# Topics in AI (CPSC 532S):
## Multimodal Learning with Vision, Language and Sound

**Lecture 11: RNNs (Part 3), Applications**

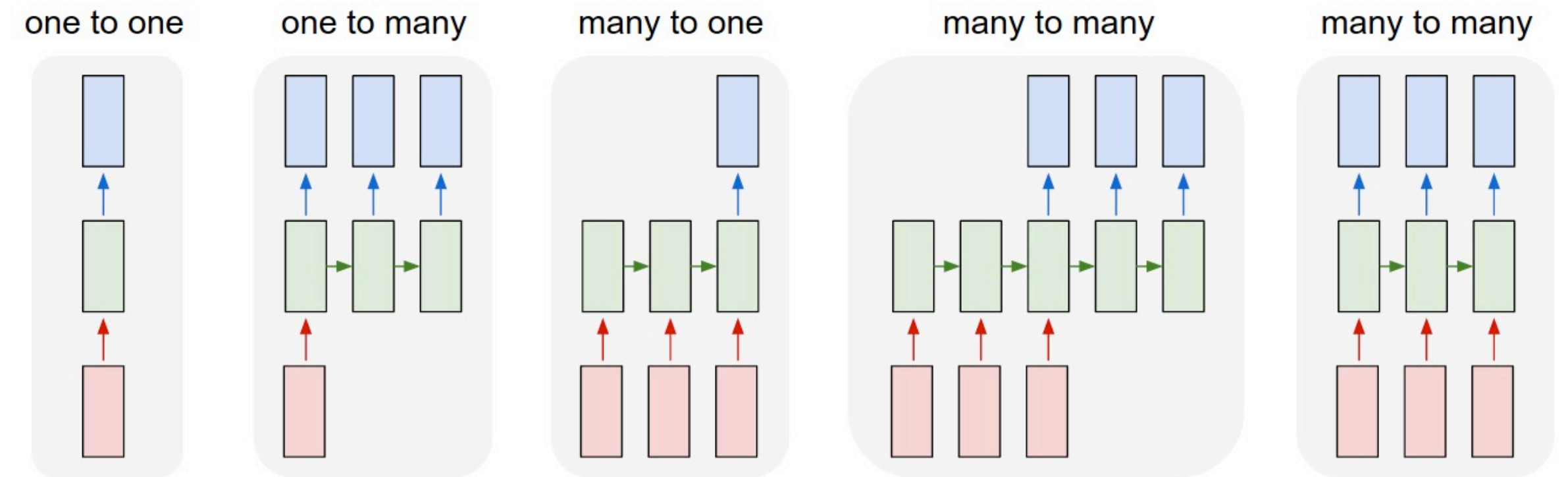# Course **Logistics**

— **Assignment 3** due date is Monday -> Wednesday

— **Assignment 4** is released Monday

— **Assignment 1** & **2** solutions are out
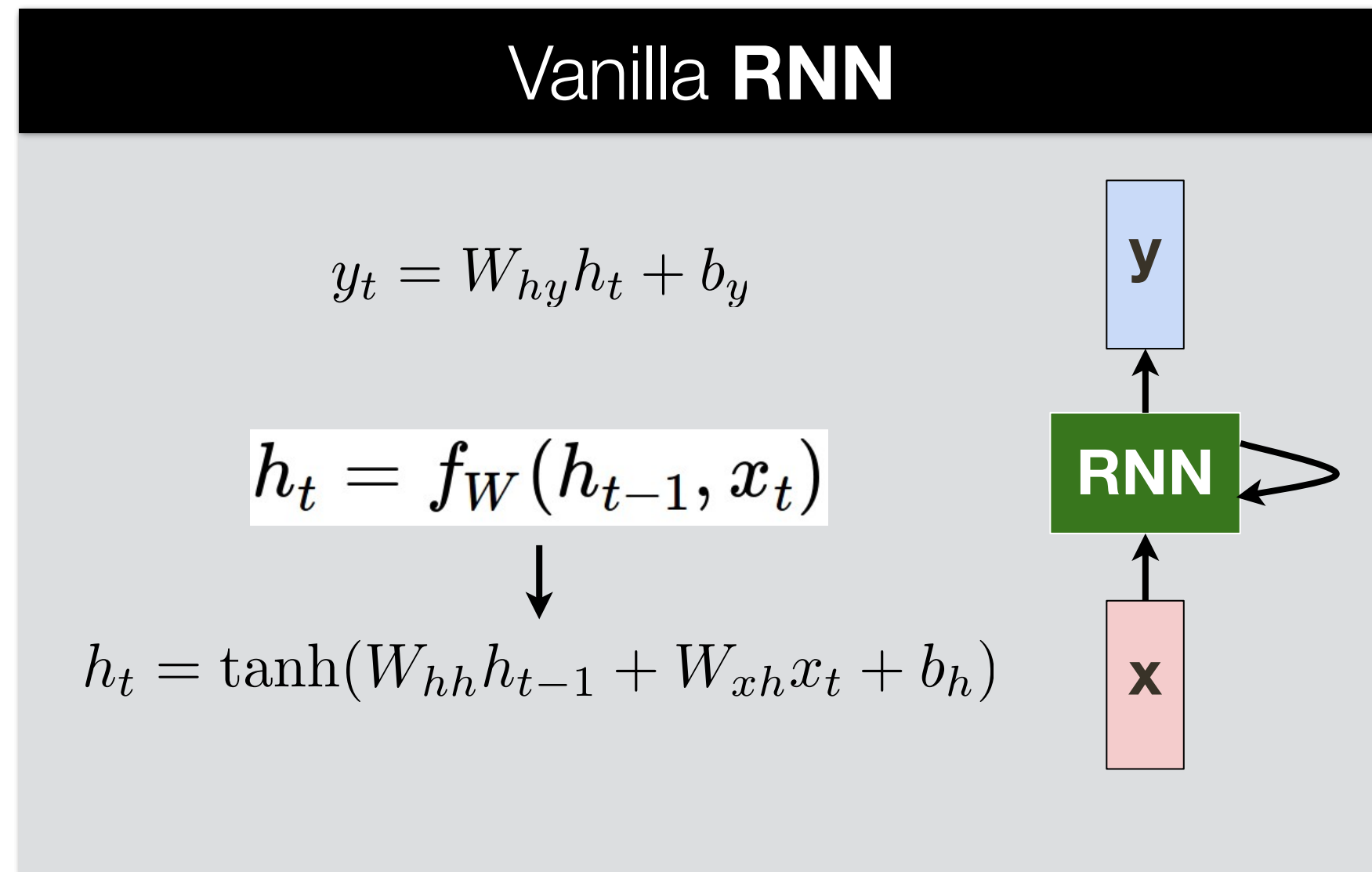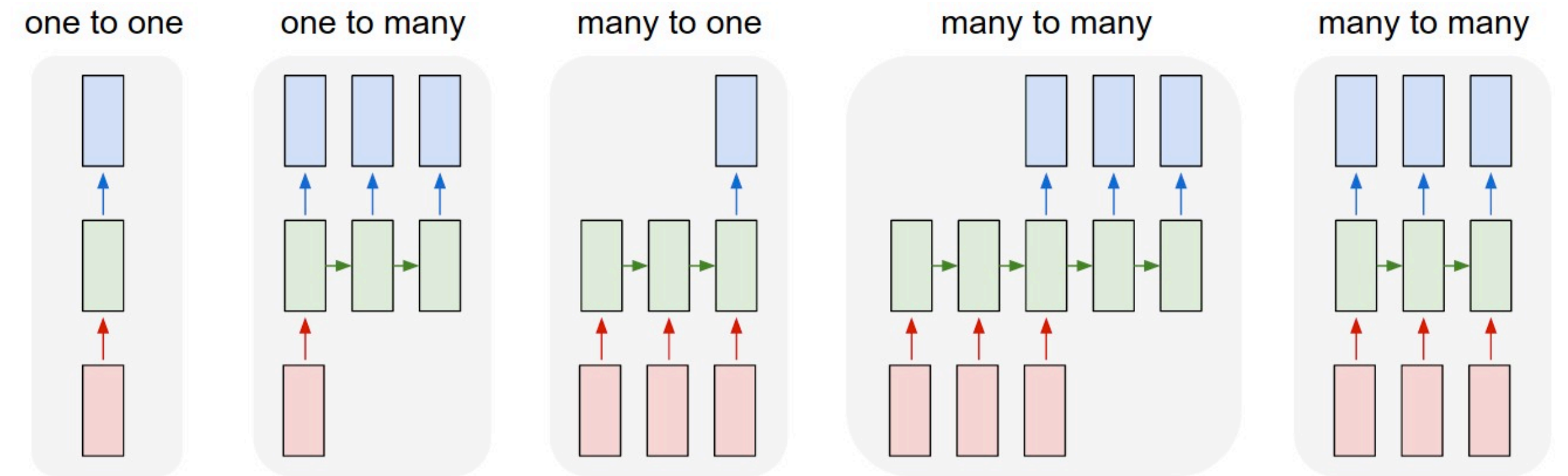
# **RNNs**: Review

## **Key Enablers:**

— Parameter sharing in computational graphs

— "Unrolling" in computational graphs

— Allows modeling **arbitrary length sequences**!



one to one    one to many    many to one    many to many    many to many
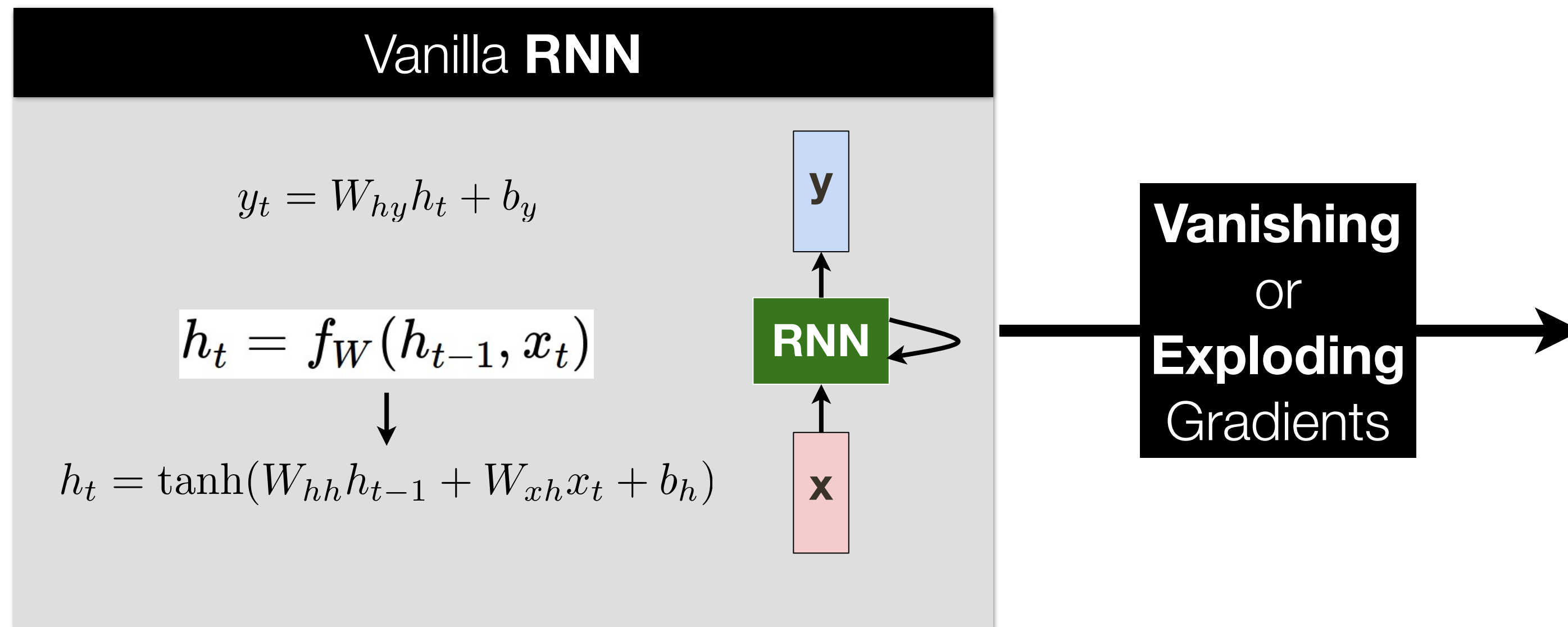
# **RNNs**: Review

## **Key Enablers:**

— Parameter sharing in computational graphs

— "Unrolling" in computational graphs

— Allows modeling **arbitrary length sequences**!



one to one    one to many    many to one    many to many    many to many

---



Vanilla **RNN**

$$y_t = W_{hy}h_t + b_y$$

$$h_t = f_W(h_{t-1}, x_t)$$

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t + b_h)$$
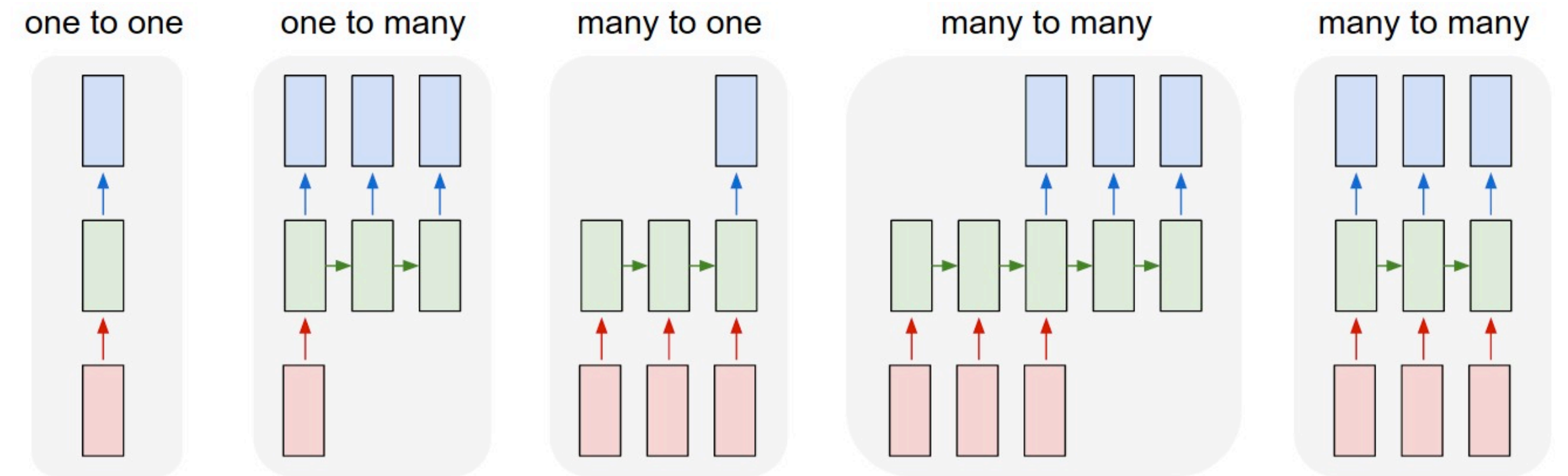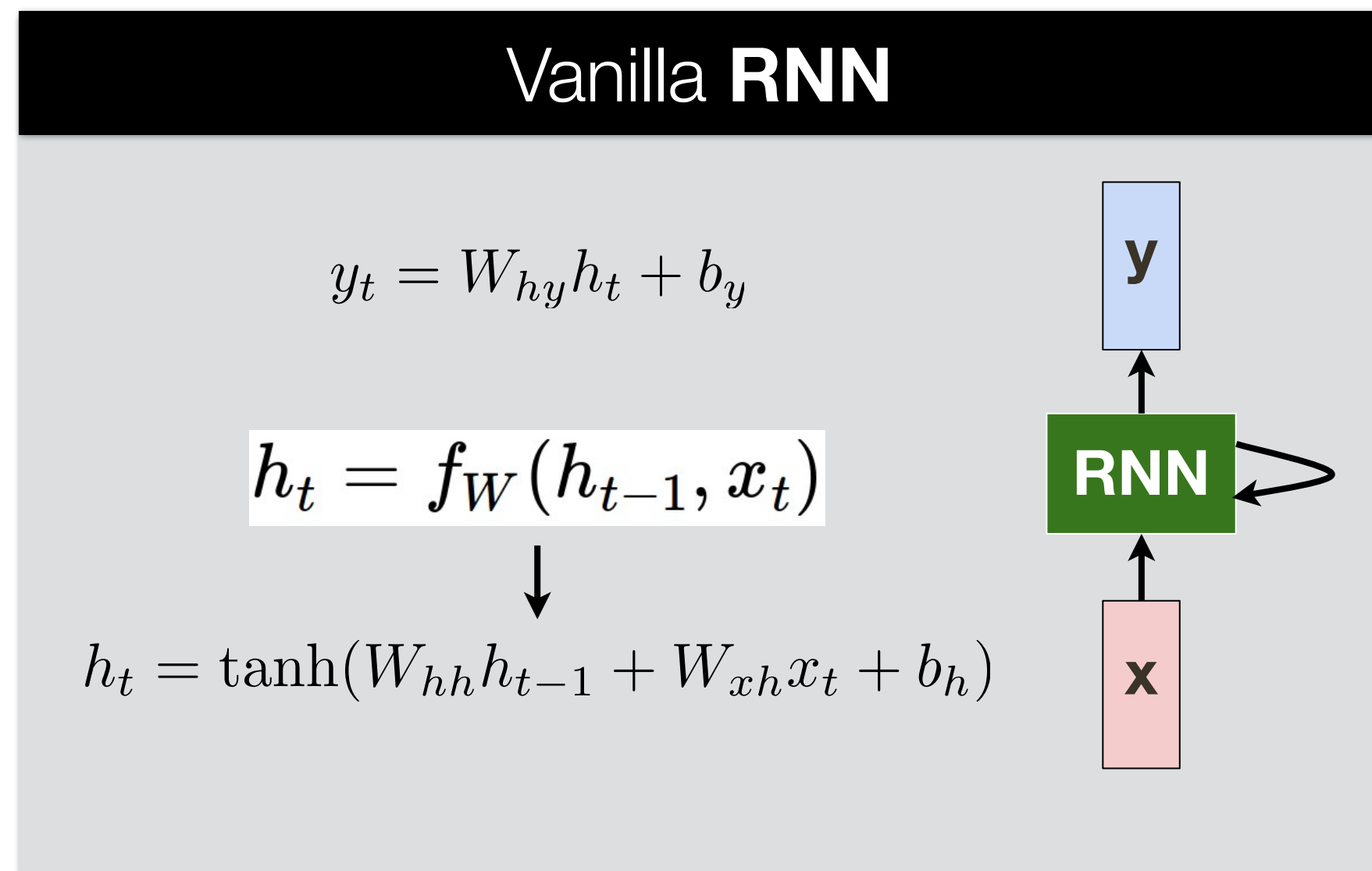
**y**

**RNN**

**x**

# **RNNs**: Review

## **Key Enablers:**

— Parameter sharing in computational graphs

— "Unrolling" in computational graphs

— Allows modeling **arbitrary length sequences**!

