# Attention Is All You Need.

Abstract

# Abstract

- Simple network architecture ... <u>Transformer</u>

- based solely on attention mechanisms!

- Experiments (machine translation):

    - superior in quality

    - more parallelizable

    - require <u>significantly</u> less time to train

- WMT 2014, best BLEU

- generalizes well!

    - *1 Constituency parsing

# 1. Introduction

## Recurrent models

- typically factor computation along the symbol position of the input and output sequences.
- Aligning the positions to steps in computation time, they generate a sequence of hidden states $h_t$, as
  - a function of previous hidden state $h_{t-1}$
  - input for position $t$.

$\downarrow$

preclude parallelization within training examples, critical at longer sequence lengths!

(improvements: factorization tricks & conditional computation)

- Attention mechanism
  : allow modeling of dependencies without regard to their distance in the input or output sentences.

## Transformer

- rely entirely on an attention mechanism to draw global dependencies between input & output
- more parallelization
- SOTA in translation quality
- 12 hrs training with 8 P100 GPUs.

## 2. Background

- to reduce sequential computation ; CNN
    - Extented Neural GPU
    - ByteNet
    - ConvS2S

↳ the number of operations required to relate signals from two arbitrary input or output positions grows in the distance between positions.
    - linearly for ConvS2S
    - logarithmically for ByteNet

☆ Transformer ☆
- reduced to a constant number of operations,
- reduced effective resolution due to averaging attention-weighted positions

    ↓

- counteract with Multi-Head Attention

- Self - Attention ( = intra - attention)
    : attention mechanism relating different positions of a single sequence in order to compute a representation of the sequence.

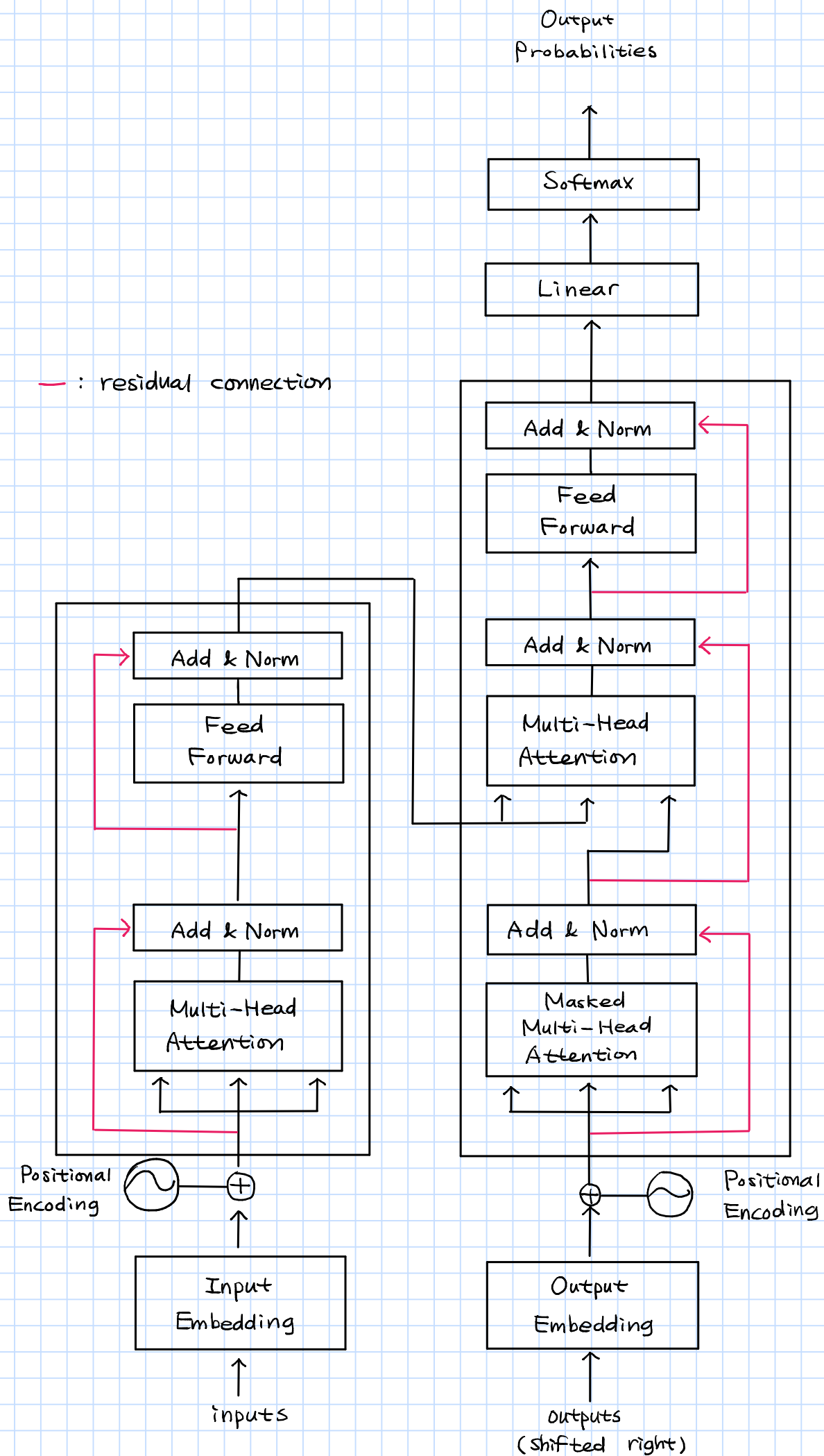- 1st transduction model relying entirely on self-attention to compute representations of its input & output without using sequence - aligned RNN or convolution.

