



DTA::CAB – a Field Spotter’s Guide

Bryan Jurish

jurish@bbaw.de

Universität Graz, Zentrum für Informationsmodellierung

19th March, 2019

The Big Picture

- The Situation
- The Problem
- The Approach

Canonicalization Methods

- Type-wise Conflation
 - ▶ Lexicon, Transliteration, Phonetization, Rewrite Cascade
- Token-wise Disambiguation
 - ▶ Dynamic Hidden Markov Model
- Problems & Workarounds

Experiments

- Generative Canonicalization Methods
- Alignment-based Lexicon

— The Big Picture —

DFG Project

(2007–2016)

- digitize ~ 1,300 print volumes, printed ~ 1600-1900
 - ▶ first editions, detailed metadata, highly accurate transcriptions
- TEI-XML corpus encoding & storage
 - ▶ DTABf dialect <http://www.deutschestextarchiv.de/doku/basisformat>
- linguistic analysis (automated)
 - ▶ tokenization, normalization, PoS-tagging, lemmatization
- online search (DDC) <http://www.ddc-concordance.org>
 - ▶ lemma-based, PoS-sensitive, spelling-tolerant

- ongoing maintenance & support through CLARIN-D

(2017–present)

In Numbers

(2019-03-11)

4,425	transcribed works
+ 603	additional works in DTAQ
214,549,119	tokens
3,305,283	surface types
958,331	digitized pages
1,918,576,860	unicode characters

<http://www.deutschestextarchiv.de>

Anmelden (DTAQ) DWDS dlexDB CLARIN-D

DTA Werke im Deutschen Textarchiv Twitter Google+ Texte ▼ Projekt ▼ Dokumentation ▼ Impressum

in den Titeldaten im Korpus in der Dokumentation [Hilfe](#)

Beispielfragen: {Gott,Herr} && "das ewige Leben #2 geben" ehelichen with \$p=VVINF "Auge um" && "Zahn um"

Deutsches Textarchiv

GRUNDLAGE FÜR EIN REFERENZKORPUS DER NEUHOCHDEUTSCHEN SPRACHE

Das Deutsche Textarchiv stellt einen disziplinen- und gattungsübergreifenden Grundbestand deutschsprachiger Texte aus dem Zeitraum von ca. 1600 bis 1900 bereit. Die Textauswahl erfolgte auf der Grundlage einer von Akademiemitgliedern erstellten und ausführlich kommentierten, umfangreichen Bibliographie. In Ergänzung wurden einschlägige Literaturgeschichten und (Fach-)Bibliographien ausgewertet. Aus der Gesamtliste der auf diesem Wege ermittelten Titel wurde von der DTA-Projektgruppe ein hinsichtlich der repräsentierten Textsorten und Disziplinen ausgewogenes Korpus zusammengestellt (weitere Informationen zur Textauswahl).

Um den historischen Sprachstand möglichst genau abzubilden, werden als Vorlage für die Digitalisierung in der Regel die Erstausgaben der Werke zugrunde gelegt. Das elektronische Volltextkorpus des DTA ist über das Internet frei zugänglich und dank seiner Aufbereitung durch (computer-)linguistische Methoden schreibweisentolerant über den gesamten jeweils verfügbaren Bestand durchsuchbar. Sämtliche Texte stehen zum Download zur Verfügung.

[mehr ...](#)

NEUIGKEITEN AUS DEM PROJEKT

Das DTA auf der LREC 2016


Das DTA in Zahlen

- 2 438 Werke
- 593 008 digitalisierte Seiten
- 139 738 217 fortlaufende Wortformen
- 991 255 987 Zeichen (Unicode)
- 447 weitere Werke in DTAQ

Das DTA am 14. Juli 2016



Am 14. Juli 1789 fand der *Sturm auf die Bastille*, Symbol für die französische Revolution, statt. Die Bastille war ein Gefängnis in Paris und galt bei den Revolutionären als Sinnbild für das Ancient Regime. Ziel der bewaffneten Demonstrantenmenge waren außerdem die dort gelagerten Munitionsvorräte. Bereits zwei Tage nach der Erstürmung begannen Abrissarbeiten, die sich über ein Jahr hinzogen. Noch heute ist der 14. Juli der Nationalfeiertag Frankreichs. Friedrich Christoph Dahlmann beschreibt die Vorgänge in seiner „Geschichte der französischen Revolution“.



