Original encoder - decoder architecture for neural machine translation



"the"  "quick"  "brown"  "fox"  <EOS>    "le"  "renard"  "brun"  "rapide"

"le"  "renard"  "brun"  "rapide"  <EOS>

code vector

encoder         decoder

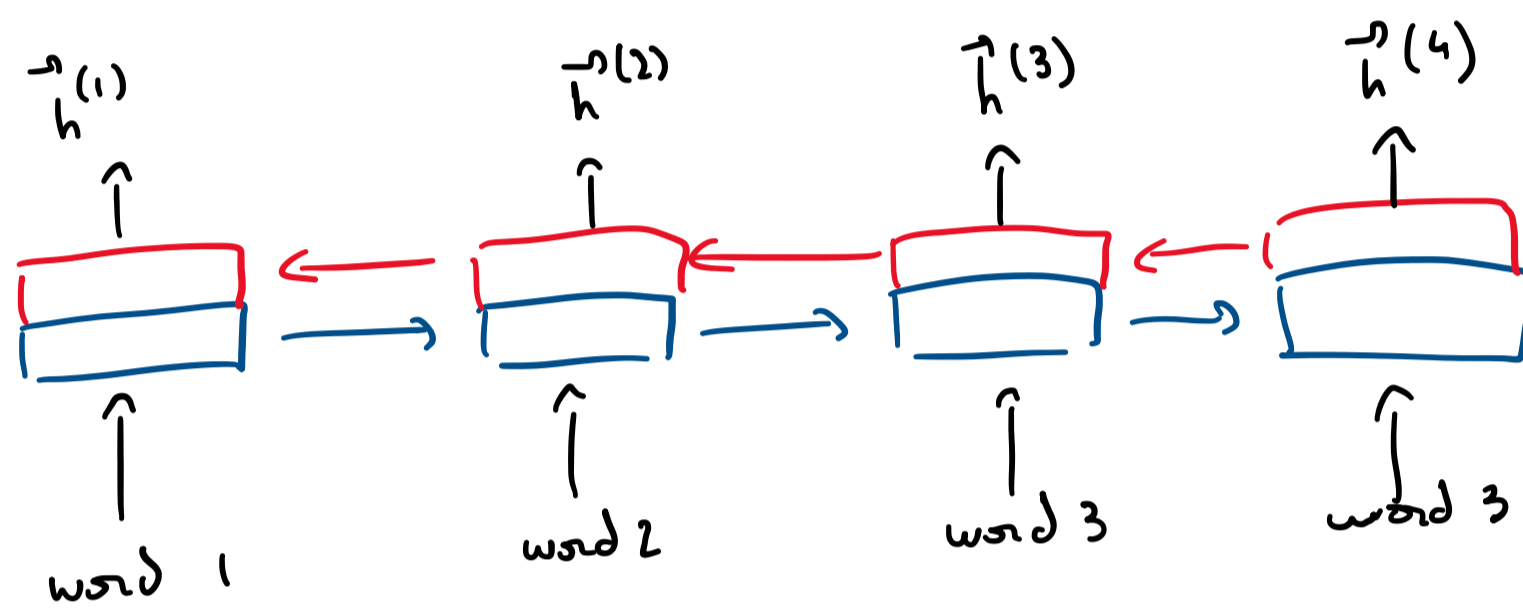memorize all relevant information about the input sentence

