# ITCS 6156/8156 Fall 2023
# Machine Learning

# Attention & Transformers

Instructor: Hongfei Xue
Email: hongfei.xue@charlotte.edu
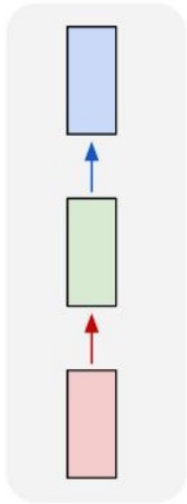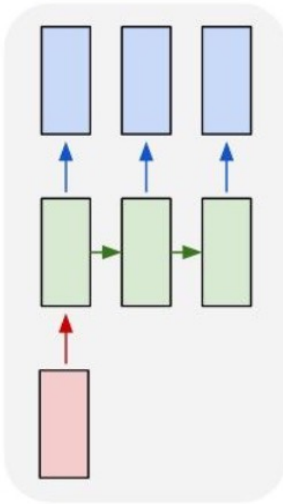Class Meeting: Mon & Wed, 4:00 PM – 5:15 PM, CHHS 376

UNIVERSITY OF NORTH CAROLINA
CHARLOTTE

Some content in the slides is based on Dr. Ruohan Gao's lectures
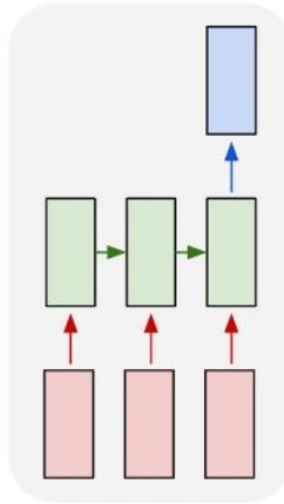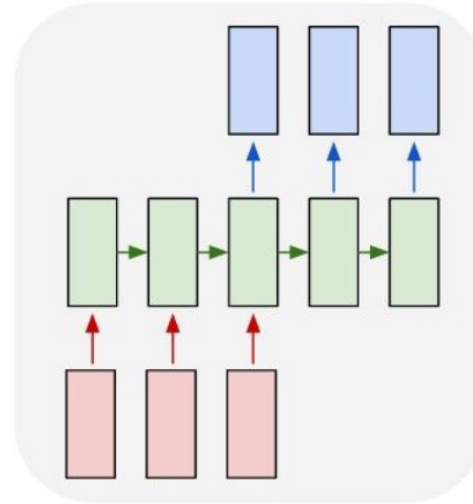
# Recurrent Neural Networks



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

# Sequence to Sequence with RNNs

**Input**: Sequence $x_1, \ldots x_T$
**Output**: Sequence $y_1, \ldots, y_{T'}$

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



Sutskever et al, "Sequence to sequence learning with neural networks", NeurIPS 2014

**Input**: Sequence $x_1, \ldots x_T$
**Output**: Sequence $y_1, \ldots, y_{T'}$

From final hidden state predict:

**Encoder**: $h_t = f_W(x_t, h_{t-1})$  **Initial decoder state** $s_0$
**Context vector** c (often $c = h_T$)



we     are     eating     bread

Sutskever et al, "Sequence to sequence learning with neural networks", NeurIPS 2014

**Input**: Sequence $x_1, \ldots x_T$
**Output**: Sequence $y_1, \ldots, y_{T'}$

**Decoder**: $s_t = g_U(y_{t-1}, s_{t-1}, c)$

estamos

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

From final hidden state predict:
**Initial decoder state** $s_0$
**Context vector** $c$ (often $c=h_T$)



Sutskever et al, "Sequence to sequence learning with neural networks", NeurIPS 2014

# Sequence to Sequence with RNNs

**Input**: Sequence $x_1, \ldots x_T$
**Output**: Sequence $y_1, \ldots, y_{T'}$

**Decoder**: $s_t = g_U(y_{t-1}, s_{t-1}, c)$

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

From final hidden state predict:
**Initial decoder state** $s_0$
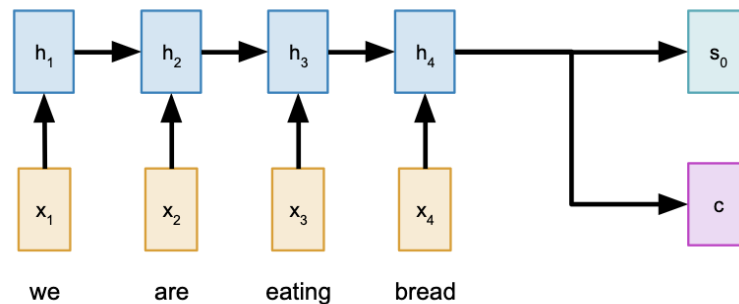**Context vector** $c$ (often $c=h_T$)



Sutskever et al, "Sequence to sequence learning with neural networks", NeurIPS 2014

