

Natural dialect processing

NLP for non-standardized language varieties

Verena Blaschke & Barbara Plank
MaiNLP lab, LMU Munich

Quantitative approaches in dialectology
and variationist sociolinguistics
December 7, 2023



Natural Language Processing

Natural Language Processing

Natural Language Processing

... but *k* \ / *W* languages?

NLP – but which “language(s)”?

- Many speakers, abundant data, standardization

NLP – but which “language(s)”?

- Many speakers, abundant data, standardization

But how do we actually use language?

NLP – but which “language(s)”?

- Many speakers, abundant data, standardization

But how do we actually use language?

- Also include minority languages, non-standard varieties

NLP – but which “language(s)”?

- Many speakers, abundant data, standardization

But how do we actually use language?

- Also include minority languages, non-standard varieties
- Tricky for NLP!
Modern methods learn from massive amounts of data –
how can we learn from sparse + heterogeneous data?

NLP – but which “language(s)”?

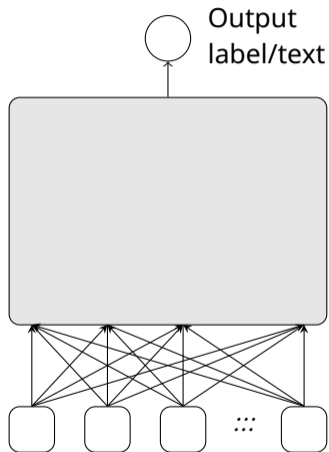
- Many speakers, abundant data, standardization

But how do we actually use language?

- Also include minority languages, non-standard varieties
- Tricky for NLP!
Modern methods learn from massive amounts of data –
how can we learn from sparse + heterogeneous data?

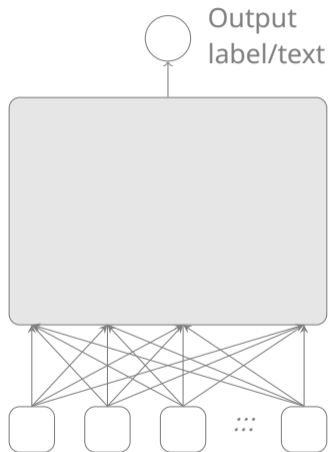
Overview of current challenges and approaches regarding
NLP & dialects

Overview



Input text sequence goes here

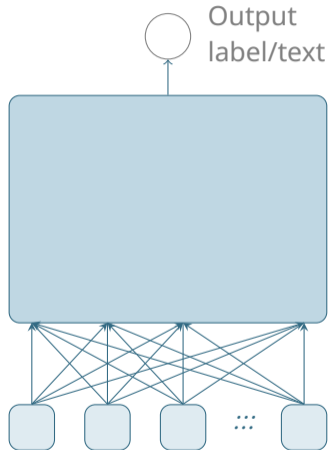
Overview



Input text sequence goes here

Available dialect data

Overview

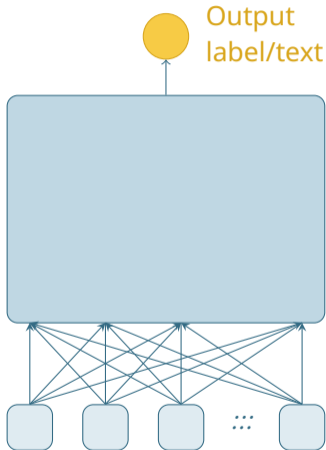


Input text sequence goes here

Modelling non-standard data

Available dialect data

Overview



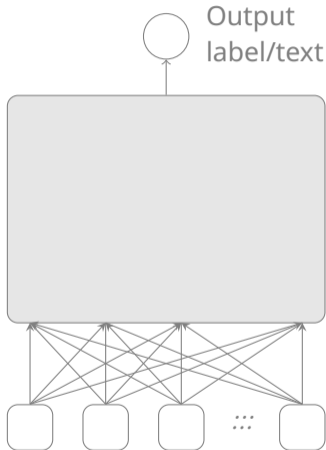
Input text sequence goes here

Human-centric NLP
(what tools and why?)

Modelling non-standard data

Available dialect data

Overview



Input text sequence goes here

Human-centric NLP
(what tools and why?)

