

# Canonicalization techniques for computer-mediated communication

Bryan Jurish,<sup>1</sup> Kay-Michael Würzner,<sup>2</sup> Maria Ermakova,<sup>1</sup> Sophie Arana<sup>1</sup>  
{jurish,wuerzner,ermakova}@bbaw.de

(1) Berlin-Brandenburgische Akademie der Wissenschaften

(2) Universität Potsdam

GSCL-2013 Workshop “*Verarbeitung . . . internetbasierter Kommunikation*”

Darmstadt, 23<sup>rd</sup> September 2013

# Overview

## The Big Picture

- The Situation
- The Problem
- The Approach

## Canonicalization Methods

- Type-wise Conflation
  - ▶ Transliteration, Phonetization, Rewrite Cascade
- Token-wise Disambiguation
  - ▶ Dynamic Hidden Markov Model
- Problems & Workarounds
  - ▶ Lexical, Morphological, Phonological, Suprasegmental

## Summary & Outlook

# — The Big Picture —

# Deutsches Textarchiv (DTA)

## Primary Goals

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

## In Numbers

*(DTA+DTAE 2013-05-02)*

1,276 transcribed works

353,245 digitized pages

567,200,587 unicode characters

81,049,777 tokens (alpha-numeric)

1,947,414 types (alpha-numeric)

**TODO: UPDATE FOR IBK CORPUS**

# The Situation

## CMC Text ≠ Orthographic Conventions

- also applies to historical text, OCR output, ...
- many non-standard graphemic forms

nicht  
“not”

nciht, nert, net, nich, nicha,  
niiiiicht, nit, nnicht, ...

tschüss  
“bye”

schüssi, tschaui, tschö,  
tschöööö, tschüssi, tschüssie,  
tschüüüüüüüüüüüüsss, ...

## Conventional NLP Tools ⇒ Strict Orthography

- IR systems, PoS taggers, stemmers, lemmatizers, morphological analyzers, parsers, ...
- **Fixed lexicon** keyed by (ortho)graphic form
- **Conventional** lexemes only (*≈ newspaper domain*)

# The Problem

$$\begin{array}{rcl} & \text{Conventional Tools} & \\ \oplus & \text{Unconventional Corpus} & \\ \hline = & & \text{Soup} \end{array}$$

- 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’, ...*

c	u	l8r	dooooood
↓	↓	↓	↓
see	you	later	dude

## In a Nutshell

- *Map* each **word**  $w$  to a unique **canonical cognate**  $\tilde{w}$
- *Defer* application analysis to canonical forms
- DTA::CAB infrastructure developed for historical text

