

IMPACT is supported by the European Community under the FP7 ICT Work Programme. The project is coordinated by the National Library of the Netherlands.

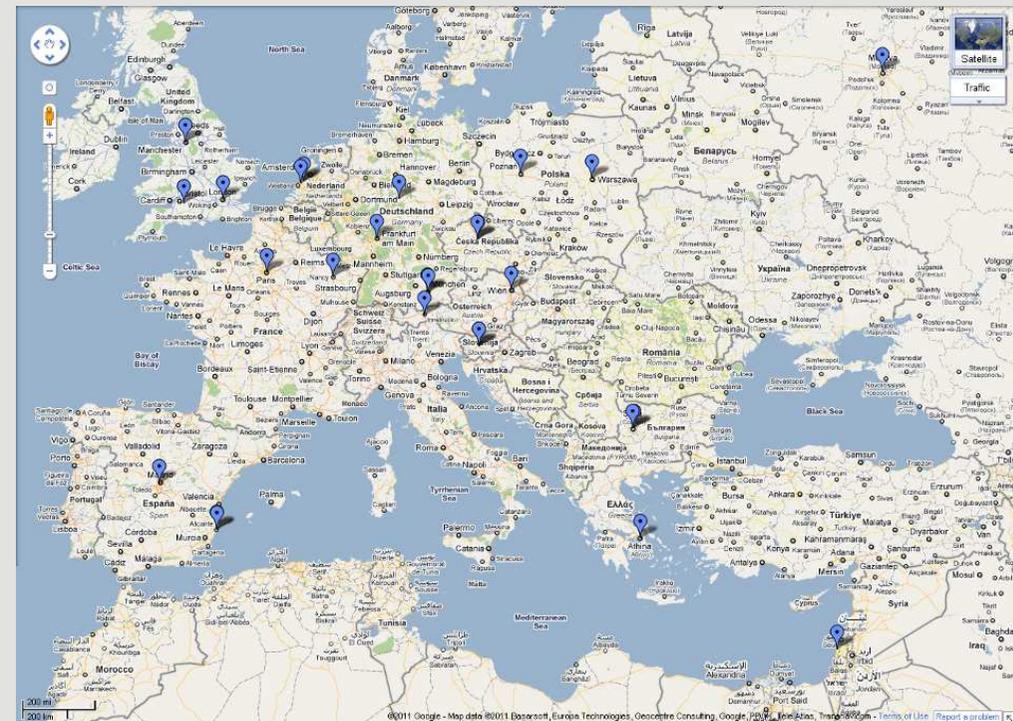
Verbesserter Zugang zu historischen Texten – Demonstration einiger Werkzeuge des Projekts IMPACT

Mark-Oliver Fischer, Münchener Digitalisierungszentrum, Bayerische Staatsbibliothek

