



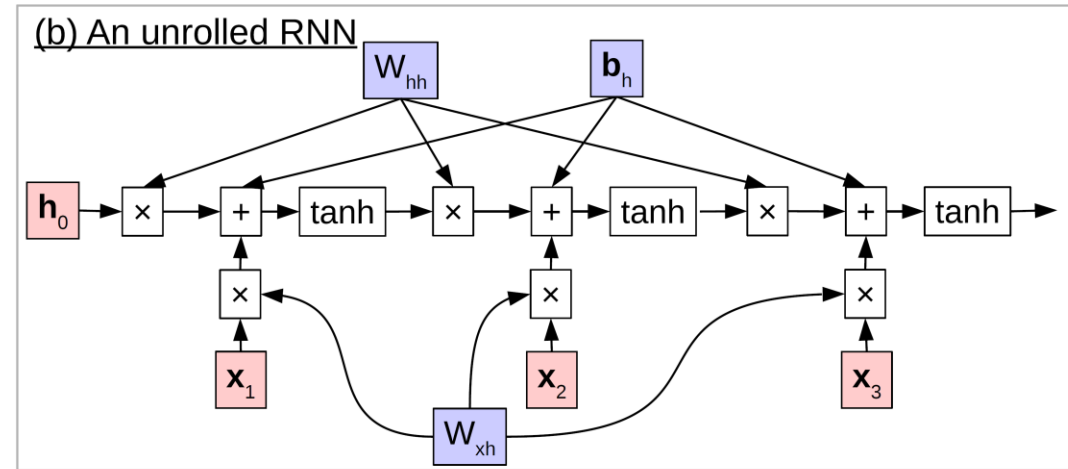
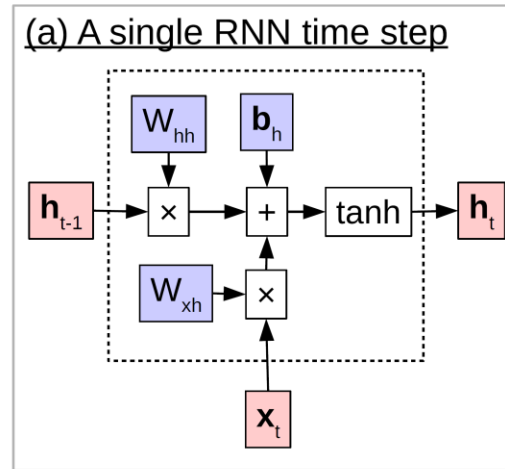
**COMPUTER SCIENCE**  
UNIVERSITY OF MARYLAND

