Tuesday • October 19, 2021

# Encoder – Decoder Networks

Yale

LING 380/780

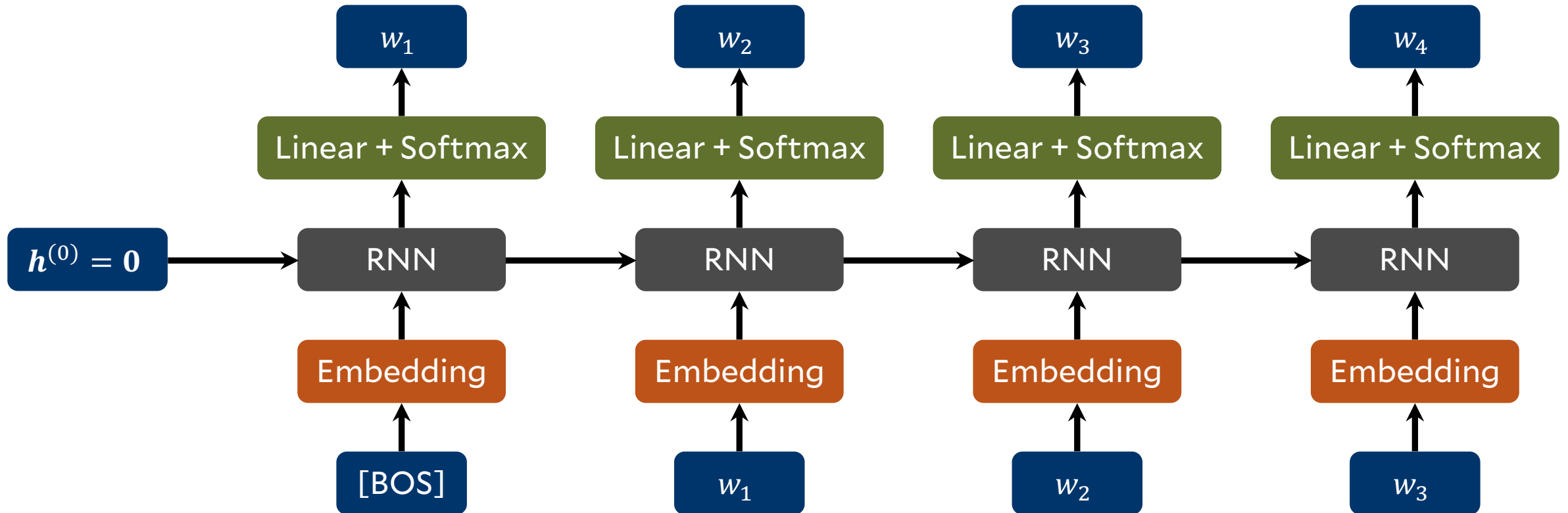*Neural Network Models of Linguistic Structure*

# Recurrent Neural Network Applications

What kinds of tasks are RNN architectures suitable for?

- Word prediction (language modeling)

- Classification

# Prediction and Generation

- Language models **generate text**.

# Greedy Generation Algorithm

- Initialize $h \leftarrow 0, w = [\text{BOS}]$.

- While $w \neq [\text{EOS}]$:

  - Feed $w$ and $h$ into the RNN.

  - Let $h$ be the new RNN state.

