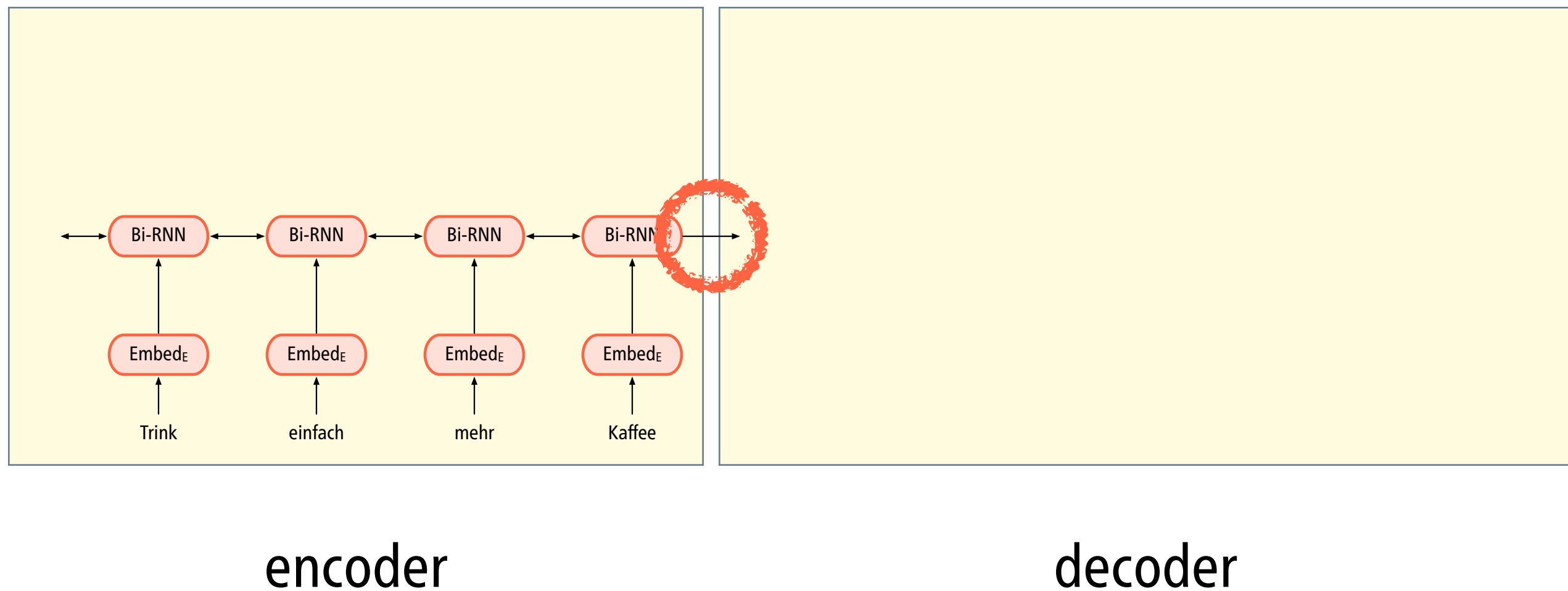


# Attention

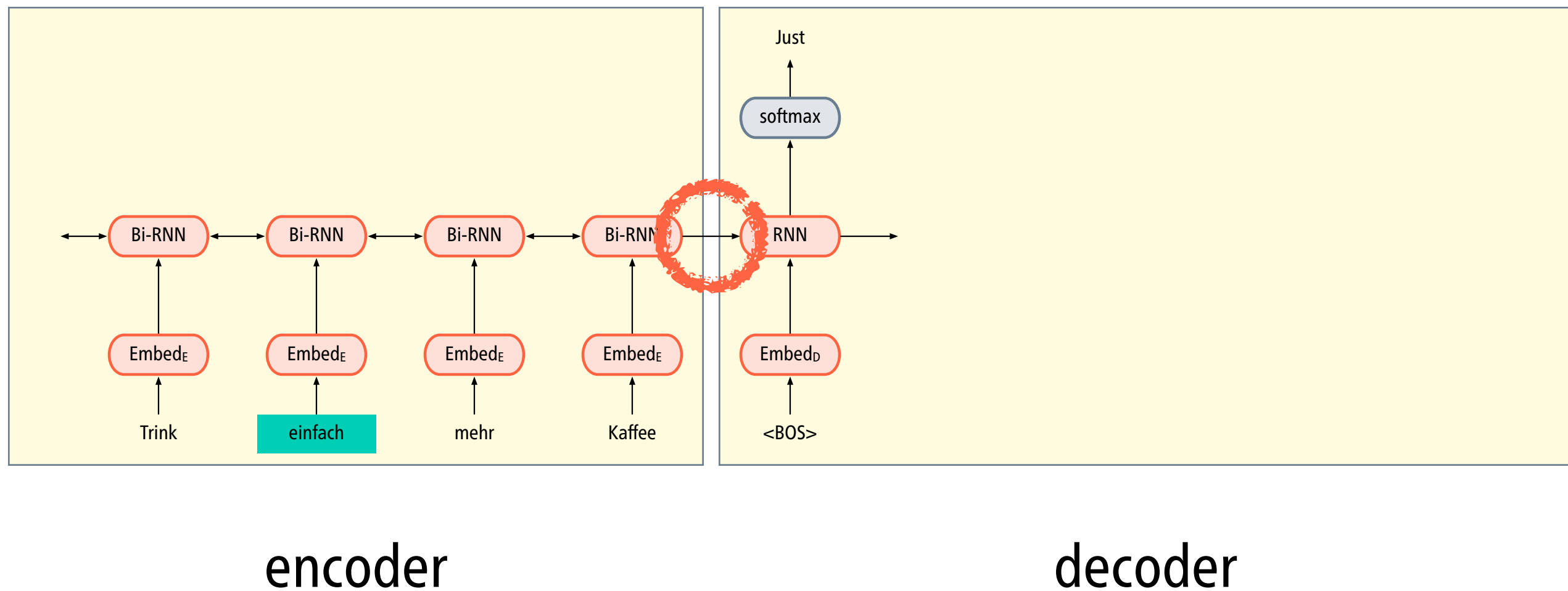
Marco Kuhlmann

