

Sequence to Sequence Models for Machine Translation

CMSC 723 / LING 723 / INST 725

Marine Carpuat

Slides & figure credits: Graham
Neubig

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

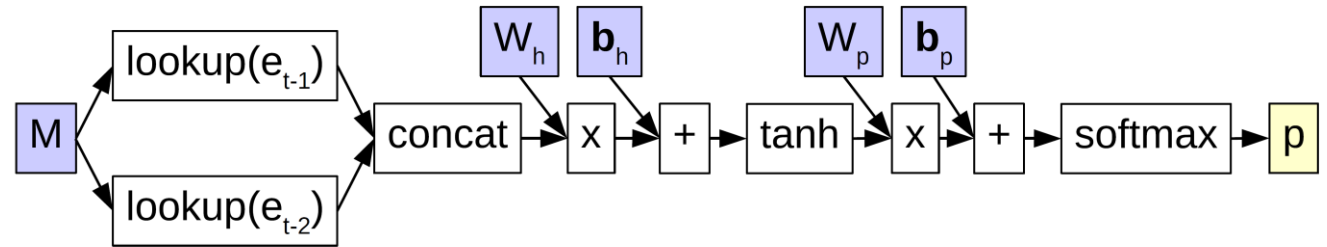
$$\hat{E} = \underset{E}{\operatorname{argmax}} P(E \mid F; \theta)$$

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

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

A feedforward neural 3-gram model



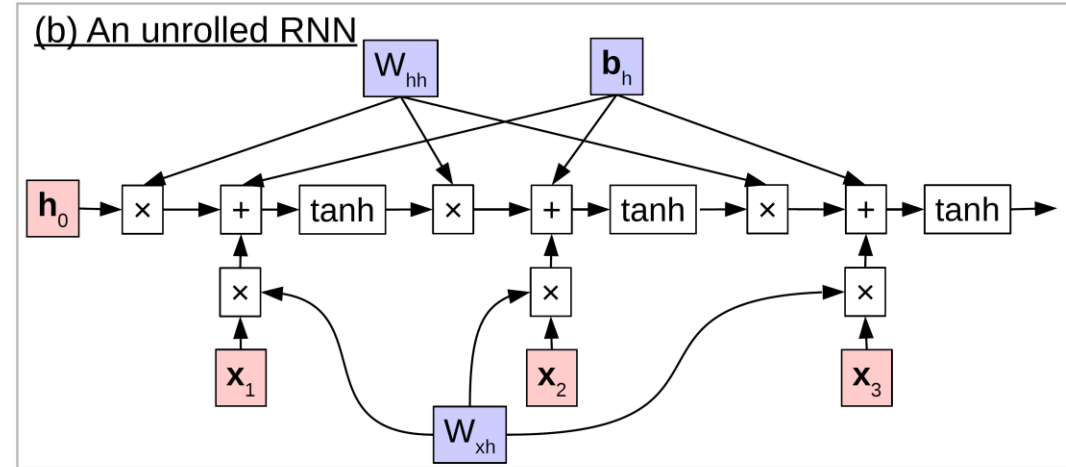
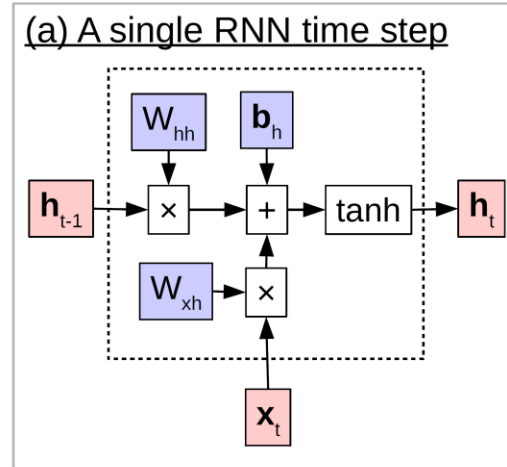
$$\mathbf{m} = \text{concat}(M_{\cdot, e_{t-2}}, M_{\cdot, e_{t-1}})$$

