# Attention Is All You Need

## The Rise of the Transformer

Jong-Cheol Lee

Department of Statistics
Pusan National University

October 29, 2024

# Contents

# Introduction

## Machine Translation



| RNN (1986) | LSTM (1997) | Seq2Seq (2014) | Attention (2015) | Tranformer (2017) | GPT-1 (2018) | BERT (2019) | GPT-3 (2020) |

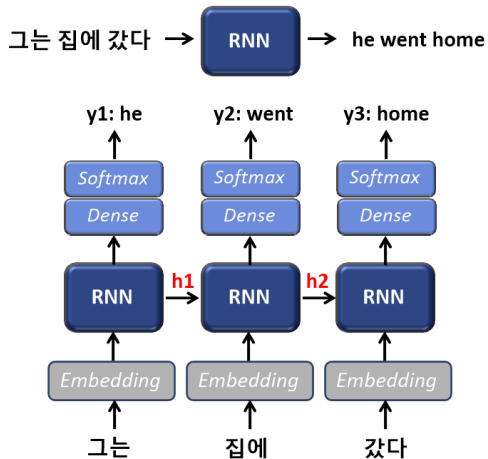*Traditional Sequential Processing* — *Attention-based Parallel Processing*

## Goal

**Key Objective**

*"Proposing the* **Transformer***, a completely new architecture."*

# RNN-based Translation

## Example



## Fundamental Concepts

- Word Embedding
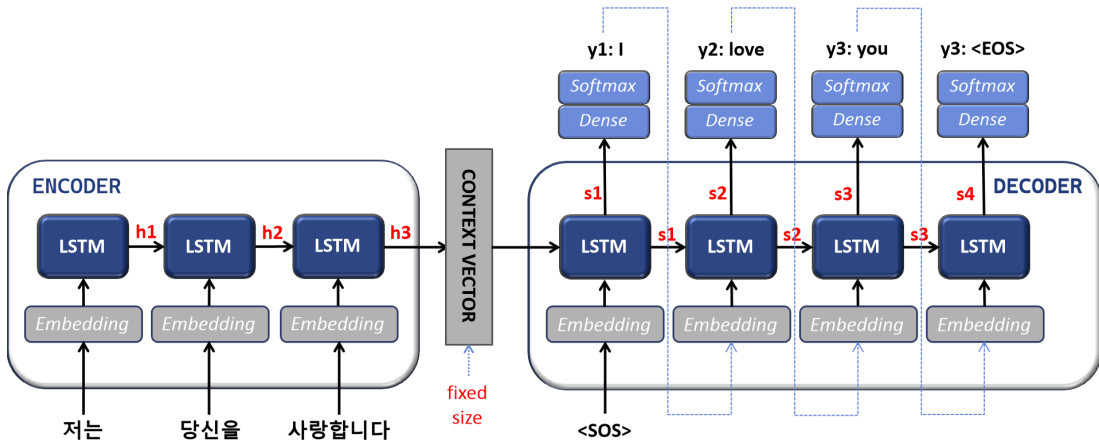
  그는 → | 0.7 | 0.3 | ... | 0.4 | 0.4 |

