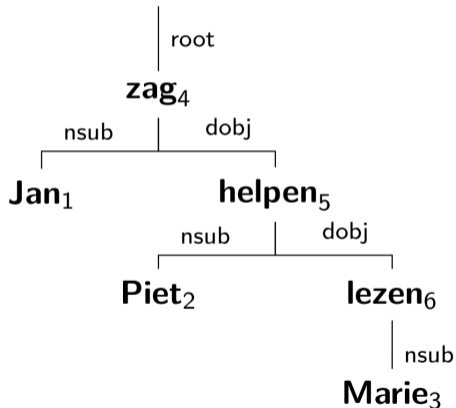


Parsing of lexicalised LCFRS via supertagging

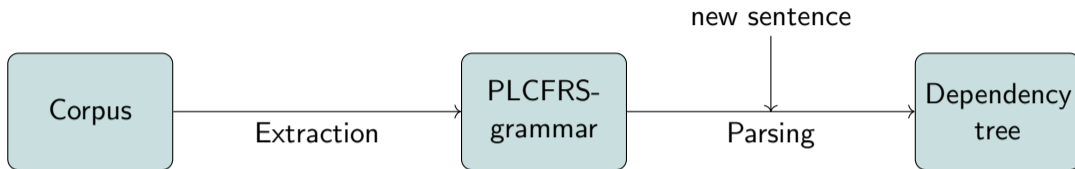
Alex Ivliev

TU Dresden

0 Jan 1 Piet 2 Marie 3 zag 4 helpen 5 lezen 6

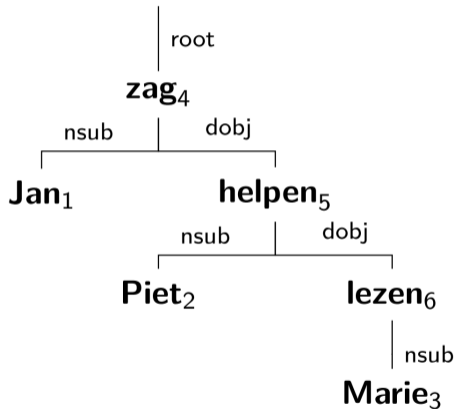


- 1 Grammar extraction and parsing
- 2 Supertagging
- 3 BERT
- 4 BERT for supertagging
- 5 Experimental results

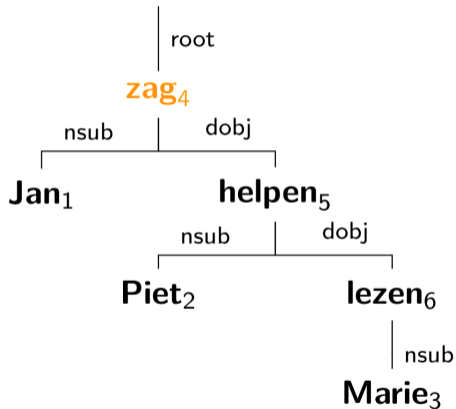


- Corpus: Hamburg Dependency Treebank
- sourced from the German news site heise.de
- 68,801 training sentences, 17,028 test sentences

Jan Piet Marie zag helpen lezen

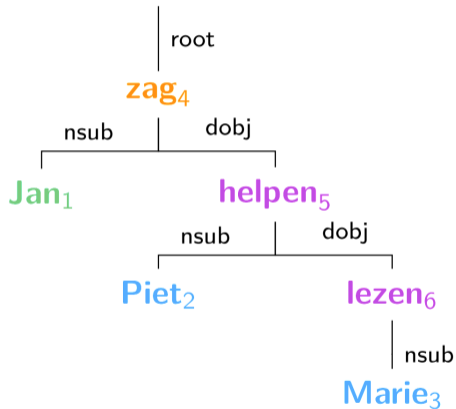


Jan Piet Marie **zag** helpen lezen



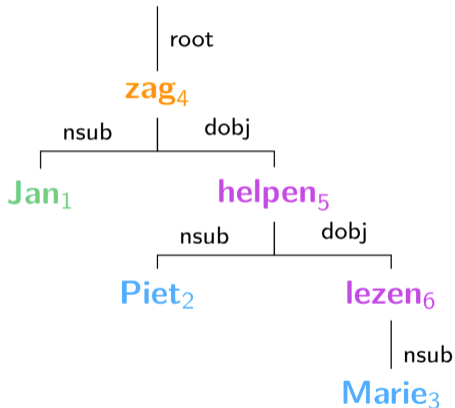
root() → nsubj()dobj()

Jan Piet Marie zag helpen lezen



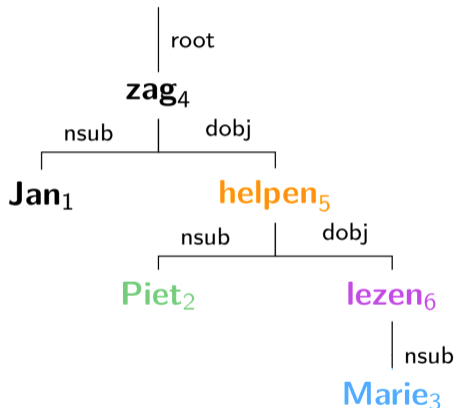
root() \rightarrow nsub(x_1)dobj(x_2, x_3)

Jan Piet Marie zag helpen lezen



$\text{root}(x_1 x_2 \text{zag} x_3) \rightarrow \text{nsub}(x_1) \text{dobj}(x_2, x_3)$

