# Self-Attention

Shusen Wang

# Self-Attention

- Self-Attention: attention beyond Seq2Seq models.

- The original self-attention paper uses LSTM.

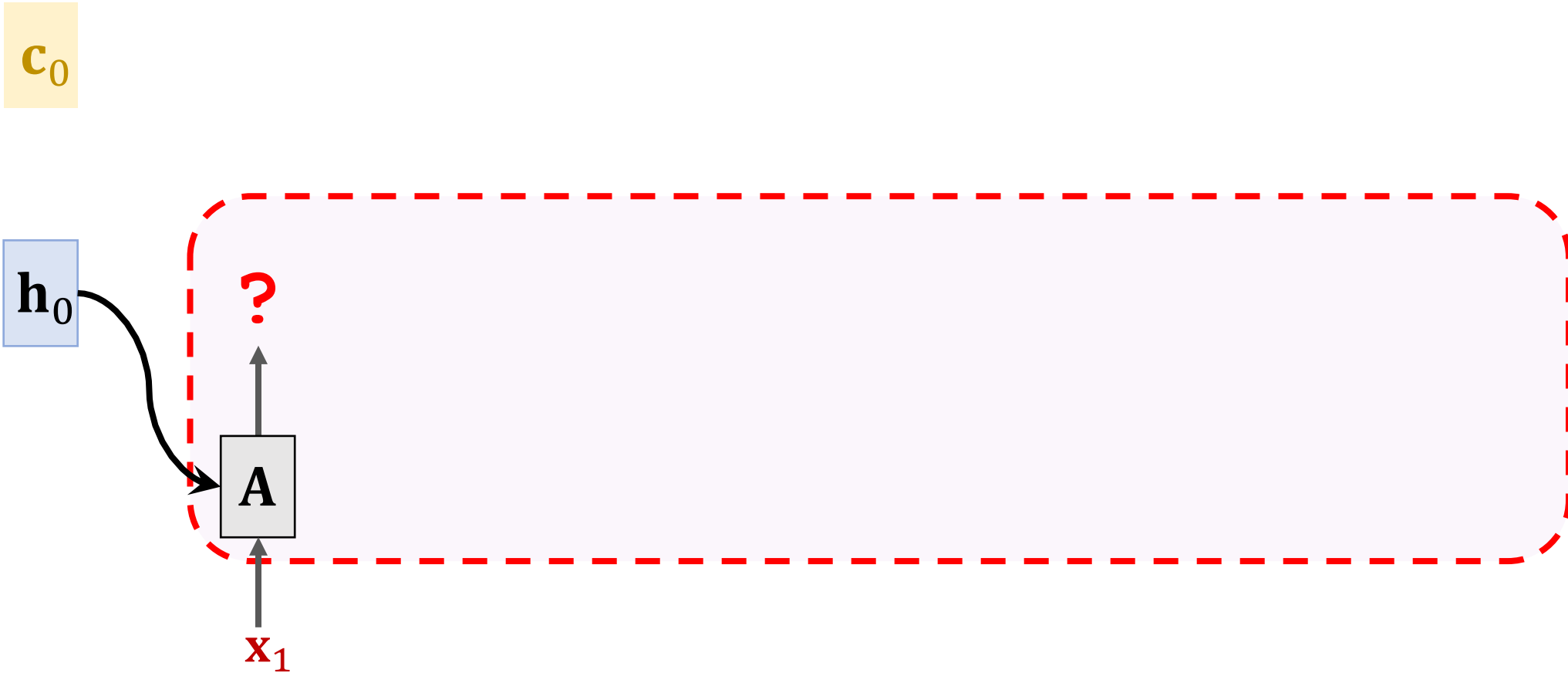- To make teaching easy, I replace LSTM by SimpleRNN.

**Original paper:**

- Cheng, Dong, & Lapata. Long Short-Term Memory-Networks for Machine Reading. In *EMNLP*, 2016.

# SimpleRNN + Self-Attention

$$\mathbf{c}_0 = \mathbf{0}$$

$$\mathbf{h}_0 = \mathbf{0}$$

# SimpleRNN + Self-Attention

$\mathbf{c}_0$

$\mathbf{h}_0$

**?**

**A**

$\mathbf{x}_1$

# SimpleRNN + Self-Attention

SimpleRNN:

$$\mathbf{h}_1 = \tanh\left(\mathbf{A} \cdot \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{h}_0 \end{bmatrix} + \mathbf{b}\right)$$

$\mathbf{c}_0$

$\mathbf{h}_0$

?

$\mathbf{A}$

$\mathbf{x}_1$

# SimpleRNN + Self-Attention
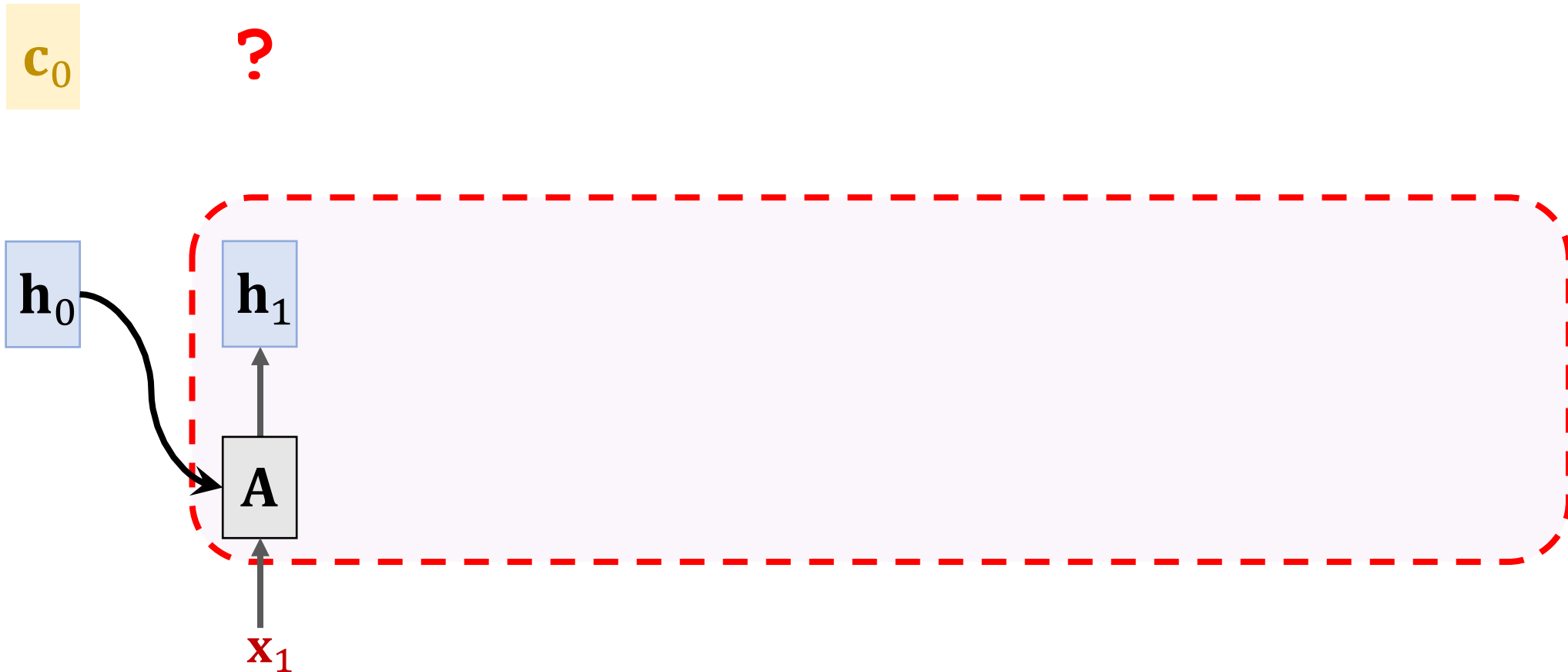
SimpleRNN + Self-Attention:

$$\mathbf{h}_1 = \tanh\left(\mathbf{A} \cdot \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{c}_0 \end{bmatrix} + \mathbf{b}\right)$$
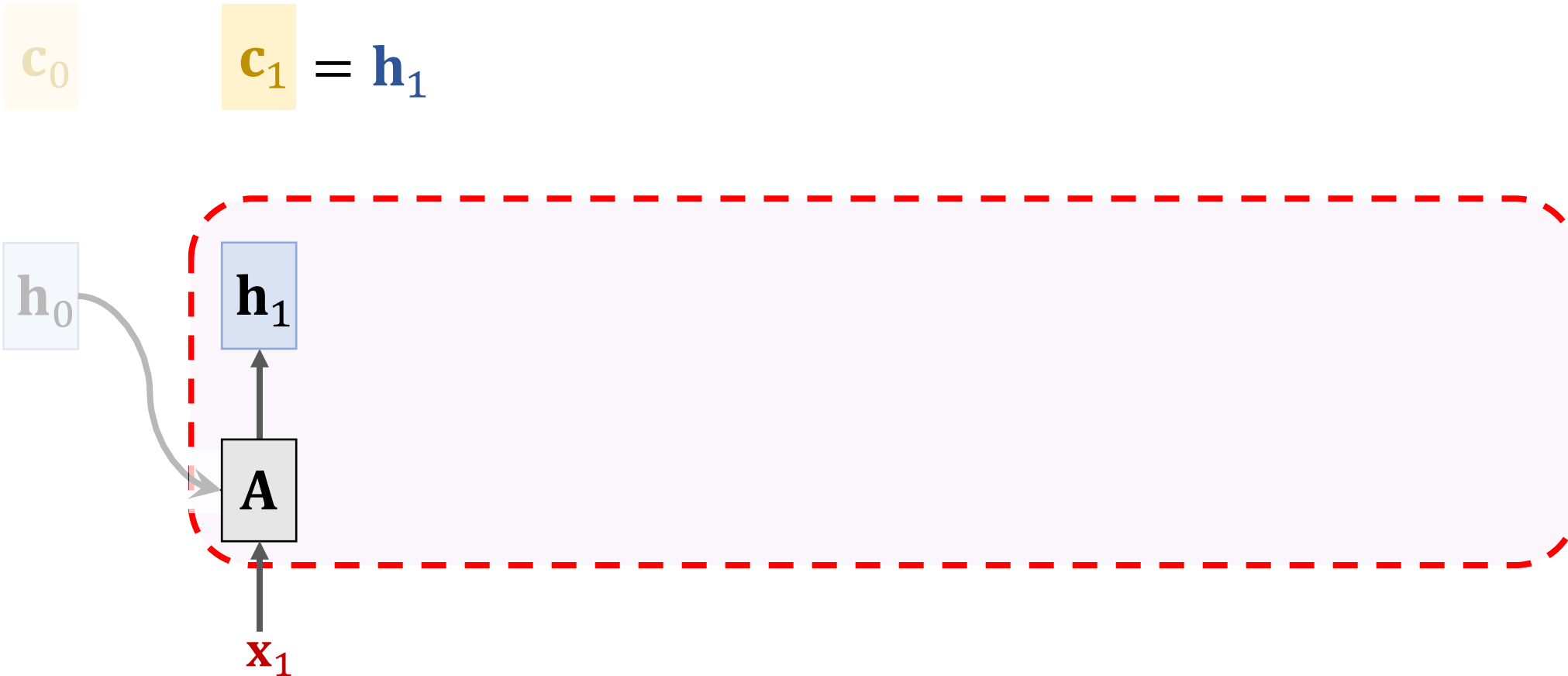
# SimpleRNN + Self-Attention

# SimpleRNN + Self-Attention

First context vector:   $c_1 = h_1$.