# From Recurrent Language Models to Machine Translation

**CMSC 470**

Marine Carpuat

# A recurrent language model

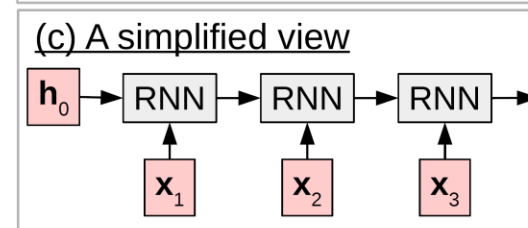
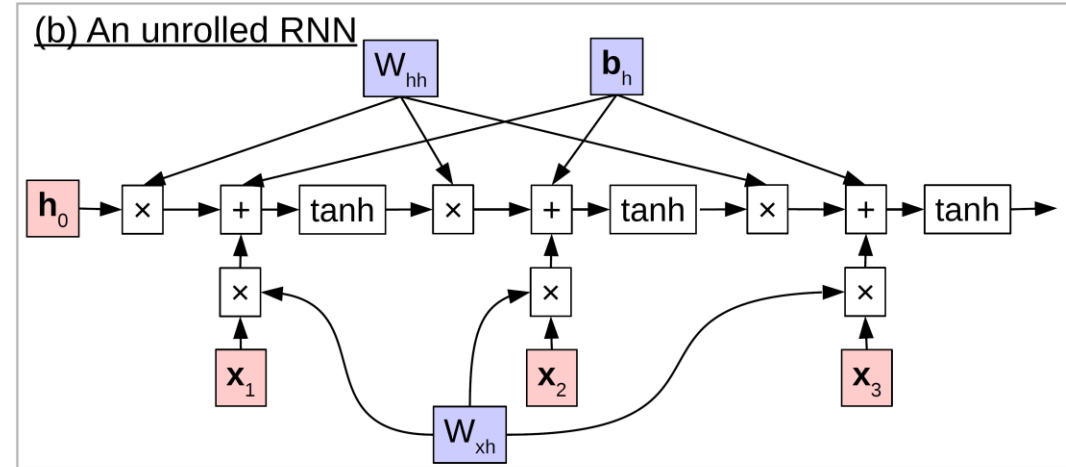
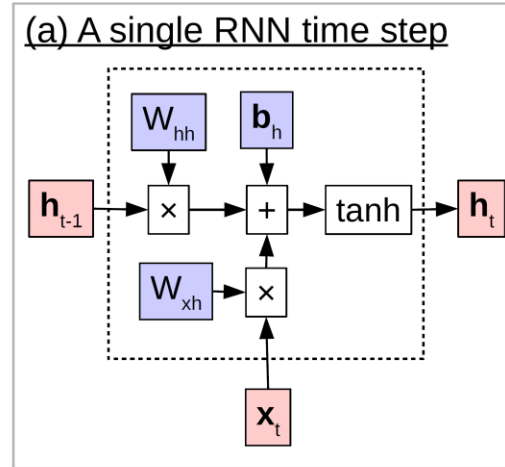


$$\mathbf{m}_t = M_{\cdot, e_{t-1}}$$

$$\mathbf{h}_t = \begin{cases} \tanh(W_{mh}\mathbf{m}_t + W_{hh}\mathbf{h}_{t-1} + \mathbf{b}_h) & t \geq 1, \\ \mathbf{0} & \text{otherwise.} \end{cases}$$

$$\mathbf{p}_t = \text{softmax}(W_{hs}\mathbf{h}_t + b_s).$$

# A recurrent language model



$$\mathbf{m}_t = M_{\cdot, e_{t-1}}$$

$$\mathbf{h}_t = \text{RNN}(\mathbf{m}_t, \mathbf{h}_{t-1})$$

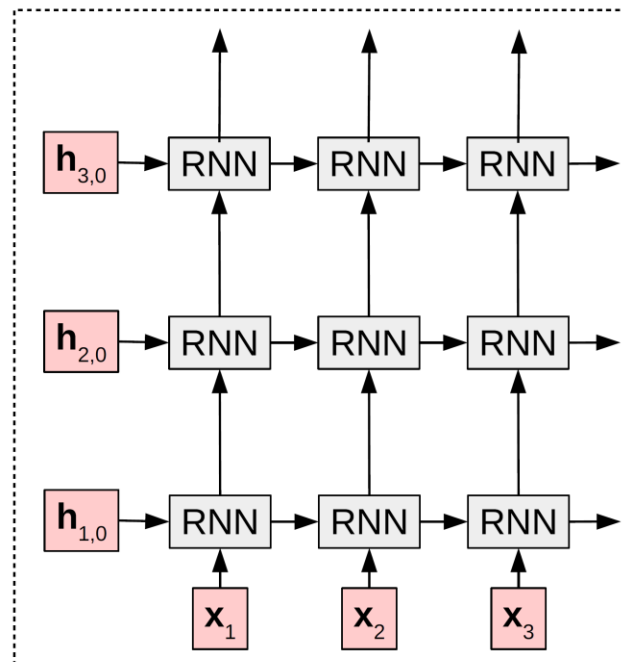
$$\mathbf{p}_t = \text{softmax}(W_{hs}\mathbf{h}_t + b_s).$$