Department of Computer and Information Science

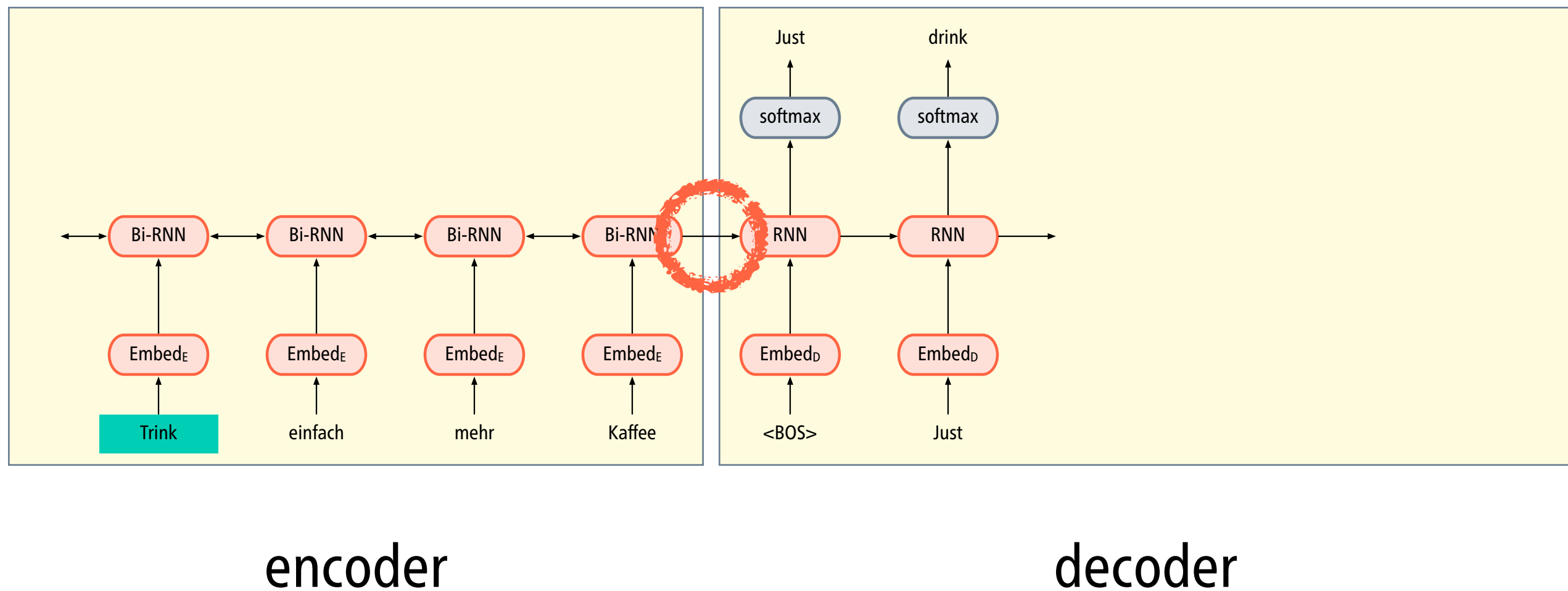
# Recency bias in recurrent neural networks