# Sequence to Sequence with RNNs

**Input**: Sequence $x_1, \ldots x_T$
**Output**: Sequence $y_1, \ldots, y_{T'}$

**Decoder**: $s_t = g_U(y_{t-1}, s_{t-1}, c)$

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

From final hidden state predict:
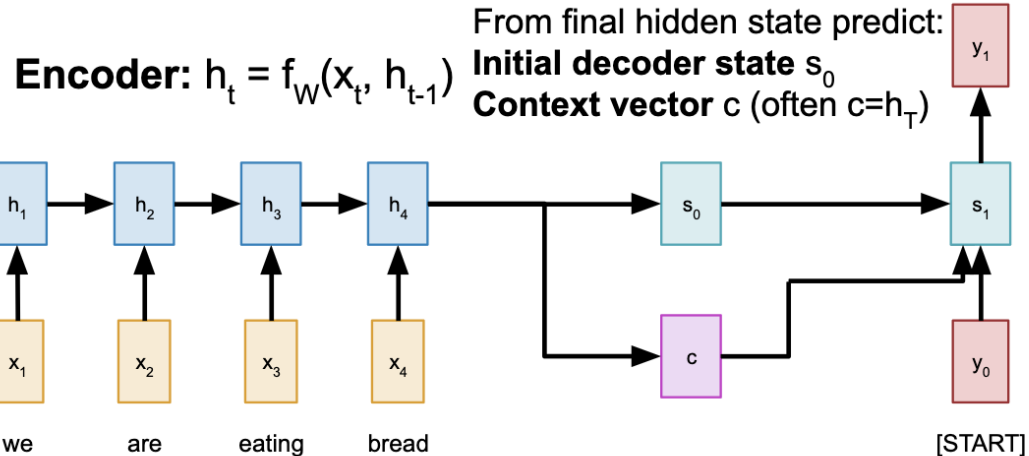**Initial decoder state** $s_0$
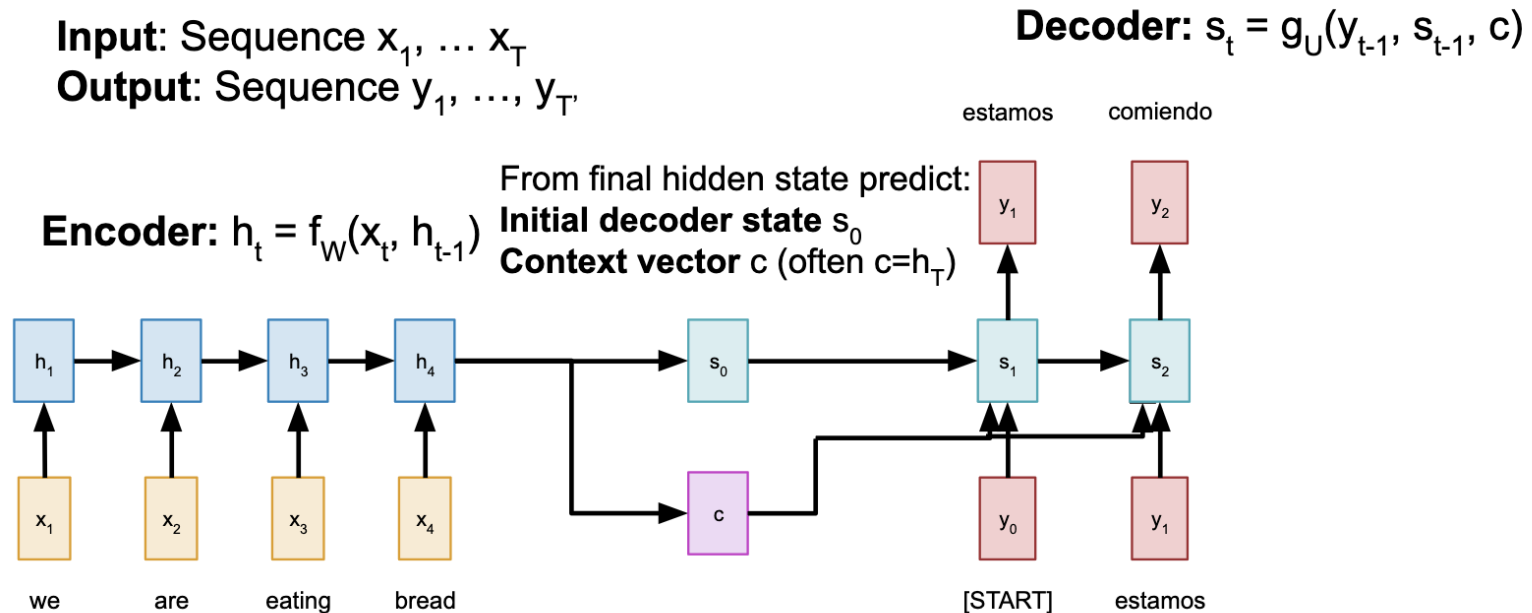**Context vector** $c$ (often $c=h_T$)
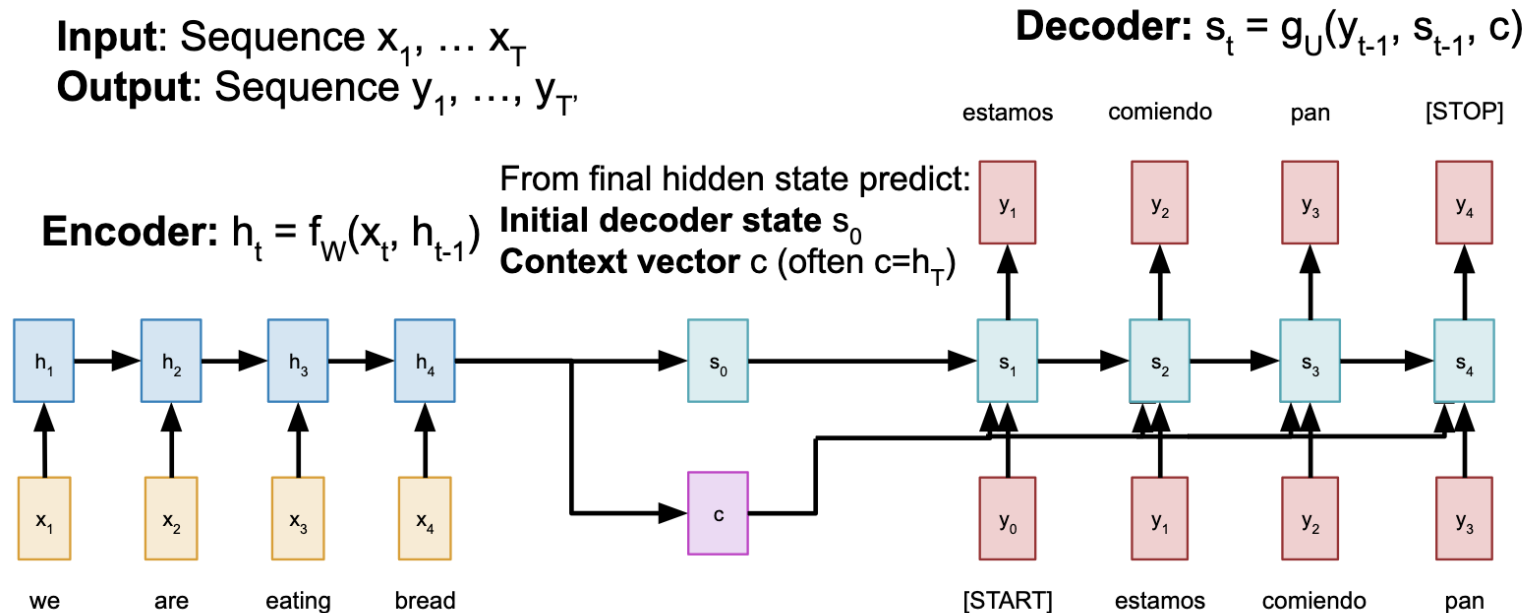