Attention - based modeling

→ look at input sequence (sentence) when generating text

Encoder : bidirectional RNN.
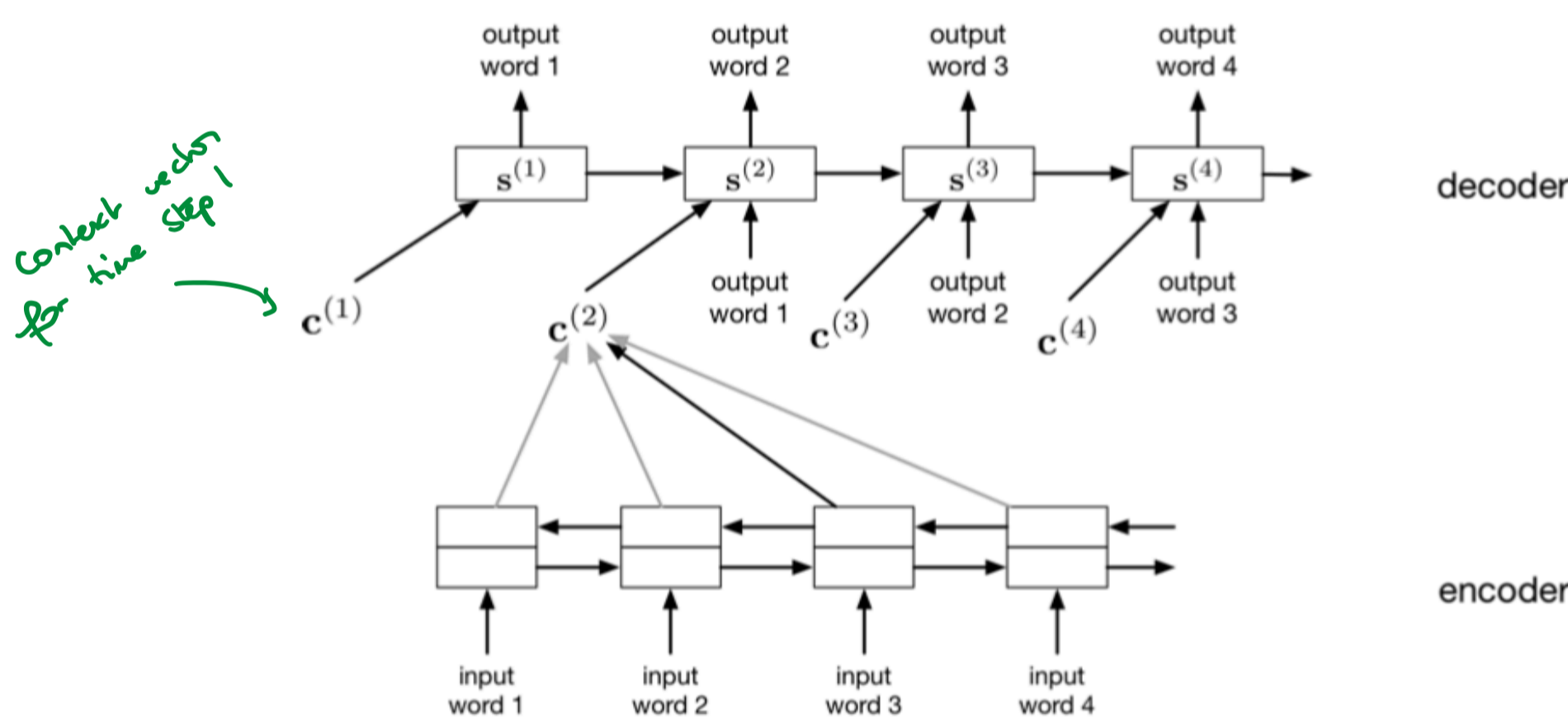
　　　Goal: compute an annotation vector
　　　　　(which the decoder will consult when it generates
　　　　　　　　the output sequence)

　　　Bidirectional : two RNNs
　　　　　　↳ (one) processes the words in the forward order
　　　　　　　↳ (second) _____ backward order



$\vec{h}^{(1)}$　　$\vec{h}^{(2)}$　　$\vec{h}^{(3)}$　　$\vec{h}^{(4)}$

word 1　　word 2　　word 3　　word 3

Annotation vector : concatenation of the hidden units of
　　　　　　　　　the 2 RNNs.

Decoder : will get a context vector in addition to
　　　　the words (generated up to current time step) as its inputs



output word 1    output word 2    output word 3    output word 4

$s^{(1)}$    $s^{(2)}$    $s^{(3)}$    $s^{(4)}$          decoder

Context vector for time step 1

$c^{(1)}$    $c^{(2)}$    output word 1  $c^{(3)}$  output word 2  $c^{(4)}$  output word 3

input word 1    input word 2    input word 3    input word 4          encoder

Soft - attention model :
　　$\vec{c}^{(i)}$ :　context vector for output time step $i$ (decoder)
　　$t$ :　index of time steps in the input (encoder)

$$\vec{c}^{(i)} = \sum_t \alpha_{it} \vec{h}^{(t)}$$

annotation vector

attention weight

$$\alpha_{it} = \frac{e^{z_{it}}}{\sum_{t'} e^{z_{it'}}} \quad \text{where} \quad z_{it} = a\left(\vec{s}^{(i-1)}, \vec{h}^{(t)}\right)$$

previous decoder hidden state

feed forward neural network

annotation vector (used for content-based addressing)