



Vom Bleiletterndruck zum Forschungskorpus Neue Dienstleistungen jenseits der OCR

Elisabeth Klein | Dr. Klaus T. Weber
BID-Kongress 2016
17. April 2016

UNIVERSITÄTS
BIBLIOTHEK
MAINZ





Überblick

1 | Faktoren der Zusammenarbeit

2 | Volltextsammlungen vs. Korpora

3 | Methoden und Test-Design

4 | Ausblick: Post-Processing

Volltextsammlung

- Menge von Volltexten in Text oder XML-Format
- meist unstrukturiert
- variable Erkennungsrate
- Metadaten optional vorhanden
- keine Annotationen enthalten



eingeschränkt durchsuchbar,
nur für qualitative Forschung
geeignet (wenn überhaupt...)

Korpus

- kuratierte Sammlung strukturierter Volltexte in Text oder XML-Format
- Metadaten für alle Texte
- stabile Erkennungsrate durch Wörterbuchabgleich
- annotiert



Negativsuchen
für qualitative & quantitative
Forschung geeignet

Volltextsammlung

- Menge von Volltexten in Text oder XML-Format
- meist unstrukturiert
- variable Erkennungsrate
- Metadaten optional vorhanden
- keine Annotationen enthalten



eingeschränkt durchsuchbar,
nur für qualitative Forschung
geeignet (wenn überhaupt...)

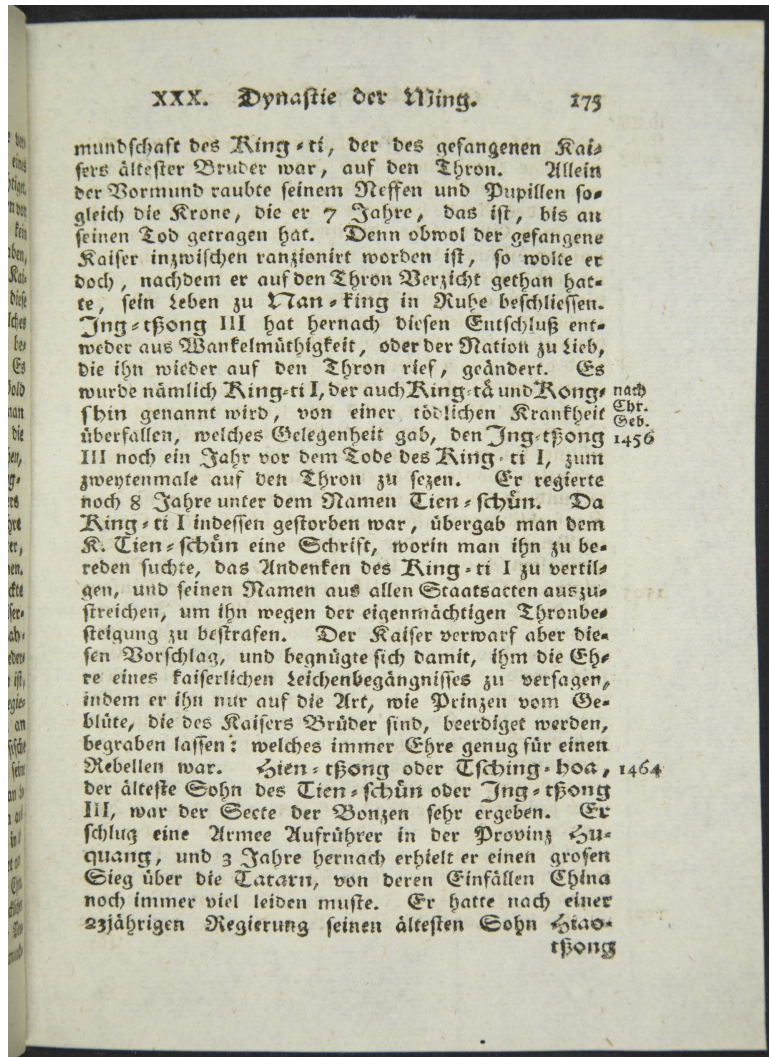
Korpus

- kuratierte Sammlung strukturierter Volltexte in Text oder XML-Format
- Metadaten für alle Texte
- stabile Erkennungsrate durch Wörterbuchabgleich
- annotiert



Negativsuchen,
für qualitative & quantitative
Methoden geeignet

Textsammlung



- Projekt „Vor der Kulturgeschichte“ (FSP „Historische Kulturwissenschaften“ der JGU)
- 11 Monographien à 200-600 Seiten
- Druckdatum vor 1850: manufakturierte Frakturtypen
- ABBYY Finereader 11