Sutskever et al, "Sequence to sequence learning with neural networks", NeurIPS 2014

# Sequence to Sequence with RNNs

**Input**: Sequence $x_1, \ldots x_T$
**Output**: Sequence $y_1, \ldots, y_{T'}$

**Decoder**: $s_t = g_U(y_{t-1}, s_{t-1}, c)$

From final hidden state predict:
**Initial decoder state** $s_0$

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

**Context vector** $c$ (often $c = h_T$)



**Problem: Input sequence bottlenecked through fixed-sized vector. What if T=1000 ?**

**Input**: Sequence $x_1, \ldots x_T$
**Output**: Sequence $y_1, \ldots, y_{T'}$

**Decoder**: $s_t = g_U(y_{t-1}, s_{t-1}, c)$

estamos    comiendo    pan    [STOP]

From final hidden state predict:
**Initial decoder state** $s_0$
**Context vector** $c$ (often $c=h_T$)

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

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| $h_1$ | $h_2$ | $h_3$ | $h_4$ | | $s_0$ | $s_1$ | $s_2$ | $s_3$ | $s_4$ |

| $y_1$ | $y_2$ | $y_3$ | $y_4$ |

| $s_1$ | $s_2$ | $s_3$ | $s_4$ |

| $x_1$ | $x_2$ | $x_3$ | $x_4$ |

| $c$ |

| $y_0$ | $y_1$ | $y_2$ | $y_3$ |

we    are    eating    bread

[START]    estamos    comiendo    pan

**Problem: Input sequence bottlenecked through fixed-sized vector. What if T=1000 ?**

**Idea: use new context vector at each step of decoder!**

9

**Input**: Sequence $x_1, \ldots x_T$
**Output**: Sequence $y_1, \ldots, y_{T'}$

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

From final hidden state:
