

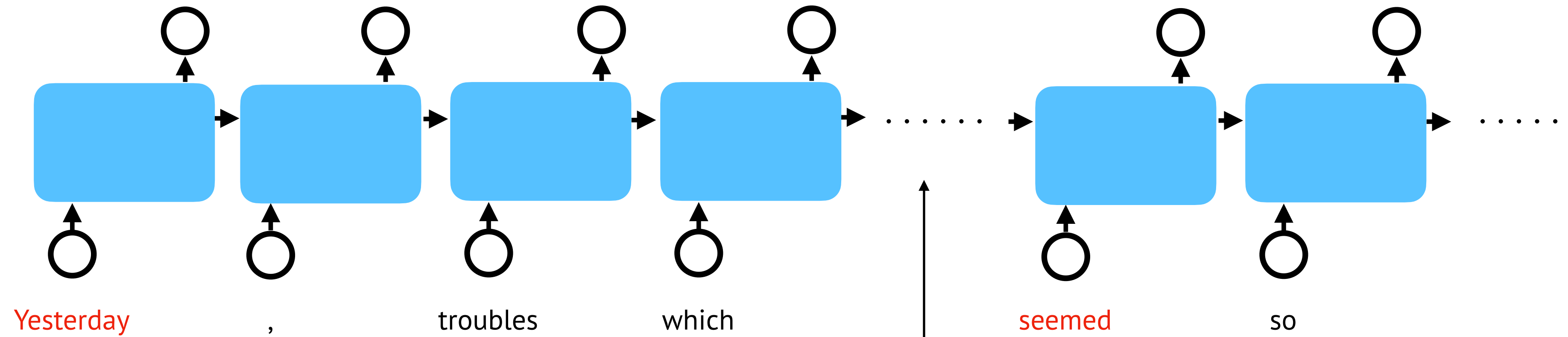
# Transformers

COMP7607 — Lecture 4

Lingpeng Kong

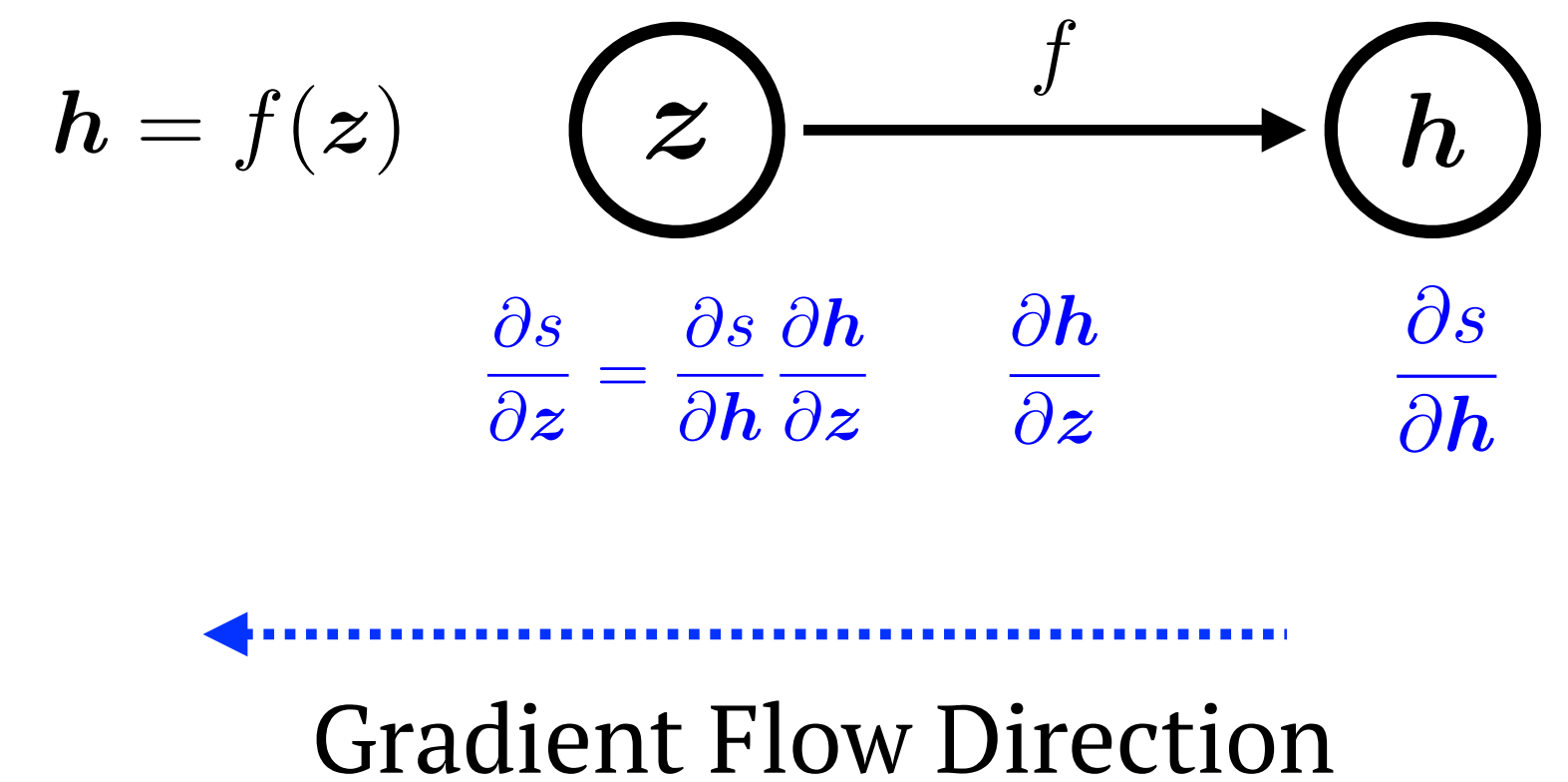
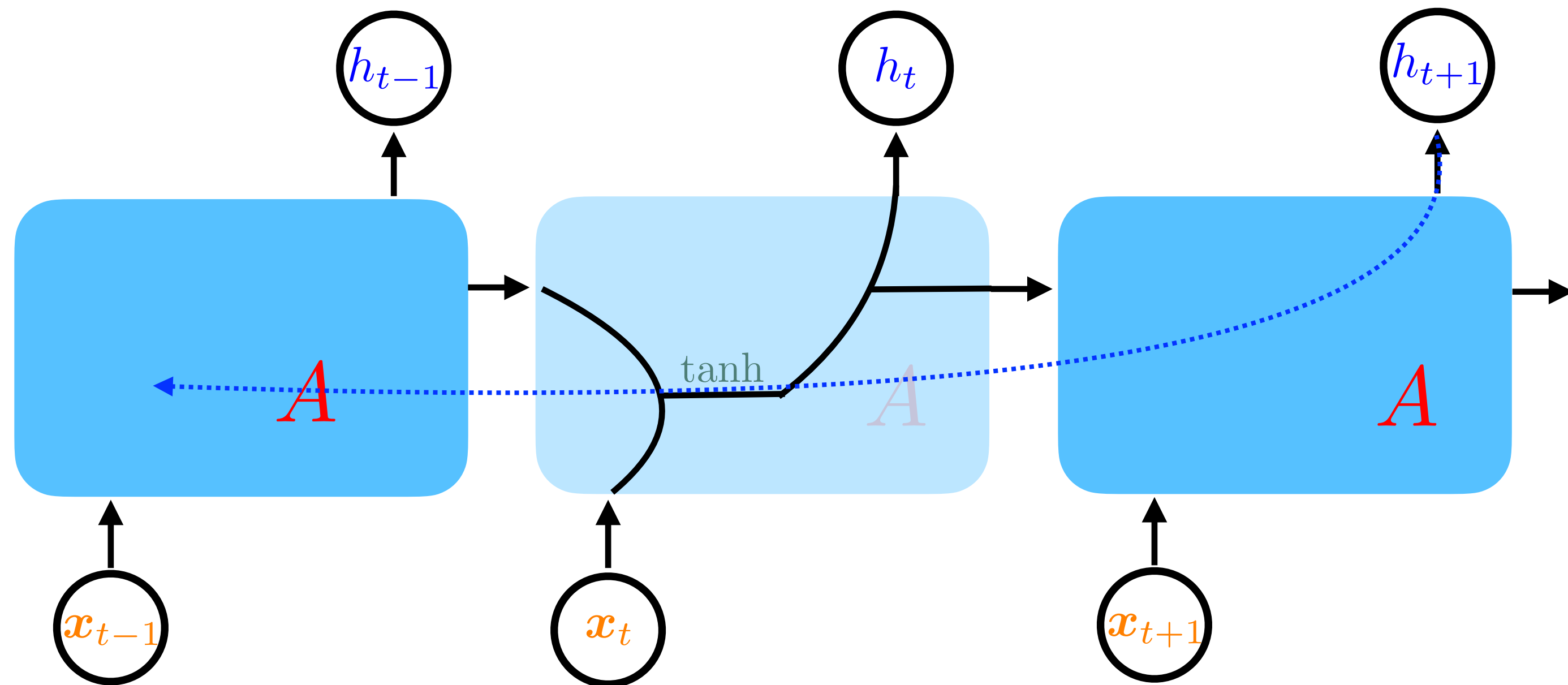
Department of Computer Science, The University of Hong Kong

# Recurrent Neural Network



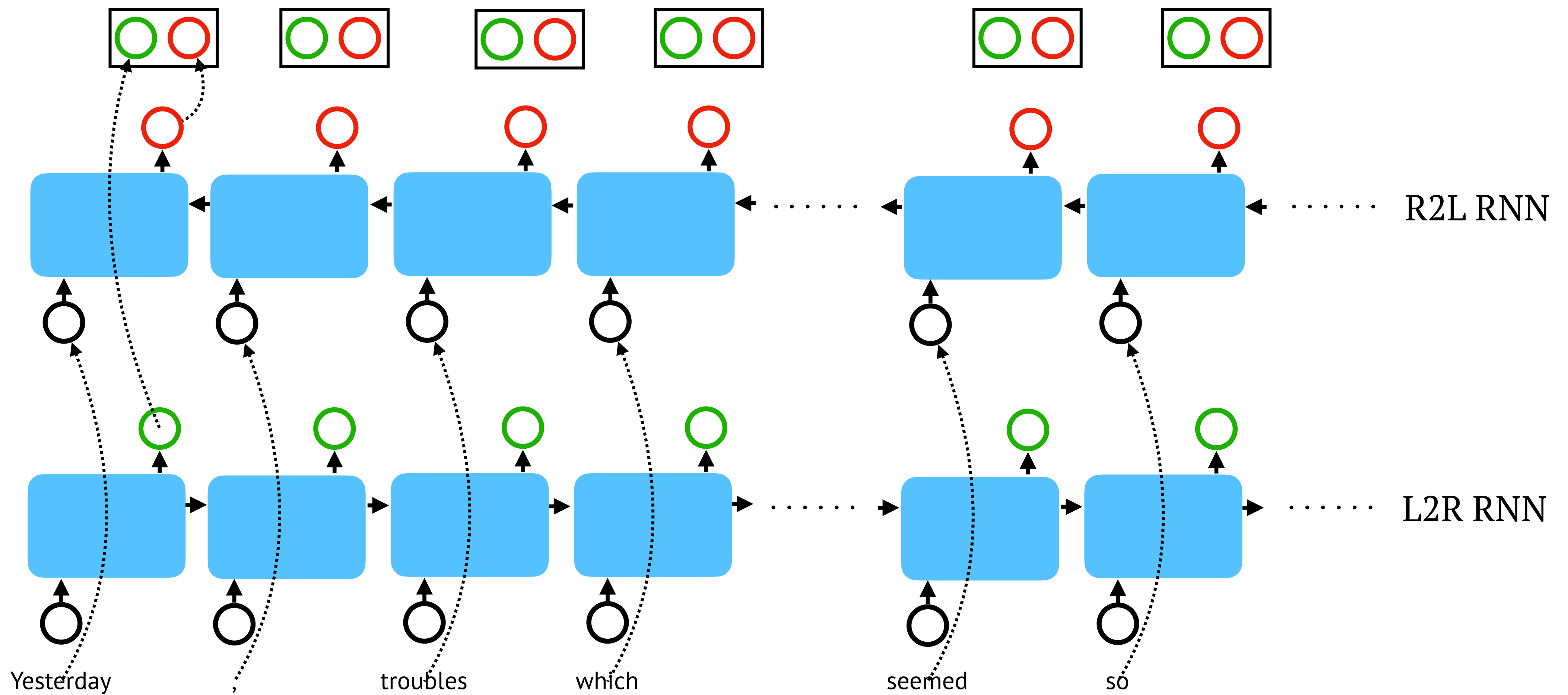
Possibly many steps  $[O(N)]$  steps before “yesterday” and “seemed” interact.

# Vanishing Gradient in RNNs

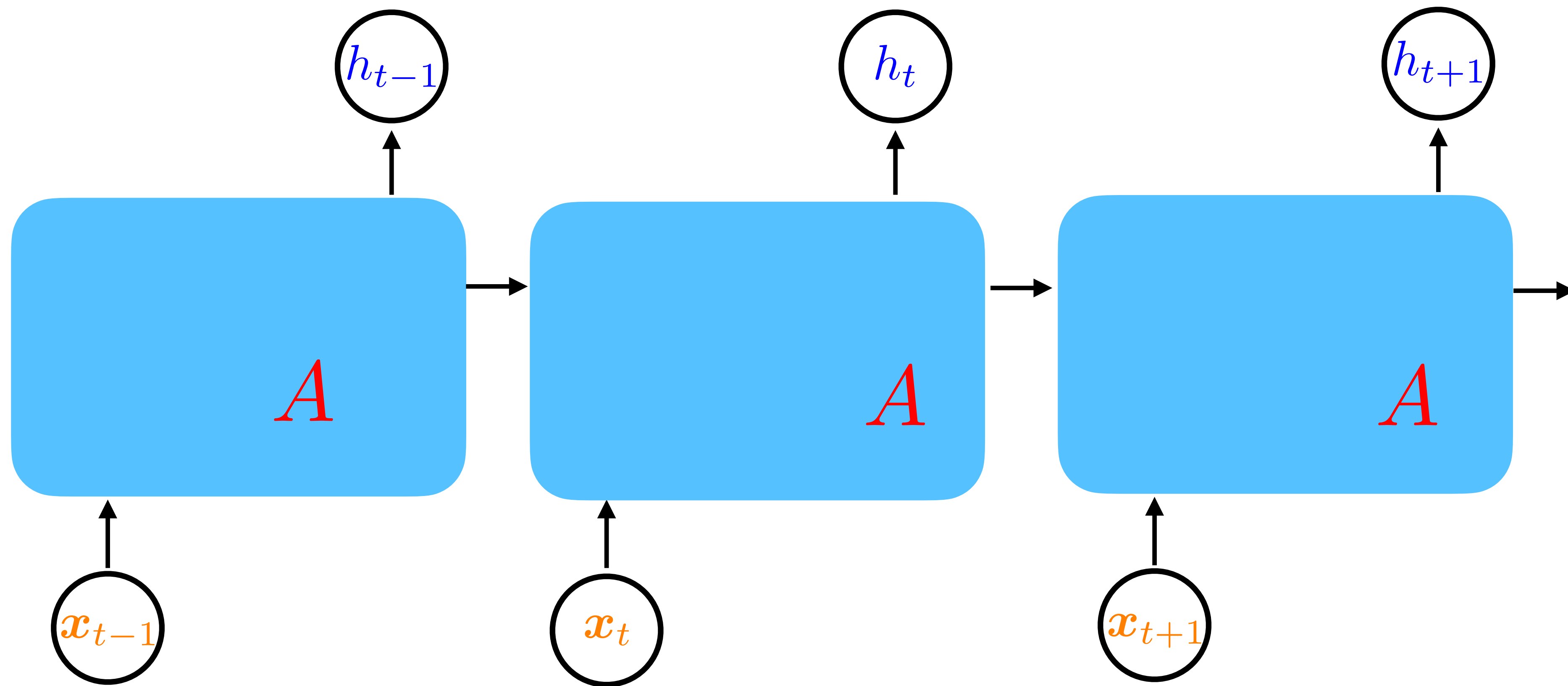


In general, the longer the path, the smaller the gradient signal.

