

Neural Natural Language Processing in WebLicht

State-of-the-Art and Future Perspectives

Daniël de Kok

What is WebLicht?

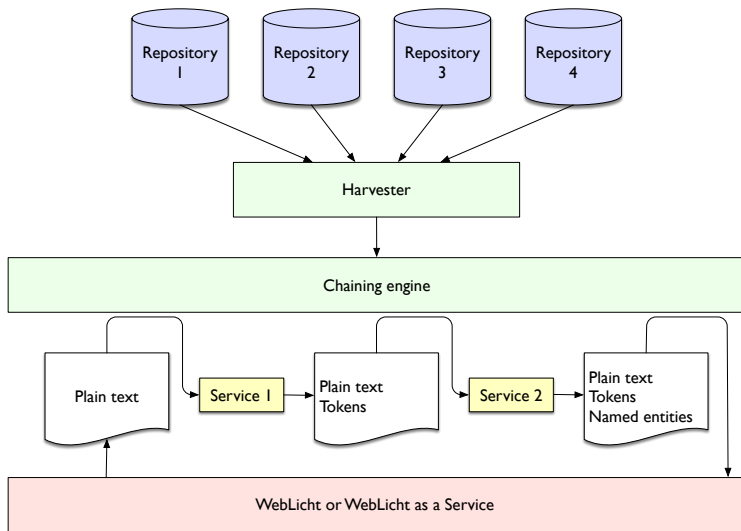
H.264

What is WebLicht?

H.264

Sherlock Holmes^{PER} ist eine 1886 vom britischen Schriftsteller Arthur Conan Doyle^{PER} geschaffene Kunstfigur .

WebLicht architecture



Text Corpus Format (TCF)

- ▶ XML data exchange format
- ▶ Annotations stored layer-wise
- ▶ Used by WebLicht webservices
- ▶ Tools for visualization, such as TüNDRA
- ▶ Tools for annotation, such as WebAnno

Text Corpus Format (TCF)

- ▶ XML data exchange format
- ▶ Annotations stored layer-wise
- ▶ Used by WebLicht webservices
- ▶ Tools for visualization, such as TüNDRA
- ▶ Tools for annotation, such as WebAnno
- ▶ Many supported annotation layers, including:
 - ▶ Tokenization
 - ▶ Lemmas
 - ▶ Constituency & Dependency Parsing
 - ▶ Morphology
 - ▶ Named Entities
 - ▶ References (anaphoric, cataphoric, coreference)
 - ▶ Geographical locations
 - ▶ Text structure

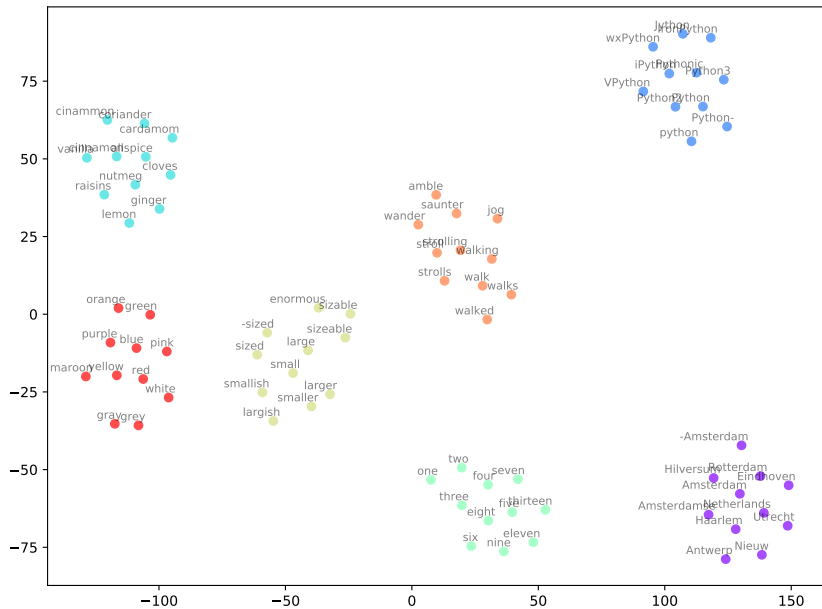
State-of-the-Art

State-of-the-Art: Neural language processing

Key ingredients:

- ▶ Vector representations of words
- ▶ Deep neural networks

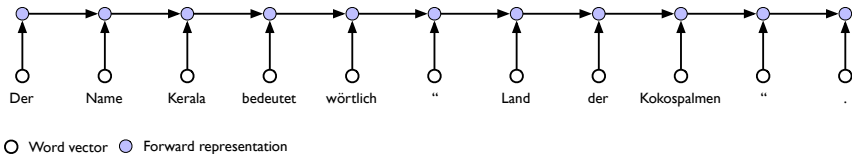
Words as vectors



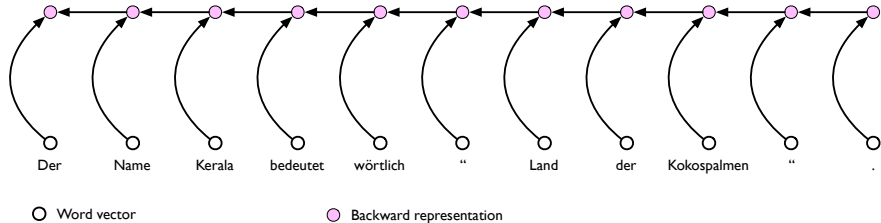
Recurrent neural networks (RNN)



