# Attention

**Shusen Wang**

# Seq2Seq Model



**Encoder** | She | is | eating | a | green | apple |

Final state of the encoder

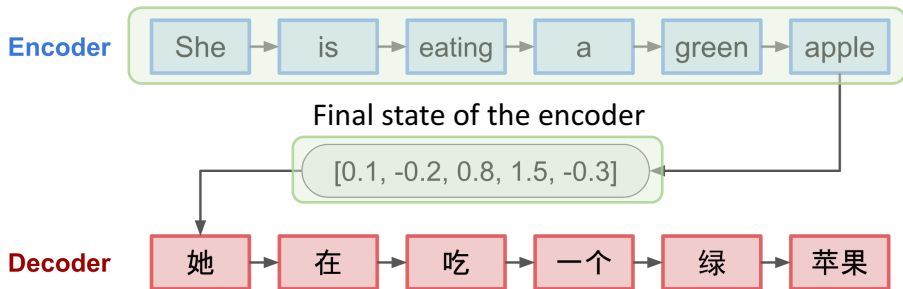[0.1, -0.2, 0.8, 1.5, -0.3]

**Decoder** | 她 | 在 | 吃 | 一个 | 绿 | 苹果 |

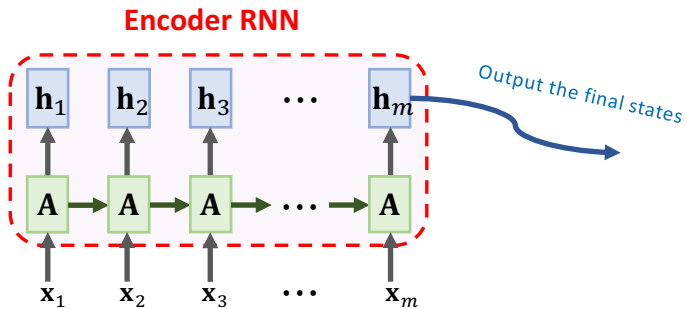The figure is from blog lilianweng.github.io

# Seq2Seq Model

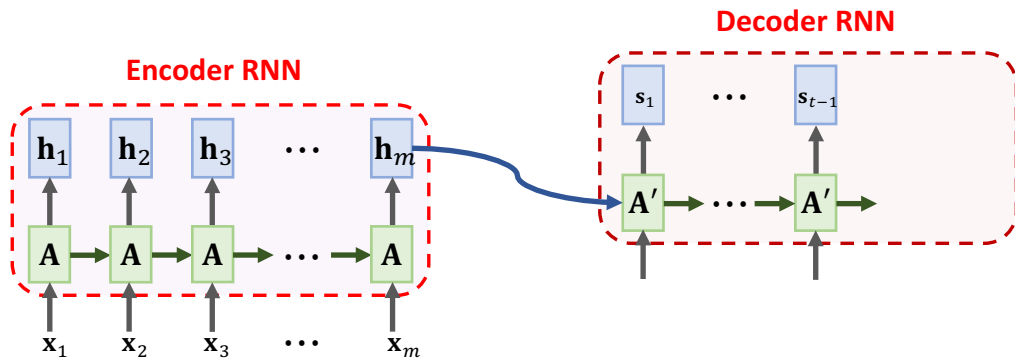**Shortcoming:** The final state is incapable of remembering a **long** sequence.



The figure is from blog lilianweng.github.io
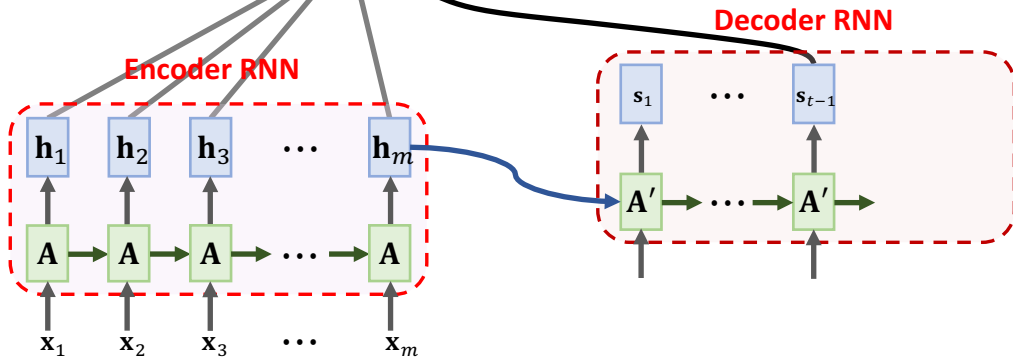
# Seq2Seq Model