| one to one | one to many | many to one | many to many | many to many |

Vanilla **RNN**

$$y_t = W_{hy}h_t + b_y$$

$$h_t = f_W(h_{t-1}, x_t)$$

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t + b_h)$$

**y**

**RNN**

**x**

**Vanishing**
or
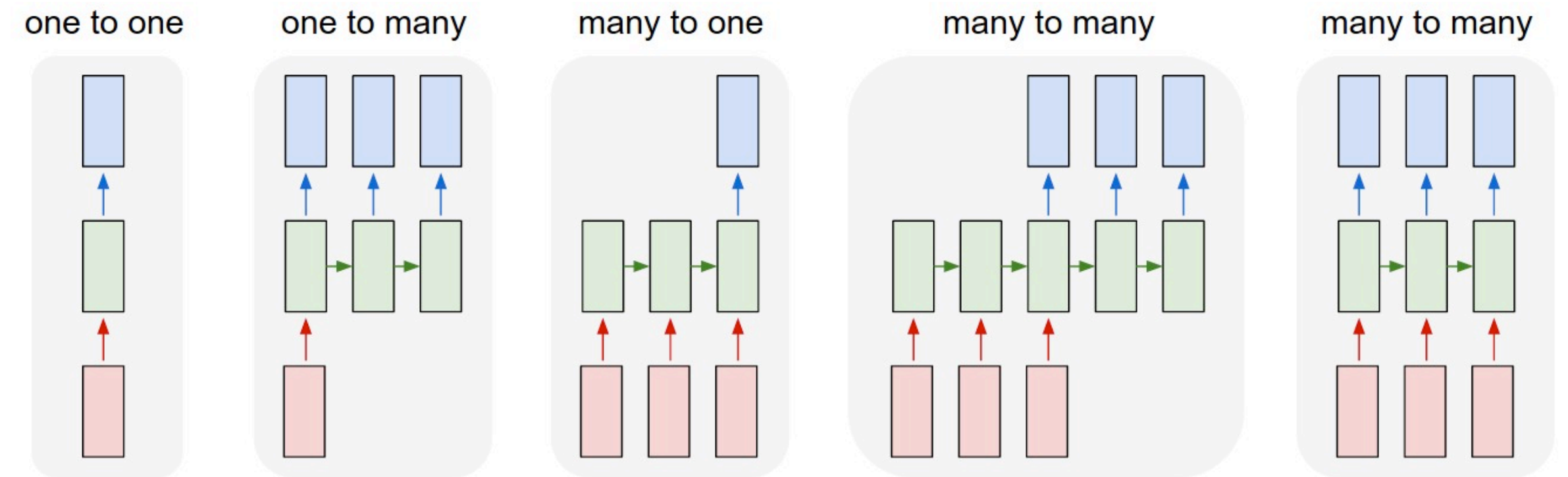**Exploding**
Gradients

# **RNNs**: Review

## **Key Enablers:**

— Parameter sharing in computational graphs

— "Unrolling" in computational graphs
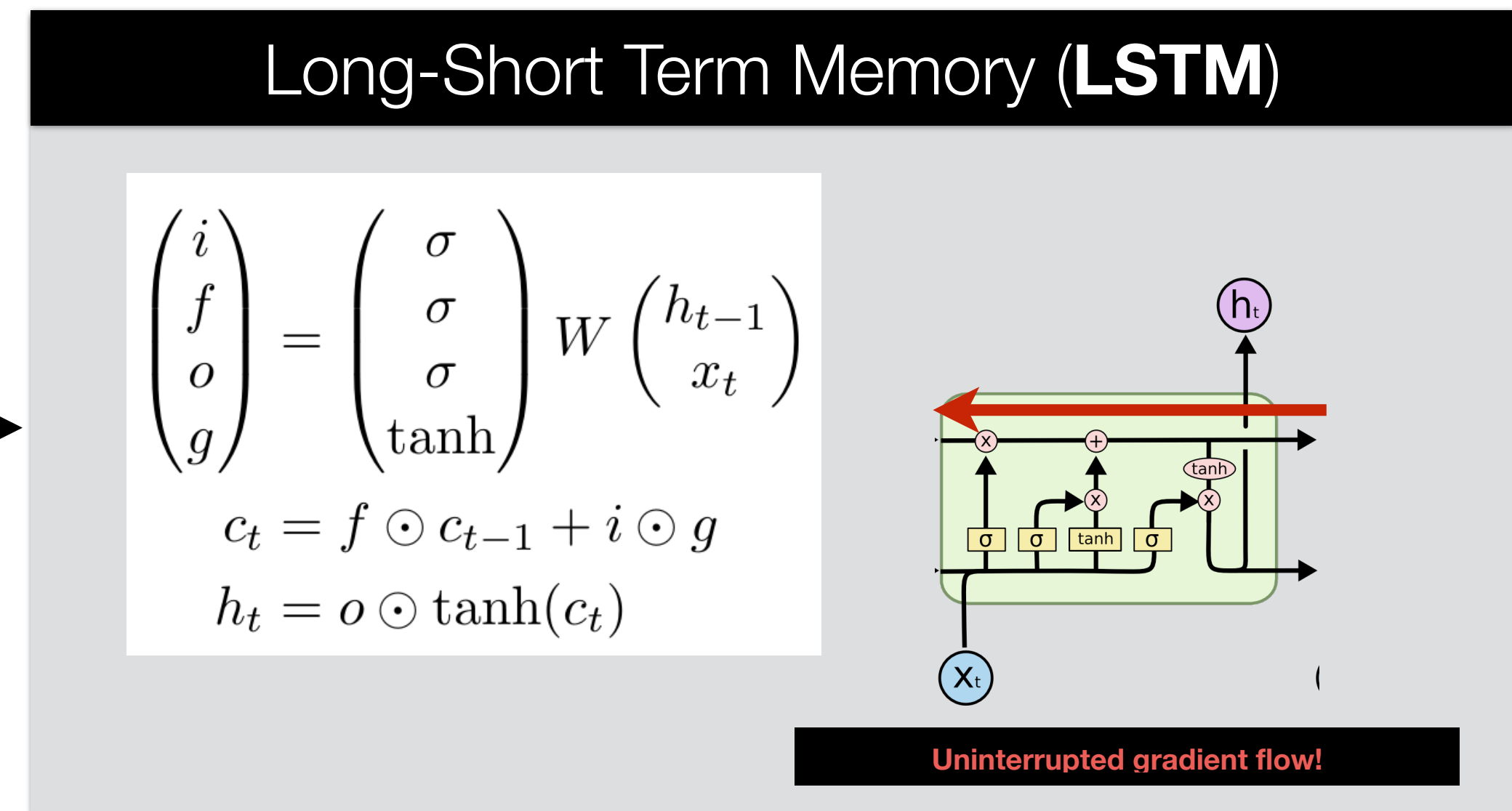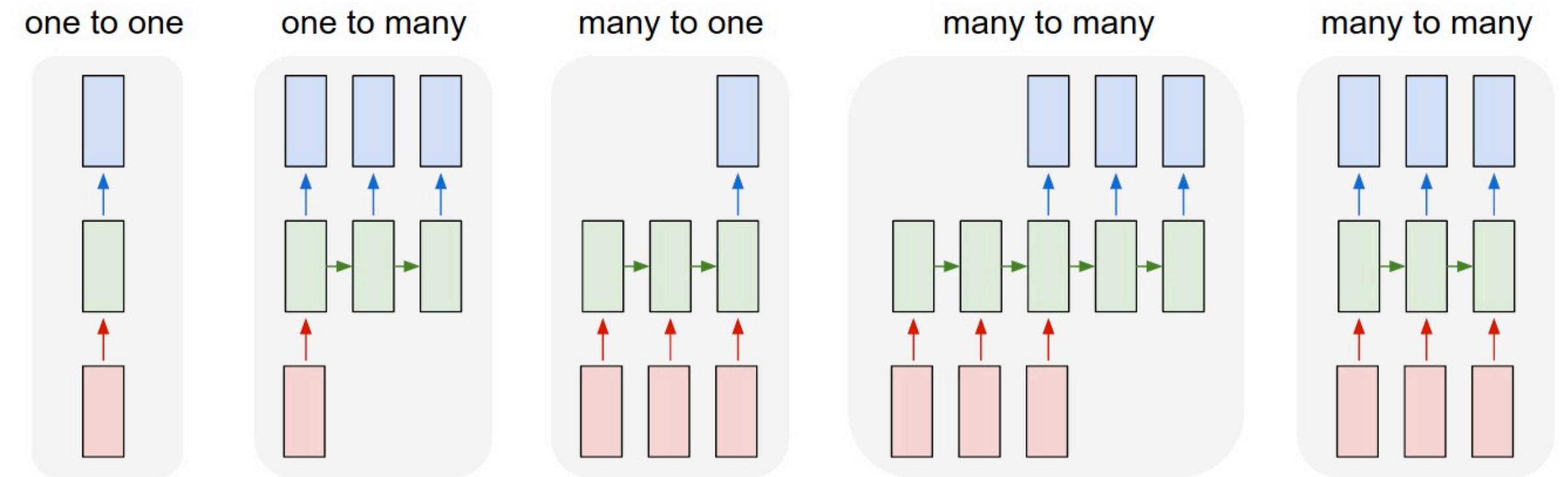
— Allows modeling **arbitrary length sequences**!



one to one     one to many     many to one     many to many     many to many

---

### Vanilla **RNN**

$$y_t = W_{hy}h_t + b_y$$

$$h_t = f_W(h_{t-1}, x_t)$$

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t + b_h)$$

**y**

**RNN**

**x**

### **Vanishing** or **Exploding** Gradients

### Long-Short Term Memory (**LSTM**)

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

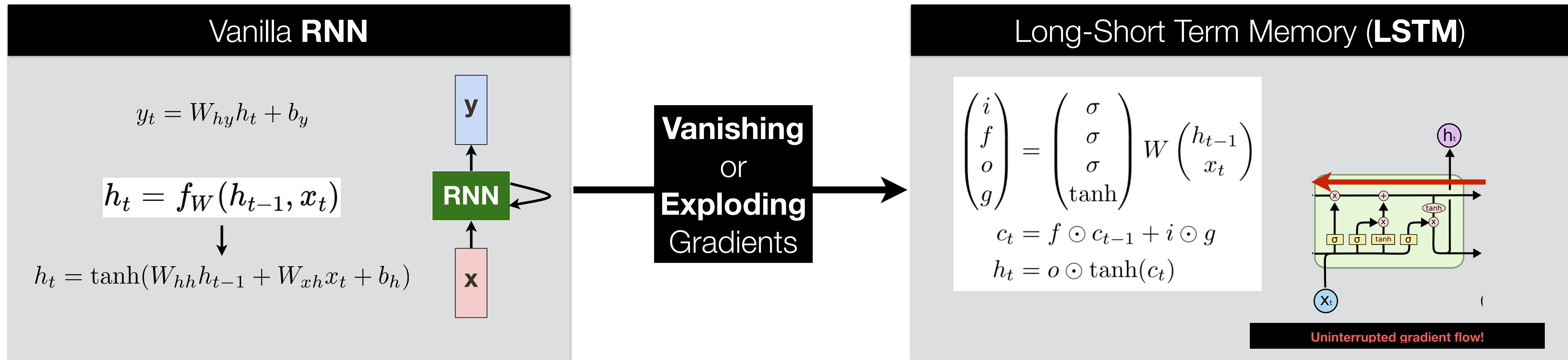**Uninterrupted gradient flow!**

# RNNs: Review

## Key Enablers:

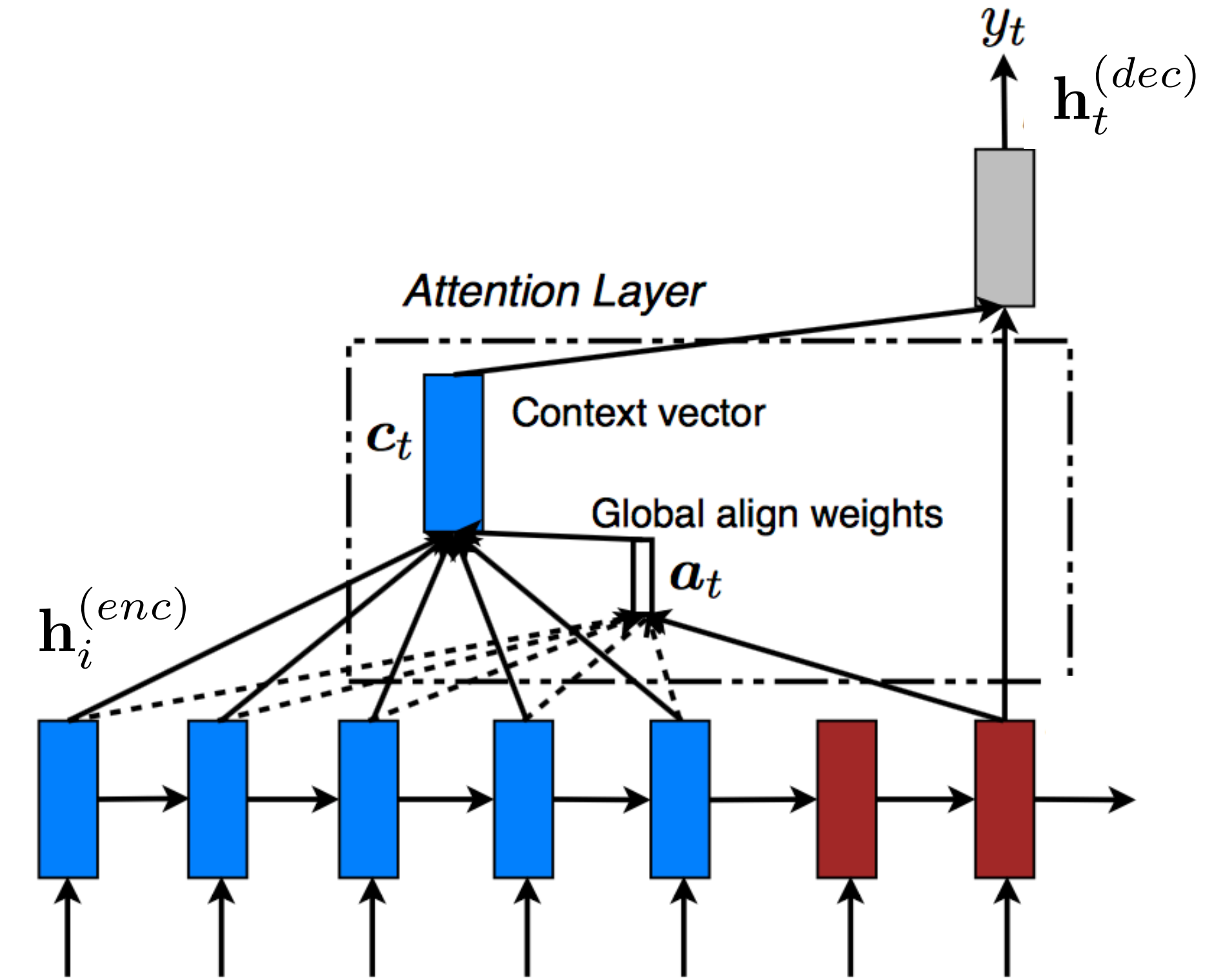— Parameter sharing in computational graphs

— "Unrolling" in computational graphs

— Allows modeling **arbitrary length sequences**!



**Loss functions:** often cross-entropy (for classification); could be max-margin (like in SVM) or Squared Loss (regression)
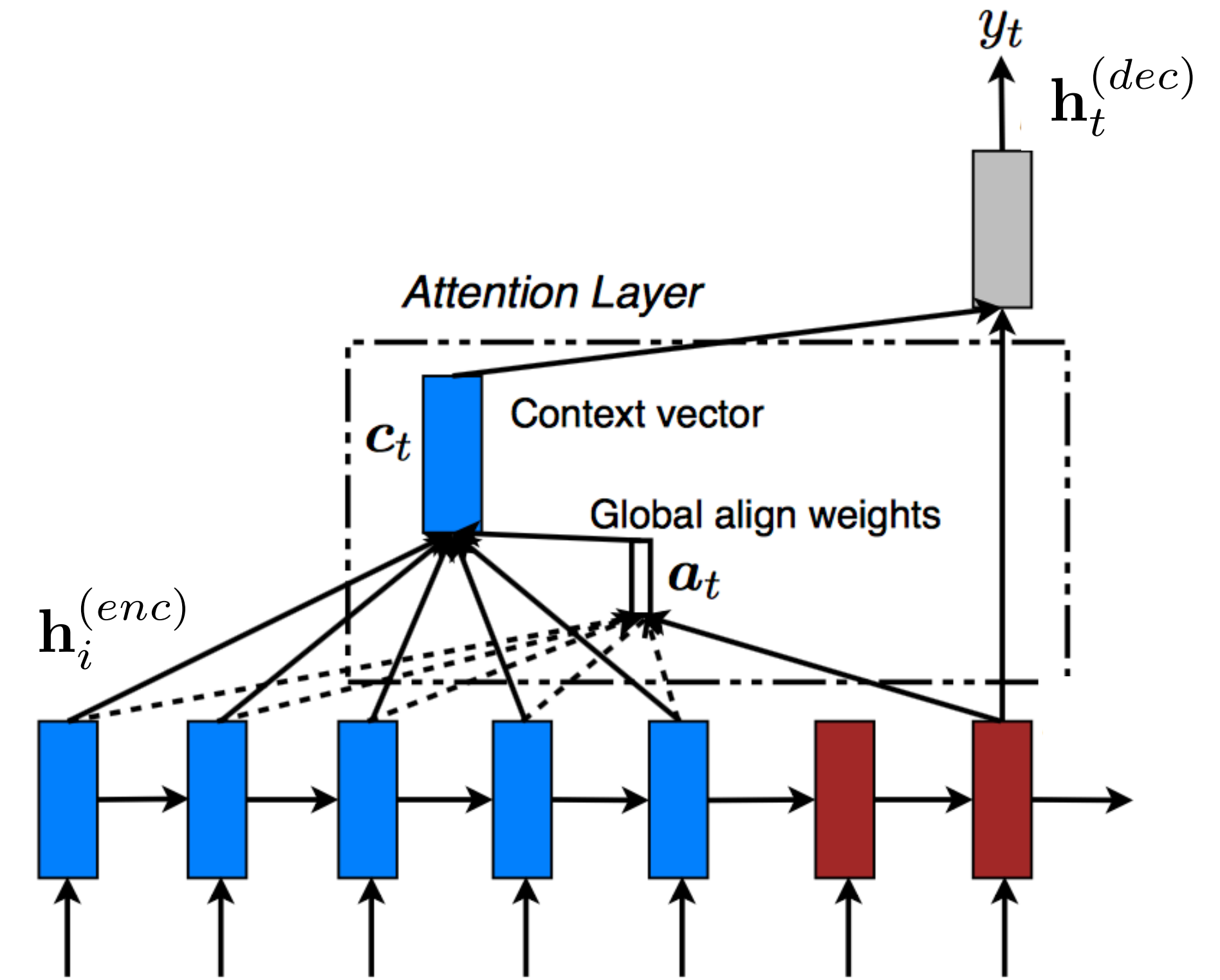
### Vanilla **RNN**

$$y_t = W_{hy} h_t + b_y$$

$$h_t = f_W(h_{t-1}, x_t)$$

$$h_t = \tanh(W_{hh} h_{t-1} + W_{xh} x_t + b_h)$$

**y**

**RNN**

**x**

**Vanishing** or **Exploding** Gradients

### Long-Short Term Memory (**LSTM**)

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

**Uninterrupted gradient flow!**

# Soft **Attention** in details

# Soft **Attention** in details

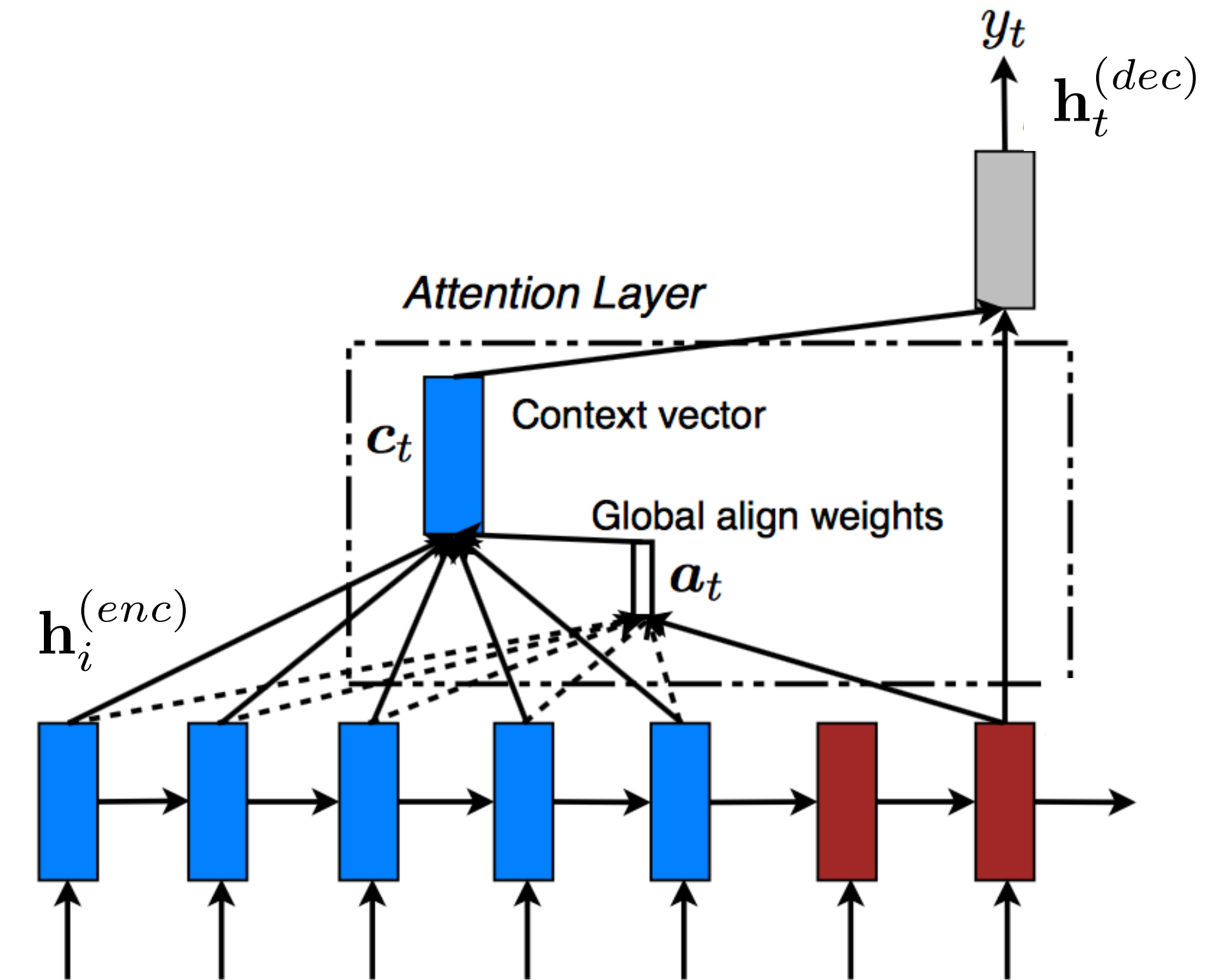$$\beta_{i,t} = score(\mathbf{h}_i^{(enc)}, \mathbf{h}_t^{(dec)})$$
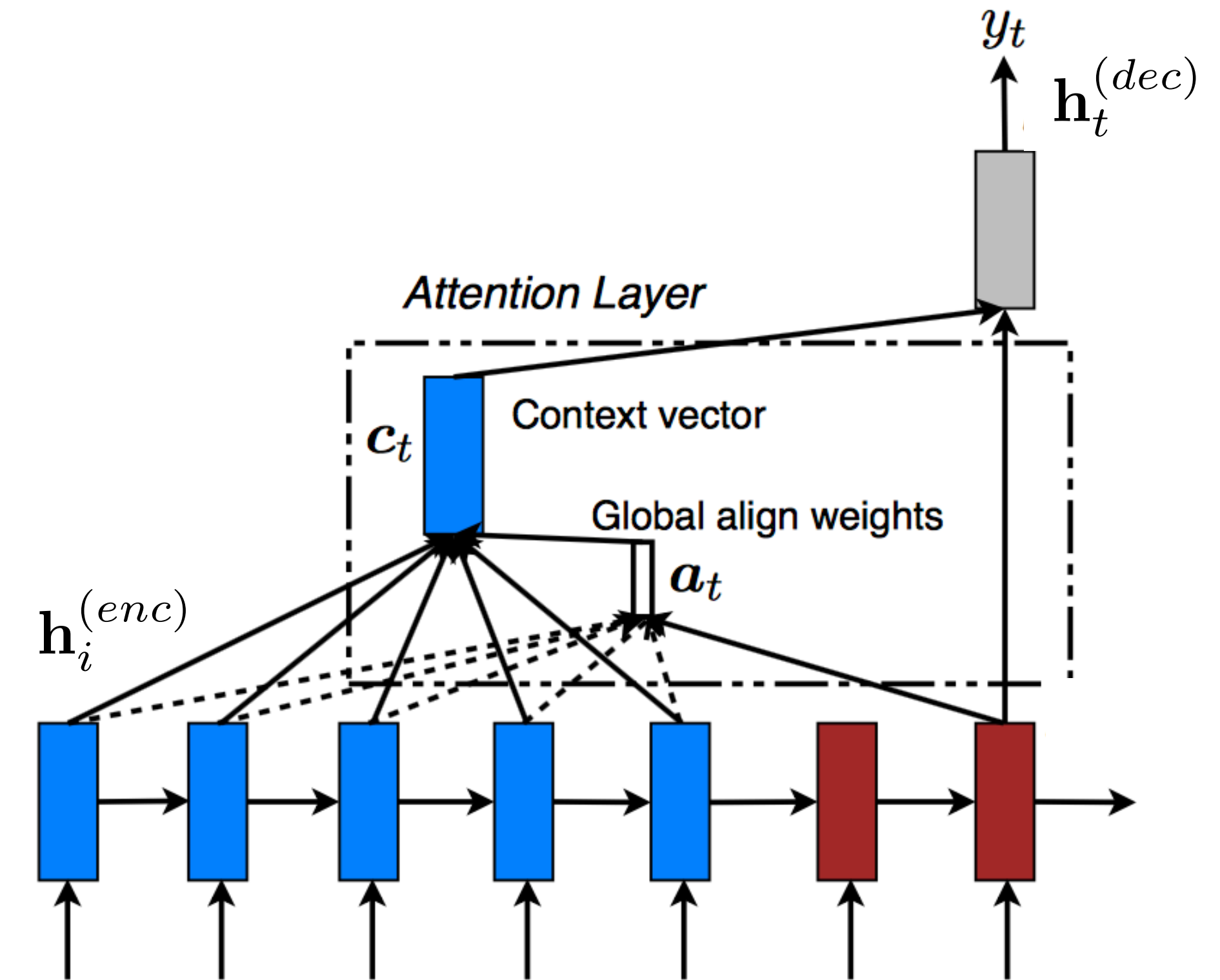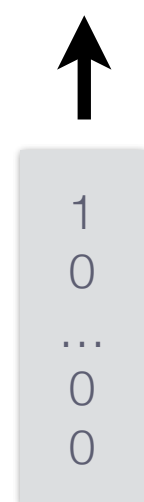
Relevance of encoding at
token i for decoding token t