- Softmax

$$\begin{bmatrix} 2.0 \\ 1.0 \\ 0.1 \end{bmatrix} \rightarrow \frac{e^z}{e^2 + e^1 + e^{0.1}} \rightarrow \begin{bmatrix} 0.67 \\ 0.23 \\ 0.1 \end{bmatrix}$$
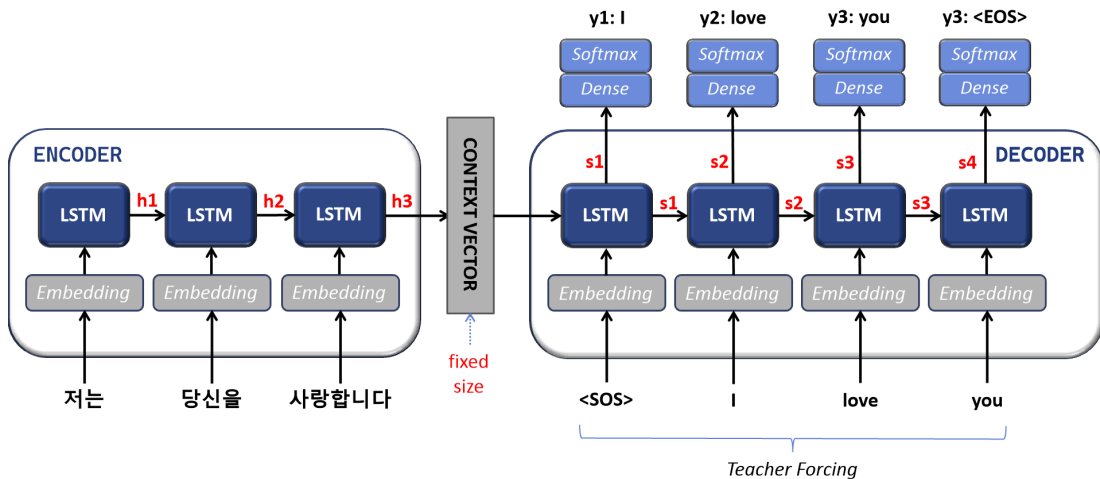
## Limitations

- Input-Output Length Constraint
- Word Order Variation Across Languages

# Seq2Seq

## Encoder-Decoder structure

# Encoder-Decoder

# Attention: Background & Main Concept

## Background

- Information loss of context vector
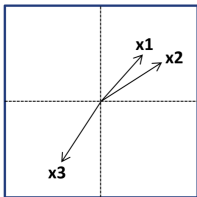- $s_1, ..., s_4$: Inefficient vector representation

## Assumption

- $h_i$ mainly captures i-th input word
- All $s_i, h_i$ are same-sized vectors

> **Main Concept**
>
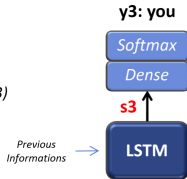> Pay **attention** to important information in each step of decoder

# Dot-Product Attention

**1) dot product ≃ similarity**



$dot(x1, x2) > dot(x1, x3)$

**2) similarity with s3 ≃ similarity with "you"**

y3: you

Softmax

Dense

s3

Previous Informations →

LSTM

**Before training**: s3 is far from **"you"**

↓

**After training**: s3 is close to **"you"**

**3)**

h1: Information of '저는'
h2: 〃 '당신을'
h3: 〃 '사랑합니다'

⟹

**dot(h1, s3)**: *similarity between* '저는' *and* s3
**dot(h2, s3)**: *similarity between* '당신을' *and* s3
**dot(h3, s3)**: *similarity between* '사랑합니다' *and* s3

**dot(h1, s3)·h1 + dot(h2, s3)·h2 + dot(h3, s3)·h3 = Attention value vector**

*weighted sum*

y3: you

Softmax

Dense

s3

# Dot-Product Attention

# Dot-Product Attention



Blend **informations** of input words

Attention Value
= <s3,h1>*h1 + <s3,h2>*h2 + <s3,h3>*h3

y3: you

# Main Idea: Query, Key, Value

## What is Q, K, V?



**Query Vector: s3**
**Key Vectors: h1 h2 h3**
**Value Vectors: h1 h2 h3**

**Query**: *Questions*
**Key**: *Tags of words*
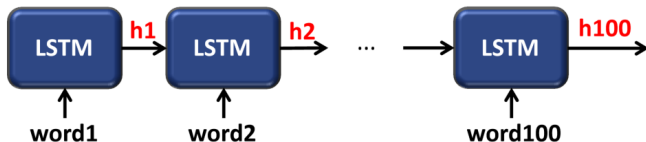**Value**: *Informations of words*