**Initial decoder state** $s_0$



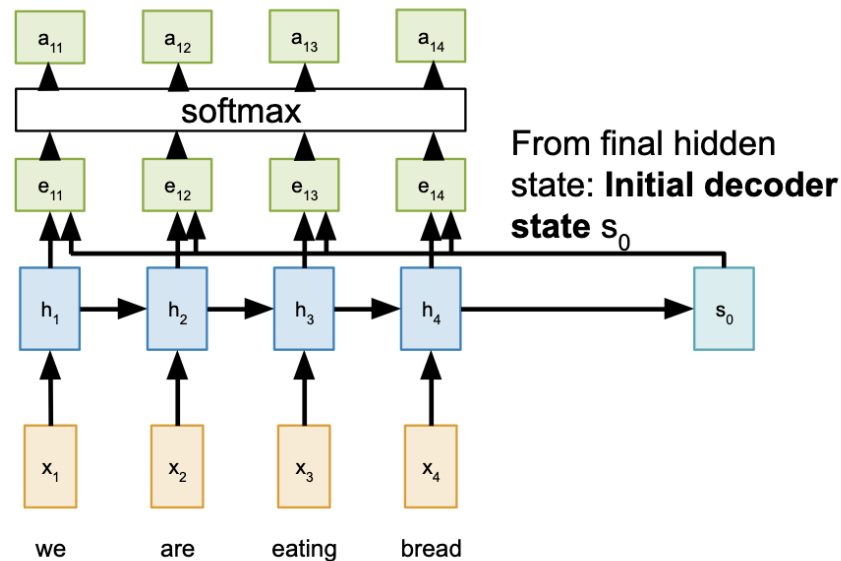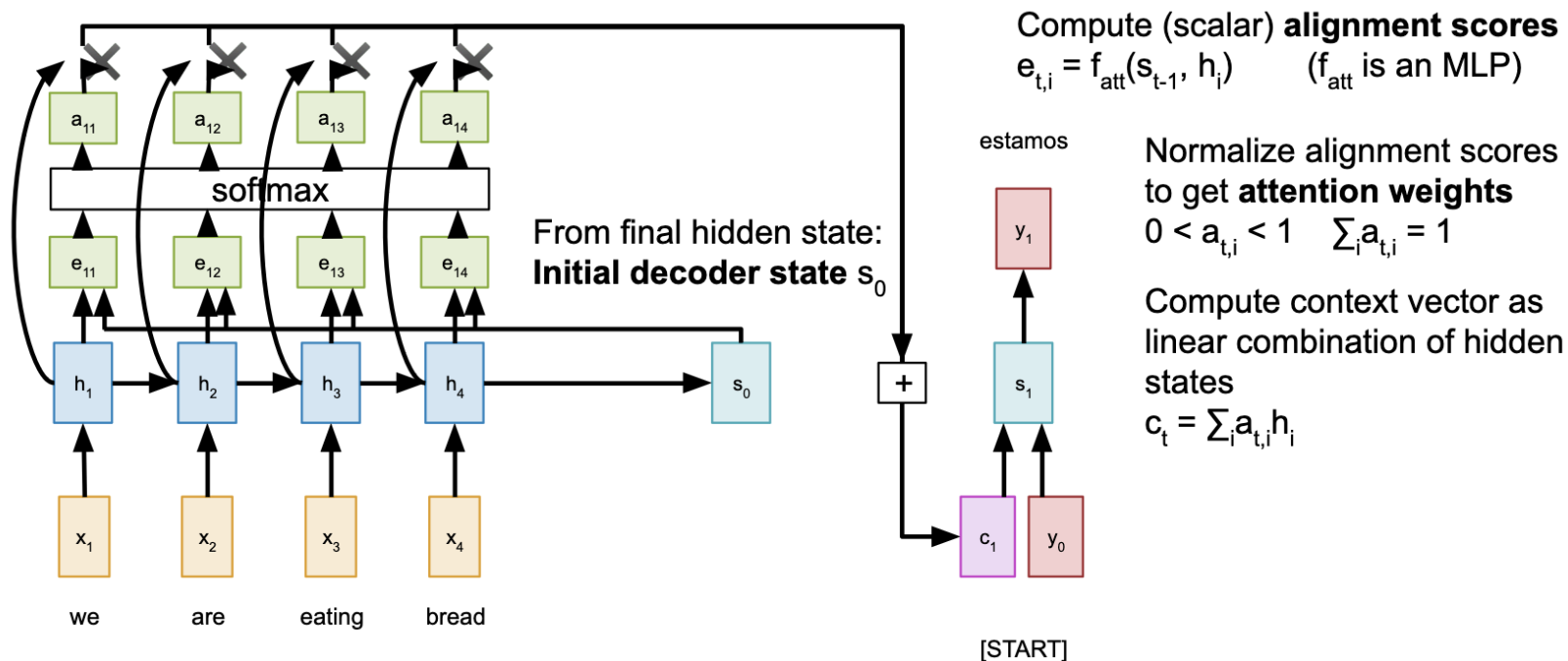Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

# Sequence to Sequence with RNNs and Attention

Compute (scalar) **alignment scores**
$e_{t,i} = f_{att}(s_{t-1}, h_i)$     ($f_{att}$ is an MLP)

From final hidden state:
**Initial decoder state** $s_0$



we     are     eating     bread

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

# Sequence to Sequence with RNNs and Attention



From final hidden state: **Initial decoder state** $s_0$

Compute (scalar) **alignment scores**
$e_{t,i} = f_{att}(s_{t-1}, h_i)$    ($f_{att}$ is an MLP)

Normalize alignment scores to get **attention weights**
$0 < a_{t,i} < 1$    $\sum_i a_{t,i} = 1$

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

Compute (scalar) **alignment scores**
$e_{t,i} = f_{att}(s_{t-1}, h_i)$     ($f_{att}$ is an MLP)

Normalize alignment scores to get **attention weights**
$0 < a_{t,i} < 1$     $\sum_i a_{t,i} = 1$

Compute context vector as linear combination of hidden states
$c_t = \sum_i a_{t,i} h_i$

From final hidden state:
**Initial decoder state** $s_0$

we     are     eating     bread

estamos

[START]

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

13

Compute (scalar) **alignment scores**
$e_{t,i} = f_{att}(s_{t-1}, h_i)$     ($f_{att}$ is an MLP)

Normalize alignment scores to get **attention weights**
$0 < a_{t,i} < 1$     $\sum_i a_{t,i} = 1$

Compute context vector as linear combination of hidden states
$c_t = \sum_i a_{t,i} h_i$

