

# Introduction to Attention Mechanism in Deep Learning

23.06.2020

# Recap

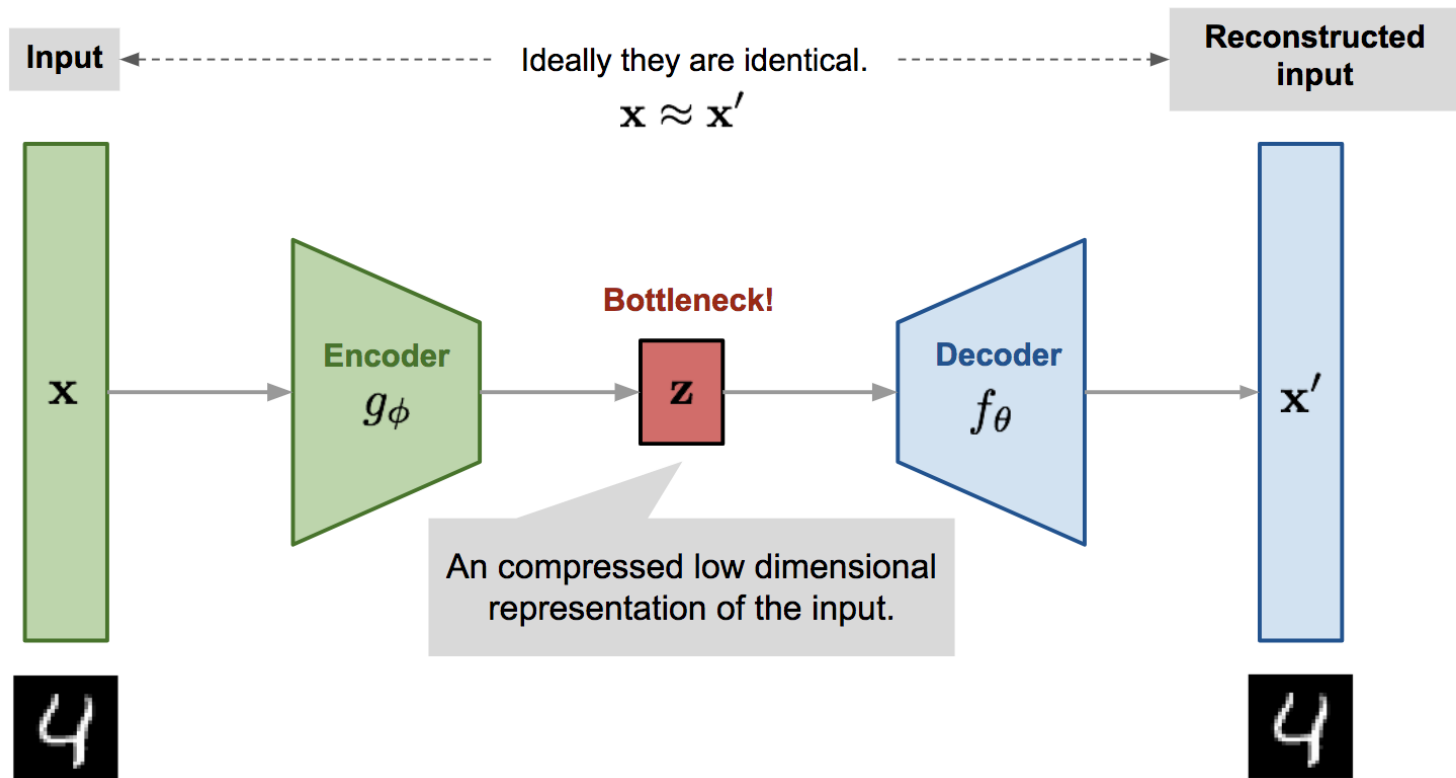


Figure 01. Autoencoder with MNIST dataset[1]

