

# Attention is all you need

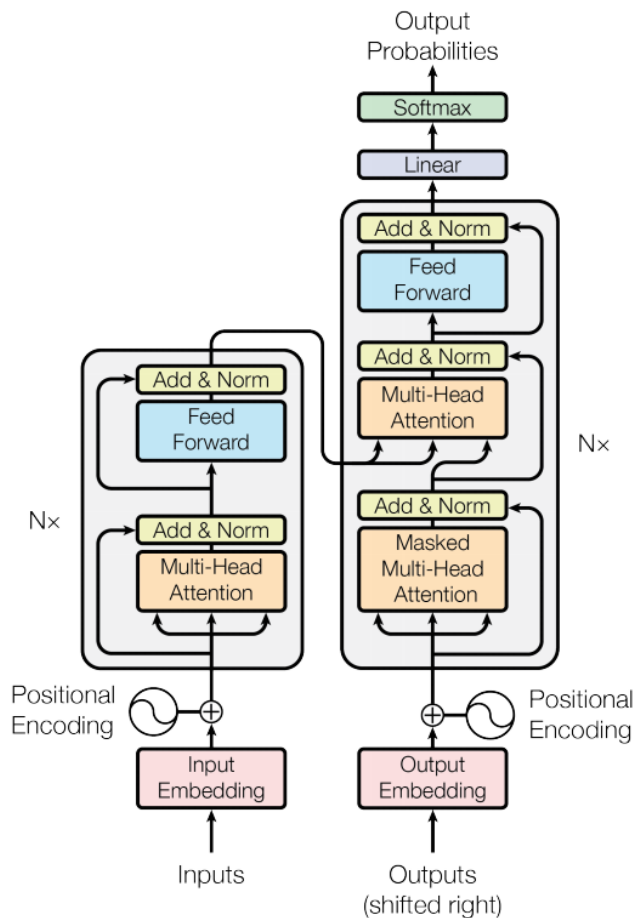


Figure 1: The Transformer - model architecture.

번역 문제에 RNN과 CNN을 사용하지 않고, Attention만을 이용하여 State-of-the-art 성능을 끌어낸 연구

# 네트워크의 특성

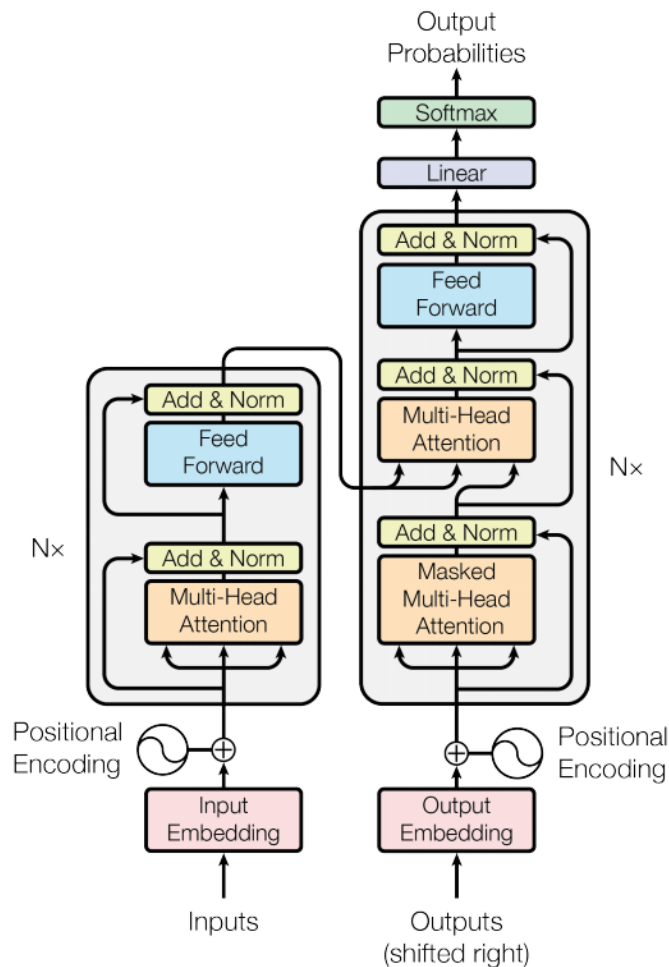


Figure 1: The Transformer - model architecture.

- Seq2seq와 유사한 Transformer 구조 사용
- 제안하는 Scaled Dot-Product Attention과, 이를 병렬로 나열한 Multi-Head Attention 블록이 알고리즘의 핵심
- RNN의 BPTT와 같은 과정이 없으므로 병렬 계산 가능
- 입력된 단어의 위치를 표현하기 위해 Positional Encoding 사용

