig erkannt, aber an der falschen Stelle in den Haupttext einsortiert wurden, ein Fehler? Bei den Stichproben wurden lediglich Buchstabengenauigkeiten, also das engere *pattern matching*, gewertet. Worttrennungen hingegen, die oft auch in der Vorlage nur „intellektuell“ zu erkennen waren, blieben unberücksichtigt. Der Einsatz von Wörterbüchern wirkte sich bei der Worttrennung positiv aus, wie sie auch die Ergebnisse nachvollziehbar verbesserten, allerdings nur in Texten, wo es keine Mischung von deutscher Fraktur und lateinischer Antiqua gab.<sup>123</sup>

Systematisch untersucht werden konnte der Einsatz nicht, jedoch scheint es eine Art absolute Grenze zu geben (bei ca. 99%), jenseits der reines *pattern matching* bzw. Training keine signifikante Verbesserung mehr zeitigt. Dies liegt vor allem daran, dass bestimmte Buchstaben nicht mehr visuell, sondern nur intellektuell differenziert werden können (wie bereits erwähnt ist dies u. a. bei c und schlecht gedrucktem e häufig der Fall). Verbesserungen jenseits dieser Grenze sind demzufolge nur noch mit Werkzeugen der *intelligent character recognition* – vor allem durch den Einsatz von Wörterbüchern – möglich.

<sup>123</sup> Hier besteht sicher Verbesserungsbedarf und dem Vernehmen nach versucht die Firma über den Ansatz der Fontserkennung auch die Wörterbuchanwendung zu optimieren. In einem ersten Schritt hatte die Mischung von Wörterbüchern positive Effekte, jedoch konnte das noch nicht systematisch untersucht werden.

Ad c) waren die Grundlage entweder hochwertige TIFFs oder hochauflösende JPEGs auf der Basis der DFG-Praxisregeln. Der Binarisierungsprozess fand also auf der Basis von guten bis sehr guten Digitalisierungsvorlagen statt (es wurden keine Vorlagen vom Film oder in bitonaler Qualität verwendet). Dennoch konnte man auch hier nicht von gleichbleibenden bzw. Normbedingungen ausgehen, denn im Digitalisierungsprozess der HAB wird Rücksicht auf die Materialität der Bücher genommen. Wenn das Buch eine starke Wölbung im Papier aufweist, so wird nicht versucht, diese durch mechanische Maßnahmen (Anpressen) zu nivellieren. Gelegentlich vorkommende leichte Schräglagen des Buchblocks werden nicht ausgeglichen, so dass die Software auf heterogene Fälle trifft, die softwareseitig ausgerichtet werden müssen und negative Einflüsse auf die Erkennungsergebnisse haben.

Noch stärker kommt dieser Effekt mit Blick auf die Materialität der Objekte selbst zum Tragen (intrinsische Effekte). Die schon erwähnten Schattenbilder durch Widerdruck oder auch Unterstreichungen im Text beeinflussen die Konversionsergebnisse massiv. Gleiches gilt für Faktoren wie die Schrifttype, Satz und Papiergüte (Bräunung, Verschmutzung, etc.).

Nach knapp zwei Jahren Erfahrung in diesem Bereich kann man zwar sagen, welche Faktoren das Ergebnis beeinflussen, nicht jedoch wie diese im Schnitt gemessen werden sollen. Je nach Scanbeschaffenheit und Vorlage erreicht man bei einem Buch 99%, bei nächsten nur noch 95%, selbst wenn alle Software-, Trainings- und Wörterbuchparameter optimiert wurden. Nötig wäre hier eine Art Qualitätsindex für das Digitalisat und die Vorlage. Dabei ist nicht nur von Interesse, ob die Software „gut“ oder „schlecht“ arbeitet, sondern auch, mit welchen Genauigkeiten man bei welcher Art von Material rechnen kann. Mit anderen Worten, um eine bessere Vorhersage zu möglichen Qualitäten zu machen und Ergebnisse überhaupt bewerten zu können, wäre auch eine Typisierung der Vorlagen bzw. eine differenzierte Herangehensweise erforderlich.

Trägt man alle diese Faktoren zusammen, so wird deutlich, dass Genauigkeit nur über eine sehr komplexe Formel, wenn überhaupt präzise zu ermitteln ist. Selbst wenn man einen Referenztext hätte, könnte man nicht zuverlässig davon ausgehen, dass die Genauigkeit, die bei diesem Text erzielt wurde, auch auf alle anderen Texte übertragbar ist. Die Daten zur Genauigkeit, die im Laufe von Helmstedt I erhoben wurden, sind daher als Tendenz, nicht als absolute Größe zu verstehen und variieren, je nachdem, welchen Aspekt man hervorhebt. In diesem Sinne kann man aus den Erfahrungen davon ausgehen, dass das Ergebnis nach der Münchner Einteilung tendenziell befriedigend bis gut ist.

