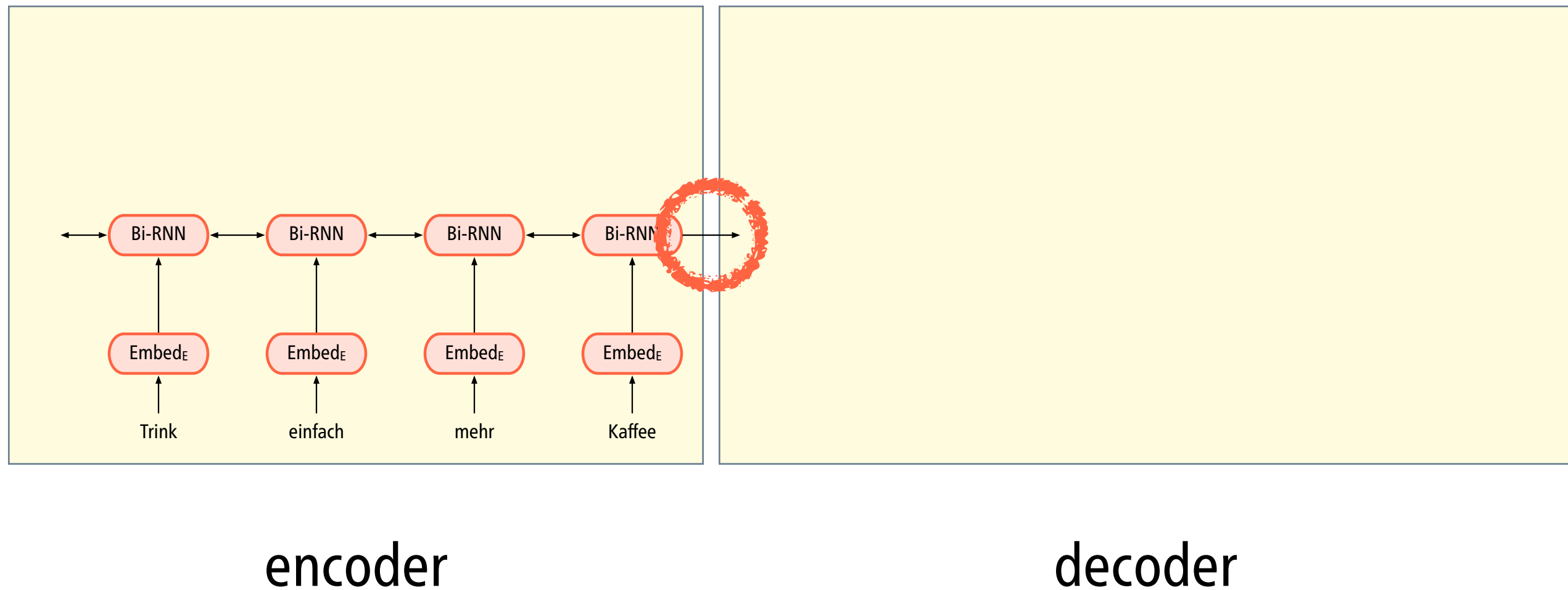Natural Language Processing

# Attention
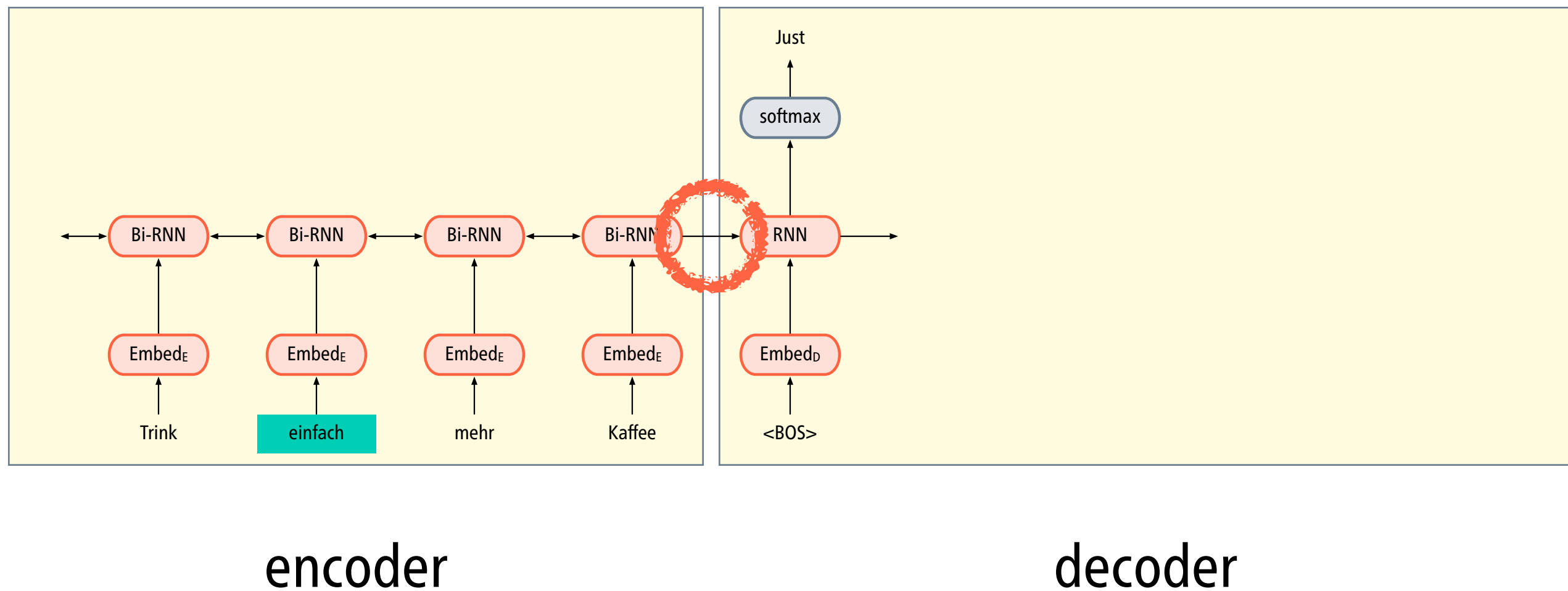
Marco Kuhlmann

Department of Computer and Information Science