Dahlmann, Friedrich Christoph: Geschichte der französischen Revolution bis auf die Stiftung der Republik. Leipzig, 1845.

Historical Text \neq Orthographic Conventions

- also applies to OCR text, E-Mail, SMS, Tweets, ...
- High variance of graphemic forms

fröhlich
"joyful"

frölich, fröhlich, vrölich, frœlich, frôlich, fröhlich,
vrölich, fröhlig, frölig, ...

Herzenleid
"heart-sorrow"

hertzenleid, herzenleit, hertzenleyd, hertzenleidt,
herzenlaid, hertzenlaid, hertzenlaydt, ...

Conventional NLP Tools \implies Strict Orthography

- IR systems, PoS taggers, lemmatizers, morphological analyzers, ...
- **Fixed lexicon** keyed by (ortho)graphic form
 - ▶ **Extant** lexemes only \leadsto **OOV failures** (tool lexicon)
 - ▶ **Attested** lexemes only \leadsto **sparsity & noise** (corpus lexicon)

	Conventional	Tools
⊕	Historical	Corpus
=		<i>Soup</i>

- Corpus variants *missing* from application lexicon
 - ▶ *low coverage* (many unknown types)
 - ▶ *poor recall* (relevant data not retrieved)
 - ▶ *spurious “noise”* (poor model fit)
 - ▶ *... and more!*

The Approach: Canonicalization

a.k.a. (orthographic) 'standardization', 'normalization', 'modernization', ...

þe	Olde	Wydgitt	Shoppe
↓	↓	↓	↓
the	old	widget	shop

In a Nutshell

- *Map* each **word** w to a unique **canonical cognate** \tilde{w}
- *Defer* application analysis to canonical forms

