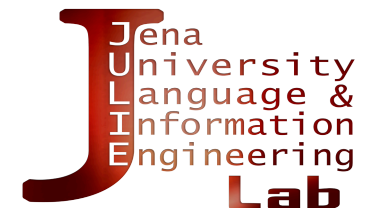


Grundlagen Kontextualisierter Wordembeddings

Erik Fäßler

Jena University Language & Information Engineering (JULIE)
Lab Friedrich-Schiller-University Jena, Germany

<http://www.julielab.de>



Language Models

Forward-Model

given the history (t_1, \dots, t_{k-1}) :

$$p(t_1, t_2, \dots, t_N) = \prod_{k=1}^N p(t_k \mid t_1, t_2, \dots, t_{k-1}).$$

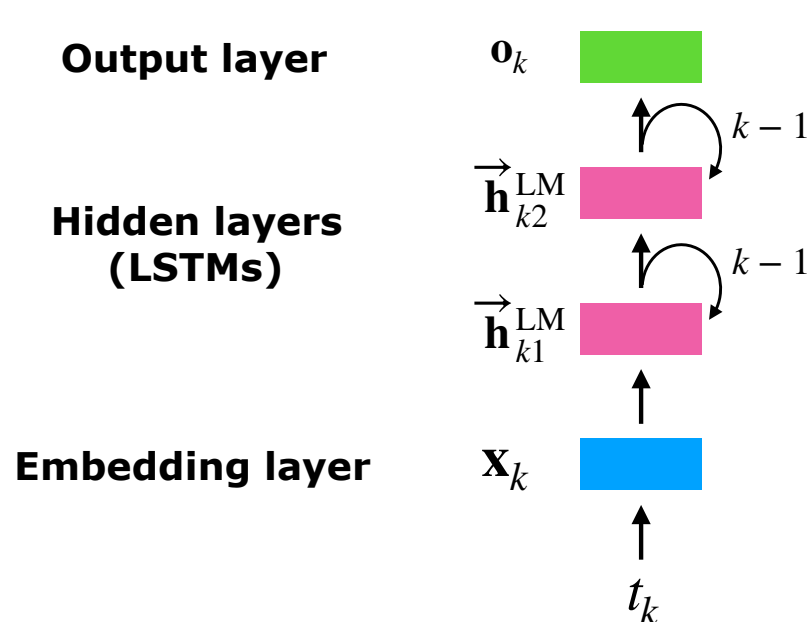
Backward-Model

$$p(t_1, t_2, \dots, t_N) = \prod_{k=1}^N p(t_k \mid t_{k+1}, t_{k+2}, \dots, t_N).$$

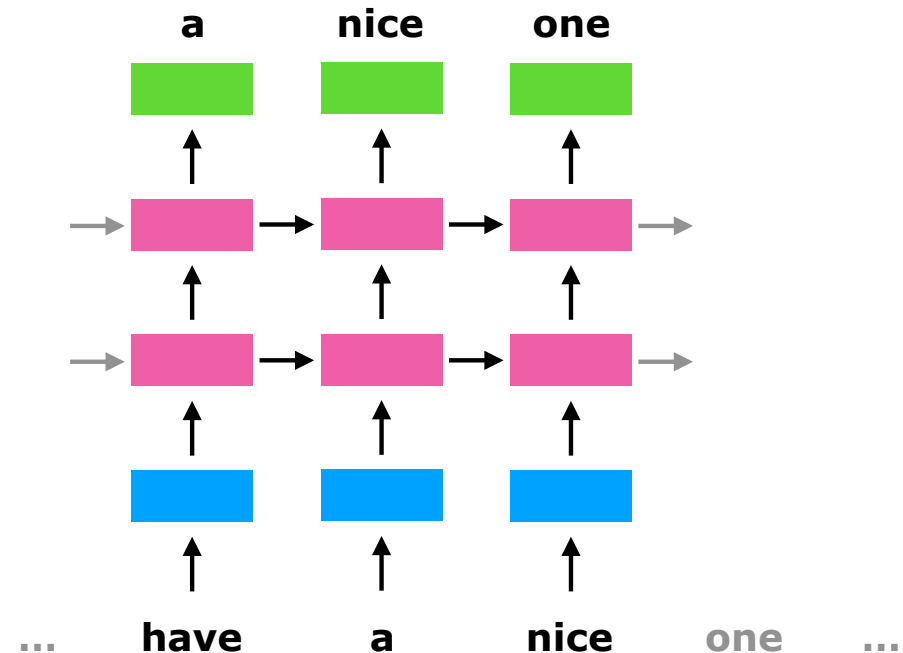
Deep LM

With long short term memory (LSTM) network,
predicting the next words in both directions to build
biLMs