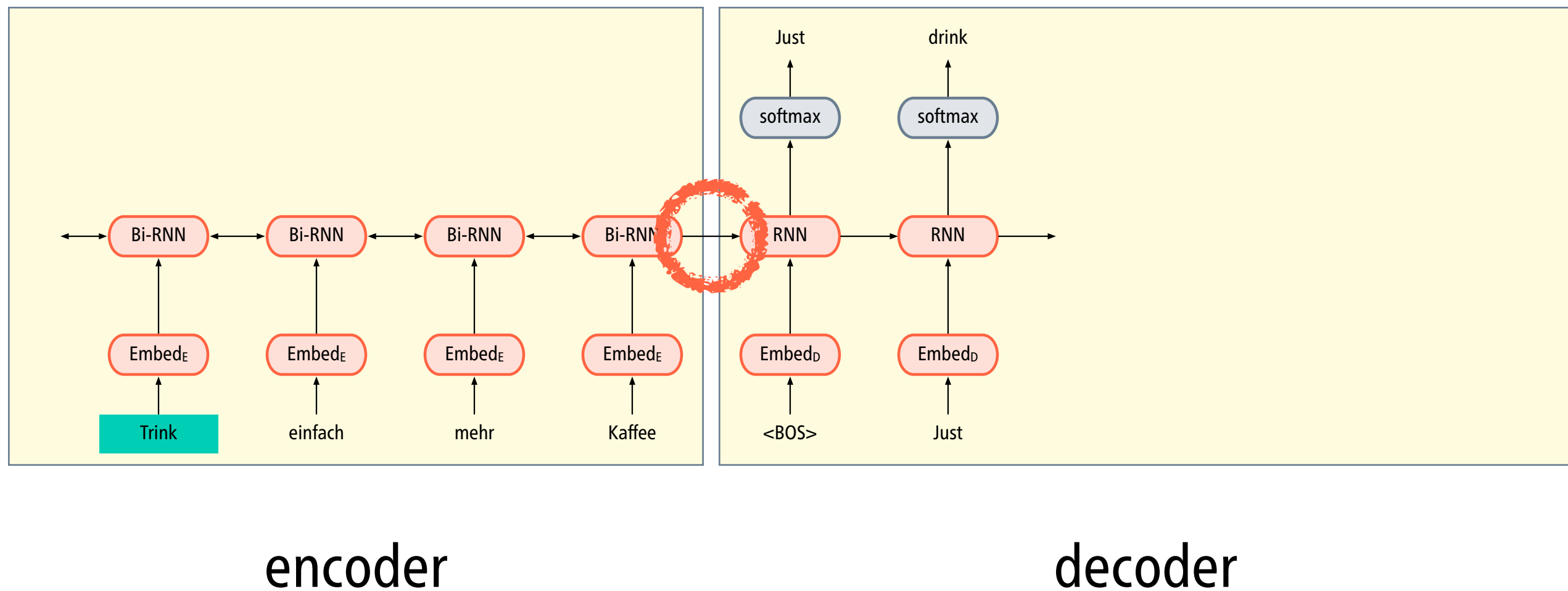
# Recency bias in recurrent neural networks



encoder                                    decoder

Sutskever et al. (2014)

# Recency bias in recurrent neural networks



encoder                    decoder

