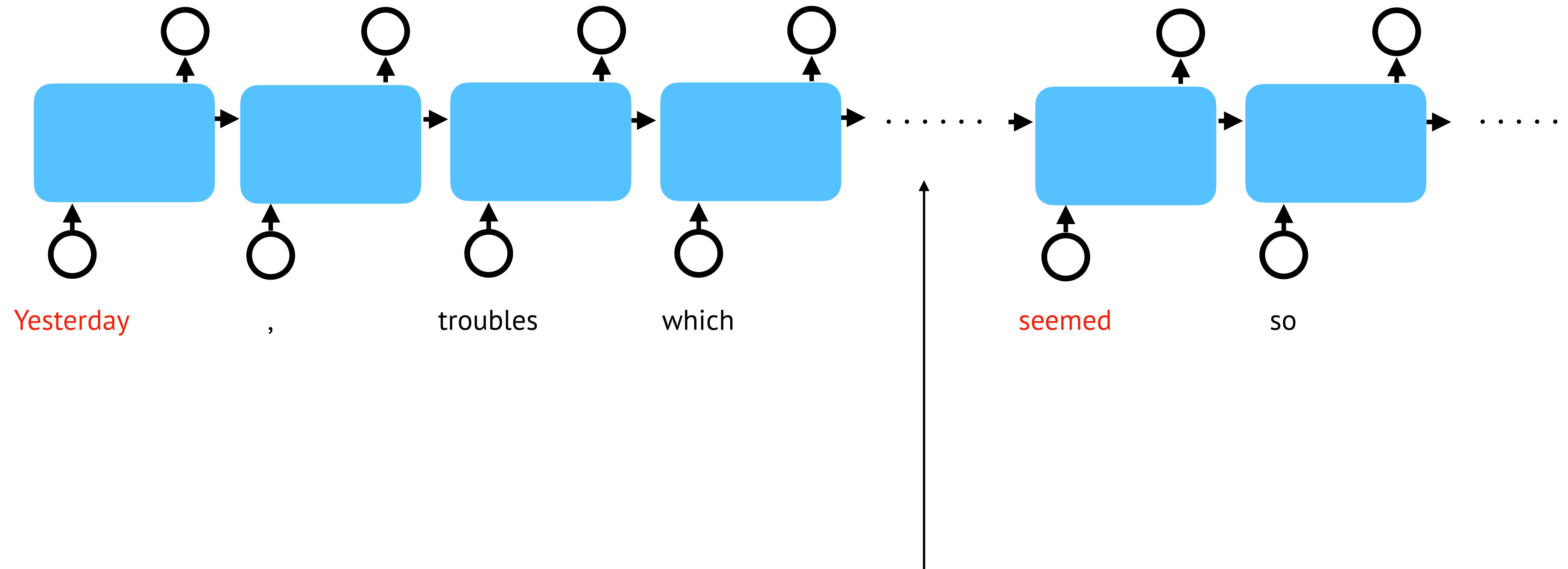# Transformers (1)

## COMP3361 — Week 4

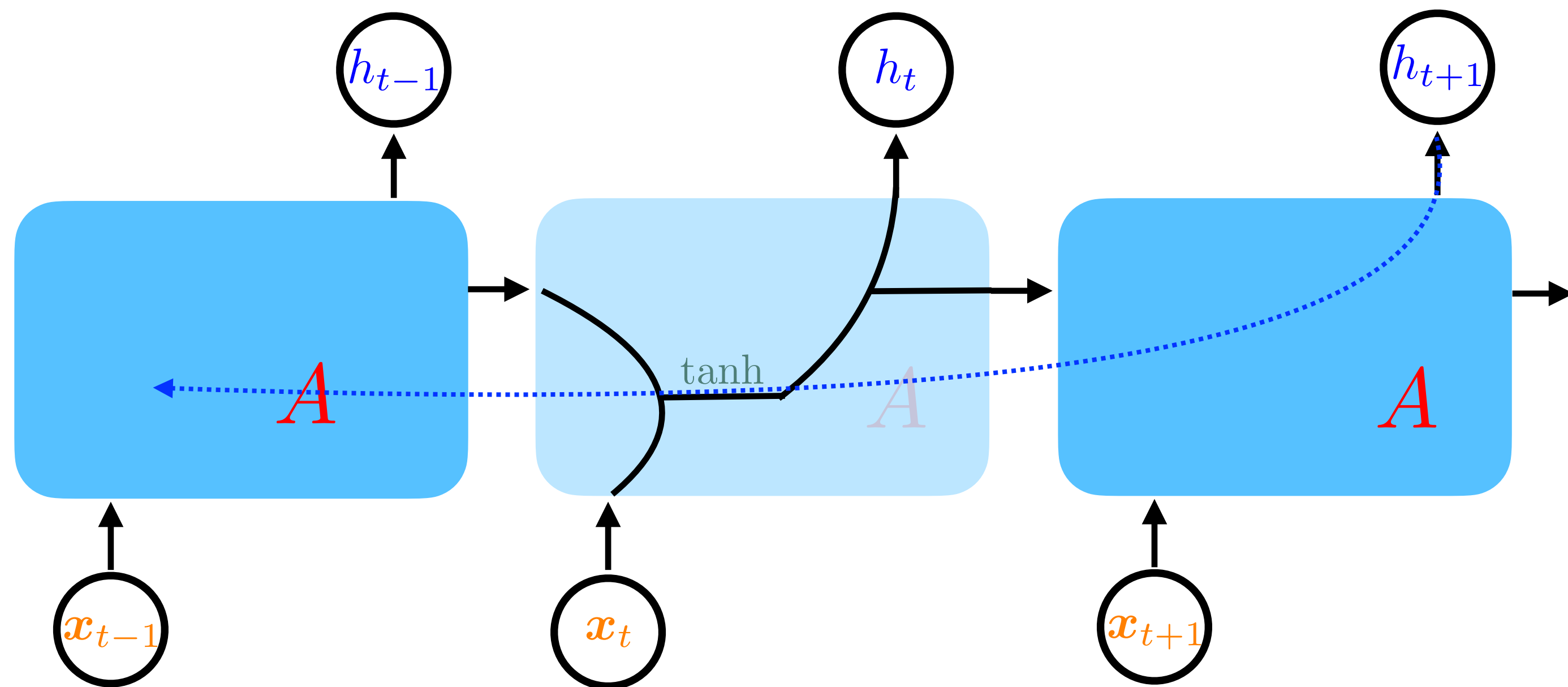**Lingpeng Kong**

**Department of Computer Science, The University of Hong Kong**

# Recurrent Neural Network



Yesterday          ,          troubles          which          seemed          so

Possibly many steps [O(N)] steps before "yesterday" and "seemed" interact.

