

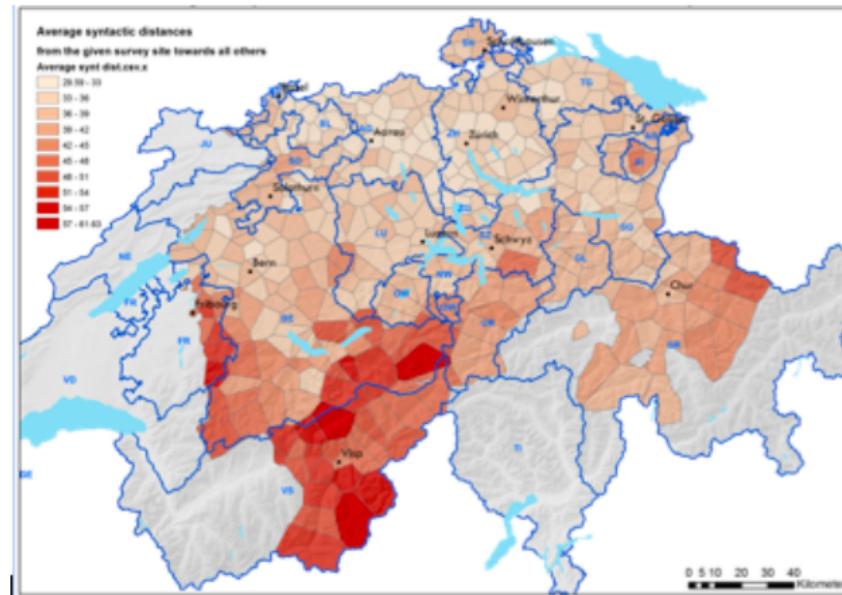


# Normalising orthographic and dialectal variants for the automatic processing of Swiss German

Tanja Samardžić, Yves Scherrer, Elvira Glaser



## Swiss German Dialects



Map by Péter Jeszenszky



## The ArchiMob corpus





## A corpus for everybody

Query **mein** 182 (1,611.86 per million)

Page  of 10 Go [Next](#) | [Last](#)

hä er dä ganz möse hööre / aber / eh jez **mine** wunsch isch dän düe noch gsii / das ich  
zwaar bin lich / i / in wolfeschiesse / **mil** vatter und mi mueter sind wolfeschiesser  
chenn ich det ässen und schlaaffe daa / salt **mi** / as du / eh / äbe z jung bisch / und  
dich nad aamälde schüsich märkli si as duu **mis** / as duu / eh / äbe z jung bisch / und  
( de ) han ich / scheen verdienet / und **mi** vatter isch glücklich gsii und gsait plib  
het er sibeze jaar gschaffet / han ich den **min** maa isch drissigi gsii / und ich abe zwaijewänzgi  
/ nuur hemmer enand hait jaare und jaare **mi** jüngsch brieder isch do achtjährig gsii  
han ich no tänkt dasch ez doch äalne wl **mi** brieder mi brieder het ganz / wißblondi  
tänkt dasch ez doch äalne wl **mi** brieder het ganz / wißblondi haarr ghaa  
chilen uis bin / dasch / glöüb ich / eh / **mi** brieder de salt er ja / het ja no vil so  
und de het mer sich doch mü ggää / und **mi** vatter isch natlierli gliklich gsii / i dére  
jeregot mer hend gad / allernöötigisch ghaa / **min** maa het esoo nach em firaabig / eppis gschaffet  
die sind den al no i de lèér gsii / und **min** maa isch fascht echli ire maischter gsii  
/ de het me gwüssst ass e chrieg chunt / **min** maa het immer gsaaat / muesch de chasch  
mobilmachung bekant gegeben / ja / und de isch **mi** maa / am zwölfi choo dee het doo vo de  
liruke de isch i de / glilice ainhalt ( wo ) **min** maa / bli Ich zu dérem und ha gsalt duu  
wäärde / isch dä kiosk offe / de isch immer **mi** maa det gsii und de ( ? ) de bek bringt  
irgend es vorlig bet oder eppis halig aso **mi** maa isch doo im dienscht / plibe / biss  
und pro maa / e franken über / und daa het **mi** maa mir zurgg-gschribre sig alles schön  
e franke ghaa han ich drii franke ghaa / **mi** maa het achzg rappe sold ghaa / und de

