

# Deep Learning for NLP

## Lecture 8: Encoder-Decoder Models

**Dr. Mohsen Mesgar**

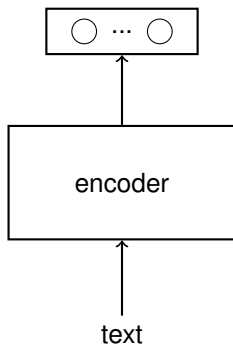
**Ubiquitous Knowledge Processing Lab (UKP Lab)**

# This lecture

- ▶ Encoder-Decoder
- ▶ Attention
- ▶ Their applications in NLP

# Encoder

A neural model to transform a text into a vector in an embedding space.

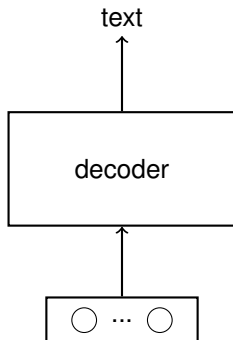


Different types of neural encoders are

- ▶ pretrained word embeddings
- ▶ MLPs, CNNs, RNNs, Transformers, ...

# Decoder

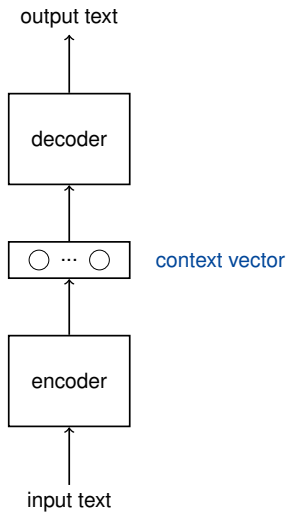
A neural model to transform a vector from an embedding space to a text.



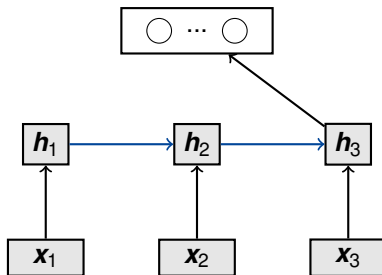
Different types of language models can be used as decoders

- ▶ RNN-based LMs
- ▶ Transformer-based LMs

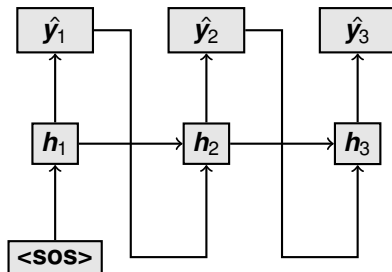
# Encoder-Decoder



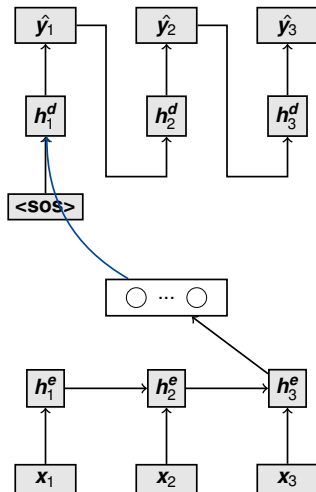
# RNN Encoder



# RNN Decoder



# RNN-based Encoder-Decoder

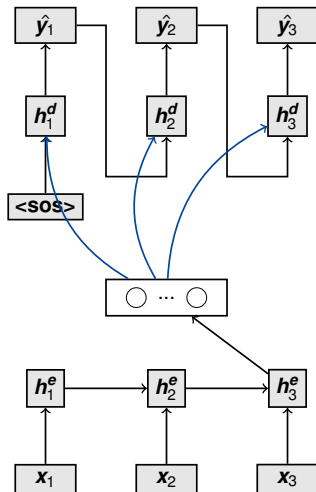




# RNN-based Encoder-Decoder

- ▶ The context vector is used to initialize the hidden state of the decoder.
- ▶ Its impact vanishes at the last steps of the decoder.

