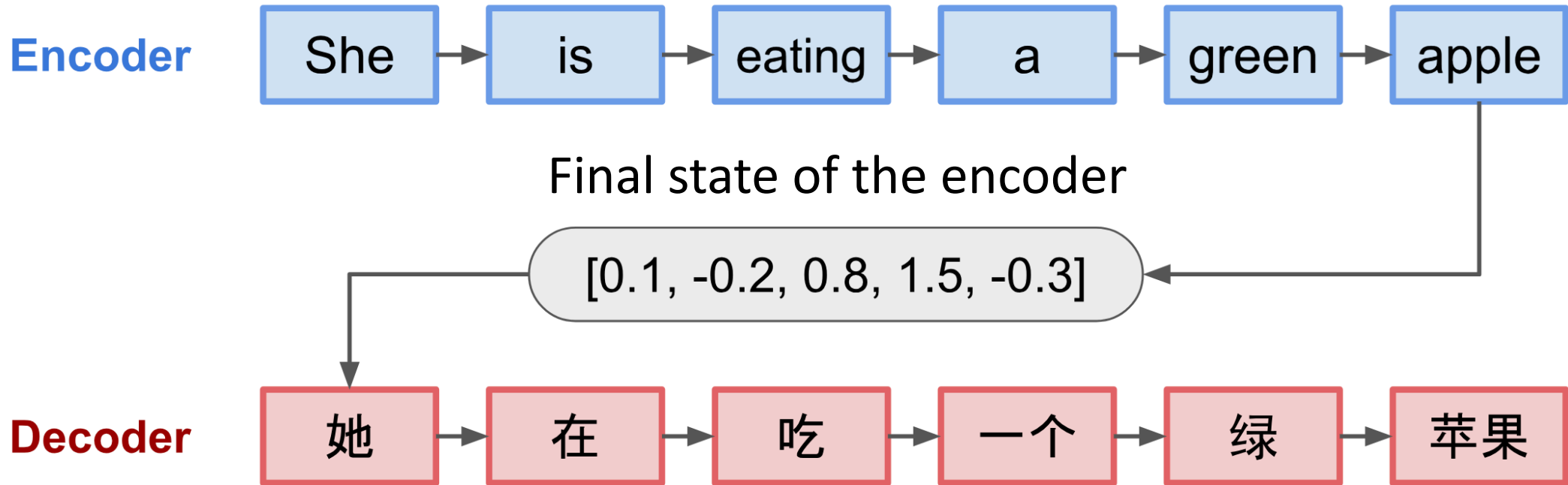


Attention

Shusen Wang

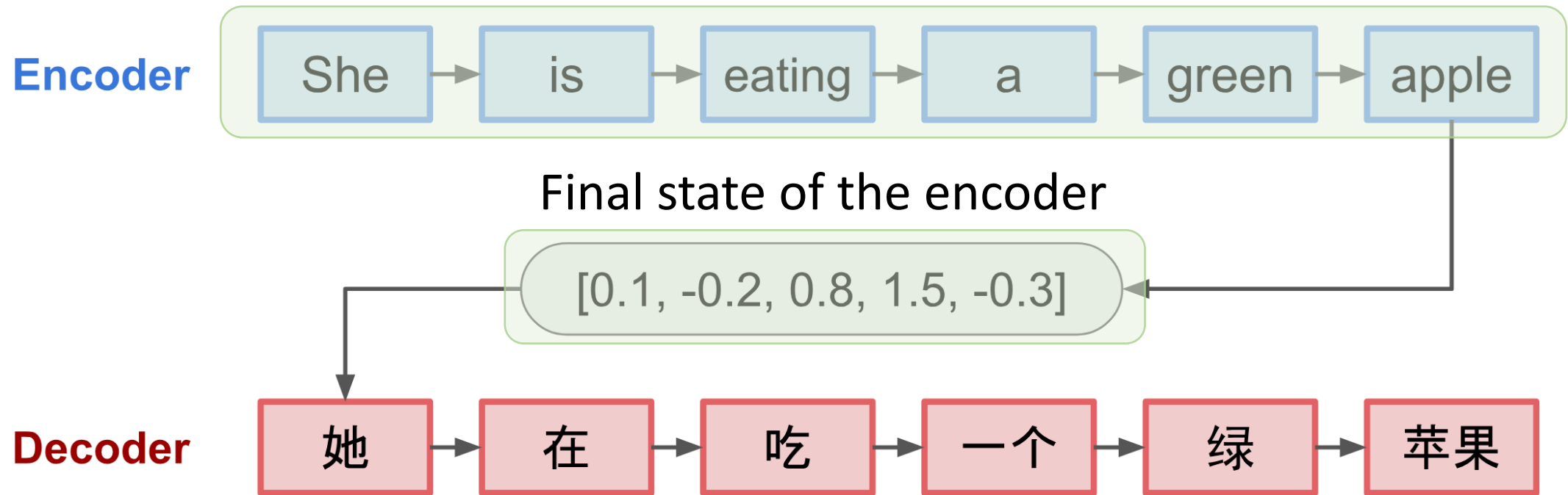
Seq2Seq Model



The figure is from blog.lilianweng.github.io