# Recency bias in recurrent neural networks



# Recency bias in recurrent neural networks



[Sutskever et al. \(2014\)](#)

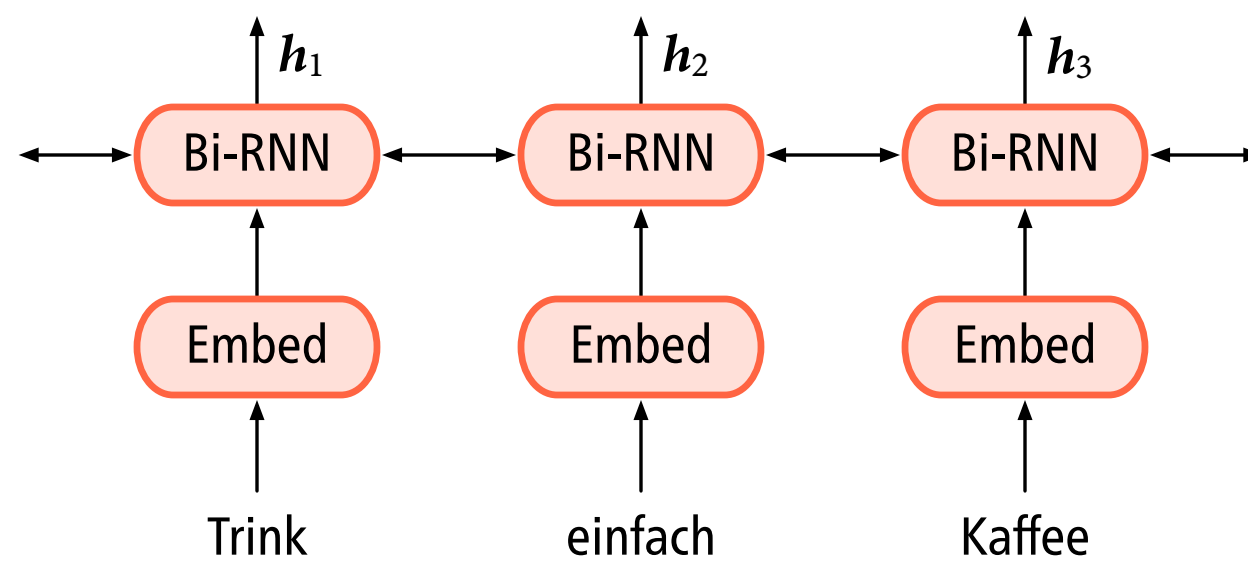
# Attention

- In the context of machine translation, **attention** enables the model to learn ‘soft’ word alignments.
- Essentially, we compute a set of weights that allow us to score words based on how much the model should ‘attend to them’.
- Attention was first proposed in the context of the sequence-to-sequence architecture, but is now used in many architectures.