$$\therefore \text{ Attention}(Q, K, V) = \text{softmax}(QK')V$$

# Transformer: Background & Main Concept

## Problems of Sequential Processing

- Difficulty to parallelize
- Long-term dependency



### Main Concept

Using only **attention mechanisms** to build an encoder-decoder

# Transformer: Overview
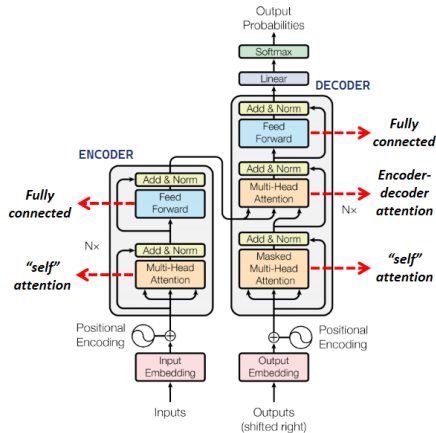
## Model Architecture



Figure 1: The Transformer - model architecture.
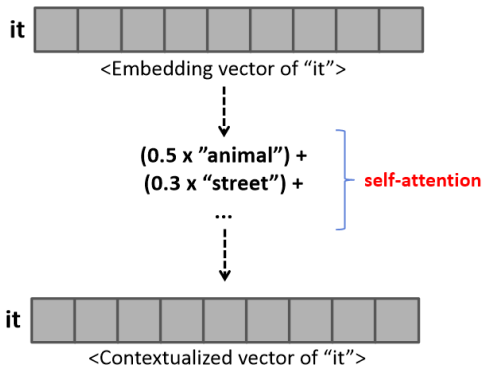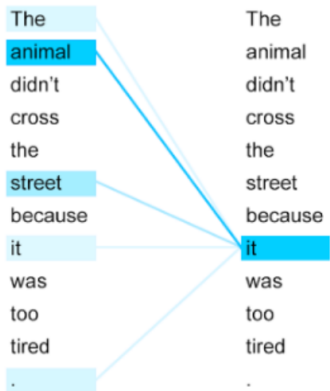
## Hyper Parameters

| Hyperparameter | Value |
|---|---|
| $d_{model}$ | 512 |
| $d_{ff}$ | 2048 |
| $N_{heads}$ | 8 |
| $N_{layers}$ | 6 |

- $d_{model}$: size of embedded vector
  (equals with size input and output)

- $d_{ff}$: hidden size of *Feed Forward*

- $N_{heads}$: # of attention heads
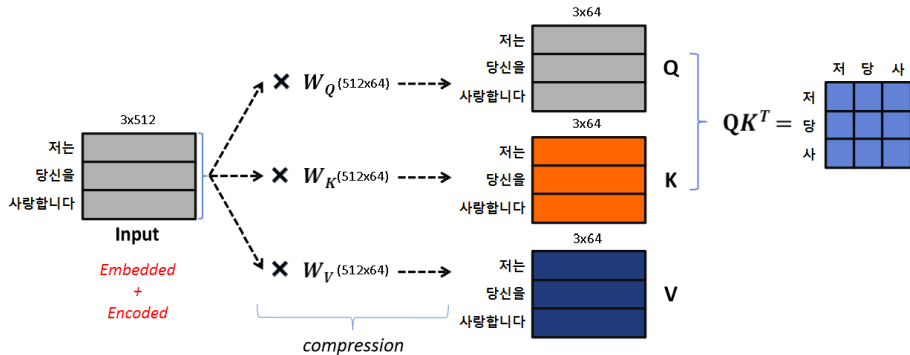
- $N_{layers}$: of layers stacked

# Why Self Attention?

## An Effect of Self-Attention

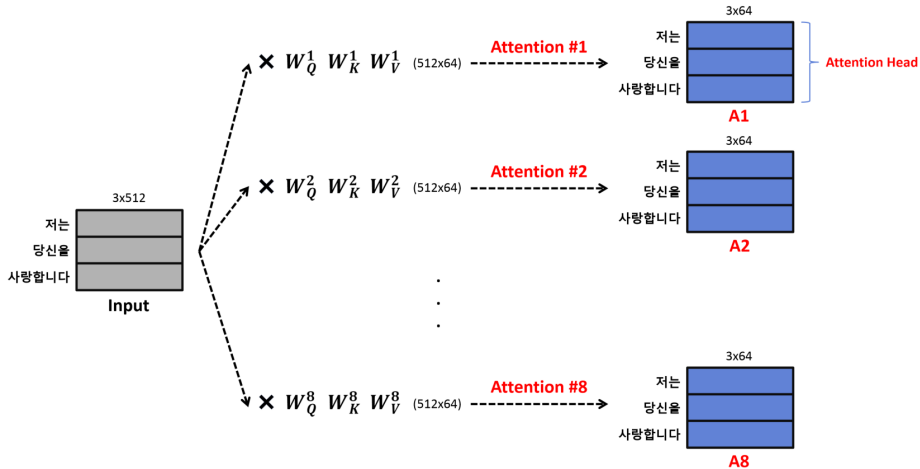Inject contextual information to embedded vector