# Vanishing Gradient in RNNs



$$\boldsymbol{h} = f(\boldsymbol{z})$$

$$\frac{\partial s}{\partial \boldsymbol{z}} = \frac{\partial s}{\partial \boldsymbol{h}} \frac{\partial \boldsymbol{h}}{\partial \boldsymbol{z}} \qquad \frac{\partial \boldsymbol{h}}{\partial \boldsymbol{z}} \qquad \frac{\partial s}{\partial \boldsymbol{h}}$$
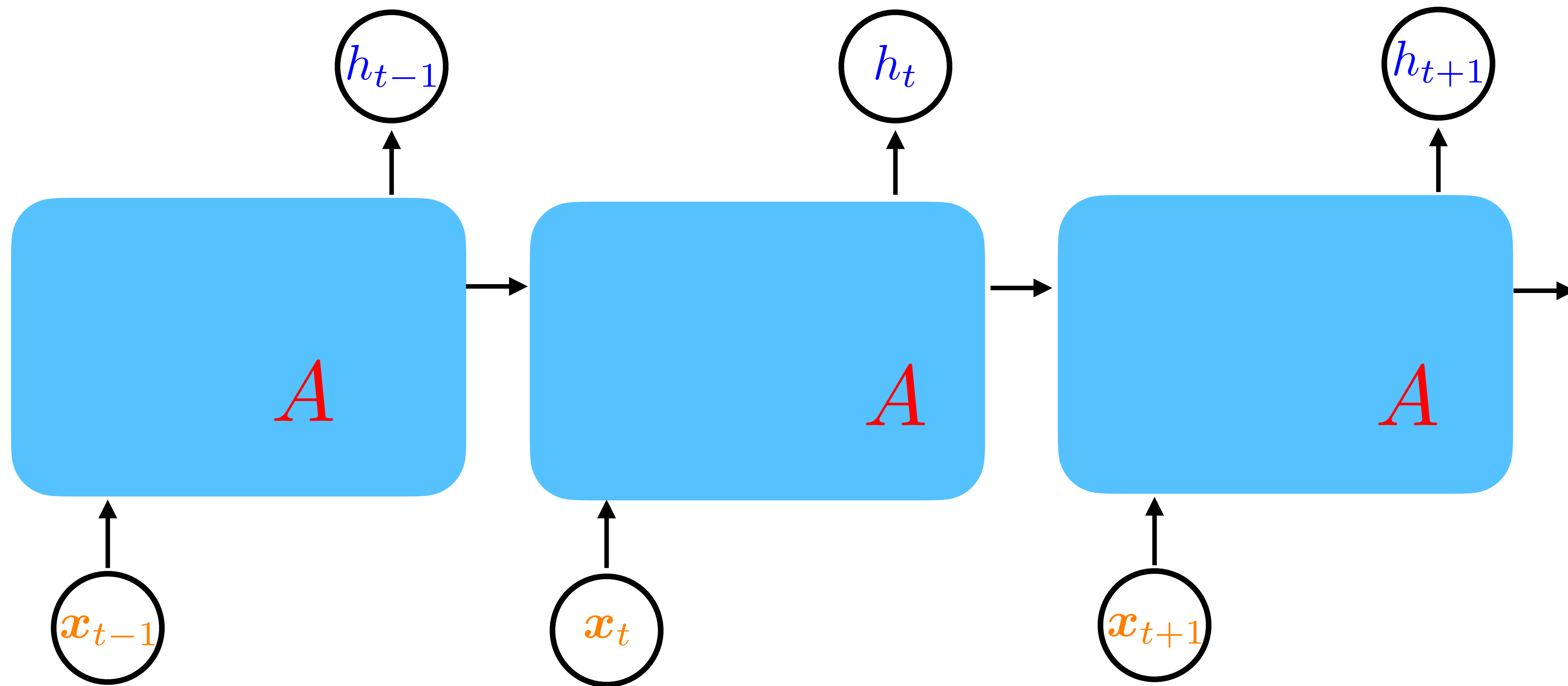
Gradient Flow Direction

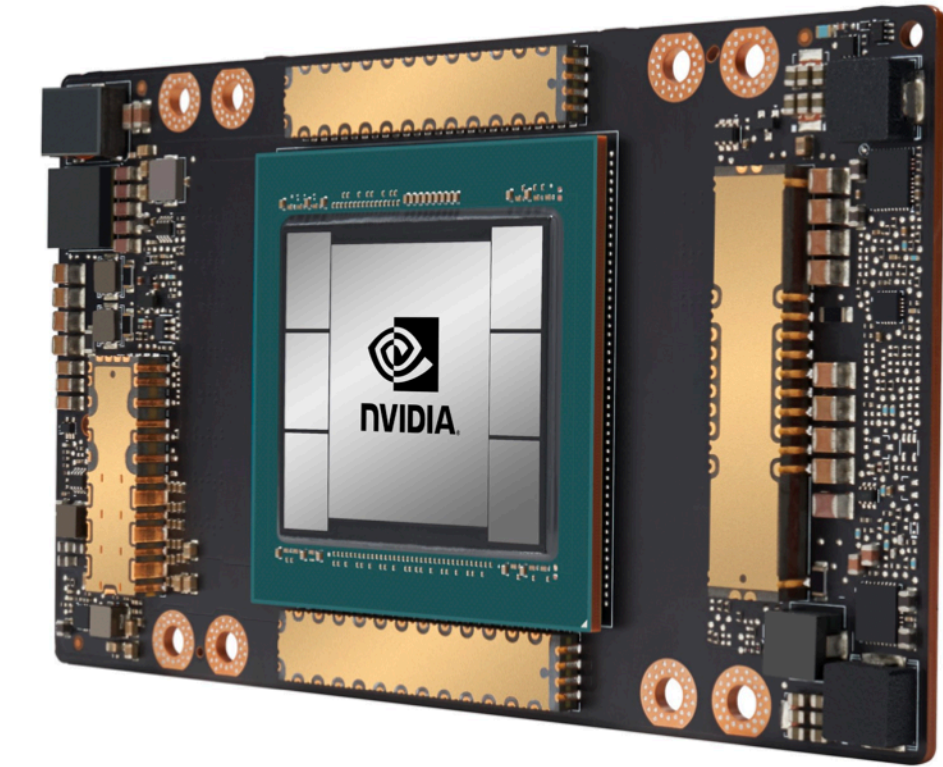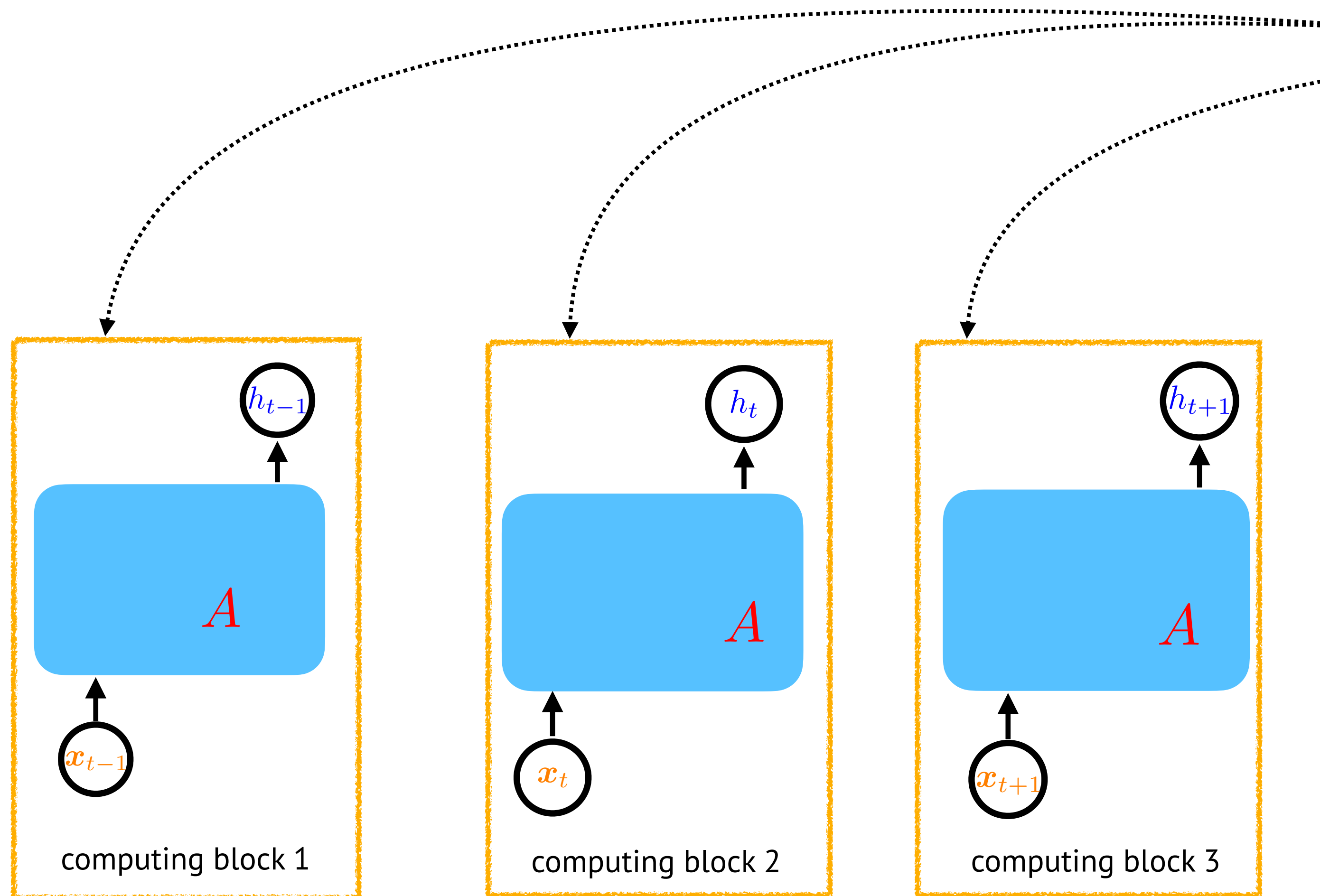In general, the longer the path, the smaller the gradient signal.