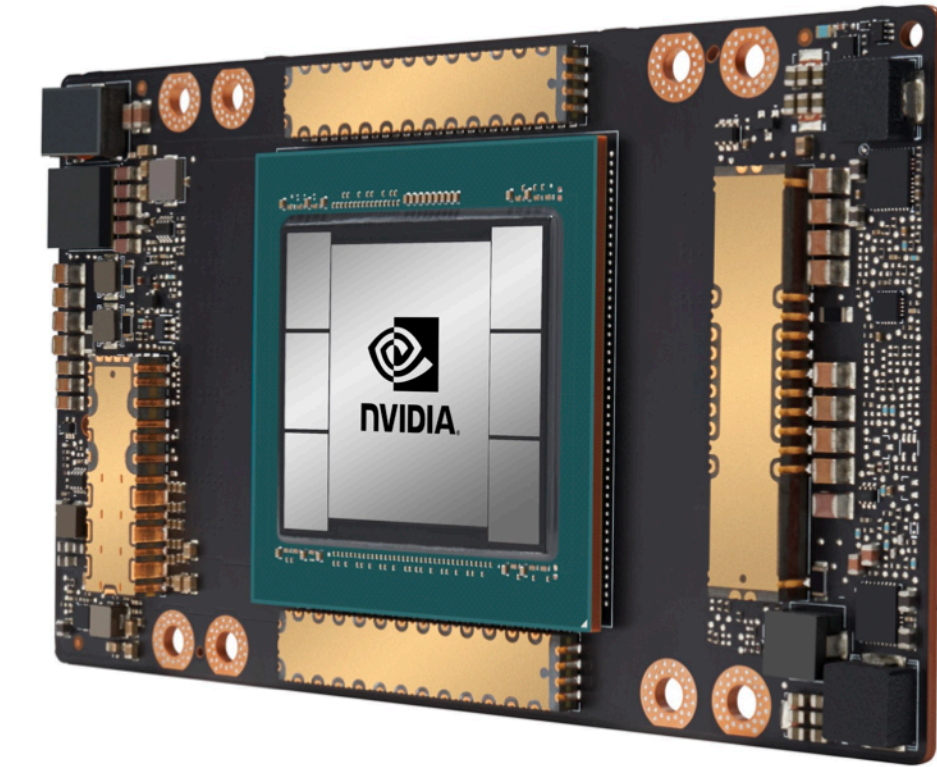
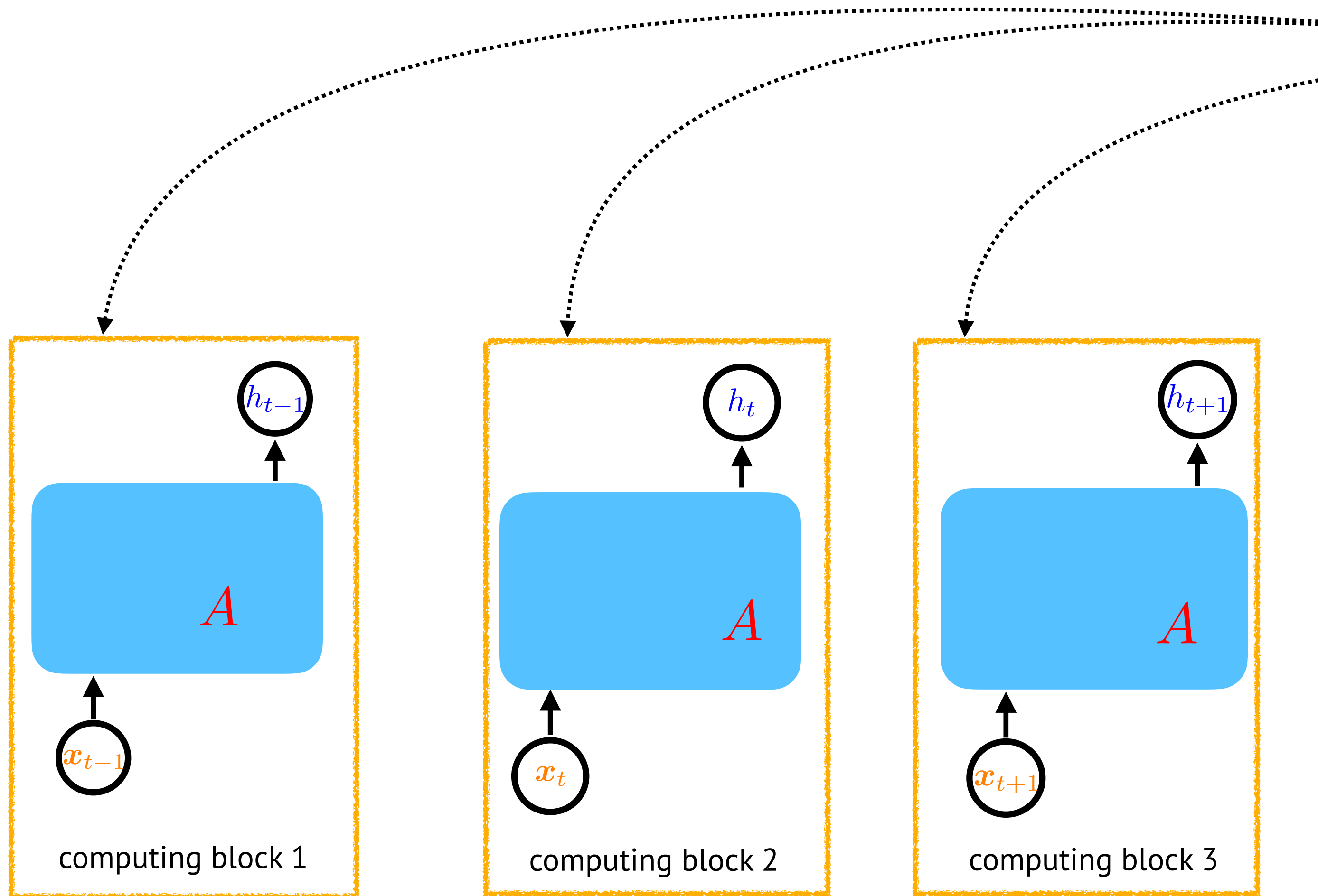
# Bidirectional Recurrent Neural Network



# Sequential Computation



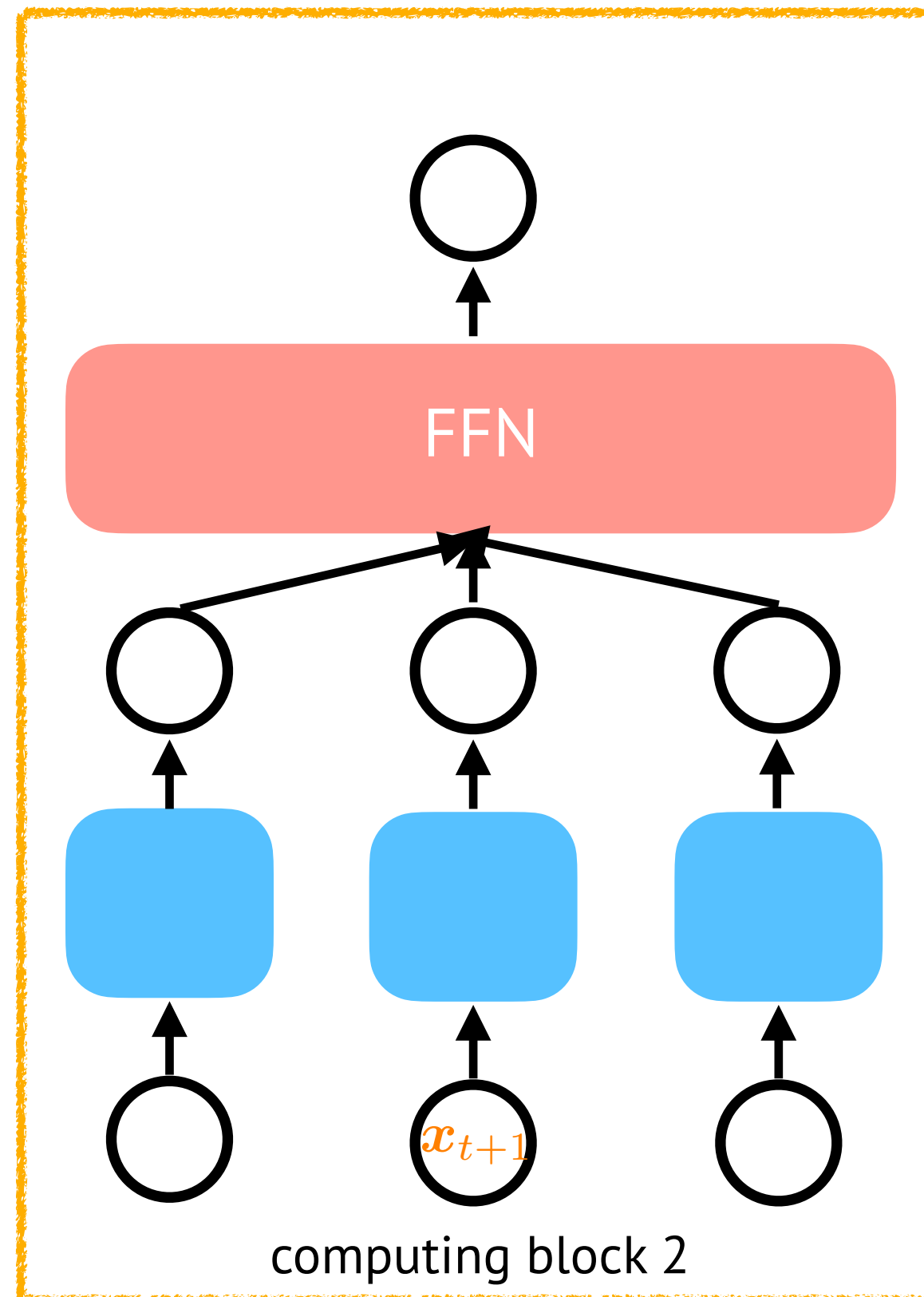
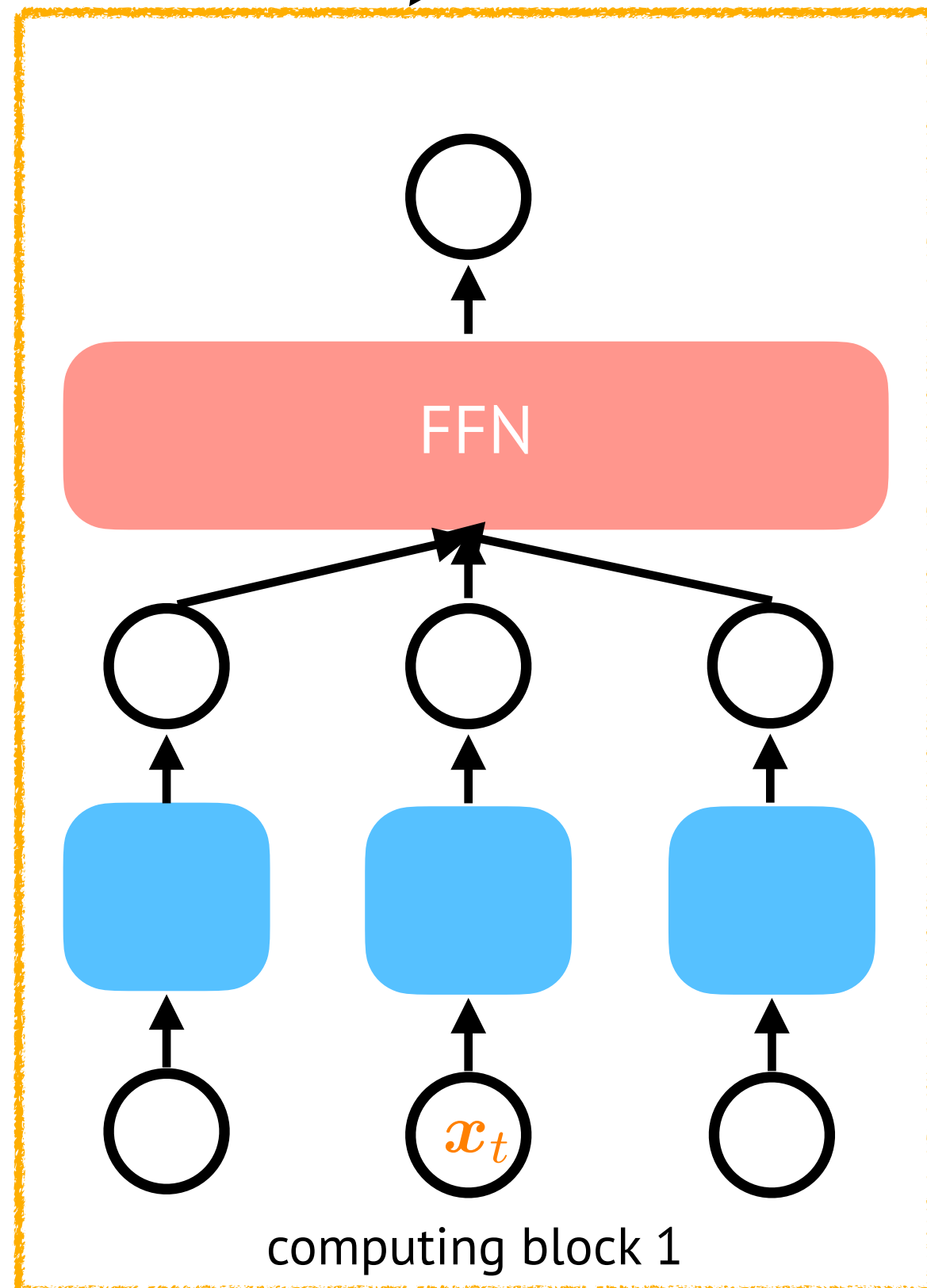
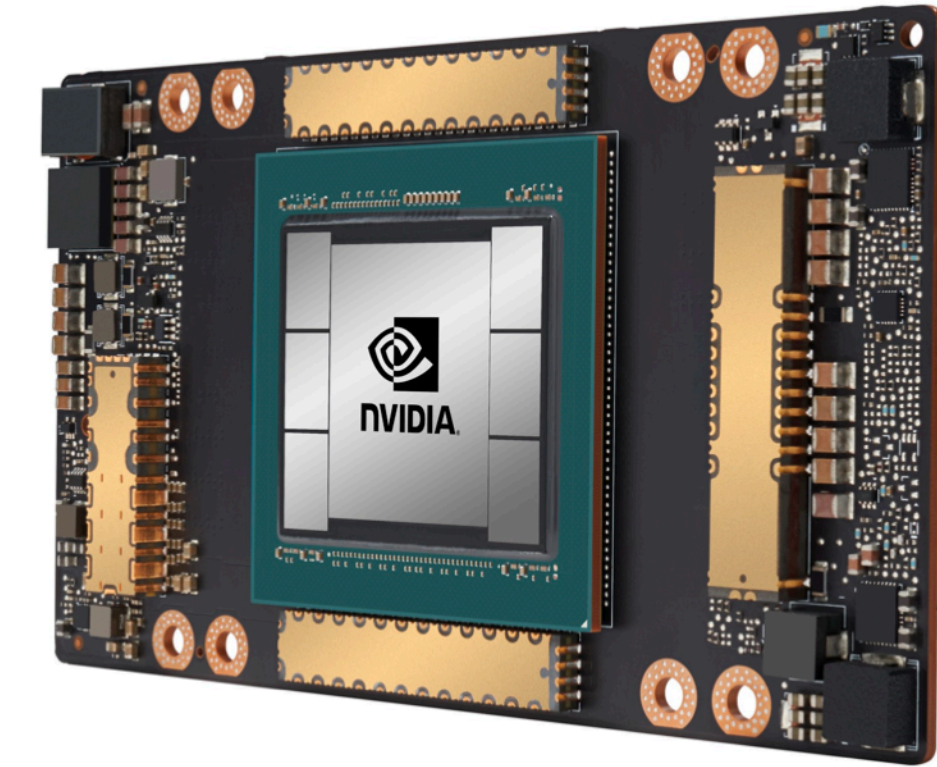
# Parallel Computing?