*What about sequential data?*

# Seq2seq Model

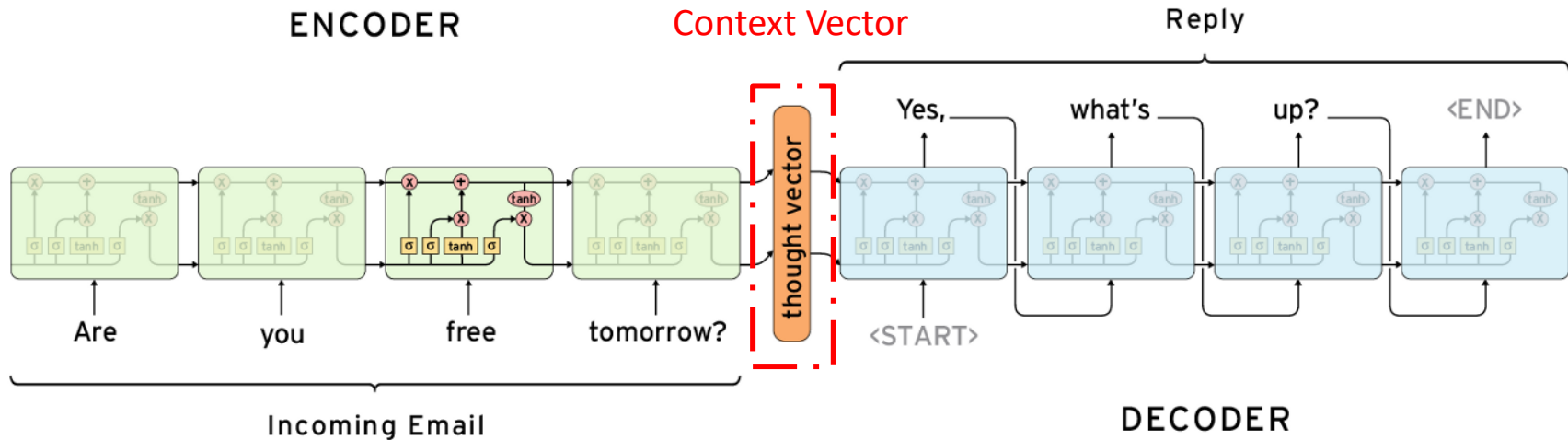


Figure 02. Seq2seq LSTMs[1]

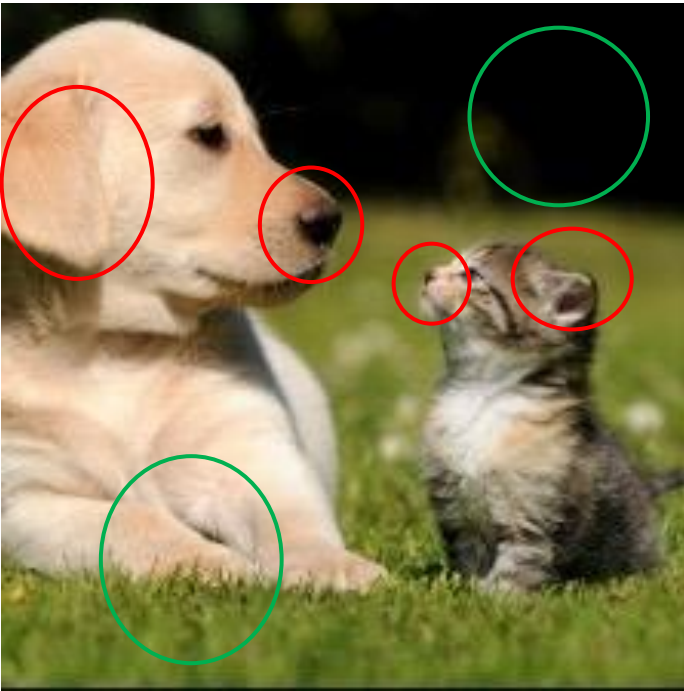
*Seq2seq was designed to encode sequential data (e.g. text, audio, speech) and decode sequential data of arbitrary length respectively.*

**Application:** Machine Translation, Speech Recognition, Parsing Sentence into Grammar Tree

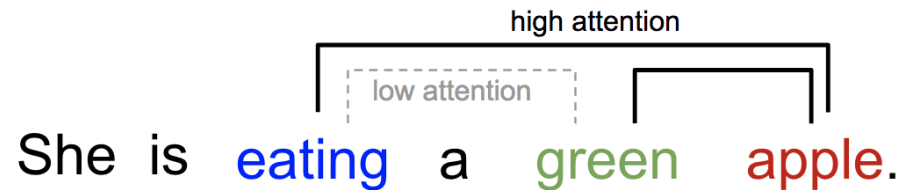
**Problem:** Context Vector of fixed length!