# Transformer vs. Seq2seq

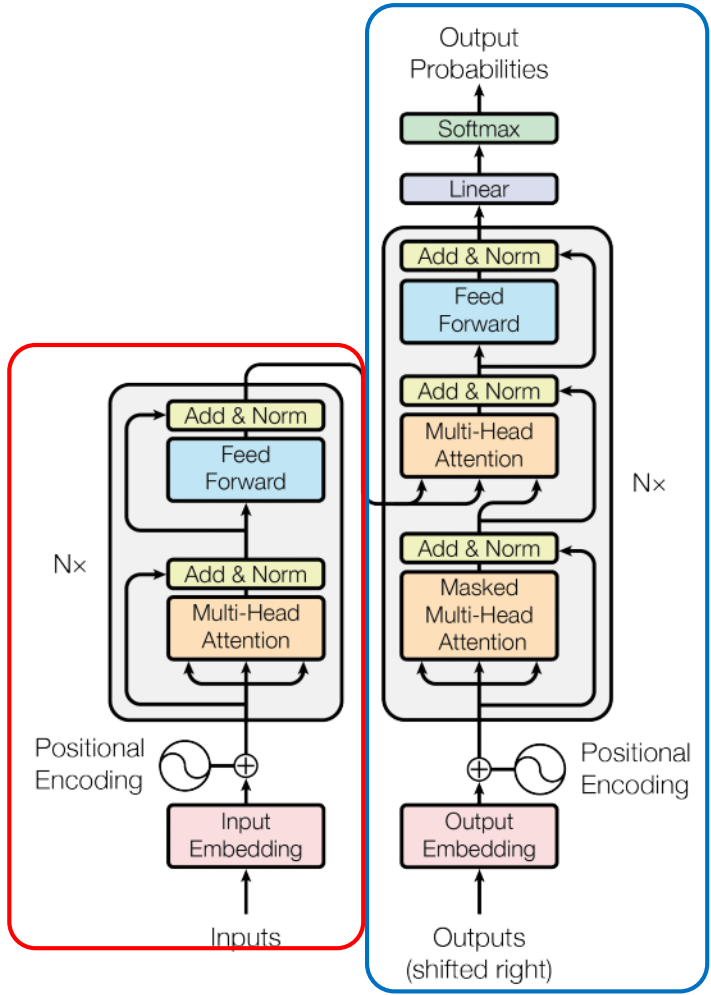
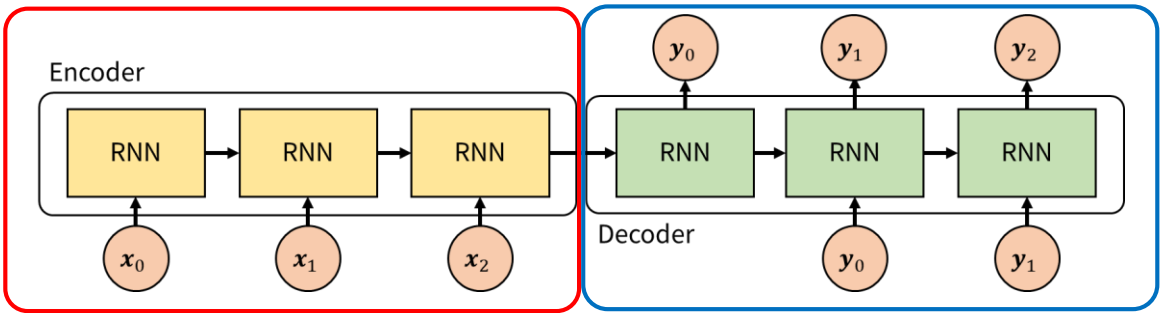


Figure 1: The Transformer - model architecture.



