

Natural Language Processing

# Attention

Marco Kuhlmann

Department of Computer and Information Science

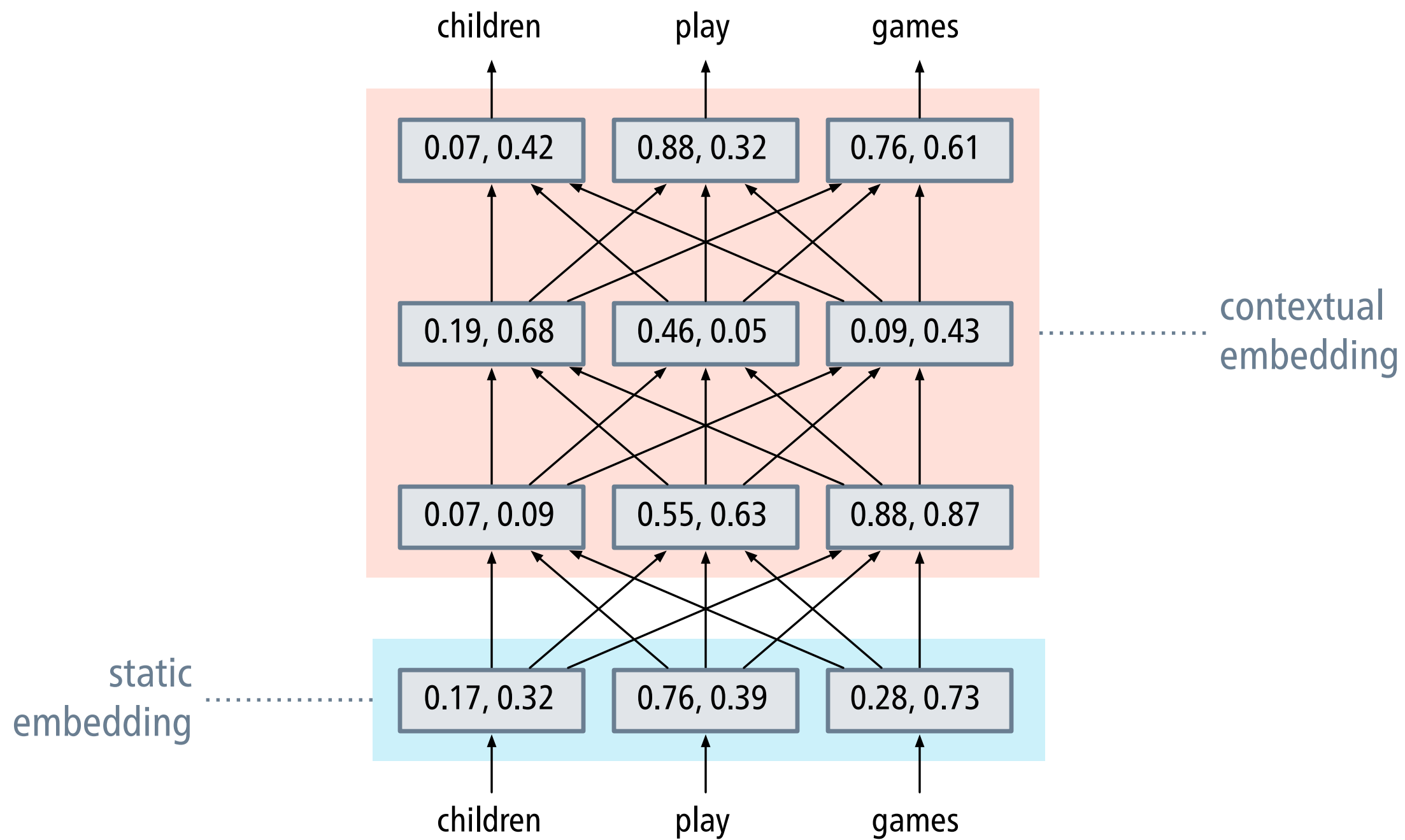
# Contextual embeddings

- Standard word embeddings map each word to a fixed vector.
- In natural language, word meaning is not static, but depends on the surrounding words.