[Bahdanau et al. \(2015\)](#)

# Attention for translation

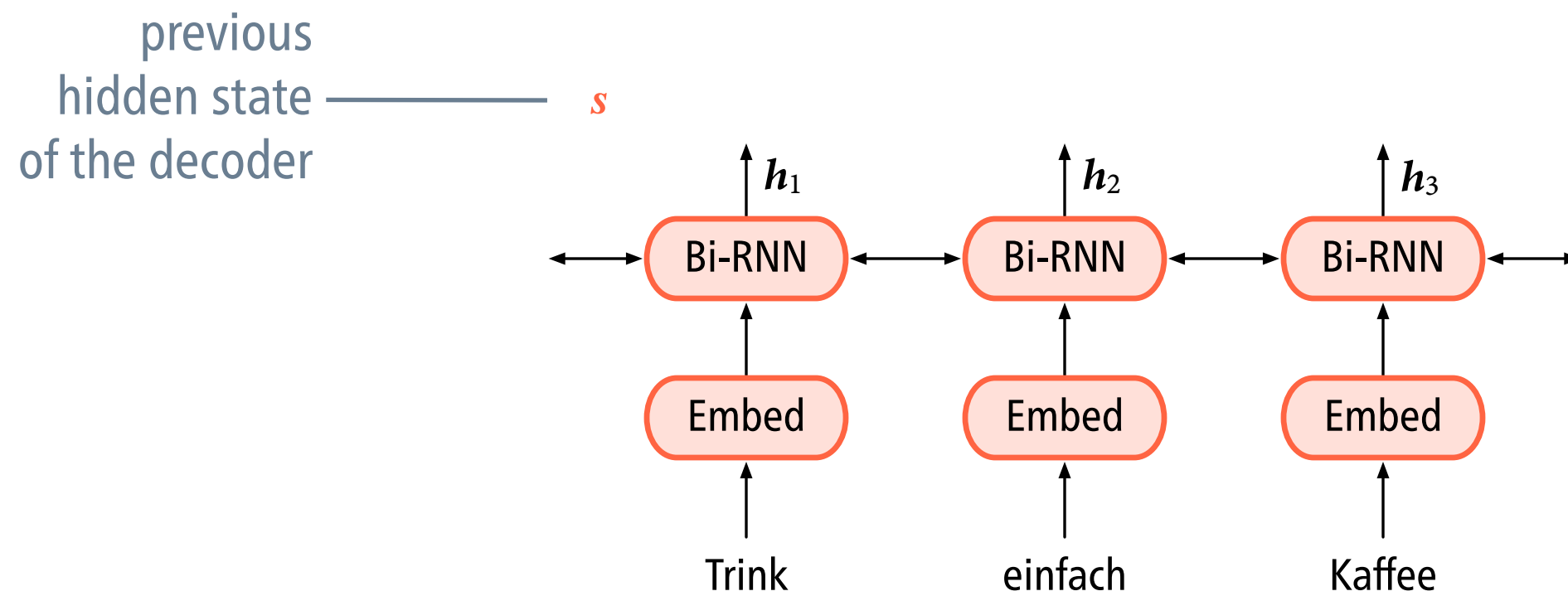
Just      drink      coffee



[Bahdanau et al. \(2015\)](#)