08.06.2011 - Deutscher Bibliothekartag,
Berlin

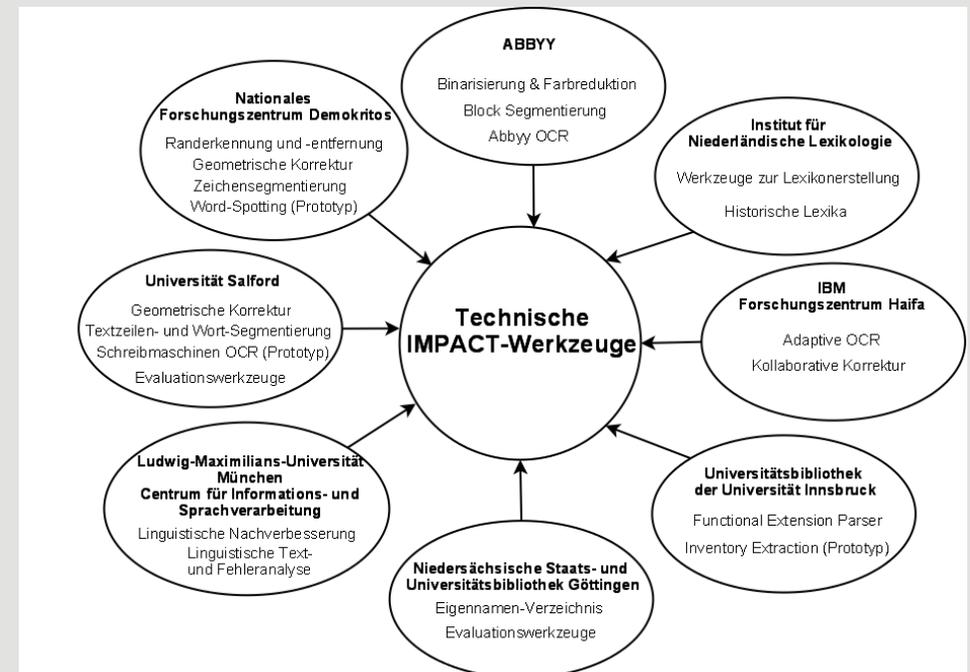
IMPACT: “IMProving ACcess to Text”

- Ziel: Elektronischen Zugang zu historischen Texten verbessern
- Unter anderem leichtere, bessere Volltexterstellung: (Weiter-)Entwicklung von OCR-Software und Software rund um die eigentliche Texterkennung
- EU-gefördertes Kooperationsprojekt
- 26 Institutionen aus 13 Ländern: Bibliotheken (z.B. BSB, DNB, BL), Forschungsinstitutionen (z.B. NCSR Demokritos, INL), kommerzielle Partner (Abbyy, IBM)



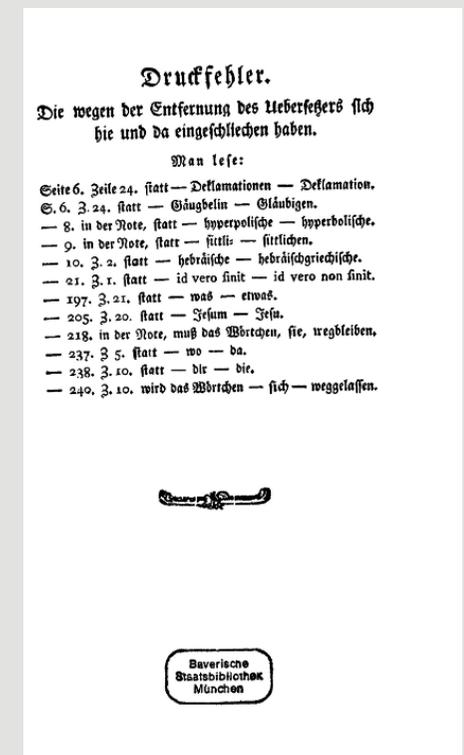
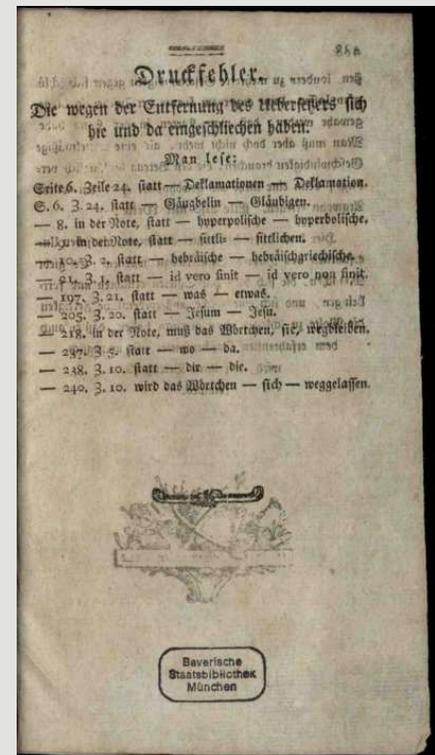
Tools für alle Prozessschritte eines OCR-Workflows