Sutskever et al. (2014)

# Recency bias in recurrent neural networks



encoder                    decoder
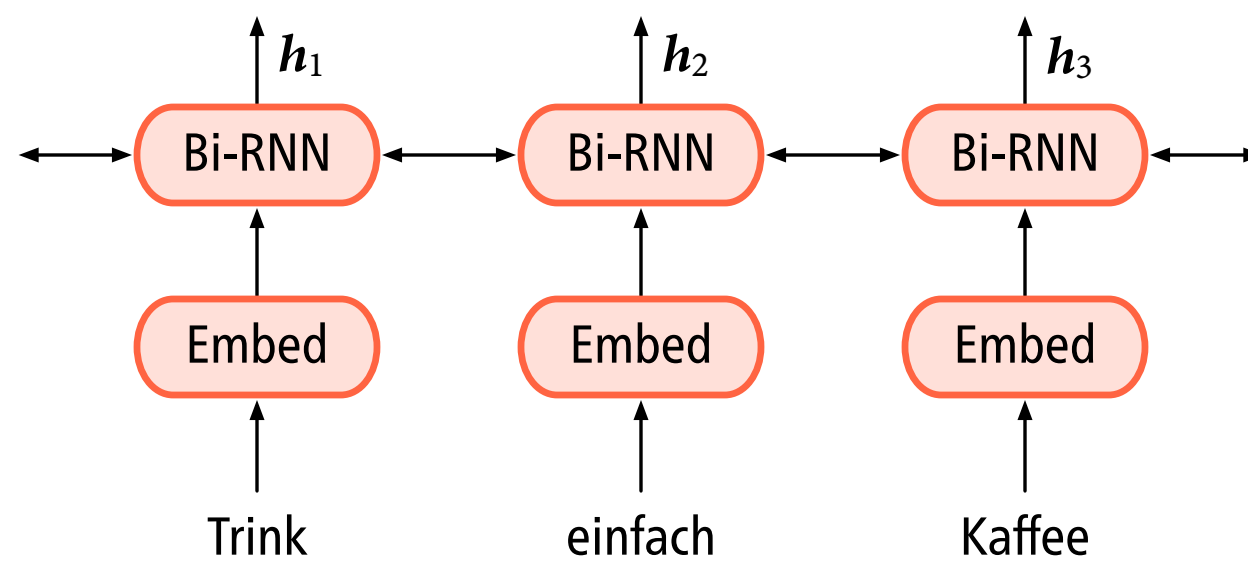
Sutskever et al. (2014)