# Soft **Attention** in details

$$\beta_{i,t} = score(\mathbf{h}_i^{(enc)}, \mathbf{h}_t^{(dec)})$$

Relevance of encoding at
token i for decoding token t



$$\alpha_{i,t} = \mathrm{Softmax}(\beta_{i,t})$$

Normalize the weights
to sum to 1

# Soft **Attention** in details

$$\beta_{i,t} = score(\mathbf{h}_i^{(enc)}, \mathbf{h}_t^{(dec)})$$

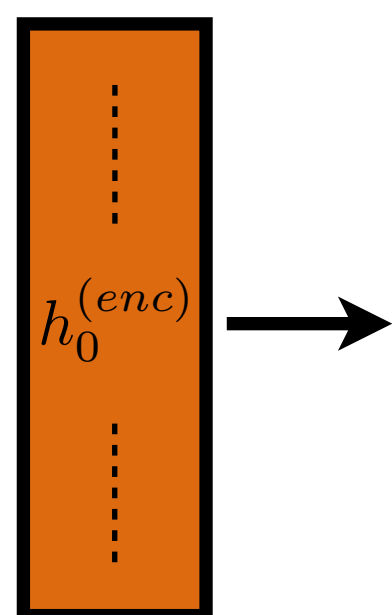Relevance of encoding at
token i for decoding token t



$$\alpha_{i,t} = \text{Softmax}(\beta_{i,t})$$

Normalize the weights
to sum to 1

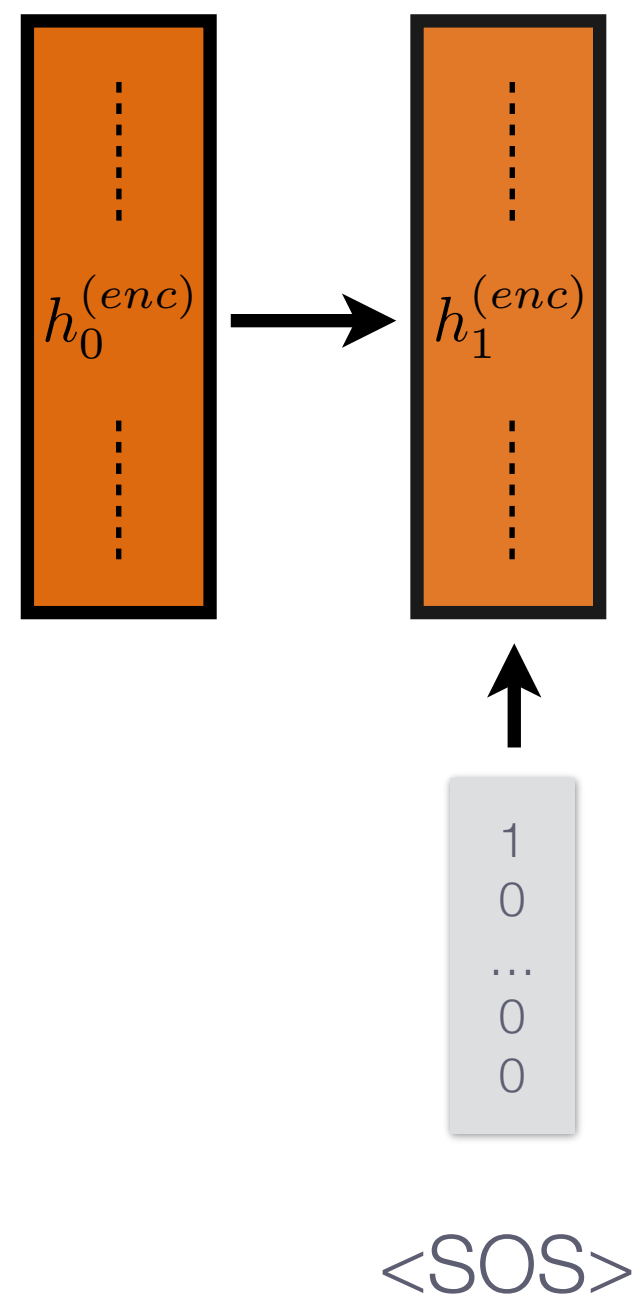$$\mathbf{c}_t = \sum_i \alpha_{i,t} \mathbf{h}_i^{(enc)}$$

Form a context vector that would simply be added to the standard decoder input
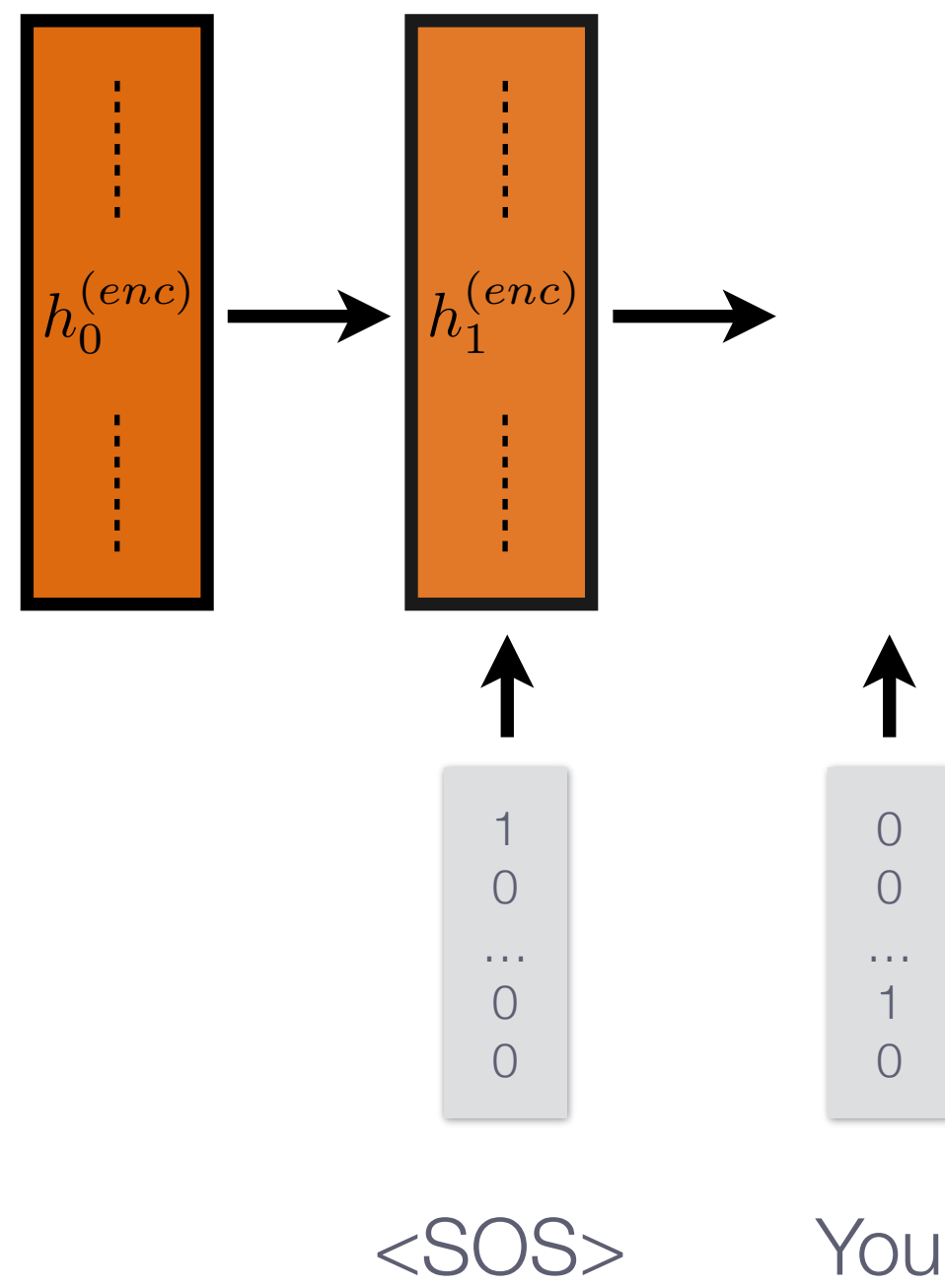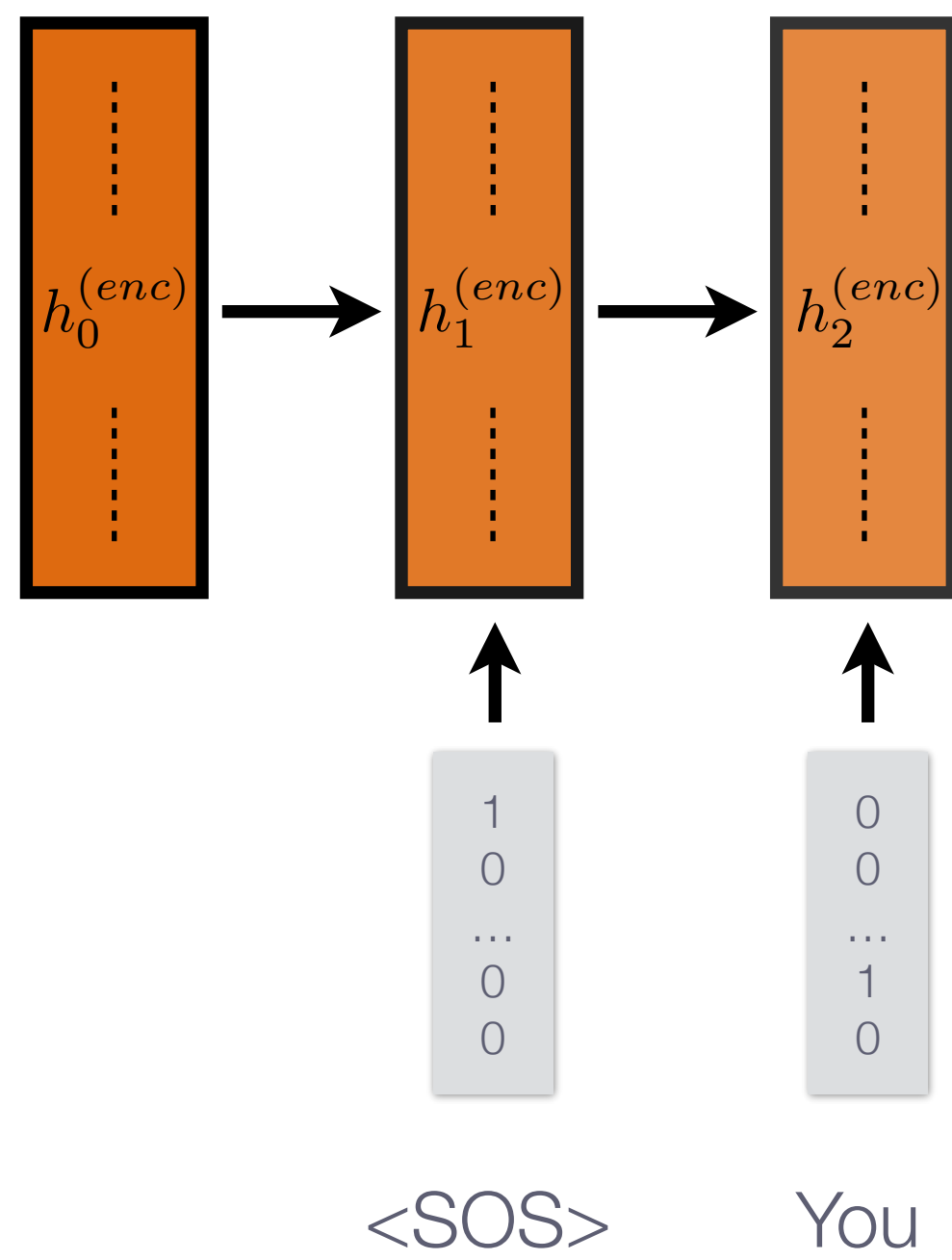
# Encoder (English)

$h_0^{(enc)}$ →

1
0
...
0

<SOS>

Encoder (English)

$h_0^{(enc)}$ → $h_1^{(enc)}$

1
0
...
0
0

<SOS>

Encoder (English)

$h_0^{(enc)}$ → $h_1^{(enc)}$ →

1
0
...
0
0

0
0
...
1
0

<SOS>          You

Encoder (English)

$h_0^{(enc)} \rightarrow h_1^{(enc)} \rightarrow h_2^{(enc)}$

1
0
...
0
0

0
0
...
1
0

<SOS>     You

Encoder (English)

**Summary Vector**

$h_0^{(enc)}$ → $h_1^{(enc)}$ → $h_2^{(enc)}$ → $h_3^{(enc)}$ → $h_4^{(enc)}$ → $h_5^{(enc)}$ → $h_6^{(enc)}$

| | 1 0 ... 0 0 | 0 0 ... 1 0 | ... | ... | 0 0 ... 0 0 | 0 1 ... 0 0 |

<SOS>    You    are    my    best    friend

Encoder (English)

**Summary Vector**

$h_0^{(enc)}$ $h_1^{(enc)}$ $h_2^{(enc)}$ $h_3^{(enc)}$ $h_4^{(enc)}$ $h_5^{(enc)}$ $h_6^{(enc)}$

1
0
...
0
0

0
0
...
1
0

...

...

0
0
...
0
0

0
1
...
0
0

<SOS>  You  are  my  best  friend

1
0
...
0

<SOS>

Decoder (Spanish)

Encoder (English)

**Summary Vector**

$h_0^{(enc)}$  $h_1^{(enc)}$  $h_2^{(enc)}$  $h_3^{(enc)}$  $h_4^{(enc)}$  $h_5^{(enc)}$  $h_6^{(enc)}$

$h_1^{(dec)}$

1
0
...
0
0

0
0
...
1
0

...

...

0
0
...
0
0

0
1
...
0
0

1
0
...
0

<SOS>  You  are  my  best  friend

<SOS>

Decoder (Spanish)

Encoder (English)

Summary Vector

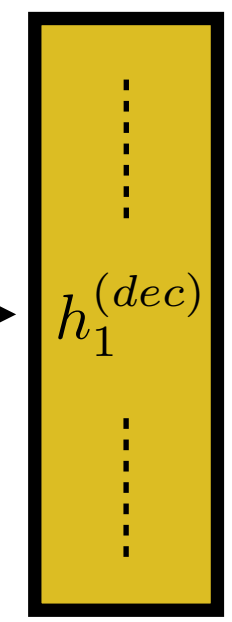Decoder (Spanish)

$h_0^{(enc)}$ $h_1^{(enc)}$ $h_2^{(enc)}$ $h_3^{(enc)}$ $h_4^{(enc)}$ $h_5^{(enc)}$ $h_6^{(enc)}$

$h_1^{(dec)}$

<SOS>    You    are    my    best    friend

<SOS>

Encoder (English)

**Summary Vector**

Eres

Softmax

$h_1^{(dec)}$

$h_0^{(enc)}$ $h_1^{(enc)}$ $h_2^{(enc)}$ $h_3^{(enc)}$ $h_4^{(enc)}$ $h_5^{(enc)}$ $h_6^{(enc)}$

<SOS> You are my best friend

<SOS>

Decoder (Spanish)

Encoder (English)

Summary Vector

Decoder (Spanish)

Eres

Softmax

$h_1^{(dec)}$

$h_0^{(enc)}$ $h_1^{(enc)}$ $h_2^{(enc)}$ $h_3^{(enc)}$ $h_4^{(enc)}$ $h_5^{(enc)}$ $h_6^{(enc)}$

<SOS> You are my best friend

<SOS> Eres

Encoder (English)

**Summary Vector**

Decoder (Spanish)

Eres

Softmax

$h_0^{(enc)}$ $h_1^{(enc)}$ $h_2^{(enc)}$ $h_3^{(enc)}$ $h_4^{(enc)}$ $h_5^{(enc)}$ $h_6^{(enc)}$

$h_1^{(dec)}$ $h_2^{(dec)}$

<SOS> You are my best friend

<SOS> Eres

Encoder (English)

$h_0^{(enc)}$ $h_1^{(enc)}$ $h_2^{(enc)}$ $h_3^{(enc)}$ $h_4^{(enc)}$ $h_5^{(enc)}$ $h_6^{(enc)}$

<SOS>  You  are  my  best  friend

Decoder (Spanish)

Eres  mi

Softmax  Softmax

$h_1^{(dec)}$ $h_2^{(dec)}$

<SOS>  Eres

Encoder (English)

$h_0^{(enc)}$ $h_1^{(enc)}$ $h_2^{(enc)}$ $h_3^{(enc)}$ $h_4^{(enc)}$ $h_5^{(enc)}$ $h_6^{(enc)}$

<SOS> You are my best friend

Decoder (Spanish)

Eres mi major amiga <EOS>

Softmax Softmax Softmax Softmax Softmax

$h_1^{(dec)}$ $h_2^{(dec)}$ $h_3^{(dec)}$ $h_4^{(dec)}$ $h_5^{(dec)}$

<SOS> Eres mi major amiga

Encoder (English)

$h_0^{(enc)}$  $h_1^{(enc)}$  $h_2^{(enc)}$  $h_3^{(enc)}$  $h_4^{(enc)}$  $h_5^{(enc)}$  $h_6^{(enc)}$

| 1 | 0 | | | 0 | 0 |
| 0 | 0 | ... | ... | 0 | 1 |
| ... | ... | | | ... | ... |
| 0 | 1 | | | 0 | 0 |
| 0 | 0 | | | 0 | 0 |

<SOS>   You   are   my   best   friend

Eres   mi   major   amiga   <EOS>

Softmax   Softmax   Softmax   Softmax   Softmax

| 0.1 | 0.3 | | | |
| -1.2 | 3.4 | ... | ... | ... |
| ... | ... | | | |
| 0.4 | -7.8 | | | |
| 0.6 | 9.1 | | | |

$h_1^{(dec)}$  $h_2^{(dec)}$  $h_3^{(dec)}$  $h_4^{(dec)}$  $h_5^{(dec)}$

| 1 |
| 0 |
| ... |
| 0 |

... ... ... ...

<SOS>   Eres   mi   major   amiga

**Encoder (English)**

$h_0^{(enc)}$ → $h_1^{(enc)}$ → $h_2^{(enc)}$ → $h_3^{(enc)}$ → $h_4^{(enc)}$ → $h_5^{(enc)}$ → $h_6^{(enc)}$

| 1 | 0 | | | 0 | 0 |
|---|---|---|---|---|---|
| 0 | 0 | ... | ... | 0 | 1 |
| ... | 1 | | | ... | ... |
| 0 | 0 | | | 0 | 0 |
| 0 | | | | 0 | 0 |

<SOS>    You    are    my    best    friend

Eres    mi    major    amiga    <EOS>

Softmax    Softmax    Softmax    Softmax    Softmax

| 0.1 | 0.3 | | | |
|-----|-----|---|---|---|
| -1.2 | 3.4 | ... | ... | ... |
| ... | ... | | | |
| 0.4 | -7.8 | | | |
| 0.6 | 9.1 | | | |

$h_1^{(dec)}$ → $h_2^{(dec)}$ → $h_3^{(dec)}$ → $h_4^{(dec)}$ → $h_5^{(dec)}$

$h_1^{(enc)}$    $h_2^{(enc)}$    $h_3^{(enc)}$    $h_4^{(enc)}$    $h_5^{(enc)}$

| 1 | | | | |
|---|---|---|---|---|
| 0 | ... | ... | ... | ... |
| ... | | | | |
| 0 | | | | |

<SOS>    Eres    mi    major    amiga

$$= \text{Score}(\quad, \quad)$$

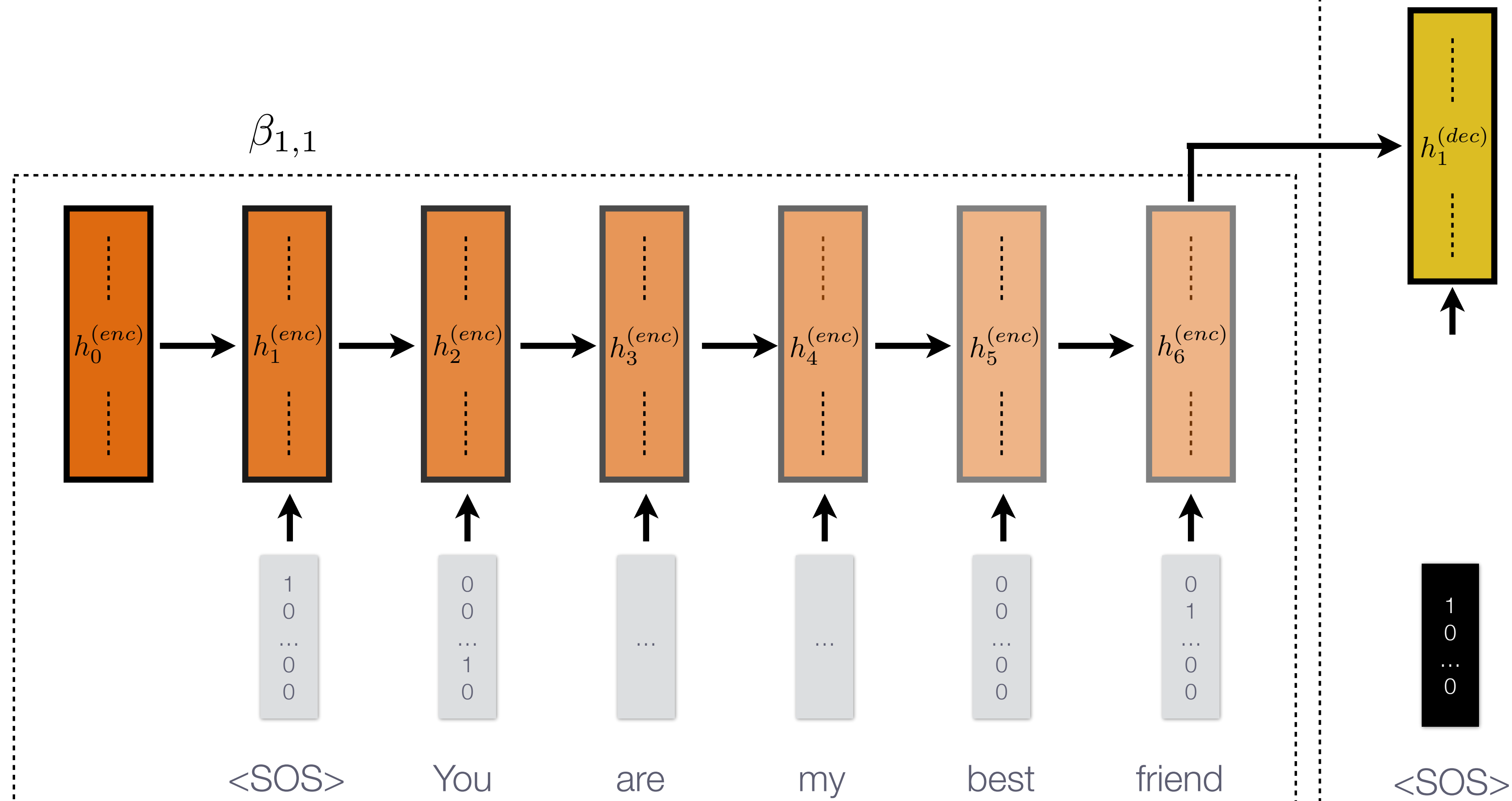$$= \text{Score}(\ h_1^{(enc)}\ ,\ h_1^{(dec)}\ )$$

$h_1^{(dec)}$

$h_0^{(enc)}$ → $h_1^{(enc)}$ → $h_2^{(enc)}$ → $h_3^{(enc)}$ → $h_4^{(enc)}$ → $h_5^{(enc)}$ → $h_6^{(enc)}$

| 1 | 0 | | | 0 | 0 |
| 0 | 0 | | | 0 | 1 |
| ... | ... | ... | ... | ... | ... |
| 0 | 1 | | | 0 | 0 |
| 0 | 0 | | | 0 | 0 |

