

# Tutorial on deep transformer

Loc Nguyen

Loc Nguyen's Academic Network, Vietnam

Email: ng\_phloc@yahoo.com

Homepage: www.locnguyen.net

## Abstract

Development of transformer is a far progressive step in the long journeys of both generative artificial intelligence (GenAI) and statistical translation machine (STM) with support of deep neural network (DNN), in which STM can be known as interesting result of GenAI because of encoder-decoder mechanism for sequence generation built in transformer. But why is transformer being preeminent in GenAI and STM? Firstly, transformer has a so-called self-attention mechanism that discovers contextual meaning of every token in sequence, which contributes to reduce ambiguousness. Secondly, transformer does not concern ordering of tokens in sequence, which allows to train transformer from many parts of sequences in parallel. Thirdly, the third reason which is result of the two previous reasons is that transformer can be trained from large corpus with high accuracy as well as highly computational performance. Moreover, transformer is implemented by DNN which is one of important and effective approaches in artificial intelligence (AI) in recent time. Although transformer is preeminent because of its good consistency, it is not easily understandable. Therefore, this technical report aims to describe transformer with explanations which are as easily understandable as possible.

**Keywords:** transformer, generative artificial intelligence, statistical translation machine, sequence generation, self-attention.

## 1. Introduction

Artificial intelligence (AI) is recent trend in technological world, especially in computer science, in which artificial neural network (ANN, NN) is one of important subjects of AI. Essentially, ANN models or implements a complicated function  $\mathbf{y} = f(\mathbf{x})$  where  $\mathbf{x} = (x_1, x_2, \dots, x_m)^T$  and  $\mathbf{y} = (y_1, y_2, \dots, y_n)^T$  are vectors so that  $\mathbf{x}$  and  $\mathbf{y}$  are imitated by input layer and output layer of ANN, respectively with note that each layer is composed of units called neurons  $x_i, y_i$ . The complication degree of function  $\mathbf{y} = f(\mathbf{x})$  is realized by hidden layers of ANN which are intermediated layers between input layer and output layer. We denote:

$$\mathbf{y} = f(\mathbf{x}|\Theta)$$

Where  $\Theta$  denotes parameters of ANN which are often weights and biases. Because  $f(\mathbf{x} | \Theta)$  is essentially vector-by-vector function whose input and output are vectors, it should have denoted as  $f(\mathbf{x} | \Theta)$  but it is still denoted as  $f(\mathbf{x} | \Theta)$  for convenience and moreover, input  $\mathbf{x}$  and output  $\mathbf{y}$  will be matrices if their elements  $x_i$  and  $y_i$  are vectors. If there are many enough hidden layers, ANN becomes a so-called deep neural network (DNN) such that DNN is cornerstone of the main subject of this report which is transformer because transformer, as its name implies, is the highly abstract and complicated version of function  $\mathbf{y} = f(\mathbf{x})$ . In other words, a transformer will make the transformation between complex and different objects if it is implemented by DNN or set of DNNs according to viewpoint of DNN. Although transformer can be applied into many areas, especially machine translation and computer vision, this report focuses on statistical machine translation (STM) because complex and different objects  $\mathbf{x}$  and  $\mathbf{y}$  in STM transformer are two sentences in two different languages where  $\mathbf{x}$  is source language sentence and  $\mathbf{y}$  is target language sentence. If

ordering of elements  $x_i / y_i$  in vector  $\mathbf{x} / \mathbf{y}$  specifying sentence is concerned as ordering of words  $x_i / y_i$  in a sentence, transformer will relate to *sequence generation*. Therefore, transformer as well as STM are inspired from sequence generation which, in turn, relates to recurrent neural network (RNN) as well as long short-term memory (LSTM) because sequence generation models are often implemented by RNN or LSTM. The most standard ANN/DNN called feedforward network (FFN) follows the one-way direction from input layer to hidden layers to output layer without reverse direction, which means that there is neither connections from output layer to hidden layers nor connections from hidden layers to input layers. In other words, there is no cycle in FFN, which cause the side-effect that it is difficult to model a sequence vector  $\mathbf{x} = (x_1, x_2, \dots, x_m)^T$  like a sentence in natural language processing (NLP) because elements / words / terms / tokens  $x_i$  in such sequence/sentence vector have the same structure and every connection  $x_i \rightarrow x_{i+1}$  of two successive words  $x_i$  and  $x_{i+1}$  is, actually, a cycle. This is the reason that recurrent neural network (RNN) is better than FFN to generate sequence. Therefore, we research transformer after researching sequence generation which is concerned after RNN is concerned. Note, sequence and sentence are two exchangeable concepts in this research.

Suppose entire FFN is reduced into a state in RNN and RNN is ordered list of neurons called sequence of neurons and moreover, output of previous neuron  $x_{i-1}$  contributes to input of current neuron  $x_i$ . Namely, for formal definition, given  $T$  time points  $t = 1, 2, \dots, T$ , then RNN is ordered sequence of  $T$  states and each state is modeled by triple  $(x_t, h_t, o_t)$  called state  $(x_t, h_t, o_t)$  where  $x_t$ ,  $h_t$ , and  $o_t$  represent input layer, hidden layer, and output layer, respectively. Without loss of generality, let  $x_t$ ,  $h_t$ , and  $o_t$  represent input neuron, hidden neuron, and output neuron, respectively when a layer is represented by one of its neurons. Please pay attention that  $x_t$ ,  $h_t$ , and  $o_t$  are represented vectors of the  $t^{\text{th}}$  word in sentence  $\mathbf{x} = (x_1, x_2, \dots, x_m)^T$  modeled by RNN in context of NLP because a word is modeled by a numeric vector in NLP. Therefore, the aforementioned sentence  $\mathbf{x} = (x_1, x_2, \dots, x_m)^T$  is a matrix indeed but  $\mathbf{x}$  is mentioned as a vector. Exactly,  $\mathbf{x}$  is vector of vectors, which leads to the convention that its elements are denoted by bold letter such as  $\mathbf{x}_i$  or  $\mathbf{x}_t$  because such elements are variable vectors representing words. Note, a word in NLP can be mentioned as term or token.

$$\begin{aligned}\mathbf{x} &= \mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m)^T \\ \mathbf{y} &= \mathbf{Y} = (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n)^T \\ \mathbf{y} &= \mathbf{Y} = f(\mathbf{x}|\Theta) = f(\mathbf{X}|\Theta)\end{aligned}$$

Note, the subscript “ $T$ ” denotes vector/matrix transposition operator. Whether the sentence / sequence is denoted as vector notation  $\mathbf{x}$  or matrix notation  $\mathbf{X}$  belongs to contextual explanations. Recall that transformer as well as STM are inspired from sequence generation which, in turn, is related to recurrent neural network (RNN) as well as long short-term memory (LSTM) because sequence generation models are often implemented by RNN or LSTM. Function  $\mathbf{y} = f(\mathbf{x} | \Theta)$  implemented by DNNs such as RNN and LSTM is also called generator because it is sequence generation model indeed. Therefore, although transformer is different from RNN and LSTM, all of them are denoted by generator  $\mathbf{y} = f(\mathbf{x} | \Theta)$  because they are sequence generation models indeed.

The  $t^{\text{th}}$  element/word in sequence/sentence  $\mathbf{x} = (x_1, x_2, \dots, x_m)^T$  is represented by the  $t^{\text{th}}$  state  $(\mathbf{x}_t, \mathbf{h}_t, \mathbf{o}_t)$  of RNN where  $\mathbf{x}_t$  is the  $t^{\text{th}}$  input word and  $\mathbf{o}_t$  is the  $t^{\text{th}}$  output word. If RNN models  $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m)^T$ , then  $T = m$  and so, if RNN models  $\mathbf{y} = (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n)^T$ , then  $T = n$ . By a convention, word and sentence are mentioned as token and sequence, respectively. Moreover,  $\mathbf{x}$  is called source sequence and  $\mathbf{y}$  is called target sequence or generated sequence. Mathematical equation to update RNN is specified as follows (Wikipedia, Recurrent neural network, 2005):

$$\begin{aligned} \mathbf{h}_t &= \sigma_h(\mathbf{W}_h \mathbf{x}_t + \mathbf{U}_h \mathbf{h}_{t-1} + \mathbf{b}_h) \\ \mathbf{o}_t &= \sigma_o(\mathbf{W}_o \mathbf{h}_t + \mathbf{b}_o) \end{aligned}$$

Where  $\mathbf{W}_h$ ,  $\mathbf{U}_h$ , and  $\mathbf{W}_o$  are weight matrices of current hidden neuron  $\mathbf{h}_t$ , previous hidden neuron  $\mathbf{h}_{t-1}$ , and current output neuron  $\mathbf{o}_t$ , respectively whereas  $\mathbf{b}_h$  and  $\mathbf{b}_o$  are bias vectors of  $\mathbf{h}_t$  and  $\mathbf{o}_t$ , respectively. Moreover,  $\sigma_h(\cdot)$  and  $\sigma_o(\cdot)$  are activation functions of  $\mathbf{h}_t$  and  $\mathbf{o}_t$ , respectively, which are vector-by-vector functions.

RNN copes with the problem of vanishing gradient when learning a long RNN of many states and so, long short-term memory (LSTM) is proposed to restrict the problem of vanishing gradient. State in RNN becomes cell in LSTM and so, given  $T$  time points  $t = 1, 2, \dots, T$ , let the pair  $(\mathbf{c}_t, \mathbf{h}_t)$  denote LSTM cell at current time point  $t$  where  $\mathbf{c}_t$  represents real information stored in memory and  $\mathbf{h}_t$  represents clear-cut information that propagates through next time points. A cell  $(\mathbf{c}_t, \mathbf{h}_t)$  has four gates such as forget gate  $\mathbf{f}_t$ , input gate  $\mathbf{i}_t$ , output gate  $\mathbf{o}_t$ , and cell gate  $\mathbf{g}_t$ . At every time point  $t$  or every iteration  $t$ , cell  $(\mathbf{c}_t, \mathbf{h}_t)$  updates its information based on these gates as follows:

$$\begin{aligned} \mathbf{f}_t &= \sigma_f(\mathbf{W}_f \mathbf{x}_t + \mathbf{U}_f \mathbf{h}_{t-1} + \mathbf{b}_f) \\ \mathbf{i}_t &= \sigma_i(\mathbf{W}_i \mathbf{x}_t + \mathbf{U}_i \mathbf{h}_{t-1} + \mathbf{b}_i) \\ \mathbf{o}_t &= \sigma_o(\mathbf{W}_o \mathbf{x}_t + \mathbf{U}_o \mathbf{h}_{t-1} + \mathbf{b}_o) \\ \mathbf{g}_t &= \sigma_g(\mathbf{W}_g \mathbf{x}_t + \mathbf{U}_g \mathbf{h}_{t-1} + \mathbf{b}_g) \\ \mathbf{c}_t &= \mathbf{f}_t \otimes \mathbf{c}_{t-1} + \mathbf{i}_t \otimes \mathbf{g}_t \\ \mathbf{h}_t &= \mathbf{o}_t \otimes \sigma_h(\mathbf{c}_t) \end{aligned}$$

Note,  $\mathbf{W}_{(\cdot)}$  and  $\mathbf{U}_{(\cdot)}$  are weight matrices whereas  $\mathbf{b}_{(\cdot)}$  are bias vectors, which are parameters. Because core information of cell  $(\mathbf{c}_t, \mathbf{h}_t)$  including  $\mathbf{c}_t$  and  $\mathbf{h}_t$  is calculated without any parameters, the problem of vanishing gradient can be alleviated when such gradient is calculated with regard to parameters such as weight matrices and bias vectors.

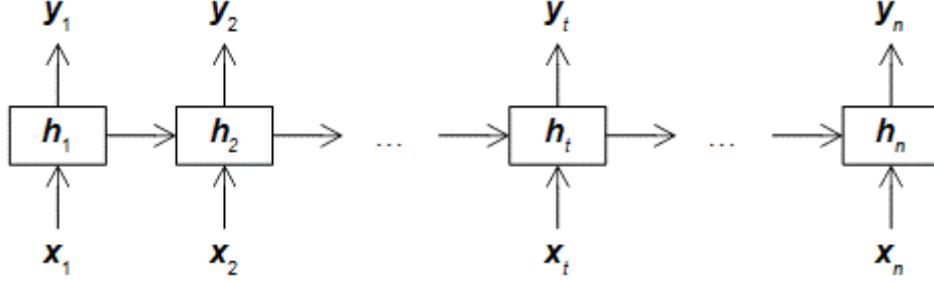
In general, when a sequence is modeled by a RNN or a LSTM, it is possible to generate a new sequence after RNN or LSTM is trained by backpropagation algorithm associated with stochastic gradient descent (SGD) algorithm. In other words, RNN and LSTM are important generation models although transformer is the main subject in this report because STM is, essentially, a sequence generation model that generates a new sentence in target language from a sentence in source language when sentence in NLP is represented by sequence. Because RNN and LSTM have the same methodological ideology, RNN is mentioned rather than LSTM because RNN is simpler one but they can be applied by exchangeable manner. For instance, given simplest case that source sequence  $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m)^T$  and target sequence also called generated sequence  $\mathbf{Y} = (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n)^T$  have the same length  $m = n$ .

$$\begin{aligned} \mathbf{X} &= (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)^T \\ \mathbf{Y} &= (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n)^T = f(\mathbf{X}|\Theta) \end{aligned}$$

Generation model  $f(\mathbf{x} | \Theta)$  is implemented by a RNN of  $n$  states  $(\mathbf{x}_t, \mathbf{h}_t, \mathbf{o}_t)$  so that  $\mathbf{o}_t = \mathbf{y}_t$  for all  $t$  from 1 to  $n$ . After RNN was trained from sample by backpropagation algorithm associated with SGD, given source sequence  $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)^T$ , target sequence  $\mathbf{Y} = (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n)^T$  is generated easily by evaluating  $n$  states of RNN.

$$\begin{aligned} \mathbf{h}_t &= \sigma_h(\mathbf{W}_h \mathbf{x}_t + \mathbf{U}_h \mathbf{h}_{t-1} + \mathbf{b}_h), \forall t = \overline{1, n} \\ \mathbf{y}_t &= \mathbf{o}_t = \sigma_o(\mathbf{W}_o \mathbf{h}_t + \mathbf{b}_o) \end{aligned}$$

Such generation process with  $n$ -state RNN is depicted by following figure:



**Figure 1.1.** RNN generation model

The next section will focus on sequence generation and attention which is a mechanism that improves generation process.

