

# **CIS 6930 Topics in Computing for Data Science**

## **Week 6a: Recurrent Neural Networks (2)**

### **Encoder-Decoder Models**

9/30/2021 & 10/5/2021  
Yoshihiko (Yoshi) Suhara

2pm-3:20pm

# This Week & Next Week

- Thu 10/7 Midterm exam (written exam **on campus**)
  - Time: 2pm-3:20pm
  - Location: LBR 252
- Fall Break

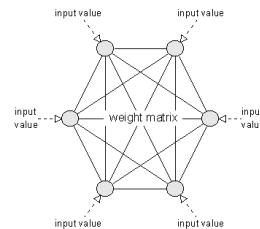
# Week 6!

- ~~Week 1: Deep Learning Basics (Thu 9/9)~~
- ~~Week 2: AutoEncoder (Tue 9/14)~~
- ~~Week 3: Convolutional Neural Networks (Thu 9/16)~~
- ~~Week 4: GAN (Tue 9/21)~~
- ~~Week 5: Word embeddings: Word2vec, GloVe (Thu 9/23)~~
- **Week 6: Recurrent Neural Networks (Tue 9/28, Thu 9/30)**
- Week 7: Review/Project pitch & Mid-term (Tue 10/5, Thu 10/7)
- Fall Break
- ...

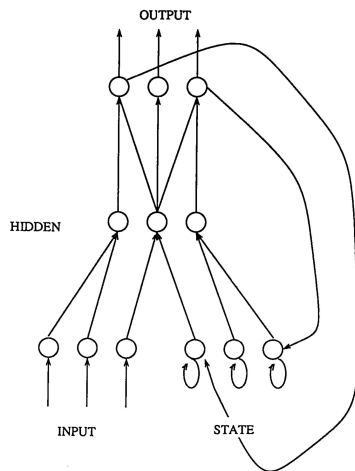
# Recap: Recurrent Neural Networks (1)

# Recurrent Neural Networks for Sequential Input

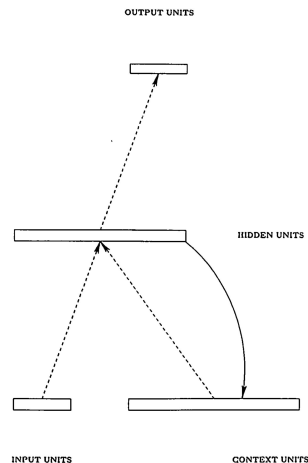
- Inspired by Hopfield Network (1982)
- “Simple” Recurrent Neural Networks
  - Jordan Net (1986)
  - **Elman Net (1990)**



**Hopfield Network**



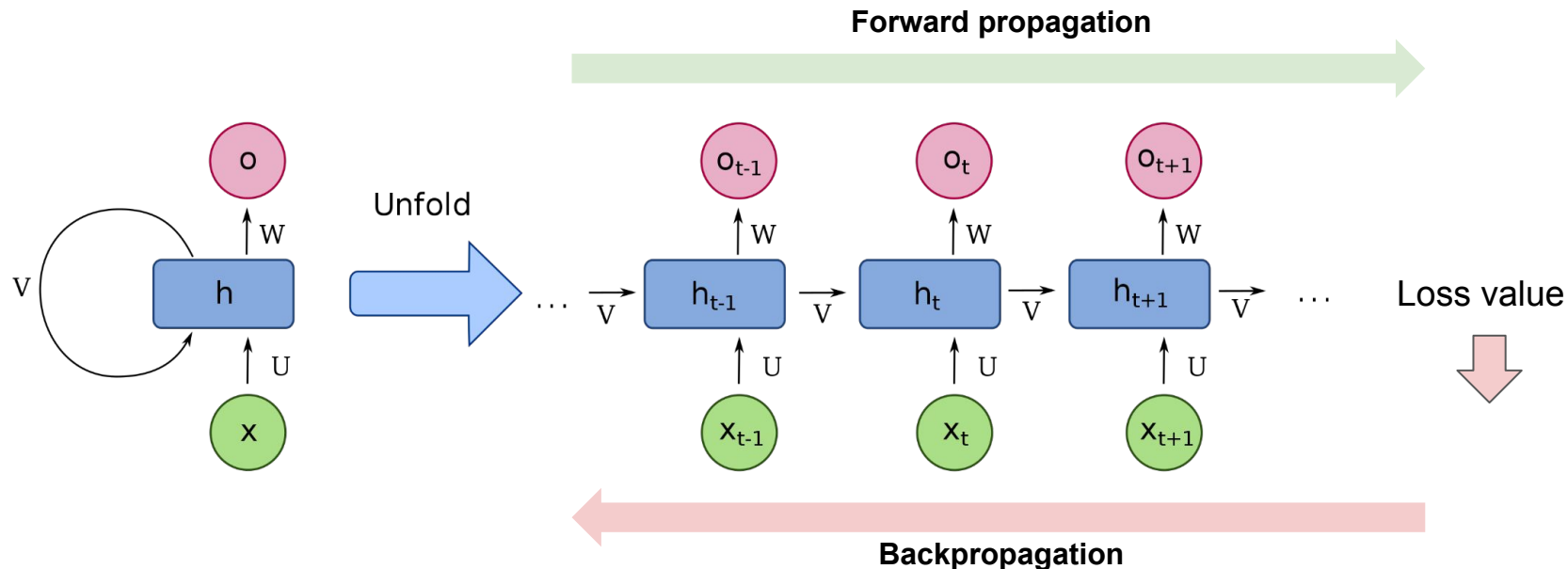
**Jordan Net**



**Elman Net**

# Unfolding RNN

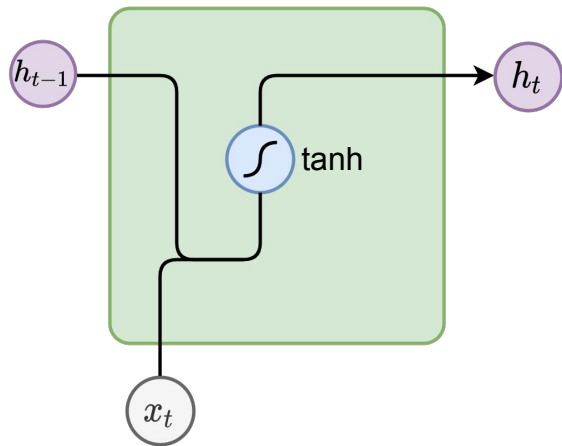
## Backpropagation Through Time (BPTT)



# RNN Cell for Sequential Data

RNN = Elman Net

- Input value + previous hidden state → Next hidden state



$$h_t = \tanh(W_{ih}x_t + b_{ih} + \underline{W_{hh}}h_{(t-1)} + b_{hh})$$

$$W_h \begin{bmatrix} x_t; h_{(t-1)} \end{bmatrix} + b_h$$

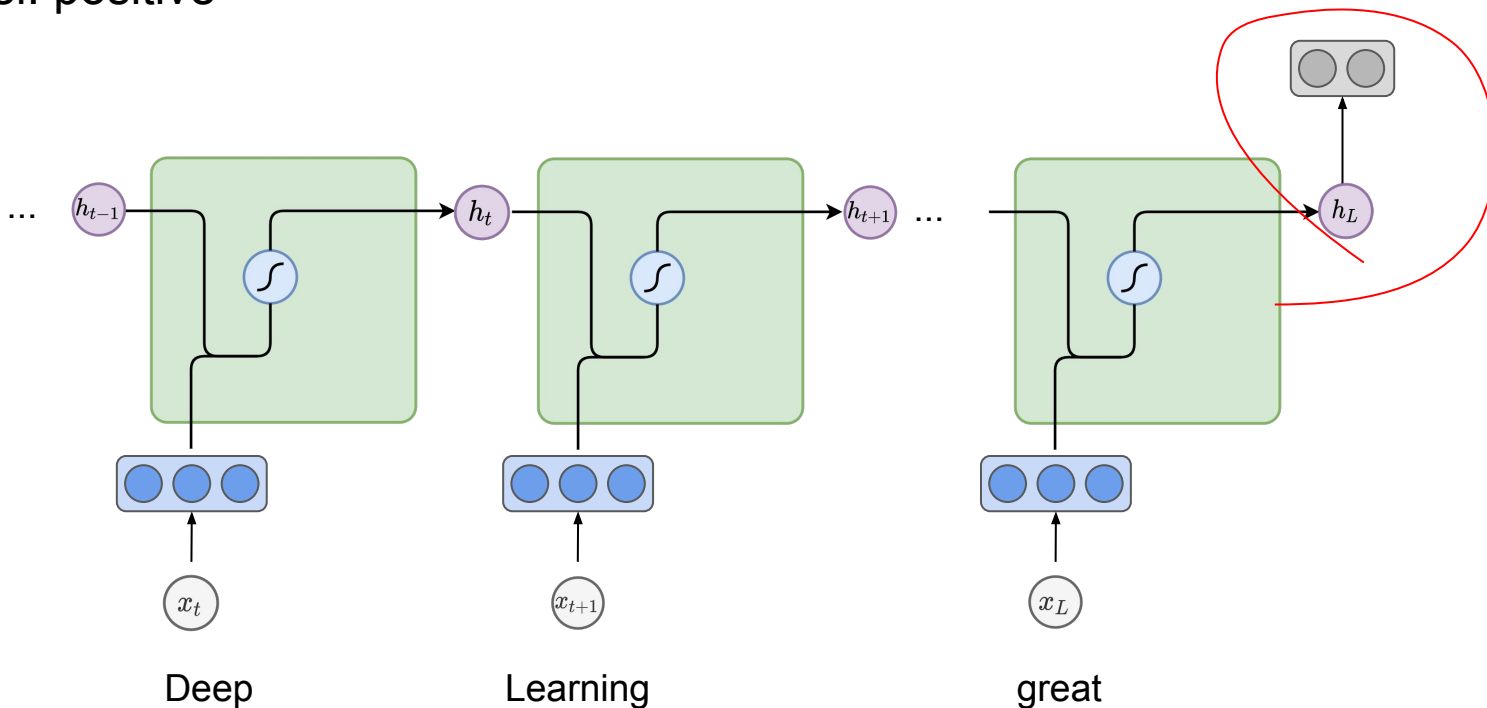
Another notation

The figures are based on the following blog article (highly recommended!)

[Illustrated Guide to LSTM's and GRU's: A step by step explanation](#)

# Text Classification Example

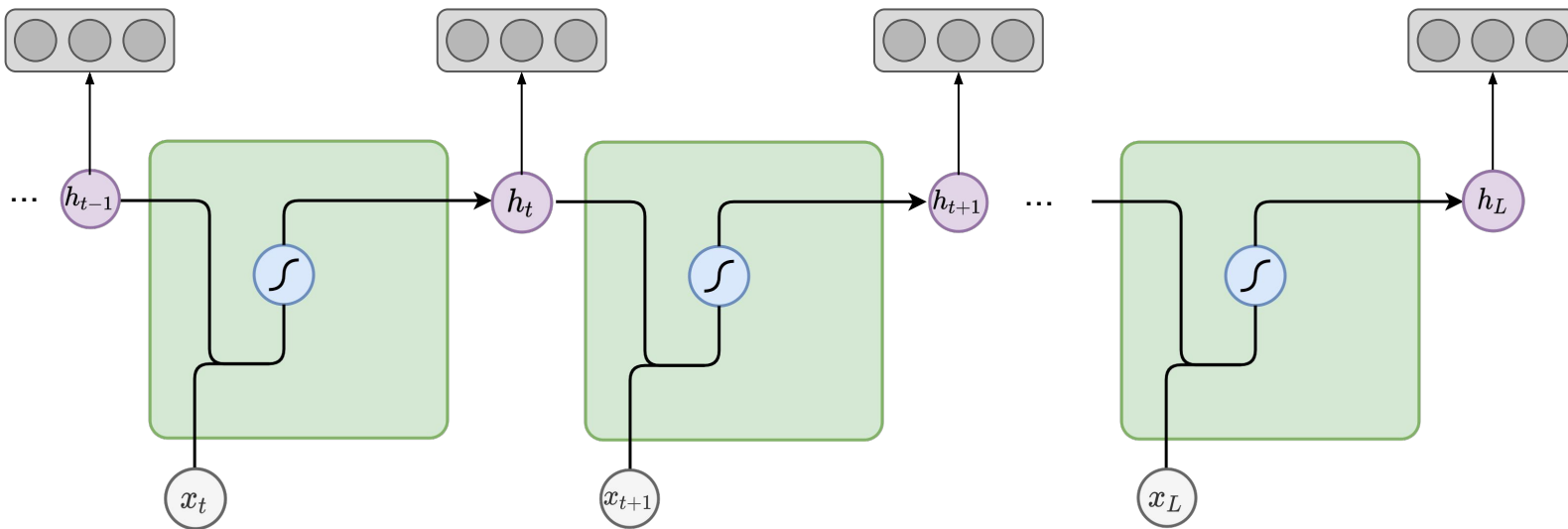
- Input: “Deep Learning is great”
- Label: positive





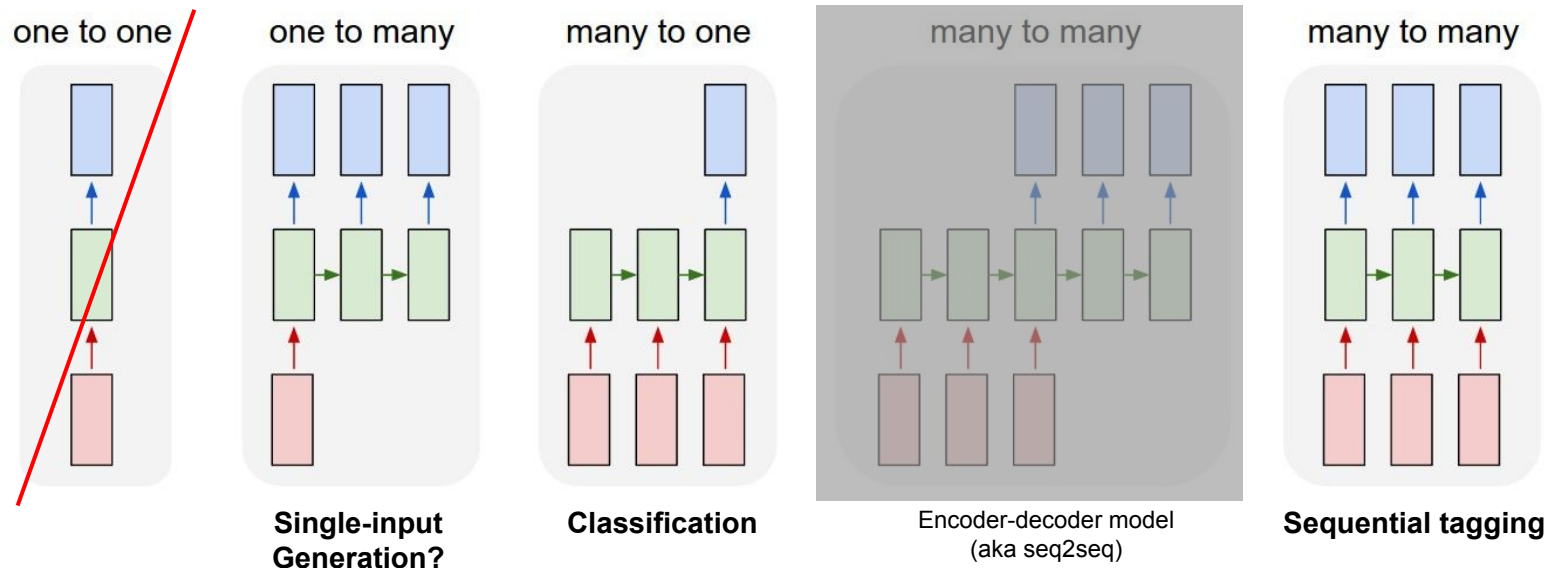
# Where to Attach Output Layer?

- Where?



