# Self-Attention

Shusen Wang

# Self-Attention

- Self-Attention: attention beyond Seq2Seq models.

- The original self-attention paper uses LSTM.

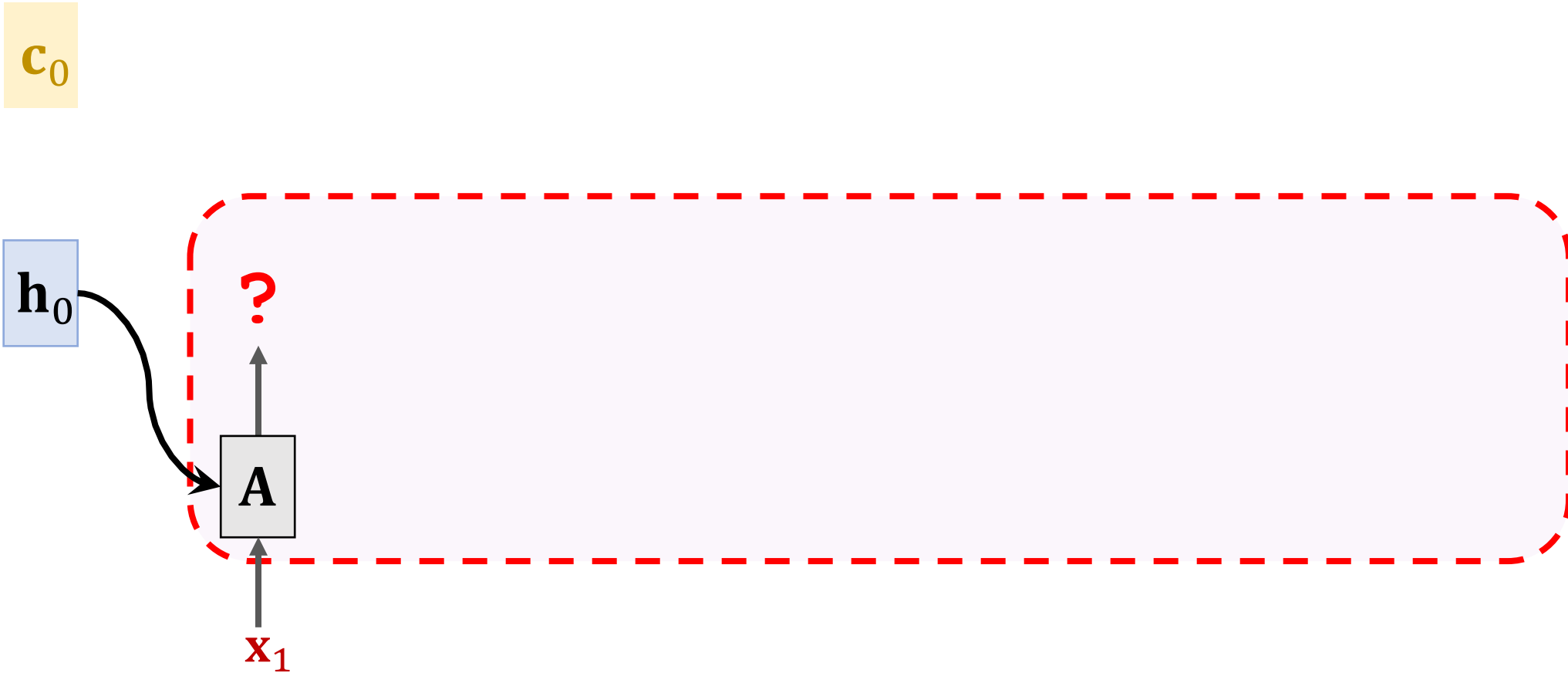- To make teaching easy, I replace LSTM by SimpleRNN.

**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

$c_0$
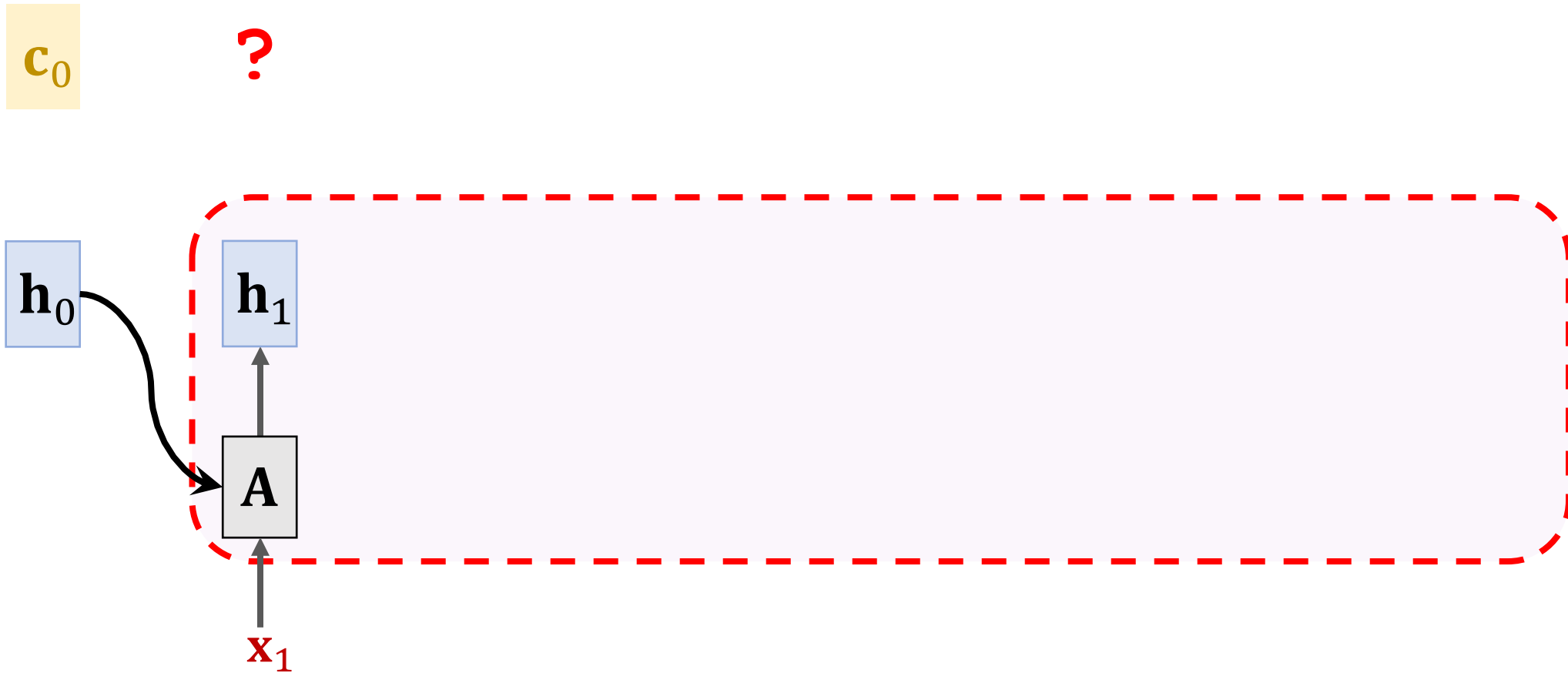
$h_0$

**?**

**A**

$x_1$

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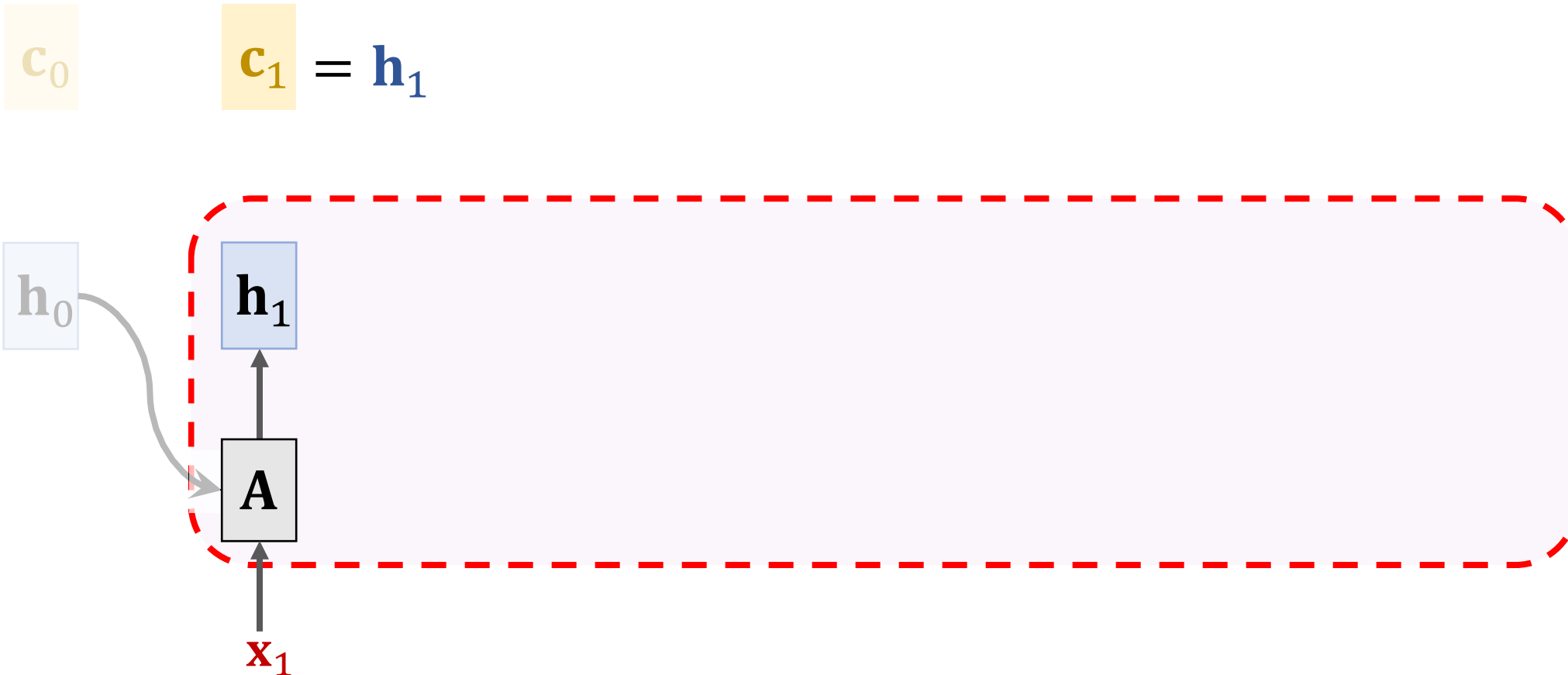
$\mathbf{c}_0$