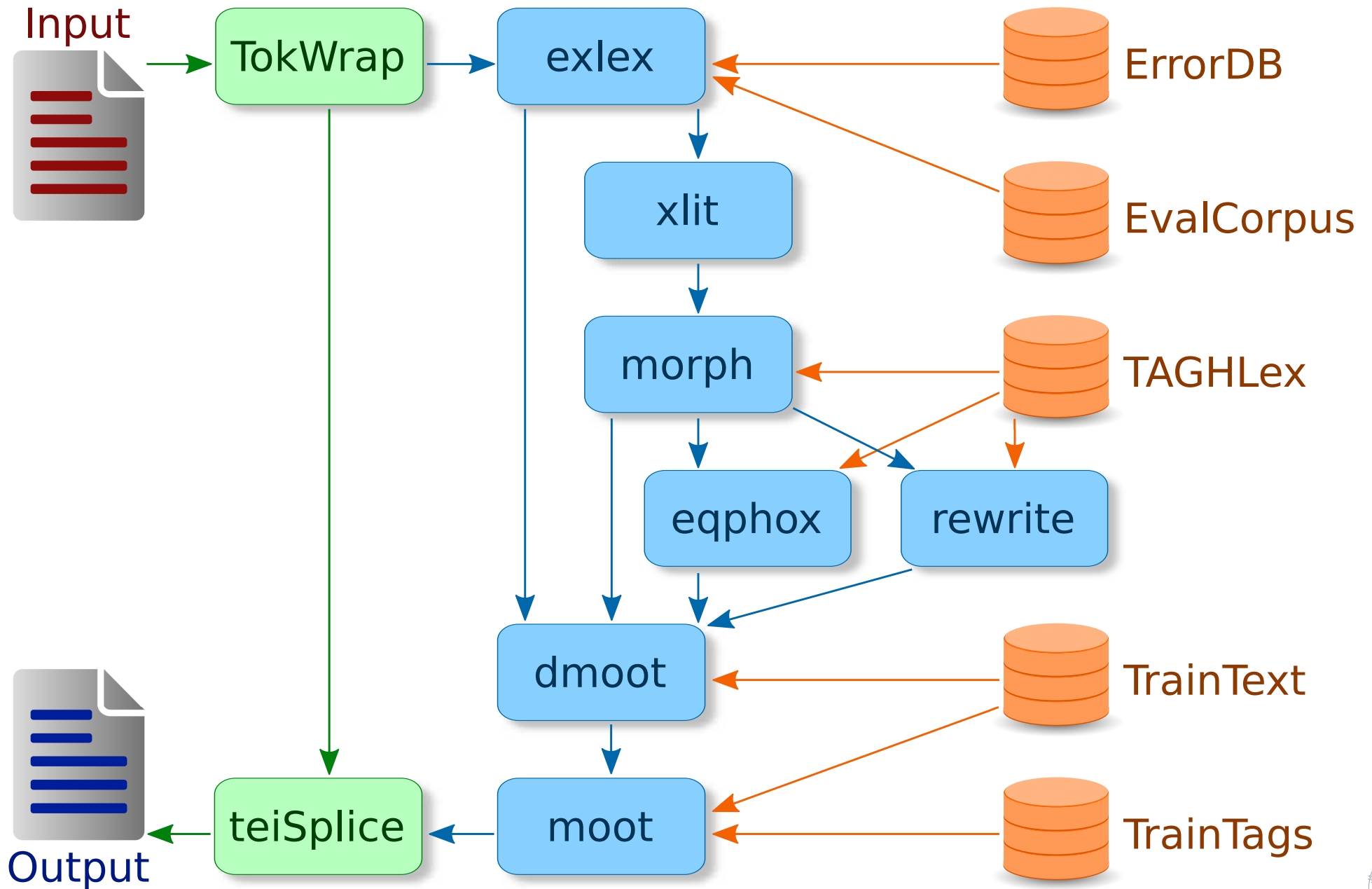
Canonical Cognates

- Synchronically active “**extant equivalent(s)**” $\tilde{w} \in \text{Lex}$
- Preserve both **root** and **relevant features** of input

Conflation Relation \sim_r

- *Binary relation* on strings (words) in \mathcal{A}^* ($\sim_r \subseteq \mathcal{A}^* \times \mathcal{A}^*$)
- Prototypically a true *equivalence relation* (reflexive, transitive, symmetric)
 - ▶ $w \sim_r v$ holds iff w and v are conflated by relation \sim_r
 - ▶ equivalence class $[w]_r = \{v \in \mathcal{A}^* \mid w \sim_r v\}$

DTA::CAB System Architecture



— Afternoon Workshop —

<http://kaskade.dwds.de/~jurish/cab-workshop/>

— Canonicalization Methods —

... and where they fail

garbage in \rightsquigarrow *garbage out*

Graphical Errors

(cf. Haaf, Wiegand & Geyken, 2013)

Print	wit \mapsto *wett	\neq mit	“with”
Transcription	vft \mapsto *fest	\neq und (v̄n)	“and”
Whitespace	daßfie \mapsto *Dash sie	\neq daß sie	“that she”

Tokenization Errors

(cf. Jurish & Würzner, 2013)

Abbrev	Durchl[.] \mapsto *torkle	\neq Durchl.	“highness”