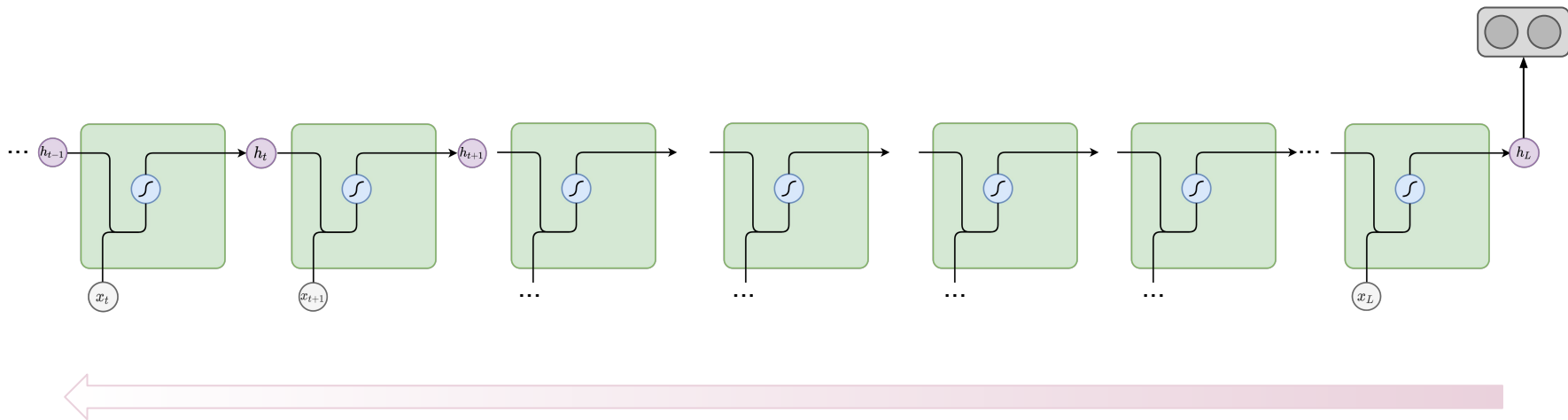
Any other options?  
→ RNN for **sequential tagging**

# RNN Application Patterns

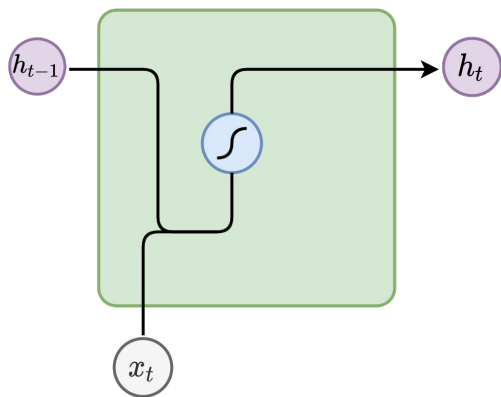


# What's the issue with RNNs?

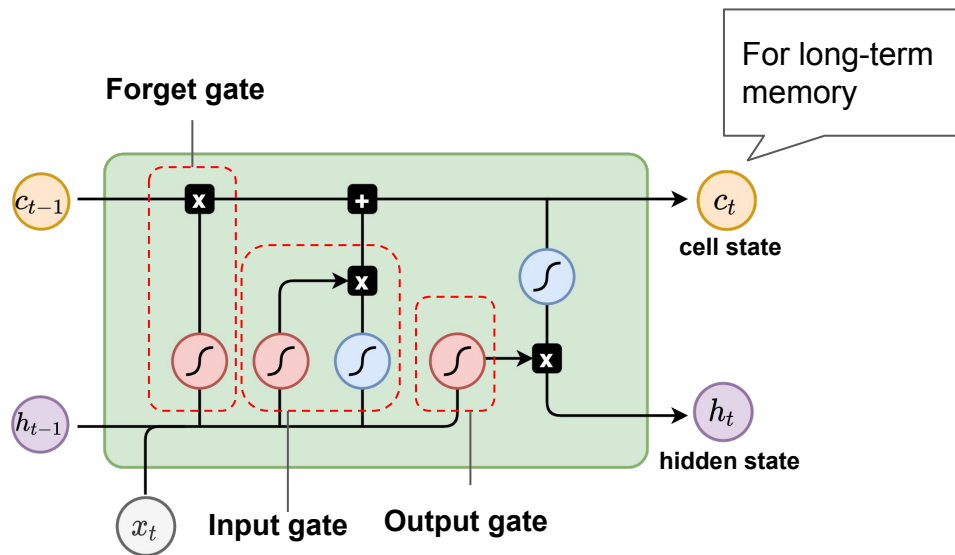
- Not a very good design for a **long sequence**
  - **Later input values** have higher impact on the last hidden state
  - Gradient vanishing problem



# LSTM: Cell State + Gating Mechanism



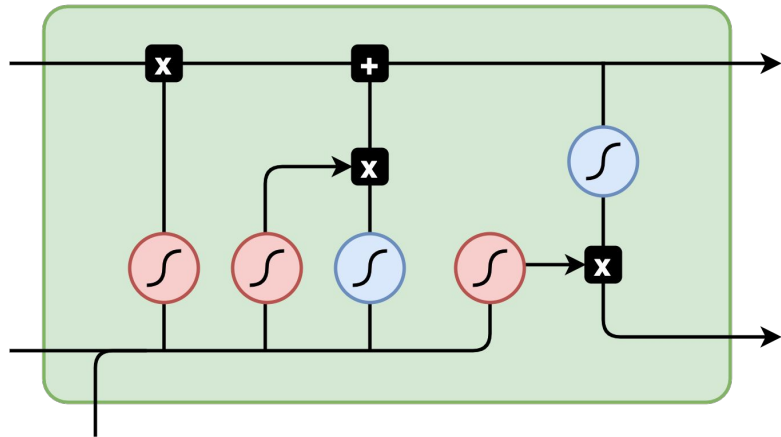
RNN Cell



LSTM Cell

# Long-short Term Memory (LSTM) Cell

[Hochreiter and Schmidhuber 1997]



$$\begin{aligned}
 i_t &= \sigma(W_{ii}x_t + b_{ii} + W_{hi}h_{t-1} + b_{hi}) \\
 f_t &= \sigma(W_{if}x_t + b_{if} + W_{hf}h_{t-1} + b_{hf}) \\
 g_t &= \tanh(W_{ig}x_t + b_{ig} + W_{hg}h_{t-1} + b_{hg}) \\
 o_t &= \sigma(W_{io}x_t + b_{io} + W_{ho}h_{t-1} + b_{ho}) \\
 c_t &= f_t \odot c_{t-1} + i_t \odot g_t \\
 h_t &= o_t \odot \tanh(c_t)
 \end{aligned}$$

# Agenda (1st Half)

- Advanced topics (contd.)
  - GRU
  - Stacked GRU/LSTM
  - Bidirectional GRU/LSTM
- Hands-on Session
- Assignment 3

# Link to the previous slide deck

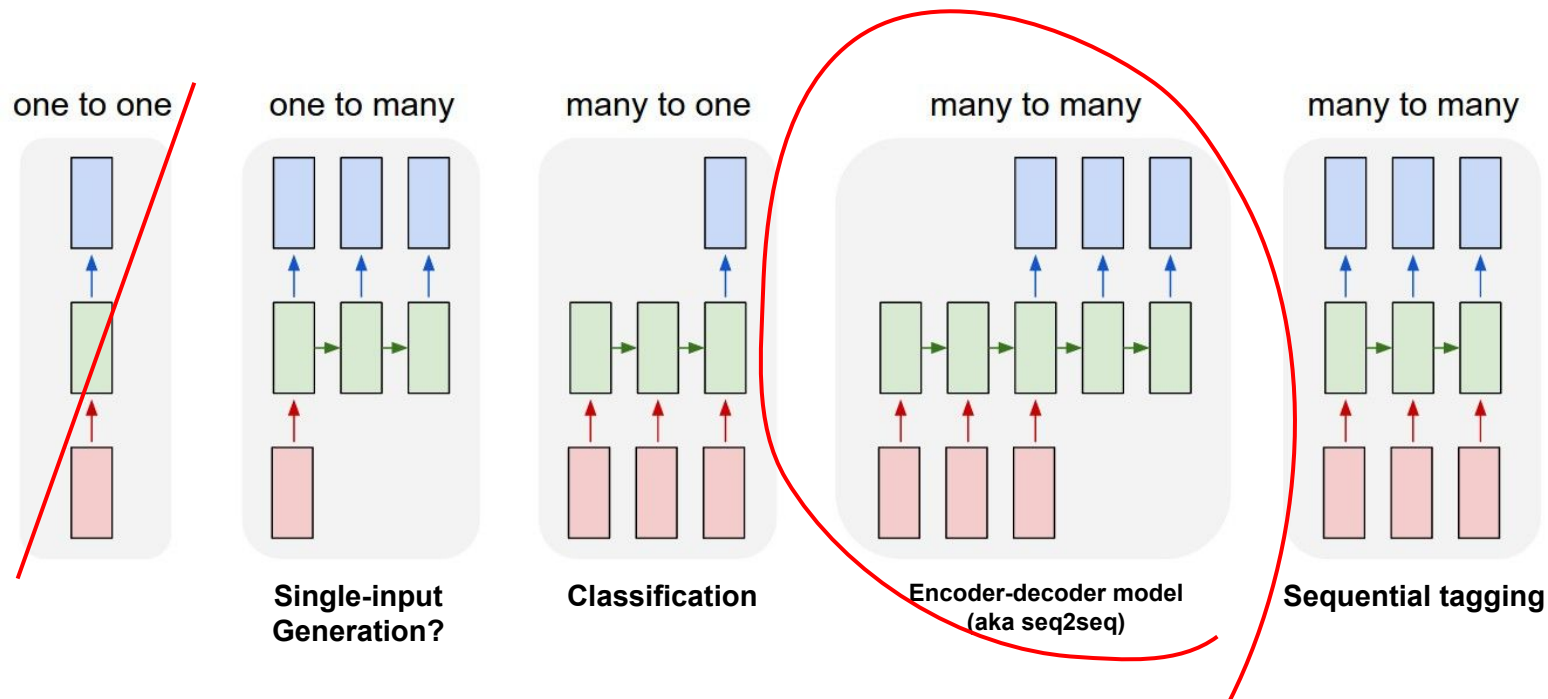
- <https://github.com/suhara/cis6930-fall2021/blob/main/slides/cis6930-week6a-recurrent-neural-networks-1.pdf>

# Agenda (2nd Half)

- Encoder-Decoder Models (aka seq2seq Models)
- Decoding Algorithms
- Attention mechanism

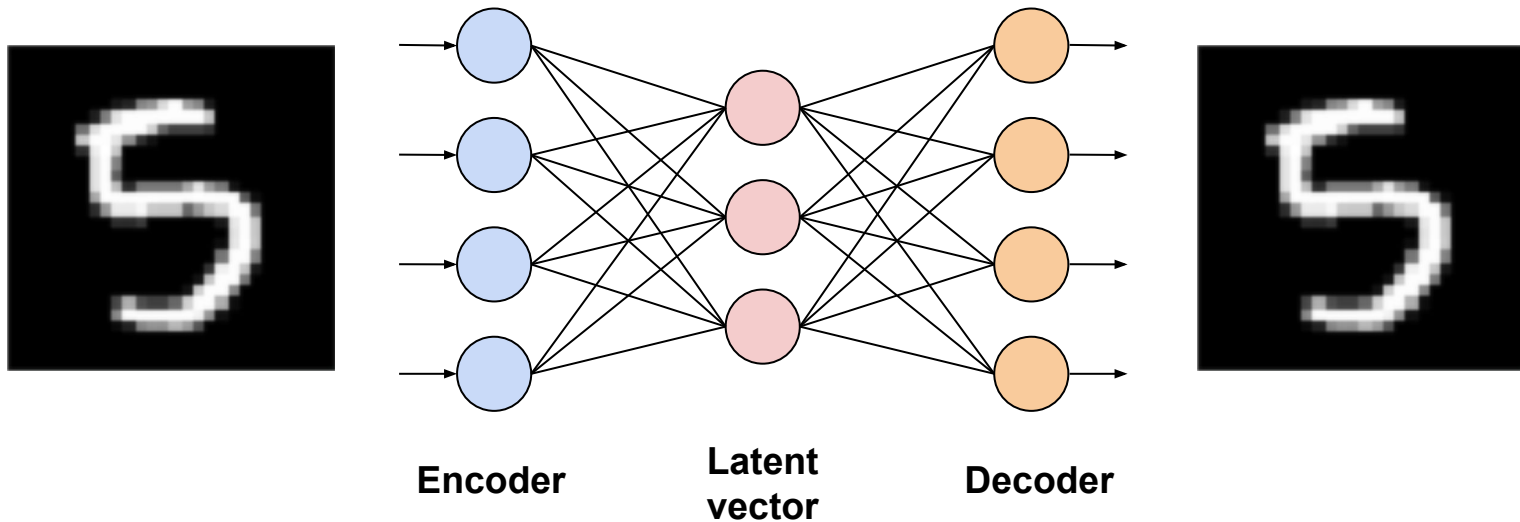


# RNN Application Patterns



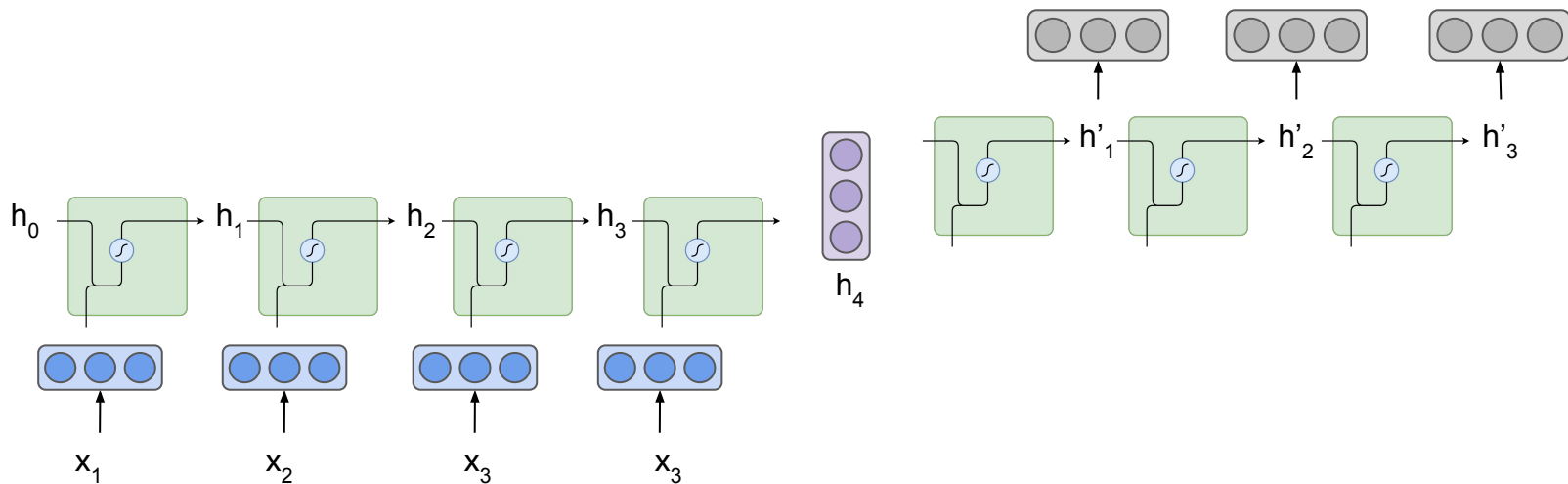
# Autoencoders

- **Encoder-decoder** models that learn to reconstruct the original data



# Encoder-Decoder Model aka Sequence-to-sequence (seq2seq) Model

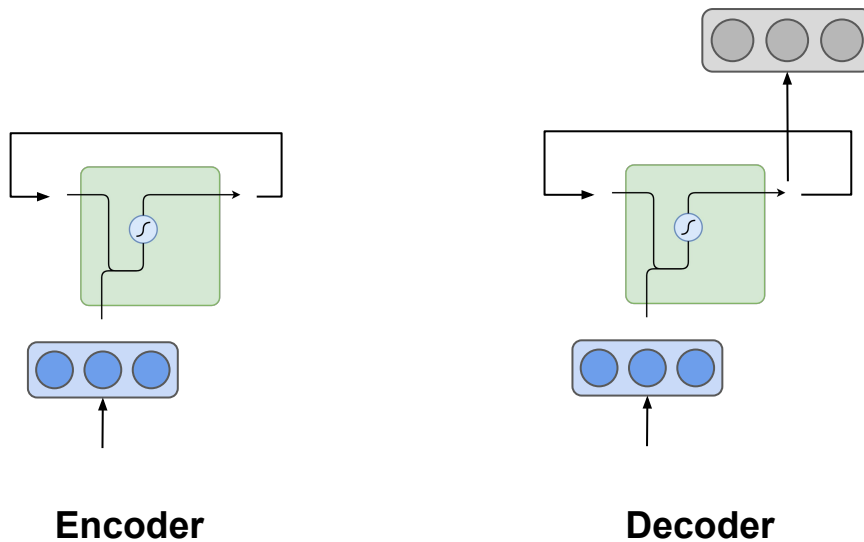
- Input: Variable length sequence
- Output: Variable length sequence



