
Applied Natural Language Processing

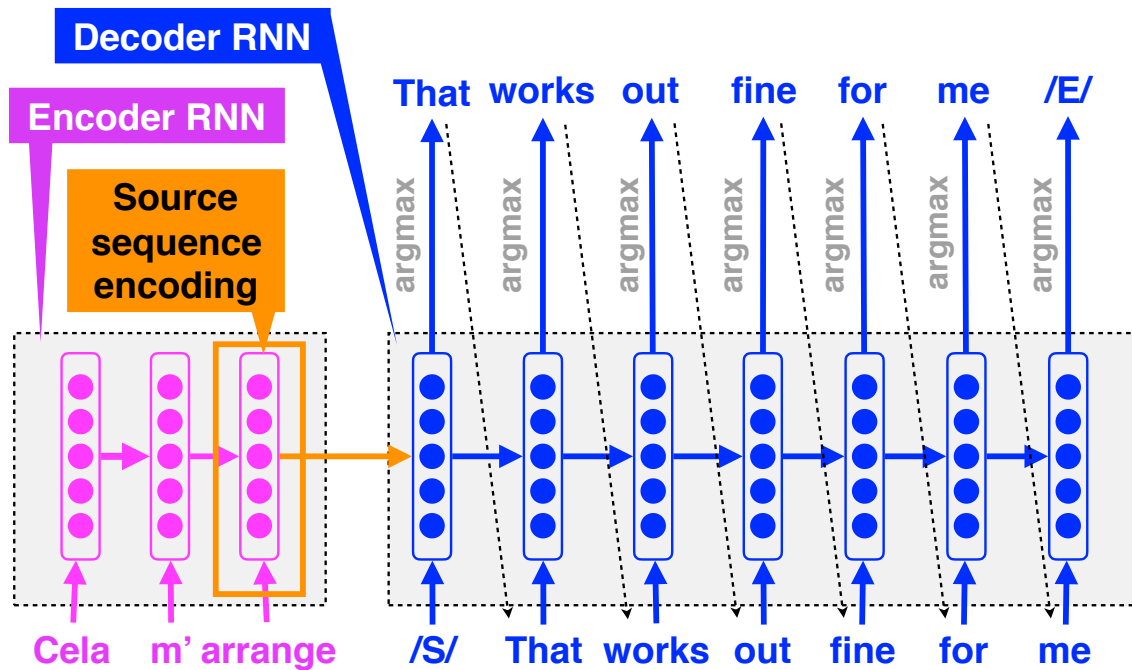
Dr. Jerome J. Braun

This Lecture: Neural Attention for NLP I

Course: Applied Natural Language Processing in Engineering
IE 7500

Unauthorized distributing, redistributing, posting, and/or reposting, of the materials of this course (including course lectures, lecture slides, slide-sets, syllabi, guidelines, assignments, quizzes, exams, presentations, software, electronic media, etc.), is prohibited.

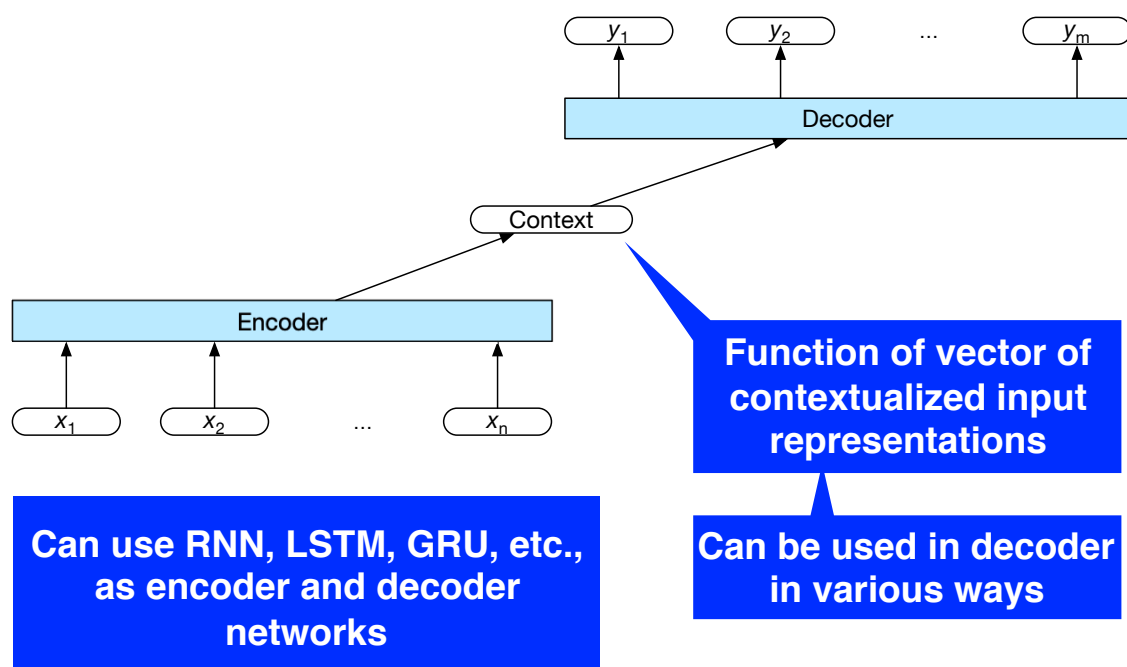
NMT Seq2seq at Test-time (Recap)