$c_0$

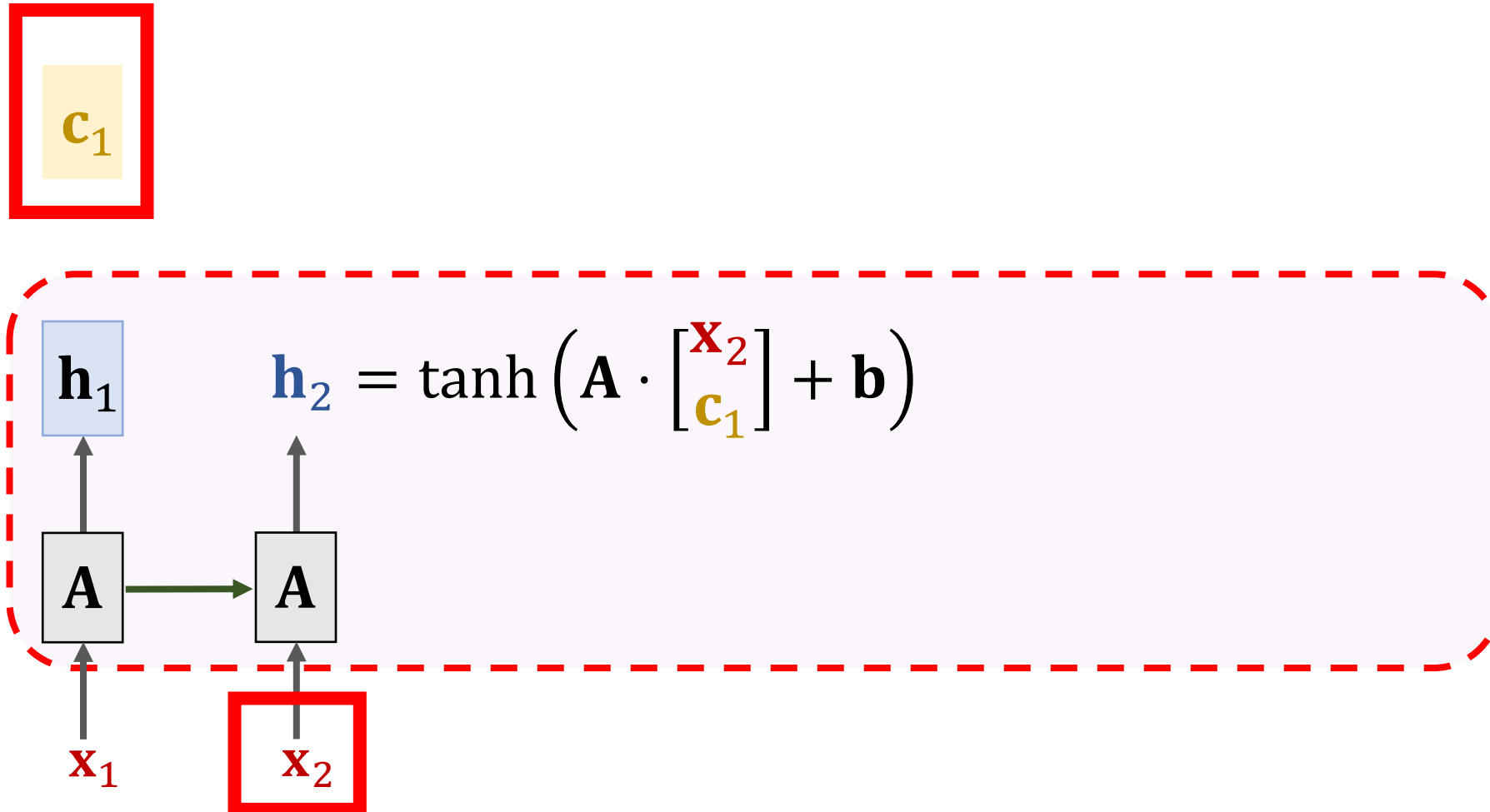$c_1 = h_1$

$h_0$

$h_1$

$A$

$x_1$

# SimpleRNN + Self-Attention

$\mathbf{c}_1$

$\mathbf{h}_1$ **?**

A → A

$\mathbf{x}_1$ $\mathbf{x}_2$

# SimpleRNN + Self-Attention



$$\mathbf{h}_2 = \tanh\left(\mathbf{A} \cdot \begin{bmatrix} \mathbf{x}_2 \\ \mathbf{c}_1 \end{bmatrix} + \mathbf{b}\right)$$

