

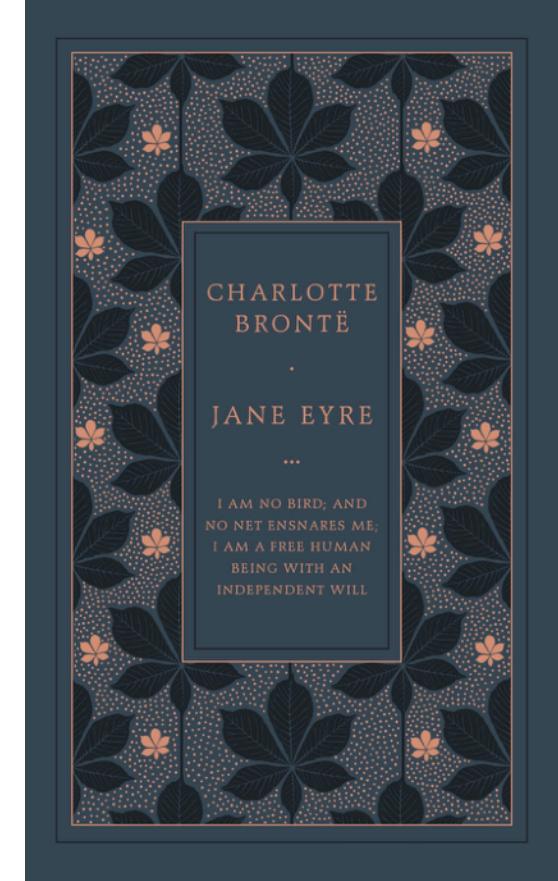
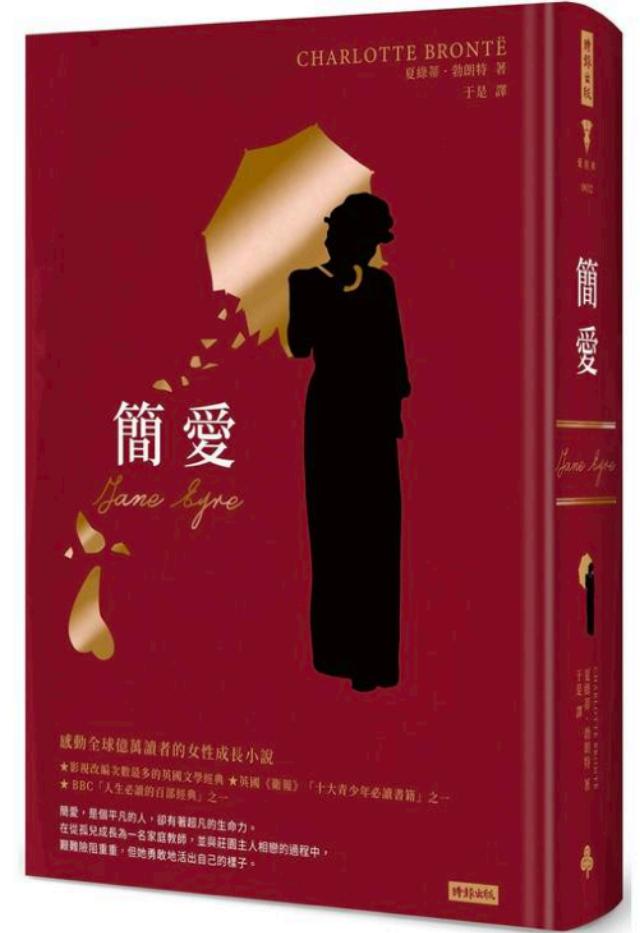
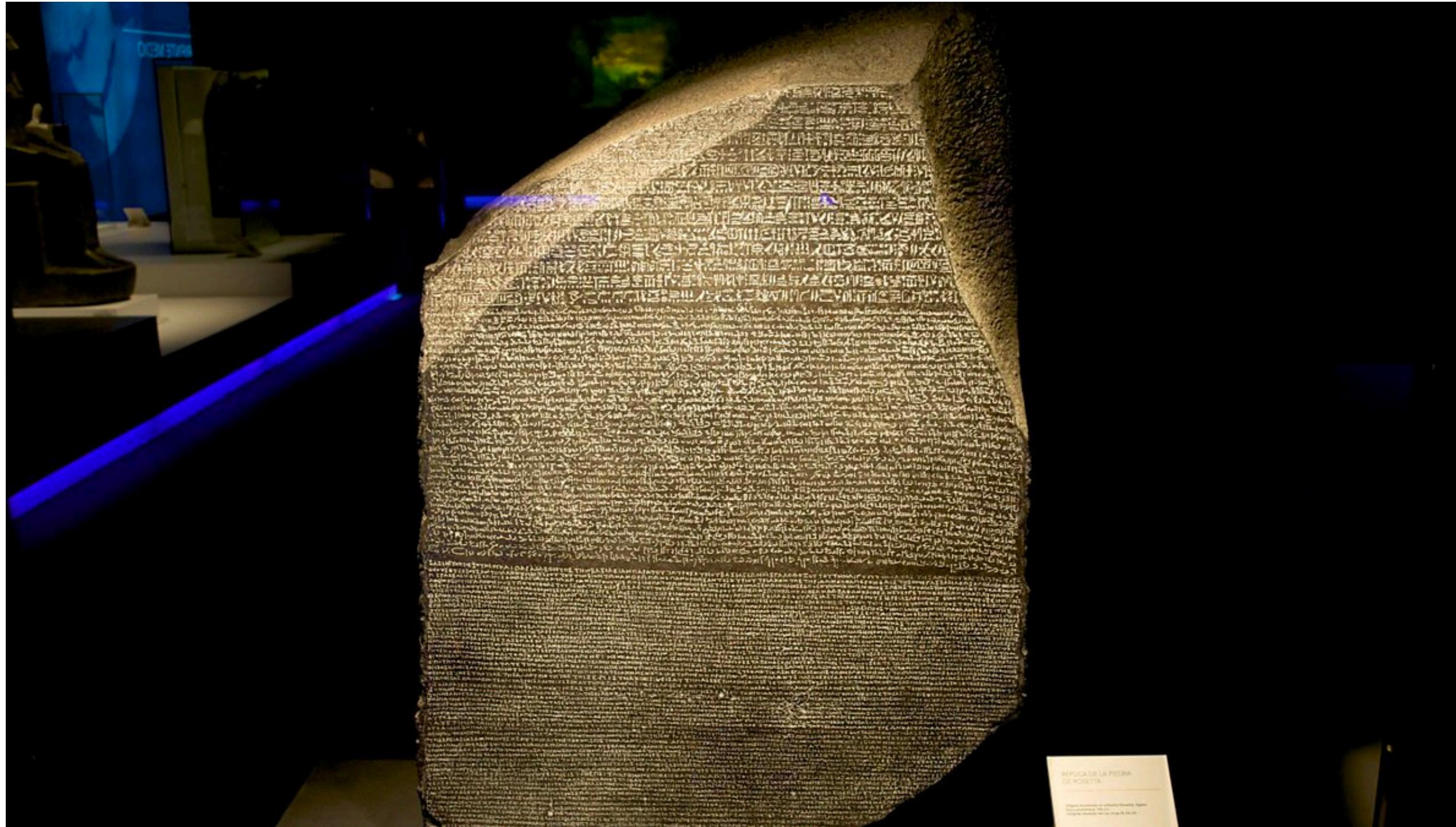
# Attention Mechanism