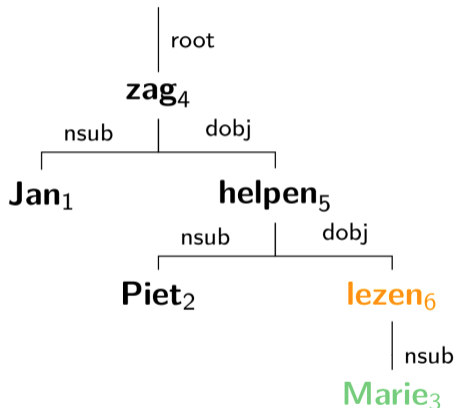
Jan Piet Marie zag helpen lezen



$$\text{root}(x_1 x_2 \text{zag} x_3) \rightarrow \text{nsub}(x_1) \text{dobj}(x_2, x_3)$$

$$\text{dobj}(x_1 x_2, \text{helpen} x_3) \rightarrow \text{nsub}(x_1) \text{dobj}(x_2, x_3)$$

Jan Piet Marie zag helpen lezen

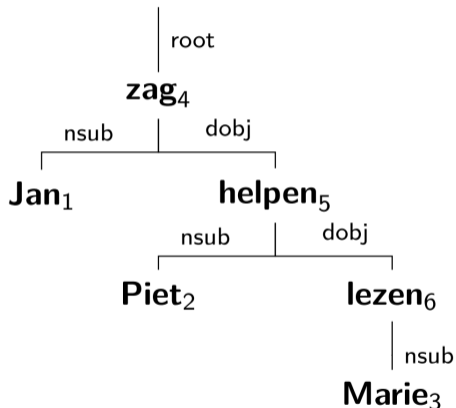


$$\text{root}(x_1 x_2 \text{zag} x_3) \rightarrow \text{nsubj}(x_1) \text{dobj}(x_2, x_3)$$

$$\text{dobj}(x_1 x_2, \text{helpen} x_3) \rightarrow \text{nsubj}(x_1) \text{dobj}(x_2, x_3)$$

$$\text{dobj}(x_1, \text{lezen}) \rightarrow \text{nsubj}(x_1)$$

Jan Piet Marie zag helpen lezen



$$\text{root}(x_1 x_2 \text{zag} x_3) \rightarrow \text{nsub}(x_1) \text{dobj}(x_2, x_3)$$

$$\text{dobj}(x_1 x_2, \text{helpen} x_3) \rightarrow \text{nsub}(x_1) \text{dobj}(x_2, x_3)$$

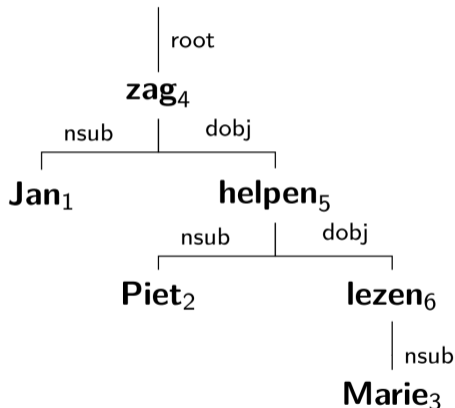
$$\text{dobj}(x_1, \text{lezen}) \rightarrow \text{nsub}(x_1)$$

$$\text{nsub}(\text{Jan}) \rightarrow \varepsilon$$

$$\text{nsub}(\text{Piet}) \rightarrow \varepsilon$$

$$\text{nsub}(\text{Marie}) \rightarrow \varepsilon$$

Jan Piet Marie zag helpen lezen



$$1.0 \text{ root}(x_1 x_2 \text{zag} x_3) \rightarrow \text{nsub}(x_1) \text{dobj}(x_2, x_3)$$

$$0.7 \text{ dobj}(x_1 x_2, \text{helpen} x_3) \rightarrow \text{nsub}(x_1) \text{dobj}(x_2, x_3)$$

$$0.3 \text{ dobj}(x_1, \text{lezen}) \rightarrow \text{nsub}(x_1)$$

$$0.1 \text{ nsub}(\text{Jan}) \rightarrow \varepsilon$$

$$0.3 \text{ nsub}(\text{Piet}) \rightarrow \varepsilon$$

$$0.6 \text{ nsub}(\text{Marie}) \rightarrow \varepsilon$$

Parseitem:

dependency	doj	0.038	weight
	$\langle (1, 3), (4, 6) \rangle$		

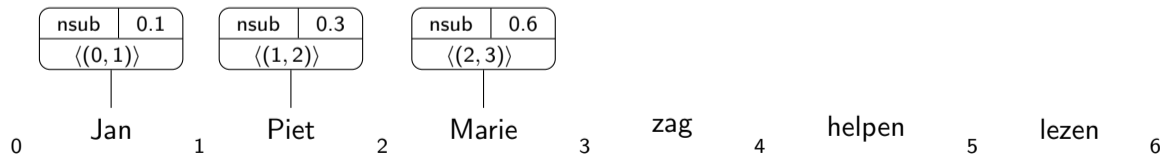
ranges

0 Jan 1 Piet 2 Marie 3 zag 4 helpen 5 lezen 6

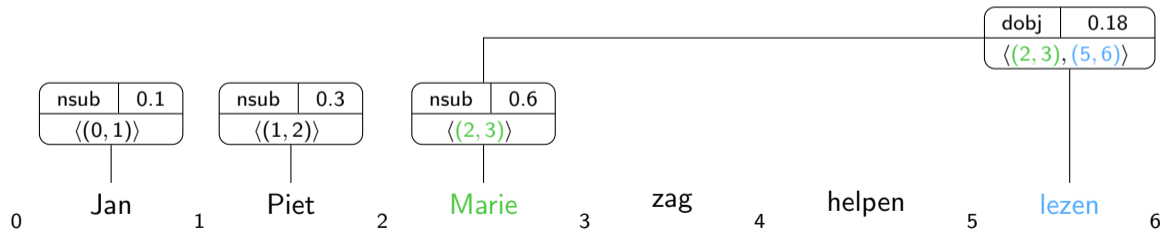
0.1: $nsub(\text{Jan}) \rightarrow \varepsilon$

