

CS 378 Lecture 18

Neural MT, seq2seq models

- Today
- Seq2seq models
 - Problems w/ seq2seq models
 - Attention

- Announcements
- A3 grading
 - A4 due Tues
 - FP

Recap Phrase-based MT

Bitext \rightarrow word alignment \rightarrow phrase table $\xrightarrow{\text{LM}}$ decoder

Alignment with IBM Model 1

n target words

$$\bar{a} = (a_1, \dots, a_n) \quad \bar{t} = (t_1, \dots, t_n)$$

$$\bar{s} = (s_1, \dots, s_m, \text{NULL}) \quad m \text{ source words}$$

Model 1: $P(\bar{t}, \bar{a} | \bar{s}) = \prod_{i=1}^n P(a_i) P(t_i | s_{a_i})$

Inference: compute $P(\bar{a} | \bar{t}, \bar{s})$ (+argmax of that)
 $P(a_i | \bar{t}, \bar{s})$ proportional to $P(t_i | s_{a_i})$

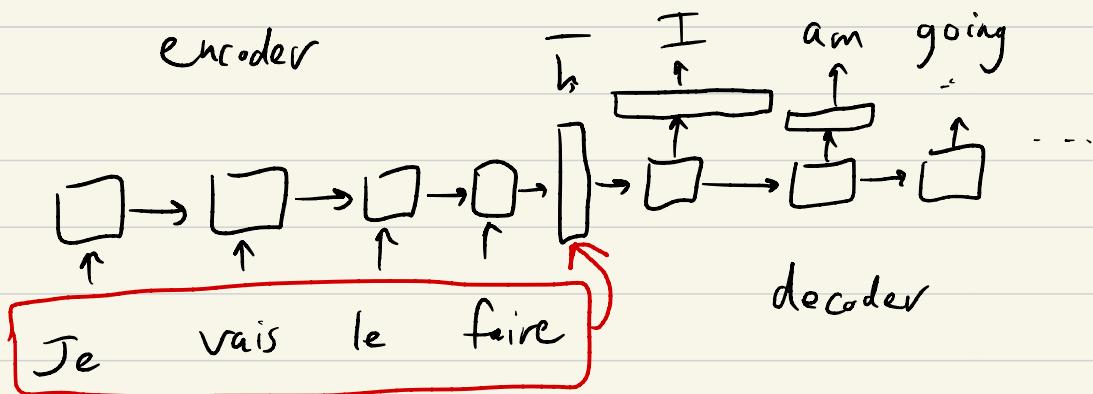
→ Look up $P(t_i | s_{a_i})$ in table for each value of a_i

→ normalize that

Seq2seq models

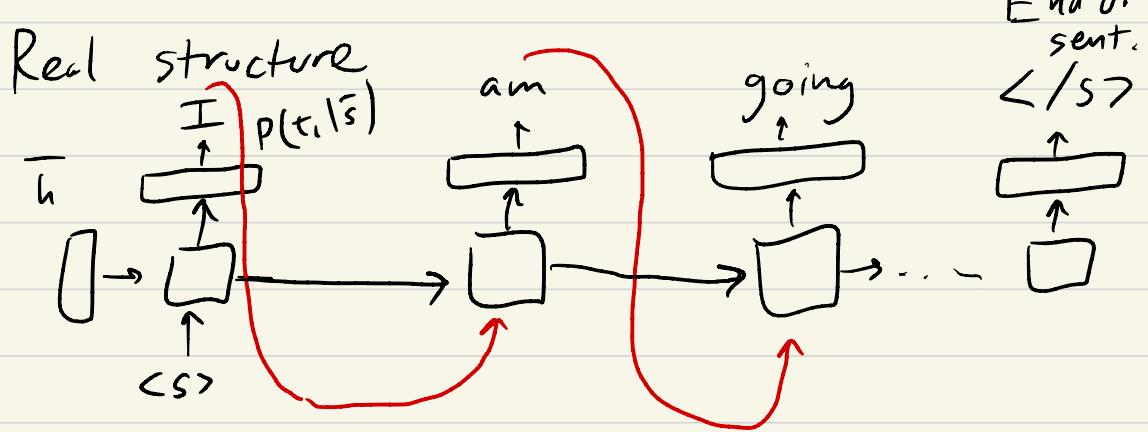
Model distribution $P(\bar{t} | \bar{s})$?

