

From Attention to Transformers

Antoine Bosselut

Attention Recap

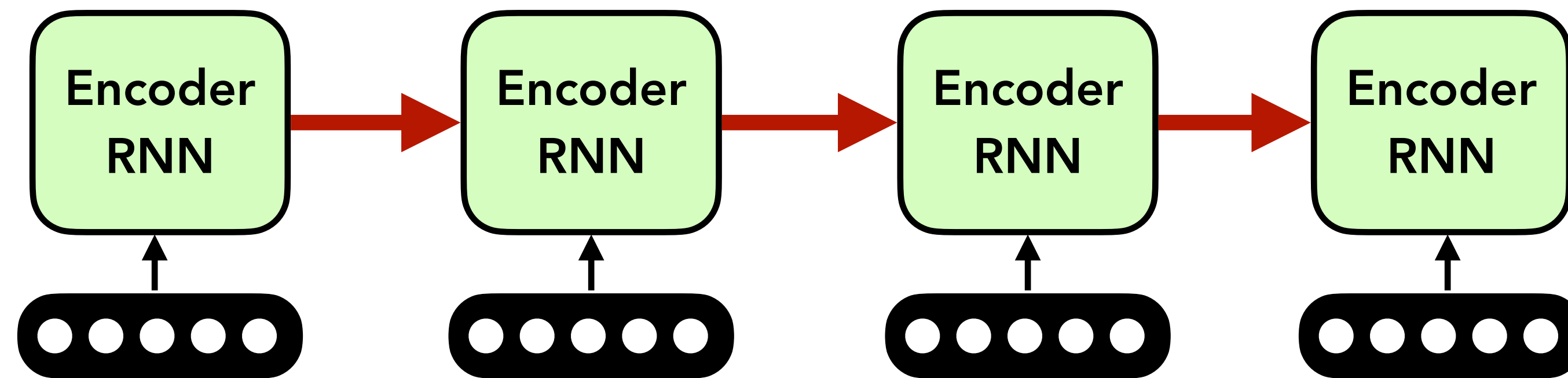
- **Main Idea:** Decoder computes a weighted sum of encoder outputs
 - Compute pairwise score between each encoder hidden state and initial decoder hidden state
- Many possible functions for computing scores (dot product, bilinear, etc.)
- **Temporal Bottleneck Fixed! Direct link** between decoder and encoder states
 - Helps with vanishing gradients and modelling long-term dependencies!
- Attention is **agnostic** to the type of RNN used in the encoder and decoder!

Question

Do any other inefficiencies remain in our sequence to sequence pipelines?

Encoder is still Recurrent

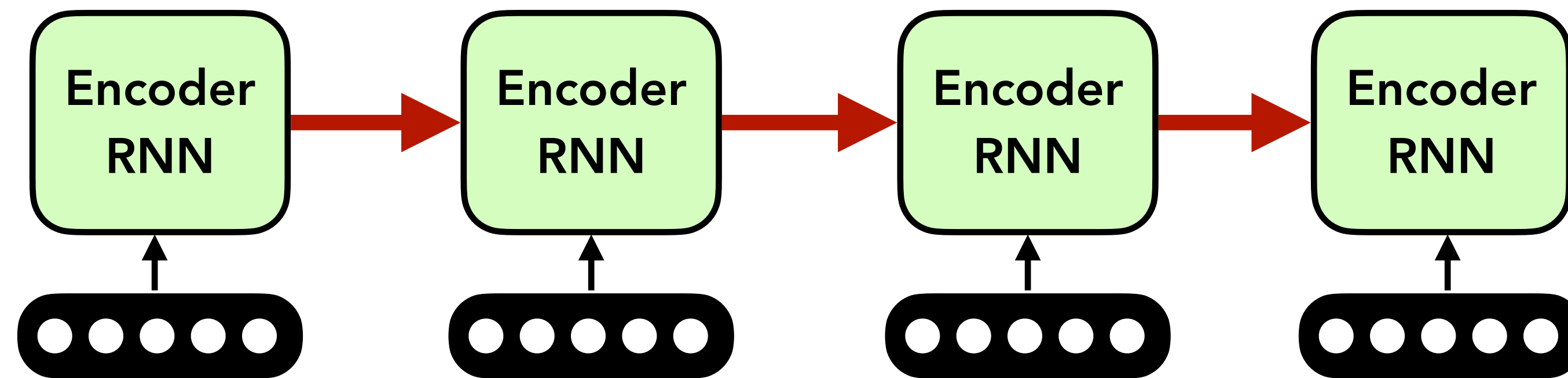
- **Encoder:** Recurrent functions can't be parallelized because previous state needs to be computed to encode next one



- **Problem:** Encoder hidden states must still be computed in series

Encoder is still Recurrent

- **Encoder:** Recurrent functions can't be parallelized because previous state needs to be computed to encode next one



- **Problem:** Encoder hidden states must still be computed in series

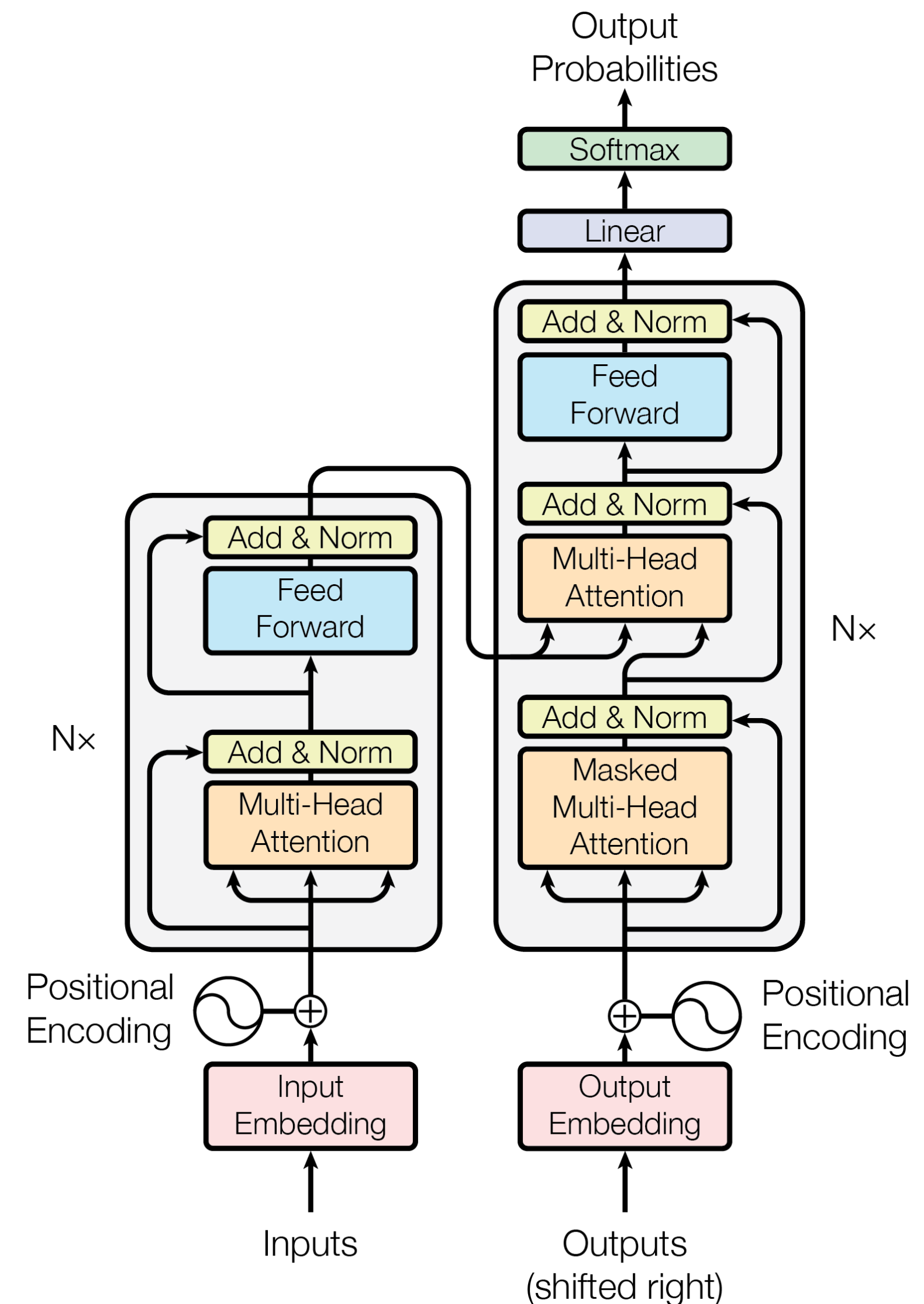
Who can think of a task where this might be a problem?

Solution:
Transformers!

Full Transformer

- Made up of encoder and decoder
- Both encoder and decoder made up of multiple cascaded transformer blocks
 - slightly different architecture in encoder and decoder transformer blocks
- Blocks generally made up **multi-headed attention** layers (self-attention) and **feedforward** layers
- No recurrent computations!

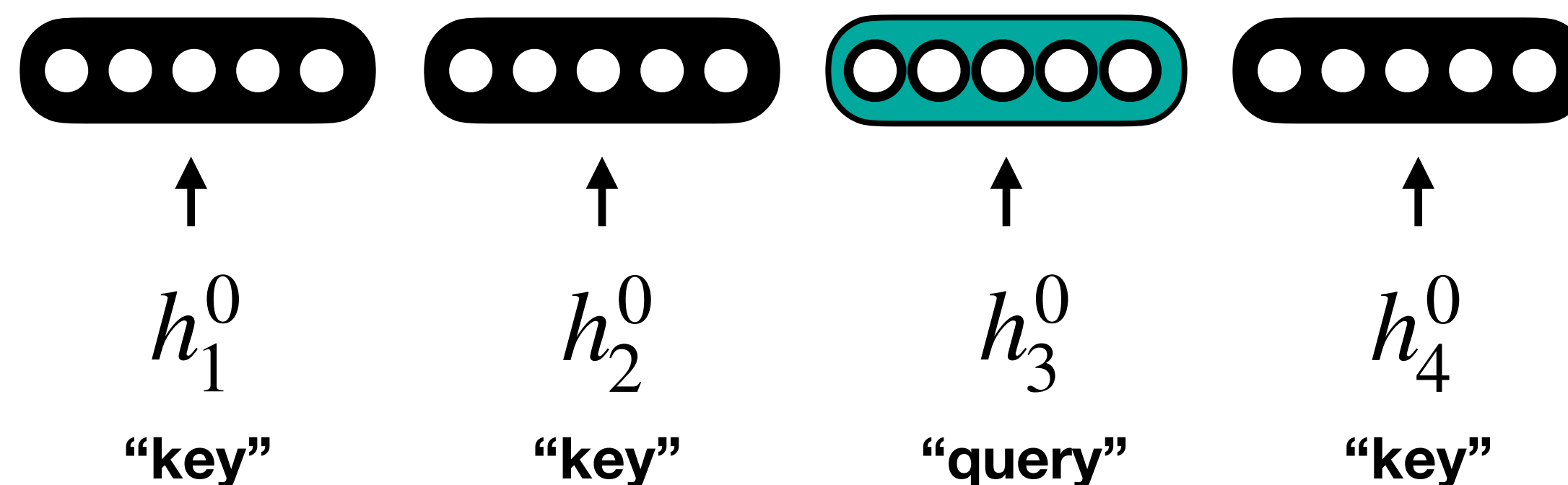
Encode sequences with self-attention



N

- **Original Idea:** Use decoder hidden state to compute attention distribution over encoder hidden states
- **New Idea:** Could we use encoder hidden states to compute attention distribution over themselves?
- **Ditch recurrence** and compute encoder state representations in parallel!

h_t^ℓ = encoder hidden state at time step t at layer ℓ



Note: Subscripts of h have switched back to t

Recap: Attention with RNNs

- **Compute** pairwise similarity between each encoder hidden state and decoder hidden state ("idea of what to decode")

$$\begin{array}{ccc} a_1 = f\left(\underbrace{\begin{array}{c} \text{key} \\ h_1^e \end{array}}_{\text{"key"}}, \underbrace{\begin{array}{c} \text{query} \\ h_1^d \end{array}}_{\text{"query"}}\right) & a_2 = f\left(\underbrace{\begin{array}{c} \text{key} \\ h_2^e \end{array}}_{\text{"key"}}, \underbrace{\begin{array}{c} \text{query} \\ h_1^d \end{array}}_{\text{"query"}}\right) & a_3 = f\left(\underbrace{\begin{array}{c} \text{key} \\ h_3^e \end{array}}_{\text{"key"}}, \underbrace{\begin{array}{c} \text{query} \\ h_1^d \end{array}}_{\text{"query"}}\right) \end{array}$$

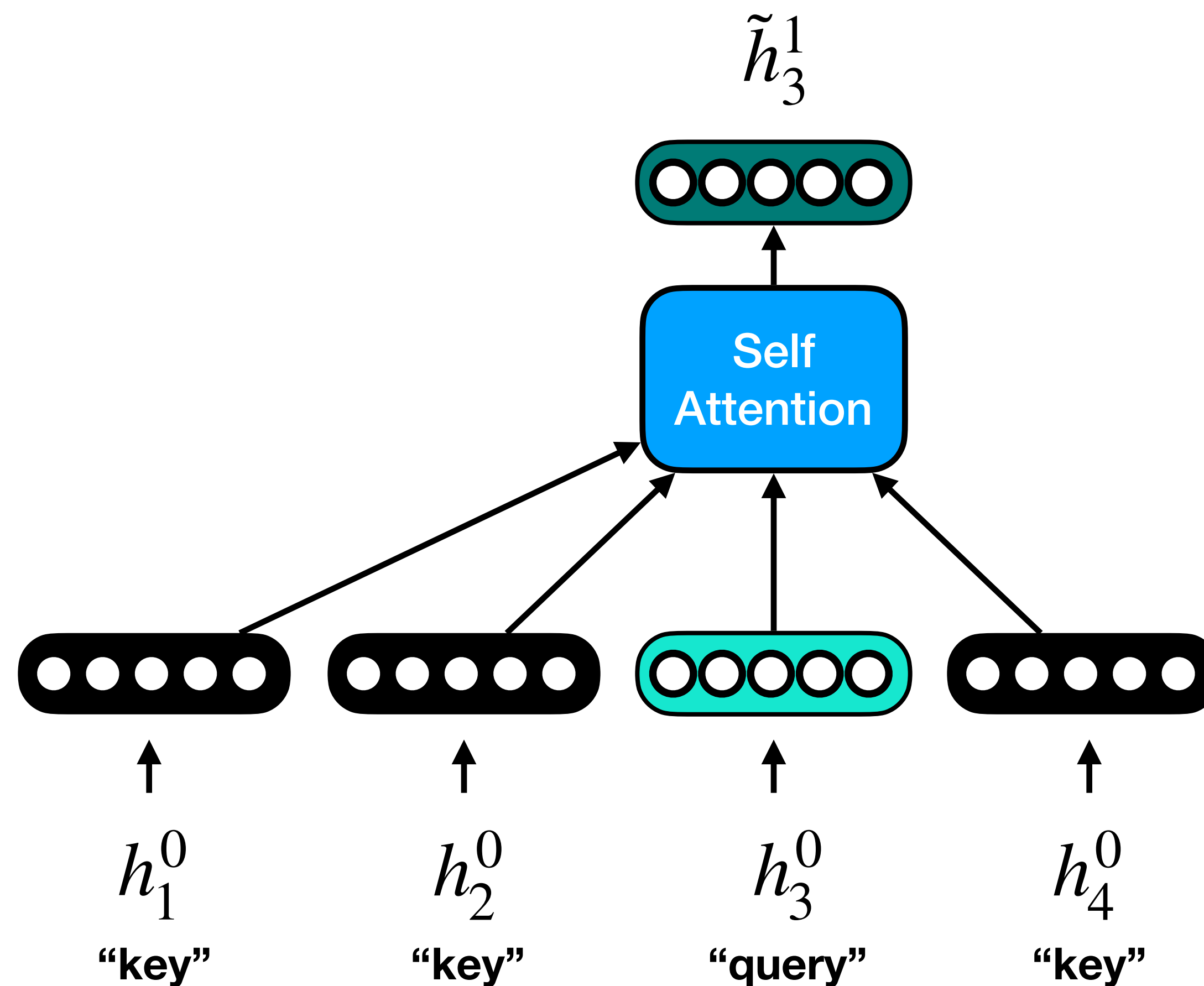
- **Convert** pairwise similarity scores to probability **distribution** (using softmax!) over encoder hidden states and compute weighted average:

Softmax!

$$\alpha_t = \frac{e^{a_t}}{\sum_j e^{a_j}} \rightarrow \begin{array}{c} \text{Bar chart of } \alpha_t \end{array} \rightarrow \tilde{h}_1^d = \sum_{t=1}^T \alpha_t h_t^e$$