# SimpleRNN + Self-Attention

$c_1$   **?**

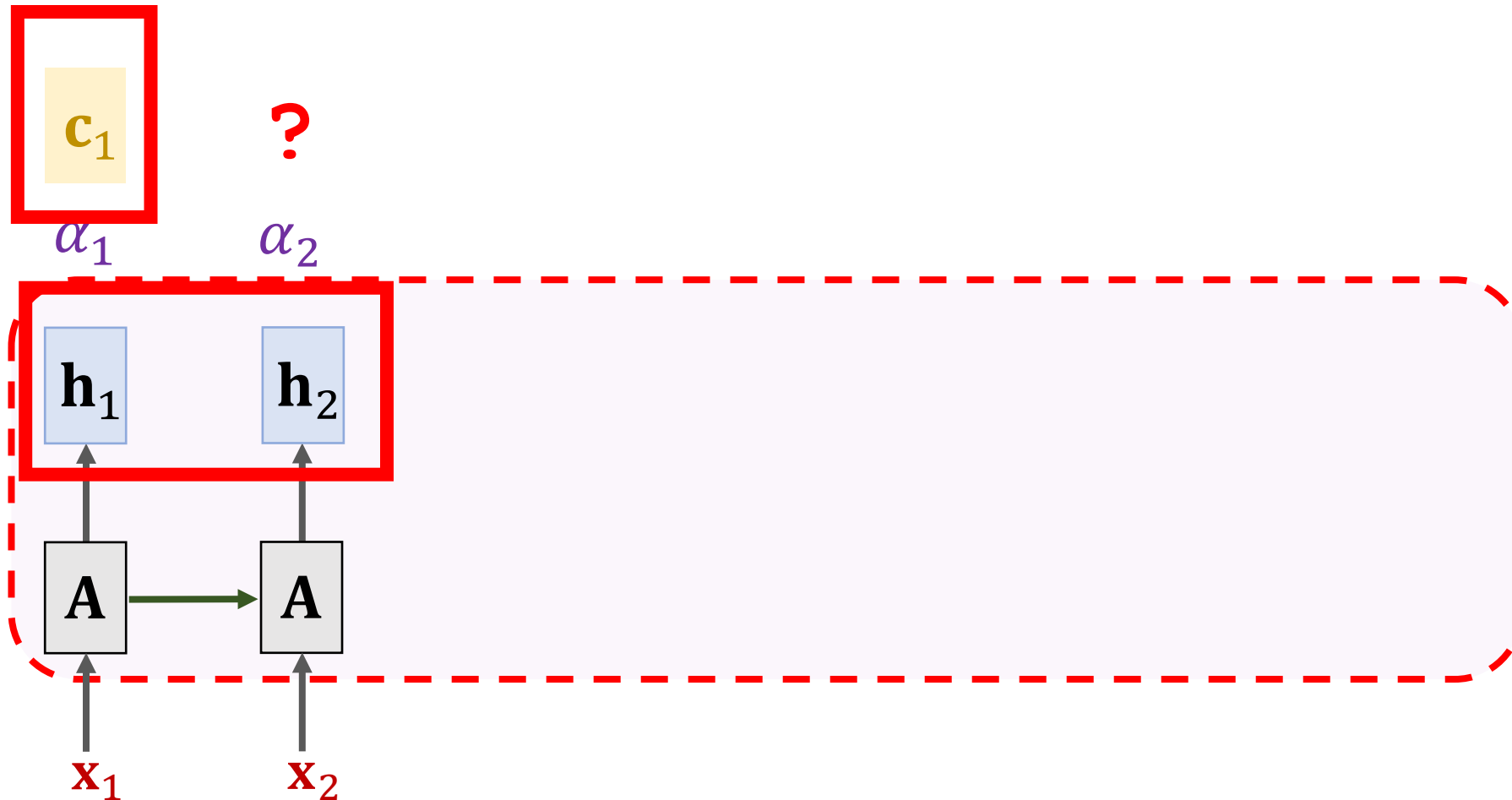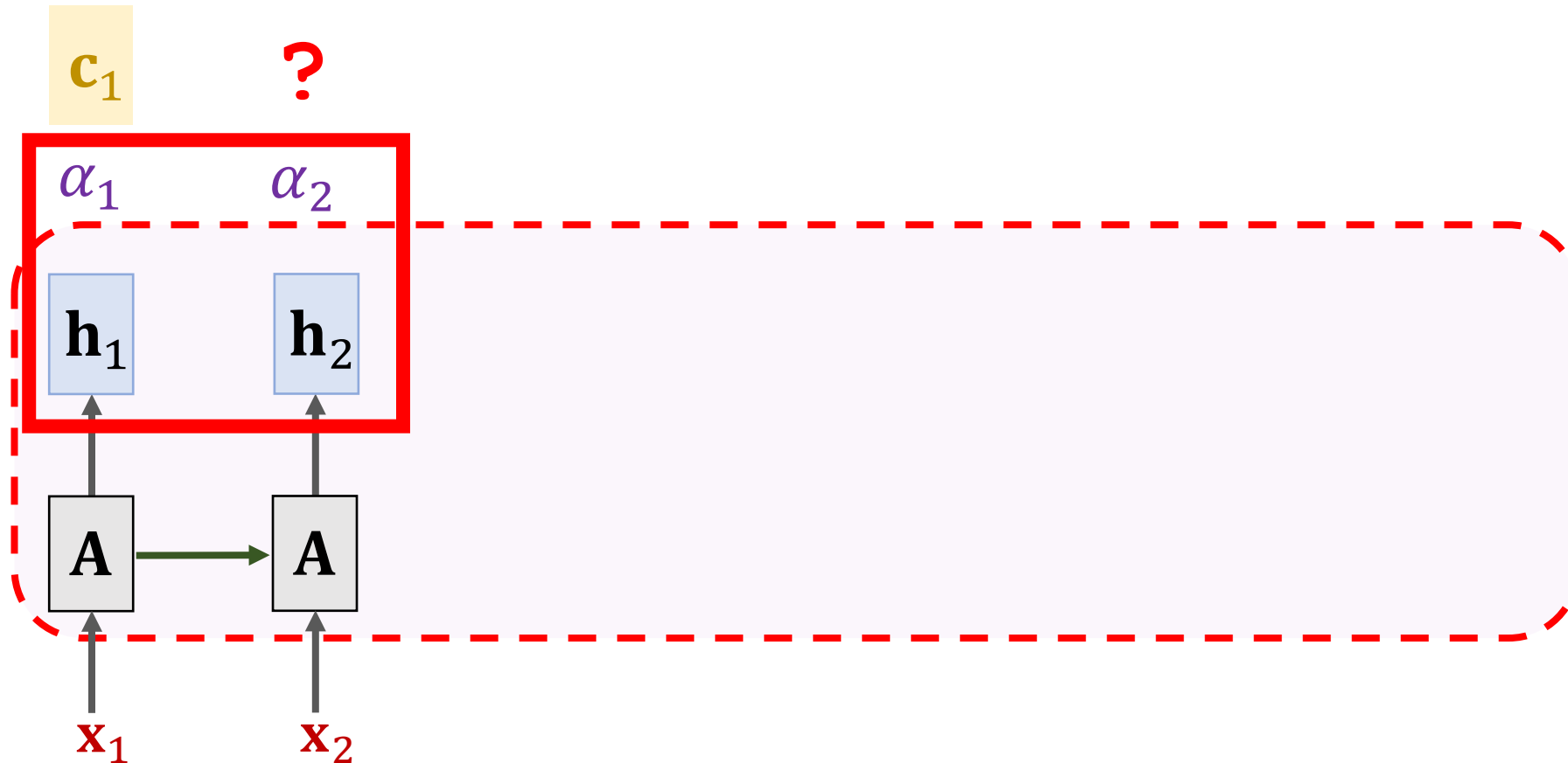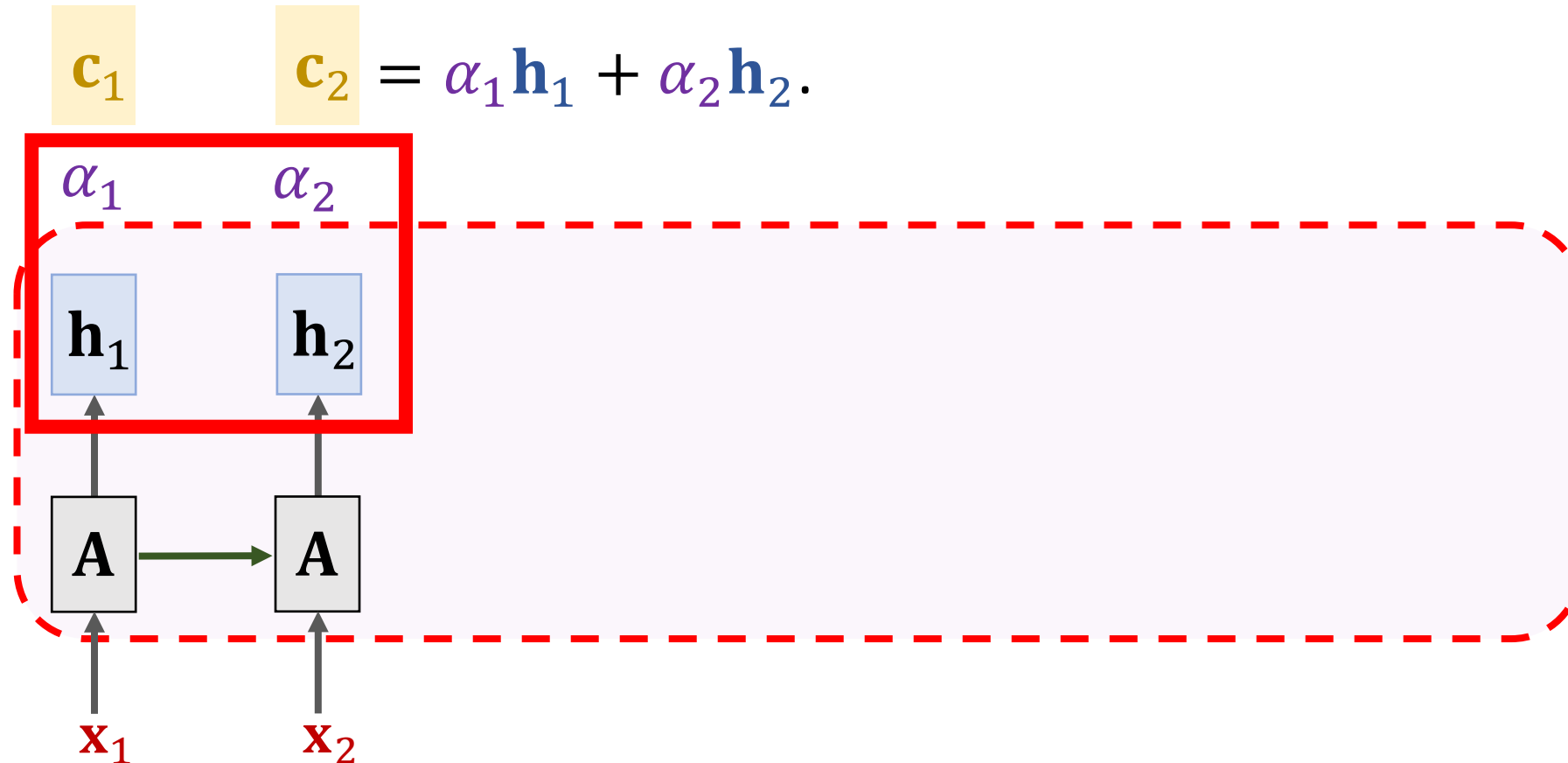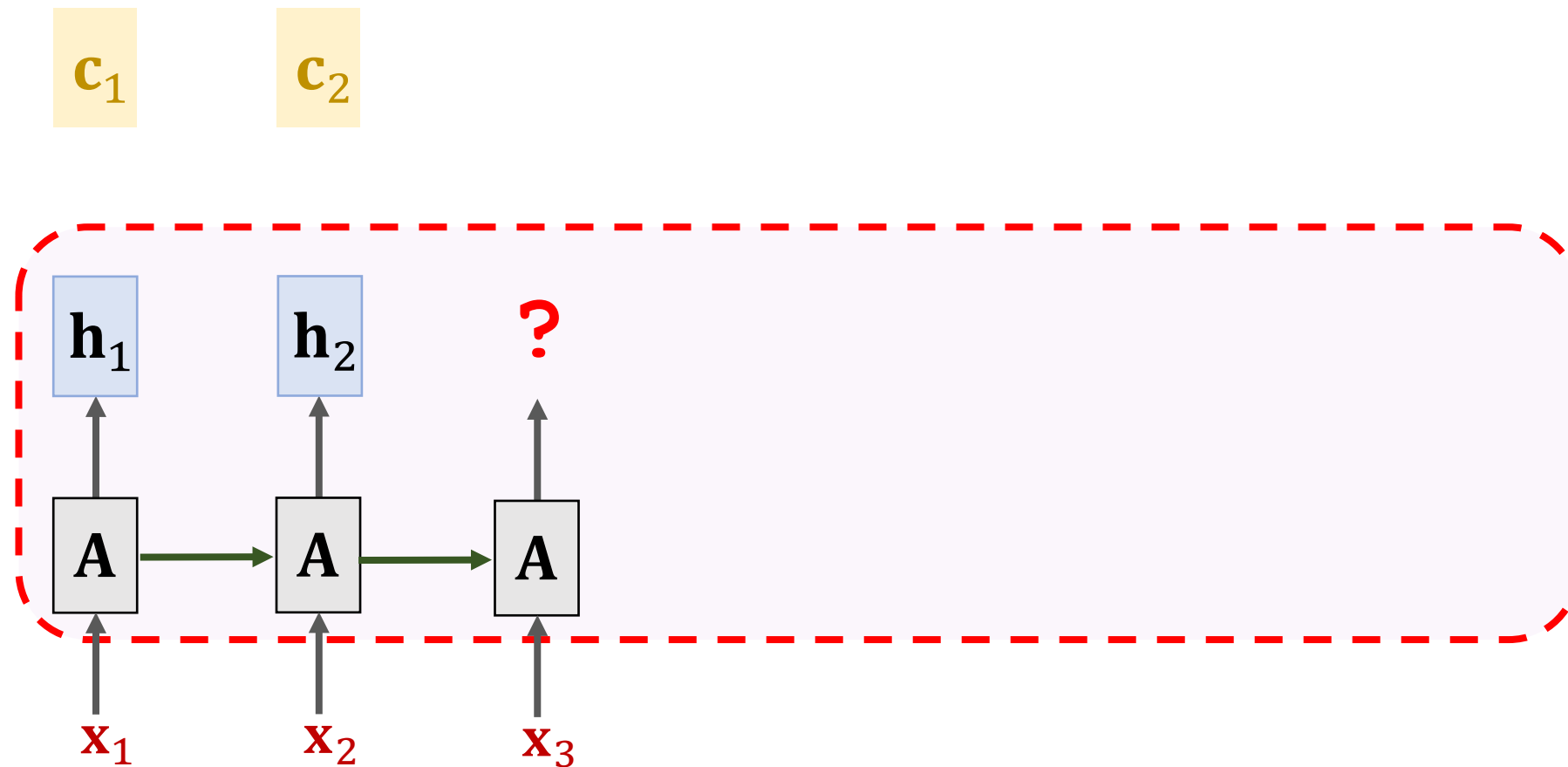$h_1$   $h_2$