## 2. Sequence generation and attention

Recall that transformer as well as statistical machine translation (STM) are inspired from sequence generation which, in turn, is related to recurrent neural network (RNN) as well as long short-term memory (LSTM) because sequence generation models are often implemented by RNN or LSTM. Function  $y = f(x | \Theta)$  implemented by DNNs such as RNN and LSTM is also called generator because it is sequence generation model indeed. Because RNN and LSTM have the same methodological ideology, RNN is mentioned rather than LSTM.

$$\begin{aligned} \mathbf{x} = \mathbf{X} &= (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m)^T \\ \mathbf{y} = \mathbf{Y} &= (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n)^T \\ \mathbf{y} = \mathbf{Y} &= f(\mathbf{x}|\Theta) = f(\mathbf{X}|\Theta) \end{aligned}$$

Note,  $\Theta$  denotes parameters of ANN which are often weights and biases whereas sequence is denoted as vector notation  $\mathbf{x}$  or matrix notation  $\mathbf{X}$  belongs to contextual explanations. This section focuses on sequence generation models such as RNN and LSTM before mentioning advanced concepts of transformer because, anyhow, transformer is next evolutionary step of sequence generation models, especially in STM and natural language processing (NLP).

Given simplest case aforementioned that source sequence  $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m)^T$  and target sequence also called generated sequence  $\mathbf{Y} = (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n)^T$  have the same length  $m = n$ .

$$\begin{aligned} \mathbf{X} &= (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)^T \\ \mathbf{Y} &= (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n)^T = f(\mathbf{X}|\Theta) \end{aligned}$$

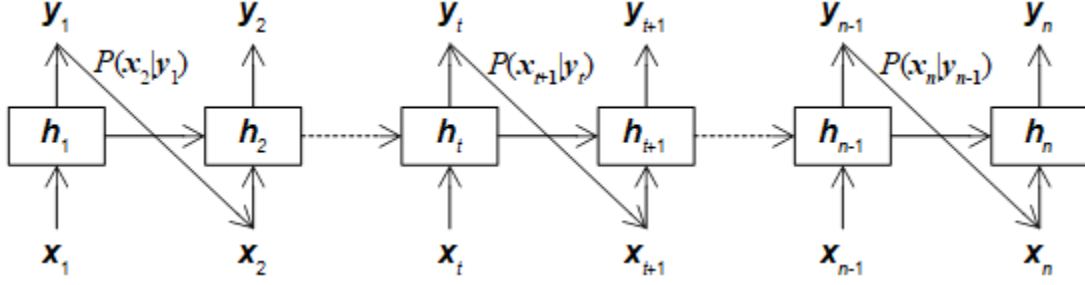
Generation model  $f(\mathbf{X} | \Theta)$  is implemented by a RNN of  $n$  states  $(\mathbf{x}_t, \mathbf{h}_t, \mathbf{o}_t)$  so that  $\mathbf{o}_t = \mathbf{y}_t$  for all  $t$  from 1 to  $n$ . After RNN was trained from sample by backpropagation algorithm associated with stochastic gradient descent (SGD) algorithm, given source sequence  $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)^T$ , target sequence  $\mathbf{Y} = (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n)^T$  is generated easily by evaluating  $n$  states of RNN.

$$\begin{aligned} \mathbf{h}_t &= \sigma_h(\mathbf{W}_h \mathbf{x}_t + \mathbf{U}_h \mathbf{h}_{t-1} + \mathbf{b}_h) \\ \mathbf{y}_t &= \mathbf{o}_t = \sigma_o(\mathbf{W}_o \mathbf{h}_t + \mathbf{b}_o) \end{aligned}$$

The simplest RNN generation needs to be extended if source sequence  $\mathbf{X}$  is incomplete, for example,  $\mathbf{X}$  has  $k$  token vectors  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k$  where  $k < n$ . When  $\mathbf{X}$  is incomplete, without loss of generality, given current output  $\mathbf{y}_t$ , it is necessary to predict the next output  $\mathbf{x}_{t+1}$  (with suppose  $t > k$ ). The prediction process, proposed by Graves (Graves, 2014), is based on estimating the predictive probability  $P(\mathbf{x}_{t+1} | \mathbf{y}_t)$  which is conditional probability of next input  $\mathbf{x}_{t+1}$  given current output  $\mathbf{y}_t$ . As a result, RNN generation model is extended as follows (Graves, 2014, p. 4):

$$\begin{aligned} \mathbf{h}_t &= \sigma_h(\mathbf{W}_h \mathbf{x}_t + \mathbf{U}_h \mathbf{h}_{t-1} + \mathbf{b}_h) \\ \mathbf{y}_t &= \sigma_o(\mathbf{W}_o \mathbf{h}_t + \mathbf{b}_o) \\ P(\mathbf{x}_{t+1} | \mathbf{y}_t) \end{aligned}$$

Following figure depicts the prediction model proposed by Graves (Graves, 2014, p. 3):



**Figure 2.1.** RNN prediction model

The problem here is how to specify predictive probability  $P(x_{t+1} | y_t)$ . In the most general form, suppose joint probability  $P(x_{t+1}, y_t)$  is parameterized by multivariate normal distribution with mean vector  $\mu$  and covariance matrix  $\Sigma$ .

$$P(x_{t+1}, y_t) = \mathcal{N}(x_{t+1}, y_t | \mu, \Sigma)$$

$$\mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}$$

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$$

It is easy to estimate  $\mu$  and  $\Sigma$  to determine  $P(x_{t+1}, y_t)$  from sample by maximum likelihood estimation (MLE) method, for instance. Consequently, predictive probability  $P(x_{t+1} | y_t)$  is determined based on joint probability  $P(x_{t+1}, y_t)$  as multivariate normal distribution with mean vector  $\mu_{12}$  and covariance matrix  $\Sigma_{12}$  specified as follows (Hardle & Simar, 2013, p. 157):

$$P(x_{t+1} | y_t) = \mathcal{N}(x_{t+1} | \mu_{12}, \Sigma_{12})$$

$$\mu_{12} = \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (y_t - \mu_2)$$

$$\Sigma_{12} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$$

Because predictive probability  $P(x_{t+1} | y_t)$  gets highest at the mean  $\mu_{12}$ , it is possible to estimate  $x_{t+1}$  given  $y_t$  by  $\mu_{12}$ .

$$h_t = \sigma_h(W_h x_t + U_h h_{t-1} + b_h)$$

$$y_t = \sigma_o(W_o h_t + b_o)$$

$$x_{t+1} = \mu_{12} = \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (y_t - \mu_2)$$

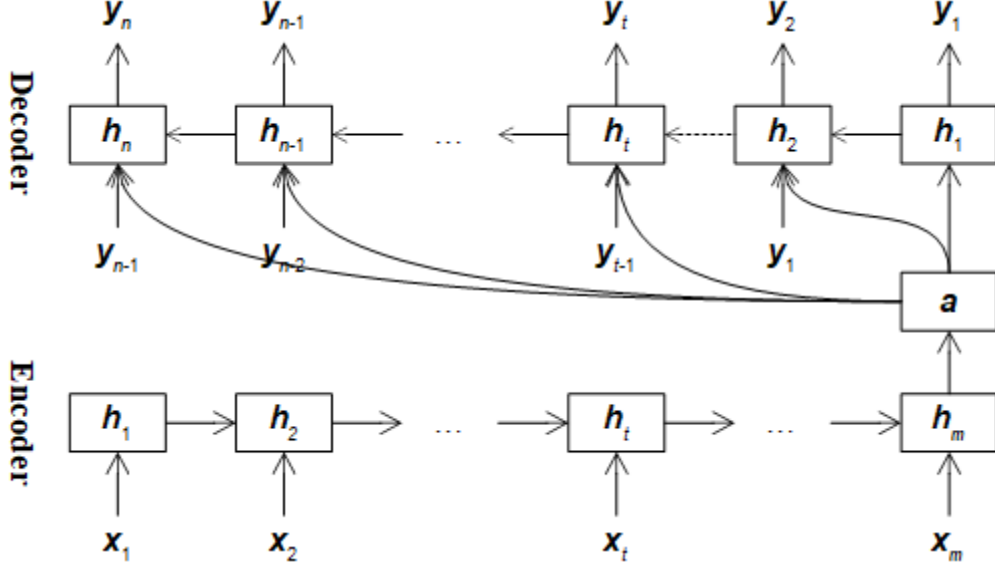
The generation model above has only one RNN because source sequence  $X$  and target sequence  $Y$  have the same length. Some real applications, especially STM applications, require that lengths of  $X$  and  $Y$  are different,  $m \neq n$ . This problem is called different-length problem.

$$X = (x_1, x_2, \dots, x_m)^T$$

$$Y = (y_1, y_2, \dots, y_n)^T$$

Solution for different-length problem is to specify two RNNs: a RNN called *encoder* for  $X$  generation and the other one called *decoder* for  $Y$  generation. Intermediate vector  $a$  is proposed to connect encoder and decoder, which is called context vector in literature (Cho, et al., 2014, p. 2). The encoder-decoder mechanism is an important progressive step in STM as well as generative artificial intelligence (GenAI) because there is no requirement of mapping token-by-token between two sequences  $X$  and  $Y$ , which is much more important than solving the different-length problem. On the other hand, sequence generation as well as its advanced development – transformer can also be classified into domain of GenAI.

According to Cho et al. (Cho, et al., 2014), context variable  $a$ , which is last output of encoder, becomes input of decoder. Following figure depicts encoder-decoder model proposed by Cho et al. (Cho, et al., 2014, p. 2) with note that context vector  $a$  has fixed length.



**Figure 2.2.** Encoder-decoder model with fixed-length context

Note, both context and current token  $t$  are inputs of next token  $t+1$ . Moreover, there is an assignment  $y_{t+1} = o_t$ . Therefore, each  $t^{\text{th}}$  state of decoder is modified as follows:

$$h_t = \sigma_h(W_h y_t + U_h h_{t-1} + V_h a + b_h)$$

$$y_{t+1} = \hat{y}_{t+1} = o_t = \sigma_o(W_o h_t + b_o)$$

Where  $V_h$  is weight matrix for context variable  $a$ . Moreover, it may be not required to calculate output for each  $t^{\text{th}}$  state of encoder. It may be only necessary to calculate hidden value of encoder.

$$h_t^{\text{encoder}} = \sigma_h^{\text{encoder}}(W_h^{\text{encoder}} x_t + U_h^{\text{encoder}} h_{t-1}^{\text{encoder}} + b_h^{\text{encoder}})$$

In STM, given source sequence  $X$  and  $t$  target tokens  $y_1, y_2, \dots, y_t$ , it is necessary to predict the next target token  $y_{t+1}$ . In other words, predictive probability  $P(y_{t+1} | \Theta, X, y_1, y_2, \dots, y_t)$  needs to be maximized so as to obtain  $y_{t+1}$ . Predictive probability  $P(y_{t+1} | \Theta, X, y_1, y_2, \dots, y_t)$  is called *likelihood* at the  $t^{\text{th}}$  state of decoder. Consequently, parameter  $\Theta$  of encoder-decoder model is maximizer of such likelihood.

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmax}} P(y_{t+1} | \Theta, X, y_1, y_2, \dots, y_t)$$

Note, parameter  $\Theta$  represents weight matrices and biases of RNN. By support of RNN and context vector  $a$  with implication of Markov property, likelihood  $P(y_{t+1} | \Theta, X, y_1, y_2, \dots, y_t)$  can become simpler:

$$P(y_{t+1} | \Theta, X, y_1, y_2, \dots, y_t) = P(y_{t+1} | \Theta, X, y_{\leq t}) = P(y_{t+1} | \Theta, a, y_t) = P(o_t | \Theta, a, y_t)$$

Likelihood  $P(y_{t+1} | \Theta, X, y_1, y_2, \dots, y_t)$ , which represents statistical language model, is object of maximum likelihood estimation (MLE) method for training encoder-decoder model (Cho, et al., 2014, p. 2). For example, the likelihood can be approximated by standard normal distribution, which is equivalent to square error function, as follows:

$$P(y_{t+1} | \Theta, X, y_1, y_2, \dots, y_t) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(f(X, y_1, y_2, \dots, y_t | \Theta) - y_{t+1})^2}{2}\right)$$

Where  $f(X, y_1, y_2, \dots, y_t | \Theta)$  denotes encoder-decoder chain.

$$f(X, y_1, y_2, \dots, y_t | \Theta) = \hat{y}_{t+1} = o_t = \sigma_o(W_o h_t + b_o)$$

Therefore, training encoder-decoder model begins with MLE associated with backpropagation algorithm and SGD from decoder back to encoder.

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmax}} \log(P(y_{t+1} | \Theta, X, y_1, y_2, \dots, y_t)) = \underset{\Theta}{\operatorname{argmin}} (f(X, y_1, y_2, \dots, y_t | \Theta) - y_{t+1})^2$$

Alternately, in STM with predefined word vocabulary, a simple but effective way to train encoder-decoder model is to replace likelihood  $P(\mathbf{y}_{t+1} \mid \Theta, \mathbf{X}, \mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_t)$  by a so-called linear component which is a feedforward network (FFN). Exactly, FNN maps the  $(t+1)^{\text{th}}$  target token specified by token vector  $\mathbf{y}_{t+1}$  to a weight vector  $\mathbf{w}$  whose each element  $w_i$  ( $0 \leq w_i \leq 1$ ) is weight of  $i^{\text{th}}$  token (Alammar, 2018).