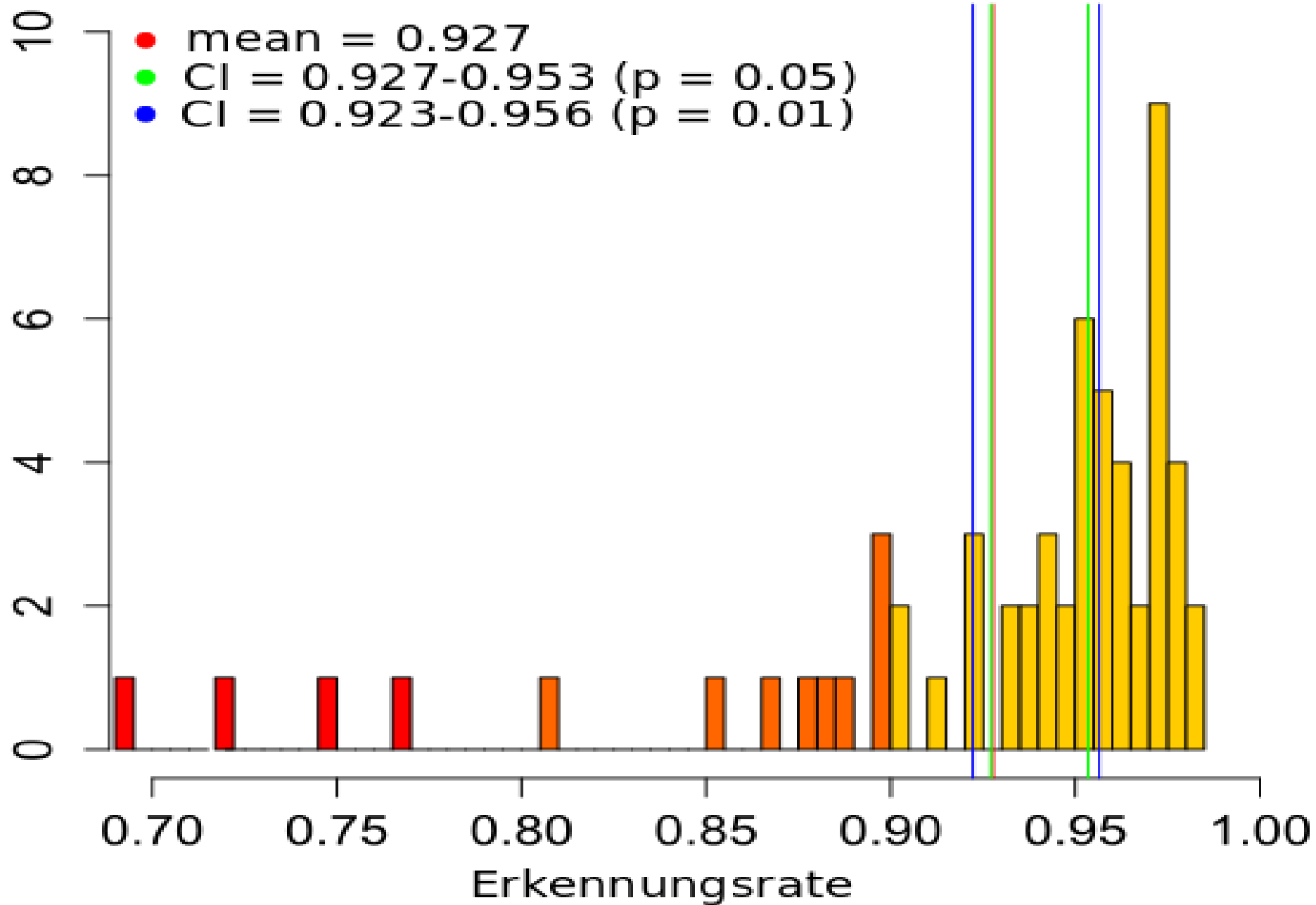
Stichprobe

- Verfahren orientiert an Stäcker (2013)
- Anpassung des Stichprobenverfahrens
 - randomisierte Stichprobe von 60 Seiten aus dem Gesamtkonvolut (83.542 Zeichen)
 - **keine** Berücksichtigung von Verschmutzungen, Durchstreichungen, Annotationen, Marginalien auf den Testseiten
 - Test unter realistischen Bedingungen

Ground Truth

- 60 Seiten manuelle Transkription durch wissenschaftliche Mitarbeiter des Projekts als Vergleichsprobe
- Herausforderung: Software zum Abgleich Stichprobe/Ground Truth finden
- Abgleich per Levenshtein-Algorithmus

Häufigkeit



Ergebnisse

veränderter Stichprobenansatz:

- Verschmutzungen, Widerdrucke, Annotationen etc. führen NICHT zu schlechteren Erkennungsraten
- fremdsprachige Abkürzungen in Kombination mit Typenwechsel und wenig bekannte Lexik/Orthographien erzeugen Erkennungsfehler
- zeigt, ob Software alltagstaugliche Ergebnisse liefert

Ergebnisse

Auswirkungen für Forschungsprojekte:

- keine phraseologischen oder Negativsuchen möglich
- Methodenmix zur wissenschaftlichen Bearbeitung des Textes nötig
 - Fehleranalyse nach der OCR zur Abschätzung kritischer Textbereiche
 - Verbindung von close & distant reading
- entstandene Volltexte sind nicht zur Bearbeitung mit quantitativen Methoden geeignet

Chancen des Post-Processing

- historische Wörterbücher
- language models
- tokenizer
- lemmatizer

} Korpuserstellung



Vielen Dank !