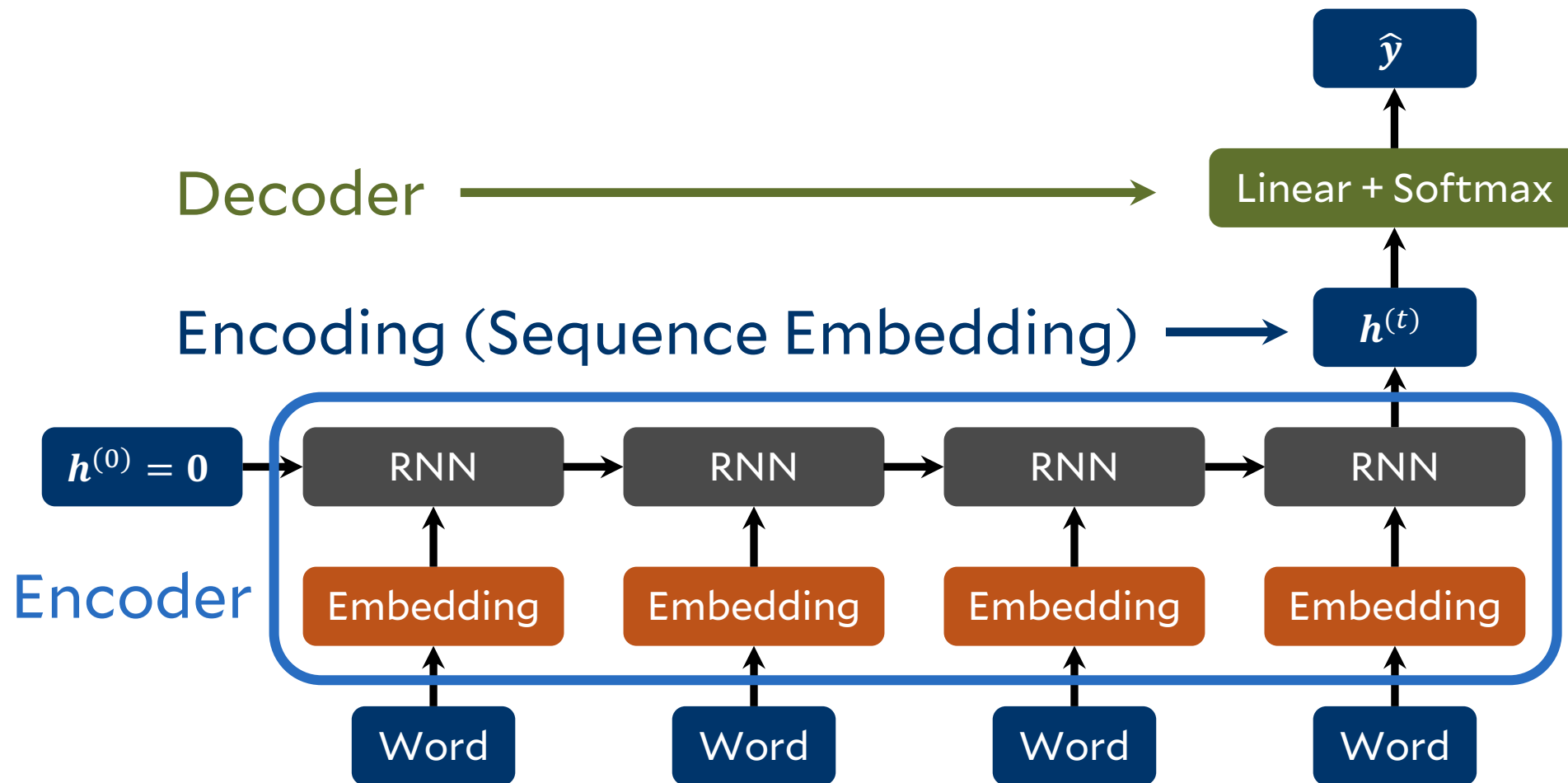  - Let $w$ be the word assigned the highest probability.

# Greedy Generation Algorithm

- "Greedy" means that at every step, the **locally best choice** is made.

- The next word is always the most likely word.