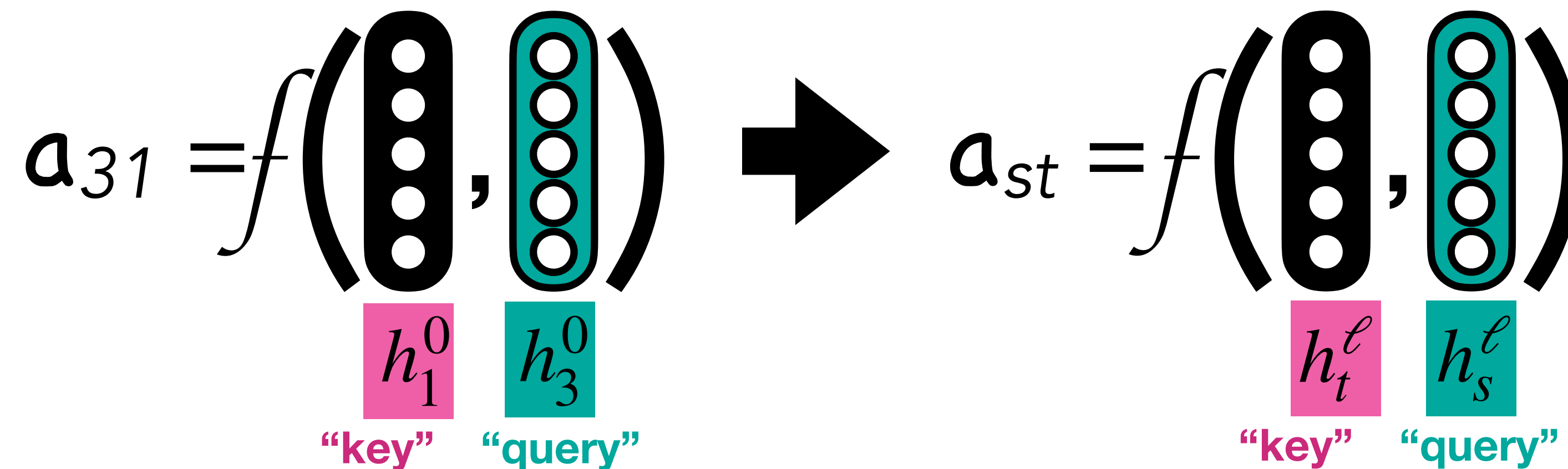
Here h_t^e is known as the "value"

Self-Attention Toy Example



Self-Attention Toy Example

h_t^ℓ = encoder hidden state at time step t at layer ℓ



$$a_{st} = \frac{(\mathbf{W}^Q h_s^\ell)^T (\mathbf{W}^K h_t^\ell)}{\sqrt{d}}$$

Compute pairwise scores

$$\alpha_{st} = \frac{e^{a_{st}}}{\sum_j e^{a_{sj}}}$$

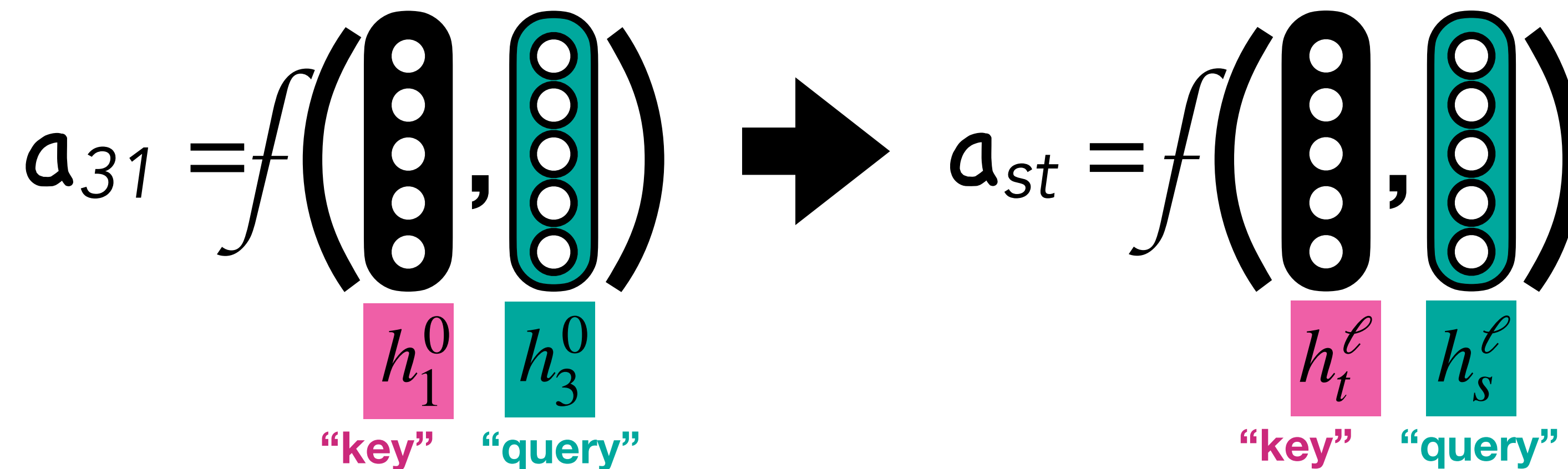
Get attention
distribution

$$\tilde{h}_s^\ell = \sum_{t=1}^T \alpha_{st} (\mathbf{W}^V h_t^\ell)$$

Attend to values to
get weighted sum

Self-Attention Toy Example

h_t^ℓ = encoder hidden state at time step t at layer ℓ



$\{1, \dots, t, \dots, T\}$
includes s !

$$a_{st} = \frac{(\mathbf{W}^Q h_s^\ell)^T (\mathbf{W}^K h_t^\ell)}{\sqrt{d}}$$

Compute pairwise scores

$$\alpha_{st} = \frac{e^{a_{st}}}{\sum_j e^{a_{sj}}}$$

Get attention
distribution

$$\tilde{h}_s^\ell = \sum_{t=1}^T \alpha_{st} (\mathbf{W}^V h_t^\ell)$$

Attend to values to
get weighted sum

Self-attention!

Self-Attention Toy Example

Compute pairwise scores

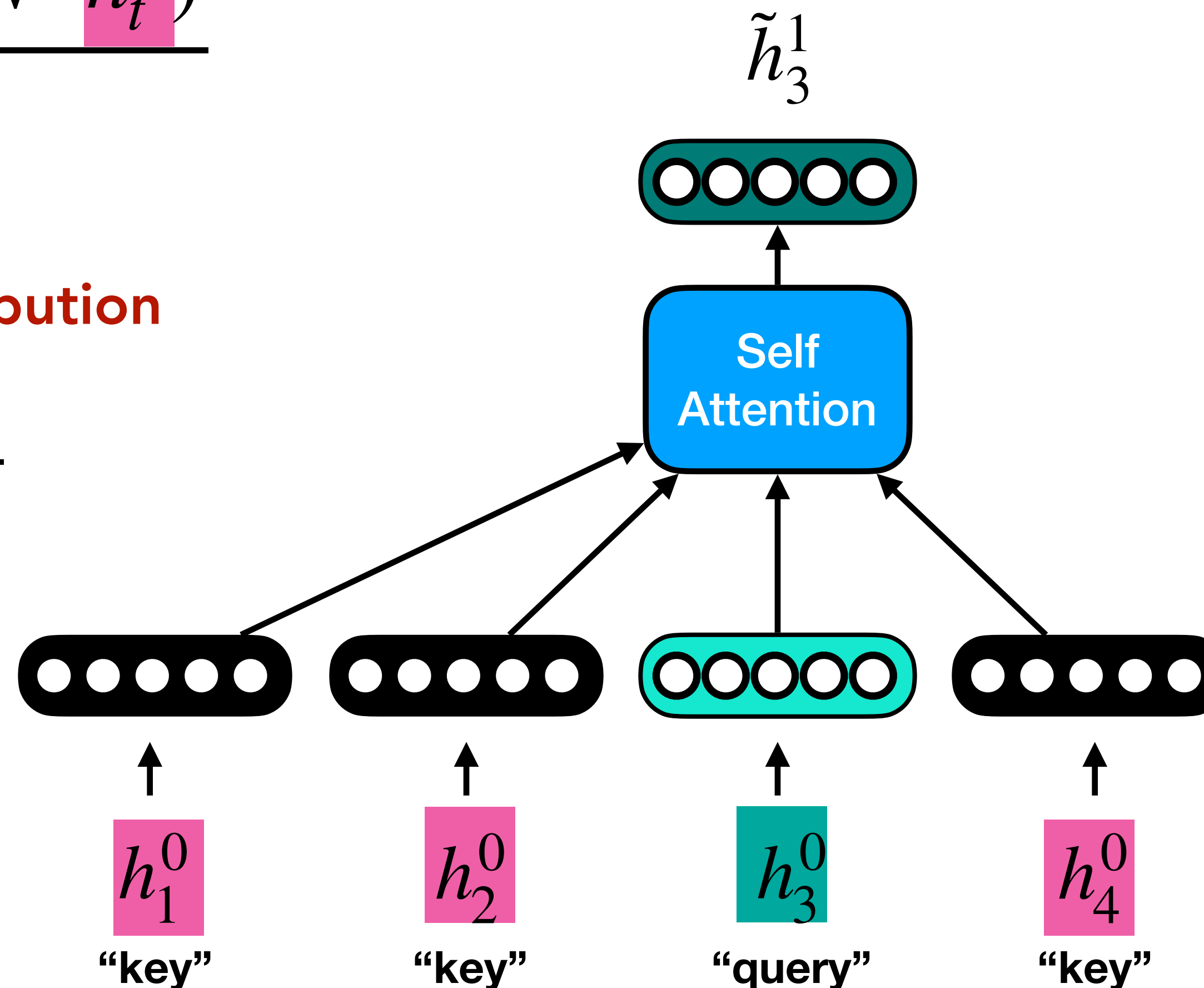
$$a_{st} = \frac{(\mathbf{W}^Q \mathbf{h}_s^\ell)^T (\mathbf{W}^K \mathbf{h}_t^\ell)}{\sqrt{d}}$$

Attend to values to get weighted sum

$$\tilde{\mathbf{h}}_s^\ell = \sum_{t=1}^T \alpha_{st} (\mathbf{W}^V \mathbf{h}_t^\ell)$$

Get attention distribution

$$\alpha_{st} = \frac{e^{a_{st}}}{\sum_j e^{a_{sj}}}$$



Self-Attention Toy Example

Compute pairwise scores

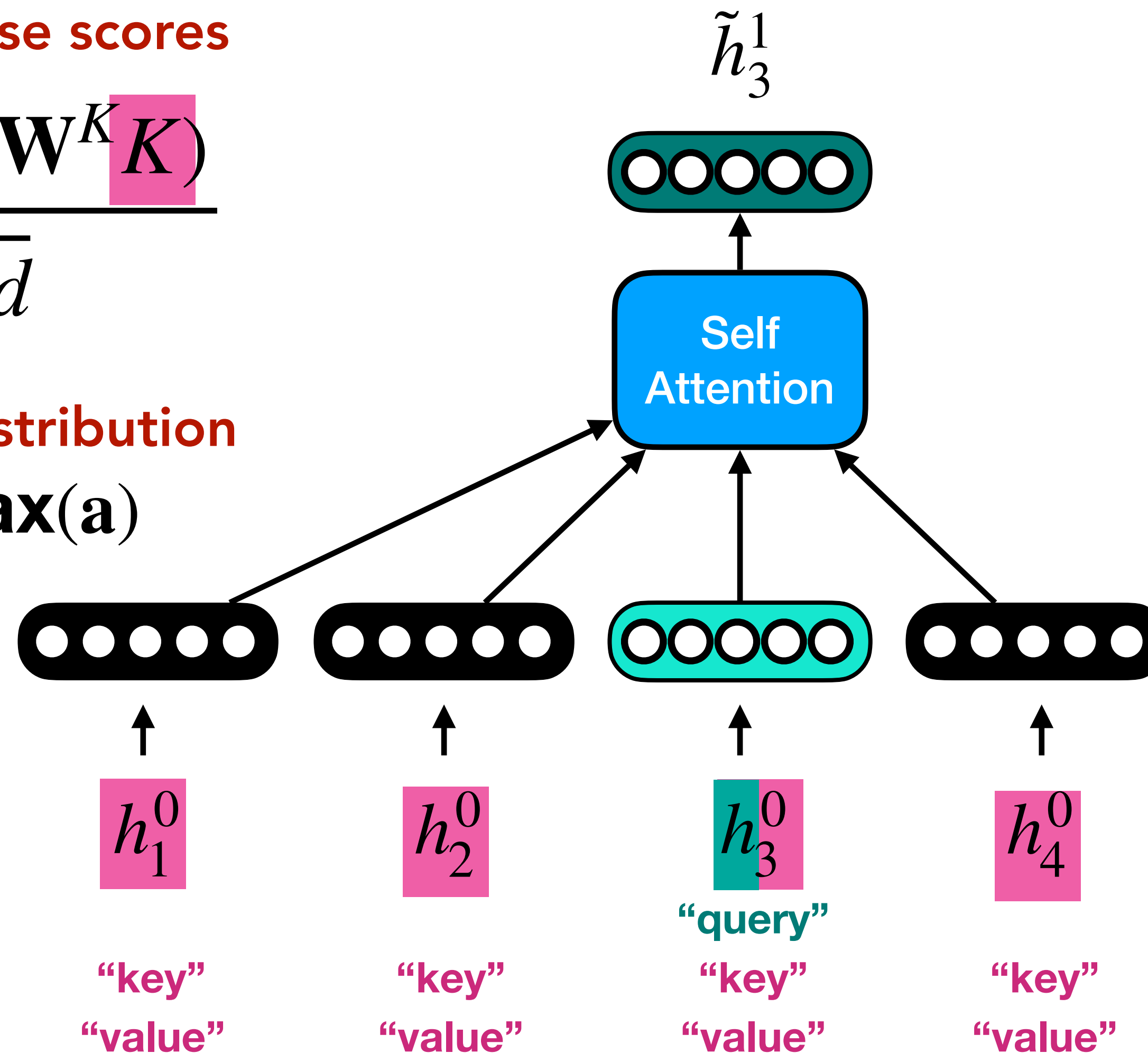
$$\mathbf{a} = \frac{(\mathbf{W}^Q \mathbf{q})(\mathbf{W}^K \mathbf{K})}{\sqrt{d}}$$

Get attention distribution

$$\alpha = \text{softmax}(\mathbf{a})$$

Attend to values to
get weighted sum

$$\tilde{h}^\ell = \mathbf{W}^O \alpha (\mathbf{V} \mathbf{W}^V)$$



“query” $\mathbf{q} = h_s^\ell$

“keys” $\mathbf{K} = \mathbf{V} = \{h_t^\ell\}_{t=0}^T$ “values”

For each attention computation, every element is a key and value, and one element is a query

Self-Attention Toy Example

Compute pairwise scores

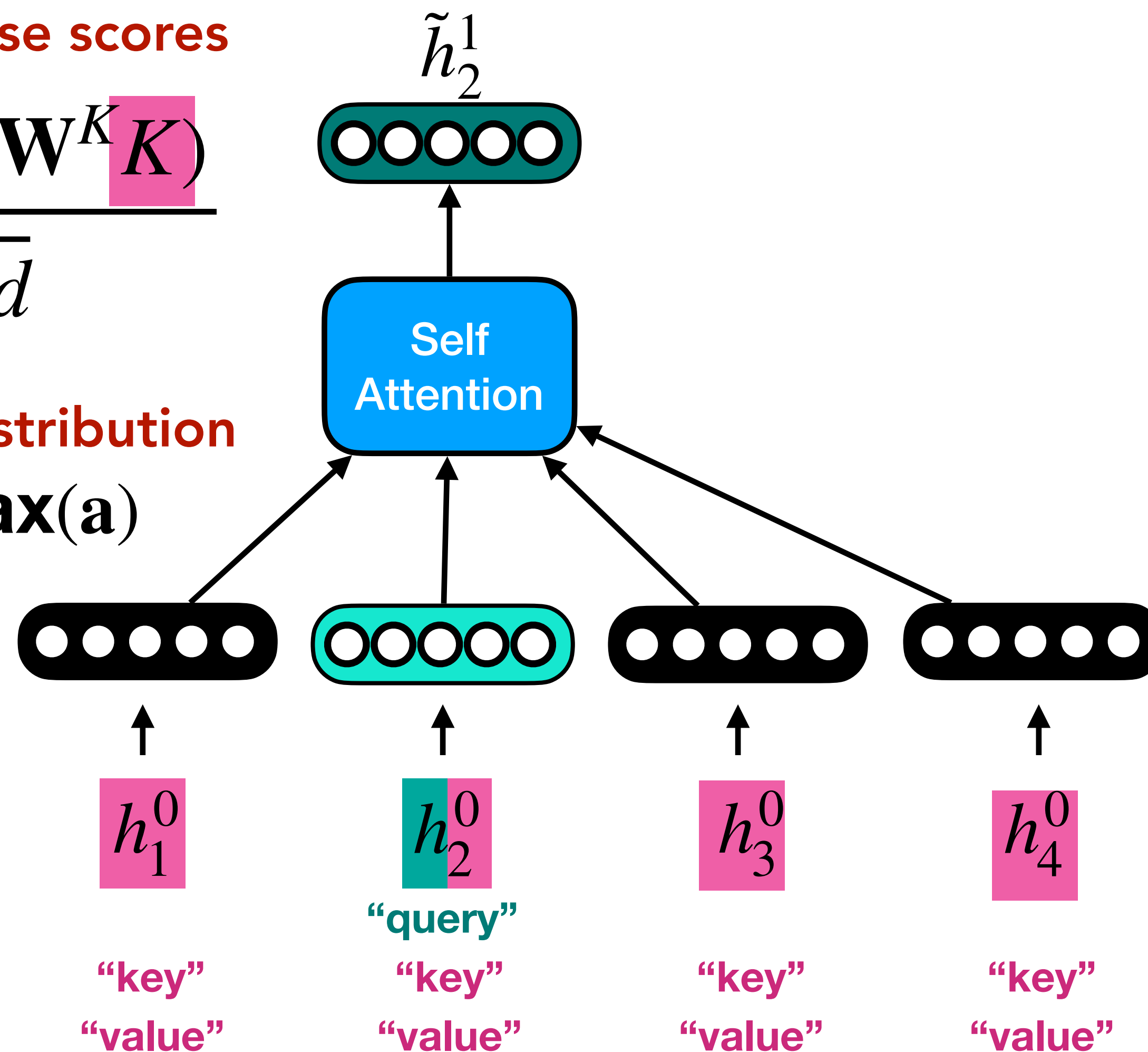
$$\mathbf{a} = \frac{(\mathbf{W}^Q \mathbf{q})(\mathbf{W}^K \mathbf{K})}{\sqrt{d}}$$