With neural models: encode \bar{s} w/RNN
"decode" \bar{t} with another RNN
initialized from the state of the first



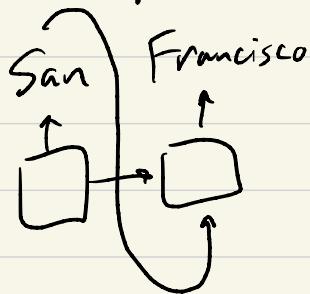
$$P(\bar{t} | \bar{s}) = P(t_1 | \bar{s}) P(t_2 | \bar{s}, t_1) \\ P(t_3 | \bar{s}, t_1, t_2) \dots$$

From \bar{h}_1 , can reconstruct the source in the target language



Feeding output t_i into step $i+1$

⇒ Model to capture word-word dependencies



bitext s^*, t^*

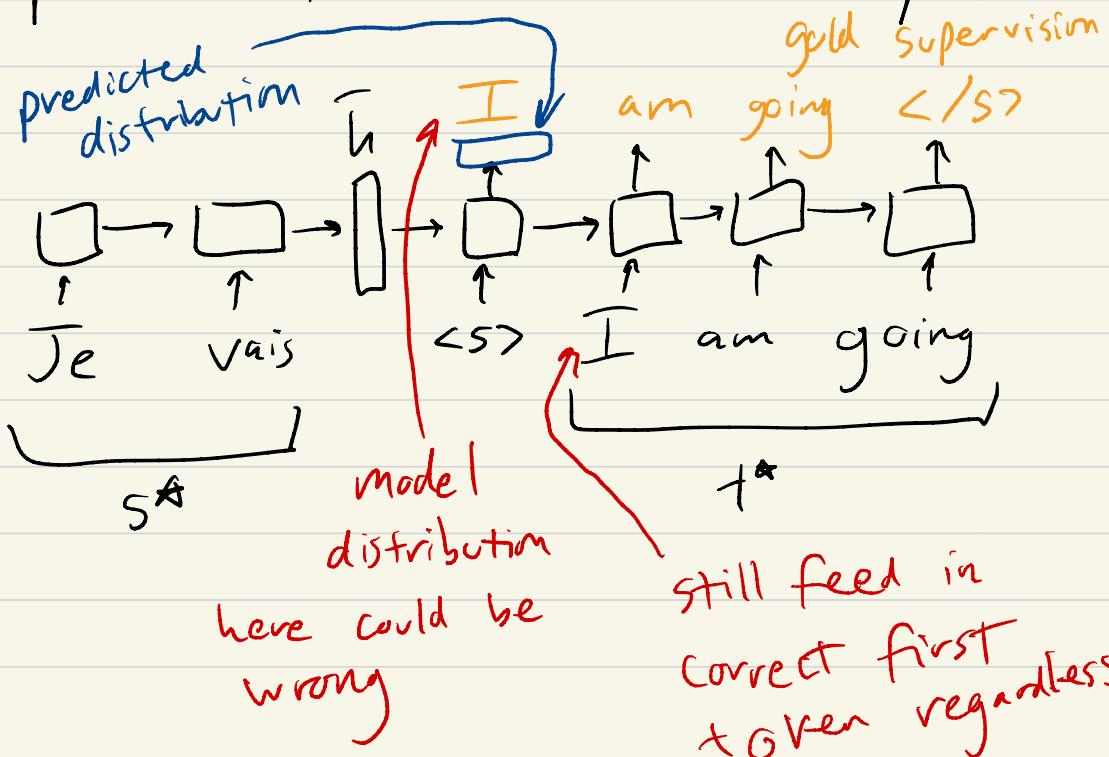
Training Given sent pairs (\bar{s}, \bar{t})

$$\text{loss} = \sum_{i=1}^n \log P(t_i^* | \bar{s}^*, t_1^*, \dots, t_{i-1}^*)$$