# Inputs & Output

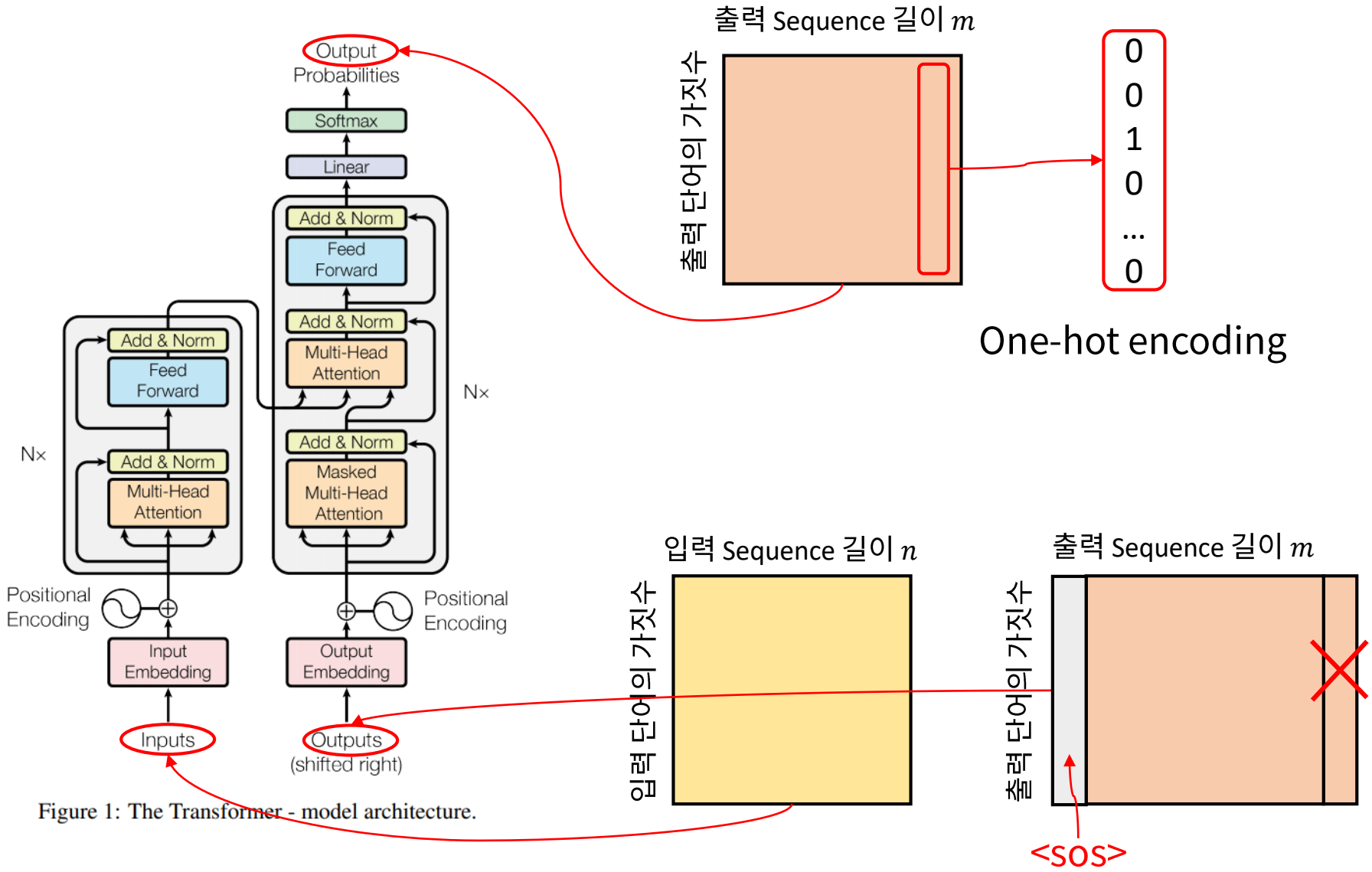
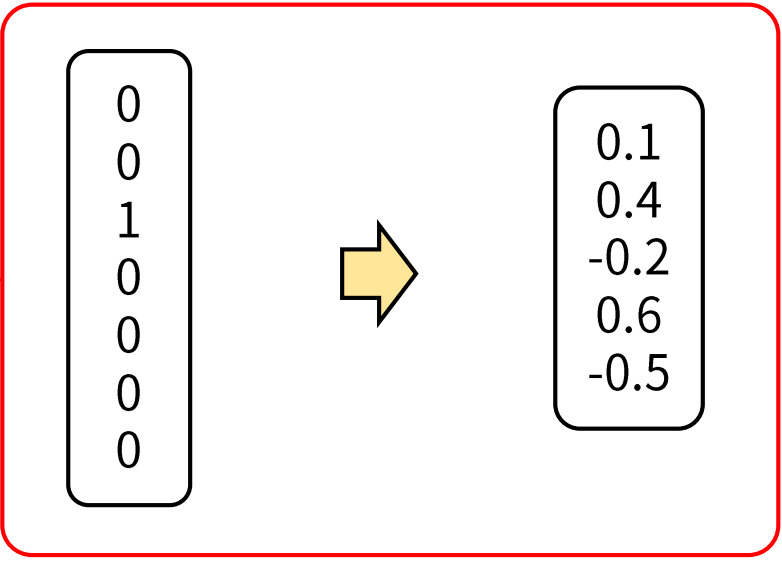
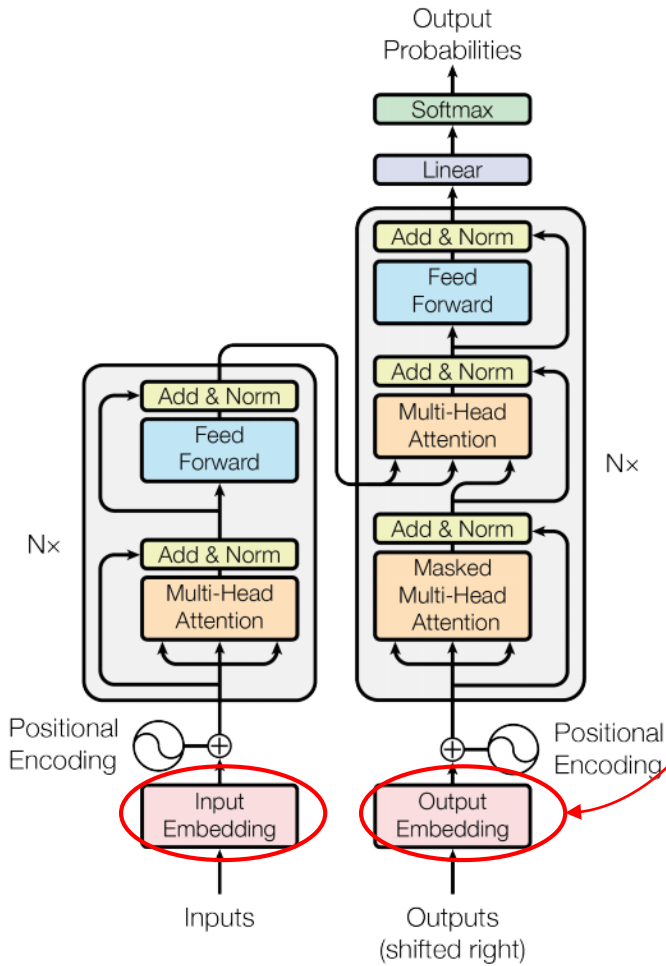


Figure 1: The Transformer - model architecture.