$\mathbf{h}_0$

**?**

$\mathbf{A}$

$\mathbf{x}_1$

# 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

# SimpleRNN + Self-Attention

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

$\mathbf{c}_1$

$\mathbf{h}_1$

$\mathbf{A}$

$\mathbf{A}$

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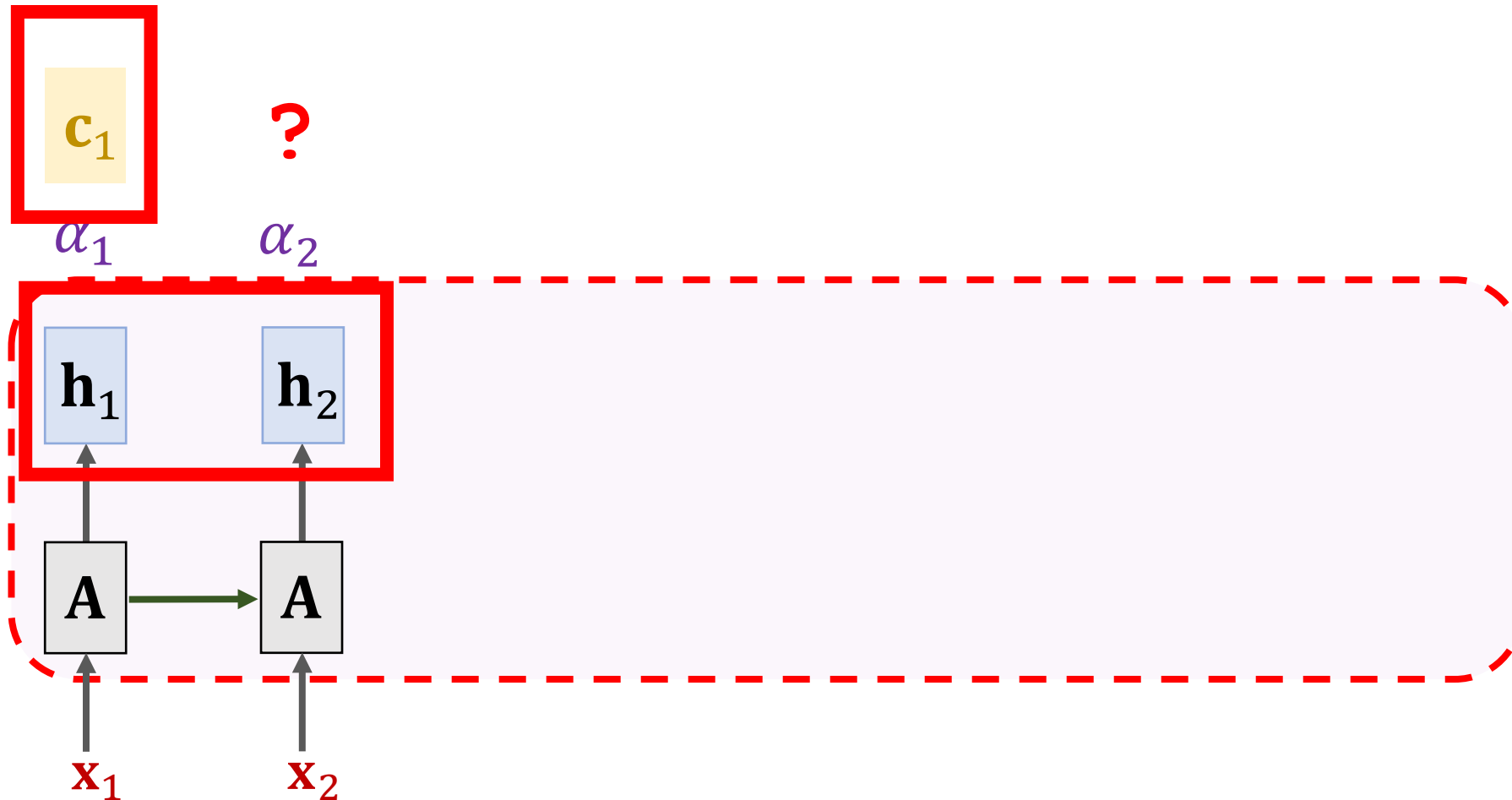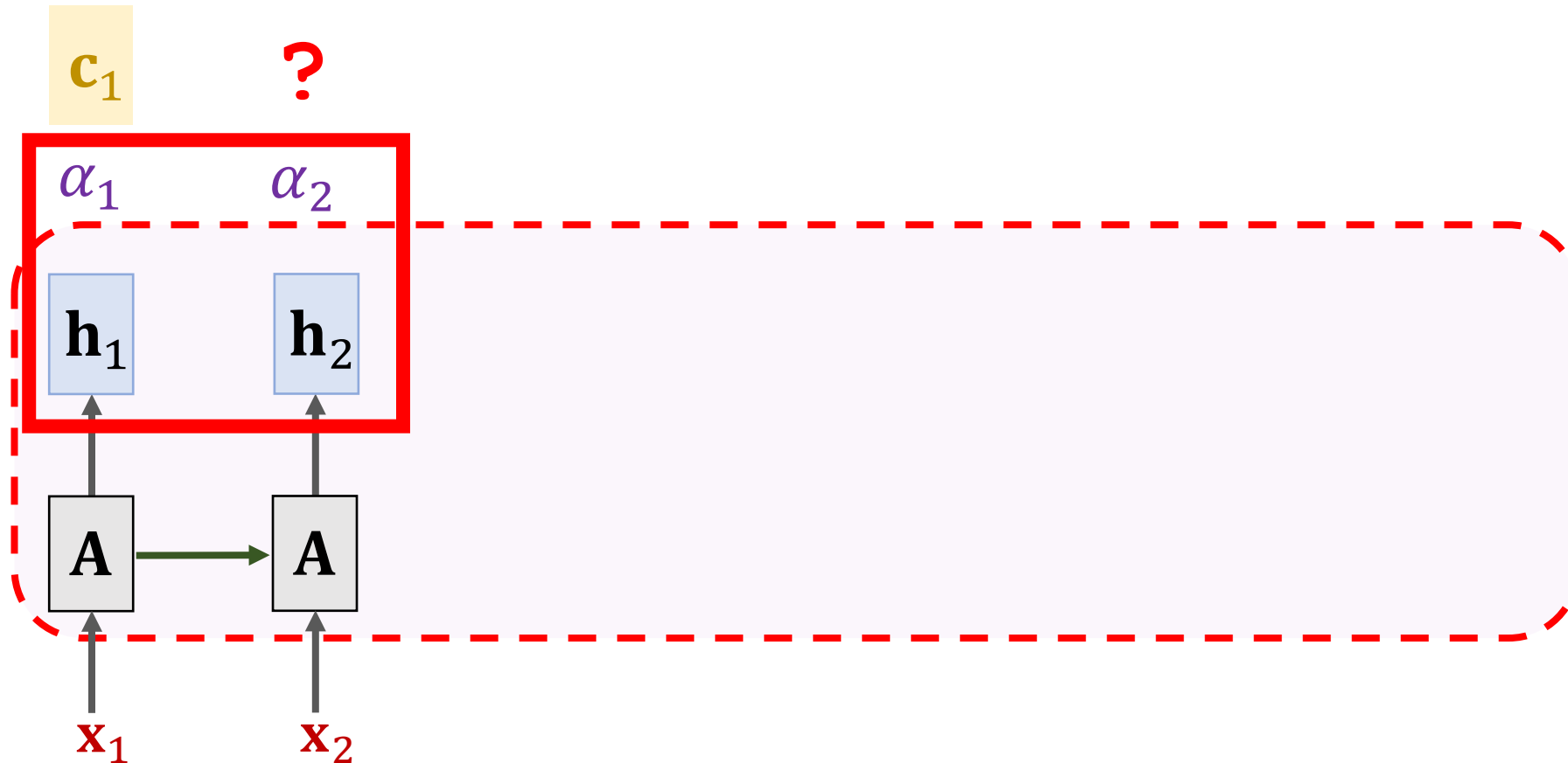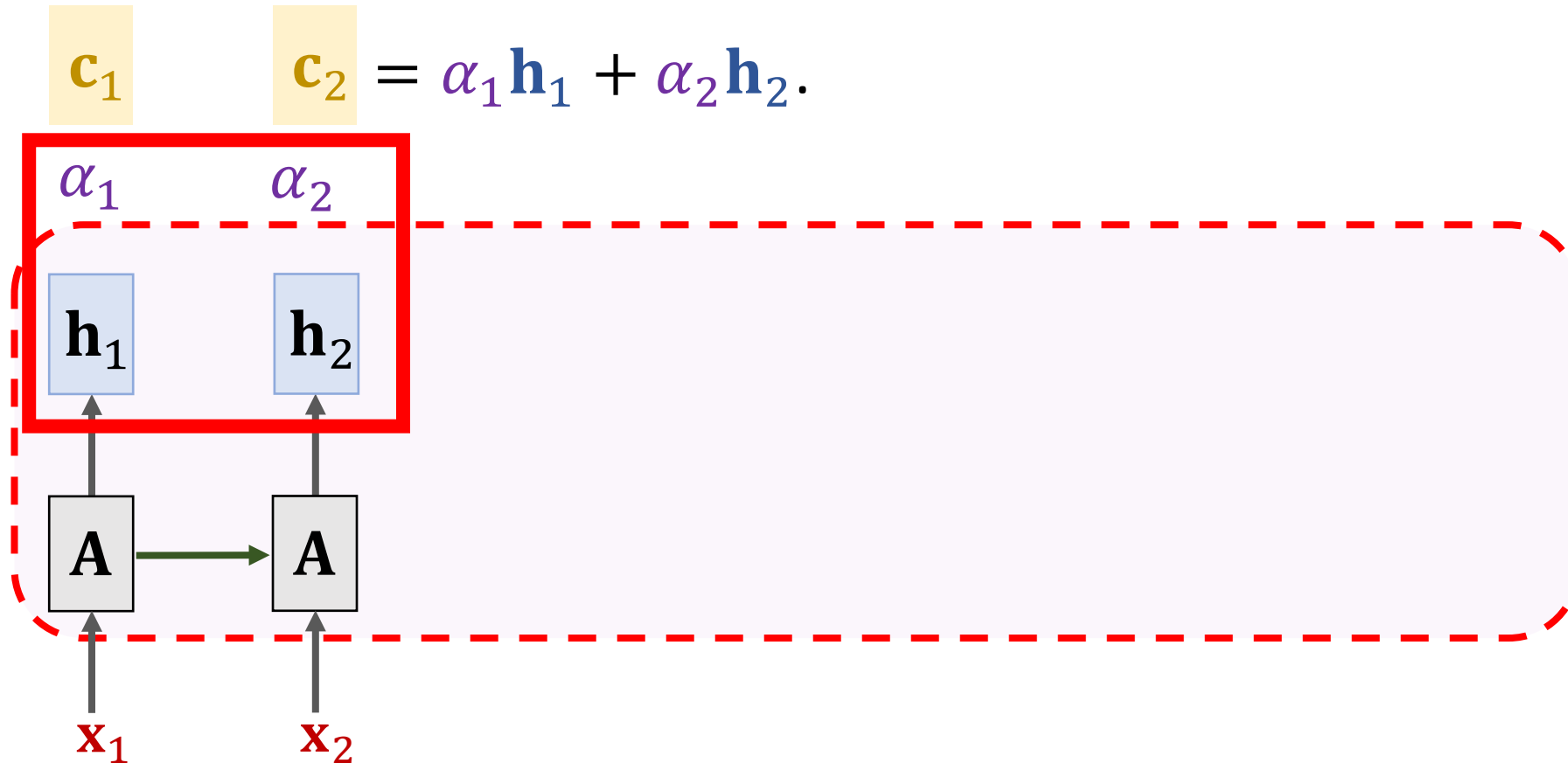
$\mathbf{x}_2$

