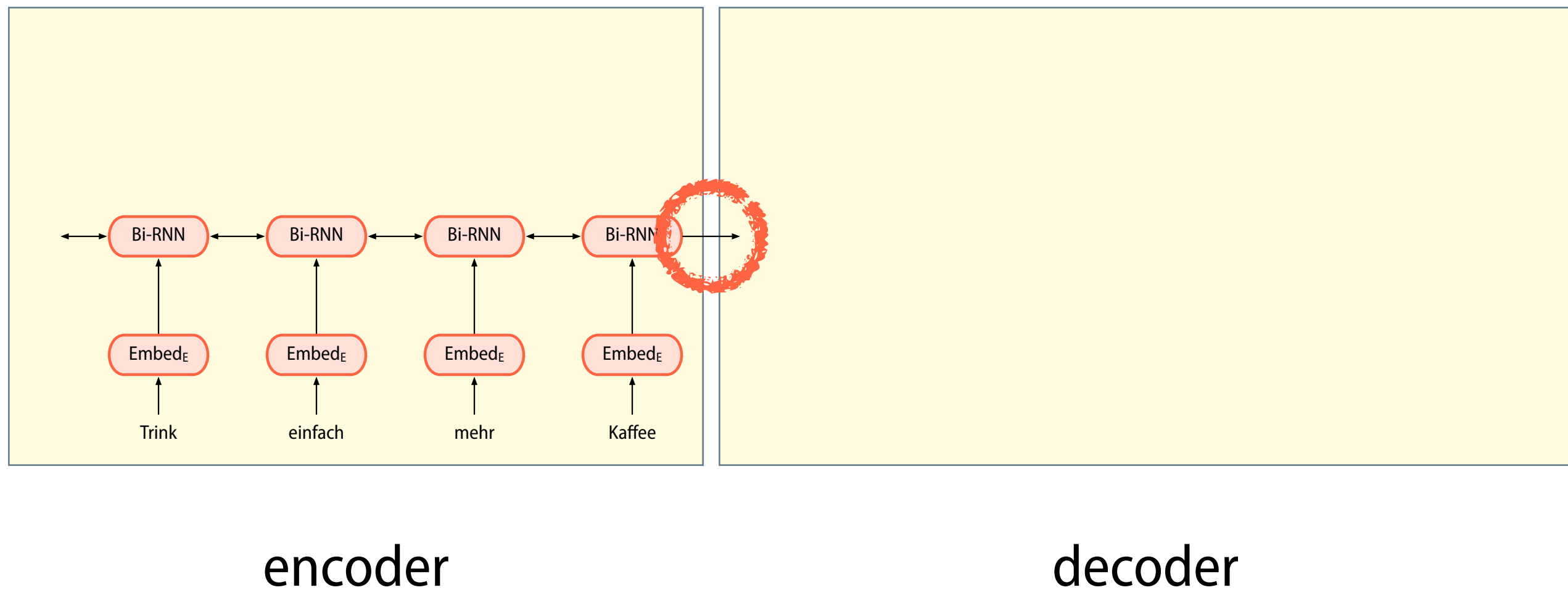


Attention

Marco Kuhlmann

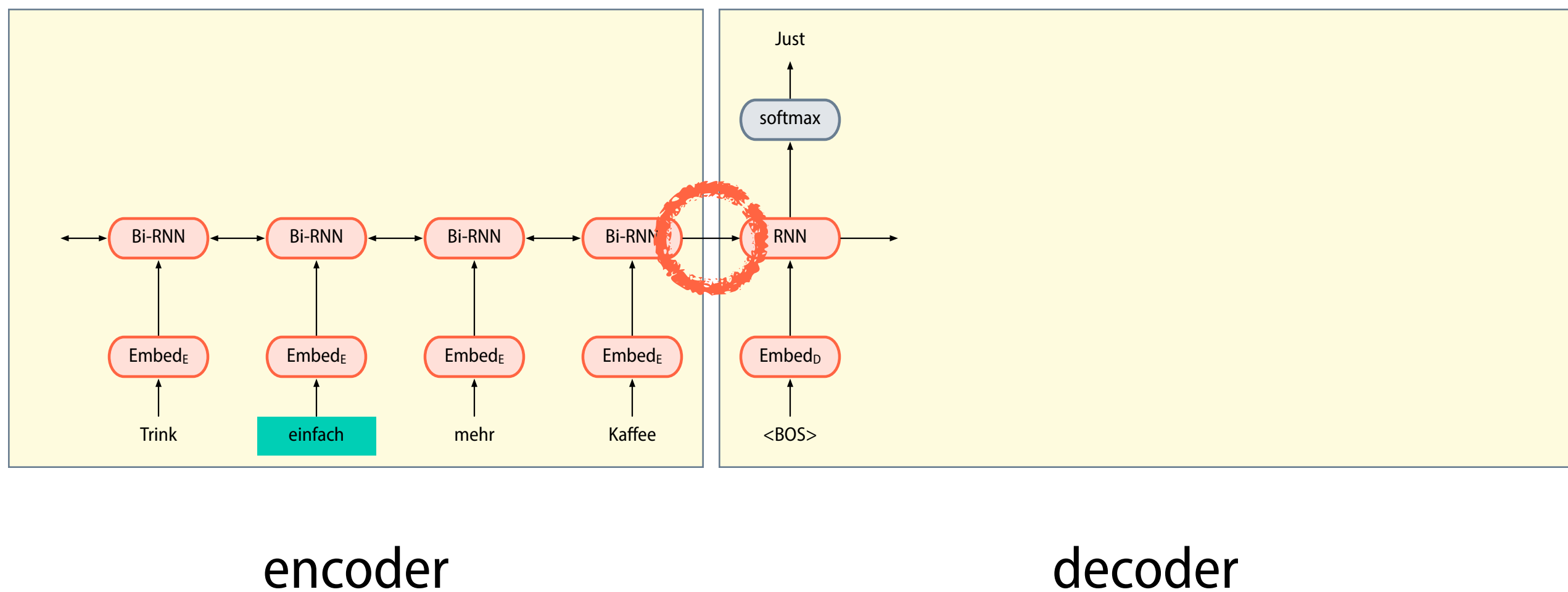
Department of Computer and Information Science