<SOS>　　You　　are　　my　　best　　friend

1
0
...
0

<SOS>

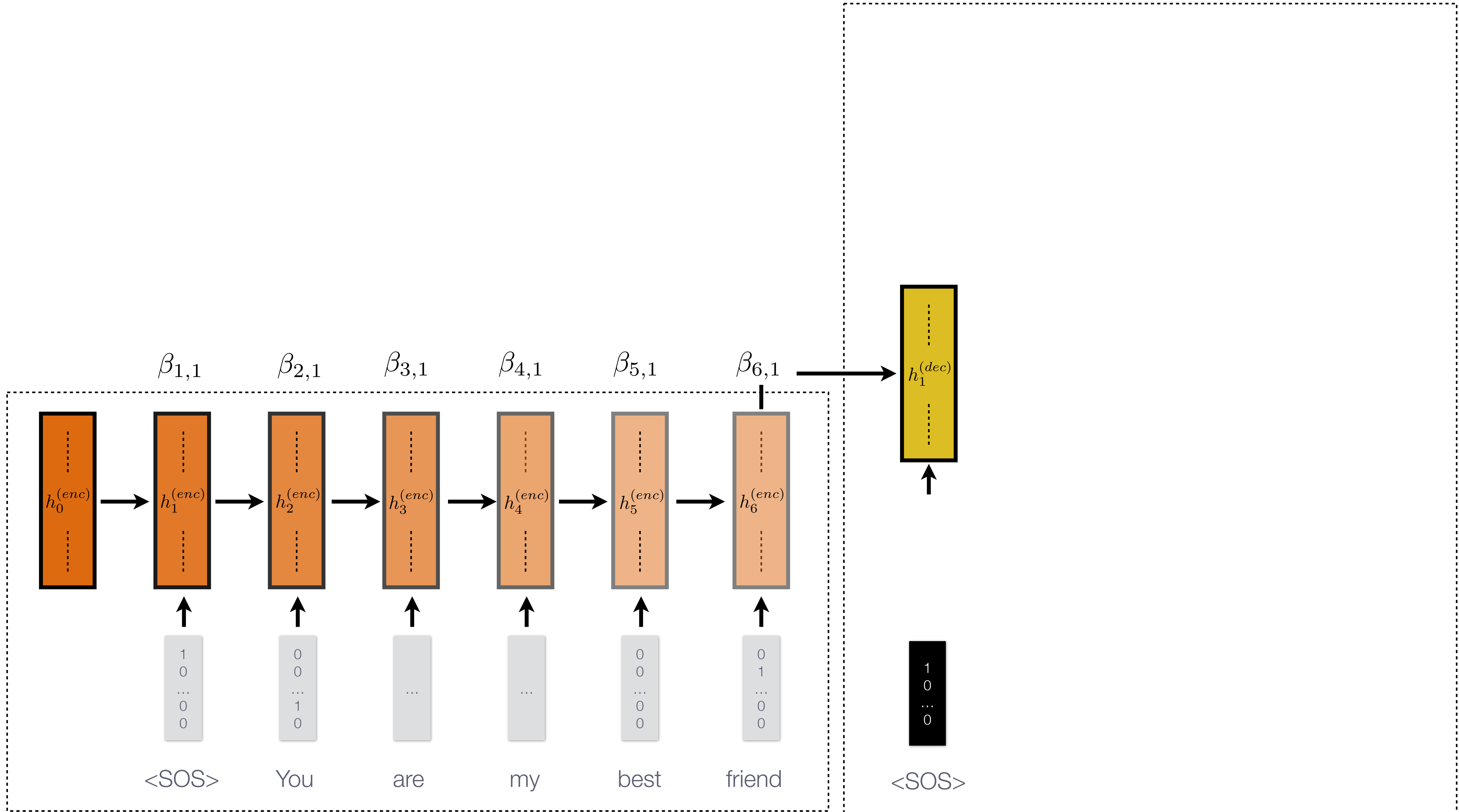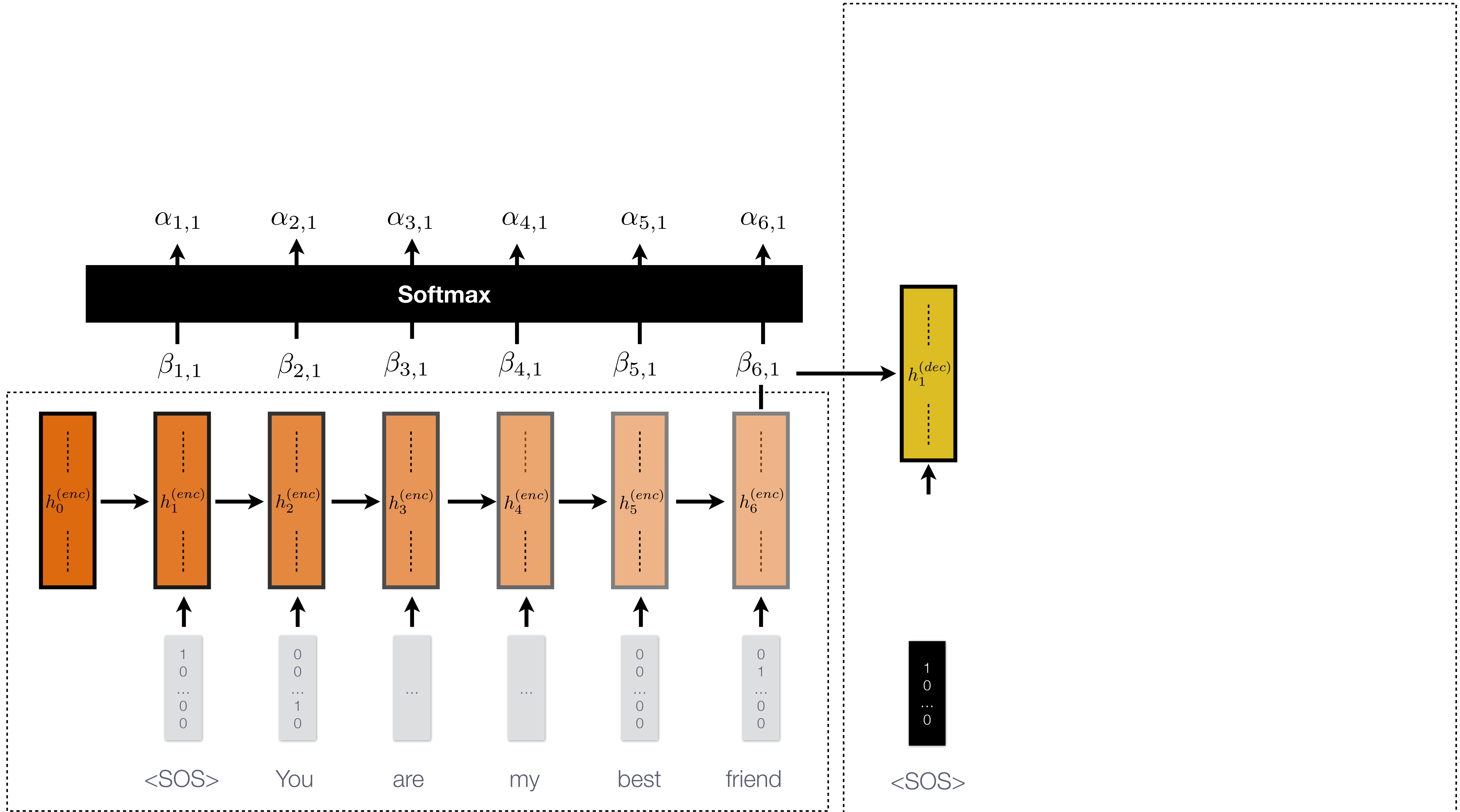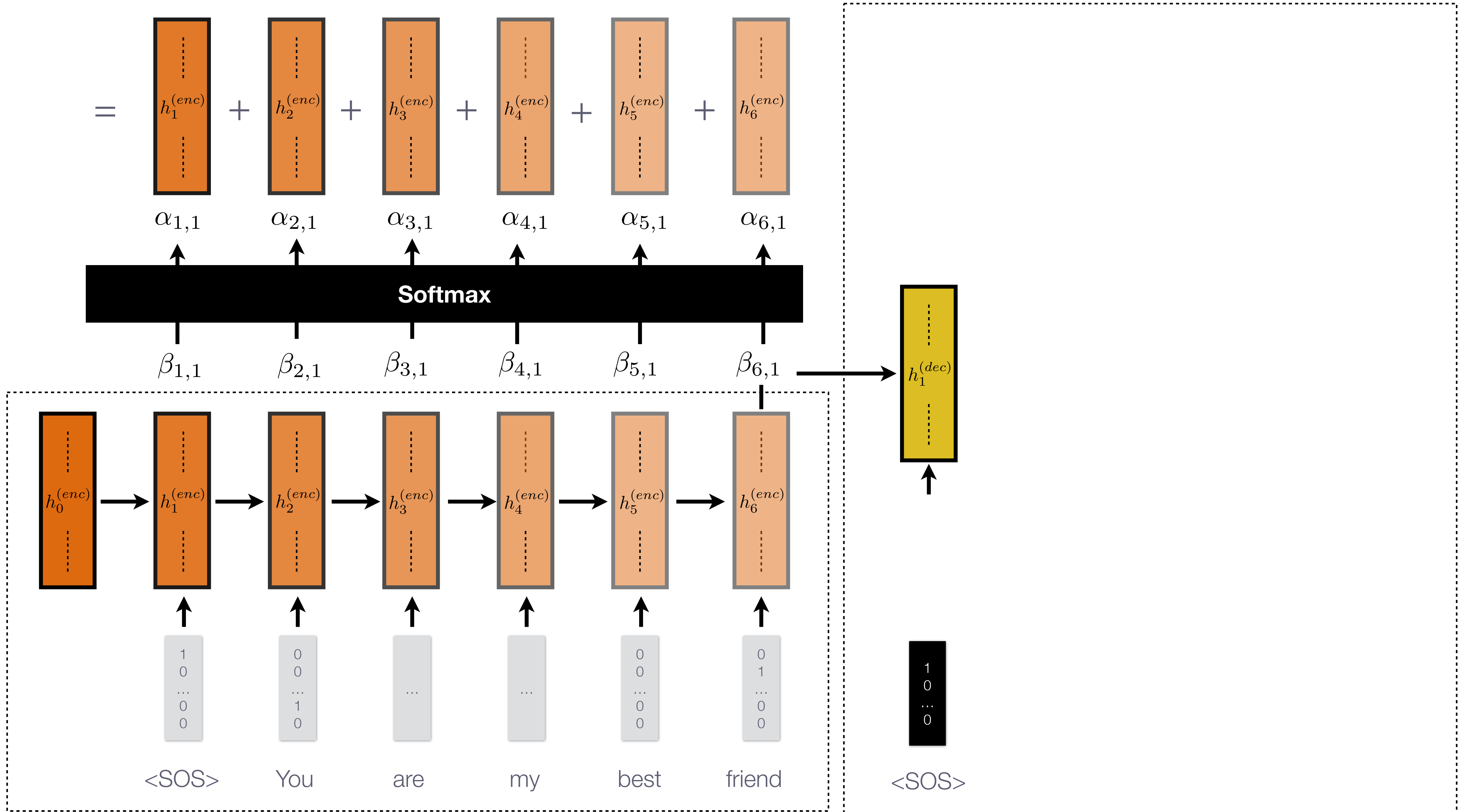$$\beta_{1,1} = \text{Score}(h_1^{(enc)}, h_1^{(dec)})$$
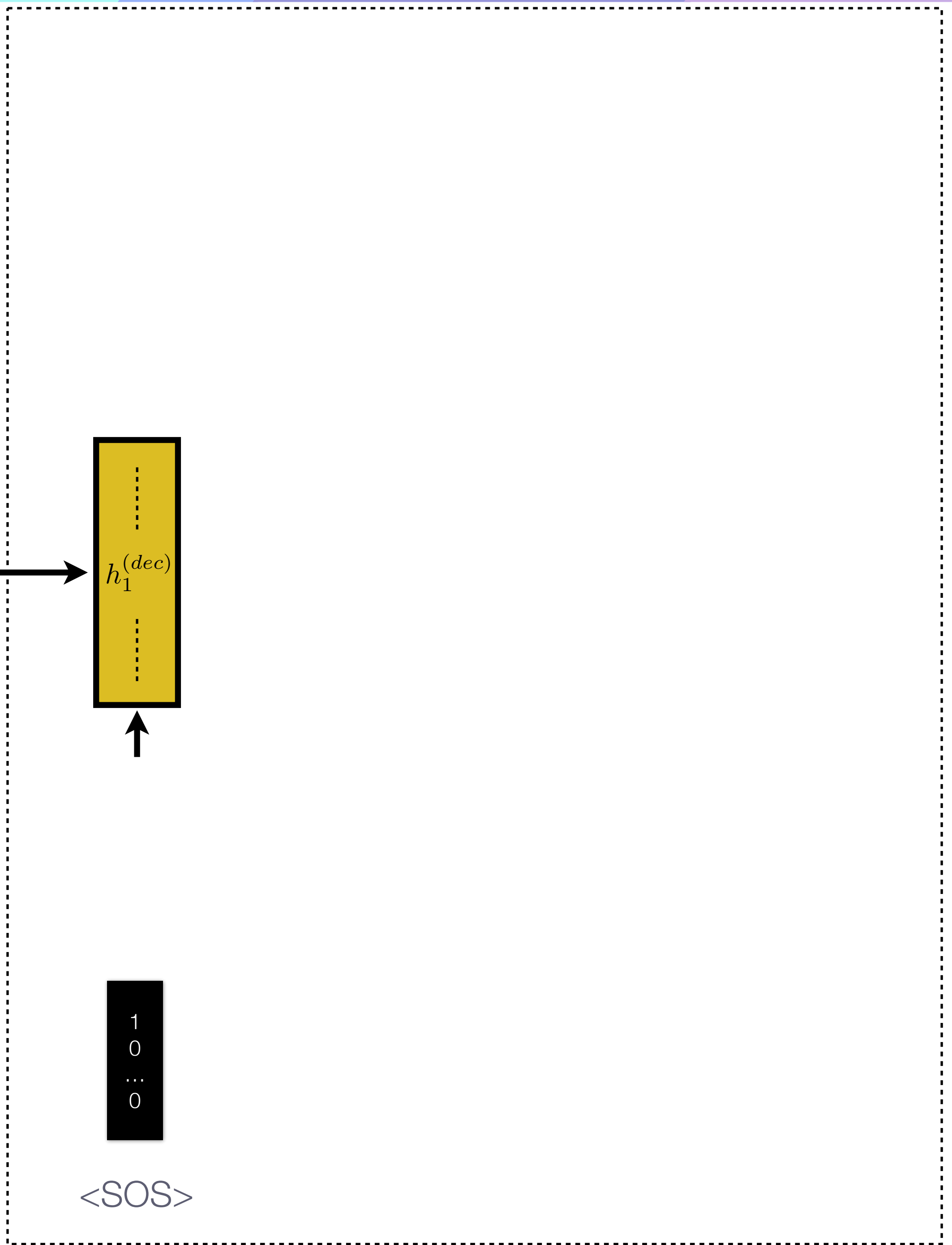
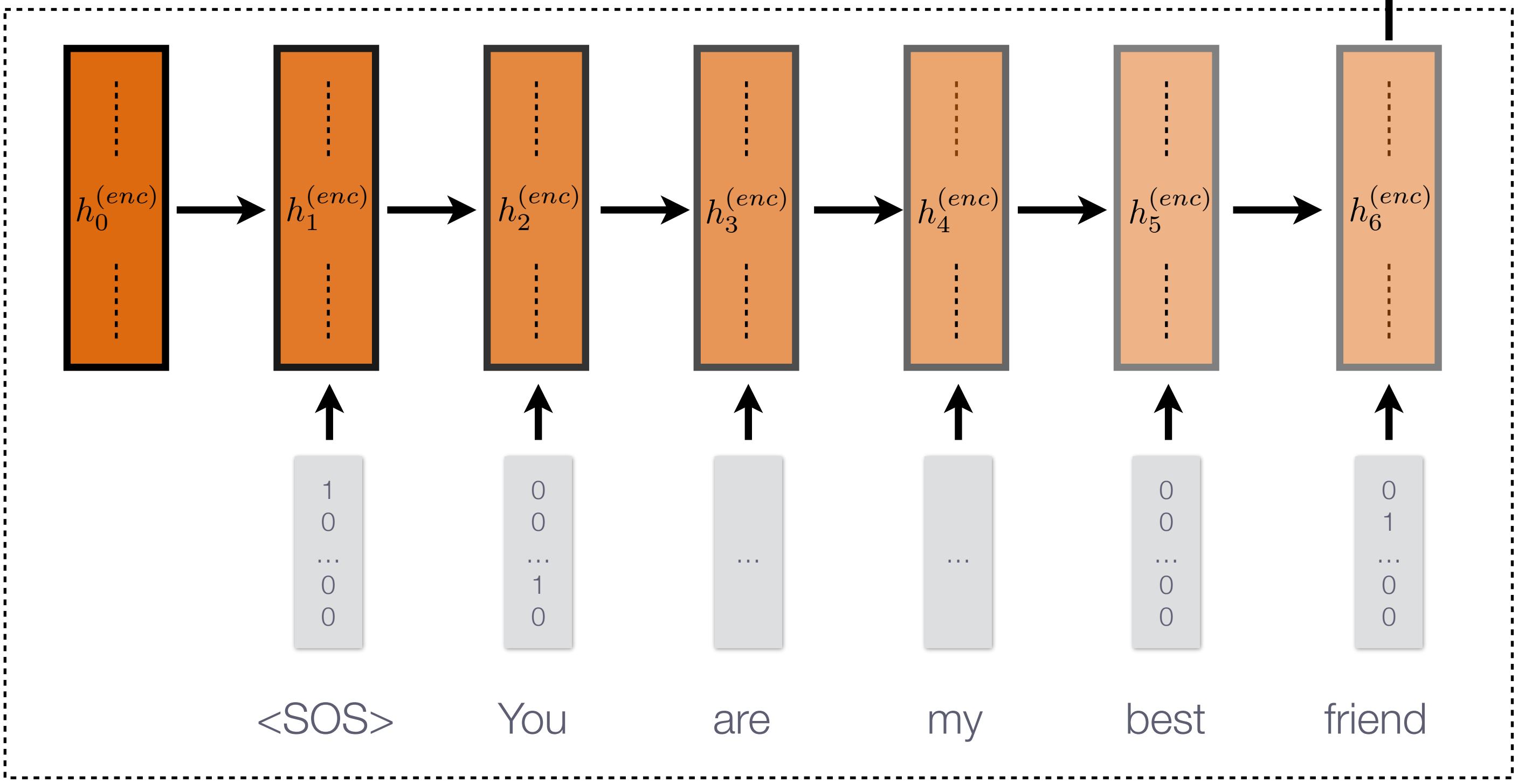$$= \text{Score}(\ h_1^{(enc)}\ ,\ h_1^{(dec)}\ )$$

$\beta_{1,1}$

$h_0^{(enc)}$ $h_1^{(enc)}$ $h_2^{(enc)}$ $h_3^{(enc)}$ $h_4^{(enc)}$ $h_5^{(enc)}$ $h_6^{(enc)}$

$h_1^{(dec)}$

| 1 | 0 | | | 0 | 0 |
| 0 | 0 | | | 0 | 1 |
| ... | ... | ... | ... | ... | ... |
| 0 | 1 | | | 0 | 0 |
| 0 | | | | 0 | 0 |

<SOS>   You   are   my   best   friend

| 1 |
| 0 |
| ... |
| 0 |

<SOS>

$$= \text{Score}(\ h_2^{(enc)}\ ,\ h_1^{(dec)}\ )$$

$\beta_{1,1}$
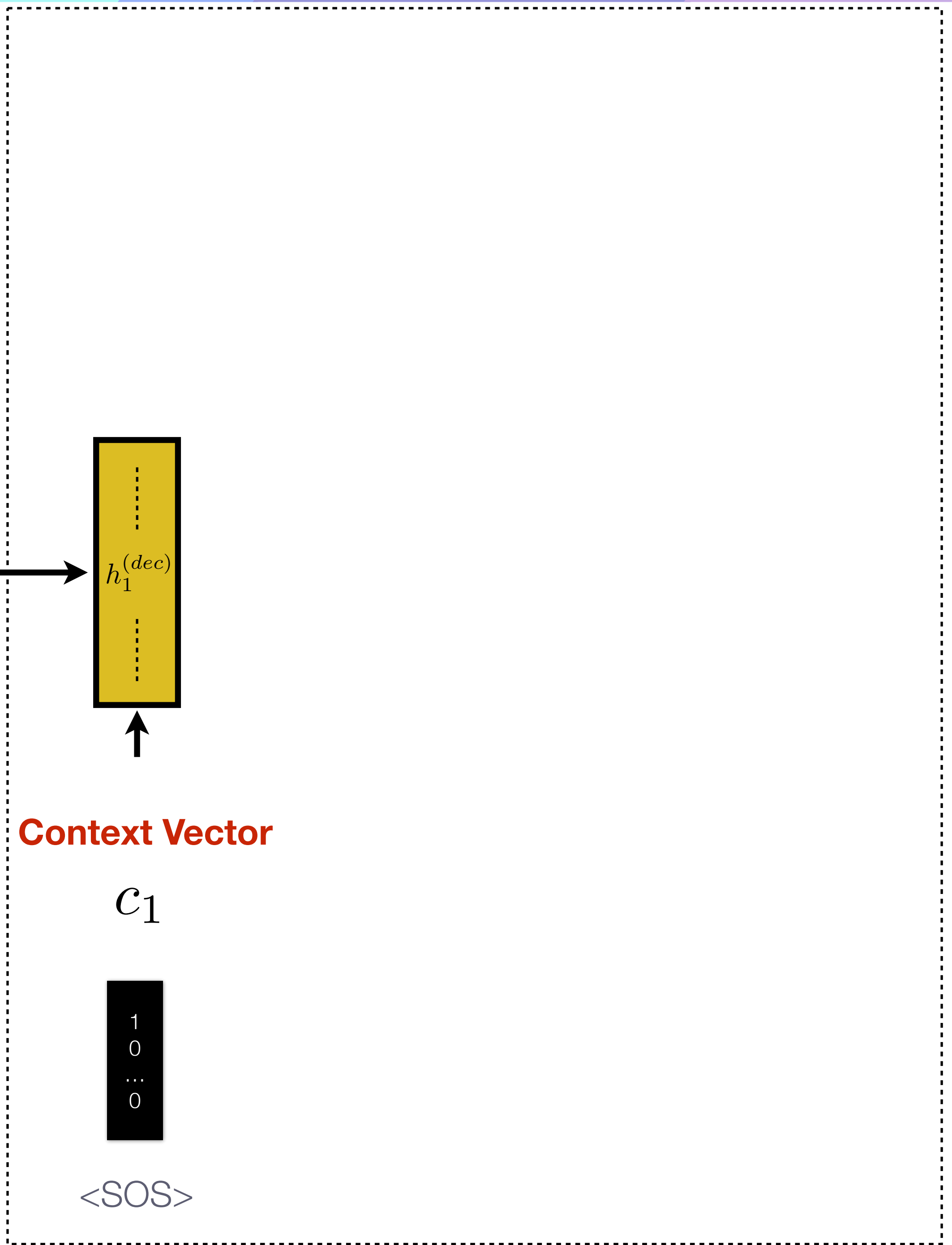
$h_0^{(enc)} \rightarrow h_1^{(enc)} \rightarrow h_2^{(enc)} \rightarrow h_3^{(enc)} \rightarrow h_4^{(enc)} \rightarrow h_5^{(enc)} \rightarrow h_6^{(enc)}$