J. Braun

3

Abstract Encoder-Decoder Network



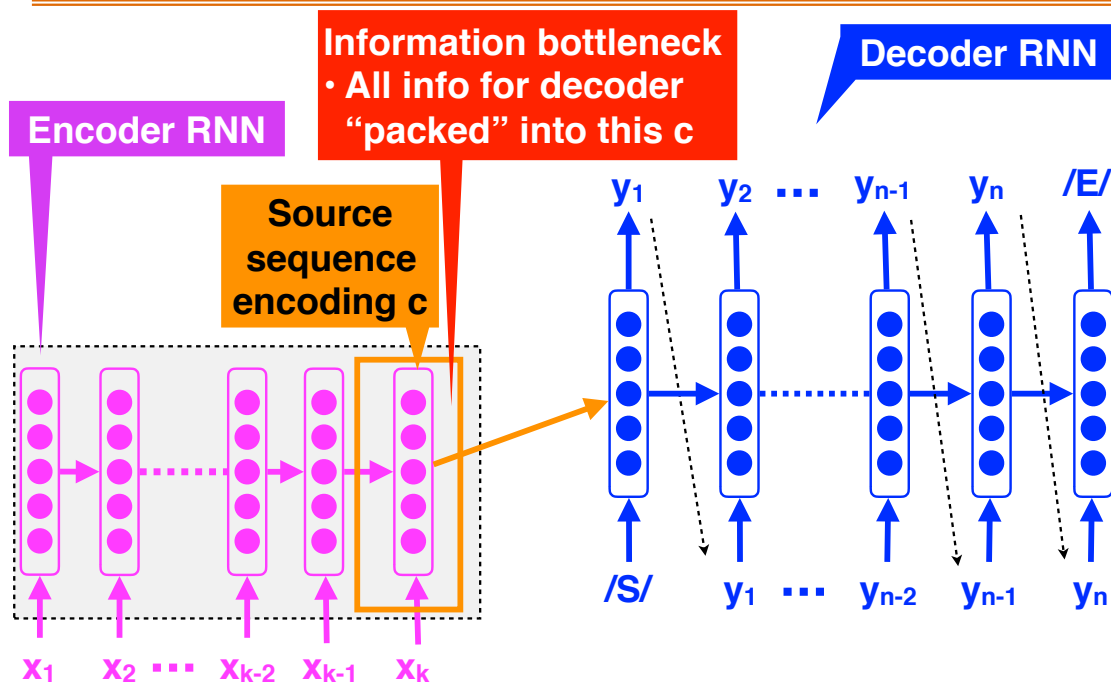
4

Context

- Context vector c should be function of encoder's hidden states, i.e., $c = f(h_1:h_n)$
 - But number of hidden states varies according to length of input
- Using encoder's final hidden state as context
 - Avoids issue of variable length of input
 - But introduces issue of excessive focus on latter parts of input sequence

5

Encoder-Decoder Information Bottleneck



Context

- **Addressing issue of excessive focus on latter parts of input sequence**
 - **Bi-RNNs — context vector formed by concatenating final states of forward and backward passes**
 - **Produce context vector by summing (or averaging) encoder hidden states**
 - ✦ **Loses info about individual encoder states that might be useful in decoding**

7

Attention Mechanism — Why

Better representation of context, for following reasons:

- **Take entire encoder context into account**
- **Update context dynamically during decoding**
- **Context remains fixed-size vector**