Modelling non-standard data

Available dialect data

(Lack of?) resources

Datasets for dialects and “small” languages

(Lack of?) resources

Datasets for dialects and “small” languages

- Two communities: variationists & NLP researchers

(Lack of?) resources

Datasets for dialects and “small” languages

- Two communities: variationists & NLP researchers
- Findable; licenses allowing re-use

(Lack of?) resources

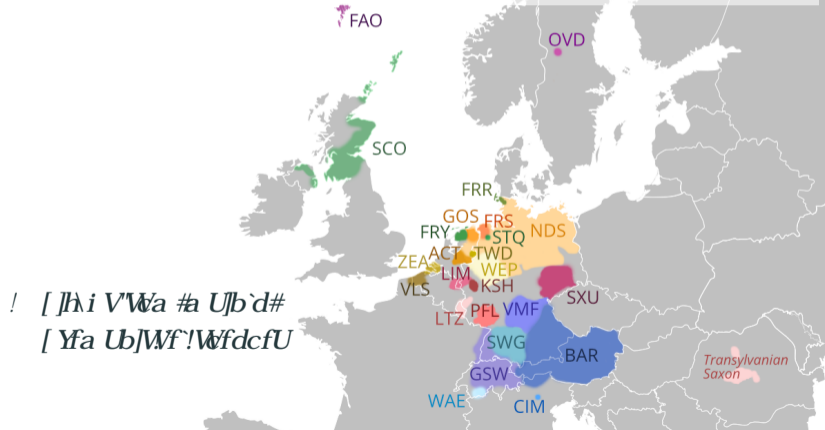
Datasets for dialects and “small” languages

- Two communities: variationists & NLP researchers
- Findable; licenses allowing re-use
- Long-term storage + accessibility

Corpus overview (small/non-std Gmc varieties)

(+spoken primarily outside Europe)

(+non-standard varieties associated with NOR, DAN, SWE, DEU)



! [ɫi vʷa #a U]b'd#
[Yfa Ub]Wf!WfdcfU

Corpus overview (small/non-std Gmc varieties)

100+ (mostly written) corpora for 30 language varieties

Blaschke et al. "A survey of corpora for Germanic low-resource languages and dialects" NoDaLiDa 2023

Corpus overview (small/non-std Gmc varieties)

100+ (mostly written) corpora for 30 language varieties

- Largely unannotated

Blaschke et al. "A survey of corpora for Germanic low-resource languages and dialects" NoDaLiDa 2023

Corpus overview (small/non-std Gmc varieties)

100+ (mostly written) corpora for 30 language varieties

- Largely unannotated
- If annotated:
 - Geolocation, dialect group
 - Morphosyntax
 - Rarely: translations, content-related annotations

Blaschke et al. "A survey of corpora for Germanic low-resource languages and dialects" NoDaLiDa 2023

Corpus overview (small/non-std Gmc varieties)

100+ (mostly written) corpora for ~30 language varieties

- Largely unannotated
- If annotated:
 - Geolocation, dialect group
 - Morphosyntax
 - Rarely: translations, content-related annotations
- Mostly dedicated, curated corpora
Recently: also uncurated, web-crawled ones

Blaschke et al. "A survey of corpora for Germanic low-resource languages and dialects" NoDaLiDa 2023

Corpus overview (small/non-std Gmc varieties)

100+ (mostly written) corpora for ~30 language varieties

- Largely unannotated
- If annotated:
 - Geolocation, dialect group
 - Morphosyntax
 - Rarely: translations, content-related annotations
- Mostly dedicated, curated corpora
Recently: also uncurated, web-crawled ones

Data exchange between research communities?

Blaschke et al. "A survey of corpora for Germanic low-resource languages and dialects" NoDaLiDa 2023

Bavarian dependency treebank

Kcf_ i bXf fy JYk

- Different dialects
(location metadata)

Bavarian dependency treebank

Kcf_ i bXf fy JYk

- Different dialects
(location metadata)
- Different text genres

Bavarian dependency treebank

Kcf_ i bXf fYj Jk

- Different dialects
(location metadata)
- Different text genres
- Freely accessible,
permissive license

Bavarian dependency treebank

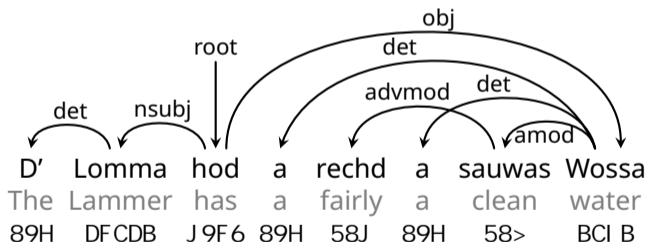
Kcf_ i bXf fy JYk

- Different dialects
(location metadata)
- Different text genres
- Freely accessible,
permissive license
- 11k words,
750 sentences