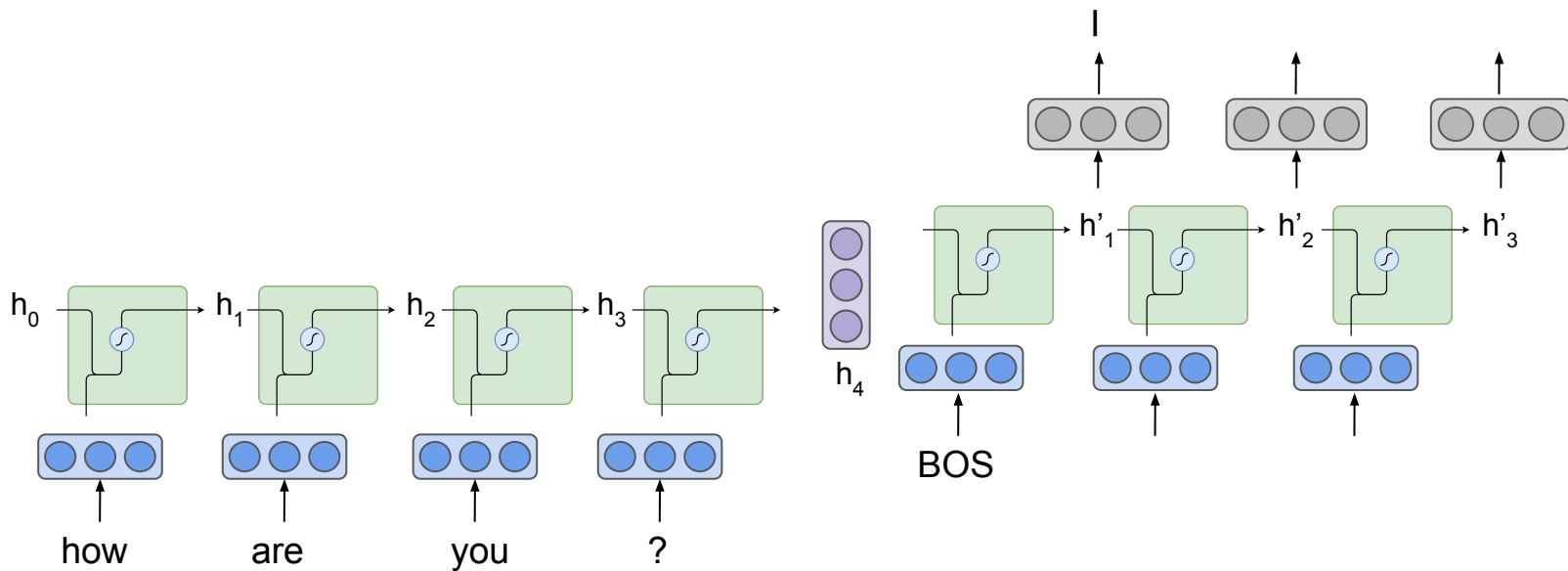
(\* RNN Cell can be LSTM/GRU Cell)