$$\mathbf{h} = \tanh(W_{mh}\mathbf{m} + \mathbf{b}_h)$$

$$\mathbf{s} = W_{hs}\mathbf{h} + \mathbf{b}_s$$

$$\mathbf{p} = \text{softmax}(\mathbf{s})$$

A recurrent language model

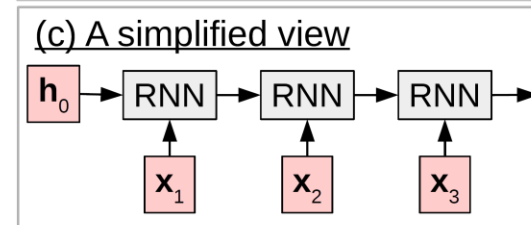
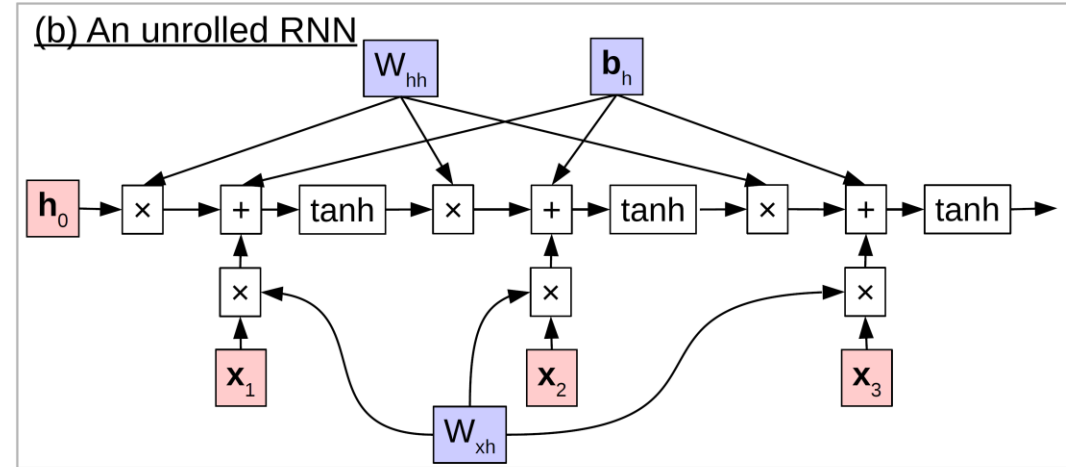
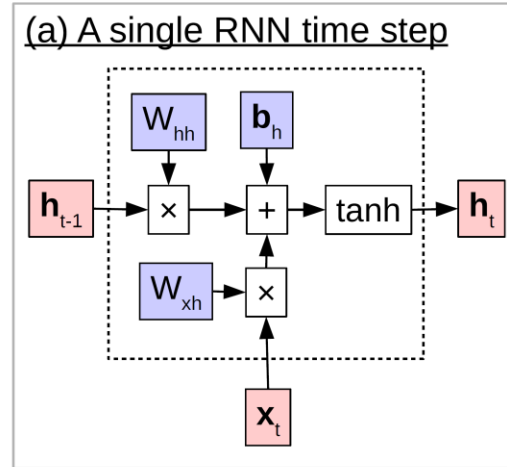


$$\mathbf{m}_t = M_{\cdot, e_{t-1}}$$

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

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

A recurrent language model



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

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

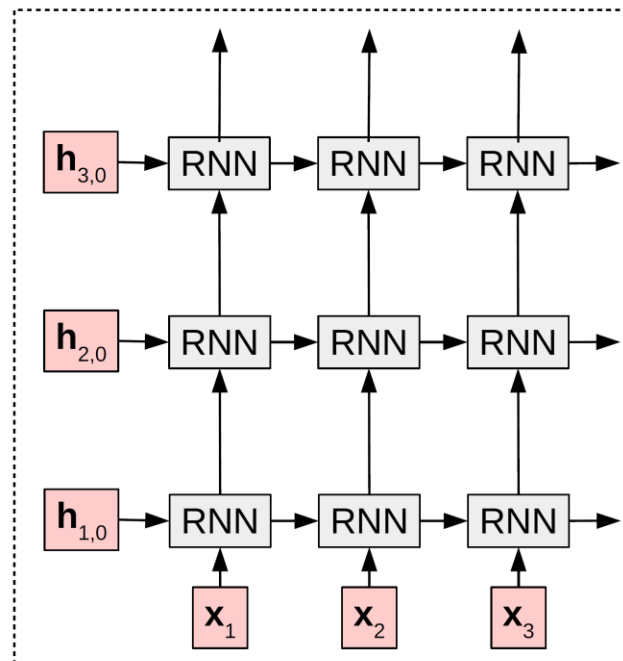
$$\mathbf{p}_t = \text{softmax}(W_{hs}\mathbf{h}_t + b_s).$$