The forward LM architecture



Expanded in the forward direction of k

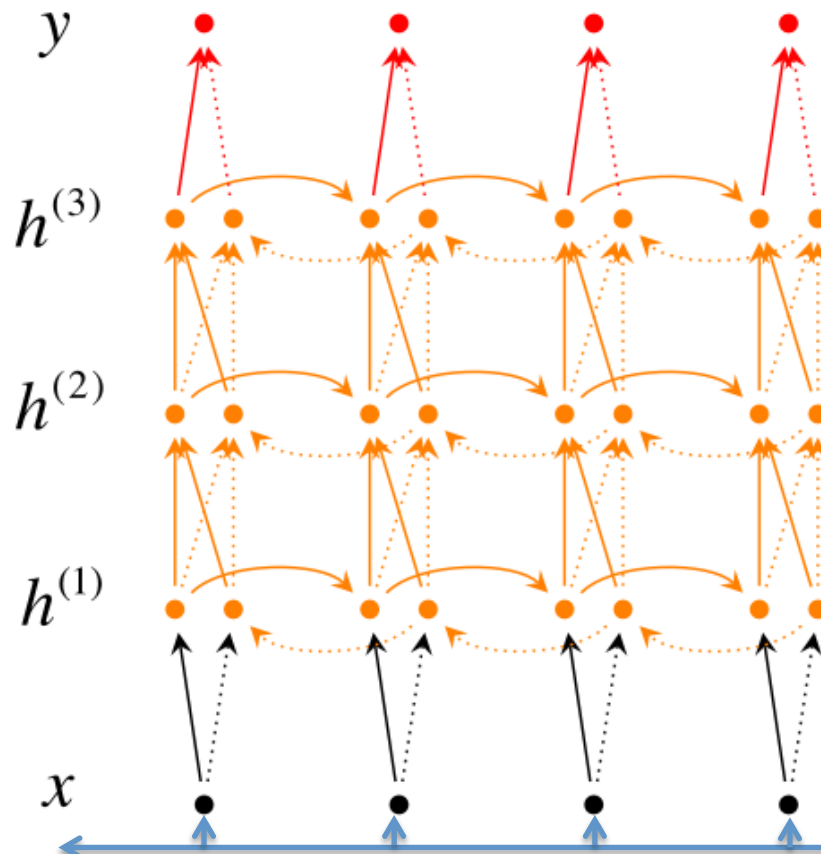


<https://de.slideshare.net/shuntaroy/a-review-of-deep-contextualized-word-representations-peters-2018>

Deep Bidirectional LM

<https://web.stanford.edu/class/archive/cs/cs224n/cs224n.1174/syllabus.html>, Lecture 8 (Feb 2), Slide 46

Deep Bidirectional RNNs



$$\vec{h}_t^{(i)} = f(\vec{W}^{(i)} h_t^{(i-1)} + \vec{V}^{(i)} \vec{h}_{t-1}^{(i)} + \vec{b}^{(i)})$$

$$\overleftarrow{h}_t^{(i)} = f(\overleftarrow{W}^{(i)} h_t^{(i-1)} + \overleftarrow{V}^{(i)} \overleftarrow{h}_{t+1}^{(i)} + \overleftarrow{b}^{(i)})$$

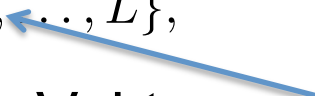
$$y_t = g(U[\vec{h}_t^{(L)}; \overleftarrow{h}_t^{(L)}] + c)$$

static word embeddings (word2vec, glove, fasttext, ...)

Each memory layer passes an intermediate sequential representation to the next.

ELMo (Peters et al.)

- BI-LSTM-LM mit L=2 Layern
- Für jeden Token-position k extrahiere
 - den statischen Embeddingvector (**CNN character n-grams!**)
 - die forward und backward biLM Zwischenrepräsentationen

$$\begin{aligned}
 R_k &= \{\mathbf{x}_k^{LM}, \overrightarrow{\mathbf{h}}_{k,j}^{LM}, \overleftarrow{\mathbf{h}}_{k,j}^{LM} \mid j = 1, \dots, L\} \\
 &= \{\mathbf{h}_{k,j}^{LM} \mid j = 0, \dots, L\},
 \end{aligned}$$