- Bildoptimierung
- Texterkennung
- Nachbearbeitung
- Nicht **eine** IMPACT-Software, sondern Zusammenspiel individueller Tools



Bildoptimierung

- Aufbereitung der Bilder für möglichst gute Texterkennung
- Nicht zwingend identisch mit optisch gutem Bild
- Tools zur geometrischen Korrektur: (NCSR Demokritos; Uni Salford)
 - Geraderücken von Seiten mit schiefem Textspiegel
 - Entzerren von Wellen, Bögen, ... im Text
 - Trennung von Doppelseiten
 - Entfernung von Bildrändern
- Verbesserte, ‚intelligente‘ Binarisierung (Abbyy)



Texterkennung

- In drei Schritten: Segmentierung, Zeichenerkennung, Wörterbuch
- IMPACT entwickelt Tools auf allen drei Ebenen
 - Verbesserte Segmentierung: Absätze, Wörter, Zeichen (Abbyy; NCSR Demokritos; Uni Salford)



- Deutliche Verbesserung der Fraktur-Erkennung (Abbyy)
- Adaptive, selbstlernende OCR (IBM)
- spezielle Maschinschrift-OCR (Uni Salford, experimentell)
- Spezifische Wörterbücher für alle Projektsprachen: Deutsch mit Schwerpunkt 16. Jahrhundert (Kooperation BSB und LMU München)

Nachbearbeitung

- Software zur Nachkontrolle, Korrektur der OCR-Ergebnisse (IBM; LMU München)
- Strukturanalyse, -auszeichnung: Druckbereich, Fußnoten, Seitenzahlen, etc. (ULB Tirol)
- Tools zur automatischen Evaluation von OCR-Ergebnissen (Uni Salford; NCSR Demokritos)

Page 14 of 42 | n --> u 1/3 x | u --> n 1/1 x

Select all | Deselect all | Correct!

was idere

föhleu -> fühlen

überdeuke -> überdenke

feinste -> feinste

Sinnen -> Sinnen

Assimiliou -> Assimilation

Der Selbstbegriff 117

...zittern, weil ich sie wahrnehme." So hängt die Objektivität dieses Satzes zunächst ab von der Objektivität der Selbstwahrnehmung. Denn wenn ich nicht erkenne, so kann ich auch nicht wahrnehmen; es fehlt also die Bedingung, unter welcher gesagt werden kann, daß Dinge existieren. Der Glaube an die Existenz von Dingen außer mir setzt also die Anerkennung meiner eigenen Existenz voraus, während ich an der Existenz der andern Dinge, wenn ich sie wahrnehme, noch immer zweifeln kann.

Über man wendet vielleicht ein: Auch die Selbstwahrnehmung hängt von der Anerkennung des Daseins anderer Dinge ab; denn mein Ich kommt sich nur als Bewußtsein des Nicht-Ich zum Bewußtsein. Dem ist zu entgegnen: Hier handelt es sich nicht um das Ich, sondern einfach um das Sein, nicht um das Was, sondern um das Daß der Existenz. Das, was ist, weiß ich nicht deshalb als Ichend, weil andere Dinge außer ihm existieren, nur daß es sich als Ich empfindet wird, hängt von der Wahrnehmung eines Nicht-Ich ab. Es wird doch niemand behaupten wollen, daß ein neugeborenes Kind, von absoluter Finsternis umgeben, und außer Berührung mit der Welt, leblos wäre, kein Bewußtsein seines Daseins hätte. Das Existenzgefühl tritt fernestwegs erst ein, wenn die Spaltung des Bewußtseins in Ich und Nicht-Ich, Subjekt und Objekt sich vollzogen hat. Das neugeborene Kind hat ohne Zweifel vom ersten Augenblick an ein Gefühl von Existenz; ob es sich aber auch vom ersten Augenblick an von der Abstraktion