COMP3361 – Week 4

Lingpeng Kong

Department of Computer Science, The University of Hong Kong

# Alignment in Machine Translation

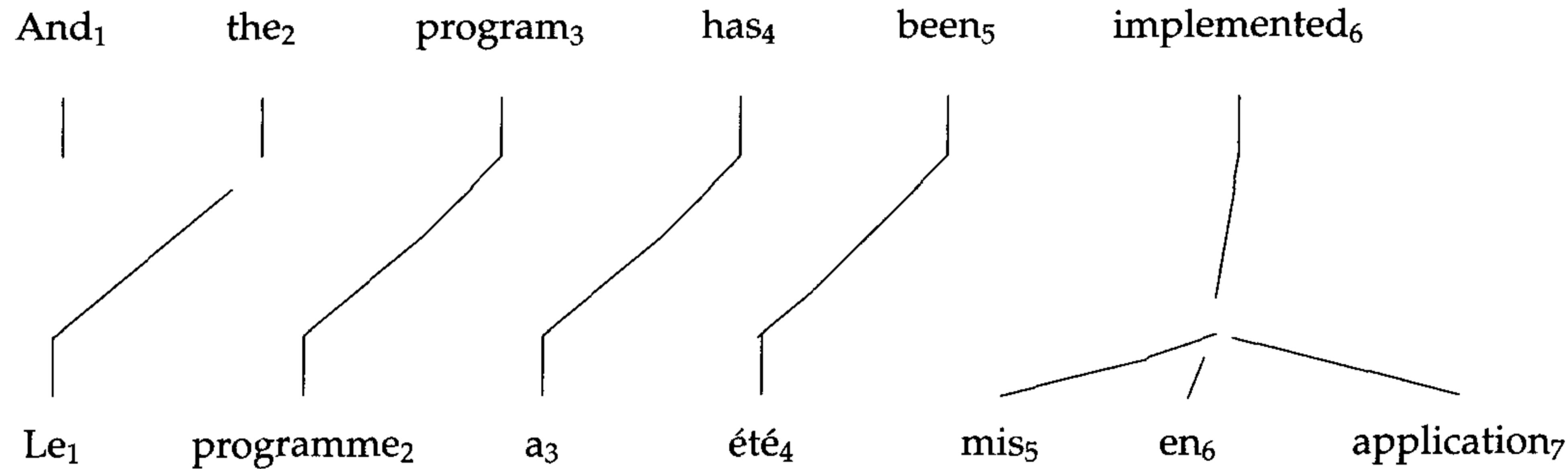


Parallel Corpus

$$p(y \mid x)$$

↑      ↑  
target   source

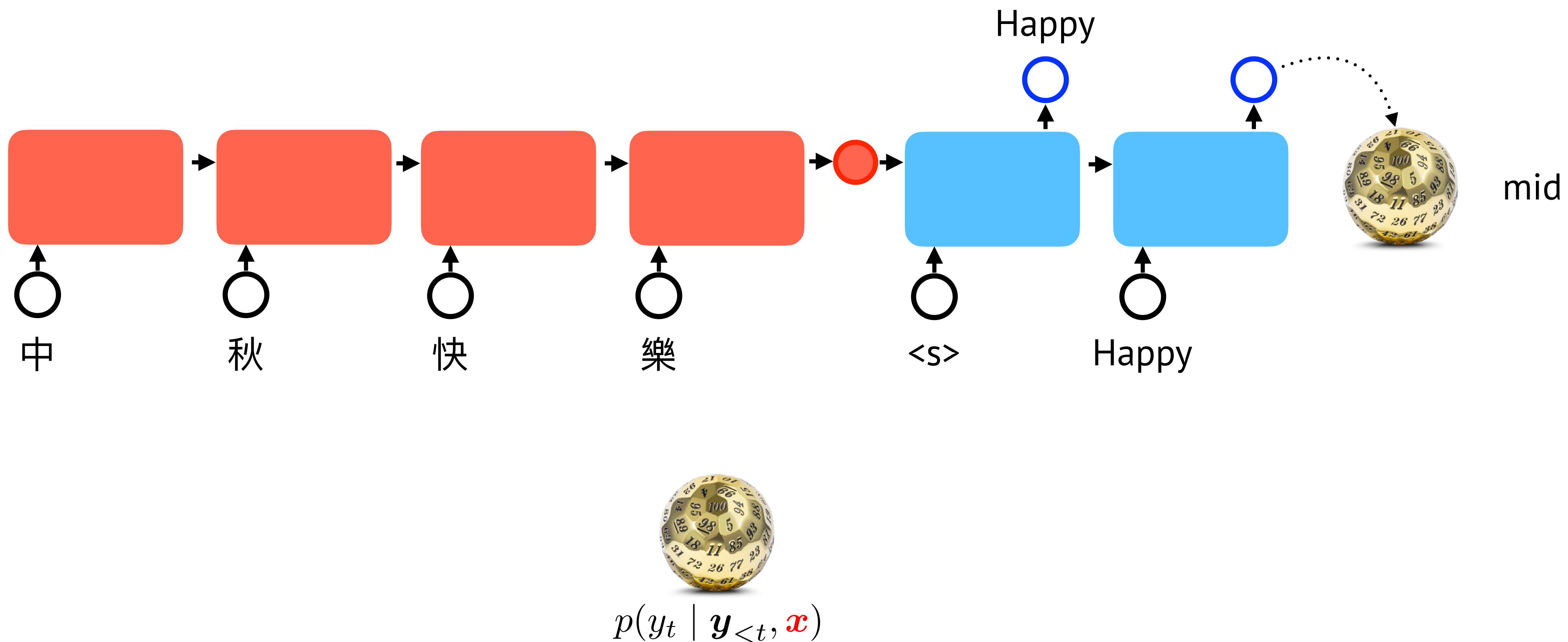
# Alignment in Machine Translation



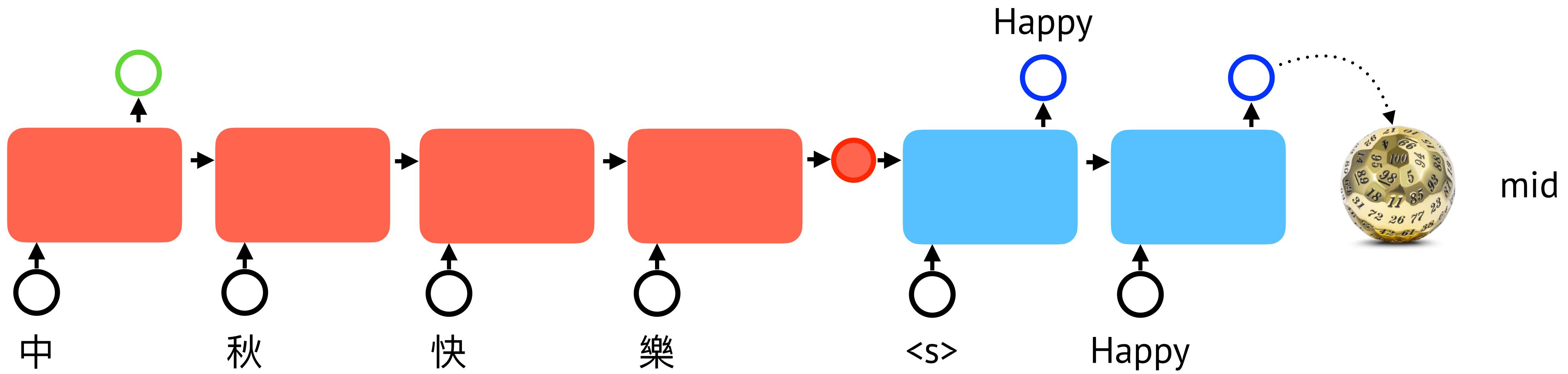
Some words might have no “counter-part”.

Alignment can be many-to-one (or one-to-many).