# Attention

- Biologically, attention existed is the visual system of a human being.
- High Attention: “High Resolution”
- Low Attention: “Low Resolution”



(a)



(b)

Figure 03. (a) Perceiving attention visually [1] and in (b) text [2]

[1] Image source: <https://distill.pub/2018/building-blocks/>

[2] Image source: <https://lilianweng.github.io/lil-log/2018/06/24/attention-attention.html>

# Attention

- Attention mechanism was design to tackle fixed context vector
- Attention create shortcuts between the context vector and the input source
- Focus more on certain regions and less on others
- Attention fells what to keep during reduction to minimize information loss salient features
- Attention provide more “fine grain” of region to be attended
- Classification: Self- Attention, Soft-Attention, Hard-Attention

# Soft vs Hard Attention

Soft Attention	Hard Attention
Different Parts, Different Subregions	Only ONE subregion
Deterministic	Stochastic
Differentiable	Non-differentiable

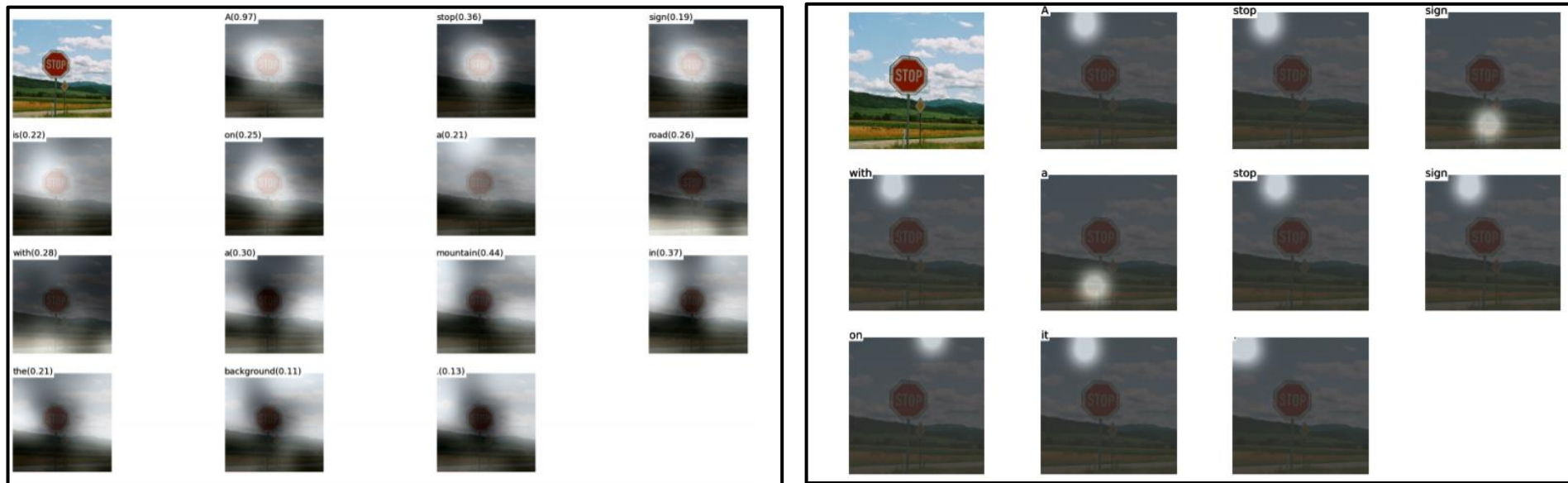


Figure 04. (Left) Soft subregions vs (Right) hard subregion [1]