Dennoch war den Projektleitern nach dieser ersten Phase klar, dass, obwohl Erfahrungen sich mit zunehmender Menge der Stichproben zwar nach dem Gesetz der großen Zahl konsolidieren lassen, eine zuverlässige Einschätzung erst auf der Grundlage methodisch zuverlässiger statistischer Verfahren möglich ist. Daher wurden die in der Folge beschriebenen statistischen Verfahren entwickelt, die, was die Überprüfung der gelieferten Genauigkeiten anlangt, im Folgeprojekt zum Einsatz kommen sollen.

## 8.7 Messung der Textgenauigkeit<sup>124</sup>

Die Frage der Messung der Textgenauigkeit ist von entscheidender Bedeutung für eine wissenschaftliche und ökonomische Bewertung der Textgüte und Textproduktion. Da wegen der Heterogenität des Materials eine Messung an Hand vorgegebener Muster und eine vollständige Überprüfung der Texte nicht möglich ist, bedarf es statistischer Verfahren, um mittels einer Stichprobe Genauigkeitswahrscheinlichkeiten bzw. Irrtumswahrscheinlichkeiten zu ermitteln.

Bei der Ermittlung von Genauigkeiten sind in Abhängigkeit von der Fragestellung zwei Verfahren in Betracht zu ziehen. Das erste Verfahren dient dazu, grundsätzlich die Güte einer Software zu testen, das zweite erlaubt die behauptete Genauigkeit eines Texts zu überprüfen.

<sup>124</sup> Die folgenden Überlegungen entstanden im Gespräch mit Herrn Burkard Rosenberger (Münster), der den Autor in den mathematischen Aspekten der Messung beraten und auch die genannten Verfahren entwickelt hat, wofür ihm herzlich gedankt sei.

## 8.7.1 Testen der Güte einer Software

Das hier beschriebene Verfahren kann dann zum Einsatz kommen, wenn es darum geht, die unbekannte Erkennungsquote einer Software abzuschätzen und dann die Güte dieser Schätzung zu bestimmen. Es wäre darin für Projektnehmer von Interesse, die den Einsatz einer bestimmten Software erwägen, oder aber auch für Hersteller von OCR-Software, um eine Erkennungsquote ihrer Software zu bestimmen.

Gegeben ist eine Grundgesamtheit von allen gescannten und OCR-erkannten Seiten aus mehreren Büchern. Die Zufallsvariable ist die Quote der pro Seite richtig erkannten Zeichen, wobei unterstellt wird, dass die Quote normalverteilt ist. Das ist realiter bei sehr heterogenen Materialien aus der Frühen Neuzeit zwar nicht immer der Fall, aber man kann von „typischen“ Materialien ausgehen und sollte für die Stichprobe nicht Seiten wählen, die durch einen ungewöhnlich hohen Grad von Verschmutzung oder Annotationen bzw. Durchstreichungen das OCR-Ergebnis massiv verfälschen würden. Das Ziel der Ermittlung ist die Quote der Erkennungsgenauigkeit durch den Mittelwert einer Stichprobe abzuschätzen und die Güte der Schätzung zu bestimmen. Die Güte der Schätzung hängt von der Stichprobengröße und der Standardabweichung der Stichprobe ab.

**1. Schritt:** Wahl einer Konfidenzzahl. Diese bestimmt die gewünschte Wahrscheinlichkeit der Aussage.

Bsp. Konfidenzzahl  $\gamma = 0.99$ , d.h. mit einer Wahrscheinlichkeit von 99% liegt die wirkliche Erkennungsquote im Intervall  $[x,y]$ . Die Konfidenzzahl gibt die Güte der Schätzung an.

**2. Schritt:** Erhebung einer Stichprobe, z. B. 10 Seiten aus einem Korpus druck- und zeittypischer Seiten. Die Auszählung der zehn Stichprobenseiten ergibt z. B. folgende Genauigkeitswerte: 0,93/0,87/0,99/0,95/0,95/0,94/0,9/0,8/0,87/0,92.

**3. Schritt:** Berechne Mittelwert und Standardabweichung dieser Stichprobe. Mittelwert hier 0,912 (Schätzgröße für die Erkennungsquote). Die Standardabweichung  $s$  beträgt 0,05411921 (je geringer, umso besser ist die Schätzung).

**4. Schritt:** Man ermittle  $c$  aus der Tabelle für die so genannte Studentsche t-Verteilung.<sup>125</sup> Dies erfolgt mittels des Stichprobenumfangs minus 1 ( $n-1$ ) und der Formel  $0.5*(1+\gamma)$ , indem  $n-1$  die Zeile, das Ergebnis aus der Berechnung von  $0.5*(1+\gamma)$  die Spalte der Tabelle der Studentischen t-Verteilung bilden. Aus der Spalte (hier  $0.5*(1+0.99)=0.995$ ) und der Zeile (hier  $n=10-1=9$ ) ist also der Wert für  $c$  zu entnehmen (hier  $c=3,25$ ).