	Superint[.] \mapsto *Super rind	\neq Superint.	“superintendent”
	vnterthän[.] \mapsto *Unter teen	\neq untertän.	“subserviently”
Space	ists \mapsto *Ists _{NN.gen}	\neq ist es	“is it”
	ers \mapsto *Airs	\neq er es	“he it”
Newline	tau[\n]sent \mapsto *Tau sehnt	\neq tausend	“thousand”

Sketch

$$\text{exlex} : \mathcal{A}^* \rightarrow \mathcal{A}^* : w \mapsto \tilde{w}$$

- Finite, deterministic type-wise mapping
- Overrides all other canonicalizers

Sources

(Jurish, Drotschmann & Ast, 2013)

- Alignment-based corpus <http://odo.dwds.de/~jurish/software/dtaec>
 - ▶ 126 volumes (1780-1901), 5.6M tokens, 212k types
- Online error database <http://kaskade.dwds.de/demo/caberr>
 - ▶ includes rudimentary inflection grammar
 - ▶ 36k database records \rightsquigarrow 437k surface mappings

Weaknesses

(cf. Kempken et al., 2006; Gotscharek et al., 2009b)

- can't handle any *ambiguity*
- can't handle *productive morphological processes*
- alignment-based bootstrapping

(Jurish & Ast, 2015)

\rightsquigarrow bogus identity alignments ($w \mapsto w$)