# Bidirectional Recurrent Neural Network

# Sequential Computation

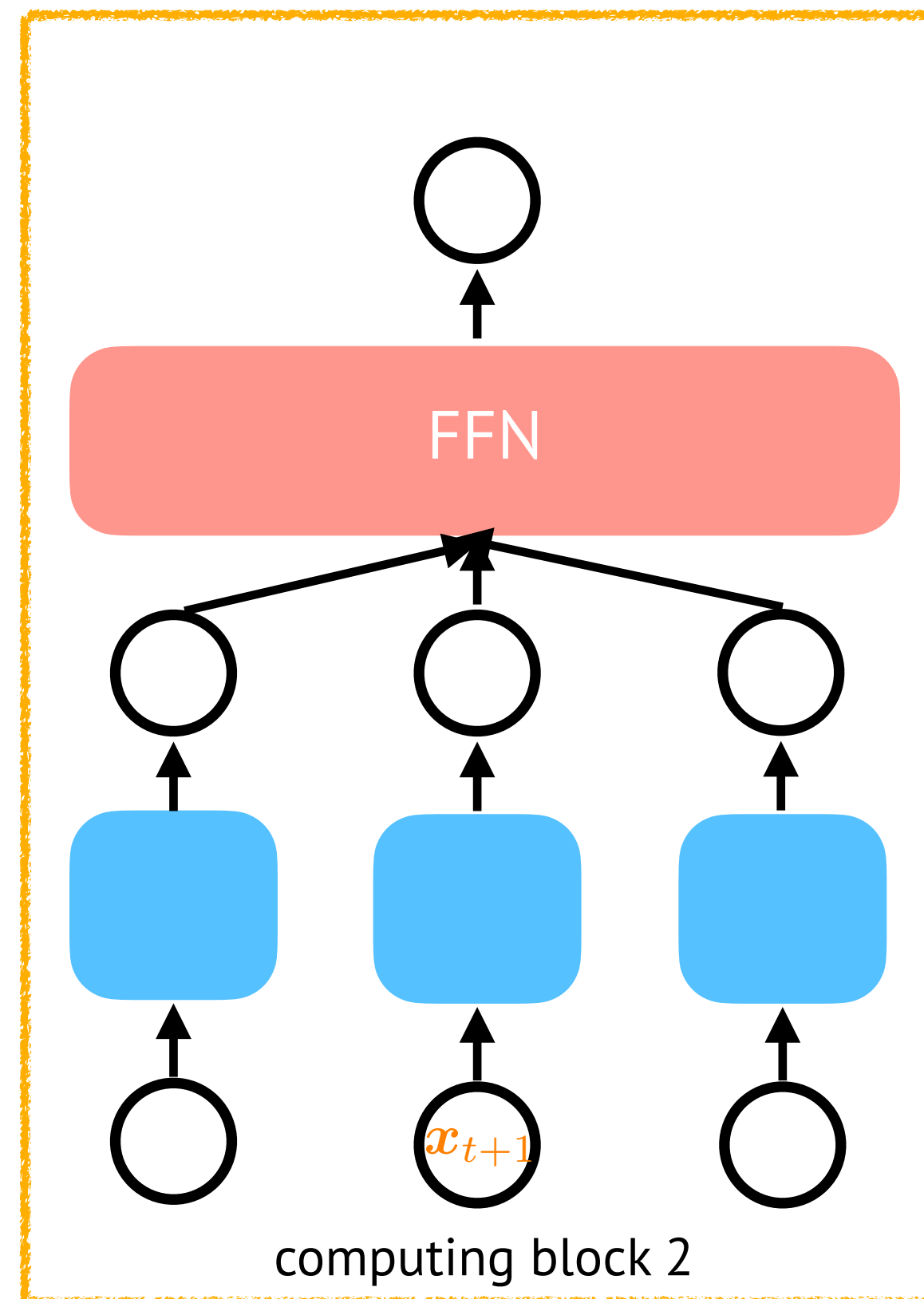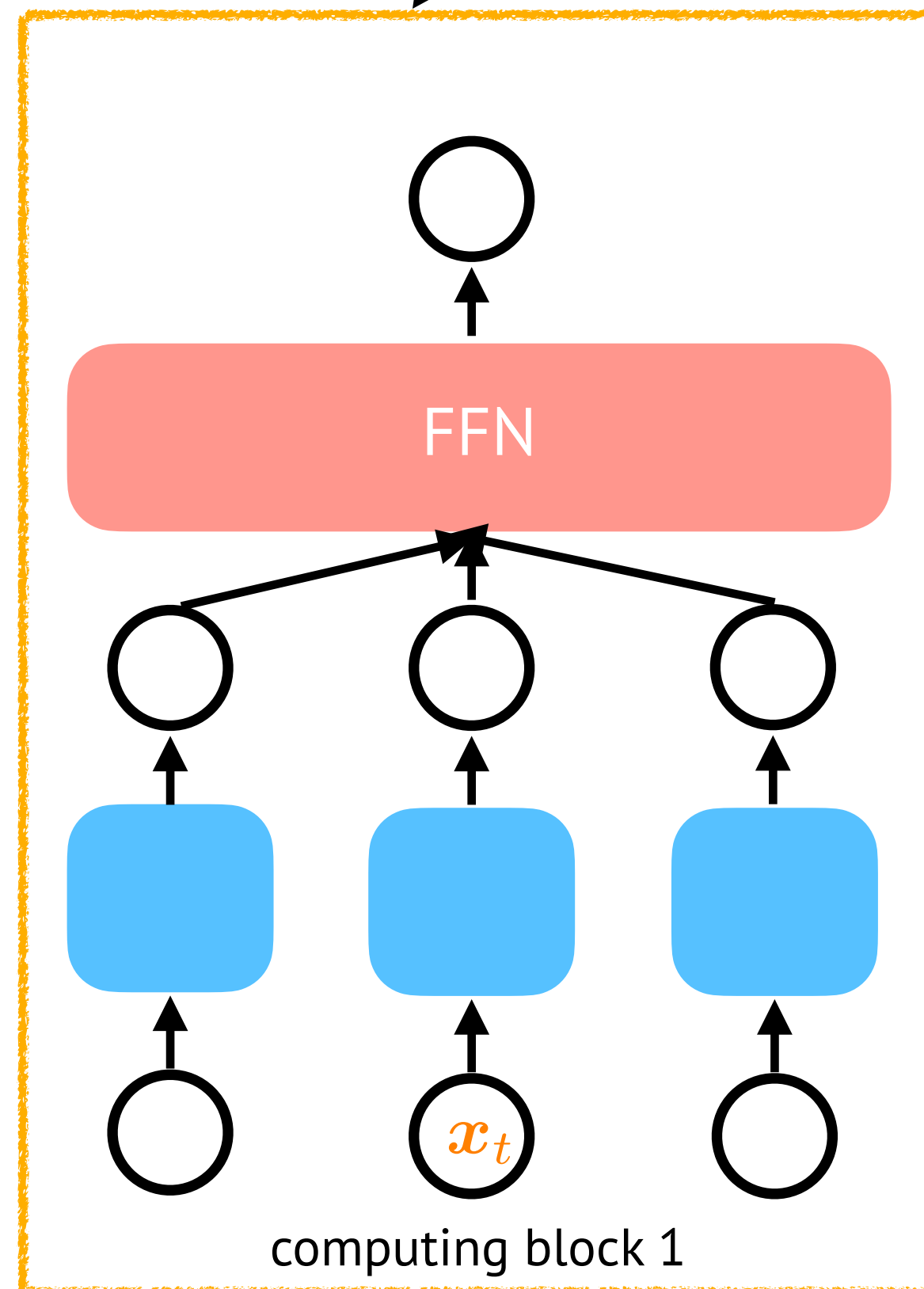