Seq2Seq Model

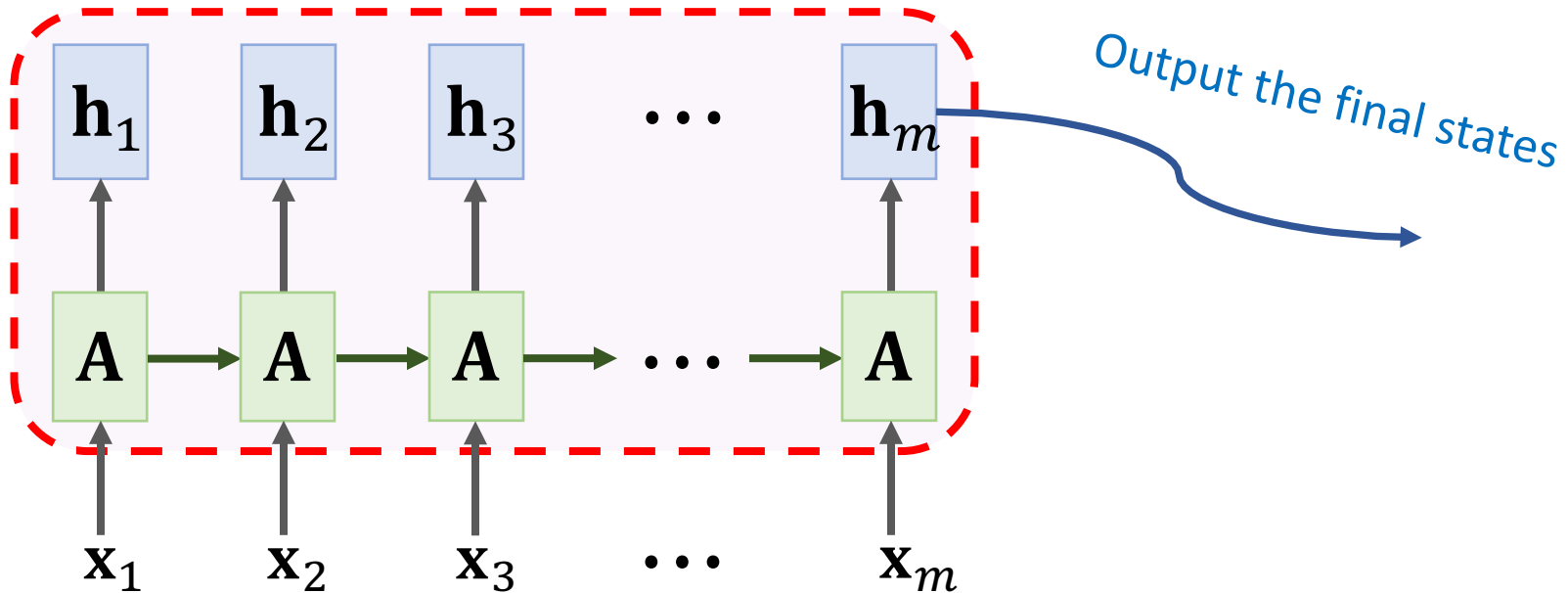
Shortcoming: The final state is incapable of remembering a **long** sequence.