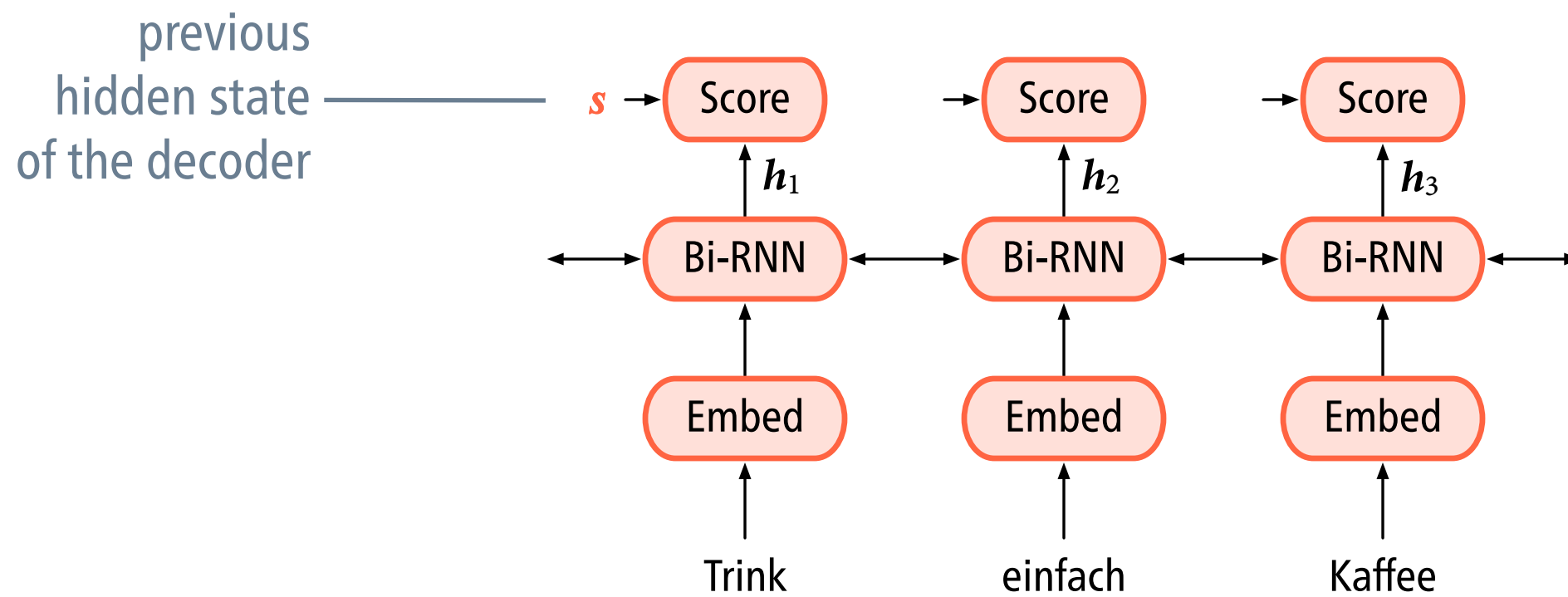
# Attention for translation

Just      drink      coffee



# Attention for translation

Just      drink      coffee









# A general characterisation of attention

- In general, attention can be described as a mapping from a query  $\mathbf{q}$  and a set of key–value pairs  $\mathbf{k}_i, \mathbf{v}_i$  to an output.
- The output is the weighted sum of the  $\mathbf{v}_i$ , where the weight of each  $\mathbf{v}_i$  is given by the affinity between  $\mathbf{q}$  and  $\mathbf{k}_i$ :

$$\text{attention}(\mathbf{q}, \mathbf{K}, \mathbf{V}) = \text{softmax}(a(\mathbf{q}, \mathbf{K}))\mathbf{V}$$

$$\mathbf{q} \in \mathbb{R}^{d_K}, \mathbf{K} \in \mathbb{R}^{n \times d_K}, \mathbf{V} \in \mathbb{R}^{n \times d_V}$$