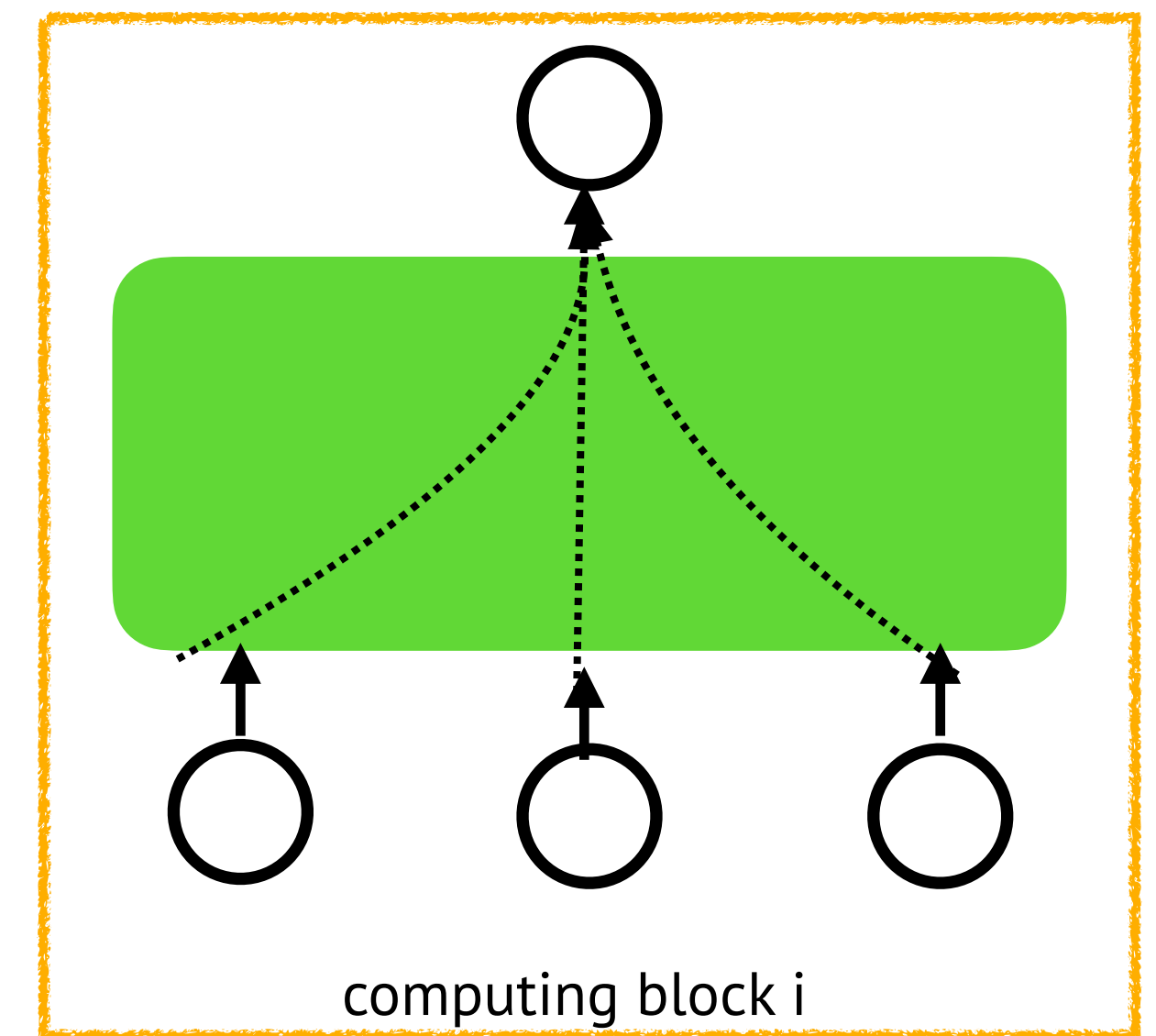
GPU loves parallel computing blocks!

.....

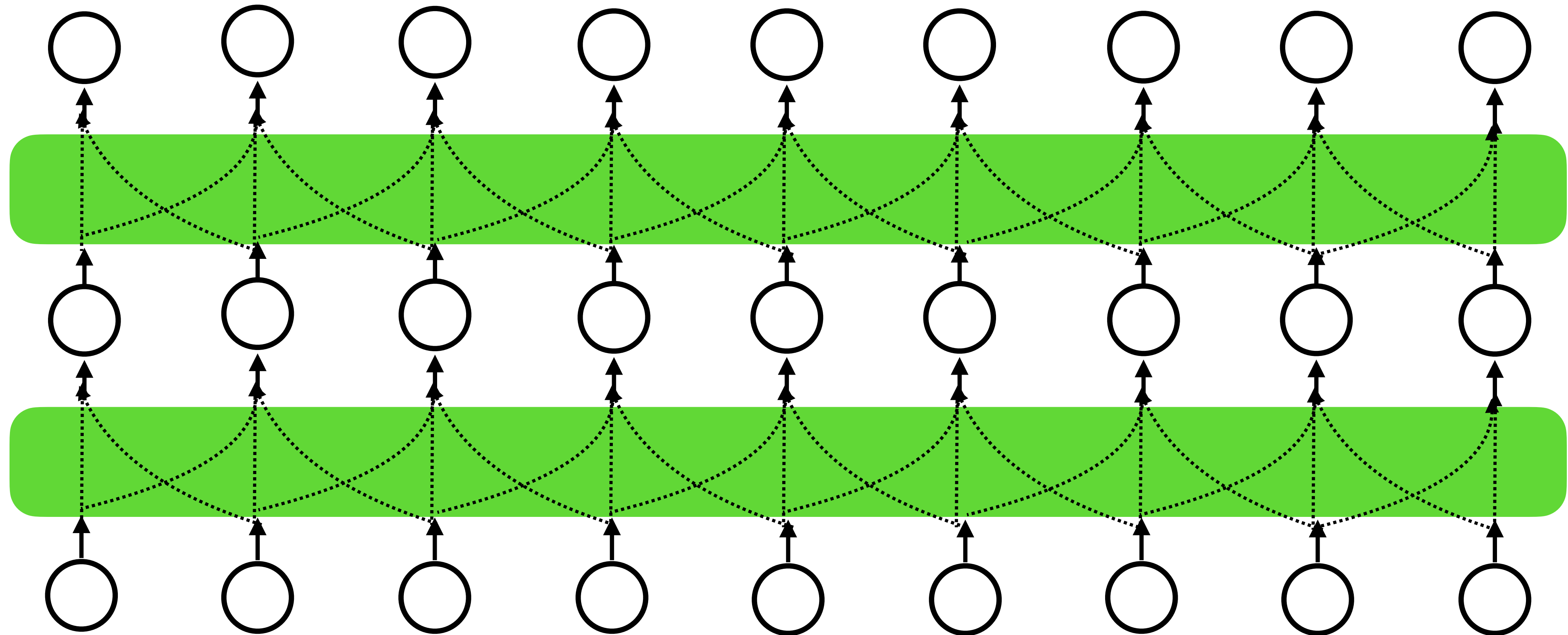
# Parallel Computing?



.....

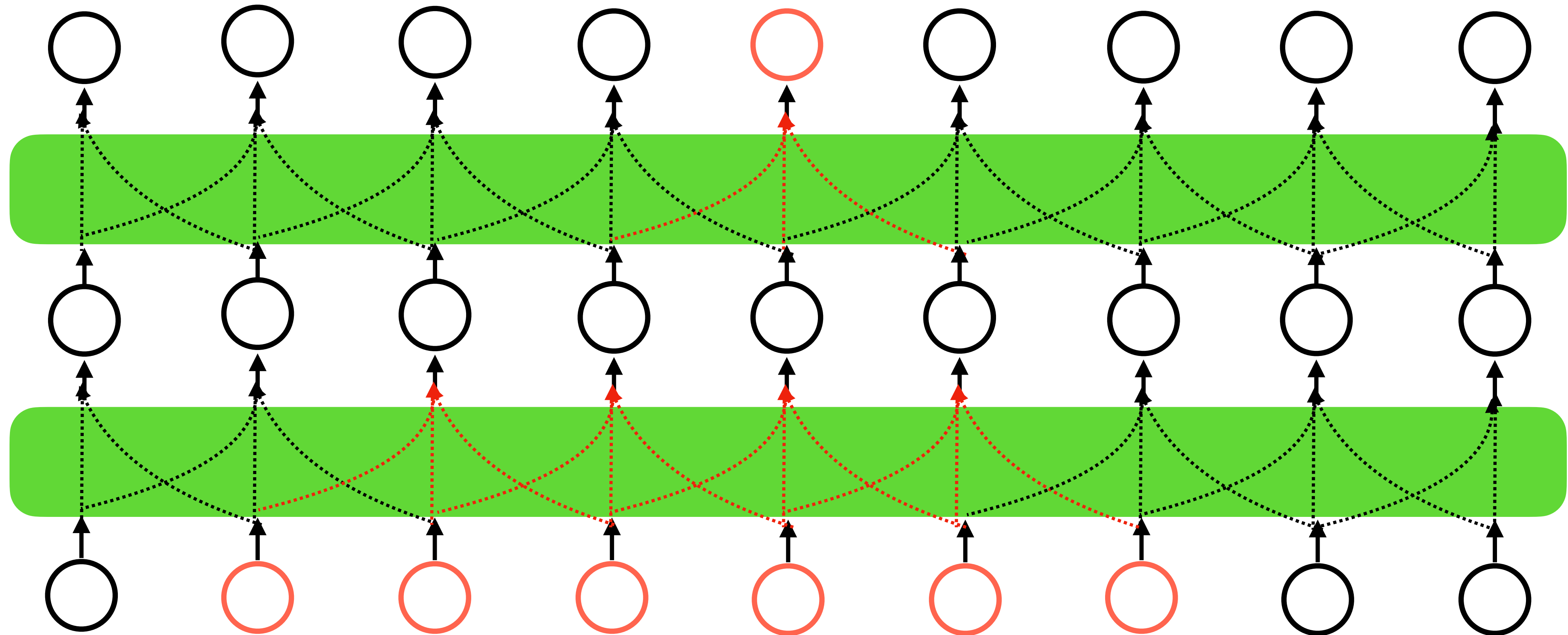


# Convolution Style Models

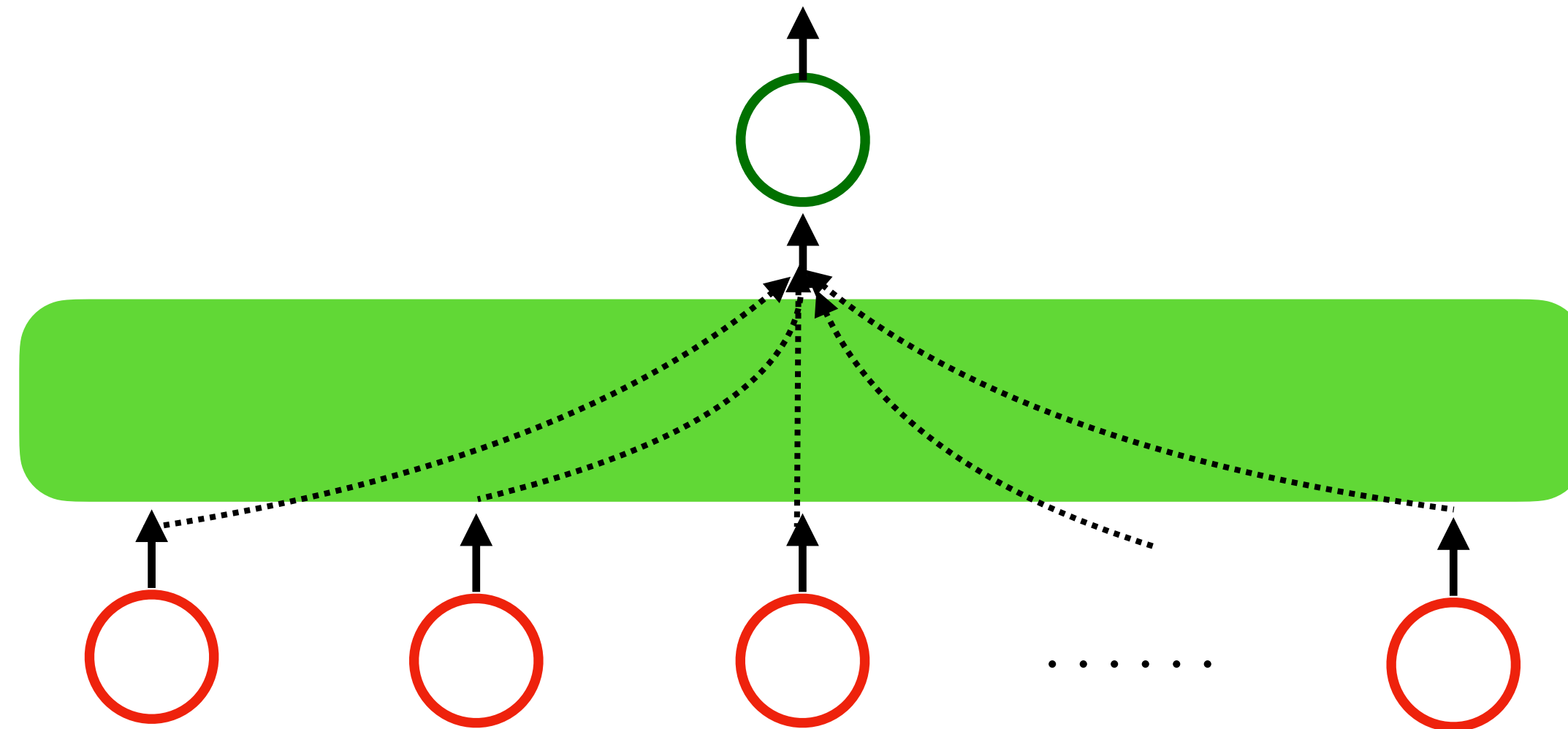




# Convolution Style Models

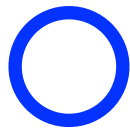
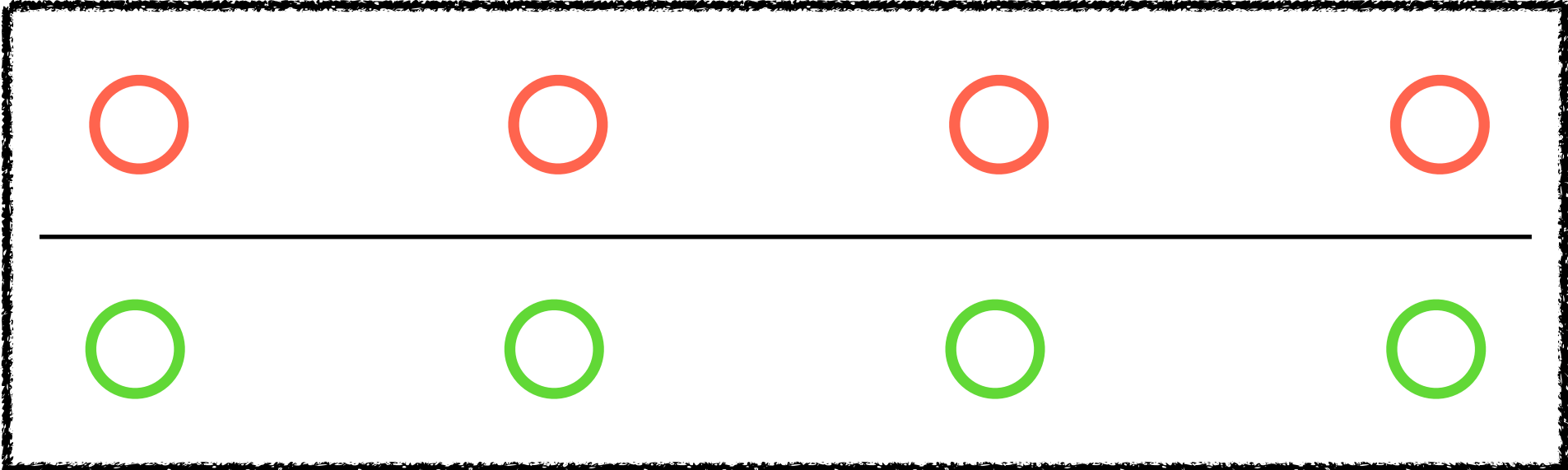


# Considering the full sequence as context



How can we achieve this?