- EvalCorpus assumption: modern edition \implies strict orthography
- Implicitly accepted **identity pairs** ($w \mapsto w$)
 - ▶ ca. 59% types, 87% tokens identical modulo transliteration
- Not always justified by the editions used

(oops)

Letter Case	bruder \mapsto *bruder \neq Bruder	trost \mapsto *trost \neq Trost	“brother” “comfort”
Extinct Forms	ward \mapsto *ward \neq wurde	däuchte \mapsto *däuchte \neq dünkte	“was” “seems”
Prosodic Foot	andre \mapsto *andre \neq andere	eigenen \mapsto *eigenen _V \neq eigenen _{ADJ}	“other” “own”
Dialect	kömmt \mapsto *kömmt \neq kommt	nich \mapsto *nich \neq nicht	“comes” “not”
Apostrophes	in’s \mapsto *in’s \neq ins	s’ist \mapsto *s’ist \neq es ist	“into the” “it is”

Sketch

$$w \sim_{\text{xlit}} v :\Leftrightarrow \text{xlit}^*(w) = \text{xlit}^*(v)$$

- **Idea:** account for *extinct characters*
- **Implementation:** $\mathcal{O}(1)$ character lookup table
- mostly useful as *preprocessor* for subsequent methods

(Jurish 2008, 2010b,c)

Successes

Long-‘s’	Abftand \mapsto Abstand	“distance”
Superscript-‘e’	n ^e ötig \mapsto nötig	“necessary”
Diacritics	Hochzît \mapsto Hochzeit	“wedding”

Failures

Diacritics	wôl \mapsto *wol	\neq wohl	“well”
Extant Characters	Thür \mapsto *Thür	\neq Tür	“door”
	vñ \mapsto *vn (\rightsquigarrow Vene)	\neq und	“and”

Sketch

(Geyken & Hanneforth, 2006)

- Model (modern) morphological processes as WFST
- Provides (infinite) *weighted target language* for subsequent methods
 - ▶ analysis cost \approx derivational complexity
- “Modern-wins” filtering

$$w \mapsto w : \Leftarrow w \in \pi_1(M_{\text{morph}})$$

Overgeneration: ‘modern’ form *shouldn’t* always win

andre	\mapsto	*André _{NE}	\neq	andere	“other”
from	\mapsto	*from _{FM.en}	\neq	fromm	“pious”
hiebei	\mapsto	*Hieb ei	\neq	hierbei	“hereby”
Zeugnuß	\mapsto	*Zeug nuß	\neq	Zeugnis	“report”
keyserlich	\mapsto	*Keys _{NE} erl _{NE} ich	\neq	kaiserlich	“imperial”
Procefs	\mapsto	*Process	\neq	Prozeß	“process”

- **Workaround:** *safety heuristics*, e.g. *_{FM}, *_{NE}, *-Ei, *-Nuß, ...

Sketch

$$w \sim_{\text{pho}} v : \Leftrightarrow \text{pho}(w) = \text{pho}(v)$$

- **Idea:** conflate words by *phonetic form*

(Jurish, 2008, 2010b)

- **Implementation:** text-to-speech rule-set

(Möhler et al., 2001)

- ▶ modified & compiled as FST M_{pho}

- ▶ online k -best equivalence cascade search

(Jurish, 2010a)

$$C_{\text{eqpho}} := M_{\text{pho}} \circ M_{\text{pho}}^{-1} \circ \pi_1(M_{\text{morph}})$$

Problems

Kopfe	↦	*Kopfweh	≠	Kopf	“head”
wes	↦	*Wehs	≠	wessen	“whose”
gewegen	↦	*Geh wegen	≠	gewogen	“weighted”
maiestat	↦	*Mais tat	≠	Majestät	“majesty”
Moyfe	↦	*Mäuse	≠	Mose(s)	“Moses”
Troglodyt	↦	*Troгло duett	≠	Troglodyt	“troglodyte”

- **Workarounds:** target language pruning, cascade lookup cutoffs

$$\begin{aligned} \text{best}_{\text{rw}}(w) &:= \arg \min_{v \in \mathcal{A}^*} \llbracket M_{\text{rw}} \circ \pi_1(M_{\text{morph}}) \rrbracket(w, v) \\ w \sim_{\text{rw}} v &:\Leftrightarrow \text{best}_{\text{rw}}(w) = \text{best}_{\text{rw}}(v) \end{aligned}$$