Examples of RNN variants

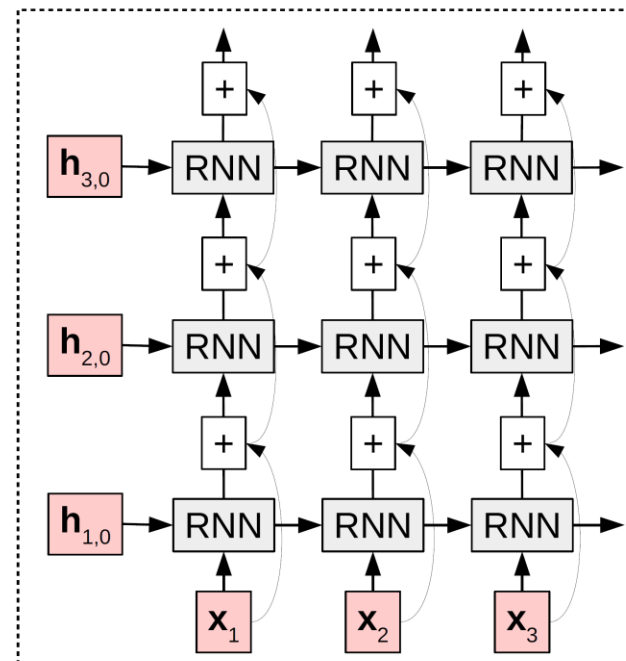
- LSTMs
 - Aim to address vanishing/exploding gradient issue