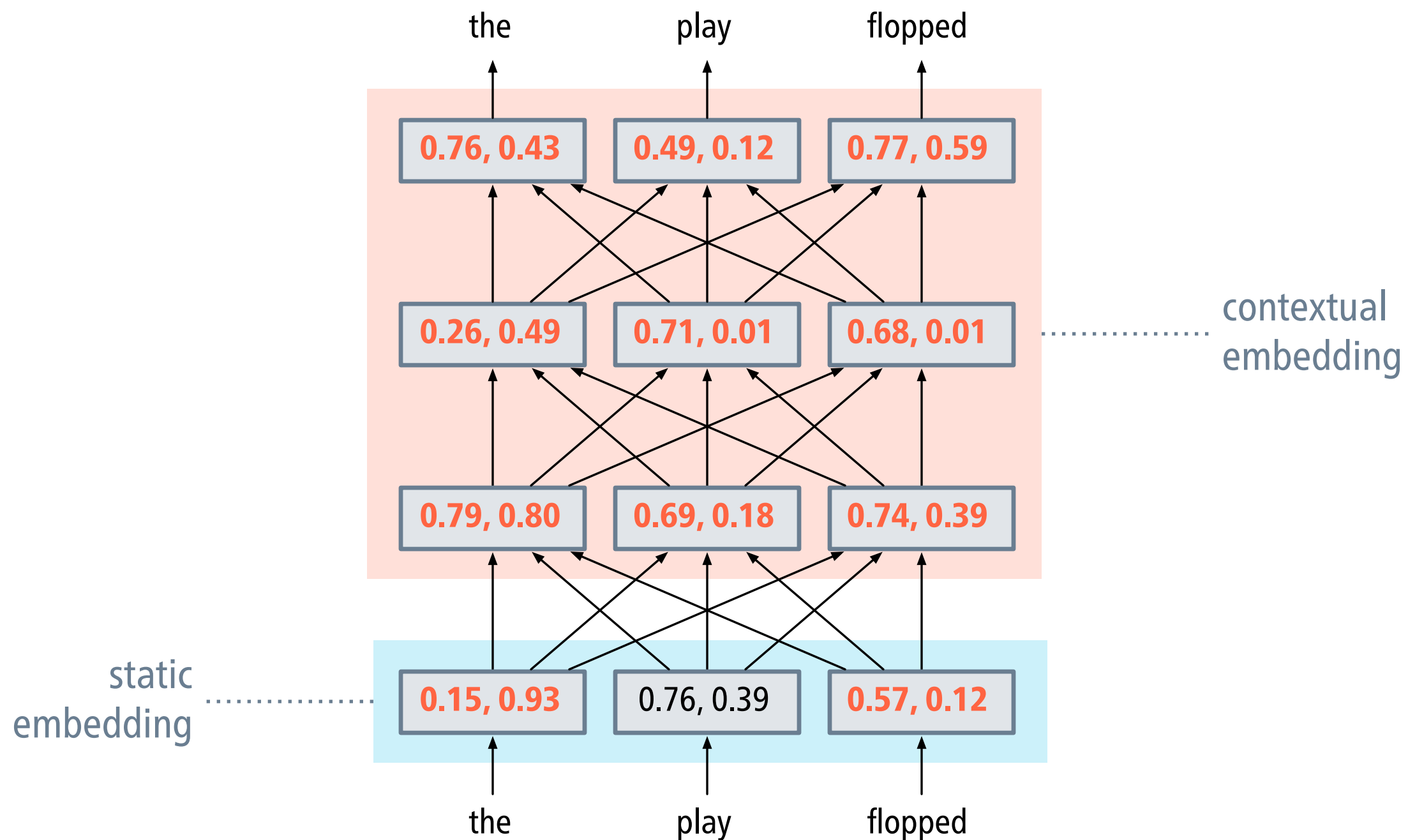
Dogs may *bark* at strangers. Trees shed *bark* in winter.

- **Contextual embeddings** assign word vectors dynamically, conditioned on context.