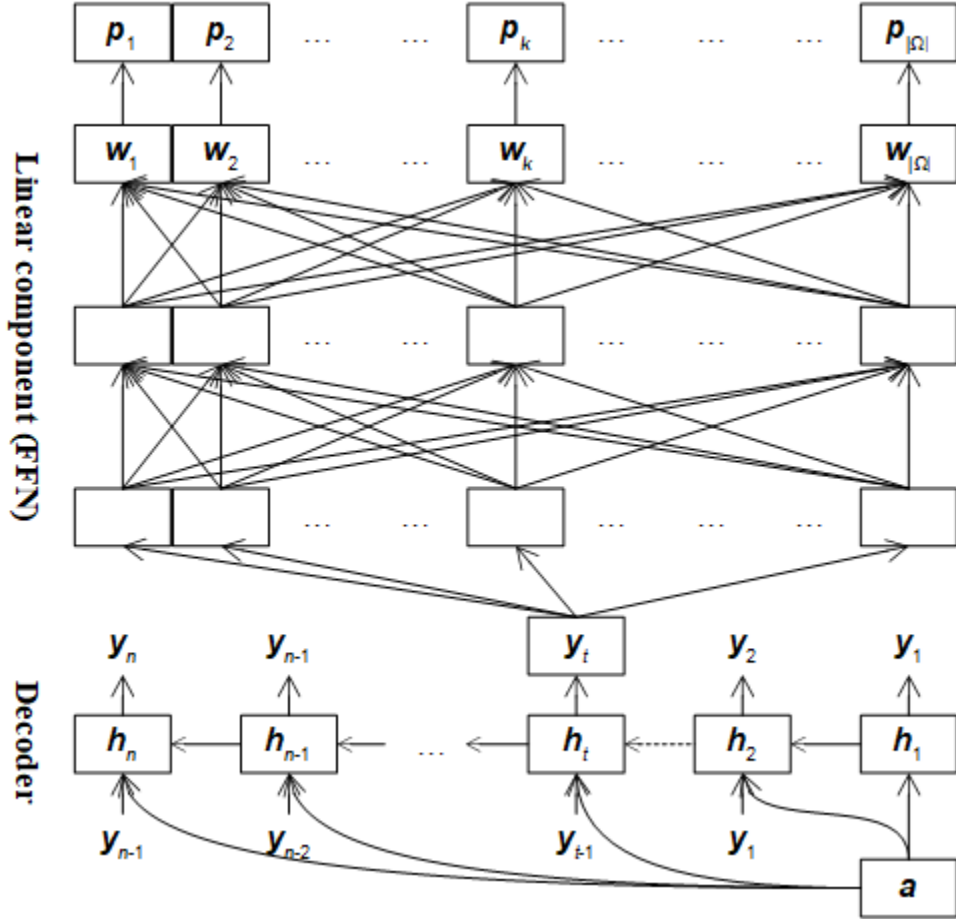
$$\mathbf{w} = (w_1, w_2, \dots, w_{|\Omega|})^T = \text{linear}(\mathbf{y}_{t+1}) = \text{FFN}(\mathbf{y}_{t+1}) = \text{FFN}(\mathbf{o}_t)$$

Length of weight vector  $\mathbf{w}$  is the cardinality  $|\Omega|$  where  $\Omega$  is the vocabulary containing all tokens. After token weight vector  $\mathbf{w}$  is determined, it is easily converted into output probability vector  $\mathbf{p} = (p_1, p_2, \dots, p_{|\Omega|})^T$  where each element  $p_i$  is probability of the  $i^{\text{th}}$  token in vocabulary given the  $(t+1)^{\text{th}}$  target token (Alammar, 2018).

$$\mathbf{p} = \text{softmax}(\mathbf{w}) = (p_1, p_2, \dots, p_{|\Omega|})^T$$

$$p_i = \frac{\exp(w_i)}{\sum_{l=1}^{|\Omega|} \exp(w_l)}$$

Following figure depicts linear component.



**Figure 2.3.** Linear component of encoder-decoder model

It is interesting that likelihood  $P(\mathbf{y}_{t+1} \mid \Theta, \mathbf{X}, \mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_t)$  can be defined as output probability vector  $\mathbf{p} = (p_1, p_2, \dots, p_{|\Omega|})^T$ . If the  $i^{\text{th}}$  token is issued, its probability  $p_i$  is 1 and other probabilities are 0.

$$P(\mathbf{y}_{t+1} \mid \Theta, \mathbf{X}, \mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_t) \stackrel{\text{def}}{=} (p_1, p_2, \dots, p_{|\Omega|})^T$$

Consequently, training encoder-decoder model begins with training linear component  $\text{FFN}(\mathbf{y}_{t+1})$  back to training decoder back to training encoder, which follows backpropagation algorithm

associated stochastic gradient descent (SGD) method. Concretely, the following *cross-entropy*  $L(\mathbf{p} | \Theta)$  is minimized so as to train  $\text{FFN}(\mathbf{y}_{t+1})$ .

$$L(\mathbf{p} | \Theta) = - \sum_{i=1}^{|\Omega|} q_i \log(p_i)$$

Where  $\Theta$  is parameter of  $\text{FFN}(\mathbf{y}_{t+1})$  and the vector  $\mathbf{q} = (q_1, q_2, \dots, q_{|\Omega|})^T$  is binary vector from sample whose each element  $q_i$  has binary values  $\{0, 1\}$  indicating whether the  $i^{\text{th}}$  token/word exists. For example, give sequence/sentence (“I”, “am”, “a”, “student”)<sup>T</sup>, if there is only one token/word “I” in sample sentence, the binary vector will be  $\mathbf{q} = (1, 0, 0, 0)^T$ . If three words “I”, “am”, and “student” are mutually existent, the binary vector will be  $\mathbf{q} = (1, 1, 0, 1)^T$ . When SGD is applied into minimizing the cross-entropy, partial gradient of  $L(\mathbf{p} | \Theta)$  with regard to  $w_j$  is:

$$\frac{\partial L(\mathbf{p} | \Theta)}{\partial w_j} = - \sum_{i=1}^{|\Omega|} q_i (\delta_{ij} - p_j)$$

Where,

$$\delta_{ij} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}$$

Proof,

$$\frac{\partial L(\mathbf{p} | \Theta)}{\partial w_j} = - \sum_{i=1}^{|\Omega|} q_i \frac{\partial \log(p_i)}{\partial w_j} = - \sum_{i=1}^{|\Omega|} q_i \frac{d \log(p_i)}{dp_i} \frac{\partial p_i}{\partial w_j} = - \sum_{i=1}^{|\Omega|} q_i \frac{1}{p_i} \frac{\partial p_i}{\partial w_j}$$

Due to:

$$\frac{\partial p_i}{\partial w_j} = \frac{\partial}{\partial w_j} \left( \frac{\exp(w_i)}{\sum_{l=1}^{|\Omega|} \exp(w_l)} \right) = p_i (\delta_{ij} - p_j)$$

We obtain:

$$\frac{\partial L(\mathbf{p} | \Theta)}{\partial w_j} = - \sum_{i=1}^{|\Omega|} q_i (\delta_{ij} - p_j) \blacksquare$$

So that gradient of  $L(\mathbf{p} | \Theta)$  with regard to  $\mathbf{w}$  is:

$$\frac{dL(\mathbf{p} | \Theta)}{d\mathbf{w}} = \left( \frac{\partial L(\mathbf{p} | \Theta)}{\partial w_1}, \frac{\partial L(\mathbf{p} | \Theta)}{\partial w_2}, \dots, \frac{\partial L(\mathbf{p} | \Theta)}{\partial w_{|\Omega|}} \right)^T$$

Therefore, parameter  $\Theta$  is updated according to SGD associated with backpropagation algorithm:

$$\Theta = \Theta - \gamma \frac{dL(\mathbf{p} | \Theta)}{d\mathbf{w}}$$

Where  $\gamma$  ( $0 < \gamma \leq 1$ ) is *learning rate*. Please pay attention that ordering of source tokens is set from the end token back to the beginning token so that null tokens specified by zero vectors are always in the opening of sequence.

When encoder-decoder model is developed, context vector  $\mathbf{a}$  becomes a so-called *attention*. The main difference between context vector and attention vector is that attention vector is calculated dynamically (customized) for each decoder state. Moreover, that context vector has fixed length restricts its prospect. Anyhow, attention mechanism fosters target sequence to pay attention to source sequence. In general, attention of a decoder state (token) is weighted sum of all encoder states (tokens) with regard to such decoder state. Suppose encoder RNN is denoted as follows:

$$\mathbf{h}_t^{\text{encoder}} = \sigma_h^{\text{encoder}} (\mathbf{W}_h^{\text{encoder}} \mathbf{x}_t + \mathbf{U}_h^{\text{encoder}} \mathbf{h}_{t-1}^{\text{encoder}} + \mathbf{b}_h^{\text{encoder}})$$



$$\mathbf{o}_t^{\text{encoder}} = \sigma_o^{\text{encoder}}(\mathbf{W}_o^{\text{encoder}}\mathbf{h}_t^{\text{encoder}} + \mathbf{b}_o^{\text{encoder}})$$

For convenience, let  $s_1, s_2, \dots, s_m$  denote  $m$  outputs of encoder such that:

$$\begin{aligned} s_1 &= \mathbf{o}_1^{\text{encoder}} \\ s_2 &= \mathbf{o}_2^{\text{encoder}} \\ &\vdots \\ s_m &= \mathbf{o}_m^{\text{encoder}} \end{aligned}$$

Let  $\text{score}(s_i, \mathbf{h}_t)$  be score of encoder output  $s_i$  and decoder hidden  $\mathbf{h}_t$  where  $\text{score}(s_i, \mathbf{h}_t)$  measures how much the  $i^{\text{th}}$  token of source sequence modeled by encoder is close to the  $t^{\text{th}}$  token of target sequence modeled by decoder. As usual,  $\text{score}(s_i, \mathbf{h}_t)$  is defined as dot product of  $s_i$  and  $\mathbf{h}_t$  (Voita, 2023).

$$\text{score}(s_i, \mathbf{h}_t) = \mathbf{s}_i^T \mathbf{h}_t$$

Where decoder hidden  $\mathbf{h}_t$  is:

$$\mathbf{h}_t = \sigma_h(\mathbf{W}_h \mathbf{y}_t + \mathbf{U}_h \mathbf{h}_{t-1} + \mathbf{b}_h)$$

Let  $\text{weight}(s_i, \mathbf{h}_t)$  be weight of encoder output  $s_i$  and decoder hidden  $\mathbf{h}_t$  over  $m$  states of encoder, which is calculated based on soft-max function (Voita, 2023):

$$\text{weight}(s_i, \mathbf{h}_t) = \frac{\exp(\text{score}(s_i, \mathbf{h}_t))}{\sum_{j=1}^m \exp(\text{score}(s_j, \mathbf{h}_t))}$$

As a result, let  $\mathbf{a}_t$  be attention of source sequence  $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)^T$  with regard to the  $t^{\text{th}}$  token of target sequence  $\mathbf{Y} = (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n)^T$ , which is weighted sum of all encoder outputs with regard to such  $t^{\text{th}}$  target token (Voita, 2023).

$$\mathbf{a}_t = \sum_{i=1}^m \text{weight}(s_i, \mathbf{h}_t) \mathbf{s}_i$$

Obviously,  $\mathbf{a}_t$  becomes one of inputs of the  $t^{\text{th}}$  token of target sequence  $\mathbf{Y} = (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n)^T$  such that:

$$\mathbf{o}_t = \sigma_o(\mathbf{W}_o \mathbf{h}_t + \mathbf{V}_o \mathbf{a}_t + \mathbf{b}_o)$$

Where  $\mathbf{V}_o$  is weight matrix of attention  $\mathbf{a}_t$ . In general, decoder RNN associated with the attention mechanism called Luong attention (Voita, 2023) is specified as follows:

$$\mathbf{h}_t = \sigma_h(\mathbf{W}_h \mathbf{y}_t + \mathbf{U}_h \mathbf{h}_{t-1} + \mathbf{b}_h)$$

$$\mathbf{a}_t = \sum_{i=1}^m \text{weight}(s_i, \mathbf{h}_t) \mathbf{s}_i$$

$$\mathbf{y}_{t+1} = \hat{\mathbf{y}}_{t+1} = \mathbf{o}_t = \sigma_o(\mathbf{W}_o \mathbf{h}_t + \mathbf{V}_o \mathbf{a}_t + \mathbf{b}_o)$$

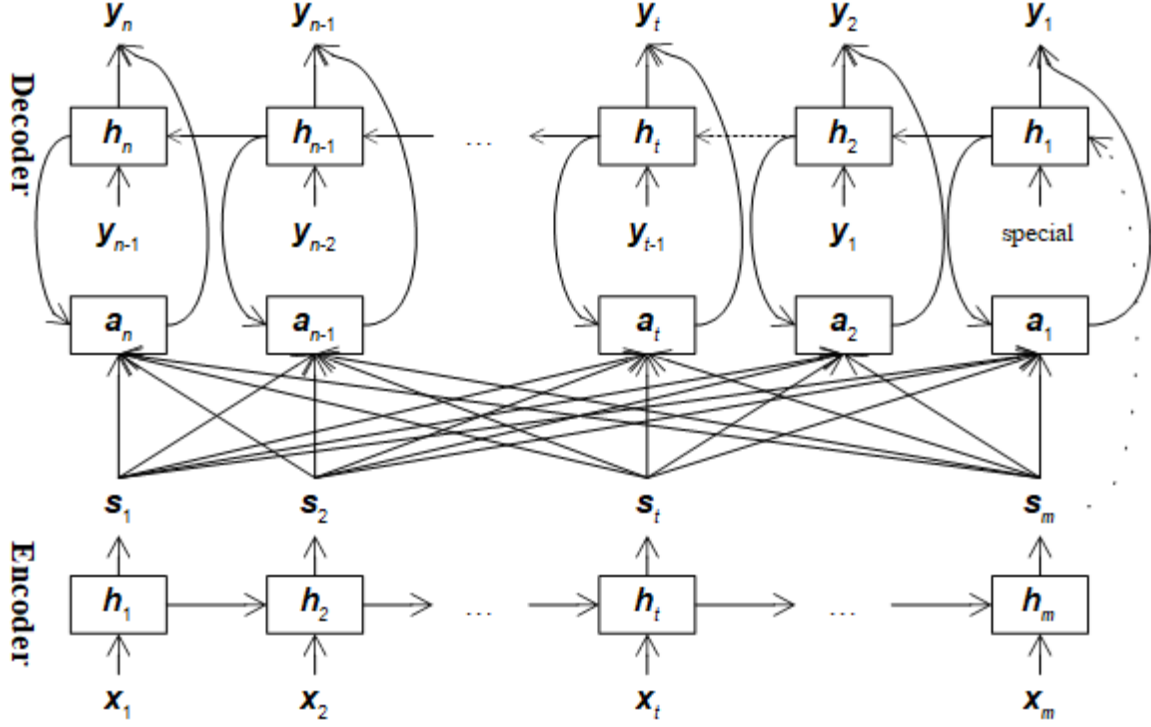
Where,

$$\mathbf{s}_i = \mathbf{o}_i^{\text{encoder}}$$

$$\text{score}(s_i, \mathbf{h}_t) = \mathbf{s}_i^T \mathbf{h}_t$$

$$\text{weight}(s_i, \mathbf{h}_t) = \frac{\exp(\text{score}(s_i, \mathbf{h}_t))}{\sum_{j=1}^m \exp(\text{score}(s_j, \mathbf{h}_t))}$$

Following figure depicts encoder-decoder model with attention (Voita, 2023):



**Figure 2.4.** Encoder-decoder model with attention

Training encoder-decoder model with support attention is still based on likelihood maximization or linear component aforementioned. Attention mechanism mentioned here does not ever concern internal meaning of every token, which only fosters target sequence to pay attention at source sequence. The attention that concerns internal meanings of tokens is called self-attention which is an advancement of attention. In other words, self-attention fosters source sequence to pay attention to itself. Transformer mentioned in the next section will implement self-attention.

