# Alternative Architectures for Neural Machine Translations

Weiting (Steven) Tan

- RNN + Attention Recap
- Convolutional Neural Net
  - How it works in encoder/decoder
  - CNN + Attention
- Transformer

#### **Feed-Forward Layer**



- Classic neural network component
- Given an input vector *x*, matrix multiplication *M* with adding a bias vector *b*

Mx + b

• Adding a non-linear activation function

 $y = \operatorname{activation}(Mx + b)$ 

Notation

$$y = FF_{\text{activation}}(x) = a(Mx + b)$$

#### Recurrent Neural Network

- propagate state *s*
- over time steps t
- receiving an input  $x_t$  at each turn

$$s_t = f(s_{t-1}, x_t)$$

(state may computed may as a feed-forward layer)

#### Recurrent Neural Network

- propagate state *s*
- over time steps t
- receiving an input  $x_t$  at each turn

$$s_t = f(s_{t-1}, x_t)$$

(state may computed may as a feed-forward layer)





#### RNN's Issue: Vanishing Gradient

 $E_0$  $E_1$  $E_2$  $E_3$  $E_4$ • Mitigation: • GRU • LSTM  $\frac{\partial E_3}{\partial s_3}$  $rac{\partial s_3}{\partial s_2}$  $\tfrac{\partial s_1}{\partial s_0}$  $rac{\partial s_2}{\partial s_1}$  $s_2$  $s_4$  $s_0$  $s_1$  $s_3$  $x_0$  $x_1$  $x_2$  $x_3$  $x_4$ 

#### RNN with Attention

- Luong et al 2015
- Encoder-Decoder Attention



#### **RNN** with Attention

- Luong et al 2015
- Encoder-Decoder Attention

Other ways to compute attention

- Dot product:  $a(s_{i-1}, h_j) = s_{i-1}^T h_j$
- Scaled dot product:  $a(s_{i-1}, h_j) = \frac{1}{\sqrt{|h_j|}} s_{i-1}^T h_j$
- General:  $a(s_{i-1}, h_j) = s_{i-1}^T W_a h_j$
- Local:  $a(s_{i-1}) = W_a s_{i-1}$







#### RNN's issue: SLOW

**Alternative Sequence Modeling** 

- Convolutional Neural Networks
- Transformer (Self-Attention)

Encoder: parallelized computation

Decoder: parallelized during training, autoregressive during inference

#### ConvNet for 2D Images



## ConvNet for Seq2Seq (Encoder)

- First end-to-end neural machine translation model of the modern era [Kalchbrenner and Blunsom, 2013]
- Encoder





#### ConvNet with Attention

- Gehring et al. 2017
- Encoder:
  - Padding from left and right (keep seq-len)
  - Parallel computation



#### ConvNet with Attention

- Decoder:
  - Autoregressive
  - Stacking of feed-forward net as decoding is processed





#### ConvNet with Attention

- Encoder-Decoder
- Attention is the simple dot product



#### Transformer

- 1. Self-Attention and Encoder-Decoder Attention
- 2. Multi-Head Attention
- 3. Positional Encoding/Embedding



Jay Alammar https://jalammar.github.io/illustrated-transformer/



#### Decoder 1

#### Transformer Encoder Attention

Scaled Dot-Product Attention







![](_page_23_Figure_1.jpeg)

×

×

Х

![](_page_23_Figure_2.jpeg)

Q

Wv		V
	=	

![](_page_24_Figure_0.jpeg)

![](_page_25_Figure_0.jpeg)

![](_page_26_Figure_0.jpeg)

### Transformer Decoder

- Self-attention with mask
- Encoder-Decoder Attention
  - Query is from decoder
  - Key and Value are from Encoder

![](_page_27_Figure_5.jpeg)

hi key and value self - attention かい

#### Positional Embedding

![](_page_29_Figure_1.jpeg)

#### Positional Embedding

$$\overrightarrow{p_t} = egin{bmatrix} \sin(\omega_1.\,t)\ \cos(\omega_1.\,t)\ \sin(\omega_2.\,t)\ \cos(\omega_2.\,t)\ arphi(\omega_2.\,t)\ arphi(\omega_{d/2}.\,t)\ arphi(\omega_{d/2}.\,t)\ \cos(\omega_{d/2}.\,t)\ arphi(\omega_{d/2}.\,t)\ arph(\omega_{d/2}.\,t)\ arphi(\omega_{d/2}.\,t)\ arph(\omega_{d/2}.\,t)\ arphi($$

$$\overrightarrow{p_t}^{(i)} = f(t)^{(i)} := egin{cases} \sin(\omega_k,t), & ext{if } i = 2k \ \cos(\omega_k,t), & ext{if } i = 2k+1 \end{cases}$$

$$\omega_k = rac{1}{10000^{2k/d}}$$

#### Positional Embedding

![](_page_31_Figure_1.jpeg)

https://github.com/jalammar/jalammar.github.io/blob/master/notebookes/transformer/transformer\_positional\_encoding\_graph.ipynb