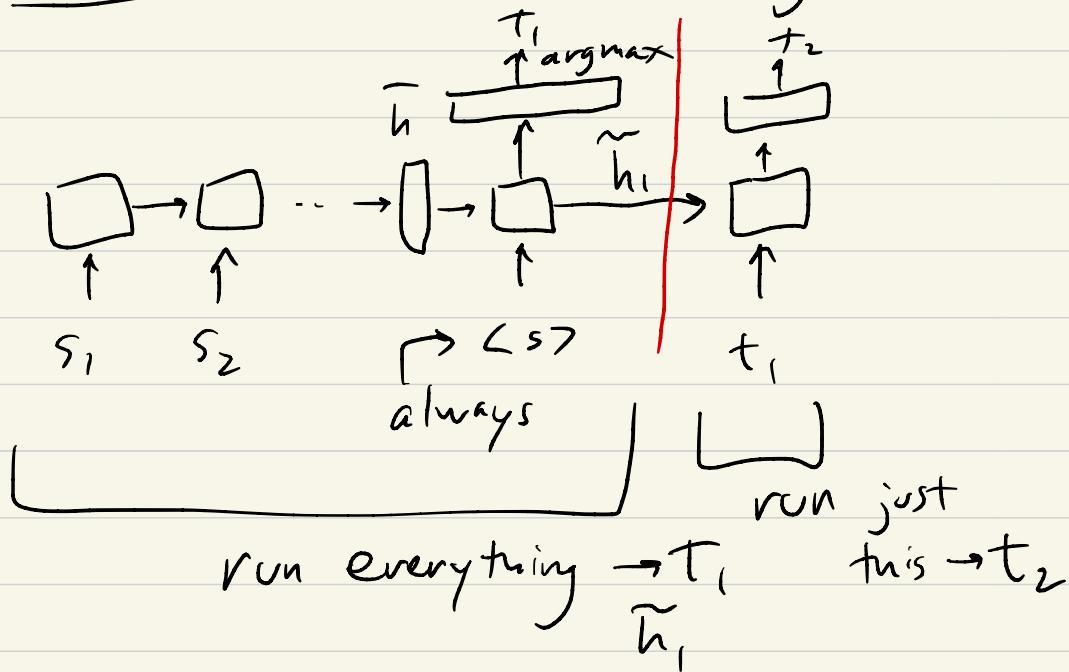
(sum over \bar{t})

Teacher Forcing

We assume everything is correct prior to the current timestep



Inference At test time: given \overline{s}



Train: single forward() + backward()

Test: n forward passes (until $\langle s \rangle$)

Save encoder state, don't rerun

Problems with seq2seq

- ① Model can repeat itself

Je fais un bureau \Rightarrow I make a desk
a desk a desk ...

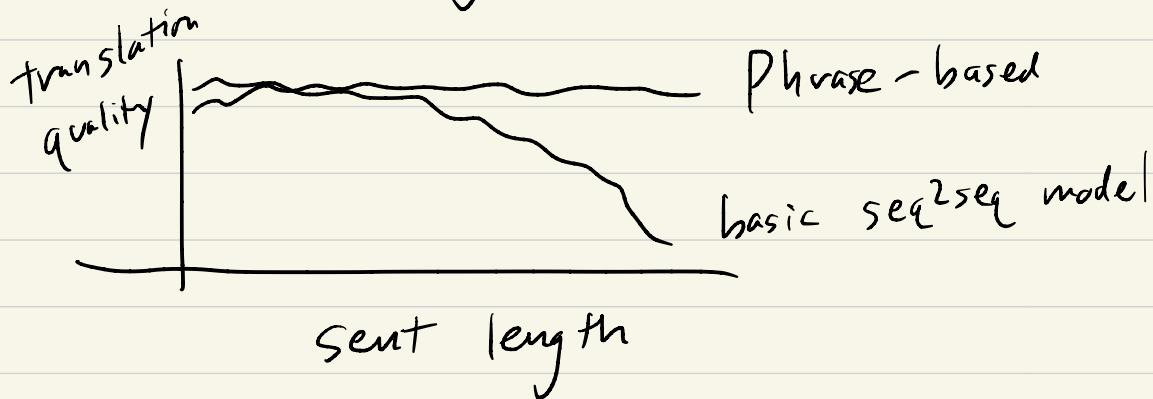
RNN doesn't have a notion of "coverage"