# SimpleRNN + Self-Attention

# SimpleRNN + Self-Attention

**Weights**: $\alpha_i = \text{align}(\mathbf{h}_i, \mathbf{c}_1)$.

# SimpleRNN + Self-Attention

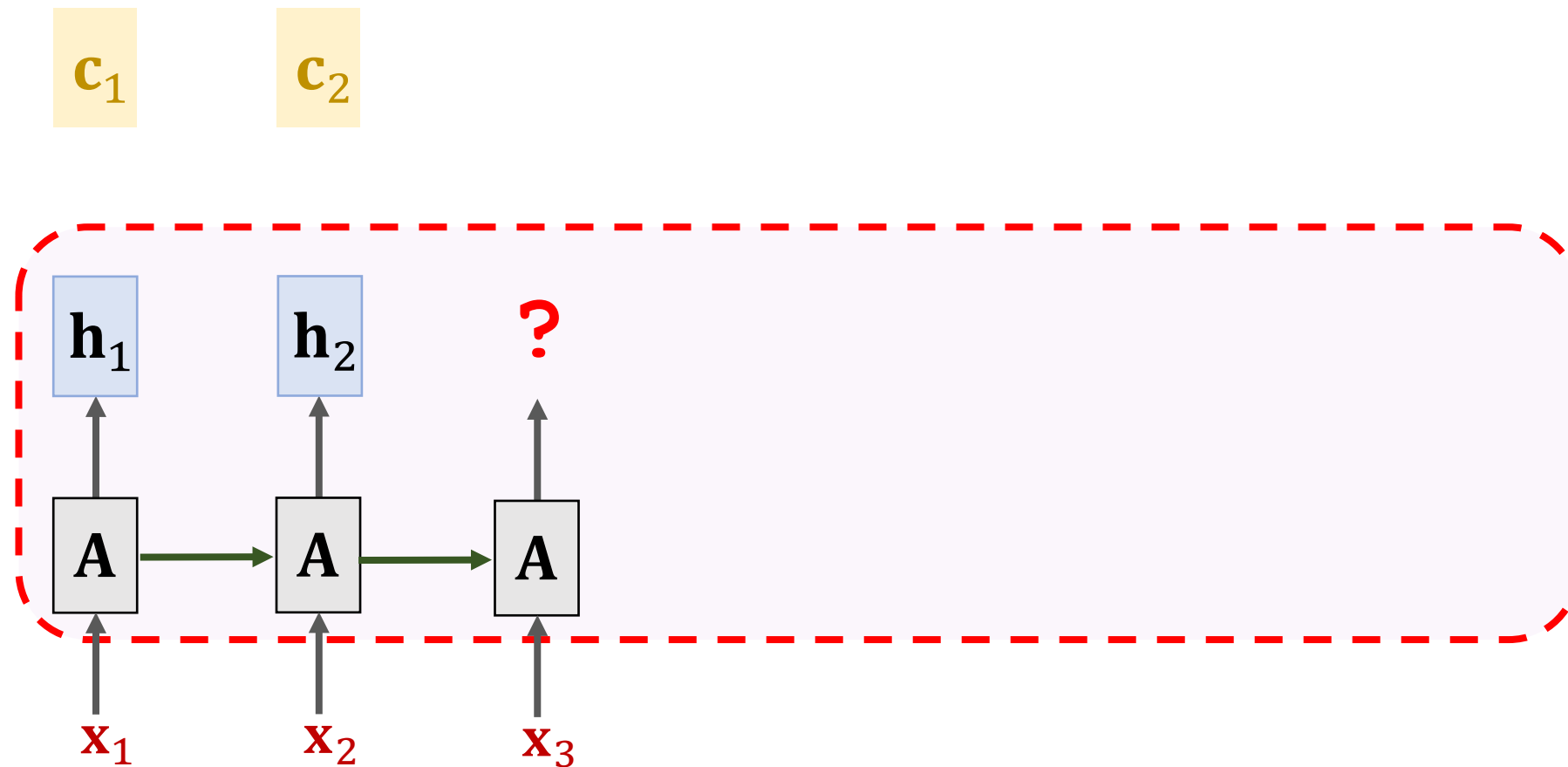**Weights**: $\alpha_i = \text{align}(\mathbf{h}_i, \mathbf{c}_1)$.
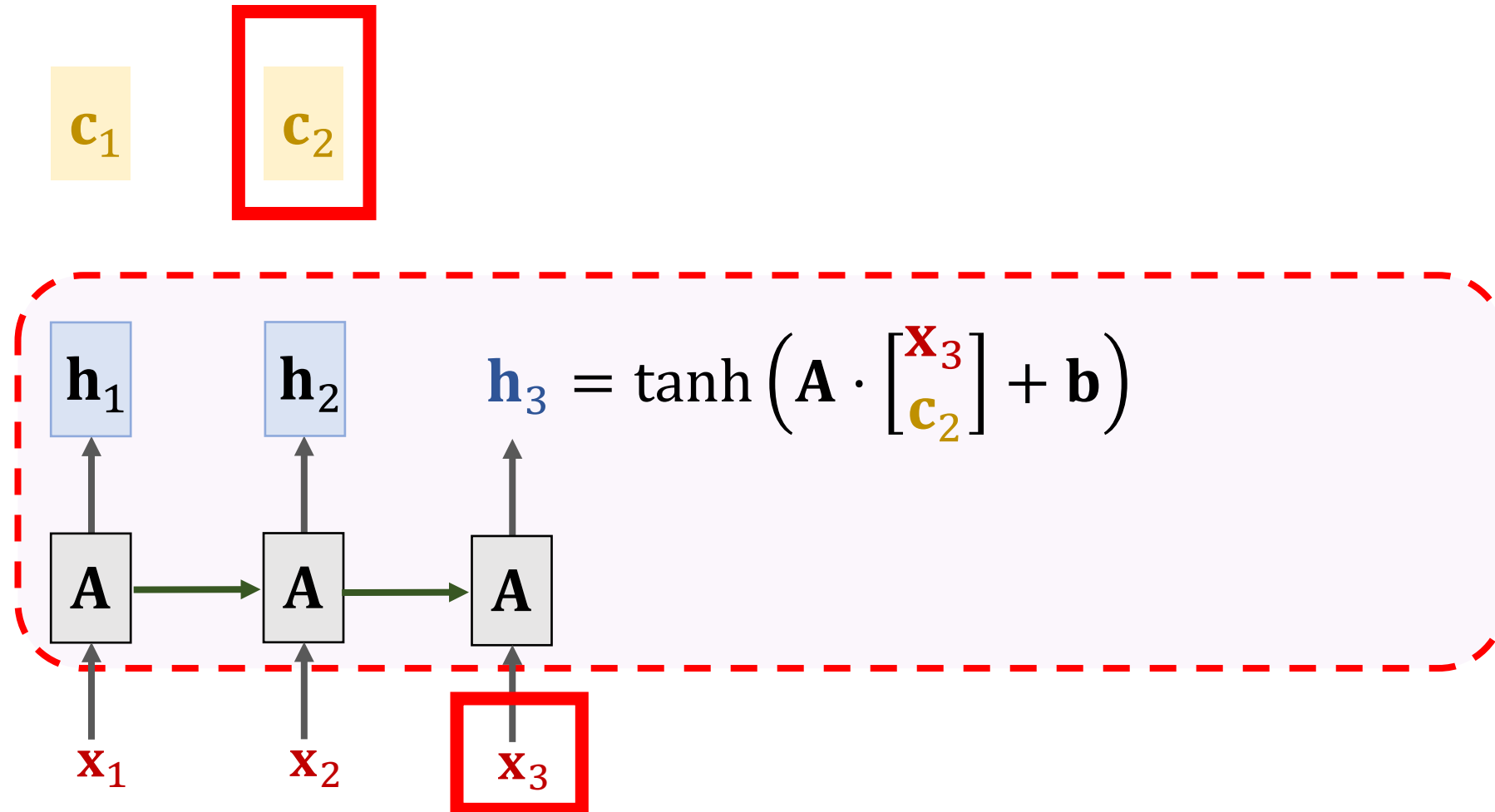
# SimpleRNN + Self-Attention