The figure is from blog.lilianweng.github.io

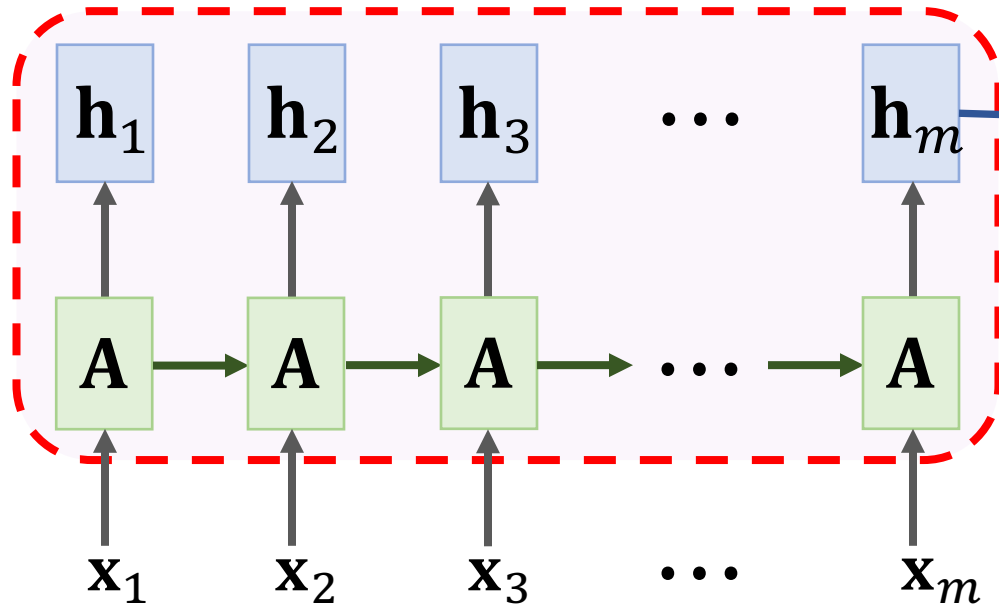
Seq2Seq Model

Encoder RNN

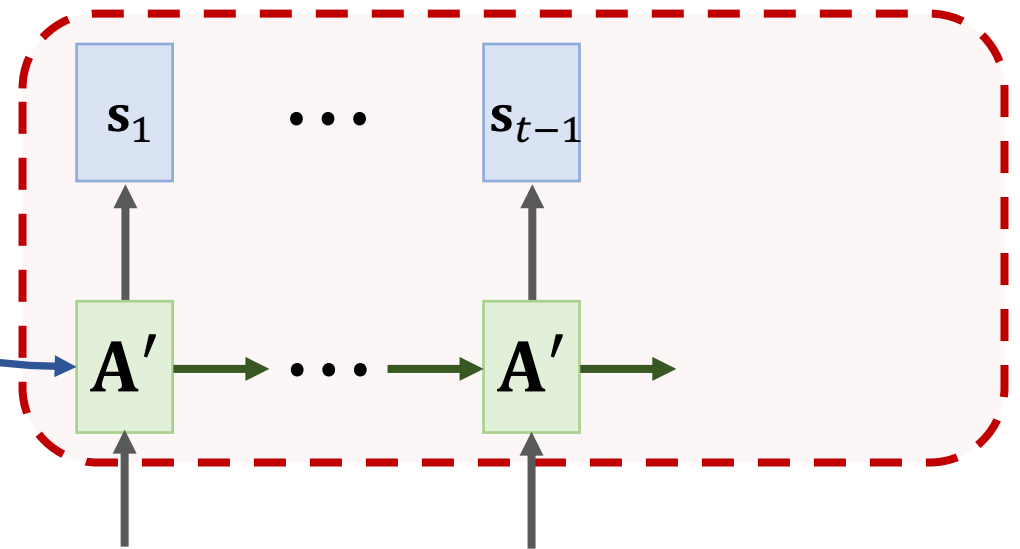


Seq2Seq Model

Encoder RNN



Decoder RNN



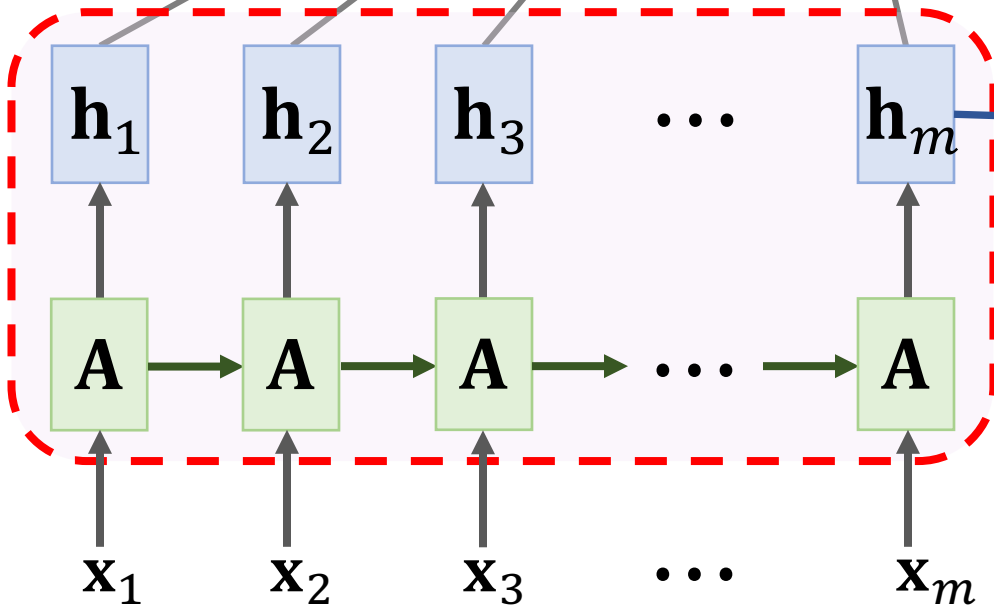
Attention

context vector

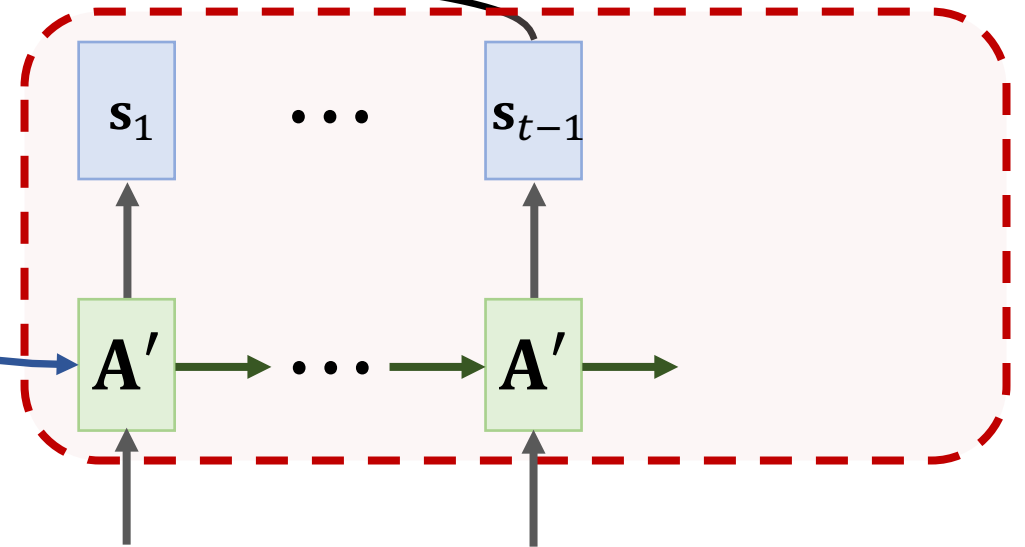
$$\mathbf{c} = \sum_{i=1}^m \alpha_i \mathbf{h}_i.$$

- α_i : similarity between \mathbf{s}_{t-1} and \mathbf{h}_i .