# Encoder Model + Decoder Model

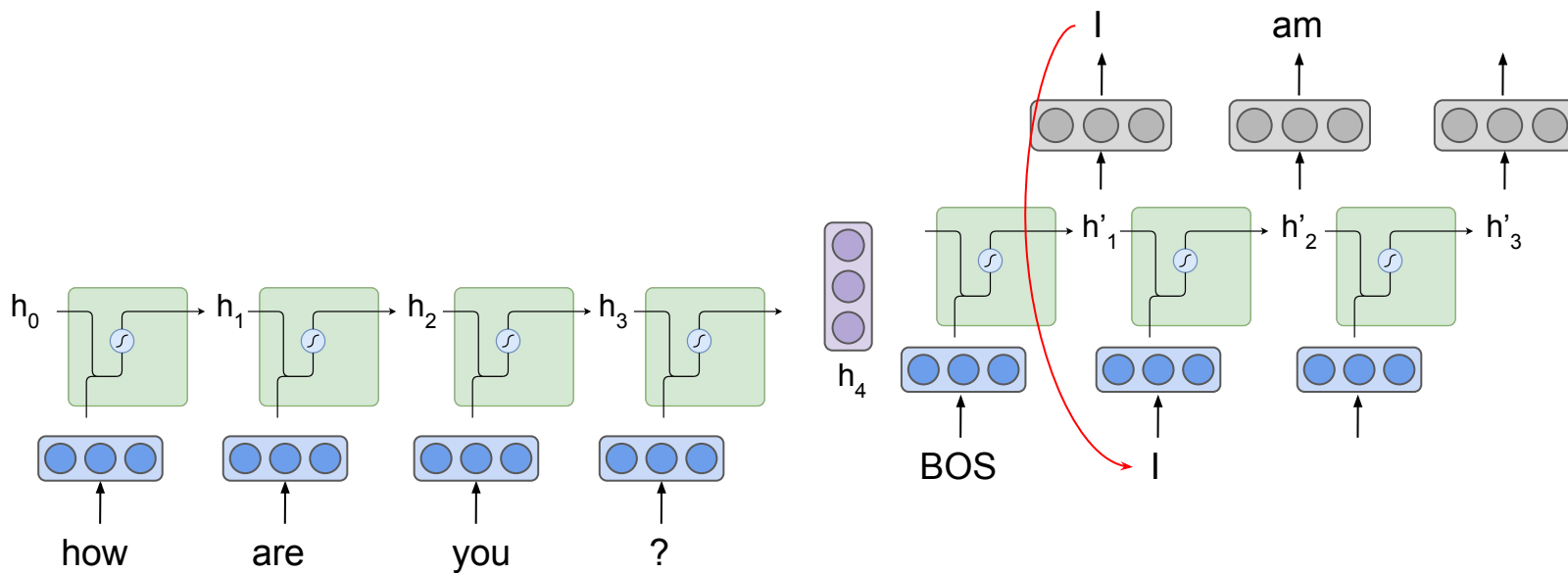
- Two different RNN models



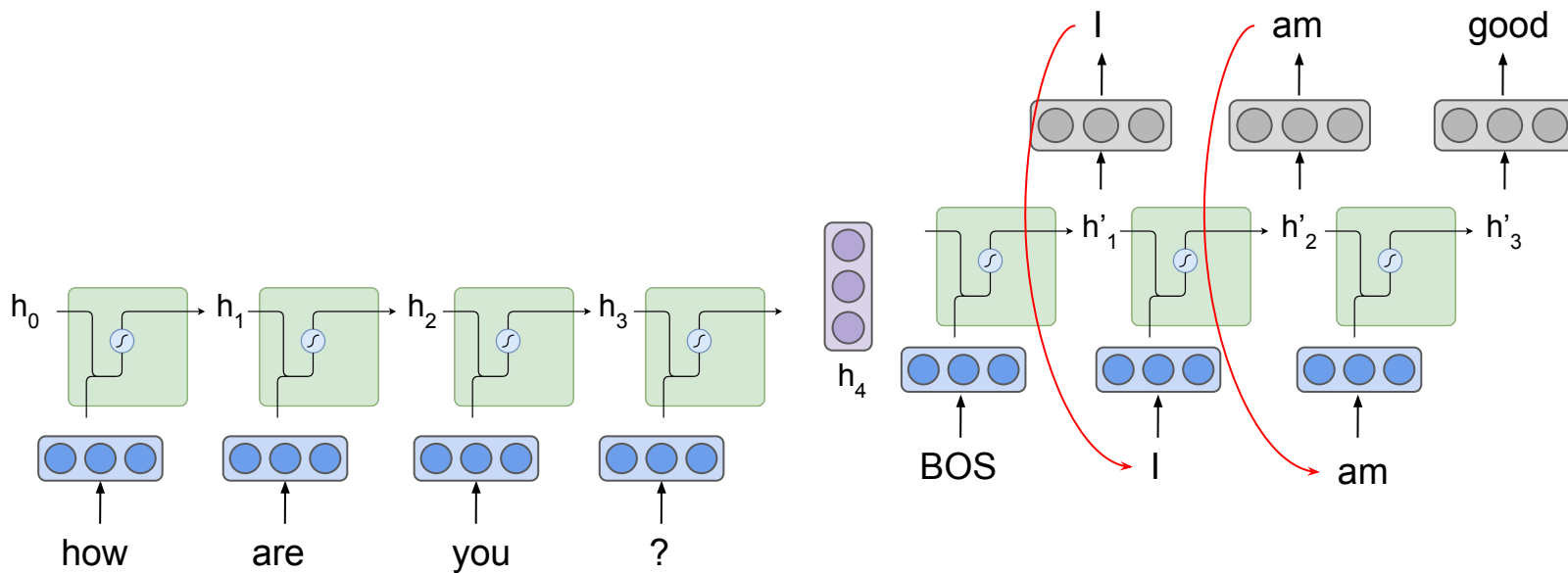
# Example 1: Response Generation



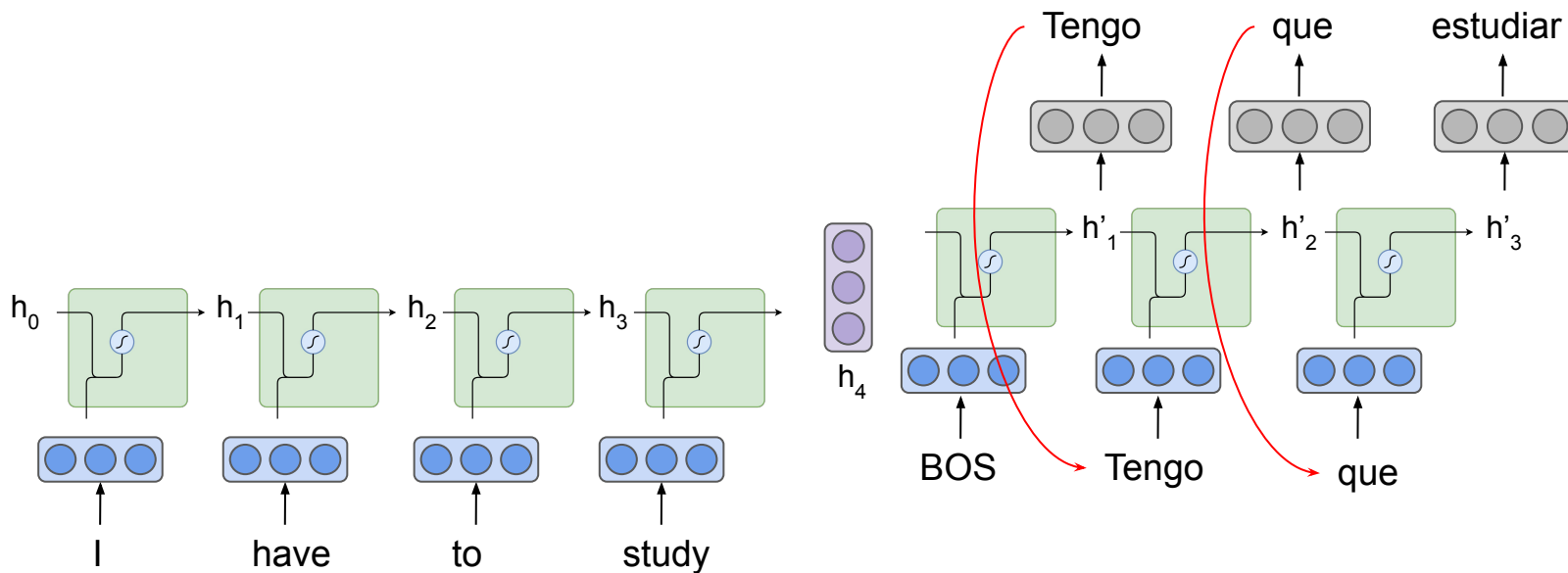
# Example 1: Response Generation



# Example 1: Response Generation

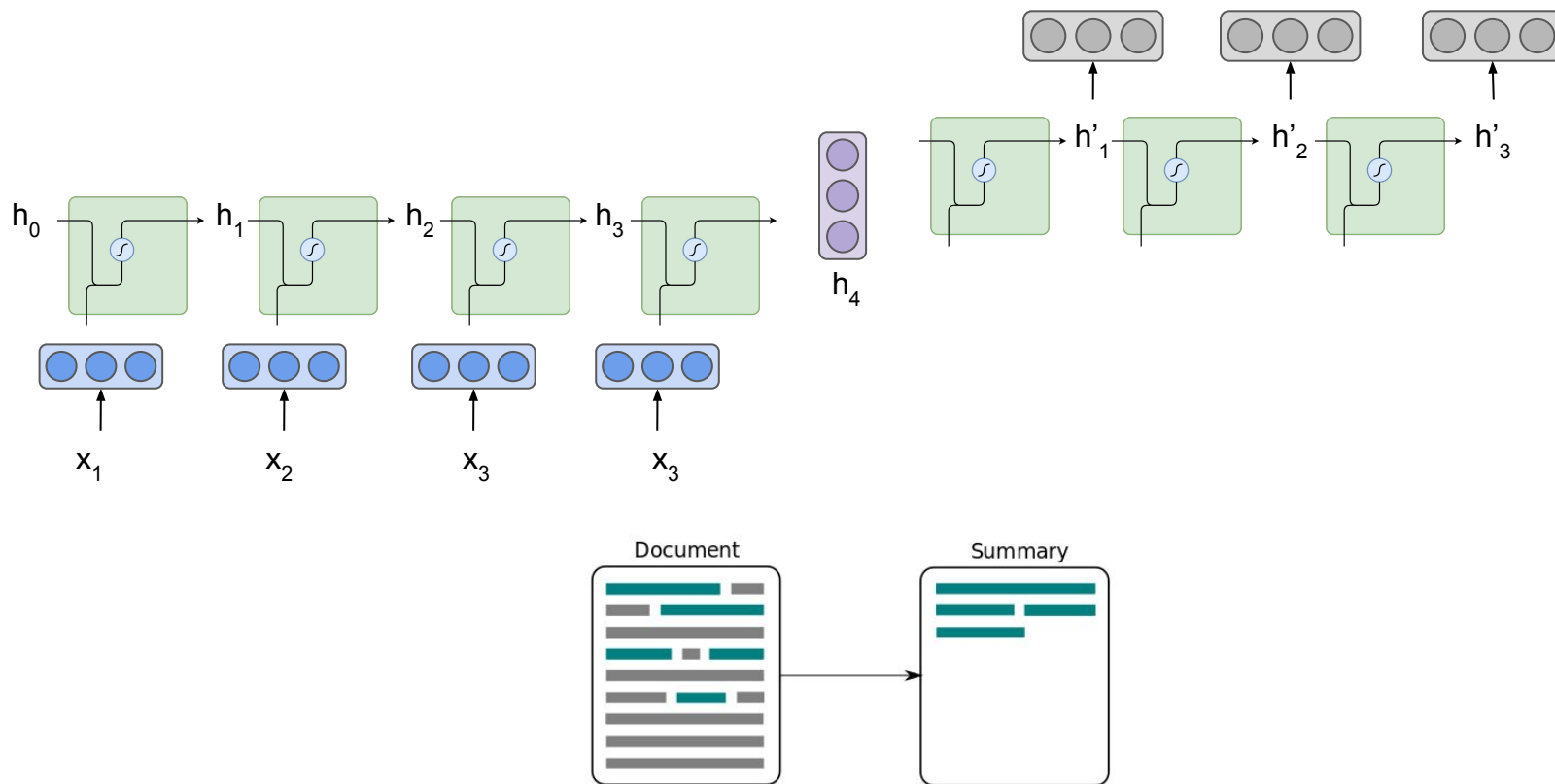


## Example 2: Machine Translation



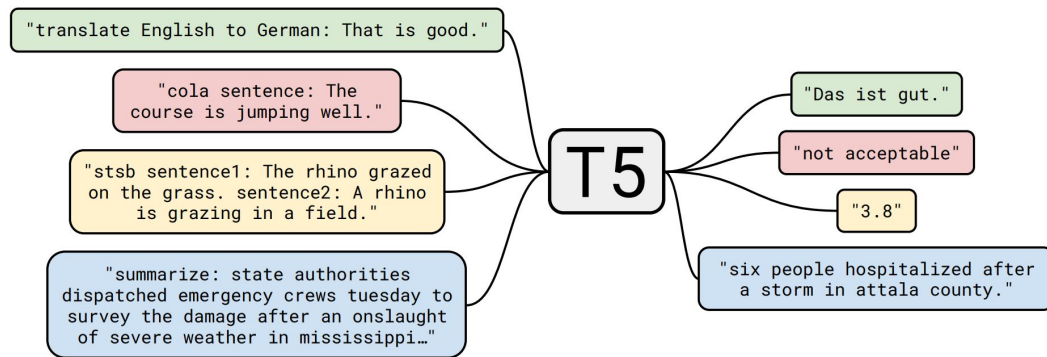


# Example 3: Text Summarization



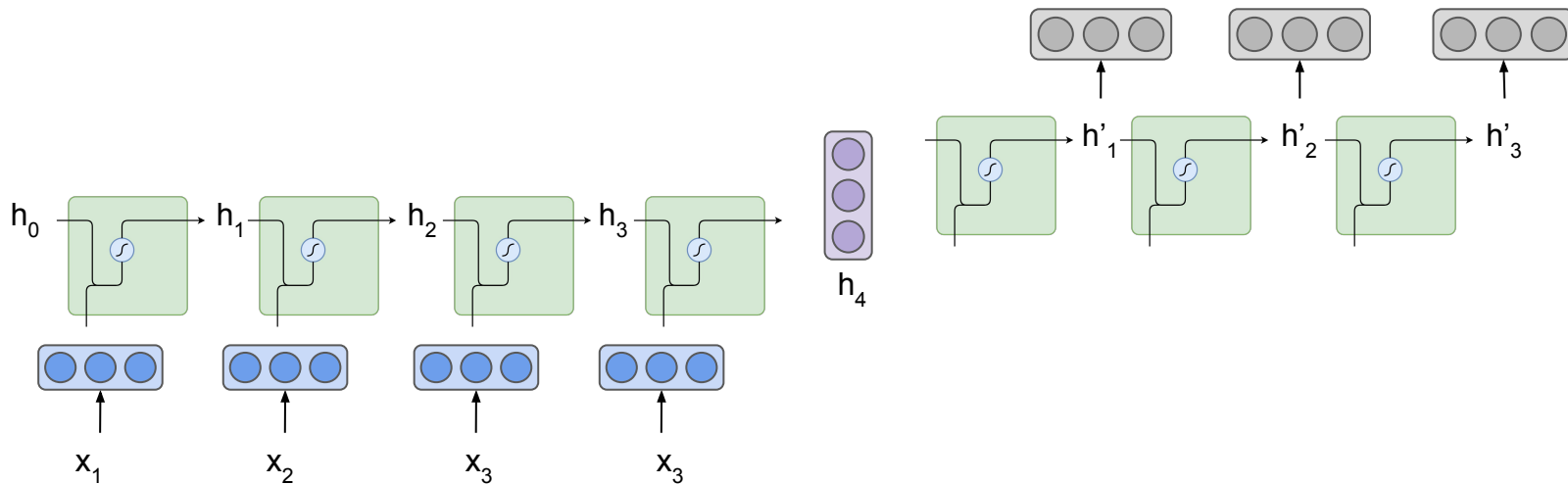
# What Can We Do with Encoder-Decoder Models?

- Almost anything! (even classification!)
- .... as long as you have a sufficient amount of **parallel data** (input-output pairs)



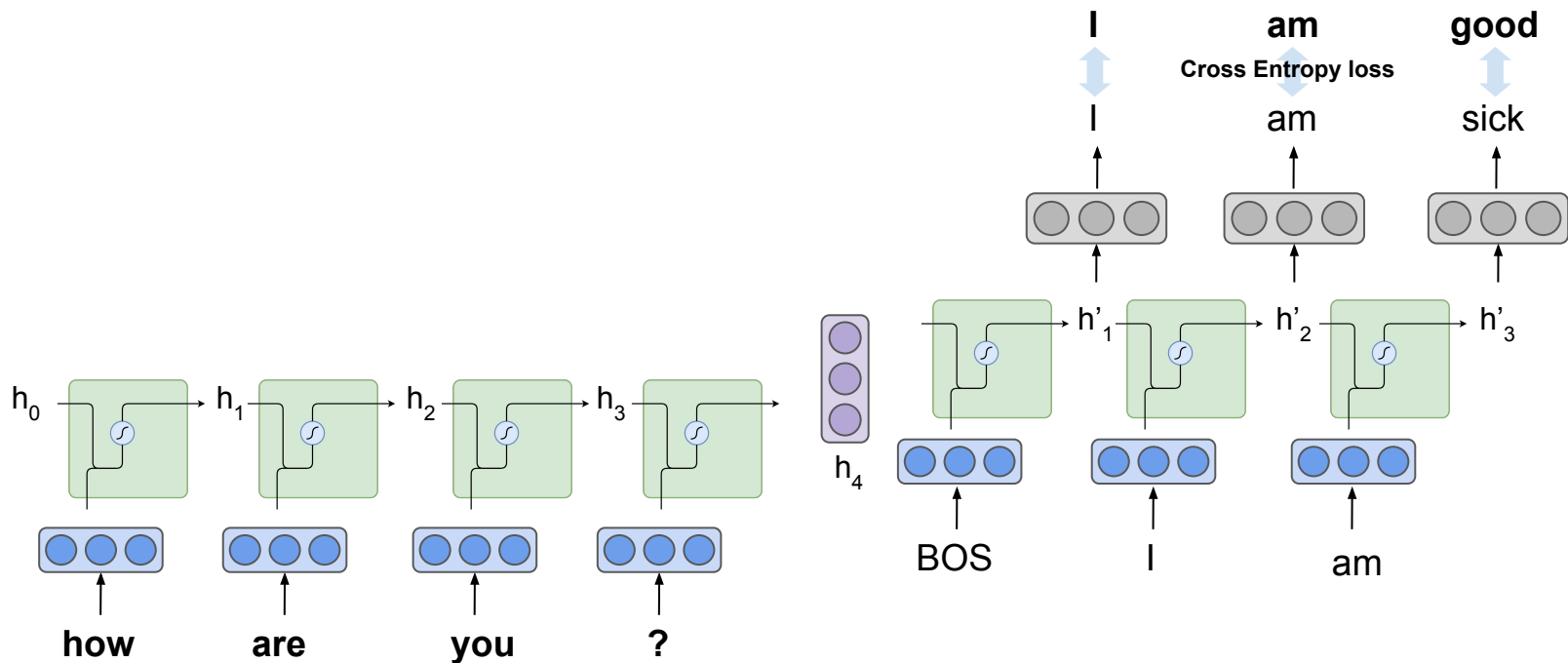
# Check: Input & Output

- **Token ID**  $\rightarrow$  One-hot vector  $\rightarrow$  Dense vector  $\rightarrow$  **RNN**
- **RNN**  $\rightarrow$  Output layer + softmax  $\rightarrow$  **Token ID**



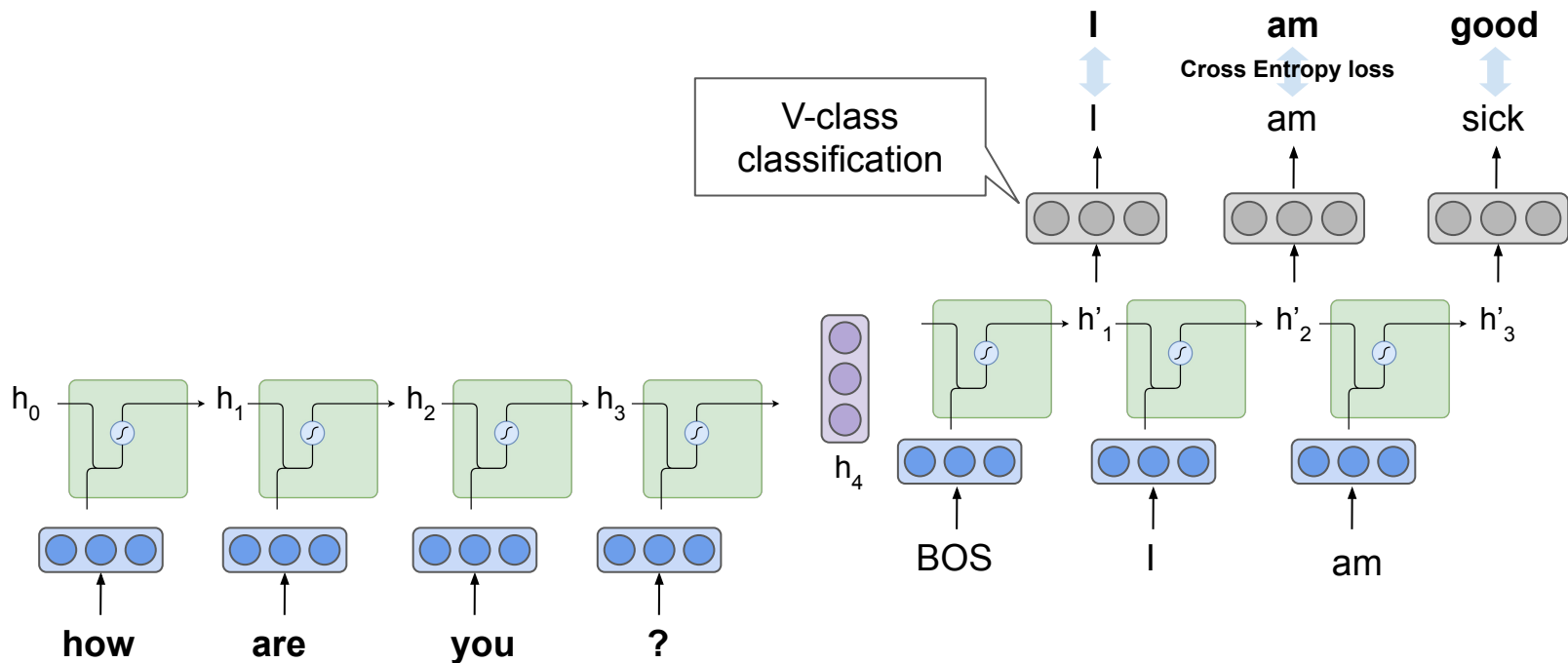
# Training Encoder-Decoder Models: Training Data

- Training data: **Pairs of text sequences** (called a **parallel corpus**)
  - e.g., (“How are you?”, “I am good”)



# Training Encoder-Decoder Models: “Teacher Forcing”

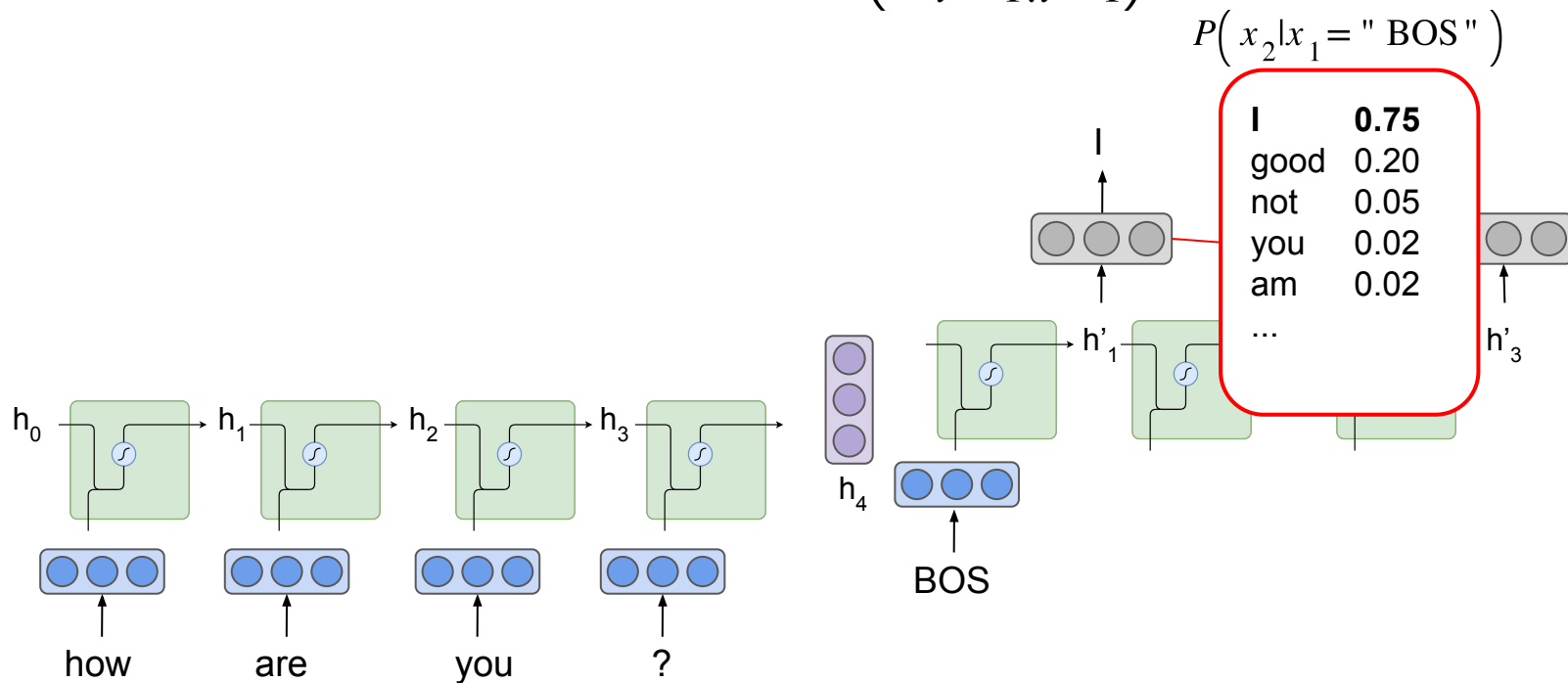
- Train the model to output **ground-truth outputs**



# Decoding Algorithms

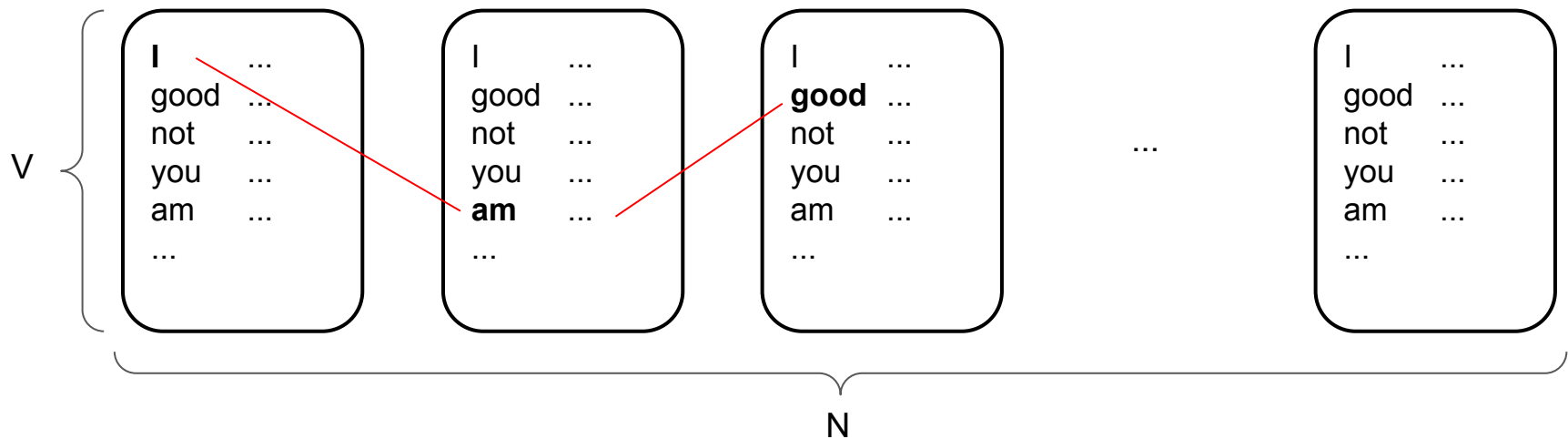
# What does the output layer do at each step?