Bavarian dependency treebank

Kcf_ i bXf fYj JYk

Universal Dependencies



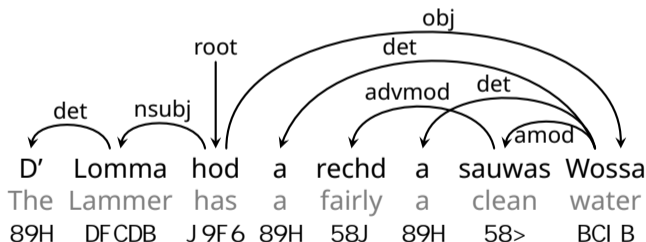
de Marneffe et al. "Universal Dependencies" *7ca di HLhcbU @b[i gHW* (2021)

Bavarian dependency treebank

Kcf_ i bXf fYj Yk

Universal Dependencies

- Cross-linguistic comparability (incl. DEU, GSW, NDS)



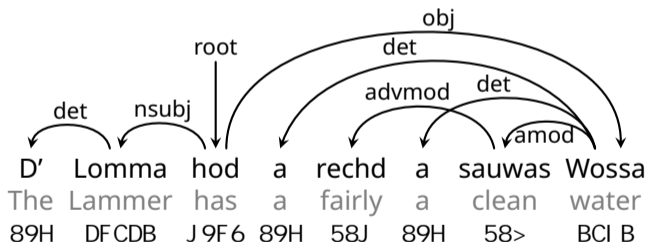
de Marneffe et al. "Universal Dependencies" *7ca di HLhcbU @b[i ghW* (2021)

Bavarian dependency treebank

Kcf_ i bXf fy Jk

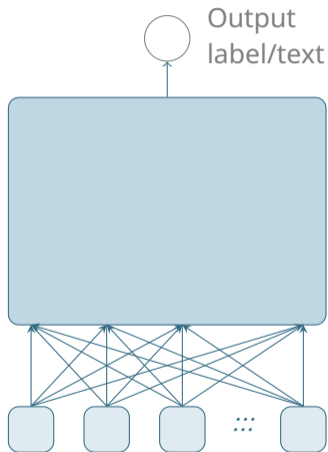
Universal Dependencies

- Cross-linguistic comparability (incl. DEU, GSW, NDS)
- Established for automatic annotation tasks



de Marneffe et al. "Universal Dependencies" *7ca di HLhcbU @b[i ghW* (2021)

Overview



Input text sequence goes here

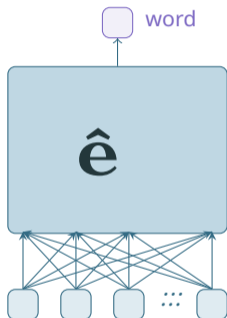
Human-centric NLP
(what tools and why?)

Modelling non-standard data

Available dialect data

Pretrain – finetune – transfer

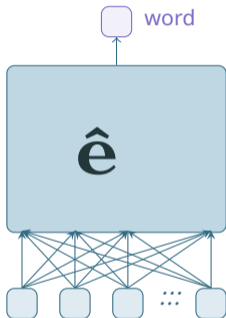
Pretraining



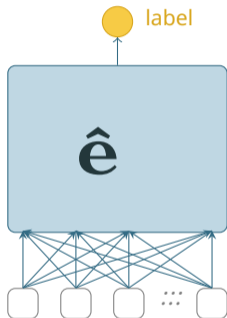
Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum. Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do

Pretrain – finetune – transfer

Pretraining



Finetuning

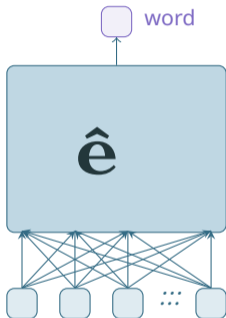


Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum. Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do

Task-specific input text

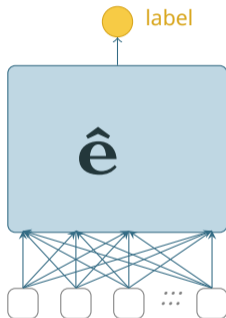
Pretrain – finetune – transfer

Pretraining



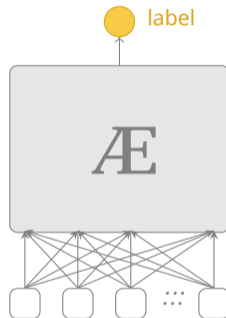
Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum. Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do

