Improving Access to Text

# iMPACT

# Constructing a Lexicon from a Historical Corpus

**Annette Gotscharek, Andreas Neumann, Ulrich Reffle,**

**Christoph Ringlstetter, Klaus U. Schulz**

*CIS, University of Munich*

**AACL 2009, Edmonton, Canada**

# Survey

- 1. Motivate of historical lexica in digitization projects
- 2. A Preliminary historical corpus for German
- 3. Infrastructure for lexicon building
- 4. Evaluation and future work

Improving Access to Text

IMPACT

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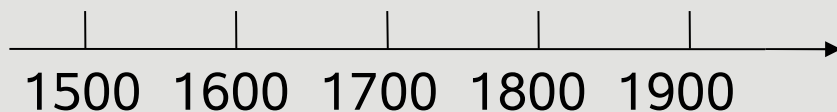
# 0. Research Context IMPACT

- IMProve ACcess to Text, Improve mass digitisation of historical books, newspapers, legal documents...
- Four year EC project started in 2008
- 15 partners: national libraries, research institutions and commercial suppliers
- (1) Image Processing
- (2) <u>Text Recognition</u>: Adaption of Optical Character Recognition to historical documents
- (3) <u>Lexicon Building</u>: Enrichment of texts to Improve IR Access to historical documents
- (4) Capacity Building for Management of Digitization Projects

Improving Access to Text

iMPACT

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

# Special Challenges for the Digitization of Historical Materials

1500  1600  1700  1800  1900

- Imaging                                                Damages on Originals
- **Optical Character Recognition**    **Rate of Recognition Errors**
- **Information Retrieval**                 **Historical Variants**
- **Human Reading**                          **Unknown Words**

# Optical Character Recognition: medium quality

Fürsten zu Gsitternwerden/wer wollte vermainen / daßwt
IhroKhurftrstl Durchl gnädigIsterHcttVatterinderpictcr
rndFrombkcltallmFürstenvorzusetzen!scyn/vnd das halst>
in^cclcQ^ vci pluz^uäzn 5accr6o5 daß tl lN KilchkN GottW
wehr als ein Priester.