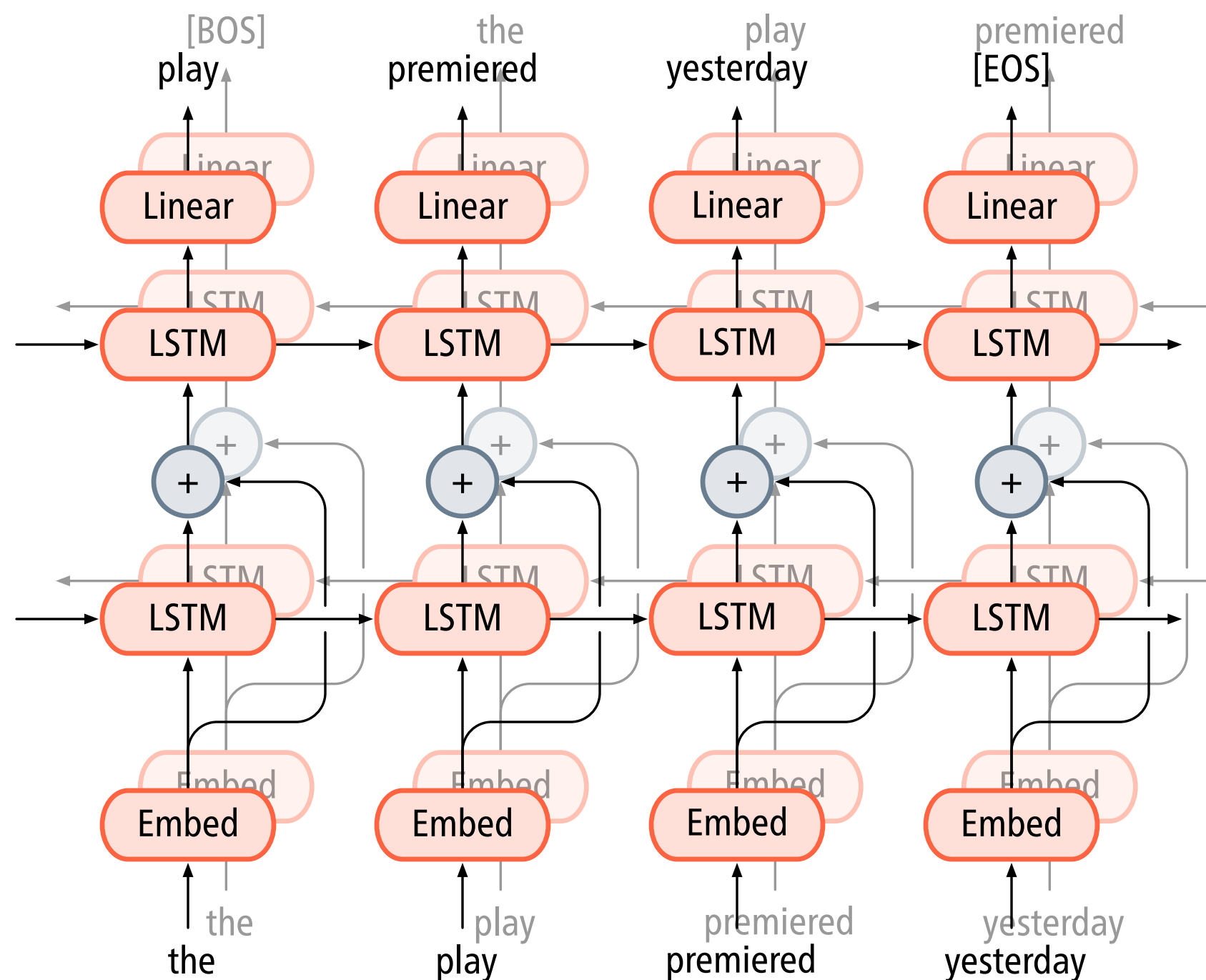
# Contextual embeddings



# Contextual embeddings



# ELMo – Embeddings from Language Models



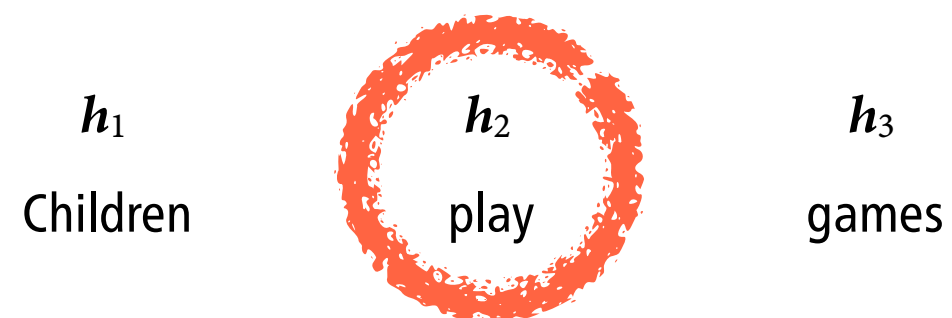
# Contextual embeddings via attention

- In attention, each word vector  $\mathbf{h}_i$  is the weighted sum of all the word vectors  $\mathbf{h}_j$  computed at the previous layer:

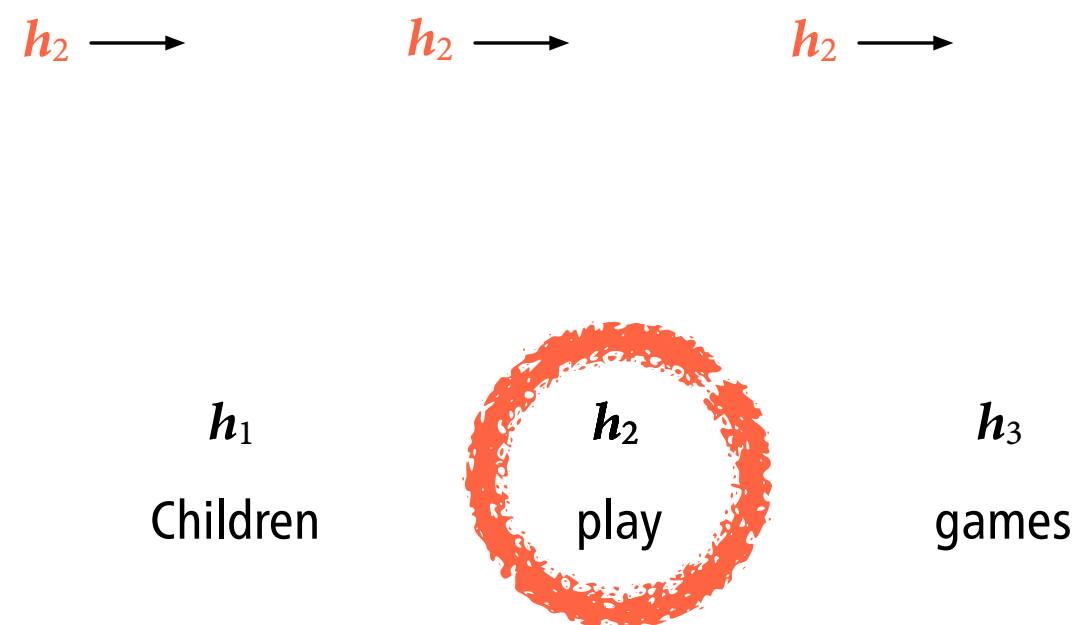
$$\mathbf{h}_i^{\ell+1} = \sum_j \alpha_{ij} \mathbf{h}_j^{\ell}$$

- The weights  $\alpha_{ij}$  express how much the model should “attend to” the context vector  $\mathbf{h}_j$  at layer  $\ell$  when forming  $\mathbf{h}_i$  at layer  $\ell + 1$ .
- The weights are proportional to the vector similarity between  $\mathbf{h}_i$  and the context vectors, and are learned during training.