Recency bias in recurrent neural networks



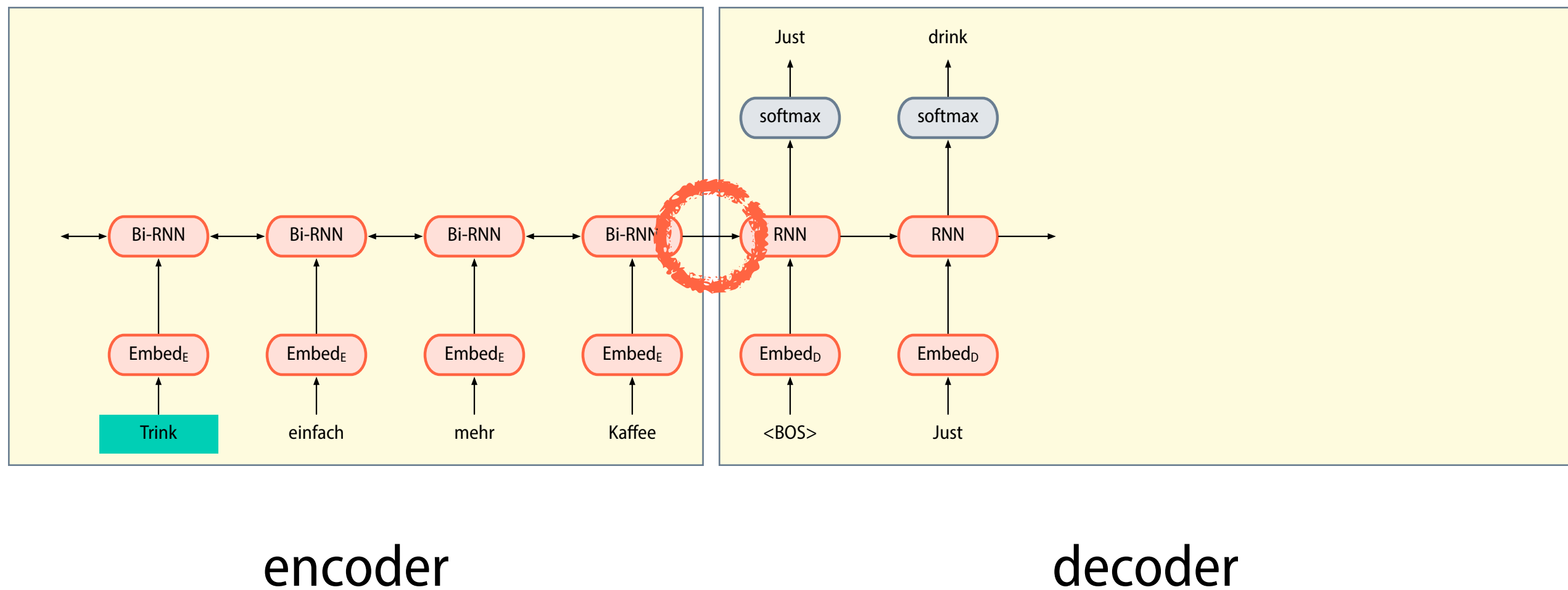
[Sutskever et al. \(2014\)](#)