- Token probability given the context  $P(x_t | x_{1:t-1})$



# Exhaustive Search is Intractable!

- We want to find the best sequence  $\operatorname{argmax}_X P(X) = \operatorname{argmax}_X \prod_{t=1}^N P(x_t | x_{1:t-1})$
- Number of candidates:  $V^N$ 
  - for a sequence length of  $N$  and a vocabulary size of  $V$



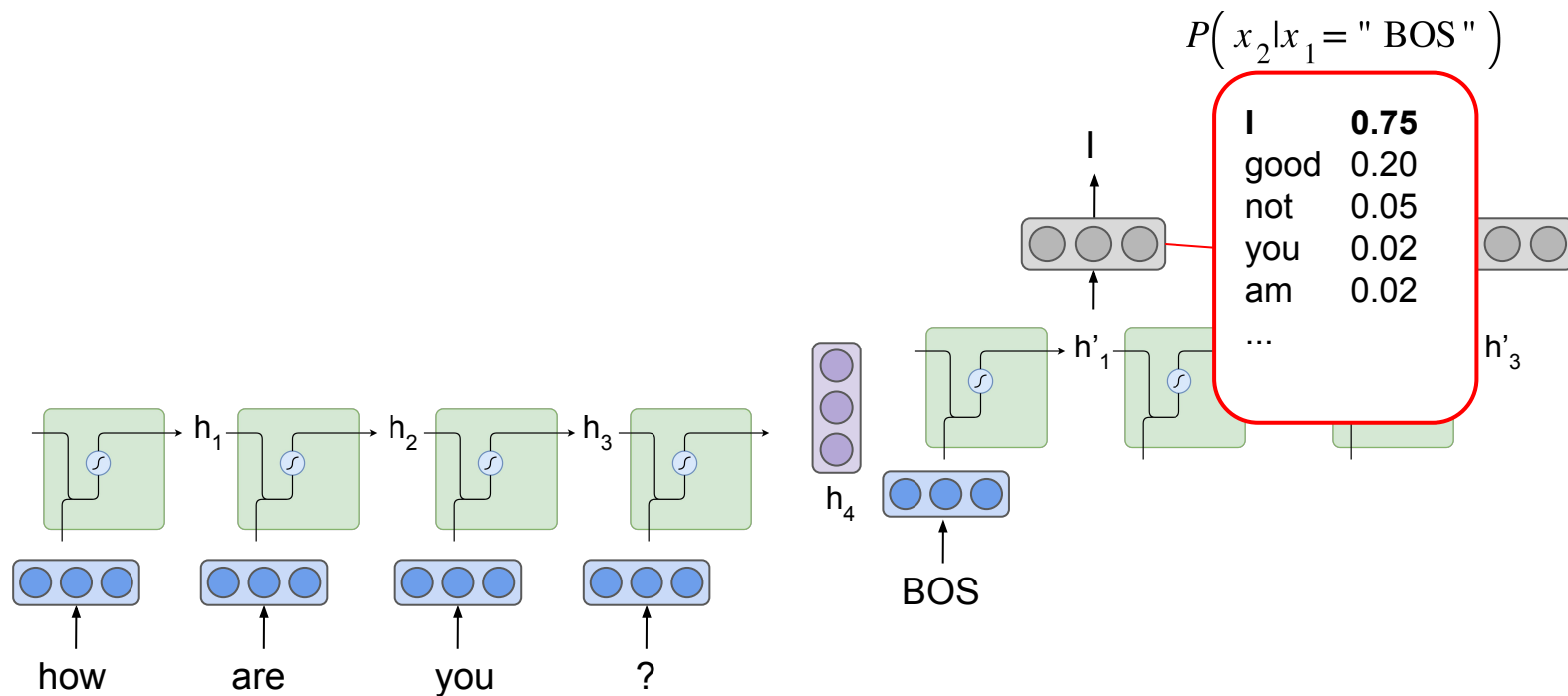


# Decoding Algorithm (i.e., Approximate Search)

- Greedy Search
- Beam Search
- Top-k Sampling (covered in Week 8-9)
- Nucleus Sampling (covered in Week 8-9)

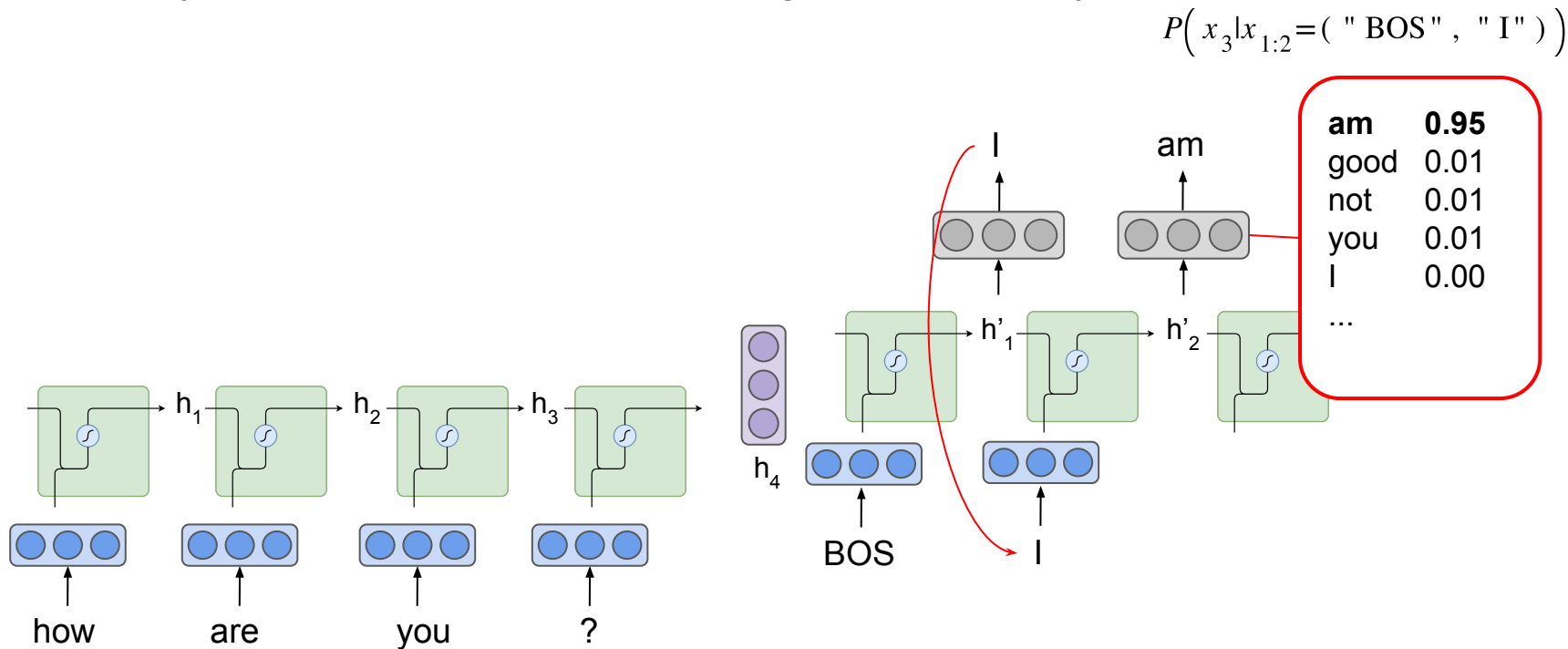
# Greedy Search

- Always choose the token with the highest probability



# Greedy Search

- Always choose the token with the highest probability

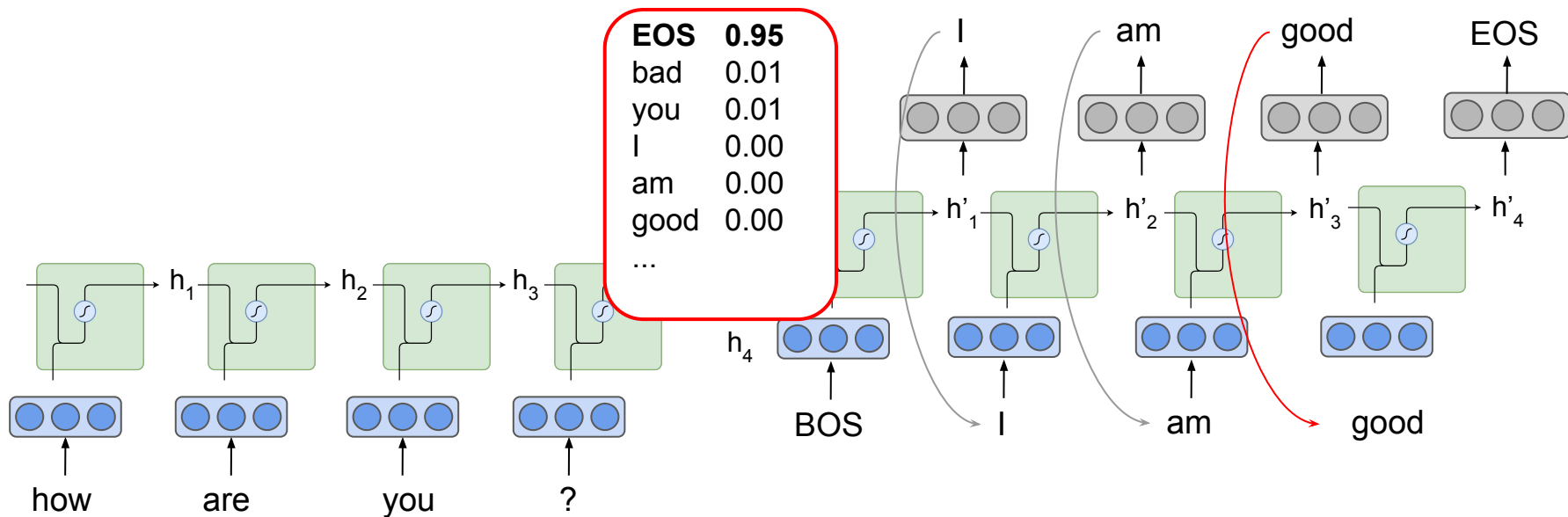




# When Does the Decoder Stop?

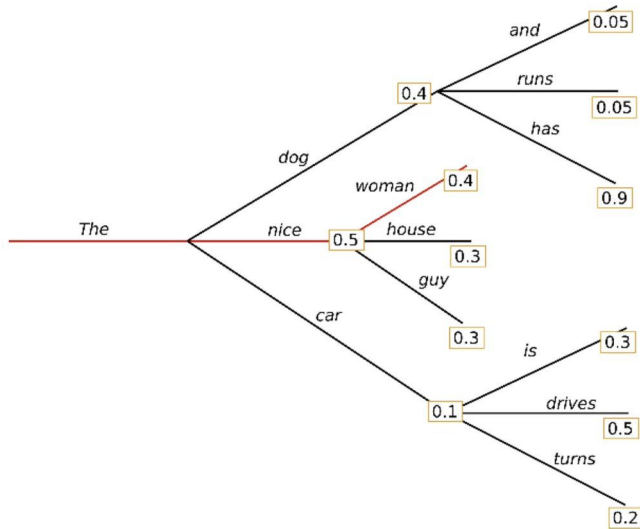
- When it selects the dummy symbol for the end of sequence (e.g., EOS)

$$P(x_5 | x_{1:4} = ( \text{" BOS "}, \text{" I "}, \text{" am "}, \text{" good "}))$$



# Beam Search

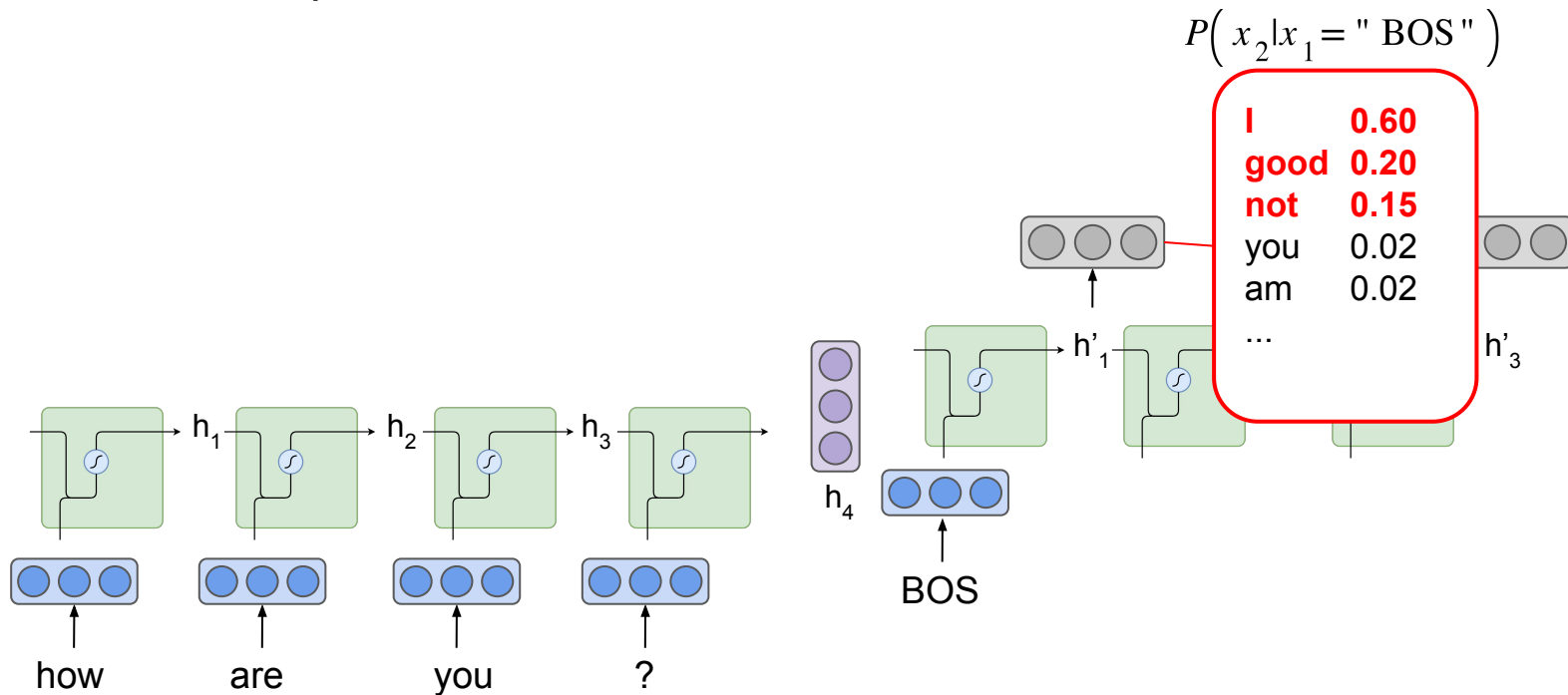
- Select top-k candidates at each step and grow the search tree