# Examples of RNN variants

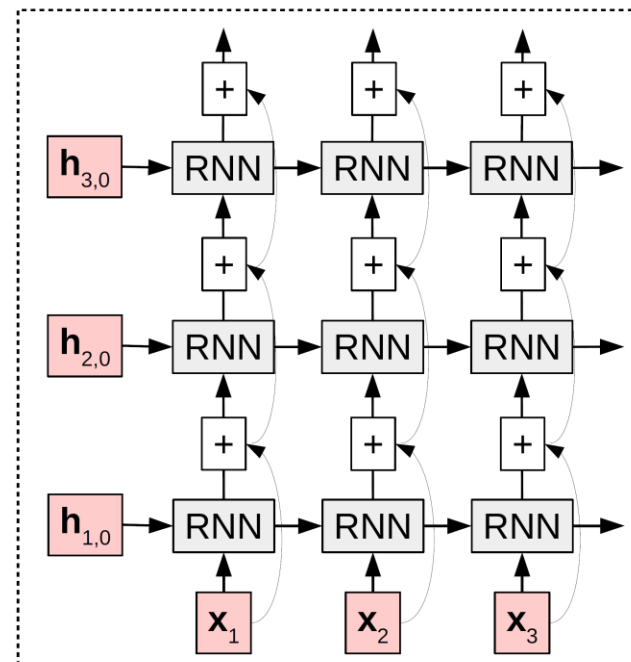
- LSTMs
  - Aim to address vanishing/exploding gradient issue

- Stacked RNNs

(a) A stacked RNN



(b) With residual connections



• ...

# Training in practice: online

---

**Algorithm 1** A fully online training algorithm

---

```
1: procedure ONLINE
2:   for several epochs of training do
3:     for each training example in the data do
4:       Calculate gradients of the loss
5:       Update the parameters according to this gradient
6:     end for
7:   end for
8: end procedure
```

---

# Training in practice: batch

---

**Algorithm 2** A batch learning algorithm

---

```
1: procedure BATCH
2:   for several epochs of training do
3:     for each training example in the data do
4:       Calculate and accumulate gradients of the loss
5:     end for
6:     Update the parameters according to the accumulated gradient
7:   end for
8: end procedure
```

---

# Training in practice: minibatch

- Compromise between online and batch
- Computational advantages
  - Can leverage vector processing instructions in modern hardware
  - By processing multiple examples simultaneously