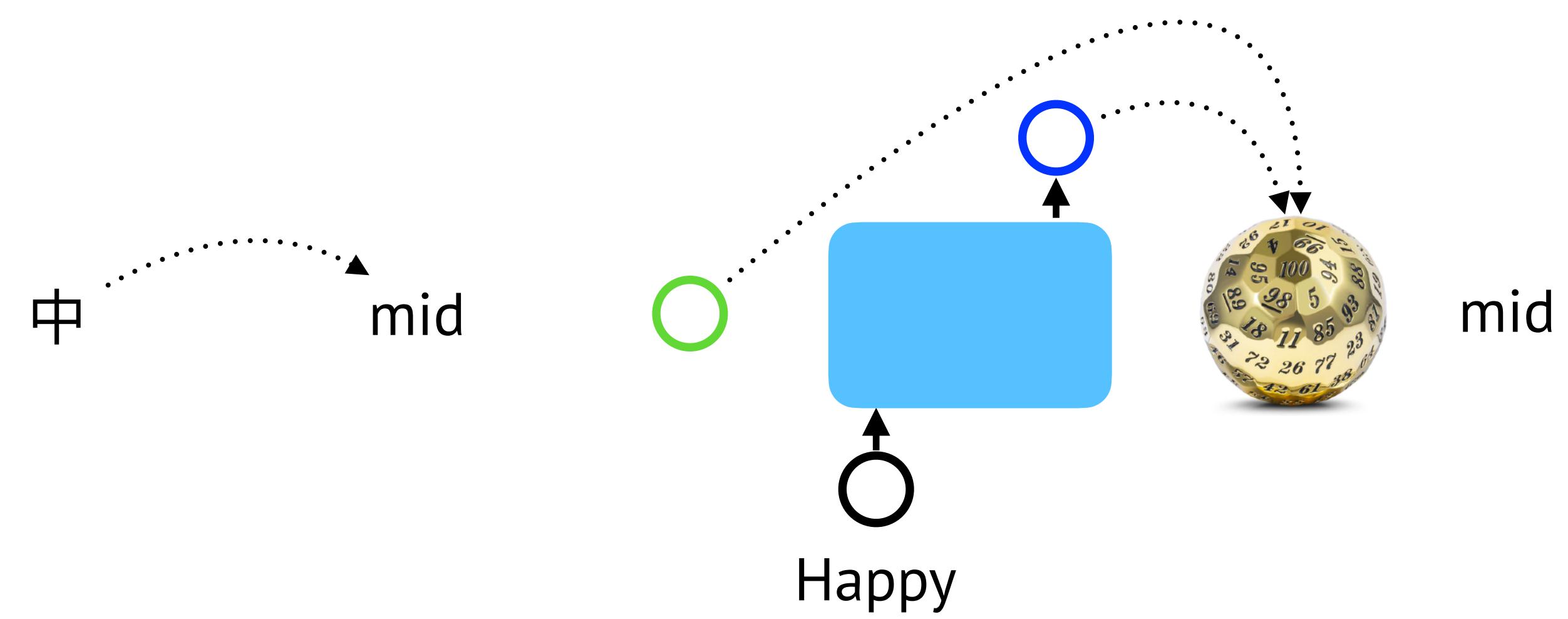
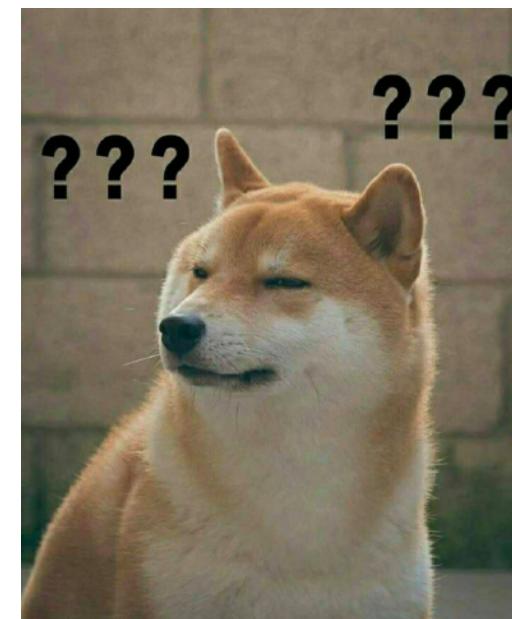
# Sequence to Sequence Model