Use context vector in decoder: $s_t = g_U(y_{t-1}, s_{t-1}, c_t)$

From final hidden state:
**Initial decoder state** $s_0$

**Intuition**: Context vector <u>attends</u> to the relevant part of the input sequence
*"estamos"* = *"we are"*
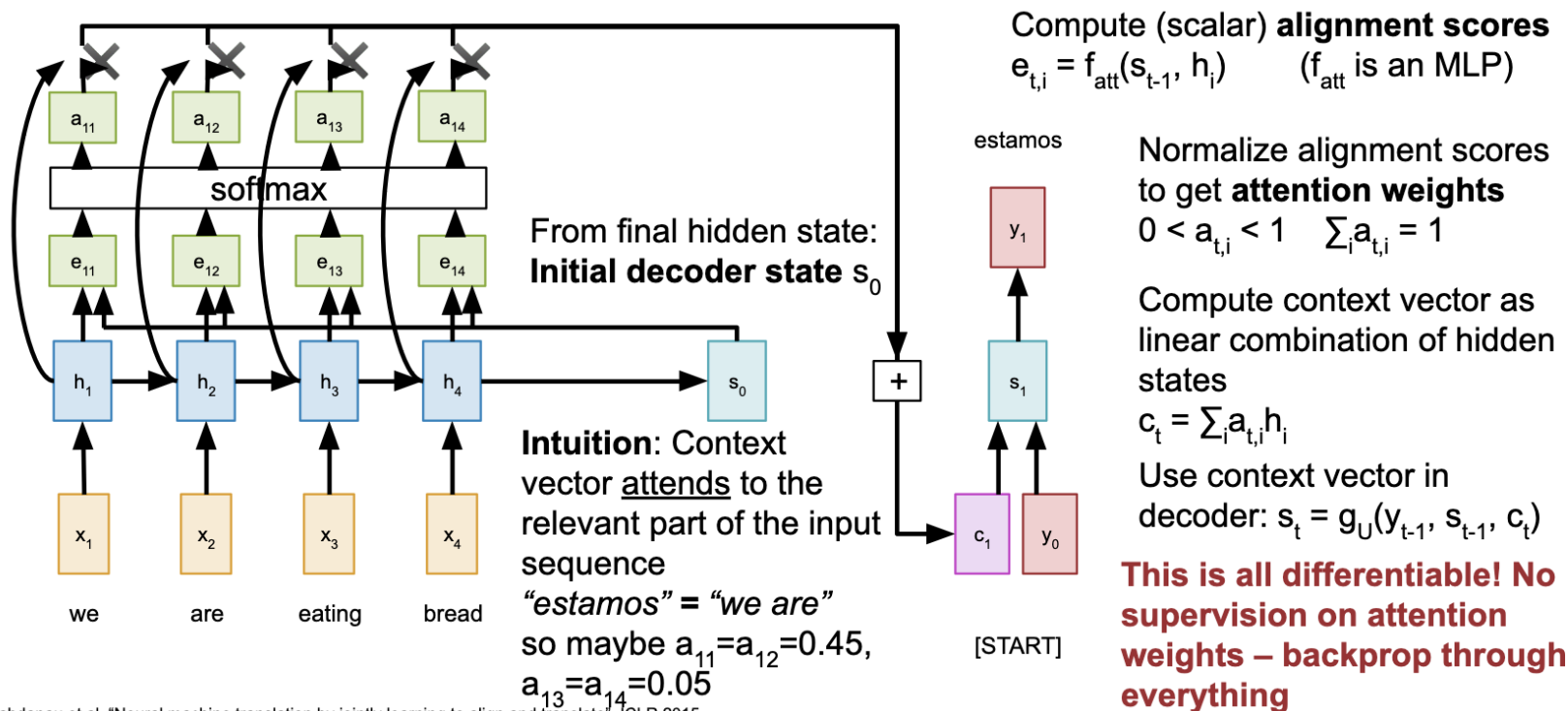so maybe $a_{11}=a_{12}=0.45$, $a_{13}=a_{14}=0.05$