A → A

$x_1$   $x_2$

# SimpleRNN + Self-Attention

# SimpleRNN + Self-Attention

**Weights**: $\quad \alpha_i = \text{similarity}(\mathbf{h}_i, \mathbf{c}_1)$
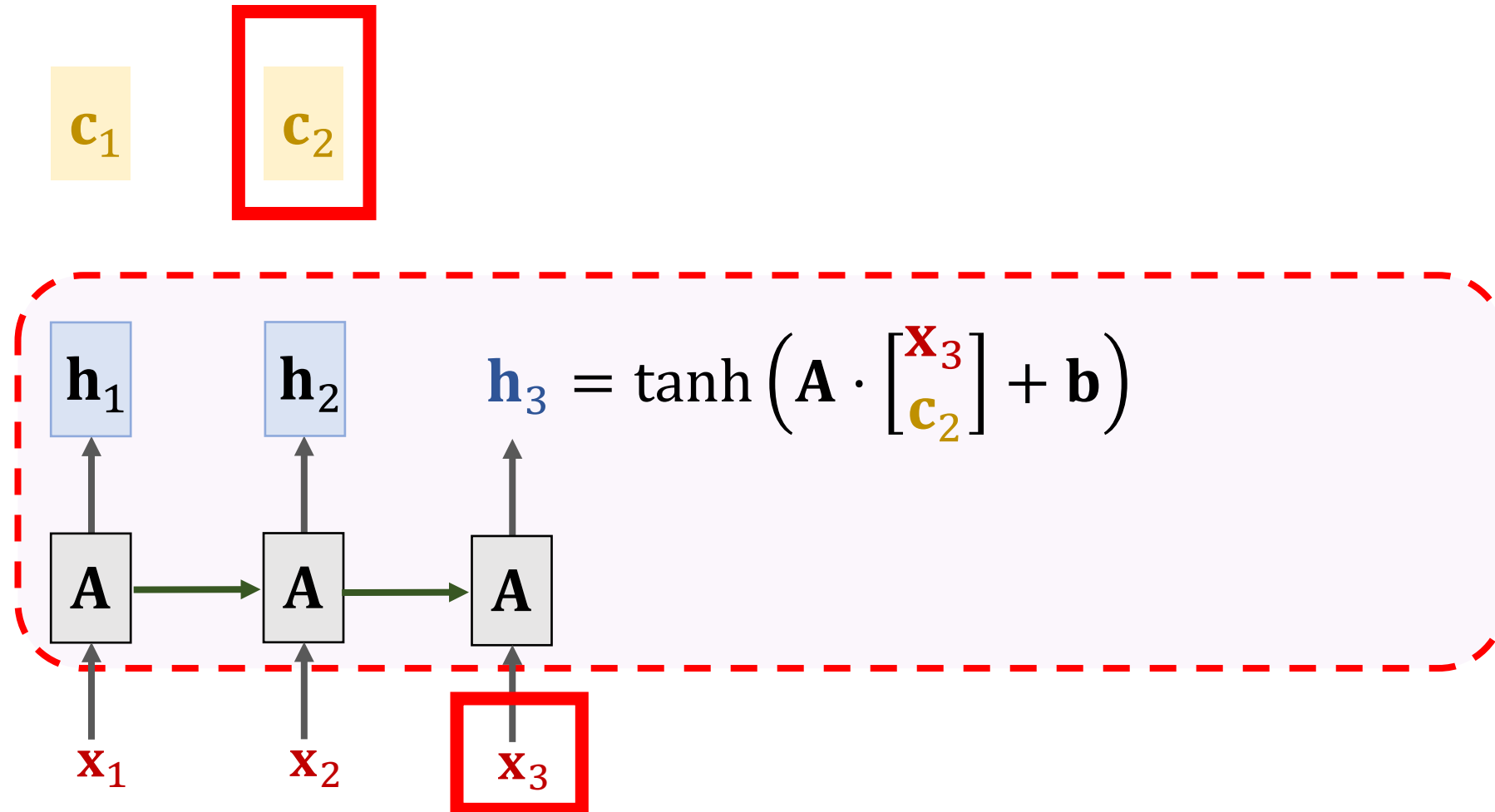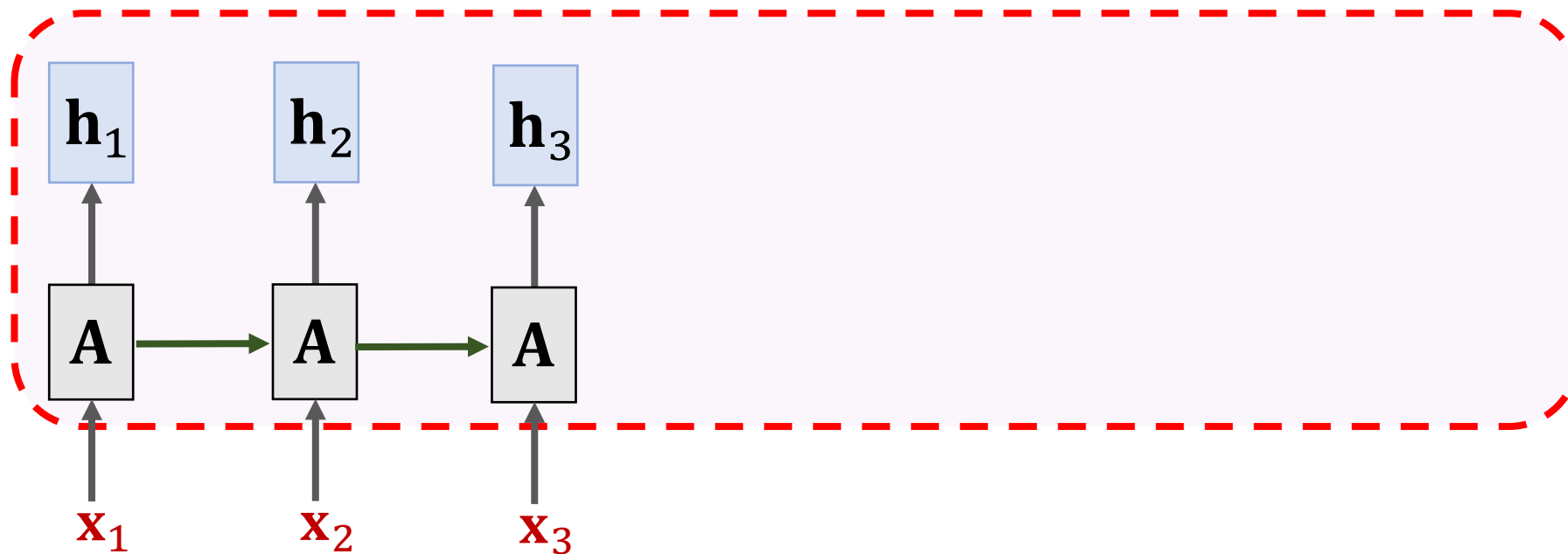
# SimpleRNN + Self-Attention

$\mathbf{c}_1 \qquad \mathbf{c}_2 = \alpha_1 \mathbf{h}_1 + \alpha_2 \mathbf{h}_2.$

$\alpha_1 \qquad \alpha_2$

$\mathbf{h}_1 \qquad \mathbf{h}_2$

$\mathbf{A} \longrightarrow \mathbf{A}$

$\mathbf{x}_1 \qquad \mathbf{x}_2$

# SimpleRNN + Self-Attention

$\mathbf{c}_1$  $\mathbf{c}_2$

$\mathbf{h}_1$  $\mathbf{h}_2$  **?**

**A** → **A** → **A**

$\mathbf{x}_1$  $\mathbf{x}_2$  $\mathbf{x}_3$

# SimpleRNN + Self-Attention

$$\mathbf{c}_1 \qquad \mathbf{c}_2$$

$$\mathbf{h}_1 \qquad \mathbf{h}_2 \qquad \mathbf{h}_3 = \tanh\left(\mathbf{A} \cdot \begin{bmatrix} \mathbf{x}_3 \\ \mathbf{c}_2 \end{bmatrix} + \mathbf{b}\right)$$