Sketch

- **Idea:** map words to *nearest* extant type (Jurish, 2010b,c)
 - ▶ generalized string edit distance (Damerau, 1964; Levenshtein, 1966)
 - ▶ computable even for infinite lexica (Mohri, 2002; Jurish, 2010a)
- **Implementation:** online ($k = 1$)–best cascade search
 - ▶ editor WFST M_{rw} compiled from ca. 300 SPE-style rules
 - ▶ EM-style weight optimization (Viterbi)
 - ▶ weight interpolation constants λ_{rw} and λ_{morph}

Problems

ehligen	↦ *Elchen	≠ ehelichen	“marital”
gewüntschten	↦ *Gewinde sekten	≠ gewünschten	“desired”
Predigamt	↦ *Birdie gambit	≠ Predigamt	“ministry”
Verdamnuß	↦ *Dominus	≠ Verdammnis	“damnation”
Weffier	↦ *Wessi	≠ Wesir	“vizier”

- **Workarounds:** rule adjustments, punitive target weights, lookup cutoffs



Idea

(Mays et al., 1991; Brill & Moore, 2000; Jurish, 2010c)

- Allow high-recall overgeneration at type level (xlit, eqpho, rw)
- Recover precision using token-level context

Implementation: Dynamic Hidden Markov Model (HMM)

- **States** are word-conflator pairs

$$Q = (\mathcal{W} \cup \{\mathbf{u}\}) \times R$$

- **Observations** are input strings

$$O_S = \bigcup_{i=1}^{n_S} \{w_i\} \subset \mathcal{A}^*$$

- **Transitions** (*static*)

$$A(\langle \tilde{w}_i, r_i \rangle_{i=1}^m) \approx p(\tilde{w}_m | \tilde{w}_1^{m-1})$$

- **Lexicon** (*dynamic*): Maxwell-Boltzmann distribution

$$B(\langle \tilde{w}, r \rangle, w) \approx \frac{b^{\beta d_r(w, \tilde{w})}}{\sum_{r' \in R} \sum_{\tilde{w}' \in \downarrow[w]_{r'}} b^{\beta d_{r'}(w, \tilde{w}')}}$$

- ▶ b, β are global model parameters ($b \geq 1, \beta \leq 0$)
- ▶ $d_r(w, \tilde{w})$ depends on conflator r

- **Lookup** (moot): Viterbi Algorithm

(Viterbi, 1967)



Sparse and/or Inappropriate Training Data

- many “normal” words treated as *unknown* (**u**)
- poor handling of historical *syntax*
 - ▶ *In deym dienst bestendig bleyben / die trubsall vnns nicht abtreiben*
M. Luther et al.: *Eyn Enchiridion oder Handbuchlein*, 1524
 - ▶ *Vnd mich deucht / das er den Hyacinthum auff einer seite lieb hatte*
G. Rollenhagen: *Vier Bücher wunderbarlicher . . . himmel*, 1603
 - ▶ *Daher newlichs einer gefagt: die tortur ist allmächtig*
F. von Spee: *Gewissens-Buch, Von Processen Gegen die Hexen*, 1647
 - ▶ *kaum sieht er hinein, so erblickt er ein Medaillon*
Goethe: *Wilhelm Meisters Lehrjahre, Bd. 4*, 1796

Edelgebohrnen	↪ *Tele gebahren	≠ edelgeborenen	“noble-born”
potz [blitz]	↪ *Potts	≠ potz _{ITJ}	“gee [whiz]”
Secretärsstelle	↪ *Secretärsstelle	≠ Sekretärsstelle	“secretary’s job”
unwifend	↪ *anwesend	≠ unwissend	“unknowing”
Phafmate	↪ *Faß matte	≠ Phasmate	“phantasms”

- **Workarounds:** language guessing, PoS heuristics, exlex, . . .
- **Cunning & Devious Plan:** incorporate PoS-tagger model (moot)

Extinct or Missing Target Lexemes

elektroskopische \mapsto *Elektro|ska|piefke “electroscopic”
Sakramentierer \mapsto *Sakrament|irrer “sacramentists”
glorwürdigen \mapsto *Chlor|würdigen_{VV} “glory-worthy”