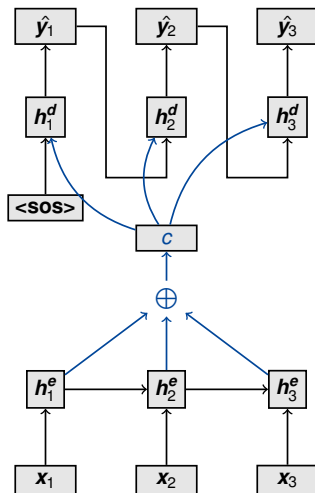
# RNN-based Encoder-Decoder



# RNN-based Encoder-Decoder

- ▶ The output of the encoder is known as the context vector.
- ▶ The dimensionality of the context vector is fixed.
- ▶ However, different input texts might have different length.
- ▶ So, considering the hidden state of the RNN encoder may not capture the entire input text.
- ▶ This is a problem especially for long input texts.

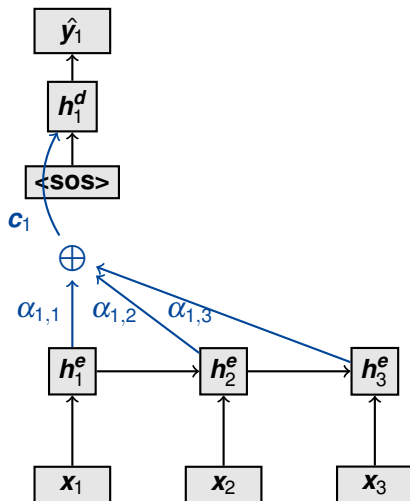
# RNN-based Encoder-Decoder



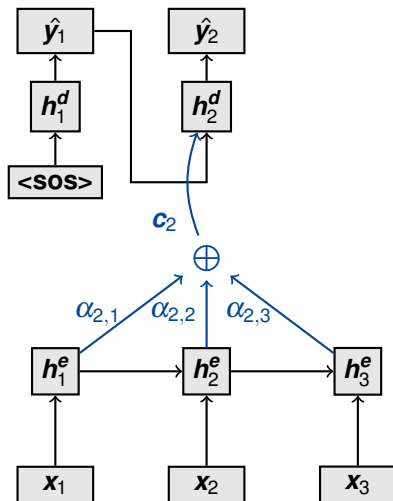
# RNN-based Encoder-Decoder

- ▶ The other problem is that context vector is unique for all decoding steps.
- ▶ The encoder treats all tokens of the input sentence equally important to produce a context vector.
- ▶ However, at any decoding step, the decoder should focus on tokens of the input sentence differently.

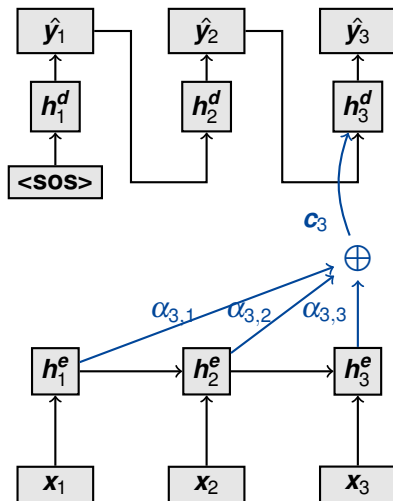
# Attention



# Attention



# Attention





# Attention

$$\mathbf{c}_t = \sum_{k=1}^N \alpha_{t,k} \mathbf{h}_k^e$$

$$\alpha_{t,k} = \frac{\exp(\text{score}(\mathbf{h}_{t-1}^d, \mathbf{h}_k^e))}{\sum_{k'=1}^N \exp(\text{score}(\mathbf{h}_{t-1}^d, \mathbf{h}_{k'}^e))}$$

# Content-based Attention (Graves et al., 2014)

$$\text{score}(\mathbf{h}_t^d, \mathbf{h}_k^e) = \text{cosine}(\mathbf{h}_t^d, \mathbf{h}_k^e)$$

# Additive Attention (Bahdanau et al., 2015)

$$\text{score}(\mathbf{h}_t^d, \mathbf{h}_k^e) = \tanh([\mathbf{h}_t^d; \mathbf{h}_k^e] \mathbf{W}^{(h)}) \mathbf{W}^{(s)}$$

# Location-based Attention (Luong et al., 2015)

$$\text{score}(\mathbf{h}_t^d, \mathbf{h}_k^e) = \text{softmax}(\mathbf{h}_t^d \mathbf{W}^{(s)})$$