Recency bias in recurrent neural networks



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

Recency bias in recurrent neural networks



Sutskever et al. (2014)

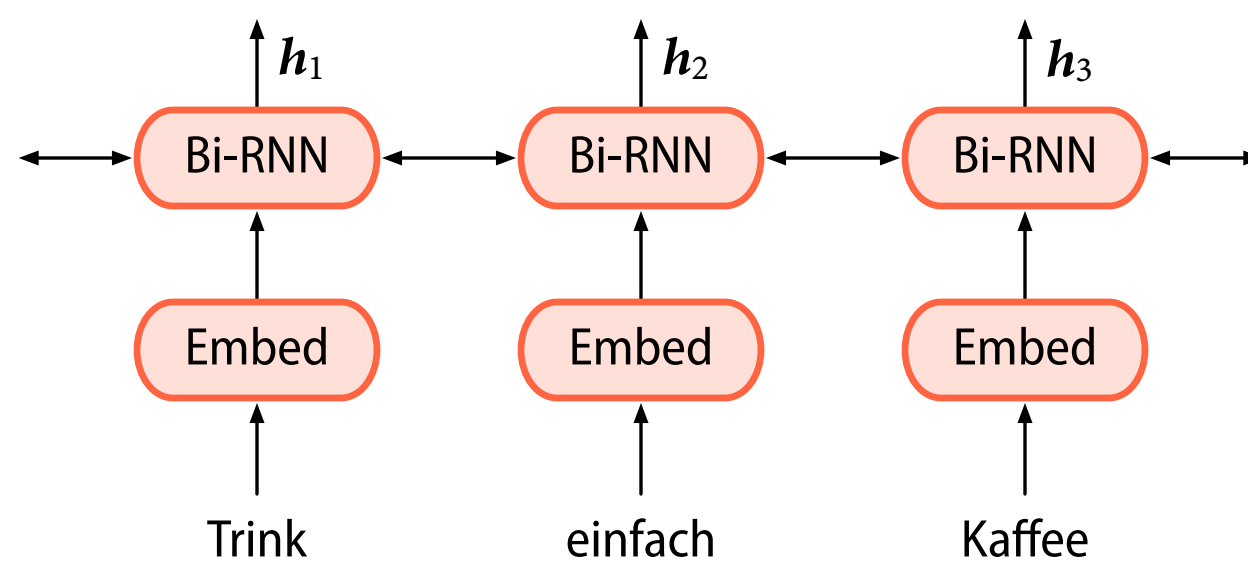
Attention

- 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. \(2015\)](#)