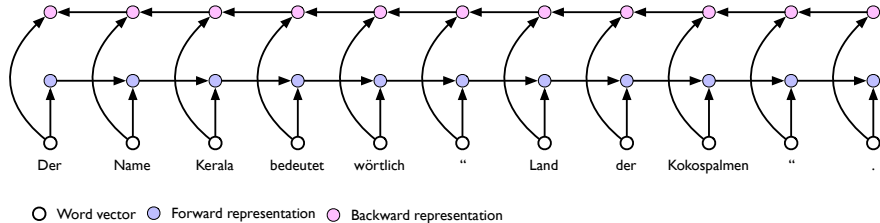
Recurrent neural networks (RNN)



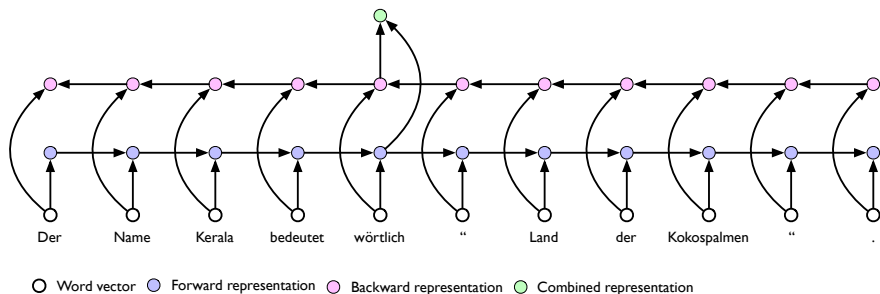
Recurrent neural networks (RNN)



Recurrent neural networks (RNN)



Recurrent neural networks (RNN)

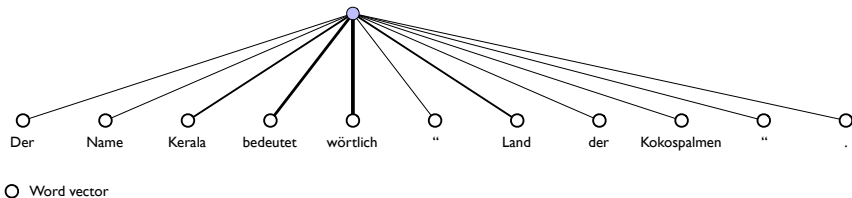


Transformer

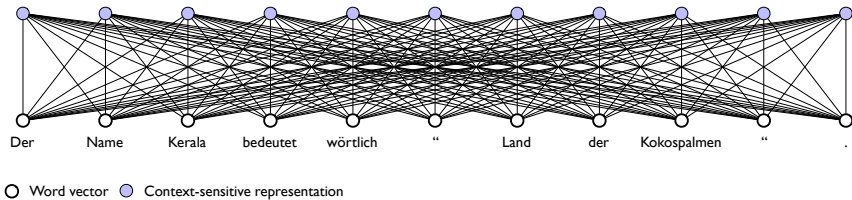
Der Name Kerala bedeutet wörtlich “ Land der Kokospalmen “ .

Word vector

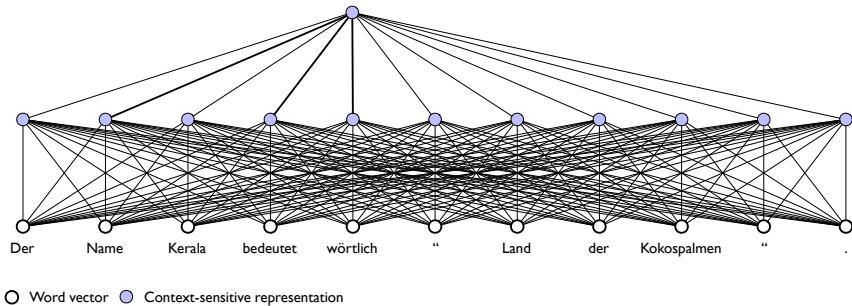
Transformer



Transformer



Transformer



Improvements in syntactic tasks

¹Joint work with Erik Schill.

²Joint work with Tobias Pütz.

Improvements in syntactic tasks

Part-of-speech tagging¹

VVFIN ART NN PTKVZ \$.
Wacht die Opposition auf ?

¹Joint work with Erik Schill.

²Joint work with Tobias Pütz.

Improvements in syntactic tasks

Part-of-speech tagging¹

VVFIN ART NN PTKVZ \$.
Wacht die Opposition auf ?

Parser	Accuracy
HMM	97.29
RNN	99.03

¹Joint work with Erik Schill.

²Joint work with Tobias Pütz.

Improvements in syntactic tasks

Part-of-speech tagging¹

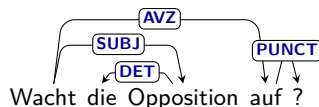
VVFIN **ART** **NN** **PTKVZ** **\$.**
Wacht die Opposition auf ?