# Beam Search (k = 3)

k is often called “beam width”

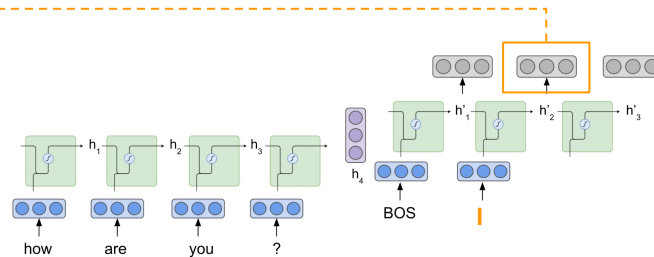
- Choose top-k tokens



## Step 2: Calculate Next Token Probability

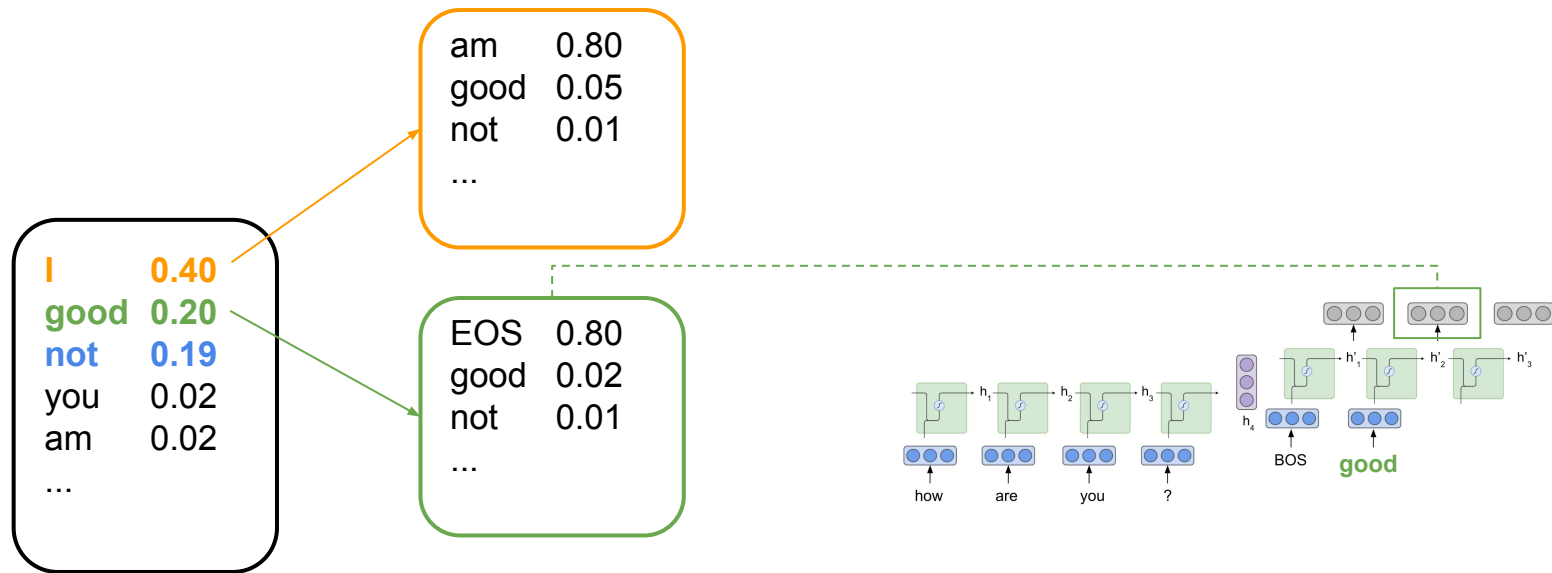
I	0.40
good	0.20
not	0.19
you	0.02
am	0.02
...	

am	0.80
good	0.05
not	0.01
...	

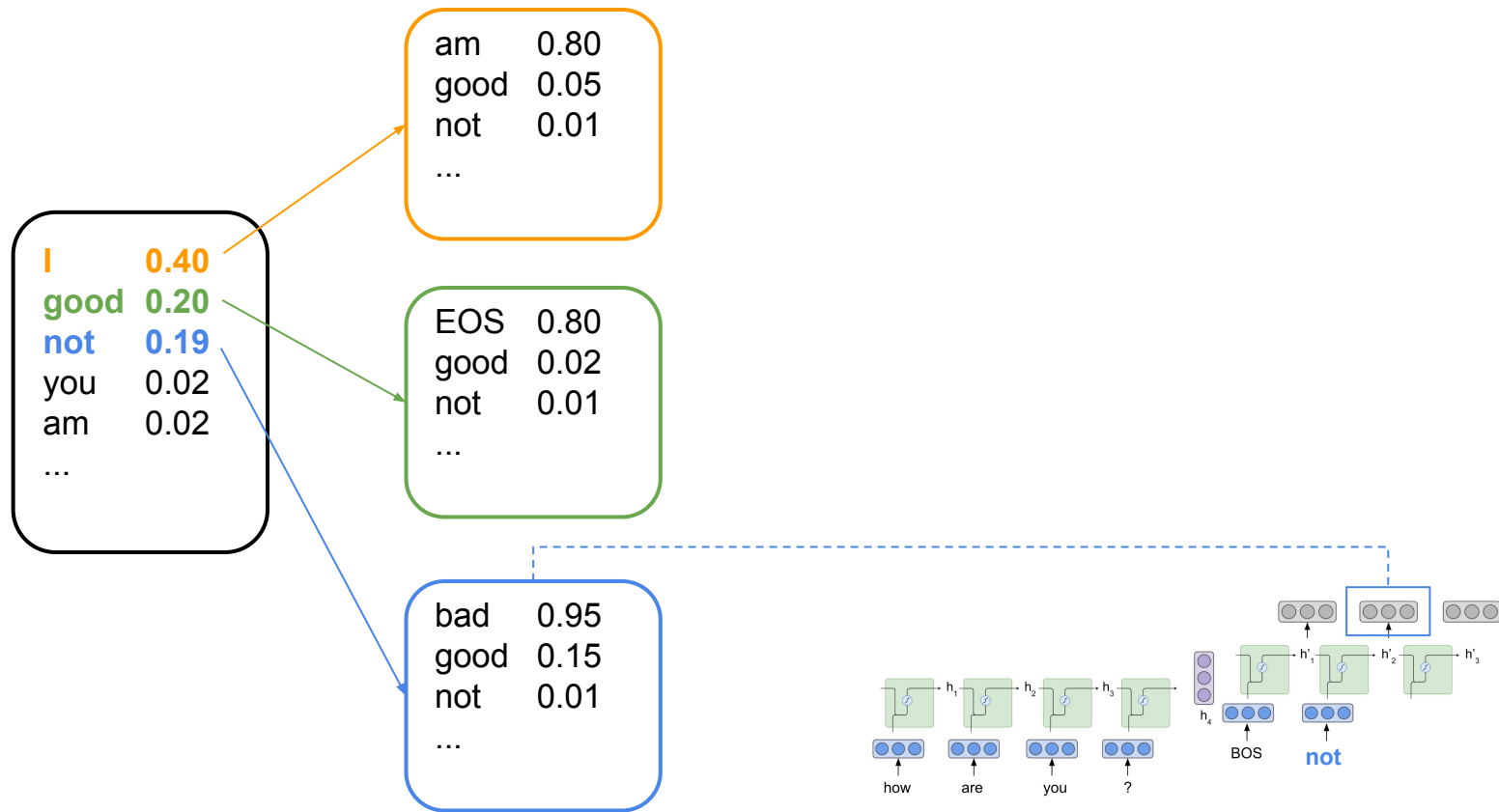




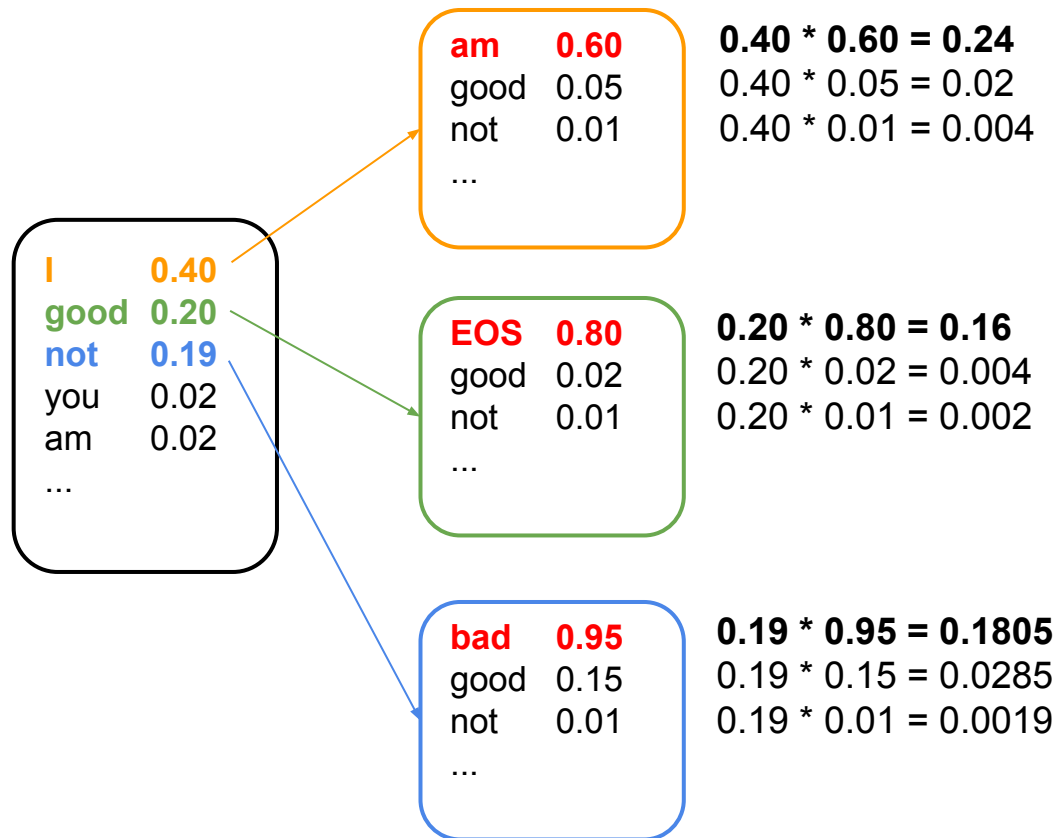
## Step 2: Calculate Next Token Probability



## Step 2: Calculate Next Token Probability



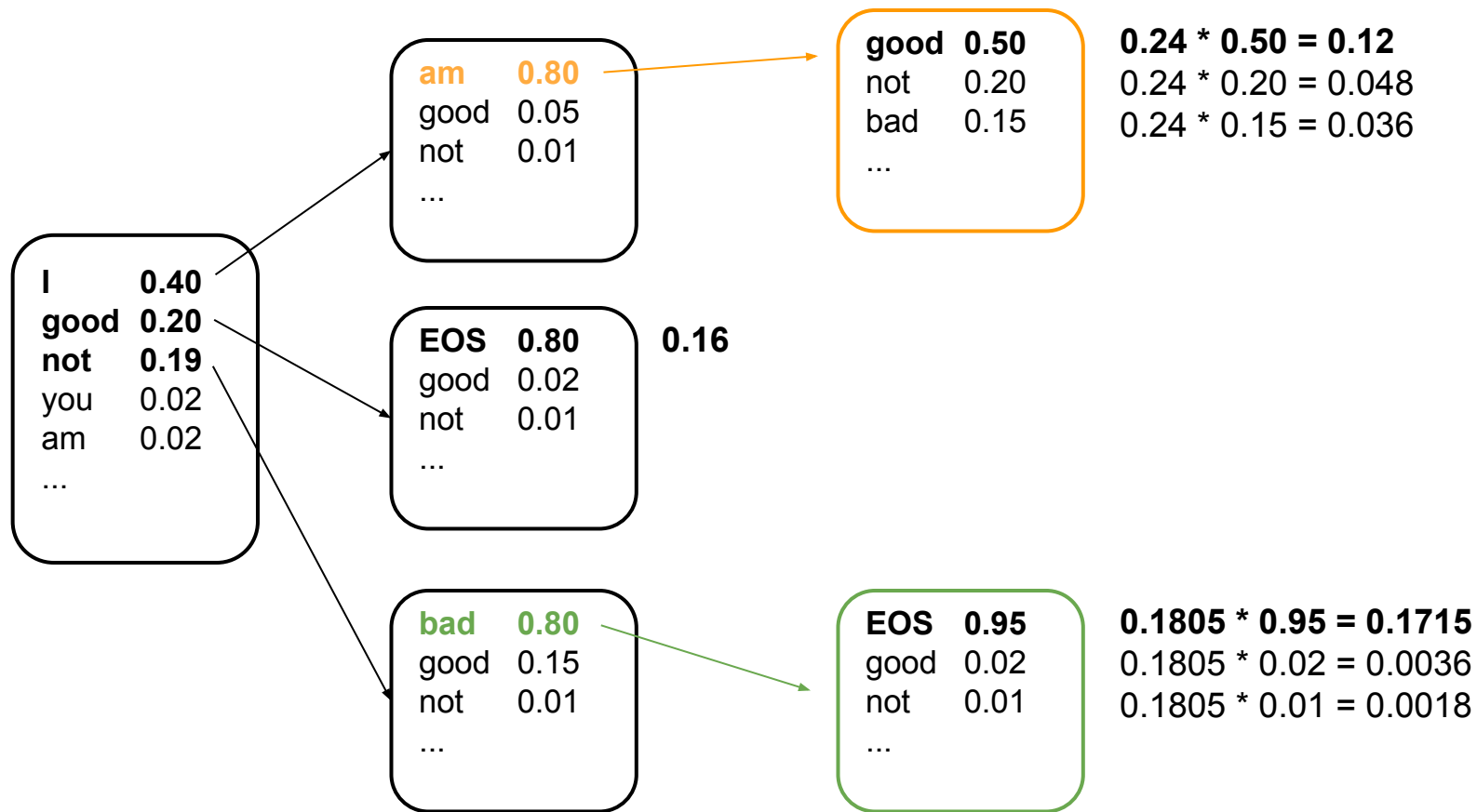
## Step 3: Choose Top-3 Candidates



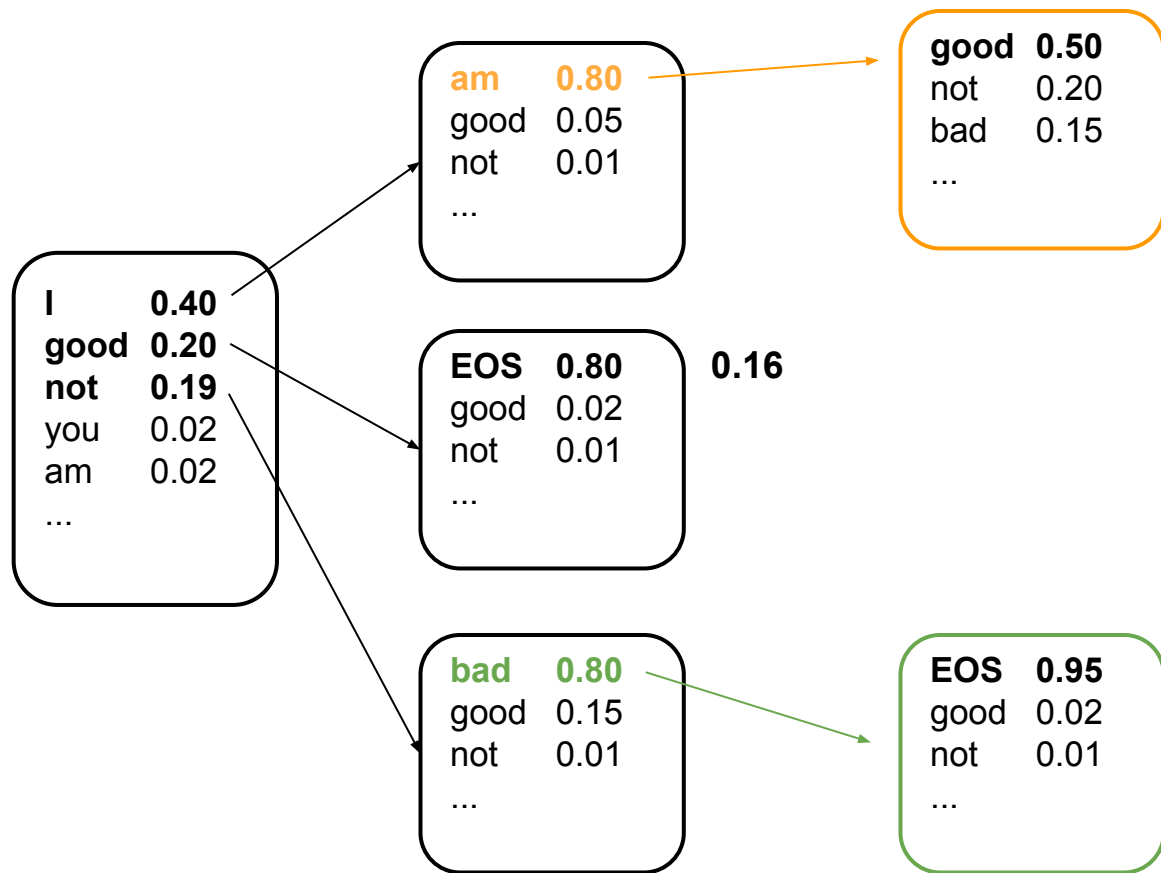
### Candidate ranking

- |                  |        |
|------------------|--------|
| 1: "I am ..."    | 0.24   |
| 2: "not bad ..." | 0.1805 |
| 3: "good"        | 0.16   |