# Scaled Dot-Product Attention

(Vaswani et al., 2017)

$$\text{score}(\mathbf{h}_t^d, \mathbf{h}_k^e) = \frac{\mathbf{h}_k^e \text{trans}(\mathbf{h}_t^d)}{\sqrt{n}}$$

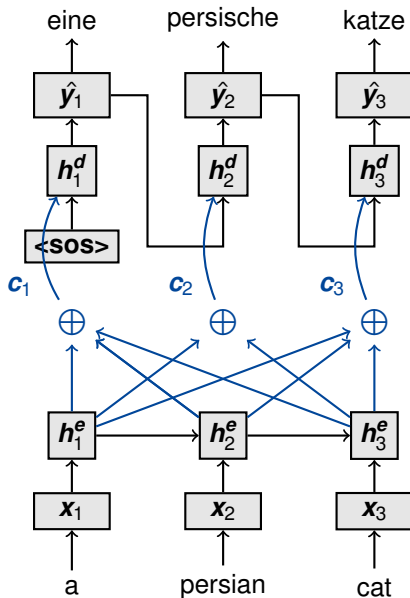
- ▶ The scaling factor  $\frac{1}{\sqrt{n}}$  is motivated by the concern when the input is large, the softmax function may have an extremely small gradient.
- ▶ Small gradients yields difficulties in learning.

# Self-Attention

- ▶ An attention mechanism to relate different tokens of an input sequence to compute a representation of the sequence itself.
- ▶ For example, the self-attention mechanism enables a model to learn the relations between a word of an input sentence and its previous words.

The FBI is chasing a criminal on the run .  
 The FBI is chasing a criminal on the run .  
 The FBI is chasing a criminal on the run .  
 The FBI is chasing a criminal on the run .  
 The FBI is chasing a criminal on the run .  
 The FBI is chasing a criminal on the run .  
 The FBI is chasing a criminal on the run .  
 The FBI is chasing a criminal on the run .  
 The FBI is chasing a criminal on the run .

# Neural Machine Translation



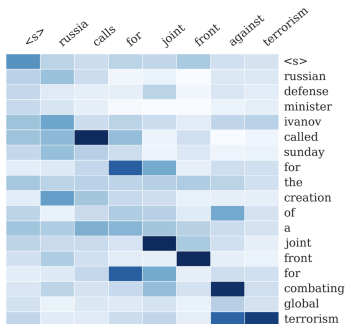
# Summarization (Rush et al. 2015)

Input ( $\mathbf{x}_1, \dots, \mathbf{x}_{18}$ ). First sentence of article:

russian defense minister ivanov called sunday for the creation of a joint front for combating global terrorism

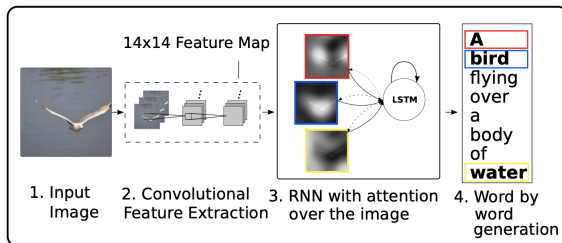
Output ( $\mathbf{y}_1, \dots, \mathbf{y}_8$ ). Generated headline:

*russia calls for joint front against **terrorism***  $\Leftarrow g(\text{terrorism}, \mathbf{x}, \text{for}, \text{joint}, \text{front}, \text{against})$





# Image Caption Generation (Xu et al. 2015)



# Other Applications

- ▶ lemmatization: `gespielt`  $\rightarrow$  `spielen`
- ▶ Spelling correction: `i_lvoe_u`  $\rightarrow$  `i_love_you`

# Auto Encoder

- ▶ An encoder-decoder model that transforms an input sequence to itself.
- ▶ It learns the identity function  $F(x) = x$ .
- ▶ It usually add some noise to the input, then the model learns to remove the noise.
- ▶ It is used for dimensionality reduction, representation learning, and unsupervised learning.
- ▶ The encoder and decoder can be used individually to solve other tasks.

# Summary

- ▶ Encoders and Decoders
- ▶ Attention mechanism
- ▶ Their applications in NLP

Thank You!