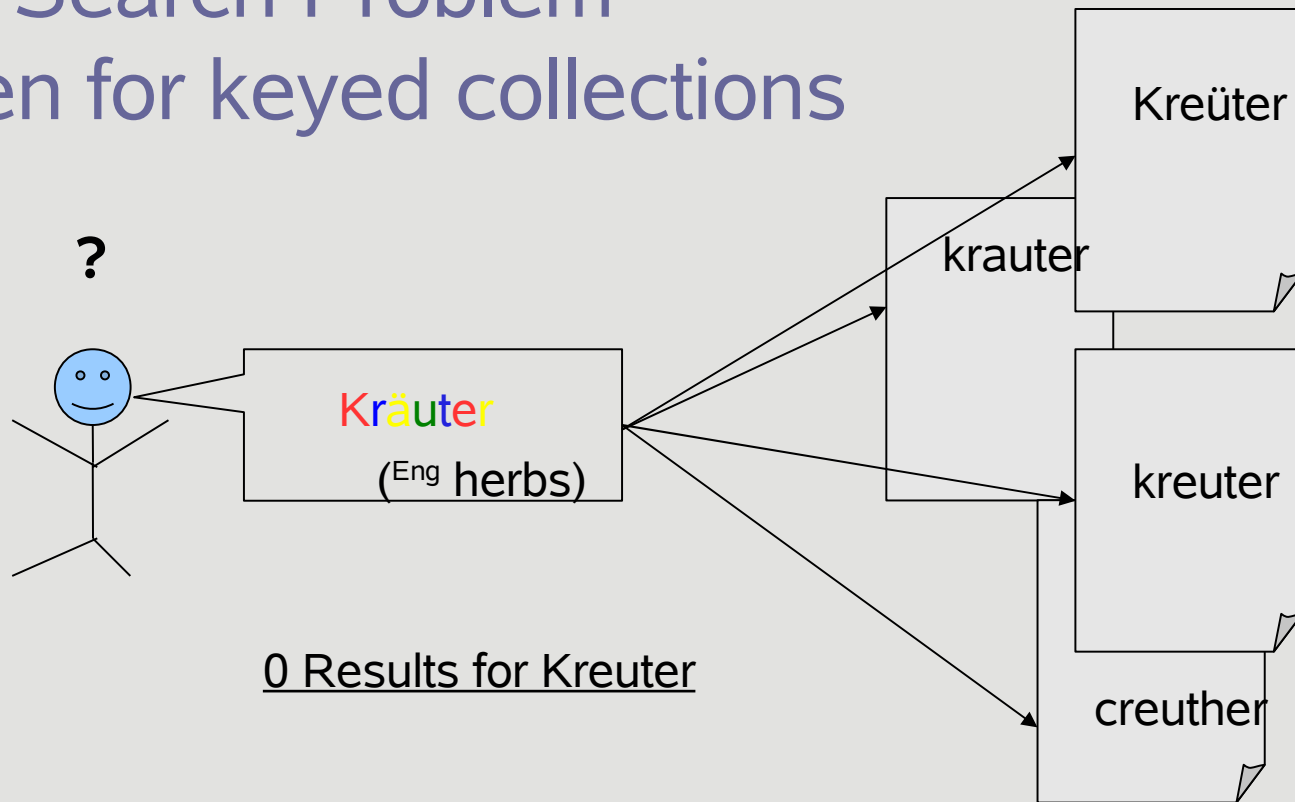# iMPACT

## Why is it so bad

(1) The image quality is challenging, further processing needed

(2) The classifiers of the OCR disregard certain type faces used in historic print

(3) **The language resources of the OCR are inappropriate: historical language**

Improving Access to Text

# iMPACT

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

# IR: Search Problem
# even for keyed collections



**?**

Kräuter

(Eng herbs)

Kreüter

krauter

kreuter

creuther

0 Results for Kreuter

− kräuter (= Engherbs) as kra،uter, kreuter, kreüter, kreuter, creuther

# Special challenge for IR and OCR on historical Texts: Spelling variation

- Missing nomalization of orthography leads to many variants in historical documents, e.g. in German texts (1500-1850):
    - teil (= [Eng]part) as theil, teyl, theyl
    - fragte (= [Eng]asked) as frug, fruk
    - statt (= [Eng]instead) as statt (=[Eng]town)
    - ringer (= [Eng]wrestler) as ringer (=[Eng] less)
- Resource: a lexicon mapping from historical variants to modern lemma

Improving Access to Text

IMPACT

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

# Language Resources need a Corpus Base

Possible sources:

- Keyed Materials on the Web
- Non Public Electronic Corpora
- Keying/corrected OCR of Image Corpora
- Noisy OCR Corpora

# iMPACT

# Status of German historical Corpora

## 1. First development corpus
- Proofread texts from 1400 to 1900: 4 institutional sources
- Medium size: 2.7 Mill. tokens
- For lexicon construction
- For preliminary diachronic analysis

## 2. IR annotated corpus for strategy testing
Texts from 16th, 17th, 18th, 19th century 31080 tokens

## 3. Project for corpus building 16th century started

## 4. Agreements on additional background materials from 1750s with 30 million tokens literature and 10 million tokens other

# Language in a historical corpus for German

Modern lexicon (CISLEX): coverage for 10 time periods

Improving Access to Text
# IMPACT

IMPACT is supported by the European Community under the FP7 ICT Work Programme. The project is coordinated by the National Library of the Netherlands.
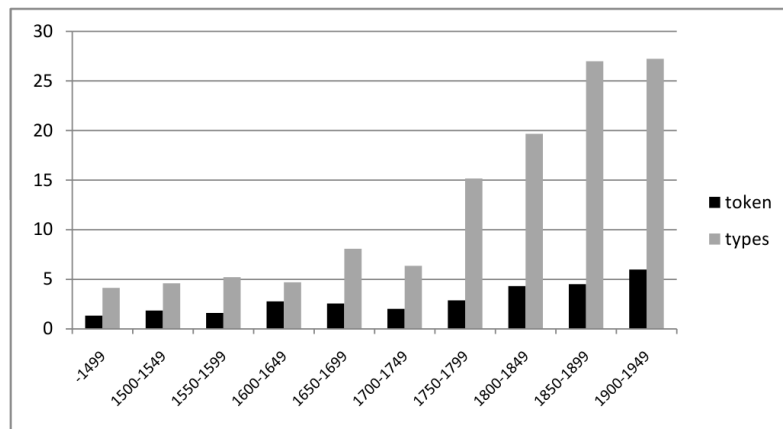
# Language in a historical corpus for German