# Repeat the Steps



# Update the Candidate List



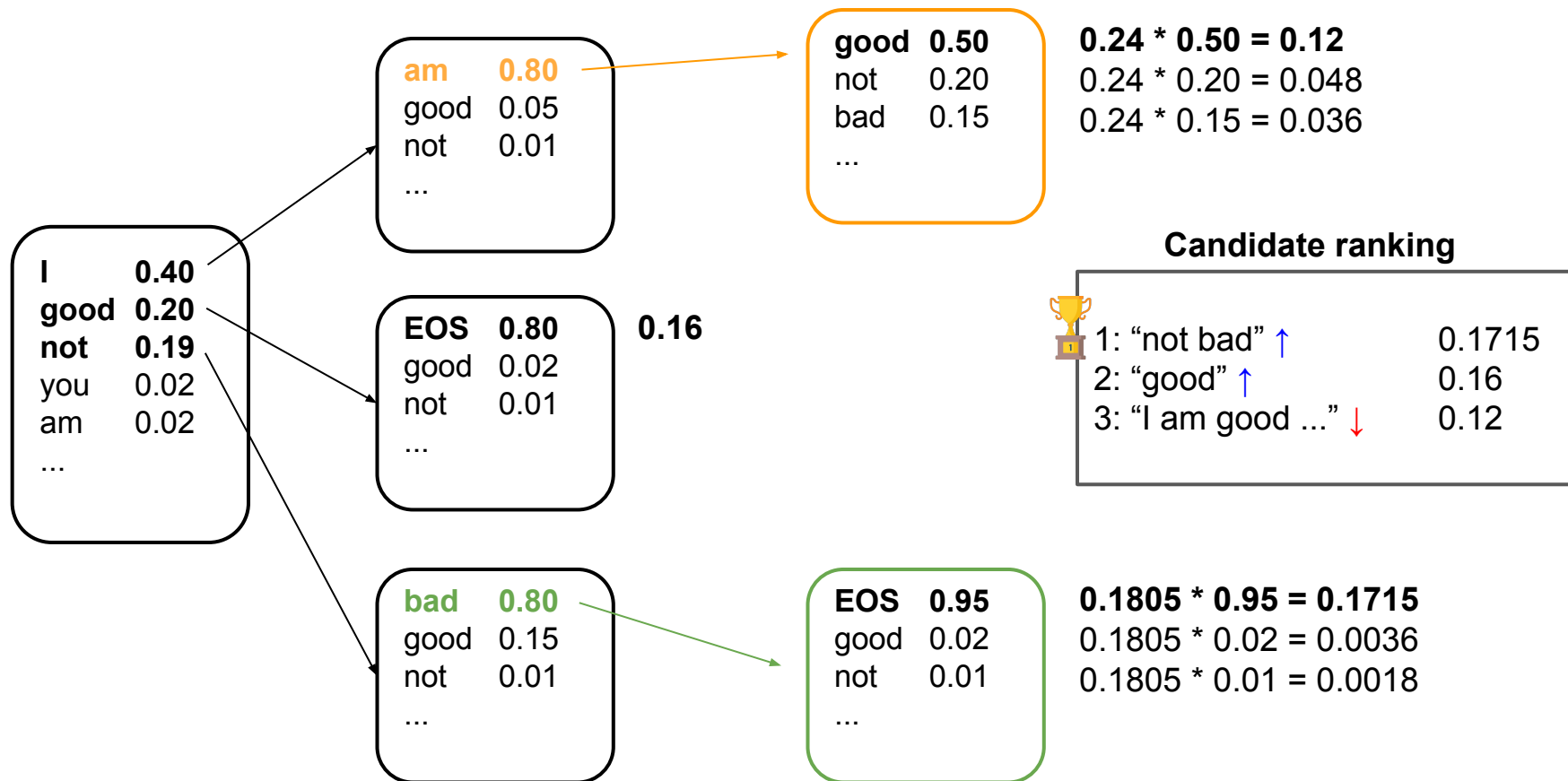
$$\begin{aligned}0.24 * 0.50 &= 0.12 \\0.24 * 0.20 &= 0.048 \\0.24 * 0.15 &= 0.036\end{aligned}$$

## Candidate ranking

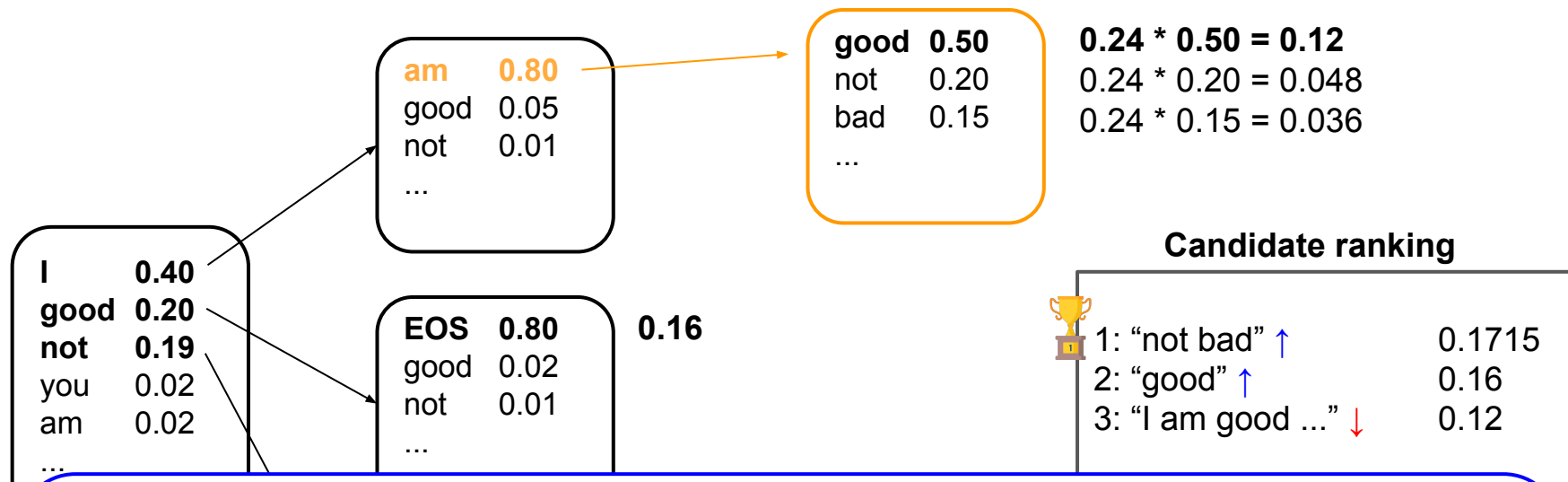
- |                      |        |
|----------------------|--------|
| 1: "not bad" ↑       | 0.1715 |
| 2: "good" ↑          | 0.16   |
| 3: "I am good ..." ↓ | 0.12   |

$$\begin{aligned}0.1805 * 0.95 &= 0.1715 \\0.1805 * 0.02 &= 0.0036 \\0.1805 * 0.01 &= 0.0018\end{aligned}$$

# Update the Candidate List



# Update the Candidate List



Key message:

The **same model** can return **a different sequence** with a different decoding algorithm

Greedy: "I am good"

Beam Search (k=2): "good"

Beam Search (k=3): "not bad"

# Summary: Decoding Algorithm

- Greedy Search
  - Time complexity:  $O(N)$
- Beam Search
  - Time complexity:  $O(k N)$
  - Common choice: Beam width = 5 ~ 10



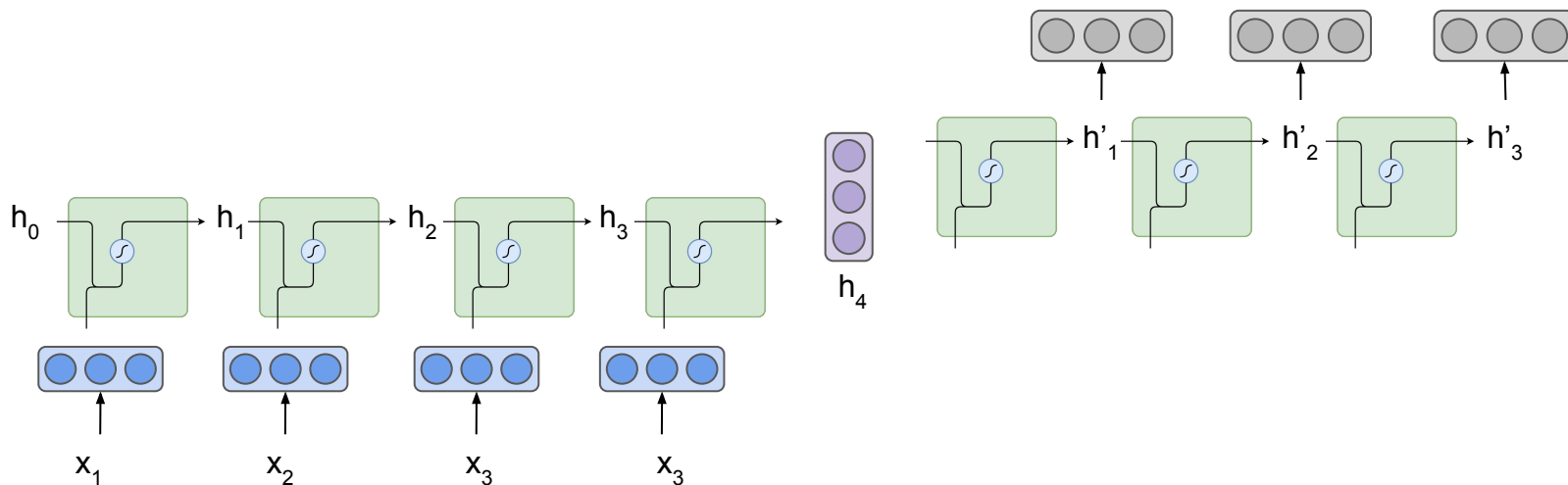
# Key Concept

- Encoder-Decoder Model

= Encoder + Decoder + Decoding algorithm

Model

Not model

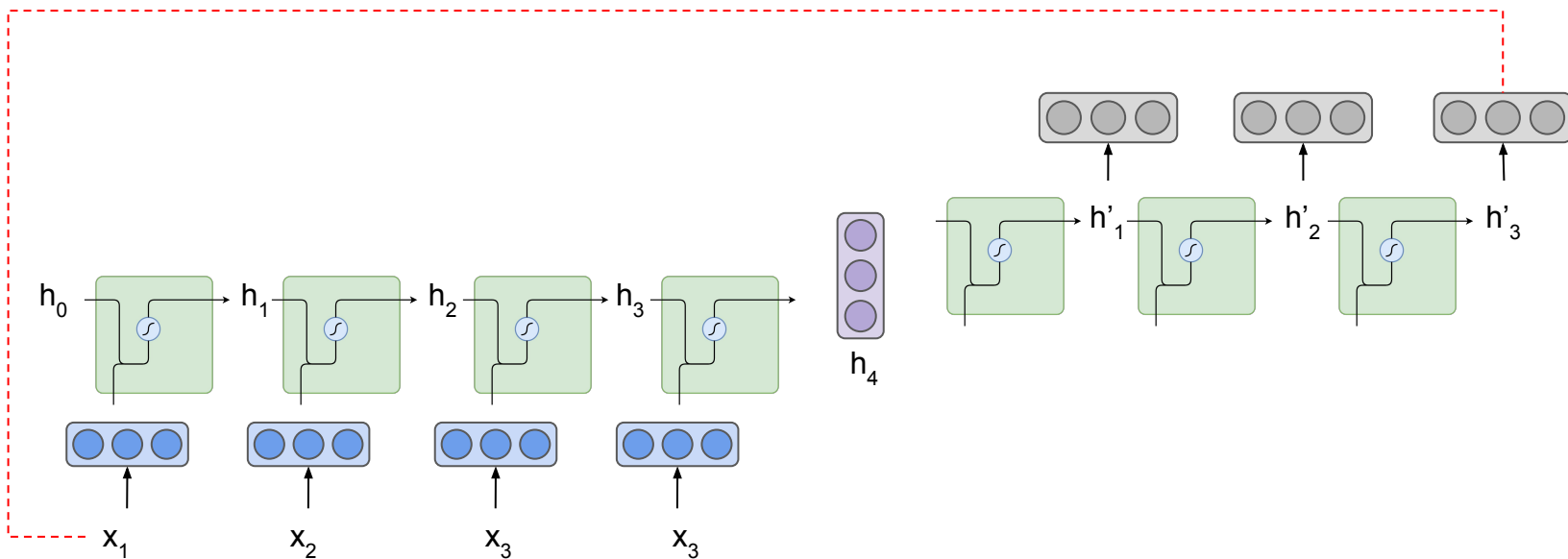


# Questions?

# Attention Mechanism

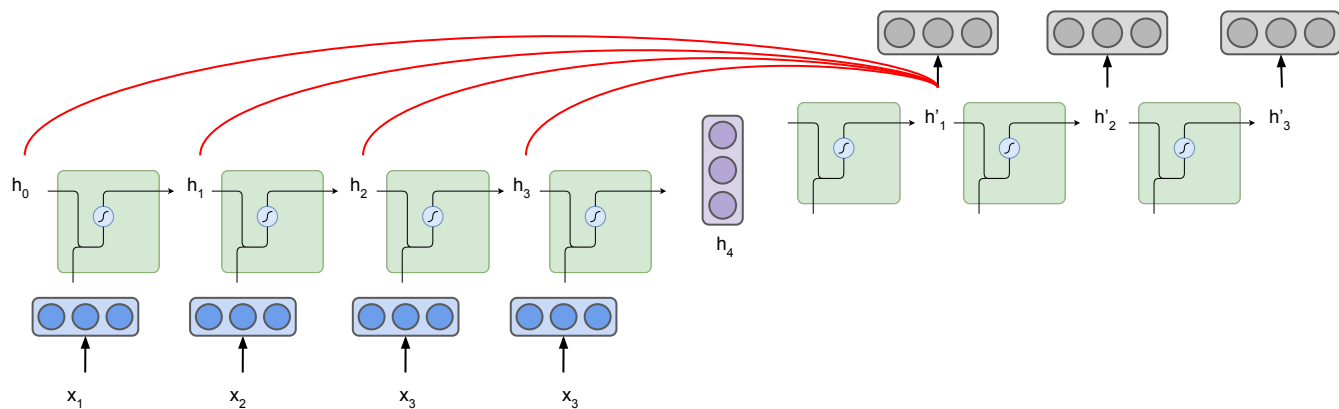
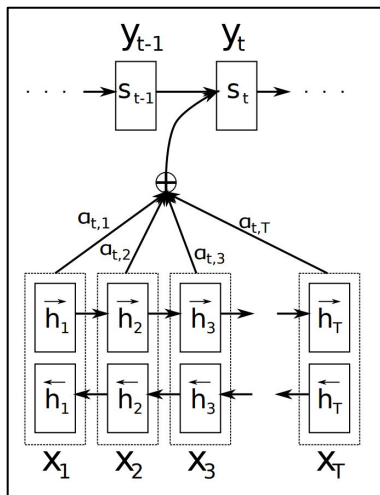
# What's the issue with Encoder-Decoder Models?