Get attention distribution

$$\alpha = \text{softmax}(\mathbf{a})$$

Attend to values to
get weighted sum

$$\tilde{h}^\ell = \mathbf{W}^O \alpha (\mathbf{V} \mathbf{W}^V)$$



"query" $\mathbf{q} = h_s^\ell$

"values" $\mathbf{K} = \mathbf{V} = \{h_t^\ell\}_{t=0}^T$
"keys"

For each attention computation, every element is a key and value, and one element is a query

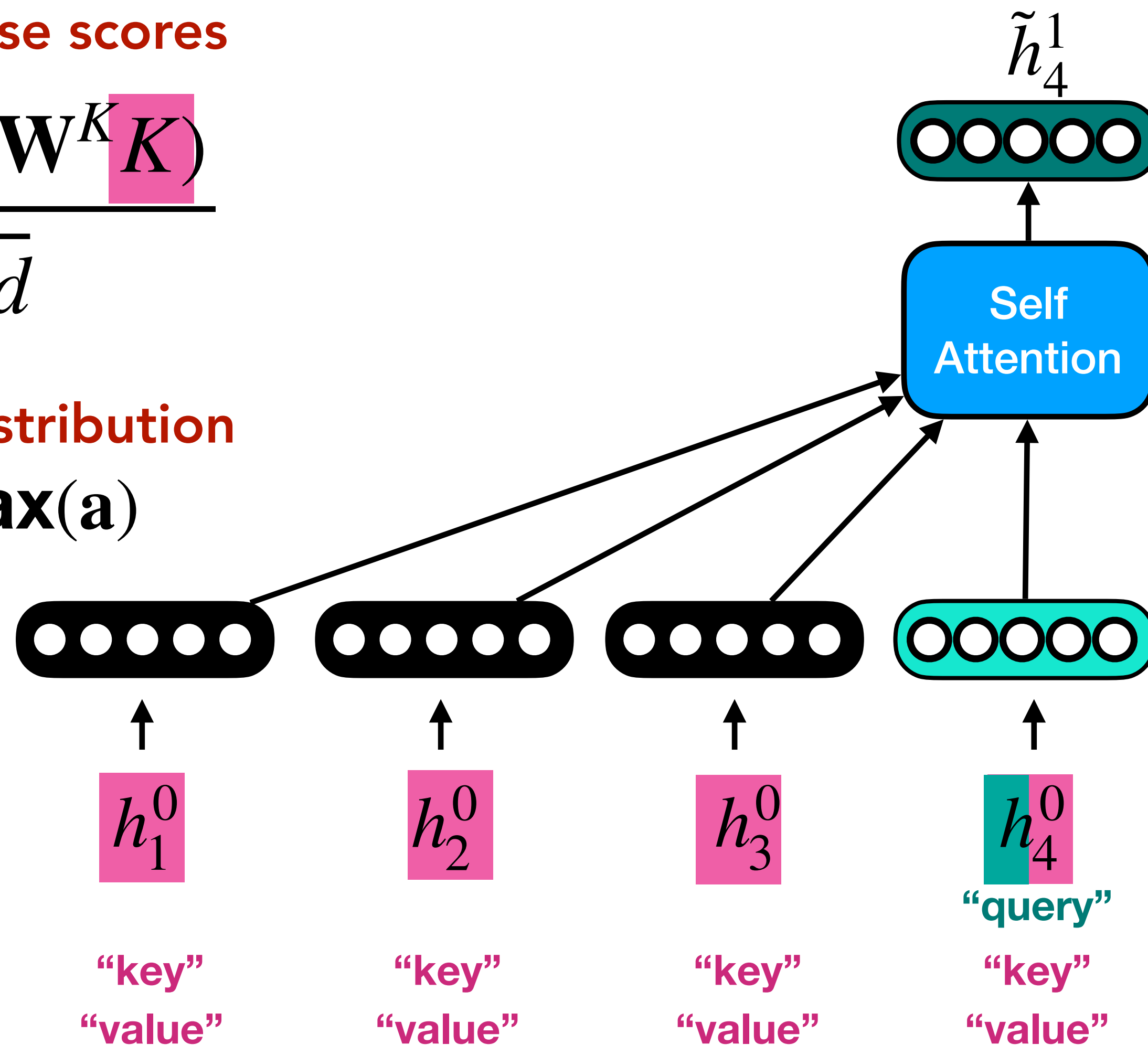
Self-Attention Toy Example

Compute pairwise scores

$$\mathbf{a} = \frac{(\mathbf{W}^Q \mathbf{q})(\mathbf{W}^K \mathbf{K})}{\sqrt{d}}$$

Get attention distribution

$$\alpha = \text{softmax}(\mathbf{a})$$



Attend to values to
get weighted sum

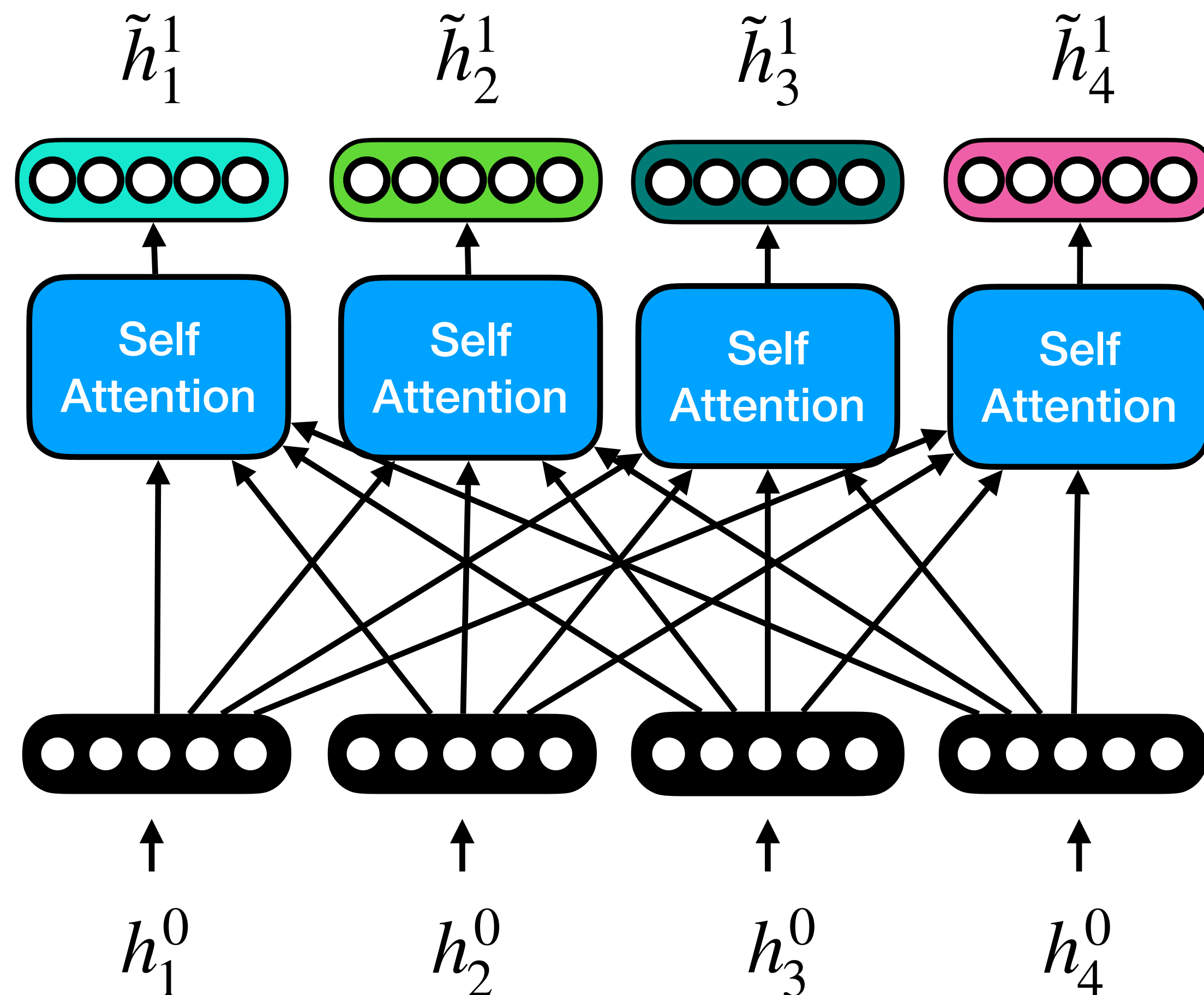
$$\tilde{h}^\ell = W^O \alpha (\mathbf{V} \mathbf{W}^V)$$

$$\text{"query"} \quad \mathbf{q} = h_s^\ell$$

$$\underset{\text{"keys"}}{\mathbf{K}} = \underset{\text{"values"}}{\mathbf{V}} = \{h_t^\ell\}_{t=0}^T$$

For each attention computation, every element is a key and value, and one element is a query

Self-Attention Toy Example



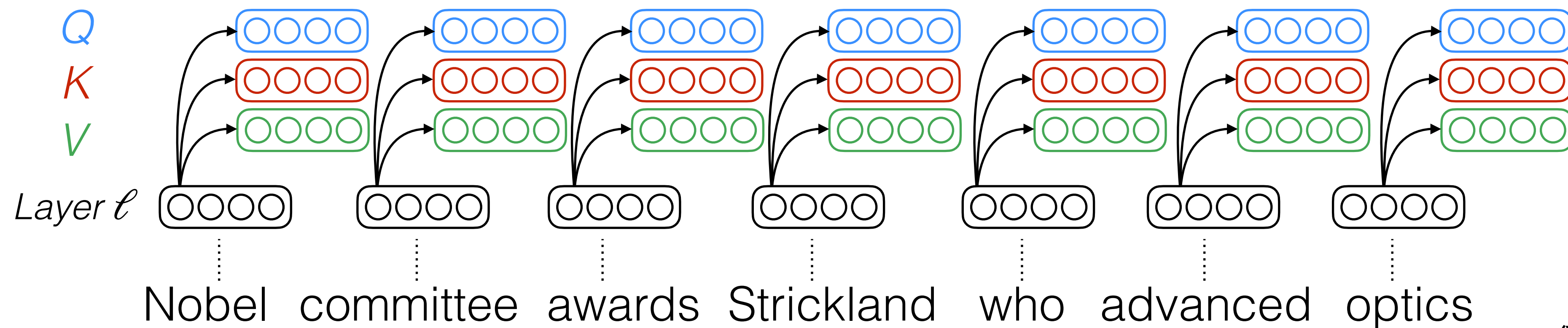
$$\tilde{h}_1^1 = \text{Attention}\left(h_1^0, \{h_t^0\}_{t=0}^{t=3}\right)$$

$$\tilde{h}_2^1 = \text{Attention}\left(h_2^0, \{h_t^0\}_{t=0}^{t=3}\right)$$

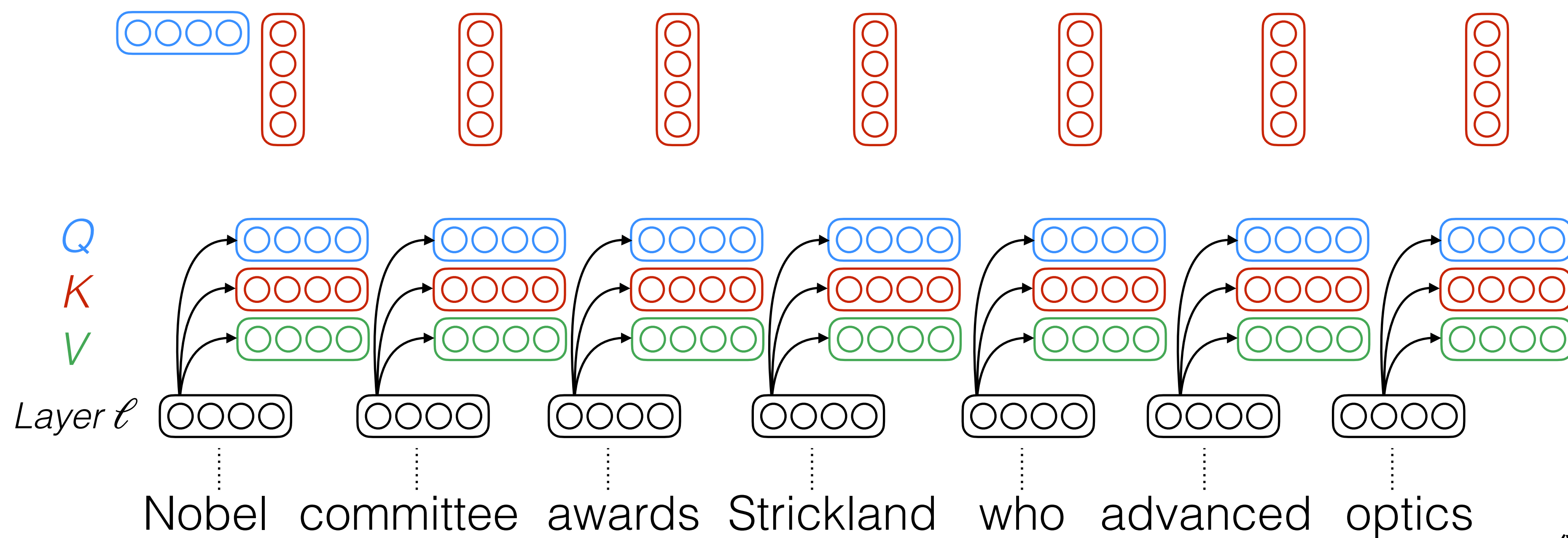
$$\tilde{h}_3^1 = \text{Attention}\left(h_3^0, \{h_t^0\}_{t=0}^{t=3}\right)$$

$$\tilde{h}_4^1 = \text{Attention}\left(h_4^0, \{h_t^0\}_{t=0}^{t=3}\right)$$

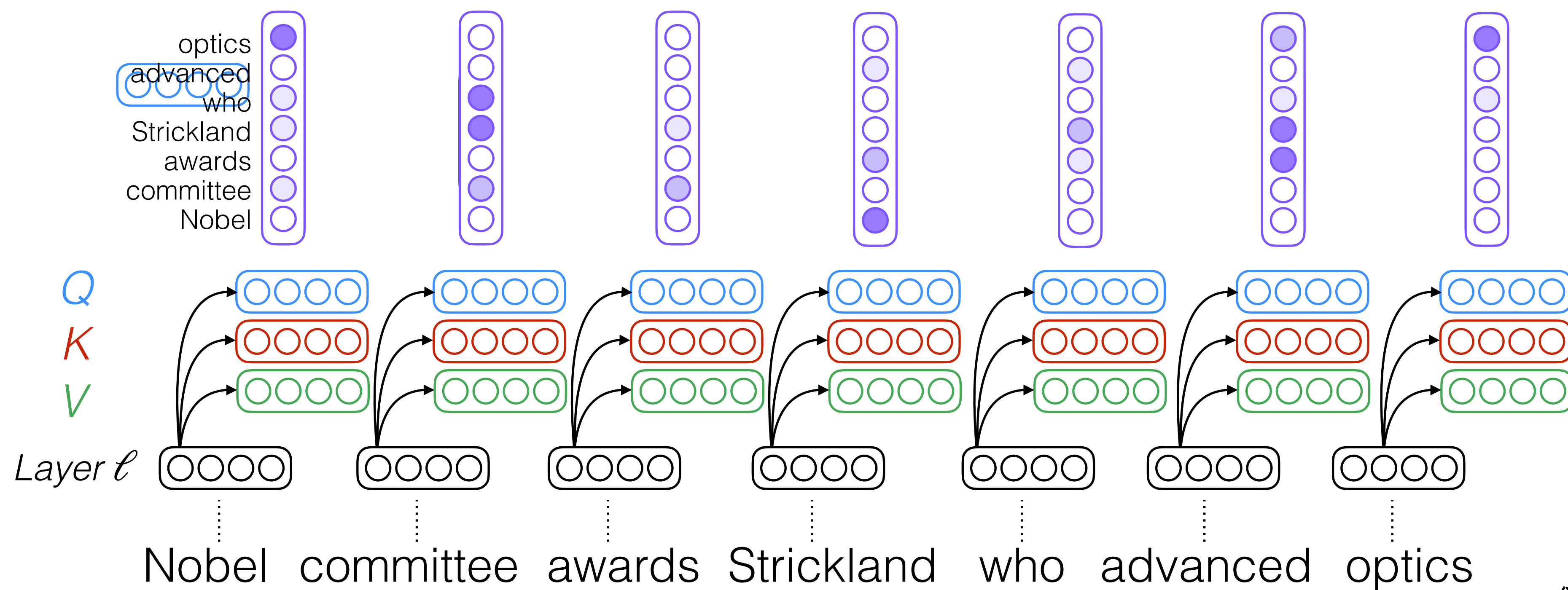
Self-attention (in encoder)



