

Transfer Learning

ELMo (Peters et al., 2018)



Learning goals

- tbd

CONTEXTUAL EMBEDDINGS



Source: *Jay Alammar*

- Bidirectional language model (LM)
- Combines a forward LM

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

and a backward LM

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

to arrive at the following loglikelihood:

$$\sum_{k=1}^N \left(\log p(t_k | t_1, \dots, t_{k-1}; \Theta_x, \vec{\Theta}_{LSTM}, \Theta_s) + \log p(t_k | t_{k+1}, \dots, t_N; \Theta_x, \overleftarrow{\Theta}_{LSTM}, \Theta_s) \right)$$

ELMO EMBEDDINGS

- Character-based (context-independent) token representations

$$\mathbf{x}_k^{LM}$$

- Two-layer biLSTM as main architecture:
 - Two context-dependent token representations *per layer*, i.e.

$$\vec{\mathbf{h}}_{k,j}^{LM} \text{ \& \; } \overleftarrow{\mathbf{h}}_{k,j}^{LM} \text{ for the } k\text{-th token in the } j\text{-th layer.}$$

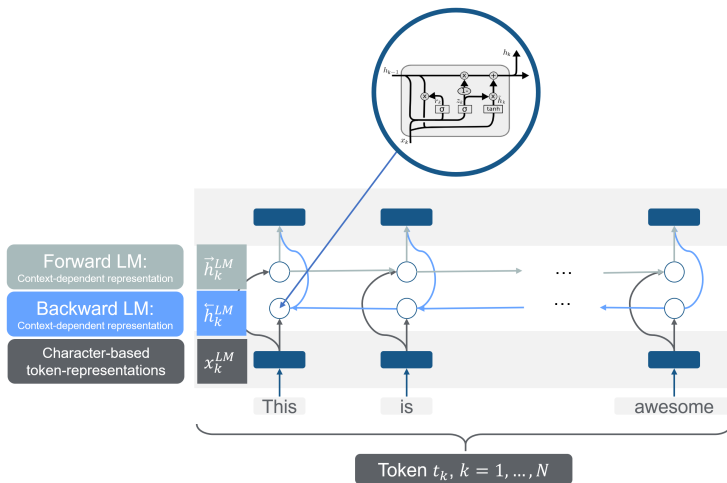
- Four context-dependent token representations in total:

$$\left\{ \vec{\mathbf{h}}_{k,j}^{LM}, \overleftarrow{\mathbf{h}}_{k,j}^{LM} \mid j = 1, 2 \right\}$$

- Five representations per token in total:

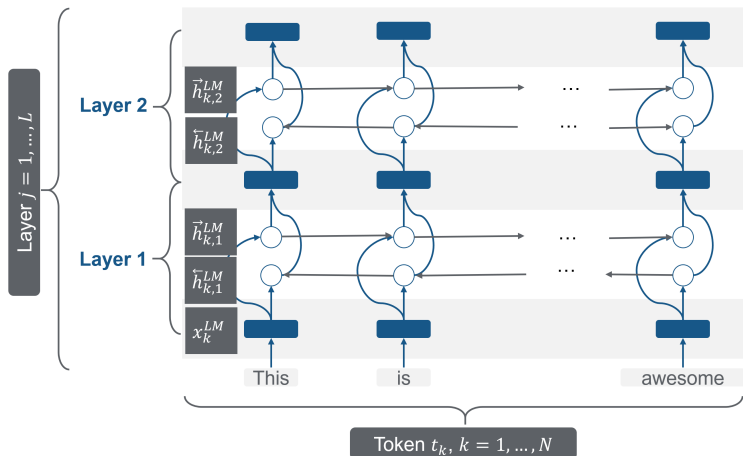
$$\begin{aligned} R_k &= \left\{ \mathbf{x}_k^{LM}, \vec{\mathbf{h}}_{k,j}^{LM}, \overleftarrow{\mathbf{h}}_{k,j}^{LM} \mid j = 1, \dots, L \right\} \\ &= \left\{ \mathbf{h}_{k,j}^{LM} \mid j = 0, 1, 2 \right\} \end{aligned}$$

GRAPHICAL REPRESENTATION



Source: *Carolyn Becker*

GRAPHICAL REPRESENTATION



Source: *Carolyn Becker*

TASK ADAPTION

Including ELMo in downstream tasks:

- Calculate task-specific weights of all five representations:

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

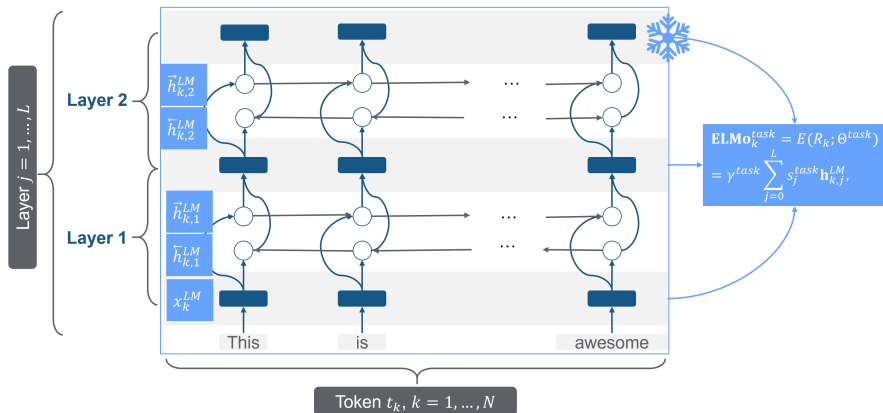
where the $\mathbf{h}_{k,j}^{\text{LM}}$ are **not trainable** anymore.

- Trainable parameters during the adaption:
 - s_j^{task} are trainable (softmax-normalized) weights
 - γ^{task} is a trainable scaling parameter

Advantages over context free-embeddings:

- Task-specific model has access to *multiple* representations of each token
- Model learns to which degree to use the different representations depending on the task at hand

TASK ADAPTION



Source: *Carolyn Becker*

FINE-TUNING APPROACH

- Pre-trained on a general domain corpus
- Embeddings are contextualized (as opposed to word2vec)
- Embeddings are not adapted to target domain/task (similar to word2vec)
- Sequential nature of LSTMs:
 - Not fully parallelizable (compared to Transformers, see next chapter)
 - Fails to capture long-range dependency during contextualization