# Dot-Product-Softmax Attention



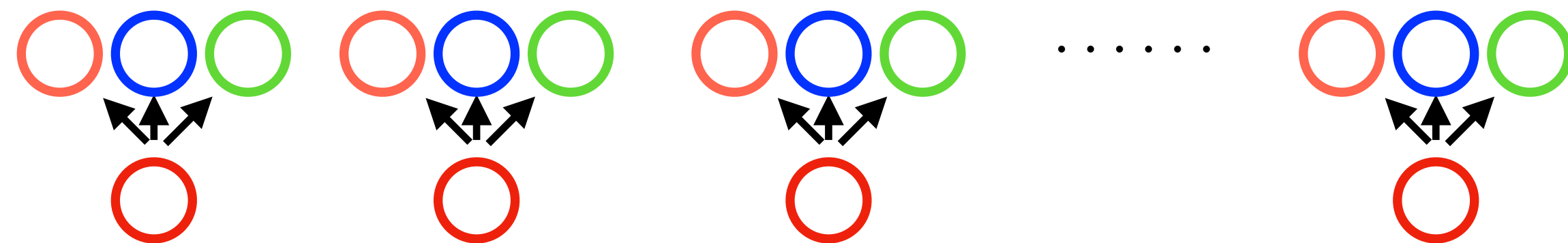
Query

Memory (key-value pairs)

$$\begin{array}{l}
 \text{○} \text{○} \quad \mathbf{q} \cdot \mathbf{k}_1 \quad \mathbf{q} \cdot \mathbf{k}_1 \\
 \text{○} \text{○} \quad \mathbf{q} \cdot \mathbf{k}_2 \quad \mathbf{q} \cdot \mathbf{k}_2 \\
 \text{○} \text{○} \quad \mathbf{q} \cdot \mathbf{k}_3 \quad \mathbf{q} \cdot \mathbf{k}_3 \\
 \text{○} \text{○} \quad \mathbf{q} \cdot \mathbf{k}_4 \quad \mathbf{q} \cdot \mathbf{k}_4
 \end{array}
 \text{softmax}(\quad) \rightarrow
 \begin{bmatrix} 0.6 \\ 0.1 \\ 0.2 \\ 0.1 \end{bmatrix}
 \begin{bmatrix} \text{○} \\ \text{○} \\ \text{○} \\ \text{○} \end{bmatrix}
 \rightarrow
 0.6 \text{○} + 0.1 \text{○} + 0.2 \text{○} + 0.1 \text{○}$$

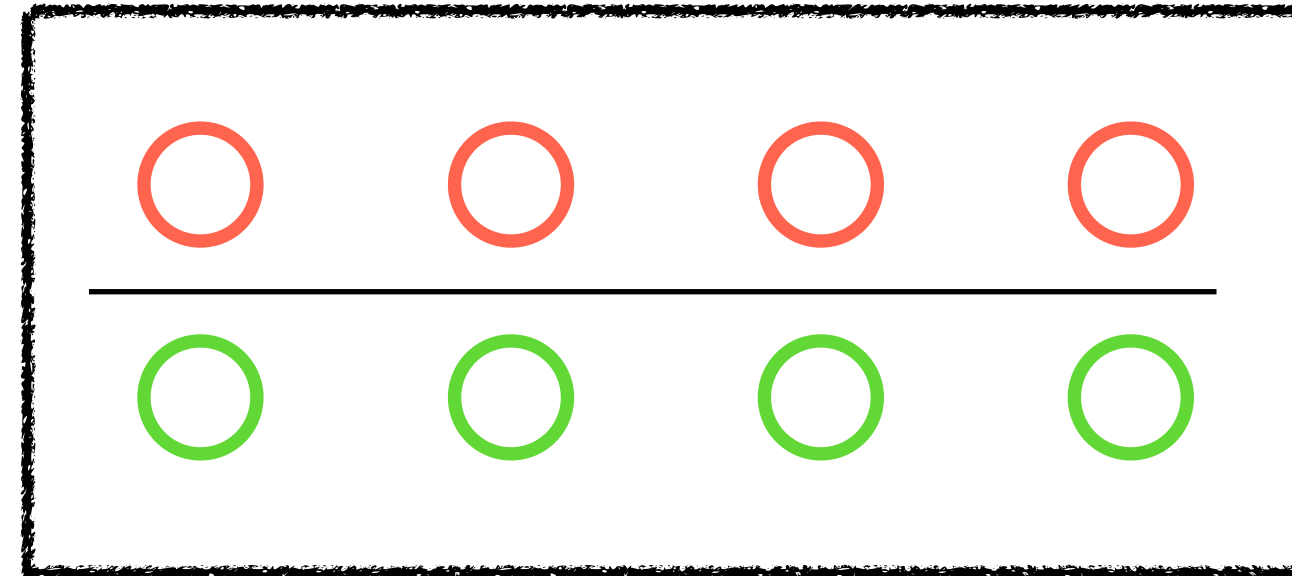
=  $\text{○}$  context vector  $\mathbf{c}$

# Considering the full sequence as context



# Attention Mechanism

○  
Query



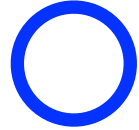
Memory (key-value pairs)

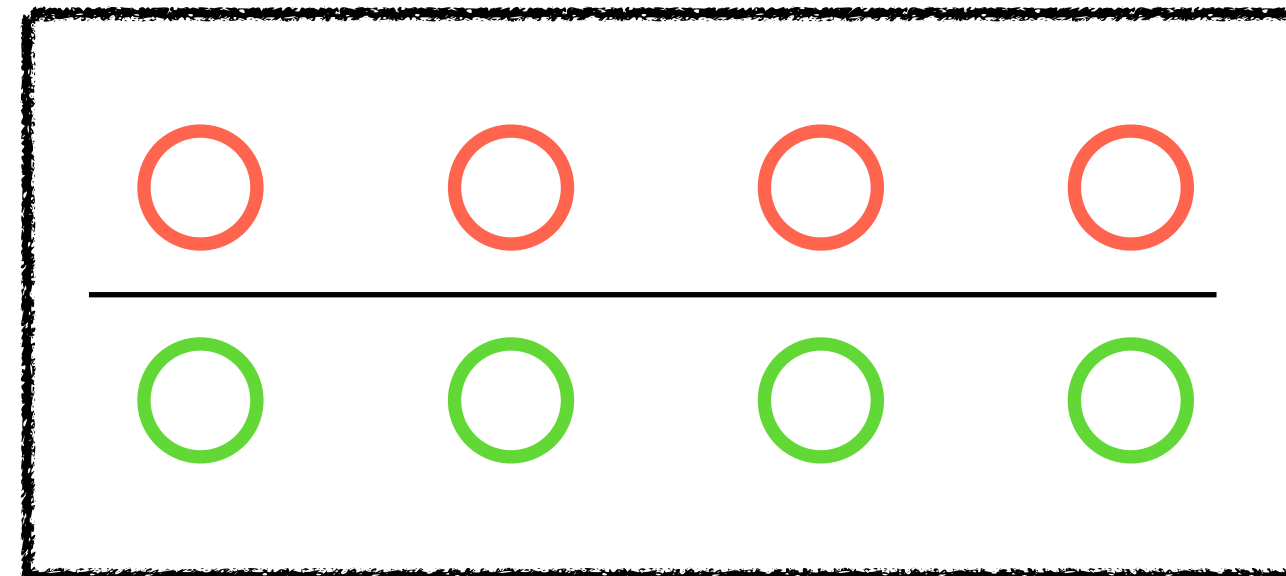


# Attention Mechanism

$$0.6 \bigcirc + 0.1 \bigcirc + 0.2 \bigcirc + 0.1 \bigcirc$$

$$= \bigcirc \text{ context vector } \mathbf{c}$$

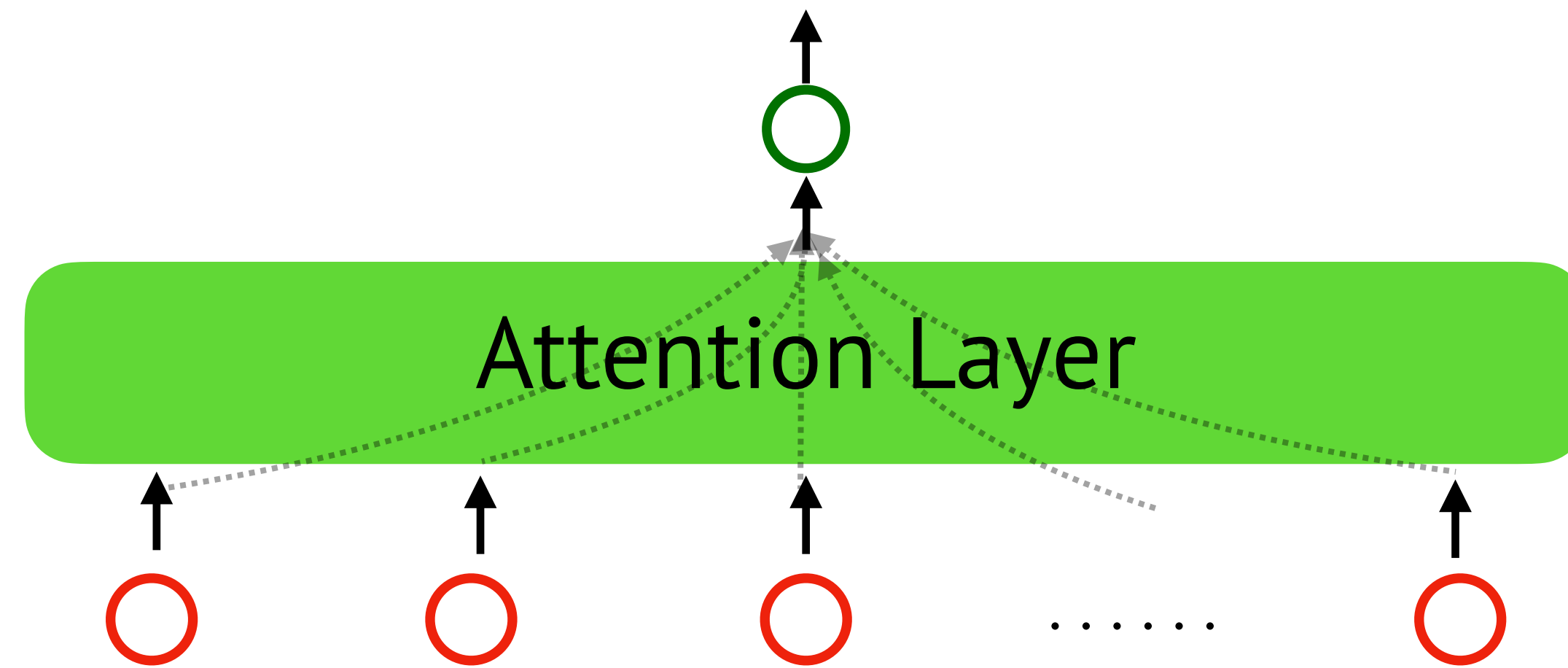
  
Query



Memory (key-value pairs)

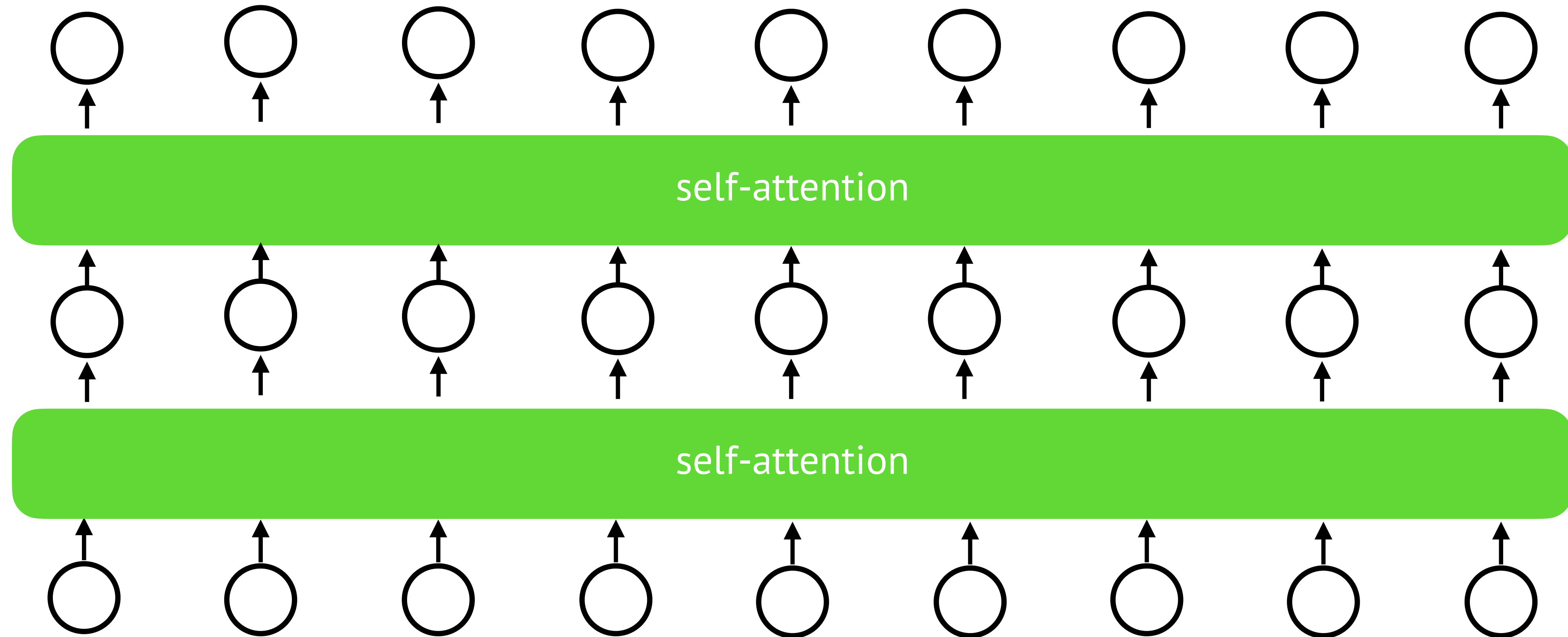


# Self-attention



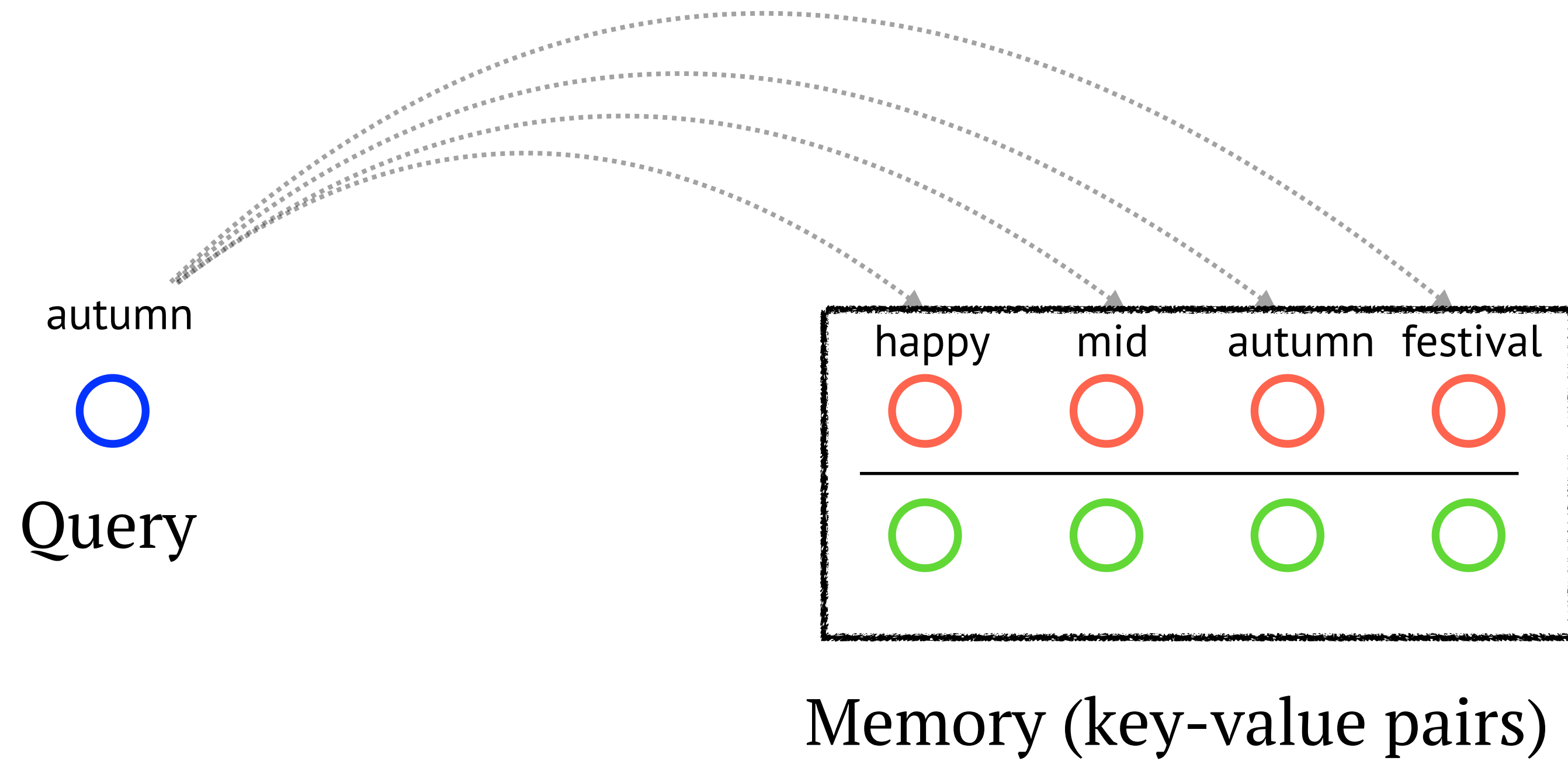
This is almost transformer — except a few things.