- 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

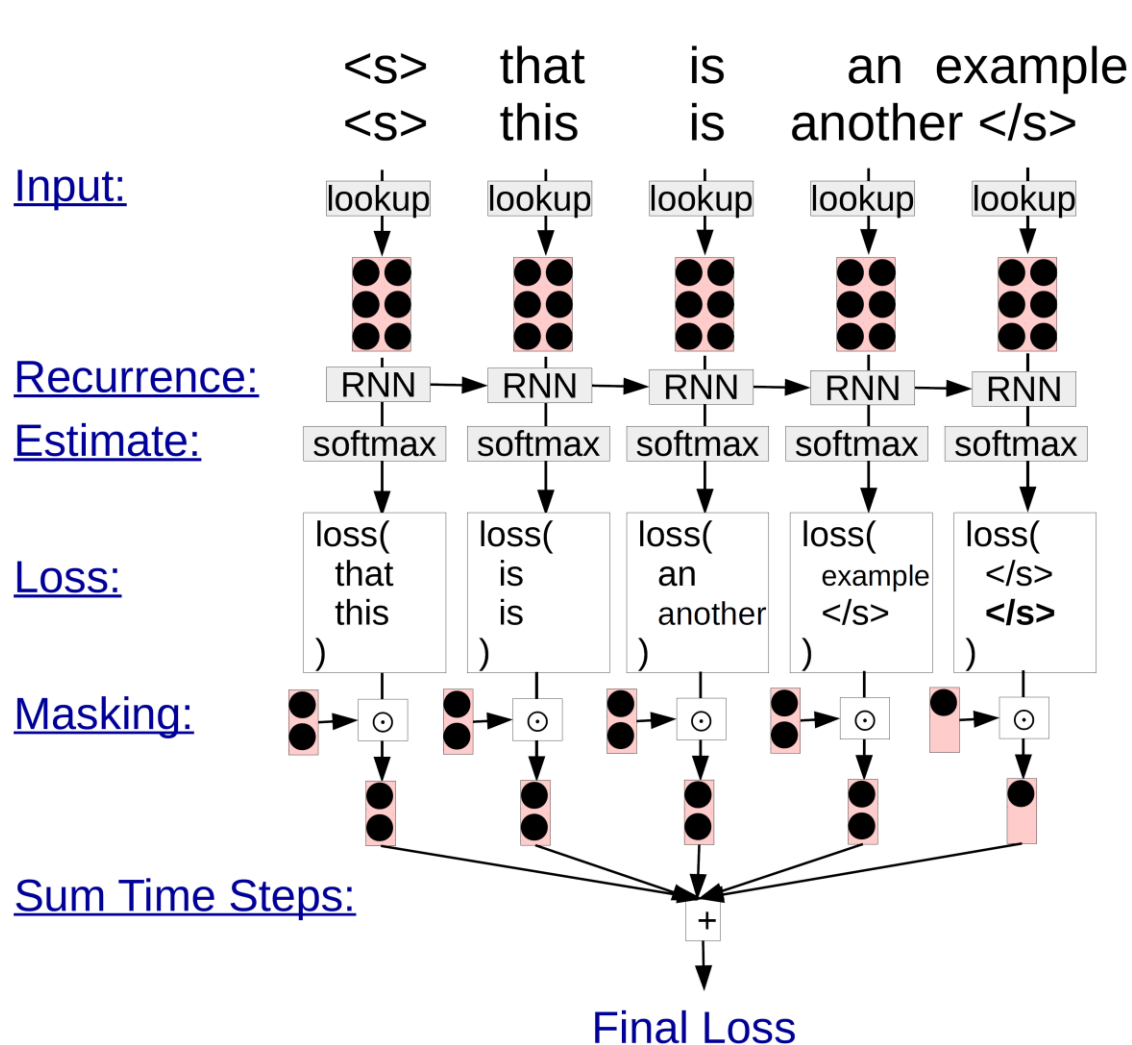
Operations w/o Minibatching

$$\tanh\left(\begin{array}{c|c|c} W & \mathbf{x}_1 & \mathbf{b} \\ \hline \begin{smallmatrix} \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{smallmatrix} & \begin{smallmatrix} \bullet \\ \bullet \\ \bullet \end{smallmatrix} & \begin{smallmatrix} \bullet \\ \bullet \\ \bullet \end{smallmatrix} \end{array}\right) \quad \tanh\left(\begin{array}{c|c|c} W & \mathbf{x}_2 & \mathbf{b} \\ \hline \begin{smallmatrix} \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{smallmatrix} & \begin{smallmatrix} \bullet \\ \bullet \\ \bullet \end{smallmatrix} & \begin{smallmatrix} \bullet \\ \bullet \\ \bullet \end{smallmatrix} \end{array}\right) \quad \tanh\left(\begin{array}{c|c|c} W & \mathbf{x}_3 & \mathbf{b} \\ \hline \begin{smallmatrix} \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{smallmatrix} & \begin{smallmatrix} \bullet \\ \bullet \\ \bullet \end{smallmatrix} & \begin{smallmatrix} \bullet \\ \bullet \\ \bullet \end{smallmatrix} \end{array}\right)$$

Operations with Minibatching

$$\begin{array}{c} \mathbf{x}_1 \quad \mathbf{x}_2 \quad \mathbf{x}_3 \rightarrow \text{concat} \rightarrow \begin{array}{c|c} W & X \\ \hline \begin{smallmatrix} \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{smallmatrix} & \begin{smallmatrix} \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{smallmatrix} \end{array} \\ \text{broadcast} \leftarrow \mathbf{b} \rightarrow \begin{array}{c|c} & B \\ \hline & \begin{smallmatrix} \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{smallmatrix} \end{array} \\ \tanh\left(\begin{array}{c|c} W & X \\ \hline \begin{smallmatrix} \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{smallmatrix} & \begin{smallmatrix} \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{smallmatrix} \end{array} + \begin{array}{c|c} & B \\ \hline & \begin{smallmatrix} \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{smallmatrix} \end{array}\right) \end{array}$$