Finetuning



Task-specific input text

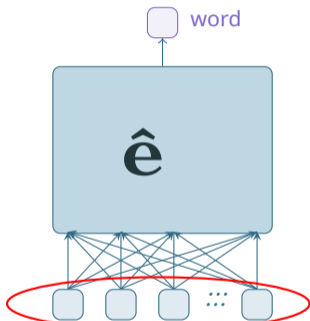
Transfer



Input text in another language

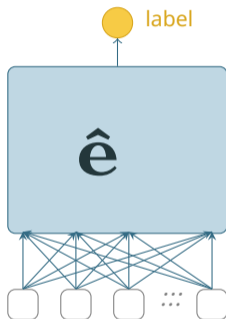
Pretrain – finetune – transfer

Pretraining



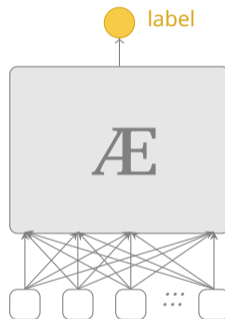
Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum. Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do

Finetuning



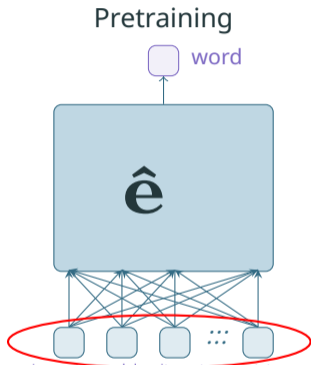
Task-specific input text

Transfer



Input text in another language

Pretrain – finetune – transfer



Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum. Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do

Encoding input text

Map character sequences (“subword tokens”) to numeric representations

Non-standard orthographies + tokenization

Die Lammer hat ein recht sauberes Wasser

D' Lomma hod a rechd a sauwas Wossa

Sentence via VUf " k] _] dYX] U" cf [#k] _] #@€aaU
Chan ea "German's next language model" 2020

Non-standard orthographies + tokenization

G Vkcfxhc_Yb]nU]cb with GBERT

Die	Lammer	hat	ein	recht	sauberes	Wasser		
[8] Y	@Jaa	-Yf	\Uh	Y] b	f YW\h	gUi VYf	-Yg	KUggYf

D'	Lomma	hod	a	rechd	a	sauwas	Wossa
----	-------	-----	---	-------	---	--------	-------

Sentence via VUf " k] _] dYX] U" cf [#k] _] #@€aaU
Chan ea "German's next language model" 2020

Non-standard orthographies + tokenization

Die Lammer hat ein recht sauberes Wasser with GBERT

Die Lammer hat ein recht sauberes Wasser
[8] Y [@Jaa -Yf [\Uh [Y] b [f YW\h [gUi VYf [-Yg [KUggYf

D' Lomma hod a rechd a sauwas Wossa
[8] fi [@ca [-aU [\c [-X [U [f YW\ [-X [U [gUi [-kUg [Kc [-ggU

Sentence via VUf " k] _] dYX] U" cf [#k] _] #@€aaU
Chan ea "German's next language model" 2020

Non-standard orthographies + tokenization

Die Lammer hat ein recht sauberes Wasser with GBERT

Die Lammer hat ein recht sauberes Wasser
[8] Y [@Jaa -Yf [\Uh [Y] b [f YW h [gUi VYf [-Yg [KUggYf

D' Lomma hod a rechd a sauwas Wossa
[8] fi [@ca [-aU [\c [-X [U [f YW [-X [U [gUi [-kUg [Kc [-ggU

Sidenote:

ChatGPT/GPT-4/etc also use such kinds of tokenization

Sentence via VUf " k] _] dYX] U" cf [#k] _] #@€aaU
Chan ea "German's next language model" 2020

Non-standard orthographies + tokenization

D' Lomma hod a rechd a sauwas Wossa
8 fi @ca -aU \c -X U fYWX -X U gUi -kUg Kc -ggU

Solutions – ongoing work in the field

Non-standard orthographies + tokenization

D' Lomma hod a rechd a sauwas Wossa
8 fi @ca -aU \c -X U fYWX -X U gUi -kUg Kc -ggU

