



OCR-D

KOORDINIERUNGSPROJEKT ZUR
WEITERENTWICKLUNG VON OCR-VERFAHREN

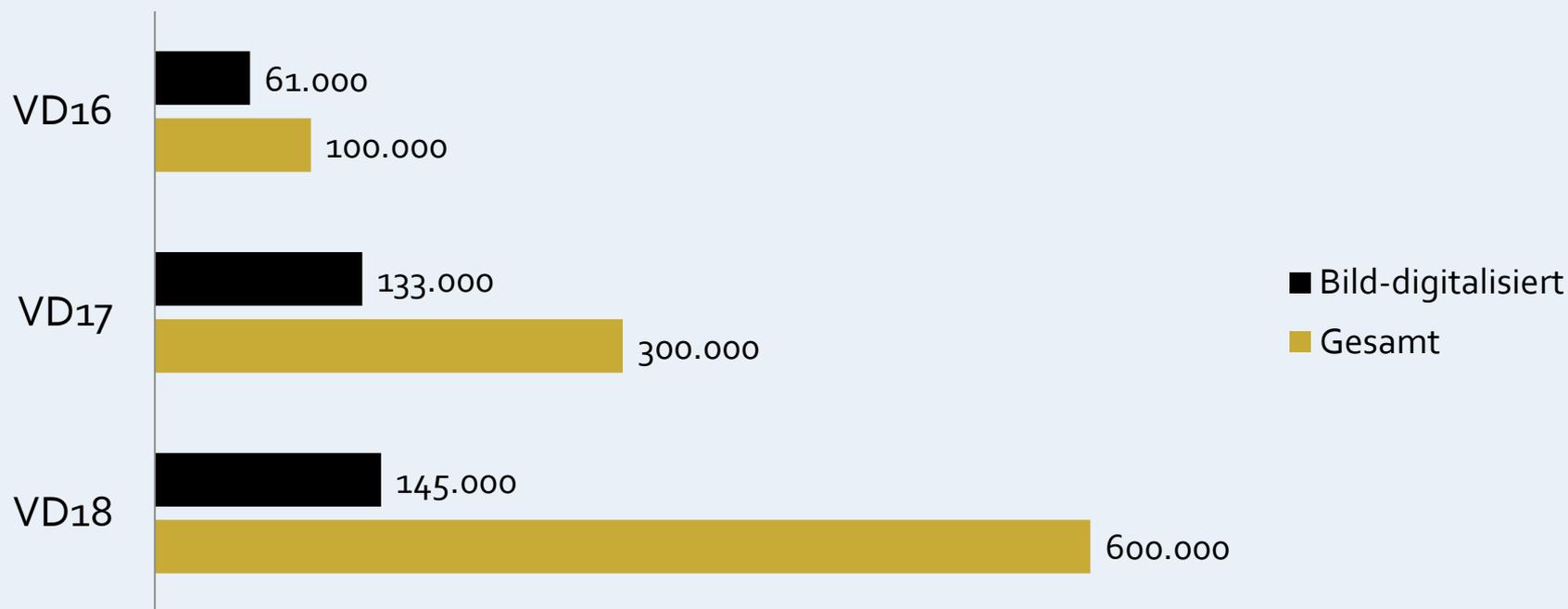
Gefördert von der Deutschen Forschungsgemeinschaft