$$\mathbf{A} \rightarrow \mathbf{A} \rightarrow \mathbf{A}$$

$$\mathbf{x}_1 \qquad \mathbf{x}_2 \qquad \mathbf{x}_3$$

# SimpleRNN + Self-Attention

$\mathbf{c}_1$  $\mathbf{c}_2$  **?**

$\mathbf{h}_1$  $\mathbf{h}_2$  $\mathbf{h}_3$

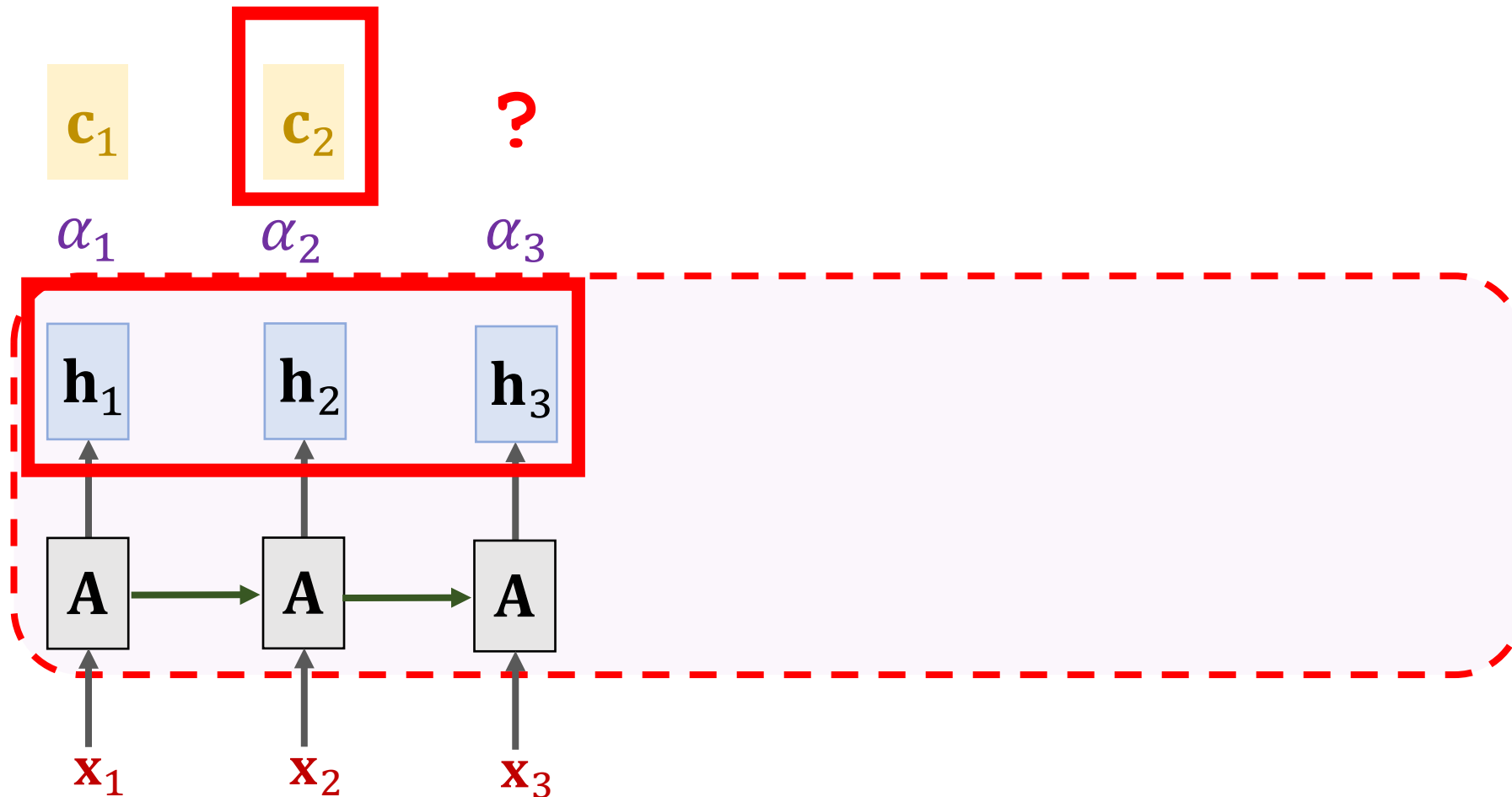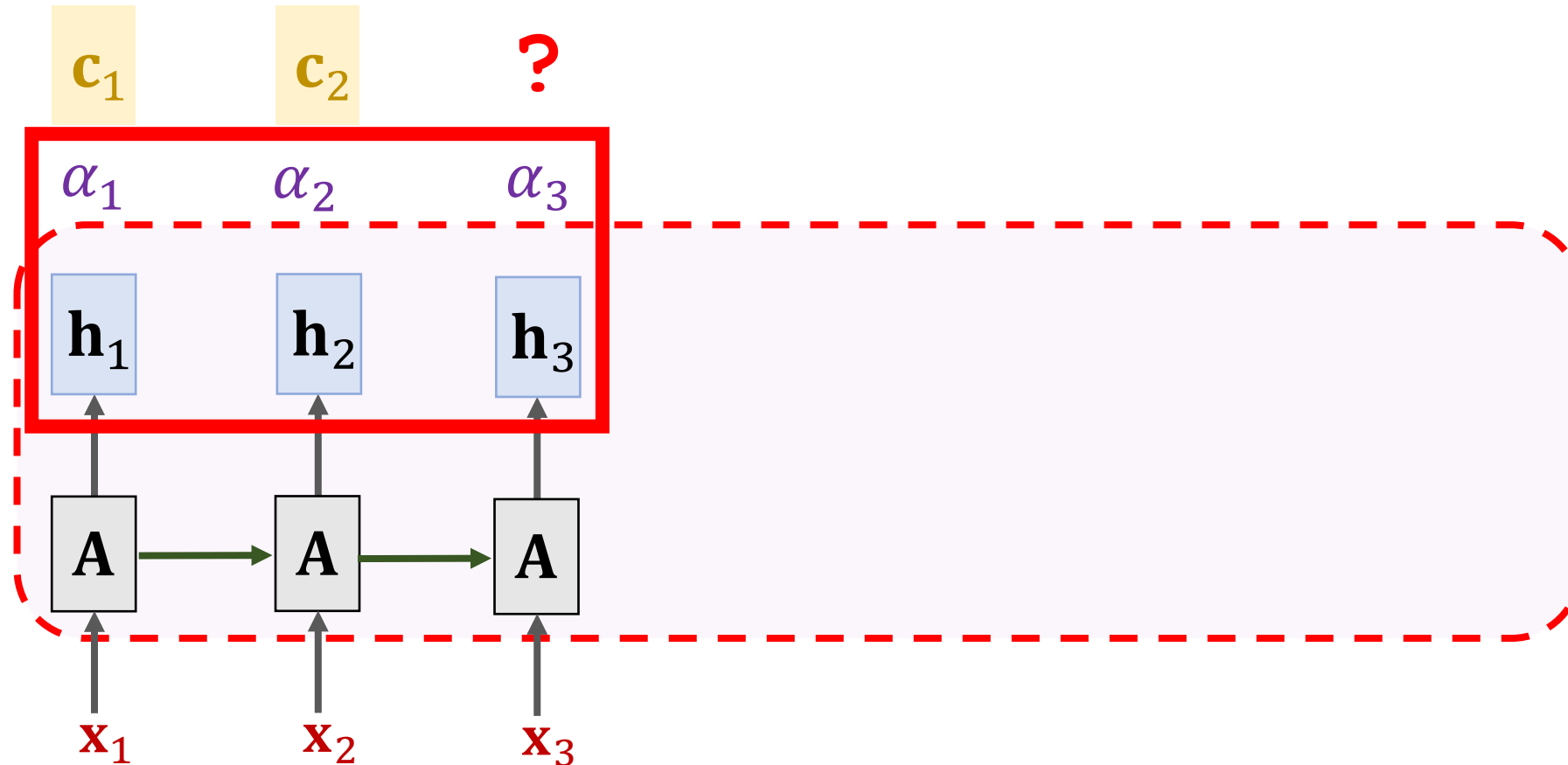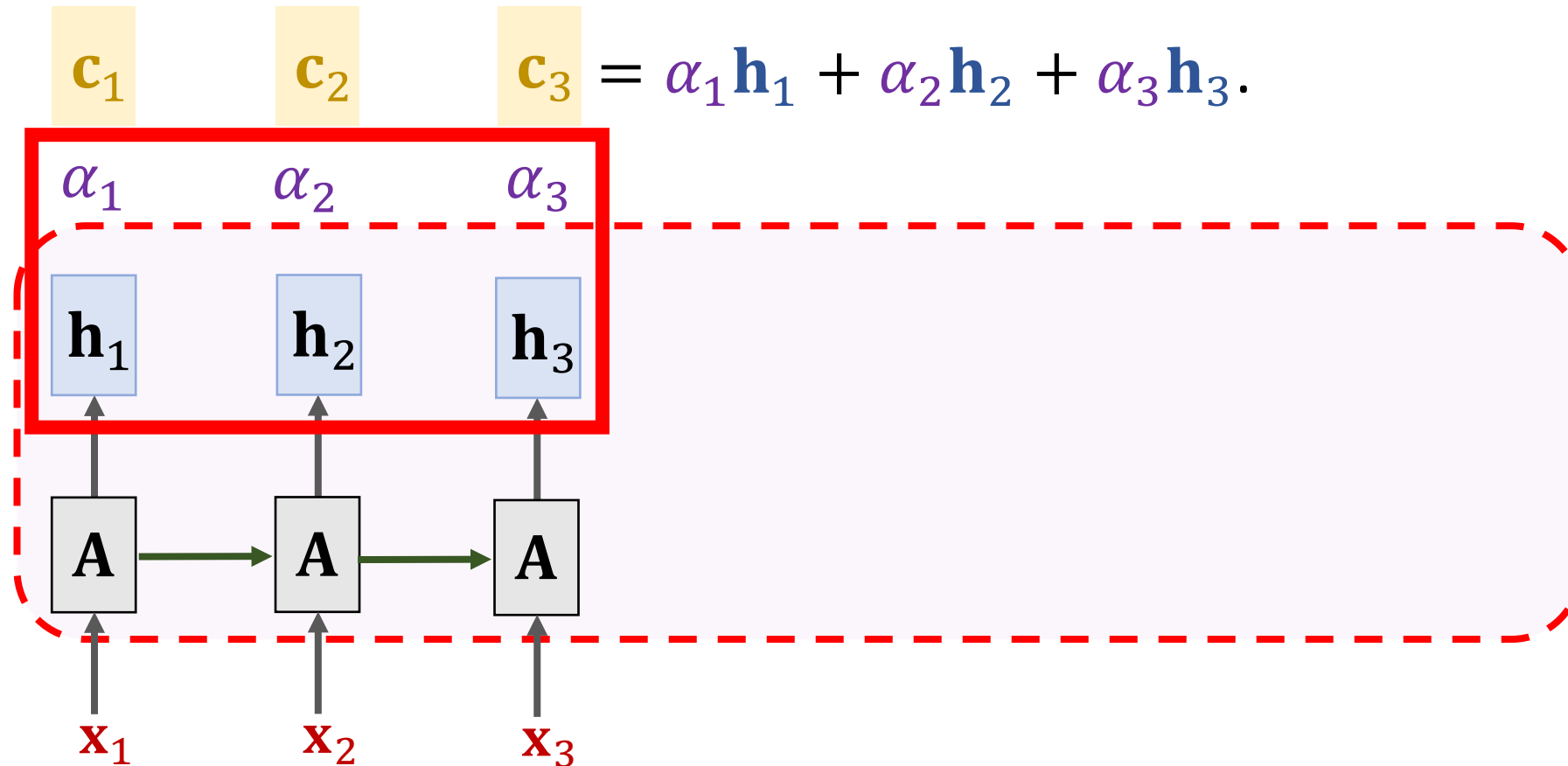**A** → **A** → **A**

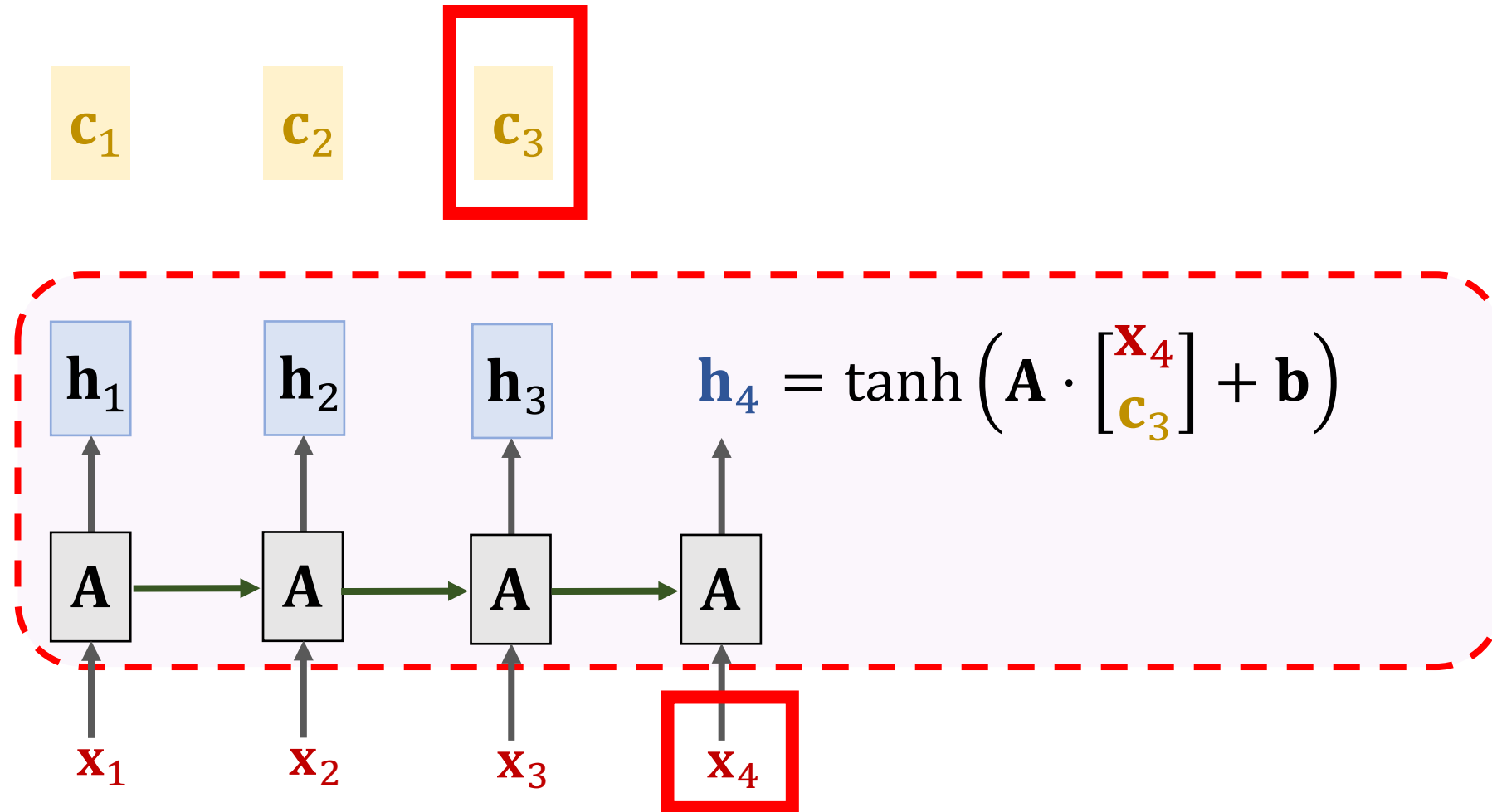$\mathbf{x}_1$  $\mathbf{x}_2$  $\mathbf{x}_3$