- Linearkombination der Vektoren 0 = Token layer = \mathbf{x}_k^{LM}
- Task-spezifische Parameter

$$\text{ELMo}_k^{task} = E(R_k; \Theta^{task}) = \gamma^{task} \sum_{j=0}^L s_j^{task} \mathbf{h}_{k,j}^{LM}.$$

- Ersetze im eigentlichen Aufgaben-KNN die statischen Embeddings durch ELMo

ELMo Vectors

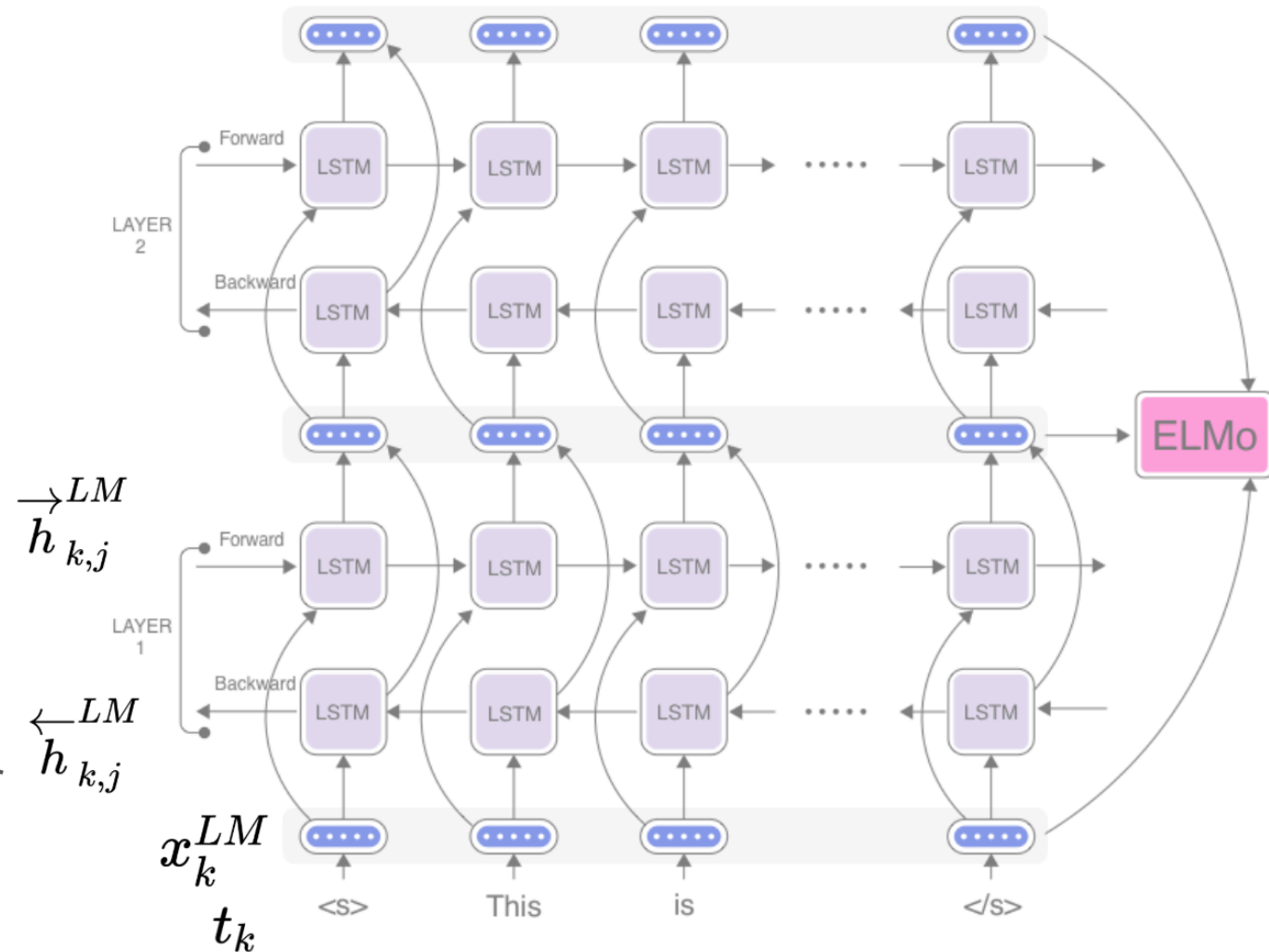
Structure

Each token t_k

L-layer biLM
computes $2L+1$
representations

k is the k -th token

j is the j -th biLM layer



<https://ireneli.eu/2018/12/17/elmo-in-practice/>

ELMo Model

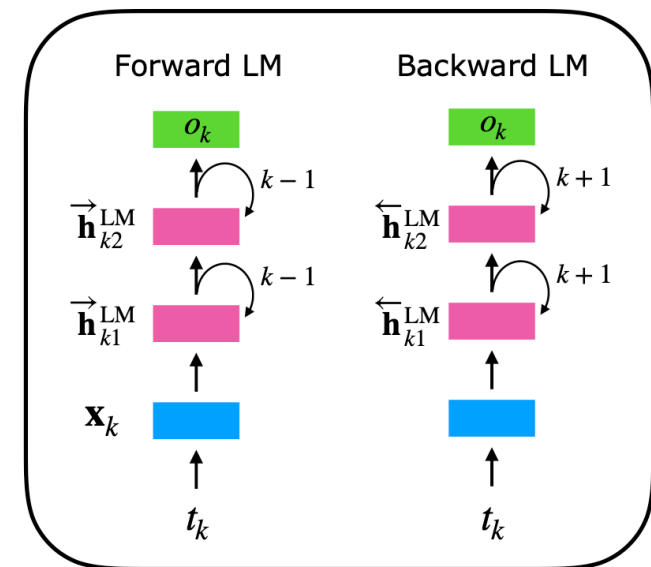
ELMo represents a word t_k as a linear combination of corresponding hidden layers (inc. its embedding)

$$\text{ELMo}_k^{\text{task}} = \gamma^{\text{task}} \times \sum \left\{ \begin{array}{l} s_2^{\text{task}} \times \mathbf{h}_{k2}^{\text{LM}} \\ s_1^{\text{task}} \times \mathbf{h}_{k1}^{\text{LM}} \\ s_0^{\text{task}} \times \mathbf{h}_{k0}^{\text{LM}} \end{array} \right. \quad \text{Concatenate hidden layers} \quad \left[\vec{\mathbf{h}}_{kj}^{\text{LM}}; \overleftarrow{\mathbf{h}}_{kj}^{\text{LM}} \right]$$

($\mathbf{x}_k; \mathbf{x}_k$)

Unlike usual word embeddings, ELMo is assigned to every *token* instead of a *type*

biLMs



37

<https://de.slideshare.net/shuntaroy/a-review-of-deep-contextualized-word-representations-peters-2018>

- Question answering (Stanford Question Answering Dataset, SQuAD)
- Textual entailment (Stanford Natural Language Inference (SNLI) corpus)
- Semantic role labeling (OntoNotes)
- Coreference solution (OntoNotes)
- Named Entity Extraction (CoNLL 2003 NER)
- Sentiment analysis