# 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
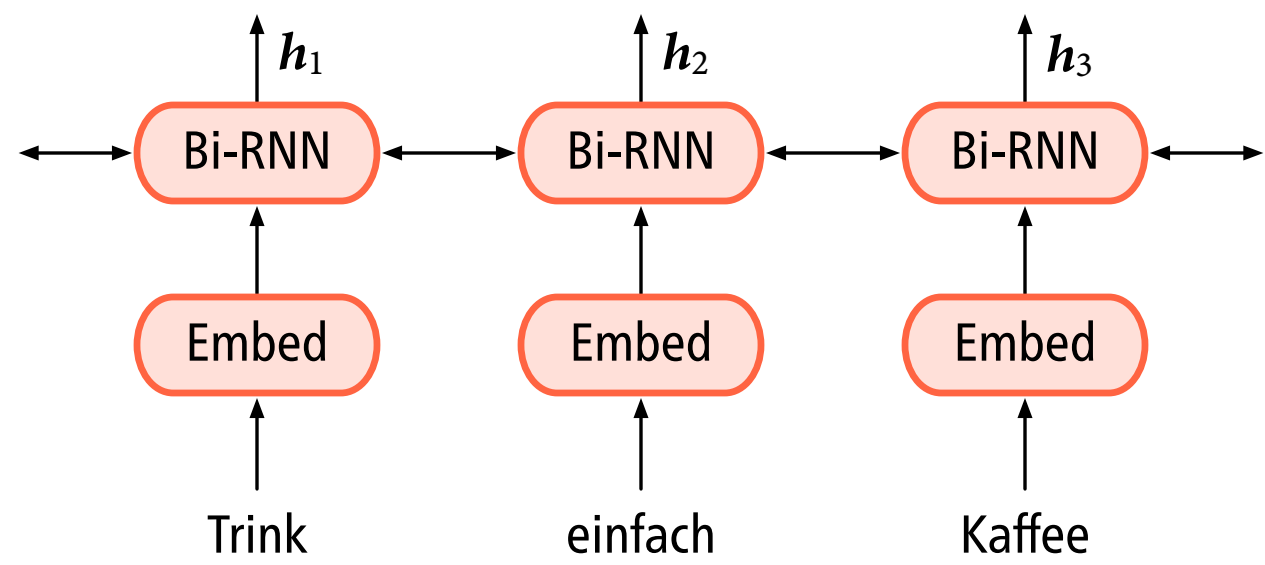


Bahdanau et al. (2015)