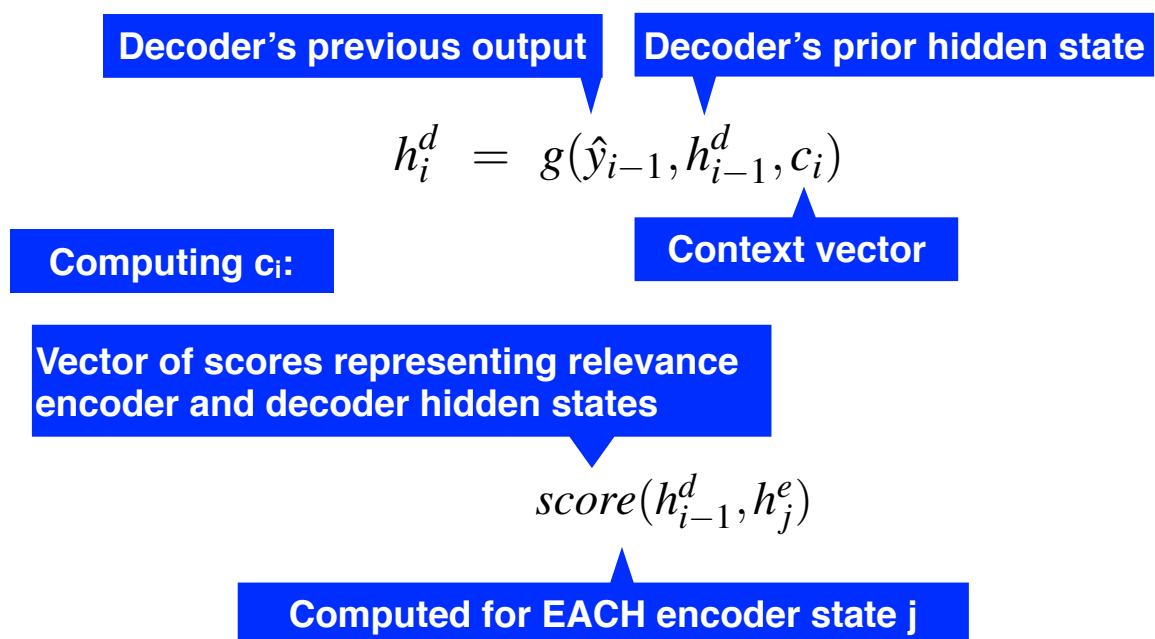
8

Attention Mechanism – How

- Static context vector replaced by dynamic construct c_i
 - Derived from encoder hidden states
 - ✦ At each point during decoding
 - New c_i generated at each decoding step i
 - ✦ This takes **all** hidden states of **encoder** into account
 - ✦ Context c_i made available during decoding
 - Current **decoder** state conditioned on all of the following
 - c_i
 - Prior hidden state
 - Previous output generated by decoder

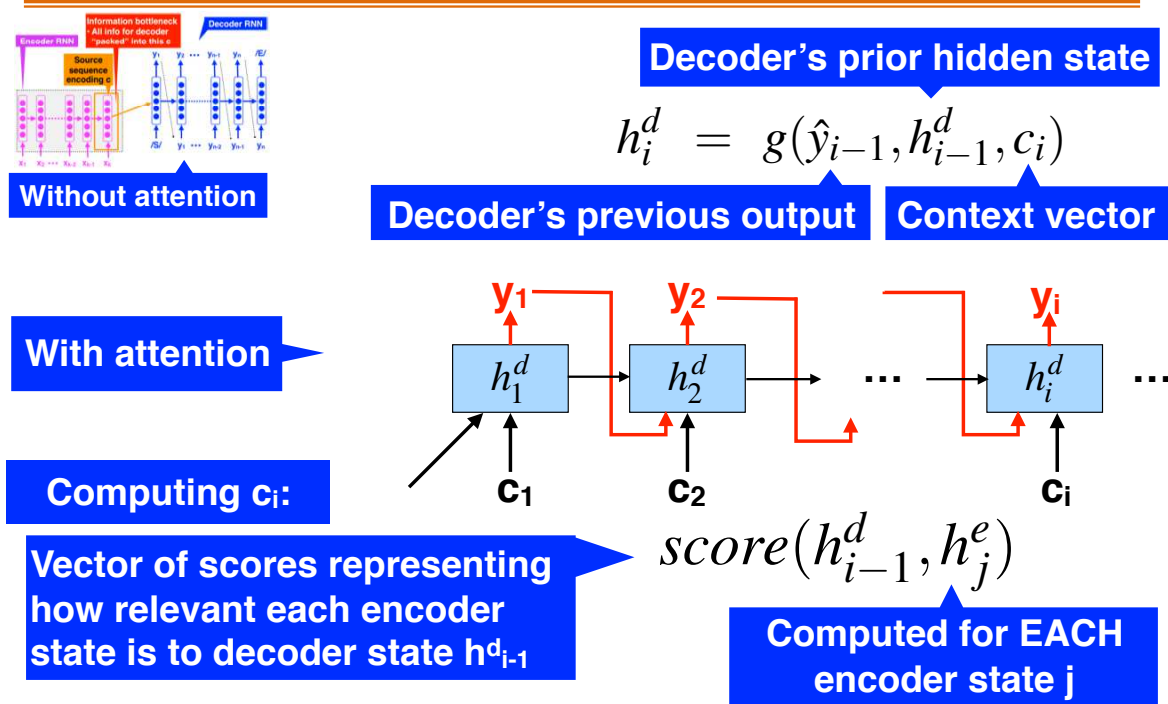
9

Attention Mechanism



10

Attention Mechanism



11

Context c_i Computation

Computing c_i :