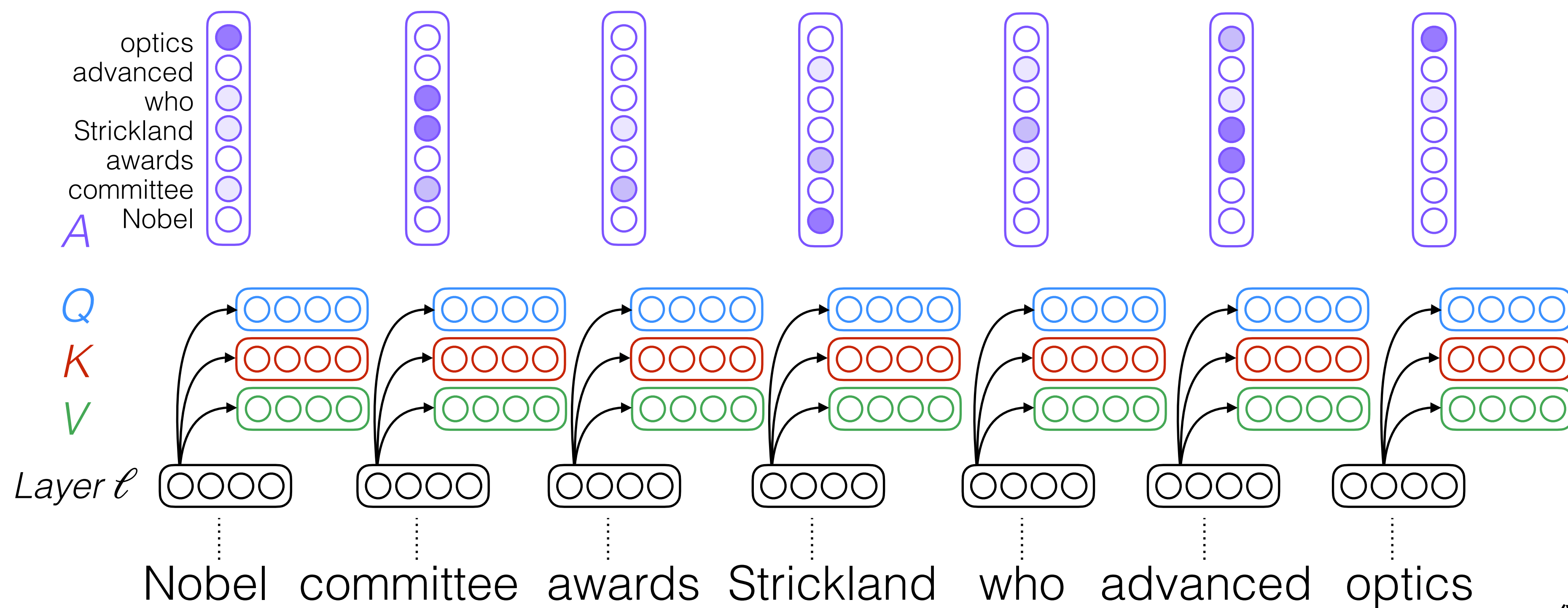
Self-attention (in encoder)



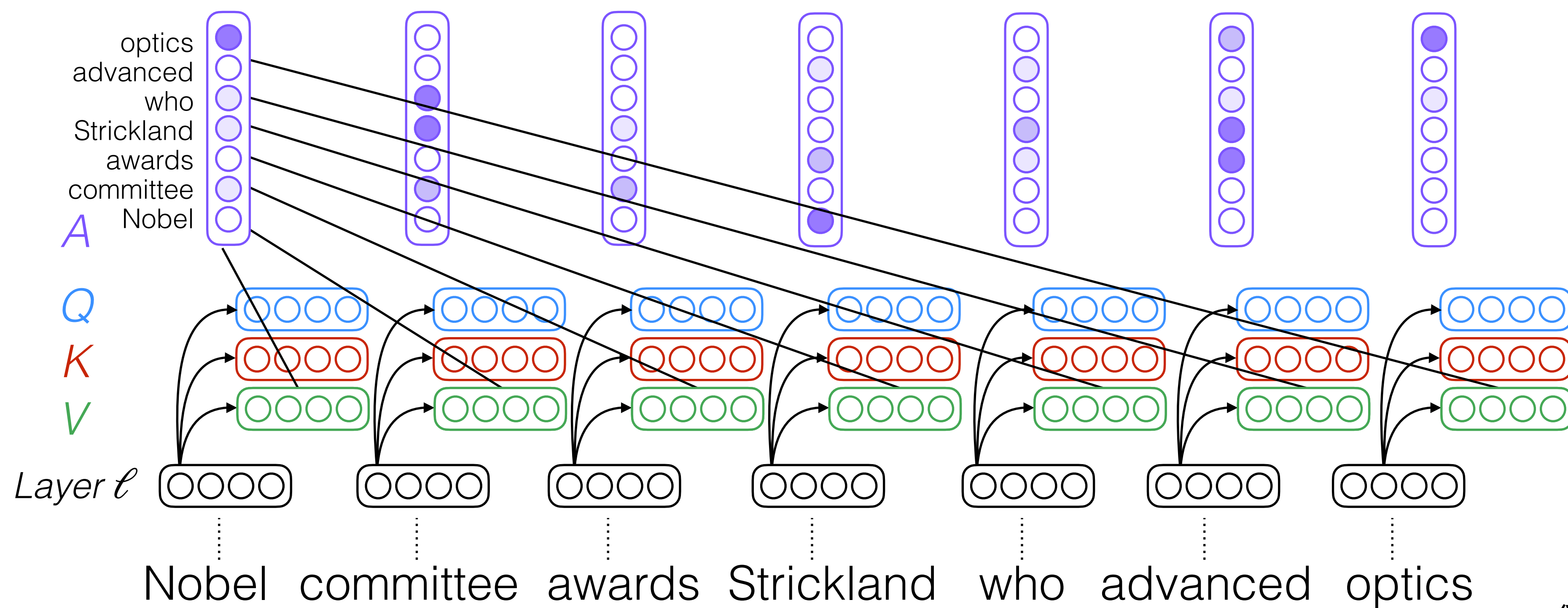
Self-attention (in encoder)



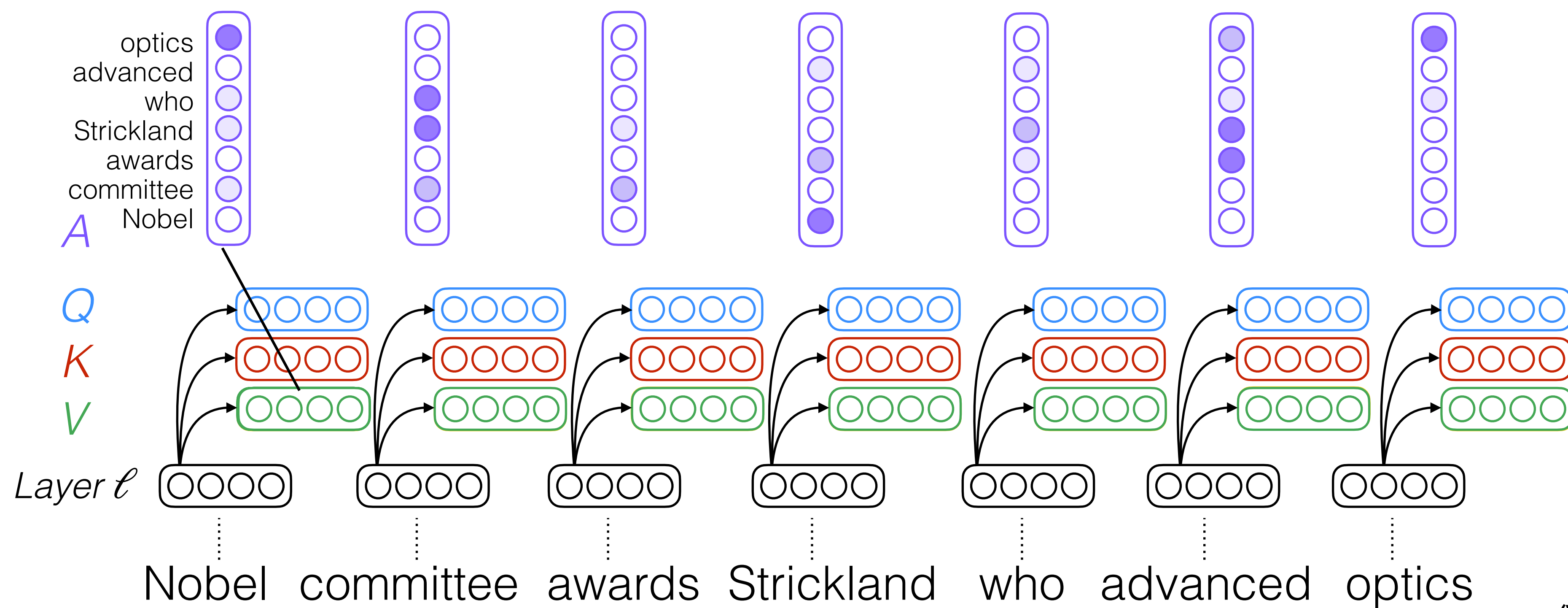
Self-attention (in encoder)



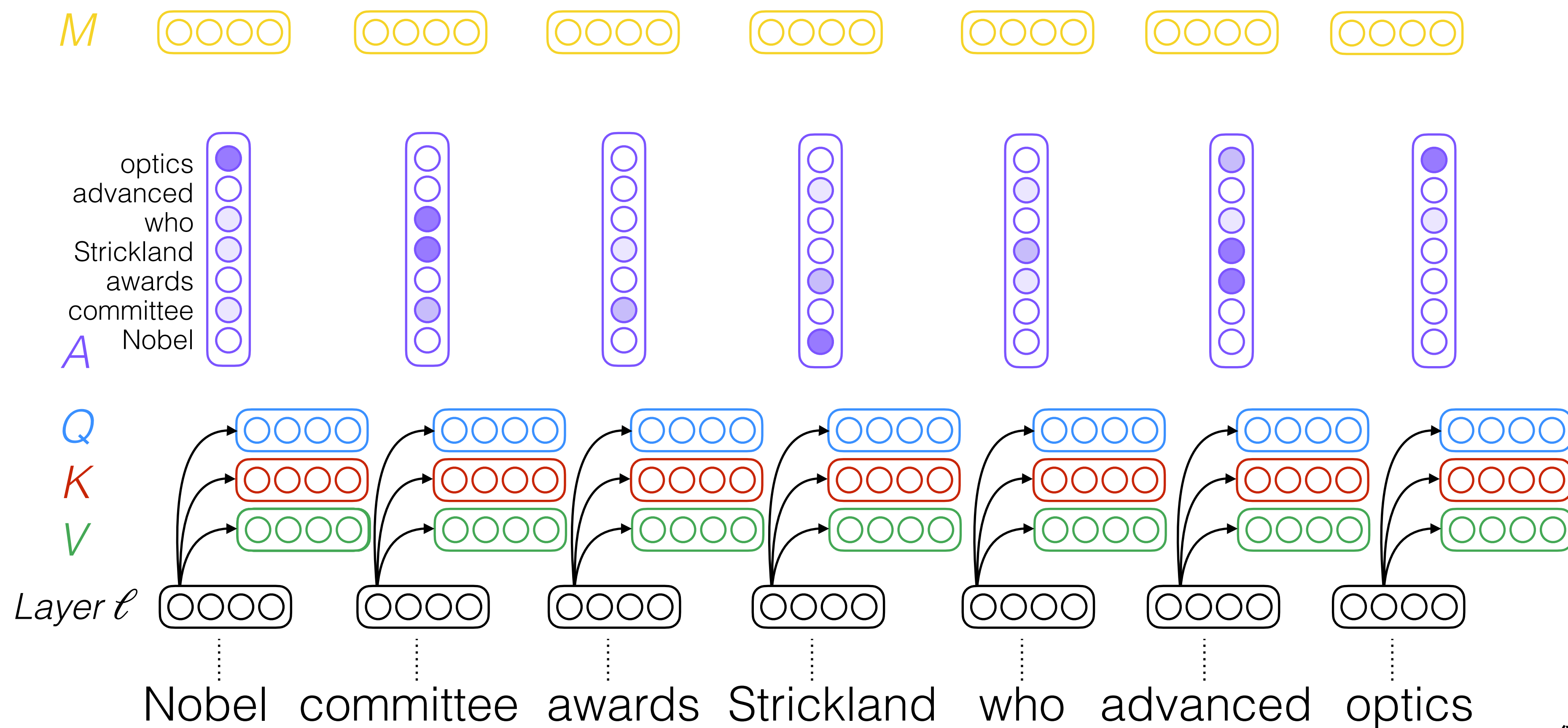
Self-attention (in encoder)



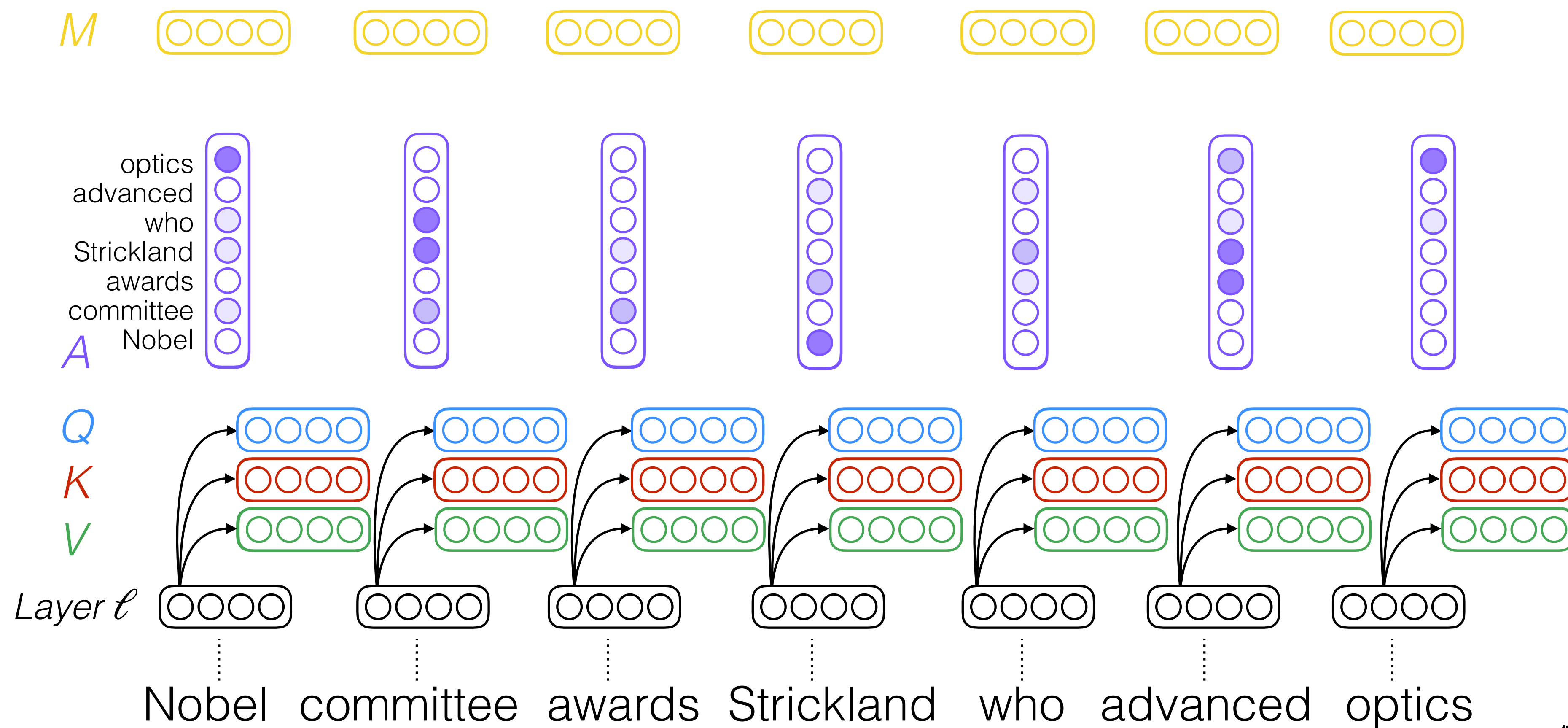
Self-attention (in encoder)