**Weights**: $\quad \alpha_i = \text{align}(\mathbf{h}_i, \mathbf{c}_1)$.

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

# SimpleRNN + Self-Attention

# SimpleRNN + Self-Attention



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

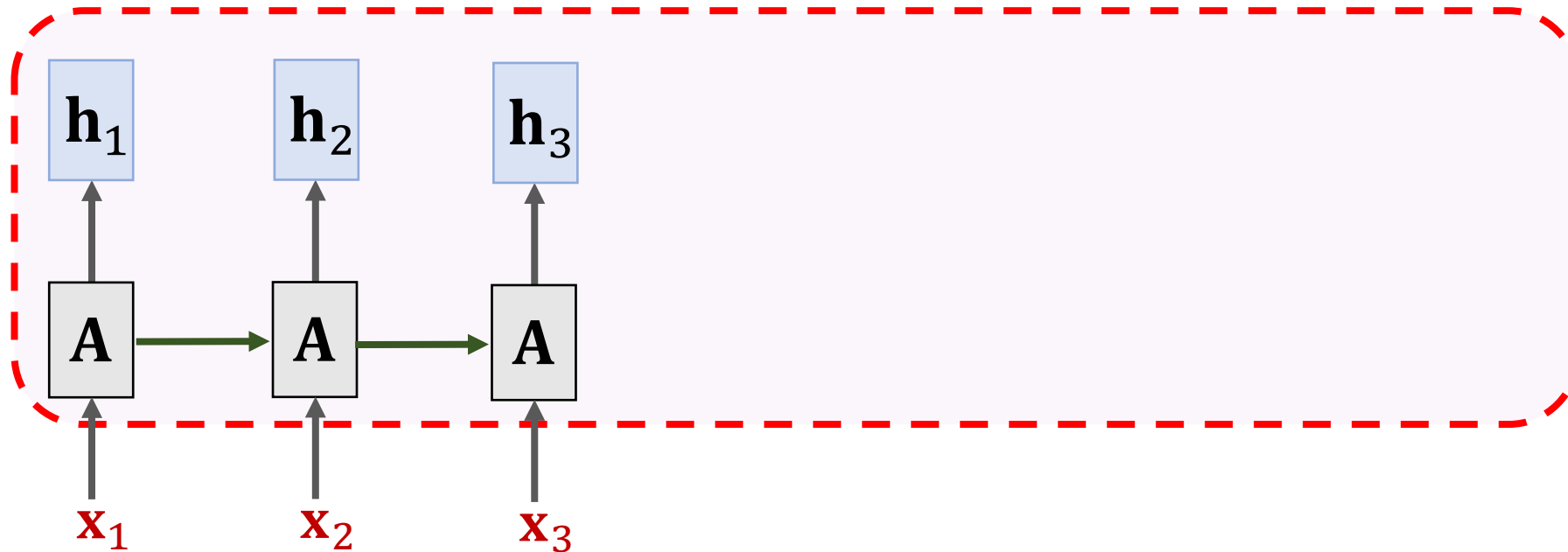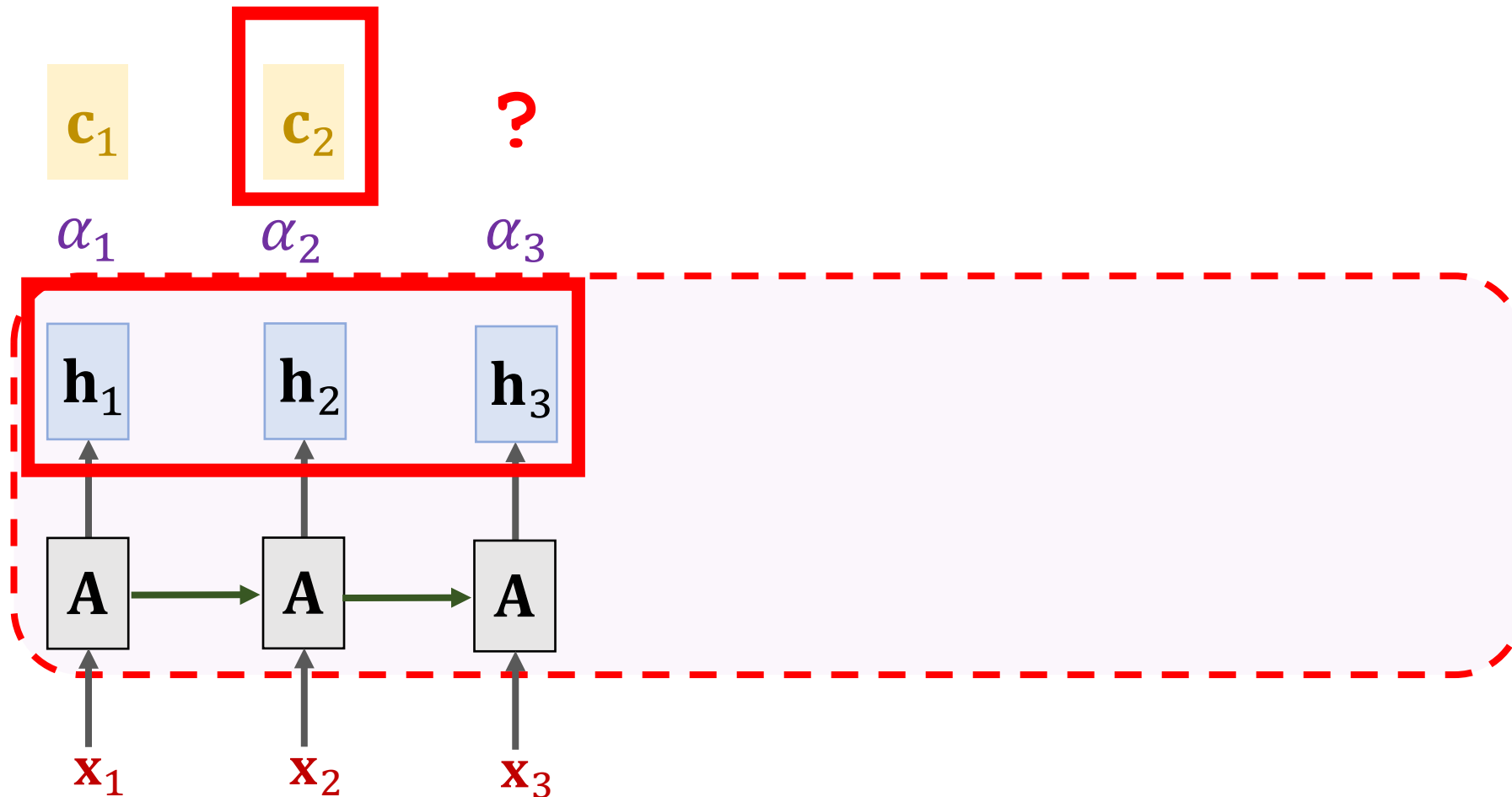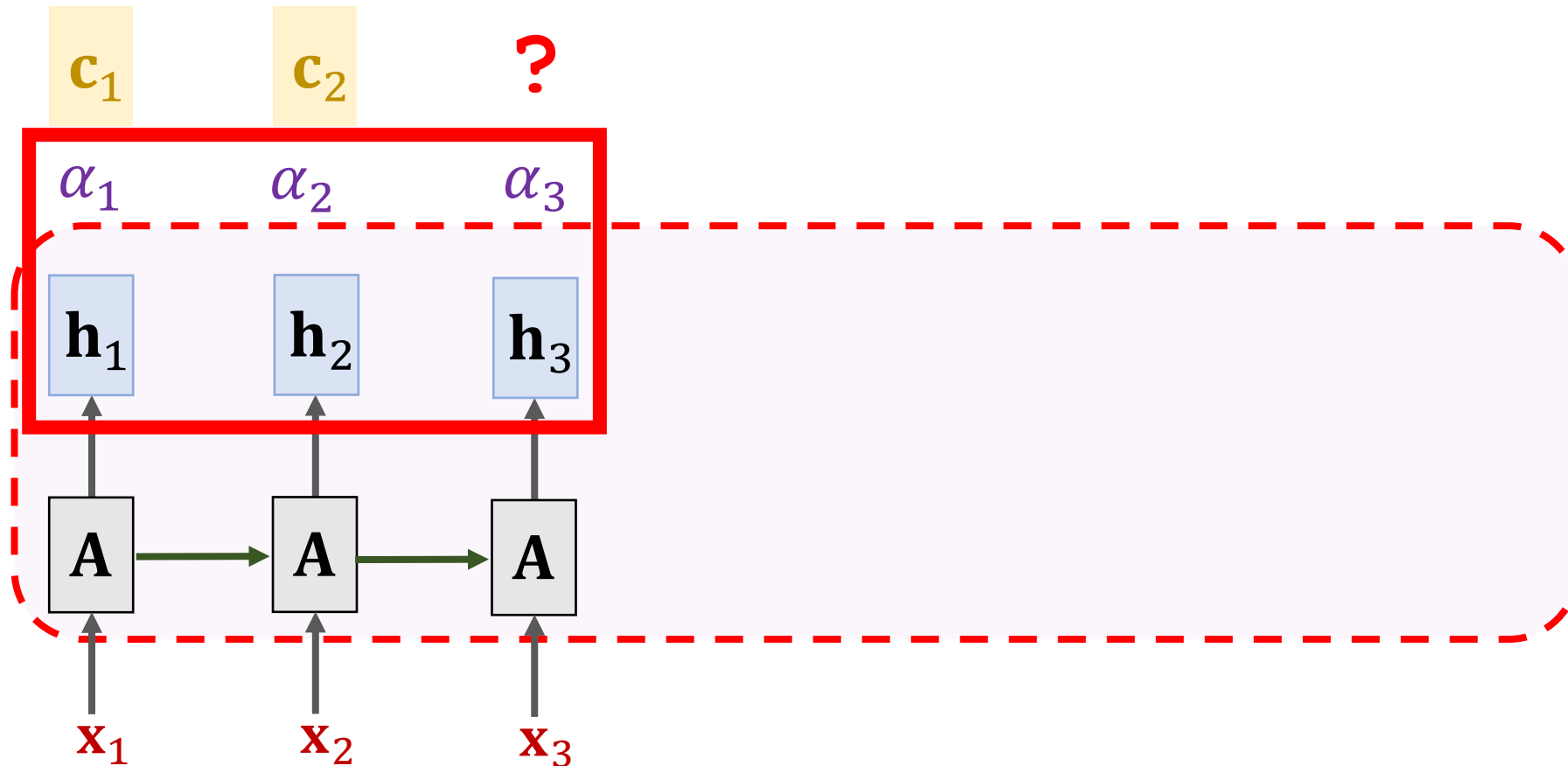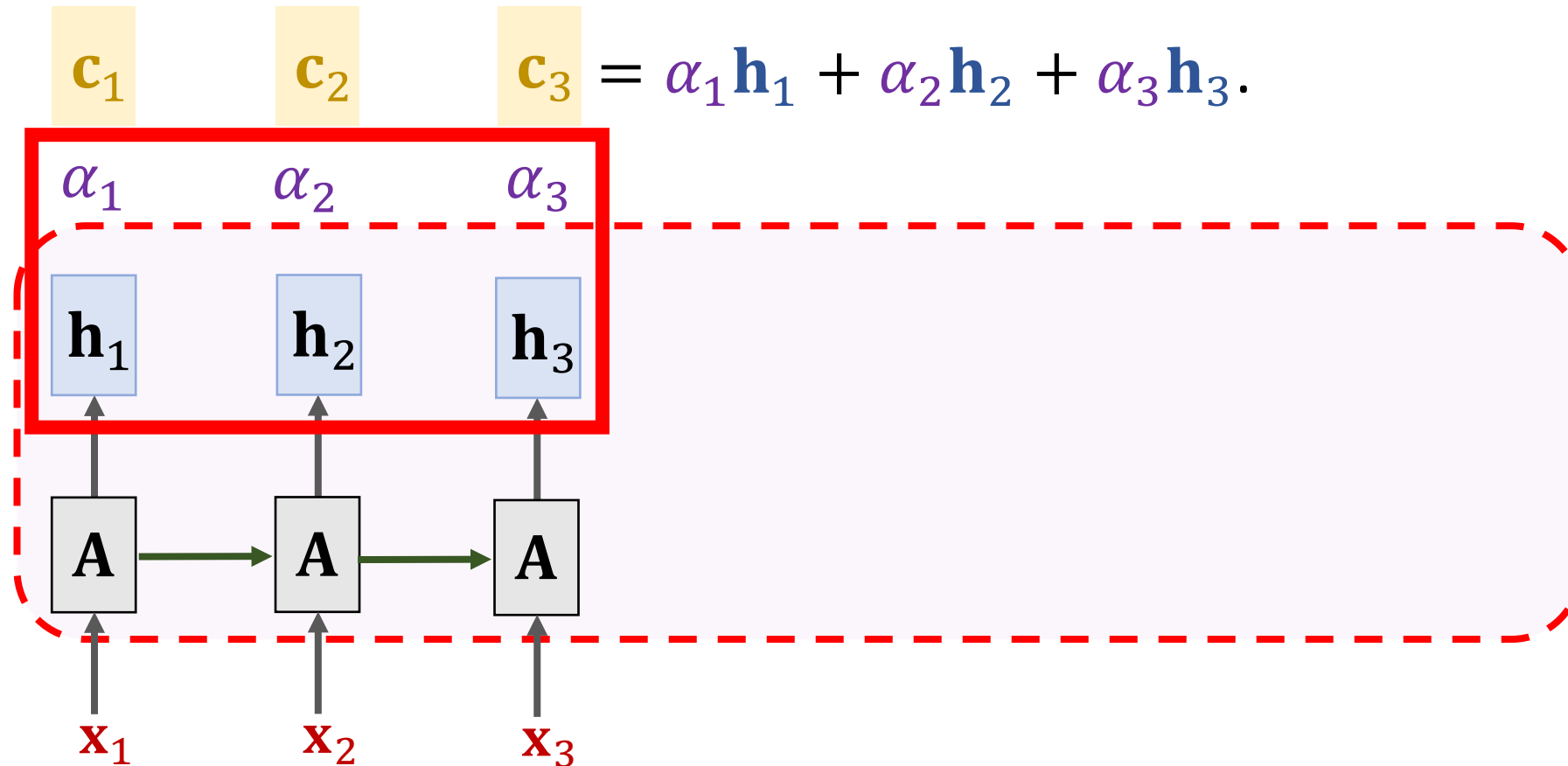# SimpleRNN + Self-Attention