## Canonical Cognates

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

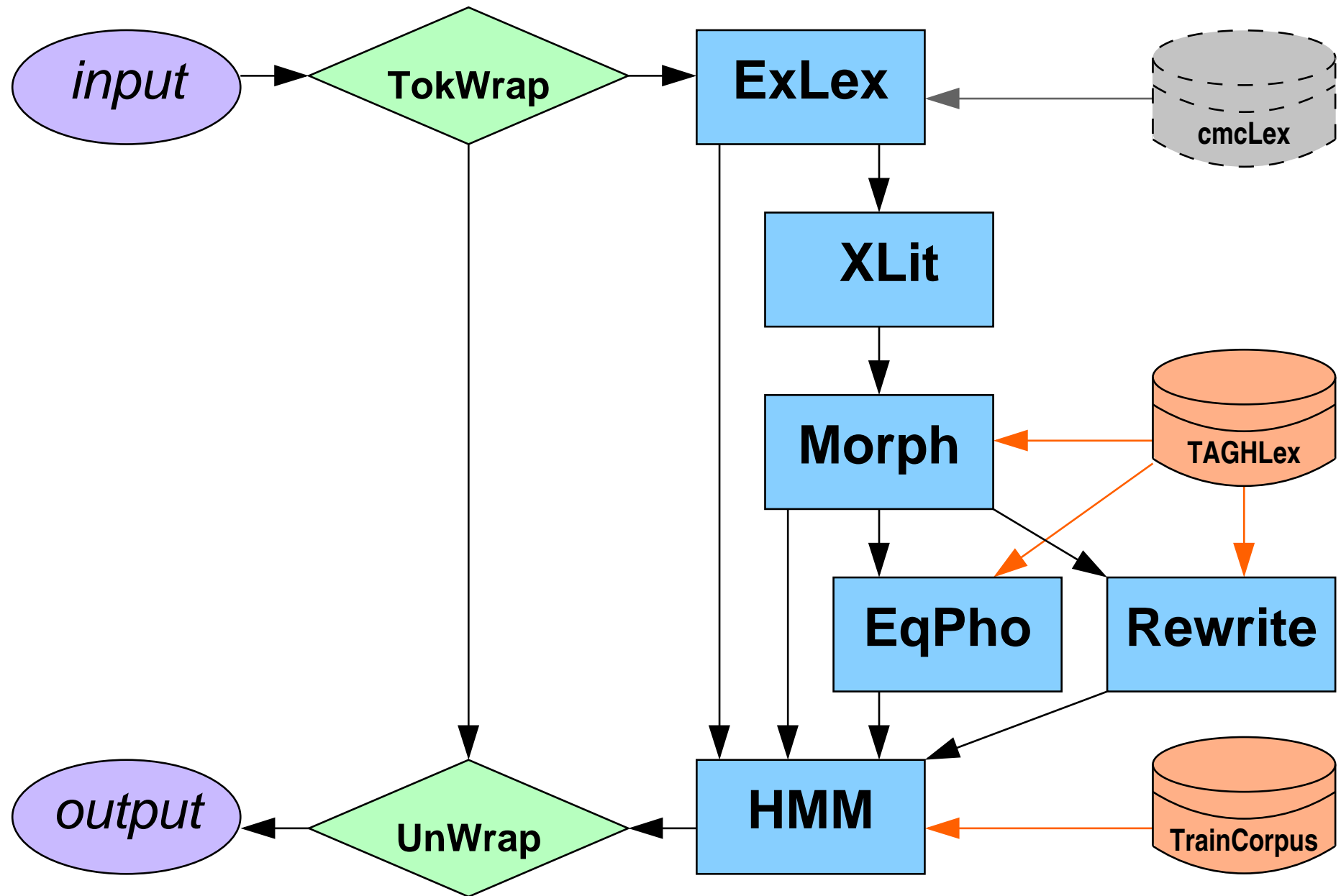
## Conflation Relation $\sim_r$

- *Binary relation* on strings (words) in  $\mathcal{A}^*$
- Prototypically a true *equivalence relation*

# — Canonicalization Methods —



# DTA::CAB System Architecture



# Deterministic Transliteration (xlit)

## Sketch

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

- **Idea:** handle *non-standard characters* (Jurish 2008, 2010b,c)
- **Implementation:**  $\mathcal{O}(1)$  character lookup table
- useful as *preprocessor* for subsequent methods

## Examples

### Foreign Scripts

北京  $\mapsto$  BeiJing

“Beijing”

Πλάτων  $\mapsto$  Platon

“Plato”

### Diacritics

schon  $\mapsto$  schon

“already”

schøn  $\mapsto$  schön

“beautiful”

### Ligatures

œde  $\mapsto$  öde

“bleak”

Ærger  $\mapsto$  Ärger

“trouble”

### (pesudo-) Allographs

hʉ@h  $\mapsto$  huch

“oops”

# TAGH Morphology (morph)

## Sketch

(Geyken & Hanneforth, 2006)