# Sequence to Sequence Model

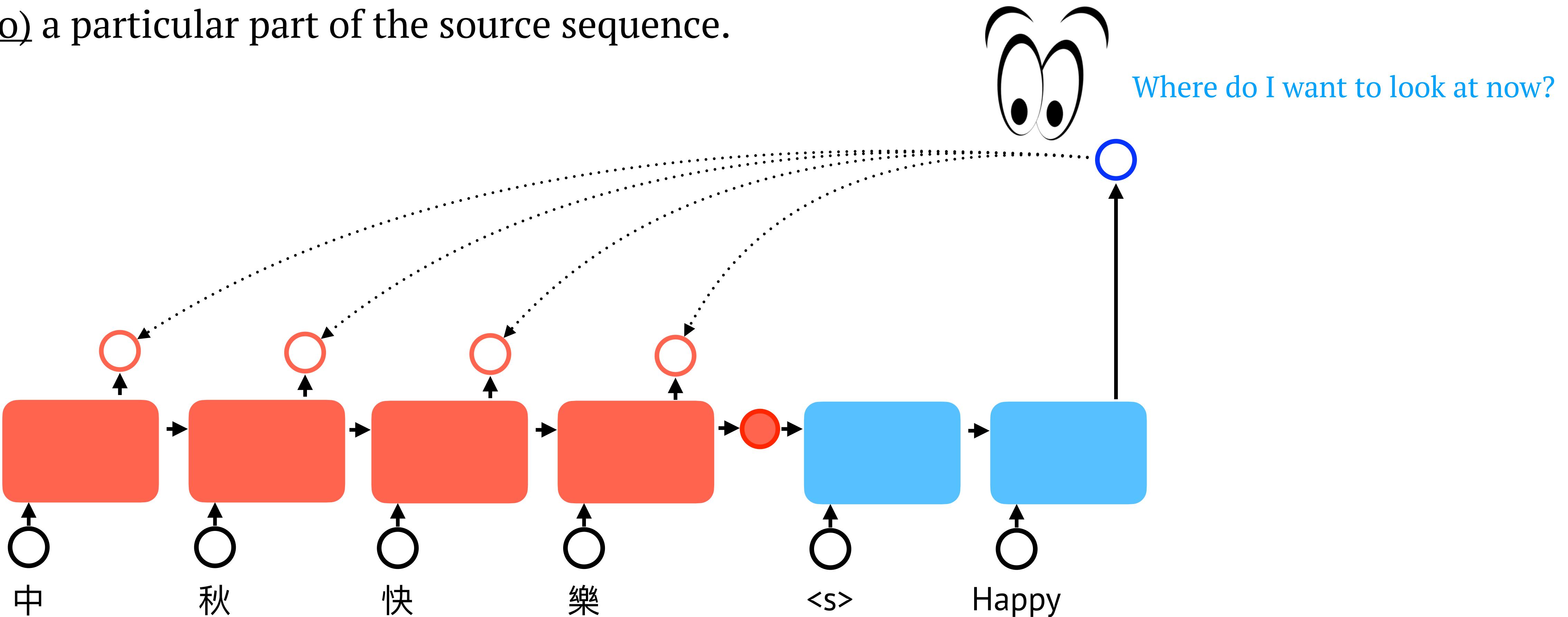


$$p(y_t \mid y_{<t}, \mathbf{x})$$



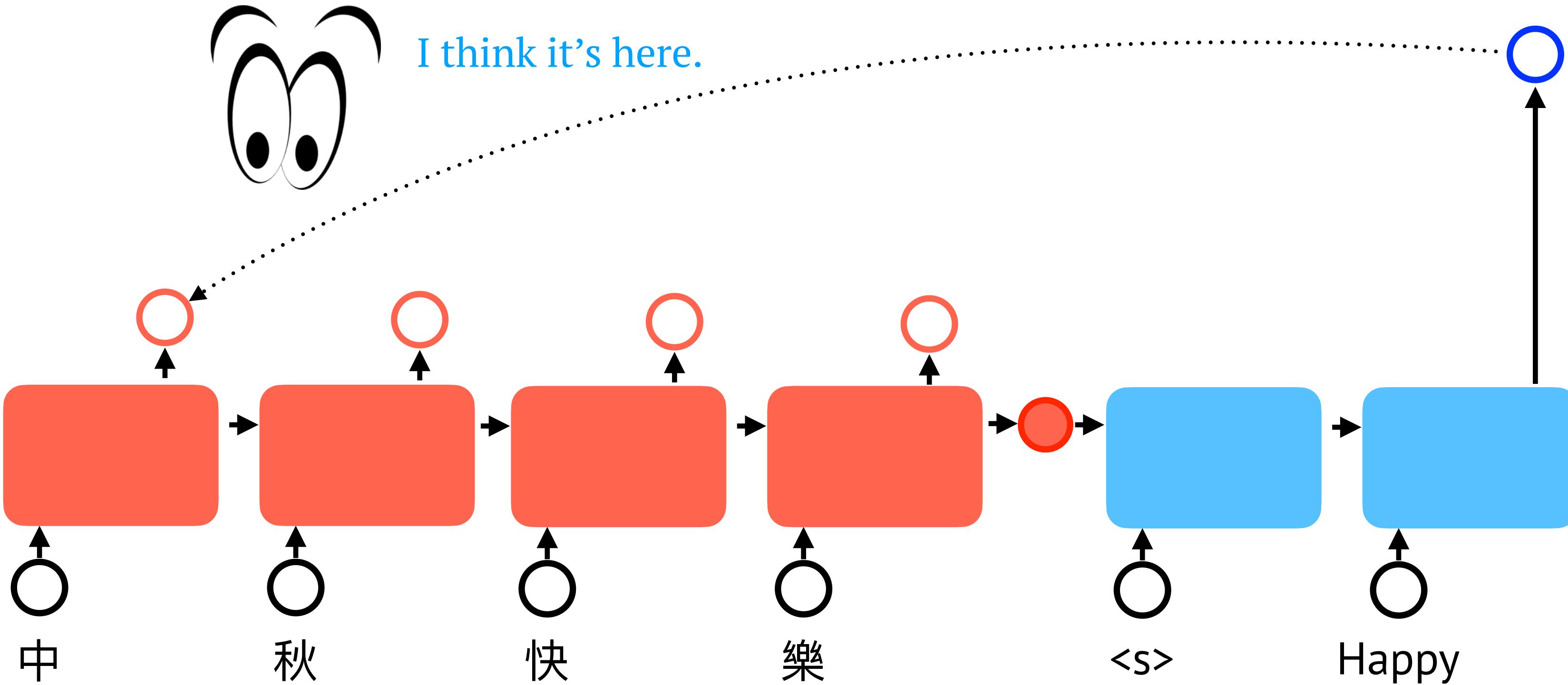
# Attention Mechanism

Use direct connection to the encoder to focus on (attend to) a particular part of the source sequence.



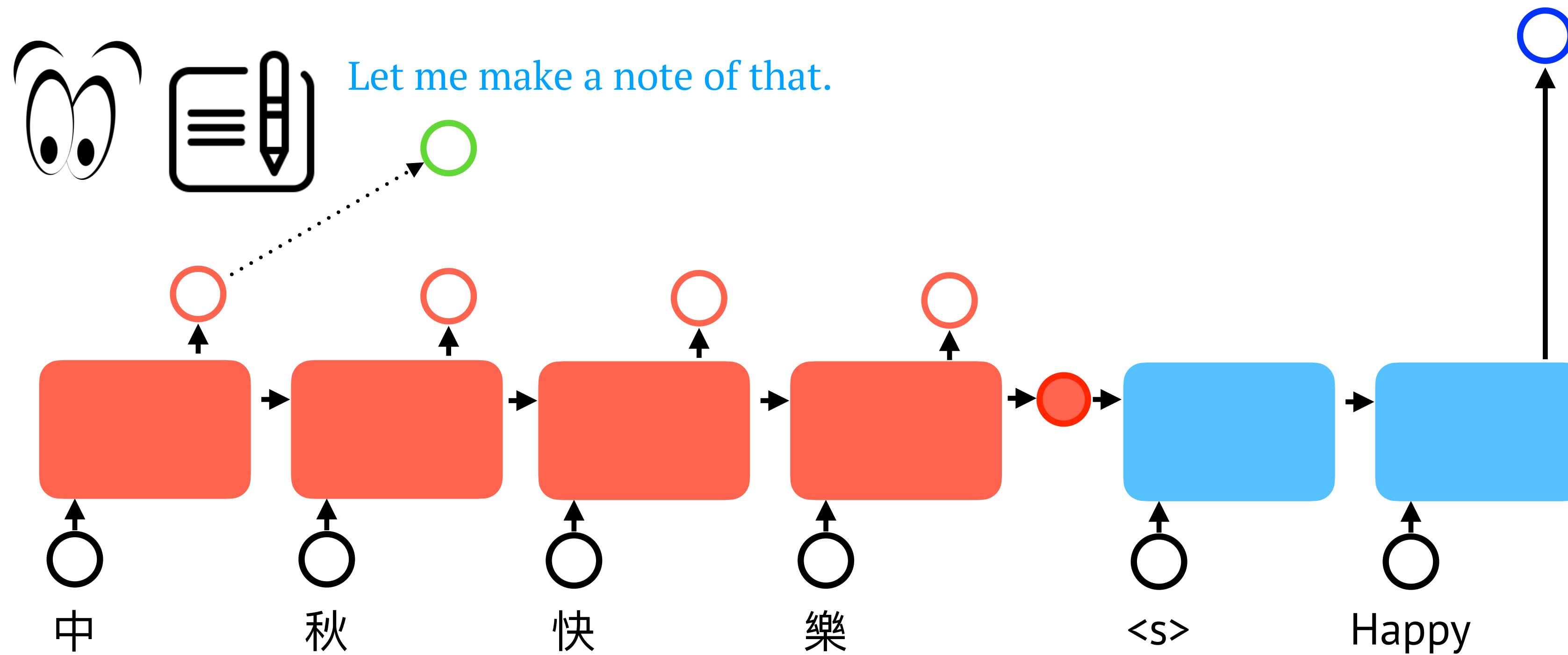
# Attention Mechanism

Use direct connection to the encoder to focus on (attend to) a particular part of the source sequence.