### 3. Transformer

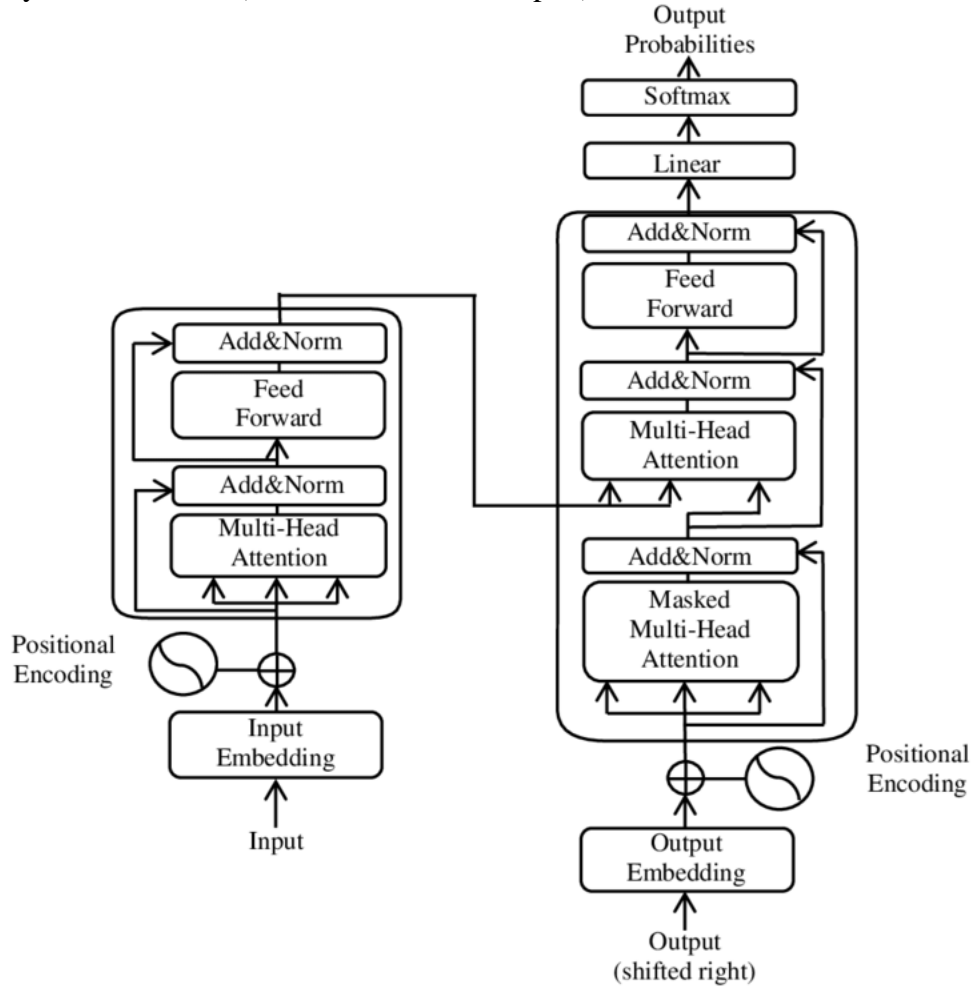
*Transformer*, developed by Vaswani et al. (Vaswani, et al., 2017) in the famous paper “Attention Is All You Need”, has also attention mechanism and encoder-decoder mechanism like the aforementioned generation model that applies recurrent neural network (RNN) and short-term memory (LSTM) but transformer does not require to process successively tokens of sequence in token-by-token ordering, which improves translation speed. Moreover, another strong point of transformer is that it has *self-attention* which is the special attention that concerns internal meanings of its own tokens. Transformer supports both attention and self-attention, which fosters target sequence to pay attention to both source sequence and target sequence and also fosters source sequence to pay attention to itself. Besides, transformer does not apply RNN / LSTM. Note that word and sentence in natural language processing (NLP) are mentioned as token and sequence, respectively by a convention, so that source sequence  $X$  is fed to encoder and target sequence  $Y$  is fed to decoder where  $X$  and  $Y$  are concerned exactly as matrices.

$$X = (x_1, x_2, \dots, x_m)^T$$

$$Y = (y_1, y_2, \dots, y_n)^T$$

Each encoder as well as each decoder in transformer are composed of some identical layers. The number of layer which is developed by Vaswani et al. (Vaswani, et al., 2017, p. 3) is 6. Each encoder layer has two sublayers which are *multi-head attention* sublayer and feedforward sublayer

whereas each decoder layer has three sublayers which are masked multi-head attention sublayer, multi-head attention sublayer, and feedforward sublayer. Every sublayer is followed by association of *residual mechanism* and *layer normalization*, denoted as  $\text{Add \& Norm} = \text{LayerNorm}(X + \text{Sublayer}(X))$ . The residual mechanism means that sublayer  $\text{Sublayer}(X)$  is added with its input as the sum  $X + \text{Sublayer}(X)$ . Note,  $\text{Sublayer}(X)$  can be attention sublayer or feedforward sublayer. The layer normalization is to normalize such sum. Following figure summarizes transformer developed by Vaswani et al. (Vaswani, et al., 2017, p. 3).



**Figure 3.1.** Architecture of transformer

Feedforward sublayer also called feedforward network (FNN) aims to fine-tune attention by increasing degree of complication.

Encoder and its attention are described firstly when multi-head attention is derived from basic concept of attention. Attention (self-attention) proposed by Vaswani et al. (Vaswani, et al., 2017) is based on three important matrices such as query matrix  $Q$ , key matrix  $K$ , and value matrix  $V$ . The number of rows of these matrices is  $m$  which is the number of tokens in sequence matrix  $X = (x_1, x_2, \dots, x_m)^T$  but the number of columns of query matrix  $Q$  and key matrix  $K$  is  $d_k$  whereas the number of columns of value matrix  $V$  is  $d_v$ . The number  $m$  of token is set according to concrete applications, which is often the number of words of the longest sentence. In literature (Vaswani, et al., 2017),  $d_k$  and  $d_v$  are called key dimension and value dimension, respectively. Dimensions of matrices  $Q$ ,  $K$ , and  $V$  are  $m \times d_k$ ,  $m \times d_k$ , and  $m \times d_v$ , respectively (Vaswani, et al., 2017), (Wikipedia, Transformer (deep learning architecture), 2019).

$$Q = \begin{pmatrix} \mathbf{q}_1 \\ \mathbf{q}_2 \\ \vdots \\ \mathbf{q}_m \end{pmatrix} = \begin{pmatrix} q_{11} & q_{12} & \cdots & q_{1d_k} \\ q_{21} & q_{22} & \cdots & q_{2d_k} \\ \vdots & \vdots & \ddots & \vdots \\ q_{m1} & q_{m2} & \cdots & q_{md_k} \end{pmatrix}, K = \begin{pmatrix} \mathbf{k}_1 \\ \mathbf{k}_2 \\ \vdots \\ \mathbf{k}_m \end{pmatrix} = \begin{pmatrix} k_{11} & k_{12} & \cdots & k_{1d_k} \\ k_{21} & k_{22} & \cdots & k_{2d_k} \\ \vdots & \vdots & \ddots & \vdots \\ k_{m1} & k_{m2} & \cdots & k_{md_k} \end{pmatrix}$$

$$V = \begin{pmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \vdots \\ \mathbf{v}_m \end{pmatrix} = \begin{pmatrix} v_{11} & v_{12} & \cdots & v_{1d_v} \\ v_{21} & v_{22} & \cdots & v_{2d_v} \\ \vdots & \vdots & \ddots & \vdots \\ v_{m1} & v_{m2} & \cdots & v_{md_v} \end{pmatrix}$$

Where,

$$\mathbf{q}_i = (q_{i1}, q_{i2}, \dots, q_{id_k})^T, \mathbf{k}_i = (k_{i1}, k_{i2}, \dots, k_{id_k})^T$$

$$\mathbf{v}_i = (v_{i1}, v_{i2}, \dots, v_{id_v})^T$$

Suppose every token vector  $\mathbf{x}_i$  in sequence matrix  $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m)^T$  has  $d_m$  elements such that  $d_m$  is called model dimension which is often 512 in NLP.

$$\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{id_m})^T$$

Query matrix  $Q$ , key matrix  $K$ , and value matrix  $V$  are determined by products of sequence matrix  $\mathbf{X}$  and query weight matrix  $W^Q$ , key weight matrix  $W^K$ , value weight matrix  $W^V$ .

$$Q = \mathbf{X}W^Q$$

$$K = \mathbf{X}W^K$$

$$V = \mathbf{X}W^V$$

Of course, dimensions of weight matrices  $W^Q$ ,  $W^K$ , and  $W^V$  are  $d_m \times d_k$ ,  $d_m \times d_k$ , and  $d_m \times d_v$ , respectively. All of them have  $d_m$  rows. Matrices  $W^Q$  and  $W^K$  have  $d_k$  columns whereas matrix  $W^V$  have  $d_v$  columns.

$$W^Q = \begin{pmatrix} w_{11}^Q & w_{12}^Q & \cdots & w_{1d_k}^Q \\ w_{21}^Q & w_{22}^Q & \cdots & w_{2d_k}^Q \\ \vdots & \vdots & \ddots & \vdots \\ w_{d_m1}^Q & w_{d_m2}^Q & \cdots & w_{d_md_k}^Q \end{pmatrix}, W^K = \begin{pmatrix} w_{11}^K & w_{12}^K & \cdots & w_{1d_k}^K \\ w_{21}^K & w_{22}^K & \cdots & w_{2d_k}^K \\ \vdots & \vdots & \ddots & \vdots \\ w_{d_m1}^K & w_{d_m2}^K & \cdots & w_{d_md_k}^K \end{pmatrix}$$

$$W^V = \begin{pmatrix} w_{11}^V & w_{12}^V & \cdots & w_{1d_v}^V \\ w_{21}^V & w_{22}^V & \cdots & w_{2d_v}^V \\ \vdots & \vdots & \ddots & \vdots \\ w_{d_m1}^V & w_{d_m2}^V & \cdots & w_{d_md_v}^V \end{pmatrix}$$

Attention is calculated based on scaled product of query matrix  $Q$ , key matrix  $K$ , and value matrix  $V$  in order to make effects on value matrix  $V$  specifying real sequence by probabilities and moreover, these probabilities are calculated by matching query matrix  $Q$  specifying query sequence and key matrix  $K$  specifying key sequence, which is similar to searching mechanism. These probabilities are also based on soft-max function, which implies weights too. Moreover, attention focuses on all tokens of sequence, which improves meaningful context of sentence in NLP. Given matrices  $Q$ ,  $K$ , and  $V$ , attention of  $Q$ ,  $K$ , and  $V$  is specified as follows:

$$\text{Attention}(\mathbf{X}) = \text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Note, the subscript “ $T$ ” denotes vector/matrix transposition operator. It is easy to recognize this attention is self-attention of only one sequence  $\mathbf{X}$  via  $Q$ ,  $K$ , and  $V$  which are essentially calculated from  $\mathbf{X}$  and weight matrices  $W^Q$ ,  $W^K$ , and  $W^V$ . Note, self-attention concerns internal meanings of

its own tokens. Transformer here fosters source sequence to pay attention to itself. The reason of dividing product  $QK^T$  by the scaling factor  $\sqrt{d_k}$  is to improve convergence speed in training transformer. Before explaining how to calculate weight / probability matrix  $\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)$ , it is necessary to skim the product  $QK^T$  of query matrix  $Q$  and key matrix  $K$  which aims to match query sequence and key sequence.

$$QK^T = \begin{pmatrix} \frac{\mathbf{q}_1 \mathbf{k}_1^T}{\sqrt{d_k}} & \frac{\mathbf{q}_1 \mathbf{k}_2^T}{\sqrt{d_k}} & \dots & \frac{\mathbf{q}_1 \mathbf{k}_m^T}{\sqrt{d_k}} \\ \frac{\mathbf{q}_2 \mathbf{k}_1^T}{\sqrt{d_k}} & \frac{\mathbf{q}_2 \mathbf{k}_2^T}{\sqrt{d_k}} & \dots & \frac{\mathbf{q}_2 \mathbf{k}_m^T}{\sqrt{d_k}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\mathbf{q}_m \mathbf{k}_1^T}{\sqrt{d_k}} & \frac{\mathbf{q}_m \mathbf{k}_2^T}{\sqrt{d_k}} & \dots & \frac{\mathbf{q}_m \mathbf{k}_m^T}{\sqrt{d_k}} \end{pmatrix}$$

The dot product  $\mathbf{q}_i \mathbf{k}_j^T$  which indicates how much the query vector  $\mathbf{q}_i$  matches or attends mutually the key vector  $\mathbf{k}_j$  is specified as follows:

$$\mathbf{q}_i \mathbf{k}_j^T = \sum_{l=1}^{d_k} q_{il} k_{jl}$$

Probability matrix  $\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)$  is specified as follows:

$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) = \begin{pmatrix} \mathbf{p}_1 \\ \mathbf{p}_2 \\ \vdots \\ \mathbf{p}_m \end{pmatrix} = \begin{pmatrix} \text{softmax}\left(\frac{\mathbf{q}_1 \mathbf{k}_1^T}{\sqrt{d_k}}, \frac{\mathbf{q}_1 \mathbf{k}_2^T}{\sqrt{d_k}}, \dots, \frac{\mathbf{q}_1 \mathbf{k}_m^T}{\sqrt{d_k}}\right) \\ \text{softmax}\left(\frac{\mathbf{q}_2 \mathbf{k}_1^T}{\sqrt{d_k}}, \frac{\mathbf{q}_2 \mathbf{k}_2^T}{\sqrt{d_k}}, \dots, \frac{\mathbf{q}_2 \mathbf{k}_m^T}{\sqrt{d_k}}\right) \\ \vdots \\ \text{softmax}\left(\frac{\mathbf{q}_m \mathbf{k}_1^T}{\sqrt{d_k}}, \frac{\mathbf{q}_m \mathbf{k}_2^T}{\sqrt{d_k}}, \dots, \frac{\mathbf{q}_m \mathbf{k}_m^T}{\sqrt{d_k}}\right) \end{pmatrix}$$