# Transformer (almost)

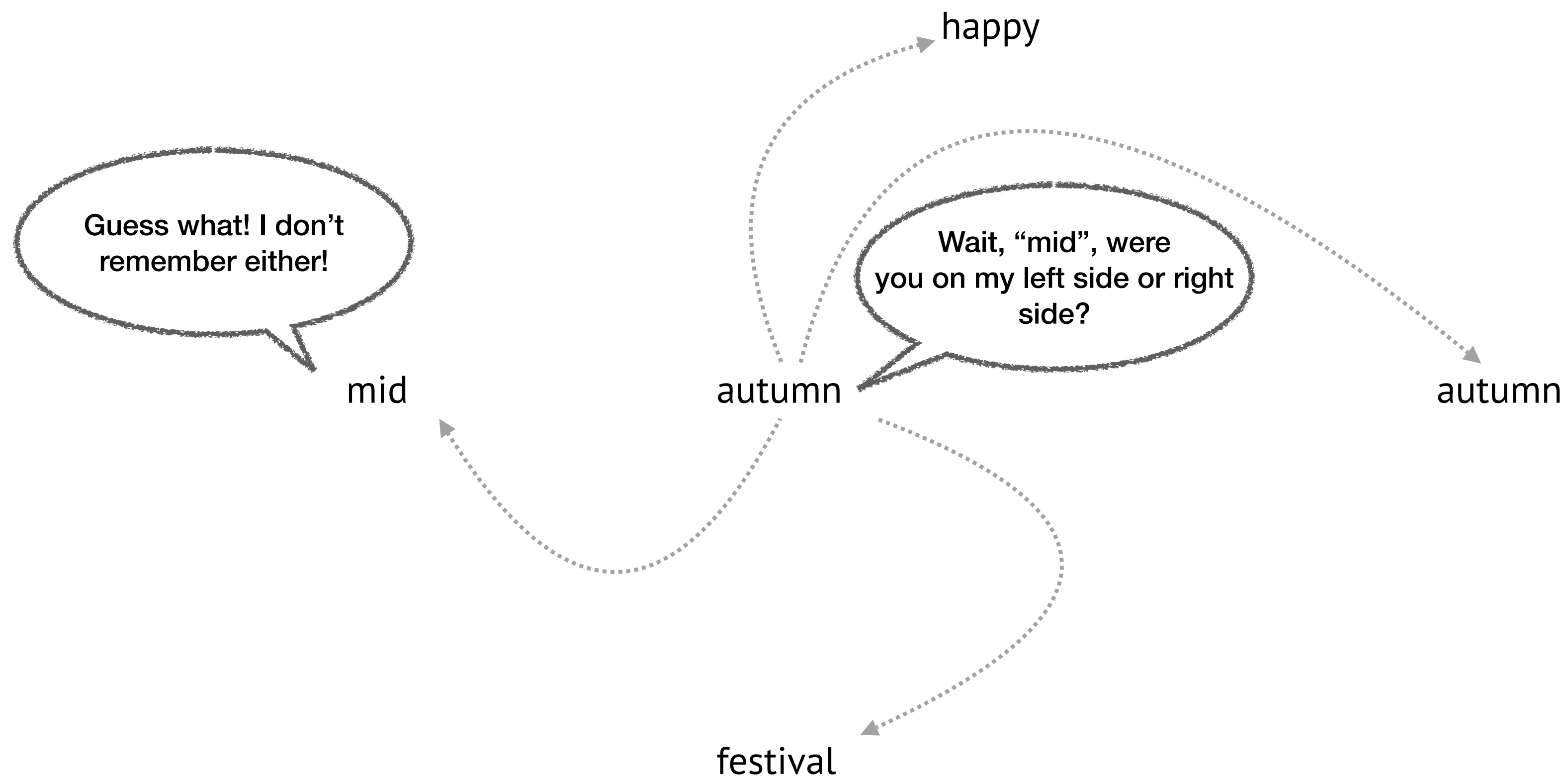




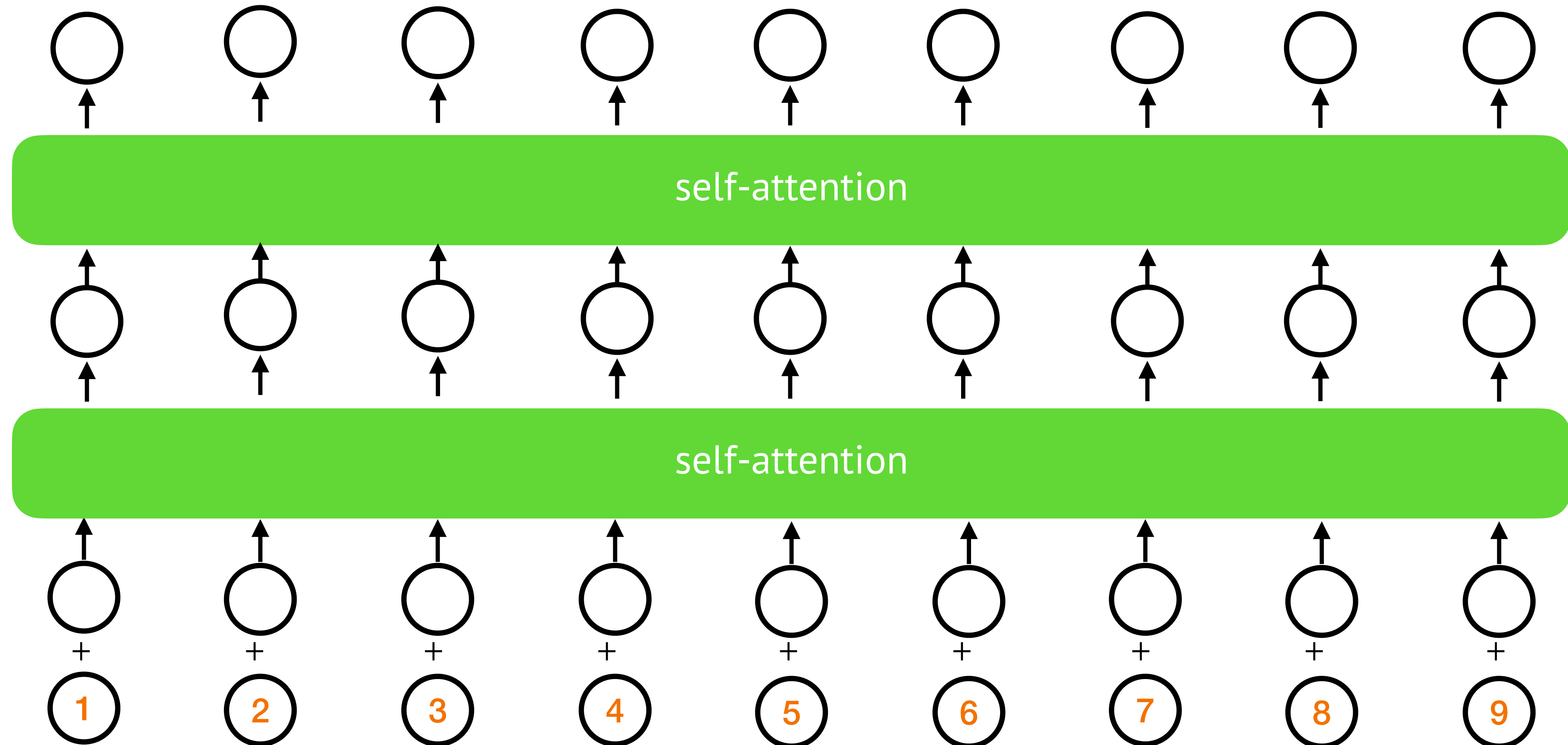
# Self-attention in Transformer



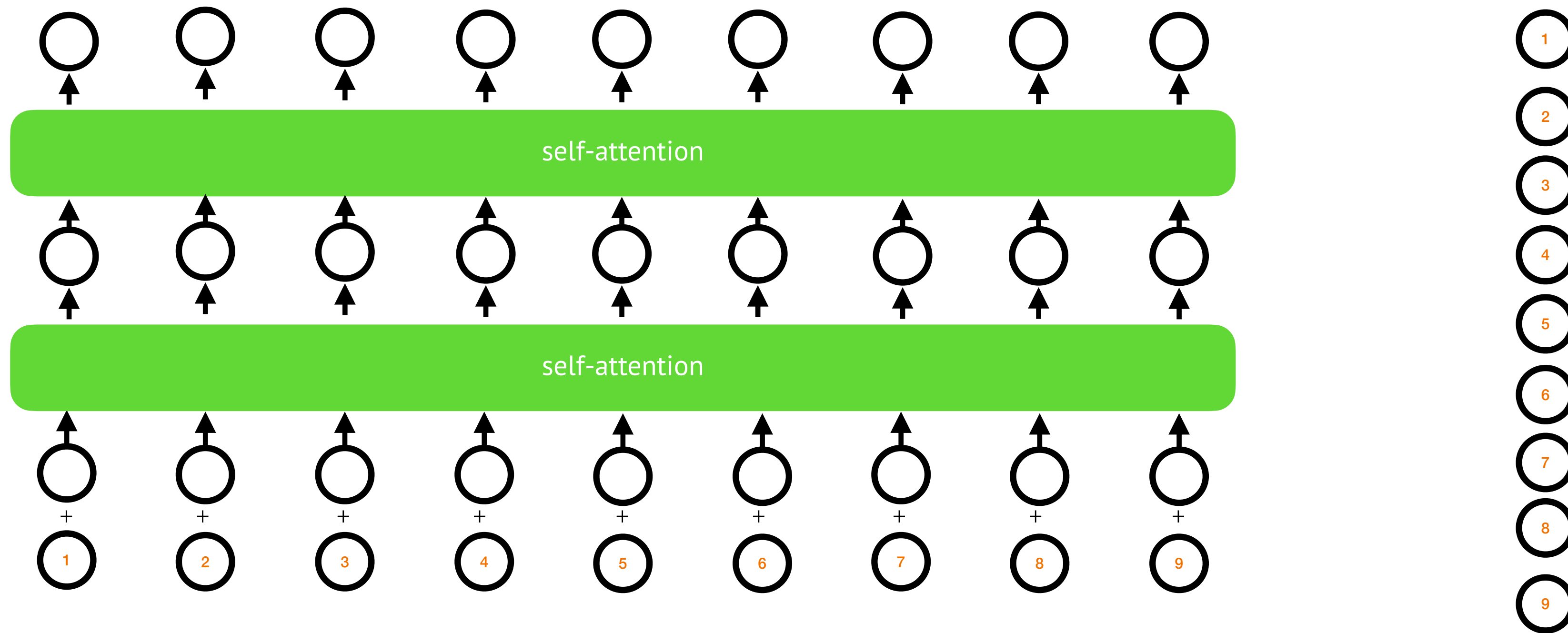
# Self-attention in Transformer



# Positional Embeddings



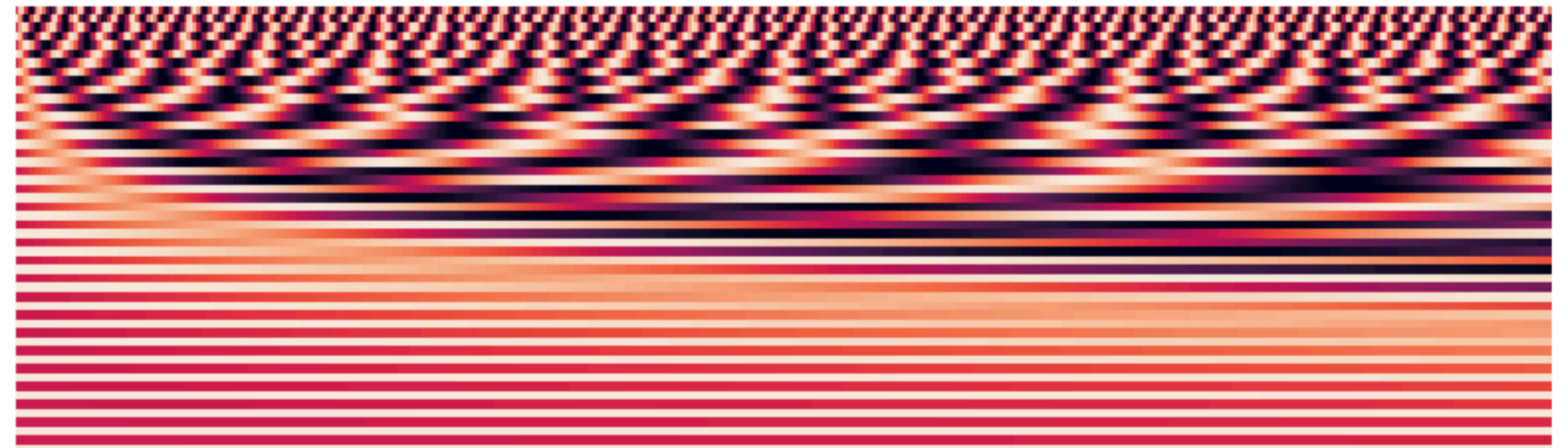
# Transformer (positional embedding)



# Positional Encoding

$$\begin{bmatrix} \sin\left(\frac{i}{10000^{2 \times \frac{1}{d}}}\right) \\ \cos\left(\frac{i}{10000^{2 \times \frac{1}{d}}}\right) \\ \vdots \\ \sin\left(\frac{i}{10000^{2 \times \frac{d/2}{d}}}\right) \\ \cos\left(\frac{i}{10000^{2 \times \frac{d/2}{d}}}\right) \end{bmatrix}$$