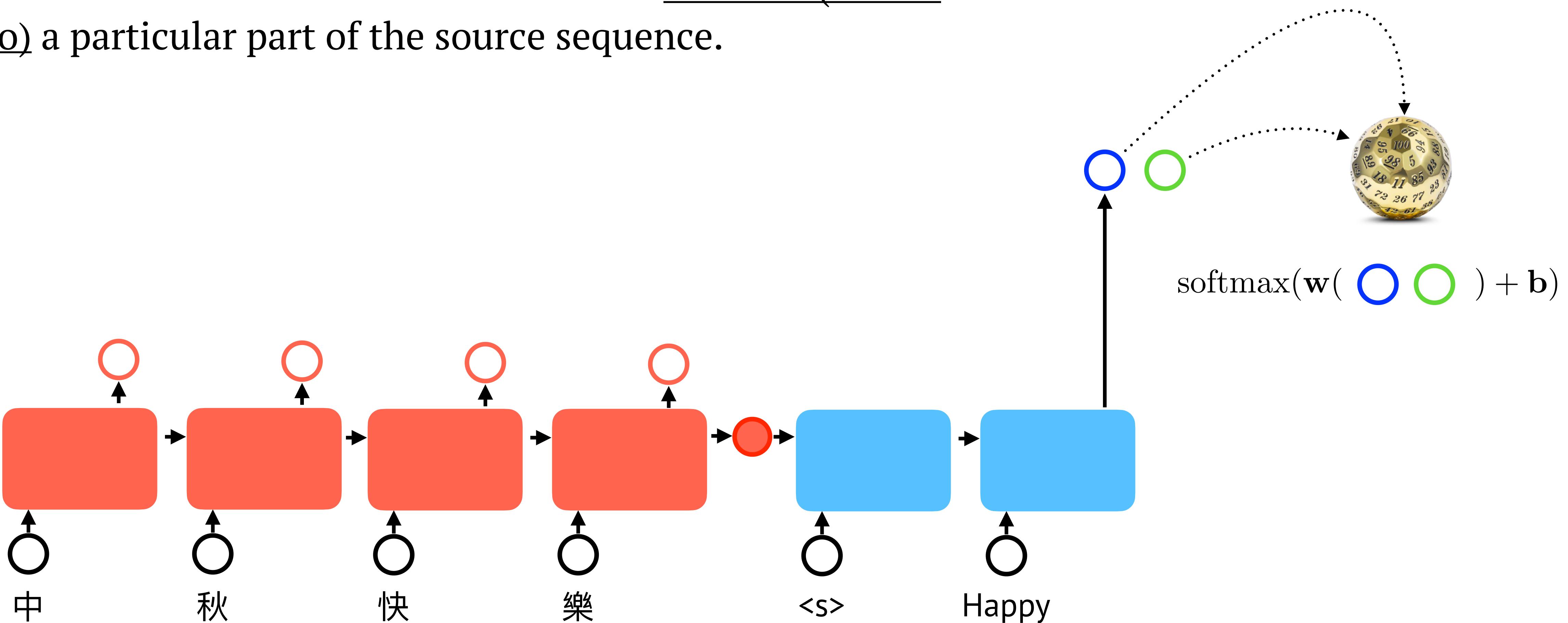
# Attention Mechanism

Use direct connection to the encoder to focus on (attend to) a particular part of the source sequence.

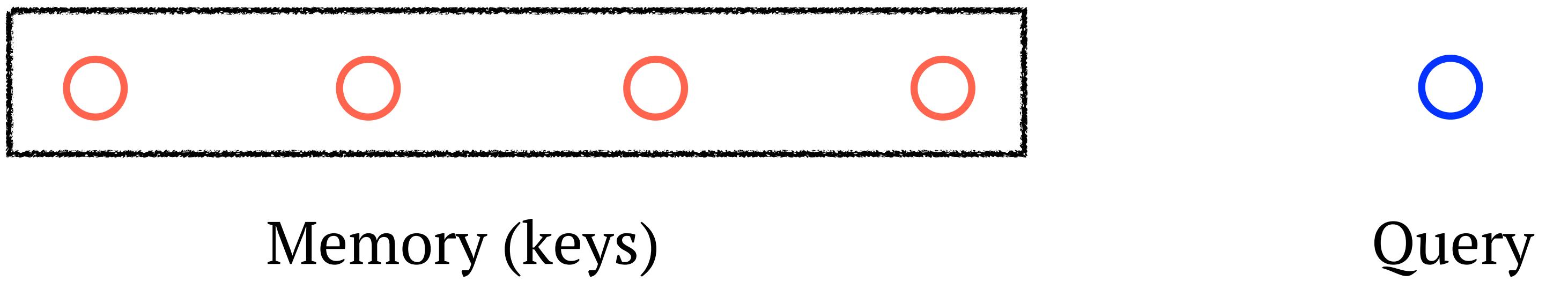


# Attention Mechanism

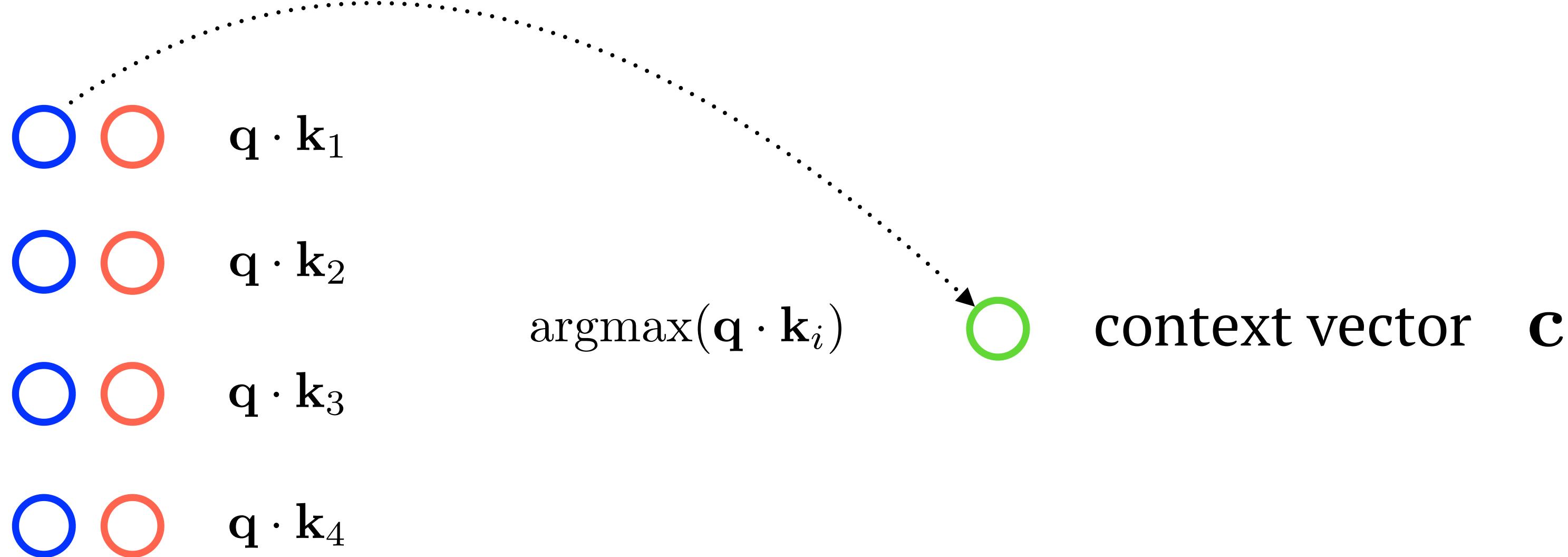
Use direct connection to the encoder to focus on (attend to) a particular part of the source sequence.



# Memory Abstraction



Task: Finding the most “relevant” item in the memory.



# Dot-Product-Softmax Attention



Query

Task: Finding the most “relevant” item in the memory.

