

Advance Deep Learning

Advance Deep Learning

Week 1: Attention

Week 2:

- Transfer Learning
- Multi-Task Learning

Week 3: Generative Adversarial Models

Week 4: Real World Applications

What is Attention?

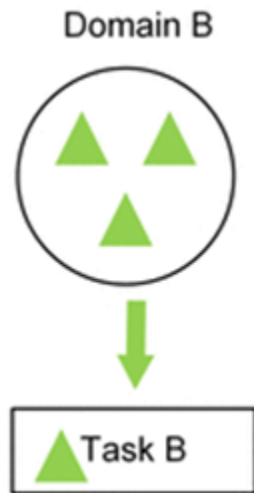
What is Attention?

It is a nice movie.

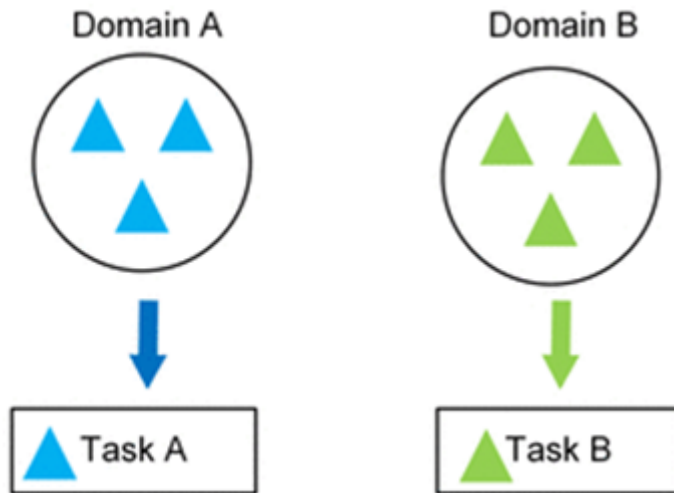
What is Attention?



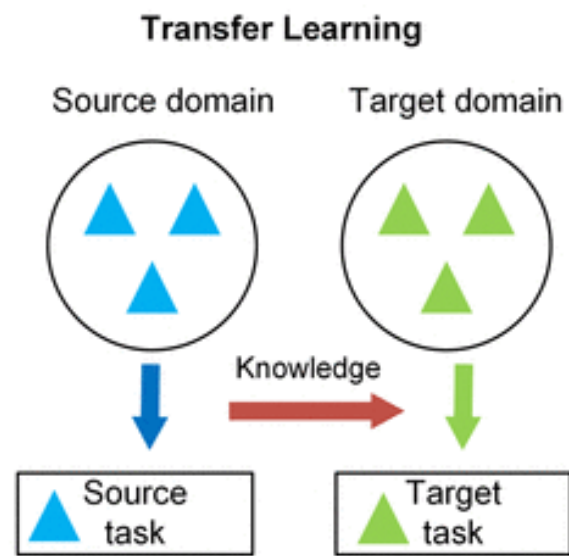
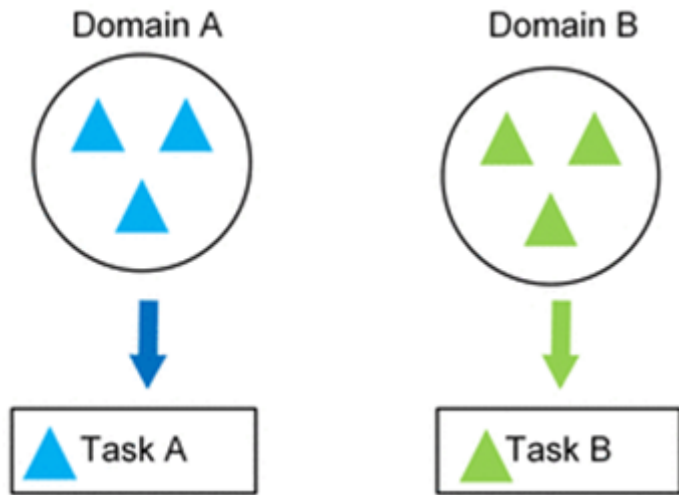
What is Transfer Learning?