0.3: $nsub(\text{Piet}) \rightarrow \varepsilon$

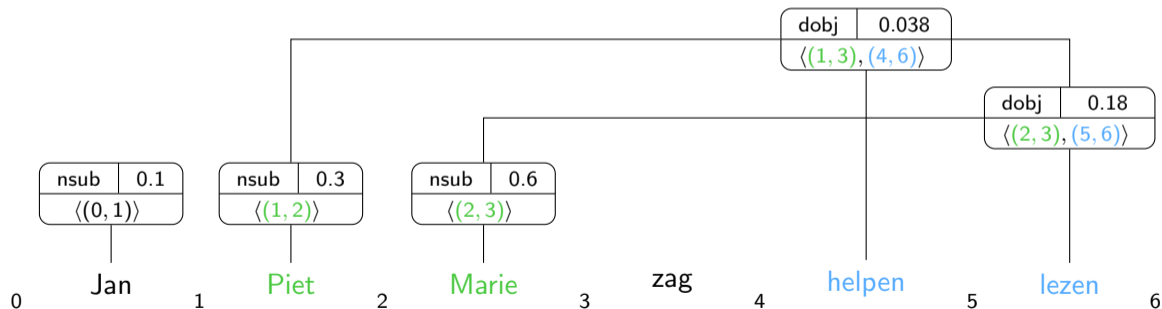
0.6: $nsub(\text{Marie}) \rightarrow \varepsilon$



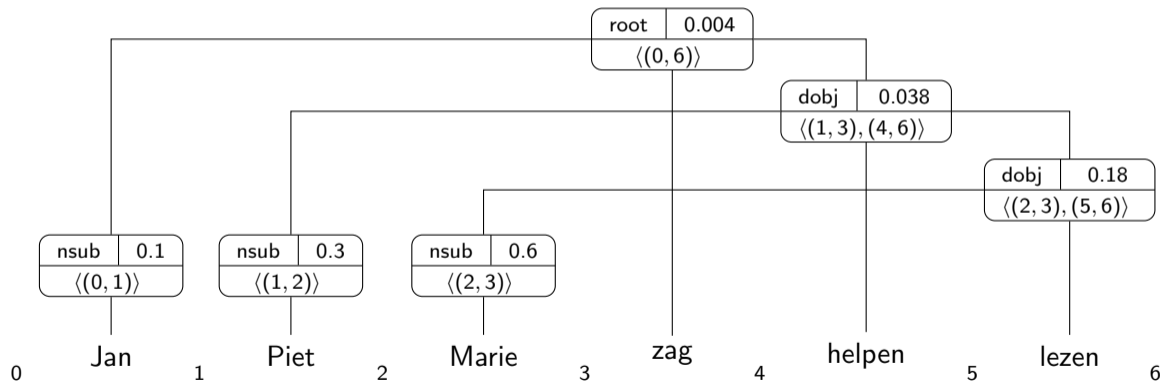
0.3: $\text{dobj}(x_1, \text{lezen}) \rightarrow \text{nsub}(x_1)$



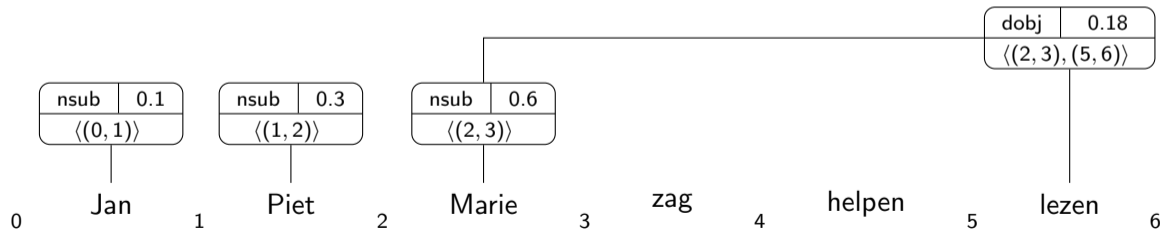
0.7: $\text{dobj}(x_1 x_2, \text{helpen} x_3) \rightarrow \text{nsub}(x_1) \text{dobj}(x_2, x_3)$



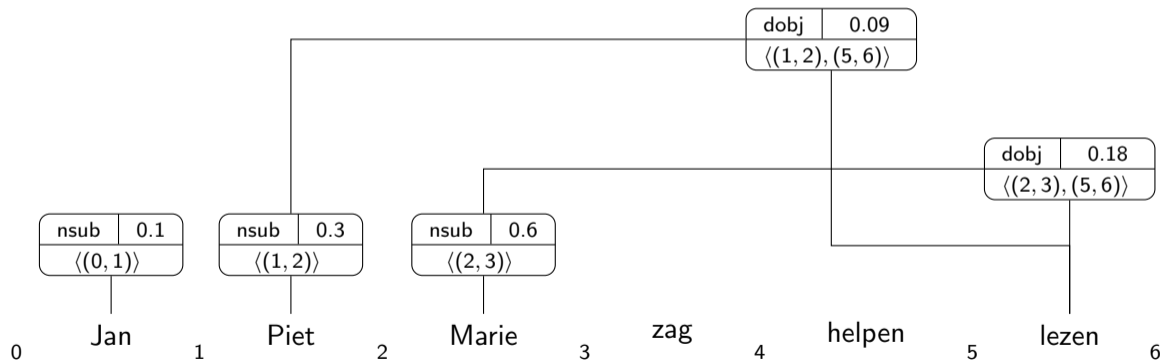
1.0: $\text{root}(x_1 x_2 \text{zag} x_3) \rightarrow \text{nsub}(x_1) \text{dobj}(x_2, x_3)$