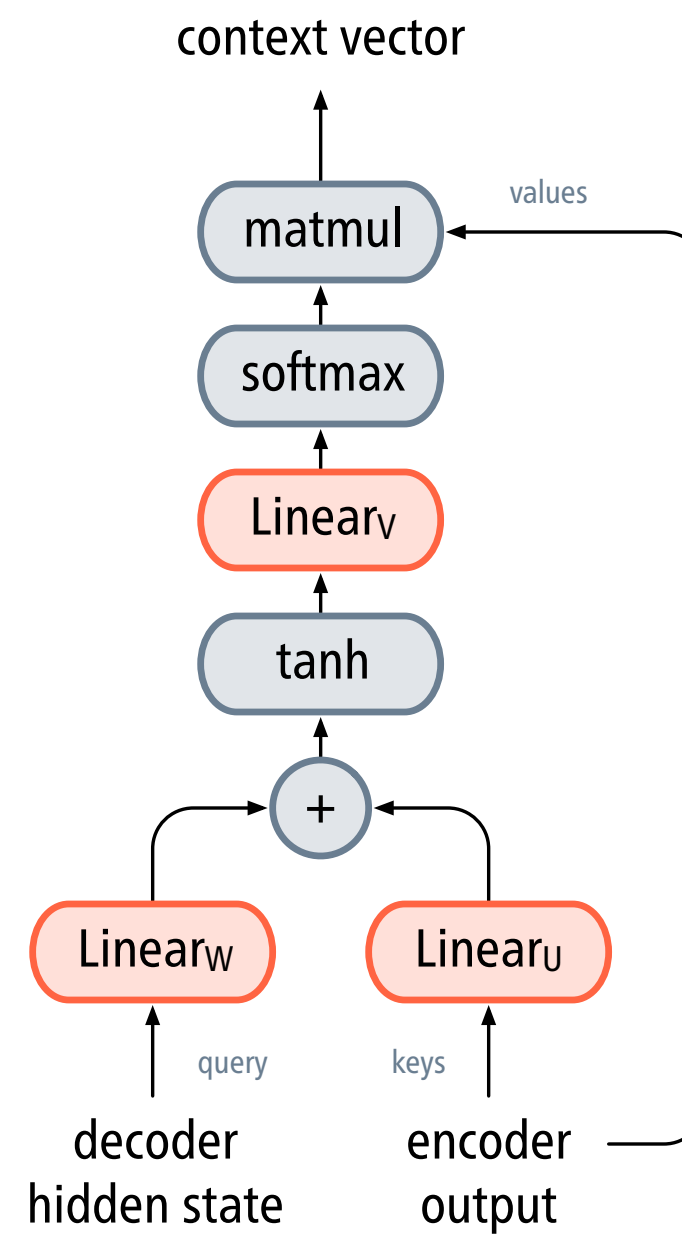
attention score

# Bahdanau attention

previous decoder  
hidden state

$$a(\mathbf{s}_{i-1}, \mathbf{h}_j) = \mathbf{v}^\top \tanh(\mathbf{W}\mathbf{s}_{i-1} + \mathbf{U}\mathbf{h}_j)$$