25.02.2016

Elisa Herrmann



VD 16-18



* Gerundete Zahlen



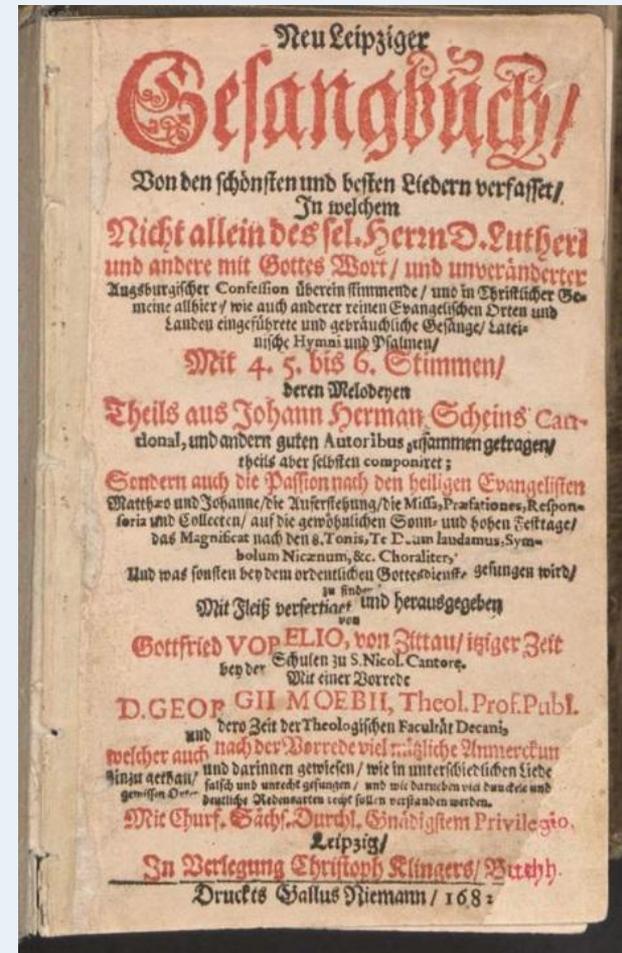
Ziel

**Konzeptionelle Vorbereitung der Transformation der VD-
Drucke (16.-18. Jh.) und der Drucke des 19. Jh. in
maschinenlesbare Form.**



Herausforderungen

- Material:
 - Sprachen: v.a. Latein, Deutsch
 - Schriftarten und -ausprägungen: u.a. Antiqua, Fraktur, Kursive
 - verschiedene Textsorten mit spezifischem Layout
- Uneinheitliche Standards
- Neuprozessierung, dynamische Datensicherung