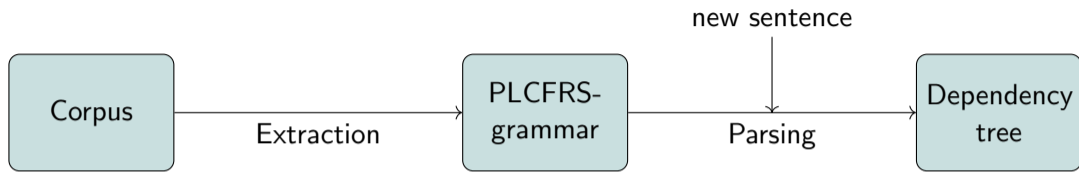


0.3: $\text{dobj}(x_1, \text{lezen}) \rightarrow \text{nsub}(x_1)$

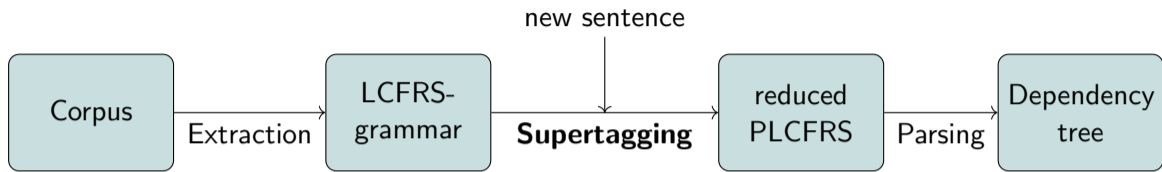


0.3: $\text{dobj}(x_1, \text{lezen}) \rightarrow \text{nsub}(x_1)$

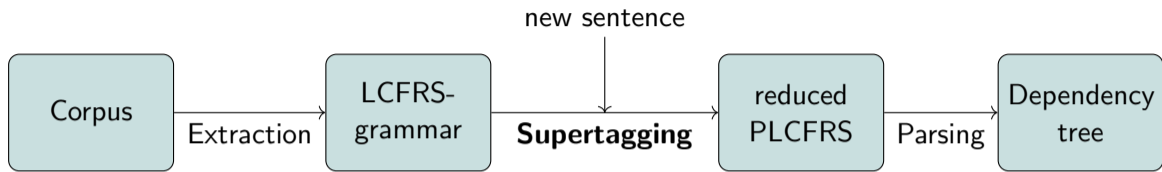




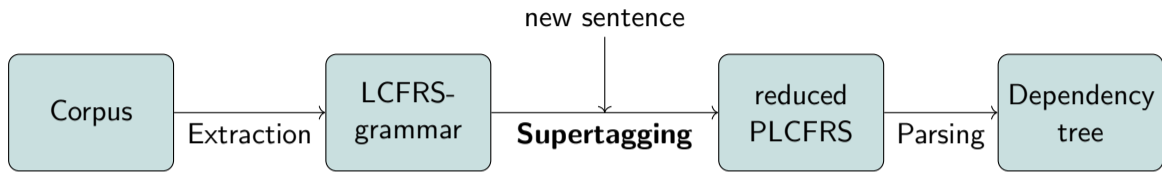
- Problem: Parsing is expensive for huge grammars
- Solution: Pre-selection of small number of promising rules



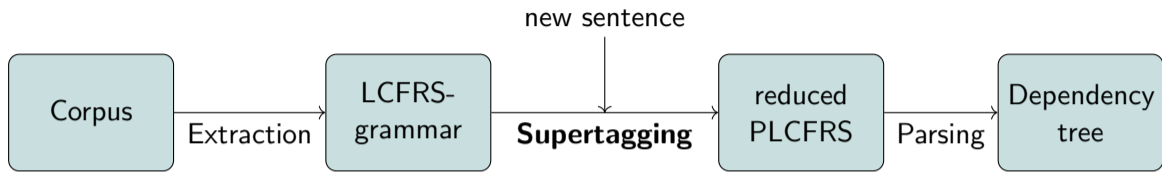
- Problem: Parsing is expensive for huge grammars
- Solution: Pre-selection of small number of promising rules



