



CS F425: Deep Learning

25

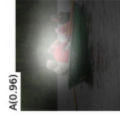
Bi-dir RNN Transformer



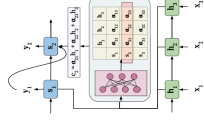
Dr. Kamlesh Tiwari
Assistant Professor, Department of CSIS,
BITS Pilani, Pilani Campus, Rajasthan-333031 INDIA
Mar 24, 2023 **ON-CAMPUS** Campus @ BITS-Pilani [Jan-May 2023]

<http://ktiwari.in/dl>

Attention Model ²



Our visual system tends to focus selectively on some parts of the image, while ignoring other irrelevant information



- Attention was first introduced in Machine Translation
- Attention weights are learned by additional feed forward NN

² cite: 27782 Bahdanau, Kyunghyun, and Yoshua. "Neural machine translation by jointly learning to align and translate." arXiv preprint arXiv:1409.0473 (2014).

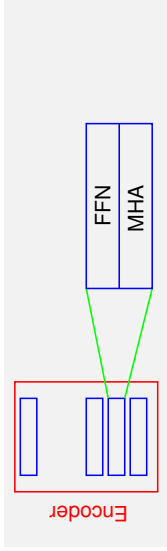
Chaudhari, Sneha, et al. "An attentive survey of attention models." ACM Tr intelligent Systems and Technology (TIST) 2021

Overview

- Has **Encoder** (self attention) and **Decoder** (self+cross attention)
- Issue with sequential and recurrent data
- Things are sequential and dependent

Dr Kalam was president. He also worked at DRDO

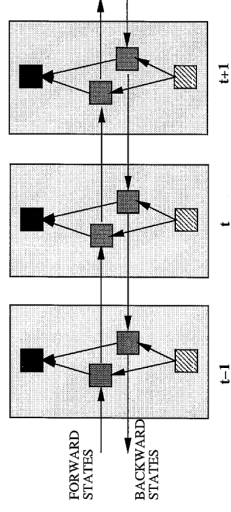
- Word2Vec⁴ or GLOVE⁵ can be used to get representations (512D)
- 6 blocks of feed forward network, and multi head attention



⁴ cite: 33782 Mikolov, Tomas, et al. "Efficient estimation of word representations in vector space." arXiv preprint arXiv:1301.3781 (2013).

⁵ cite: 33108 Pennington, Jeffrey, et al. "Glove: Global vectors for word representation." Proceedings of the conference on empirical methods in natural language processing (EMNLP), 2014.

Bidirectional LSTM/RNN ¹



- No limitation on using input just up to a preset future frame
- Simultaneous training in positive and negative time direction
- Modified one can estimate conditional posterior probability of complete symbol sequences without assumption on distribution
- Useful to real data, classification experiments

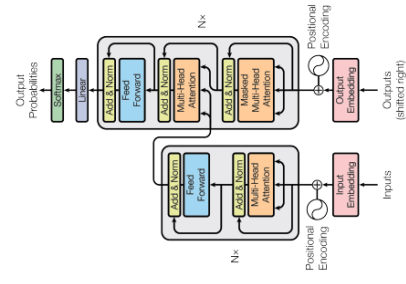
¹ cite: 8304 Schuster, Mike, and Kulkup K. Paliwal. "Bidirectional recurrent neural networks." IEEE transactions on Signal Processing 45.11 (1997): 2673-2681.

Transformer ³

Sequence transduction models include an encoder and a decoder that are connected through attention

- Transformer is simple, solely based on attention mechanism
- It takes significantly less time to train
- Simultaneous input of sentence handles issues related to exploding/vanishing gradients

Training is sequential and end-to-end



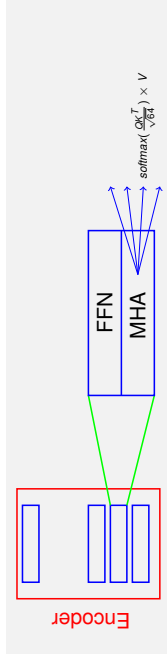
³ cite: 69087 Vaswani, et al. "Attention is all you need." Advances in neural information processing systems (NIPS), 2017.

Multi Head Attention (MHA)

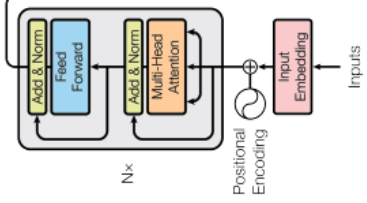
- **Multiple** set of three dedicated NN circuit Query, Key, and Value
- Converts 512D vector to 64D size
- Attention? connection strength between different words

The dog eat food because it was hungry.

- Self attention and cross attention
- Consider vectors X_q^i, X_k^k , and X_v^l for I am good.
- Word embedding $\rightarrow \text{softmax}(\frac{QK^T}{\sqrt{64}}) \times V$
- Concatenate **multiple** ones and apply FC to get 64D vector



Transformer Multi-Head Attention (MHA)



Output of MHA would be same for different permutation of the same sentence

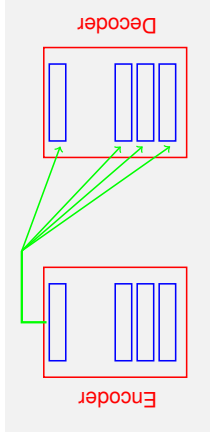
- Positional embedding is added

$$X = embedding + P$$
- P being sinusoidal encoding
- Residual connections are added and layer normalization is applied

Thank You!

Thank you very much for your attention!

Decoder



- Masking is needed while training for self attention
- Not for cross attention
- Residual connections and layer normalization are added