$$\text{blue circle} \text{ } \text{red circle} \quad q \cdot k_1$$

$$\text{blue circle} \text{ } \text{red circle} \quad q \cdot k_2$$

$$\text{blue circle} \text{ } \text{red circle} \quad q \cdot k_3$$

$$\text{blue circle} \text{ } \text{red circle} \quad q \cdot k_4$$

$$q \cdot k_1$$

$$q \cdot k_2$$

$$q \cdot k_3$$

$$q \cdot k_4$$

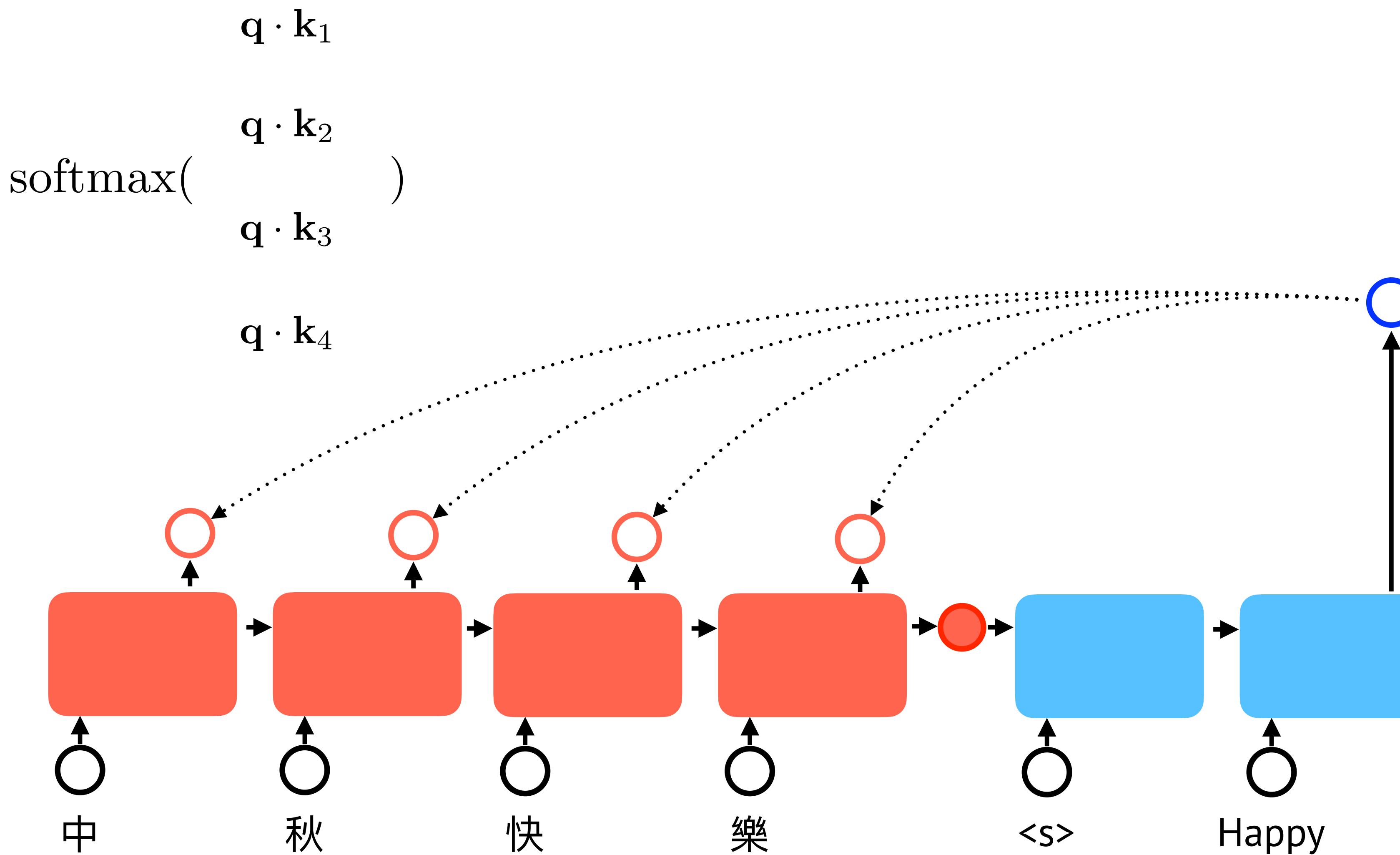
$$\text{softmax}(\quad \quad \quad ) \rightarrow$$

$$\begin{bmatrix} 0.6 \\ 0.1 \\ 0.2 \\ 0.1 \end{bmatrix}$$

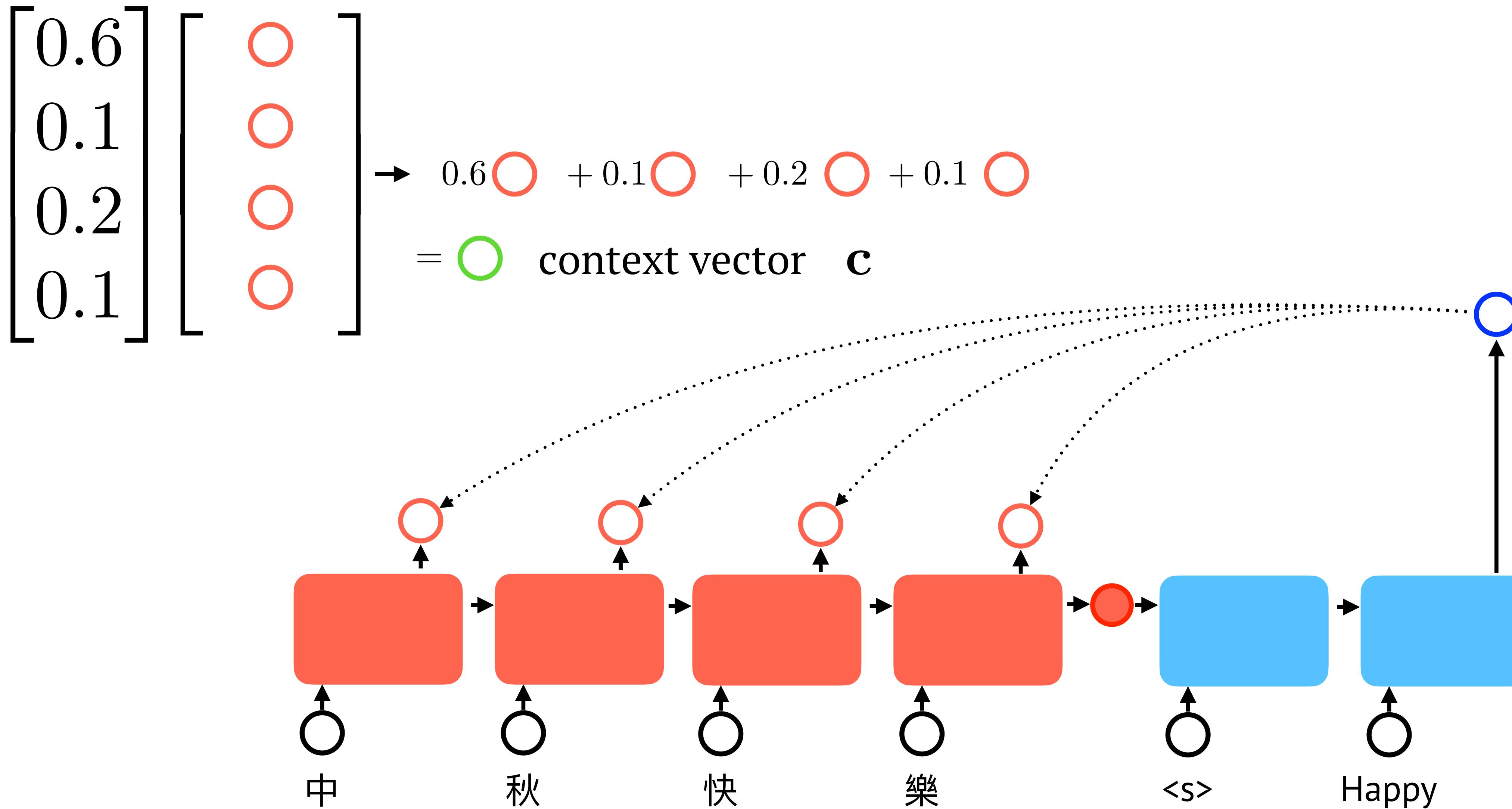
$$\begin{bmatrix} \text{red circle} \\ \text{red circle} \\ \text{red circle} \\ \text{red circle} \end{bmatrix}$$

$$\rightarrow 0.6 \text{ } \text{red circle} + 0.1 \text{ } \text{red circle} + 0.2 \text{ } \text{red circle} + 0.1 \text{ } \text{red circle}$$
$$= \text{green circle} \text{ context vector } \mathbf{c}$$