- ② Fixed vocabulary

Elle est allée à Pont-de-Buis \Rightarrow She
went to
decoder has $|V|$ words in it UNK
 $(\sim 30,000)$

No notion of "copy-paste"

③ Bad at long sentences



Two reasons:

- ① T_h is a fixed size
- ② Long distances are hard in LSTMs

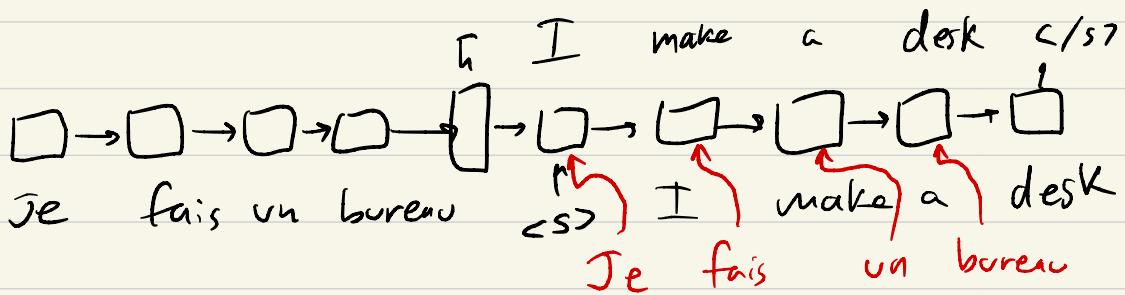
(trick: reverse the input)

$s_m \dots s_1 \rightarrow t_1 t_2 \dots$



do well on this

Attention



This is a word-by-word translation

What if...