The  $i^{\text{th}}$  row of probability matrix  $\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)$  includes weights / probabilities that the  $i^{\text{th}}$  token is associated with all tokens including itself with note that  $\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)$  is  $m \times m$  matrix, specified by weight/probability vector  $\mathbf{p}_i$ . It is necessary to explain the  $i^{\text{th}}$  row of probability matrix  $\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)$  which is the following row vector:

$$\text{softmax}\left(\frac{\mathbf{q}_i \mathbf{k}_1^T}{\sqrt{d_k}}, \frac{\mathbf{q}_i \mathbf{k}_2^T}{\sqrt{d_k}}, \dots, \frac{\mathbf{q}_i \mathbf{k}_m^T}{\sqrt{d_k}}\right) = \mathbf{p}_i = (p_{i1}, p_{i2}, \dots, p_{im})$$

Each probability  $p_{ij}$ , which is weight indeed, is calculated by soft-max function as follows:

$$p_{ij} = \frac{\exp\left(\frac{\mathbf{q}_i \mathbf{k}_j^T}{\sqrt{d_k}}\right)}{\sum_{l=1}^m \exp\left(\frac{\mathbf{q}_i \mathbf{k}_l^T}{\sqrt{d_k}}\right)}, \forall j = \overline{1, m}$$

Where  $\exp(\cdot)$  is natural exponential function. Therefore, probability matrix  $\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)$  is totally determined:

$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) = \begin{pmatrix} \mathbf{p}_1 \\ \mathbf{p}_2 \\ \vdots \\ \mathbf{p}_m \end{pmatrix} = \begin{pmatrix} p_{11} & p_{12} & \cdots & p_{1m} \\ p_{21} & p_{22} & \cdots & p_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ p_{m1} & p_{m2} & \cdots & p_{mm} \end{pmatrix}$$

Where,

$$\mathbf{p}_i = (p_{i1}, p_{i2}, \dots, p_{im})^T, \forall i = \overline{1, m}$$

Self-attention of  $Q$ ,  $K$ , and  $V$  is totally determined as follows:

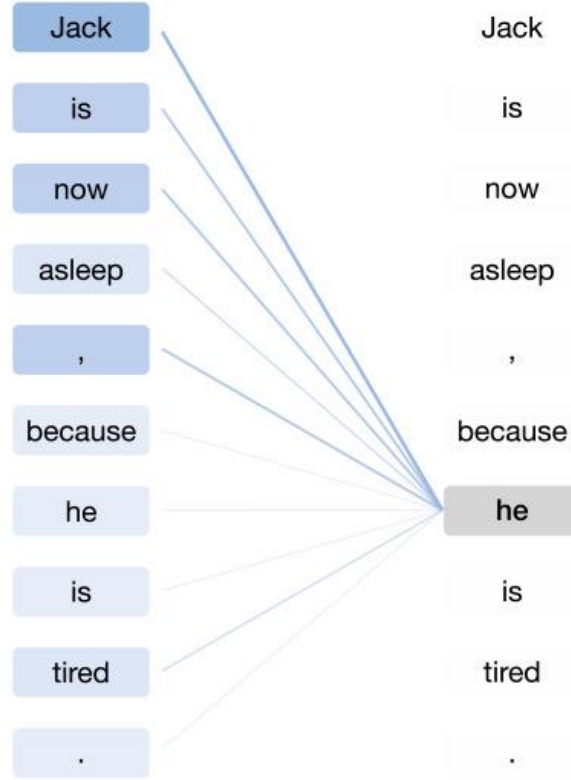
$$\begin{aligned} \text{Attention}(\mathbf{X}) &= \text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V = \begin{pmatrix} \mathbf{p}_1 \\ \mathbf{p}_2 \\ \vdots \\ \mathbf{p}_m \end{pmatrix} (\bar{\mathbf{v}}_1, \bar{\mathbf{v}}_2, \dots, \bar{\mathbf{v}}_{d_v}) \\ &= \begin{pmatrix} \mathbf{p}_1 \bar{\mathbf{v}}_1^T & \mathbf{p}_1 \bar{\mathbf{v}}_2^T & \cdots & \mathbf{p}_1 \bar{\mathbf{v}}_{d_v}^T \\ \mathbf{p}_2 \bar{\mathbf{v}}_1^T & \mathbf{p}_2 \bar{\mathbf{v}}_2^T & \cdots & \mathbf{p}_2 \bar{\mathbf{v}}_{d_v}^T \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{p}_m \bar{\mathbf{v}}_1^T & \mathbf{p}_m \bar{\mathbf{v}}_2^T & \cdots & \mathbf{p}_m \bar{\mathbf{v}}_{d_v}^T \end{pmatrix} = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1d_v} \\ a_{21} & a_{22} & \cdots & a_{2d_v} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{md_v} \end{pmatrix} \end{aligned}$$

Where,

$$\bar{\mathbf{v}}_j = \begin{pmatrix} v_{1j} \\ v_{2j} \\ \vdots \\ v_{mj} \end{pmatrix}$$

$$a_{ij} = \mathbf{p}_i \bar{\mathbf{v}}_j^T = \sum_{k=1}^m p_{ik} v_{kj}$$

Note,  $\bar{\mathbf{v}}_j$  is the  $j^{\text{th}}$  column vector of value matrix  $V$ . Of course, dimension of self-attention  $\text{Attention}(Q, K, V)$  is  $m \times d_v$  having  $m$  rows and  $d_v$  columns. Attention  $\text{Attention}(Q, K, V)$  is also called scaled dot product attention because of dot product  $\mathbf{q}_i \mathbf{k}_j^T$  and scaling factor  $\sqrt{d_k}$ . Each row  $\mathbf{a}_i = (a_{i1}, a_{i2}, \dots, a_{id_v})^T$  of  $\text{Attention}(Q, K, V)$ , which is a  $d_v$ -length vector, is self-attention of the  $i^{\text{th}}$  token which is contributed by all tokens via scaled dot products  $QK^T$ . Therefore, the preeminence of self-attention is that self-attention concerns all tokens in detail instead of concerning only sequence and the self-attention  $\mathbf{a}_i = (a_{i1}, a_{i2}, \dots, a_{id_v})^T$  of the  $i^{\text{th}}$  token is attended by all tokens. For example, given sentence “Jack is now asleep, because he is tired.”, the word “he” is strongly implied to the word “Jack” by self-attention of the word “he” although the word “he” is ambiguous. Following figure (Han, et al., 2021, p. 231) illustrates the self-attention of the word “he” in which each strength of implication of another word (accept itself “he”) to the word “he” is indicated by strong degree of connection color.



**Figure 3.2.** Self-attention example

Vaswani et al. (Vaswani, et al., 2017) proposed an improvement of attention called multi-head attention which is concatenation of many attentions. The existence of many attentions aims to discover as much as possible different meanings under attentions and the concatenation mechanism aims to unify different attentions into one self-attention. Following equation specifies multi-head attention with note that the multi-head attention here is self-attention.

$$\text{MultiheadAttention}(\mathbf{X}) = \text{concatenate}(\text{head}_1, \text{head}_2, \dots, \text{head}_h) \mathbf{W}^O$$

Where,

$$\text{head}_i = \text{Attention}(Q_i, K_i, V_i)$$

$$Q_i = \mathbf{X} \mathbf{W}_i^Q$$

$$K_i = \mathbf{X} \mathbf{W}_i^K$$

$$V_i = \mathbf{X} \mathbf{W}_i^V$$

Of course,  $\mathbf{W}_i^Q$ ,  $\mathbf{W}_i^K$ , and  $\mathbf{W}_i^V$  are query weight matrix, key weight matrix, and value weight matrix for the  $i^{\text{th}}$  head, respectively whereas  $\mathbf{W}^O$  is the entire weight matrix whose dimension is often set as  $hd_v \times d_m$  so that multi-head attention  $\text{MultiheadAttention}(\mathbf{X})$  is  $m \times d_m$  matrix which is the same to dimension of input sequence matrix  $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m)^T$ . Note that the concatenation mechanism follows horizontal direction so that the concatenation  $\text{concatenate}(\text{head}_1, \text{head}_2, \dots, \text{head}_h)$  is  $m \times hd_v$  matrix when each head  $\text{head}_i = \text{Attention}(Q_i, K_i, V_i)$  is  $m \times d_v$  matrix. There are  $h$  heads (attentions) in the equation above. In practice,  $h$  is set so that  $hd_v = d_m$  which is model dimension. Recall that  $d_m$  is often 512 in NLP. For easy illustration, the concatenation of  $h$  attentions is represented as  $m \times hd_v$  as follows:

$$\text{concatenate}(\text{head}_1, \dots, \text{head}_i, \dots, \text{head}_h)$$

$$= \begin{pmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \dots & a_{1d_v}^{(1)} & \dots & a_{11}^{(i)} & a_{12}^{(i)} & \dots & a_{1d_v}^{(i)} & \dots & a_{11}^{(h)} & a_{12}^{(h)} & \dots & a_{1d_v}^{(h)} \\ a_{21}^{(1)} & a_{22}^{(1)} & \dots & a_{2d_v}^{(1)} & \dots & a_{21}^{(i)} & a_{22}^{(i)} & \dots & a_{2d_v}^{(i)} & \dots & a_{21}^{(h)} & a_{22}^{(h)} & \dots & a_{2d_v}^{(h)} \\ \vdots & \vdots & \ddots & \vdots & \dots & \vdots & \vdots & \ddots & \vdots & \dots & \vdots & \vdots & \ddots & \vdots \\ a_{m1}^{(1)} & a_{m2}^{(1)} & \dots & a_{md_v}^{(1)} & \dots & a_{m1}^{(i)} & a_{m2}^{(i)} & \dots & a_{md_v}^{(i)} & \dots & a_{m1}^{(h)} & a_{m2}^{(h)} & \dots & a_{md_v}^{(h)} \end{pmatrix}$$

Obviously, weight matrix  $W^O$  is  $hd_v \times d_m$  matrix so that multi-head attention  $\text{MultiheadAttention}(X)$  is  $m \times d_m$  matrix, as follows:

$$W^O = \begin{pmatrix} \omega_{11}^{(1)} & \omega_{12}^{(1)} & \dots & \omega_{1d_m}^{(1)} \\ \omega_{21}^{(1)} & \omega_{22}^{(1)} & \dots & \omega_{2d_m}^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ \omega_{d_v1}^{(1)} & \omega_{d_v2}^{(1)} & \dots & \omega_{d_vd_m}^{(1)} \\ \vdots & \vdots & \vdots & \vdots \\ \omega_{11}^{(i)} & \omega_{12}^{(i)} & \dots & \omega_{1d_m}^{(i)} \\ \omega_{21}^{(i)} & \omega_{22}^{(i)} & \dots & \omega_{2d_m}^{(i)} \\ \vdots & \vdots & \ddots & \vdots \\ \omega_{d_v1}^{(i)} & \omega_{d_v2}^{(i)} & \dots & \omega_{d_vd_m}^{(i)} \\ \vdots & \vdots & \vdots & \vdots \\ \omega_{11}^{(h)} & \omega_{12}^{(h)} & \dots & \omega_{1d_m}^{(h)} \\ \omega_{21}^{(h)} & \omega_{22}^{(h)} & \dots & \omega_{2d_m}^{(h)} \\ \vdots & \vdots & \ddots & \vdots \\ \omega_{d_v1}^{(h)} & \omega_{d_v2}^{(h)} & \dots & \omega_{d_vd_m}^{(h)} \end{pmatrix}$$

After multi-head attention goes through residual mechanism and layer normalization of attention sublayer, it is fed to feedforward sublayer or feedforward network (FFN) to finish the processing of encoder. Let  $\text{EncoderAttention}(X)$  be output of encoder which is considered as attention:

$$\text{MultiheadAttention}(X) = \text{LayerNorm}(X + \text{MultiheadAttention}(X))$$

$$\text{EncoderAttention}(X)$$

$$= \text{LayerNorm}(\text{MultiheadAttention}(X) + \text{FFN}(\text{MultiheadAttention}(X)))$$

If there is a stack of  $N$  encoders, the process above is repeated  $N$  times. In literature (Vaswani, et al., 2017),  $N$  is set to be 6. Without loss of generality, we can consider  $N = 1$  as simplest case for easy explanations.

Now it is essential to survey how decoder applies encoder attention  $\text{EncoderAttention}(X)$  into its encoding task. Essentially, decoder has two multi-head attentions such as masked multi-head attention and multi-head attention whereas encoder has only one multi-head attention. Their attentions are similar to encoder's attention but there is a slight difference. Firstly, decoder input sequence  $Y = (y_1, y_2, \dots, y_n)^T$  is fed to masked multi-head attention sublayer with note that  $Y$  is  $n \times d_m$  matrix with support that model dimension  $d_m$ , which is often set to be 512 in natural language processing (NLP), may not be changed with regard to decoder. Because masked multi-head attention is composed by concatenation of masked head attentions by the same way of encoder, we should concern masked head attention. Sequence  $Y$  should have  $n = m$  tokens like sequence  $X$  in practice. This is necessary because the length  $m = n$  is the largest number of possible tokens in



any sequence. For shorter sentences in NLP, redundant tokens are represented by zeros. Moreover, most of parameters (weight matrices) of encoder and decoder are independent from  $m$  and  $n$ , especially in the case  $m = n$ .