Parser	Accuracy
---------------	-----------------

HMM	97.29
-----	-------

RNN	99.03
-----	-------

Dependency parsing²



¹Joint work with Erik Schill.

²Joint work with Tobias Pütz.

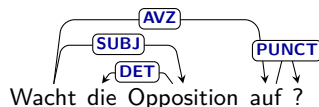
Improvements in syntactic tasks

Part-of-speech tagging¹

VVFIN **ART** **NN** **PTKVZ** **\$.**
Wacht die Opposition auf ?

Parser	Accuracy
HMM	97.29
RNN	99.03

Dependency parsing²



Tagger	Accuracy
Linear SVM	89.9
Feed forward NN	91.8
RNN	94.4
Transformer	94.7
RNN + pretraining	96.1
Transformer + pretraining	96.3

¹Joint work with Erik Schill.

²Joint work with Tobias Pütz.

Conclusion

- ▶ WebLicht allows you to build text analysis pipelines:
<https://weblicht.sfs.uni-tuebingen.de>
- ▶ WebLicht can be integrated in other applications using WebLicht as a Service:
<https://weblicht.sfs.uni-tuebingen.de/WaaS/>
- ▶ State-of-the-Art: annotations tools using neural networks

Thank you!

Future perspectives

Pretraining (Peters et al., 2018 & Devlin et al., 2019)

Pretraining

Finetuning

Deep neural model

Was sollte Daewoo gegen den betriebswirtschaftlich günstigeren Standort
[MASK] haben ? [SEP] Oder ist Bremerhaven nicht günstiger?

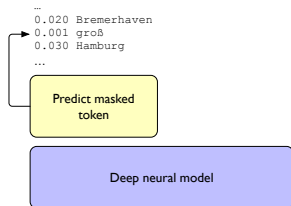
Wikipedia

Taz

Common Crawl

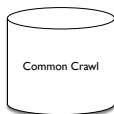
Pretraining (Peters et al., 2018 & Devlin et al., 2019)

Pretraining



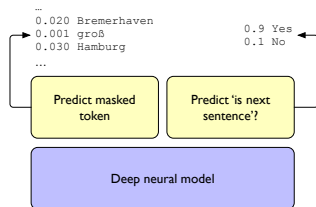
Finetuning

Was sollte Daewoo gegen den betriebswirtschaftlich günstigeren Standort [MASK] haben ? [SEP] Oder ist Bremerhaven nicht günstiger?



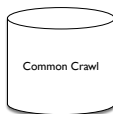
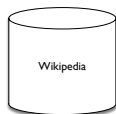
Pretraining (Peters et al., 2018 & Devlin et al., 2019)

Pretraining



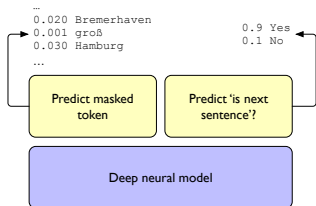
Finetuning

Was sollte Daewoo gegen den betriebswirtschaftlich günstigeren Standort [MASK] haben ? [SEP] Oder ist Bremerhaven nicht günstiger?

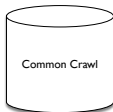


Pretraining (Peters et al., 2018 & Devlin et al., 2019)

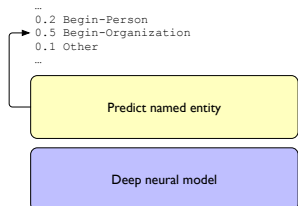
Pretraining



Was sollte Daewoo gegen den betriebswirtschaftlich günstigeren Standort [MASK] haben ? [SEP] Oder ist Bremerhaven nicht günstiger?



Finetuning



Was sollte Daewoo gegen den betriebswirtschaftlich günstigeren Standort Bremerhaven haben ?



Pretraining holds promise for text analysis

- ▶ Pretraining models seem to pick up a certain amount of syntax (Tenney et al., 2019 & Clark et al., 2019)

Pretraining holds promise for text analysis

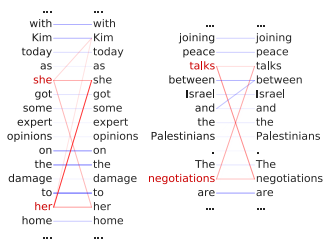
- ▶ Pretraining models seem to pick up a certain amount of syntax (Tenney et al., 2019 & Clark et al., 2019)
- ▶ Pretraining models with fine tuning shown to work well on textual entailment, paraphrase detection, and question-answer pair classification. (Devlin et al., 2019)

Pretraining holds promise for text analysis

- ▶ Pretraining models seem to pick up a certain amount of syntax (Tenney et al., 2019 & Clark et al., 2019)
- ▶ Pretraining models with fine tuning shown to work well on textual entailment, paraphrase detection, and question-answer pair classification. (Devlin et al., 2019)

Head 5-4

- **Coreferent** mentions attend to their antecedents
- 65.1% accuracy at linking the head of a coreferent mention to the head of an antecedent



(Source: Clark et al., 2019)