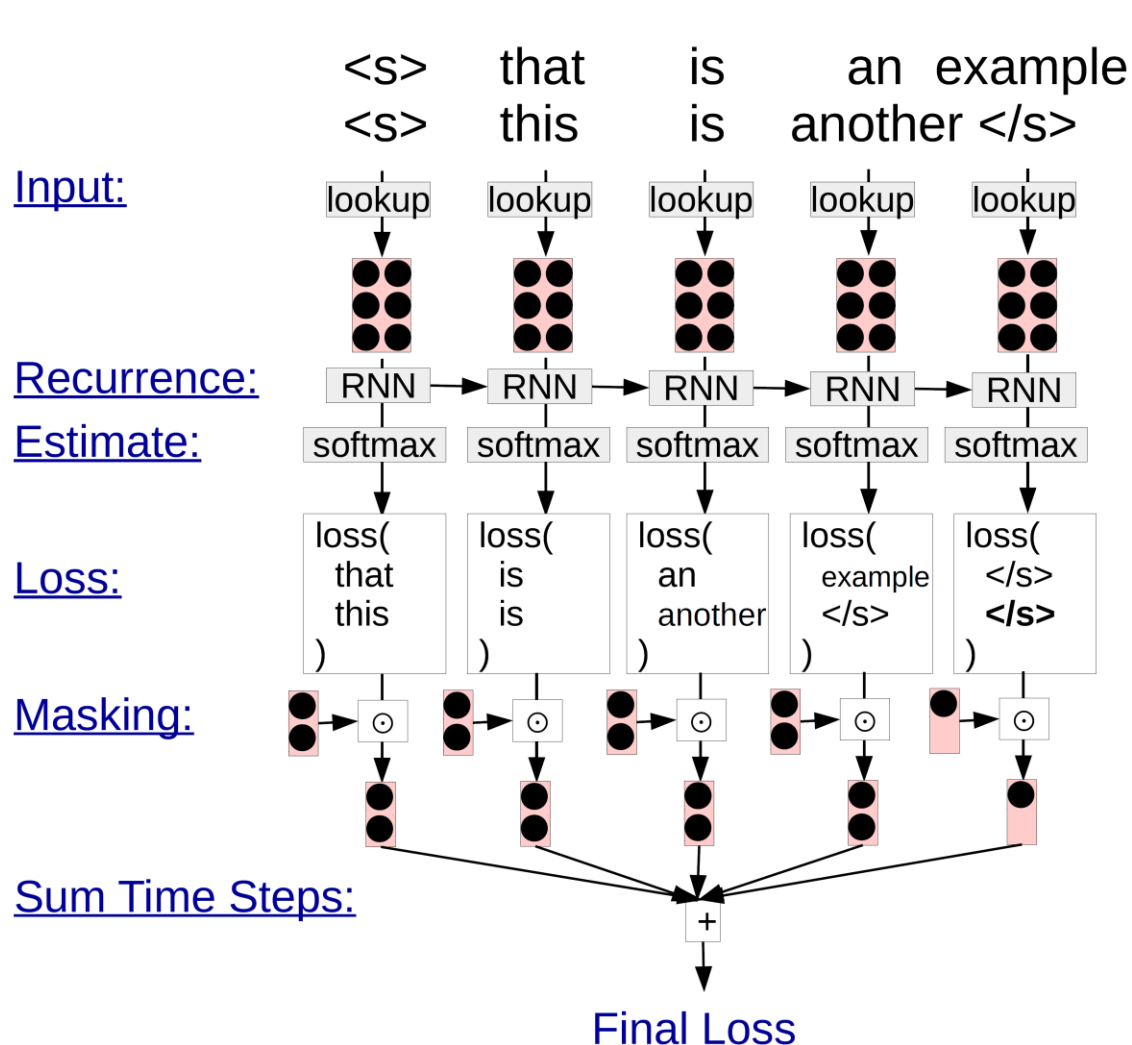
Operations w/o Minibatching

$$\tanh\left(\begin{array}{c|c|c} W & x_1 & b \\ \hline \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{array}\right) + \tanh\left(\begin{array}{c|c|c} W & x_2 & b \\ \hline \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{array}\right) + \tanh\left(\begin{array}{c|c|c} W & x_3 & b \\ \hline \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{array}\right)$$

Operations with Minibatching

$$\begin{array}{c} x_1 \ x_2 \ x_3 \rightarrow \text{concat} \rightarrow \\ \text{broadcast} \leftarrow b \\ \begin{array}{c|c|c} W & X & B \\ \hline \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{array} \end{array}$$
$$\tanh\left(\begin{array}{c|c|c} W & X & B \\ \hline \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{array}\right)$$

# Problem with minibatches: in language modeling, examples don't have the same length



## • 3 tricks

### • Padding

- Add </s> symbol to make all sentences same length

### • Masking

- Multiply loss function calculated over padded symbols by zero

- + sort sentences by length



# Machine Translation

- Translation system
  - Input: source sentence  $F$
  - Output: target sentence  $E$
  - Can be viewed as a function