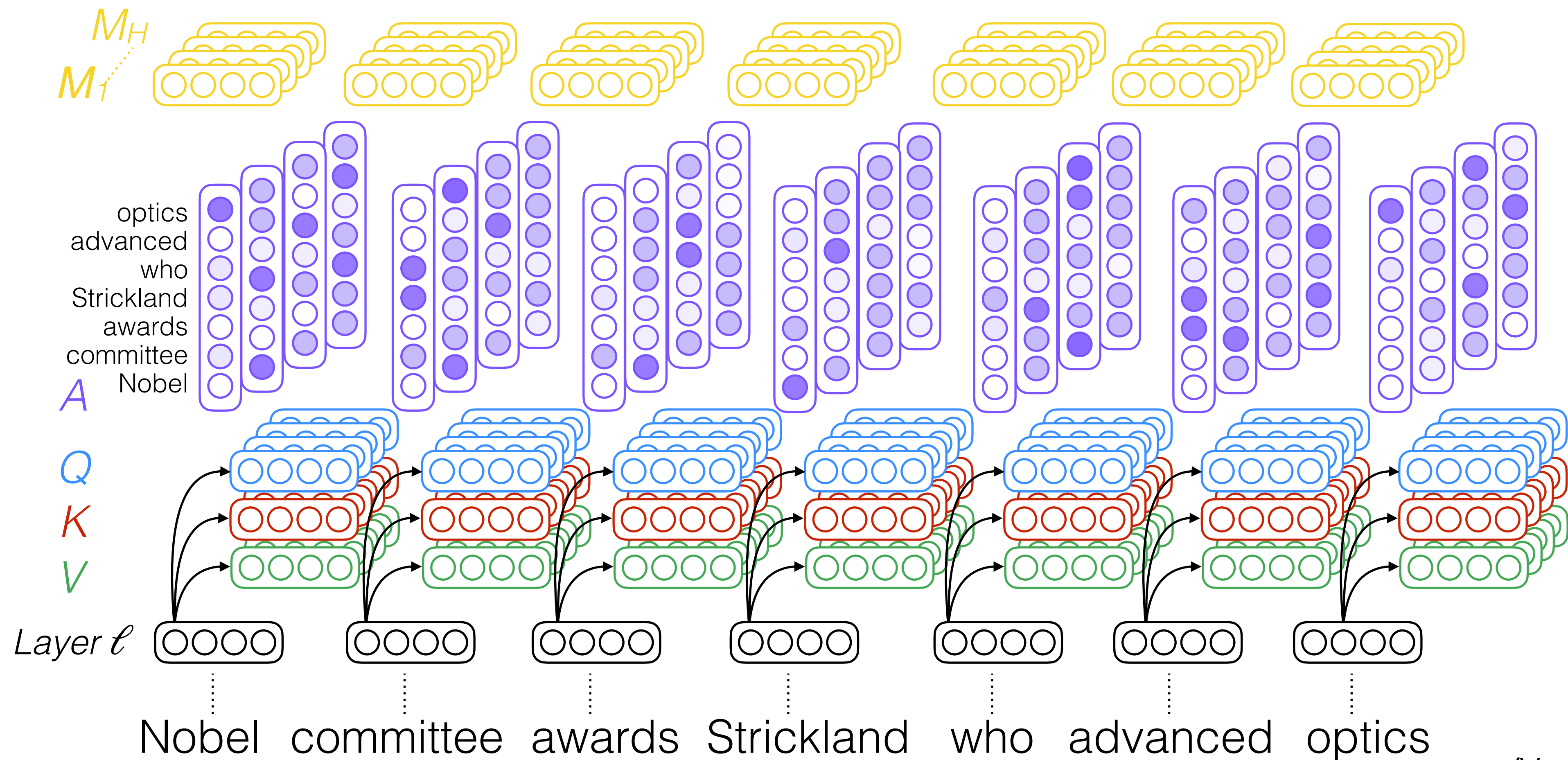
Self-attention (in encoder)



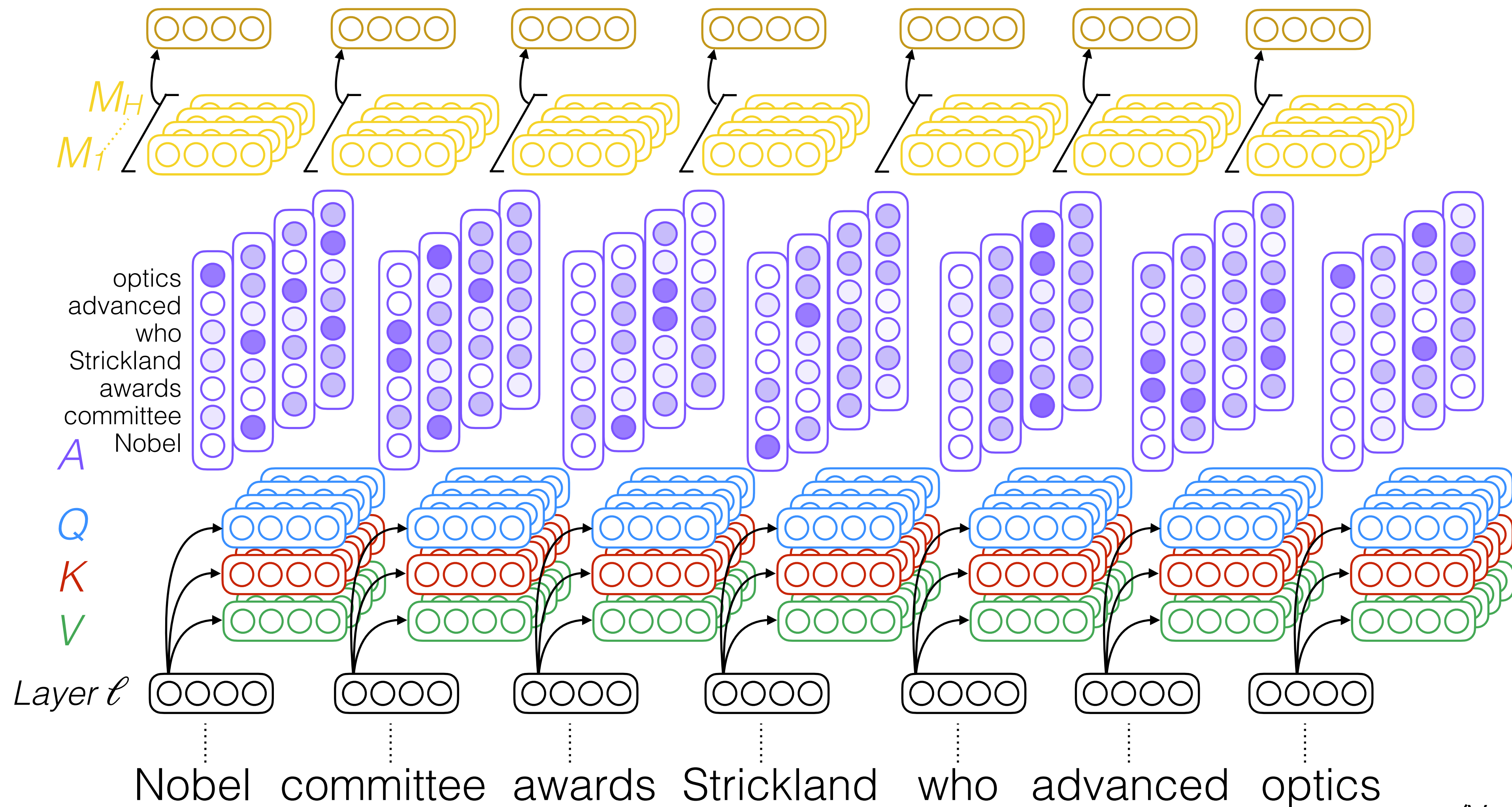
Self-attention (in encoder)



Multi-head self-attention



Multi-head self-attention

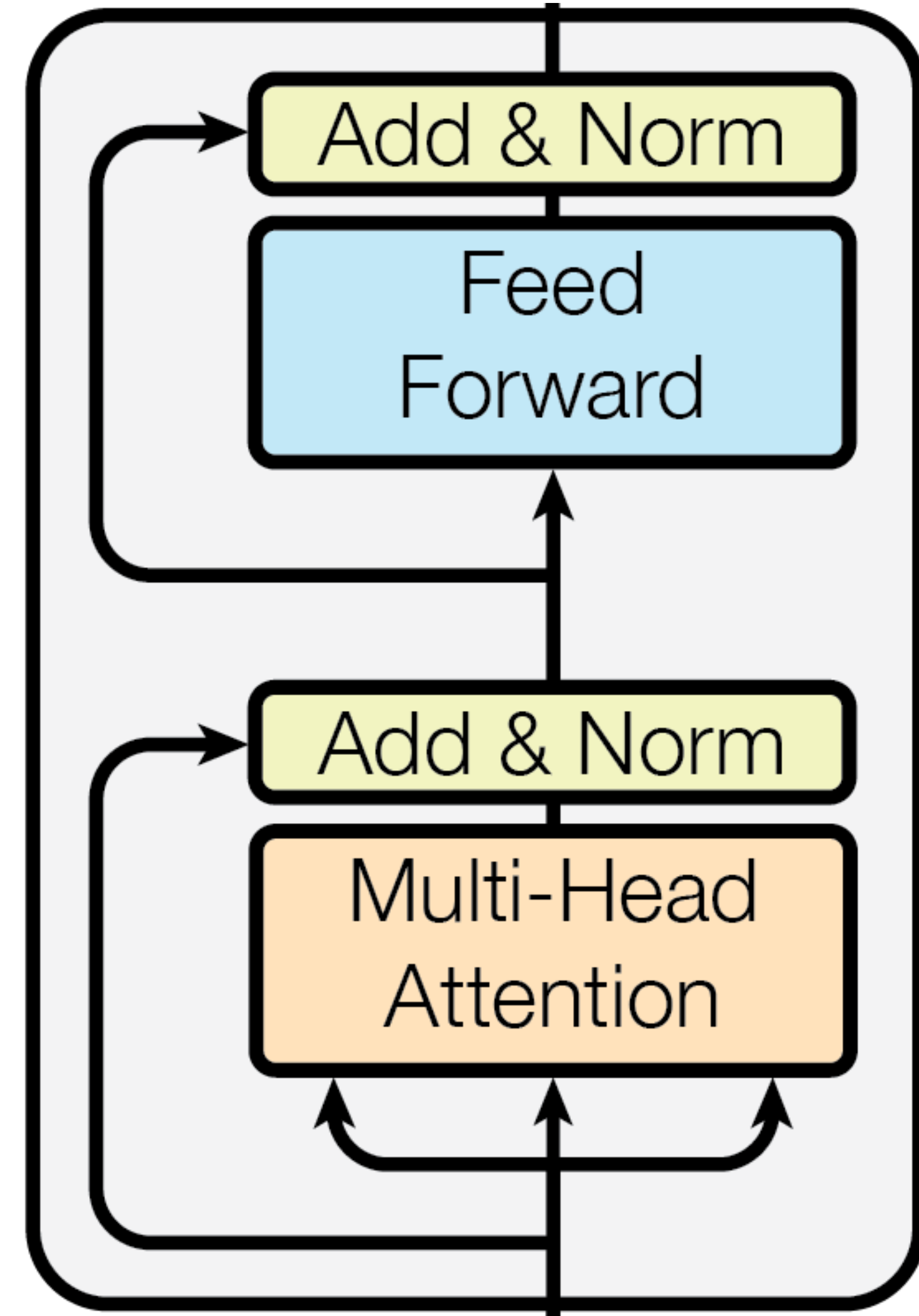


Question

What are two advantages of self-attention over recurrent models?

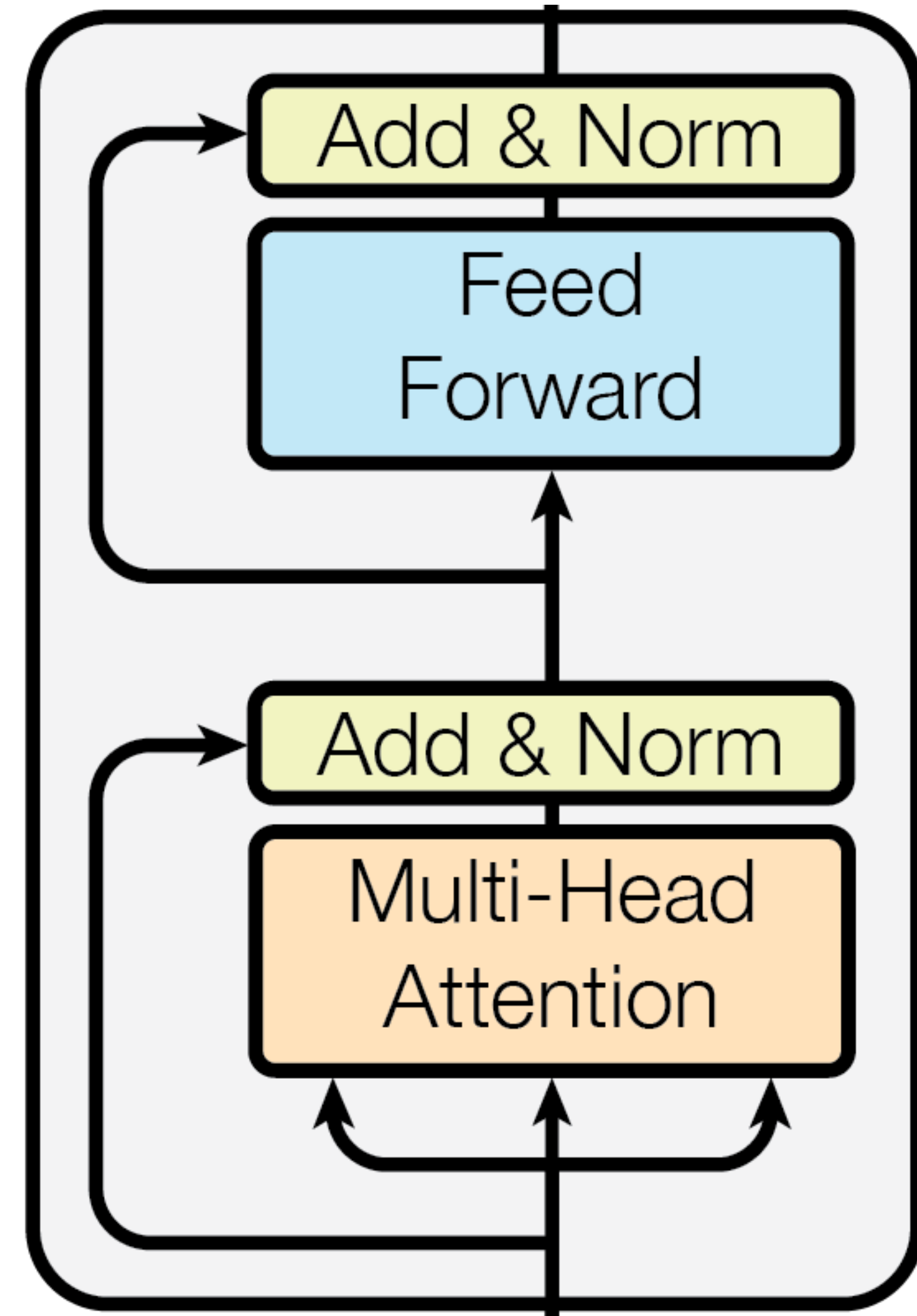
Transformer Block

- Multi-headed attention is the main innovation of the transformer model!