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

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

## Overgeneration: ‘known’ form *shouldn’t* always win

hab	$\mapsto$	*hab <sub>VVIMP</sub>	$\neq$	habe	“have”
Andre	$\mapsto$	*André <sub>NE</sub>	$\neq$	Andere	“other”
muste	$\mapsto$	*mus te	$\neq$	muss te	“must”
Grad	$\mapsto$	*Grad	$\neq$	Gerade	“directly”
schaftsinnig	$\mapsto$	*schaf sinnig	$\neq$	scharfsinnig	“perceptive”

- **Workaround:** *safety heuristics*, e.g. \*V.imp, \*NE, \*f = 0

# Phonetic Equivalence (eqpho)

## 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  $M_{\text{pho}}$
  - ▶ online  $k$ –best equivalence cascade search (Jurish, 2010a)

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

## Examples

<b>Successes</b>	betrib	$\mapsto$	Betrieb	“operation”
	eindoitig	$\mapsto$	eindeutig	“univocal”
	hallloooo	$\mapsto$	hallo	“hello”
<b>Failures</b>	mirsch	$\nrightarrow$	mich	“me”
	nischt	$\nrightarrow$	nicht	“not”

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

# Rewrite Cascade ( $_{rw}$ )

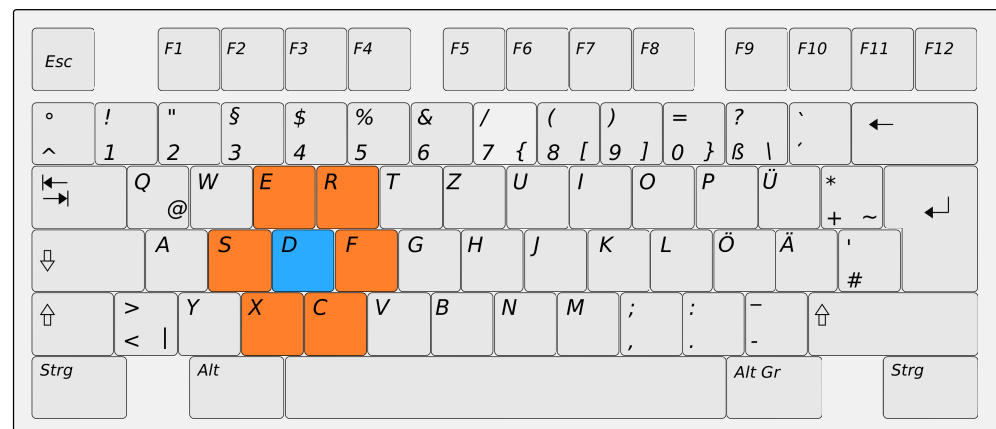
## Sketch

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

- **Idea:** map words to *nearest* conventional type (Jurish, 2010b,c)
  - ▶ generalized string edit distance (Levenshtein, 1966)
  - ▶ computable even for infinite lexica (Mohri, 2002; Jurish, 2010a)
- **Implementation:** online ( $k = 1$ )–best cascade search
  - ▶ editor WFST  $M_{rw}$  compiled from *keyboard adjacency model*
  - ▶ weight interpolation constants  $\lambda_{rw}$  and  $\lambda_{\text{morph}}$

## Keyboard Adjacency

- limit search space
- attenuate edit costs
- insert, delete, replace  
... but **NOT** transpose!



# Error Model $M_{rw}$

## Generic Edit Operations

Deletion	gu <u>e</u> t $\mapsto$ gut	<i>“good”</i>
Insertion	is $\mapsto$ is <u>t</u>	<i>“is”</i>
Substitution	g <u>o</u> ht $\mapsto$ g <u>e</u> ht	<i>“goes”</i>
Transposition	u <u>dn</u> $\mapsto$ u <u>nd</u>	<i>“and”</i>
Doubling	da <u>s</u> $\mapsto$ da <u>ss</u>	<i>“that”</i>
Un-doubling	neu <u>ees</u> $\mapsto$ neu <u>e</u> s	<i>“new”</i>
Un-potentialiation	juh <u>uuuu</u> $\mapsto$ juh <u>u</u>	<i>“woo-hoo”</i>