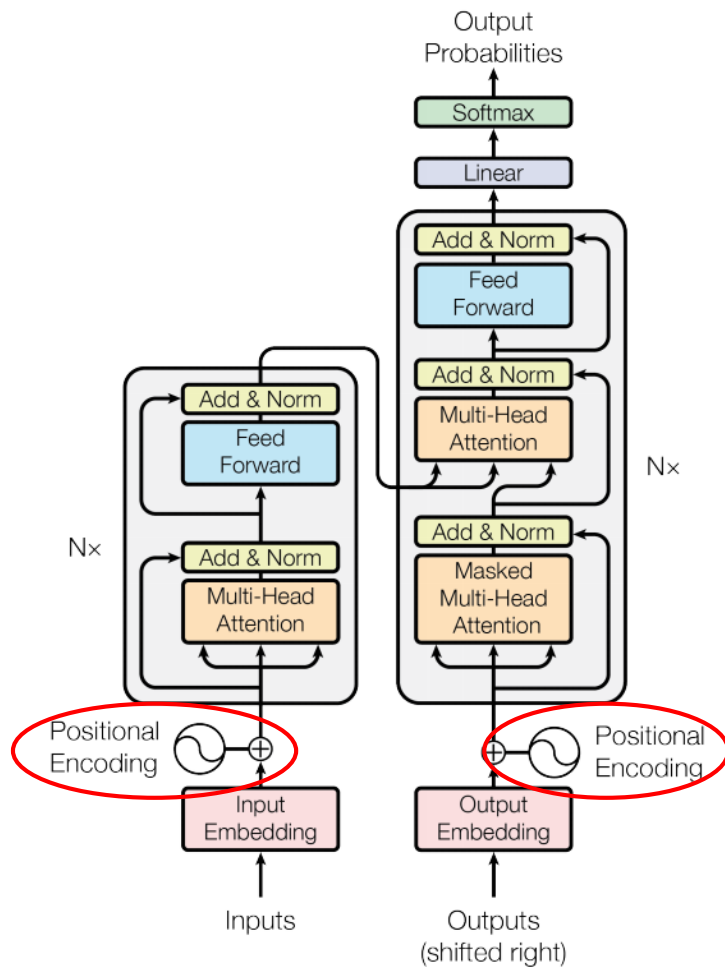
# Word Embedding



One-Hot Encoding → Embedding

One-Hot Encoding된 단어를 실수 형태로 변경하면서 차원의 수를 줄이는 방법

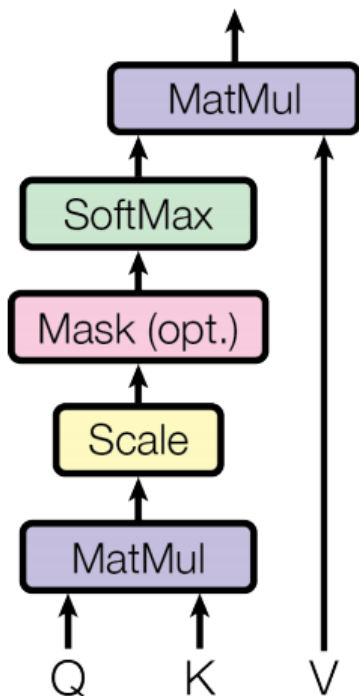
# Positional Encoding



- 시간적 위치별로 **고유의 Code**를 생성하여 더하는 방식
- 전체 Sequence의 길이 중 상대적 위치에 따라 고유의 벡터를 생성하여 Embedding된 벡터에 더해줌

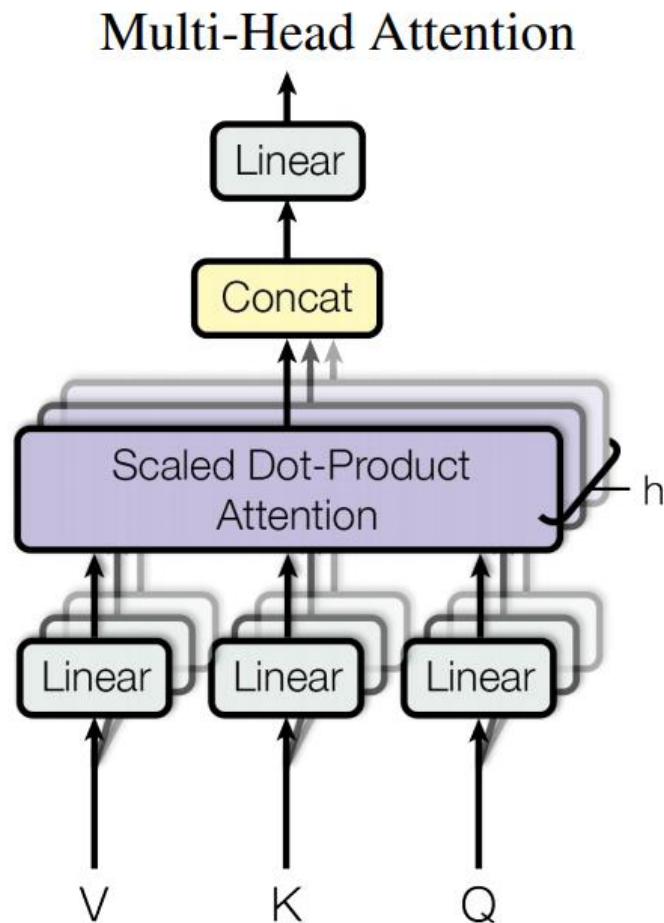
# Scaled Dot-Product Attention

## Scaled Dot-Product Attention



- Query, Key-Value의 구조를 띄고 있음
- Q와 K의 비교 함수는 **Dot-Product**와 **Scale**로 이루어짐
- Mask를 이용해 Illegal connection의 attention을 금지
- Softmax로 유사도를 0 ~ 1의 값으로 Normalize
- 유사도와 V를 결합해 **Attention value** 계산