encoder hidden  
state at position j

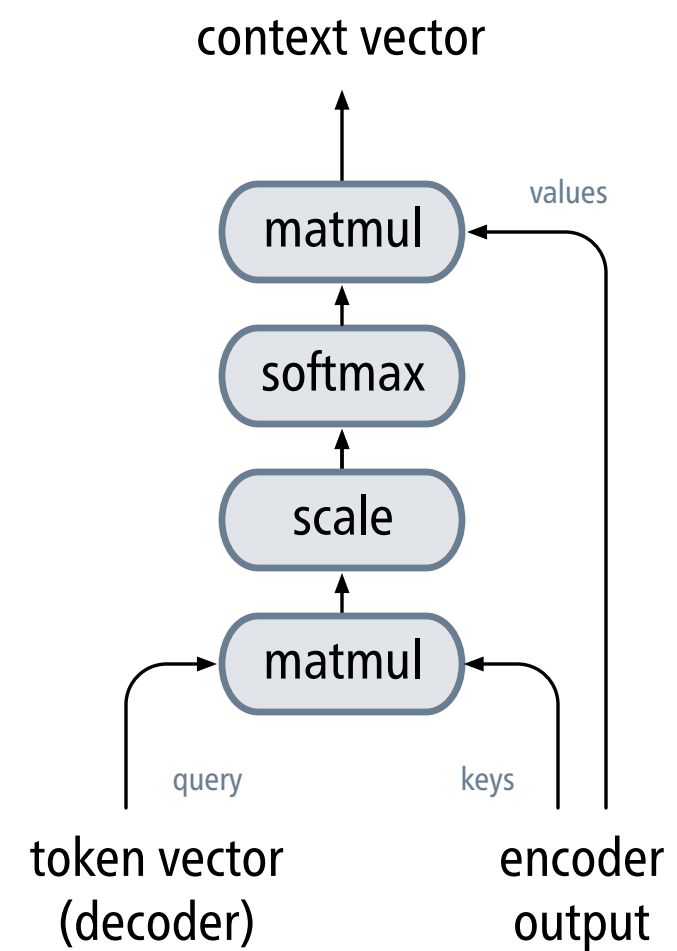


# Scaled dot-product attention

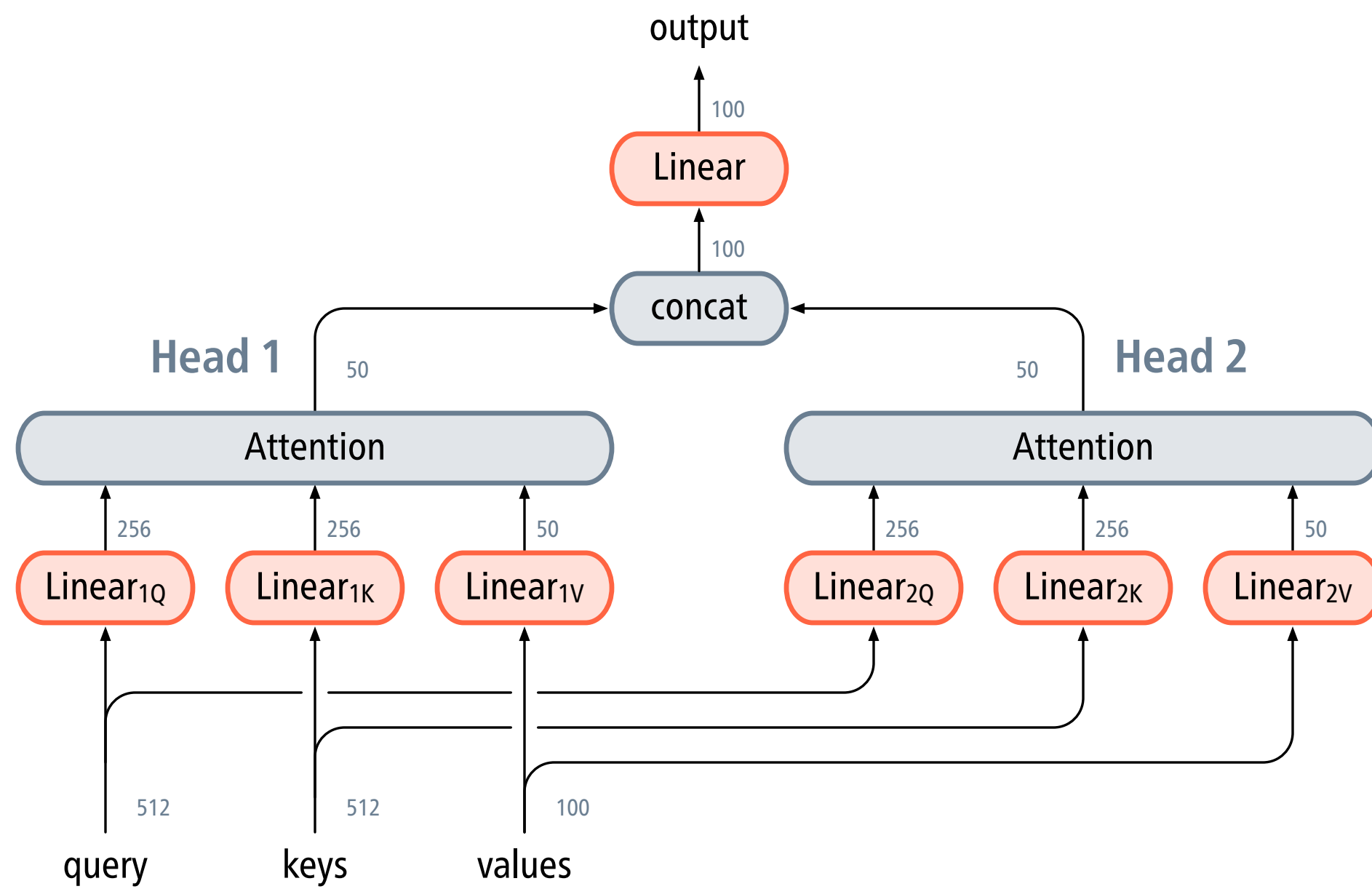
token vector  
at position i

$$a(\mathbf{h}_i, \mathbf{h}_j) = \frac{\mathbf{h}_i \mathbf{h}_j^T}{\sqrt{d_K}}$$