# Contextual embeddings via attention

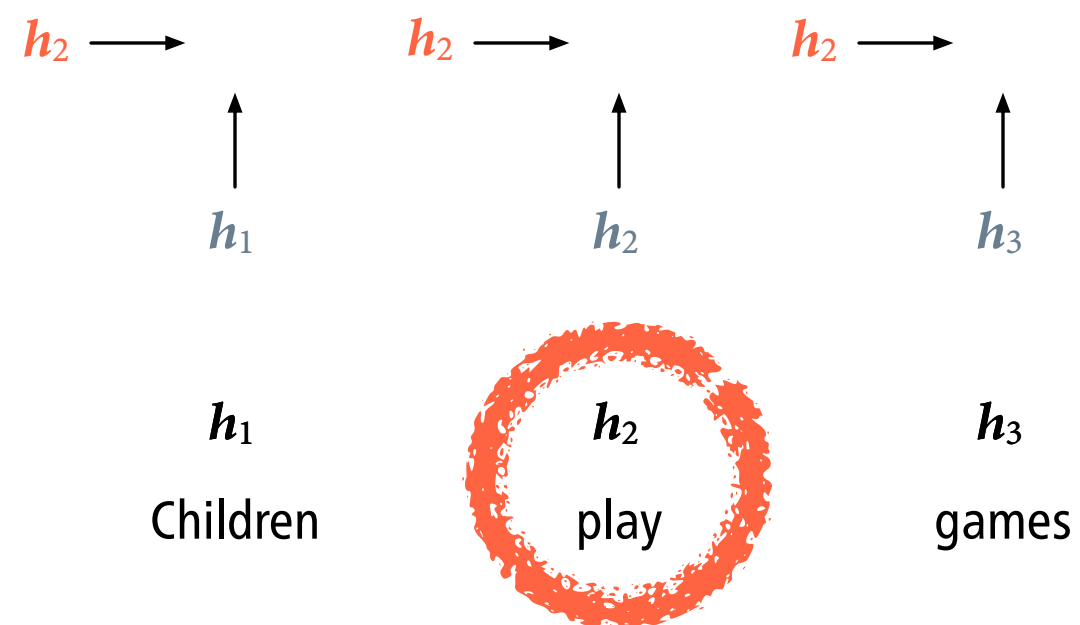


# Contextual embeddings via attention

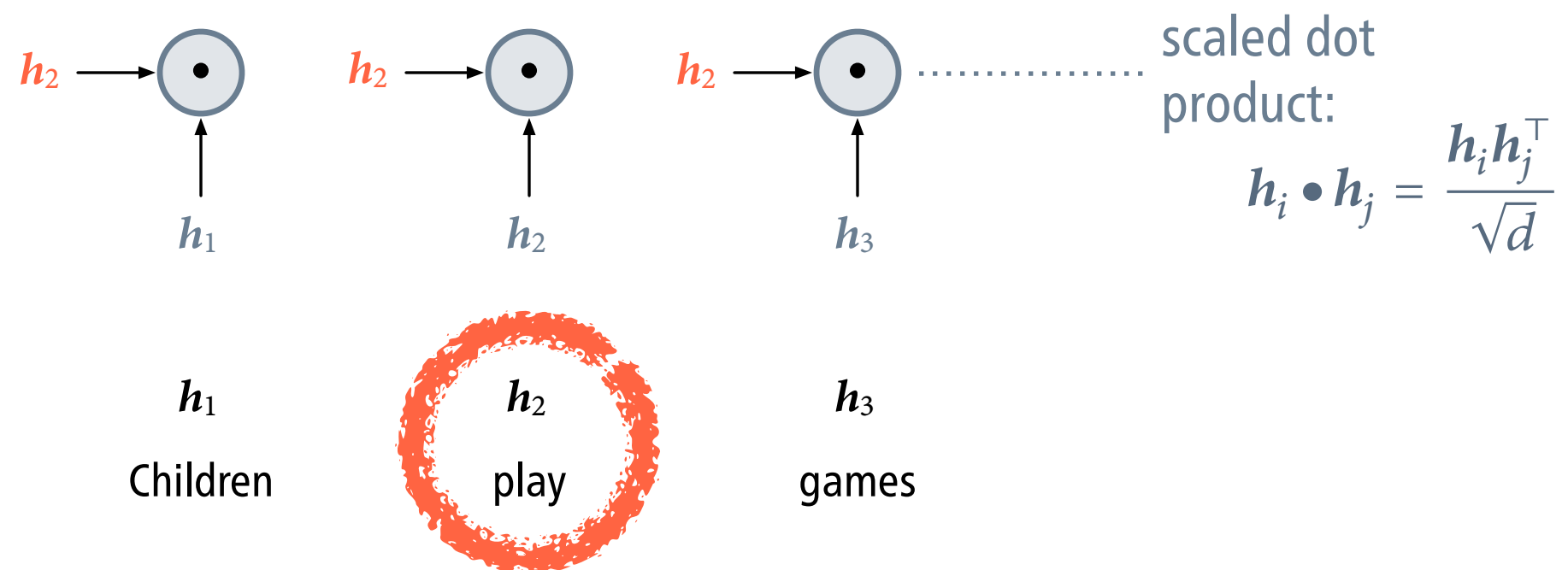




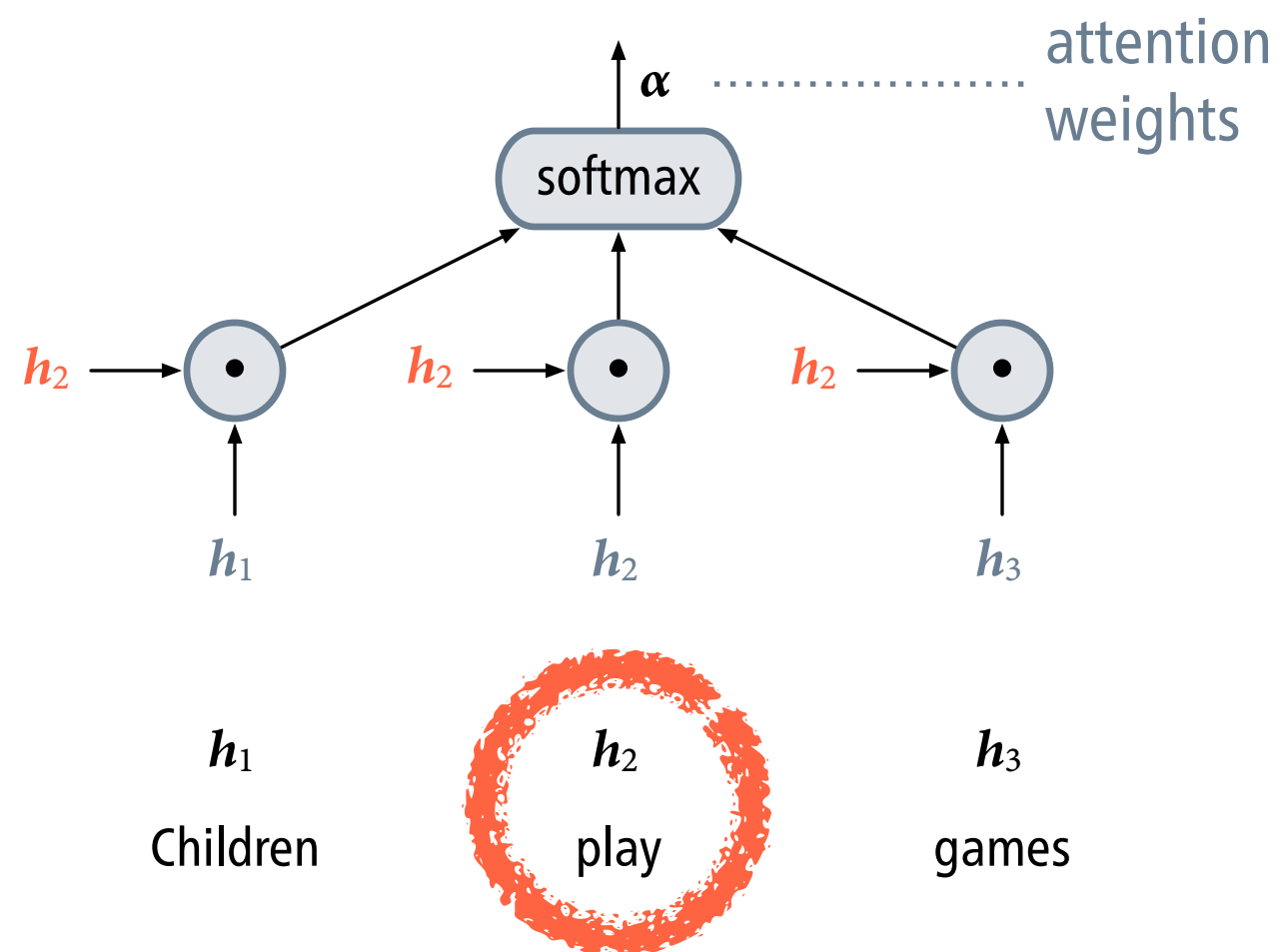
# Contextual embeddings via attention



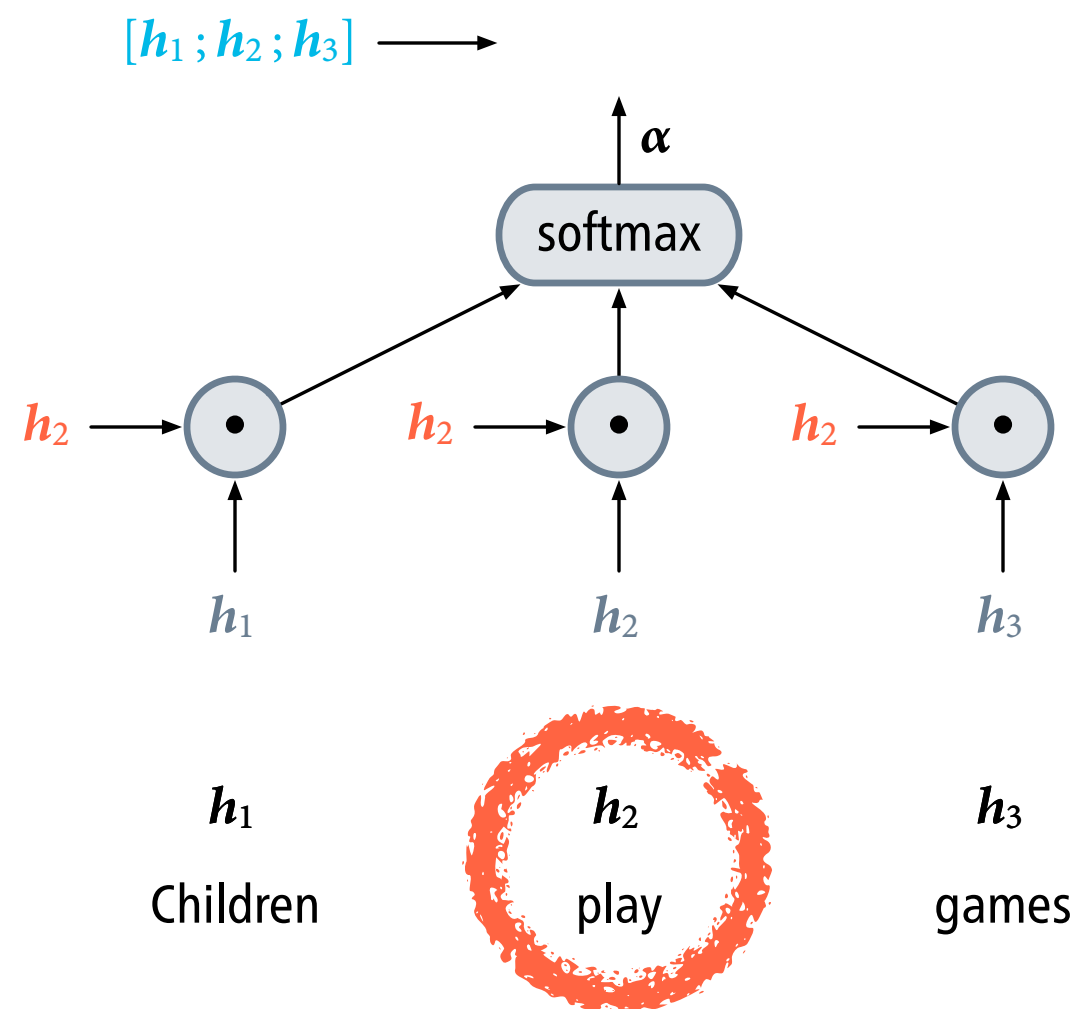