Solutions – ongoing work in the field

- Move away from character-sequence-based representations altogether – requires training new models, experimental (overall performance might be worse)

Non-standard orthographies + tokenization

D' Lomma hod a rechd a sauwas Wossa
8 fi @ca -aU \c -X U fYWX -X U gUi -kUg Kc -ggU

Solutions – ongoing work in the field

- Move away from character-sequence-based representations altogether – requires training new models, experimental (overall performance might be worse)
- (How to) make existing models more robust to variation

Brittleness towards uncommon structures

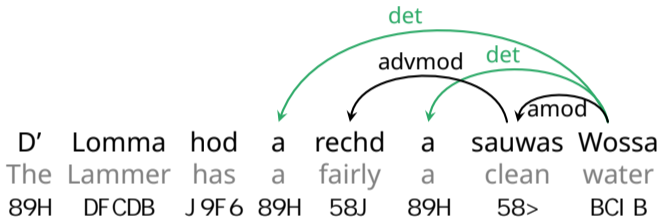
Kcf_ i bXf fy Jk

D'	Lomma	hod	a	rechd	a	sauwas	Wossa
The	Lammer	has	a	fairly	a	clean	water
89H	DFCDB	J9F6	89H	58J	89H	58>	BCI B

Sentence via VUf" k]_] dYX] U" cf [#k]_] #@€aaU

Brittleness towards uncommon structures

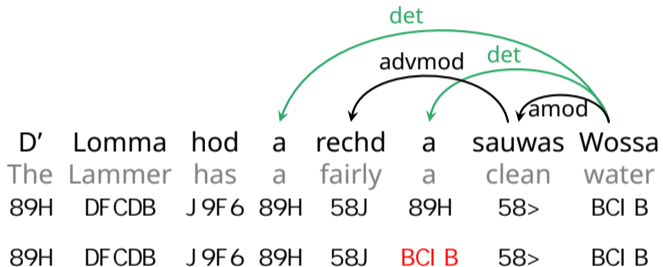
Kcf_ i bXf fy JYk



Sentence via VUf " k] _] dYX] U" cf [#k] _] # @€aaU

Brittleness towards uncommon structures

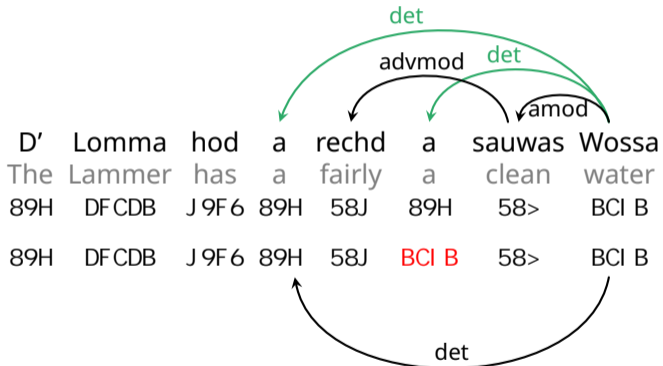
Kcf_ i bXf fy Jk



Sentence via VUf " k] _] dYX] U" cf [#k] _] #@€aaU

Brittleness towards uncommon structures

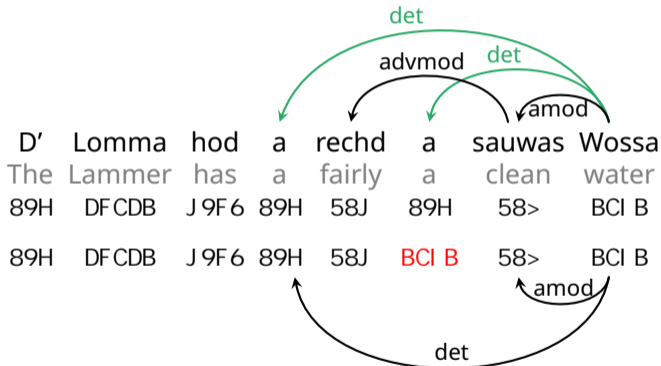
Kcf_ i bXf fy JYk