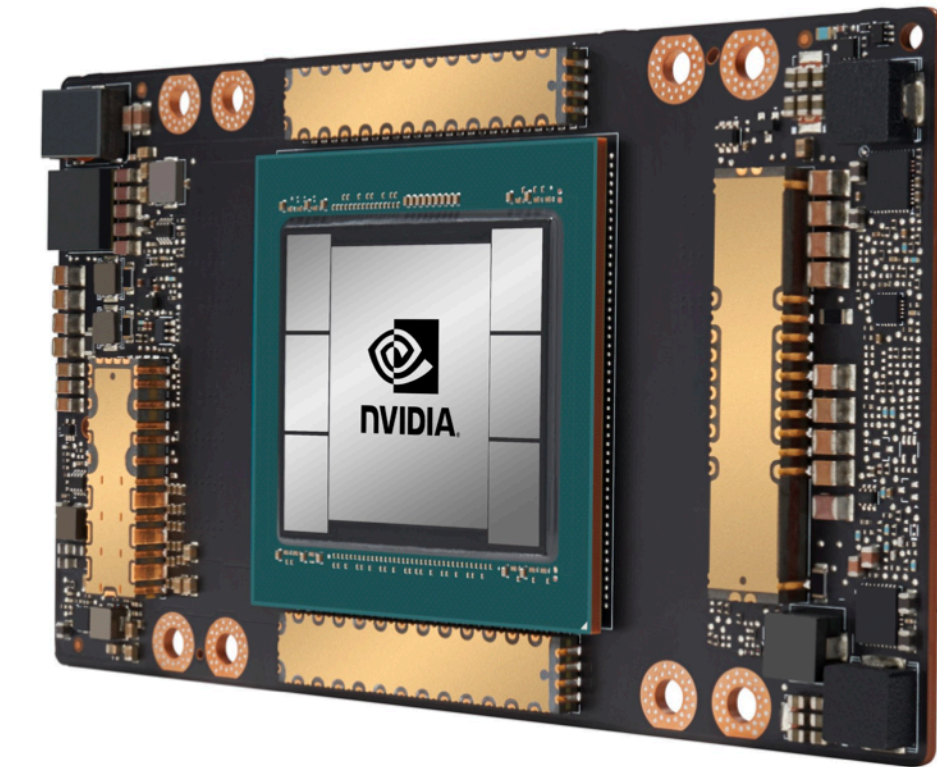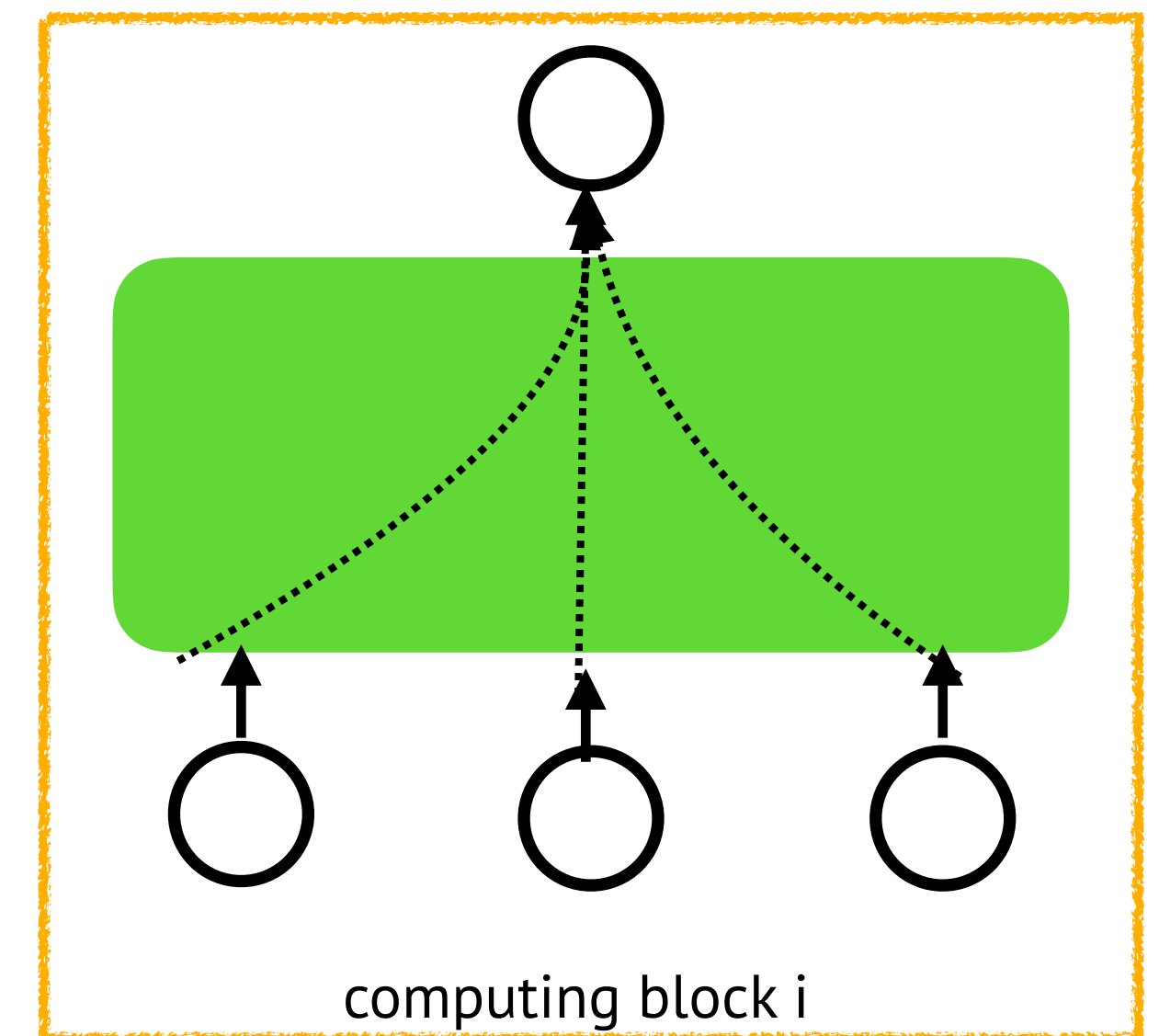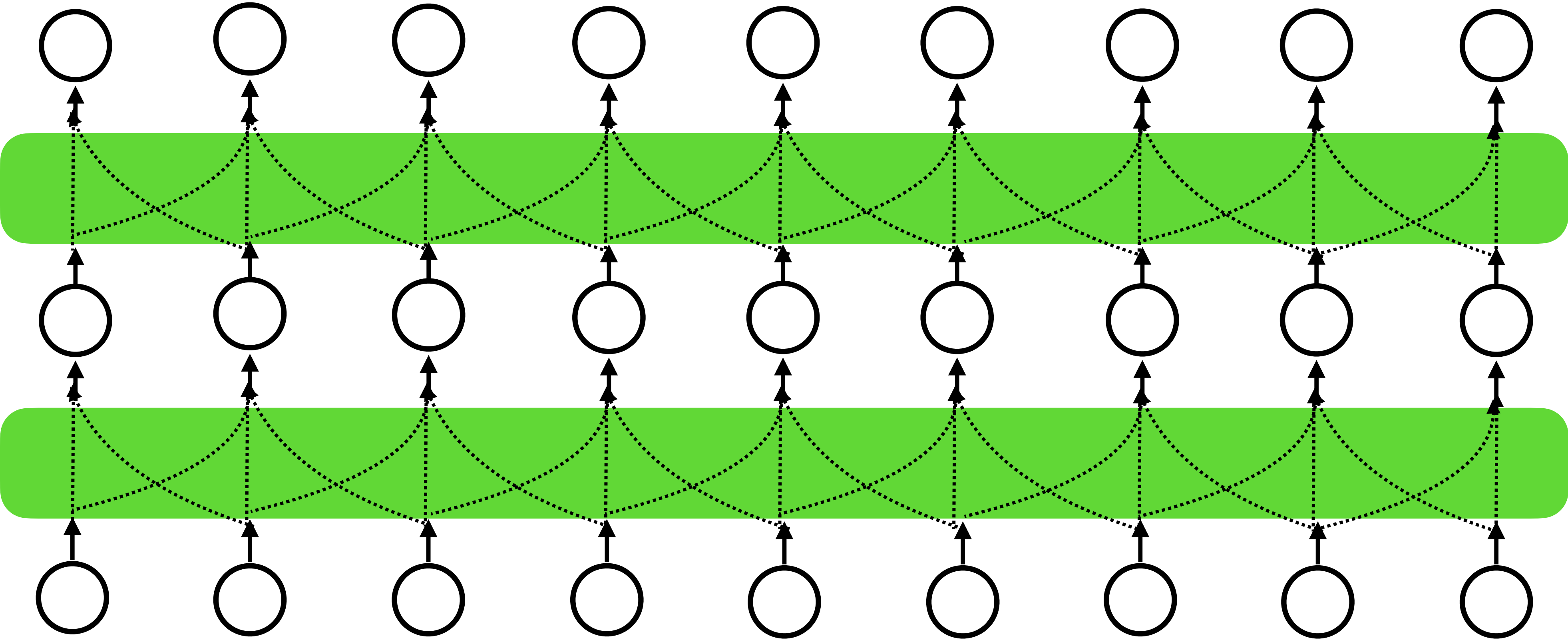# Parallel Computing?