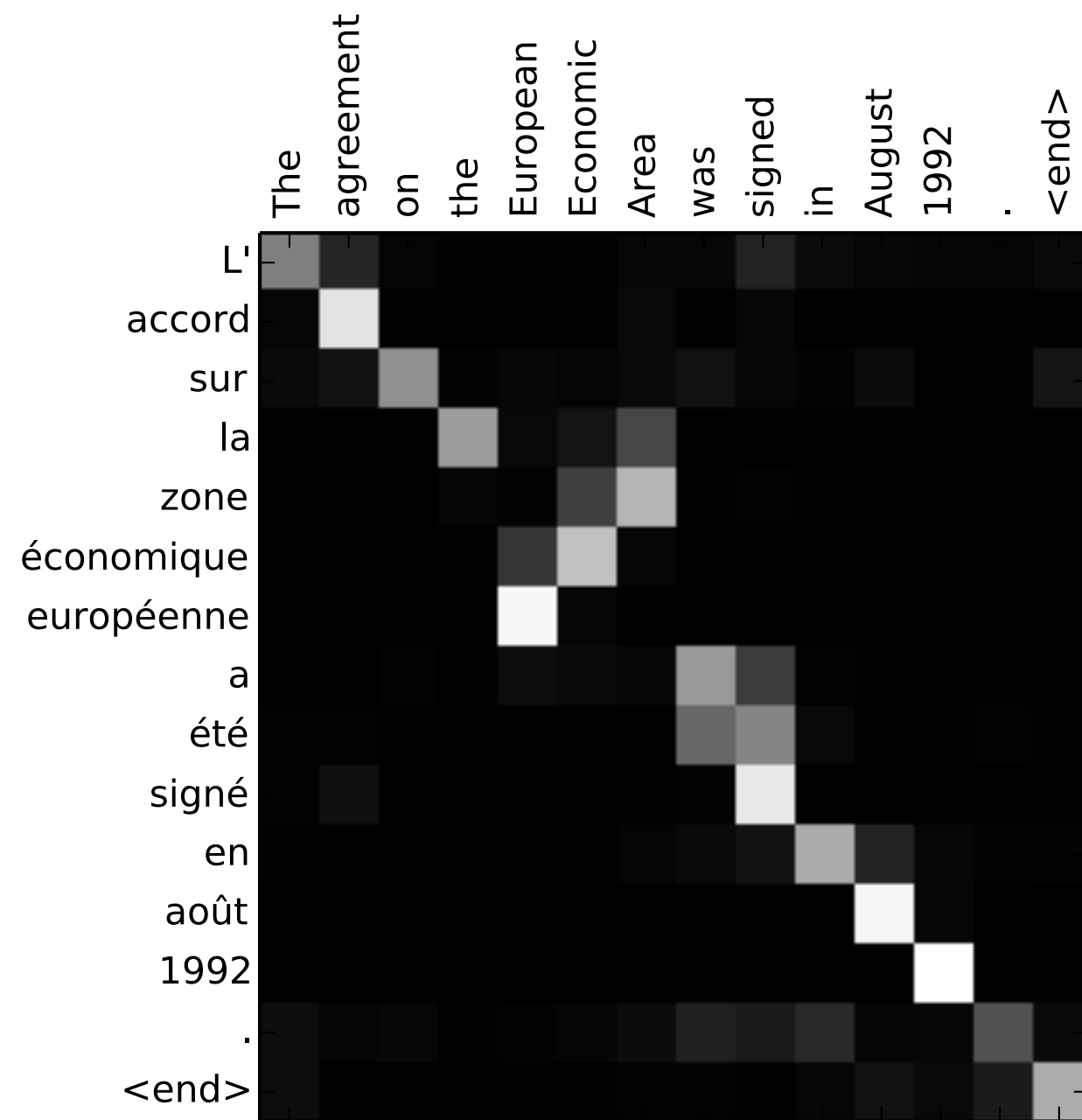
encoder output  
at position j



# Multi-head attention



# Attention as word alignments



In the context of the encoder–decoder architecture for neural machine translation, attention weights resemble soft word alignments.

Image source: [Bahdanau et al. \(2015\)](#)