Sentence via VUf " k] _] dYX] U" cf [#k] _] #@€aaU

Brittleness towards uncommon structures

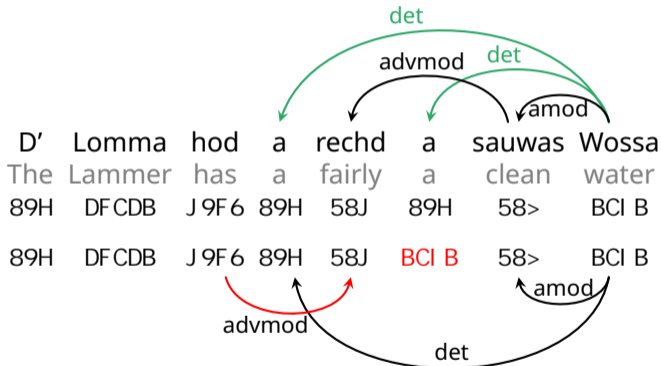
Kcf_ i bXf fy JYk



Sentence via VUf " k] _] dYX] U" cf [#k] _] #@€aaU

Brittleness towards uncommon structures

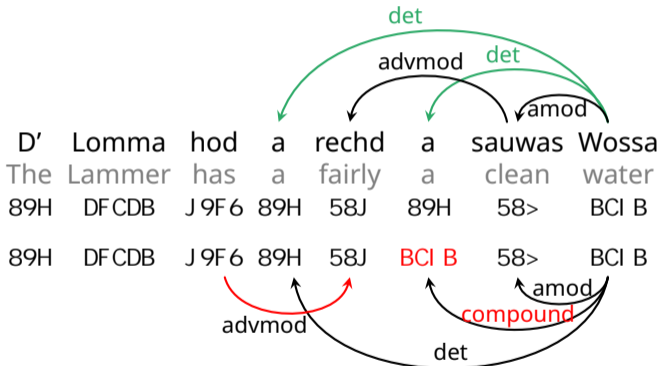
Kcf_ i bXf fy JYk



Sentence via VUf " k] _] dYX] U" cf [#k] _] #@€aaU

Brittleness towards uncommon structures

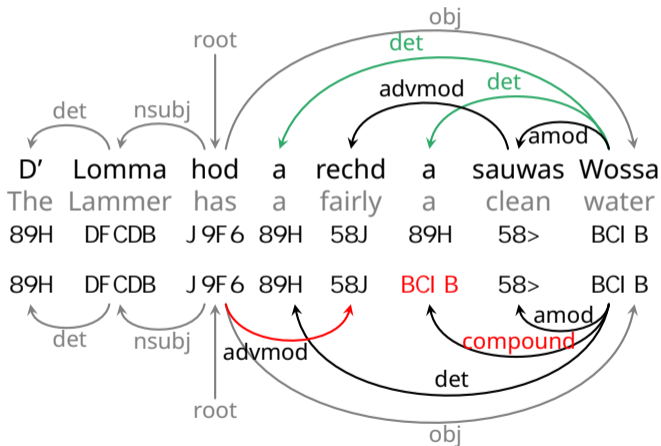
Kcf_ i bXf fy JYk



Sentence via VUf " k] _] dYX] U" cf [#k] _] # @€aaU

Brittleness towards uncommon structures

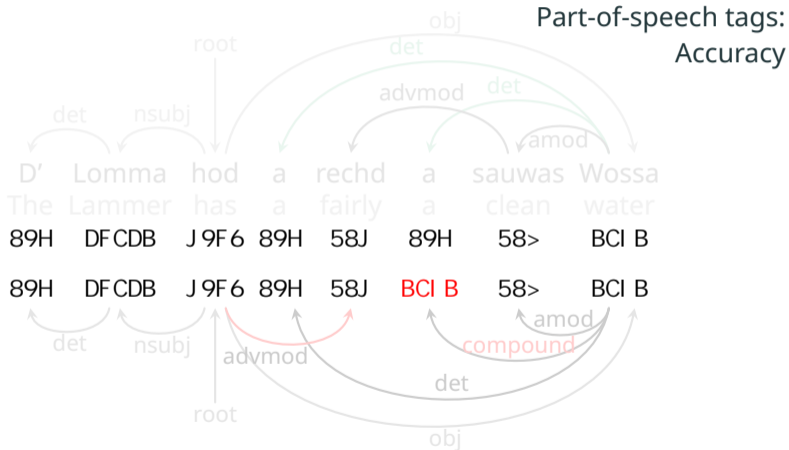
Kcf_ i bXf fy Jk



Sentence via VUf " k] _] dYX] U" cf [#k] _] # @ € aaU

Brittleness towards uncommon structures

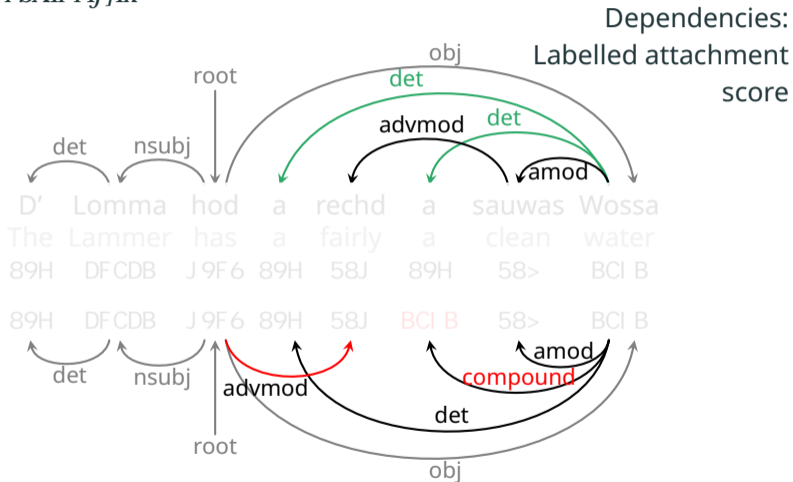
Kcf_ i bXf fy JYk



Sentence via VUf " k] _] dYX] U" cf [#k] _] # @€aaU

Brittleness towards uncommon structures

Kcf_ i bXf fy JYk



Sentence via VUf " k] _] dYX] U" cf [#k] _] #@€aaU

Automatic tagging/parsing

Kcf_ i bXf fy Jk

Train on German data (there is no Bavarian training data!),
test on German vs. Bavarian

Model	Test lang	Acc (%)	LAS (%)
Stanza	DEU	95.9	83.7
GBERT	DEU	96.8	83.1
UDPipe	DEU	96.5	84.9

Acc = accuracy (part-of-speech tags); LAS = labelled attachment score