GPU loves parallel computing blocks!

$h_{t-1}$

$h_t$

$h_{t+1}$

$A$

$A$

$A$

$\boldsymbol{x}_{t-1}$

$\boldsymbol{x}_t$

$\boldsymbol{x}_{t+1}$

computing block 1

computing block 2

computing block 3

. . . . . .

# Parallel Computing?



computing block 1

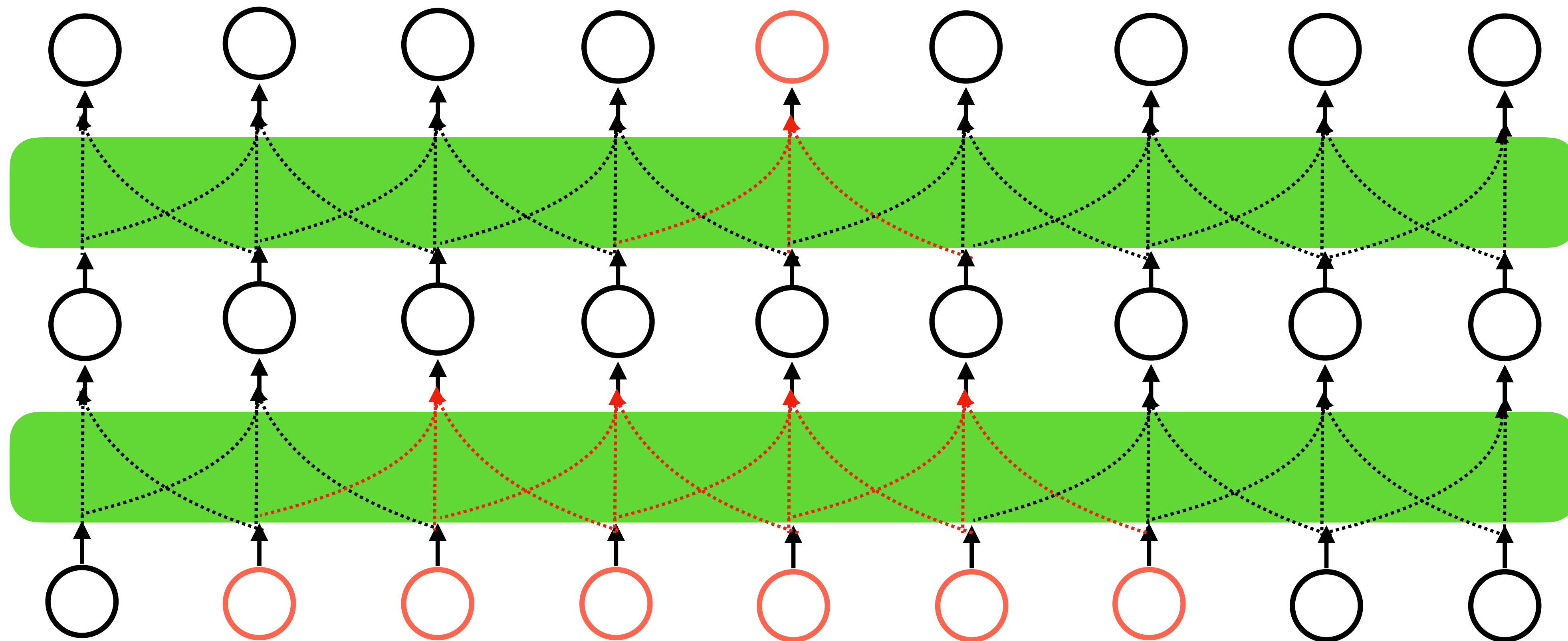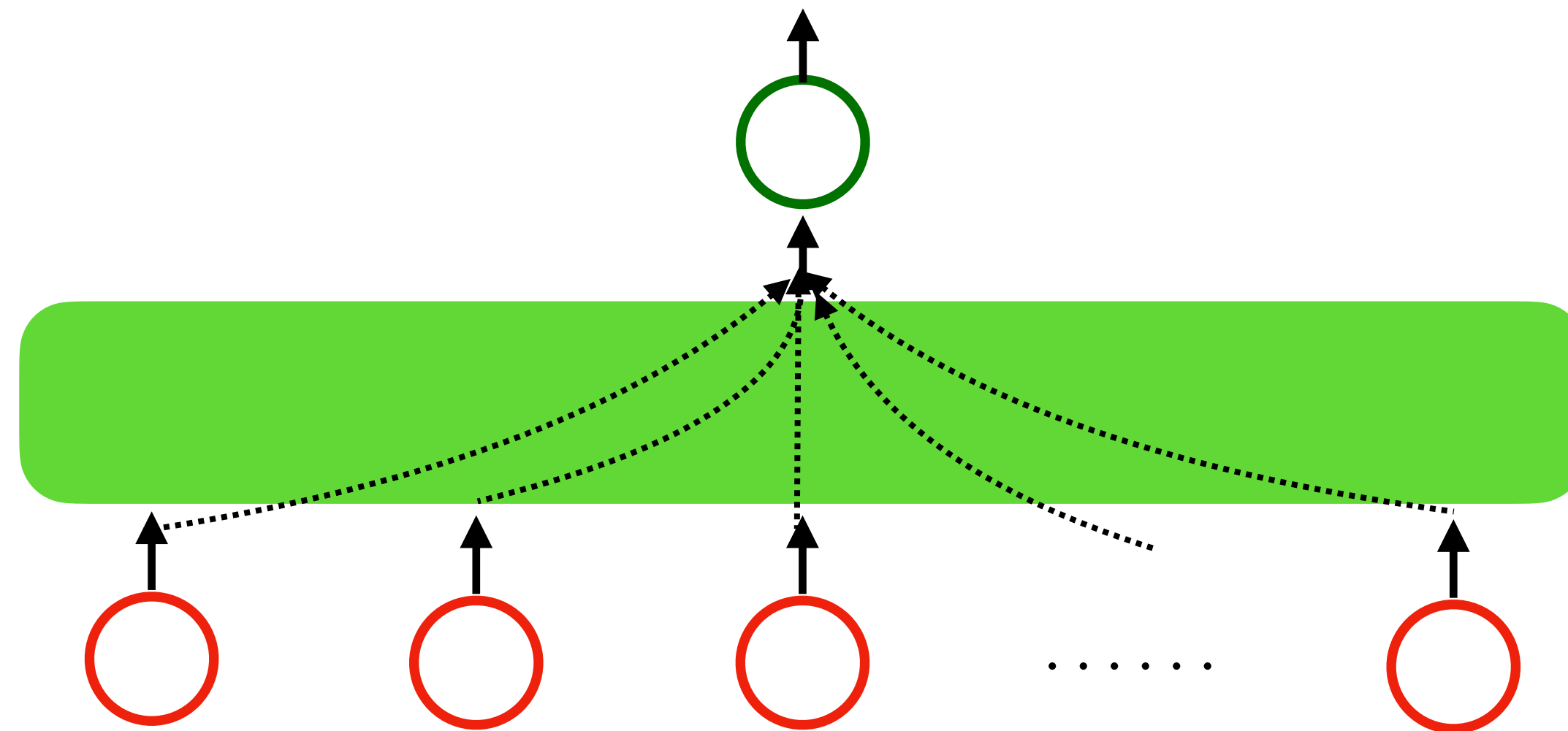computing block 2