Encoder RNN



Decoder RNN



Attention

context vector

$$\mathbf{c} = \sum_{i=1}^m \alpha_i \mathbf{h}_i.$$

- α_i : similarity between \mathbf{s}_{t-1} and \mathbf{h}_i .

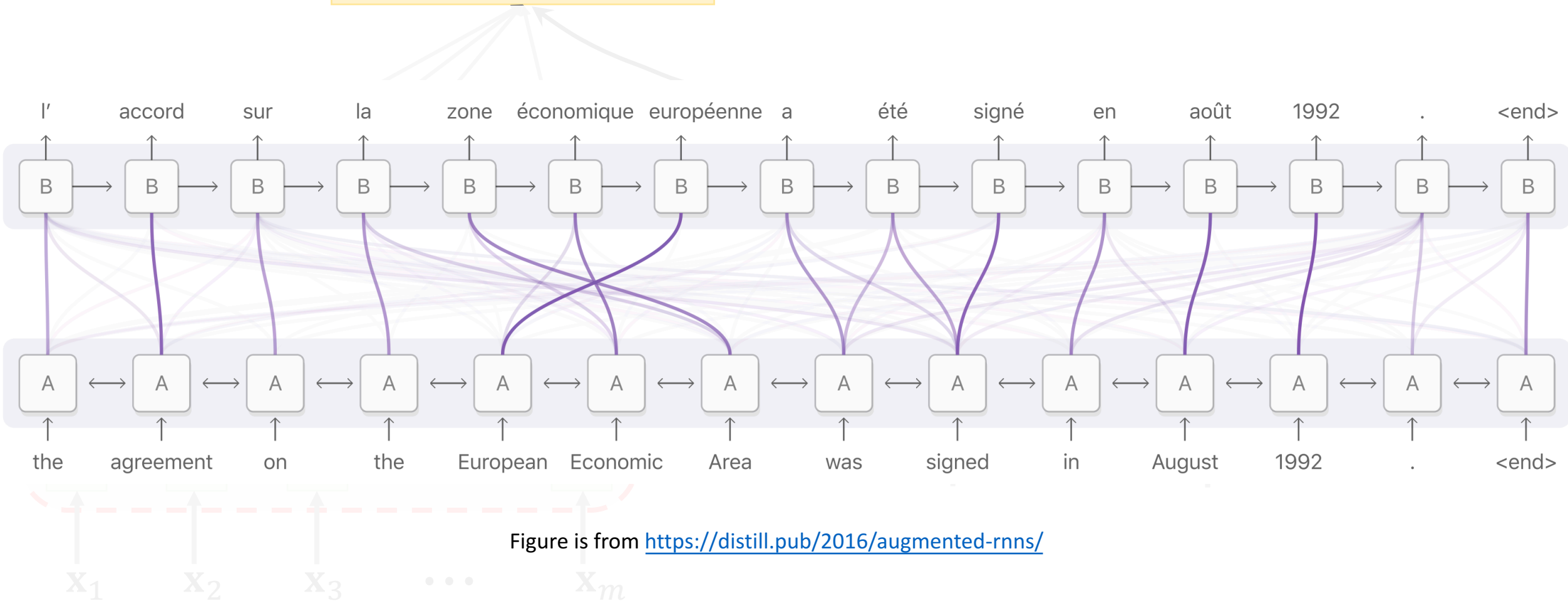


Figure is from <https://distill.pub/2016/augmented-rnns/>

Attention

context vector

$$\mathbf{c} = \sum_{i=1}^m \alpha_i \mathbf{h}_i.$$

- α_i : similarity between \mathbf{s}_{t-1} and \mathbf{h}_i .