Automatic tagging/parsing

Kcf_ i bXf fy Jk

Train on German data (there is no Bavarian training data!),
test on German vs. Bavarian

Model	Test lang	Acc (%)	LAS (%)
Stanza	DEU	95.9	83.7
GBERT	DEU	96.8	83.1
UDPipe	DEU	96.5	84.9
Stanza	BAR	40.9	23.1
GBERT	BAR	57.4	30.1
UDPipe	BAR	80.5	67.3

Acc = accuracy (part-of-speech tags); LAS = labelled attachment score

Automatic tagging/parsing

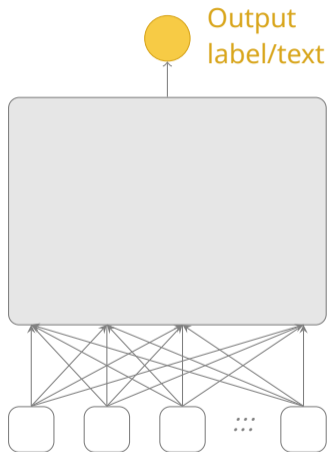
Kcf_ i bXf fy Jk

Train on German data (there is no Bavarian training data!),
test on German vs. Bavarian

Model	Test lang	Acc (%)	LAS (%)	Input representation
Stanza	DEU	95.9	83.7	
GBERT	DEU	96.8	83.1	
UDPipe	DEU	96.5	84.9	
Stanza	BAR	40.9	23.1	Full words
GBERT	BAR	57.4	30.1	Subword tokens
UDPipe	BAR	80.5	67.3	Subword tok. + characters

Acc = accuracy (part-of-speech tags); LAS = labelled attachment score

Overview



Input text sequence goes here

Human-centric NLP
(what tools and why?)

Modelling non-standard data

Available dialect data

What NLP tools and why?

Computational linguistics & machine learning research

- Quantitative patterns
- How to learn from sparse + heterogeneous data?

Blokland et al. "Language documentation meets language technology" *Computational Linguistics* (2015)

What NLP tools and why?