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



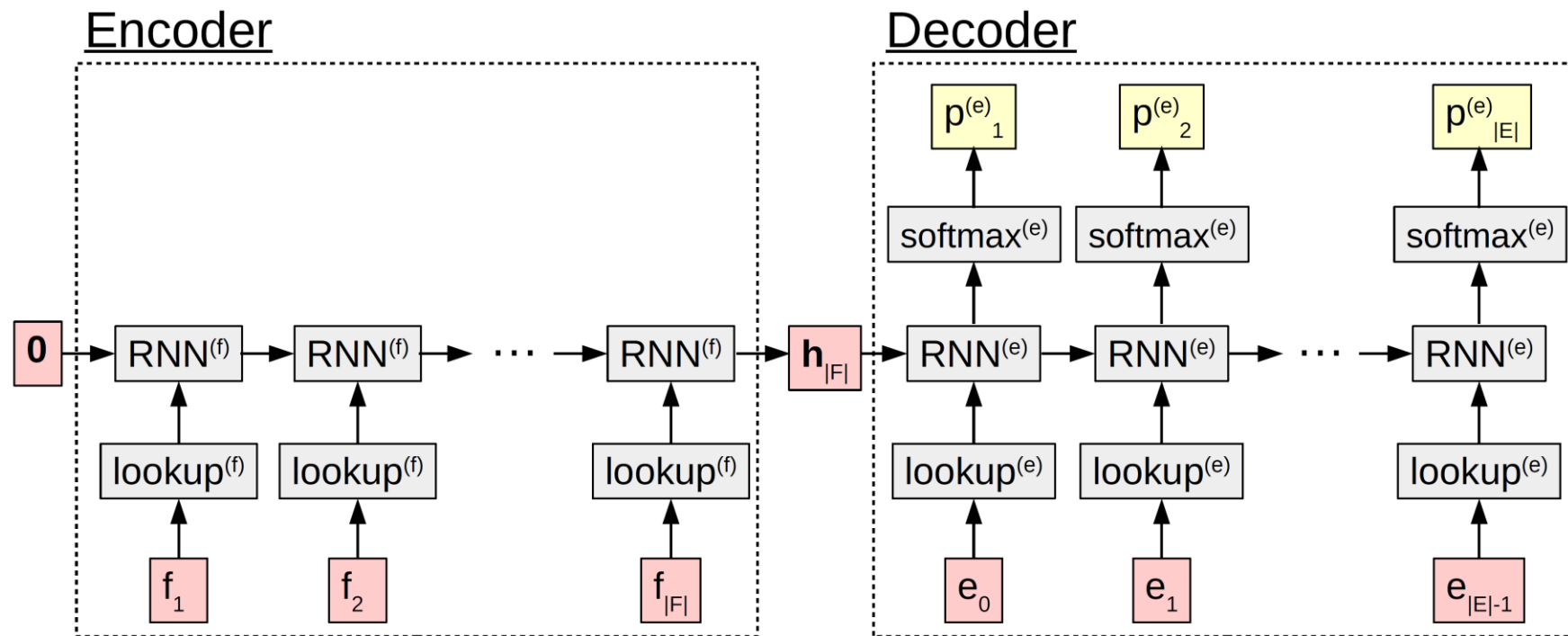
• 3 tricks

- Padding
 - Add `</s>` symbol to make all sentences same length
- Masking
 - Multiply loss function calculated over padded symbols by zero
- + sort sentences by length

