

# Constructing a Canonicalized Corpus of Historical German by Text Alignment

Bryan Jurish,  
Marko Drotschmann,  
Henriette Ast

*Deutsches Textarchiv*  
*Berlin-Brandenburgische Akademie der Wissenschaften*  
<http://deutsches-textarchiv.de>

*New Methods in Historical Corpora*  
Manchester, 29<sup>th</sup>-30<sup>th</sup> April, 2011

# Overview

## The Big Picture

- Canonicalization
- Desiderata
- Proposal

## Construction

- Sources
- Text Alignment
- Manual Annotation

## Applications

- Test Corpus
- Canonicalization Lexicon

## Conclusion

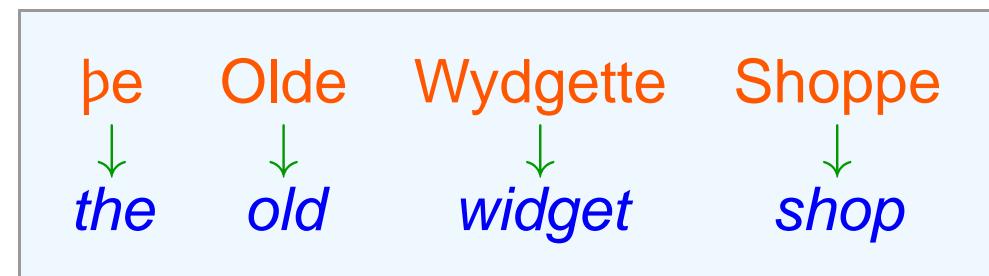
# — The Big Picture —

# Canonicalization

a.k.a. (orthographic) ‘standardization’, ‘normalization’, ‘modernization’, ...

## The Problem

- Historical text  $\not\models$  orthographic conventions
- Conventional NLP tools  $\Rightarrow$  strict orthography
  - ▶ *Fixed lexicon* keyed by **orthographic form**
  - ▶ *Extant* lexemes only



## The Approach

- *Map* each word  $w$  to a unique **canonical cognate**  $\tilde{w}$ 
  - ▶ Synchronously active “extant equivalents”
  - ▶ Preserve both *root* and **relevant features** of input
- *Defer* application analysis to canonical forms

# Desiderata

## Evaluation

- **Compare** various canonicalization functions  $c(\cdot)$
- **Task:** information retrieval  $\implies (\text{precision}, \text{recall})$
- **Retrieval via canonical equivalence:**  
 $\text{retrieved}_c(w) := [c \circ c^{-1}](w) = \{v : c(v) = c(w)\}$
- **Relevance requires manual verification!**  
 $\text{relevant}(w) := ?$

## Ground-Truth Corpus

- Manually verified canonicalization pairs  $(w, \tilde{w})$
- “Gold standard”  $\widehat{c}(\cdot)$  for training & evaluation  
 $\text{relevant}(w) := \text{retrieved}_{\widehat{c}}(w) = \{v : \tilde{v} = \tilde{w}\}$
- Minimize manual annotation effort

*... but how?*

# Proposal

## Intuitions

- Contemporary editions of historical works  
     $\Rightarrow$  **already standardized**
- Expect mostly identity canonicalizations  $w = \tilde{w}$   
*(at least for 18<sup>th</sup>-19<sup>th</sup> century German)*

## Construction (Sketch)

- **Align** historical text with a contemporary edition
  - ▶ maximize identity alignments
- **Confirm** or **Reject** type-wise alignments
  - ▶ exploit Heaps' Law
- **Manually annotate** only unconfirmed tokens
  - ▶ don't lose "interesting" anomalous material

# — Construction —

## Text Resources

- Source texts: *Deutsches Textarchiv (DTA)*
  - ▶ *Belles lettres*, drama, verse, philosophy
- Target texts: [gutenberg.org](http://gutenberg.org), [zeno.org](http://zeno.org)

## Prototype Corpus

- 13 volumes, published 1780–1880
- ca. 350k tokens ~ 28k types *(words only)*

## Ongoing Construction ('full' corpus)

- 129 volumes, published 1780–1901
- ca. 5.2M tokens ~ 219k types *(words only)*

# Text Alignment

## Preprocessing

- Tokenization (1 word / line)
- Transliteration e.g. ( $f \mapsto s$ ), ( $\ddot{o} \mapsto \ddot{o}$ )

## Basic Alignment

- Token-wise LCS (GNU diff)
- > 77% identity, > 94% transliterated identity

## Heuristic Alignment

- For each change change hunk
  - ▶ multi-token alignments e.g. (*zwei und vierzig*  $\mapsto$  *zweiundvierzig*)
  - ▶ character-wise ‘best’ match (Levenshtein)