Dimension



Index in the sequence

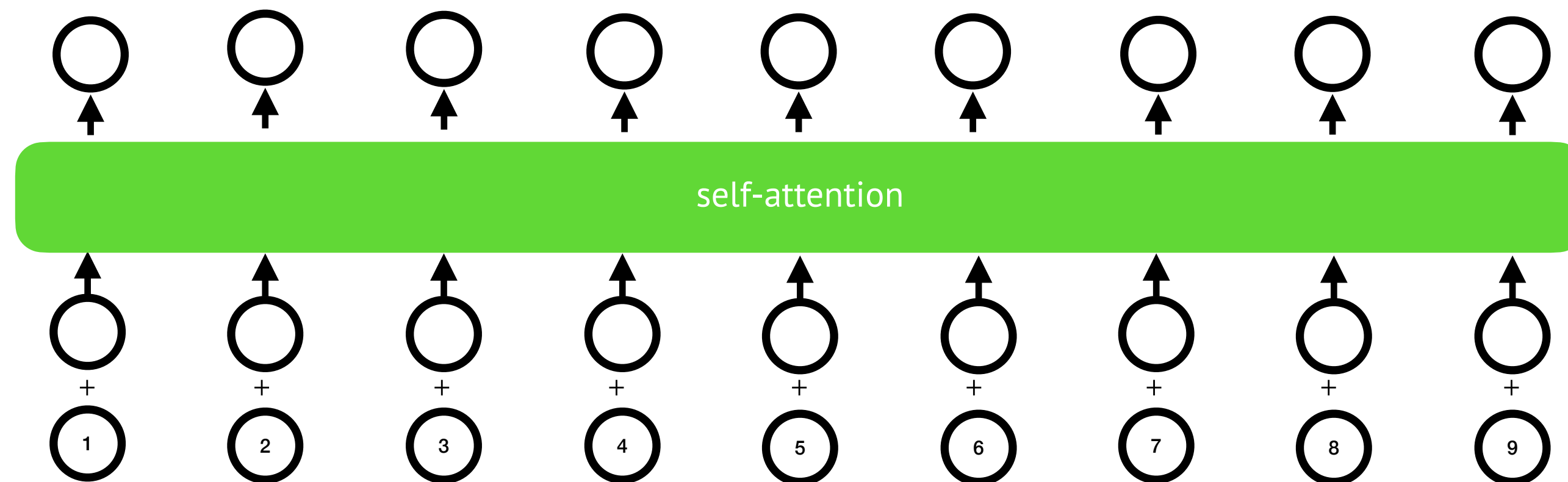


The idea of relative position

# Positional Encoding

$$\begin{bmatrix} \sin\left(\frac{i}{10000^{2 \times \frac{1}{d}}}\right) \\ \cos\left(\frac{i}{10000^{2 \times \frac{1}{d}}}\right) \\ \vdots \\ \sin\left(\frac{i}{10000^{2 \times \frac{d/2}{d}}}\right) \\ \cos\left(\frac{i}{10000^{2 \times \frac{d/2}{d}}}\right) \end{bmatrix}$$

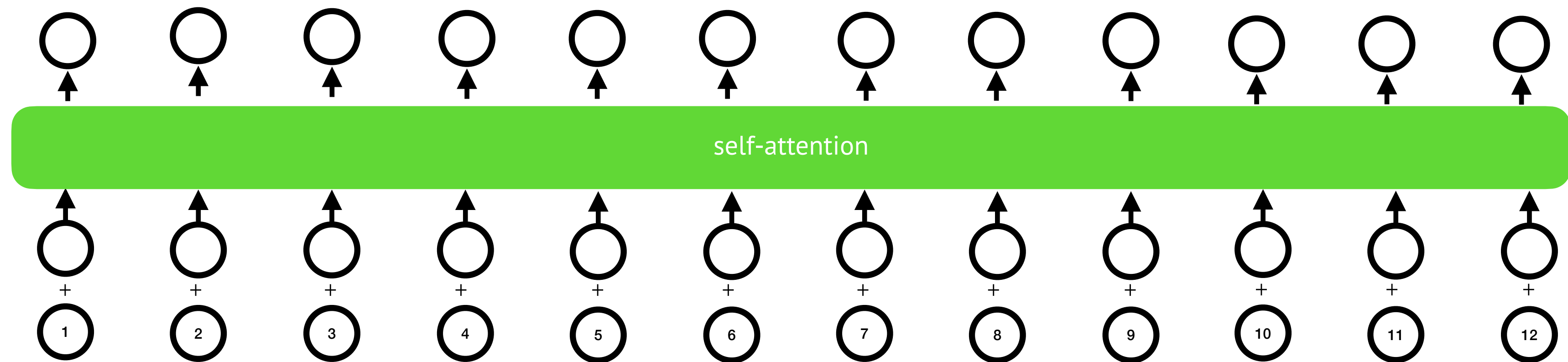
Periodic: Hope this will work in extrapolation. (No)



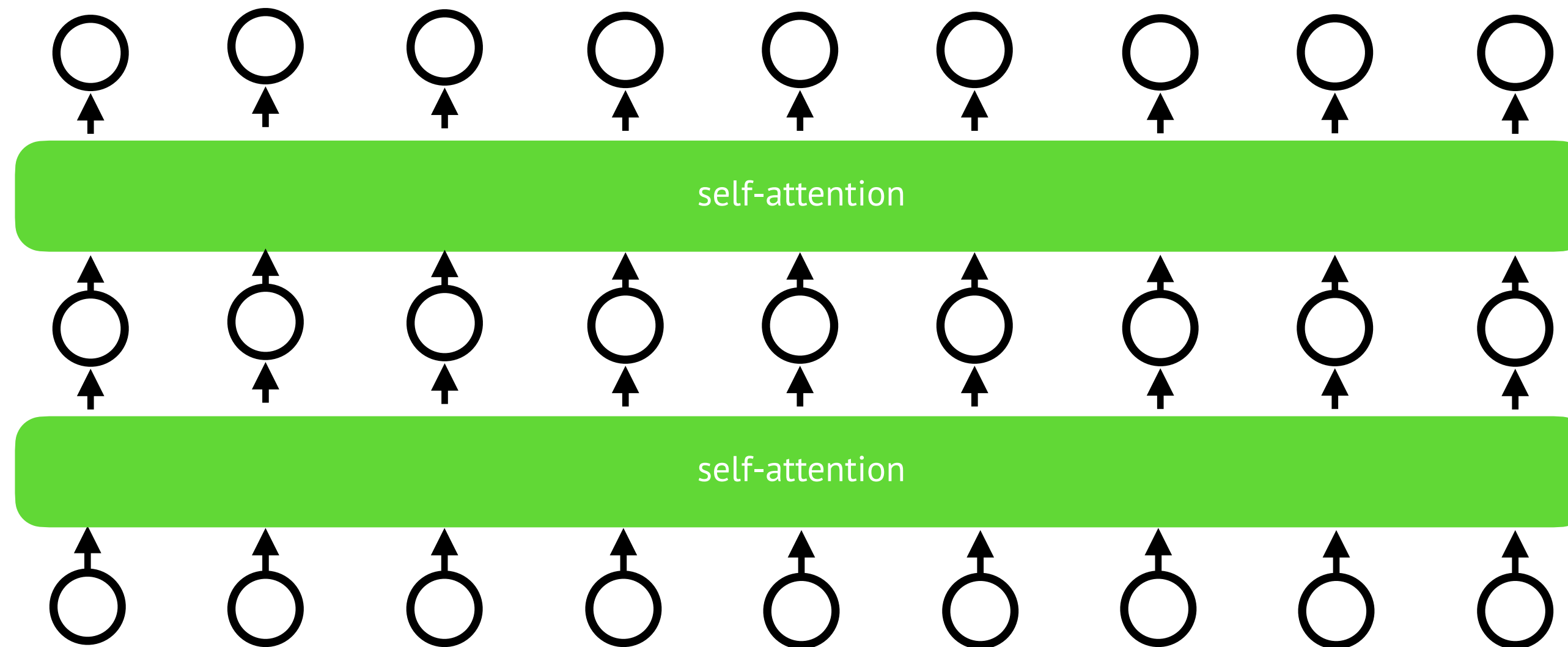
# Positional Encoding

$$\begin{bmatrix} \sin\left(\frac{i}{10000^{2 \times \frac{1}{d}}}\right) \\ \cos\left(\frac{i}{10000^{2 \times \frac{1}{d}}}\right) \\ \vdots \\ \sin\left(\frac{i}{10000^{2 \times \frac{d/2}{d}}}\right) \\ \cos\left(\frac{i}{10000^{2 \times \frac{d/2}{d}}}\right) \end{bmatrix}$$

Periodic: Hope this will work in extrapolation. (No)

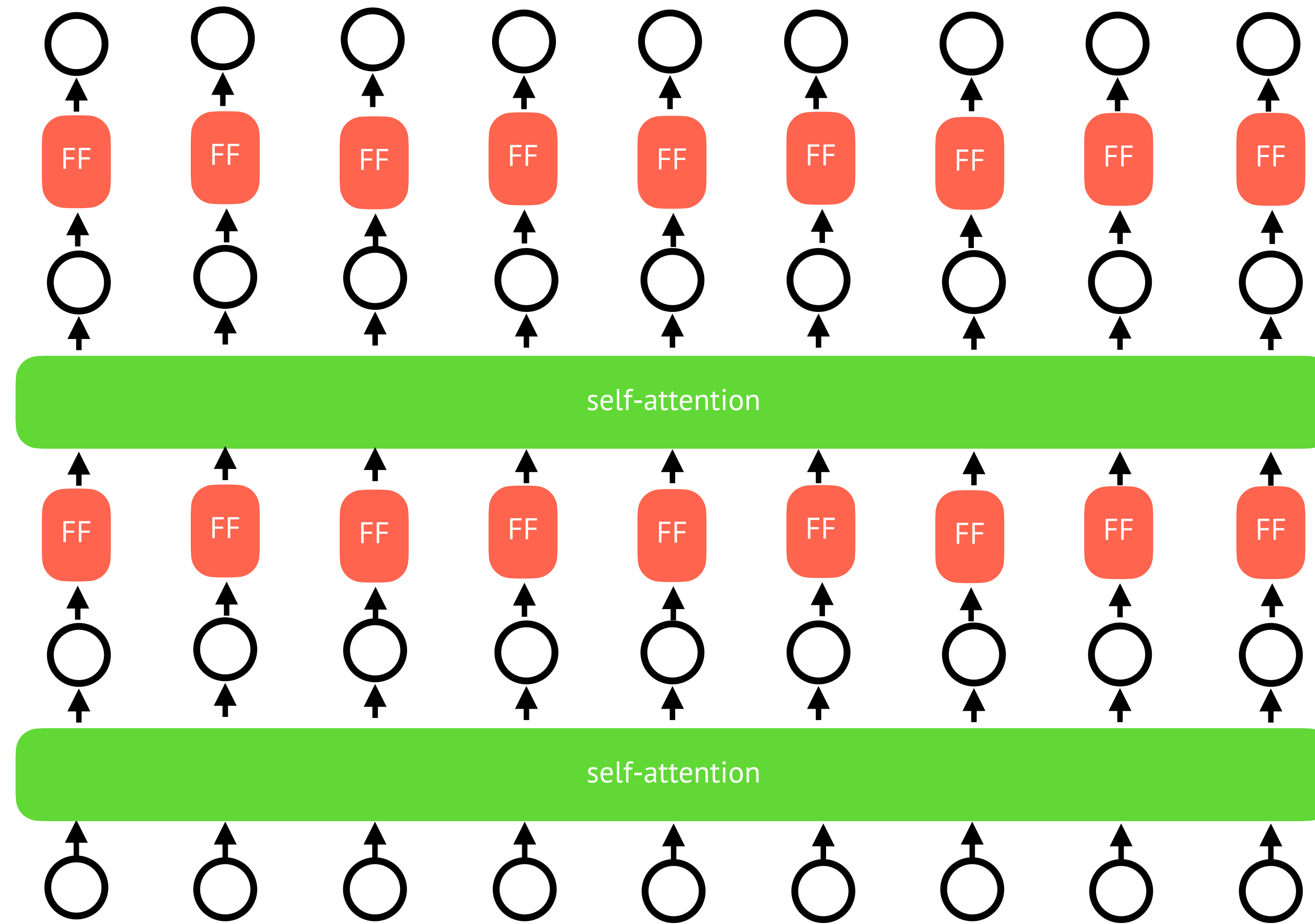


# Feed Forward Layer





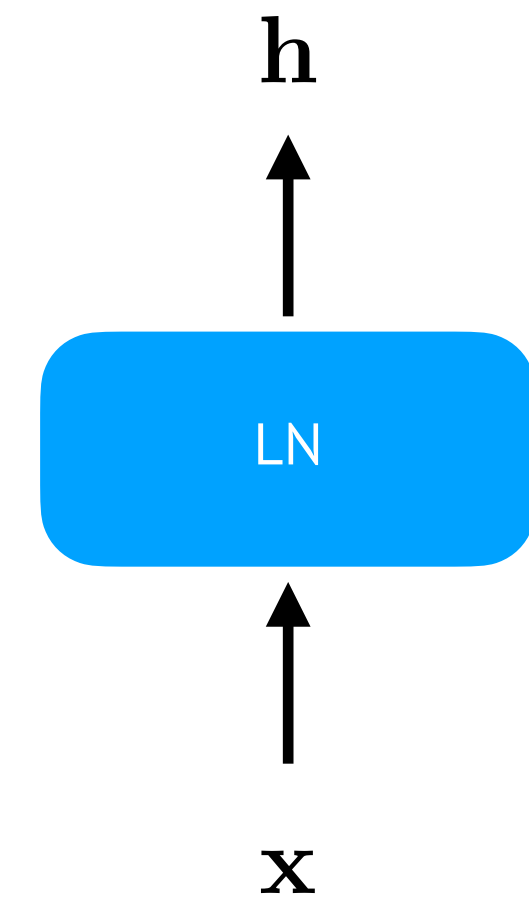
# Feed Forward Layer



# Layer Normalization (Ba et al, 2016)

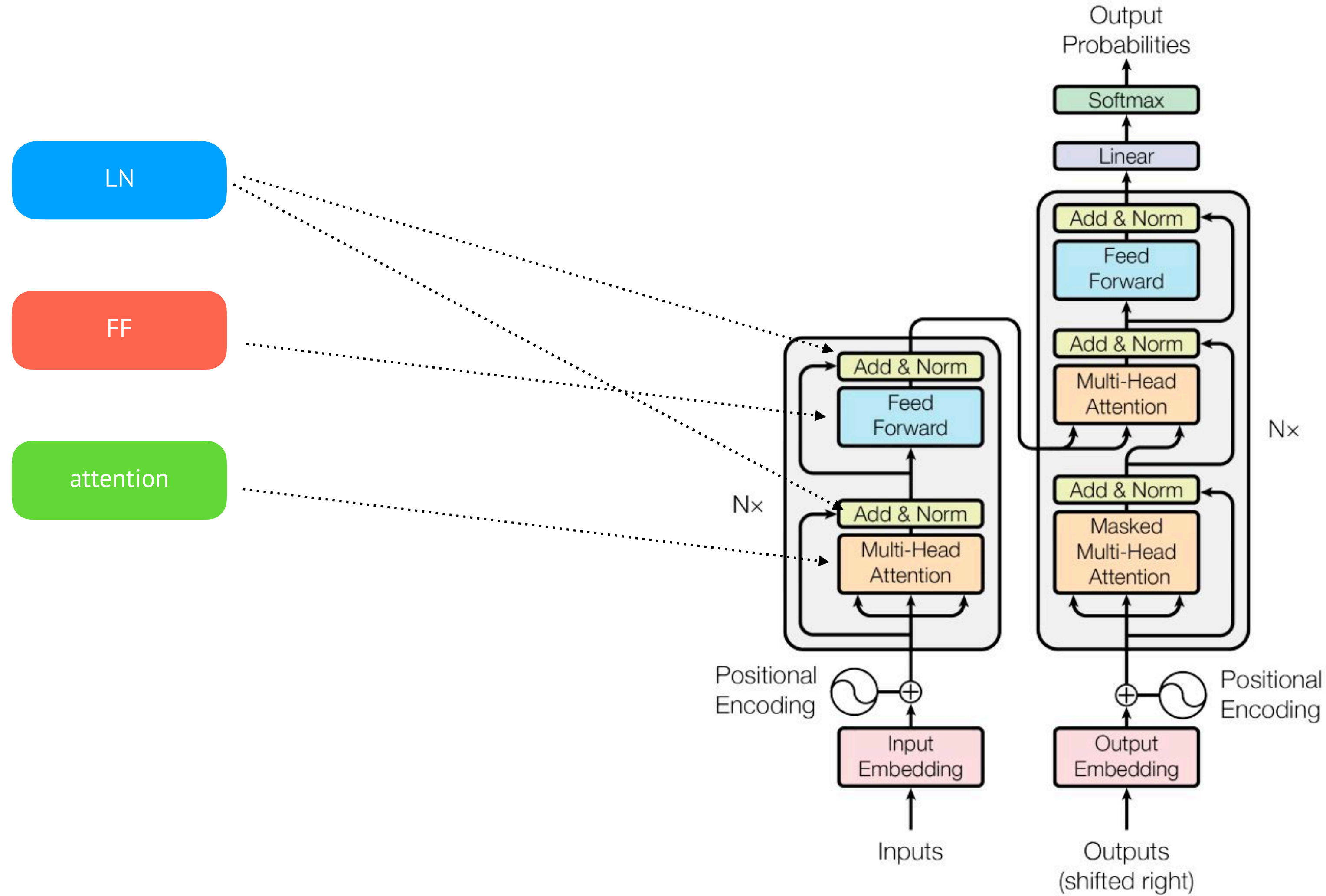
$$\mathbf{h} = \mathbf{g} \odot N(\mathbf{x}) + \mathbf{b}$$

$$N(\mathbf{x}) = \frac{\mathbf{x} - \mu}{\sigma} \quad \mu = \frac{1}{H} \sum_{i=1}^H x_i \quad \sigma = \sqrt{\frac{1}{H} \sum_{i=1}^H (x_i - \mu)^2}$$