There is a principle that a token  $y_i$  in sequence  $Y$  does not know its successive tokens  $y_{i+1}, y_{i+2}, \dots, y_n$  with note that these tokens are called unknown tokens for token  $y_i$ , which causes that soft-max function needs to be added a mask matrix  $M$  whose unknown positions are removed by setting them to be negative infinities because evaluation of negative infinite by exponential function is zero. *Masked attention* is self-attention too.

$$\text{MaskedAttention}(Y) = \text{MaskedAttention}(Q, K, V) = \text{softmax}\left(M + \frac{QK^T}{\sqrt{d_k}}\right)V$$

Where masked matrix  $M$  is triangle matrix with negative infinities on upper part and zeros on lower part as follows:

$$M = \begin{pmatrix} 0 & -\infty & -\infty & \dots & -\infty \\ 0 & 0 & -\infty & \dots & -\infty \\ 0 & 0 & 0 & \dots & -\infty \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 0 \end{pmatrix}$$

Note,

$$Q = YW^Q$$

$$K = YW^K$$

$$V = YW^V$$

Where  $W^Q$ ,  $W^K$ , and  $W^V$  are weight matrices with note that they are different from the ones of encoder. Dimensions of weight matrices  $W^Q$ ,  $W^K$ , and  $W^V$  are  $d_m \times d_k$ ,  $d_m \times d_k$ , and  $d_m \times d_v$ , respectively. Dimensions of matrices  $Q$ ,  $K$ , and  $V$  are  $n \times d_k$ ,  $n \times d_k$ , and  $n \times d_v$ , respectively whereas dimension of masked matrix  $M$  is  $n \times d_m$ . We have  $QK^T$  is  $n \times n$  matrix:

$$QK^T = \begin{pmatrix} \frac{q_1 k_1^T}{\sqrt{d_k}} & -\infty & -\infty & \dots & -\infty \\ \frac{q_2 k_1^T}{\sqrt{d_k}} & \frac{q_2 k_2^T}{\sqrt{d_k}} & -\infty & \dots & -\infty \\ \frac{q_3 k_1^T}{\sqrt{d_k}} & \frac{q_3 k_2^T}{\sqrt{d_k}} & \frac{q_3 k_3^T}{\sqrt{d_k}} & \dots & -\infty \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{q_n k_1^T}{\sqrt{d_k}} & \frac{q_n k_2^T}{\sqrt{d_k}} & \frac{q_n k_3^T}{\sqrt{d_k}} & \dots & \frac{q_n k_n^T}{\sqrt{d_k}} \end{pmatrix}$$

Recall that the purpose of masked matrix  $M$  is to remove the affections of current token from its after tokens such that:

$$\text{softmax}\left(M + \frac{QK^T}{\sqrt{d_k}}\right) = \begin{pmatrix} p_1 \\ p_2 \\ p_3 \\ \vdots \\ p_n \end{pmatrix} = \begin{pmatrix} p_{11} & 0 & 0 & \dots & 0 \\ p_{21} & p_{22} & 0 & \dots & 0 \\ p_{31} & p_{32} & p_{33} & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ p_{n1} & p_{n2} & p_{n3} & \dots & p_{nn} \end{pmatrix}$$

Where,

$$p_{ij} = \frac{\exp\left(\frac{\mathbf{q}_i \mathbf{k}_j^T}{\sqrt{d_k}}\right)}{\sum_{l=1}^n \exp\left(\frac{\mathbf{q}_i \mathbf{k}_l^T}{\sqrt{d_k}}\right)}, \forall i \geq j$$

$$p_{ij} = \frac{\exp(-\infty)}{\sum_{l=1}^n (\dots)} = 0, \forall i < j$$

Therefore, masked attention is determined as follows:

$$\text{MaskedAttention}(\mathbf{Y}) = \text{MaskedAttention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\mathbf{M} + \frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$$

$$= \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1d_v} \\ a_{21} & a_{22} & \cdots & a_{2d_v} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nd_v} \end{pmatrix}$$

Where attention element  $a_{ij}$  is calculated by the aforementioned way:

$$a_{ij} = \mathbf{p}_i \bar{\mathbf{v}}_j^T = \sum_{k=1}^n p_{ik} v_{kj}$$

Dimension of masked attention  $\text{MaskedAttention}(\mathbf{Y})$  is  $n \times d_v$  having  $n$  rows and  $d_v$  columns. Following equation specifies masked multi-head attention which is concatenation of some masked attentions.

$$\text{MaskedMultiheadAttention}(\mathbf{Y}) = \text{concatenate}(\text{head}_1, \text{head}_2, \dots, \text{head}_h) \mathbf{W}^O$$

Where,

$$\begin{aligned} \text{head}_i &= \text{MaskedAttention}(\mathbf{Q}_i, \mathbf{K}_i, \mathbf{V}_i) \\ \mathbf{Q}_i &= \mathbf{Y} \mathbf{W}_i^Q \\ \mathbf{K}_i &= \mathbf{Y} \mathbf{W}_i^K \\ \mathbf{V}_i &= \mathbf{Y} \mathbf{W}_i^V \end{aligned}$$

Please pay attention that weights matrices  $\mathbf{W}_i^Q$ ,  $\mathbf{W}_i^K$ ,  $\mathbf{W}_i^V$ , and  $\mathbf{W}^O$  are different from the ones of encoder. Dimensions of  $\mathbf{W}_i^Q$ ,  $\mathbf{W}_i^K$ ,  $\mathbf{W}_i^V$ , and  $\mathbf{W}^O$  are  $d_m \times d_k$ ,  $d_m \times d_k$ ,  $d_m \times d_v$ , and  $hd_v \times d_m$  so that dimension of masked multi-head attention  $\text{MaskedMultiheadAttention}(\mathbf{Y})$  is  $n \times d_m$ . Residual mechanism and layer normalization are applied into masked multi-head attention too:

$$\text{MaskedMultiheadAttention}(\mathbf{Y}) = \text{LayerNorm}(\mathbf{Y} + \text{MaskedMultiheadAttention}(\mathbf{Y}))$$

Because mechanism of multi-head attention of decoder is relatively special, it is called complex multi-head attention for convention. Because complex multi-head attention is composed by concatenation of some complex attentions by the same way of encoder, we should concern complex attention.

$$\text{Attention}(\mathbf{X}, \mathbf{Y}) = \text{ComplexAttention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$$

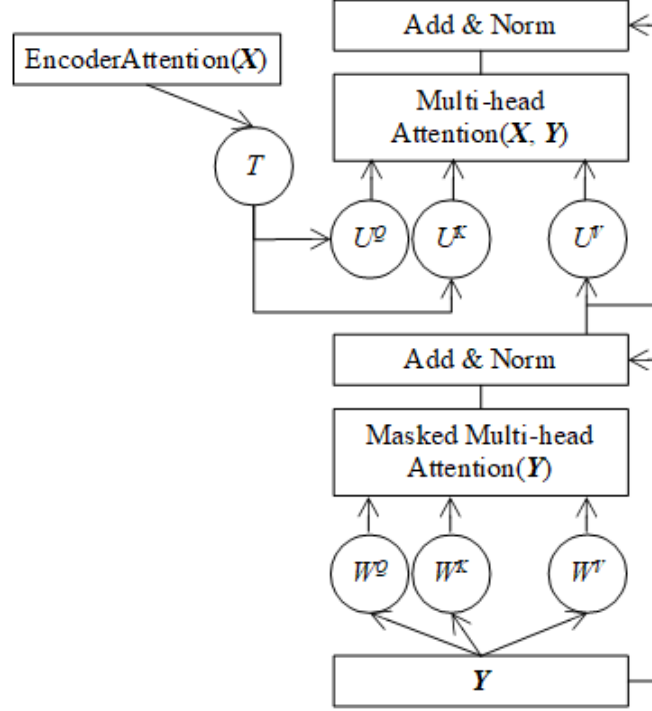
Query matrix  $\mathbf{Q}$  and key matrix  $\mathbf{K}$  of complex attention are products of encoder attention  $\text{EncoderAttention}(\mathbf{X})$  and query weight matrix  $\mathbf{U}^Q$  and key weight matrix  $\mathbf{U}^K$ , respectively.

$$\begin{aligned} \mathbf{Q} &= \text{TEncoderAttention}(\mathbf{X}) \mathbf{U}^Q \\ \mathbf{K} &= \text{TEncoderAttention}(\mathbf{X}) \mathbf{U}^K \end{aligned}$$

Where  $T$  is transformation matrix whose dimension is  $n \times m$ . If  $n = m$ , matrix  $T$  will be removed. Value matrix  $V$  of complex attention is product of masked multi-head attention and value weight matrix  $U^V$ .

$$V = \text{MaskedMultiheadAttention}(Y)U^V$$

Dimensions of weight matrices  $U^Q$ ,  $U^K$ , and  $U^V$  are  $d_m \times d_k$ ,  $d_m \times d_k$ , and  $d_m \times d_v$ , respectively. Following figure depicts  $\text{Attention}(X, Y)$  in general view.



**Figure 3.3.** Decoder attention  $\text{Attention}(X, Y)$  in general view

Transformer here fosters target sequence to pay attention to itself and source sequence by masked self-attention and encoder attention. Of course, after complex attention is calculated, multi-head attention of decoder (complex multi-head attention) is totally determined.

$$\text{MultiheadAttention}(X, Y) = \text{concatenate}(\text{head}_1, \text{head}_2, \dots, \text{head}_h)U^O$$

Where,

$$\text{head}_i = \text{ComplexAttention}(Q_i, K_i, V_i)$$

$$Q_i = T\text{EncoderAttention}(X)U_i^Q$$

$$K_i = T\text{EncoderAttention}(X)U_i^K$$

$$V_i = \text{MaskedMultiheadAttention}(Y)U_i^V$$

Of course,  $U_i^Q$ ,  $U_i^K$ , and  $U_i^V$  are query weight matrix, key weight matrix, and value weight matrix of the  $i^{\text{th}}$  head, respectively whereas  $U^O$  is entire weight matrix and  $T$  is transformation matrix. Because encoder attention  $\text{EncoderAttention}(X)$  is  $m \times d_m$  matrix, dimension of transformation matrix  $T$  is  $n \times m$ . If  $n = m$ , matrix  $T$  will be removed. In practice, it is necessary to set  $n = m$ . Dimensions of  $U_i^Q$ ,  $U_i^K$ ,  $U_i^V$ , and  $U^O$  are  $d_m \times d_k$ ,  $d_m \times d_k$ ,  $d_m \times d_v$ , and  $hd_v \times d_m$  so that dimension of multi-head attention  $\text{MultiheadAttention}(X, Y)$  is  $n \times d_m$ . Residual mechanism and layer normalization are applied into decoder multi-head attention too:

$$\text{MultiheadAttention}(X, Y)$$

$$= \text{LayerNorm}(\text{MaskedMultiheadAttention}(Y) + \text{MultiheadAttention}(X, Y))$$

Let  $Z$  be output of decoder which is decoder attention too, we obtain:

$$\begin{aligned}\mathbf{Z} &= \text{DecoderAttention}(\mathbf{X}, \mathbf{Y}) \\ &= \text{LayerNorm} \left( \text{MultiheadAttention}(\mathbf{X}, \mathbf{Y}) + \text{FFN}(\text{MultiheadAttention}(\mathbf{X}, \mathbf{Y})) \right)\end{aligned}$$

Where FFN denotes feedforward network or feedforward sublayer. If there is a stack of  $N$  decoders, the process above is repeated  $N$  times. In literature (Vaswani, et al., 2017),  $N$  is set to be 6. Without loss of generality, we can consider  $N = 1$  as simplest case for easy explanations. Note, dimension of  $\mathbf{Z}$  is  $n \times d_m$ . Model dimension  $d_m$  is often set to be 512 in NLP.

In context of statistical translation machine (STM), it is necessary to calculate probabilities of words (tokens) in vocabulary  $\Omega$ . Because these probabilities are calculated based on soft-max function, it is first to map decoder output matrix  $\mathbf{Z}$  into weight vector  $\mathbf{w} = (w_1, w_2, \dots, w_{|\Omega|})^T$  where every element  $w_i$  of vector  $\mathbf{w}$  is weight of the  $i^{\text{th}}$  word in vocabulary  $\Omega$ . The mapping is implemented by a feedforward network (FNN) called linear component in literature (Vaswani, et al., 2017, p. 3). In other words, input of linear component is sequence matrix  $\mathbf{Z}$  whereas its output is weight vector  $\mathbf{w}$  (Alammar, 2018). Please pay attention that the length of  $\mathbf{w}$  is the number of words (tokens) in vocabulary  $\Omega$  and so,  $\mathbf{w}$  is also called token/word weight vector.

$$\mathbf{w} = (w_1, w_2, \dots, w_{|\Omega|})^T = \text{linear}(\mathbf{Z}) = \text{FFN}(\mathbf{Z})$$

In practice,  $\mathbf{Z}$  is flattened into long vector because  $\mathbf{w}$  is vector too so that FNN can be implemented. After token weight vector  $\mathbf{w}$  is determined, it is easily converted into output probability vector  $\mathbf{p} = (p_1, p_2, \dots, p_{|\Omega|})^T$  where each element  $p_i$  is probability of the  $i^{\text{th}}$  word/token in vocabulary when sentence/sequence  $\mathbf{Z}$  is raised (Alammar, 2018). If the  $i^{\text{th}}$  word is issued, its probability  $p_i$  is 1 and other probabilities are 0.

$$\begin{aligned}\mathbf{p} &= \text{softmax}(\mathbf{w}) = (p_1, p_2, \dots, p_{|\Omega|})^T \\ p_i &= \frac{\exp(w_i)}{\sum_{l=1}^{|\Omega|} \exp(w_l)}\end{aligned}$$

Consequently, the next token which is predicted in STM for example is the one whose probability is highest, which means that the largest element in  $\mathbf{p}$  need to be found for STM translation after linear component  $\mathbf{w}$  and output probability  $\mathbf{p}$  are evaluated given  $\mathbf{Z}$  which in turn determined based on source sequence  $\mathbf{X}$  and target sequence  $\mathbf{Y}$  via mechanism encoder/decoder and attention.

It is not difficult to learn linear component  $\text{FFN}(\mathbf{Z})$  by backpropagation algorithm associated stochastic gradient descent (SGD) method. Concretely, the following cross-entropy  $L(\mathbf{p} | \Theta)$  is minimized so as to train  $\text{FFN}(\mathbf{Z})$ .

$$L(\mathbf{p} | \Theta) = - \sum_{i=1}^{|\Omega|} q_i \log(p_i)$$

Where  $\Theta$  is parameter of  $\text{FFN}(\mathbf{Z})$  and the vector  $\mathbf{q} = (q_1, q_2, \dots, q_{|\Omega|})^T$  is binary vector from sample whose each element  $q_i$  has binary values  $\{0, 1\}$  indicating whether the  $i^{\text{th}}$  token/word exists. For example, give sequence/sentence ("I", "am", "a", "student")<sup>T</sup>, if there is only one token/word "I" in sample sentence, the binary vector will be  $\mathbf{q} = (1, 0, 0, 0)^T$ . If three words "I", "am", and "student" are mutually existent, the binary vector will be  $\mathbf{q} = (1, 1, 0, 1)^T$ . When SGD is applied into minimizing the cross-entropy, partial gradient of  $L(\mathbf{p} | \Theta)$  with regard to  $w_j$  is:

$$\frac{\partial L(\mathbf{p} | \Theta)}{\partial w_j} = - \sum_{i=1}^{|\Omega|} q_i (\delta_{ij} - p_j)$$

Where,

$$\delta_{ij} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}$$

Proof,

$$\frac{\partial L(\mathbf{p}|\Theta)}{\partial w_j} = - \sum_{i=1}^{|\Omega|} q_i \frac{\partial \log(p_i)}{\partial w_j} = - \sum_{i=1}^{|\Omega|} q_i \frac{d \log(p_i)}{dp_i} \frac{\partial p_i}{\partial w_j} = - \sum_{i=1}^{|\Omega|} q_i \frac{1}{p_i} \frac{\partial p_i}{\partial w_j}$$

Due to:

$$\frac{\partial p_i}{\partial w_j} = \frac{\partial}{\partial w_j} \left( \frac{\exp(w_i)}{\sum_{l=1}^{|\Omega|} \exp(w_l)} \right) = p_i (\delta_{ij} - p_j)$$

We obtain:

$$\frac{\partial L(\mathbf{p}|\Theta)}{\partial w_j} = - \sum_{i=1}^{|\Omega|} q_i (\delta_{ij} - p_j) \blacksquare$$

So that gradient of  $L(\mathbf{p} | \Theta)$  with regard to  $\mathbf{w}$  is:

$$\frac{dL(\mathbf{p}|\Theta)}{d\mathbf{w}} = \left( \frac{\partial L(\mathbf{p}|\Theta)}{\partial w_1}, \frac{\partial L(\mathbf{p}|\Theta)}{\partial w_2}, \dots, \frac{\partial L(\mathbf{p}|\Theta)}{\partial w_{|\Omega|}} \right)^T$$

Therefore, parameter  $\Theta$  is updated according to SGD associated with backpropagation algorithm:

$$\Theta = \Theta - \gamma \frac{dL(\mathbf{p}|\Theta)}{d\mathbf{w}}$$

Where  $\gamma$  ( $0 < \gamma \leq 1$ ) is learning rate.