# SimpleRNN + Self-Attention

**Weights**: $\alpha_i = \text{align}(\mathbf{h}_i, \mathbf{c}_2).$

# SimpleRNN + Self-Attention

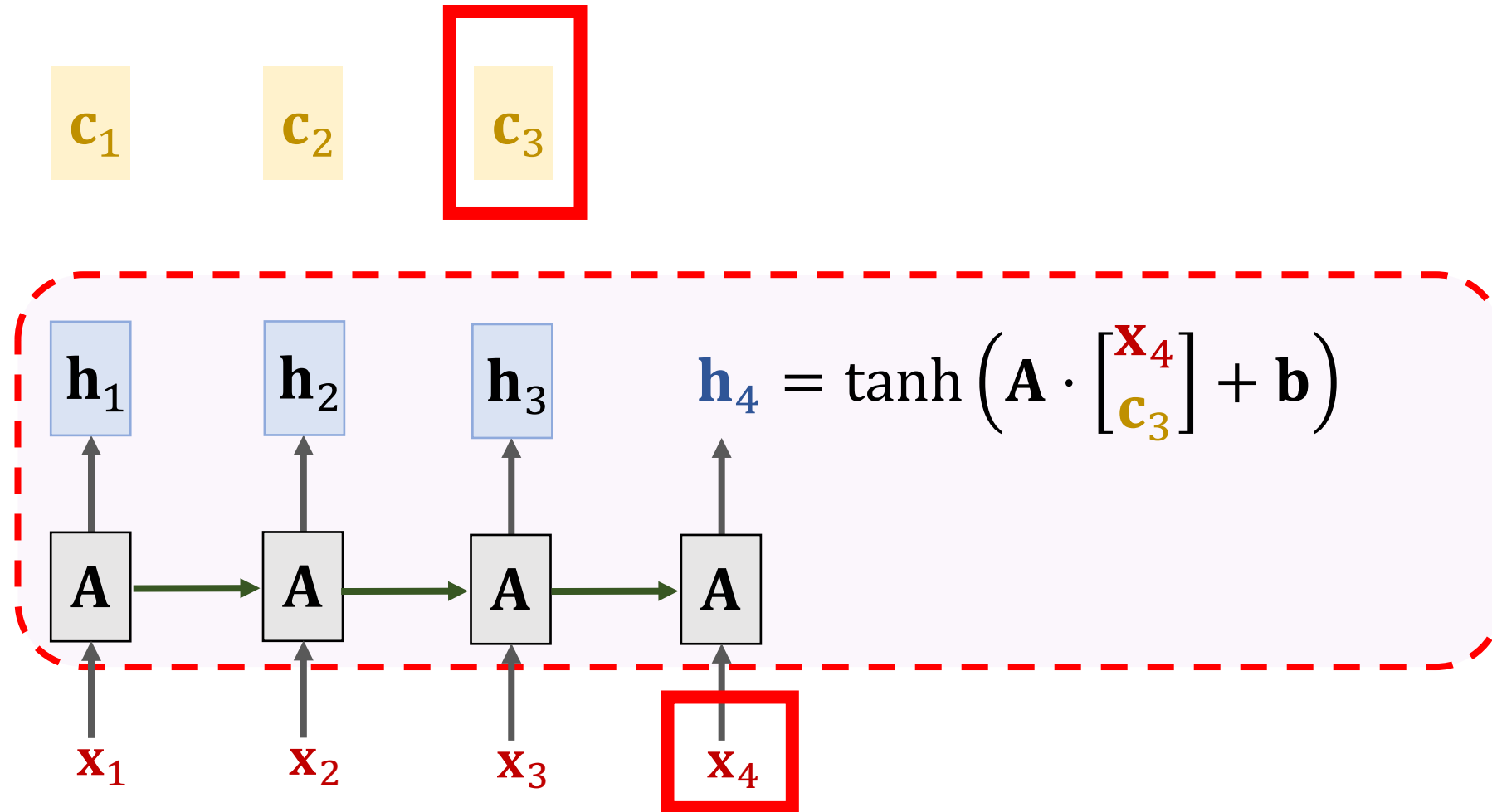**Weights**: $\alpha_i = \text{align}(\mathbf{h}_i, \mathbf{c}_2)$.