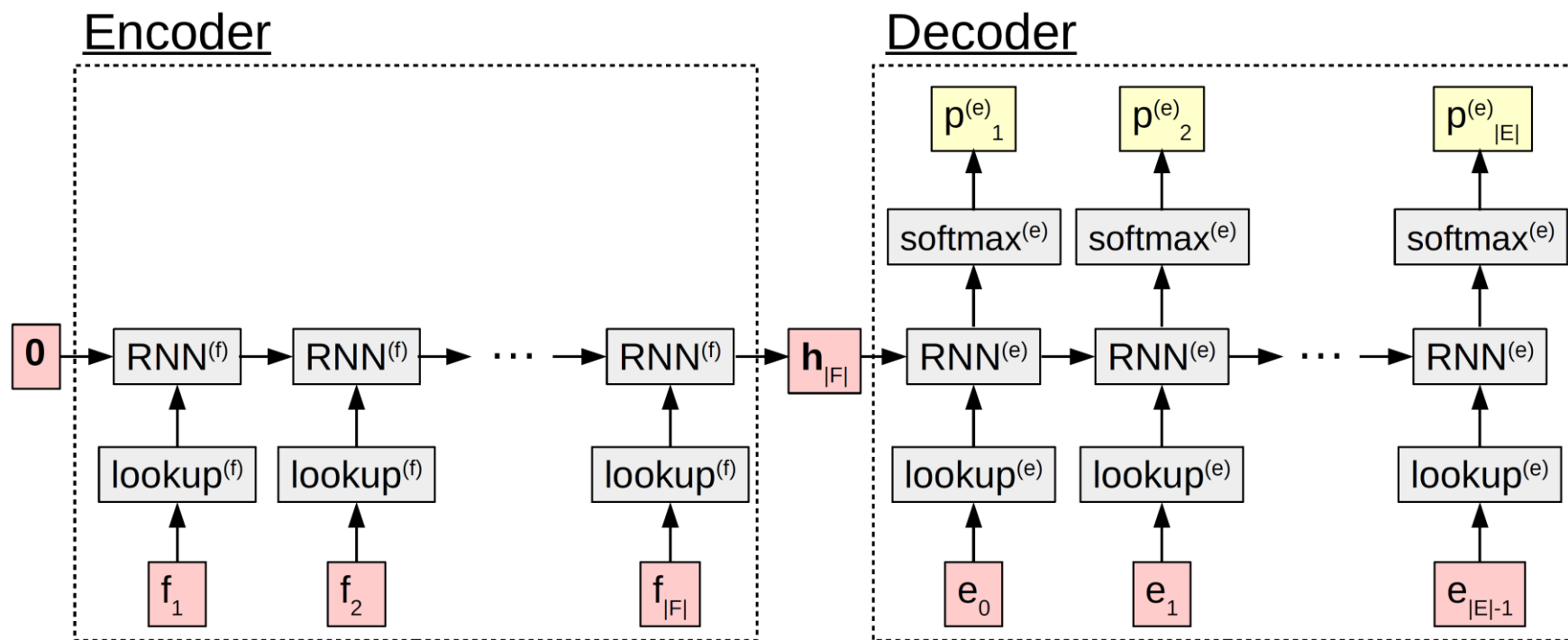
$$\hat{E} = \text{mt}(F)$$

- Statistical machine translation systems

$$\hat{E} = \underset{E}{\operatorname{argmax}} P(E | F; \theta)$$

- 3 problems
  - Modeling
    - how to define  $P(\cdot)$ ?
  - Training/Learning
    - how to estimate parameters from parallel corpora?
  - Search
    - How to solve  $\operatorname{argmax}$  efficiently?