# SimpleRNN + Self-Attention

**Weights**: $\quad \alpha_i = \text{similarity}(\mathbf{h}_i, \ \mathbf{c}_2)$
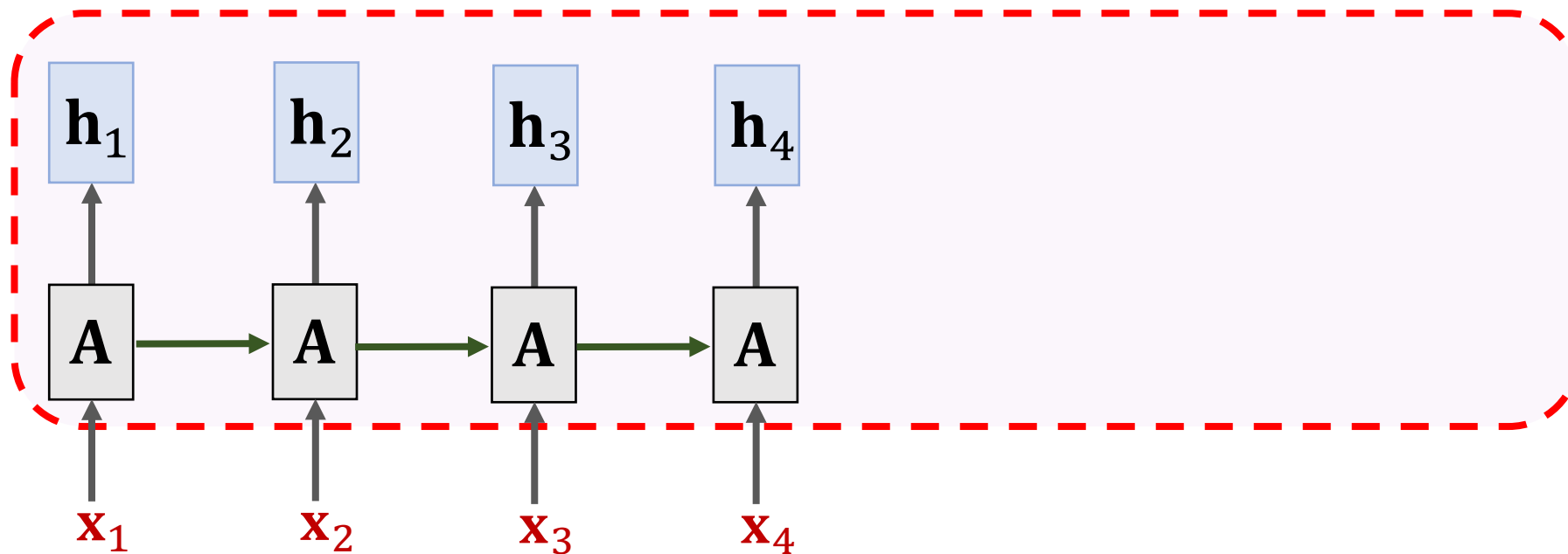
# SimpleRNN + Self-Attention

# SimpleRNN + Self-Attention
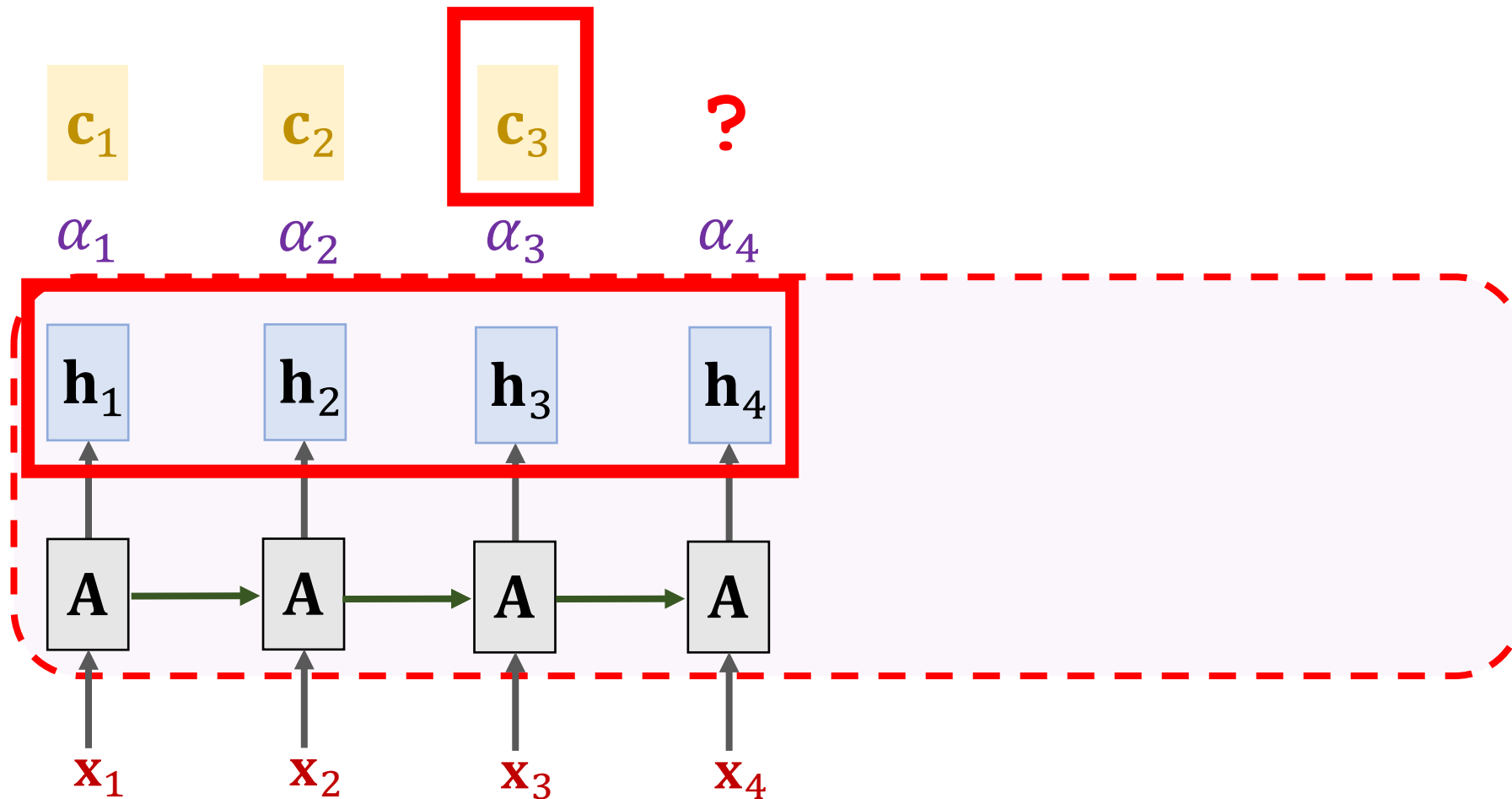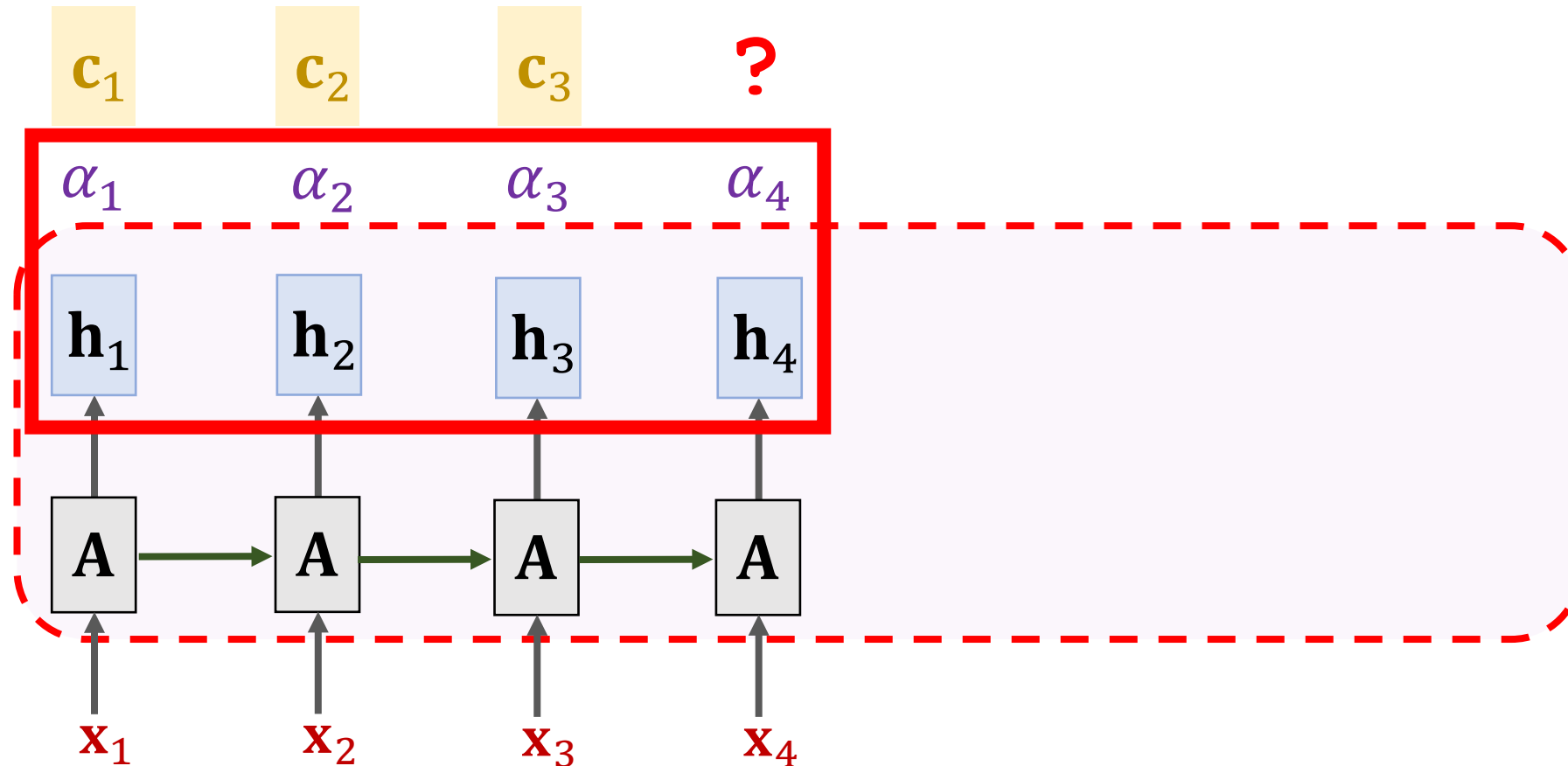
$\mathbf{c}_1 \qquad \mathbf{c}_2 \qquad \mathbf{c}_3 = \alpha_1 \mathbf{h}_1 + \alpha_2 \mathbf{h}_2 + \alpha_3 \mathbf{h}_3.$

$\alpha_1 \qquad \alpha_2 \qquad \alpha_3$

$\mathbf{h}_1 \qquad \mathbf{h}_2 \qquad \mathbf{h}_3$

$\mathbf{A} \qquad \mathbf{A} \qquad \mathbf{A}$

$\mathbf{x}_1 \qquad \mathbf{x}_2 \qquad \mathbf{x}_3$

# SimpleRNN + Self-Attention



$$\mathbf{h}_4 = \tanh\left(\mathbf{A} \cdot \begin{bmatrix} \mathbf{x}_4 \\ \mathbf{c}_3 \end{bmatrix} + \mathbf{b}\right)$$