Extinct or Missing Morphological Processes

aller- aller|größte \mapsto *grüßest “greatest of all”
hier- hier|nächst \mapsto *hin|nagst “subsequent”
ob- ob|angeregter \mapsto *Oboen|gerechter “aforementioned”
-e Kopf|e \mapsto *Kopf|weh “head”
-ig standhaft|ig \mapsto *stoned|hastig “steadfast”
-lich auflös|lich \mapsto *auslöschlich “soluble”

Lexical Shift / Codification

bälder \mapsto *Bälde_{NN} “sooner”
neulicher \mapsto *neulich_{ADV} “recent”

Proper Names

NE \mapsto ***Lex** Carlowitz \mapsto *Gorilla|witz_{NN} \neq Carlowitz
Philipson \mapsto *File|bison_{NN} \neq Philipson
Moÿfe \mapsto *Mäuse_{NN} \neq Moses

NE +flect Jef|um \mapsto *Jass|ohm \neq Jesus
Mathilde|n \mapsto *modelten \neq Mathilde
George|ns \mapsto *Chargen \neq Georg

Non-Lexical Material

from \mapsto *from_{FM.en} \neq fromm “pious”
medicinæ \mapsto *Medizin|ah \neq medicinae “medical” (FM.lat)
mon Dieu \mapsto *Mohn Dia \neq mon Dieu “my God” (FM.fr)
Zn \mapsto *Zen \neq Zn “zinc” (chem.)

Ambiguity / Divergence

wit \mapsto mit, wie, weit, wir, wit
widder \mapsto wider, wieder, weder, Widder

— Experiments —
*... or, “how bad is it **really**?”*

DTA Evaluation Corpus

(Jurish, Drotschmann & Ast, 2013)

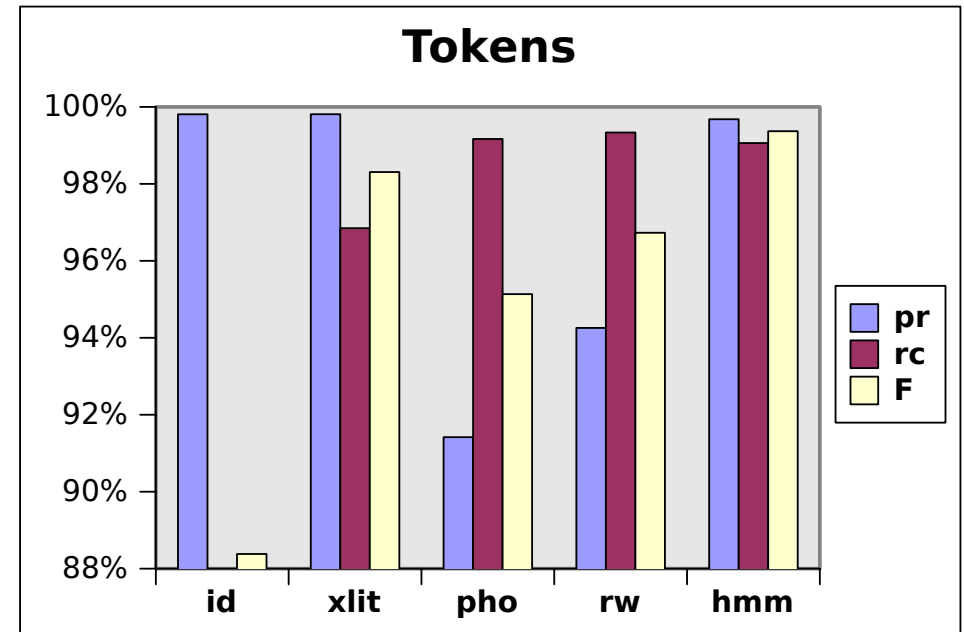
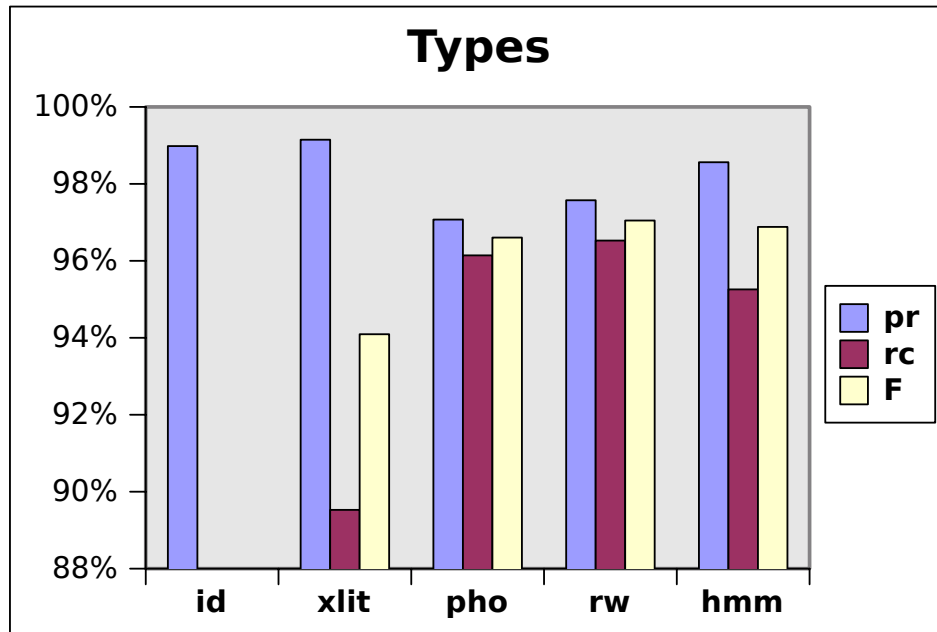
- Gold standard subset of *DTA* Phase 1
 - ▶ 13 volumes *(1780 – 1880)*
 - ▶ 114,542 tokens; 10,070 types *(alphabetic only)*
- *Canonical cognate* assigned to each token
 - ▶ automatic **alignment** with conventional edition
 - ▶ 3-pass manual **review**
 - Type-wise *(conservative)*
 - Token-wise *(unverified tokens only)*
 - “Suspicious” pairs *(heuristic selection)*