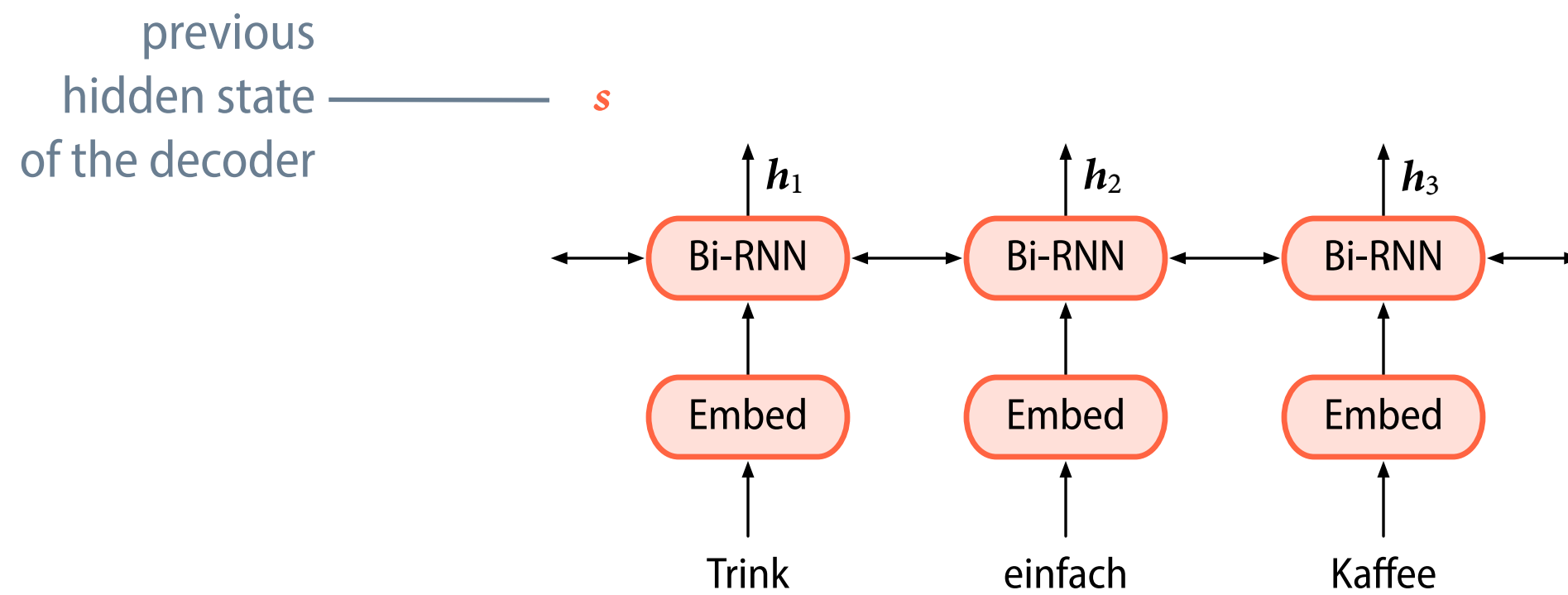
Attention for translation

Just drink coffee



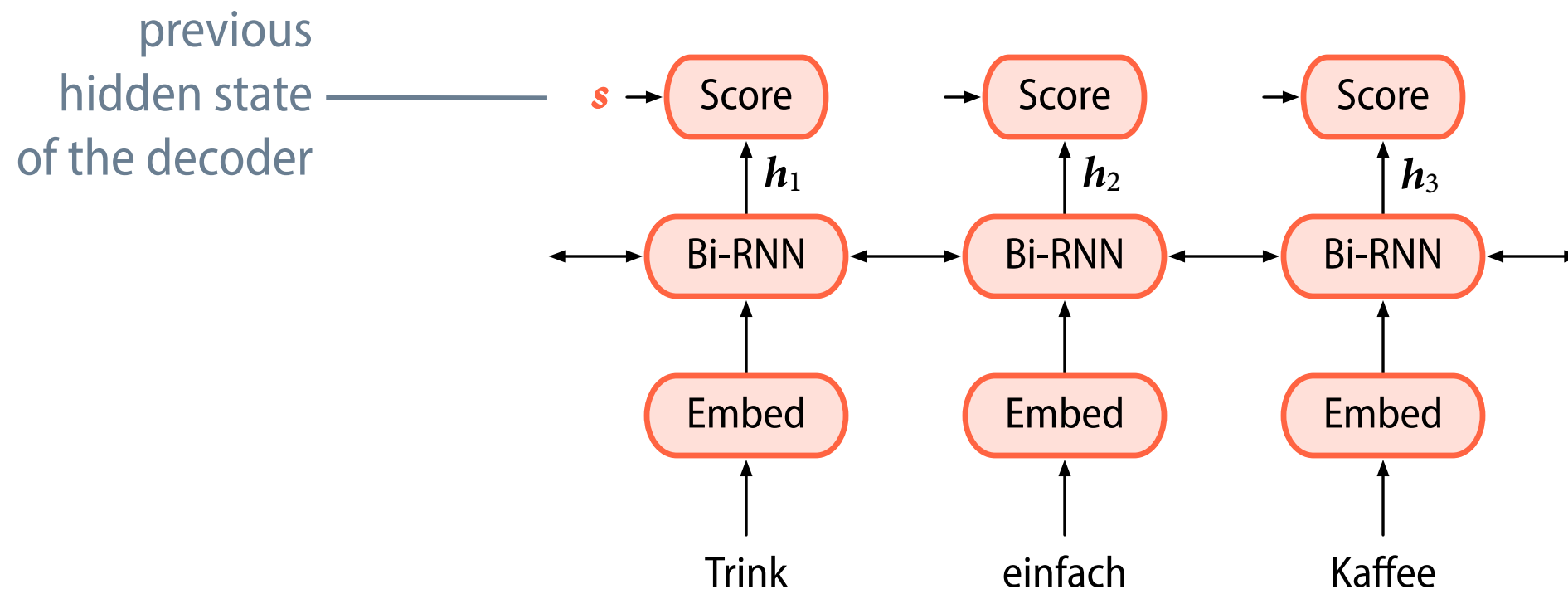
Attention for translation

Just drink coffee

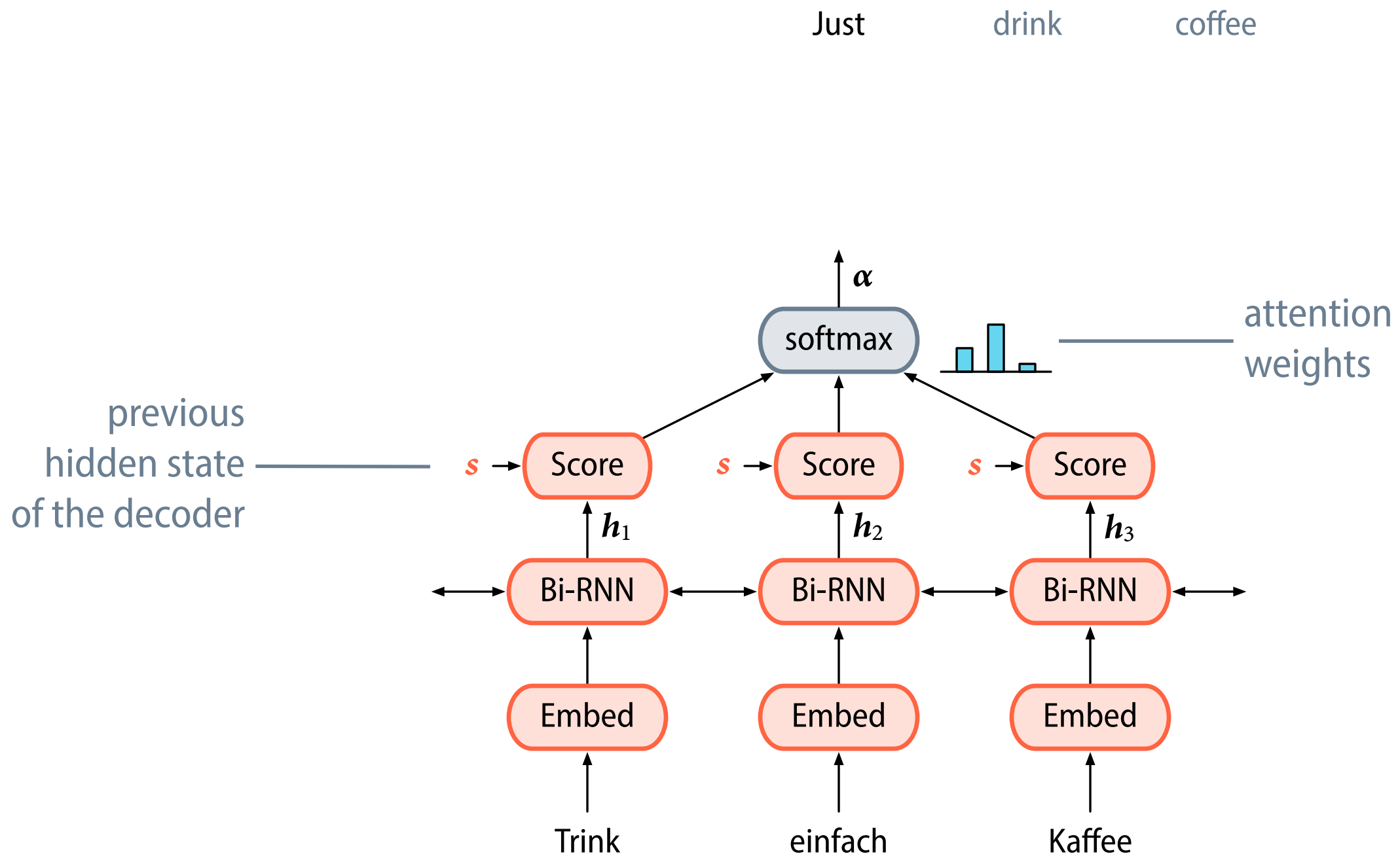


Attention for translation

Just drink coffee