Introduction to Neural Machine Translation

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

Encoder-decoder model



Encoder-decoder model

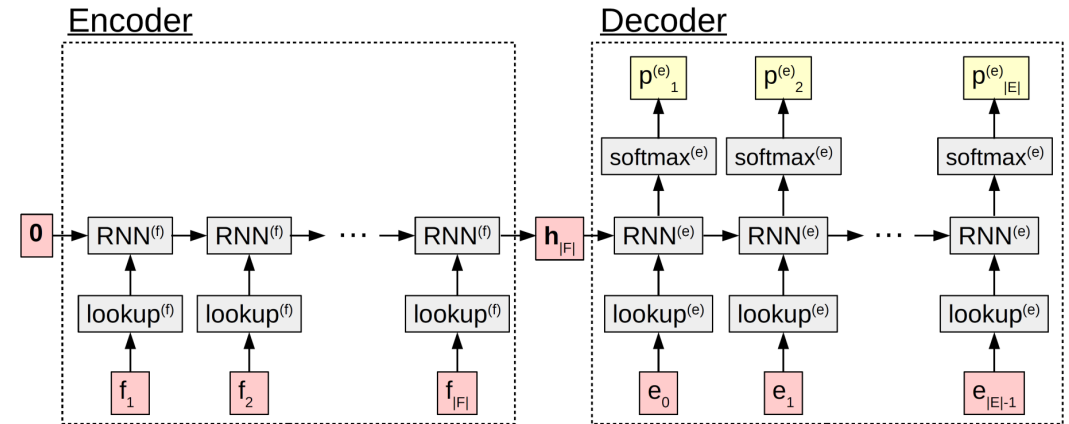
$$\mathbf{m}_t^{(f)} = M_{\cdot, f_t}^{(f)}$$