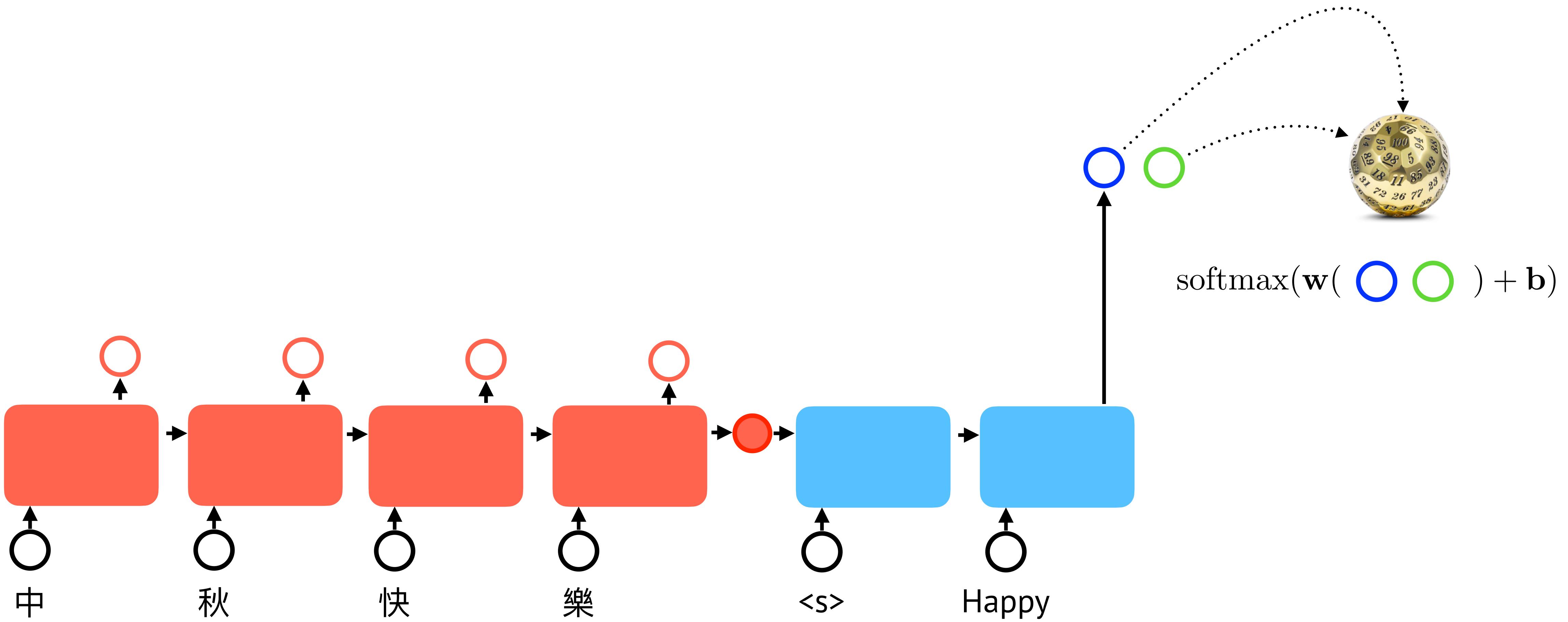
# Attention Mechanism



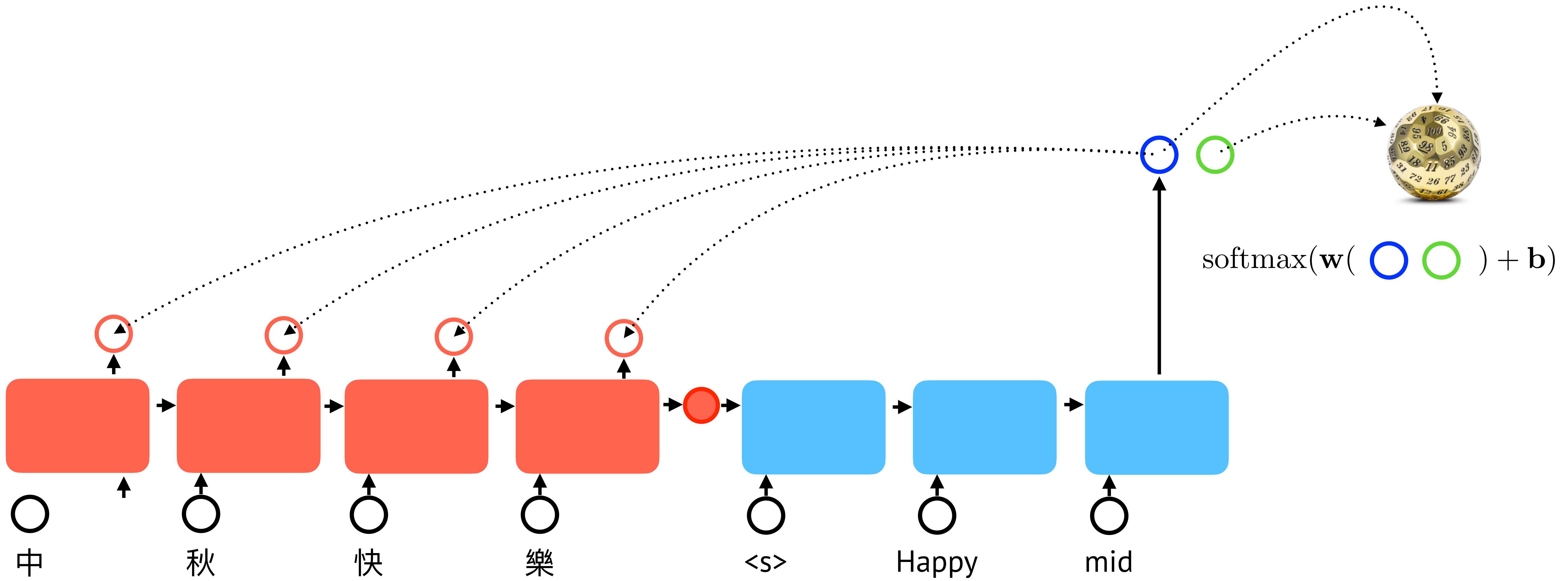
# Attention Mechanism



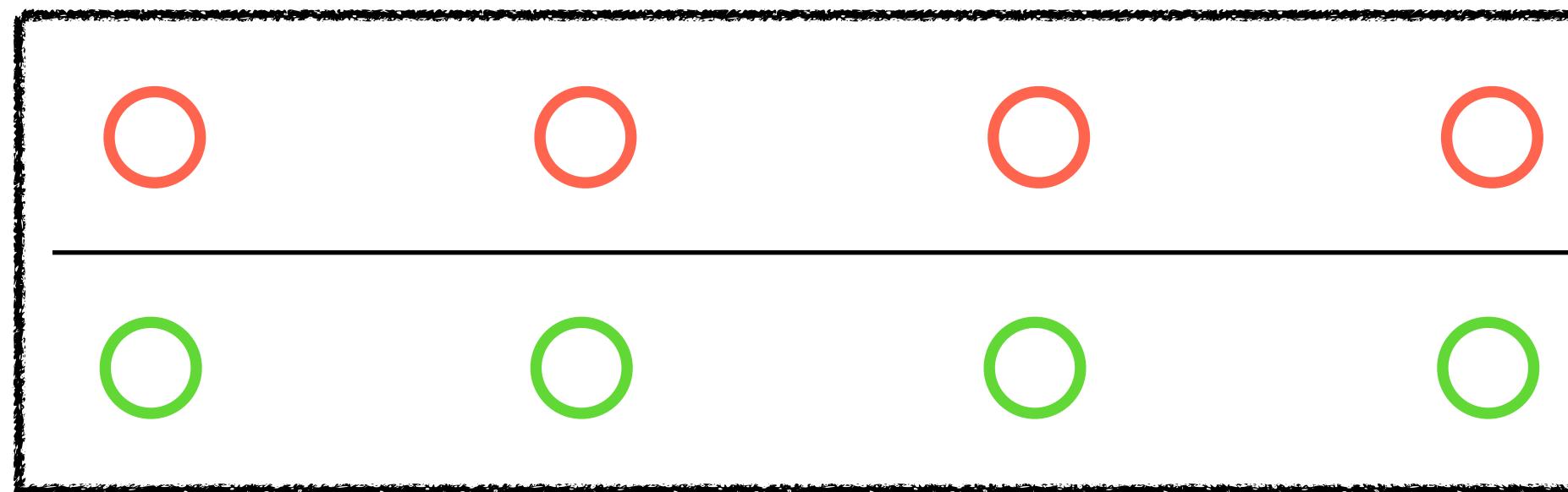
# Attention Mechanism



# Attention Mechanism



# Dot-Product-Softmax Attention



○  
Query

Memory (key-value pairs)