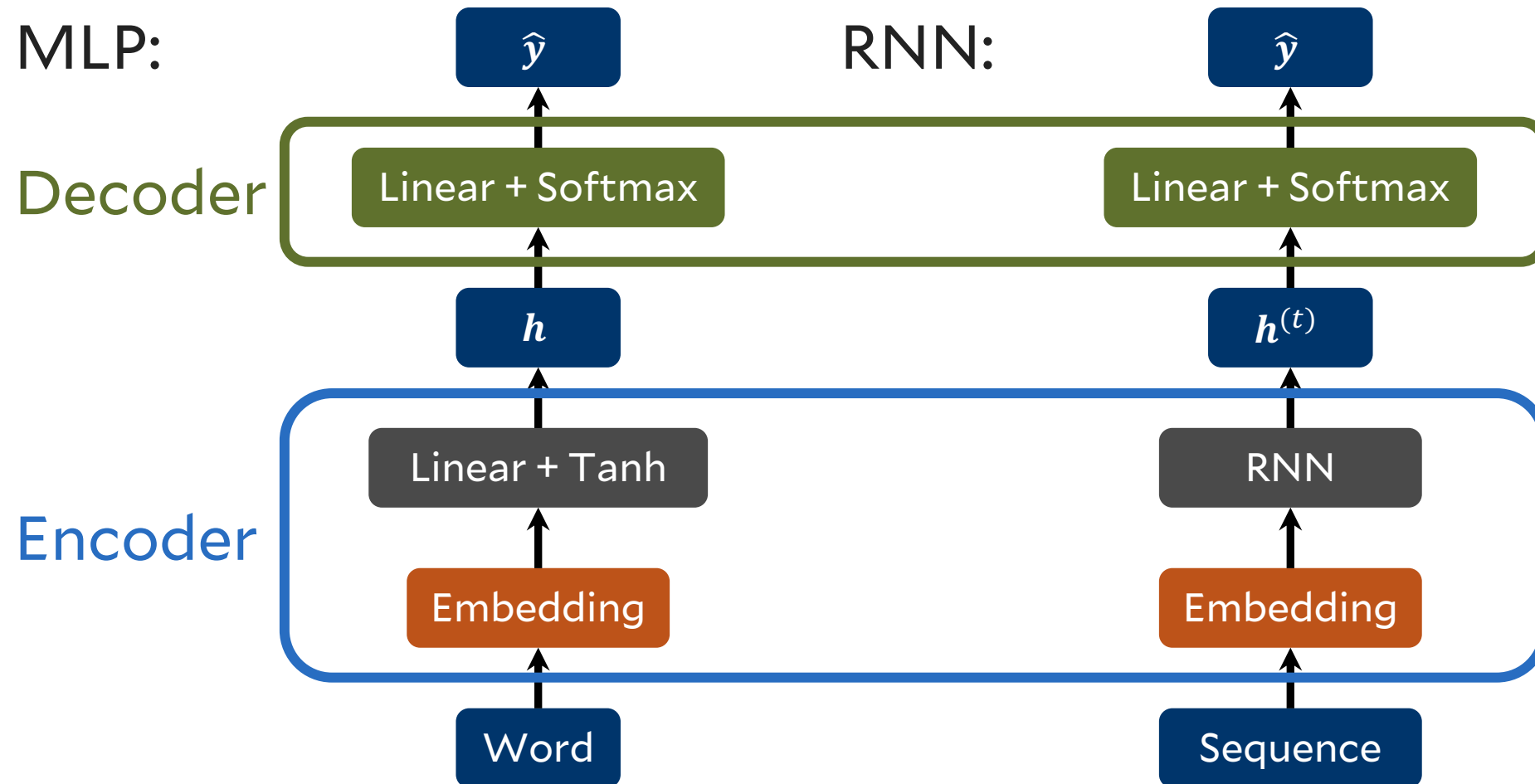- But this doesn't necessarily give the **best sequence**!

# Beam Search Generation Algorithm

- Initialize $s^{(i)} = [\text{BOS}]$ for $i \in \{1, 2, \ldots, k\}$.

- While none of the $s^{(i)}$s end with $[\text{EOS}]$:

  - For each $i$, let $s^{(i,j)}$ be the sequence obtained by adding the $j$th most likely next token to $s^{(i)}$.