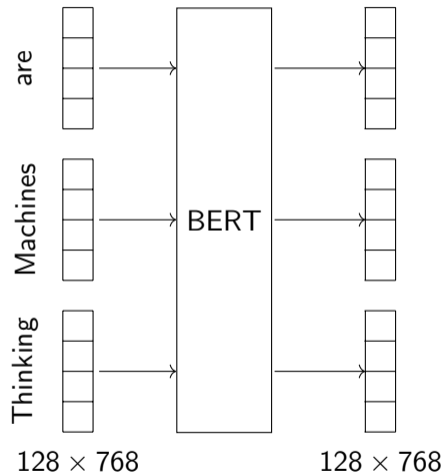
Jan	$\text{nsub}(\text{Jan}) \rightarrow \varepsilon$	$\text{subj}(\text{Jan}) \rightarrow \varepsilon$
Piet	$\text{dobj}(\text{Piet}) \rightarrow \varepsilon$	$\text{nsub}(\text{Piet}) \rightarrow \varepsilon$
Marie	$\text{subj}(x_1, \text{Marie}) \rightarrow \text{nsub}(x_1)$	$\text{nsub}(\text{Marie}) \rightarrow \varepsilon$
zag	$\text{root}(x_1 x_2 \text{zag} x_3) \rightarrow \text{nsub}(x_1) \text{dobj}(x_2, x_3)$	$\text{aux}(x_1 \text{zag} x_2) \rightarrow \text{dobj}(x_1, x_2)$
helpen	$\text{dobj}(x_1 x_2, \text{helpen} x_3) \rightarrow \text{nsub}(x_1) \text{dobj}(x_2, x_3)$	$\text{dobj}(x_1, \text{helpen}) \rightarrow \text{nsub}(x_1)$
lezen	$\text{dobj}(\text{lezen}, x_1) \rightarrow \text{nsub}(x_1)$	$\text{dobj}(x_1, \text{lezen}) \rightarrow \text{nsub}(x_1)$



Jan	$\text{nsub}(\text{David}) \rightarrow \varepsilon$	$\text{subj}(\text{Jan}) \rightarrow \varepsilon$
Piet	$\text{dobj}(\text{Piet}) \rightarrow \varepsilon$	$\text{nsub}(\text{Piet}) \rightarrow \varepsilon$
Marie	$\text{subj}(x_1, \text{Marie}) \rightarrow \text{nsub}(x_1)$	$\text{nsub}(\text{Marie}) \rightarrow \varepsilon$
zag	$\text{root}(x_1 x_2 \text{zag} x_3) \rightarrow \text{nsub}(x_1) \text{dobj}(x_2, x_3)$	$\text{aux}(x_1 \text{zag} x_2) \rightarrow \text{dobj}(x_1, x_2)$
helpen	$\text{dobj}(x_1 x_2, \text{helpen} x_3) \rightarrow \text{nsub}(x_1) \text{dobj}(x_2, x_3)$	$\text{dobj}(x_1, \text{helpen}) \rightarrow \text{nsub}(x_1)$
lezen	$\text{dobj}(\text{lezen}, x_1) \rightarrow \text{nsub}(x_1)$	$\text{dobj}(x_1, \text{lezen}) \rightarrow \text{nsub}(x_1)$



Jan	$nsub(\langle T \rangle) \rightarrow \varepsilon$	$subj(\langle T \rangle) \rightarrow \varepsilon$
Piet	$dobj(\langle T \rangle) \rightarrow \varepsilon$	$nsub(\langle T \rangle) \rightarrow \varepsilon$
Marie	$subj(x_1, \langle T \rangle) \rightarrow nsub(x_1)$	$nsub(\langle T \rangle) \rightarrow \varepsilon$
zag	$root(x_1 x_2 \langle T \rangle x_3) \rightarrow nsub(x_1) dobj(x_2, x_3)$	$aux(x_1 \langle T \rangle x_2) \rightarrow dobj(x_1, x_2)$
helpen	$dobj(x_1 x_2, \langle T \rangle x_3) \rightarrow nsub(x_1) dobj(x_2, x_3)$	$dobj(x_1, \langle T \rangle) \rightarrow nsub(x_1)$
lezen	$dobj(\langle T \rangle, x_1) \rightarrow nsub(x_1)$	$dobj(x_1, \langle T \rangle) \rightarrow nsub(x_1)$



Network is trained on

1 Masked Language Model

The man went to the [REDACTED] to buy a [REDACTED] of milk.

Network is trained on

1 Masked Language Model

The man went to the store to buy a gallon of milk.

Network is trained on

1 Masked Language Model

The man went to the store to buy a gallon of milk.

2 Next Sentence Prediction

A: The man went to the store.

B: He bought a gallon of milk.

A: The man went to the store.

B: Penguins are flightless.

Network is trained on

1 Masked Language Model

The man went to the store to buy a gallon of milk.

2 Next Sentence Prediction

A: The man went to the store.

B: He bought a gallon of milk.

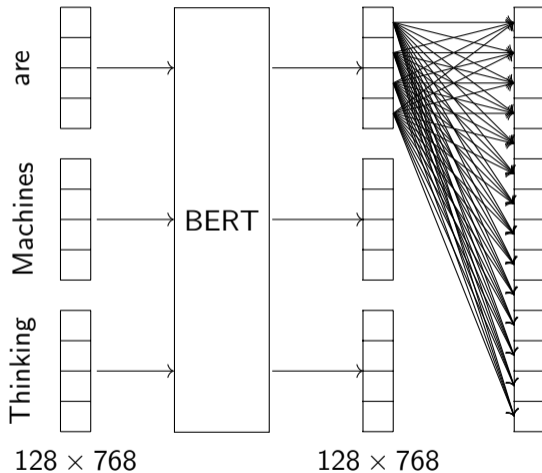
A: The man went to the store.

B: Penguins are flightless.