Evaluation Measures

- Simulated information retrieval task
- Type- and token-wise precision (pr), recall (rc), and F



Experiment 1: Results



	% Types			% Tokens		
	pr	rc	F	pr	rc	F
id	99.0	59.2	74.1	99.8	79.3	88.4
xlit	99.1	89.5	94.1	99.8	96.8	98.3
pho	97.1	96.1	96.6	91.4	99.2	95.1
rw	97.6	96.5	97.0	94.3	99.3	96.7
hmm	98.6	95.3	96.9	99.7	99.1	99.4

‘Evaluation’ Corpus \rightsquigarrow Ground-Truth Relevance

$$\text{relevant}(w, \tilde{w}) := \{(v, \tilde{v}) : \tilde{v} = \tilde{w}\}$$

- Most thoroughly annotated corpus subset
- 13 volumes \sim 320k tokens \sim 28k types (words only)

‘Training’ Corpus \rightsquigarrow Canonicalization Lexicon (lex)

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

- Strictly disjoint from test corpus (by author)
- 101 volumes \sim 3.5M tokens \sim 158k types (words only)

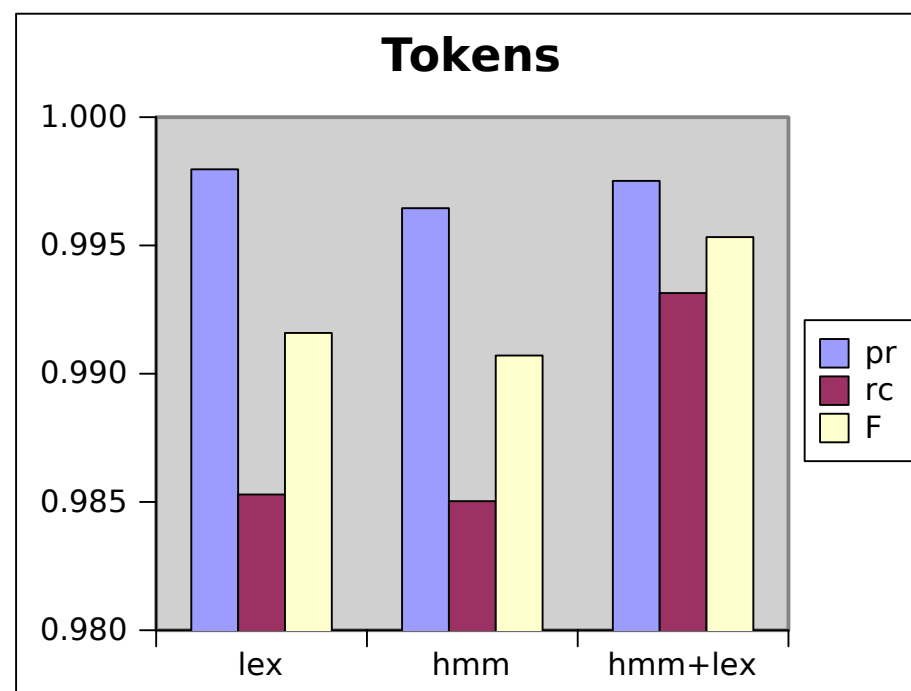
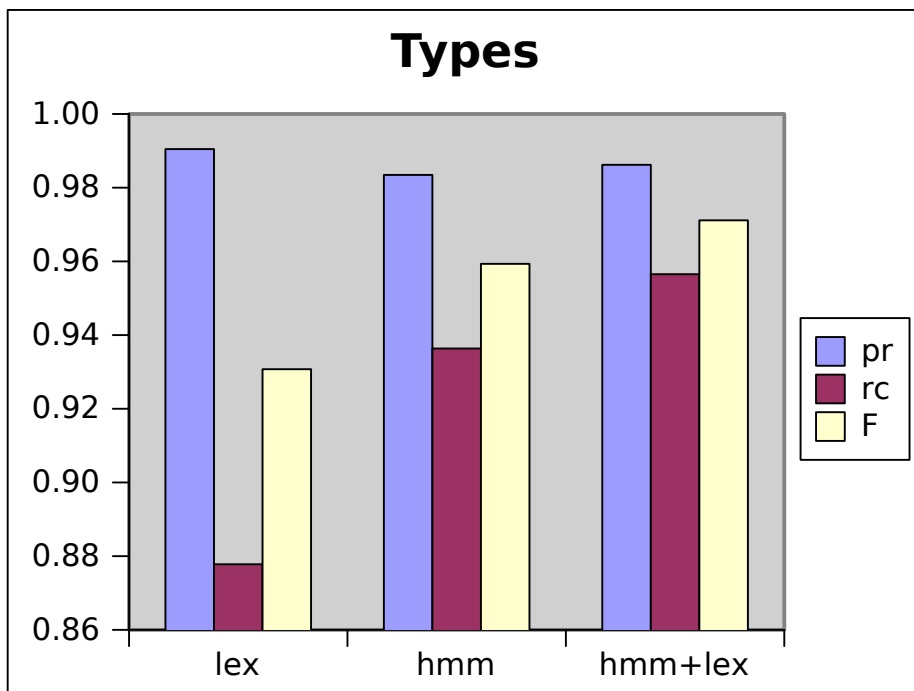
Evaluation

- Simulated information retrieval (pr, rc, F)
- Tested methods: id, lex, hmm, hmm+lex

(van Rijsbergen, 1979)



Experiment 2: Results



	% Types			% Tokens		
	pr	rc	F	pr	rc	F
id	99.1	55.7	71.3	99.8	78.5	87.9
lex	99.0	87.8	93.1	99.8	98.5	99.2
hmm	98.3	93.6	95.9	99.6	98.5	99.1
hmm+lex	98.6	95.7	97.1	99.8	99.3	99.5

Historical Text and Conventional Tools

don't play together nicely "out of the box"

Canonicalization Methods

- static lexicon
- transliteration
- phonetic equivalence
- rewrite cascade
- HMM disambiguator

~ effective but sparse

~ quick and dirty

~ elegant but coarse

~ flexible but costly

~ precision recovery

Advanced Topics

- parameter optimization
- rewrite editor
- HMM smoothing

≥ 27 free parameters

induction from training pairs

PoS, semantics, morphs, ...

þe Olde Lafft Slynde ("The End")

<http://www.deutschestextarchiv.de/demo/cab/>