# SimpleRNN + Self-Attention

**Weights**: $\quad \alpha_i = \text{align}(\mathbf{h}_i, \ \mathbf{c}_2).$



$\mathbf{c}_1 \qquad \mathbf{c}_2 \qquad \mathbf{c}_3 = \alpha_1 \mathbf{h}_1 + \alpha_2 \mathbf{h}_2 + \alpha_3 \mathbf{h}_3.$

# 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**: $\alpha_i = \text{align}(\mathbf{h}_i, \mathbf{c}_3)$.

# SimpleRNN + Self-Attention

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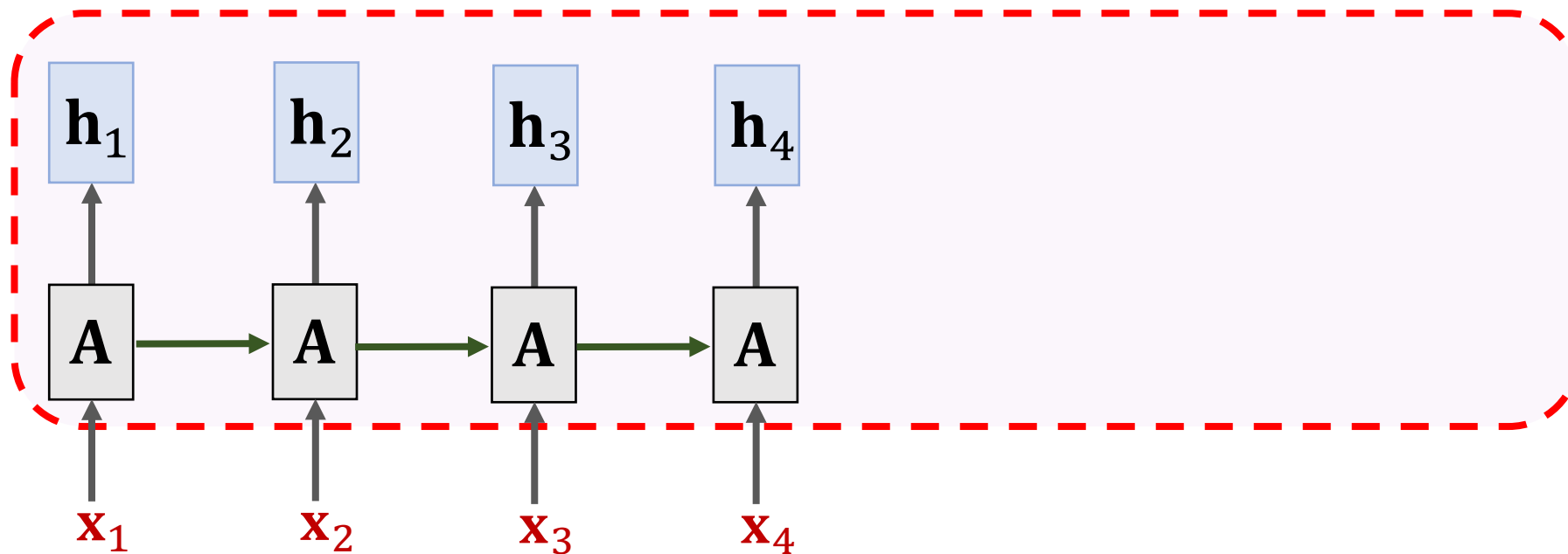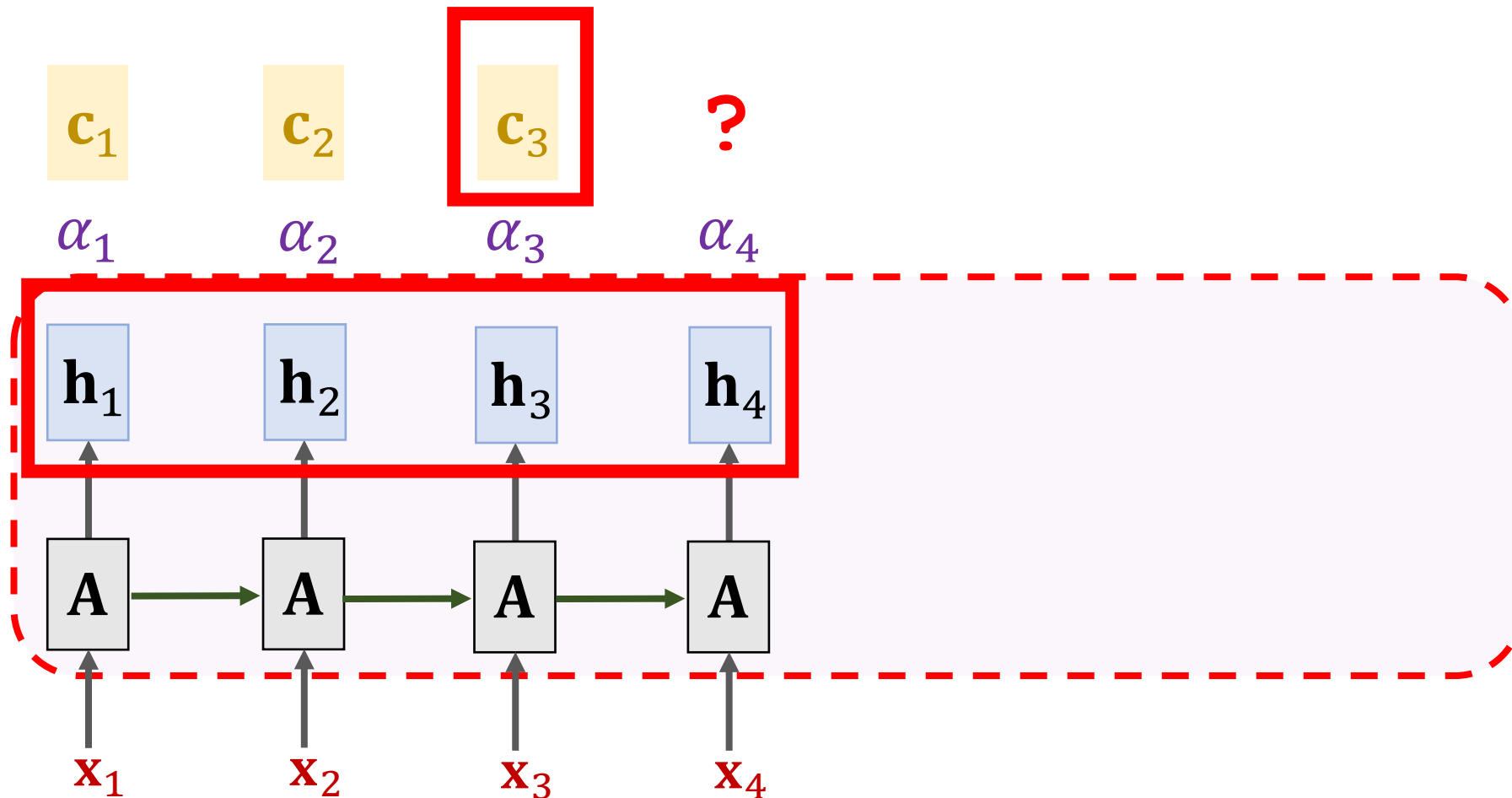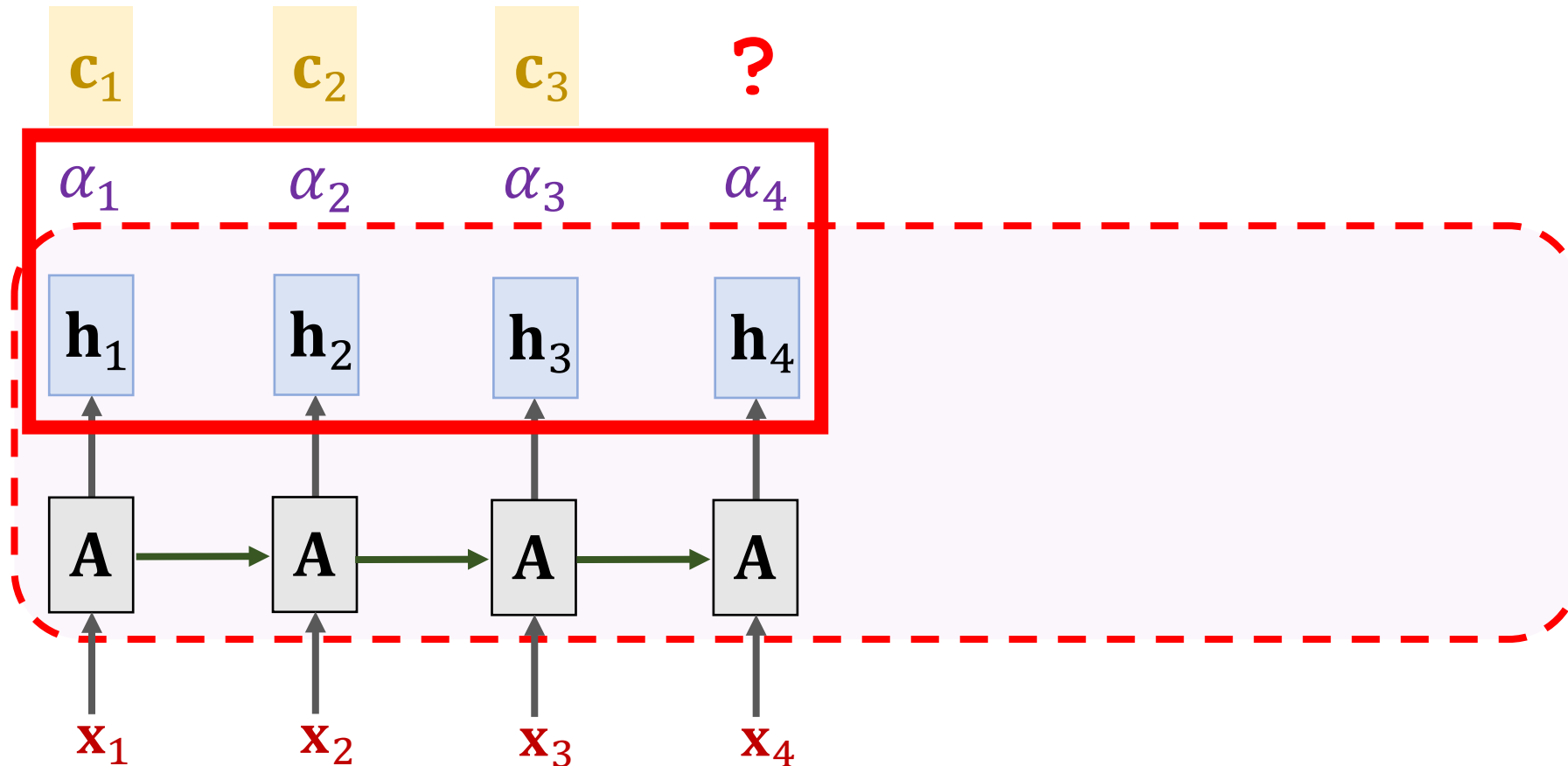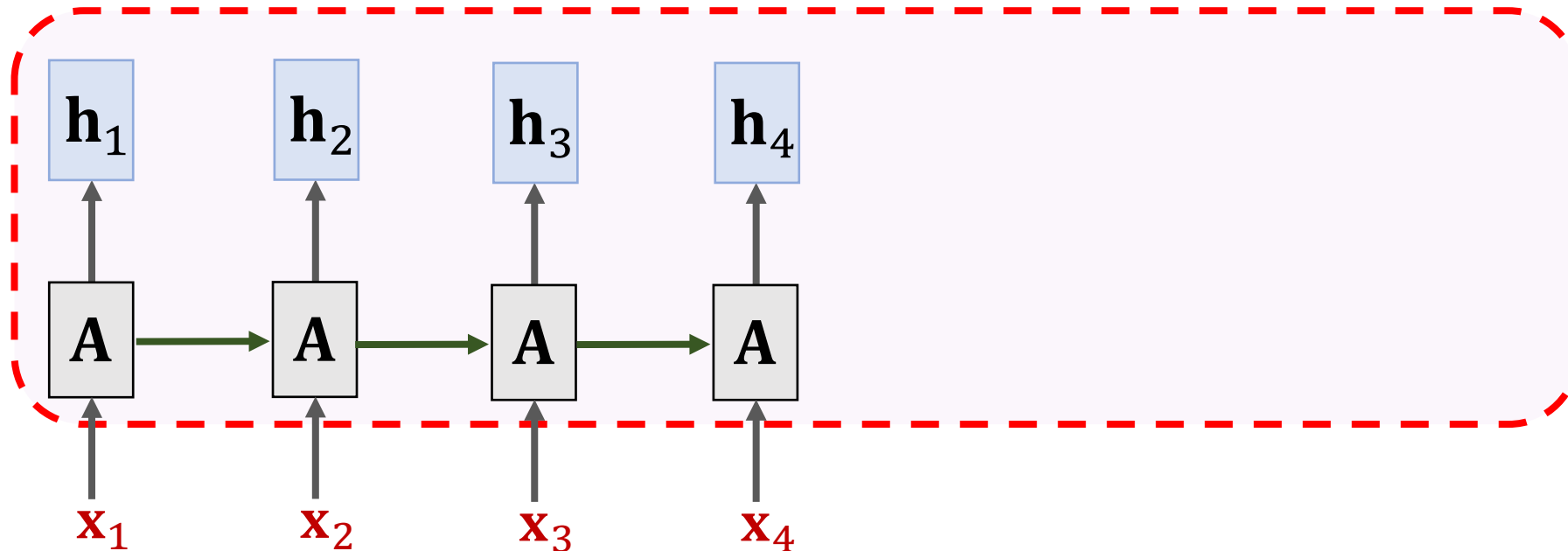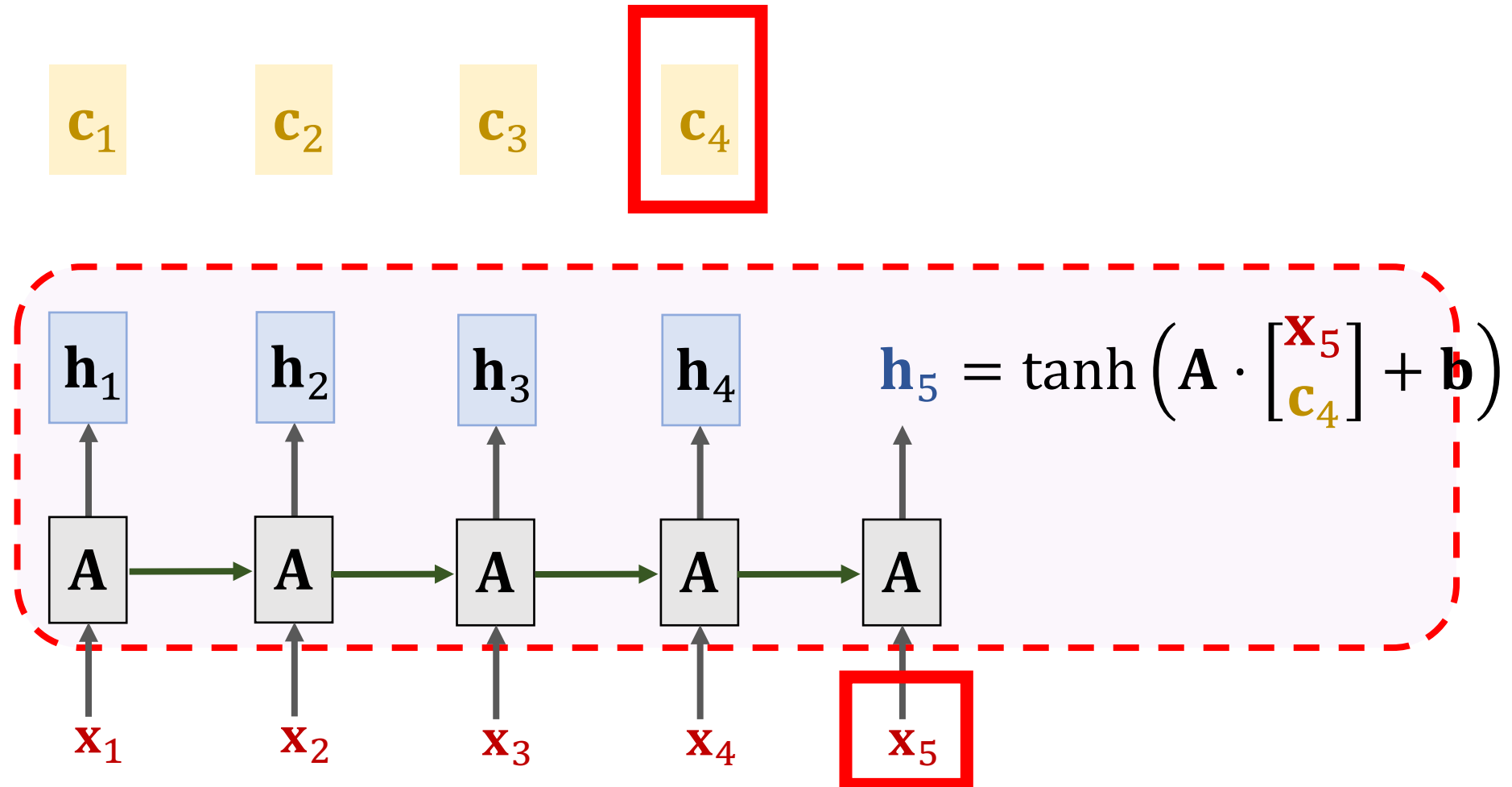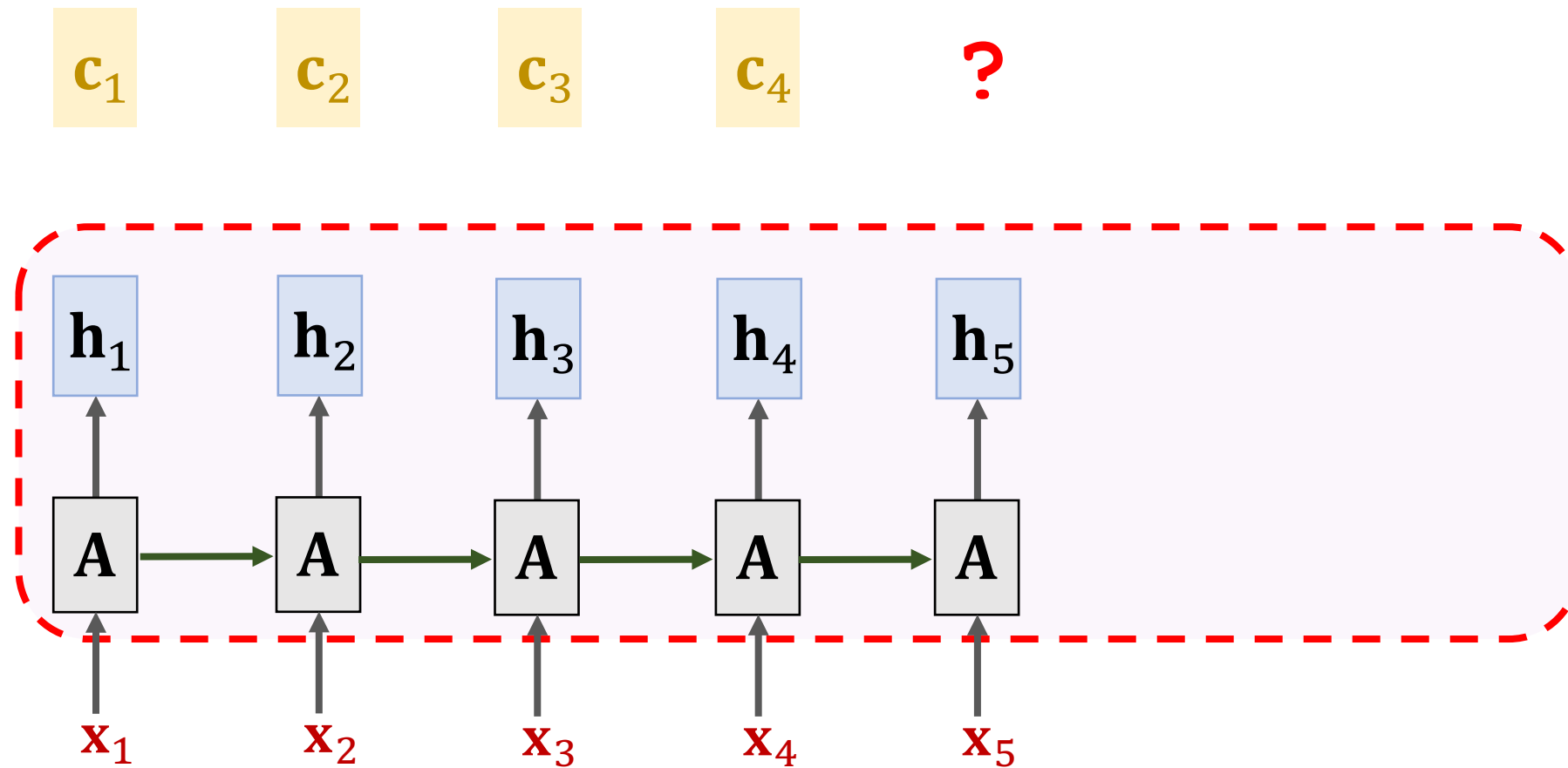
# SimpleRNN + Self-Attention