## Keyboard-Adjacent Edit Operations

Deletion	<u>i</u> ommer $\mapsto$ <u>i</u> mmmer	<i>“always”</i>
	h <u>oi</u> $\mapsto$ h <u>i</u>	<i>“hi”</i>
Insertion	<u>z</u> $\mapsto$ <u>zu</u>	<i>“to”</i>
	scha <u>f</u> $\mapsto$ schar <u>f</u>	<i>“sharp”</i>
Substitution	gu <u>v</u> kt $\mapsto$ gu <u>b</u> kt	<i>“looks”</i>

# Token-wise Disambiguation (hmm)

## Idea

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

- Allow high-recall overgeneration at type level
- Recover precision using token-level context

## Implementation: Dynamic Hidden Markov Model (HMM)

- **States** are word-conflator pairs

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

- **Observations** are input strings

$$\mathcal{O}_S = \bigcup_{i=1}^{n_S} \{\mathbf{w}_i\} \subset \mathcal{A}^*$$

- **Transitions** (*static*)

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

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

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

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

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

(Viterbi, 1967)

# HMM Example

**Input** hallooo wida muste kuz w€g



# HMM Example

<b>Input</b>	hallooo	wida	muste	kuz	wEg
--------------	---------	------	-------	-----	-----

---

**xlit**

*hallooo*

*wida*

*muste*

*kuz*

*wEg*

# HMM Example

Input	hallooo	wida	muste	kuz	wEg
-------	---------	------	-------	-----	-----

---

**xlit**

*hallooo*

*wida*

*muste*

*kuz*

*wEg*

**pho**

{*hallo*}

$\left\{ \begin{array}{l} \textit{wider,} \\ \textit{wieder} \end{array} \right\}$

{*muste*}

$\left\{ \begin{array}{l} \textit{Cuts,} \\ \textit{Kuts} \end{array} \right\}$

$\left\{ \begin{array}{l} \textit{weck,} \\ \textit{weg} \end{array} \right\}$

# HMM Example

Input	hallooo	wida	muste	kuz	wEg
<b>xlit</b>	<i>hallooo</i>	<i>wida</i>	<i>muste</i>	<i>kuz</i>	<i>wEg</i>
<b>pho</b>	{ <i>hallo</i> }	$\left\{ \begin{array}{l} \textit{wider,} \\ \textit{wieder} \end{array} \right\}$	{ <i>muste</i> }	$\left\{ \begin{array}{l} \textit{Cuts,} \\ \textit{Kuts} \end{array} \right\}$	$\left\{ \begin{array}{l} \textit{weck,} \\ \textit{weg} \end{array} \right\}$
<b>rw</b>	<i>hallo</i>	<i>Weda</i>	<i>musste</i>	<i>kurz</i>	<i>weg</i>

# HMM Example

Input	hallooo	wida	muste	kuz	wEg
xlit	<i>hallooo</i>	<i>wida</i>	<i>muste</i>	<i>kuz</i>	<i>wEg</i>
pho	{ <u>hallo</u> }	$\left\{ \begin{array}{l} \text{wider,} \\ \text{wieder} \end{array} \right\}$	{ <i>muste</i> }	$\left\{ \begin{array}{l} \text{Cuts,} \\ \text{Kuts} \end{array} \right\}$	$\left\{ \begin{array}{l} \text{weck,} \\ \text{weg} \end{array} \right\}$
rw	<u>hallo</u>	<i>Weda</i>	<u>musste</u>	<u>kurz</u>	<u>weg</u>
hmm	<i>hallo</i>	<i>wieder</i>	<i>musste</i>	<i>kurz</i>	<i>weg</i>