**5. Schritt:** Berechnung des Konfidenzintervalls  $a$ : Zusammen mit den Angaben zur Stichprobe ( $n$ ) und der Standardabweichung ( $s$ ) kann nun mit Hilfe der Formel  $a = c*s / \sqrt{n}$  das Konfidenzintervall berechnet werden ( $a = 3,25*0,05411921/\sqrt{10} = 0,05562049$ )

Als Ergebnis der Berechnung kann festgestellt werden, dass mit einer 99%igen Wahrscheinlichkeit das Ergebnis im Intervall  $[\text{Mittelwert}-a, \text{Mittelwert}+a]$ , (hier  $[0,912-0,05562049; 0,912+0,05562049]$ ) also im Bereich zwischen 0,85637951 und 0,96762049 liegt. Wir haben hier also noch eine recht große Schwankungsbreite, was vor allem an der relativ geringen Stichprobengröße liegt. Mit anderen Worten, je größer die Stichprobe und je geringer die Standardabweichung, umso schmaler wird dieser Schwankungsbereich werden. Natürlich kann man auch mit einer abweichenden Konfidenzzahl, z. B. 80%, rechnen, doch scheint dies für eine möglichst genaue Bewertung des OCR-Ergebnisses nicht zielführend.

<sup>125</sup> [http://de.wikipedia.org/wiki/Studentsche\\_t-Verteilung](http://de.wikipedia.org/wiki/Studentsche_t-Verteilung) [Stand: 02/2013]

## 8.7.2 Überprüfen der Güte eines Texts

Neben der Frage nach der unbekanntem Güte der Erkennungsquote der Software ist vor allem die Frage nach der behaupteten Güte des Texts von Bedeutung, wenn Arbeitsergebnisse von Dienstleistern überprüft werden müssen. Auch hier verbietet die Textmenge aus pragmatischen Gründen eine Auszählung bzw. Überprüfung und erfordert ein statistisches Verfahren zur Ermittlung der Genauigkeit. Es geht dabei darum, anhand einer Stichprobe zu überprüfen, ob die vom Dienstleister behauptete Erkennungsquote stimmt, wobei man die Wahrscheinlichkeit für einen eigenen Irrtum möglichst gering halten möchte. Das ist ein sog. Bernoulli-Experiment, und die zugehörige diskrete Verteilung ist die Binomialverteilung. Dafür gibt es einige markante „Grenzen“ in Abhängigkeit von Stichprobengröße und akzeptierter Irrtumswahrscheinlichkeit.

Beispiel: Ein Anbieter behauptet, dass seine Software mindestens 90% der Zeichen korrekt erkennt. Wenn man nun eine Stichprobe erhebt und darin exakt 90% der Zeichen korrekt erkannt werden, kann man nicht ohne weiteres davon ausgehen, dass die Aussage des Anbieters auch wirklich stimmt, denn jedes Ergebnis einer Stichprobe streut um den errechenbaren Mittelwert. Es könnte also durchaus sein, dass die Erkennungsquote bei 85% liegt und nur zufällig in der Stichprobe ein Ausreißer nach oben vorlag. Tatsächlich liegt die Wahrscheinlichkeit, dass diese Aussage zutrifft, bei lediglich 50%.

Es erhebt sich nun die Frage, wie selten solche Ausreißer in Abhängigkeit von der Stichprobengröße sind, d. h. welche Irrtumswahrscheinlichkeit man bereit ist zu akzeptieren und welche Stichprobengröße mindestens einzuplanen ist, um ein akzeptables Ergebnis zu erlangen. Es leuchtet unmittelbar ein, dass die Irrtumswahrscheinlichkeit abnimmt, je größer die Stichprobe wird, jedoch möchte man die Stichprobe so klein wie möglich und nötig halten, um den Aufwand des Auszählens zu optimieren bzw. überhaupt einen pragmatischen Zugang zu einer soliden Abschätzung zu eröffnen. Dabei muss man eine plausible Balance finden zwischen akzeptierter Irrtumswahrscheinlichkeit und Größe der Stichprobe.