$h_1^{(dec)}$

| | |
|---|---|
| 1 | 0 |
| 0 | 0 |
| ... | ... |
| 0 | 1 |
| 0 | 0 |

<SOS>    You    are    my    best    friend

<SOS>

$$\beta_{1,1} \quad \beta_{2,1} \quad \beta_{3,1} \quad \beta_{4,1} \quad \beta_{5,1} \quad \beta_{6,1}$$

$h_0^{(enc)} \quad h_1^{(enc)} \quad h_2^{(enc)} \quad h_3^{(enc)} \quad h_4^{(enc)} \quad h_5^{(enc)} \quad h_6^{(enc)}$

$h_1^{(dec)}$

<SOS> You are my best friend

<SOS>

$$\alpha_{1,1} \quad \alpha_{2,1} \quad \alpha_{3,1} \quad \alpha_{4,1} \quad \alpha_{5,1} \quad \alpha_{6,1}$$

**Softmax**

$$\beta_{1,1} \quad \beta_{2,1} \quad \beta_{3,1} \quad \beta_{4,1} \quad \beta_{5,1} \quad \beta_{6,1}$$

$h_0^{(enc)}$  $h_1^{(enc)}$  $h_2^{(enc)}$  $h_3^{(enc)}$  $h_4^{(enc)}$  $h_5^{(enc)}$  $h_6^{(enc)}$
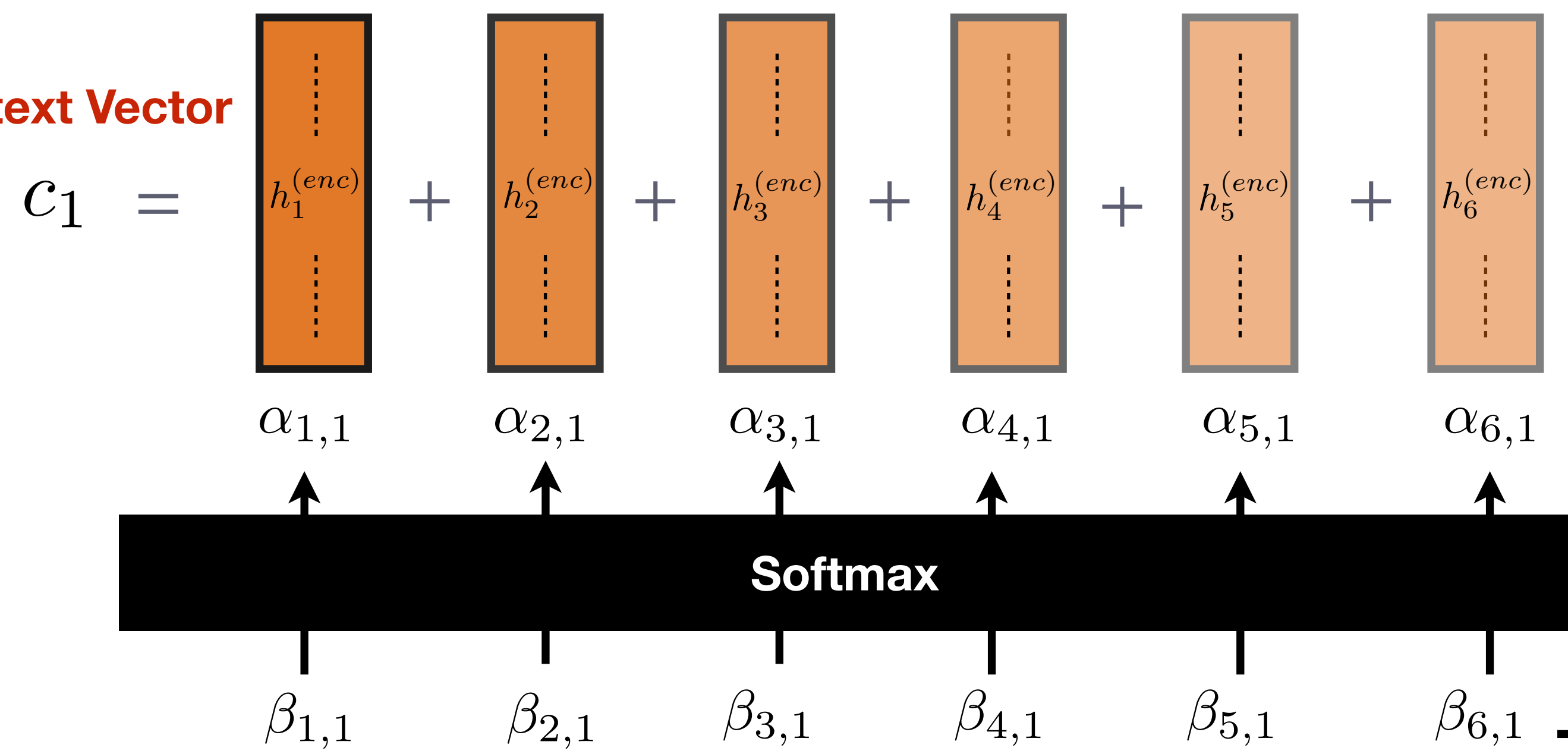
$h_1^{(dec)}$

<SOS>  You  are  my  best  friend
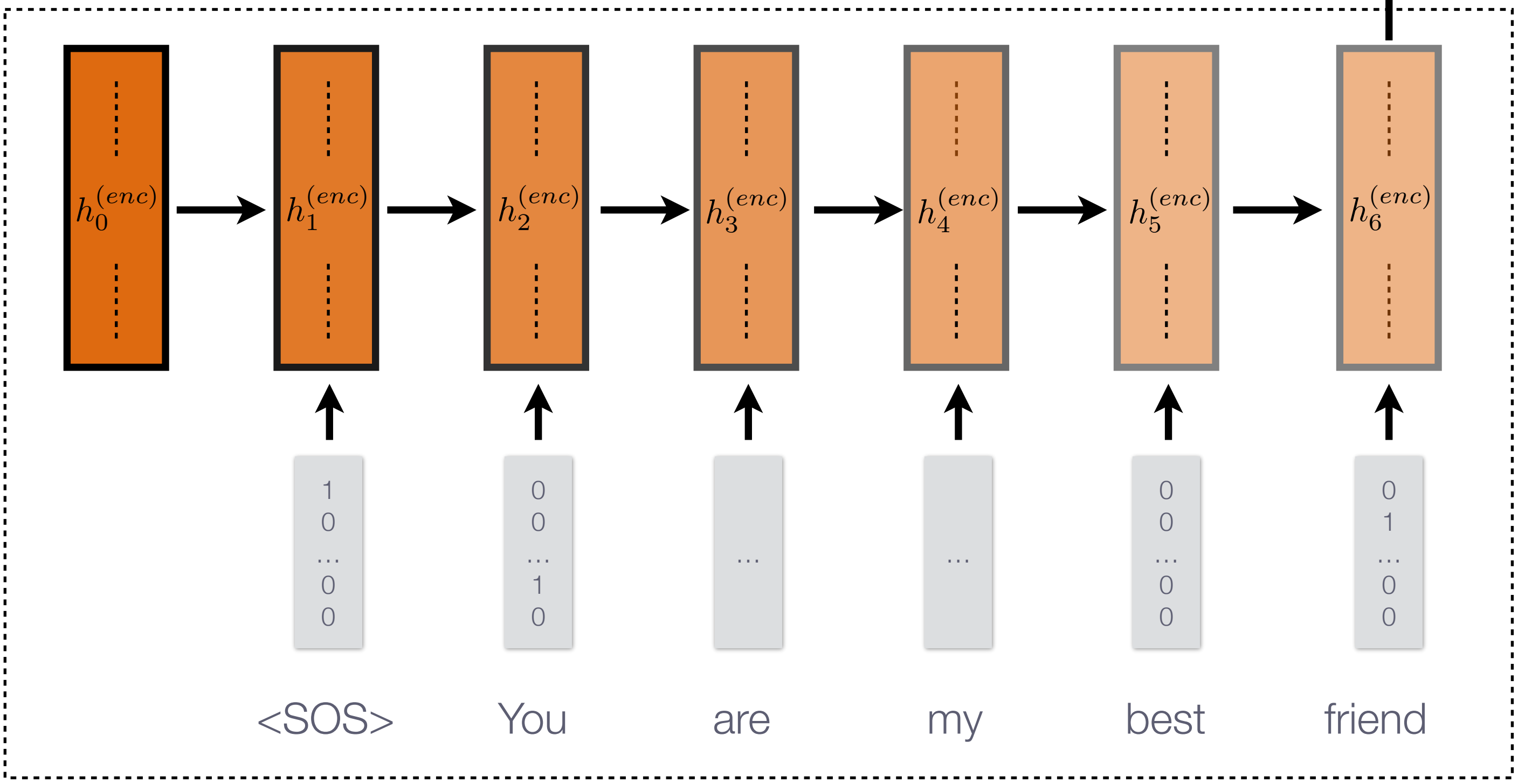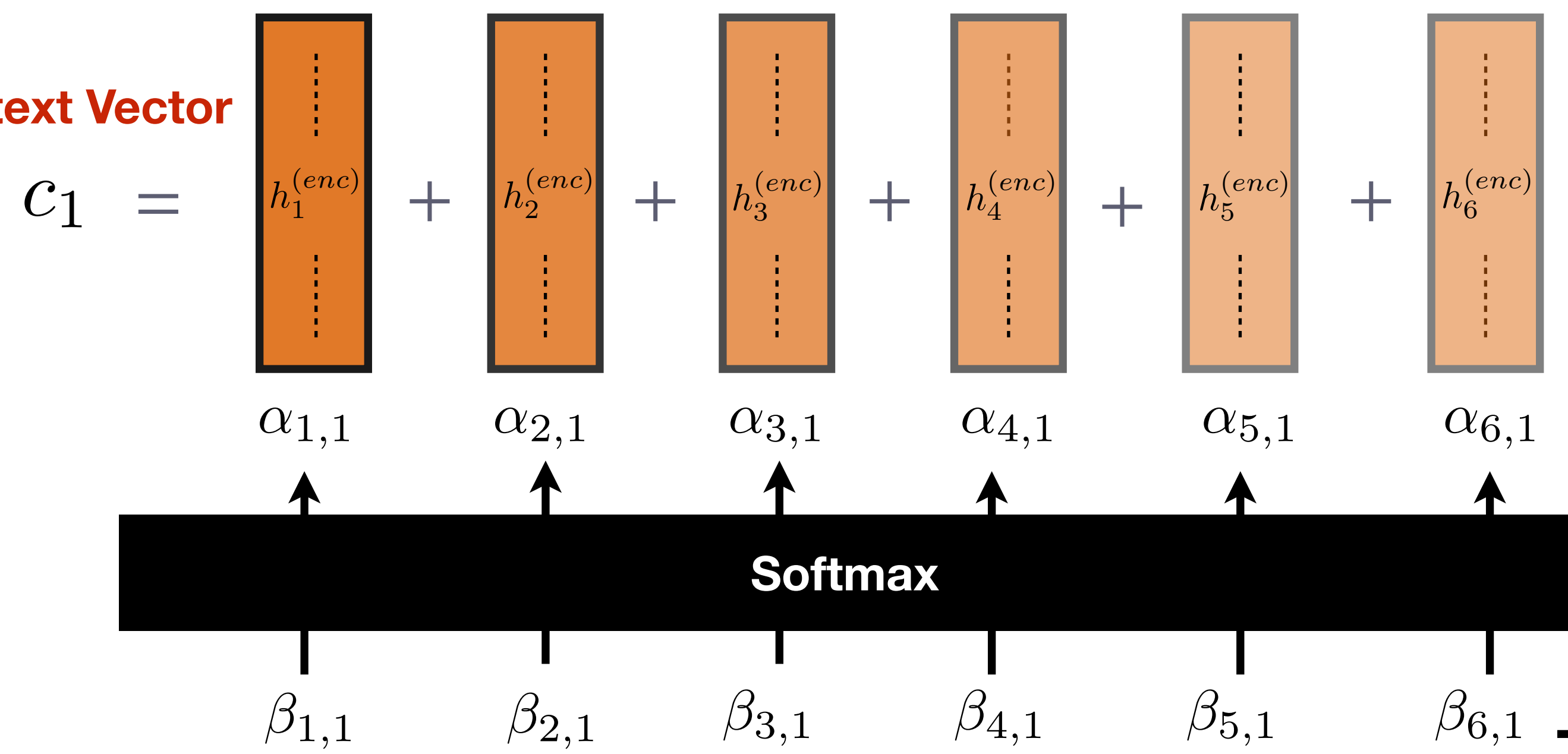
<SOS>

**Context Vector**

$$c_1 = h_1^{(enc)} + h_2^{(enc)} + h_3^{(enc)} + h_4^{(enc)} + h_5^{(enc)} + h_6^{(enc)}$$

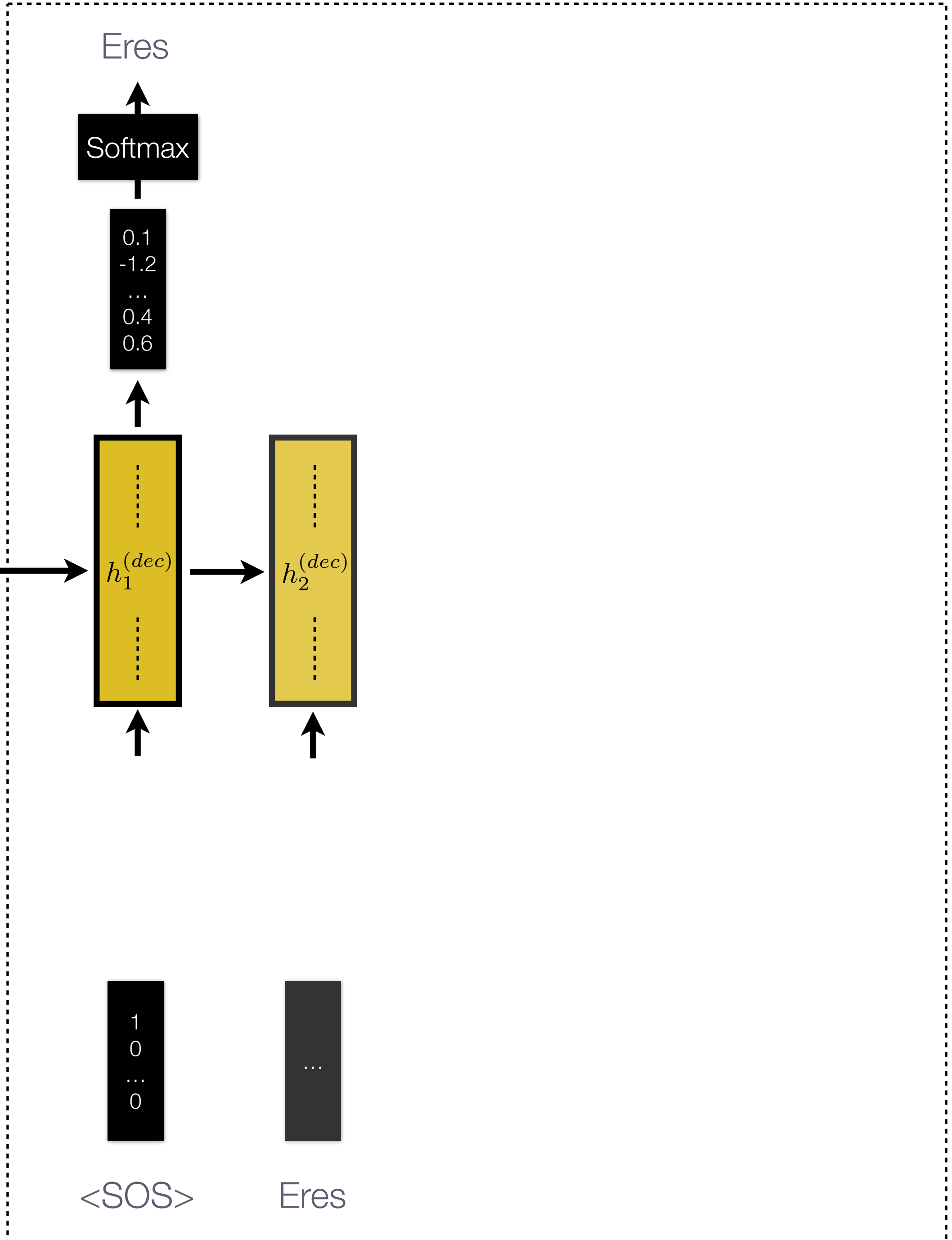$$\alpha_{1,1} \qquad \alpha_{2,1} \qquad \alpha_{3,1} \qquad \alpha_{4,1} \qquad \alpha_{5,1} \qquad \alpha_{6,1}$$

**Softmax**

$$\beta_{1,1} \qquad \beta_{2,1} \qquad \beta_{3,1} \qquad \beta_{4,1} \qquad \beta_{5,1} \qquad \beta_{6,1}$$

$h_0^{(enc)} \rightarrow h_1^{(enc)} \rightarrow h_2^{(enc)} \rightarrow h_3^{(enc)} \rightarrow h_4^{(enc)} \rightarrow h_5^{(enc)} \rightarrow h_6^{(enc)}$

$h_1^{(dec)}$

| 1 | 0 | | | 0 | 0 |
| 0 | 0 | | | 0 | 1 |
| ... | ... | ... | ... | ... | ... |
| 0 | 1 | | | 0 | 0 |
| 0 | 0 | | | 0 | 0 |

<SOS>    You    are    my    best    friend

1
0
...
0

<SOS>

**Context Vector**

$$c_1 \; = \; h_1^{(enc)} \; + \; h_2^{(enc)} \; + \; h_3^{(enc)} \; + \; h_4^{(enc)} \; + \; h_5^{(enc)} \; + \; h_6^{(enc)}$$

$\alpha_{1,1}$  $\alpha_{2,1}$  $\alpha_{3,1}$  $\alpha_{4,1}$  $\alpha_{5,1}$  $\alpha_{6,1}$

**Softmax**

$\beta_{1,1}$  $\beta_{2,1}$  $\beta_{3,1}$  $\beta_{4,1}$  $\beta_{5,1}$  $\beta_{6,1}$

$h_1^{(dec)}$

$h_0^{(enc)}$  $h_1^{(enc)}$  $h_2^{(enc)}$  $h_3^{(enc)}$  $h_4^{(enc)}$  $h_5^{(enc)}$  $h_6^{(enc)}$

| 1 | 0 | | | 0 | 0 |
| 0 | 0 | | | 0 | 1 |
| ... | ... | ... | ... | ... | ... |
| 0 | 1 | | | 0 | 0 |
| 0 | | | | 0 | 0 |

<SOS>   You   are   my   best   friend

**Context Vector**

$$c_1$$

| 1 |
| 0 |
| ... |
| 0 |

<SOS>

**Context Vector**

$$c_1 \;=\; h_1^{(enc)} \;+\; h_2^{(enc)} \;+\; h_3^{(enc)} \;+\; h_4^{(enc)} \;+\; h_5^{(enc)} \;+\; h_6^{(enc)}$$

$\alpha_{1,1}$  $\alpha_{2,1}$  $\alpha_{3,1}$  $\alpha_{4,1}$  $\alpha_{5,1}$  $\alpha_{6,1}$

**Softmax**

$\beta_{1,1}$  $\beta_{2,1}$  $\beta_{3,1}$  $\beta_{4,1}$  $\beta_{5,1}$  $\beta_{6,1}$

$h_0^{(enc)}$  $h_1^{(enc)}$  $h_2^{(enc)}$  $h_3^{(enc)}$  $h_4^{(enc)}$  $h_5^{(enc)}$  $h_6^{(enc)}$

1 0 ... 0 0

0 0 ... 1 0

...

...

0 0 ... 0 0

0 1 ... 0 0

<SOS>  You  are  my  best  friend

Eres

Softmax

0.1
-1.2
...
0.4
0.6

$h_1^{(dec)}$

**Context Vector**

$c_1$

1
0
...
0

<SOS>

# Context Vector

$$c_1 = h_1^{(enc)} + h_2^{(enc)} + h_3^{(enc)} + h_4^{(enc)} + h_5^{(enc)} + h_6^{(enc)}$$

$\alpha_{1,1}$    $\alpha_{2,1}$    $\alpha_{3,1}$    $\alpha_{4,1}$    $\alpha_{5,1}$    $\alpha_{6,1}$

**Softmax**

$\beta_{1,1}$    $\beta_{2,1}$    $\beta_{3,1}$    $\beta_{4,1}$    $\beta_{5,1}$    $\beta_{6,1}$

Eres

Softmax

0.1
-1.2
...
0.4
0.6

$h_1^{(dec)}$   $h_2^{(dec)}$

$h_0^{(enc)}$   $h_1^{(enc)}$   $h_2^{(enc)}$   $h_3^{(enc)}$   $h_4^{(enc)}$   $h_5^{(enc)}$   $h_6^{(enc)}$

1   0   ...   ...   0   0
0   0       0   1
...   ...       ...   ...
0   1       0   0
0   0       0   0

<SOS>    You    are    my    best    friend

# Context Vector
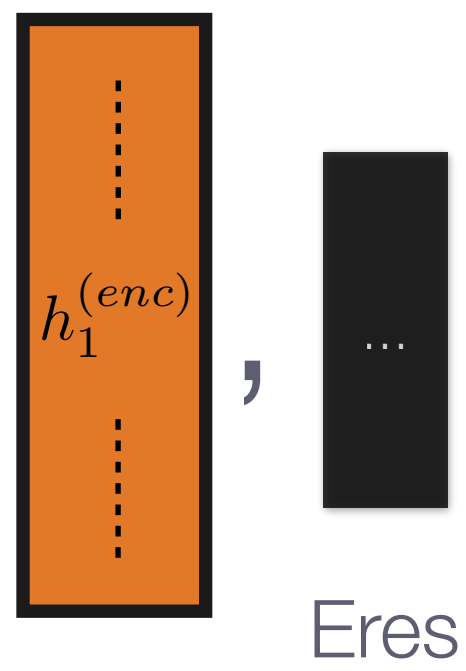
$$c_1$$

1
0
...
0

...

<SOS>    Eres

$$\beta_{1,2} = \text{Score}(\ h_1^{(enc)}\ ,\ h_2^{(dec)}\ )$$

Eres

Softmax

0.1
-1.2
...
0.4
0.6

$h_1^{(dec)}$ → $h_2^{(dec)}$

$h_0^{(enc)}$ → $h_1^{(enc)}$ → $h_2^{(enc)}$ → $h_3^{(enc)}$ → $h_4^{(enc)}$ → $h_5^{(enc)}$ → $h_6^{(enc)}$

1
0
...
0
0

0
0
...
1

...

...

0
0
...
0

0
1
...
0
0

<SOS>    You    are    my    best    friend

1
0
...
0

...