Pattern based  lexicon; coverage for 10 time periods
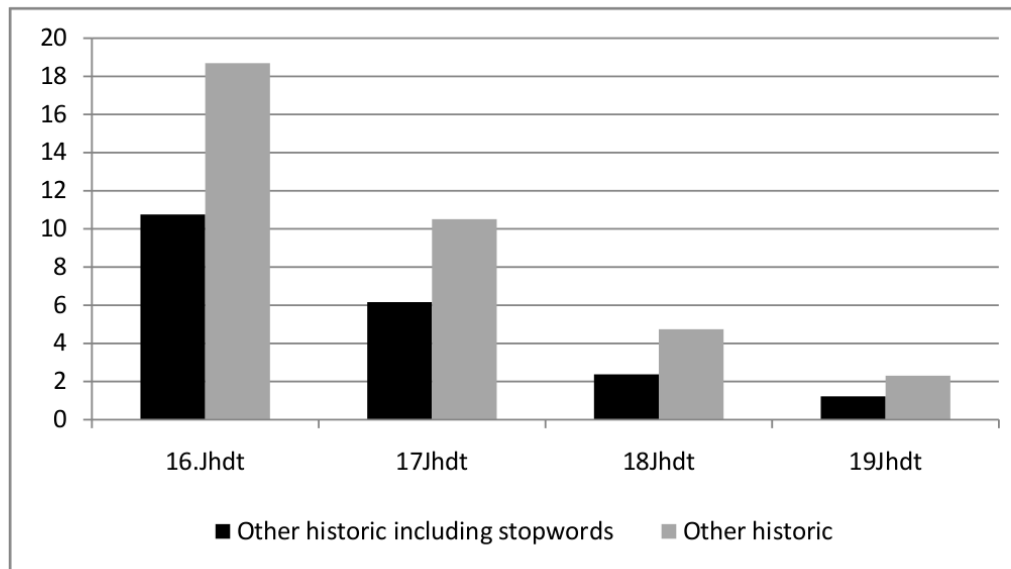
# Language in a historical corpus for German

Compounds (modern components); coverage for 10 time periods

Improving Access to Text

# iMPACT

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

Test corpus:

potential of  additional non-pattern based historical lexicon

Improving Access to Text

# iMPACT

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

# Tool for collaborative web based lexicon building

- Goal: lexicon that assigns modern lemmas to historical full forms
- Corpus management component
- Two kinds of historical words: pattern based and other
- Modern lemmas and POS information for historical full forms automatically suggested by system
- Collaborative handling of difficult cases
- Workflows: corpus frequencies or on single documents

# Main Background Ressource

# Approximate matching procedure

**Modern lexicon**

| **Inflected forms** | **Lemmatizing information** |
|---|---|
| … | … |
| teile | teil (= part) |
|  | teilen (= to share) |
| ... |  |
| taille | taille (= waist) |
| fragte | fragen (= to ask) |
| … | … |

Improving Access to Text

IMPACT

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

# Approximate matching procedure

| ~ 140 Patterns | Modern lexicon | |
|---|---|---|
| | **Inflected forms** | **Lemmatizing information** |
| … | … | … |
| th → t | teile → teil (= part) | |
| ei → ai | teile ↘ teilen (= to share) | |
| ey → ei | ... | |
| l → ll | taille → taille (= waist) | |
| … | fragte → fragen (= to ask) | |
| | … | … |

# Approximate matching procedure

| Spelling variation | ~ 140 Patterns | Modern lexicon |
|---|---|---|
| | | **Inflected forms** / **Lemmatizing information** |
| | … | … / … |
| | th → t | |
| theile | ei → ai | teile → teil (= part) |
| | ey → ei | → teilen (= to share) |
| | l → ll | ... |
| | … | taille → taille (= waist) |
| | | fragte → fragen (= to ask) |
| | | … / … |

Improving Access to Text
**IMPACT**

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

# Approximate matching procedure



**Spelling variation**

**~ 140 Patterns**

…

th → t

ei → ai

ey → ei

l → ll

…

theile

**Modern lexicon**

**Inflected forms**

**Lemmatizing information**

…                    …

teile ⟶ teil (= part)

…        teilen (= to share)

taille ⟶ taille (= waist)

fragte ⟶ fragen (= to ask)

…                    …

# Approximate matching procedure



**Spelling variation**

**~ 140 Patterns**

…

th → t
ei → ai
ey → ei
l → ll

…

theile

**Modern lexicon**

**Inflected forms**

**Lemmatizing information**

…

teile ———→ teil (= part)
———→ teilen (= to share)

…

taille ———→ taille (= waist)

fragte ———→ fragen (= to ask)

…

…

Improving Access to Text

IMPACT

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

# Lexicon Tool Corpus Mode: User selects a word

# Pattern Variants created automatically: User selects

## Interpretations for string "theile" :

☑ **theile** matches **teile** . Applied Patterns: t→th at position 0

☑ **theile** matches **taille** . Applied patterns: t→th at position 0; ai→ei at Ppsition 1; ll→l at position 3

[ Confirm matches ]

●

## Add token to special list

Classify **theile** as:

Historic word without modern equivalent [ Add ]

Historic abbreviation [ Add ]

Pattern matcher failed [ Add ]

Named Entity [ Add ]

Missing in modern lexicon [ Add ]