Transformer Block

- Multi-headed attention is the main innovation of the transformer model!
- Each block also composed of:
 - a layer normalisations
 - a feedforward network
 - residual connections



LayerNorm & Residual Connections

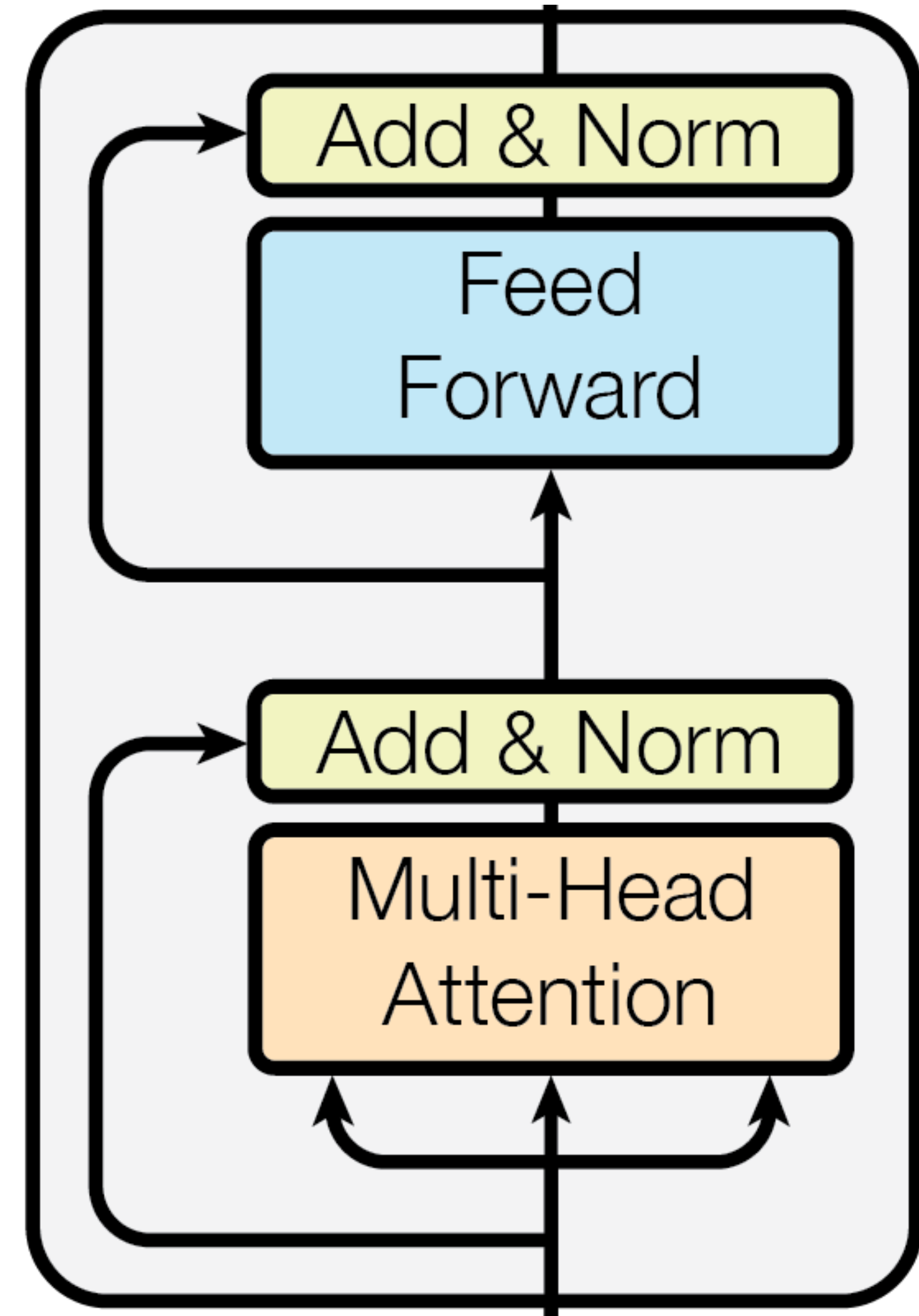
- **Layer Normalisation**

- Normalize the outputs of different modules

$$y = \frac{x - \mathbb{E}[x]}{\sqrt{\text{Var}[x] + \epsilon}} * \gamma + \beta$$

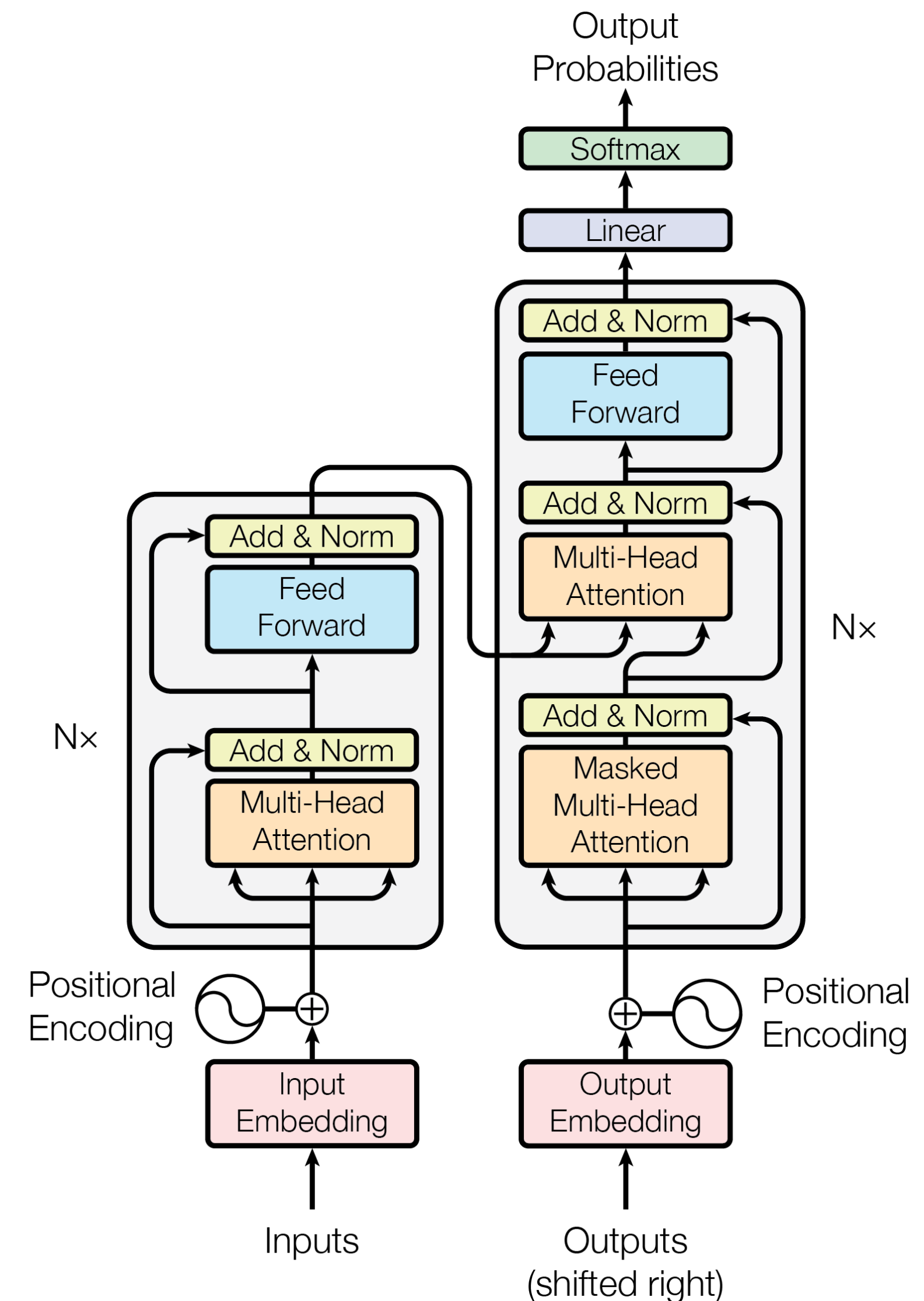
- **Residual Connections**

- Add the input of a module to its output
- $\text{LayerNorm}(x + \text{Sublayer}(x))$



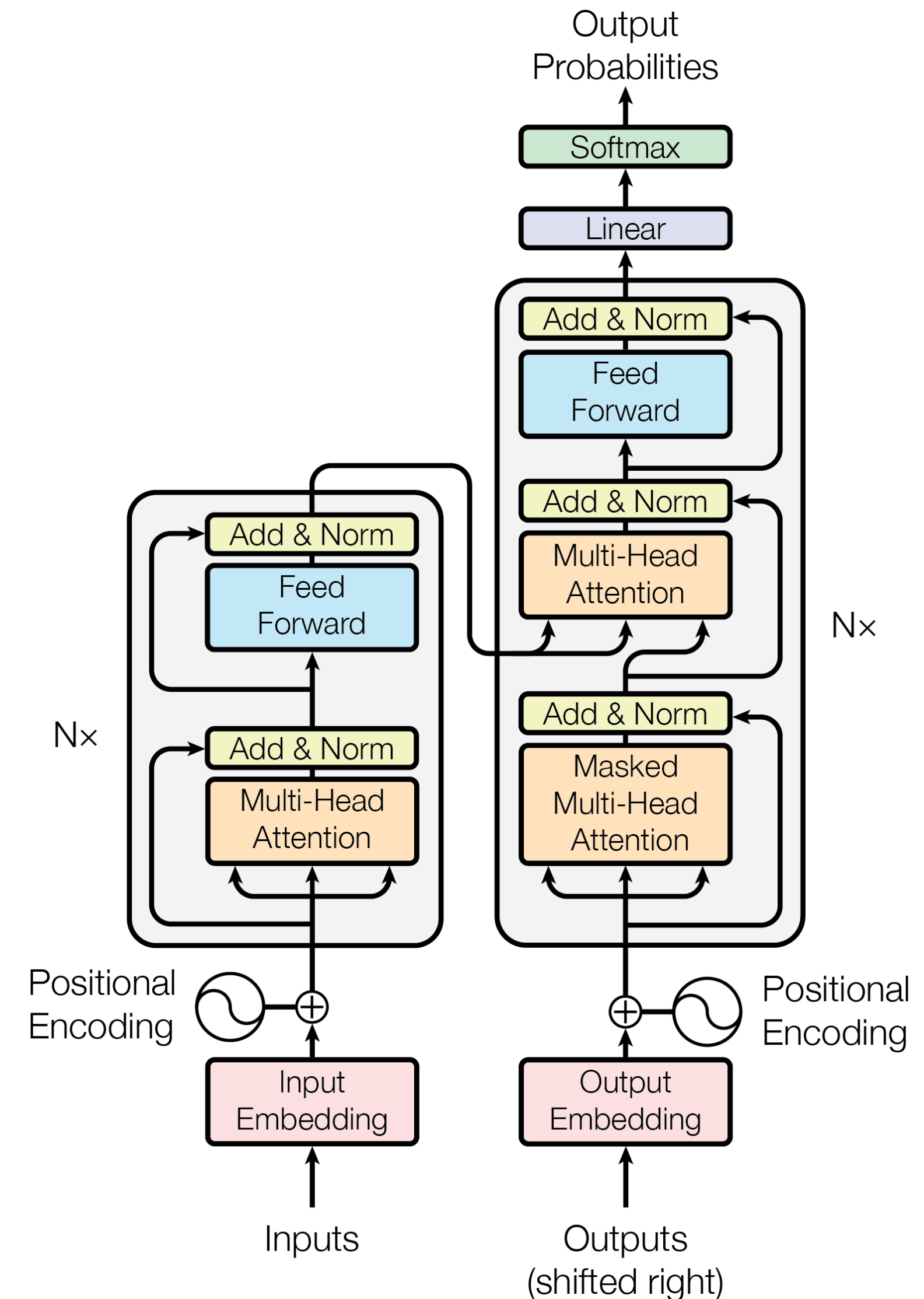
Full Transformer

- Full transformer encoder is multiple cascaded transformer blocks
 - **build up compositional representations of inputs**



Full Transformer

- Full transformer encoder is multiple cascaded transformer blocks
 - **build up compositional representations of inputs**
- Transformer decoder (right) similar to encoder
 - First layer of block is **masked** multi-headed attention
 - Second layer is multi-headed attention over *final-layer* encoder outputs (**cross-attention**)
 - Third layer is feed-forward network



Question

**What is an issue with self-attention
for the decoder?**

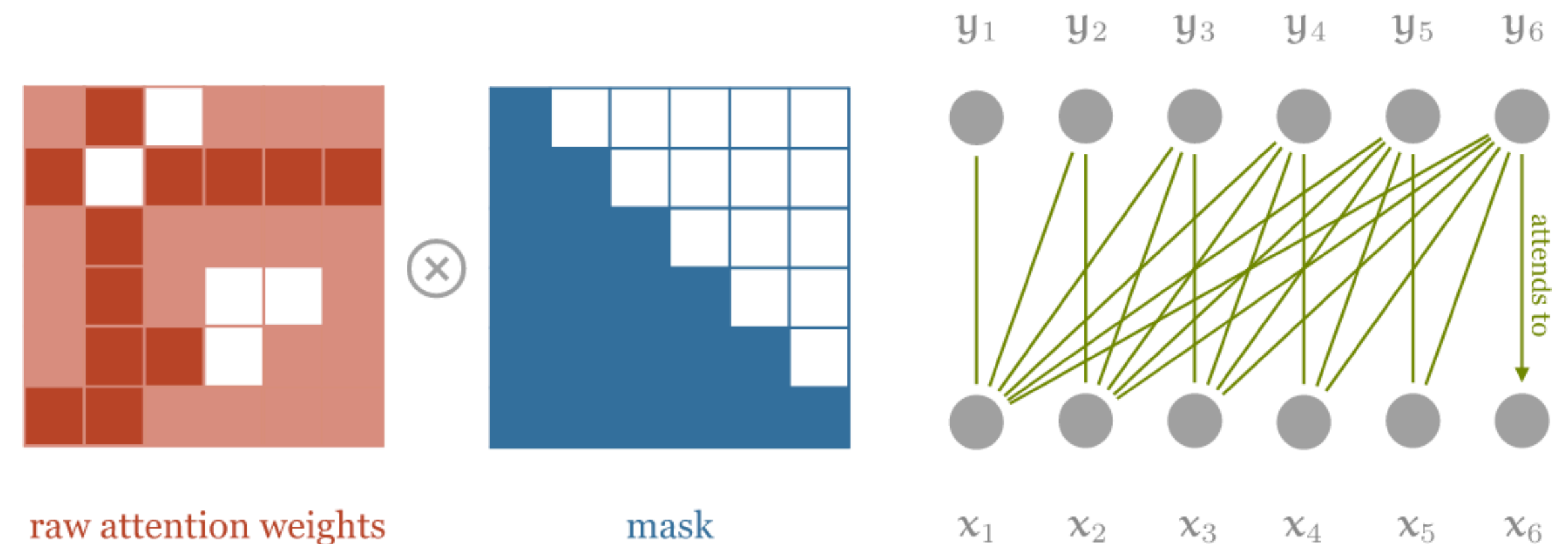
Masked Multi-headed Attention

- Self-attention can attend to any token in the sequence
- For the decoder, **you don't want tokens to attend to future tokens**
 - Decoder used to generate text (i.e., machine translation)

Masked Multi-headed Attention

- Self-attention can attend to any token in the sequence
- For the decoder, **you don't want tokens to attend to future tokens**
 - Decoder used to generate text (i.e., machine translation)

Mask the attention scores of future tokens so their attention = 0

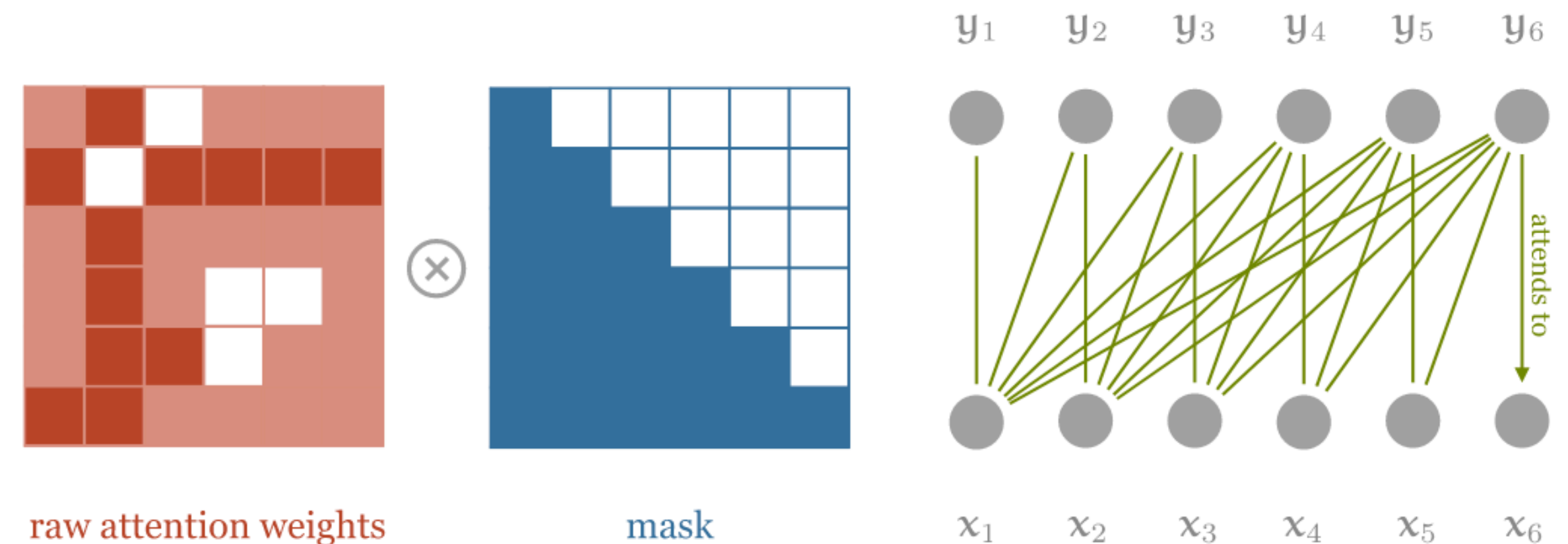


$$a_{st} = \frac{(\mathbf{W}^Q h_s^\ell)^T (\mathbf{W}^K h_t^\ell)}{\sqrt{d}} \quad \rightarrow \quad a_{st} := a_{st} - \infty ; s < t$$

Masked Multi-headed Attention

- Self-attention can attend to any token in the sequence
- For the decoder, **you don't want tokens to attend to future tokens**
 - Decoder used to generate text (i.e., machine translation)