computing block i

# Convolution Style Models

# Convolution Style Models

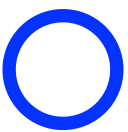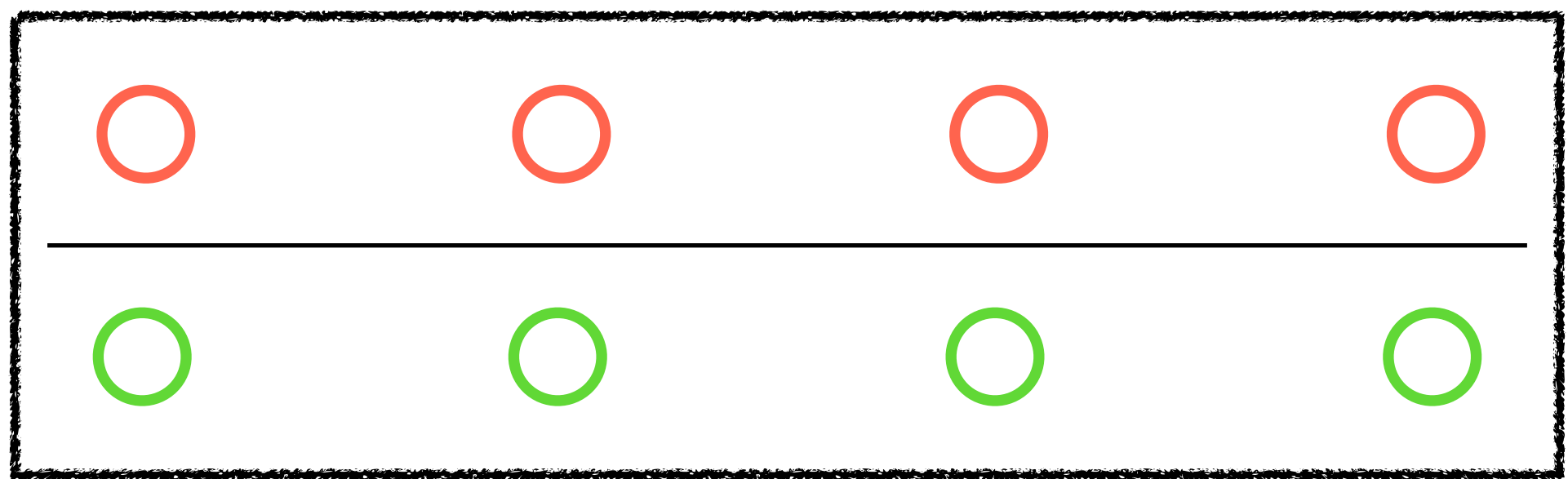# Considering the full sequence as context



How can we achieve this?

# Dot-Product-Softmax Attention

Query

Memory (key-value pairs)

$$\text{softmax}\left(\begin{array}{c} \mathbf{q} \cdot \mathbf{k}_1 \\ \mathbf{q} \cdot \mathbf{k}_2 \\ \mathbf{q} \cdot \mathbf{k}_3 \\ \mathbf{q} \cdot \mathbf{k}_4 \end{array}\right) \rightarrow \begin{bmatrix} 0.6 \\ 0.1 \\ 0.2 \\ 0.1 \end{bmatrix}$$