  - Let $s^{(1)}, s^{(2)}, \ldots, s^{(k)}$ be the $k$ most likely sequence among the $s^{(i,j)}$s.

# RNN Classification

# MLP vs. RNN Classifiers

MLP:

RNN:

**Decoder**

**Encoder**

| | |
|---|---|
| $\hat{y}$ | $\hat{y}$ |
| Linear + Softmax | Linear + Softmax |
| $h$ | $h^{(t)}$ |
| Linear + Tanh | RNN |
| Embedding | Embedding |
| Word | Sequence |

# Types of RNNs

- **RNN Generators:** Autoregressively generate text from $h^{(0)} = 0$ and $w_0 = [\text{BOS}]$.

- **RNN Classifiers:** "Encode" an embedding for the sequence, then "decode" it using a linear + softmax.

- **Sequence-to-Sequence RNNs:** A third type of RNN!

# Sequence-to-Sequence (Seq2Seq) Tasks

**Sequence-to-sequence tasks** (seq2seq) are ones that take a **sequence** as input and **produce a sequence** as output.

**Examples:**

- Machine translation

- Text summarization

- Question answering
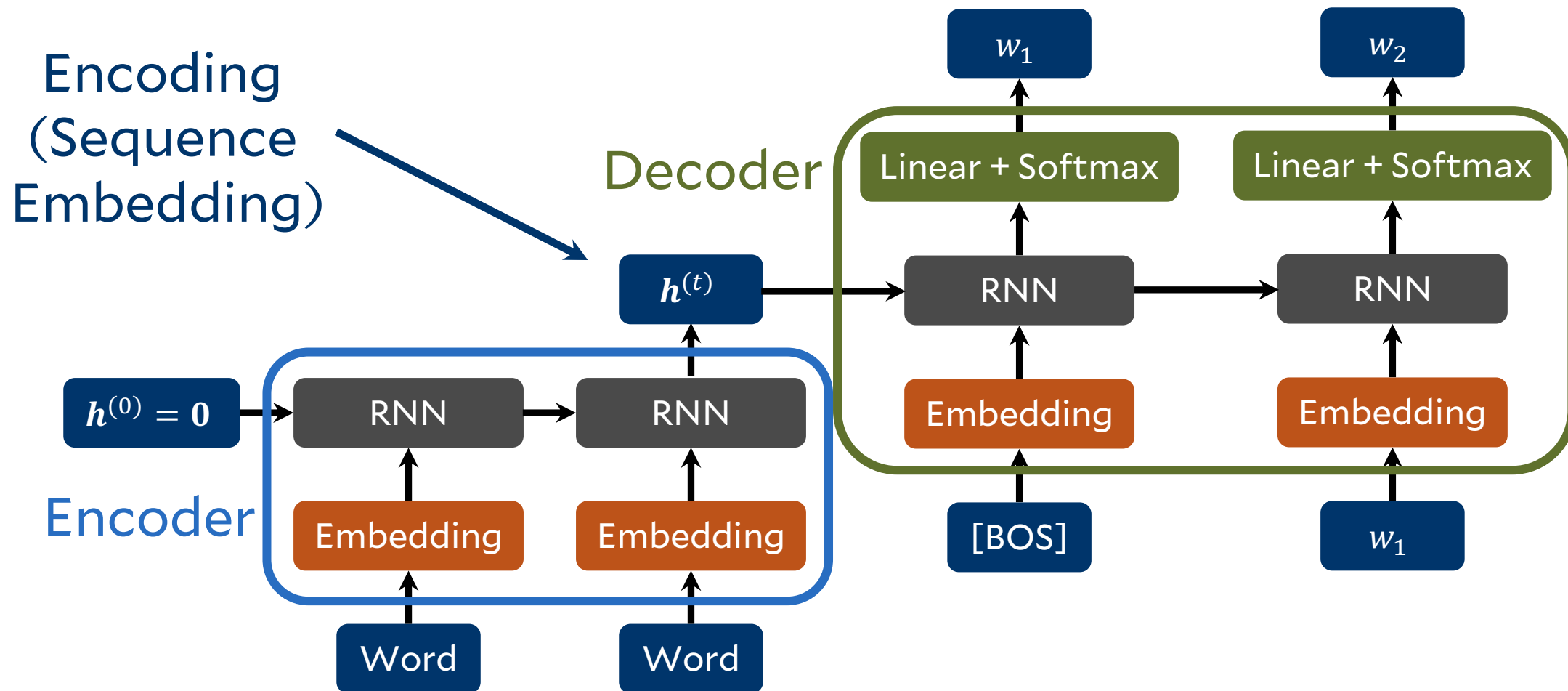
# Sequence-to-Sequence (Seq2Seq) Tasks