# Type-wise Confirmation

## Idea

- Manually **confirm** or **reject** non-identity alignments
- Exploit Heaps' Law
  - ▶ vocabulary grows *logarithmically* with corpus size
- Conservative acceptance only

## Results (prototype corpus)

- Available: 18k tokens ~ 5.8k types
- Confirmed: 16k tokens (90%) ~ 4.5k types (77%)

## Throughput

- ca. 3.95 seconds / pair  $\approx$  15 words / second

# Token-wise Annotation

## Idea

- Resolve remaining uncanonicalized tokens (ca. 2%)
- Retain anomalous canonicalization patterns

## Preprocessing Filters

- Block pruning ( $\approx 2.2\%$ )
- Closed-class lexicon

## Annotations

- Canonical form + administrative flags
- Expert review for problematic cases
- Throughput (total)  $\approx 1.3$  words / second

# — Experiments —

# Materials

## Prototype Corpus $\rightsquigarrow$ Ground-Truth Relevance

- Most thoroughly annotated corpus subset
- 340k tokens; 29k types (words only)

## Full Corpus $\rightsquigarrow$ Canonicalization Lexicon (LEX)

$$\text{LEX}(\textcolor{orange}{w}) = \begin{cases} \arg \max_{\tilde{w}} f(\textcolor{orange}{w}, \tilde{w}) & \text{if } f(\textcolor{orange}{w}) > 0 \\ \textcolor{orange}{w} & \text{otherwise} \end{cases}$$

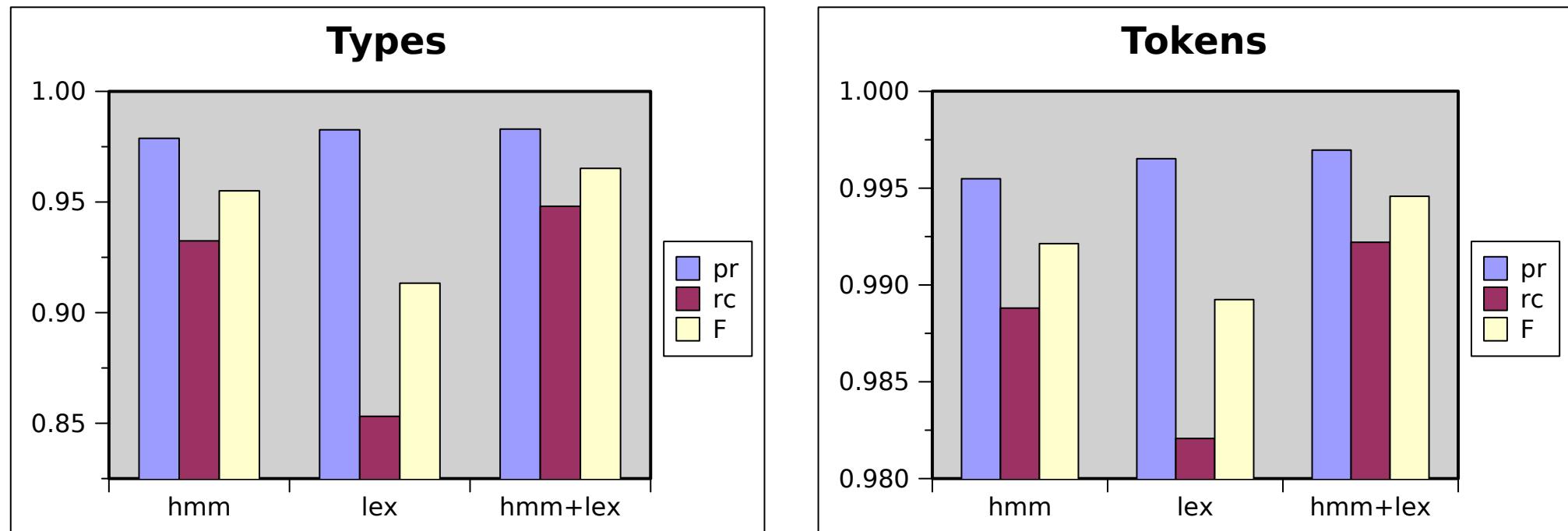
- Strictly disjoint from test corpus (by author)
- Partially annotated (no expert review)
- 2.4M tokens; 140k types (words only)

## HMM Canonicalization Cascade

(Jurish, 2010c)