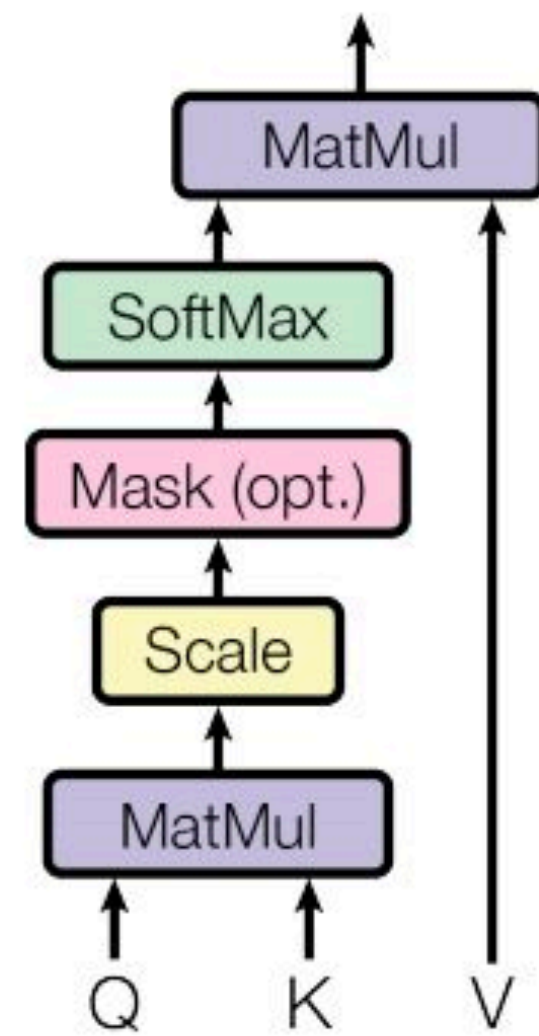
Smoother gradients, faster training and better generalization accuracy. (Xu et al, Neurips 2019)

# Layer Normalization



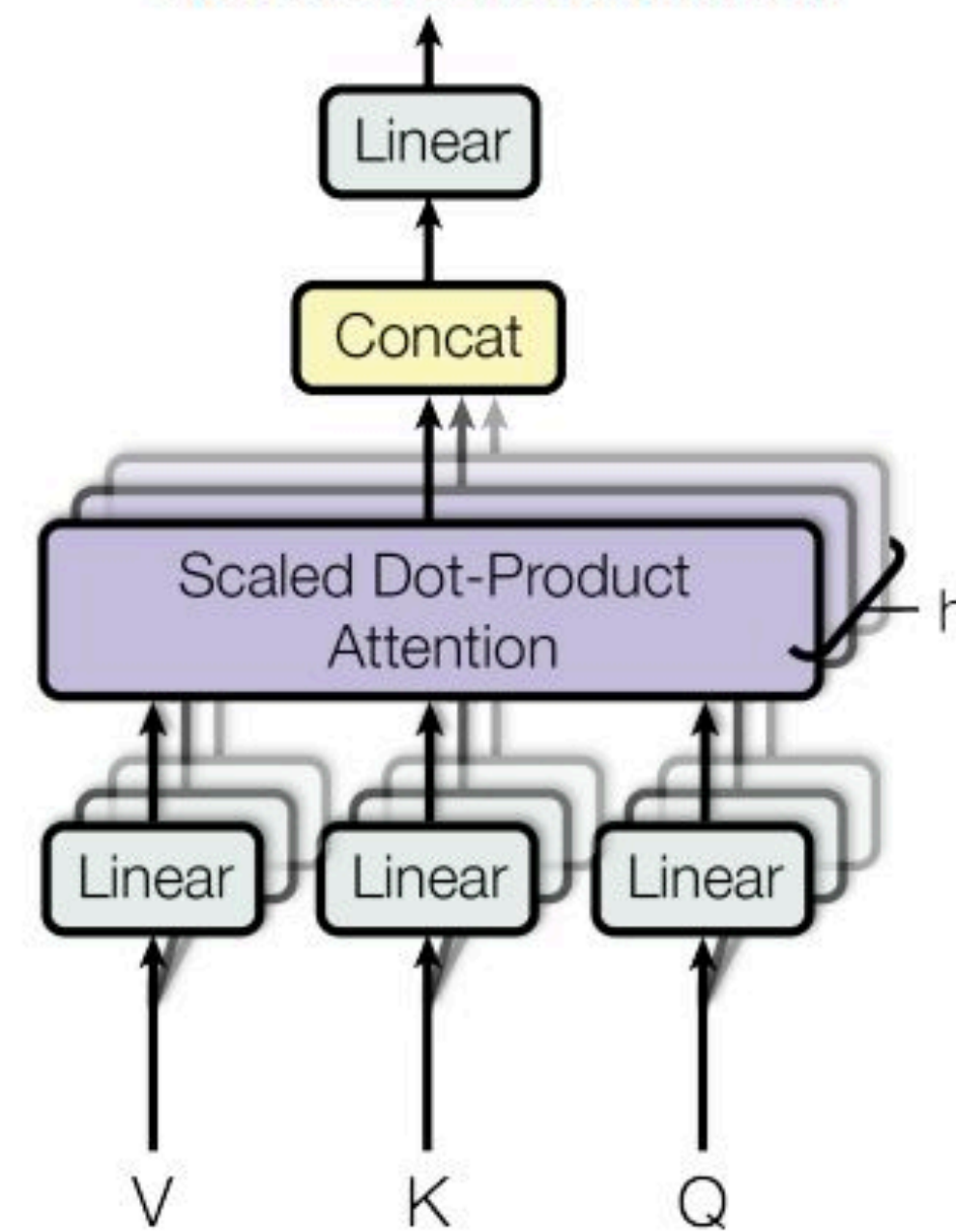
# Multi-head Attention

Scaled Dot-Product Attention



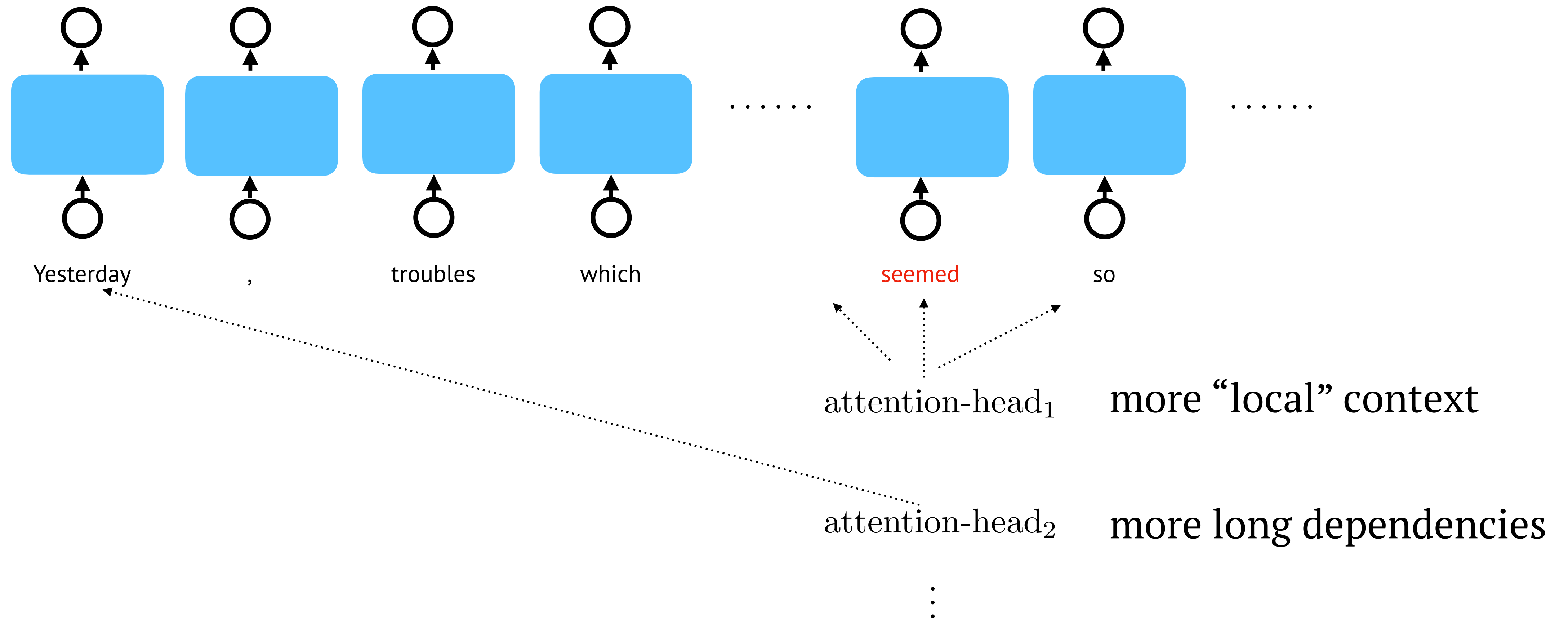
$$\text{score}(q, k) = \frac{q^T k}{\sqrt{d_k}}$$

Multi-Head Attention



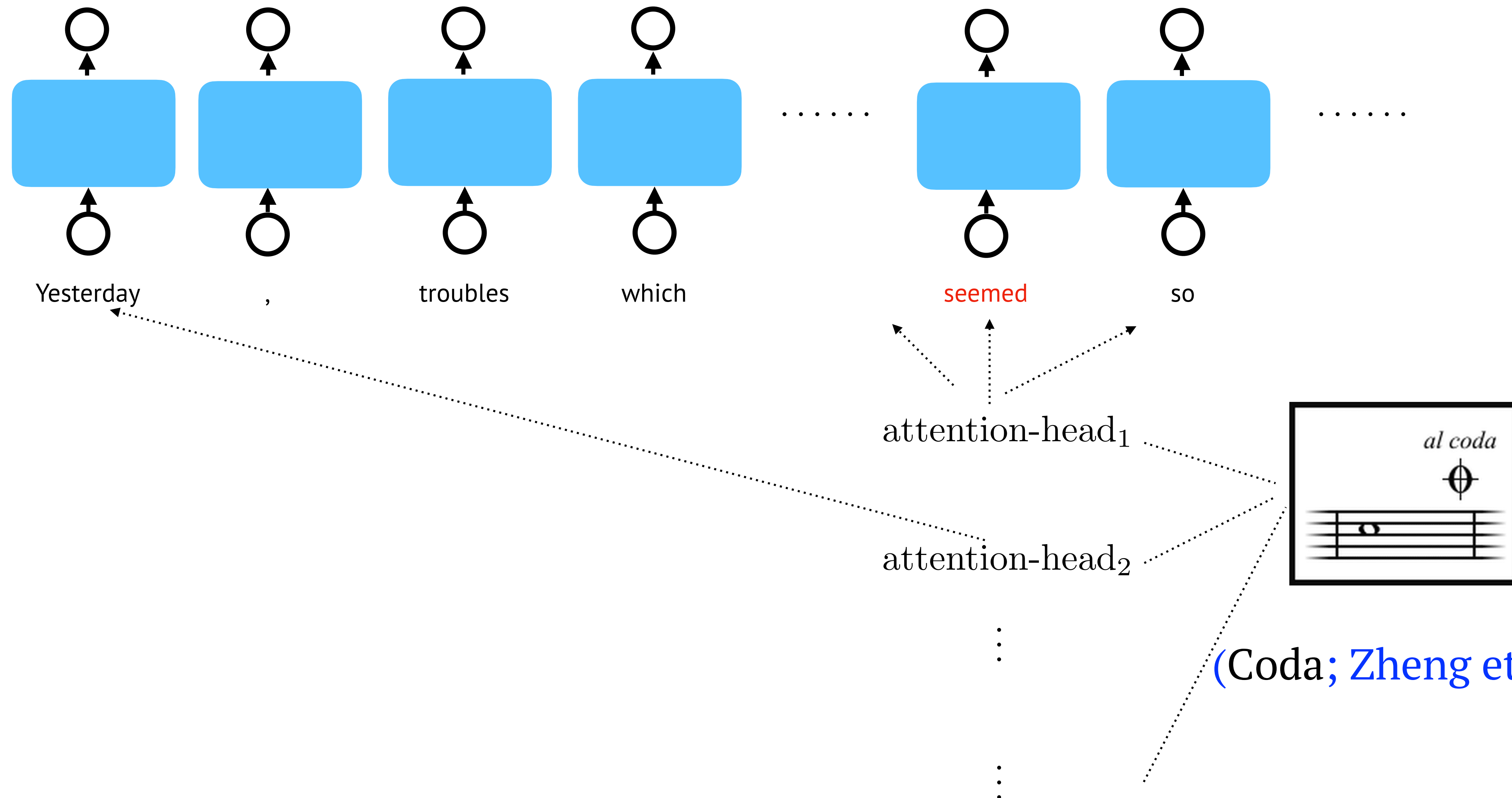
multiple copies

# Multi-head Attention



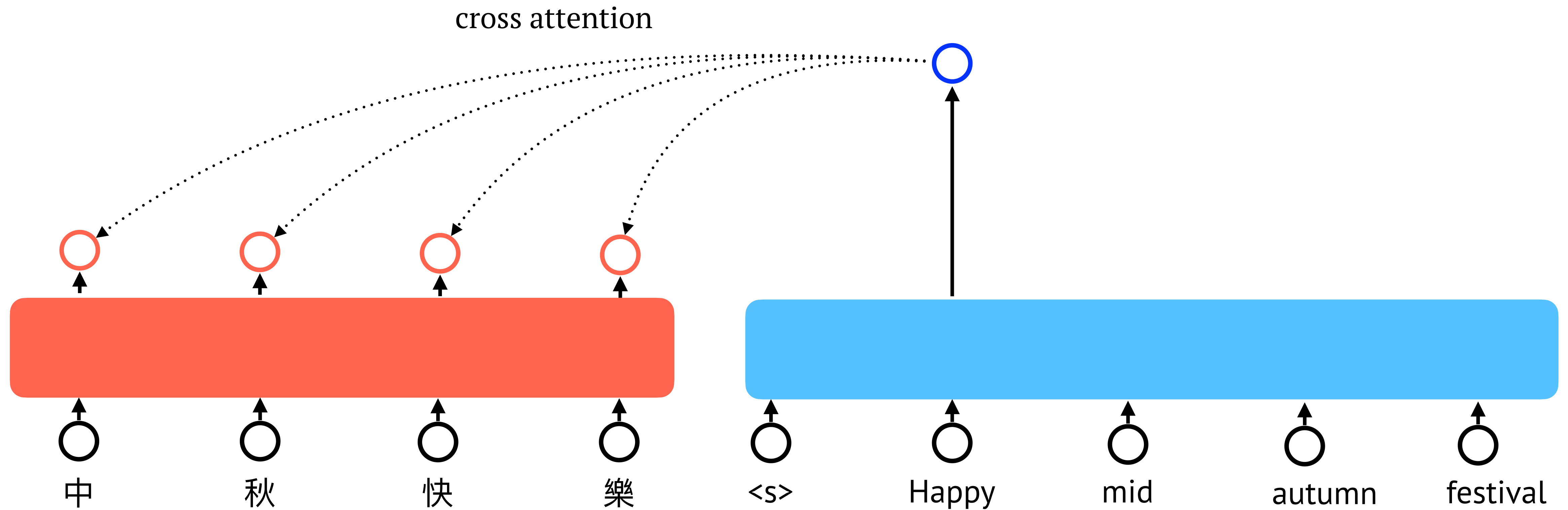
Improve the “resolution” of the attention mechanism.