How do you apply RNNs to seq2seq tasks?

- **Classification:** RNN encoder, 1-layer MLP decoder

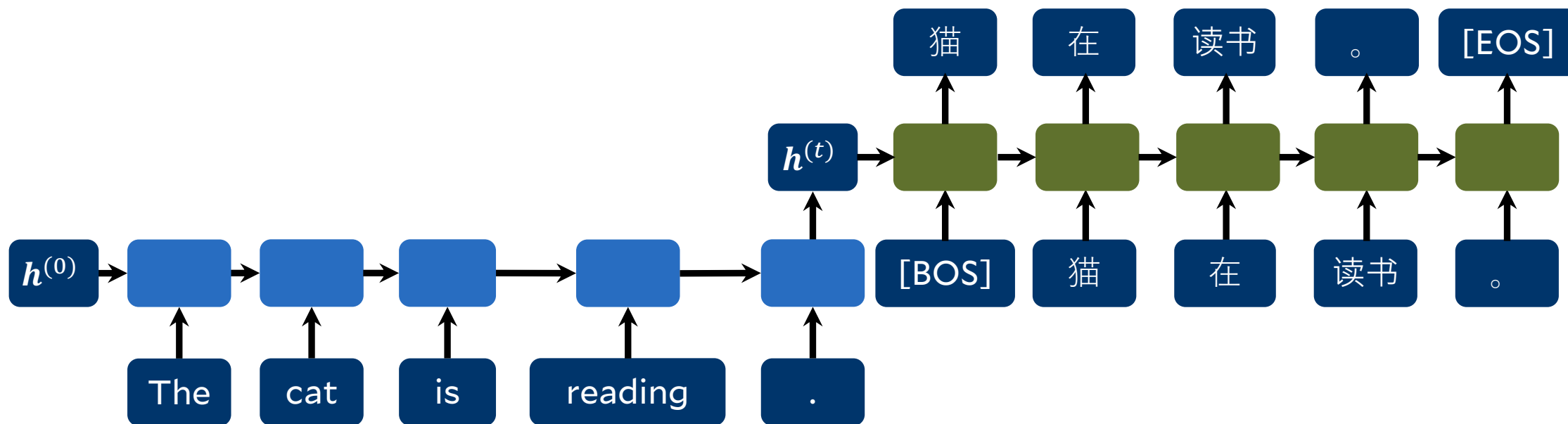- **Seq2Seq:** RNN encoder, autoregressive RNN decoder

# Types of Neural Networks

|  | Single Input | Sequence Input |
|---|---|---|
| Single Output | MLP | RNN Classifier |
| Sequence Output | RNN Generator | RNN Encoder-Decoder |

# Encoder–Decoder Architecture



Encoding (Sequence Embedding)

Decoder

$w_1$    $w_2$

Linear + Softmax    Linear + Softmax

$\boldsymbol{h}^{(t)}$

RNN    RNN

Embedding    Embedding

[BOS]    $w_1$

$\boldsymbol{h}^{(0)} = \boldsymbol{0}$

RNN    RNN

Encoder

Embedding    Embedding

Word    Word

# Example: Machine Translation

# Problem

- In the encoder–decoder architecture, all information about the sequence must be stored in the encoding vector.

