

Transformer: Multi-head Attention

Ki Hyun Kim

nlp.with.deep.learning@gmail.com

Attention: Query Generation

- Example:



강남역에서 가장 회식 하기 좋은 오리고기 맛있는 집은 어디야?



Google Search

I'm Feeling Lucky

Google offered in: [한국어](#)



강남역 오리고기 회식장소 맛집



Google Search

I'm Feeling Lucky

Google offered in: [한국어](#)

마음의 상태(state)를 잘 반영하면서
좋은 검색 결과를 이끌어내는 쿼리를 얻기 위함



Attention: Query Generation

- Example:

Google

강남역에서 가장 회식 하기 좋은 오리고기 맛있는 집은 어디야?



만약 검색을 다양하게 할 수 있다면?

Google Search

I'm Feeling Lucky

Google offered in: 한국어

Google

강남역 오리고기 회식장소 맛집



Google Search

I'm Feeling Lucky

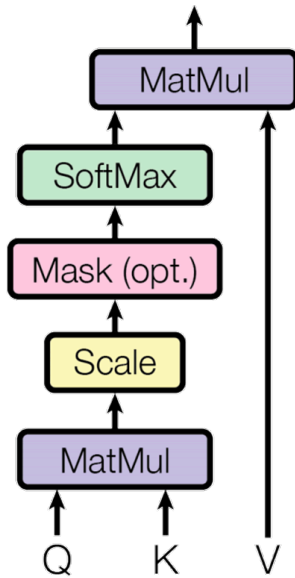
Google offered in: 한국어

마음의 상태(state)를 잘 반영하면서
좋은 검색 결과를 이끌어내는 쿼리를 얻기 위함

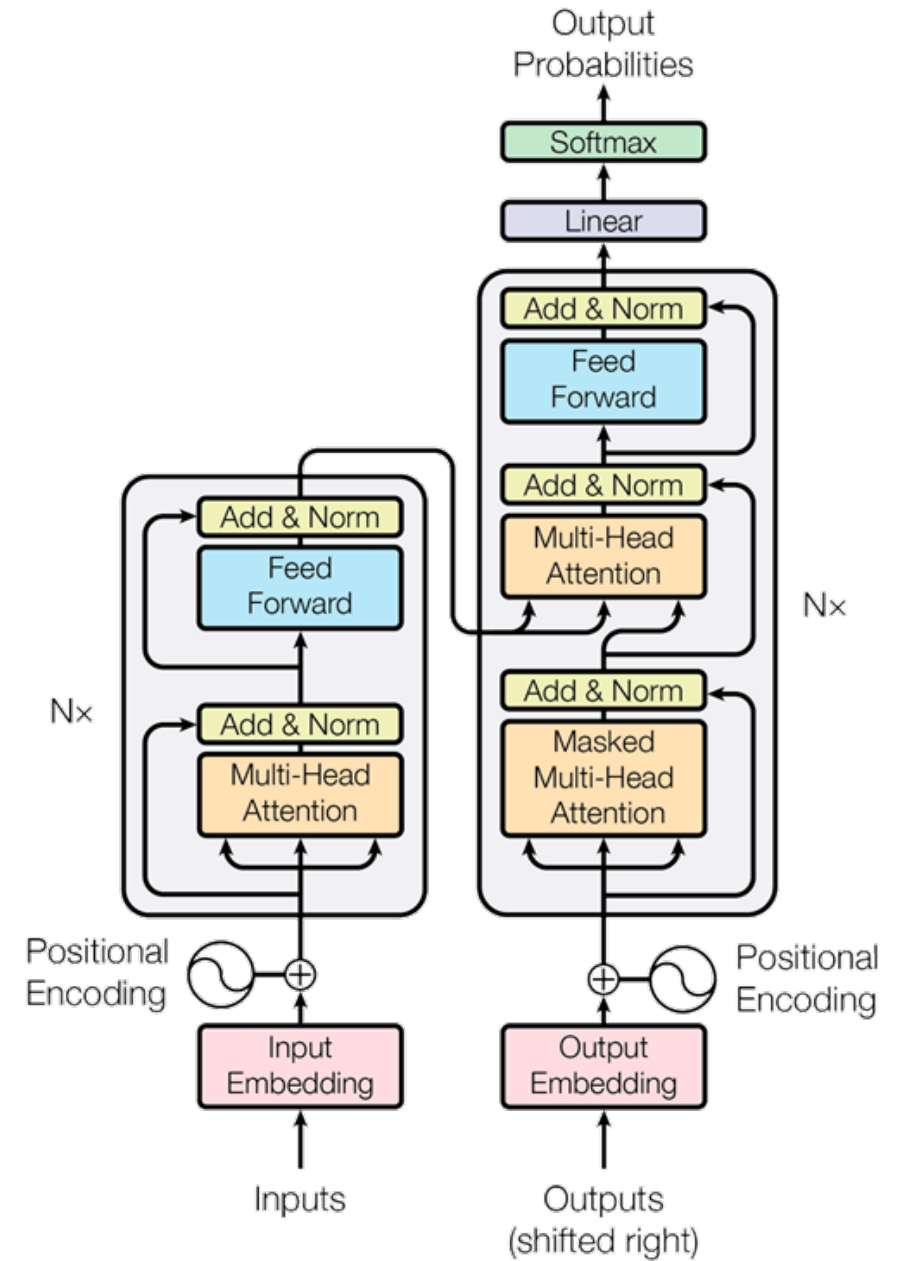
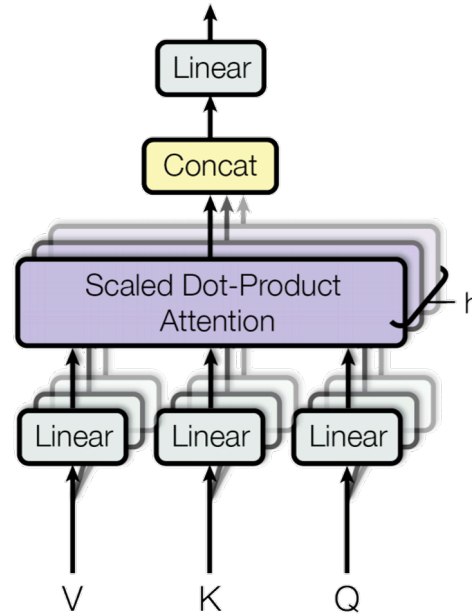