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015
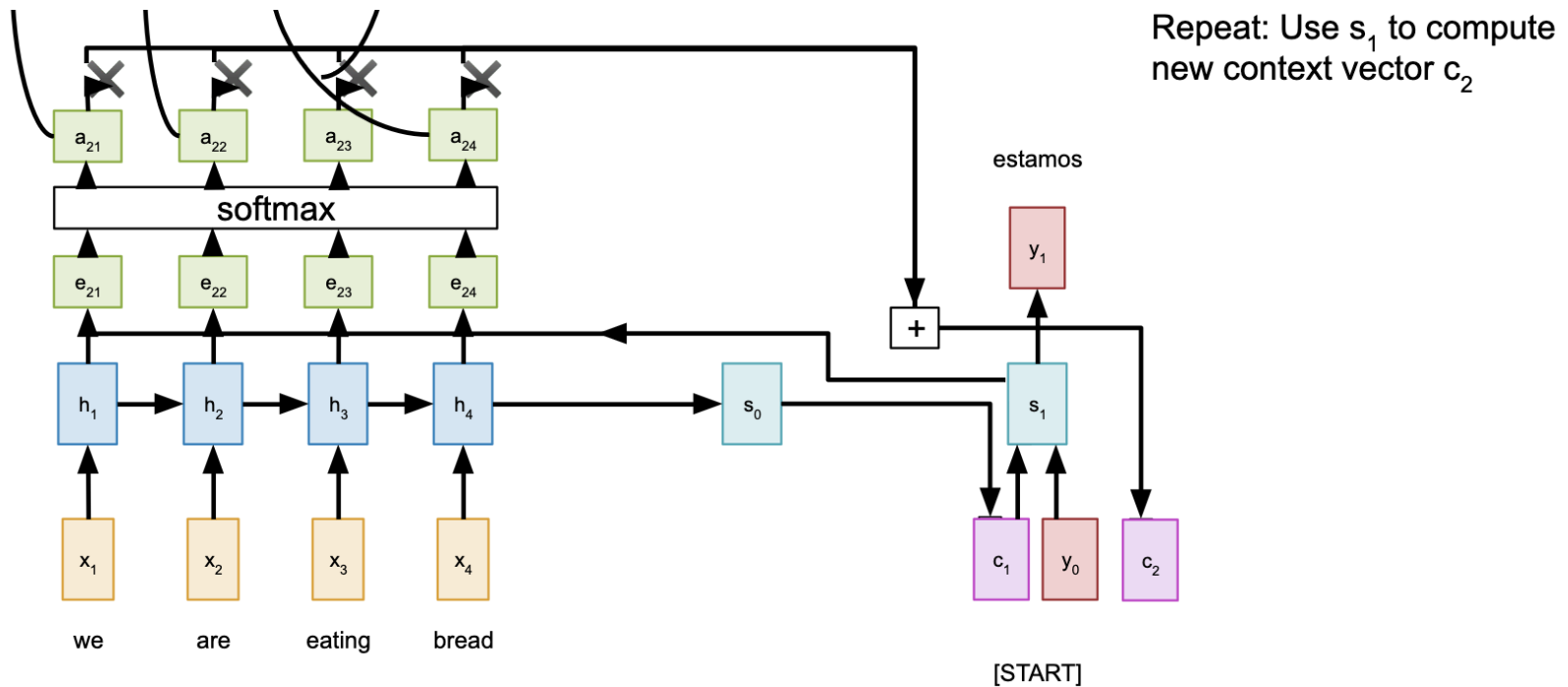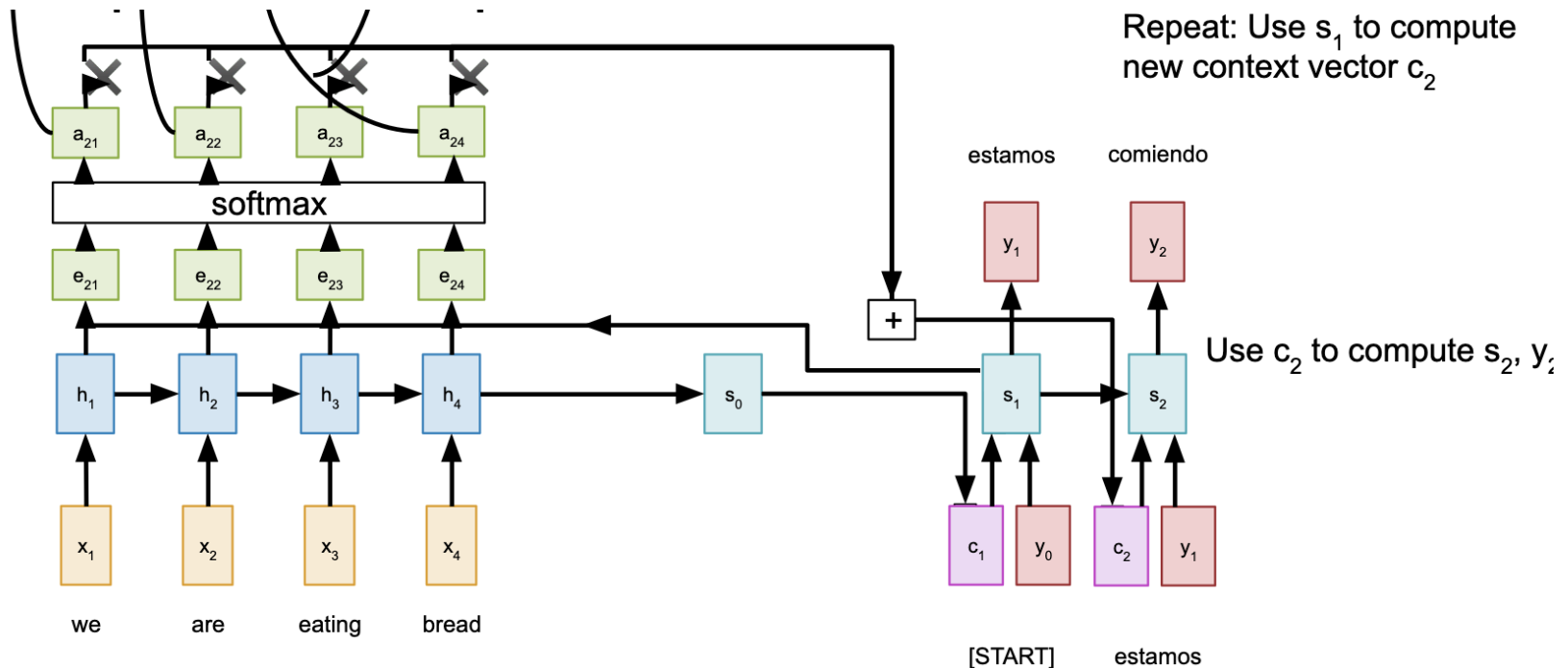
14

# Sequence to Sequence with RNNs and Attention



Compute (scalar) **alignment scores**
$e_{t,i} = f_{att}(s_{t-1}, h_i)$        ($f_{att}$ is an MLP)

Normalize alignment scores to get **attention weights**
$0 < a_{t,i} < 1$    $\sum_i a_{t,i} = 1$

Compute context vector as linear combination of hidden states
$c_t = \sum_i a_{t,i} h_i$