●

Improving Access to Text

# IMPACT

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

# Lemma Readings created automatically: User confirms

Improving Access to Text

# IMPACT

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

# Lexicon Tool Lemma Attestations: User selects



## Choose attestations for wordform "theile" (Noun,neut.)

Add Attestaions

theile

1556 : Vonn warer / wesenlicher / vnd pleibēder Gegenwertigkeit des Leybs und Blůts Christi (...). durch Johannem Gropperum d. Archidiacō der h. Kirchen zu Cȏllen [display text]

☐ 59780: der einigen Person CHRISTI mehe Personen/ oder ye fill verscheiden theile. Nů wissen wir aber vß dem lieben H. Johanne / das welcher Christ

1668 : Spiegel der Ehren des (...) Erzhauses Oesterreich (...) 1212 anfahend (...) 1519 sich endend. Erstlich vor mer als C Jahren verfasset durch (...) Johann Jacob Fugger (...) nunmehr aber (...) aus dem Original neu-ůblicher ůmgesetzet (...) erweitert (...) durch Sigmund von Birken [display text]

☐ 43264: re/ gienge A. 1272 der Lármen wieder an/ und hatten sich beyde theile aufs neue zum Krieg gerůstet. Der Bischoff/ die Macht Rudolphi zuschw

1752 : Die Sitten der americanischen Wilden im Vergleich zu den Sitten der Frühzeit [Text zeigen]

☐ 42663: und Trauer zu handeln. Die Arzeneikunst für ihre Krankheiten theile ich in zween Theile, nemlich in die natürliche und unnatürliche, weni

1757 : Vorkritische Schriften II 1757-1777 [display text]

☐ 471216: merkt bald, daß diese ehrwürdige Gesellschaft sich in zwei Logen theile, in die der Grillenfänger und die der Gecken. Ein gelehrter Grillenfä

☐ 597633: Vorsorge annimmt und um welcher willen sie Verfügungen macht. Ich theile diese Krankheiten zwiefach ein, in die der Ohnmacht und in die der Verk

| theile (Noun, neut.) | | |
| --- | --- | --- |
| PreContext | Wordform | frequency |
| das | theil | 1673 |
| des | theils | 668 |
| des | theiles | 17 |
| dem | theil | 1673 |
| dem | theile | 641 |
| das | theil | 1673 |
| die | theile | 641 |
| der | theile | 641 |
| den | theilen | 372 |
| die | theile | 641 |
| Score: 3371 | | |

Improving Access to Text

**iMPACT**

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

# Resulting Lexicon Structure

**Improving Access to Text**

# iMPACT

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

# Additional View: Lexicon Tool Document Mode

Improving Access to Text

# IMPACT

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

# Resources built

- **Validated lexicon**
  15,000 entries ("kernel lexicon for hist. German")
- **Improved hypothetical lexicon**
  Set of observed transition patterns with time stamps and frequency information from automatic attestations enables a targeted hypothetical lexicon and quantitative investigations on language change.

# Pattern based IR lexicon: precision and recall

Results on IR corpus, manual evaluation of 31,745 tokens

|            | 16th century | 17th century | 18th century | 19th century |
|------------|--------------|--------------|--------------|--------------|
| Prec       | 71.84%       | 84.22%       | 94.76%       | 97.68%       |
| **Prec(-stop)** | **62.09%** | **76.49%** | **91.09%** | **95.78%** |
| Rec        | 84.70%       | 86.60%       | 94.25%       | 94.95%       |
| **Rec(-stop)** | **76.16%** | **78.69%** | **89.69%** | **90.26%** |
| FF         | 6.07%        | 4.80%        | 1.57%        | 0.29%        |
| **FF(-stop)** | **8.47%** | **6.92%** | **1.52%** | **0.55%** |

Improving Access to Text

**IMPACT**

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

# Project Status

- Lexicon work continues

- Benefits for IR and OCR already proven

- New corpus especially 16$^{th}$ century integrated by end of 2009

- CL: study on historical transformation patterns

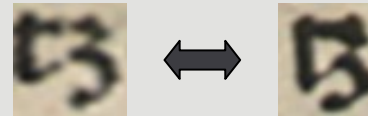- CL: study on German compound words
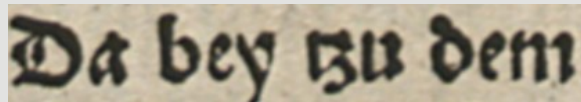
- Work on English starts 2010

# Corpus Extension

*Iterative Process with Bavarian State Library to*
*Create More Resources for Early High German*

(1) A random selection of 200 pages from 100 sources

(2) Selection of usable sources

(3) Specification of keying by BSB/CIS for 70 complete books usable for both presentation and linguistic resources building: 500,000 tokens

(4) Contract with service providers

Improving Access to Text

IMPACT

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

# Next Step Corpus Extension: difficulties

From the service provider:
„Our question is whether there exists a difference between tz that looks normal and the form that looks more similar to ß.“