# Soft Attention-Additive Attention (1)

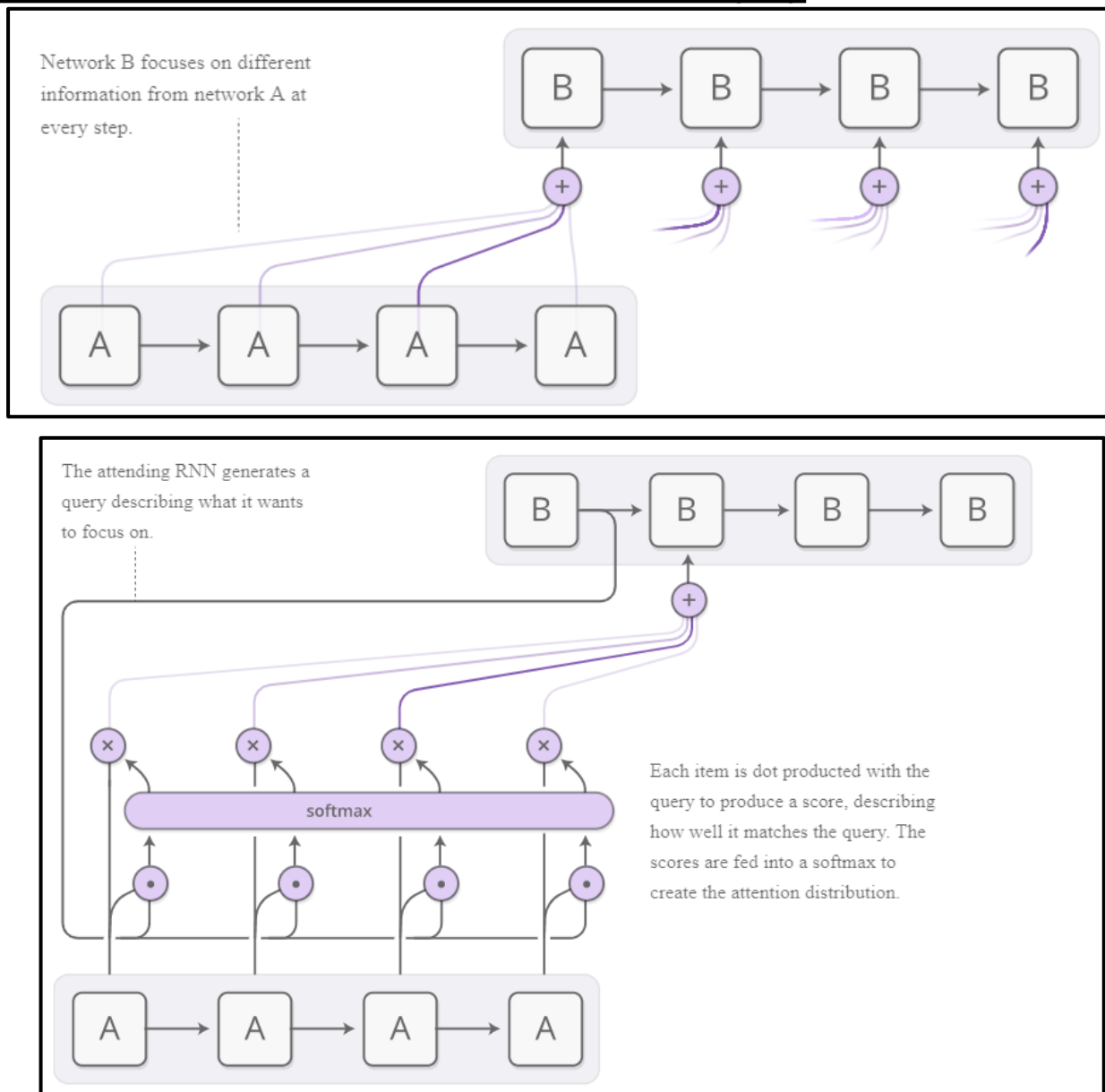


Figure 05. (Top) Aggregated additive attention mechanism. (Bottom) Unfold additive attention mechanism [1]