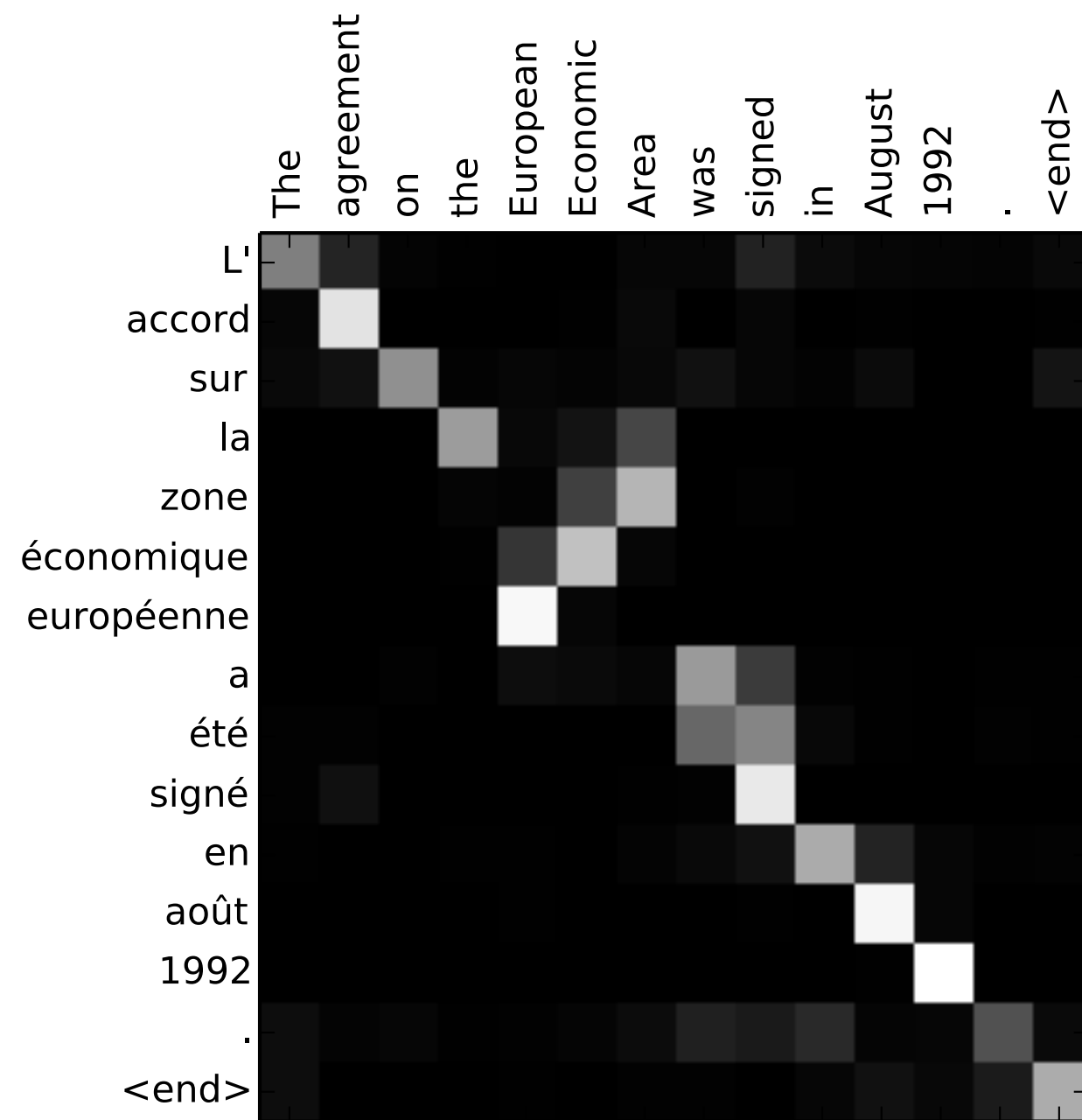
Attention for translation



A more 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 compatibility between \mathbf{q} and \mathbf{k}_i .
- In the translation architecture, the query \mathbf{q} corresponds to the hidden state of the decoder; and the keys and values correspond to the hidden states of the encoder, \mathbf{h}_i .

Attention as word alignments



In the context of the encoder–decoder architecture for neural machine translation, attention can be interpreted as word alignments.

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