# Multi-Head Attention



- Linear 연산 (Matrix Mult)를 이용해  $Q, K, V$ 의 차원을 감소  
Q와 K의 차원이 다른 경우 이를 이용해 동일하게 맞춤
- $h$ 개의 Attention Layer를 병렬적으로 사용 – 더 넓은 계층
- 출력 직전 Linear 연산을 이용해 Attention Value의 차원을 필요에 따라 변경
- 이 메커니즘을 통해 병렬 계산에 유리한 구조를 가지게 됨



# Masked Multi-Head Attention

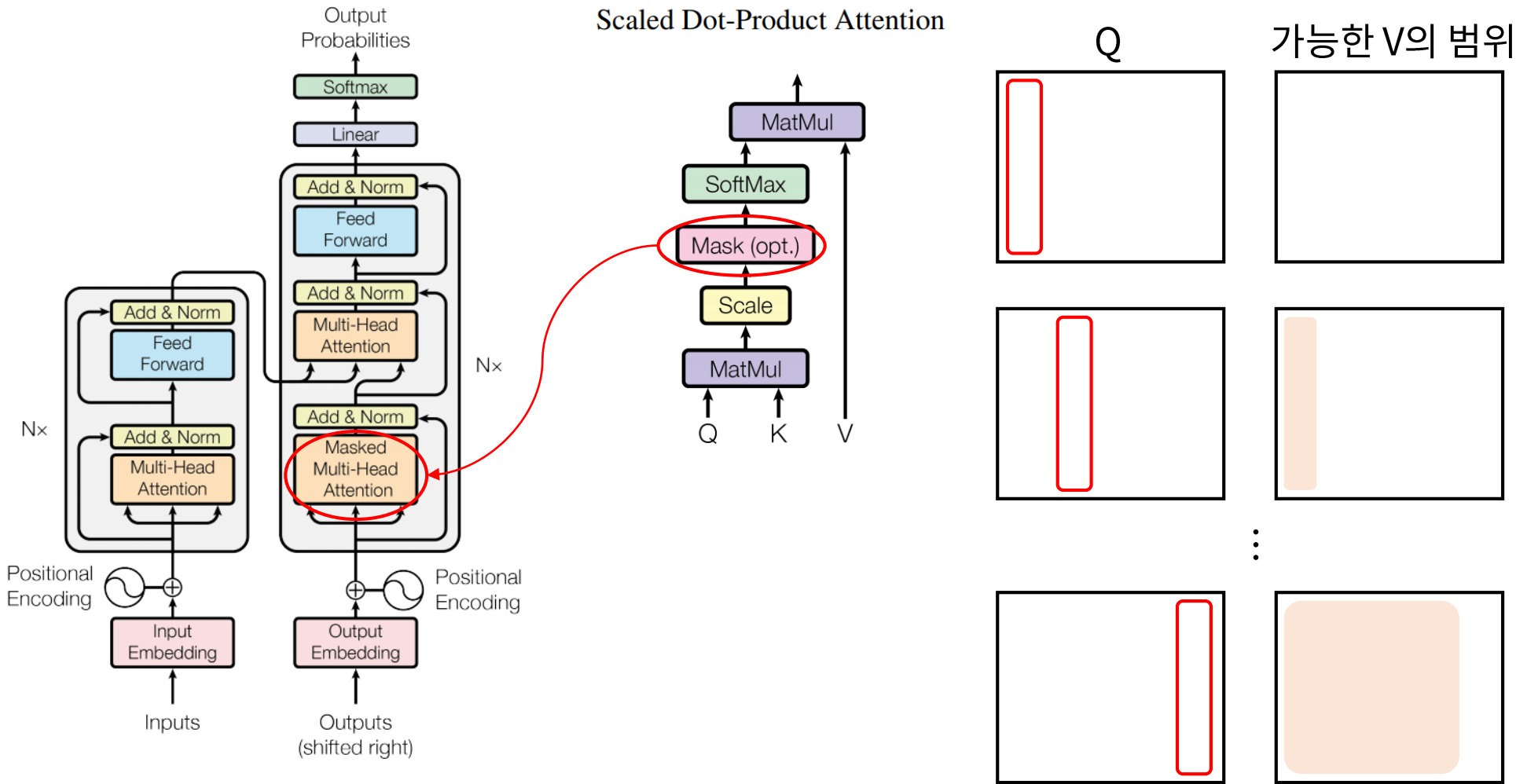


Figure 1: The Transformer - model architecture.

Self-Attention에서 자기 자신을 포함한 미래의 값과는 Attention을 구하지 않기 때문에, Masking을 사용한다.