Transformer & Attention

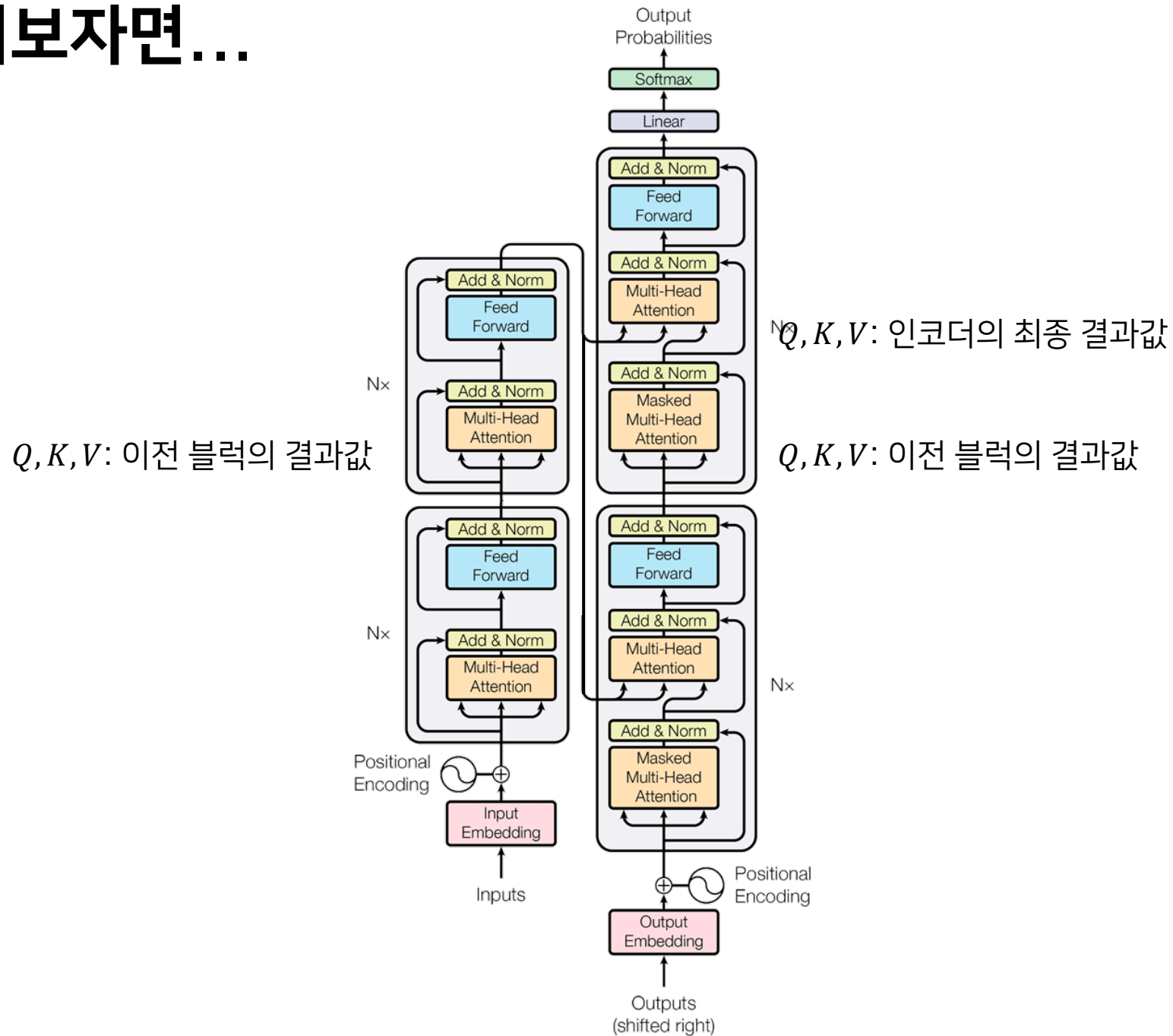
Scaled Dot-Product Attention



Multi-Head Attention



굳이 그려보자면...



Equations

In case of Q, K and V come from same origin,
 $|Q| = |K| = |V| = (\text{batch_size}, n \text{ or } m, \text{hidden_size})$.

In case of Q and K, V come from different origin,
 $|Q| = (\text{batch_size}, n, \text{hidden_size})$
 $|K| = |V| = (\text{batch_size}, m, \text{hidden_size})$.

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{Q \cdot K^\top}{\sqrt{d_{\text{head}}}}\right) \cdot V$$

$$\text{MultiHead}(Q, K, V) = [\text{head}_1; \dots; \text{head}_h] \cdot W^O$$

where $\text{head}_i = \text{Attention}(Q \cdot W_i^Q, K \cdot W_i^K, V \cdot W_i^V)$,

and $W_i^Q \in \mathbb{R}^{d_{\text{model}} \times d_{\text{head}}}$, $W_i^K \in \mathbb{R}^{d_{\text{model}} \times d_{\text{head}}}$, $W_i^V \in \mathbb{R}^{d_{\text{model}} \times d_{\text{head}}}$, $W^O \in \mathbb{R}^{(h \times d_{\text{head}}) \times d_{\text{model}}}$.

$$d_{\text{head}} = d_{\text{model}} / h = 64$$

$$h = 8, d_{\text{model}} = 512$$

Summary

- Previous method: attention in sequence to sequence
 - Query를 잘 만들어 key-value를 잘 matching 시키자
- Multi-head Attention
 - 여러 개의 query를 만들어 다양한 정보를 잘 얻어오자
- Attention 자체로도 정보의 encoding과 decoding이 가능함을 보여줌