# HMM Example

Input	hallooo	wida	muste	kuz	wEg
xlit	<i>hallooo</i>	<i>wida</i>	<i>muste</i>	<i>kuz</i>	<i>wEg</i>
pho	{ <u>hallo</u> }	$\left\{ \begin{array}{l} \textit{wider,} \\ \textit{wieder} \end{array} \right\}$	{ <i>muste</i> }	$\left\{ \begin{array}{l} \textit{Cuts,} \\ \textit{Kuts} \end{array} \right\}$	$\left\{ \begin{array}{l} \textit{weck,} \\ \textit{weg} \end{array} \right\}$
rw	<u>hallo</u>	<i>Weda</i>	<u>musste</u>	<u>kurz</u>	<u>weg</u>
hmm	<i>hallo</i>	<i>wieder</i>	<i>musste</i>	<i>kurz</i>	<i>weg</i>
Output	<b>hallo</b>	<b>wieder</b>	<b>musste</b>	<b>kurz</b>	<b>weg</b>

# Some Anticipated Difficulties

## Lexical

- neologisms & “missing” lexemes
- Workaround:** static exception lexicon  $M_{\text{exlex}}$

ne  $\mapsto$  eine “a”  
ma  $\mapsto$  mal “once”

## Morphological\*

- “missing” productive processes
- Workaround:**  $M_{\text{morph}}$  plug-in(s)

glaub  $\mapsto$  glaube “believe”  
welch  $\mapsto$  welchen “which”

\* ... or typographical simulations of phonological processes acting on inflectional morphemes?

## Phonetic / Phonological

- identity criterion is too strict
- Workaround:**  $M_{\text{rw}}$  plug-in(s)

sehn  $\sim$  sehen “see”  
zehn  $\not\sim$  Zehen “ten”

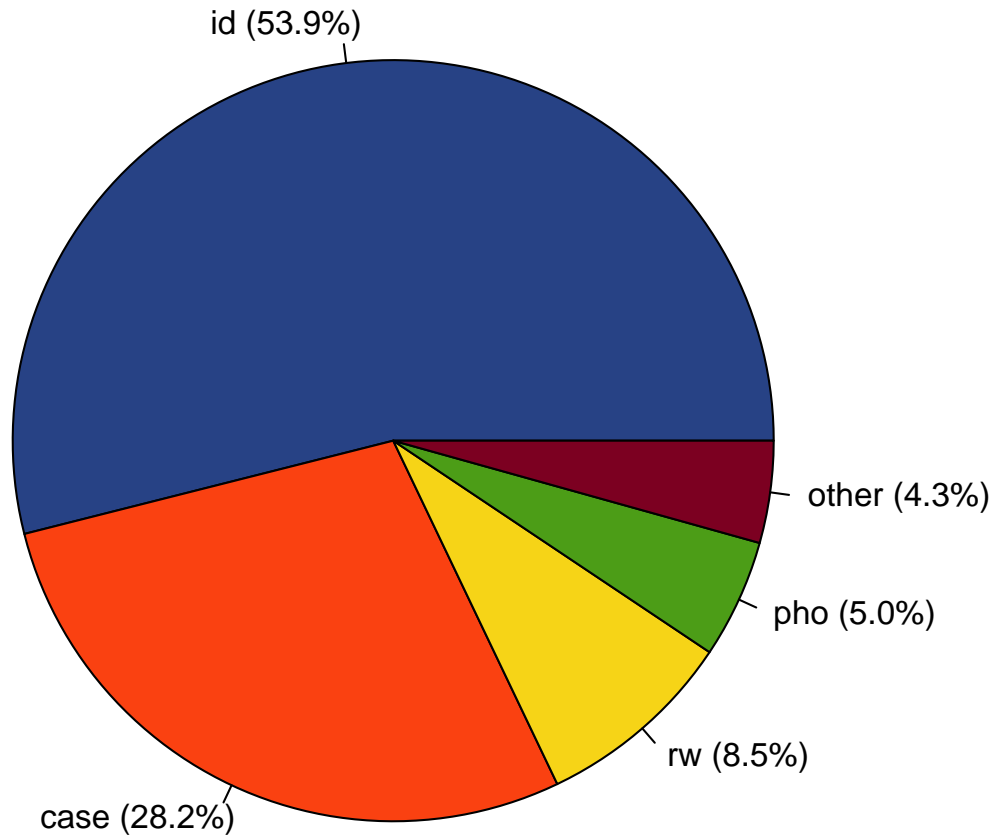
## Suprasegmental

- unclear token boundaries
- Workaround:** ???