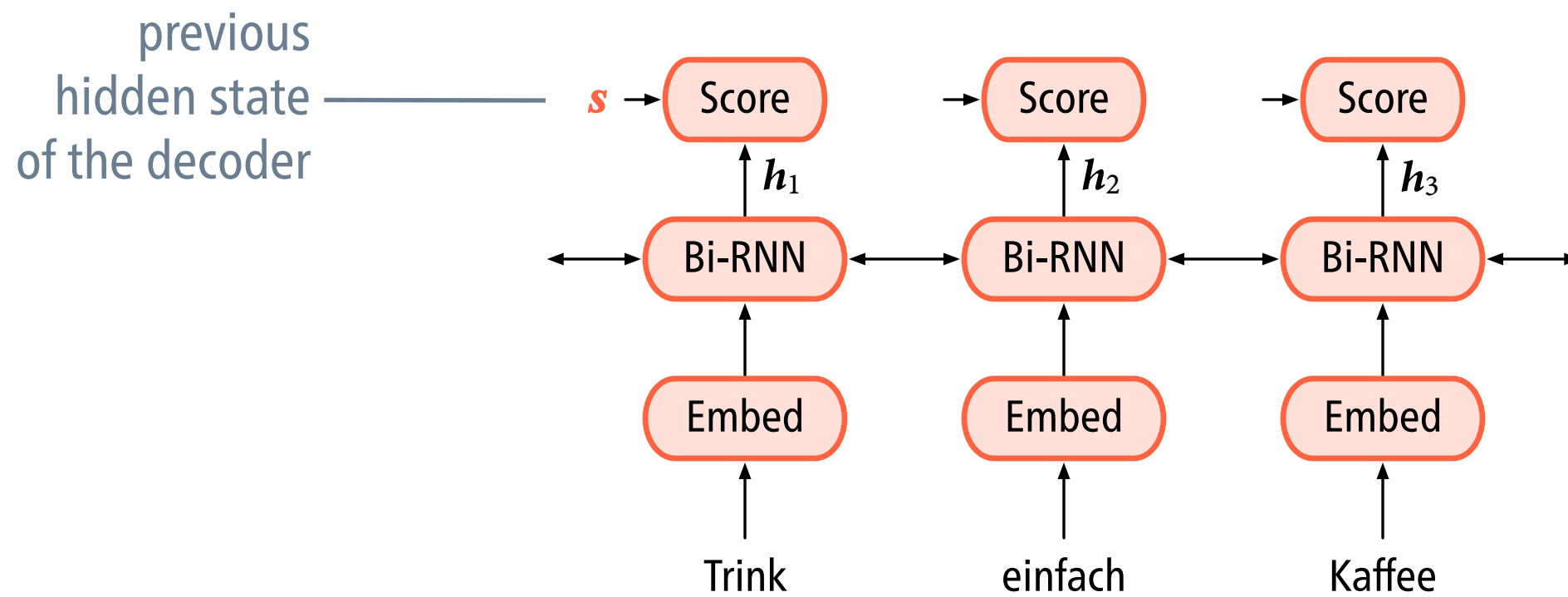# Attention for translation

Just        drink        coffee
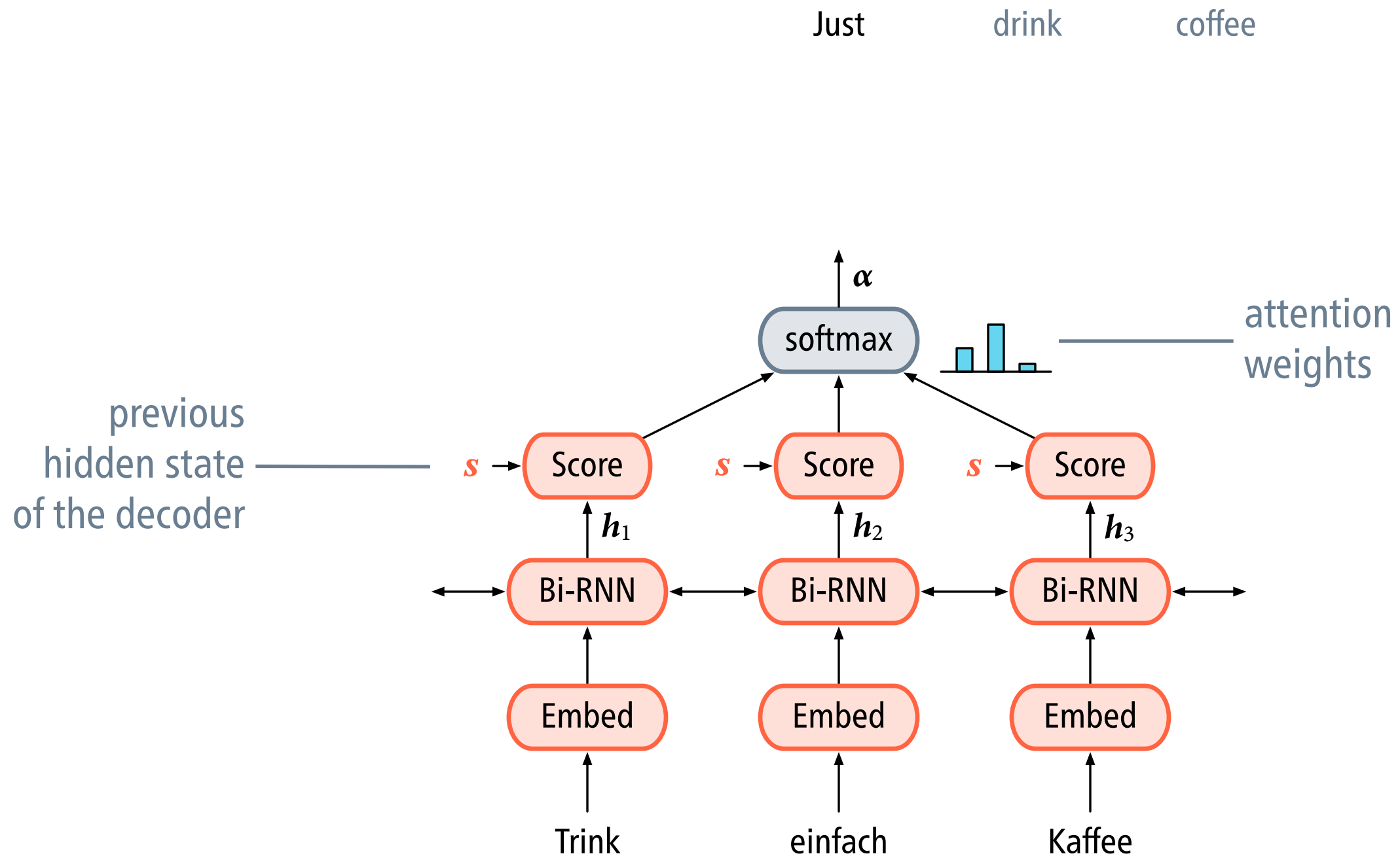
previous
hidden state —————— $s$
of the decoder

$h_1$        $h_2$        $h_3$

←  Bi-RNN  →  Bi-RNN  →  Bi-RNN  →

Embed      Embed      Embed

Trink        einfach      Kaffee

Bahdanau et al. (2015)

# Attention for translation

Just     drink     coffee

previous
hidden state
of the decoder

$s$

| Score | | Score | | Score |

Bi-RNN ↔ Bi-RNN ↔ Bi-RNN

$h_1$     $h_2$     $h_3$

Embed    Embed    Embed

Trink     einfach     Kaffee

Bahdanau et al. (2015)

# Attention for translation

Just      drink      coffee



Bahdanau et al. (2015)

# Attention for translation



Bahdanau et al. (2015)

# A general characterisation of attention

- In general, attention can be described as a mapping from a query $q$ and a set of key–value pairs $k_i$, $v_i$ to an output.

- The output is the weighted sum of the $v_i$, where the weight of each $v_i$ is given by the affinity between $q$ and $k_i$:

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

$$q \in \mathbb{R}^{d_K}, K \in \mathbb{R}^{n \times d_K}, V \in \mathbb{R}^{n \times d_V} \qquad \text{attention score}$$

Vaswani et al. (2017)