# Soft Attention-Additive Attention(2)

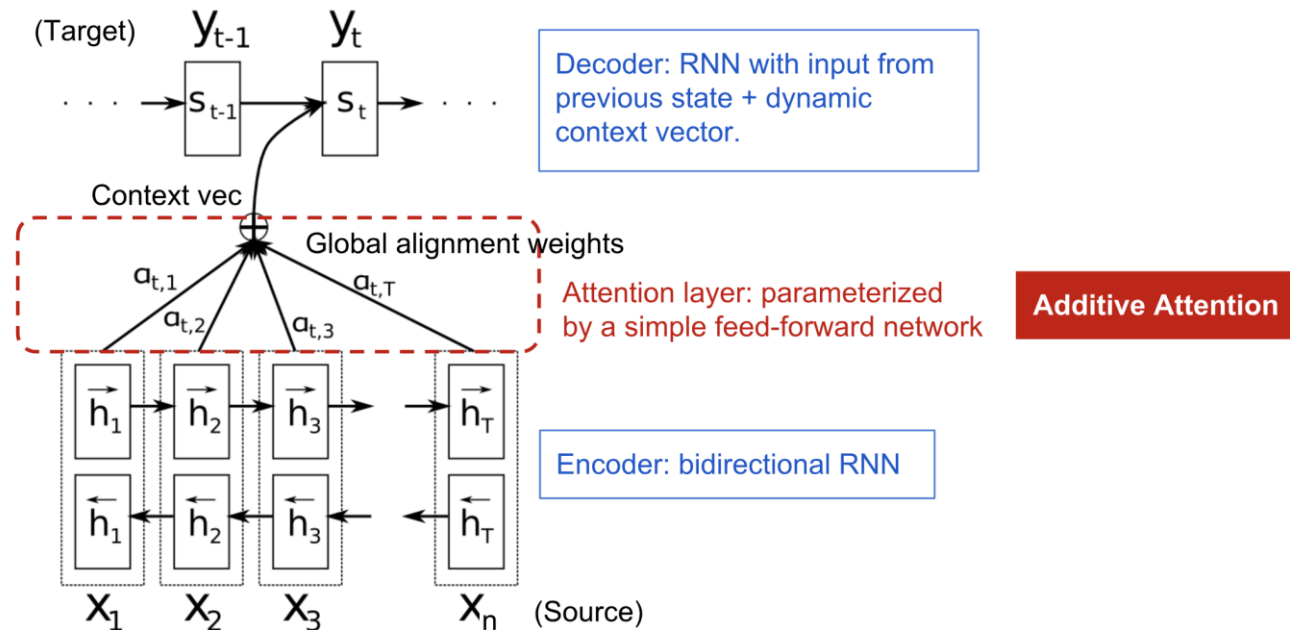


Figure 06. Diagram derived from Bahdanau et al., 2015 with addition information from [1]

$$\mathbf{c}_t = \sum_{i=1}^n \alpha_{t,i} \mathbf{h}_i \quad ; \text{ Context vector for output } y_t$$

$$\alpha_{t,i} = \text{align}(y_t, x_i) \quad ; \text{ How well two words } y_t \text{ and } x_i \text{ are aligned.}$$

$$= \frac{\exp(\text{score}(s_{t-1}, \mathbf{h}_i))}{\sum_{i'=1}^n \exp(\text{score}(s_{t-1}, \mathbf{h}_{i'}))} \quad ; \text{ Softmax of some predefined alignment score..}$$

$$\text{score}(s_t, \mathbf{h}_i) = \mathbf{v}_a^\top \tanh(\mathbf{W}_a [s_t; \mathbf{h}_i])$$