Digitalisierung gefördert durch die Deutsche Forschungsgemeinschaft · DFG



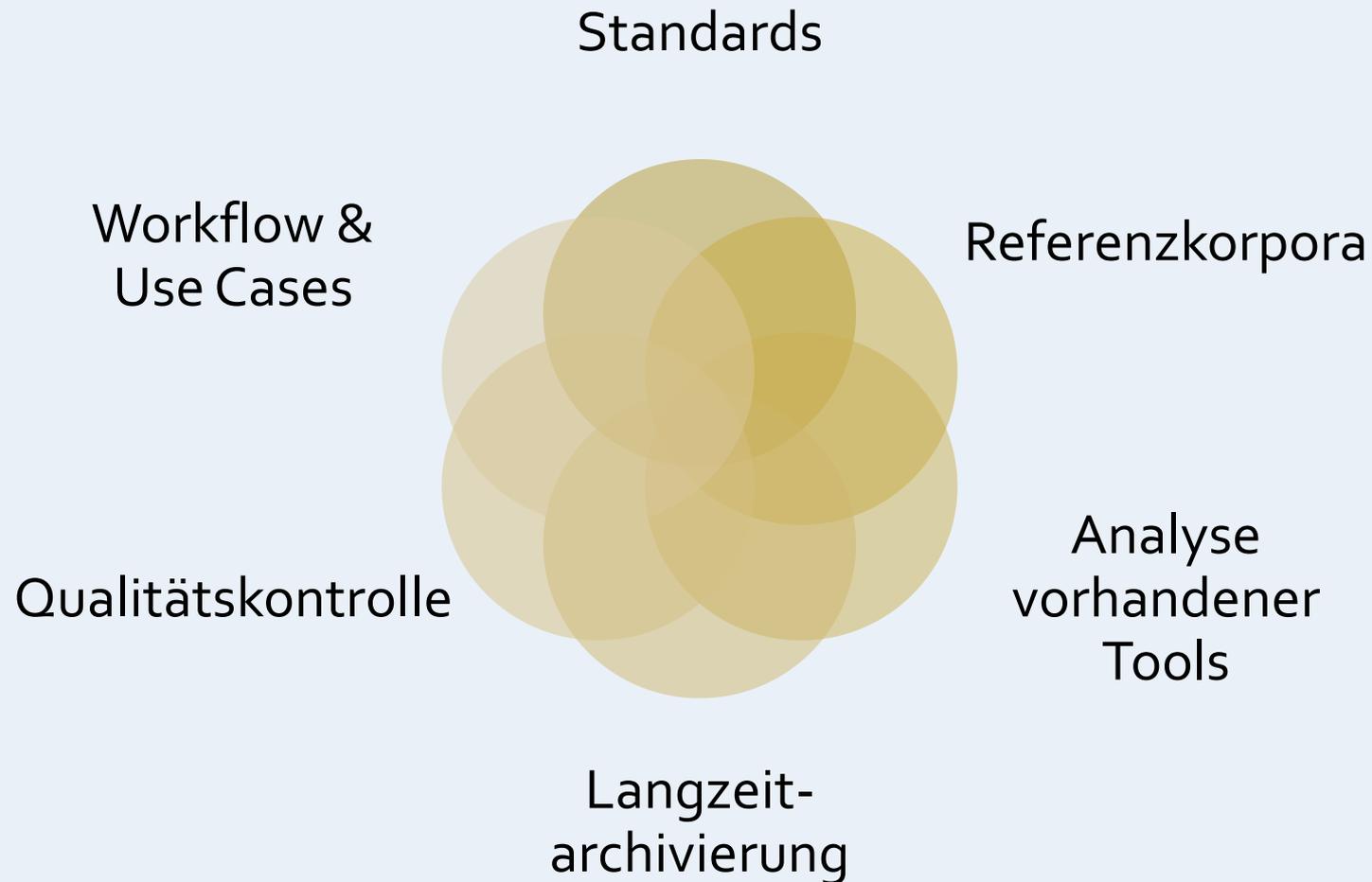
Das Projekt

- Projektpartner
 - Herzog August Bibliothek Wolfenbüttel
 - Berlin-Brandenburgische Akademie der Wissenschaften, insb. Deutsches Textarchiv (DTA)
 - Bayerische Staatsbibliothek in München

- 2 Phasen:
 1. Aufbau der Koordinierungsstruktur und Konzeption der Projektphase
 2. Ausschreibung und konzeptionelle Begleitung der Pilotprojekte



Arbeitspakete





Funktionsmodell

- Einbindung neuer Erkenntnisse aus den einzelnen Arbeitspaketen
- Einbindung bestehender Werkzeuge
- Modularer Aufbau
 - Adaptive Anpassung an verschiedene Bedürfnisse
 - Anpassung an zukünftige Entwicklungen



Funktionsmodell - I



- Bild-Digitalisat bereitstellen



- Splitting (optional)

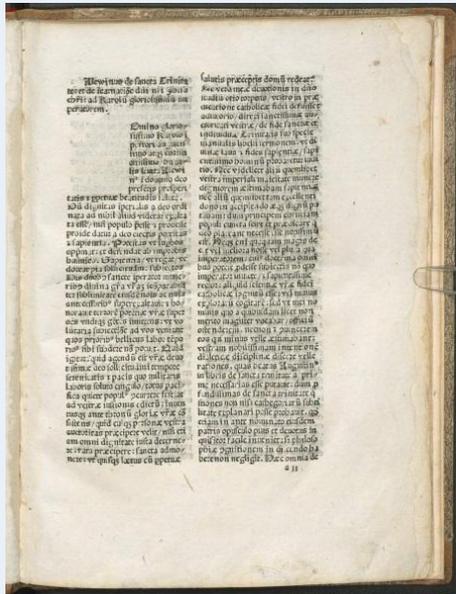


- Bildvorsortierung

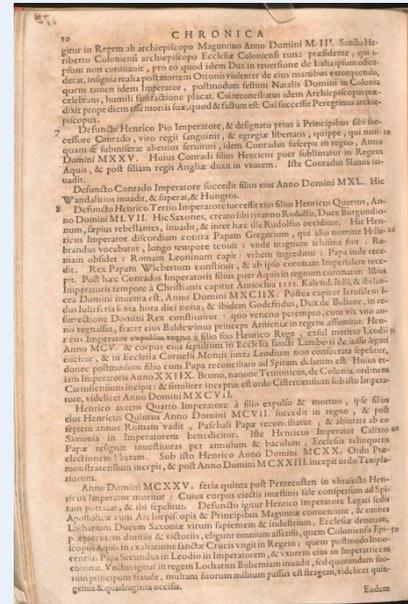


Bildvorsortierung

- Bildvorsortierung
 - Z.B. Layoutanalyse



http://daten.digitalle-sammlungen.de/bsb00005199/image_9



http://daten.digitalle-sammlungen.de/bsb00010079/image_16



Funktionsmodell - II

- 
- Preprocessing Page Level
 - (z.B. Cropping, Deskewing, Binarisierung)