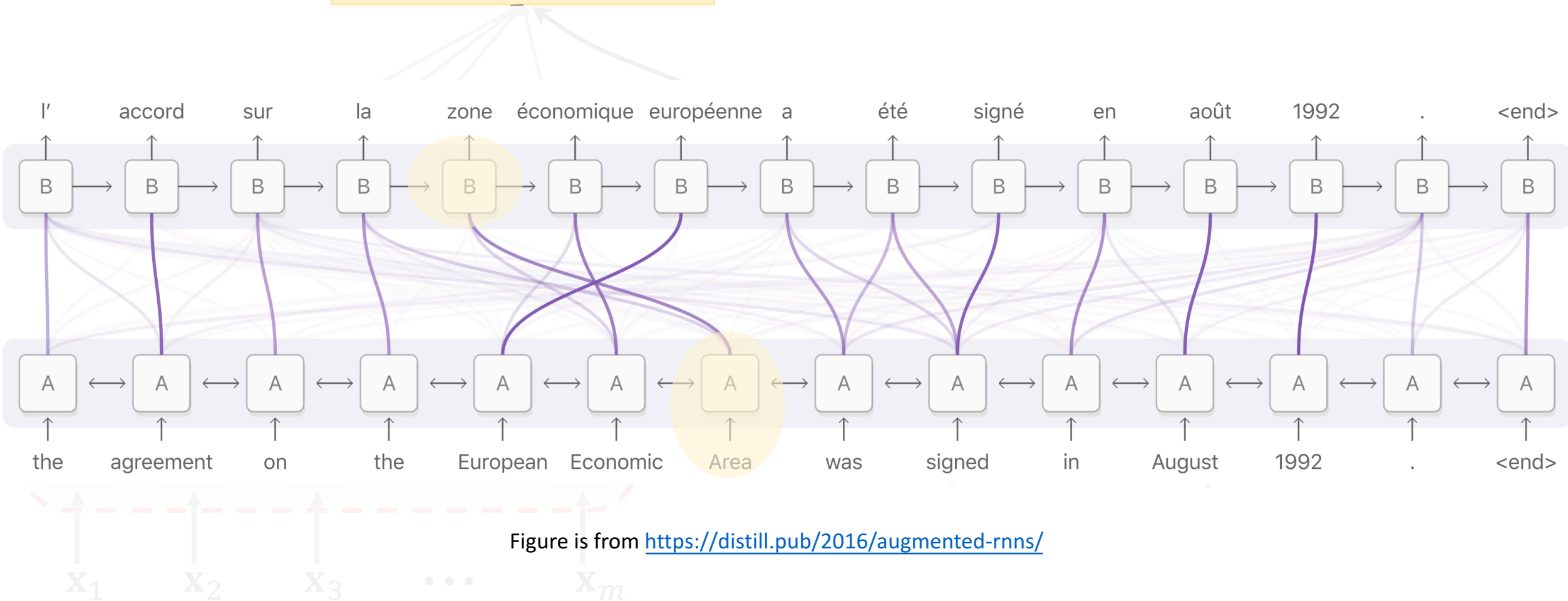


Figure is from <https://distill.pub/2016/augmented-rnns/>

Attention

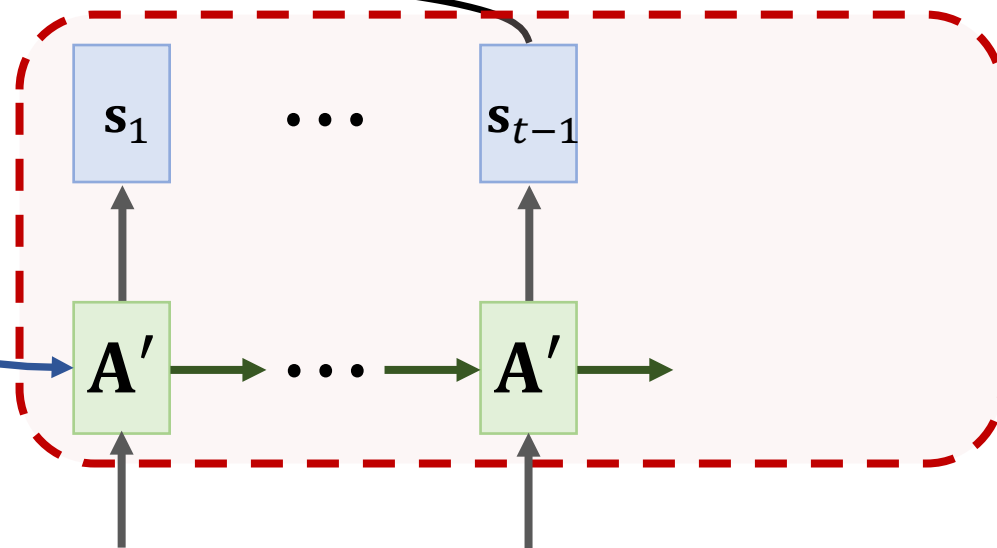
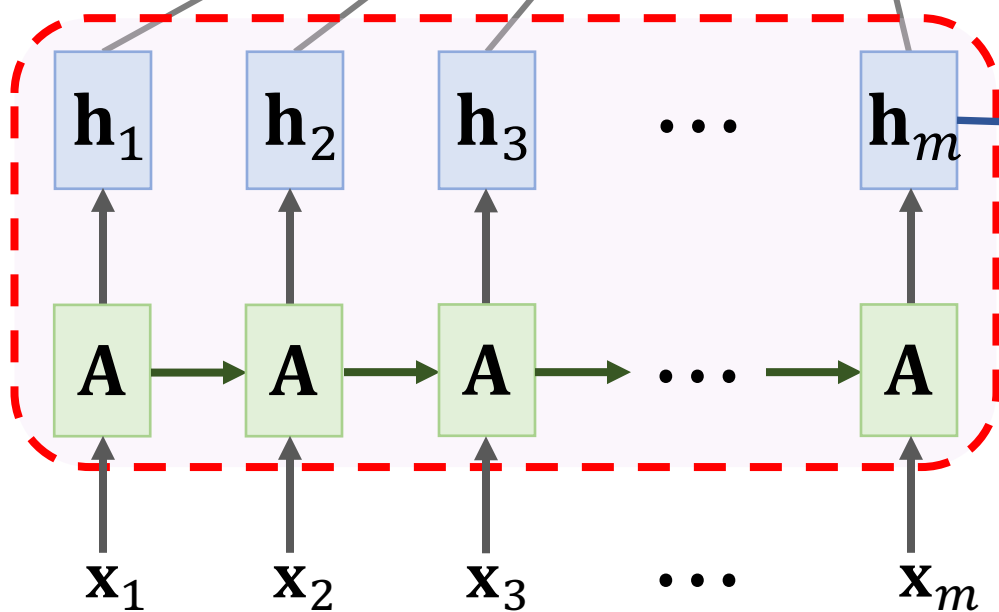
context vector