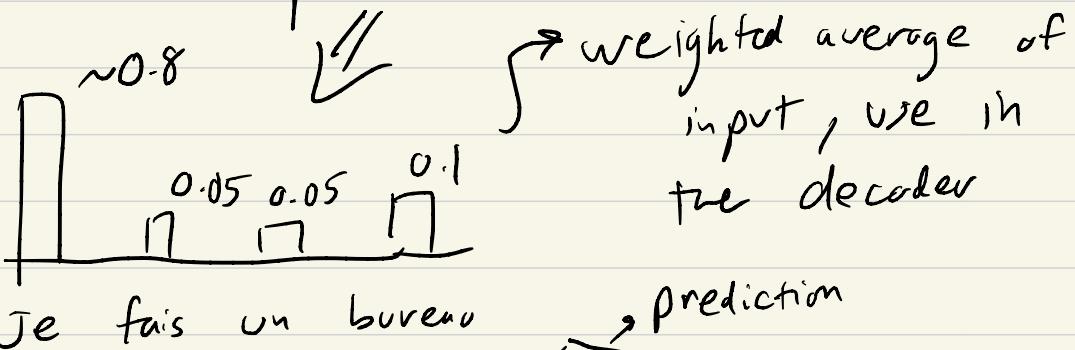
modify LSTM to take 2 inputs

Now it's easy, scales to long seqs

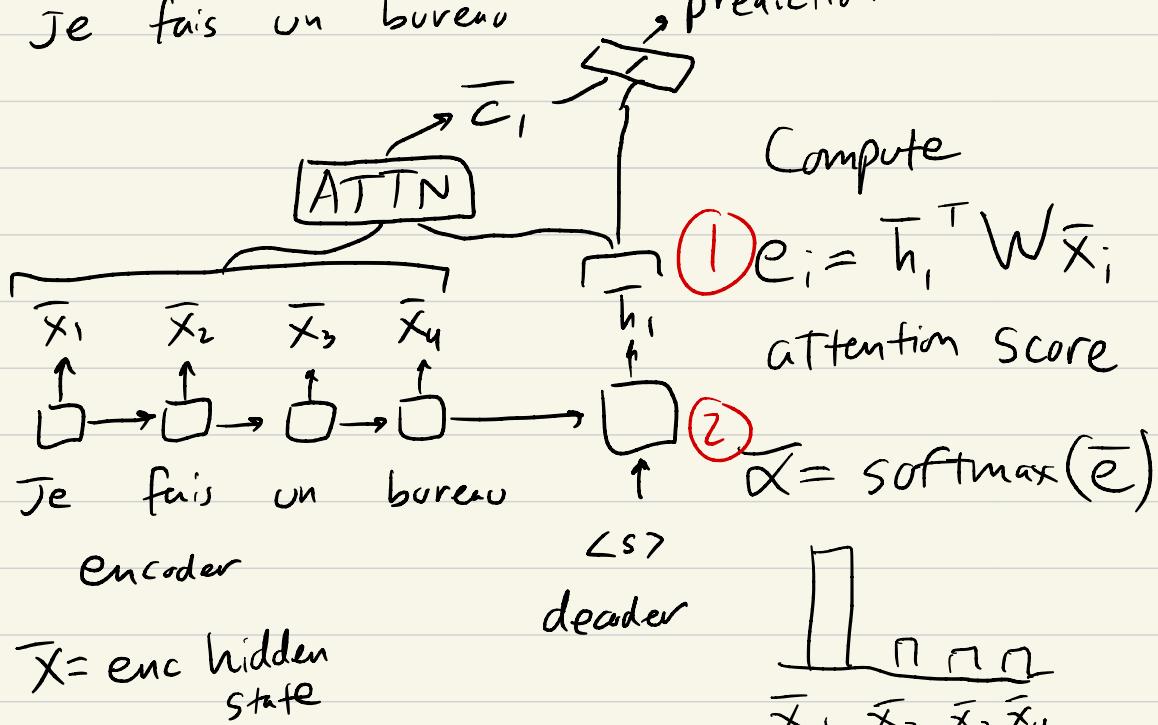
Problem: not always word-by-word
different parts of input may
be relevant

Solution: softly pick where we look
in the input

Attention = predict

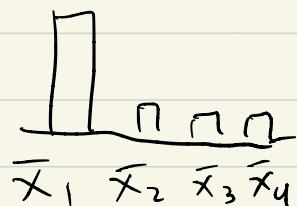


Je fais un bureau



\bar{X} = enc hidden state

decoder
↳



$$\textcircled{3} \quad \bar{c}_i = \sum_{i=1}^m \alpha_i \bar{x}_i$$

weighted sum of \bar{x}_i

Then:

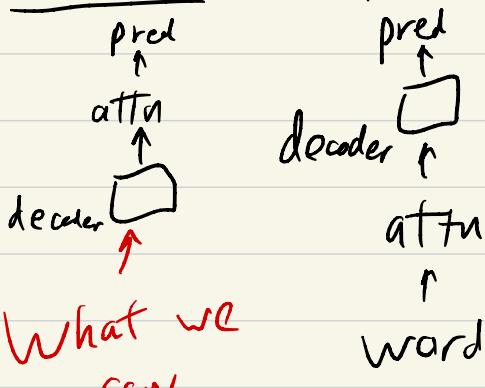
$$P(t_1 | \bar{s}) = \text{FFNN}(\bar{h}_1, \bar{c}_1)$$

\nearrow
concatenate

$\bar{c}_1 \approx "je"$ \rightarrow makes $t_1 = "I"$ easy to predict

Training: random init + backprop

Details Many ways to do this



Bahdanau

Luong