# Multi-Head Attention in Action

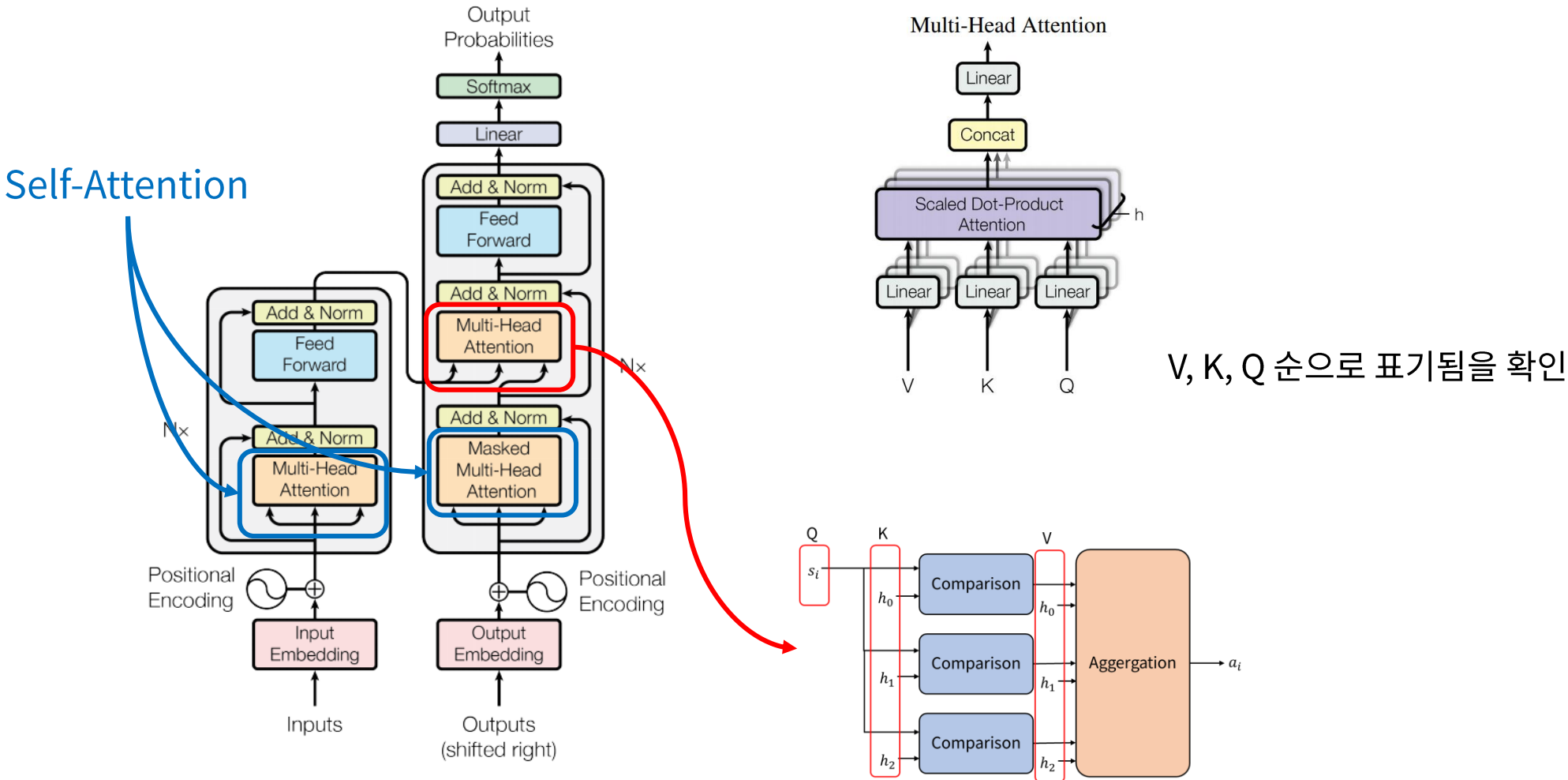


Figure 1: The Transformer - model architecture.