<SOS>    Eres

$$\beta_{1,2} = \text{Score}(\ h_1^{(enc)}\ ,\ h_2^{(dec)}\ )$$

We don't have this
(we need a proxy)

Eres

Softmax

0.1
-1.2
...
0.4
0.6

$h_1^{(dec)}$ → $h_2^{(dec)}$

$h_0^{(enc)}$ → $h_1^{(enc)}$ → $h_2^{(enc)}$ → $h_3^{(enc)}$ → $h_4^{(enc)}$ → $h_5^{(enc)}$ → $h_6^{(enc)}$

1
0
...
0
0

0
0
...
1
0

...

...

0
0
...
0
0

0
1
...
0
0

<SOS>    You    are    my    best    friend

1
0
...
0

...

<SOS>    Eres

$$\beta_{1,2} = \text{Score}( \; h_1^{(enc)} , \quad )$$

**We don't have this
(we need a proxy)**

Eres

Softmax

0.1
-1.2
...
0.4
0.6

$h_1^{(dec)}$

$h_0^{(enc)}$  $h_1^{(enc)}$  $h_2^{(enc)}$  $h_3^{(enc)}$  $h_4^{(enc)}$  $h_5^{(enc)}$  $h_6^{(enc)}$

1      0      ...    ...    0      0
0      0                    0      1
...    ...                  ...    ...
0      1                    0      0
0      0                    0      0

<SOS>   You    are    my    best   friend

1
0
...
0

...

<SOS>    Eres

$$\beta_{1,2} = \text{Score}(\ h_1^{(enc)}\ ,\ h_1^{(dec)}\ )$$

Eres

Softmax

0.1
-1.2
...
0.4
0.6

$h_1^{(dec)}$

$h_0^{(enc)}$ → $h_1^{(enc)}$ → $h_2^{(enc)}$ → $h_3^{(enc)}$ → $h_4^{(enc)}$ → $h_5^{(enc)}$ → $h_6^{(enc)}$

1
0
...
0
0

0
0
...
1
0

...

...

0
0
...
0
0

0
1
...
0
0

<SOS>    You    are    my    best    friend

1
0
...
0

...

<SOS>    Eres

# **Additive** Attention

$$\beta_{1,2} = \text{Score}(\ h_1^{(enc)}\ ,\ h_1^{(dec)}\ ) = \text{NN}(\ \substack{h_1^{(enc)} \\ h_1^{(dec)}}\ )$$

key     query

# **Additive** Attention

$$\beta_{1,2}, \beta_{2,2}, \beta_{3,2}, \beta_{4,2}, \beta_{5,2}, \beta_{6,2} = \text{NN}(\quad)$$

$h_1^{(dec)}$ $h_1^{(dec)}$ $h_1^{(dec)}$ $h_1^{(dec)}$ $h_1^{(dec)}$ $h_1^{(dec)}$

$h_1^{(enc)}$ $h_2^{(enc)}$ $h_3^{(enc)}$ $h_4^{(enc)}$ $h_5^{(enc)}$ $h_6^{(enc)}$

keys

# **Additive** Attention

query (replicated)



$$\beta_{1,2}, \beta_{2,2}, \beta_{3,2}, \beta_{4,2}, \beta_{5,2}, \beta_{6,2} = \text{NN}(\qquad)$$

keys

# **Dot-product** Attention

$$\beta_{1,2} = \text{Score}(\ h_1^{(enc)}\ ,\ h_1^{(dec)}\ ) = h_1^{(enc)} \quad h_1^{(dec)}$$

key  query

# **Dot-product** Attention

$$\beta_{1,2} = \text{Score}(\; h_1^{(enc)} \;,\; h_1^{(dec)} \;) = h_1^{(enc)} \quad \boxed{\text{------} \; h_1^{(dec)} \; \text{------}}$$

key     query

$$\beta_{1,2}, \beta_{2,2}, \beta_{3,2}, \beta_{4,2}, \beta_{5,2}, \beta_{6,2} \; = \; h_1^{(enc)} \; h_2^{(enc)} \; h_3^{(enc)} \; h_4^{(enc)} \; h_5^{(enc)} \; h_6^{(enc)} \quad \boxed{\text{------} \; h_1^{(dec)} \; \text{------}}$$

keys

# **General** Dot-product Attention

$$\beta_{1,2} = \text{Score}( \quad \mathbf{W}_k \quad h_1^{(enc)} , \quad \mathbf{W}_q \quad h_1^{(dec)} \quad )$$

# **Scaled** General Dot-product Attention

$$\beta_{1,2} = \text{Score}( \quad \mathbf{W}_k \quad h_1^{(enc)} \quad , \quad \mathbf{W}_q \quad h_1^{(dec)} \quad )$$

$$n = \text{Length}\left(h_i^{(enc)}\right)$$

$$\hat{\beta}_{1,2} = \frac{\beta_{1,2}}{\sqrt{n}}$$

# **Scaled** General Dot-product Attention

$$\beta_{1,2} = \text{Score}(\; \mathbf{W}_k \;\; h_1^{(enc)} \;,\; \mathbf{W}_q \;\; h_1^{(dec)} \;)$$

$$\hat{\beta}_{1,2} = \frac{\beta_{1,2}}{\sqrt{n}}$$

$$n = \text{Length}\left(h_i^{(enc)}\right)$$

$$n = \text{Length}\left(h_i^{(dec)}\right)$$

# Soft **Attention** in details

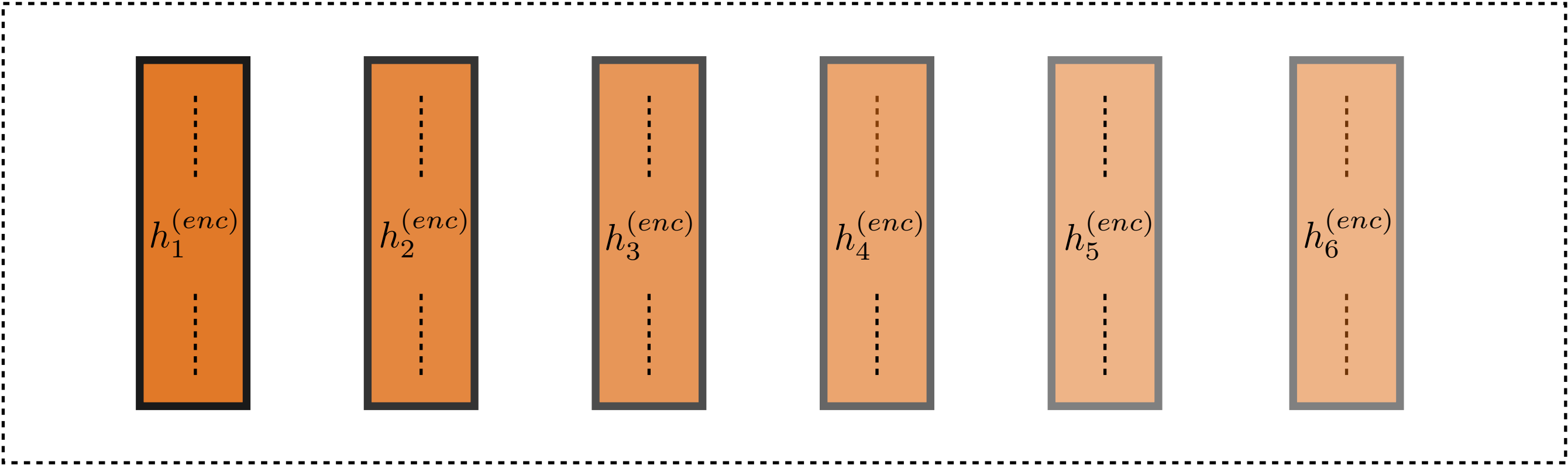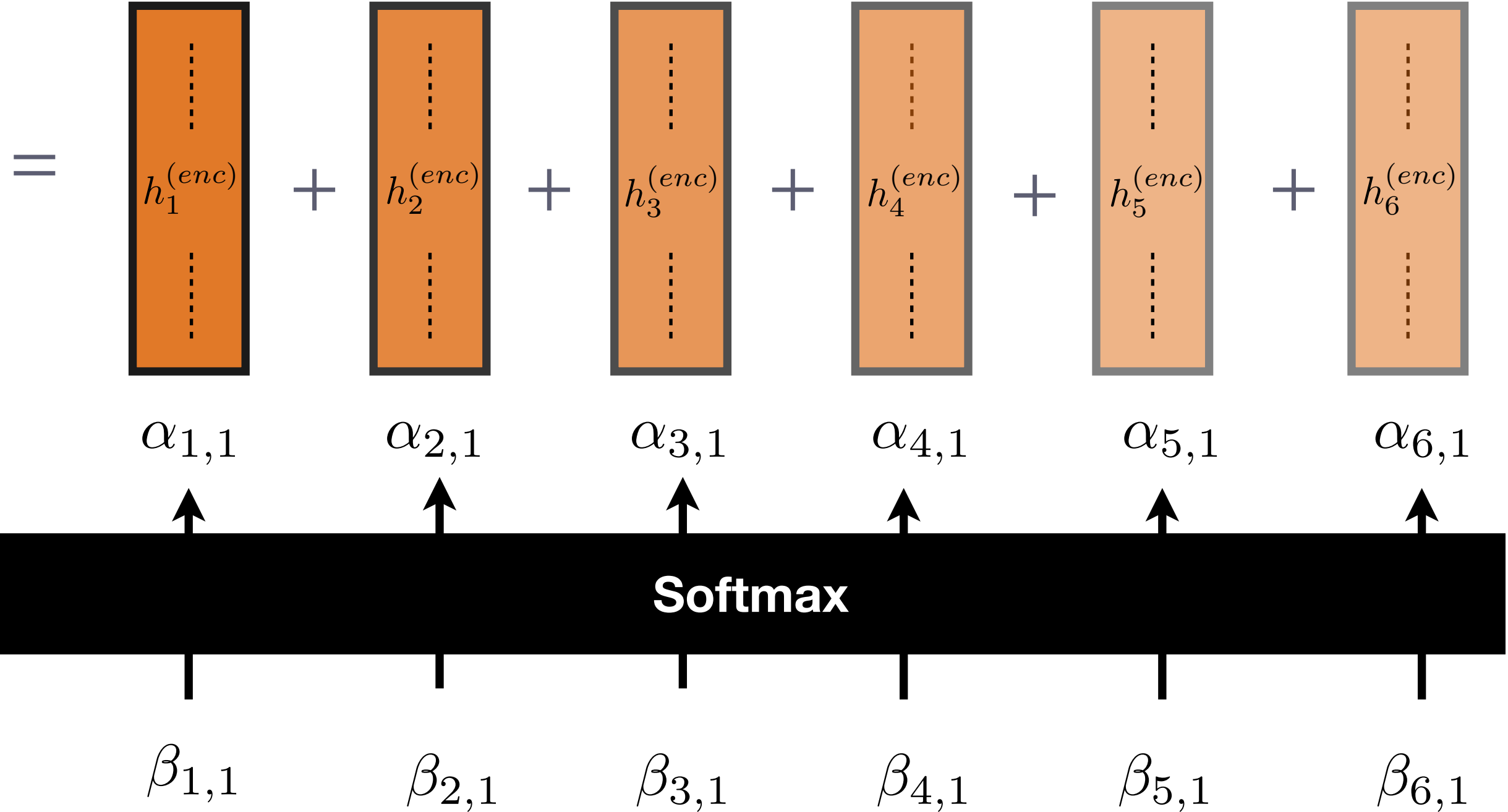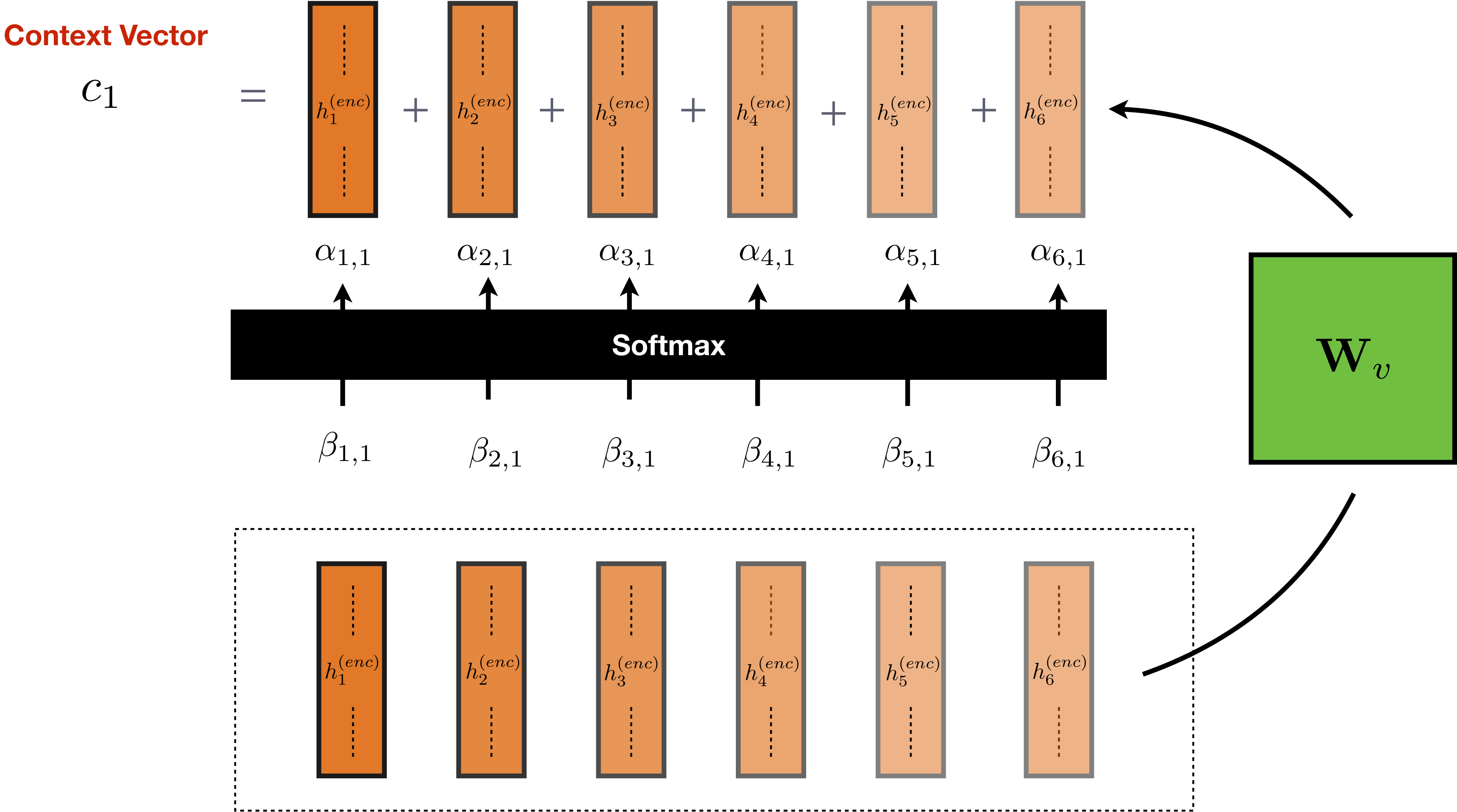| Name | Alignment score function | Citation |
|---|---|---|
| Content-base attention | $\text{score}(\boldsymbol{s}_t, \boldsymbol{h}_i) = \text{cosine}[\boldsymbol{s}_t, \boldsymbol{h}_i]$ | Graves2014 |
| Additive(*) | $\text{score}(\boldsymbol{s}_t, \boldsymbol{h}_i) = \mathbf{v}_a^\top \tanh(\mathbf{W}_a[\boldsymbol{s}_t; \boldsymbol{h}_i])$ | Bahdanau2015 |
| Location-Base | $\alpha_{t,i} = \text{softmax}(\mathbf{W}_a \boldsymbol{s}_t)$ <br> Note: This simplifies the softmax alignment to only depend on the target position. | Luong2015 |
| General | $\text{score}(\boldsymbol{s}_t, \boldsymbol{h}_i) = \boldsymbol{s}_t^\top \mathbf{W}_a \boldsymbol{h}_i$ <br> where $\mathbf{W}_a$ is a trainable weight matrix in the attention layer. | Luong2015 |
| Dot-Product | $\text{score}(\boldsymbol{s}_t, \boldsymbol{h}_i) = \boldsymbol{s}_t^\top \boldsymbol{h}_i$ | Luong2015 |
| Scaled Dot-Product(^) | $\text{score}(\boldsymbol{s}_t, \boldsymbol{h}_i) = \dfrac{\boldsymbol{s}_t^\top \boldsymbol{h}_i}{\sqrt{n}}$ <br> Note: very similar to the dot-product attention except for a scaling factor; where n is the dimension of the source hidden state. | Vaswani2017 |

# Forming a **Context** Vector

**Context Vector**

$$c_1 \quad =$$

$\alpha_{1,1} \qquad \alpha_{2,1} \qquad \alpha_{3,1} \qquad \alpha_{4,1} \qquad \alpha_{5,1} \qquad \alpha_{6,1}$

**Softmax**

$\beta_{1,1} \qquad \beta_{2,1} \qquad \beta_{3,1} \qquad \beta_{4,1} \qquad \beta_{5,1} \qquad \beta_{6,1}$

$h_1^{(enc)} \qquad h_2^{(enc)} \qquad h_3^{(enc)} \qquad h_4^{(enc)} \qquad h_5^{(enc)} \qquad h_6^{(enc)}$

# Forming a **Context** Vector

**Context Vector**

$$c_1 = \quad h_1^{(enc)} + h_2^{(enc)} + h_3^{(enc)} + h_4^{(enc)} + h_5^{(enc)} + h_6^{(enc)}$$

$\alpha_{1,1} \qquad \alpha_{2,1} \qquad \alpha_{3,1} \qquad \alpha_{4,1} \qquad \alpha_{5,1} \qquad \alpha_{6,1}$

**Softmax**

$\beta_{1,1} \qquad \beta_{2,1} \qquad \beta_{3,1} \qquad \beta_{4,1} \qquad \beta_{5,1} \qquad \beta_{6,1}$

$h_1^{(enc)} \quad h_2^{(enc)} \quad h_3^{(enc)} \quad h_4^{(enc)} \quad h_5^{(enc)} \quad h_6^{(enc)}$

# Forming a **Context** Vector

**Context Vector**

$$c_1 = h_1^{(enc)} + h_2^{(enc)} + h_3^{(enc)} + h_4^{(enc)} + h_5^{(enc)} + h_6^{(enc)}$$

$\alpha_{1,1} \quad \alpha_{2,1} \quad \alpha_{3,1} \quad \alpha_{4,1} \quad \alpha_{5,1} \quad \alpha_{6,1}$

**Softmax**

$\beta_{1,1} \quad \beta_{2,1} \quad \beta_{3,1} \quad \beta_{4,1} \quad \beta_{5,1} \quad \beta_{6,1}$

$\mathbf{W}_v$

$h_1^{(enc)} \quad h_2^{(enc)} \quad h_3^{(enc)} \quad h_4^{(enc)} \quad h_5^{(enc)} \quad h_6^{(enc)}$

# Forming a **General Context** Vector

**Context Vector**

$$c_1 = \boxed{h_1^{(enc)}} + \boxed{h_2^{(enc)}} + \boxed{h_3^{(enc)}} + \boxed{h_4^{(enc)}} + \boxed{h_5^{(enc)}} + \boxed{h_6^{(enc)}}$$

$\alpha_{1,1}$     $\alpha_{2,1}$     $\alpha_{3,1}$     $\alpha_{4,1}$     $\alpha_{5,1}$     $\alpha_{6,1}$

**Softmax**

$\mathbf{W}_v$

$\beta_{1,1}$     $\beta_{2,1}$     $\beta_{3,1}$     $\beta_{4,1}$     $\beta_{5,1}$     $\beta_{6,1}$

$h_1^{(enc)}$     $h_2^{(enc)}$     $h_3^{(enc)}$     $h_4^{(enc)}$     $h_5^{(enc)}$     $h_6^{(enc)}$

# Soft **Attention** in details

$y_t$

$\mathbf{h}_t^{(dec)}$

$\beta_{i,t} = score(\mathbf{h}_i^{(enc)}, \mathbf{h}_t^{(dec)})$ ~~(crossed out)~~

Relevance of encoding at token i for decoding token t

**Query**: $\mathbf{Q}_t$

$\beta_{i,t} = score(\mathbf{h}_i^{(enc)}, \mathbf{x}_t^{(dec)})$

$\beta_{i,t} = score(\mathbf{h}_i^{(enc)}, \mathbf{h}_{t-1}^{(dec)})$

$\beta_{i,t} = score(\mathbf{h}_i^{(enc)}, [\mathbf{x}_t^{(dec)}, \mathbf{h}_{t-1}^{(dec)}])$

**Key:** $\mathbf{K}_i$

Attention Layer

$c_t$ Context vector

Global align weights

$a_t$

$\mathbf{h}_i^{(enc)}$

$\alpha_{i,t} = \text{Softmax}(\beta_{i,t})$

Normalize the weights to sum to 1

$\mathbf{c}_t = \sum_i \alpha_{i,t} \mathbf{h}_i^{(enc)}$

**Value:** $\mathbf{V}_i$

Form a context vector that would simply be added to the standard decoder input

# Generalized Soft **Attention** in details

$\beta_{i,t} = score(\mathbf{h}_i^{(enc)}, \mathbf{h}_t^{(dec)})$ ~~(crossed out)~~

Relevance of encoding at token i for decoding token t

**Query**: $\mathbf{Q}_t$

$\beta_{i,t} = score(\mathbf{W}_k \mathbf{h}_i^{(enc)}, \mathbf{W}_q \mathbf{x}_t^{(dec)})$

$\beta_{i,t} = score(\mathbf{W}_k \mathbf{h}_i^{(enc)}, \mathbf{W}_q \mathbf{h}_{t-1}^{(dec)})$

$\beta_{i,t} = score(\mathbf{W}_k \mathbf{h}_i^{(enc)}, \mathbf{W}_q [\mathbf{x}_t^{(dec)}, \mathbf{h}_{t-1}^{(dec)}])$

**Key:** $\mathbf{K}_i$

$\alpha_{i,t} = \mathrm{Softmax}(\beta_{i,t})$

Normalize the weights to sum to 1

$\mathbf{c}_t = \sum_i \alpha_{i,t} \mathbf{W}_v \mathbf{h}_i^{(enc)}$

**Value:** $\mathbf{V}_i$

Form a context vector that would simply be added to the standard decoder input

$y_t$

$\mathbf{h}_t^{(dec)}$

*Attention Layer*

$\mathbf{c}_t$ — Context vector

Global align weights

$\boldsymbol{a}_t$

$\mathbf{h}_i^{(enc)}$

# **Attention** Mechanisms and RNNs

[ Cho et al., 2015 ]

# **Self** Attention



(Source-Target-Attention)

σ ( Query   Key ) Value

f(u)

Target
u

g(x)    h(x)

Source
x

(Self-Attention)

σ ( Query   Key ) Value

f(x)   g(x)   h(x)

Self
x

# **Self** Attention



Encoder (English)

$h_0^{(enc)}$ → $h_1^{(enc)}$ → $h_2^{(enc)}$ → $h_3^{(enc)}$ → $h_4^{(enc)}$ → $h_5^{(enc)}$ → $h_6^{(enc)}$

| 1 | | 0 | | | | 0 | | 0 |
| 0 | | 0 | | | | 0 | | 1 |
| ... | | ... | | ... | | ... | | ... |
| 0 | | 1 | | | | 0 | | 0 |
| 0 | | 0 | | | | 0 | | 0 |