Da die Berechnung mathematisch kompliziert ist,<sup>126</sup> seien im Folgenden mehrere Tabellen mit fest vorgegebenen Werten aufgeführt, die typische Angaben enthalten. Vorgeschlagen wird eine Stichprobengröße von 500 beliebigen Zeichen, die idealerweise durch Benutzung eines Zufallsgenerators ausgewählt werden. Unter dieser Voraussetzung gilt folgende Tabelle. In der linken Spalte ist die behauptete Erkennungsquote angegeben, in der rechten die Zahl der in der Stichprobe mindestens korrekt erkannten Zeichen, die vorliegen muss, um die Behauptung des Dienstleisters bei einer Irrtumswahrscheinlichkeit von 2,5% als korrekt einstufen zu können. Wenn also ein Dienstleister behauptet, dass ein Text eine Genauigkeit von 96% hat, müssen in der Stichprobe von 500 Zeichen mindestens 489 Zeichen korrekt erkannt werden, damit bei einer Irrtumswahrscheinlichkeit von 2,5% die Behauptung des Dienstleisters akzeptiert werden kann.

<sup>126</sup> Zur Binomialverteilung vgl. das Skript von Albers/Yanik: [http://www.math.uni-bremen.de/didaktik/ma/ralbers/Veranstaltungen/Stochastik11/Material/ScriptK7\\_110703.pdf](http://www.math.uni-bremen.de/didaktik/ma/ralbers/Veranstaltungen/Stochastik11/Material/ScriptK7_110703.pdf) [Stand: 02/2013]

Behauptete Erkennungsquote	Mindestzahl der korrekt erkannten Zeichen (Stichprobengröße = 500)
95%	485 (96,9%)
96%	489 (97,8%)
97%	493 (98,5%)
98%	496 (99,3%)
99%	499 (99,9%)
>99%	500 (100%)

Zum Vergleich eine Tabelle mit einer Stichprobengröße von 2.000 Zeichen:

Behauptete Erkennungsquote	Mindestzahl der korrekt erkannten Zeichen (Stichprobengröße = 2000)
95%	1920 (96%)
96%	1938 (96,9%)
97%	1956 (97,8%)
98%	1972 (98,6%)
99%	1988 (99,4%)
>99%	2000 (100%)

Bei einer behaupteten Erkennungsquote von über 99% scheint es sinnvoll, ggf. die Stichprobengröße zu erhöhen. Nachstehend zwei Tabellen, die beispielhaft zeigen, wie hoch die ermittelte Erkennungsquote in Abhängigkeit von einer bestimmten Stichprobengröße sein muss, wenn Texte eine behauptete Erkennungsquote von 99,5% bzw. 99,7% haben. Die Irrtumswahrscheinlichkeit liegt jeweils bei 2,5%:

Behauptete Genauigkeit: 99,5%	
Stichprobengröße	Mindestzahl der korrekt erkannten Zeichen
500	500
1000	999
2000	1996
5000	4985
10000	9960

Behauptete Genauigkeit: 99,7%	
Stichprobengröße	Mindestzahl der korrekt erkannten Zeichen
500	500
1000	1000
2000	1998
5000	4995
10000	9990

Abschließend eine Tabelle mit einer höheren Irrtumswahrscheinlichkeit von 5% statt 2,5%:

Behauptete Erkennungsquote	Mindestzahl der korrekt erkannten Zeichen (Stichprobengröße = 2000)
95%	1916 (95,8%)
96%	1934 (96,7%)
97%	1952 (97,6%)
98%	1970 (98,5%)
99%	1988 (99,4%)
>99%	2000 (100%)

Die Frage, welche Irrtumswahrscheinlichkeit man sinnvollerweise zugrunde legt, kann nicht verbindlich beantwortet werden. Erfahrungswerte zeigen, dass 5% eine sinnvolle Größe darstellt, die nicht überschritten werden sollte. Da der Mehraufwand der Abschätzung bei 2,5% gegenüber 5% gering ist, sollte eine Irrtumswahrscheinlichkeit von 2,5% zugrunde gelegt werden, um ein verlässliches Ergebnis zu erhalten.

Ein Excelsheet, mit dem man weitere Angaben durch Eingabe alternativer Parameter berechnen kann, findet sich auf der Projektseite der Herzog August Bibliothek.<sup>127</sup> Eine Schwierigkeit besteht darin, ein pragmatisches Verfahren zur Auszählung zu finden. Zunächst muss die Stichprobe per Zufallsprinzip erhoben werden. Nutzen kann man dafür die in Programmiersprachen meist vorhandenen Zufalls-generatoren. Um echte Zufallszahlen<sup>128</sup> zu generieren, lassen sich geeignete Webservices nutzen.<sup>129</sup> Um die Auszählung zu erleichtern, entwickelt die Herzog August Bibliothek derzeit für die Messung von Genauigkeiten einen Webservice, der im zweiten Teil des Projekts, Helmstedter Drucke online II, zum Einsatz kommen soll.