ELMo Evaluation

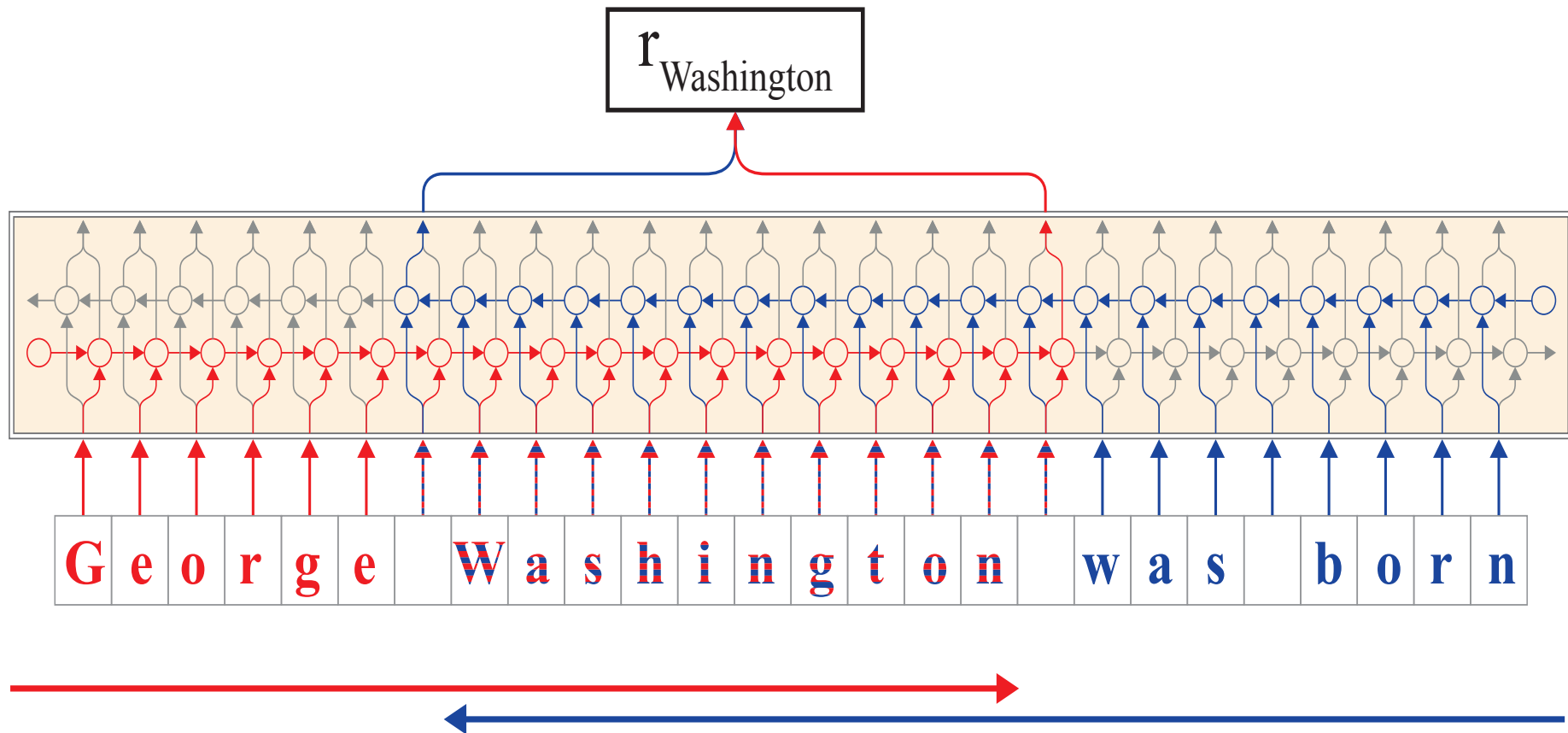
nicht-triviale KNN
Architekturen



TASK	PREVIOUS SOTA		OUR BASELINE	ELMo + BASELINE	INCREASE (ABSOLUTE/ RELATIVE)
SQuAD	Liu et al. (2017)	84.4	81.1	85.8	4.7 / 24.9%
SNLI	Chen et al. (2017)	88.6	88.0	88.7 \pm 0.17	0.7 / 5.8%
SRL	He et al. (2017)	81.7	81.4	84.6	3.2 / 17.2%
Coref	Lee et al. (2017)	67.2	67.2	70.4	3.2 / 9.8%
NER	Peters et al. (2017)	91.93 \pm 0.19	90.15	92.22 \pm 0.10	2.06 / 21%
SST-5	McCann et al. (2017)	53.7	51.4	54.7 \pm 0.5	3.3 / 6.8%

Flair (Akbik et al.)

- Basiert ebenfalls auf biLMs
- Allerdings nicht auf Token- sondern auf **Zeichenebene**
- Bisläng nur Sequenzentagging
- Im Flair Framework gibt es auch Dokumentenklassifikation
- Flair embedding für token k:
 - Forward: LM Hidden State nach letztem Token Zeichen
 - Backward: LM Hidden State vor erstem Zeichen
- Im Flair Paper und –Framework werden die unterschiedlichen Embeddings concateniert (**gestackt**)
 - Flair Forward + Backward + (GloVe | word2vec | fasttext | ...)



ELMo



Task	Language	Dataset	Flair	Previous best
Named Entity Recognition	English	Conll-03	93.18 (F1)	92.22 (<i>Peters et al., 2018</i>)
Named Entity Recognition	English	Ontonotes	89.3 (F1)	86.28 (<i>Chiu et al., 2016</i>)
Emerging Entity Detection	English	WNUT-17	49.49 (F1)	45.55 (<i>Aguilar et al., 2018</i>)
Part-of-Speech tagging	English	WSJ	97.85	97.64 (<i>Choi, 2016</i>)
Chunking	English	Conll-2000	96.72 (F1)	96.36 (<i>Peters et al., 2017</i>)
Named Entity Recognition	German	Conll-03	88.27 (F1)	78.76 (<i>Lample et al., 2016</i>)
Named Entity Recognition	German	Germeval	84.65 (F1)	79.08 (<i>Hänig et al., 2014</i>)
Named Entity Recognition	Dutch	Conll-03	90.44 (F1)	81.74 (<i>Lample et al., 2016</i>)
Named Entity Recognition	Polish	PolEval-2018	86.6 (F1) (<i>Borchmann et al., 2018</i>)	85.1 (<i>PolDeepNer</i>)