<SOS>    You    are    my    best    friend

# **Self** Attention



Encoder (English)

$h_1^{(enc)}$    $h_2^{(enc)}$    $h_3^{(enc)}$    $h_4^{(enc)}$    $h_5^{(enc)}$    $h_6^{(enc)}$

$c_1$    $c_2$    $c_3$    $c_4$    $c_5$    $c_6$

$h_1^{(enc)}$    $h_2^{(enc)}$    $h_3^{(enc)}$    $h_4^{(enc)}$    $h_5^{(enc)}$    $h_6^{(enc)}$

&lt;SOS&gt;    You    are    my    best    friend

# Transformers: Attention is all you need

# Transformers: Attention is all you need (Encoder)



**Note**: for assignment you are <u>not</u> implementing transformer encoder

# **Transformers**: Attention is all you need  (Encoder)



**Note**: for assignment you are <u>not</u> implementing transformer encoder

# **Transformers**: Attention is all you need  (Encoder)



$$PE_{(pos,2i)} = sin(pos/10000^{2i/d_{model}})$$

$$PE_{(pos,2i+1)} = cos(pos/10000^{2i/d_{model}})$$

**Note**: for assignment you are <u>not</u> implementing transformer encoder

# **Transformers**: Attention is all you need  (Encoder)



**Note**: for assignment you are <u>not</u> implementing transformer encoder

# **Transformers**: Attention is all you need (Encoder)



**Note**: for assignment you are <u>not</u> implementing transformer encoder

# **Transformers**: Attention is all you need (Encoder)



**Note**: for assignment you are <u>not</u> implementing transformer encoder

# Transformers: Attention is all you need (Encoder)

Normalize by sqrt of dimensionality
(leads to more stable gradients)

# Transformers: Attention is all you need (Encoder)

# Transformers: Attention is all you need (Encoder)

# Transformers: Attention is all you need (Encoder)

In practice, we use multiple self-attention heads

# **Transformers**: Attention is all you need  (Encoder)



In practice, we use multiple self-attention heads
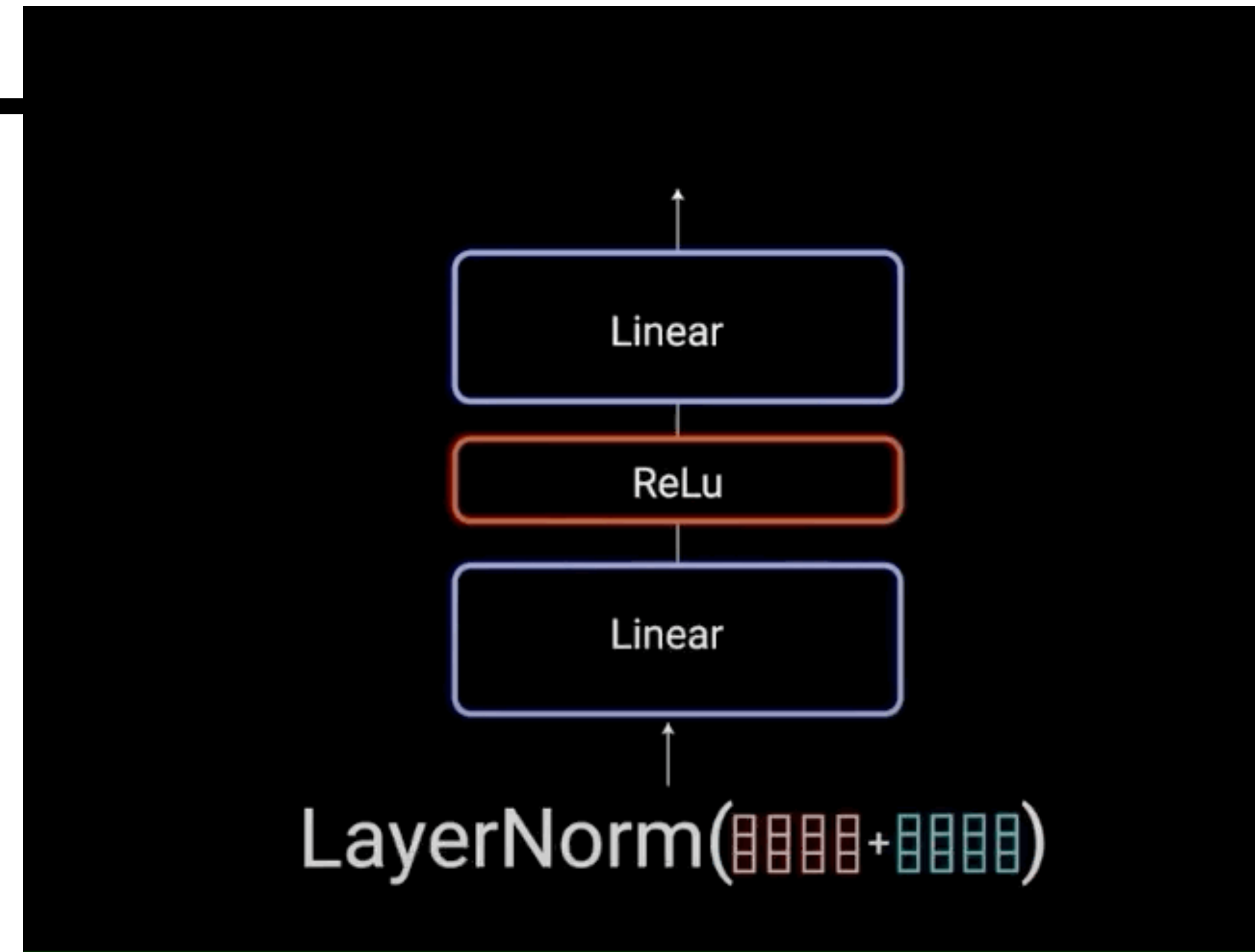
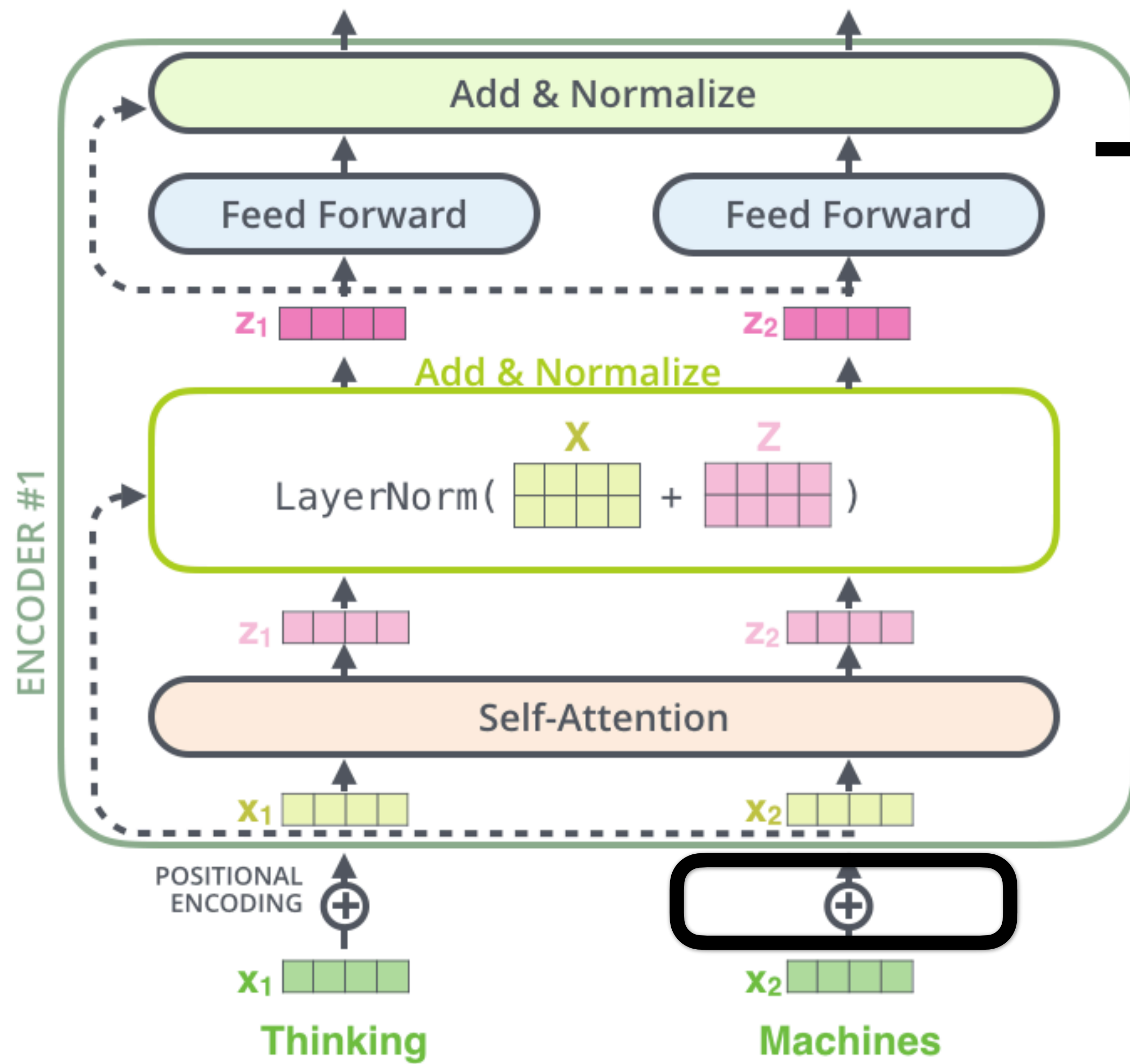# Transformers: Attention is all you need  (Encoder)

Residual connection with LayerNorm

# **Transformers**: Attention is all you need  (Encoder)
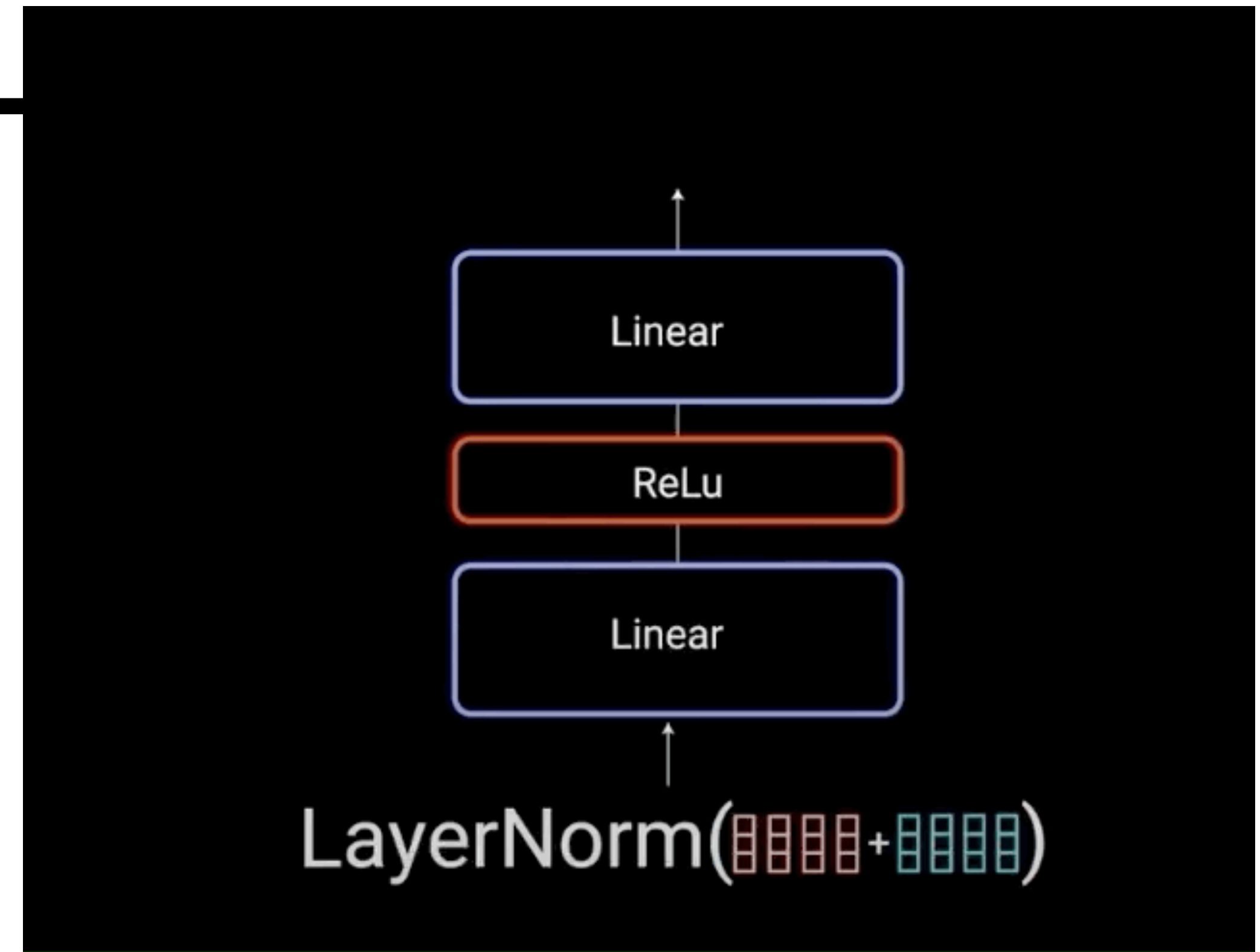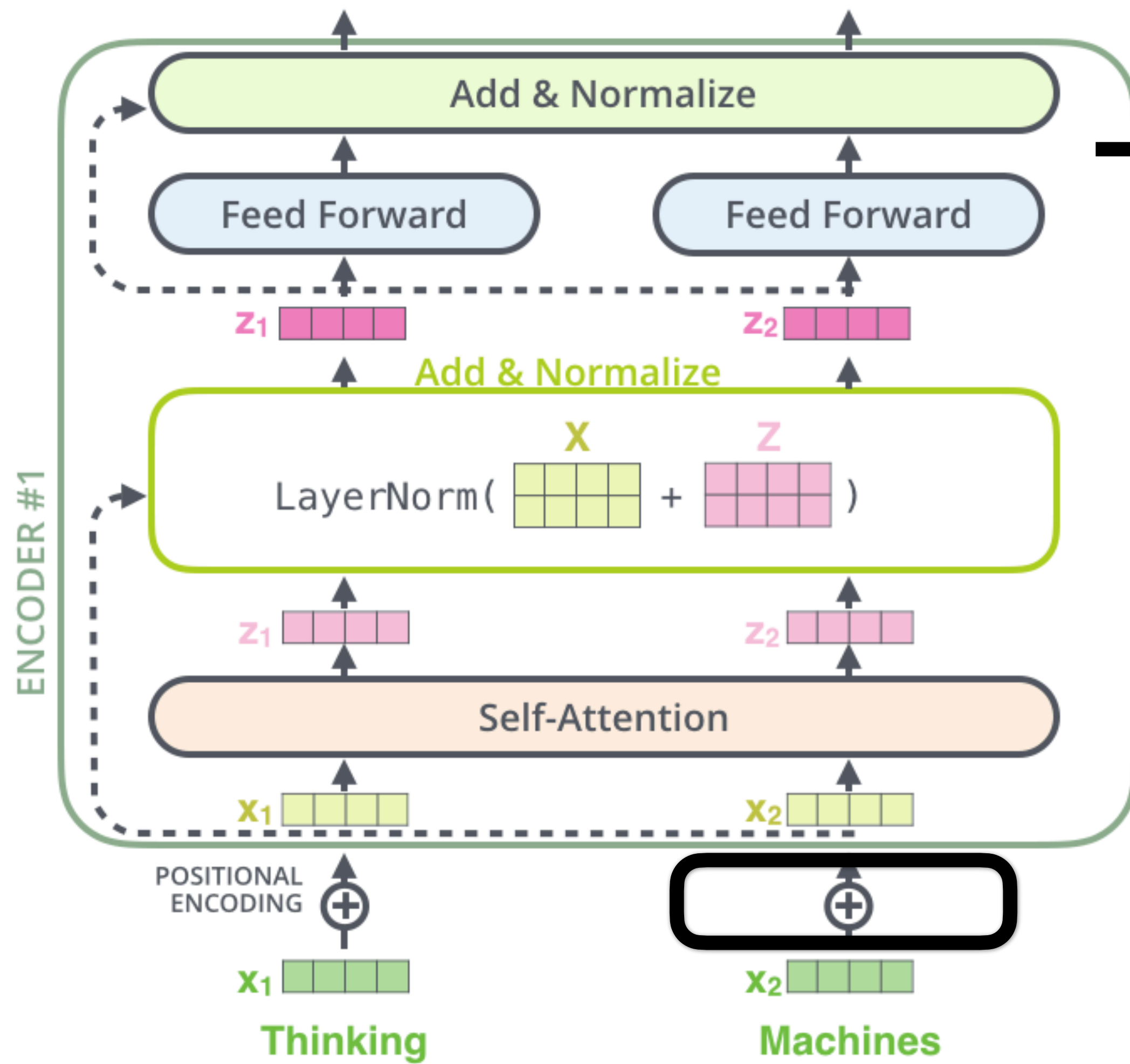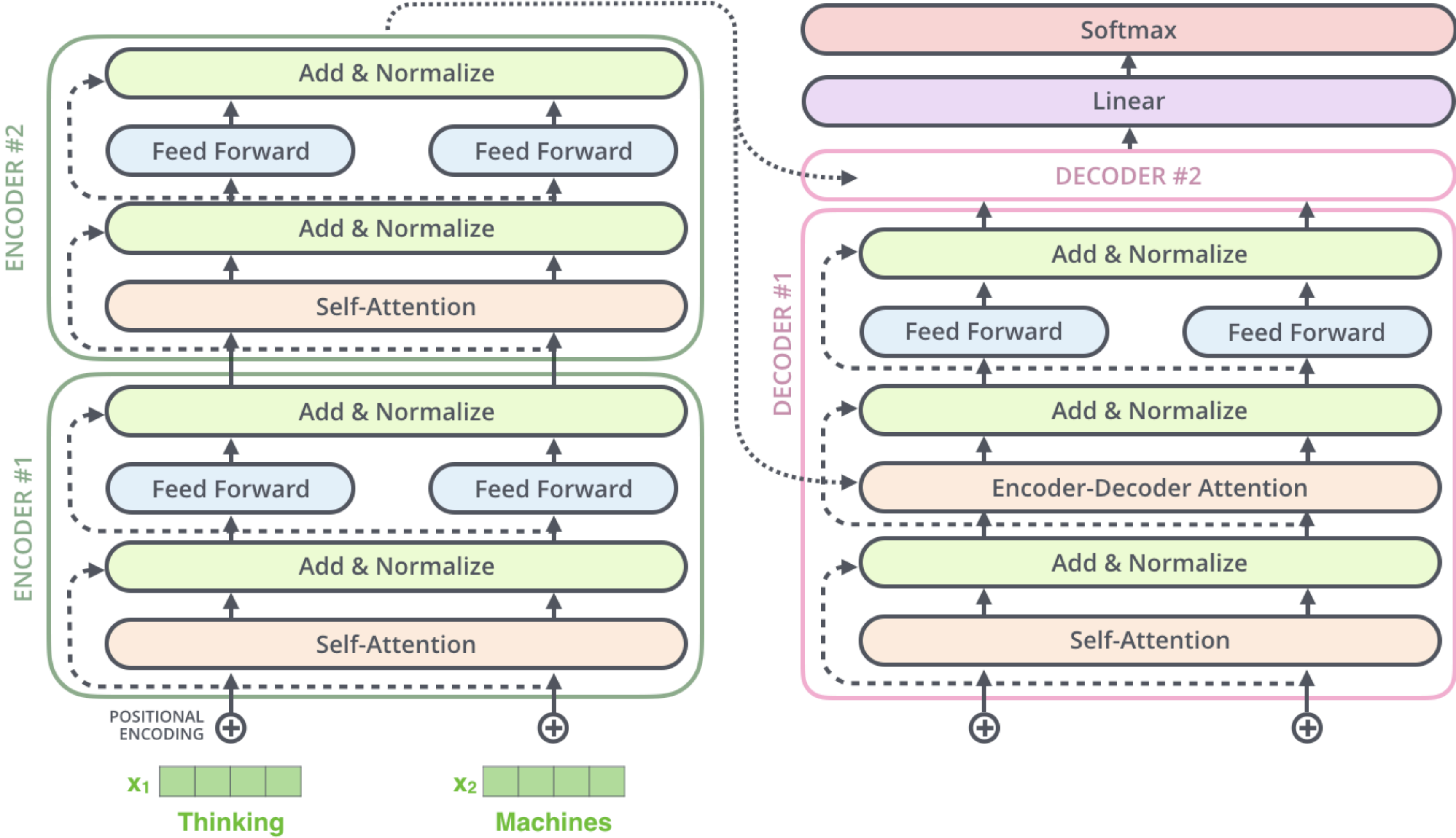


Residual connection with LayerNorm

# Transformers: Attention is all you need (Encoder)

# **Transformers**: Attention is all you need  (Encoder)

# Transformers: Attention is all you need (Encoder)

# Transformers: Attention is all you need

# **Transformers**: Attention is all you need

# **Transformers**: Attention is all you need

# **Self** Attention



The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

# **Benefits** of Transformers

1. Tokens are processes in **parallel** in both encoder and decoder, which is much faster than RNN or LSTM

2. Can (in principle) **model infinite history**, unlike RNN or LSTM that typically only carries context for relatively small number of steps

3. **No gradient flow issues**, due to residual architecture design of Transformer layers — similar to LSTM in some sense.
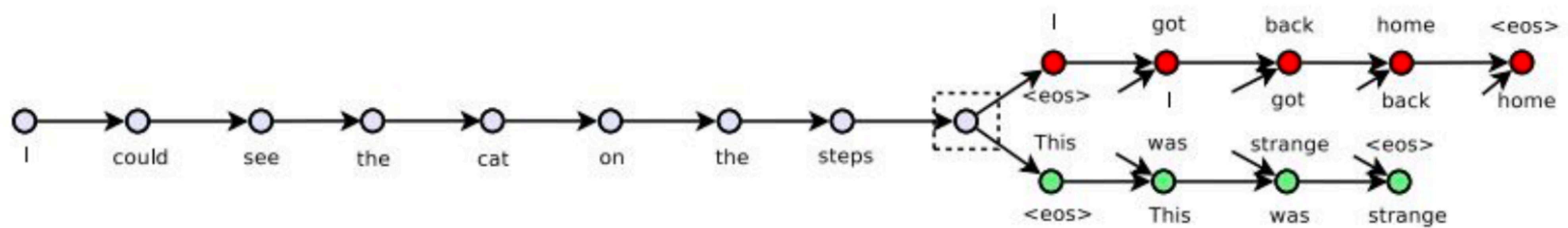
# **Benefits** of Transformers

**Note:** In principle Transformer can model RNN-line or LSTM-like recursion by using causal mask and computing relevance based on "positional" information stored in a token representation and context based on "content" information stored in a token