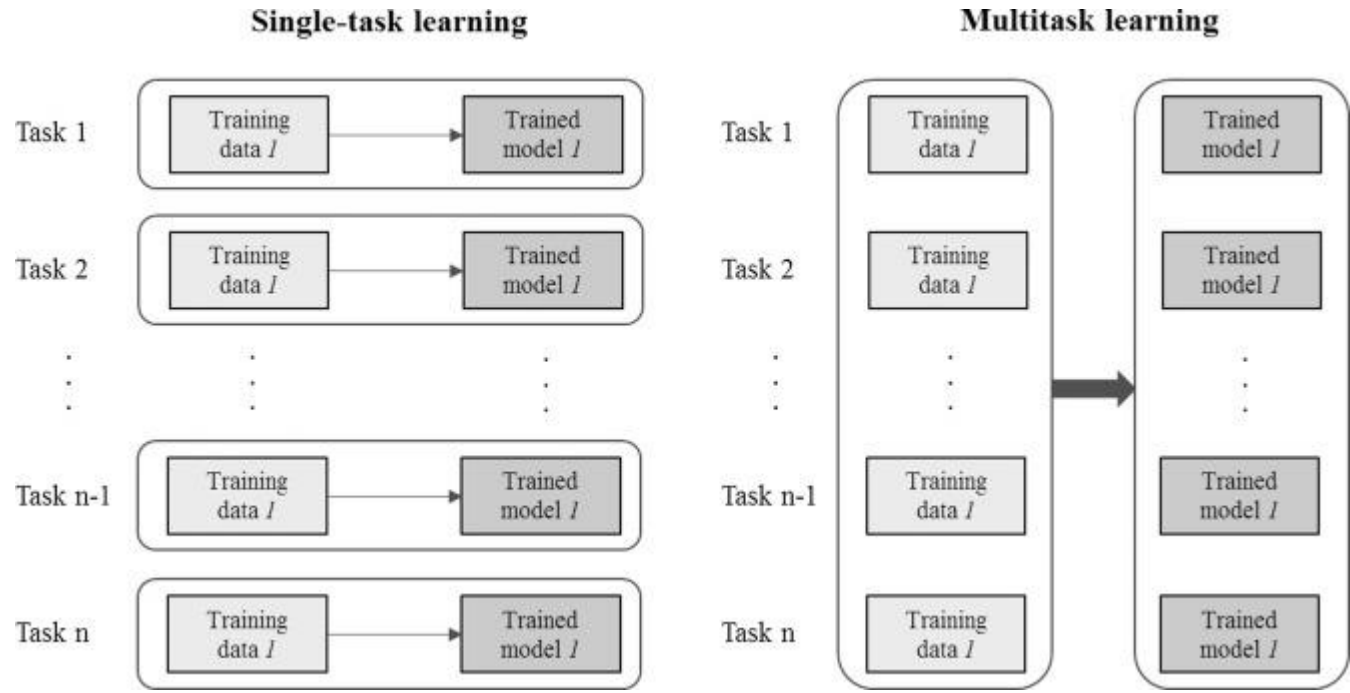
What is Transfer Learning?



What is Transfer Learning?

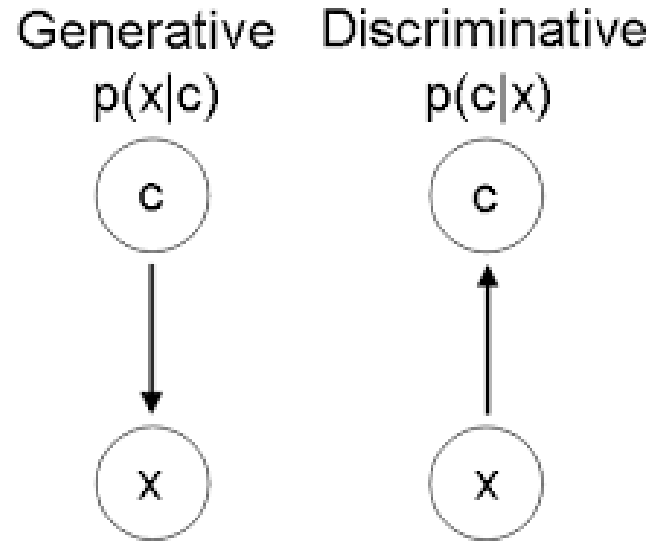


What is Multi-task Learning?