## 8.8 Ökonomische Betrachtungen

Die Kosten für die Konversion einer Seite können beträchtlich schwanken. Bei der Nutzung von Standard-OCR-Software liegen die sächlichen Kosten nur in der Beschaffung der Software. Die Personalkosten für die Auswahl, Vorbereitung und Prozessverwaltung müssen jedoch berücksichtigt werden. Gerade bei Massenkonzessionsprozessen ist die Etablierung fester Workflows wichtig und der Einsatz von leistungsfähiger Hardware, denn der OCR-Prozess ist sehr rechenintensiv und daher bei schlechter Hardwareausstattung sehr zeitaufwendig und mit Kosten verbunden. Insofern ist es eine irri-ge Annahme, man könne mit Standardsoftware nahezu kostenneutral arbeiten.

Dies vorangeschickt, ist an dieser Stelle die grundsätzliche Frage zu stellen, welche Qualität erforderlich ist und wie viel sie kosten darf. Denn ohne eine Qualitätsanforderung zu stellen, sind auch differenzierte ökonomische Betrachtungen wertlos; die Entscheidung wäre dann einfach und müsste immer für die günstigste Variante ausfallen. Wenn aber die günstigste Variante nicht die Qualität zu liefern vermag, die mindestens erforderlich ist, dann ist das hier aufgewendete Geld schlicht verschwendet und selbst die günstigste Variante noch zu teuer. Sich diese wirtschaftlichen Gemeinplätze immer wieder zu vergegenwärtigen und Aspekte von Effektivität und Effizienz in die Diskussion um die Qualität von OCR-Software einzubringen, ist dringend erforderlich, wenn man realistische Szenarien für eine massenhafte Konversion von Drucken der Frühen Neuzeit in maschinenlesbare Texte entwickeln will.

Den Ausgang für solche Überlegungen sollte man bei einer Festlegung der Zwecke nehmen, für die maschinenlesbare Texte verwendet werden sollen und für die bestimmte Qualitäten notwendig sind.

<sup>127</sup> Siehe: <http://www.hab.de/de/home/wissenschaft/projekte/helmstedter-drucke-online.html> [Stand: 02/2013]

<sup>128</sup> <http://de.wikipedia.org/wiki/Zufallszahl> [Stand: 02/2013]

<sup>129</sup> Z. B. <https://www.random.org/> [Stand: 02/2013]

Leider gibt es bislang hierzu in der wissenschaftlichen *community* keine normativen Aussagen. Sicher kann man davon ausgehen, dass ein wissenschaftlich brauchbarer Text, also ein Text der sowohl als Lesefassung als auch für zuverlässige Recherchen *ex silentio*<sup>130</sup> taugt, eine Genauigkeit von mindestens 99,95%<sup>131</sup> besitzen muss, d. h. man nimmt pro Seite im Schnitt etwa 1-2 Buchstabenfehler in Kauf. Texte dieser Güte lassen sich für Materialien der Handpressenzeit (15. Jh. bis 18. Jh.) derzeit nur durch Abschreiben (*double keying*) erreichen. Die Kosten schwanken zwischen 1,20 und 2,20 € pro Seite, wobei ggf. weitere Dienstleistungen wie die XML-Grundkodierungen der Dokumente hinzukommen. Ginge man davon aus, dass z. B. im 17. Jh. etwa 200.000 Drucke konvertiert werden sollten und legte man eine Durchschnittszahl von 140 Seiten zugrunde,<sup>132</sup> müssten für die 28 Mio Seiten mindestens 33,6 Mio Euro aufgewendet werden. Selbst wenn eine solche Qualität wünschbar wäre, scheint sie doch angesichts der hohen Kosten unrealistisch, ja nicht einmal nötig, da man bei ausgewählten Materialien auch mit einer geringeren Qualität arbeiten könnte. Doch welche Qualität ist ausreichend?

Standard-OCR-Software führt, wie man weiss, bei Materialien der Frühen Neuzeit zu hohen Fehlerquoten und die Qualität sinkt bei den hier betrachteten Materialien deutlich unter den Wert, den man nach der Münchner Skala als „schlecht“ bezeichnen könnte. Da es im Bereich des „Schlechten“ aber durchaus noch einen Bereich des „Nützlichen“ gibt, bleibt diese Qualifizierung abhängig von der Frage, welchen Gebrauch man von dieserart Texten machen möchte. Auch Aspekte der nachträglichen Verbesserungen (semantische Analyse, Wörterbücher, *crowdsourcing*) stehen unter diesem Vorbehalt.<sup>133</sup> Darauf nach dem derzeitigen Stand des Wissens eine angemessene und für alle wissenschaftlichen Felder gleichermaßen zufriedenstellende Antwort zu finden, dürfte schwer fallen. Dennoch kann man Tendenzen feststellen, die zumindest eine grobe Orientierung bieten.