68 OAN. KEPL. DE STEL. CYGNI.

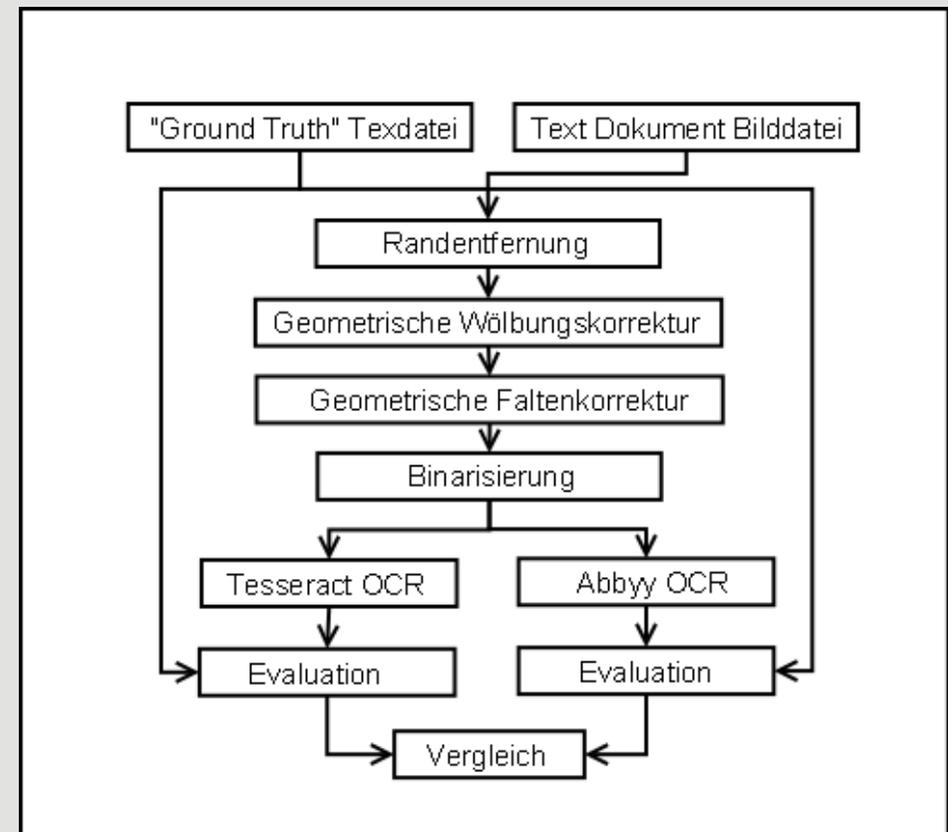
Ex quibus distantijs extruxerunt Braheani, circumspiculis omnibus locum 26. 18/ Aquarij, latit: 55. 30/ Bor. Hinc invenitur ascensio 26. 18/ Nova 300. 46/ declinatio 36. 52/ Borealis. Culminatigitur cum 28. 17/ Capric.

In Hispania parte Andalusia, in Sicilia, Peloponneso, Italia, Cilicia, Syria, caeterisque locis Terrarum, sub hoc eodem parallelo sitis, per Verticem quotidie transit. Quibus vero est altitudo Poli 57. 8/ ij horizontem stringit in Septentrione; ut Anglia, Hollandia, Brunsvigo, Marchia, Livonia, Moscovia. Ulterius versus Septentrionem non occidit.

Species Oloris polii accessum Nova. N. Novam detinet.

Interoperabilität und Modularisierung

- Ziel: Alle IMPACT-Tools können zusammenwirken, sind miteinander einsetzbar
- Dazu stehen fast alle Tools als Webservices zur Verfügung, die zu individuellen ‚Workflows‘ kombiniert werden können





Demonstration

- Stand der Deutschen Nationalbibliothek: G11
- Mittwoch und Donnerstag, je 13 bis 15 Uhr
- Weitere Informationen: <http://www.impact-project.eu/>