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

$\mathbf{c}_1 \qquad \mathbf{c}_2 \qquad \mathbf{c}_3 \qquad \mathbf{c}_4 = \alpha_1 \mathbf{h}_1 + \alpha_2 \mathbf{h}_2 + \alpha_3 \mathbf{h}_3 + \alpha_4 \mathbf{h}_4$.
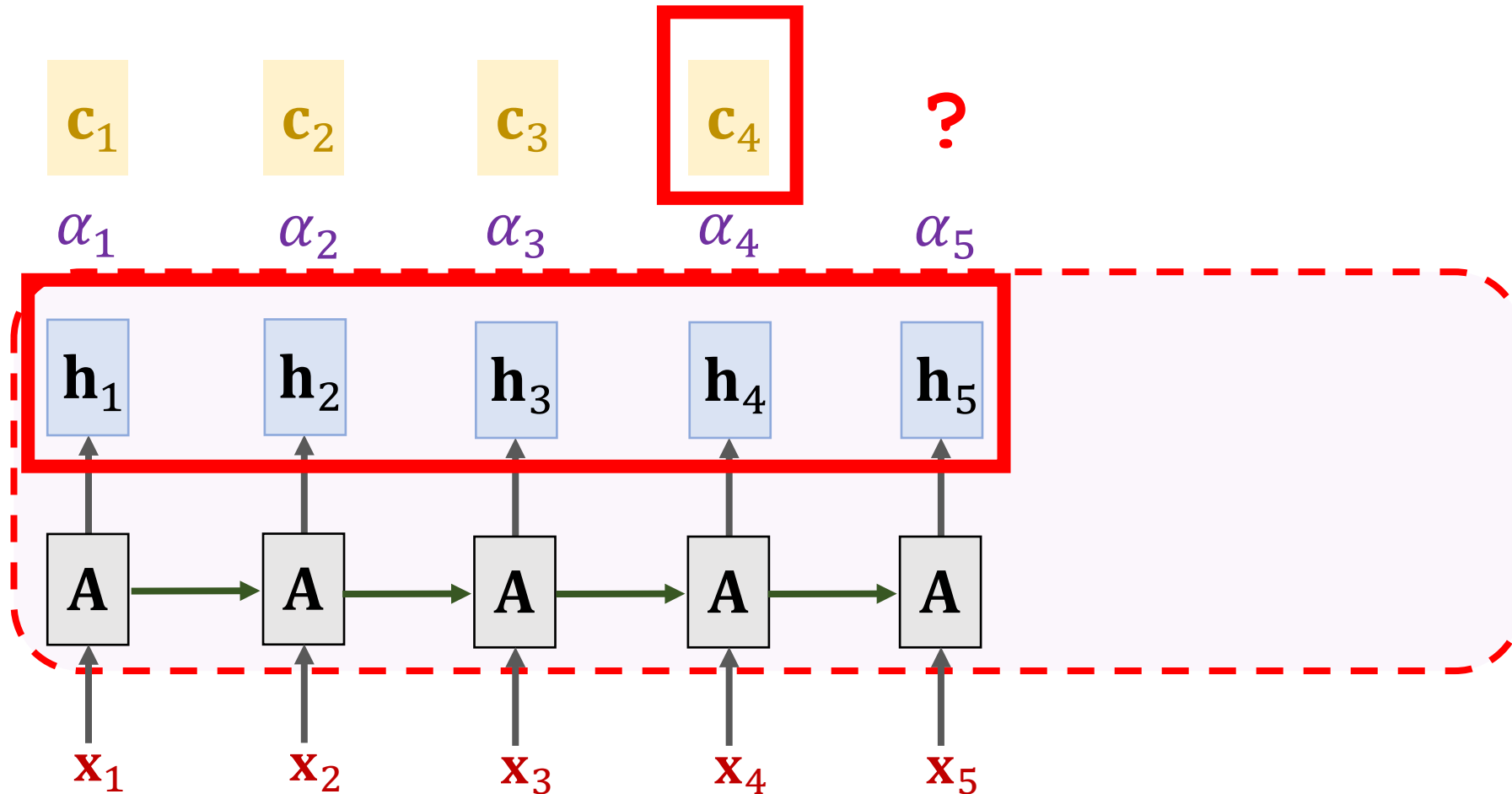
# SimpleRNN + Self-Attention

$c_1$  $c_2$  $c_3$  $c_4$
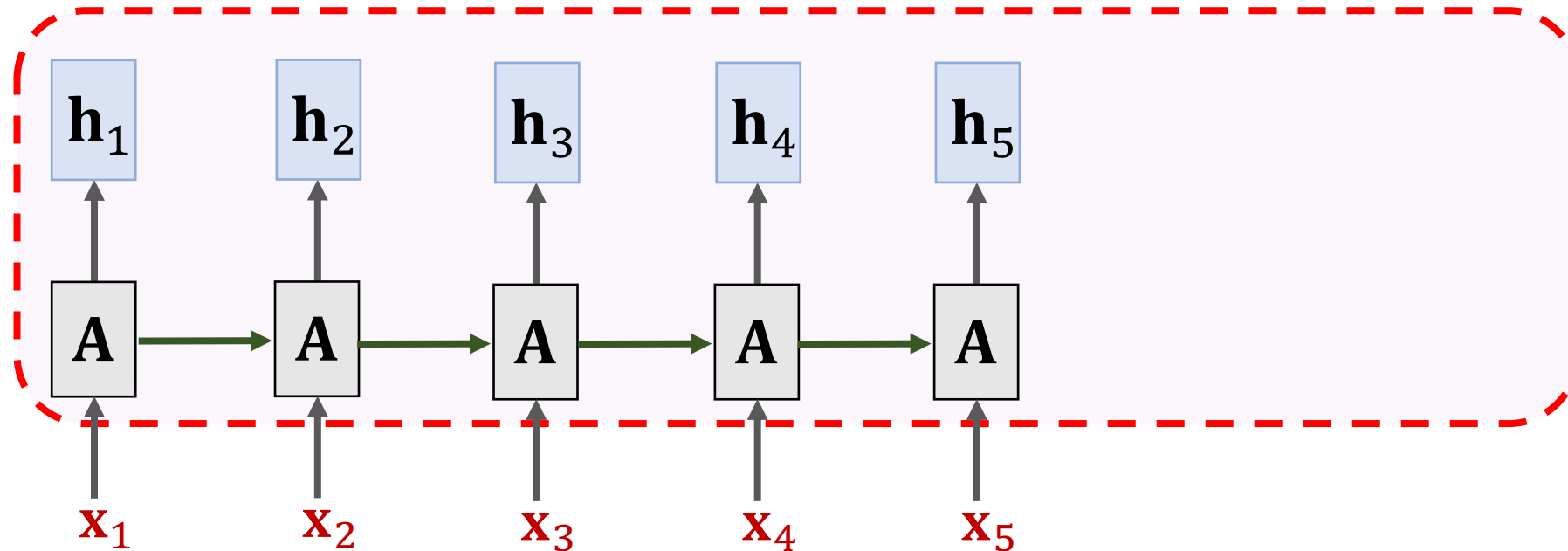
$h_1$  $h_2$  $h_3$  $h_4$  $h_5 = \tanh\left(A \cdot \begin{bmatrix} x_5 \\ c_4 \end{bmatrix} + b\right)$

A   A   A   A   A

$x_1$  $x_2$  $x_3$  $x_4$  $x_5$

# SimpleRNN + Self-Attention

$\mathbf{c}_1$     $\mathbf{c}_2$     $\mathbf{c}_3$     $\mathbf{c}_4$     **?**

$\mathbf{h}_1$     $\mathbf{h}_2$     $\mathbf{h}_3$     $\mathbf{h}_4$     $\mathbf{h}_5$

**A** → **A** → **A** → **A** → **A**

$\mathbf{x}_1$     $\mathbf{x}_2$     $\mathbf{x}_3$     $\mathbf{x}_4$     $\mathbf{x}_5$

# SimpleRNN + Self-Attention

# SimpleRNN + Self-Attention

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

# 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!