# Multi-head Attention

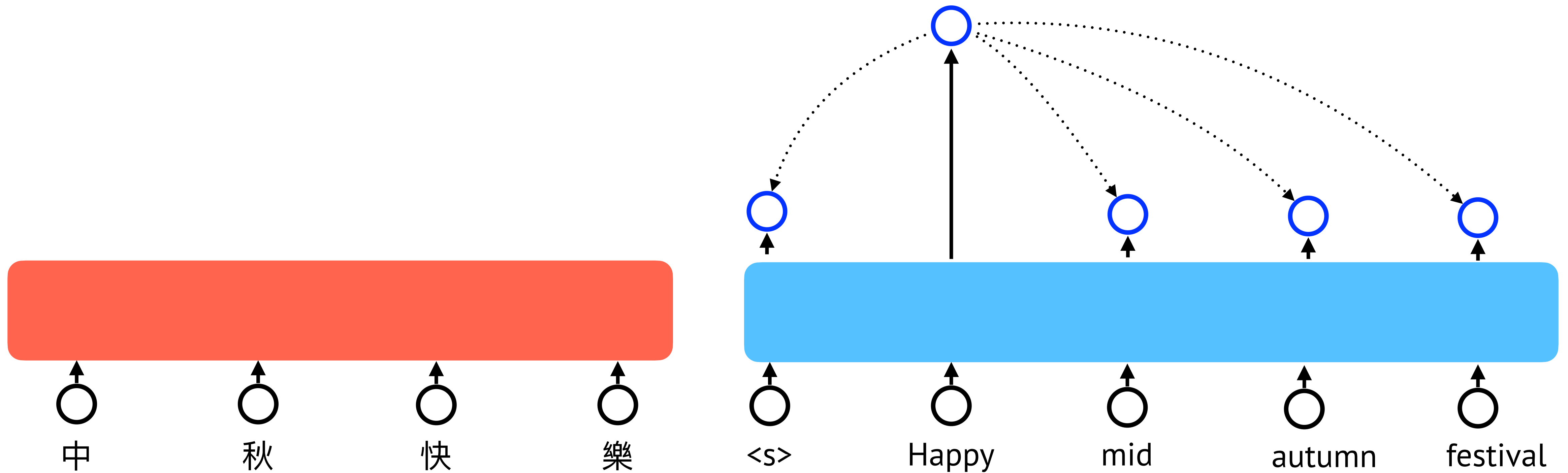


(Coda; Zheng et al, ACL 2021)

# Transformer as Decoder

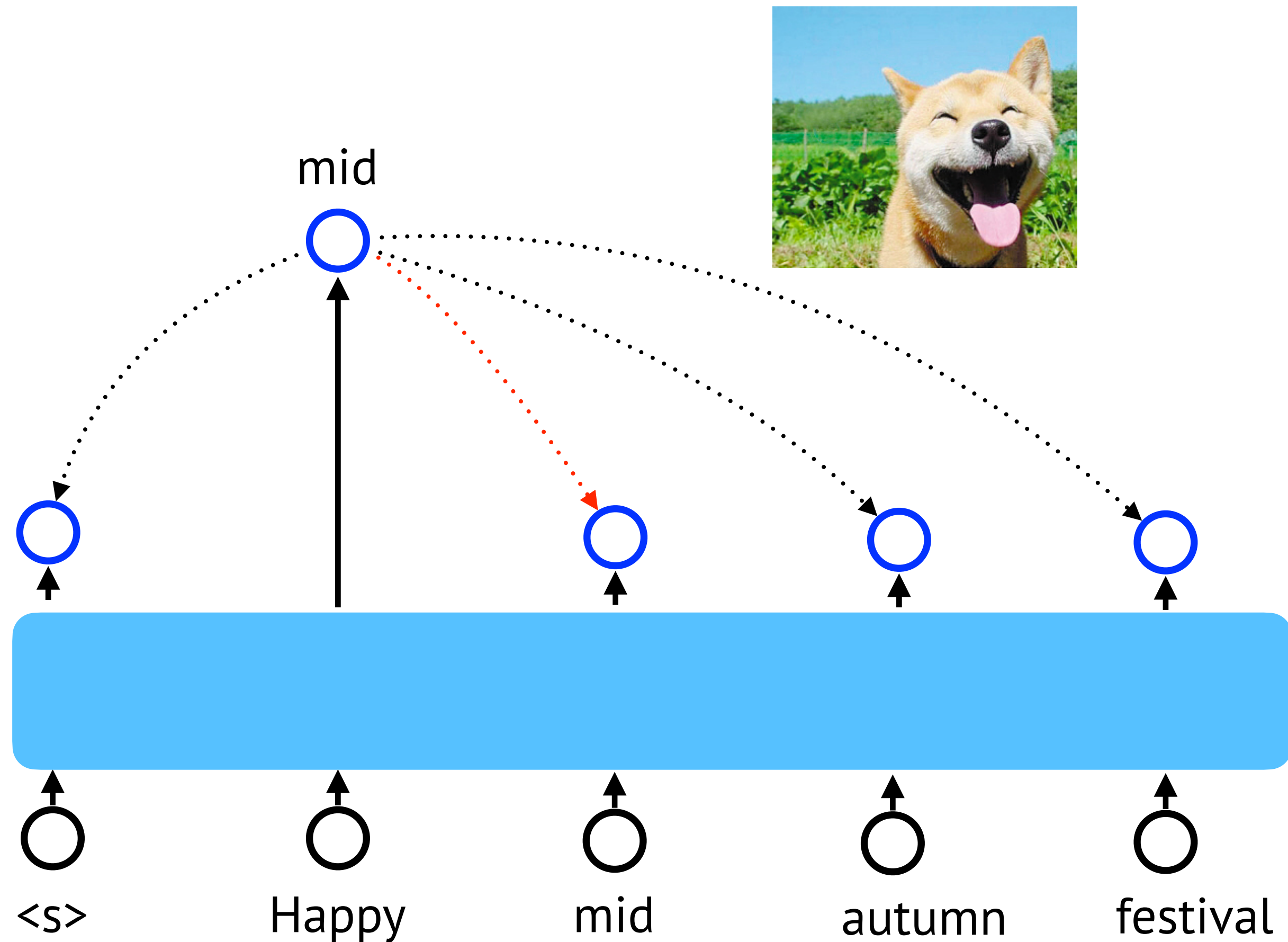


# Transformer as Decoder





# Transformer as Decoder

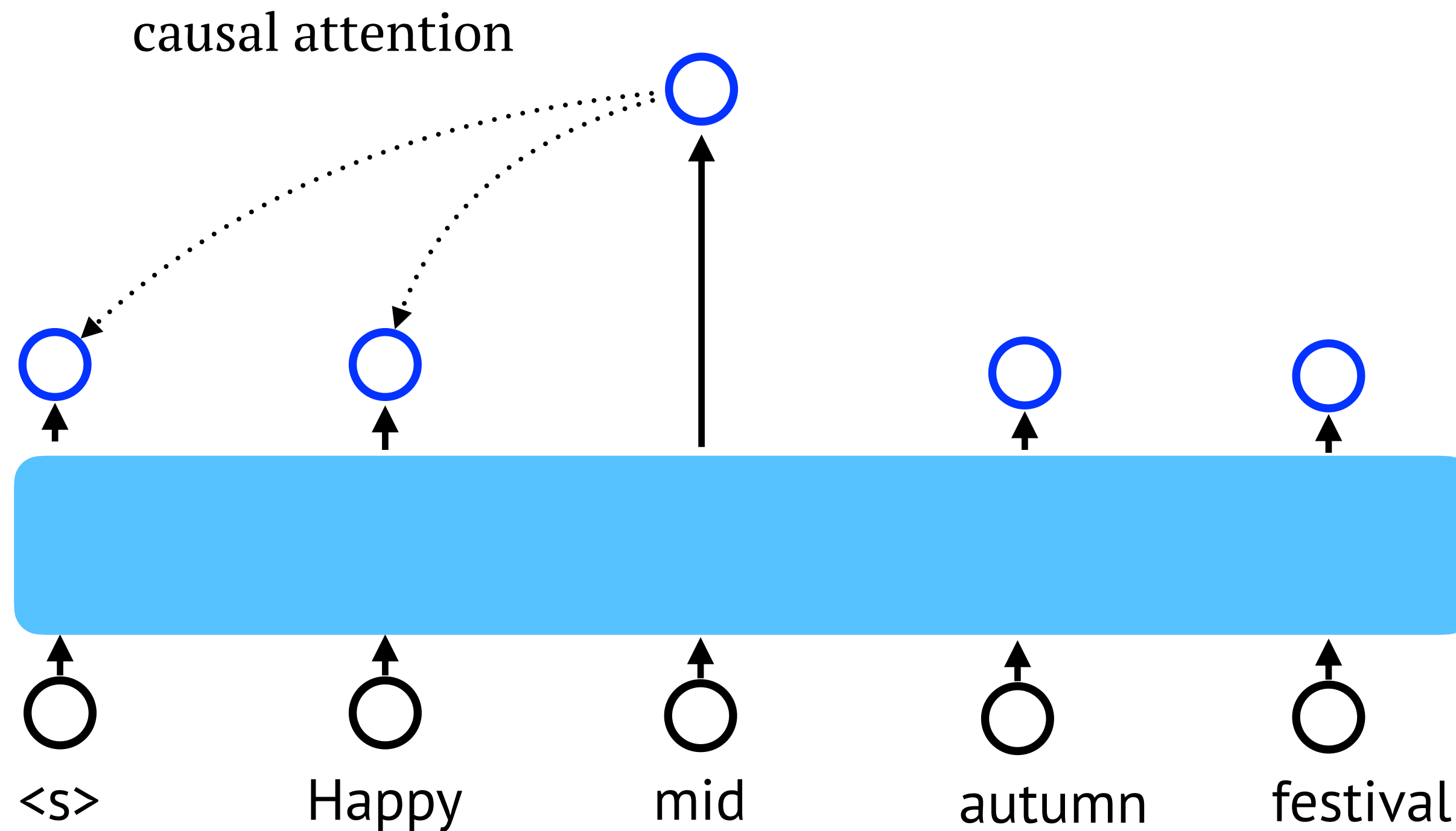


Need to prevent the attention the future words.

	Happy	mid	autumn	festival
Happy	$-\infty$	$-\infty$	$-\infty$	$-\infty$
mid		$-\infty$	$-\infty$	$-\infty$
autumn			$-\infty$	$-\infty$
festival				$-\infty$

$$e_{ij} = \begin{cases} q_i^\top k_j, & j < i \\ -\infty, & j \geq i \end{cases}$$

# Transformer as Decoder

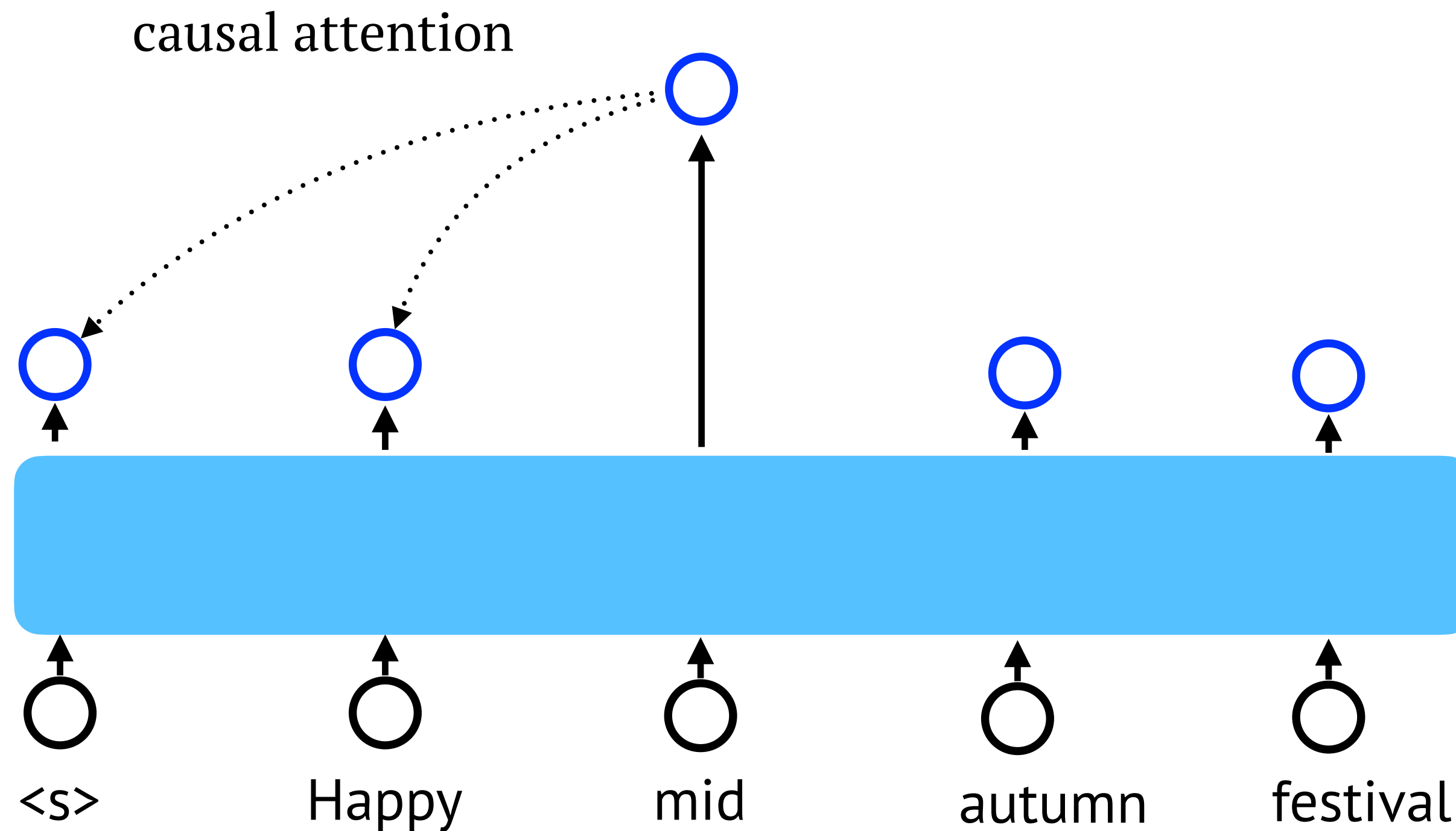


Need to prevent the attention the future words.

	Happy	mid	autumn	festival
Happy	$-\infty$	$-\infty$	$-\infty$	$-\infty$
mid		$-\infty$	$-\infty$	$-\infty$
autumn			$-\infty$	$-\infty$
festival				$-\infty$

$$e_{ij} = \begin{cases} q_i^\top k_j, & j < i \\ -\infty, & j \geq i \end{cases}$$

# Transformer as Decoder



Need to prevent the attention the future words.

	Happy	mid	autumn	festival
Happy	$-\infty$	$-\infty$	$-\infty$	$-\infty$
mid		$-\infty$	$-\infty$	$-\infty$
autumn			$-\infty$	$-\infty$
festival				$-\infty$

$$e_{ij} = \begin{cases} q_i^\top k_j, & j < i \\ -\infty, & j \geq i \end{cases}$$