Use context vector in decoder: $s_t = g_U(y_{t-1}, s_{t-1}, c_t)$

**This is all differentiable! No supervision on attention weights – backprop through everything**

From final hidden state:
**Initial decoder state** $s_0$

**Intuition**: Context vector <u>attends</u> to the relevant part of the input sequence
*"estamos"* = *"we are"*
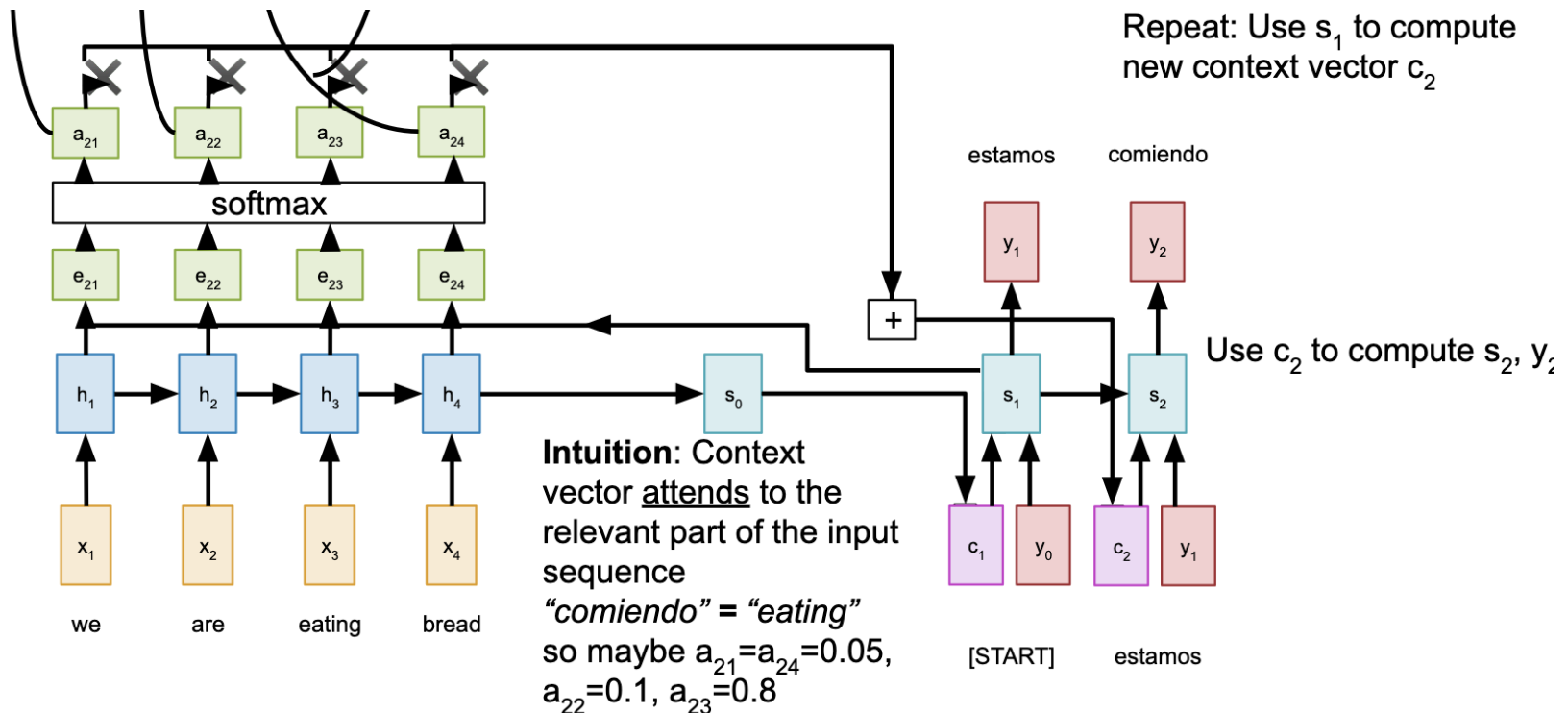so maybe $a_{11}=a_{12}=0.45$, $a_{13}=a_{14}=0.05$

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

15

# Sequence to Sequence with RNNs and Attention



Repeat: Use $s_1$ to compute new context vector $c_2$

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

Repeat: Use $s_1$ to compute new context vector $c_2$

Use $c_2$ to compute $s_2$, $y_2$

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015
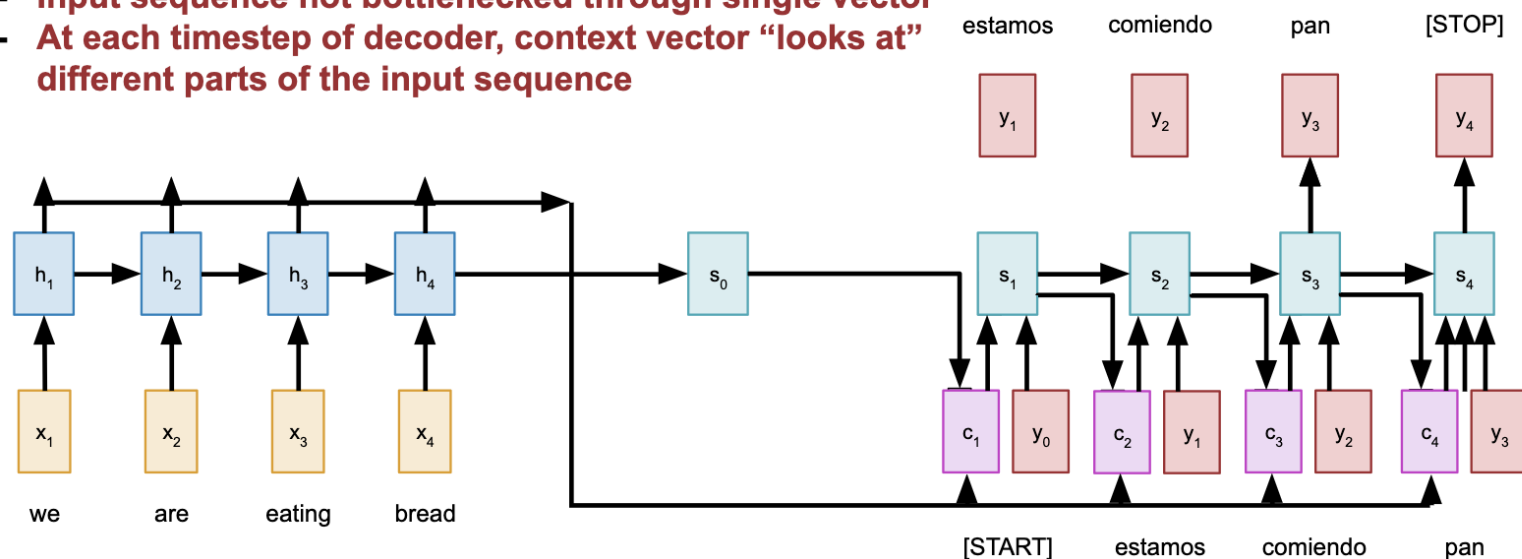
# Sequence to Sequence with RNNs and Attention