- Robust finite-state canonicalizer
- Tested methods: ID, LEX, HMM, HMM+LEX

# Results



	% Types			% Tokens		
	pr	rc	F	pr	rc	F
ID	98.3	57.1	72.2	99.7	79.1	88.2
HMM	97.9	93.2	95.5	99.5	98.9	99.2
LEX	98.3	85.3	91.3	99.7	98.2	98.9
HMM+LEX	98.3	94.8	96.5	99.7	99.2	99.5

# Conclusion

## Construction

- Alignment with contemporary edition
- Type-wise confirmation
- Token-wise annotation

~*minimal-effort corpus bootstrapping*

## Applications

- Simple corpus-based lexicon ⇒ surprisingly effective
    - ▶ very high precision
    - ▶ mediocre recall for unknown types (sparse data)
  - ‘Exception’ lexicon for HMM canonicalizer
    - ▶ best overall performance
    - ▶ corpus-based and generative techniques
- complement one another*

# þe Olde Laste Slÿde ("The End")

***Thank you for listening!***

# — Addenda —

# Token Annotation GUI

DTA::EvalCorpus::Editor::Gtk2: phase3.corpus\_01.xml \*

File Edit Go History View Help

Open Save Back Forward Previous Next Apply Quit

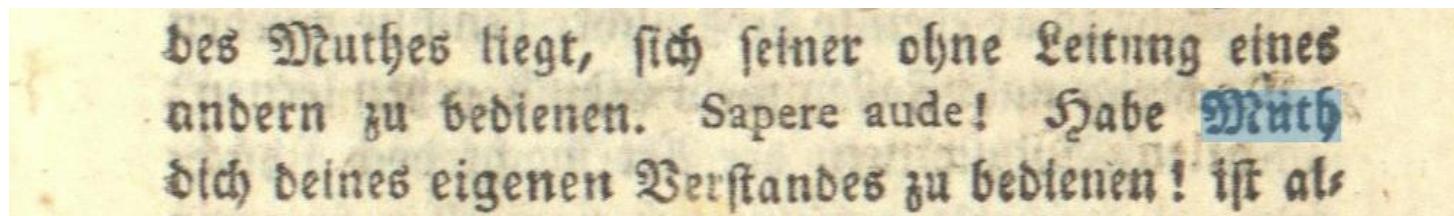
Context  
Sapere aude !  
Habe **Muth** dich deines eigenen Verstandes zu bedienen !  
ift alſo der Wahlspruch der Aufklärung .

General  
**src:** kant\_aufklaerung\_1784.phase3.q001.xml  
**wid:** w838  
**sid:** s138  
**page:** 17 (56%)  
**u:** Muth  
**old:** Muth  
**cur:** Mut  
**new:** Mut  
**ext:**

Flags  
LEX  
 Valid Sentence  
 Valid Token  
 Review Required  
 Notes

Actions  
CAB Query  
Browse Source  
View XML  
Copy Original  
Apply

Changed to document data\_01/kant\_aufklaerung\_1784.phase3.q001.xml 100%



des Muthes liegt, sich seiner ohne Leitung eines  
anderen zu bedienen. Sapere aude! Habe **Muth**  
dich deines eigenen Verstandes zu bedienen! ist als .

# GUI: Batch Editor Window

DTA::EvalCorpus::Editor::Gtk2: Batch Editor

Search XPath: `./w[@old="Muth"]`

Global Search

n	class	left	old	right	base
1	LEX	... durch Vernunft , frohen	<b>Muth</b>	und guten Willen zu über...	goethe_lehrjahre04_1796
2	LEX	...Alter einen fo geruhigen	<b>Muth</b>	und eine fo schöne Nacht...	brentano_kasperl_1838
3	LEX	Habe	<b>Muth</b>	dich deines eigenen Verf...	kant_aufklaerung_1784
4	LEX	... Sprache , und hatte den	<b>Muth</b>	nicht ihren Augen zu beg...	goethe_lehrjahre04_1796
5	LEX	Er sprach ihnen	<b>Muth</b>	ein , und feine Gründe w...	goethe_lehrjahre02_1795
6	LEX	...flecktes Ziel Mit frohem	<b>Muth</b>	und strengem Fleiß errei...	goethe_torquato_1790
7	LEX	... , ihm ganz ängstlich zu	<b>Muth</b>	zu werden , daß es ihr n...	spyri_heidi_1880
8	LEX	" - " - " - "	<b>Muth</b>	" - " - " - "	" - " - " - "

Eval Perl Code:

# Administrivia

Class	N	%Edited
LEX	2684	59.22 %
NE	874	19.29 %
JOIN	792	17.48 %
GRAPH	101	2.23 %
SPLIT	72	1.59 %
BUG	40	0.88 %
GONE	8	0.18 %
FM	1	0.02 %