# Seq2Seq Model

# Attention



context vector

$$\mathbf{c}_t = \sum_{i=1}^{m} \alpha_{ti} \ \mathbf{h}_i.$$

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

**Decoder RNN**

**Encoder RNN**

# Attention

context vector

$$\mathbf{c}_t = \sum_{i=1}^{m} \alpha_{ti} \ \mathbf{h}_i.$$

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



Figure is from https://distill.pub/2016/augmented-rnns/
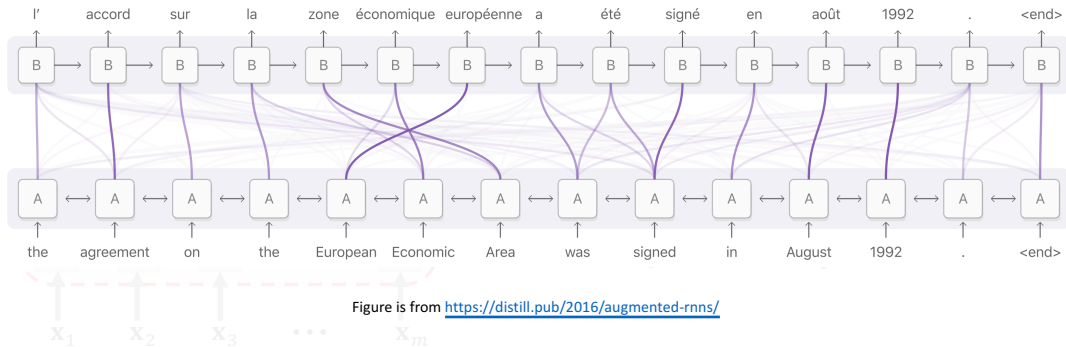
# Attention



context vector

$$\mathbf{c}_t = \sum_{i=1}^{m} \alpha_{ti} \; \mathbf{h}_i.$$

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

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

# Attention



context vector

$$\mathbf{c}_t = \sum_{i=1}^{m} \alpha_{ti} \ \mathbf{h}_i.$$

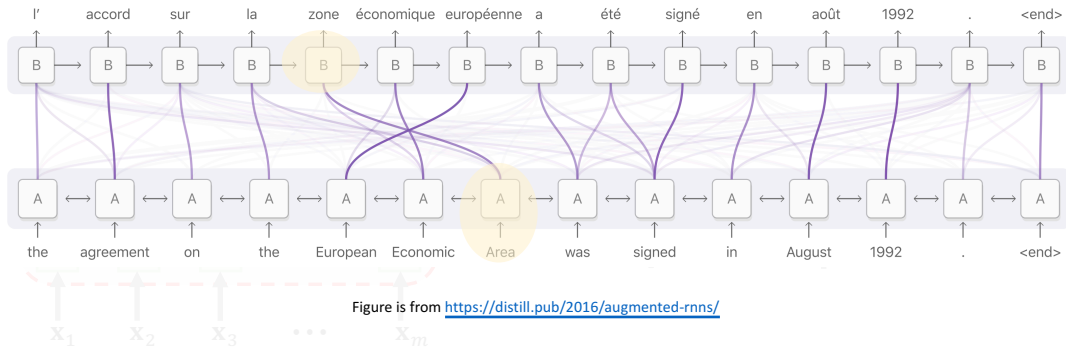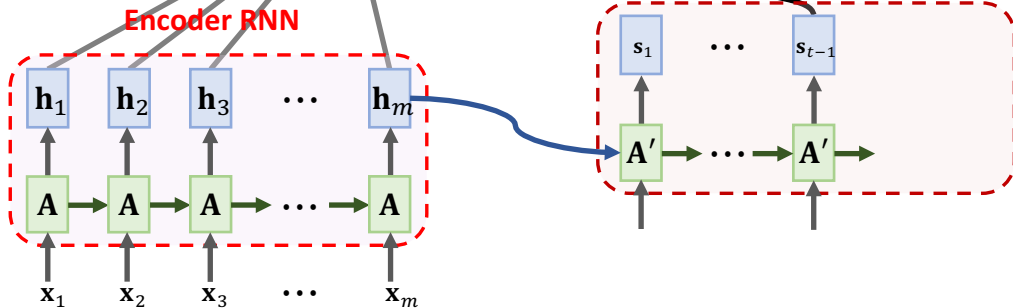- $\alpha_{ti}$: similarity between $\mathbf{s}_{t-1}$ and $\mathbf{h}_i$.
- $\alpha_{ti}$ is computed by a neural network taking $\mathbf{s}_{t-1}$ and $\mathbf{h}_i$ as input.

Decoder RNN

Encoder RNN

**Attention**

$$\mathbf{c}_t = \sum_{i=1}^{m} \alpha_{ti} \ \mathbf{h}_i.$$

Encoder RNN

Decoder RNN