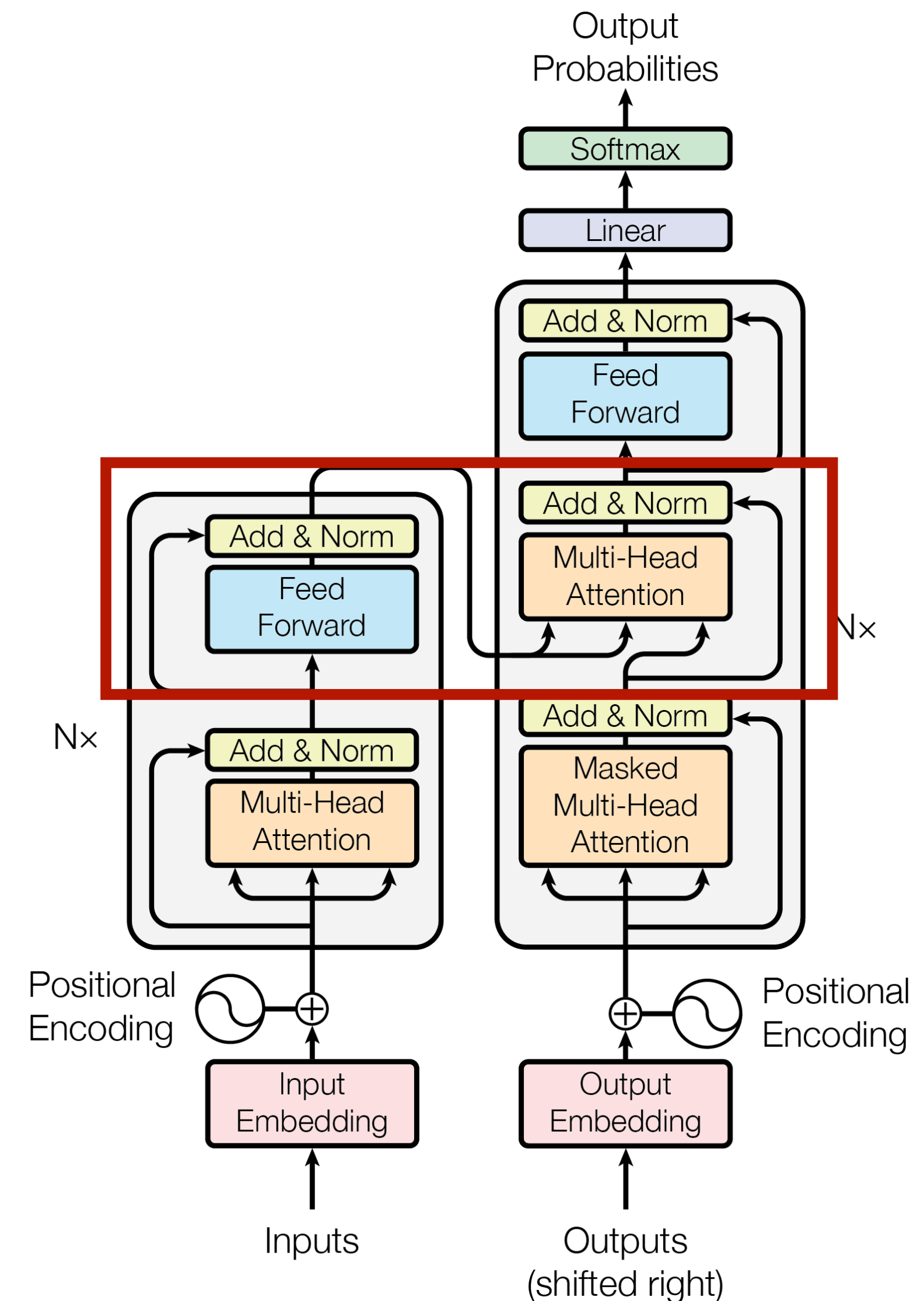
Mask the attention scores of future tokens so their attention = 0



$$a_{st} = \frac{(\mathbf{W}^Q h_s^\ell)^T (\mathbf{W}^K h_t^\ell)}{\sqrt{d}} \quad \rightarrow \quad a_{st} := a_{st} - \infty ; s < t \quad \rightarrow \quad \alpha_{st} = \frac{e^{a_{st}}}{\sum_j e^{a_{sj}}} = 0$$

Cross-attention

- **Cross attention** is the same classical attention as in the RNN encoder-decoder model
- The query to the attention function is the output of the masked multi-headed attention in the decoder (i.e., a decoder state)
- The keys and values are the output of the **final** encoder transformer
- Once again, a representation from the decoder is used to **attend** to the encoder outputs



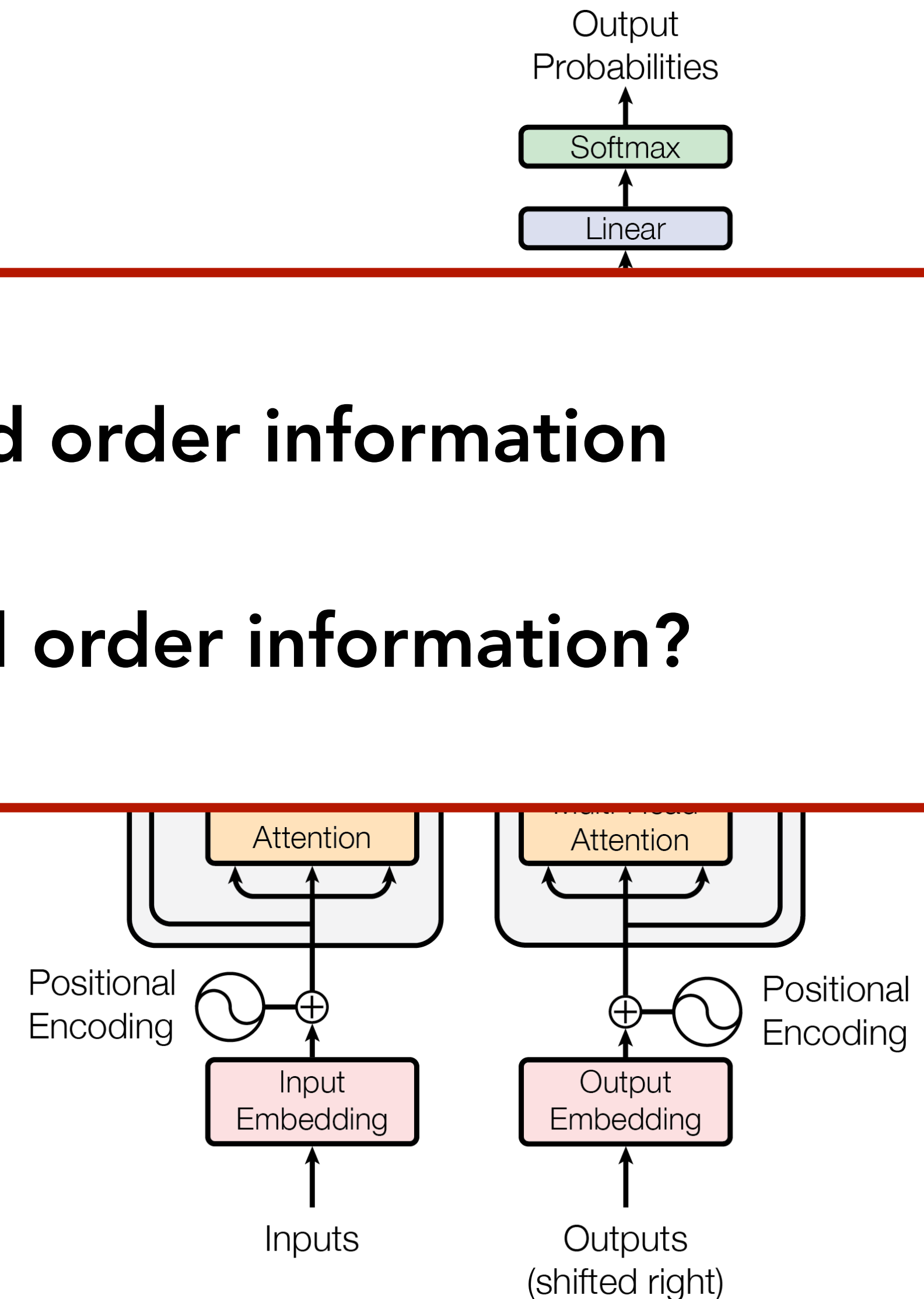
Full Transformer

- Full transformer encoder is multiple cascaded transformer blocks

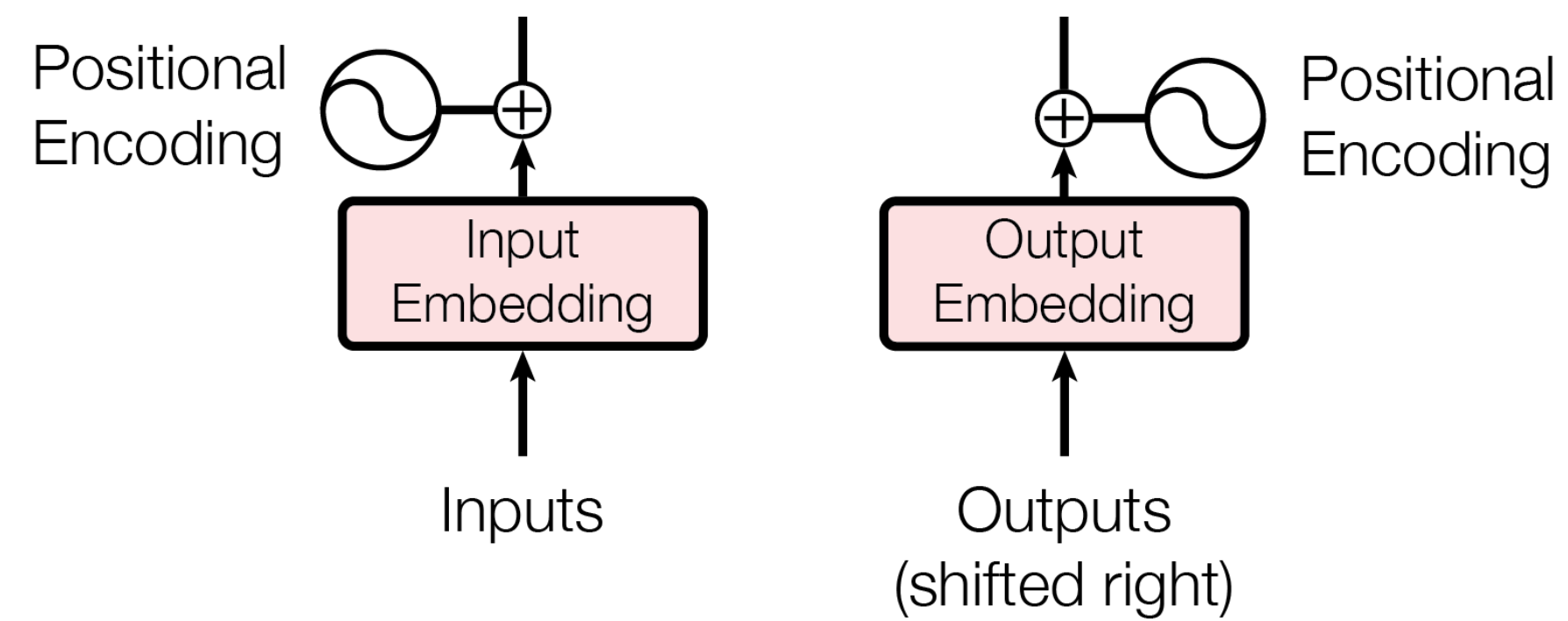
Recurrent models provided word order information

Does self-attention provide word order information?

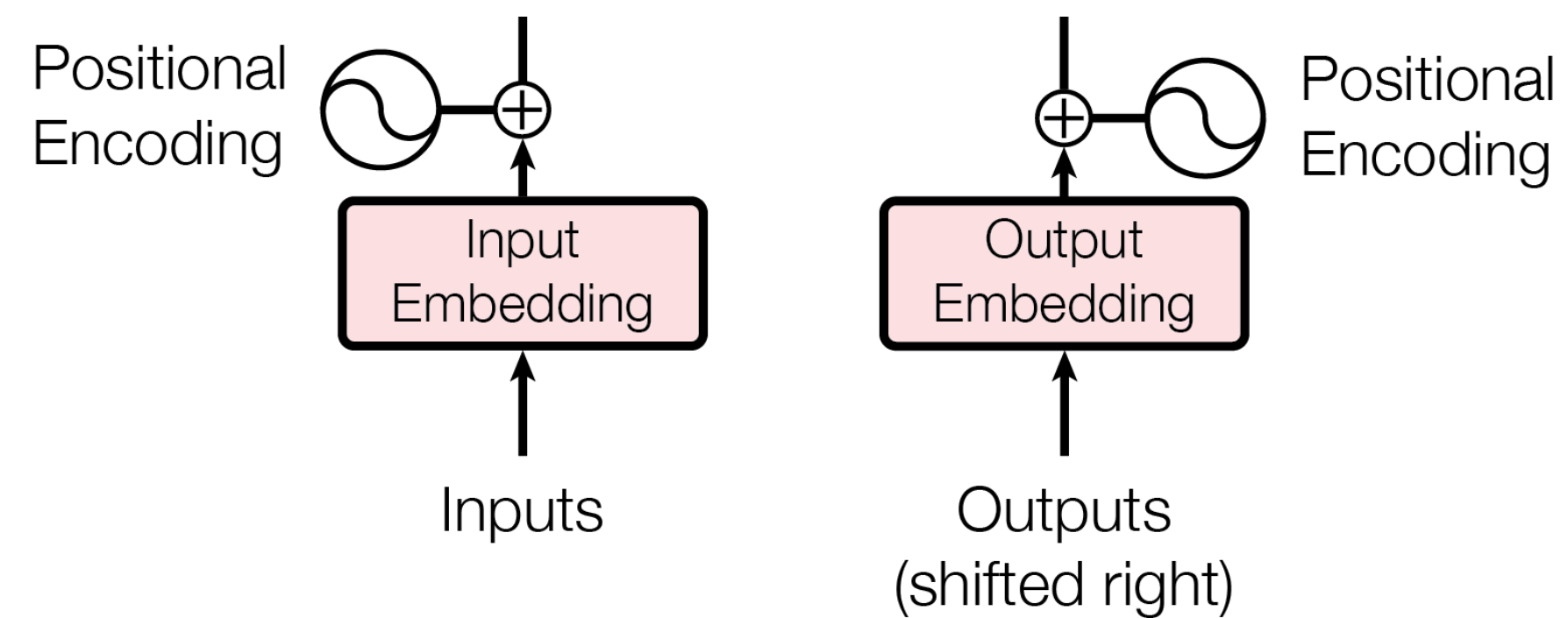
- Second layer is multi-headed attention over encoder outputs (**cross-attention**)
- Third layer is feed-forward network



Position Embeddings

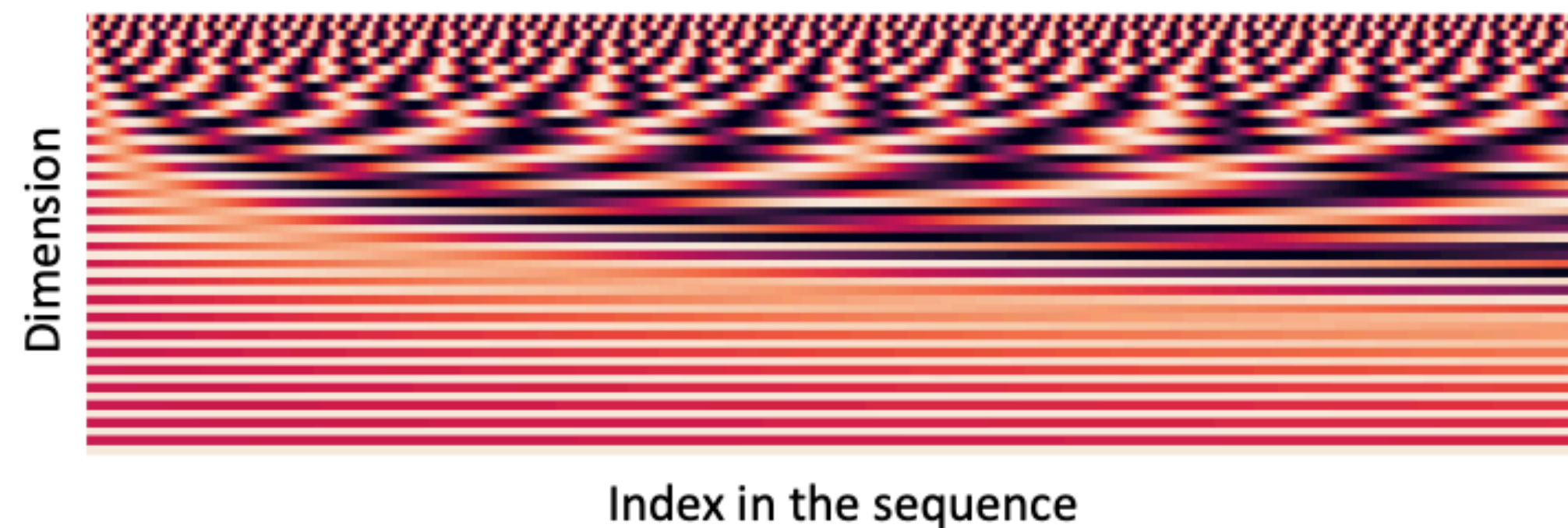


Position Embeddings

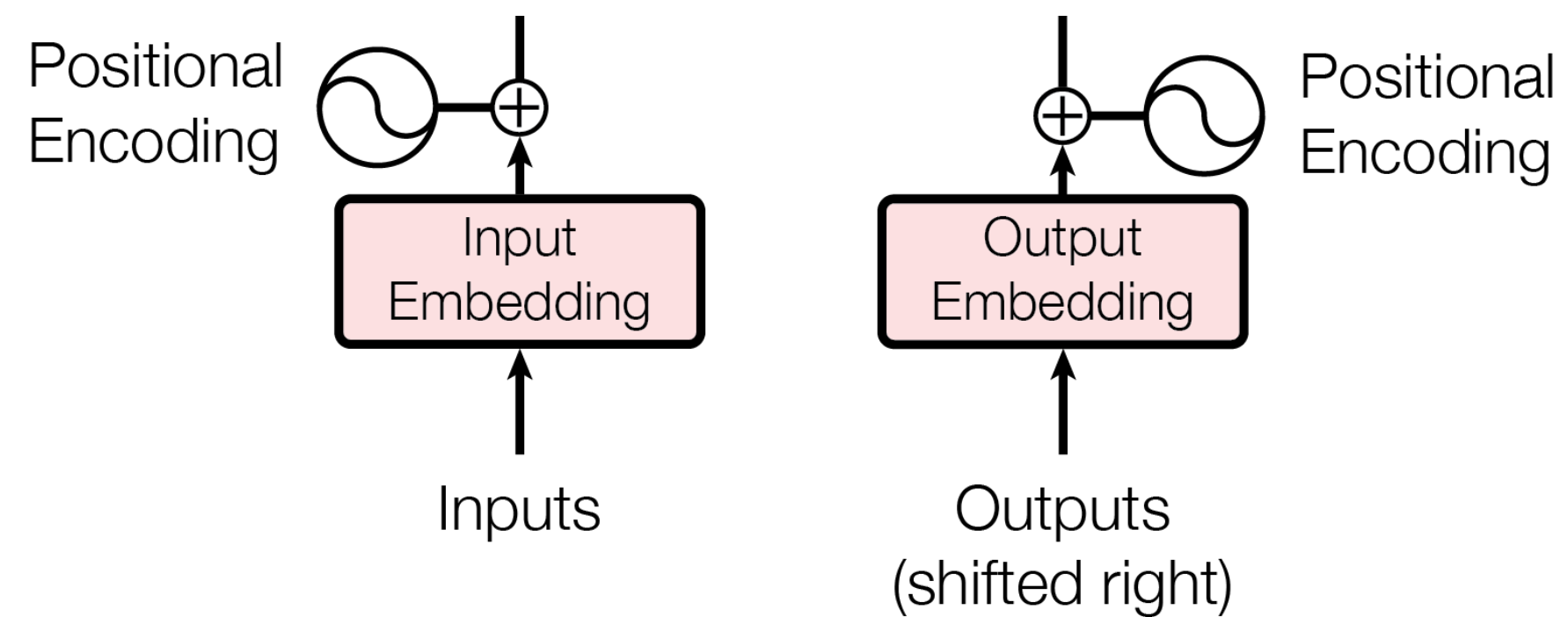


- Early position embeddings encoded a sinusoid function that was offset by a phase shift proportional to sequence position