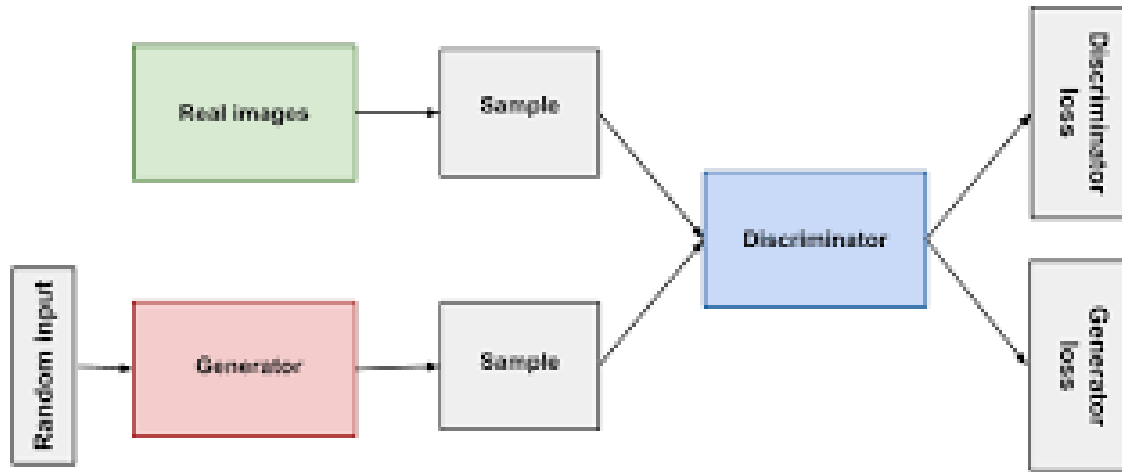


What is Generative Adversarial Network?

What are Generative and Discriminative Models?



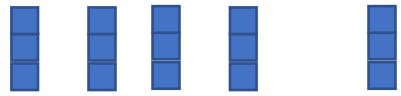
Generative Adversarial Network?