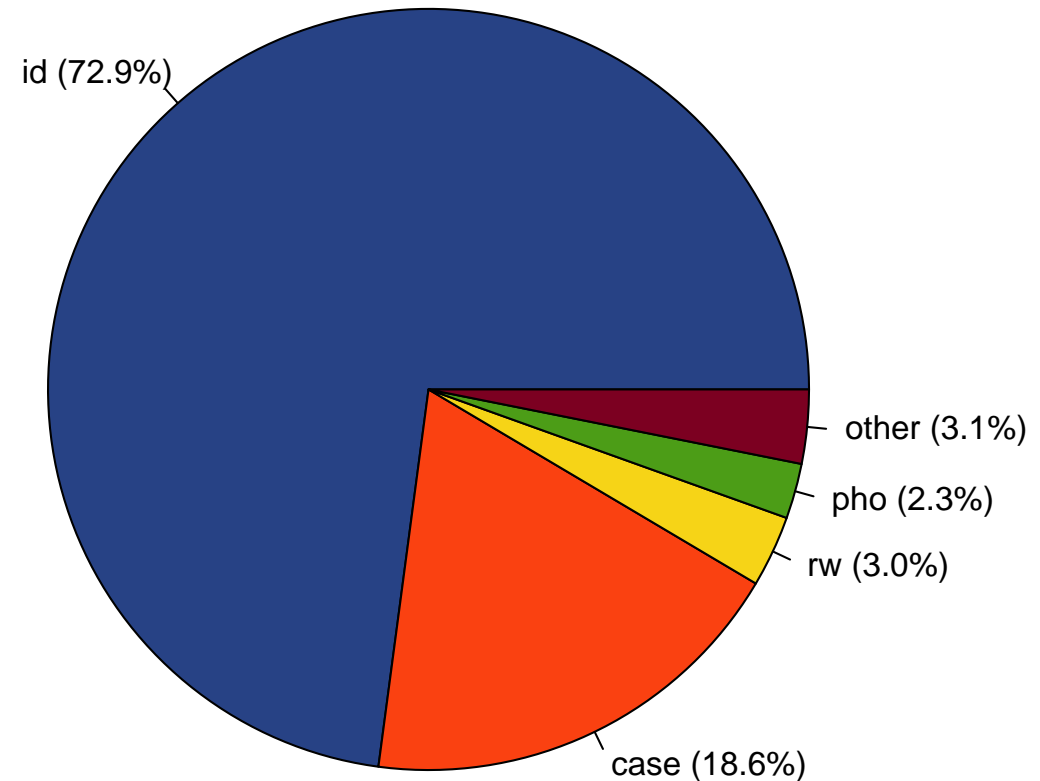
hat's  $\mapsto$  hat es “has it”  
kannste  $\mapsto$  kannst du “can you”

# Observed Mappings

**by Type** ( $N = 4340$ )



**by Token** ( $N = 13611$ )



# Concluding Remarks

## Unconventional Text and Conventional Tools

*won't play together nicely "out of the box"*

## Canonicalization Methods

- transliteration *~ quick and dirty*
- phonetic equivalence *~ elegant but coarse*
- rewrite cascade *~ flexible but costly*
- HMM disambiguator *~ precision recovery*

## Future Work

- part-of-speech annotations *STTS*
- plug-in resources  *$M_{\text{exlex}}, M_{\text{rw}}, M'_{\text{morph}}$*
- model parameter optimization  *$\geq 27$  free parameters*



C U L8R DØØD\$  
 (“The End”)

*Thank you for listening!*