$$\mathbf{c} = \sum_{i=1}^m \alpha_i \mathbf{h}_i.$$

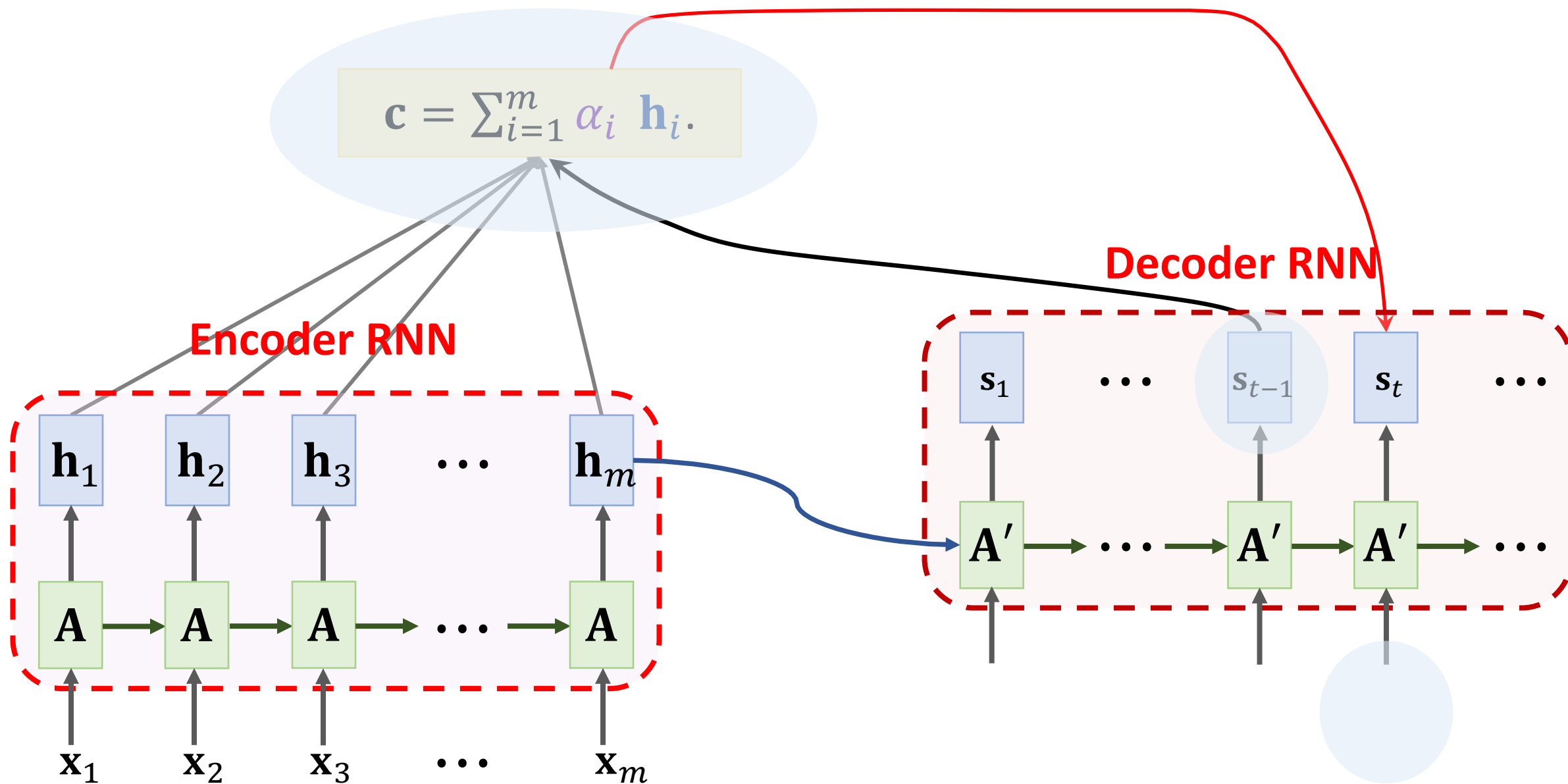
- α_i : similarity between \mathbf{s}_{t-1} and \mathbf{h}_i .
- α_i is computed by a neural network taking \mathbf{s}_{t-1} and \mathbf{h}_i as input.