Page  of 10 Go [Next](#) | [Last](#)

Goal: 700 000 words (44 docs)   Currently : 500 000 words (34 docs)



## The problem of variation

Other docs	mine	man	hed	ime	gsäit
	mini		hèd	imer	gsääit
	määin		hét	emmer	
	mäin		heet	iiimer	
	main		haa		
			händ		
			hüt		
Same doc	mi	ma	hat		gsait
	mii		hät		
	miin				
	mis				
	miis				
Transcription	min	maa	het	immer	gsaait
Normalised	<b>mein</b>	<b>mann</b>	<b>hat</b>	<b>immer</b>	<b>gesagt</b>
Translation	My husband has always said...				

Table : A segment of a transcribed and normalised text with corresponding variants found in the same document and in other documents.



## Word-level annotation

jaa	ja	ITJ
de	dann	ADV
het	hat	VAFIN
me	man	PIS
no	noch	ADV
gluegt	gelugt	VVPP
tänkt	gedacht	VVPP
dasch	das_ist	PDS+
ez	jetzt	ADV
de	der	ART
genneraal	general	NN
jaa	ja	ITJ
das	das	PDS
ischsch	ist	VAFIN
en	ihn	PPER
ez	jetzt	ADV



## The approach

- Manual normalisation of 6 documents with VARD 2 and IGT
- Developing an automatic approach using the manually annotated sample for training and testing
- Semi-automatic processing of the remaining documents:
  1. Train the system on all the manually annotated data
  2. Process a new document
  3. Evaluate and correct the output
  4. Add the new document to the training set
  5. Repeat 1-4 for all the remaining documents



## Automatic normalisation as translation

jaa	ja
de	dann
het	hat
me	man
no	noch
gluegt	gelugt
tänkt	gedacht
dasch	das_ist
ez	jetzt
de	der
genneraal	general

jaa	ja
das	das
ischsch	ist
en	ihn
ez	jetzt



## The approach

The problem: little data, a lot of variation

Dividing the task:

**Unique**: no ambiguity, only one normalisation possible

**Ambig1**: multiple possible normalisations, one predominant

**Ambig2**: multiple possible normalisations, none predominant

**New**: unknown words (no information about possible normalisations)



## Automatic methods

**W-b-w:** Assign to each test word the most probable normalisation based on word frequency

**CSMT:** Generate for each test word the most probable normalisation based on Character-Based Statistical Machine Translation

**ALM:** Rerank the candidates using the ArchiMob language model

**ELM:** Rerank the candidates using an extended language model (TüBa-D/S + ArchiMob)



## Results

	Proportion (%)	Accuracy (%)					
		No	W-b-w	CSMT	ALM	ELM	Best
Unique	44.53	23.06	<b>87.83</b>	87.83	87.83	87.83	87.83
Ambig1	43.68	21.92	<b>80.63</b>	80.63	70.11	67.12	80.63
Ambig 2	0.52	4.17	<b>38.43</b>	<b>39.81</b>	<b>46.06</b>	<b>49.54</b>	49.54
New	11.27	3.64	3.64	<b>23.88</b>	23.88	23.88	23.88
All	100	20.28	74.94	77.23	72.67	71.38	77.28

Table : Accuracies of the automatic normalisation methods for the different word classes, averaged over the 5 folds.



## Conclusions

- Combined methods are better than any individual method
- Advanced methods help with hard cases (ambiguous and new words), but more data needed
- Data from standard German do not help much