For STM example, given French source sentence “Je suis étudiant” (Alammar, 2018) is translated into English target sentence “I am a student” (Alammar, 2018) by transformer which is trained with corpus before (transformer was determined), which goes through following rounds:

Round 1:

- French source sentence “Je suis étudiant” coded by sentence/sequence matrix  $\mathbf{X} = (\mathbf{x}_1 = c("<bos>"), \mathbf{x}_2 = c("je"), \mathbf{x}_3 = c("suis"), \mathbf{x}_4 = c("étudiant"), \mathbf{x}_5 = c("<eos>"))^T$  where  $c(.)$  is embedding numeric vector of given word with note that words “<bos>” and “<eos>” are special predefined words indicating the beginning of sentence and the end of sentence, respectively. As a convention,  $c(.)$  is called word/token vector whose dimension can be  $d_m=512$ . If predefined sentence length is longer, redundant word vectors are set to be zeros, for example, let  $\mathbf{x}_6 = \mathbf{0}, \mathbf{x}_7 = \mathbf{0}, \dots, \mathbf{x}_{100} = \mathbf{0}$  given the maximum number words in sentence is 100. These zero vectors do not affect decoder evaluation and training parameters.
- Source sequence  $\mathbf{X}$  is fed to encoder so as to produce encoder attention  $\text{EncoderAttention}(\mathbf{X})$ .

Round 2:

- English target sentence is coded by sequence/matrix  $\mathbf{Y} = (\mathbf{y}_1 = c("<bos>"))^T$ . If predefined sentence length is longer, redundant word vectors are set to be zeros.
- Target sequence  $\mathbf{Y} = (\mathbf{y}_1 = c("<bos>"))^T$  and encoder attention  $\text{EncoderAttention}(\mathbf{X})$  are fed to decoder so as to produce decode output  $\mathbf{Z}$ .
- Output  $\mathbf{Z}$  goes through linear component  $\mathbf{w} = \text{linear}(\mathbf{Z})$  and soft-max function component  $\mathbf{p} = \text{softmax}(\mathbf{w})$  so as to find out the maximum probability  $p_i$  so that the  $i^{\text{th}}$  associated word in vocabulary is “i”. As a result, the embedding numeric vector of the word “i” is added to target sequence so that we obtain  $\mathbf{Y} = (\mathbf{y}_1 = c("<bos>"), \mathbf{y}_2 = c("i"))^T$ .

Round 3:

- Both target sequence  $\mathbf{Y} = (\mathbf{y}_1 = c("<bos>"), \mathbf{y}_2 = c("i"))^T$  and encoder attention  $\text{EncoderAttention}(\mathbf{X})$  are fed to decoder so as to produce decode output  $\mathbf{Z}$ .
- Output  $\mathbf{Z}$  goes through linear component  $\mathbf{w} = \text{linear}(\mathbf{Z})$  and soft-max function component  $\mathbf{p} = \text{softmax}(\mathbf{w})$  so as to find out the maximum probability  $p_i$  so that the  $i^{\text{th}}$  associated word in vocabulary is "am". As a result, the embedding numeric vector of the word "am" is added to target sequence so that we obtain  $\mathbf{Y} = (\mathbf{y}_1 = c("<bos>"), \mathbf{y}_2 = c("i"), \mathbf{y}_3 = c("am"))^T$ .

Similarly, rounds 4, 5, and 6 are processed by the same way so as to obtain final target sequence  $\mathbf{Y} = (\mathbf{y}_1 = c("<bos>"), \mathbf{y}_2 = c("i"), \mathbf{y}_3 = c("am"), \mathbf{y}_4 = c("a"), \mathbf{y}_5 = c("student"), \mathbf{y}_6 = c("<eos>"))^T$  which is the English sentence "I am a student" translated from the French sentence "Je suis étudiant". Note, the translation process is stopped when the end-of-sentence word "<eos>" is met.

Main ideas of transformer were described but there are two improvements such as positional encoding and normalization. Firstly, *positional encoding* is that sequences  $\mathbf{X}$  and  $\mathbf{Y}$  were added by their corresponding position vectors:

$$\begin{aligned}\mathbf{X} &= \mathbf{X} + (\text{pos}(\mathbf{x}_1), \text{pos}(\mathbf{x}_2), \dots, \text{pos}(\mathbf{x}_m))^T \\ \mathbf{Y} &= \mathbf{Y} + (\text{pos}(\mathbf{y}_1), \text{pos}(\mathbf{y}_2), \dots, \text{pos}(\mathbf{y}_n))^T\end{aligned}$$

Without loss of generality, let  $\text{POS}(\mathbf{X}) = (\text{pos}(\mathbf{x}_1), \text{pos}(\mathbf{x}_2), \dots, \text{pos}(\mathbf{x}_m))^T$  be position vector whose each element is position  $\text{pos}(\mathbf{x}_i)$  of token  $\mathbf{x}_i$ . It is necessary to survey  $\text{pos}(\mathbf{x}_i)$ .

$$\begin{aligned}\text{pos}(\mathbf{x}_i) &= \text{pos}((x_{i1}, x_{i2}, \dots, x_{ij}, \dots, x_{id_m})^T) \\ &= (\text{pos}(x_{i1}), \text{pos}(x_{i2}), \dots, \text{pos}(x_{ij}), \dots, \text{pos}(x_{id_m}))^T\end{aligned}$$

This implies how to calculate position vector  $\text{POS}(\mathbf{X})$  is how to calculate position  $\text{pos}(x_{ij})$  where  $i$  is position of the  $i^{\text{th}}$  token and  $j$  is position of the  $j^{\text{th}}$  numeric value of such token vector. We have:

$$\text{pos}(x_{ij}) = \text{pos}(i, j)$$

Suppose two successive numeric values such as  $j^{\text{th}}$  numeric value and  $(j+1)^{\text{th}}$  numeric value such that  $j = 2k$  and  $j+1 = 2k+1$ , we need to calculate two kinds of positions as follows:

$$\begin{aligned}\text{pos}(i, 2k) \\ \text{pos}(i, 2k + 1)\end{aligned}$$

Fortunately, these positions are easily calculated by sine function and cosine function as follows (Vaswani, et al., 2017, p. 6):

$$\begin{aligned}\text{pos}(i, 2k) &= \sin\left(\frac{i}{10000^{2k/d_m}}\right) \\ \text{pos}(i, 2k + 1) &= \cos\left(\frac{i}{10000^{2k/d_m}}\right)\end{aligned}$$

Recall that  $d_m$  is model dimension which is the length of token vector  $\mathbf{x}_i$ . It is often set to be 512 in NLP. As a result, we have:

$$\begin{aligned}\mathbf{X} &= \mathbf{X} + \text{POS}(\mathbf{X}) \\ \text{POS}(\mathbf{X}) &= (\text{pos}(\mathbf{x}_1), \text{pos}(\mathbf{x}_2), \dots, \text{pos}(\mathbf{x}_m))^T \\ \text{pos}(\mathbf{x}_i) &= (\text{pos}(i, 1), \text{pos}(i, 2), \dots, \text{pos}(i, d_m))^T \\ \text{pos}(i, j) &= \begin{cases} \sin\left(\frac{i}{10000^{2k/d_m}}\right) & \text{if } j = 2k \\ \cos\left(\frac{i}{10000^{2k/d_m}}\right) & \text{if } j = 2k + 1 \end{cases}\end{aligned}$$

Please pay attention that target sequence  $\mathbf{Y}$  is added by position vector  $\text{POS}(\mathbf{Y})$  by the same way too. There may be a question that why sequences  $\mathbf{X}$  and  $\mathbf{Y}$  are added by their position vectors before

they are fed into encoder/decoder when tokens in a sequence have their own orders because a sequence is an ordered list of tokens indeed. The answer depends on computational effectiveness as well as flexibility. For example, when sequences are added by their position vectors, transformer can be trained by incomplete French source sequence “<bos> Je suis” and incomplete English target sequence “a student <eos>” because there is no requirement of token ordering. Moreover, sequences can be split into many parts and these parts are trained parallel. This improvement is necessary in case of training huge corpus.

The second improvement is layer (network) normalization:

$$\begin{aligned} & \text{LayerNorm}(\mathbf{X} + \text{Sublayer}(\mathbf{X})) \\ & \text{LayerNorm}(\mathbf{Y} + \text{Sublayer}(\mathbf{Y})) \end{aligned}$$

Because residual mechanism is implemented by the sum  $\mathbf{X} + \text{Sublayer}(\mathbf{X})$  or  $\mathbf{Y} + \text{Sublayer}(\mathbf{Y})$ , it is necessary to survey the following normalization without loss of generality:

$$\text{LayerNorm}(\mathbf{x})$$

Where  $\mathbf{x} = (x_1, x_2, \dots, x_n)^T$  is layer of  $n$  neuron  $x_i$  with note that each neuron  $x_i$  is represented by a number. Suppose  $\mathbf{x}$  as a sample conforms normal distribution, its sample mean and variance are calculated as follows:

$$\begin{aligned} \bar{x} &= \frac{1}{n} \sum_{i=1}^n x_i \\ s^2 &= \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \end{aligned}$$

As a result, layer normalization is distribution normalization:

$$\begin{aligned} \text{LayerNorm}(\mathbf{x}) &= (\text{norm}(x_1), \text{norm}(x_2), \dots, \text{norm}(x_n))^T \\ \text{norm}(x_i) &= \frac{x_i - \bar{x}}{\sqrt{s^2}} \end{aligned}$$

In literature, layer normalization aims to improve convergence speed in training.

It is not difficult to train transformer from corpus which can be a huge set of pairs of source/target sequences. Backpropagation algorithm associated with stochastic gradient descent (SGD) is a simple and effective choice. Feedforward sublayer represented by feedforward network (FFN) is easily trained by backpropagation algorithm associated SGD, besides attention sublayers can be trained by backpropagation algorithm associated SGD too. For instance, attention parameters for encoder such as weight matrices  $W_i^Q$ ,  $W_i^K$ ,  $W_i^V$ , and  $W^O$  can be learned by backpropagation algorithm associated with SGD. Attention parameters for decoder such as weight matrices  $W_i^Q$ ,  $W_i^K$ ,  $W_i^V$ ,  $W^O$ ,  $T$ ,  $U_i^Q$ ,  $U_i^K$ ,  $U_i^V$ , and  $U^O$  can be learned by backpropagation algorithm associated SGD too. Note, starting point for backpropagation algorithm to train transformer is to make comparison of target sequence (for example, the English target sentence “I am a student” given the French source sentence “Je suis étudiant”) and evaluated sequence (for example, the English evaluated sentence “We are scholars” given the same French source sentence “Je suis étudiant”) at decoder, which goes backward encoder. Moreover, please pay attention that zero vectors representing redundant tokens do not affect updating these weight matrices when training transformer.

#### 4. Pre-trained model

AI models cope with two problems of model learning: 1) it is impossible to preprocess or annotate (label) huge data so as to make the huge data better for training, and 2) huge data is often come

with data stream rather than data scratch. Note, the first problem is most important. *Transfer learning* (Han, et al., 2021, pp. 226-227) can solve the two problems by separating the training process by two stages: 1) *pre-training stage* aims to draw valuable knowledge from data stream / data scratch, and 2) *fine-tuning stage* later will take advantages of knowledge from pre-training stage so as to apply the knowledge into solving task-specific problem just by fewer samples or smaller data. As its name hints, transfer learning draws knowledge from pre-training stage and then transfers such knowledge to fine-tuning stage for doing some specific task. Capturing knowledge in pre-training stage is known as source task and doing some specific task is known as target task (Han, et al., 2021, p. 227). Source task and target task may be essentially similar like GPT model and BERT model for token generation mentioned later but these tasks can be different or slightly different. The fine-tuning stage is dependent on concrete application and so, pre-training stage is focused in this section. The purpose of pre-training stage is to build a *large-scale pre-trained model* called *PTM* which must have ability to process huge data or large-scale data. If large-scale data is come from data stream called downstream data, PTM will need to reach the strong point that is parallel computation. If large-scale is too huge, PTM will need to reach the strong point that is efficient computation. When efficient computation can be reached by good implementation, parallel computation requires an improvement of methodology. In order to catch knowledge inside data without human interference with restriction that such knowledge represented by label, annotation, context, meaning, etc. is better than cluster and group, self-supervised learning is often accepted as a good methodology for PTM (Han, et al., 2021, pp. 227-229). Essentially, *self-supervised learning* tries to draw pseudo-supervised information from unannotated/unlabeled data so that such pseudo-supervised information plays the role of supervised information like annotation and label that fine-tuning stage applies into supervised learning tasks for solving specific problem with limited data. The pseudo-supervised information is often relationships and contexts inside data structure. Anyhow, self-supervised learning is often associated with transfer learning because, simply, annotating entirely huge data is impossible. Self-supervised learning associated with pre-training stage is called self-supervised pre-training. Although self-supervised pre-training is preeminent, pre-training stage can apply other learning approaches such as supervised learning and unsupervised learning.

That the essentially strong point of transformer is self-attention makes transformer appropriate to be a good PTM when self-attention follows essentially ideology of self-supervised learning because self-attention mechanism tries to catch contextual meaning of every token inside its sequence. Moreover, transformer supports parallel computation based on its other aspect that transformer does not concern token ordering in sequence. Anyhow, transformer is suitable to PTM for transfer learning and so this section tries to explain large-scaled pre-trained model (PTM) via transformer as an example of PTM. Note, fine-tuning stage of transfer learning will take advantages of PTM for solving task-specific problem; in other words, fine-tuning stage will fine-tune or retrain PTM with downstream data, smaller data, or a smaller group of indications. When fine-tuning stage is not focused in description, PTM is known as transfer learning model which includes two stages such as pre-training stage and fine-tuning stage. In this case, source task and target task of transfer learning have the same model architecture (model backbone) which is the same PTM architecture. Large-scale PTM implies its huge number of parameters as well as huge data from which it is trained.