Seq2seq의 Attention과 동일한 구조

# Position-wise Feed-Forward

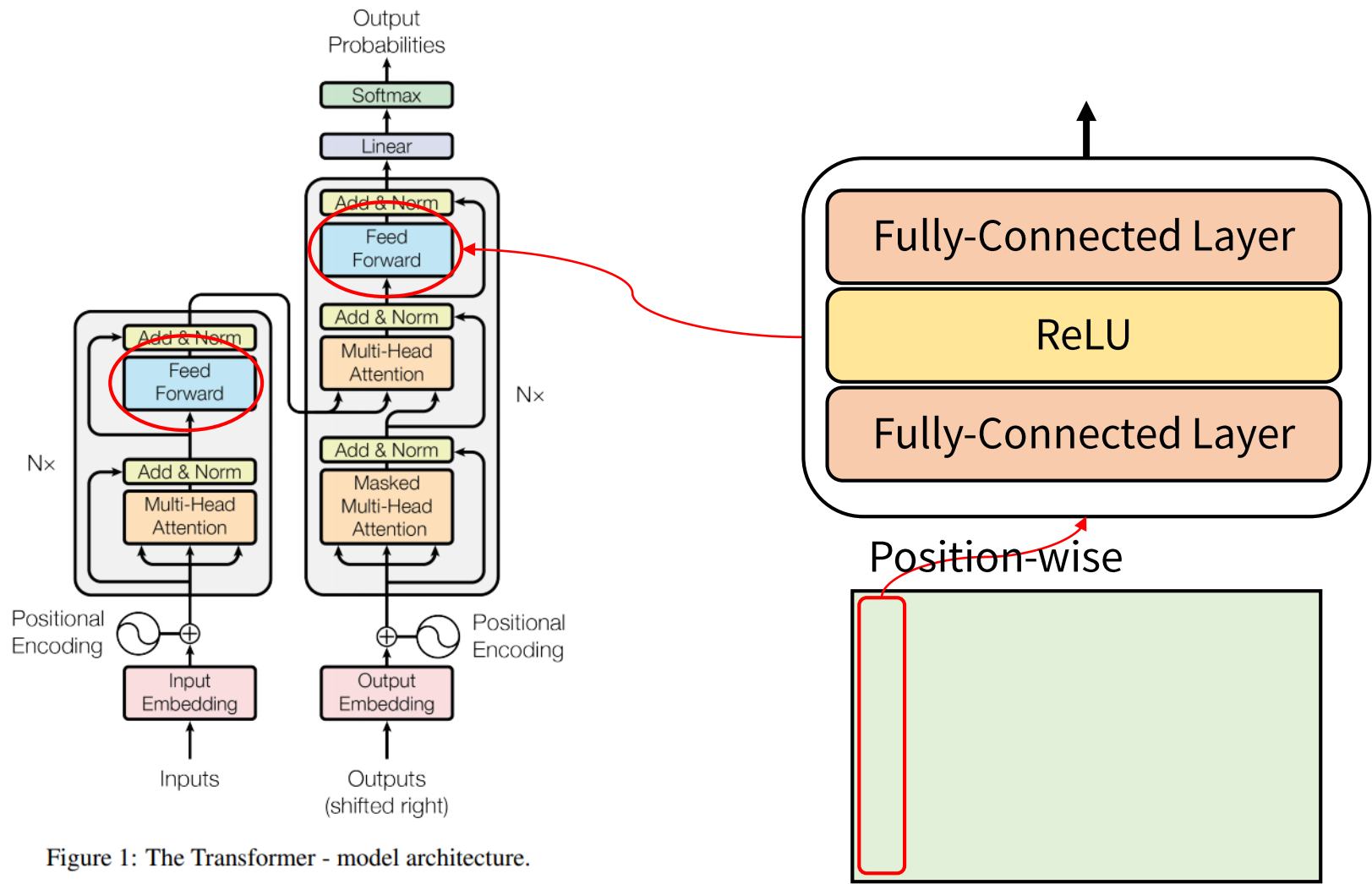


Figure 1: The Transformer - model architecture.

# Add & Norm

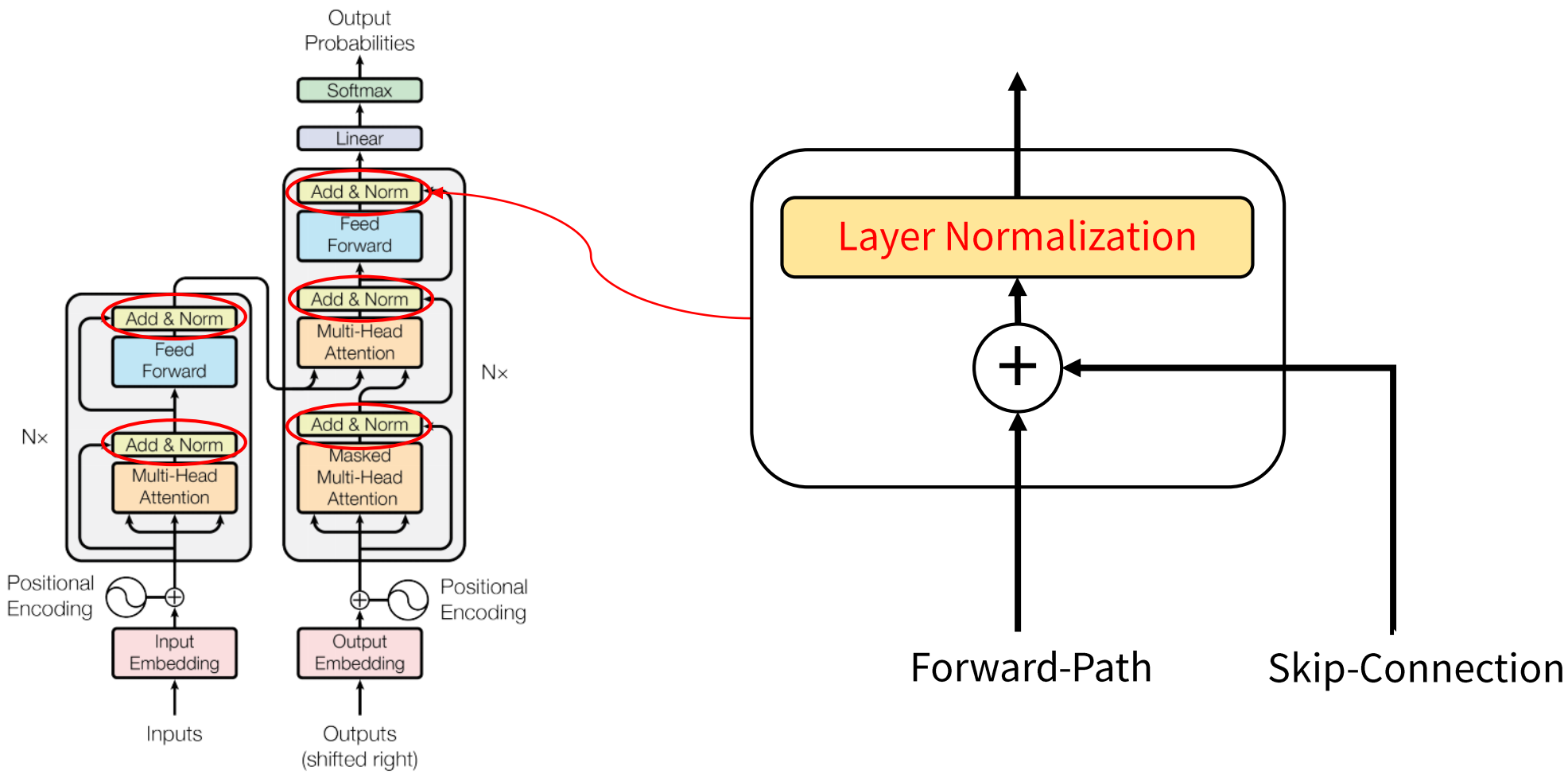


Figure 1: The Transformer - model architecture.