Computational linguistics & machine learning research

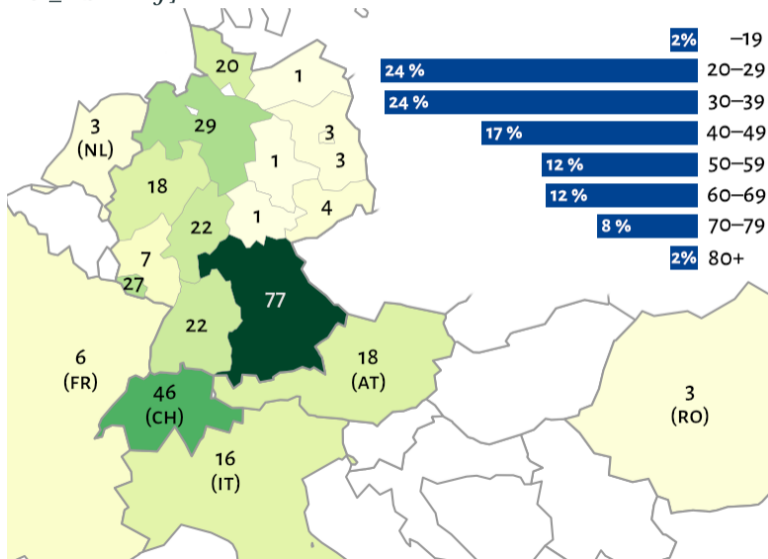
- Quantitative patterns
- How to learn from sparse + heterogeneous data?

NLP tools for linguists – also require dialogue between the communities!

Blokland et al. "Language documentation meets language technology" *Computational Linguistics* (2015)

Language technology for dialect speakers

Kcf_i bXXf fy JYk



Language technology for dialect speakers

Kcf_ i bXf fy Jk

Self-identified dialect speakers

- How often / in what contexts do people speak/write their dialects?
- What do they think about various language-related technologies for their dialects? (e.g., spellcheckers, machine translation, digital assistants...)

Language technology for dialect speakers

Kcf_ i bXf fy Jk

Self-identified dialect speakers

- How often / in what contexts do people speak/write their dialects?
- What do they think about various language-related technologies for their dialects? (e.g., spellcheckers, machine translation, digital assistants...)

Results

- Dialect input > dialect output

Language technology for dialect speakers

Kcf_ i bXf fy Jk

Self-identified dialect speakers

- How often / in what contexts do people speak/write their dialects?
- What do they think about various language-related technologies for their dialects? (e.g., spellcheckers, machine translation, digital assistants...)

Results

- Dialect input > dialect output
- Processing speech > processing text

Language technology for dialect speakers

Kcf_ i bXf fy Jk

Self-identified dialect speakers

- How often / in what contexts do people speak/write their dialects?
- What do they think about various language-related technologies for their dialects? (e.g., spellcheckers, machine translation, digital assistants...)

Results

- Dialect input > dialect output
- Processing speech > processing text
- Speaker communities vary in their attitudes!

Language technology for dialect speakers

Kcf_ i bXf fy Jk

Self-identified dialect speakers

- How often / in what contexts do people speak/write their dialects?
- What do they think about various language-related technologies for their dialects? (e.g., spellcheckers, machine translation, digital assistants...)

Results

- Dialect input > dialect output
- Processing speech > processing text
- Speaker communities vary in their attitudes!
 - Between subgroups (e.g., GSW vs. NDS speakers)

Language technology for dialect speakers

Kcf_ i bXf fy Jk

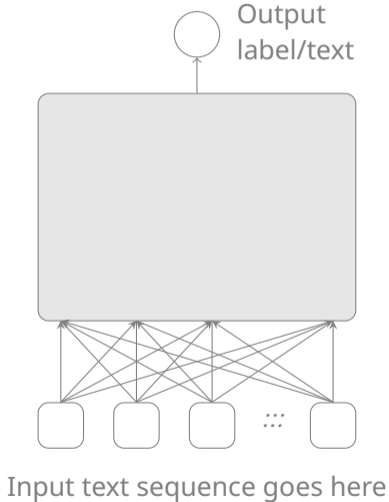
Self-identified dialect speakers

- How often / in what contexts do people speak/write their dialects?
- What do they think about various language-related technologies for their dialects? (e.g., spellcheckers, machine translation, digital assistants...)

Results

- Dialect input > dialect output
- Processing speech > processing text
- Speaker communities vary in their attitudes!
 - Between subgroups (e.g., GSW vs. NDS speakers)
 - Within subgroups

Summary



Reflecting on
what tools we build

Representing/modelling
non-standard data

Data availability
! []h i V W a # a U] b ` d #
[Y f a U b] W f ! W f d c f U

Summary



Reflecting on
what tools we build

Representing/modelling
non-standard data

Data availability
! [h i v w # a U] b d #
[Y f a U b] W f ! W f d c f U