- 
- Qualitätskontrolle

- 
- OLR₁: Page Segmentation

- 
- Preprocessing Segment Level

- 
- Qualitätskontrolle



Page Segmentation



- Segmentierung des Bilddigitalisats in
 - Textzonen
 - Nichttextzonen



Funktionsmodell - III



- OCR



- Qualitätskontrolle



- OLR 2: Region Classification



- Qualitätskontrolle



OCR & OLR

- OCR
 - Vgl. bzw. Kombination **klassischer Zeichenerkennungsverfahren** auf Glyphenebene mit **segmentierungsfreien Ansätzen** (z.B. "Deep Learning"-Verfahren auf Basis neuronaler Netze)
- OLR
 - **Region Classification:** Bestimmung der layout-semantischen Funktion der einzelnen Regionen (Überschrift, Marginalie, Fußnote etc.)
 - **Document Analysis:** Extrapolierung der Dokumentstruktur aus den entsprechenden Strukturelementen (Überschrift)



Qualitätsprüfung

- Einbindung verschiedener Qualitätskontrollen im Prozess
 - Binarisierung
 - Layout Analyse
 - Text Segmentation
 - OCR
- Einordnung der Ergebnisse nach Use Cases



Qualitätsprüfung

EPIS J-OLJE
JOANNEM AKEPPLERUM
MATHEMATICUM C'SAREUM
S C R I P T Æ;
INSERTIS AD EASDEM
RESPONSIONIBUS KEPPLERIANIS,
QUOTQUOT HACTENUS REPERIRI
OPUS NOVUM, Qy6'RECONDITA KEPPLERIAN'E
DOCTRIN'E CAPITA DILUCIDE EXPLICANTUR, ET HISTORIA
LITERARIA IN UNIVERSUM MIRIHCE ILLUSTRATUR,
NUNC PRIMUM
CUM PR'ÆFATIONE DE MERITIS GERMANORUM IN MATHESIN,
INTRODUCTIONE IN HISTORIAM LITERARIAM S./ECULORUM
XVI. ET XVII. ET JO. KEPPLERI VITA



Le texte affiché peut comporter un certain nombre d'erreurs.
En effet, le mode texte de ce document a été généré de façon
automatique par un programme de reconnaissance optique de
caractères (OCR). Le taux de reconnaissance estimé pour ce
document est de **86.43** %.
En savoir plus sur l'OCR



- Grundmethode: Six Sigma





Funktionsmodell - IV



- Export



- Nachkorrektur/ Crowdsourcing



- Qualitätskontrolle



- LZA



- Antworten auf technische, informationswissenschaftliche & organisatorische Probleme

Konsolidiertes Verfahren zur OCR-Verarbeitung von Digitalisaten des schriftlichen deutschen Kulturerbes des 16.-19. Jh.



Kontakt Daten

- Webseite: www.ocr-d.de

OCR-D
Koordinierungsprojekt zur Weiterentwicklung von
Verfahren der Optical Character Recognition (OCR)

Home Projektbeschreibung **▼** Kontakt Impressum

Projektziele
Technology Watch

HISTORISCHE DRUCKE (16.-19. JH.)

Das OCR-Verfahren wird insbesondere für Drucke aus dem deutschsprachigen Raum des 16.-19. Jahrhunderts optimiert. Herausforderungen sind v.a. die verschiedenen Sprachen und Schrifttypen, wie bspw. Fraktur und Antiqua.

Das Projekt OCR-D

- Elisa Herrmann, elisa.herrmann@hab.de, +49 5331 808-306