# SimpleRNN + Self-Attention

# SimpleRNN + Self-Attention

**Weights**: $\alpha_i = \text{similarity}(\mathbf{h}_i, \mathbf{c}_3)$
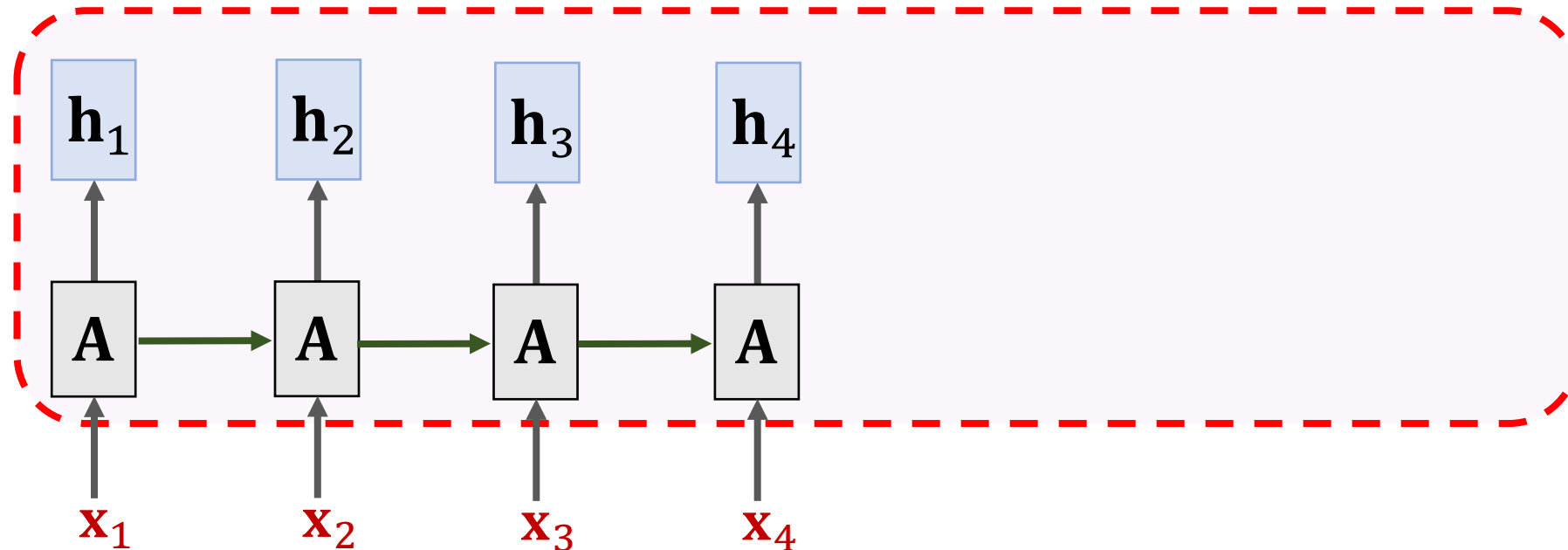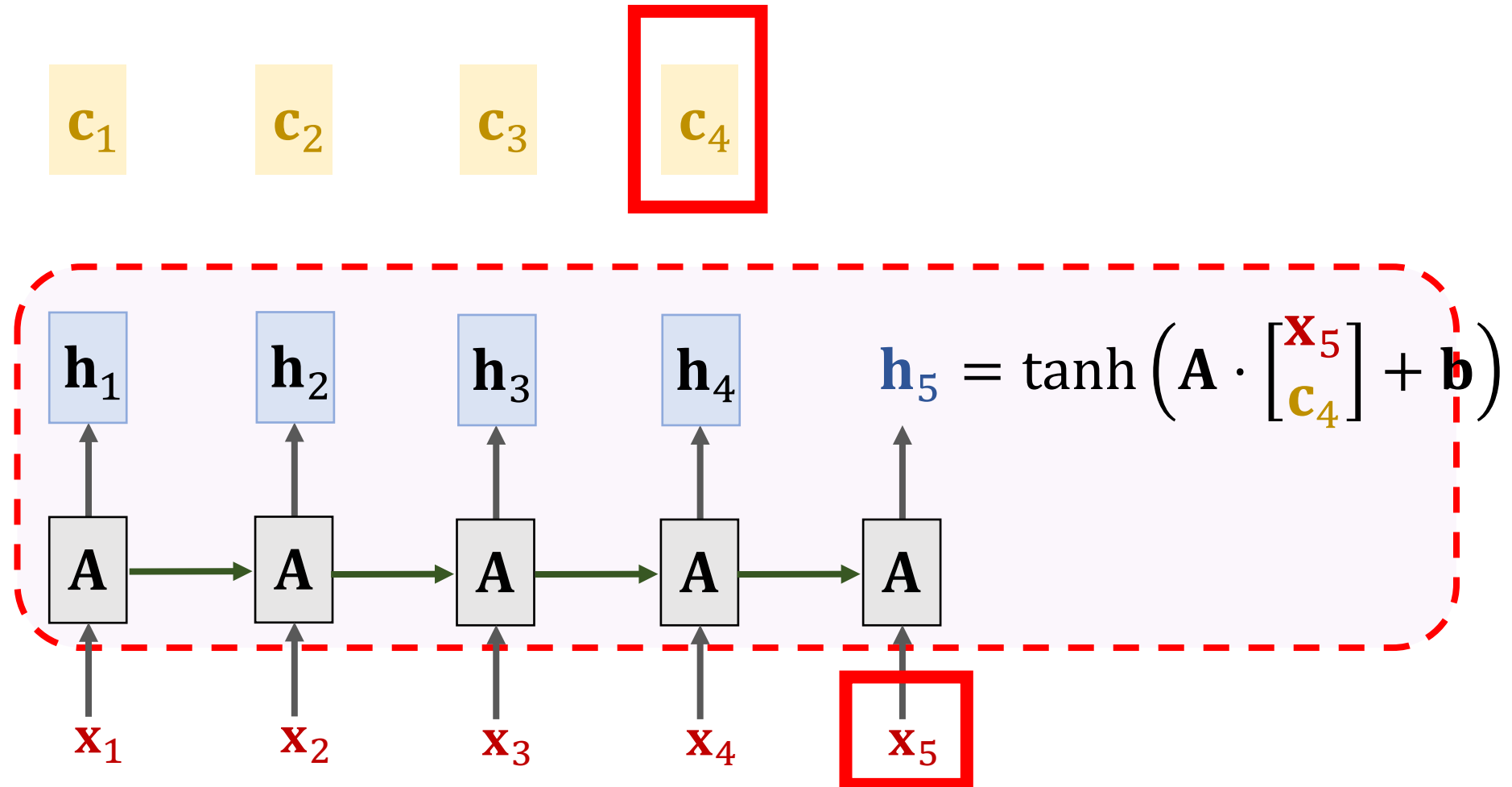
# SimpleRNN + Self-Attention

# SimpleRNN + Self-Attention

**Weights**: $\quad \alpha_i = \text{similarity}(\mathbf{h}_i, \mathbf{c}_3)$
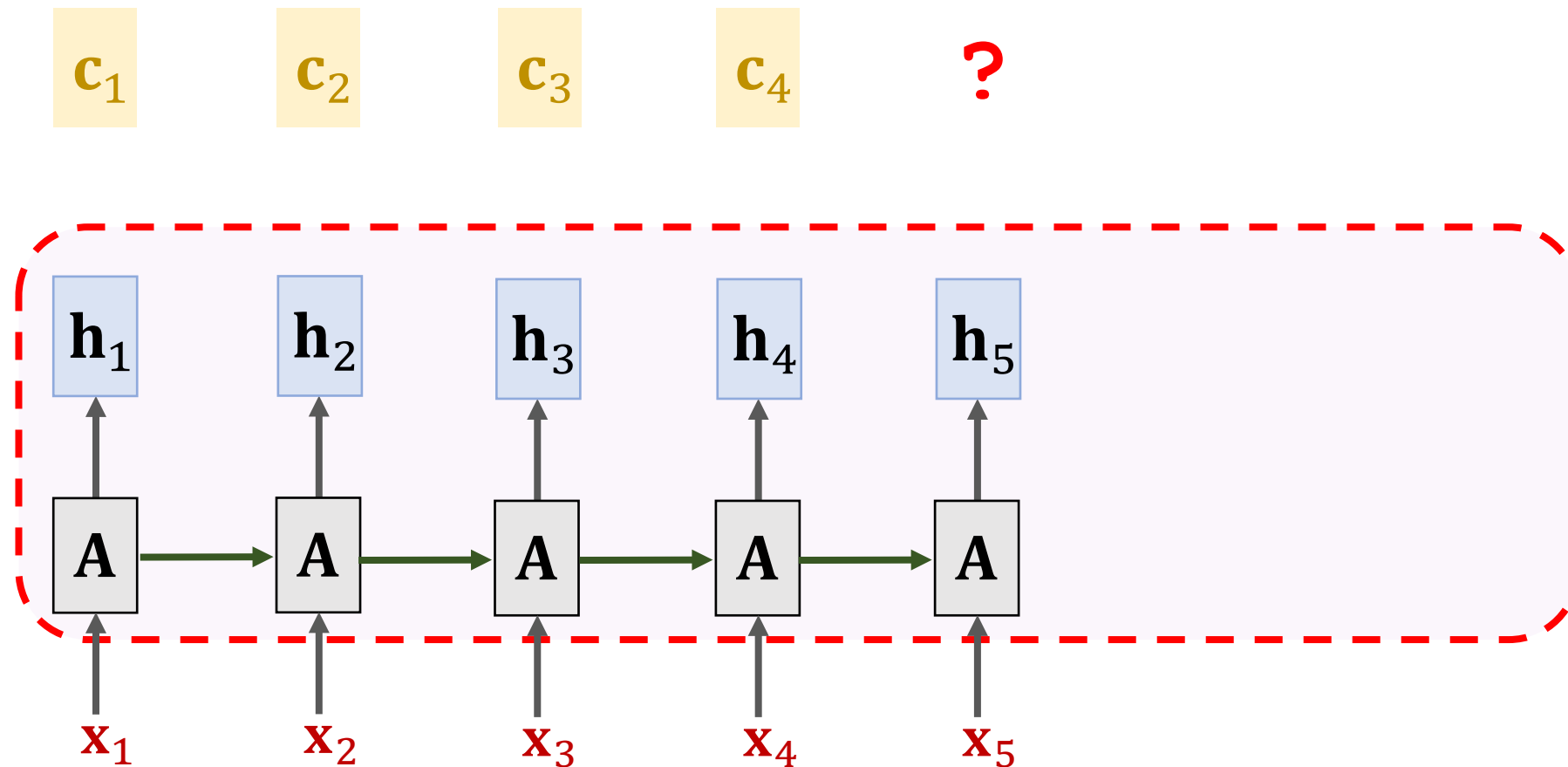
$\mathbf{c}_1 \qquad \mathbf{c}_2 \qquad \mathbf{c}_3 \qquad \mathbf{c}_4 = \alpha_1 \mathbf{h}_1 + \alpha_2 \mathbf{h}_2 + \alpha_3 \mathbf{h}_3 + \alpha_4 \mathbf{h}_4.$
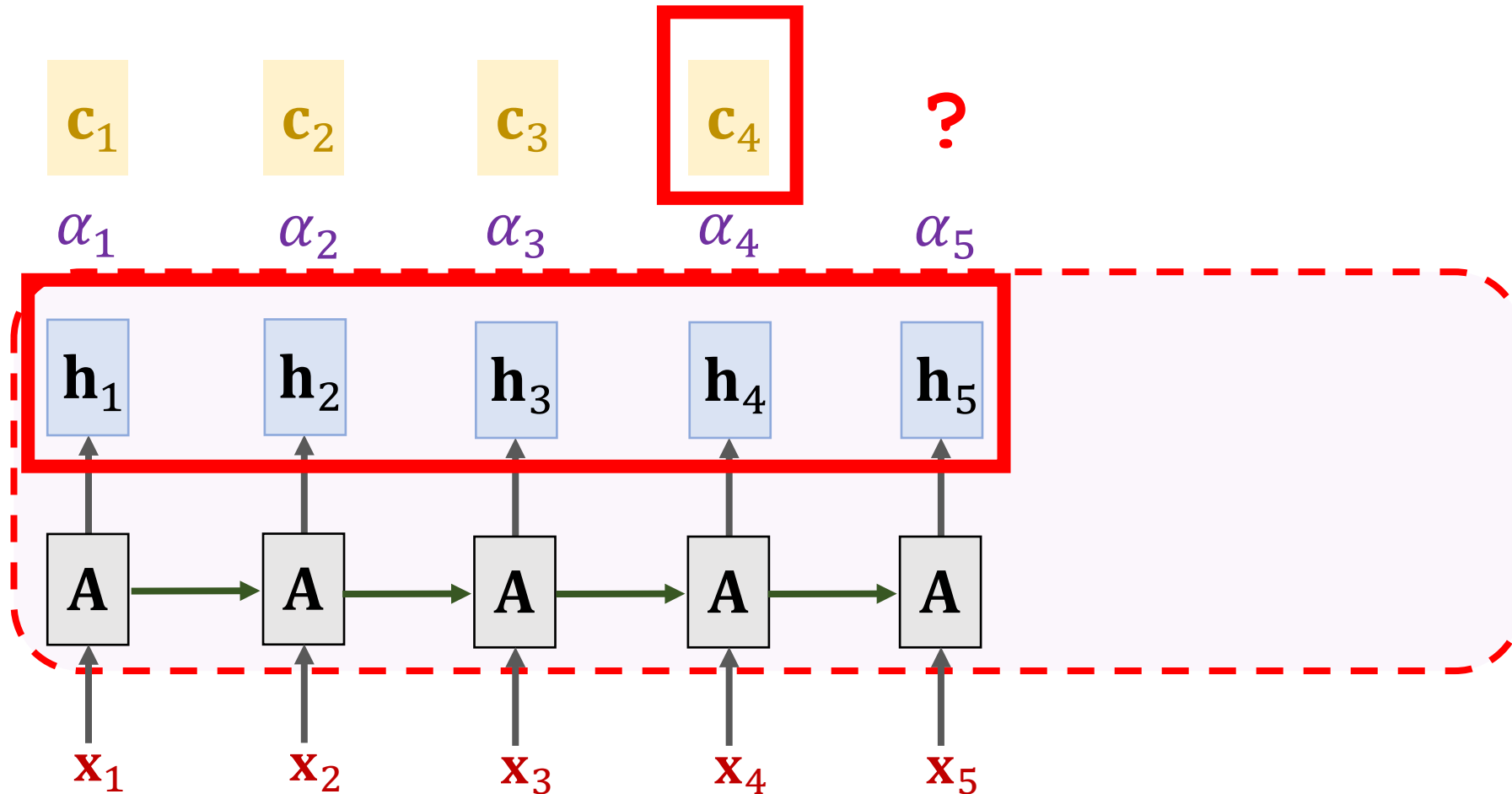
# SimpleRNN + Self-Attention