Klammert man editorische Ansprüche aus, die eigene Anforderungen mitbringen,<sup>134</sup> dürfte der maßgebliche Gebrauch, der von Volltexten in wissenschaftlichem Kontext gemacht wird, die Suche sein. Demgegenüber scheint die Lesbarkeit eines Texts minder wichtig, zumindest dann, wenn die Alignierung des digitalen Faksimiles mit dem Volltext gewährleistet ist. Das digitale Faksimile übernimmt hierbei die Funktion der Lesefassung. Ein in den Geisteswissenschaften noch wenig wahrgenommener Nutzen liegt in quantitativen Analyseverfahren und der semantischen Verarbeitung maschinenlesbarer Texte, die innerhalb der voranschreitenden *digital humanities* entwickelt werden. Auch diese bringen einen eigenen Qualitätsanspruch mit.

Wenn man die Suche als das derzeit wichtigste Instrument herausgreift und fragt, welche Textgüte einer wissenschaftlichen Anforderung genügt, lassen sich zwei Typen von hauptsächlichen Suchverfahren unterscheiden, nämlich phraseologische und Suchen nach Entitäten oder Stichwörtern. Phraseologische Suchen nach bestimmten sprachlichen, aus mehreren Wörtern bzw. aus längeren Phrasen bestehenden Textteilen können bei einer Qualität, die nicht bei 99,95% liegt, vermutlich nicht sinnvoll durchgeführt werden oder die Ergebnisse sind stark zufallsbehaftet. Anders liegt der Fall bei Suchen nach Entitäten oder einzelnen Stichwörtern. Das ist auch typischerweise das, was Forscher bei Google und ähnlichen Suchmaschinen suchen. Entitäten- oder in einem weiteren Sinn Einwortsuchen – sind auch in Texten möglich, die man als „schlecht“ bezeichnen könnte, indem das gesuchte Wort „zufällig“ richtig erkannt wurde. Allerdings lauern hier Gefahren der Fehleinschätzung des Suchergebnisses. Die Wahrscheinlichkeit bei Texten der Frühen Neuzeit einen Treffer zu erzielen, ist bei handelsüblicher Software stark eingeschränkt, denn eine Erkennungsgüte von unter 90%, die hier nicht untypisch ist, bedeutet, dass ca. jedes vierte Wort falsch ist. Das klingt zunächst nicht dramatisch, weil es ja auch bedeutet, dass drei Wörter richtig sind und man z. B. einen Namen sehr wohl finden könnte, jedoch muss man in Betracht ziehen, dass diese Fehler sich nicht statistisch antizipieren lassen und möglicherweise nur „uninteressante“ Wörter korrekt sind. Problematisch in Texten der Frühen Neuzeit sind Mischtexte, in denen bei dieserart Software nur Antiqua oder nur Fraktur korrekt erkannt werden, so dass gerade die

<sup>130</sup> D. h. dass man ein negatives Rechercheergebnis so werten kann, dass ein bestimmter Begriff nicht vorkommt.

<sup>131</sup> Hier und im Folgenden ist von Buchstabengenauigkeit die Rede.

<sup>132</sup> Vgl. Anm. 118

<sup>133</sup> Vgl. dazu die Ergebnisse des IMPACT Projekts (<http://www.digitisation.eu/>) [Stand: 02/2013]

<sup>134</sup> Zu den Wünschen von Nutzern, die an eine elektronische Edition gestellt werden, vgl. Sören A. Steding, *Computer-based scholarly editions : context, concept, creation, clientele*. Berlin 2002, S. 234.

interessierenden deutschen oder lateinischen Namensformen nicht gefunden werden. Beide Umstände lassen die Wahrscheinlichkeit eines Treffers mit handelsüblicher Software weiter sinken. Die Gefahr der Fehleinschätzung des Suchergebnisses besteht darin, dass der Suchende entweder glaubt, dass ein Text einen Namen nur singular nennt, wo er in Wahrheit allgegenwärtig ist, oder er meint, dass bei einem negativen Resultat der Name, wenn er vorkommt, doch nur selten vorkommt, denn andernfalls hätte man ihn doch wenigstens einmal finden müssen. Um diese Effekte abzufedern, hilft die Forderung der Praxisregeln der DFG, schmutziges OCR nicht zu verbergen, weil sichtbar wird, wie gut oder schlecht der Text tatsächlich ist, doch kann sie dieserart Fehleinschätzungen nicht grundsätzlich verhindern bzw. der Suchende kann eigentlich nie zuverlässig davon ausgehen, dass das, was ihm als Ergebnis vorliegt, tatsächlich so zu interpretieren ist, wie es sich zeigt. Mit anderen Worten, der Fund ist eigentlich nutzlos oder nur in einem sehr oberflächlichen Sinne hilfreich, indem man argumentiert, dass es immer besser ist, etwas zu haben, als nichts.