- If the text is long, the vector might be too small to store the entire sequence.

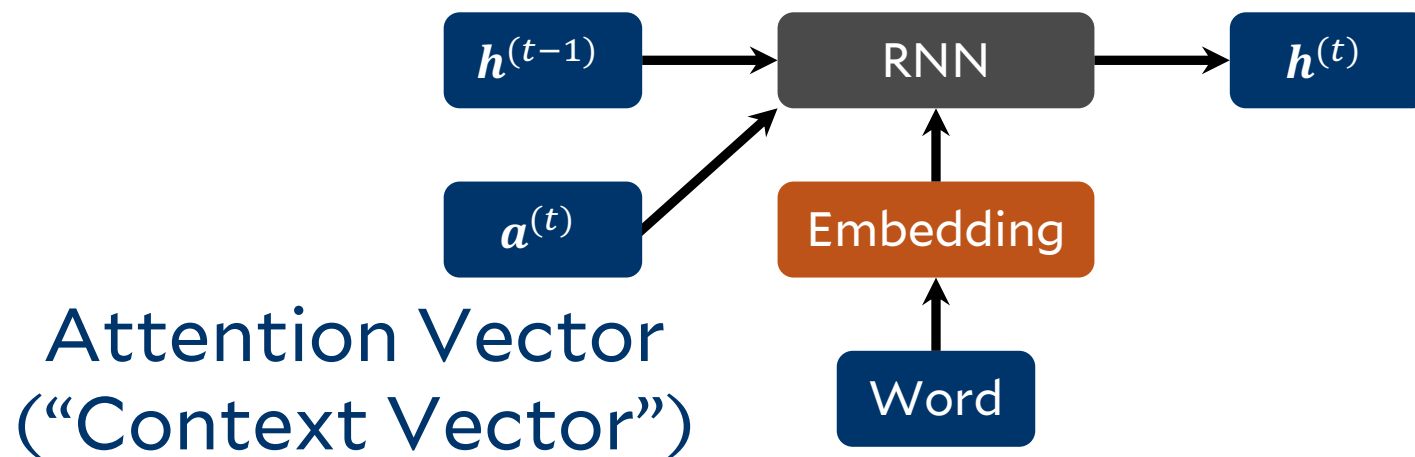- **Solution:** allow the decoder to peek at the input sequence!

窥视

# Attention

Two methods for incorporating the input sequence in the decoder:

- Bahdanau attention

- Luong attention

# Bahdanau Attention

Add a third input to the RNN decoder unit:



$h^{(t-1)}$ → RNN → $h^{(t)}$

$a^{(t)}$

Embedding

Word

Attention Vector
("Context Vector")

# Bahdanau Attention

- At each decoding step, the decoder looks at ("attends to") parts of the encoder output and retrieves an attention vector $\boldsymbol{a}^{(t)}$.

$$\text{Linear}\left(\boldsymbol{x}^{(t)}, \boldsymbol{h}^{(t-1)}, \boldsymbol{a}^{(t)}\right) = \boldsymbol{W} \begin{bmatrix} \boldsymbol{x}^{(t)} \\ \boldsymbol{h}^{(t-1)} \\ \boldsymbol{a}^{(t)} \end{bmatrix} + \boldsymbol{b}$$

$$= \boldsymbol{W}^{(x)} \boldsymbol{x}^{(t)} + \boldsymbol{W}^{(h)} \boldsymbol{h}^{(t-1)} + \boldsymbol{W}^{(a)} \boldsymbol{a}^{(t)} + \boldsymbol{b}$$