# Encoder-decoder model



# Encoder-decoder model

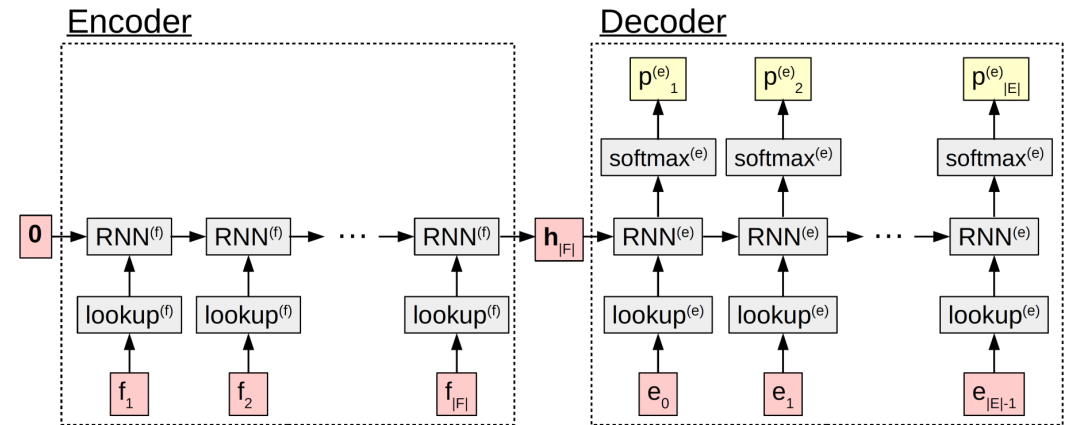
$$\mathbf{m}_t^{(f)} = M_{\cdot, f_t}^{(f)}$$

$$\mathbf{h}_t^{(f)} = \begin{cases} \text{RNN}^{(f)}(\mathbf{m}_t^{(f)}, \mathbf{h}_{t-1}^{(f)}) & t \geq 1, \\ \mathbf{0} & \text{otherwise.} \end{cases}$$

$$\mathbf{m}_t^{(e)} = M_{\cdot, e_{t-1}}^{(e)}$$

$$\mathbf{h}_t^{(e)} = \begin{cases} \text{RNN}^{(e)}(\mathbf{m}_t^{(e)}, \mathbf{h}_{t-1}^{(e)}) & t \geq 1, \\ \mathbf{h}_{|F|}^{(f)} & \text{otherwise.} \end{cases}$$

$$\mathbf{p}_t^{(e)} = \text{softmax}(W_{hs} \mathbf{h}_t^{(e)} + b_s)$$



# Generating Output

- We have a model  $P(E|F)$ , how can we generate translations?
- 2 methods
  - **Sampling**: generate a random sentence according to probability distribution
  - **Argmax**: generate sentence with highest probability

# Ancestral Sampling

- Randomly generate words one by one
- Until end of sentence symbol
- Done!

```
while  $y_{j-1} \neq \text{"</s>"}$ :  
   $y_j \sim P(y_j \mid X, y_1, \dots, y_{j-1})$ 
```