# Soft Attention-Additive Attention(3)

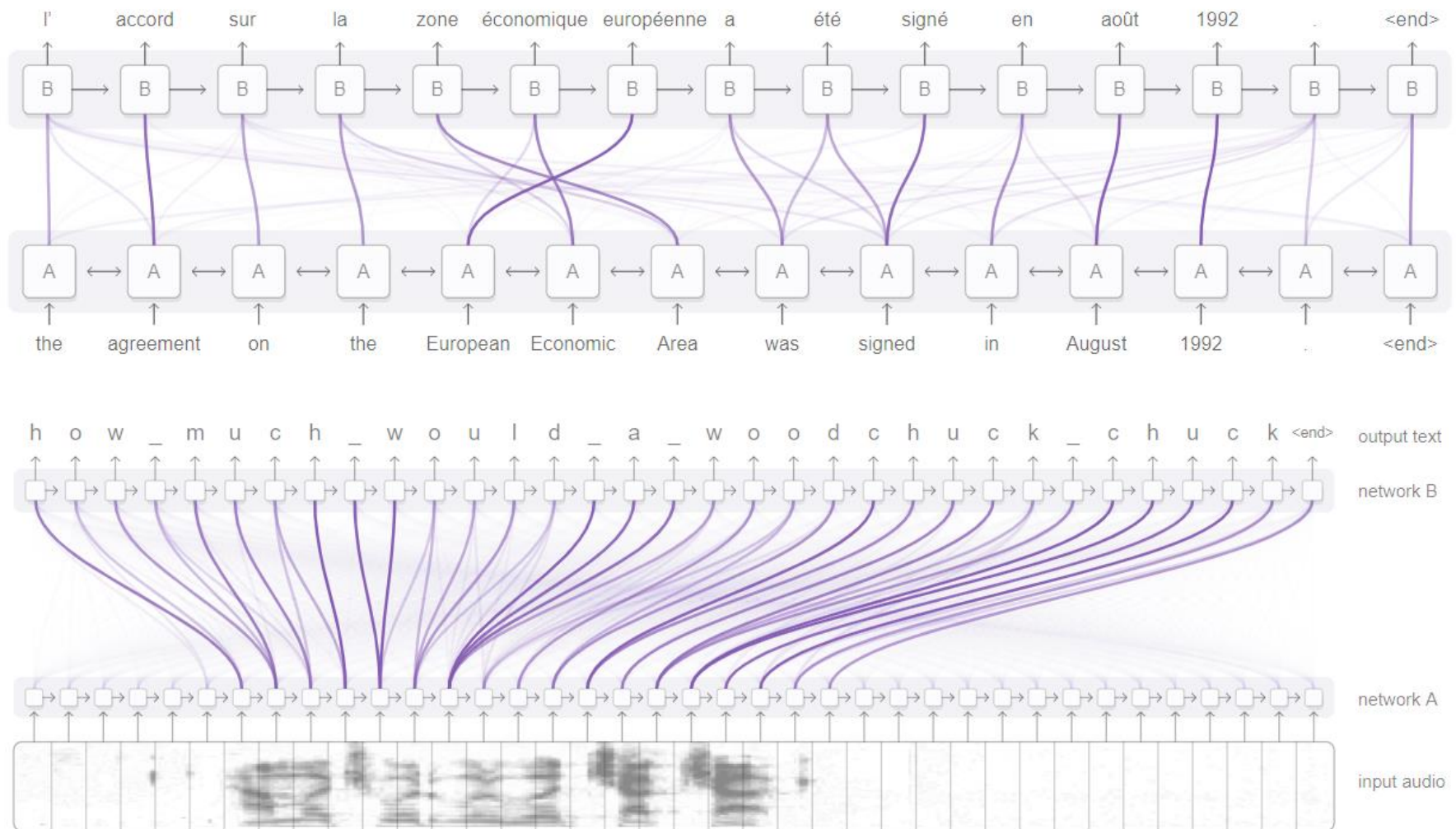


Figure 06. (Top) Text2text machine translation, (Bottom) Speech2text machine translation [1]

# Summary

Name	Definition	Citation
Self-Attention(&)	Relating different positions of the same input sequence. Theoretically the self-attention can adopt any score functions above, but just replace the target sequence with the same input sequence.	<a href="#">Cheng2016</a>
Global/Soft	Attending to the entire input state space.	<a href="#">Xu2015</a>
Local/Hard	Attending to the part of input state space; i.e. a patch of the input image.	<a href="#">Xu2015</a> ; <a href="#">Luong2015</a>

Name	Alignment score function	Citation
Content-base attention	$\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \text{cosine}[\mathbf{s}_t, \mathbf{h}_i]$	<a href="#">Graves2014</a>
Additive(*)	$\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \mathbf{v}_a^\top \tanh(\mathbf{W}_a[\mathbf{s}_t; \mathbf{h}_i])$	<a href="#">Bahdanau2015</a>
Location-Base	$\alpha_{t,i} = \text{softmax}(\mathbf{W}_a \mathbf{s}_t)$ Note: This simplifies the softmax alignment to only depend on the target position.	<a href="#">Luong2015</a>
General	$\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \mathbf{s}_t^\top \mathbf{W}_a \mathbf{h}_i$ where $\mathbf{W}_a$ is a trainable weight matrix in the attention layer.	<a href="#">Luong2015</a>
Dot-Product	$\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \mathbf{s}_t^\top \mathbf{h}_i$	<a href="#">Luong2015</a>
Scaled Dot-Product(^)	$\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \frac{\mathbf{s}_t^\top \mathbf{h}_i}{\sqrt{n}}$ Note: very similar to the dot-product attention except for a scaling factor; where n is the dimension of the source hidden state.	<a href="#">Vaswani2017</a>

Thank you for your kind  
attention!