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

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

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

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



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

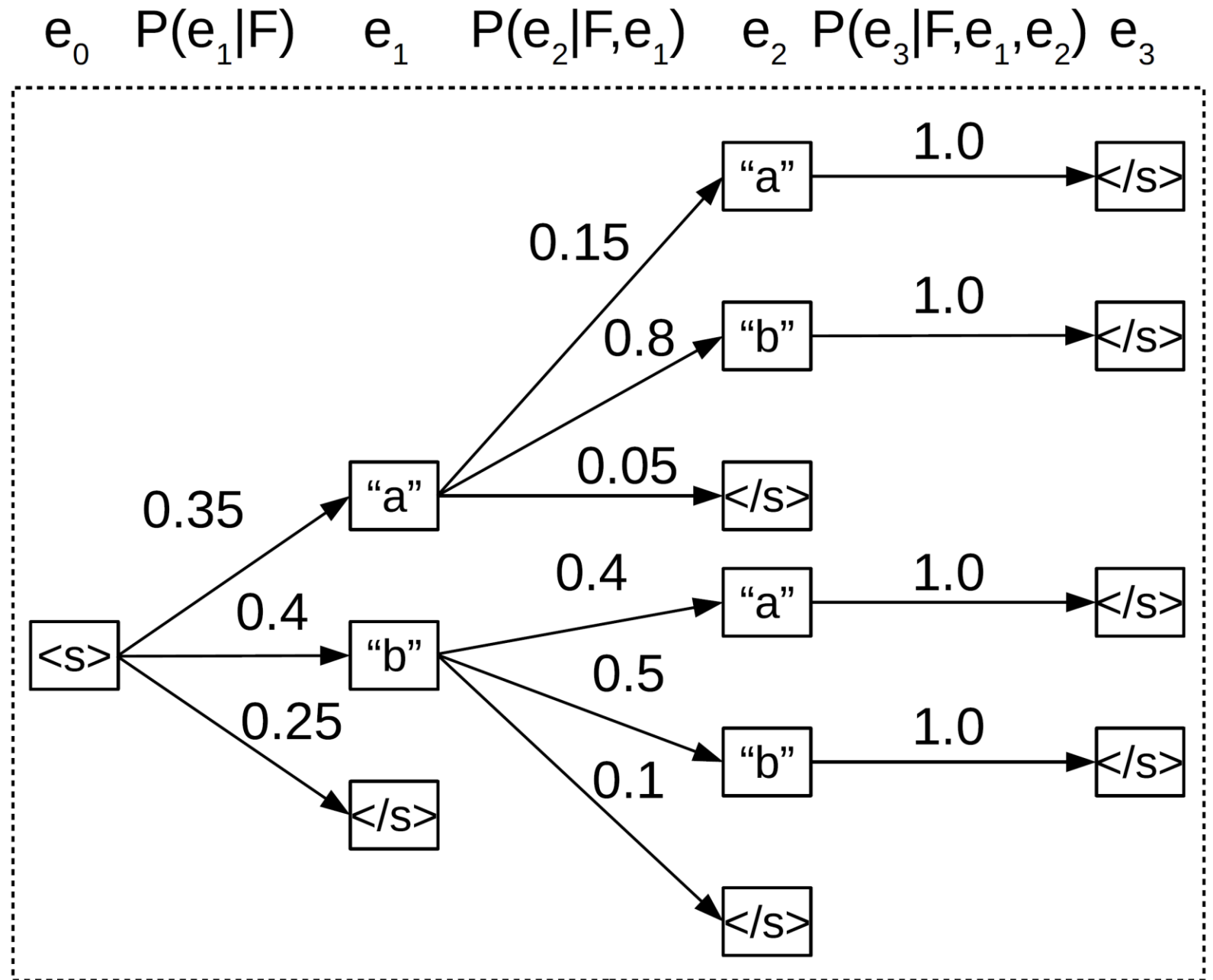
- One by one, pick single highest probability word

```
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

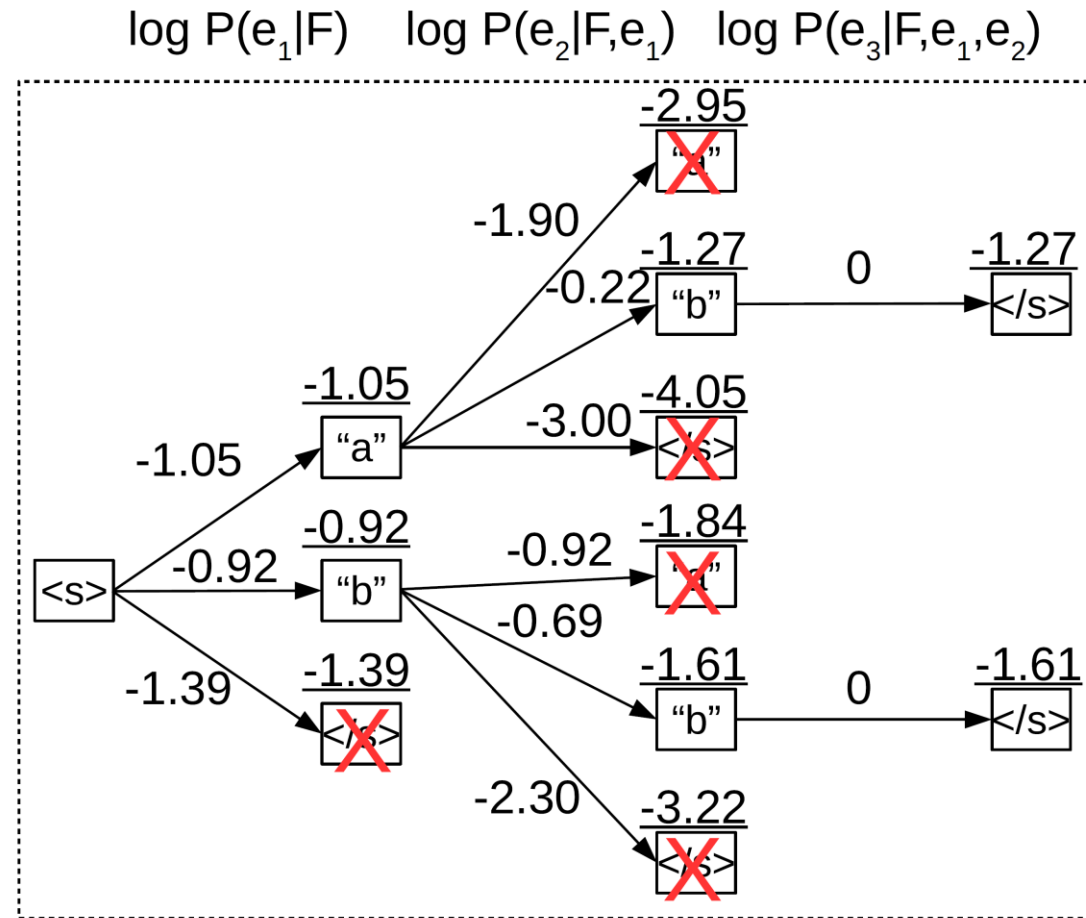
Example



Beam Search

Example with beam size $b = 2$

We consider b top hypotheses at each time step



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