*Zusammenfassend scheint daher, dass der Nutzen einer wissenschaftlich motivierten Volltextsuche mit schmutzigem OCR oder „schlechtem Text“ (< 90% nach der Münchner Skala) gering bzw. sogar fragwürdig ist, weil es nahezu unmöglich ist, Treffer zu kontextualisieren und korrekt zu bewerten, – wenn man sich nicht die Mühe macht oder aus zeitlichen Gründen machen kann, das Dokument, aus dem der Fund stammt, zu lesen. Positive Treffer bleiben zufällig und das Risiko der Fehlinterpretation ist hoch, auch wenn natürlich ein Tauchgang die eine oder andere Perle zutage zu fördern vermag. Angesichts dieser Grundüberlegungen scheint die Bereitstellung von Volltext frühneuzeitlicher Texte unterhalb eines bestimmten Qualitätsniveaus als nicht sinnvoll, zumindest für den Bereich der direkten Suche. Statistische oder probabilistische Verfahren, wie sie sich mit Konzepten aus den digital humanities verbinden (z. B. Google ngram Viewer), folgen dabei anderen Gesetzmäßigkeiten, weil hier Mängel im Detail durch statistische Masse ausgeglichen werden. Welche Qualität in dieserart quantitativen Projekten ausreichend ist, liesse sich mit Hilfe statistischer Verfahren sicher ermitteln, möglicherweise wäre ein Nutzen auch noch bei „schlechten“ Texten gegeben. Erfahrungen geschweige denn Studien dazu gibt es bisher kaum, so dass an dieser Stelle nicht näher darauf eingegangen werden kann.*

Obzwar deutlich scheint, dass eine Qualität unterhalb von 90% für die Suche unzureichend ist, bleibt schwer zu entscheiden, welche Qualität für welche Zwecke zwischen 90% und 99,95% benötigt wird und welche Kosten angemessen sind. Mit den von B.I.T. Tomasi nach händischen Ermittlungsmethoden erhobenen 96% ist sicher ein sinnvolles Mittelmaß zwischen Qualitätsanspruch und bereitgestellter Menge für wissenschaftliche Zwecke erreicht, die preislich etwa in der Mitte zwischen einer Standard-OCR-Lösung<sup>135</sup> und einer händischen Transkription liegt. In der Projektplanung kommt es also darauf an zu entscheiden, ob man Texte im Umfang mit einer Qualität von 96% oder alternativ weniger Texte im Umfang von ca. 1/3 auf einem sehr hohen Niveau per händischer Abschrift zur Verfügung stellt. Im Falle der Helmstedter Drucke, bei denen es sich überwiegend um serielles Material (serielles Universitätschrifttum mit rekursiven Phänomenen) handelt, wäre eine solche Umfangsreduktion nicht zu rechtfertigen, mag es auch in anderen Fällen, z. B. herausragende Lexika, literaturwissenschaftlich oder historisch zentrale Werke u.ä., unabweislich erforderlich sein. Hier muss gelten, das richtige, meint effizienteste Werkzeug zu gebrauchen, wo also Kosten und Nutzen miteinander in Einklang stehen. Dieser Mittelweg scheint mit der formulierten Qualitätsanforderung und dem hier verfolgten Konversionsverfahren gefunden worden zu sein, auch wenn perspektivisch die Kosten für das Verfahren noch sinken müssen, wenn es in noch größerem Maße wirtschaftlich eingesetzt werden soll.

<sup>135</sup> Vgl. z. B. das Angebot der VZG in Göttingen: <http://www.gbv.de/Verbundzentrale/serviceangebote/ocr-service-der-vzg> [Stand: 02/2013]