Decoder RNN

Encoder RNN

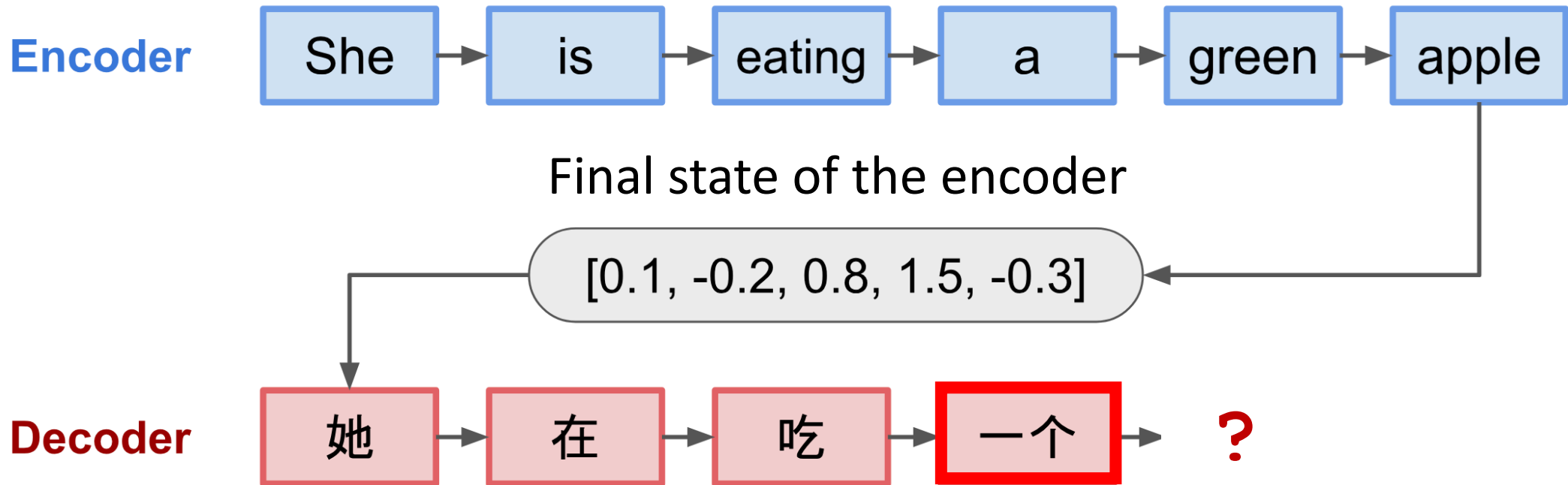


Attention



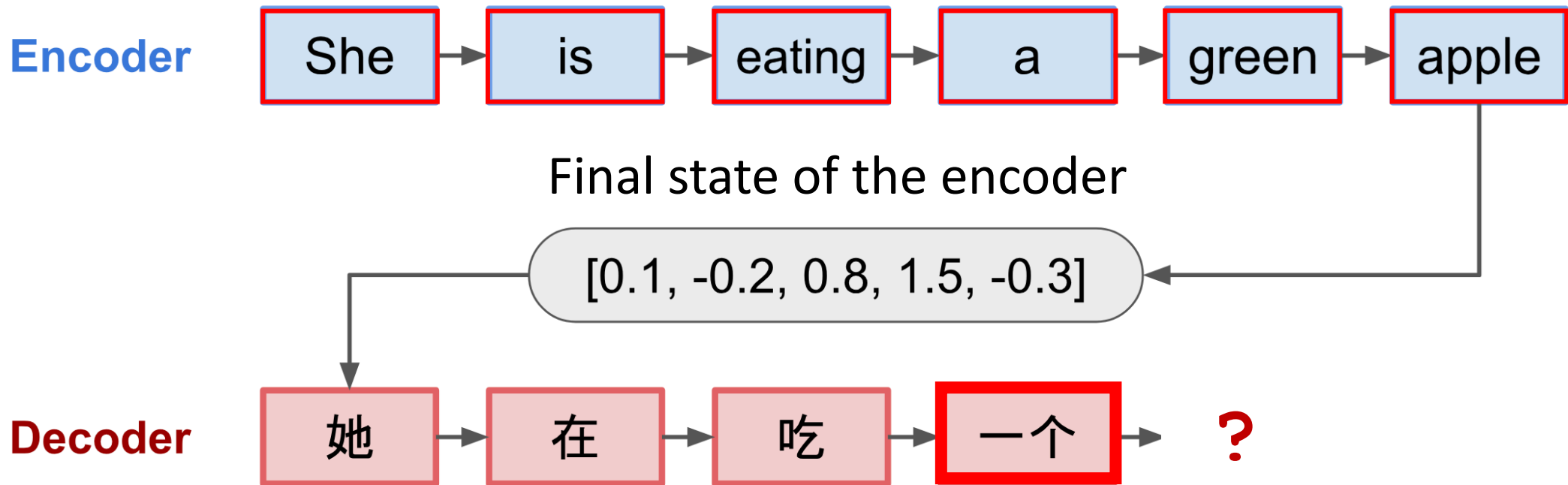
Summary

- Standard Seq2Seq model: the decoder looks at only **its current state**.



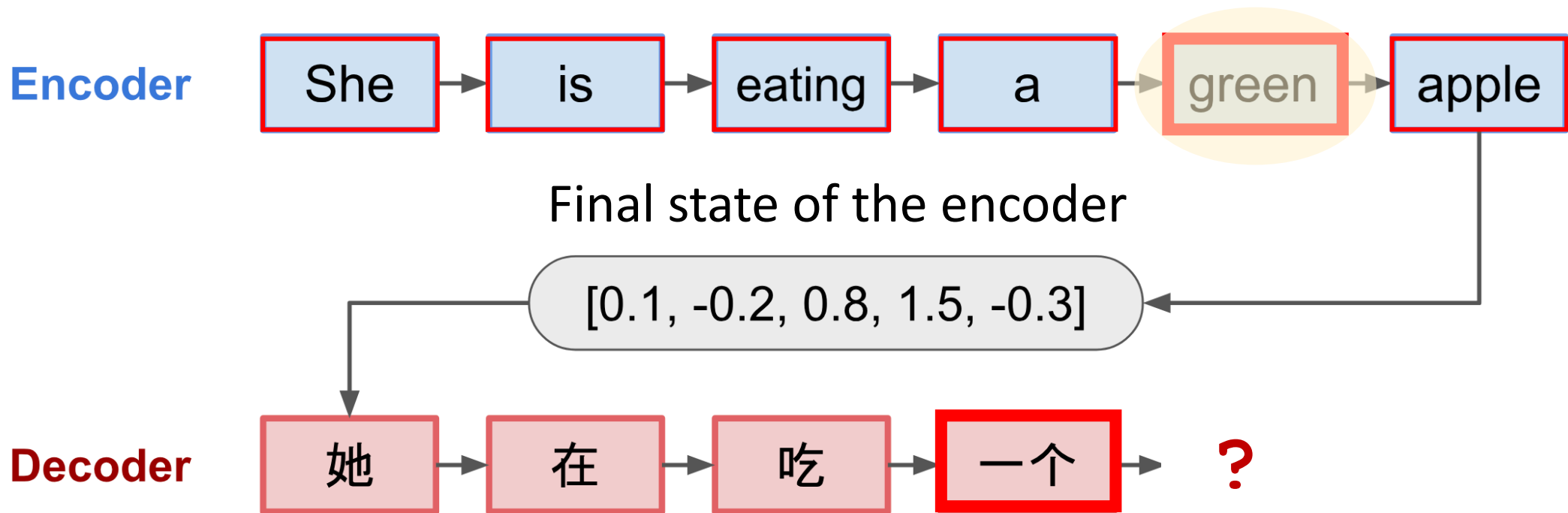
Summary

- Standard Seq2Seq model: the decoder looks at only its current state.
- Attention: decoder additionally looks at **all the states of the encoder**.



Summary

- Standard Seq2Seq model: the decoder looks at only its current state.
- Attention: decoder additionally looks at **all the states of the encoder**.



Summary

- Standard Seq2Seq model: the decoder looks at only its current state.
- Attention: decoder additionally looks at all the states of the encoder.
- Downside: higher time complexity.
 - l_1 : input sequence length
 - l_2 : target sequence length
 - Standard Seq2Seq: $O(l_1 + l_2)$ time complexity
 - Seq2Seq + attention: $O(l_1 l_2)$ time complexity