Repeat: Use $s_1$ to compute new context vector $c_2$

Use $c_2$ to compute $s_2$, $y_2$

**Intuition**: Context vector <u>attends</u> to the relevant part of the input sequence *"comiendo" = "eating"* so maybe $a_{21}=a_{24}=0.05$, $a_{22}=0.1$, $a_{23}=0.8$

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

18

**Use a different context vector in each timestep of decoder**

- **Input sequence not bottlenecked through single vector**
- **At each timestep of decoder, context vector "looks at" different parts of the input sequence**



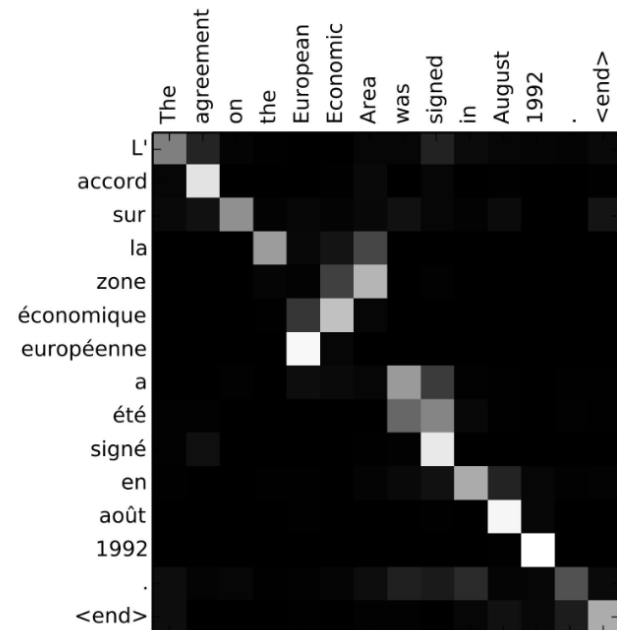Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

19

**Example**: English to French translation

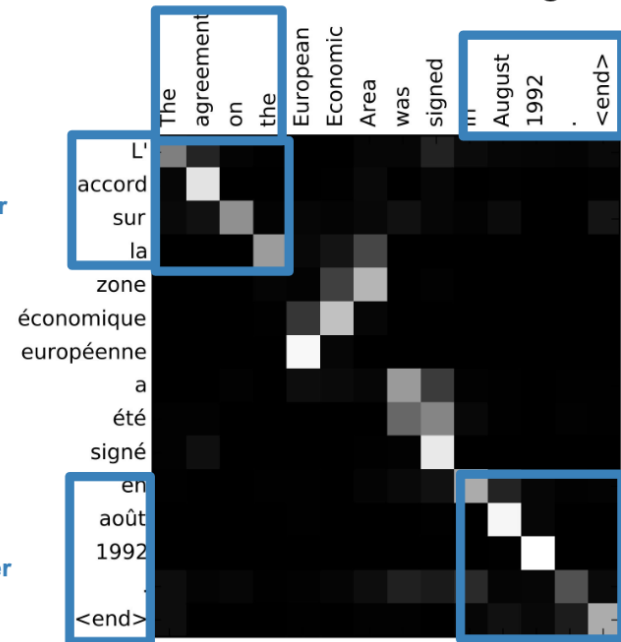**Input**: "The agreement on the European Economic Area was signed in August 1992."

**Output**: "L'accord sur la zone économique européenne a été signé en août 1992."

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

Visualize attention weights $a_{t,i}$

# Sequence to Sequence with RNNs and Attention

**Example**: English to French translation

**Input**: "**The agreement on the** European Economic Area was signed **in August 1992**."

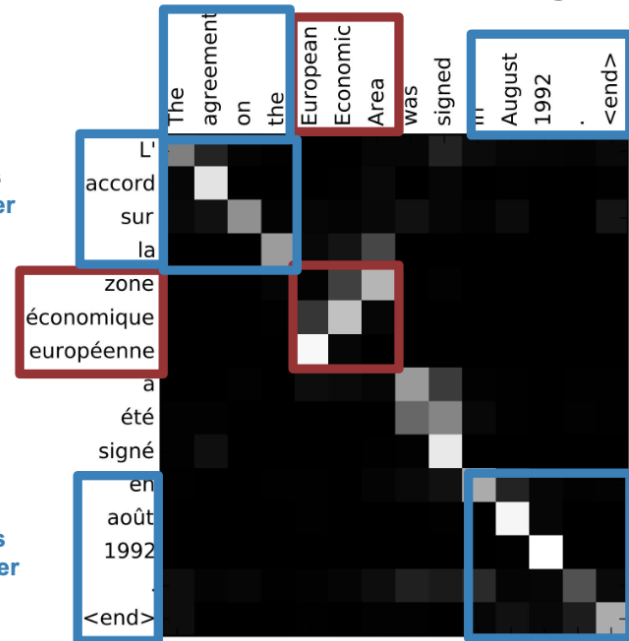**Output**: "**L'accord sur la** zone économique européenne a été signé **en août 1992**."

**Diagonal attention means words correspond in order**

**Diagonal attention means words correspond in order**

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

Visualize attention weights $a_{t,i}$

**Example**: English to French translation

**Input**: "**The agreement on the European Economic Area** was signed **in August 1992**."

**Output**: "**L'accord sur la zone économique européenne** a été signé **en août 1992**."

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

## Visualize attention weights $a_{t,i}$



**Diagonal attention means words correspond in order**

**Attention figures out different word orders**

**Diagonal attention means words correspond in order**

# Questions?

University of North Carolina CHARLOTTE