# Greedy search

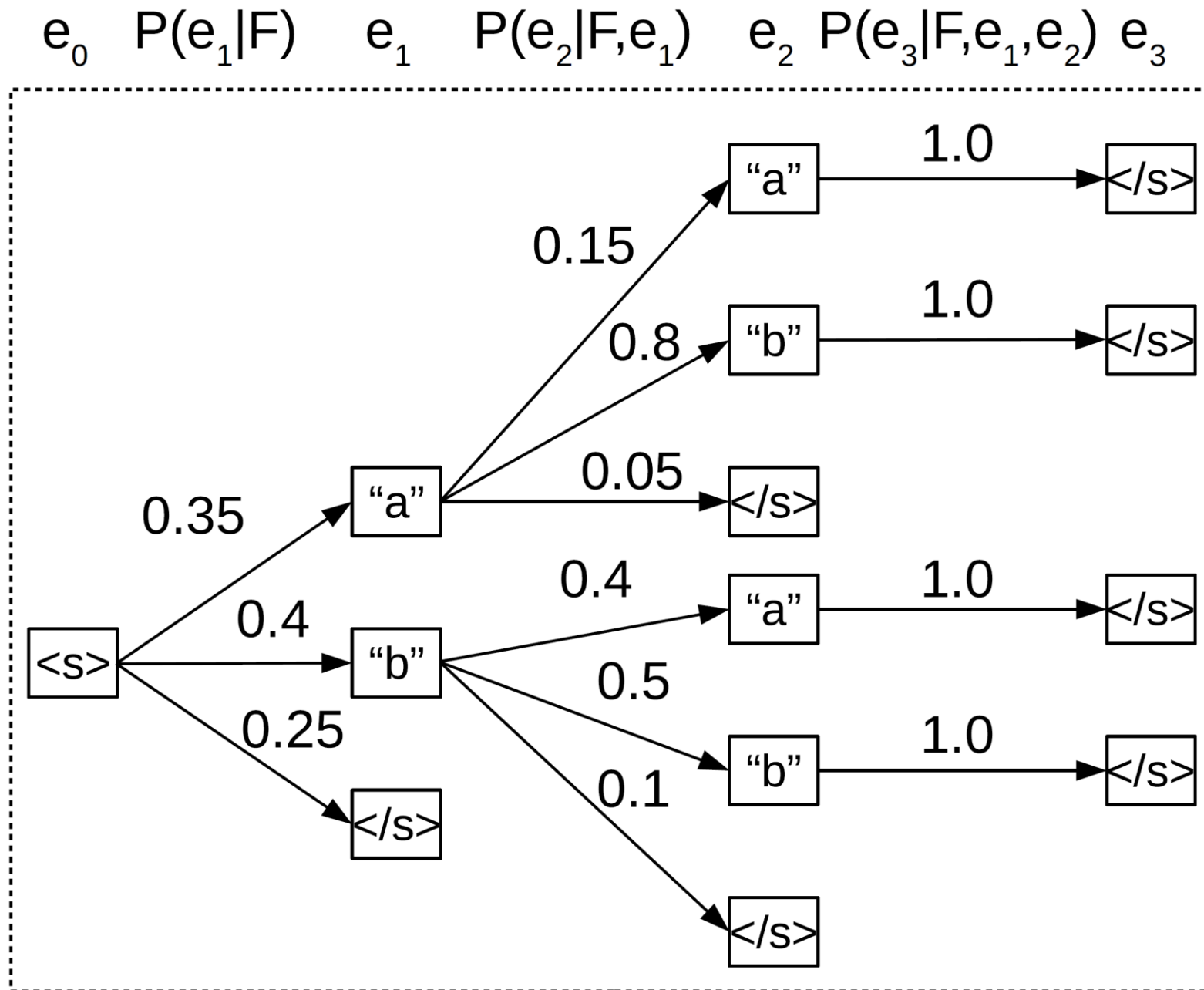
- One by one, pick single highest probability word

```
while  $y_{j-1} \neq \text{"</s>"}$ :  
   $y_j = \operatorname{argmax} P(y_j \mid X, y_1, \dots, y_{j-1})$ 
```

- Problems
  - Often generates easy words first
  - Often prefers multiple common words to rare words