Attention

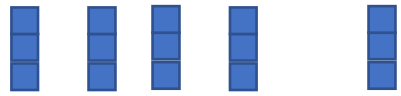
It is a nice movie

It is a nice movie

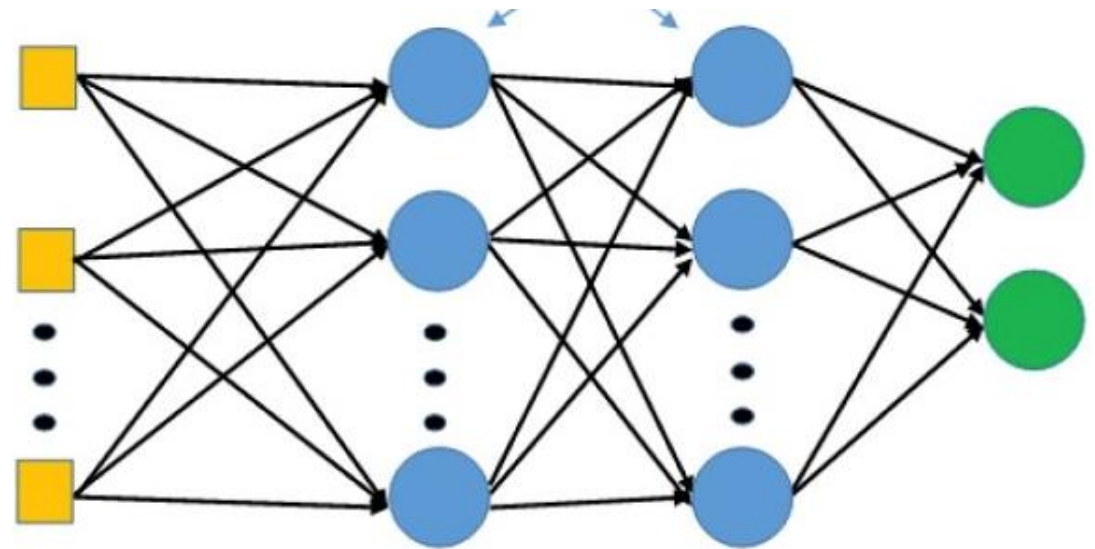


Word Embedding

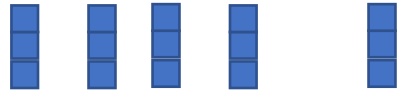
It is a nice movie



Word Embedding



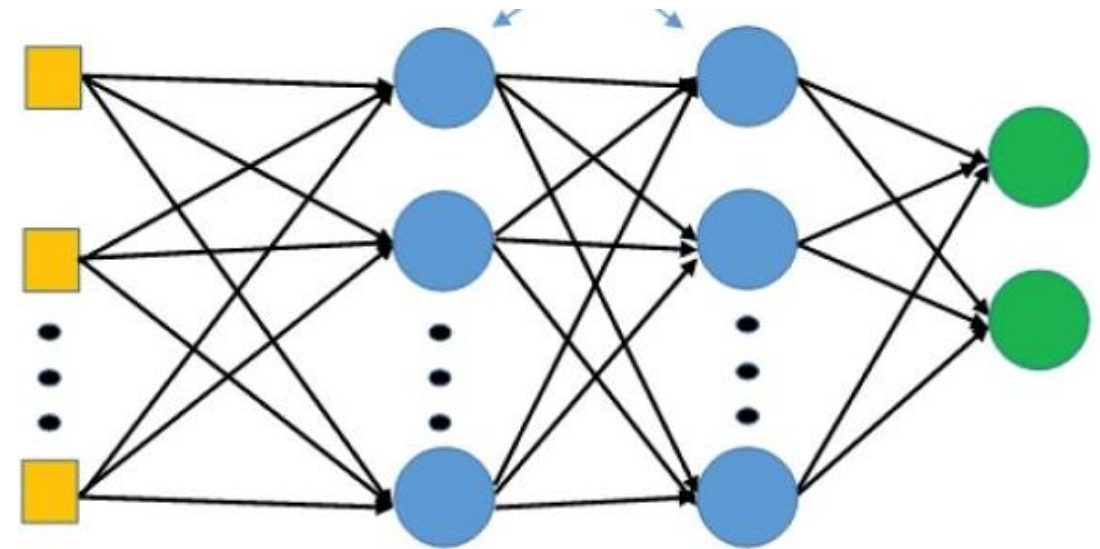
It is a nice movie



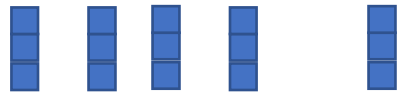
Word Embedding

Generate Embedding of the sentence

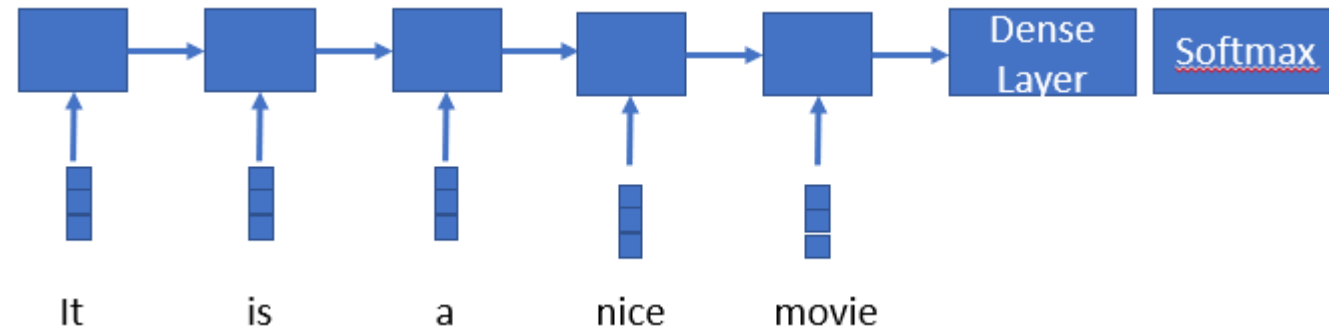
- **Sum of the embedding vector**
- **Average of the embedding vector**
- **Concatenation of the vectors**
- **Or, any sentence embedding**



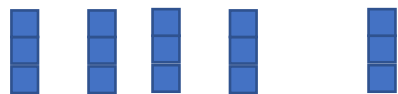
It is a nice movie



Word Embedding



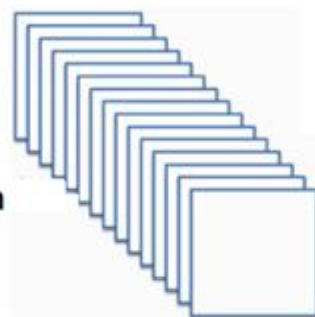
It is a nice movie



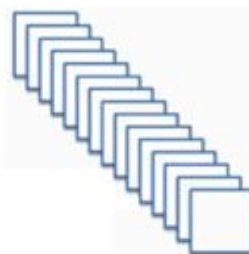
Word Embedding

0	0	0	0	0	0	0
0	1	0	0	0	1	0
0	0	0	0	0	0	0
0	0	0	1	0	0	0
0	1	0	0	0	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