# 3. Model Architecture

- Encoder - decoder structure
- Encoder: maps an input sequence of symbol representations $(x_1, ..., x_n)$

  to a sequence of continuous representations $z = (z_1, ..., z_n)$

  $\downarrow$

- Decoder: generates an output sequence $(y_1, ..., y_m)$ of symbols

  one element at a time.

- auto - regressive at each step.

  (consuming previously generated symbols as additional input

  when generating the next)


- Stacked self - attention
- point - wise, fully connected layers for both encoder & decoder

Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Add & Norm

Multi-Head
Attention

Add & Norm

Feed
Forward

Add & Norm

Masked
Multi-Head
Attention

— : residual connection

Positional
Encoding

Positional
Encoding

Input
Embedding

Output
Embedding

inputs

outputs
(shifted right)

## 3.1 Encoder & Decoder Stacks

- Encoder

  - stack of N = 6 identical layers

  - 1st: multi-head self-attention mechanism

  - 2nd: position-wise fully connected feed-forward network

  - residual connection around each of the two sub-layers, followed by layer normalization.

  - the output of each sub-layer is LayerNorm ($x$ + <u>Sublayer($x$)</u>)

    <span style="color:red">function implemented by the sub-layer itself.</span>

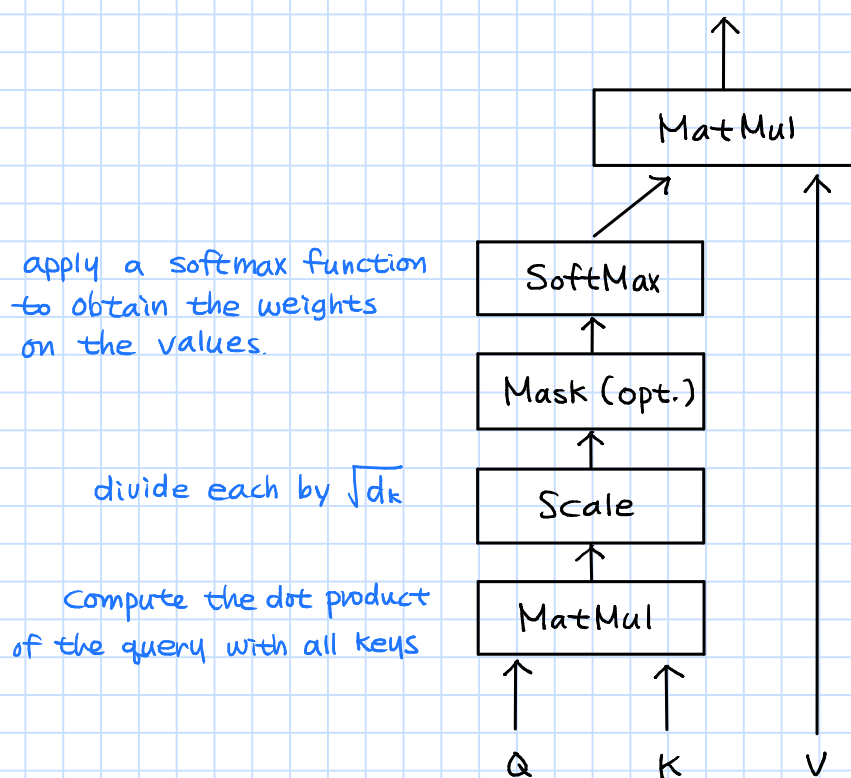  - $d_{model}$ = 512

- Decoder

  - stack of N = 6 identical layers

  - masking: prevent positions from attending to subsequent positions.

## 3.2 Attention

- Mapping a query & a set of key-value pairs to an output (Q, K, V, output are all vectors)
- output: weighted sum of the values
- the weight assigned to each value is computed by a compatibility function of the query with the corresponding key.

### 3.2.1 Scaled Dot-Product Attention

apply a softmax function to obtain the weights on the values.

divide each by $\sqrt{d_K}$

Compute the dot product of the query with all keys



- input: queries, keys of dimension $d_K$, values of dimension $d_V$
- In practice, compute attention function simultaneously, packed together into a matrix Q. K(keys), V(Values)

- $\text{Attention}(Q, K, V) = \text{softmax}\left(\dfrac{QK^T}{\sqrt{d_K}}\right) V$

2 most commonly used attention functions:

① additive attention

② dot-product (multiplicative) attention


① additive attention

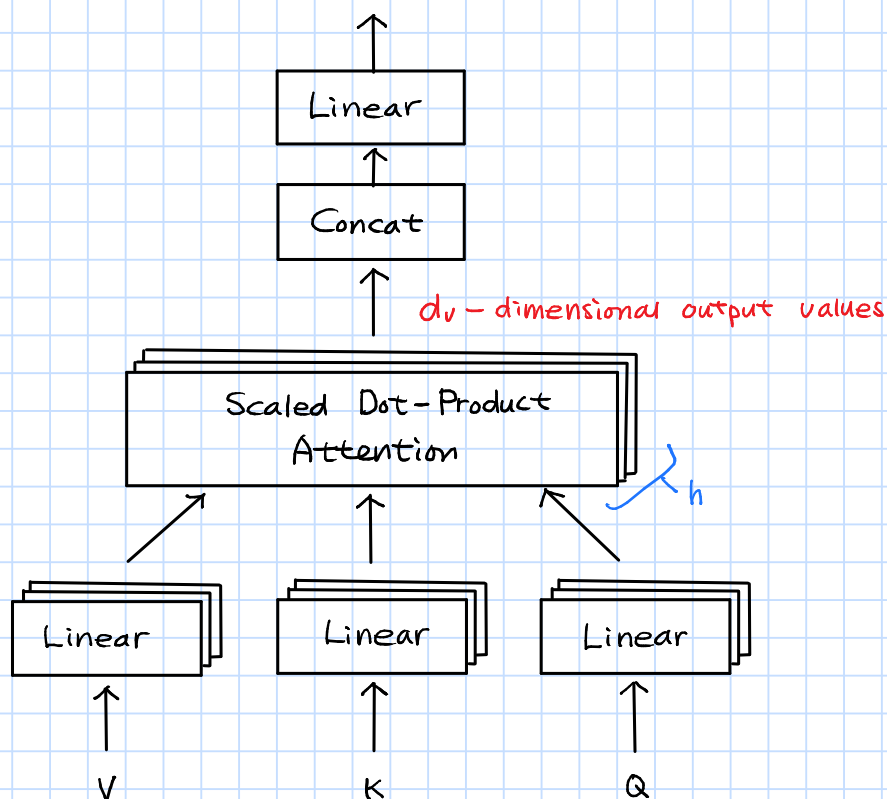- computes compatibility function using a feed-forward network with a single hidden layer.