$\mathbf{q} \cdot \mathbf{k}_1$

$\mathbf{q} \cdot \mathbf{k}_2$

$\mathbf{q} \cdot \mathbf{k}_3$

$\mathbf{q} \cdot \mathbf{k}_4$

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

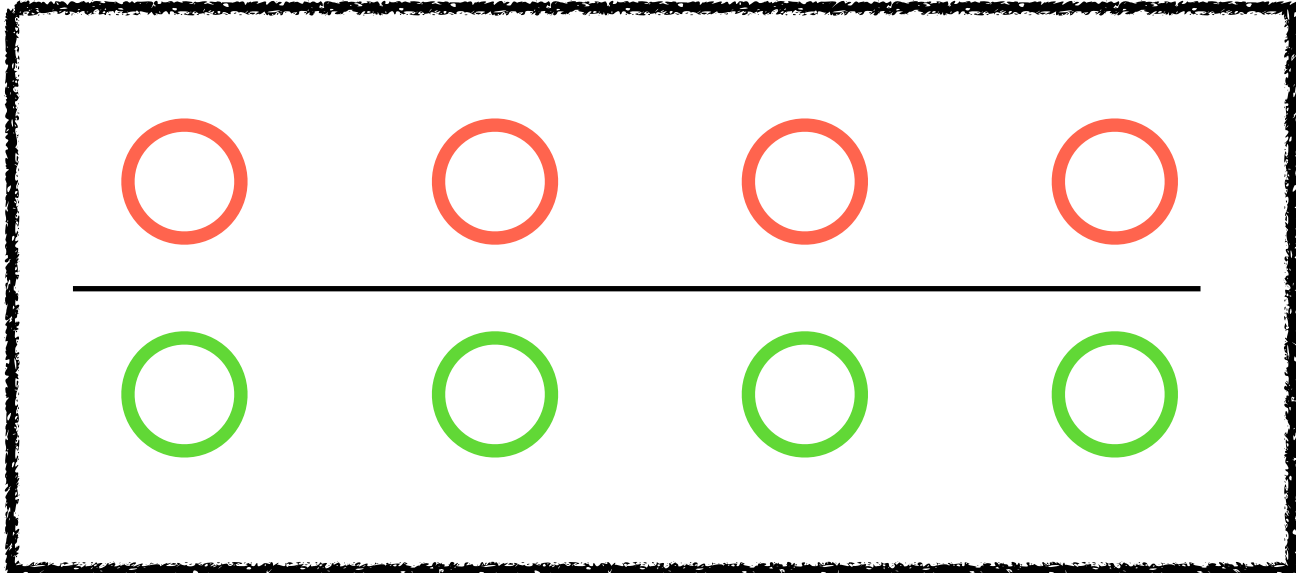$= \bigcirc$ context vector $\quad \mathbf{c}$

# Considering the full sequence as context

# Attention Mechanism

Query

Memory (key-value pairs)

. . . . . .