## Q, K, V in Self Attention



$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{64}})\mathbf{V} =$$
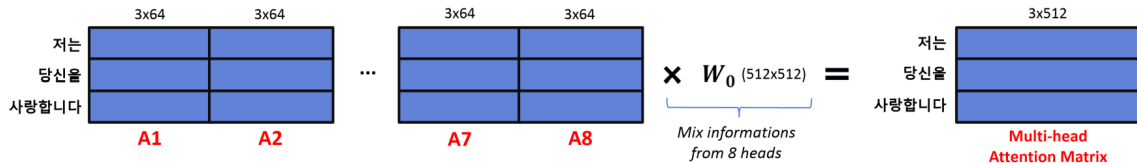
# Multi-head Attention

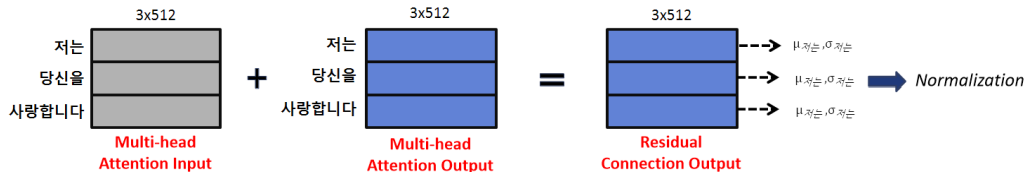## Calculate $N_{heads}$ Attention Matrix