# Output Softmax

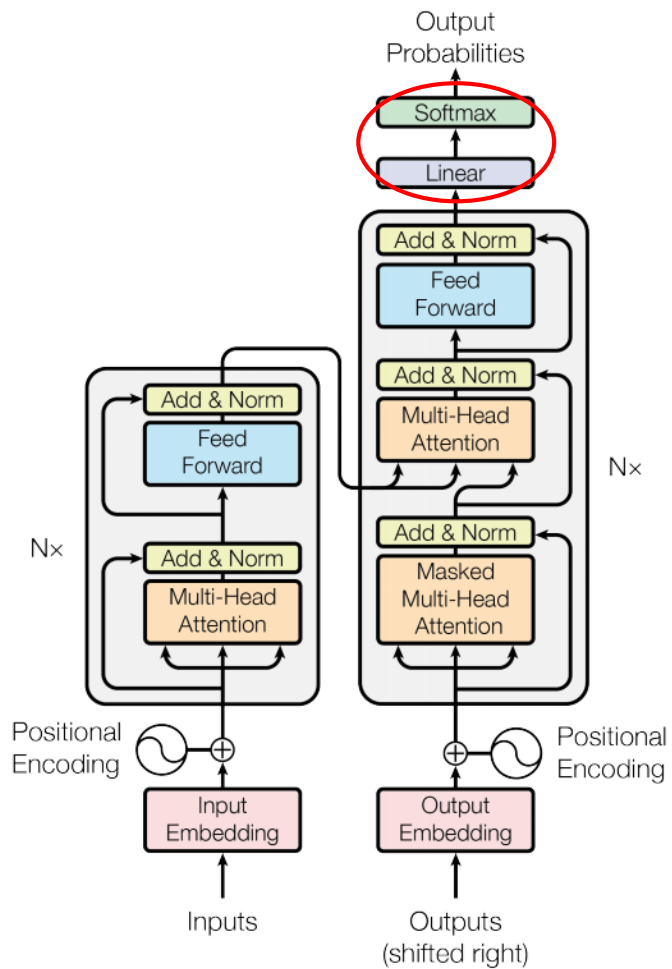


Figure 1: The Transformer - model architecture.

- Linear 연산을 이용해 출력 단어 종류의 수에 맞춤
- Softmax를 이용해 어떤 단어인지 Classification문제 해결

# Attention is really all you need

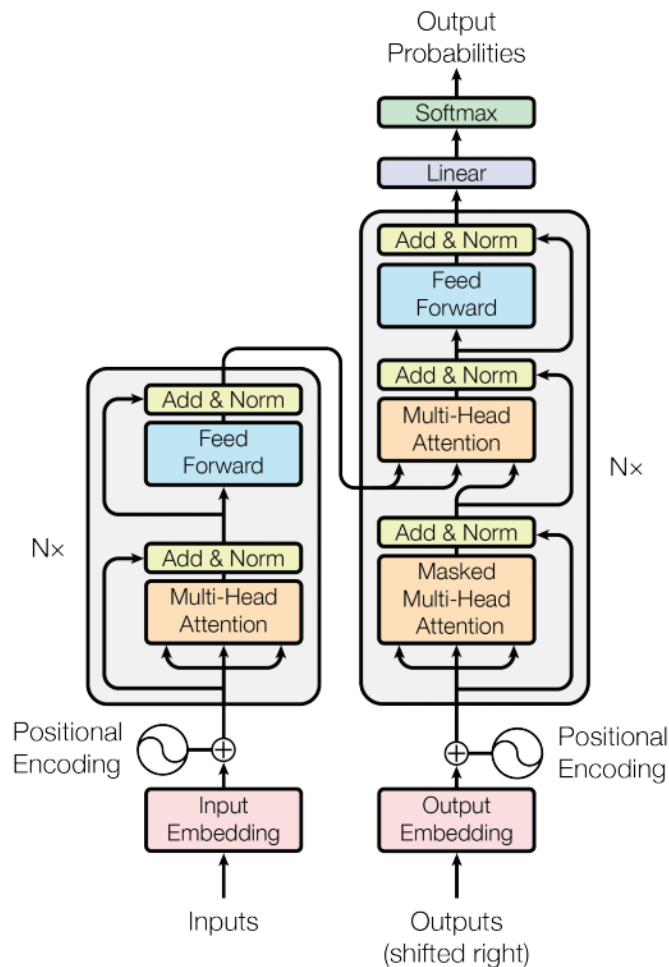


Figure 1: The Transformer - model architecture.

다시 한번 아래 내용을 곱씹어 봅시다.

- Seq2seq와 유사한 Transformer 구조 사용
- 제안하는 Scaled Dot-Product Attention과, 이를 병렬로 나열한 Multi-Head Attention 블록이 알고리즘의 핵심
- RNN의 BPTT와 같은 과정이 없으므로 병렬 계산 가능
- 입력된 단어의 위치를 표현하기 위해 Positional Encoding 사용