$$p_i = \begin{pmatrix} \sin(i/10000^{2*1/d}) \\ \cos(i/10000^{2*1/d}) \\ \vdots \\ \sin(i/10000^{2*\frac{d}{2}/d}) \\ \cos(i/10000^{2*\frac{d}{2}/d}) \end{pmatrix}$$

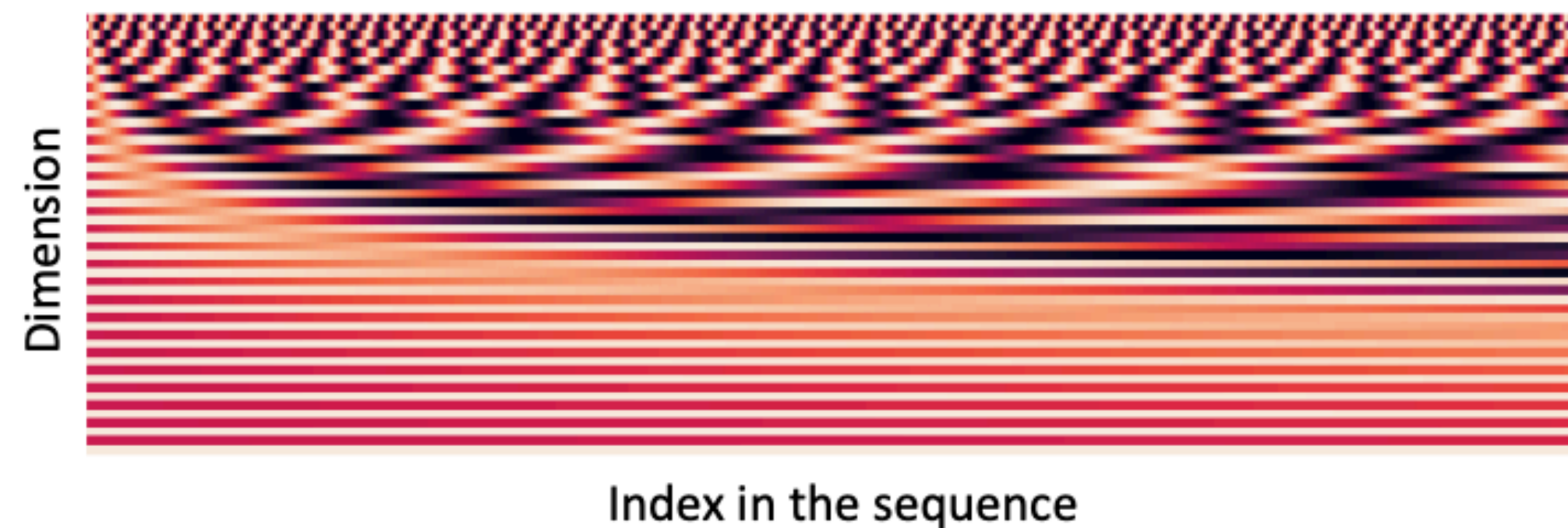


Position Embeddings



- Early position embeddings encoded a sinusoid function that was offset by a phase shift proportional to sequence position
- **In practice, easiest is to learn position embeddings from scratch**

$$p_i = \begin{pmatrix} \sin(i/10000^{2*1/d}) \\ \cos(i/10000^{2*1/d}) \\ \vdots \\ \sin(i/10000^{2*\frac{d}{2}/d}) \\ \cos(i/10000^{2*\frac{d}{2}/d}) \end{pmatrix}$$



Question

What might be a disadvantage of using learned position embeddings?

Poor generalisation to sequences longer than the maximum position embedding you have learned

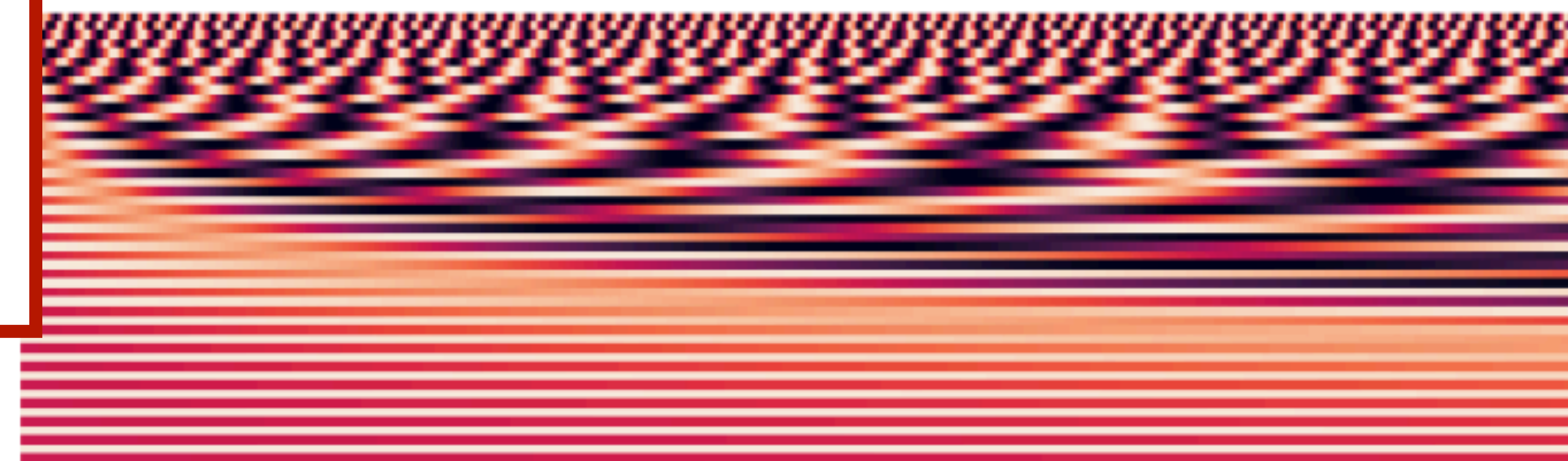
Position Embeddings

Lots of potential for new methods that generalise to longer sequences

Position embeddings remain an active area of research

- Early position embeddings encoded a sinusoid function that was offset by a phase shift proportional to sequence position
- **In practice, easiest is to learn position embeddings from scratch**

$$\begin{pmatrix} \sin(i/10000^{2*\frac{u}{2}/d}) \\ \cos(i/10000^{2*\frac{d}{2}/d}) \end{pmatrix}$$



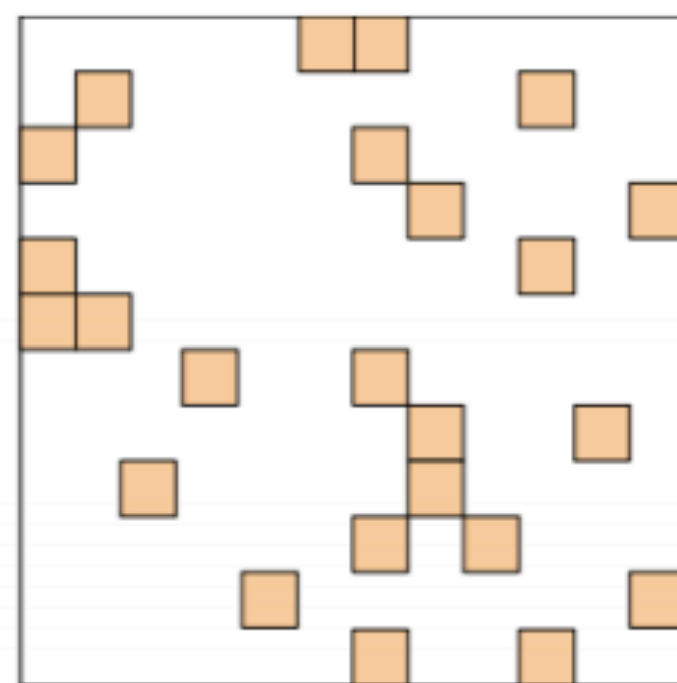
Index in the sequence

Performance: Machine Translation

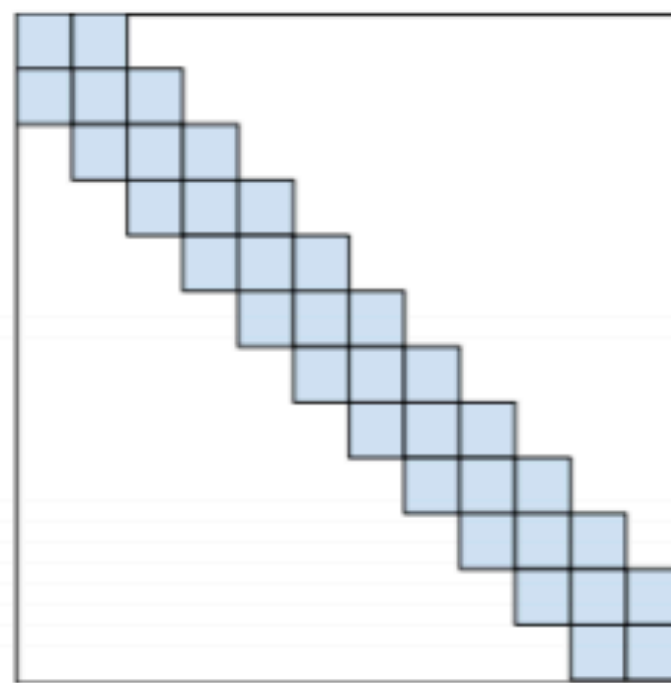
| Model | BLEU | | Training Cost (FLOPs) | |
|---------------------------------|-------------|--------------|---------------------------------------|---------------------|
| | EN-DE | EN-FR | EN-DE | EN-FR |
| ByteNet [15] | 23.75 | | | |
| Deep-Att + PosUnk [32] | | 39.2 | | $1.0 \cdot 10^{20}$ |
| GNMT + RL [31] | 24.6 | 39.92 | $2.3 \cdot 10^{19}$ | $1.4 \cdot 10^{20}$ |
| ConvS2S [8] | 25.16 | 40.46 | $9.6 \cdot 10^{18}$ | $1.5 \cdot 10^{20}$ |
| MoE [26] | 26.03 | 40.56 | $2.0 \cdot 10^{19}$ | $1.2 \cdot 10^{20}$ |
| Deep-Att + PosUnk Ensemble [32] | | 40.4 | | $8.0 \cdot 10^{20}$ |
| GNMT + RL Ensemble [31] | 26.30 | 41.16 | $1.8 \cdot 10^{20}$ | $1.1 \cdot 10^{21}$ |
| ConvS2S Ensemble [8] | 26.36 | 41.29 | $7.7 \cdot 10^{19}$ | $1.2 \cdot 10^{21}$ |
| Transformer (base model) | 27.3 | 38.1 | $3.3 \cdot 10^{18}$ | |
| Transformer (big) | 28.4 | 41.0 | $2.3 \cdot 10^{19}$ | |

Question

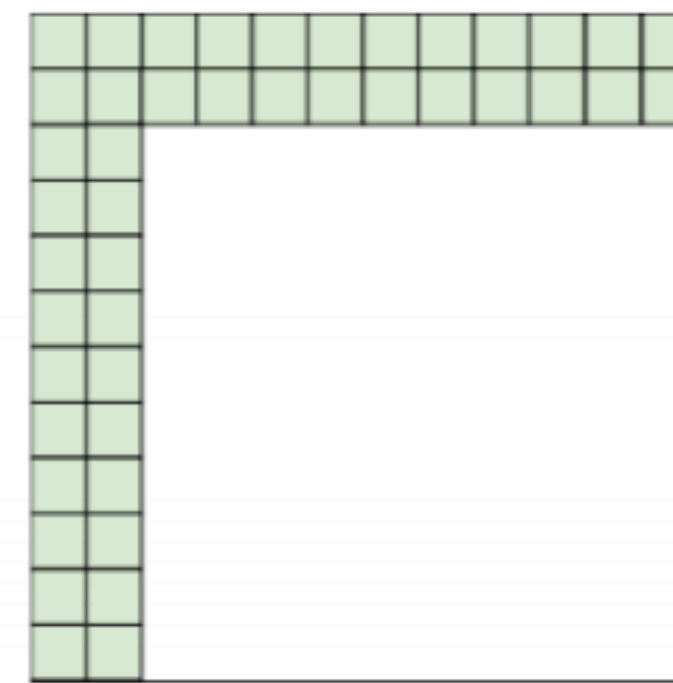
What could be a disadvantage of transformers over RNNs?



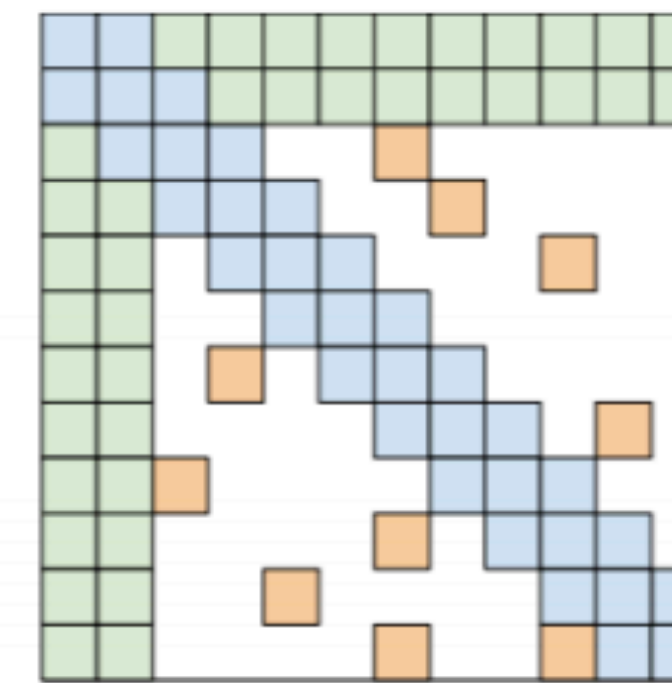
(a) Random attention



(b) Window attention



(c) Global Attention



(d) BIGBIRD

Other Resources of Interest

- The Annotated Transformer
 - <https://nlp.seas.harvard.edu/2018/04/03/attention.html>
- The Illustrated Transformer
 - <https://jalammar.github.io/illustrated-transformer/>
- Only basics presented here today! Many modifications to initial transformers exist

Recap

- **Temporal Bottleneck:** **Vanishing gradients** stop many RNN architectures from learning **long-range dependencies**
- **Parallelisation Bottleneck:** RNN states depend on previous time step hidden state, so must be **computed in series**
- **Attention:** Direct connections between output states and inputs (solves temporal bottleneck)
- **Self-Attention:** Remove recurrence, allowing parallel computation
- Modern **Transformers** use attention, but require position embeddings to capture sequence order

References

- Bahdanau, D., Cho, K., & Bengio, Y. (2014). Neural Machine Translation by Jointly Learning to Align and Translate. *CoRR*, *abs/1409.0473*.
- Vaswani, A., Shazeer, N.M., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., & Polosukhin, I. (2017). Attention is All you Need. *ArXiv*, *abs/1706.03762*.
- Wu et al., Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation. *arxiv* 2016