(in other words, it is more or less strict generalization)
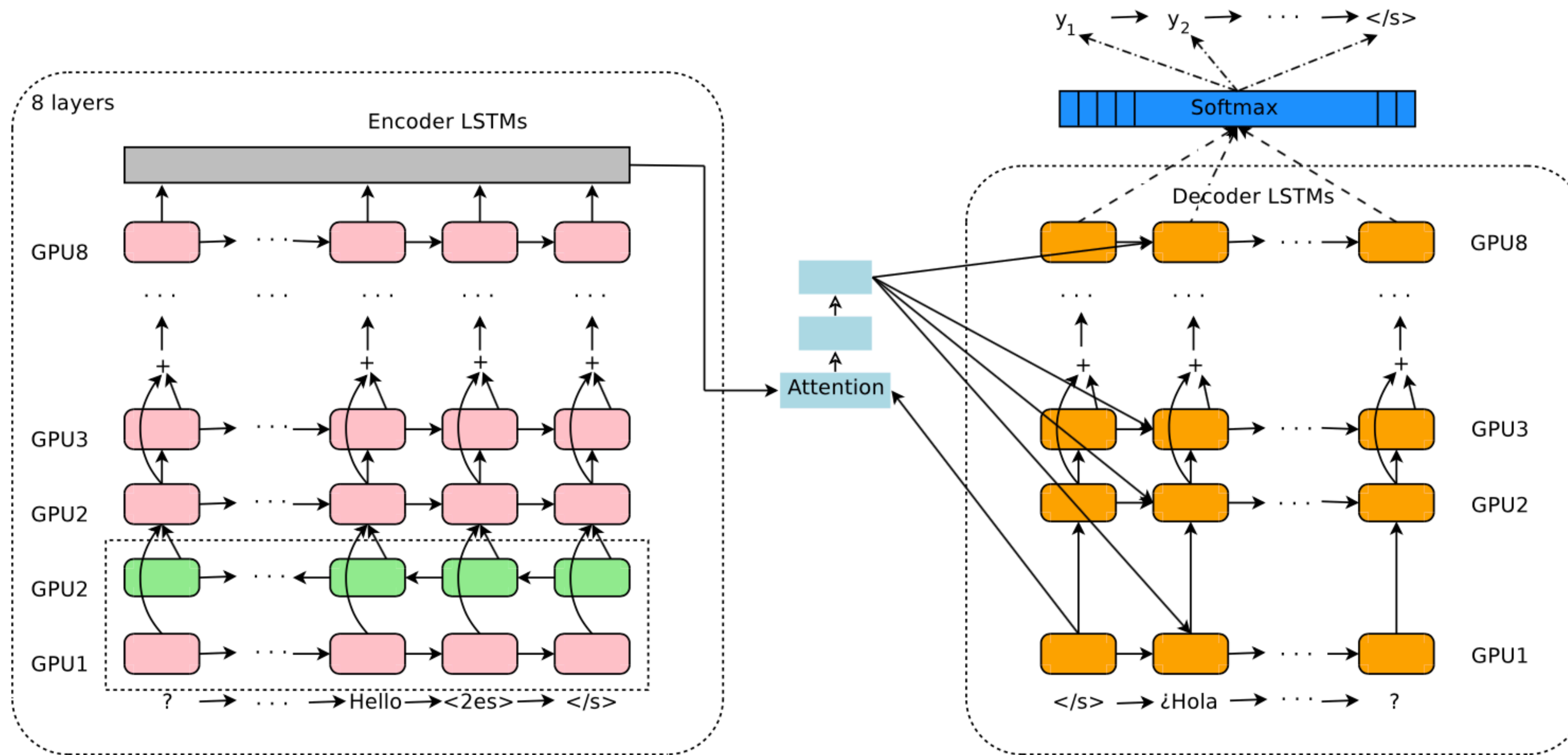
Let us look at some actual practical uses of RNNs

# Applications: Skip-thought Vectors

word2vec but for sentences, where each sentence is processed by an LSTM



[ Kiros et al., 2015 ]

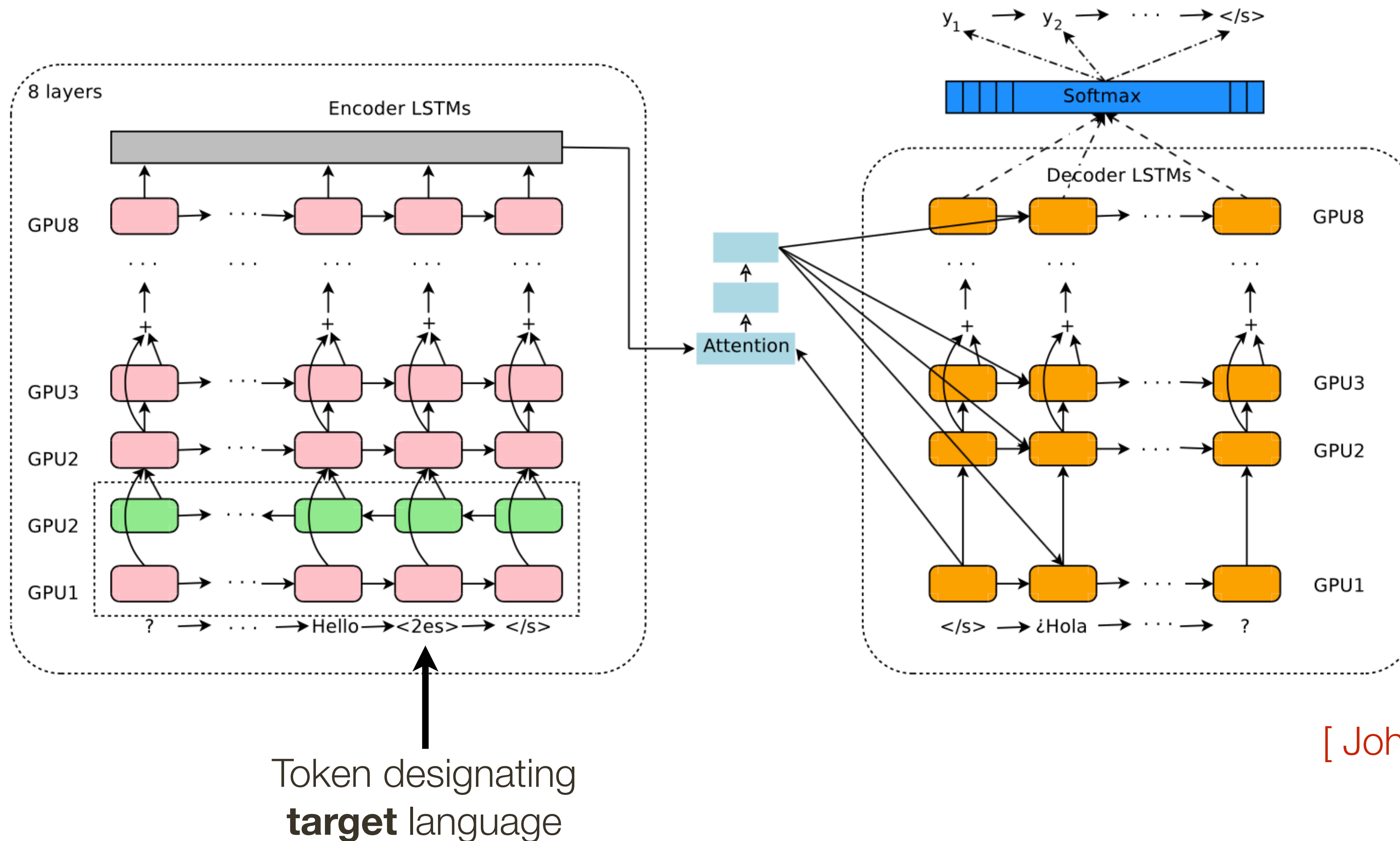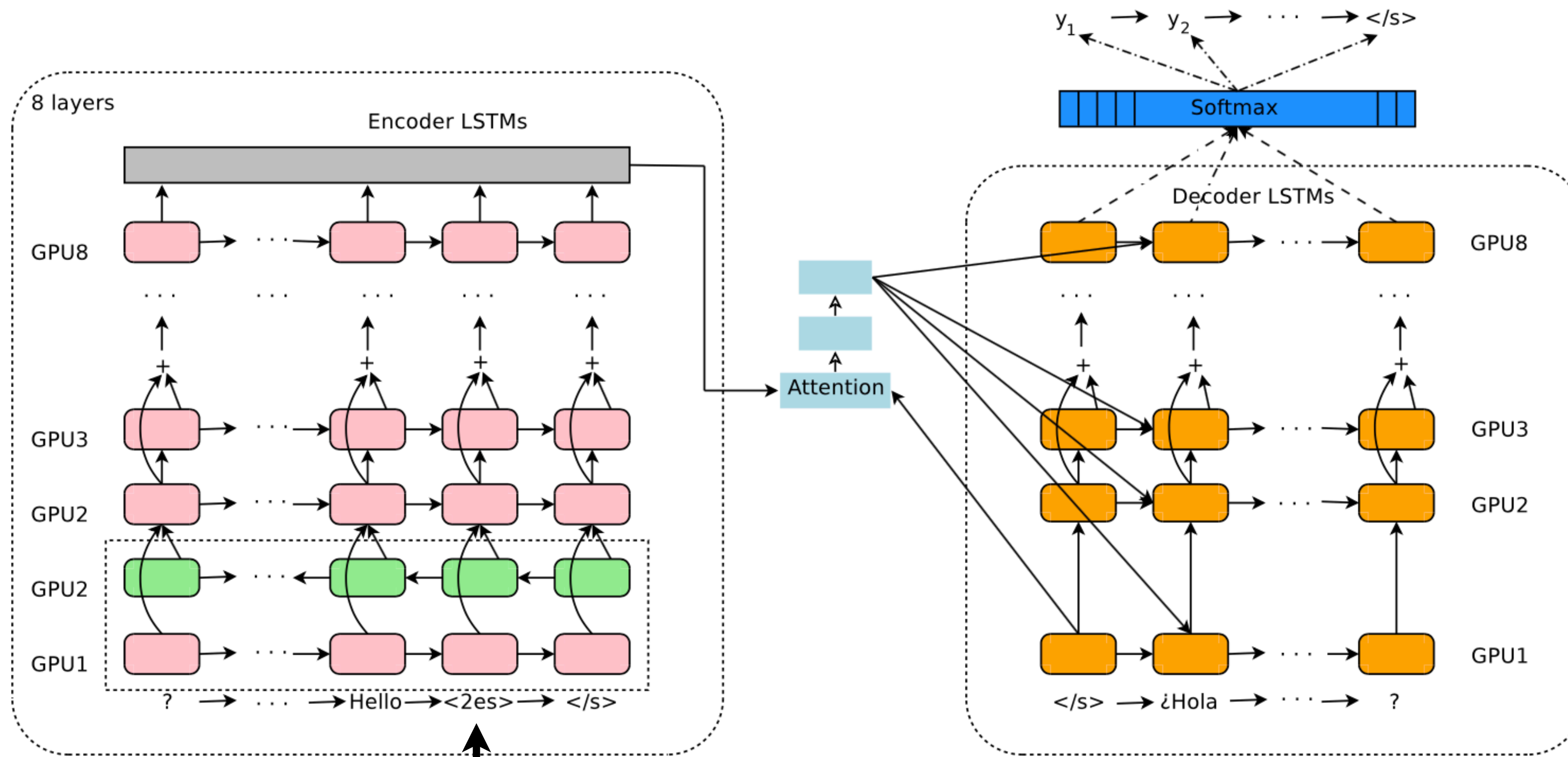# **Applications:** Google Language Translation

One model to translate from **any language** to any other language



[ Johnson et al., 2017 ]

# **Applications:** Google Language Translation

One model to translate from **any language** to any other language



Token designating
**target** language

[ Johnson et al., 2017 ]

# Applications: Google Language Translation

One model to translate from **any language** to any other language
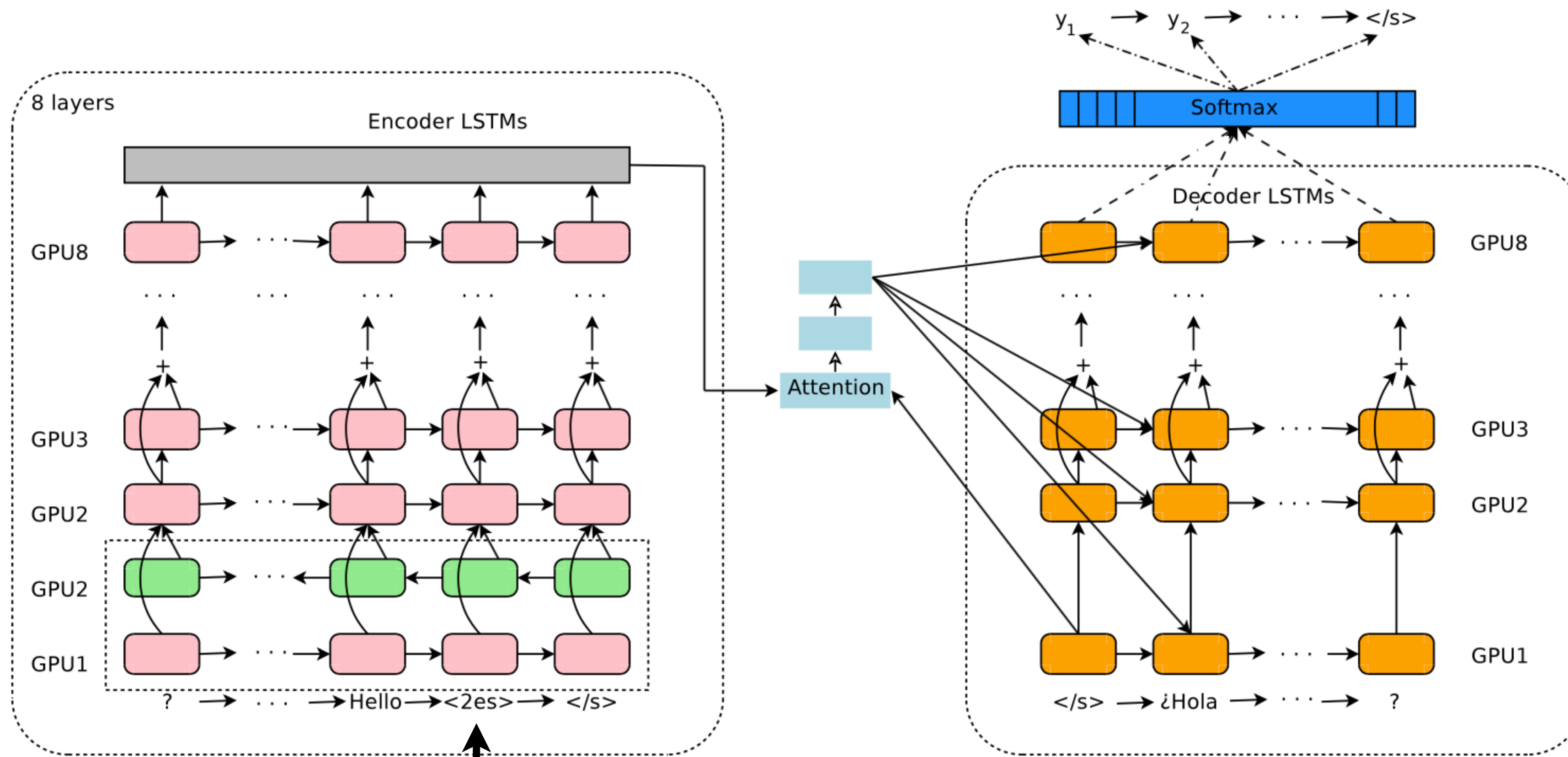


**Flipped** order encoding

Token designating
**target** language

[ Johnson et al., 2017 ]

# **Applications:** Google Language Translation

One model to translate from **any language** to any other language



**Flipped** order encoding

**Why?**

Token designating
**target** language

[ Johnson et al., 2017 ]

# **Applications:** Google Language Translation

One model to translate from **any language** to any other language
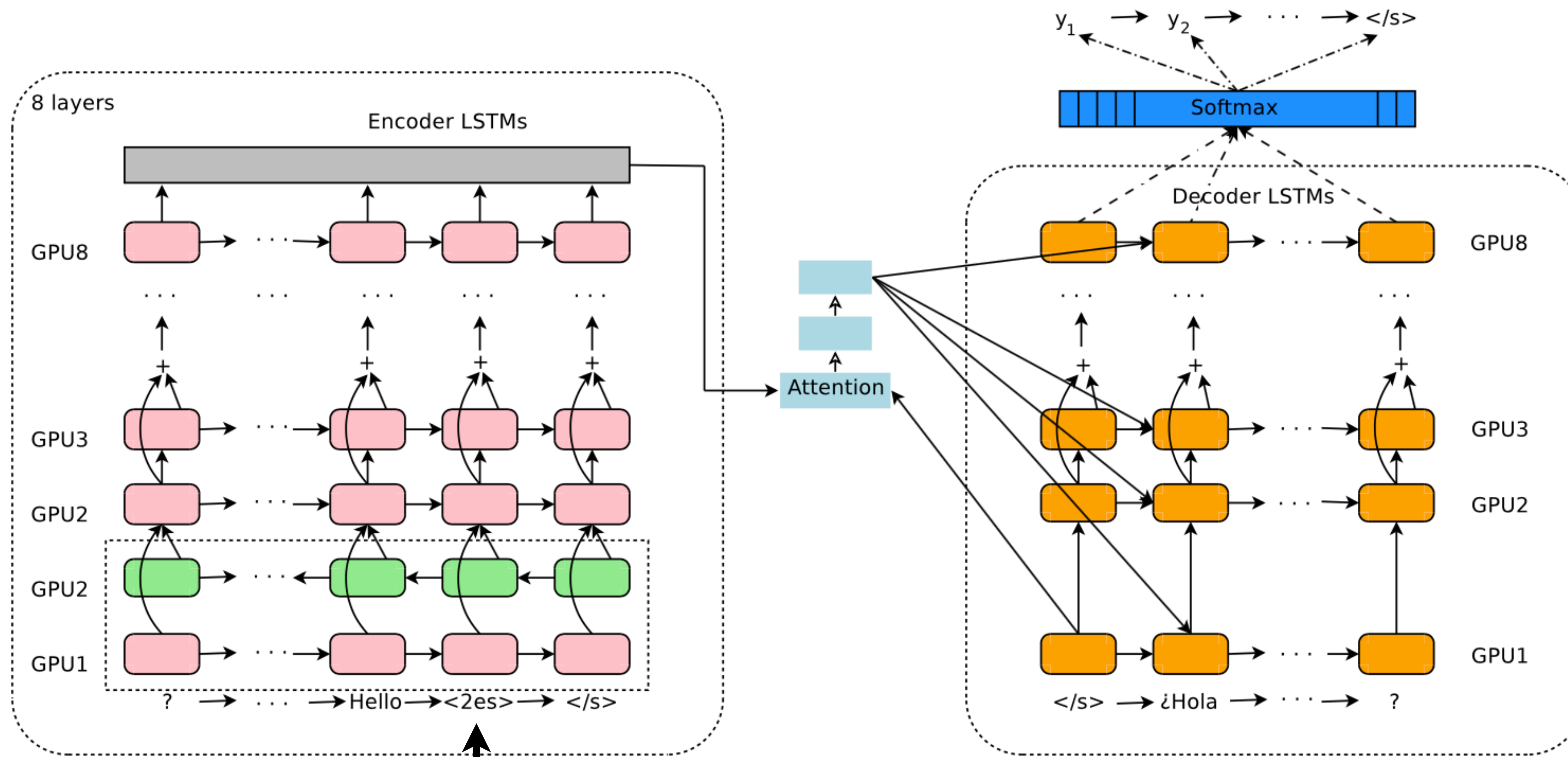


**Flipped** order encoding

Token designating
**target** language

[ Johnson et al., 2017 ]

8! layer LSTM decoder and encoder

# **Applications:** Google Language Translation

One model to translate from **any language** to any other language



**Residual** at
other layers
(ResNet style)

**Bi-directional**
at lower layers

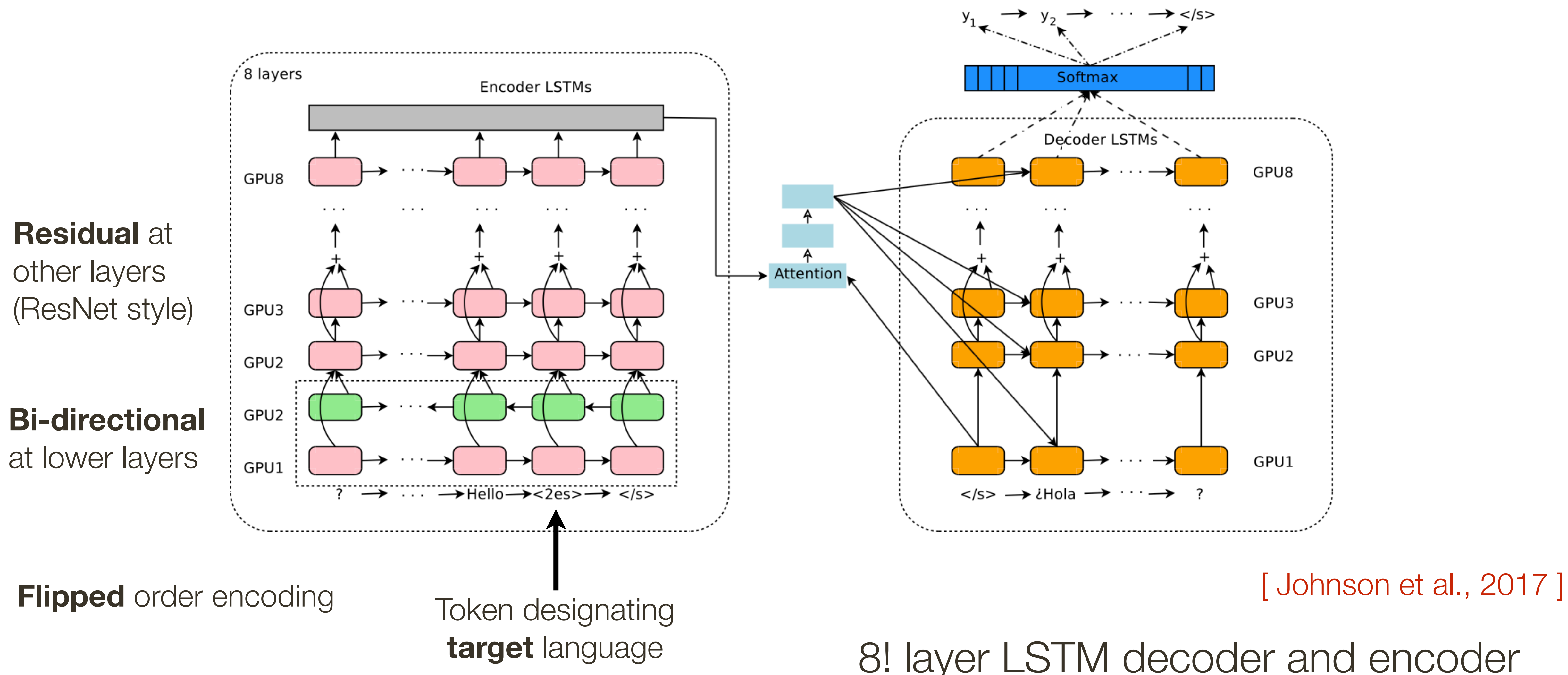**Flipped** order encoding

Token designating
**target** language

[ Johnson et al., 2017 ]

8! layer LSTM decoder and encoder

# **Applications:** Google Language Translation

One model to translate from **any language** to any other language



**Residual** at other layers (ResNet style)

**Bi-directional** at lower layers

**Flipped** order encoding

Token designating **target** language

[ Johnson et al., 2017 ]

8! layer LSTM decoder and encoder