$$\text{○ ○ } q \cdot k_1$$

$$\text{○ ○ } q \cdot k_2$$

$$\text{○ ○ } q \cdot k_3$$

$$\text{○ ○ } q \cdot k_4$$

$$q \cdot k_1$$

$$q \cdot k_2$$

$$q \cdot k_3$$

$$q \cdot k_4$$

$$\text{softmax}(\quad \quad \quad ) \rightarrow$$

$$\begin{bmatrix} 0.6 \\ 0.1 \\ 0.2 \\ 0.1 \end{bmatrix}$$

$$\begin{bmatrix} \text{○} \\ \text{○} \\ \text{○} \\ \text{○} \end{bmatrix}$$

$$\rightarrow 0.6 \text{○} + 0.1 \text{○} + 0.2 \text{○} + 0.1 \text{○} = \text{○} \text{ context vector } \mathbf{c}$$

# Dot-Product-Softmax Attention

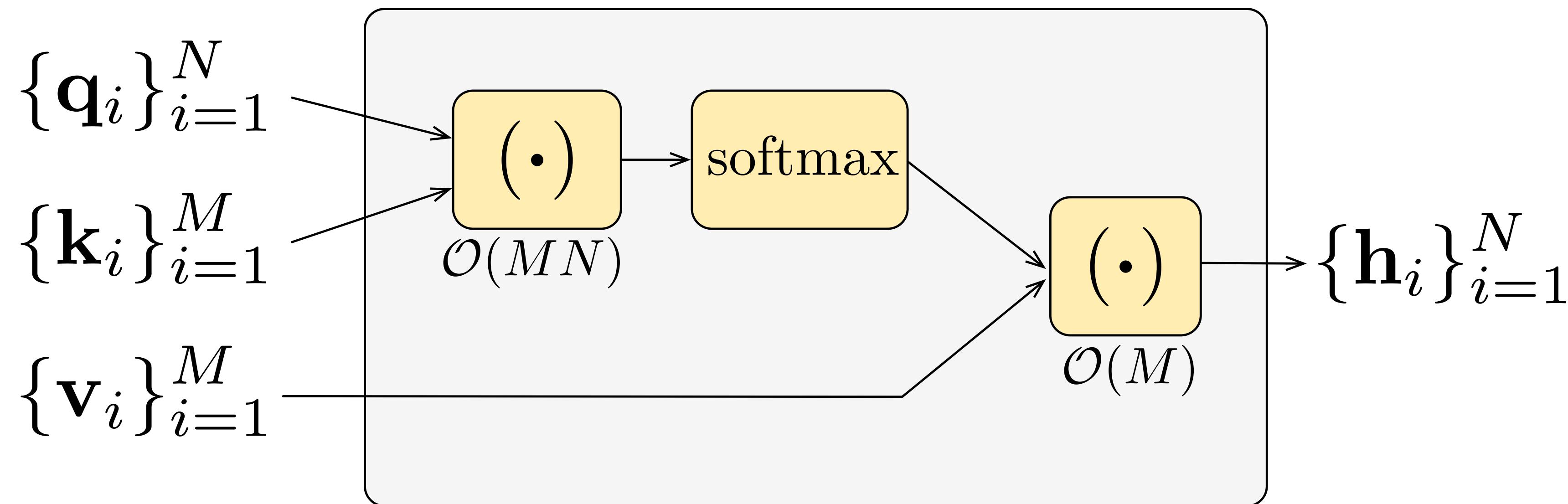
$$\sum_{m=1}^M \frac{\exp(q_n k_m)}{\sum_{m'=1}^M \exp(q_n k_{m'})} v_m^\top = \mathbf{V}^\top \text{softmax}(\mathbf{K} q_n)$$

similarity

normalized similarity

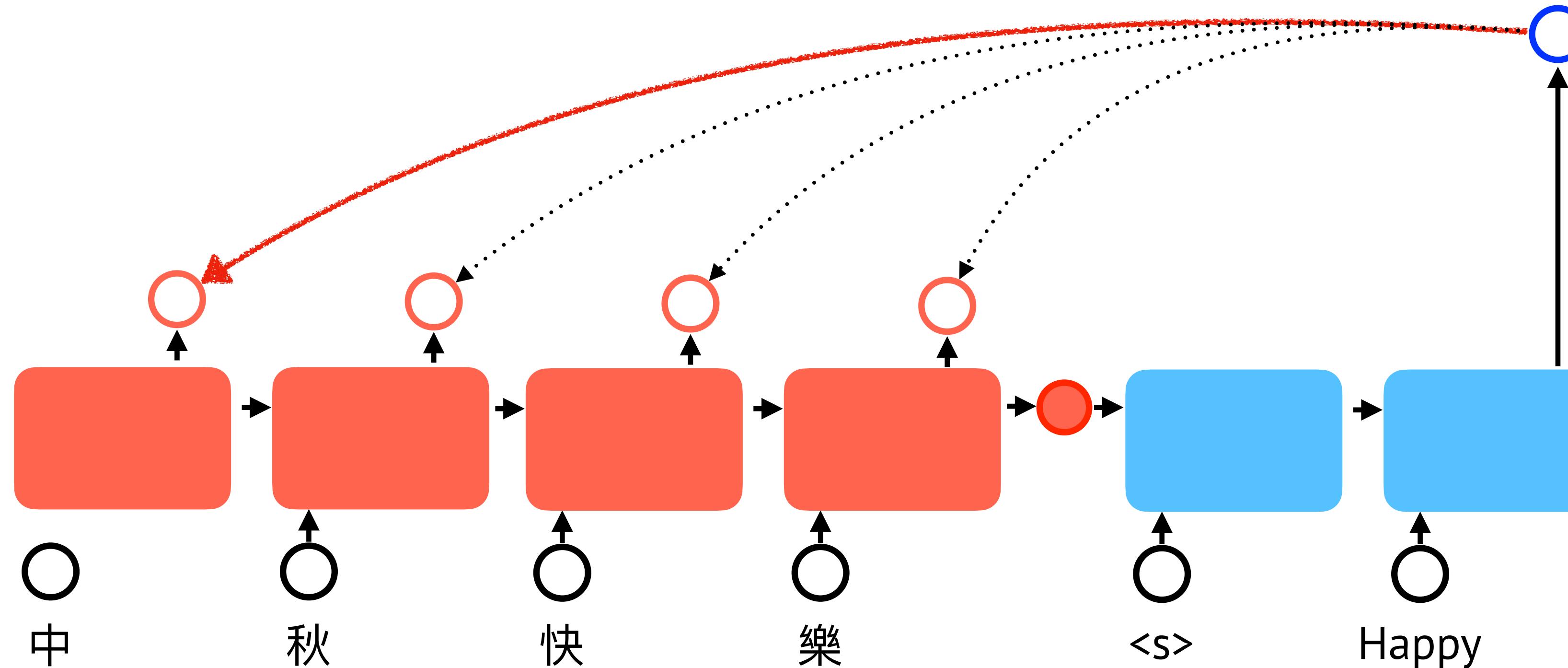
weighted sum

# Computational Complexity



# Attention Mechanism

It helps with vanishing gradient problem.



# Attention Mechanism

It offers some interpretability.