# Multi-head Attention
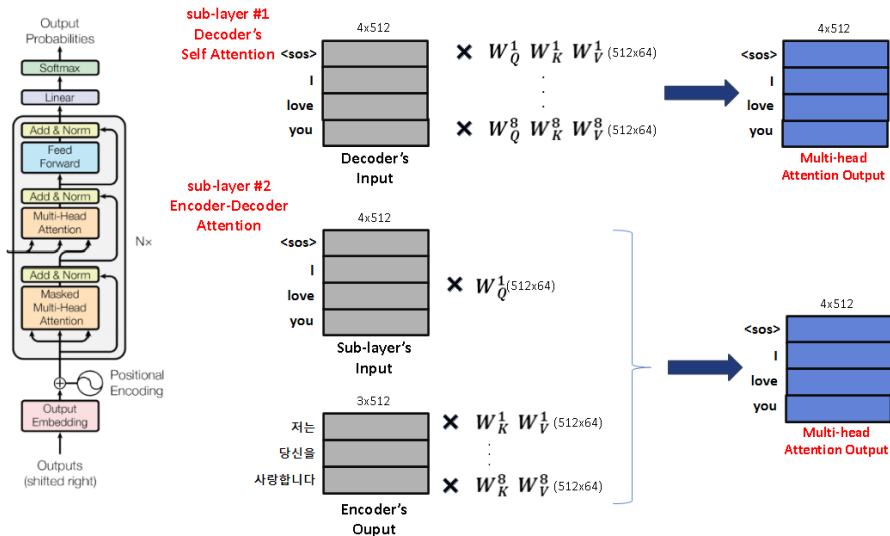
## Multi-head Attention Matrix



## Residual Connection & Layer Normalization

# Attention in Decoder

# Positional Encoding
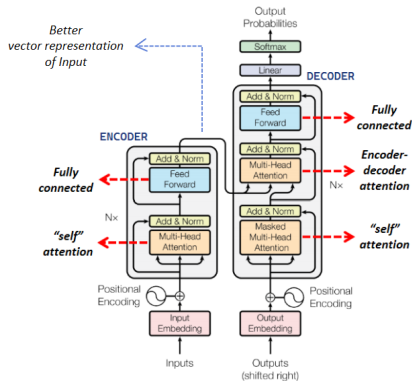
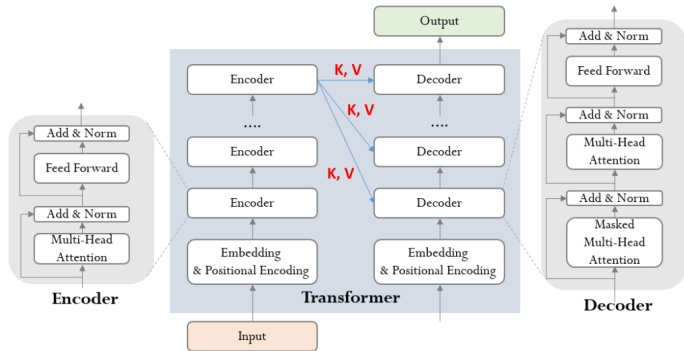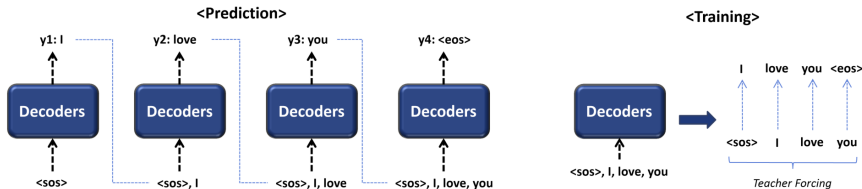# Stacking Encoders and Decoders



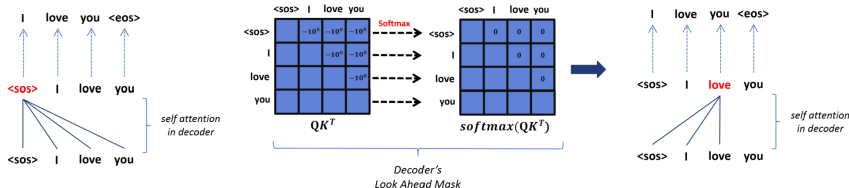Figure 1: The Transformer - model architecture.

# How the Transformer Works?

## Prediction & Training



## Masked Self Attention

# Evaluation

Table: Comparison with English-to-German and English-to-French newstest2014 tests

| Model | Base | BLEU | | Training Cost (FLOPs) | |
|---|---|---|---|---|---|
| | | EN-DE | EN-FR | EN-DE | EN-FR |
| ByteNet [18] | CNN | 23.75 | 39.2 | $1.0 \times 10^{20}$ | - |
| Deep-Att + PosUnk [39] | RNN + Attention | - | 39.2 | - | $1.0 \times 10^{20}$ |
| GNMT + RL [38] | RNN | 24.6 | 39.92 | $2.3 \times 10^{19}$ | $1.4 \times 10^{20}$ |
| ConvS2S [9] | CNN | 25.16 | 40.46 | $9.6 \times 10^{18}$ | $1.5 \times 10^{20}$ |
| MoE [32] | MoE | 26.03 | 40.56 | $2.0 \times 10^{19}$ | $1.2 \times 10^{20}$ |
| Deep-Att + PosUnk Ensemble [39] | RNN + Attention | - | 40.4 | - | $8.0 \times 10^{20}$ |
| GNMT + RL Ensemble [38] | RNN | 26.30 | 41.16 | $1.8 \times 10^{20}$ | $1.1 \times 10^{21}$ |
| ConvS2S Ensemble [9] | CNN | 26.36 | 41.29 | $7.7 \times 10^{19}$ | $1.2 \times 10^{21}$ |
| Transformer (base model) | Transformer | 27.3 | 38.1 | $3.3 \times 10^{18}$ | - |
| Transformer (big) | Transformer | **28.4** | **41.8** | $2.3 \times 10^{19}$ | - |

# References

📄 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin.
**Attention is All You Need**
In *Advances in Neural Information Processing Systems*, 2017.

📄 Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio.
**Neural Machine Translation by Jointly Learning to Align and Translate**
In *International Conference on Learning Representations (ICLR)*, 2015.

📄 Ilya Sutskever, Oriol Vinyals, and Quoc V. Le.
**Sequence to Sequence Learning with Neural Networks**
In *Advances in Neural Information Processing Systems*, 2014.