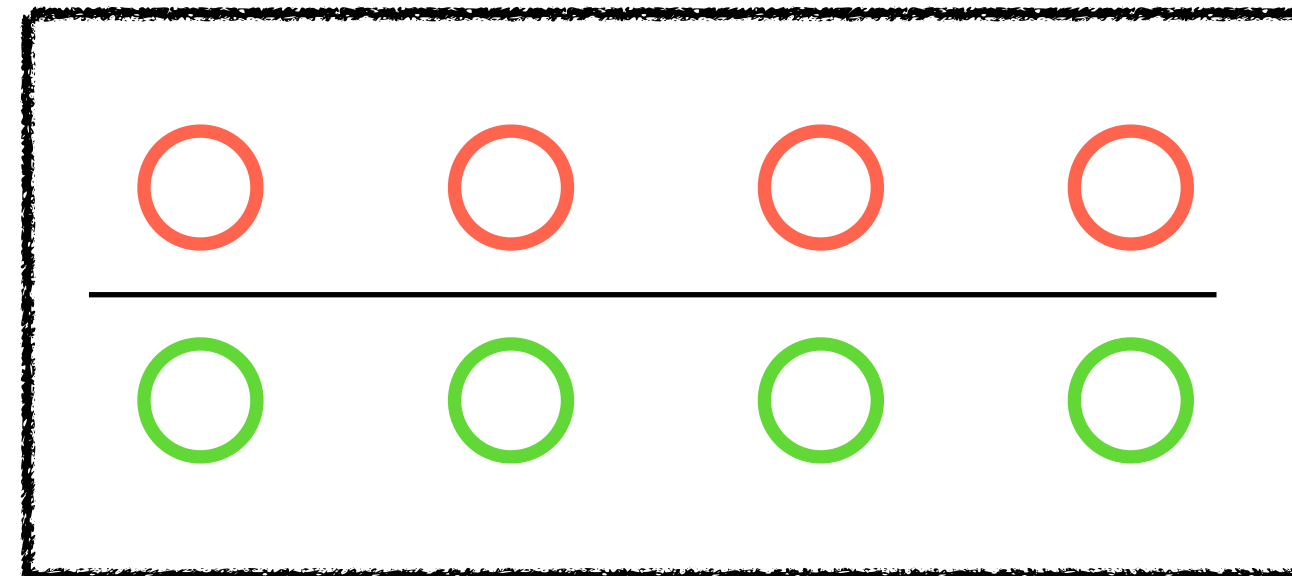# Attention Mechanism

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

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

Query

Memory (key-value pairs)
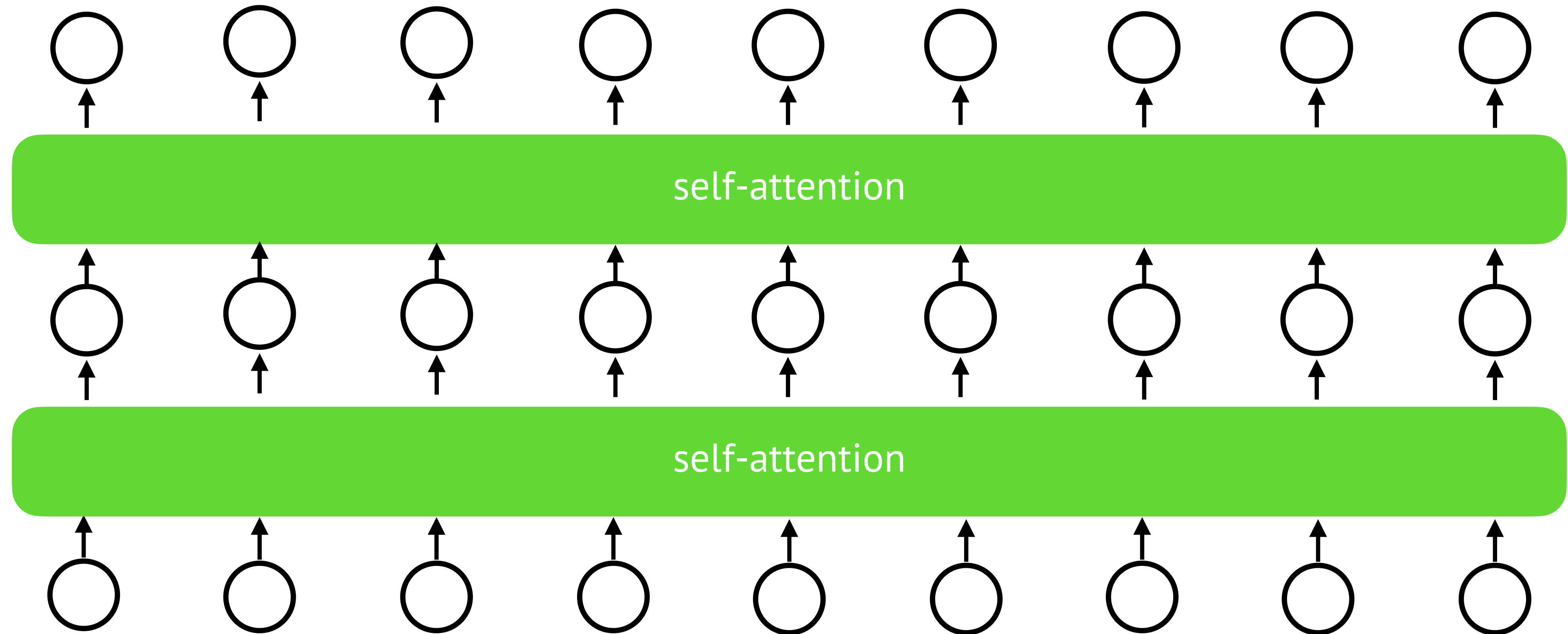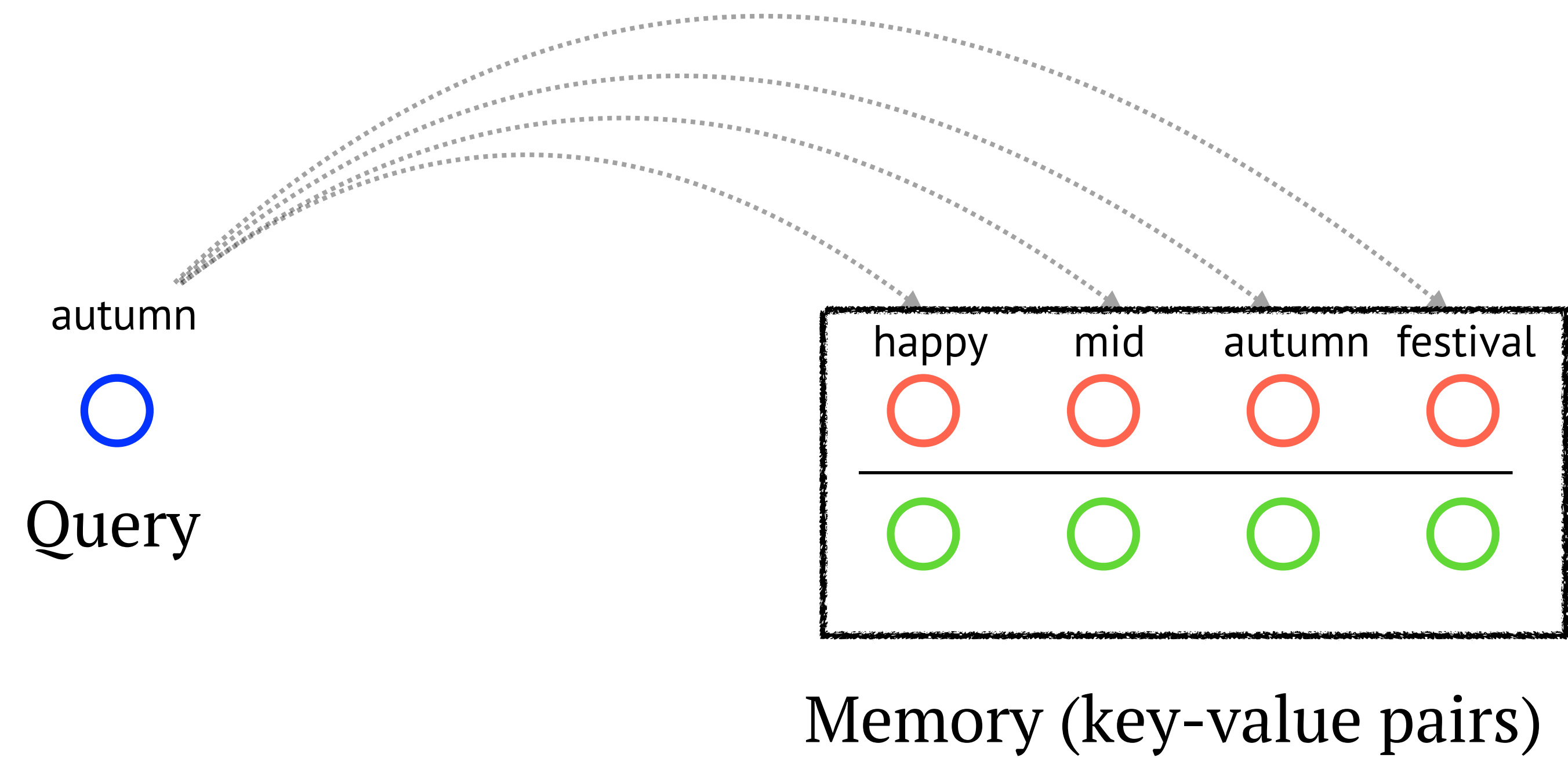
· · · · · ·

# Self-attention



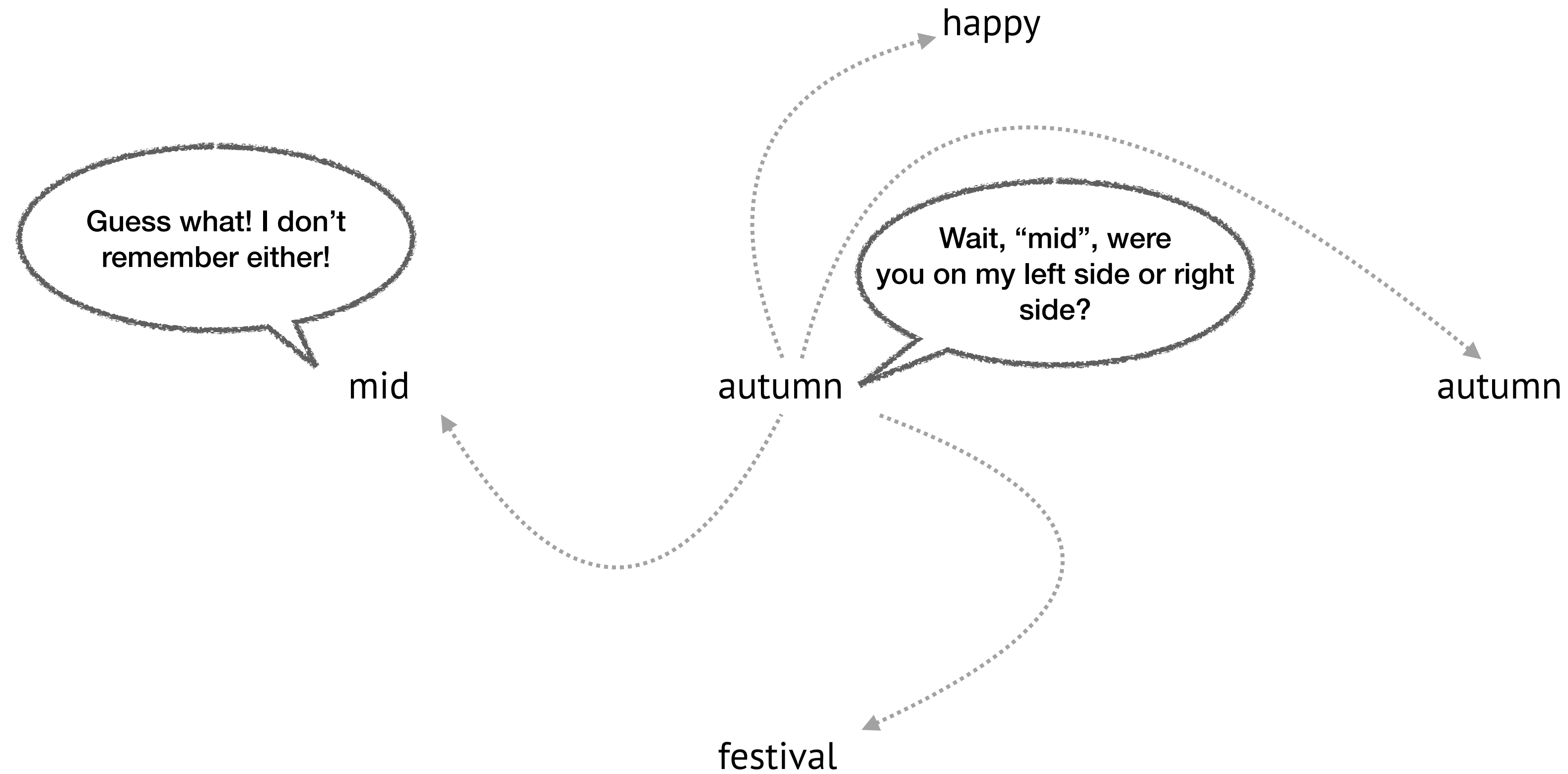This is almost transformer — except a few things.

# Transformer (almost)

# Self-attention in Transformer

# Self-attention in Transformer

# Positional Embeddings