- Not a very good design for **long input/output sequence** even with LSTM/GRU
  - e.g., dependency b/w the beginning of input & the end of output



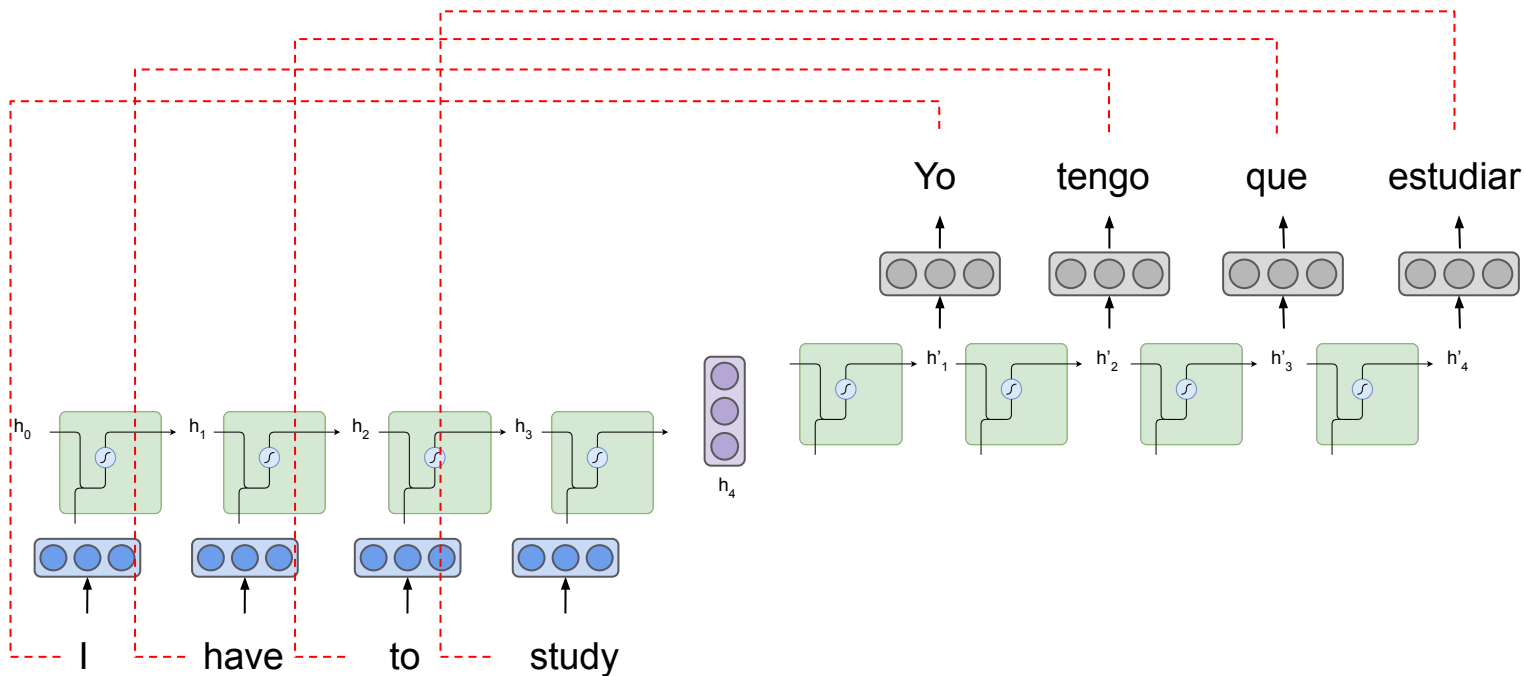
# Attention Mechanism [Bahdanau, Cho, and Bengio ICLR 2016]

- Additional input to the decoder based on **the alignments b/w encoder and decoder steps**

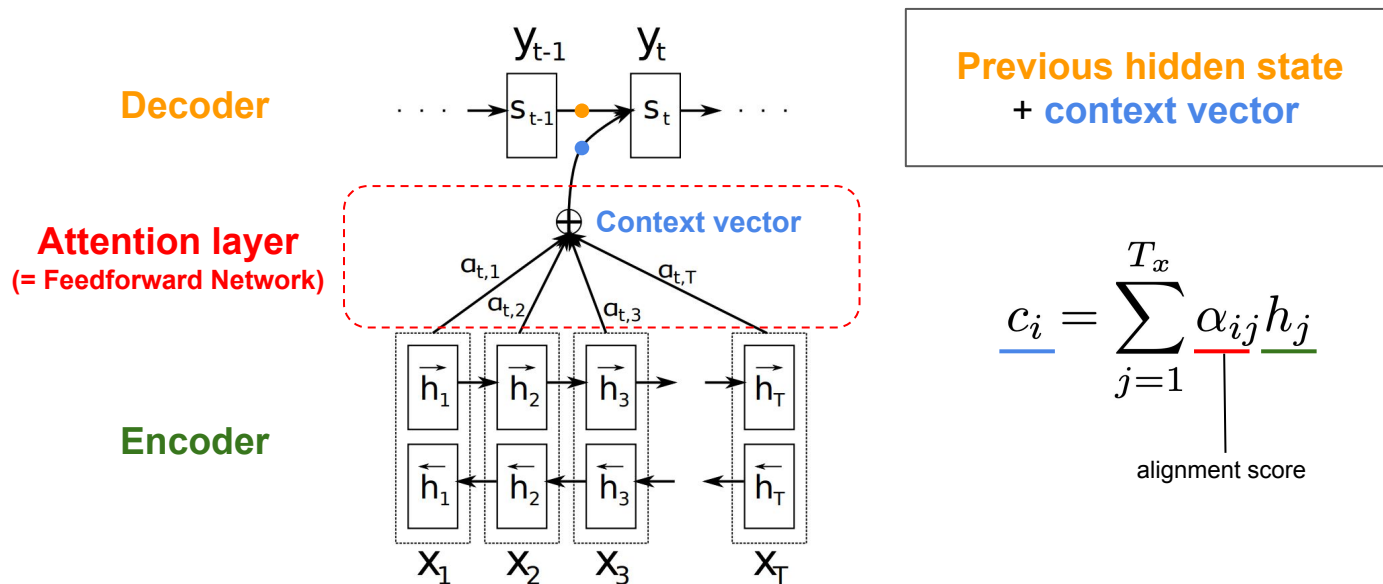


# Attention Mechanism: Intuition

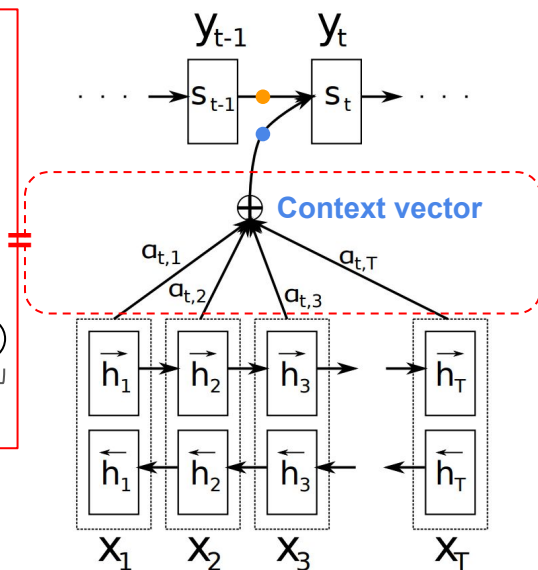
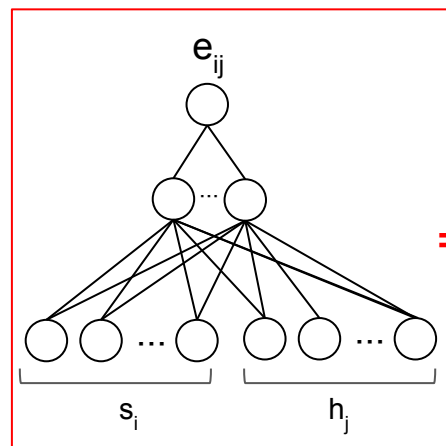
- “Direct” connections between any steps in the encoder/decoder model



# Attention Mechanism: Attention Layer



# Attention Layer in Depth



$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j$$

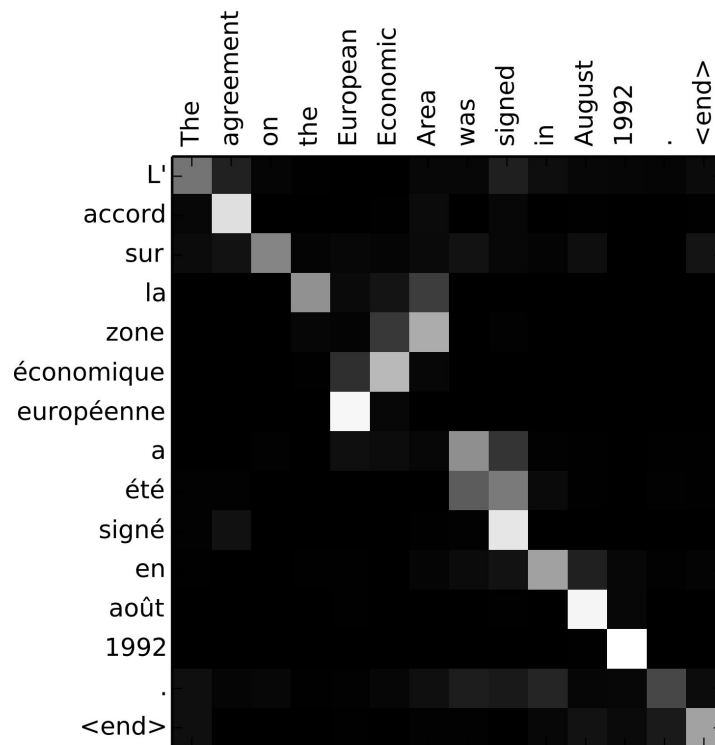
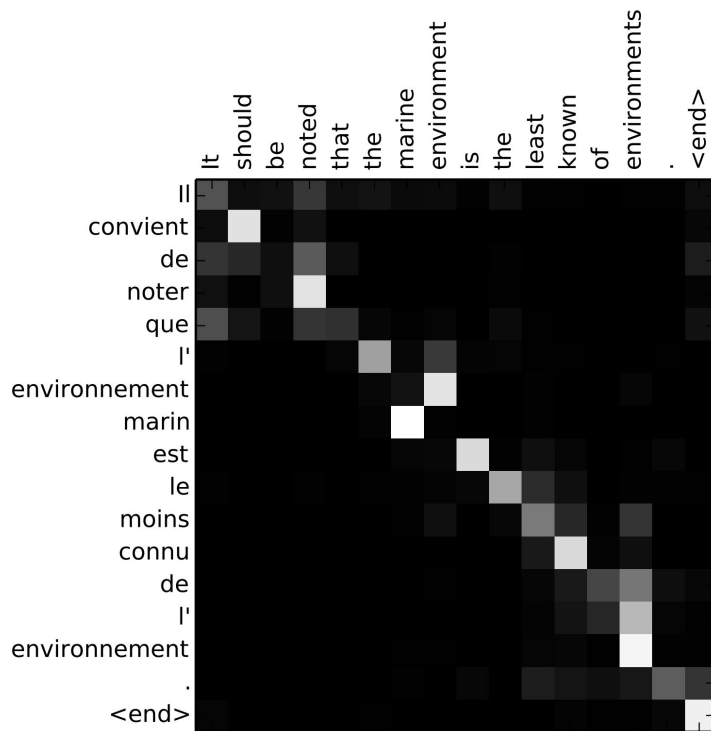
$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})}$$

$$e_{ij} = \text{score}(s_i, h_j)$$

$$\text{score}(s_i, h_j) = \mathbf{v}_a^T \tanh(\mathbf{W}_a [s_i; h_j])$$



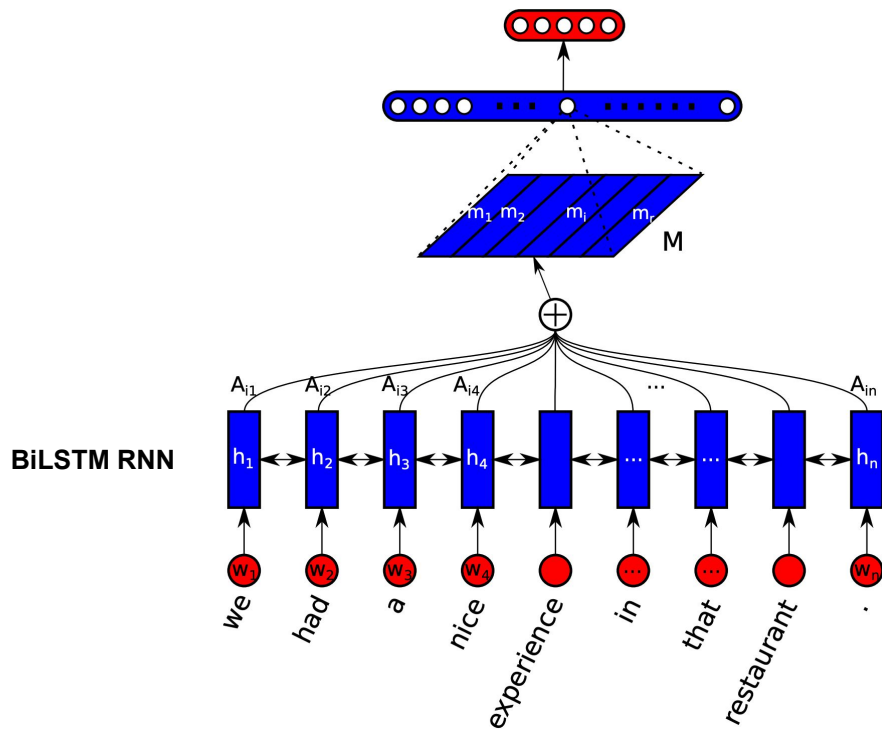
# Attention Weight Analysis



# Alignment Score Options

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

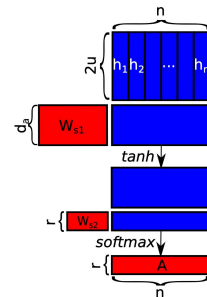
# Self-Attention Mechanism for Text Classification



$$M = AH$$

$$A = \text{softmax} (W_{s2} \tanh (W_{s1} H^T))$$

$$H = (h_1, h_2, \dots, h_n)$$



# Takeaway

- Encoder-Decoder Models (aka seq2seq Models)
  - Many applications: Response Generation, Machine Translation, Text Summarization, etc.
- Decoding Algorithms
  - Greedy Search
  - Beam Search
- Attention mechanism
  - Alignment scores
  - Self-attention mechanism

# Next Week

- Tue 10/5 Review session & initial discussion about term projects
  - Quick survey: Do you have any ideas in your mind? (Text? Image? or else?)
- Thu 10/7 Midterm exam (written exam on campus)
  - Time: 2pm-3:20pm
  - Location: TBA
- Fall Break!