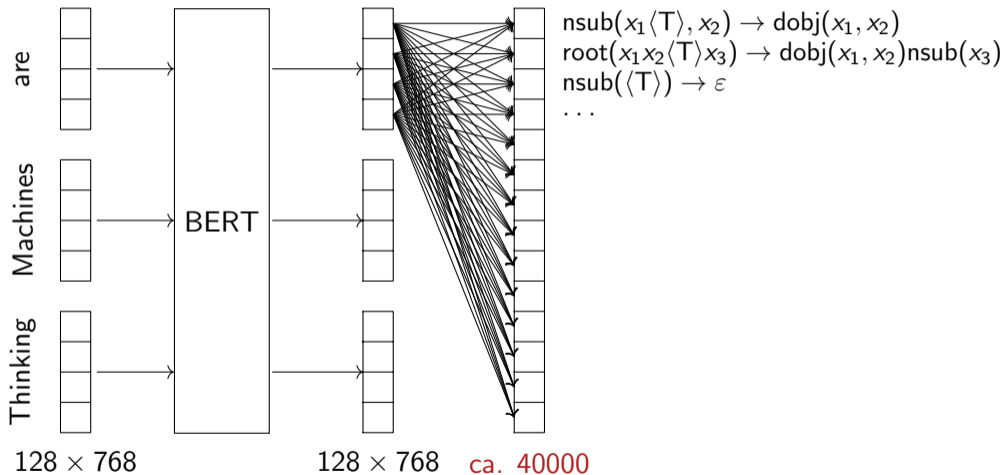
→ Pretrained model has an understanding of language

→ Pretrained model is available online

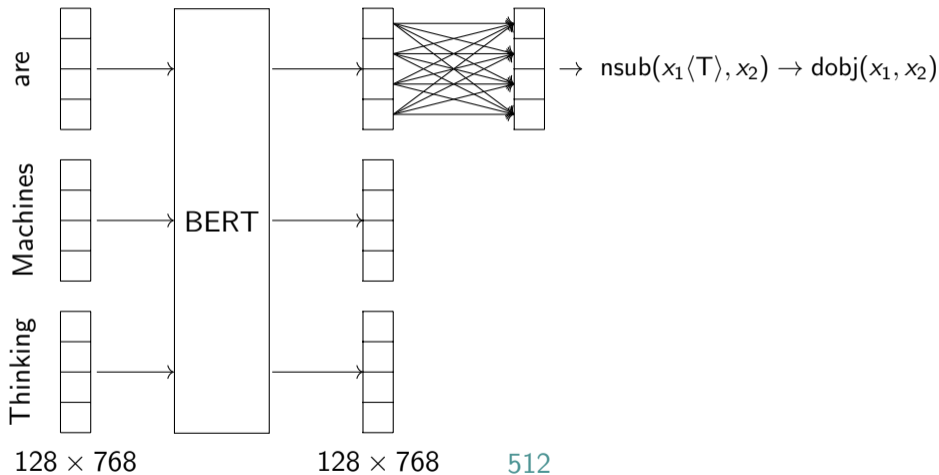
BERT - Finetuning



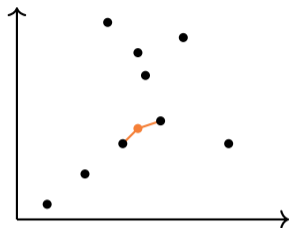
BERT - Finetuning



BERT - Finetuning



Embedding of LCFRS-rules



- LCFRS rules
- Output BERT

- Assign a vector of real numbers to each LCFRS-rule
- BERT outputs a vector
- Choose the n closest vectors

→ Nearby vectors should represent “similar” rules

Similarity measures - Tokens

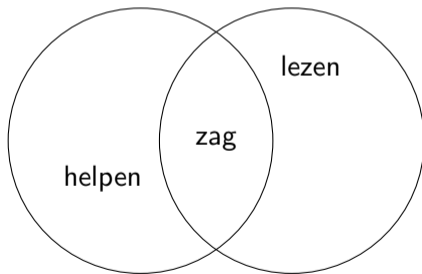
$$r = \text{dobj}(x_1 \langle T \rangle, x_2) \rightarrow \text{nsup}(x_1, x_2)$$

$\rightarrow T(r) = \{\text{zag}, \text{helpen}, \text{lezen}\}$

$\rightarrow \text{dobj}(x_1 \text{zag}, x_2) \rightarrow \text{nsup}(x_1, x_2)$

$\rightarrow \text{dobj}(x_1 \text{helpen}, x_2) \rightarrow \text{nsup}(x_1, x_2)$

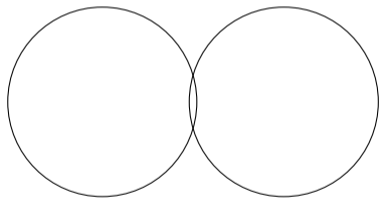
$\rightarrow \text{dobj}(x_1 \text{lezen}, x_2) \rightarrow \text{nsup}(x_1, x_2)$



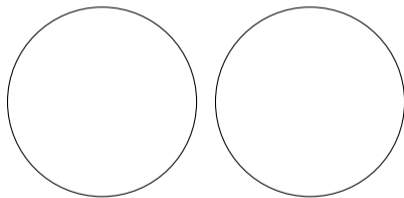
Similarity measures - Tokens

$r = \text{dobj}(x_1 \langle T \rangle, x_2) \rightarrow \text{nsub}(x_1, x_2)$
 $\rightarrow T(r) = \{\text{zag}, \text{helpen}, \text{lezen}\}$

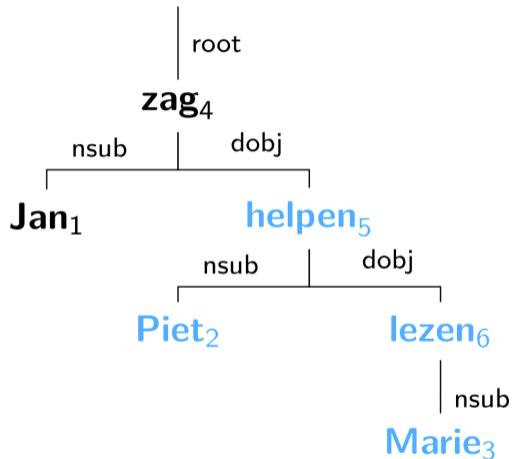
$\rightarrow \text{dobj}(x_1 \text{zag}, x_2) \rightarrow \text{nsub}(x_1, x_2)$
 $\rightarrow \text{dobj}(x_1 \text{helpen}, x_2) \rightarrow \text{nsub}(x_1, x_2)$
 $\rightarrow \text{dobj}(x_1 \text{lezen}, x_2) \rightarrow \text{nsub}(x_1, x_2)$