$$\mathbf{h}_5 = \tanh\left(\mathbf{A} \cdot \begin{bmatrix} \mathbf{x}_5 \\ \mathbf{c}_4 \end{bmatrix} + \mathbf{b}\right)$$
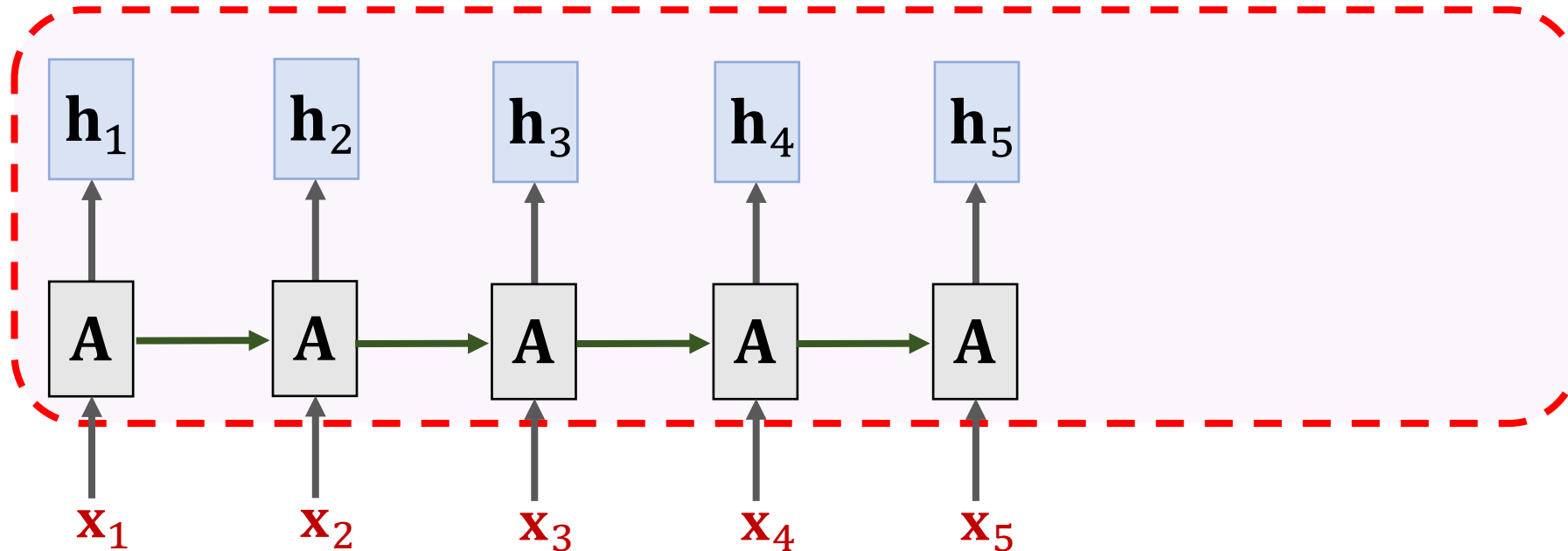
# SimpleRNN + Self-Attention

# SimpleRNN + Self-Attention

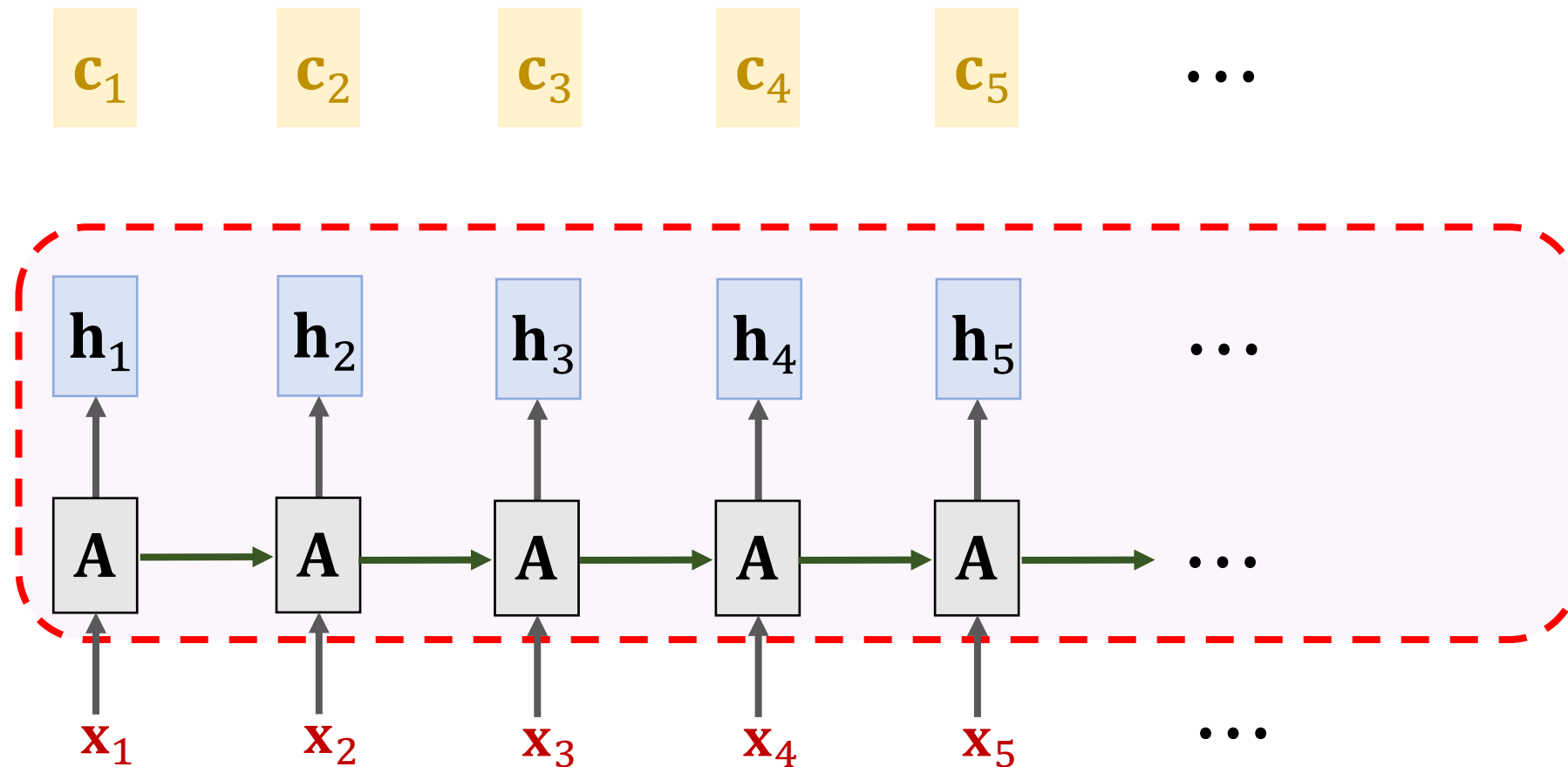**Weights**: $\alpha_i = \text{similarity}(\mathbf{h}_i,\ \mathbf{c}_3)$

# SimpleRNN + Self-Attention

$\mathbf{c}_1$ $\quad$ $\mathbf{c}_2$ $\quad$ $\mathbf{c}_3$ $\quad$ $\mathbf{c}_4$ $\quad$ $\mathbf{c}_5 = \alpha_1 \mathbf{h}_1 + \alpha_2 \mathbf{h}_2 + \cdots + \alpha_5 \mathbf{h}_5.$

# SimpleRNN + Self-Attention

# Summary

- With self-attention, RNN is less likely to forget.

# Summary

- With self-attention, RNN is less likely to forget.

- Pay attention to the context relevant to the new input.



Figure is from the paper " Long Short-Term Memory-Networks for Machine Reading."

# Thank you!