Assume this score is:

- Measure of similarity of (between):
- Decoder hidden state
- Each hidden state of encoder

Use dot product as similarity measure

$$score(h_{i-1}^d, h_j^e) = h_{i-1}^d \cdot h_j^e$$

12

More Robust Similarity Score (1)

- Dot product is static measure — does not adapt during training
- More robust similarity score

Parametrize score with dedicated set of weights (W_s)

$$\text{score}(h_{i-1}^d, h_j^e) = h_{i-1}^d W_s h_j^e$$

Allows network to learn WHICH ASPECTS of SIMILARITY between decoder and encoder states are important to current task

13

More Robust Similarity Score (2)

Normalize parametrized scores with softmax

$$\alpha_{ij} = \text{softmax}(\text{score}(h_{i-1}^d, h_j^e) \quad \forall j \in e)$$

$$= \frac{\exp(\text{score}(h_{i-1}^d, h_j^e))}{\sum_k \exp(\text{score}(h_{i-1}^d, h_k^e))}$$

Results in vector of weights α_{ij} that express proportional relevance of each encoder hidden state j to current decoder state i

14

More Robust Similarity Score (3)

Compute fixed-length context vector

For current decoder state

Computed as weighted average over all encoder hidden states

$$c_i = \sum_j \alpha_{ij} h_j^e$$

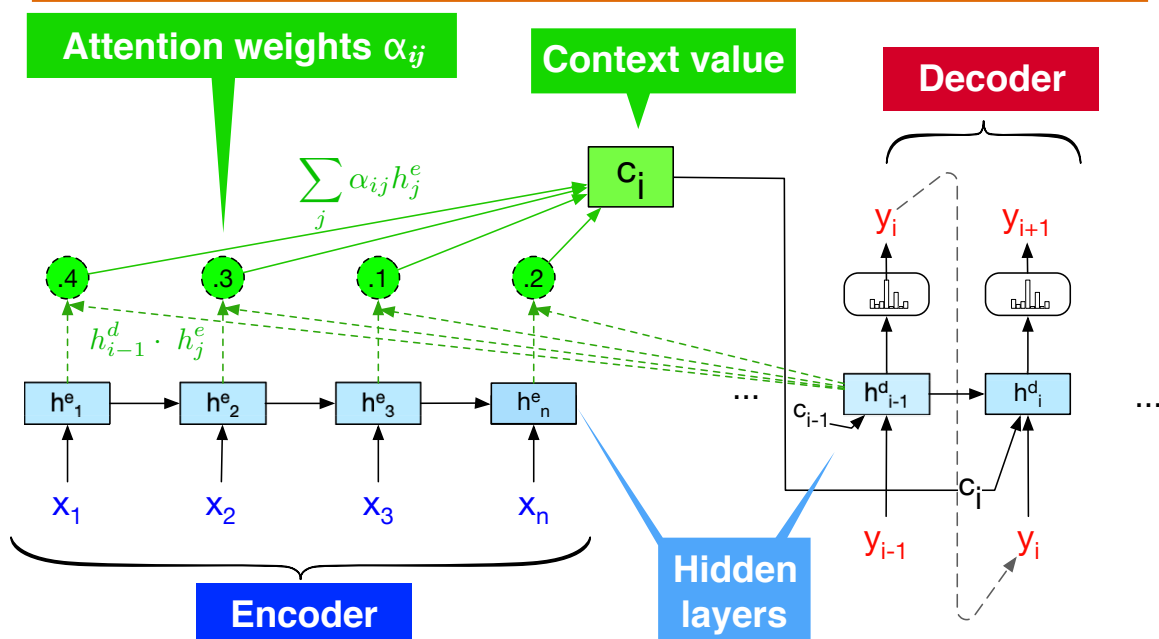
Results in vector of weights α_{ij} that express proportional relevance of each encoder hidden state j to current decoder state i

$$\alpha_{ij} = \text{softmax}(\text{score}(h_{i-1}^d, h_j^e) \quad \forall j \in e)$$

$$= \frac{\exp(\text{score}(h_{i-1}^d, h_j^e))}{\sum_k \exp(\text{score}(h_{i-1}^d, h_k^e))}$$

15

Encoder-Decoder with Attention



16



-

Questions?