1%



99%



$\rightarrow r = \text{dobj}(x_1 x_2, \langle T \rangle x_3) \rightarrow \text{nsub}(x_1) \text{dobj}(x_2, x_3)$

$C(r) \supseteq \{\text{helpen, Piet, lezen, Marie}\}$

Similarity measures - String distances

$\text{dobj}(x_1, \langle T \rangle x_2 x_3) \rightarrow \text{nsub}(x_1) \text{dobj}(x_2, x_3)$

$\text{dobj}(x_1, x_2 \langle T \rangle x_3) \rightarrow \text{nsub}(x_1) \text{dobj}(x_2, x_3)$

vs.

$\text{root}(x_1 x_2 x_3, \langle T \rangle) \rightarrow \text{aux}(x_1, x_2) \text{nsub}(x_3)$

$\text{dobj}(\langle T \rangle x_1, x_2) \rightarrow \text{nsub}(x_1) \text{det}(x_2)$

Similarity measures - String distances

$\text{dobj}(x_1, \langle T \rangle x_2 x_3) \rightarrow \text{nsub}(x_1) \text{dobj}(x_2, x_3)$

vs.

$\text{root}(x_1 x_2 x_3, \langle T \rangle) \rightarrow \text{aux}(x_1, x_2) \text{nsub}(x_3)$

$\text{dobj}(x_1, x_2 \langle T \rangle x_3) \rightarrow \text{nsub}(x_1) \text{dobj}(x_2, x_3)$

$\text{dobj}(\langle T \rangle x_1, x_2) \rightarrow \text{nsub}(x_1) \text{det}(x_2)$

String of rule: $A(x_1 x_2 \langle T \rangle, x_3) \rightarrow B(x_1) C(x_2, x_3) \rightsquigarrow BC \langle T \rangle, C$

Similarity measures - String distances

$\text{dobj}(x_1, \langle T \rangle x_2 x_3) \rightarrow \text{nsub}(x_1) \text{dobj}(x_2, x_3)$
 $\text{dobj}(x_1, x_2 \langle T \rangle x_3) \rightarrow \text{nsub}(x_1) \text{dobj}(x_2, x_3)$

vs.

$\text{root}(x_1 x_2 x_3, \langle T \rangle) \rightarrow \text{aux}(x_1, x_2) \text{nsub}(x_3)$
 $\text{dobj}(\langle T \rangle x_1, x_2) \rightarrow \text{nsub}(x_1) \text{det}(x_2)$

String of rule: $A(x_1 x_2 \langle T \rangle, x_3) \rightarrow B(x_1) C(x_2, x_3) \rightsquigarrow BC \langle T \rangle, C$

Bigram	BC⟨T⟩,C	⟨T⟩,CBCB	Δ
BC	1	1	0
C⟨T⟩	1	0	1
⟨T⟩,	1	1	0
,C	1	1	0
CB	0	2	2

Similarity measures - String distances

$\text{dobj}(x_1, \langle T \rangle x_2 x_3) \rightarrow \text{nsub}(x_1) \text{dobj}(x_2, x_3)$
 $\text{dobj}(x_1, x_2 \langle T \rangle x_3) \rightarrow \text{nsub}(x_1) \text{dobj}(x_2, x_3)$

vs. $\text{root}(x_1 x_2 x_3, \langle T \rangle) \rightarrow \text{aux}(x_1, x_2) \text{nsub}(x_3)$
 $\text{dobj}(\langle T \rangle x_1, x_2) \rightarrow \text{nsub}(x_1) \text{det}(x_2)$

String of rule: $A(x_1 x_2 \langle T \rangle, x_3) \rightarrow B(x_1) C(x_2, x_3) \rightsquigarrow BC \langle T \rangle, C$

Bigram	BC⟨T⟩,C	⟨T⟩,CBCB	Δ
BC	1	1	0
C⟨T⟩	1	0	1
⟨T⟩,	1	1	0
,C	1	1	0
CB	0	2	2

Damerau–Levenshtein distance

- edit distance between two sequences
- minimum number of operations to change one word into the other
- insertions, deletions or substitutions, transposition

Training of the embedding

Gradient descent algorithm

$$\text{Loss} = \sum_{r_1 \in P} \sum_{r_2 \in P} \mu(r_1, r_2) \cdot \text{dist}(e(r_1), e(r_2)) + \lambda(\|e(r_1)\|_2^2 + \|e(r_2)\|_2^2) \rightarrow \min$$

Training of the embedding

Gradient descent algorithm

$$\text{Loss} = \sum_{r_1 \in P} \sum_{r_2 \in P} \mu(r_1, r_2) \cdot \text{dist}(e(r_1), e(r_2)) + \lambda(\|e(r_1)\|_2^2 + \|e(r_2)\|_2^2) \rightarrow \min$$

$\mu: P \times P \rightarrow [-1, 1]$... Similarity measure

Training of the embedding

Gradient descent algorithm

$$\text{Loss} = \sum_{r_1 \in P} \sum_{r_2 \in P} \mu(r_1, r_2) \cdot \text{dist}(e(r_1), e(r_2)) + \lambda(\|e(r_1)\|_2^2 + \|e(r_2)\|_2^2) \rightarrow \min$$