# Contextual embeddings via attention



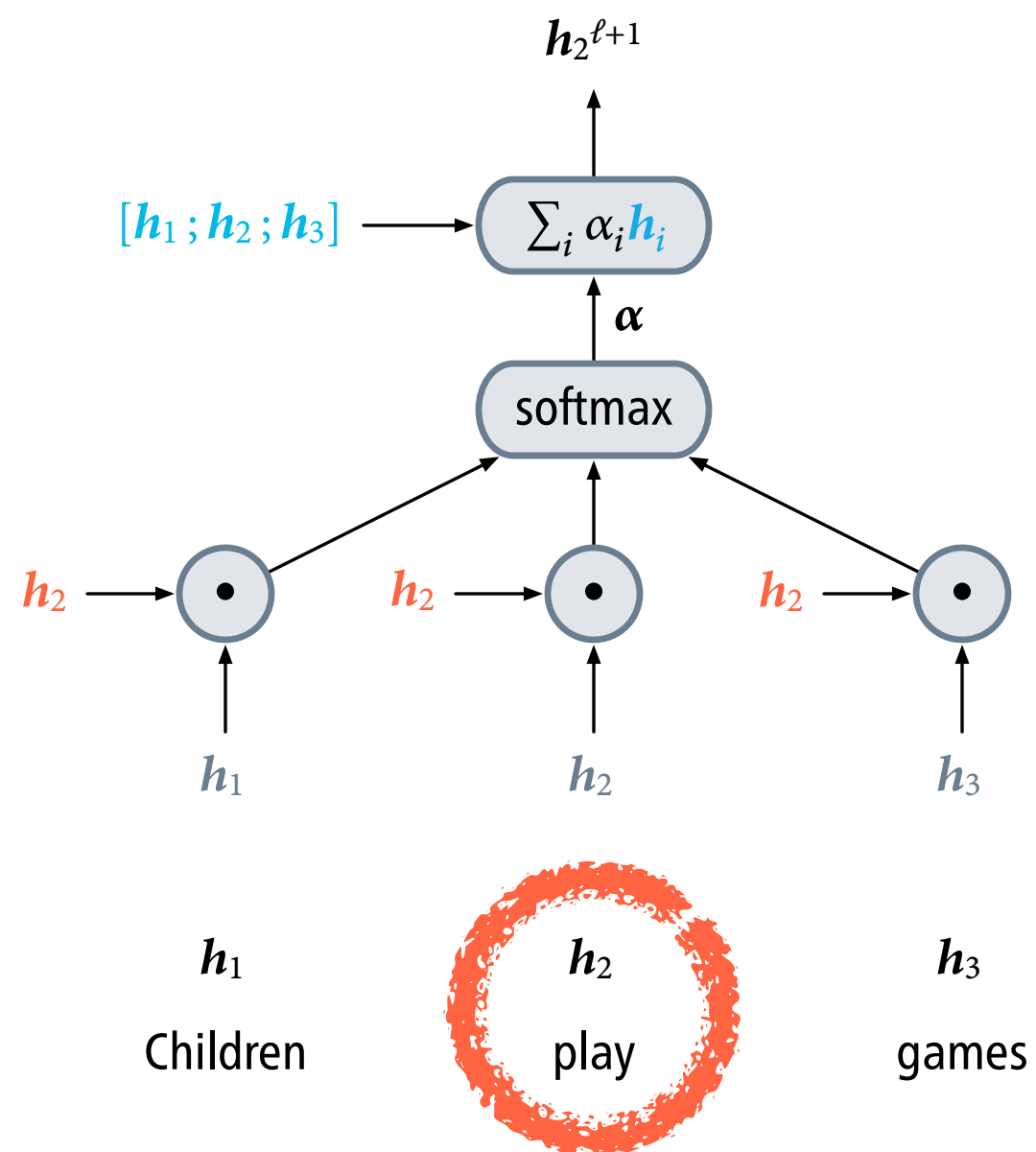
# Contextual embeddings via attention



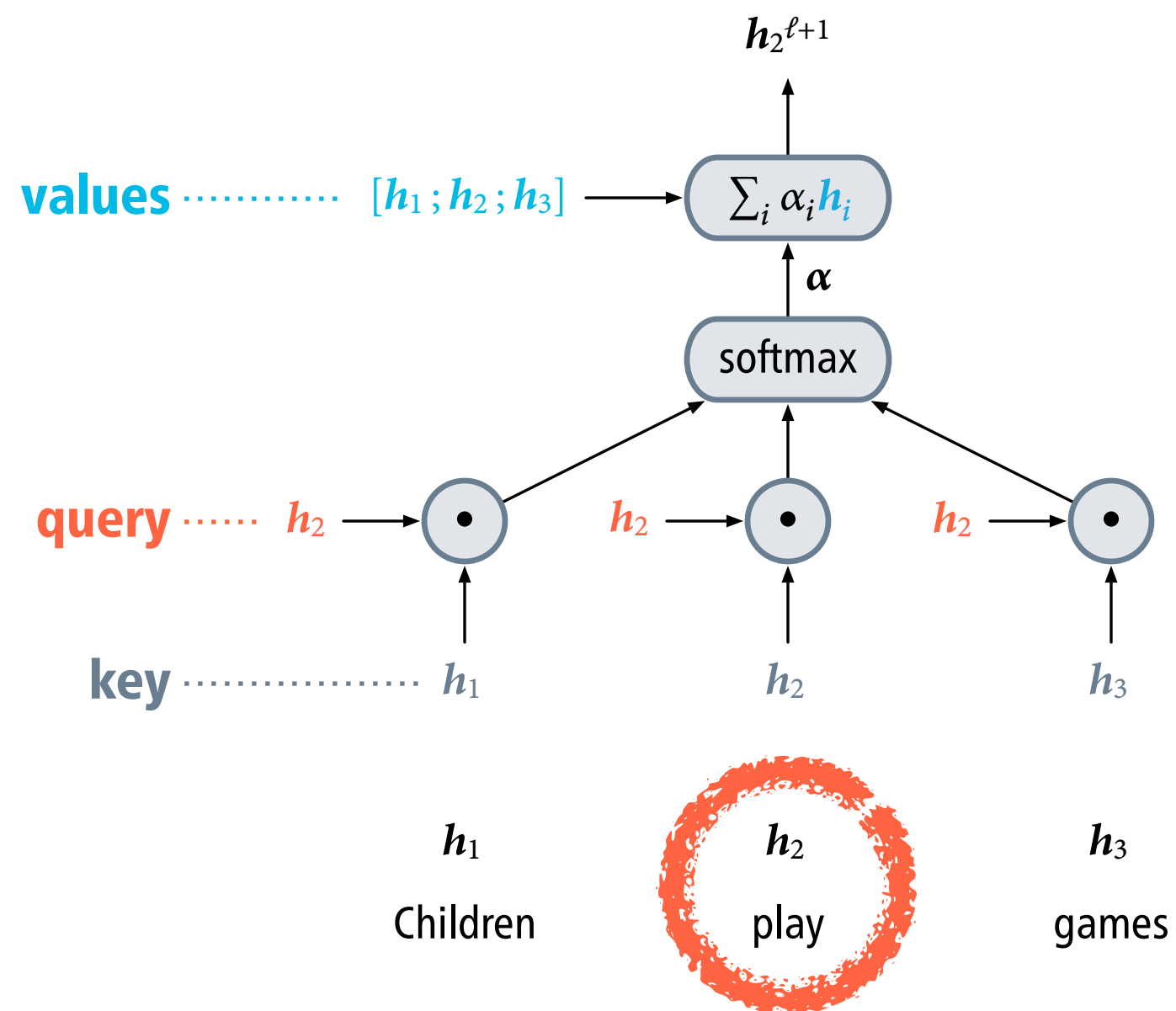
# Contextual embeddings via attention



# Contextual embeddings via attention



# Queries, keys, and values



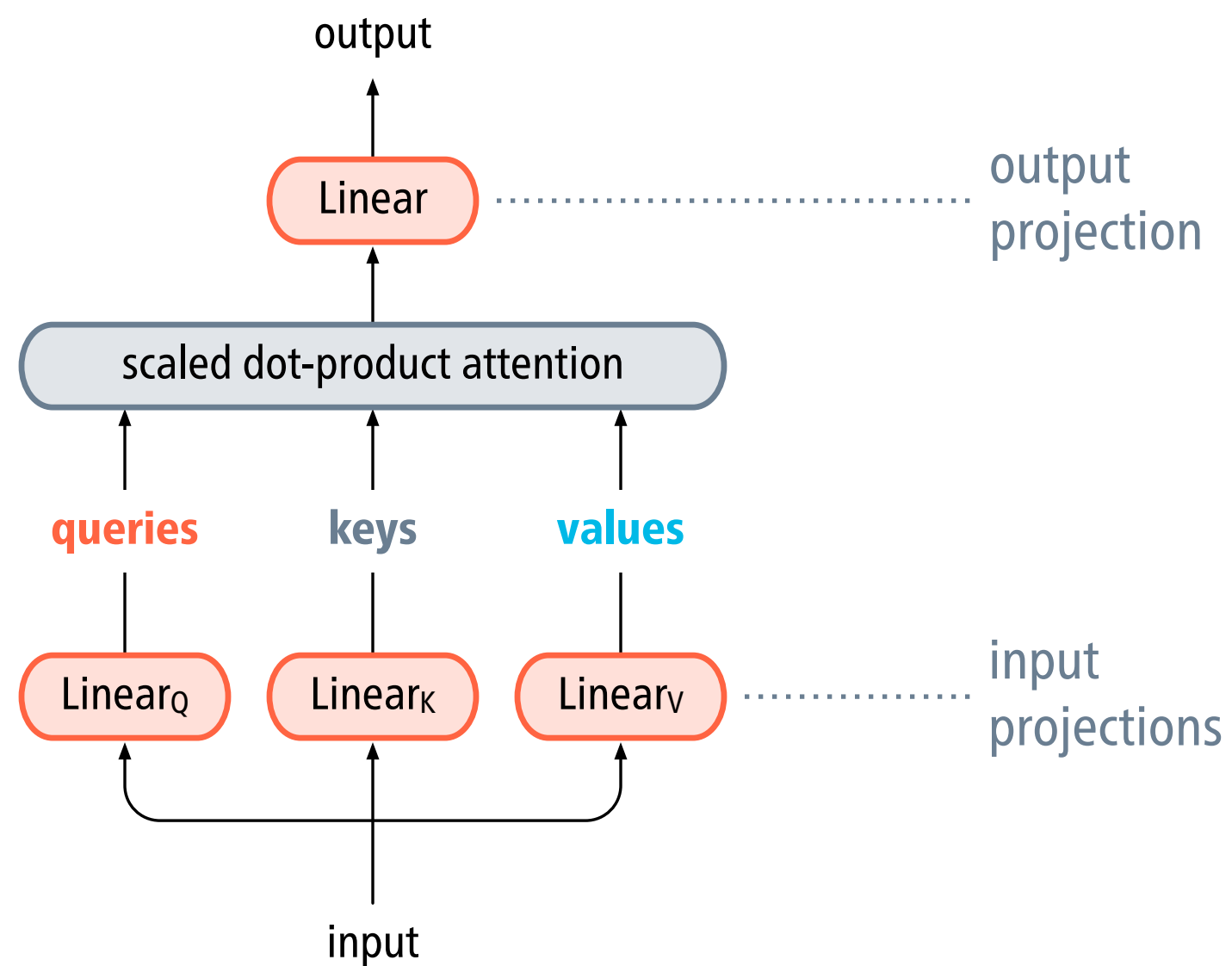
# A general characterisation of attention

- In general, attention can be described as a mapping from a query  $\mathbf{q}$  and a set of key–value pairs  $\langle \mathbf{k}_i, \mathbf{v}_i \rangle$  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 attention score between  $\mathbf{q}$  and  $\mathbf{k}_i$ :

$$\alpha_i = \text{softmax}(\text{score}(\mathbf{q}, \mathbf{K})) \mathbf{V}$$

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

# Attention in the Transformer





# Scaled dot-product attention in PyTorch

```
# Input: queries q, keys k, values v
# shape of q, k, v: [num_words, d]

# Compute the attention scores (scaled dot product)
scores = q @ k.transpose(-1, -2) / hidden_dim**0.5
# shape of scores: [num_words, num_words]

# Normalise the attention scores
alphas = F.softmax(scores, dim=-1)
# shape of alphas: [num_words, num_words]

# The output is the alpha-weighted sum of the values
result = alphas @ v
# shape of result: [num_words, d]
```

# Multi-head attention