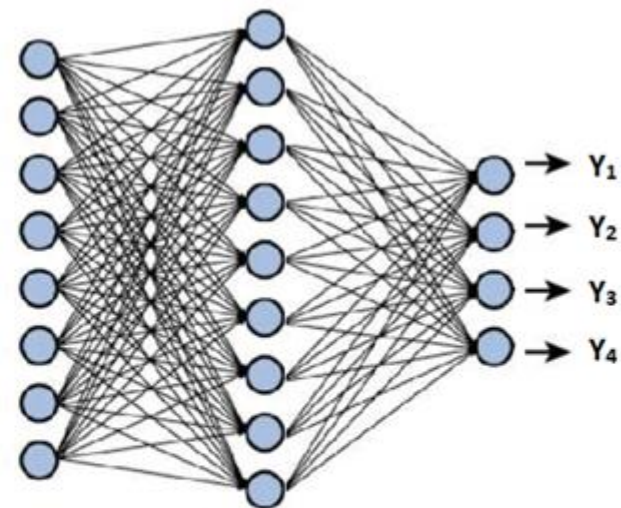
Convolution



Pooling



Flattening



Sentence



Convolutional layer



Pooling layer



Dense layers



Output layer

Sequence-to-Sequence Problem

- POS
- NER
- Machine Translation
- Machine Transliteration
- so on

Sequence-to-Sequence (Seq2Seq) Model

i love you



मैं तुमसे प्यार करता हूँ

main tumase pyaar karata hoon

Sequence-to-Sequence (Seq2Seq) Model

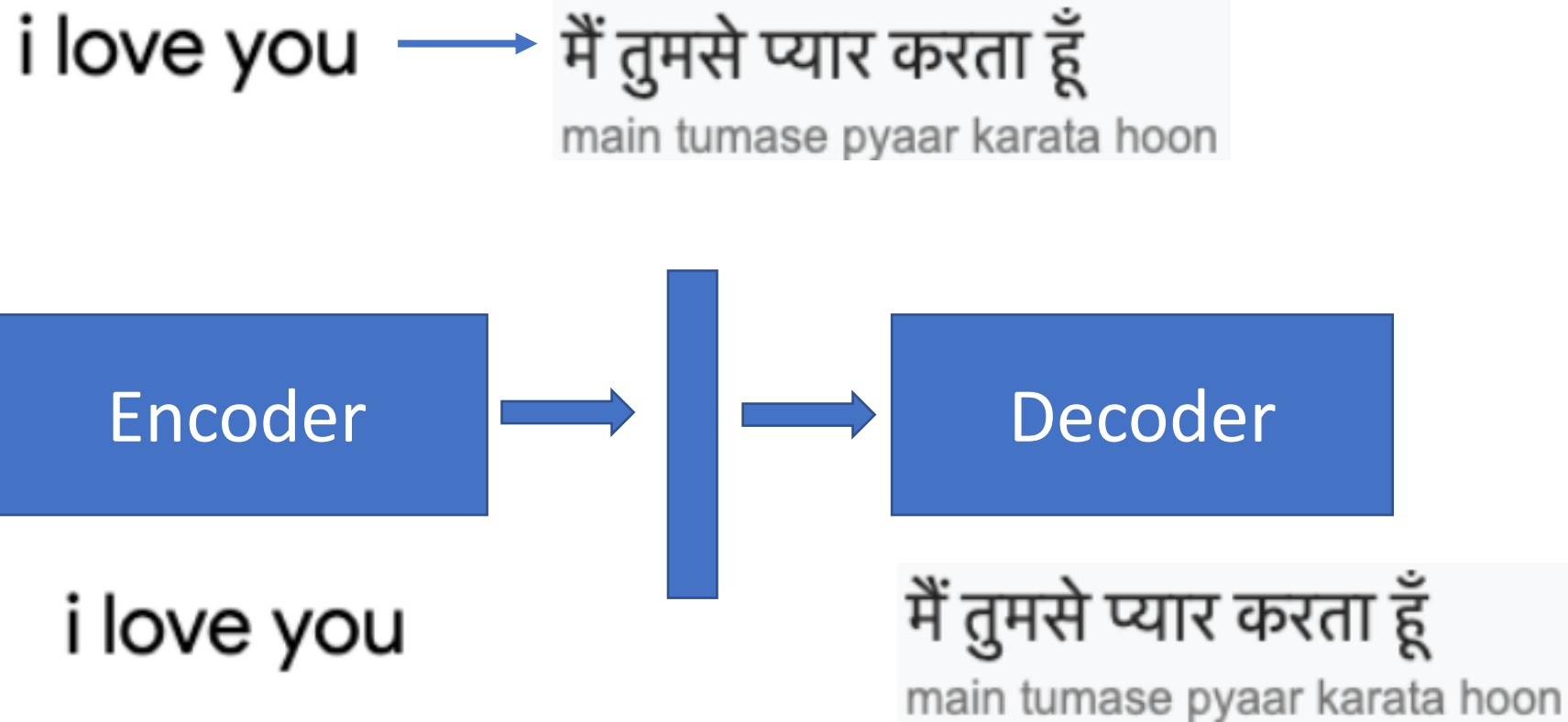
i love you → मैं तुमसे प्यार करता हूँ
main tumase pyaar karata hoon