$\mu: P \times P \rightarrow [-1, 1]$... Similarity measure

$e: P \times P \rightarrow \mathbb{R}^n$... Embedding function

Training of the embedding

Gradient descent algorithm

$$\text{Loss} = \sum_{r_1 \in P} \sum_{r_2 \in P} \mu(r_1, r_2) \cdot \text{dist}(e(r_1), e(r_2)) + \lambda(\|e(r_1)\|_2^2 + \|e(r_2)\|_2^2) \rightarrow \min$$

$\mu: P \times P \rightarrow [-1, 1]$... Similarity measure

$e: P \times P \rightarrow \mathbb{R}^n$... Embedding function

$\|\cdot\|_2^2 = \sum_{i=1}^n (x_i)^2$... Square of L2-Norm

Training of the embedding

Gradient descent algorithm

$$\text{Loss} = \sum_{r_1 \in P} \sum_{r_2 \in P} \mu(r_1, r_2) \cdot \text{dist}(e(r_1), e(r_2)) + \lambda(\|e(r_1)\|_2^2 + \|e(r_2)\|_2^2) \rightarrow \min$$

$\mu: P \times P \rightarrow [-1, 1]$... Similarity measure

$e: P \times P \rightarrow \mathbb{R}^n$... Embedding function

$\|\cdot\|_2^2 = \sum_{i=1}^n (x_i)^2$... Square of L2-Norm

$\lambda \in \mathbb{R}$... Regularization factor

Training of the embedding

Gradient descent algorithm

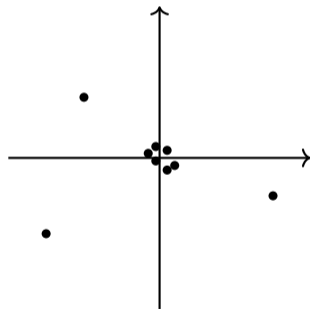
$$\text{Loss} = \sum_{r_1 \in P} \sum_{r_2 \in P} \mu(r_1, r_2) \cdot \text{dist}(e(r_1), e(r_2)) + \lambda(\|e(r_1)\|_2^2 + \|e(r_2)\|_2^2) \rightarrow \min$$

$\mu: P \times P \rightarrow [-1, 1]$... Similarity measure

$e: P \times P \rightarrow \mathbb{R}^n$... Embedding function

$\|\cdot\|_2^2 = \sum_{i=1}^n (x_i)^2$... Square of L2-Norm

$\lambda \in \mathbb{R}$... Regularization factor



Gradient descent algorithm

$$\text{Loss} = \sum_{r_1 \in P} \sum_{r_2 \in P} (1 - \mu(r_1, r_2) - \text{dist}(e(r_1), e(r_2)))^2 \rightarrow \min$$

$\mu: P \times P \rightarrow [0, 1]$... Similarity measure

$e: P \times P \rightarrow \mathbb{R}^n$... Embedding function

$\|\cdot\|_2^2 = \sum_{i=1}^n (x_i)^2$... Square of L2-Norm

	$N = 25$	$N = 50$	$N = 75$	$N = 100$
Bigram	66.61/0.00	13.15/0.98	1.33/8.6	0.07/19.63
Trigram	95.76/0.00	10.09/0.12	0.75/1.52	0.38/3.69
Levenshtein	99.31/0.26	90.00/8.84	71.72/26.16	52.56/44.53

- Format: Parse failure / timeout in % (Timeout = 60s)
- Levenshtein is unable to parse most sentences
- Larger N decrease parse failure % but increase timeout %

	$N = 25$	$N = 50$	$N = 75$	$N = 100$
Bigram	42/19	1965/189	8468/422	16725/1815
Trigram	27/15	922/191	3344/540	4923/705
Levenshtein	165/1	5401/3	15691/4	27170/14

- Format: Time average / median in ms
- Time increases with larger N
- Trigram is faster than Bigram
- Median much lower than average

	$N = 25$	$N = 50$	$N = 75$	$N = 100$
Bigram	76.14/82.00	60.74/67.54	60.44/67.95	61.98/69.76
Trigram	44.46/64.32	8.82/21.51	10.77/23.59	12.49/25.15
Levenshtein	90.79/96.65	90.42/96.64	90.10/97.00	89.01/96.35

- Format: unabled / labeled in %
- Trigram much worse than Bigram
- Labeled and unlabeled precision similar

- Supertagging: Pre-selection of small number of rules to increase performance
- Neural network BERT can be used
- requires vector encoding of LCFRS-rules
- training embedding based on similarity measures
- of the ones tested, bigram performed the best



DEVLIN, J. ; CHANG, M.-W. ; LEE, K. ; TOUTANOVA, K. :

BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding.

In: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*.

Association for Computational Linguistics, 4171–4186



JOSHI, A. K. ; SRINIVAS, B. :

Disambiguation of Super Parts of Speech (or Supertags): Almost Parsing.

In: *COLING 1994 Volume 1: The 15th International Conference on Computational Linguistics*



KUHLMANN, M. ; SATTÀ, G. :

Treebank Grammar Techniques for Non-Projective Dependency Parsing.

In: *Proceedings of the 12th Conference of the European Chapter of the ACL (EACL 2009)*.

Association for Computational Linguistics, 478–486



NEDERHOF, M. :

Weighted Deductive Parsing and Knuth's Algorithm.

In: *Computational Linguistics* 29 (2003), Nr. 1, S. 135–143