② dot-product attention

- much faster, space-efficient in practice
- for large values of $d_k$, the dot products grow large in magnitude, pushing the softmax function into regions where it has extremely small gradients.

↓

scale the dot product by $\frac{1}{\sqrt{d_k}}$

## 3.2.2 Multi-Head Attention



$d_v$ – dimensional output values

- linearly project the queries, keys, and values  $h$  times with different, learned linear projections  to  $d_k$, $d_k$, $d_v$  dimensions, respectively.
- allows models to jointly attend to information from different representation subspace at different positions.

- $\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, ..., \text{head}_n) W^O$

   where $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^k, VW_i^v)$

   - $W_i^Q \in \mathbb{R}^{d_{model} \times d_k}$
   - $W_i^k \in \mathbb{R}^{d_{model} \times d_k}$
   - $W_i^v \in \mathbb{R}^{d_{model} \times d_v}$
   - $W^O \in \mathbb{R}^{h d_v \times d_{model}}$

- $h = 8$ parallel attention layers, or heads.
- $d_k = d_v = d_{model} / h = 64$
- Due to the reduced dimension of each head, the total computational cost is similar to that of single-head attention with full dimensionality.

3.2.3  Application of Attention in our Model

   — Transformer uses multi-head attention in 3 ways:

① In encoder-decoder attention layers

   — Queries: Come from the previous decoder layer.

   — Memory keys & values : Come from the output of the encoder.

   — allows every position in the decoder to attend over all positions
     in the input sequence.

② Self-attention layers in encoder

   — In a self-attention layer, all of the keys, values, and queries
     Come from the output of the previous layer in the encoder.

   — Each position in the encoder can attend to all positions
     in the previous layer of the encoder.

③ Self-attention layers in decoder

   — allow each position in the decoder to attend to all positions in
     the decoder up to & including that position

   — need to prevent leftward information flow in the decoder
     to preserve the auto-regressive property.

     ↓

   — inside of scaled dot-product attention by masking out
     all values in the input of the softmax which correspond to
     illegal connections.

## 3.3 Position-wise Feed-Forward Networks

- Each of the layers in the encoder & decoder contains a fully connected feed-forward network
- Consists of 2 linear transformations with a ReLU activation in between

- $FFN(x) = \max(0, xW_1 + b_1) W_2 + b_2$

- Linear transformations : same across different positions
  use different parameters from layer to layer.
- Dimensionality :
  - input & output : $d_{model} = 512$
  - inner-layer : $d_{ff} = 2048$

## 3.4 Embeddings & Softmax
- Embedding : to convert the input ~~tokens~~ & output ~~tokens~~
  to vectors of dimension $d_{model}$
- learned linear transformation & softmax function
  : to convert the decoder output to predict next-token probabilities.

- Share the same weight matrix between 2 embedding layers & the pre-softmax linear transformation
- In the embedding layers, those weights are multiplied by $\sqrt{d_{model}}$

## 3.5 Positional Encoding

$$PE_{(pos, 2i)} = \sin(pos / 10000^{2i/d_{model}})$$

$$PE_{(pos, 2i+1)} = \cos(pos / 10000^{2i/d_{model}})$$

position     dimension

- each dimension of the positional encoding corresponds to a sinusoid.

- The wavelengths form a geometric progress from $2\pi$ to $10000 \cdot 2\pi$

# *1. Constituency parsing

: 문장이 구 단위로 묶여가면서 구조를 이뤄가는 방법

어순이 고정적인 영어에서 쓰임.

: 문장이 구 단위로 묶여가면서 구조를 이뤄가는 방법

어순이 고정적인 영어에서 쓰임.