# Greedy Search

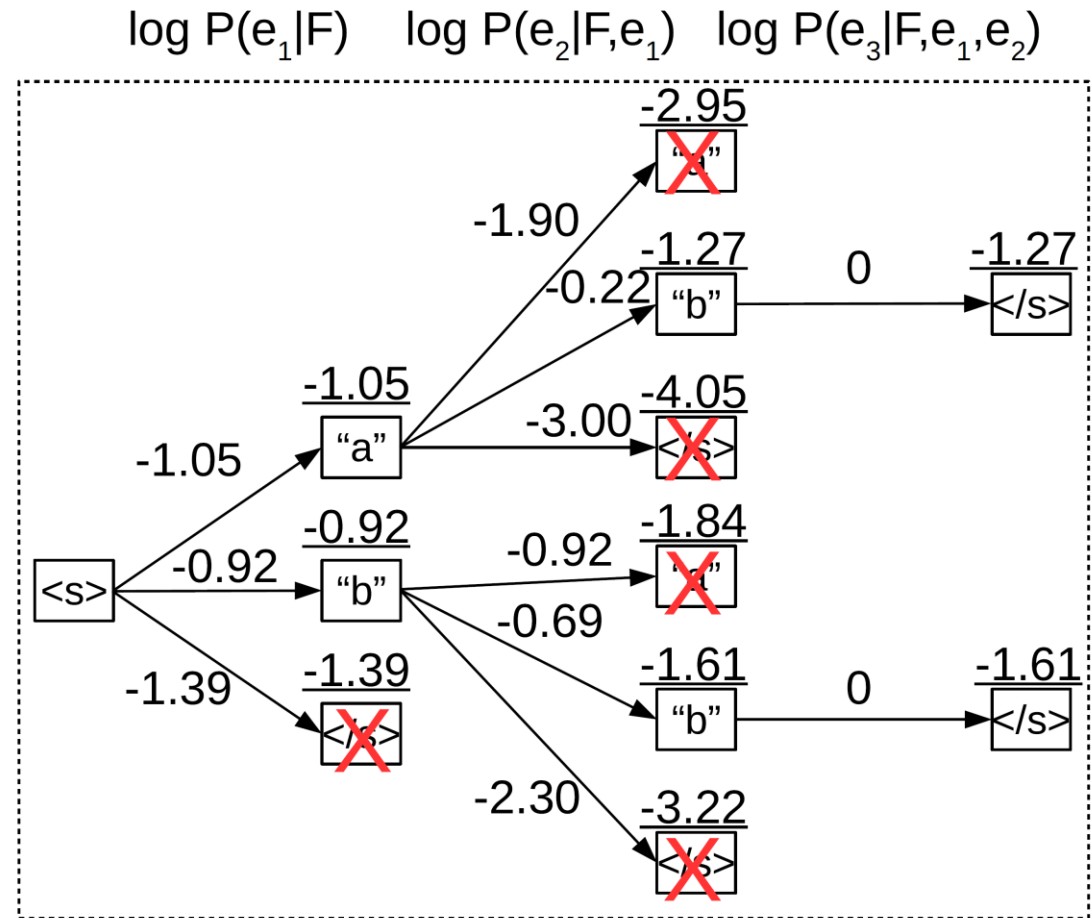
Example



# Beam Search

Example with beam size  $b = 2$

We consider  $b$  top hypotheses at each time step





# Other encoder structures: Bidirectional encoder

$$\vec{\mathbf{h}}_t^{(f)} = \begin{cases} \overrightarrow{\text{RNN}}^{(f)}(\mathbf{m}_t^{(f)}, \vec{\mathbf{h}}_{t+1}^{(f)}) & t \geq 1, \\ \mathbf{0} & \text{otherwise.} \end{cases}$$

$$\overleftarrow{\mathbf{h}}_t^{(f)} = \begin{cases} \overleftarrow{\text{RNN}}^{(f)}(\mathbf{m}_t^{(f)}, \overleftarrow{\mathbf{h}}_{t+1}^{(f)}) & t \leq |F|, \\ \mathbf{0} & \text{otherwise.} \end{cases}$$

Motivation:

- Help bootstrap learning
- By shortening length of dependencies

$$\mathbf{h}_0^{(e)} = \tanh(W_{\vec{f}e} \vec{\mathbf{h}}_{|F|} + W_{\overleftarrow{f}e} \overleftarrow{\mathbf{h}}_1 + \mathbf{b}_e)$$

Motivation:

- Take 2 hidden vectors from source encoder
- Combine them into a vector of size required by decoder

# Introduction to Neural Machine Translation

- Neural language models review
- Sequence to sequence models for MT
  - Encoder-Decoder
  - Sampling and search (greedy vs beam search)
  - **Practical tricks**
- **Sequence to sequence models for other NLP tasks**
- **Attention mechanism**