## 9 Register

- Abkürzungen** 20, 46, 53  
**Antiqua** 21 f., 43, 69, 124, 126  
**Bedienoberflächen** 55, 80, 116 ff.  
**Batchlauf** 36 ff., 59 ff., 98 f.  
**Benutzerschnittstelle** 44 f., 50, 58, 60, 80 ff.  
**Berufsbezeichnungen** 14  
**Betriebsvoraussetzungen** 61  
**Bildvorverarbeitung** 28 f., 42 f.  
**Binarisierung** 28 f., 42 f., 64, 67, 72 f., 102 f., 130 f.  
**Binominalverteilung** 133 ff.  
**Buchgestaltung** 14 ff.  
**Desiderata** 41, 42, 43, 45, 48, 49, 51, 58, 60  
**digital humanities** 129 ff.  
**Distanzmaß** 35, 37  
**Erkennungsgüte** 60, 72 ff., 85, 129 ff.  
**Evaluation der Ergebnisse** 44 f., 50, 58, 60, 81 f., 85  
**Exportformat** 10, 48, 51 ff., 83, 87 ff., 126  
**Exportschnittstelle** 8, 51 ff.  
**Fehlerquellen** 26, 28, 29, 31, 33, 37, 38  
**Fraktur** 21 ff., 43, 46, 69, 124, 126  
**Funeralschriftensammlung** 12 ff.  
**Gebrochene Schriften** 21 ff.  
**Good Practice** 78  
**Grafische Benutzeroberfläche** 55, 57, 80  
**Granularität** 60  
**Ground-Truth-Daten** 25, 34, 73, 130  
**Gruppierung** 22, 25, 59, 64, 67, 72  
**Handbuch** 43 ff., 102 ff.  
**Helmstedter Drucke** 123 ff.  
**Hilfsdateien** 39 ff., 79 f.  
**I / J-Gebrauch** 19 ff.  
**Importschnittstelle** 79 ff.  
**Informationsverlust** 26, 28, 29, 31, 33, 37, 38  
**Installationsvoraussetzungen** 61 f.  
**Interaktivität** 81  
**Konfidenzwert** 30 ff., 59, 65, 85, 132 f.  
**Konfigurationsmöglichkeiten** 39 ff., 102 ff.  
**Korrekturfehler** 34 ff., 48, 66 f., 74, 128  
**Krankheitsbezeichnungen** 14  
**Kurzanleitung** 98 ff.  
**Layoutanalyse** 29 f., 36, 44 ff., 128  
**Leichenpredigt** 11 f.  
**Lexikalische Nachkorrektur** 32 ff., 48 ff.  
**Linguistische Nachkorrektur** 32 ff.  
**Lizenzmodell** 40, 47, 61  
**Marginalien** 16, 75, 128, 130  
**Modularität der Software** 43, 45, 49 f., 58, 82, 85  
**Musterbibliothek** 25, 30, 39 f., 52 ff., 62, 67 ff., 73 f., 100 f.  
**Mustereditor** 53 f., 56 f.  
**Named Entities** 10, 36, 50  
**OCR** 30 ff., 46 ff., 98 f.  
**OCR-Muster** 30 f., 52 ff., 67 ff., 73, 100 f.  
**OCR-Optimierungsmöglichkeiten** 62 ff.  
**OCR-Verarbeitungsstufen** 26 ff., 42 ff.  
**Orthographische Besonderheiten** 19 ff.  
**Parameterdateien** 39 ff., 79 f.  
**Rundes „r“** 20  
**Schrifttypen** 23 ff., 67, 69 f.  
**Schwabacher** 21 f., 25  
**Segmentierung** 29 ff., 44 ff., 72, 81, 105 f., 128  
**Semantische Erschließung** 36, 50, 136  
**Sequencer** 46, 52, 53, 55, 78, 126, 128  
**Software-Weiterentwicklung** 79 ff.  
**Sortierung** 25, 26, 64, 72, 85 f.  
**Speichermodell** 39 ff., 79 f.  
**Sprachengruppe** 40, 46 f., 56, 70, 74, 99  
**Stapelverarbeitung** 36 ff., 59 ff., 98 f.  
**Strukturerkennung** 27, 50  
**Studentische t-Verteilung** 131  
**Textgenauigkeit** 129 ff.  
**Training** 52 ff., 64 f., 67 ff., 73 f., 83 f., 100 f., 126, 128  
**Typografie** 14 ff.  
**U / V-Gebrauch** 19  
**Validierungsmöglichkeiten** 49  
**Wörterbuch** 32 ff., 48 ff., 66 ff., 74, 128, 130  
**Workflow** 26 f., 38, 79, 85 f., 135  
**Wortbibliothek** 32 ff., 48 ff., 66 ff., 74, 128, 130  
**Wortkoordinaten** 32, 35, 84  
**Wortkorrektur** 32 ff., 48 ff., 78, 84

**Der vorliegende Werkstattbericht aus der Abteilung Historische Drucke der Staatsbibliothek zu Berlin – Preußischer Kulturbesitz beschreibt Testszenarien zur automatischen Texterkennung (OCR) von frühneuzeitlichen Funeralschriften. In einem von der Deutschen Forschungsgemeinschaft (DFG) geförderten Pilotprojekt wurden Optimierungs- und Konfigurationsmöglichkeiten zweier ausgewählter Softwareprodukte getestet, um gattungsspezifische Optionen der Volltextgenerierung auszuloten. Neben den Ergebnissen dieser Tests werden die Vor- und Nachteile der verschiedenen Softwarelösungen dargestellt sowie allgemeine Anforderungen an die OCR Alter Drucke formuliert. Ein Erfahrungsbericht aus dem Projekt „Helmstedter Drucke Online“ der Herzog August Bibliothek Wolfenbüttel ergänzt die Abhandlung. Thomas Stäcker beschreibt darin auch Methoden der Ermittlung von Textgenauigkeiten und geht auf ökonomische Aspekte der Volltexterstellung Alter Drucke ein.**