Generative Pre-trained Transformer (*GPT*), developed in 2018 with GPT-1 by OpenAI ([www.openai.com](http://www.openai.com)) whose product is ChatGPT launched in 2022, is a PTM that applies only decoder of transformer into sequence generation. In pre-training stage, GPT trains its decoder from



huge data over internet and available sources so as to predict next word  $y_{t+1}$  from previous words  $y_1, y_2, \dots, y_t$  by maximizing likelihood  $P(y_{t+1} | \Theta, y_1, y_2, \dots, y_t)$  and taking advantages of self-attention mechanism aforementioned (Han, et al., 2021, p. 231). Maximization of likelihood  $P(y_{t+1} | \Theta, y_1, y_2, \dots, y_t)$  belongs to autoregressive language model.

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmax}} \log(P(y_{t+1} | \Theta, y_1, y_2, \dots, y_t))$$

$$\hat{y}_{t+1} = \underset{y_{n+1}}{\operatorname{argmax}} P(y_{t+1} | \hat{\Theta}, y_1, y_2, \dots, y_t)$$

Where,

$$P(y_{t+1} | \Theta, y_1, y_2, \dots, y_t) \stackrel{\text{def}}{=} \mathbf{p}$$

$$\mathbf{p} = \operatorname{softmax}(\mathbf{w}) = (p_1, p_2, \dots, p_{|\Omega|})^T$$

$$p_i = \frac{\exp(w_i)}{\sum_{l=1}^{|\Omega|} \exp(w_l)}$$

$$\mathbf{w} = (w_1, w_2, \dots, w_{|\Omega|})^T = \operatorname{linear}(\mathbf{Z}) = \operatorname{FFN}(\mathbf{Z})$$

And,

$$\operatorname{MultiheadAttention}(\mathbf{X}, \mathbf{Y})$$

$$= \operatorname{LayerNorm}(\operatorname{MaskedMultiheadAttention}(\mathbf{Y}) + \operatorname{MultiheadAttention}(\mathbf{X}, \mathbf{Y}))$$

$$\mathbf{Z} = \operatorname{DecoderAttention}(\mathbf{X}, \mathbf{Y})$$

$$= \operatorname{LayerNorm}(\operatorname{MultiheadAttention}(\mathbf{X}, \mathbf{Y}) + \operatorname{FFN}(\operatorname{MultiheadAttention}(\mathbf{X}, \mathbf{Y})))$$

Because GPT has only one decoder, sequence  $\mathbf{X}$  is null in GPT.

$$\operatorname{MultiheadAttention}(\mathbf{X}, \mathbf{Y}) = \operatorname{MultiheadAttention}(\mathbf{Y})$$

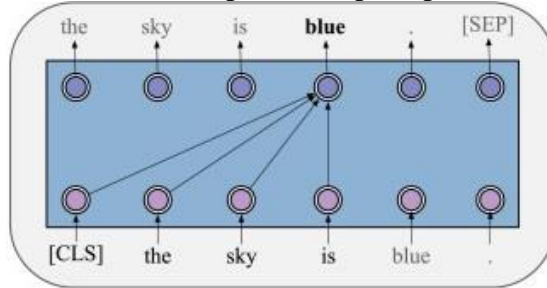
$$\operatorname{DecoderAttention}(\mathbf{X}, \mathbf{Y}) = \operatorname{DecoderAttention}(\mathbf{Y})$$

Likelihood  $P(y_{t+1} | \Theta, y_1, y_2, \dots, y_t)$  is simplified for easy explanation. Exactly, given sequence  $\mathbf{Y} = (y_1, y_2, \dots, y_{n+1})^T$ , GPT aims to maximize log-likelihood  $L(\Theta | \mathbf{Y})$  as follows (Han, et al., 2021, p. 231):

$$L(\Theta | \mathbf{Y}) = \sum_{i=1}^n \log(P(y_{i+1} | \Theta, y_1, y_2, \dots, y_i))$$

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmax}} L(\Theta | \mathbf{Y})$$

Later on, GPT improves its pre-trained decoder in fine-tuning stage by re-training the decoder with annotated data, high-quality data, and domain-specific data so as to improve pre-trained parameters. Moreover, GPT adds extra presentation layers in fine-tuning stage (Han, et al., 2021, p. 231). Following figure (Han, et al., 2021, p. 232) depicts prediction process of GPT.



**Figure 4.1.** Prediction process of GPT

Bidirectional Encoder Representations from Transformers (*BERT*), developed in 2018 by Google, is a PTM that applies only encoder of transformer into sequence generation. In pre-training stage,

BERT trains its encoder from huge data over internet and available sources. Given  $(t+1)$ -length sequence  $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{t+1})^T$ , BERT applies masked language model to randomize an unknown token at random position denoted  $masked$  where the random index  $masked$  is randomized in  $t+1$  indices  $\{1, 2, \dots, t+1\}$  with note that the randomization process can be repeated many times. Such unknown token, which is called masked token denoted  $\mathbf{y}_{masked}$ , will be predicted given  $t$ -length sequence  $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t)^T$  without loss of generality. In order words, masked words  $\mathbf{x}_{masked}$  is predicted from other words  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t$  by maximizing likelihood  $P(\mathbf{x}_{masked} | \Theta, \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t)$  and taking advantages of self-attention mechanism aforementioned (Han, et al., 2021, p. 232).

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmax}} \log(P(\mathbf{x}_{masked} | \Theta, \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t))$$

$$\hat{\mathbf{y}}_{masked} = \underset{\mathbf{y}_{masked}}{\operatorname{argmax}} P(\mathbf{x}_{masked} | \hat{\Theta}, \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t)$$

Where,

$$P(\mathbf{x}_{masked} | \Theta, \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t) \stackrel{\text{def}}{=} \mathbf{p}$$

$$\mathbf{p} = \operatorname{softmax}(\mathbf{w}) = (p_1, p_2, \dots, p_{|\Omega|})^T$$

$$p_i = \frac{\exp(w_i)}{\sum_{l=1}^{|\Omega|} \exp(w_l)}$$

$$\mathbf{w} = (w_1, w_2, \dots, w_{|\Omega|})^T = \operatorname{linear}(\mathbf{Z}) = \operatorname{FFN}(\mathbf{Z})$$

And,

$$\operatorname{MultiheadAttention}(\mathbf{X}) = \operatorname{LayerNorm}(\mathbf{X} + \operatorname{MultiheadAttention}(\mathbf{X}))$$

$$\mathbf{Z} = \operatorname{EncoderAttention}(\mathbf{X})$$

$$= \operatorname{LayerNorm}(\operatorname{MultiheadAttention}(\mathbf{X}) + \operatorname{FFN}(\operatorname{MultiheadAttention}(\mathbf{X})))$$

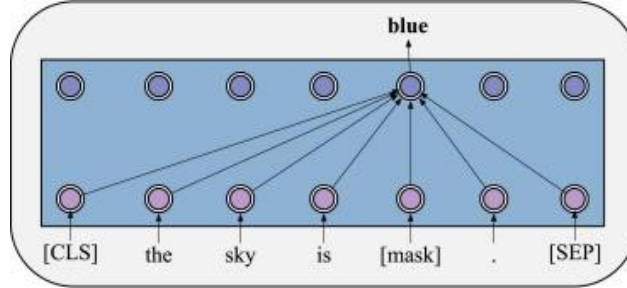
Likelihood  $P(\mathbf{y}_{masked} | \Theta, \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m)$  is simplified for easy explanation, thereby it is necessary to explain more how BERT defines and maximizes likelihood with support of masked language model. Given sequence  $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m)^T$ , let  $R = \{r_1, r_2, \dots, r_k\}$  be the set of indices whose respective tokens are initially masked, for instance, token  $\mathbf{x}_{r_j}$  will be initially masked if  $r_j$  belongs to mask set  $R$ . Let  $\mathbf{X}_{r_j}$  be the set of  $r_j-1$  tokens which are unmasked later, for instance, the tokens  $\mathbf{x}_{r_1}, \mathbf{x}_{r_2}, \dots, \mathbf{x}_{r_{j-1}}$  which were initially masked before are now unmasked (known) at current iteration. Note, the set  $R$  is called mask set or mask pattern and  $\mathbf{X}_{r_j}$  does not include token  $\mathbf{x}_{r_j}$ . BERT randomizes  $k$  masked indices so as to establish mask set  $R$ . Let  $S$  be the set of indices whose tokens are always known, which is the complement of mask set  $R$  with regard to all indices so that union of  $R$  and  $S$  is  $\{1, 2, \dots, m\}$ . Thereby, let  $\mathbf{S}$  be the set of tokens whose indices are in  $S$ . In other words,  $\mathbf{S}$  contains tokens which are always known. BERT aims to maximize log-likelihood  $L(\Theta | \mathbf{X})$  as follows (Han, et al., 2021, p. 232):

$$L(\Theta | \mathbf{X}) = \sum_{r_j \in R} \log(P(\mathbf{x}_{r_j} | \Theta, \mathbf{X}_{r_j}, \mathbf{S}))$$

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmax}} L(\Theta | \mathbf{X})$$

Later on, BERT improves its pre-trained encoder in fine-tuning stage by re-training the encoder with annotated data, high-quality data, and domain-specific data so as to improve pre-trained parameters. By support of masked language model (autoencoding language model) for masking tokens, BERT can predict a token at any position in two directions given a list of other tokens while GPT only predicts a token at next position given previous tokens. The name ‘‘BERT’’, which

is abbreviation of “Bidirectional Encoder Representations from Transformers”, hints that BERT can generate tokens/words in bidirectional way at any positions. Therefore, GPT is appropriate to language generation and BERT is appropriate to language understanding (Han, et al., 2021, p. 231). BERT also adds extra presentation layers in fine-tuning stage (Han, et al., 2021, p. 232). Following figure depicts prediction process of BERT.



**Figure 4.2.** Prediction process of BERT

Recall that given a transfer model, capturing knowledge in pre-training stage is known as source task and doing some specific task is known as target task (Han, et al., 2021, p. 227), thereby there is a question that how source task transfers knowledge to target task or how PTM makes connection between source task and knowledge task. The answer is that there are two transferring approaches such as feature transferring and parameter transferring (Han, et al., 2021, p. 227). *Feature transferring* converts coarse data like unlabeled data into fine data like labeled data so that fine data considered as feature is fed to fine-tuning stage. *Parameter transferring* transfers parameters learned at pre-training stage to fine-tuning stage. If pre-training stage and fine-tuning stage share the same model architecture which is the same PTM architecture, parameter transferring will always occur in PTM. Both GPT and BERT apply parameter transferring because they will initialize or set up their models such as GPT decoder and BERT encoder by billions of parameters that were learned in pre-training stage with the same model architecture (model backbone) such as GPT decoder and BERT encoder before they perform fine-tuning task in fine-tuning stage. Self-supervised learning which trains unlabeled data is appropriate to pre-training stage because unlabeled data is much more popular than labeled data, thereby parameter transferring is often associated with self-supervised learning. Because transformer is suitable to self-supervised learning due to its self-attention mechanism, parameter transferring is suitable to PTMs like GPT and BERT. Moreover, if they apply transformer into annotating or creating task-specific data / fine data for improving their decoder and encoder in fine-tuning stage, they will apply feature transferring too. In general, within parameter transferring and same architecture, PTM itself is backbone for both pre-training stage and fine-tuning stage.

## 5. Conclusions

As the paper title “Attention is all you need” (Vaswani, et al., 2017) hints, attention-awarded transformer is the important framework for generative artificial intelligence and statistical translation machine whose applications are not only large but also highly potential. For instance, it is possible for transformer to generate media content like sound, image, video from texts, which is very potential for cartoon industry and movie making applications (film industry). The problem of difference in source data and target data, which can be that, for example, source sequence is text sentence and target sequence is raster data like sound and image, can be solved effectively and smoothly because of two aforementioned strong points of transformer such as self-attention and not concerning token ordering. Moreover, transformer’s methodology is succinct with support of

encoder-decoder mechanism and deep neural network. Therefore, it is possible to infer that applications of transformer can go beyond some recent pre-trained models and/or pre-trained models based on transformer can be improved more.

## References

- Alammar, J. (2018, June 27). *The Illustrated Transformer*. (GitHub) Retrieved June 2024, from Jay Alammar website: <https://jalammar.github.io/illustrated-transformer>
- Cho, K., Merrienboer, B. v., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., & Bengio, Y. (2014, September 3). Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation. *arXiv preprint*, 1-15. doi:10.48550/arXiv.1406.1078
- Graves, A. (2014, June 5). Generating Sequences With Recurrent Neural Networks. *arXiv preprint*, 1-43. doi:10.48550/arXiv.1308.0850
- Han, X., Zhang, Z., Ding, N., Gu, Y., Liu, X., Huo, Y., . . . Zhu, J. (2021, August 26). Pre-trained models: Past, present and future. *AI Open*, 2(2021), 225-250. doi:10.1016/j.aiopen.2021.08.002
- Hardle, W., & Simar, L. (2013). *Applied Multivariate Statistical Analysis*. Berlin, Germany: Research Data Center, School of Business and Economics, Humboldt University.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., . . . Polosukhin, I. (2017). Attention Is All You Need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, & S. Vishwanathan (Ed.), *Advances in Neural Information Processing Systems (NIPS 2017)*. 30. Long Beach: NeurIPS. Retrieved from <https://arxiv.org/abs/1706.03762>
- Voita, L. (2023, November 17). *Sequence to Sequence (seq2seq) and Attention*. (GitHub) Retrieved June 2024, from Elena (Lena) Voita website: [https://lena-voita.github.io/nlp\\_course/seq2seq\\_and\\_attention.html](https://lena-voita.github.io/nlp_course/seq2seq_and_attention.html)
- Wikipedia. (2005, April 7). *Recurrent neural network*. (Wikimedia Foundation) Retrieved from Wikipedia website: [https://en.wikipedia.org/wiki/Recurrent\\_neural\\_network](https://en.wikipedia.org/wiki/Recurrent_neural_network)
- Wikipedia. (2019, August 25). *Transformer (deep learning architecture)*. (Wikimedia Foundation) Retrieved from Wikipedia website: [https://en.wikipedia.org/wiki/Transformer\\_\(deep\\_learning\\_architecture\)](https://en.wikipedia.org/wiki/Transformer_(deep_learning_architecture))