# Bahdanau attention

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

previous decoder
hidden state

encoder hidden
state at position j

context vector

values

matmul

softmax

Linear$_V$

tanh

+

Linear$_W$

Linear$_U$

query

keys

decoder
hidden state

encoder
output

# Scaled dot-product attention

context vector

values

matmul

softmax

scale

matmul

query · keys

token vector
(decoder)

encoder
output

token vector
at position i

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

encoder output
at position j

Vaswani et al. (2017)
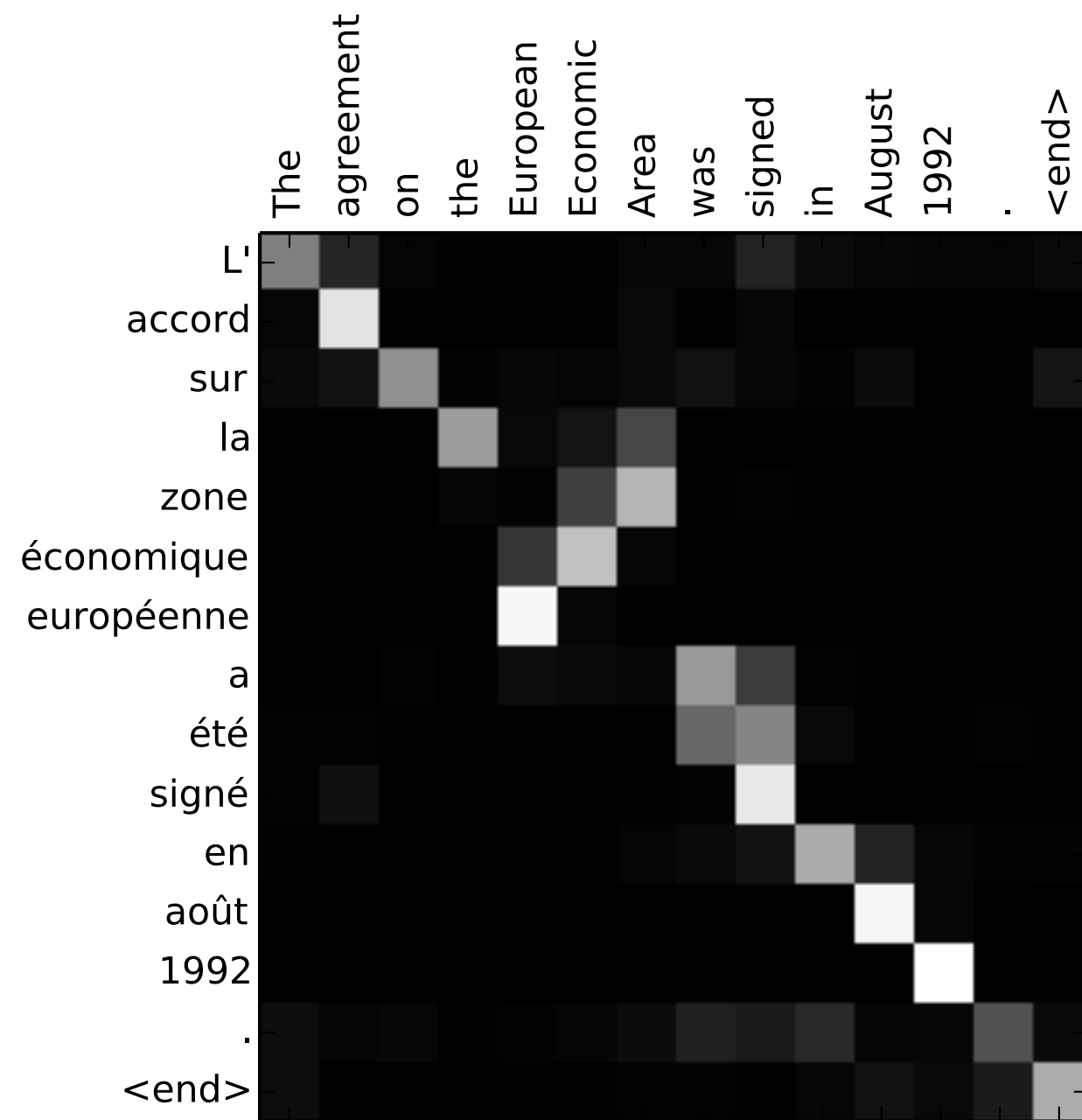
# Multi-head attention

# 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)