# SRN with Attention

$$h^{(t)} = \tanh(W^{(x)}x^{(t)} + W^{(h)}h^{(t-1)} + W^{(a)}a^{(t)} + b)$$

# LSTM with Attention

$$f^{(t)} = \sigma(W^{(f,x)}x^{(t)} + W^{(f,h)}h^{(t-1)} + W^{(f,a)}a^{(t)} + b^{(f)})$$

$$i^{(t)} = \sigma(W^{(i,x)}x^{(t)} + W^{(i,h)}h^{(t-1)} + W^{(i,a)}a^{(t)} + b^{(i)})$$

$$o^{(t)} = \sigma(W^{(o,x)}x^{(t)} + W^{(o,h)}h^{(t-1)} + W^{(o,a)}a^{(t)} + b^{(o)})$$

$$\tilde{c}^{(t)} = \tanh(W^{(c,x)}x^{(t)} + W^{(c,h)}h^{(t-1)} + W^{(c,a)}a^{(t)} + b^{(c)})$$

$$c^{(t)} = f^{(t)} \odot c^{(t-1)} + i^{(t)} \odot \tilde{c}^{(t)}$$

$$h^{(t)} = o^{(t)} \odot c^{(t)}$$

# GRU with Attention

$$r^{(t)} = \sigma(W^{(r,x)}x^{(t)} + W^{(r,h)}h^{(t-1)} + W^{(r,a)}a^{(t)} + b^{(r)})$$

$$z^{(t)} = \sigma(W^{(z,x)}x^{(t)} + W^{(z,h)}h^{(t-1)} + W^{(z,a)}a^{(t)} + b^{(z)}))$$

$$\widetilde{h}^{(t)} = \tanh(W^{(h,x)}x^{(t)} + W^{(h,a)}a^{(t)} + b^{(h)}$$

$$+ r^{(t)} \odot (W^{(r,h)}h^{(t-1)} + b^{(h,r)}))$$

$$h^{(t)} = z^{(t)} \odot h^{(t-1)} + (1 - z^{(t)}) \odot \widetilde{h}^{(t)}$$

# Computing the Attention Vector

First, use an MLP (2 layers, tanh activation, no softmax) to calculate an **attention score**

$$s_i^{(t)} = \text{MLP}\left(\begin{bmatrix} \boldsymbol{h}^{(e,i)} \\ \boldsymbol{h}^{(d,t-1)} \end{bmatrix}\right)$$

where

- $\boldsymbol{h}^{(e,i)}$ is the $i$th encoder hidden state

- $\boldsymbol{h}^{(d,t)}$ is the $(t-1)$st decoder hidden state

# Computing the Attention Vector

Then, convert the attention scores into **attention weights**:

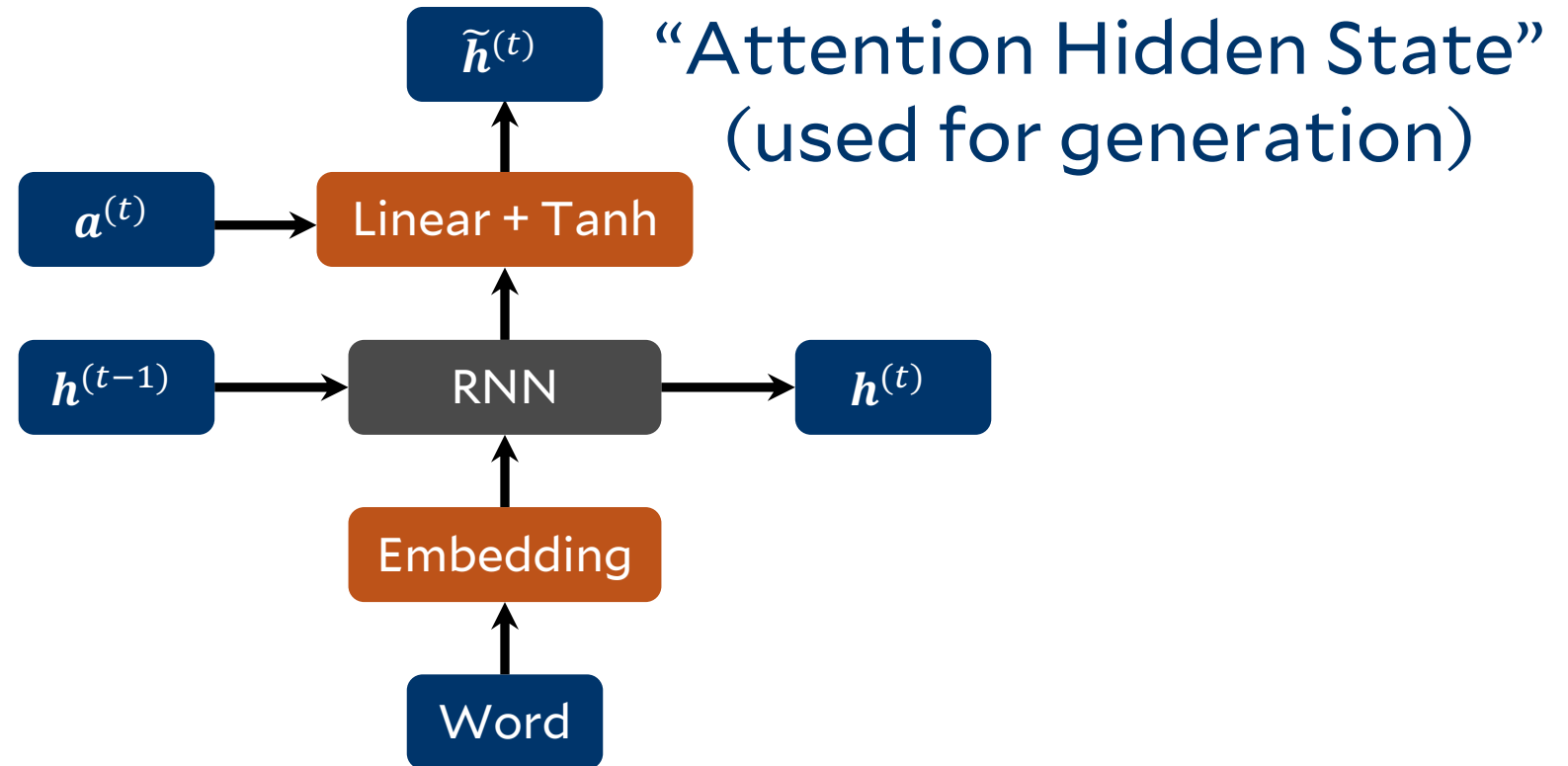$$\boldsymbol{\alpha}^{(t)} = \text{softmax}(\boldsymbol{s}^{(t)})$$

Use the attention weights to take a weighted average of encoder hidden states:

$$\boldsymbol{a}^{(t)} = \sum_i \alpha_i^{(t)} \, \boldsymbol{h}^{(e,i)}$$

# Luong Attention

Add the attention vector to the RNN decoder hidden state:

$$\widetilde{\boldsymbol{h}}^{(t)}$$

"Attention Hidden State"
(used for generation)

$\boldsymbol{a}^{(t)}$ → Linear + Tanh

$\boldsymbol{h}^{(t-1)}$ → RNN → $\boldsymbol{h}^{(t)}$

Embedding

Word

# Luong Attention

Computing attention scores

- "Dot": $s_i^{(t)} = \left( \boldsymbol{h}^{(d,t)} \right)^\top \boldsymbol{h}^{(e,i)}$

- "General": $s_i^{(t)} = \left( \boldsymbol{h}^{(d,t)} \right)^\top \boldsymbol{W} \boldsymbol{h}^{(e,i)}$

- "Concat" (similar to Bahdanau): $s_i^{(t)} = \mathrm{MLP} \left( \begin{bmatrix} \boldsymbol{h}^{(e,i)} \\ \boldsymbol{h}^{(d,t)} \end{bmatrix} \right)$

# Visualizing Attention Weights

Bahdanau et al. (2015):
English to French