i → मैं

love → प्यार

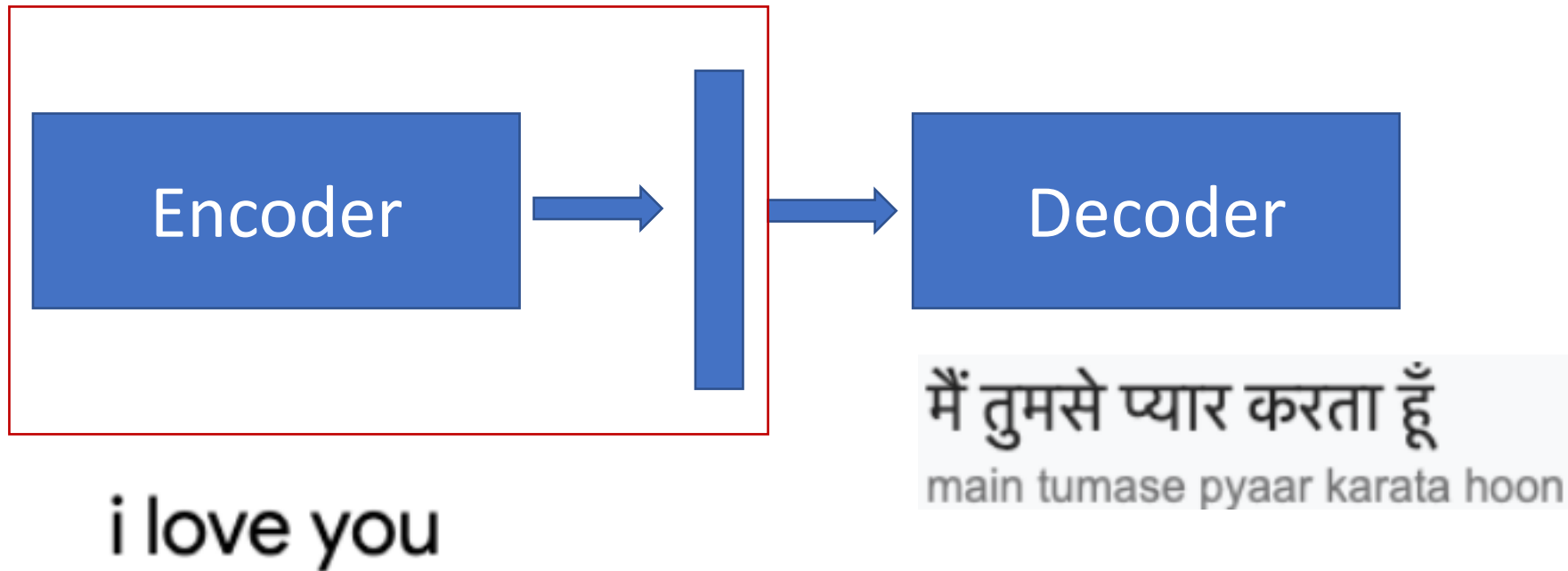
you → तुम

Sequence-to-Sequence (Seq2Seq) Model

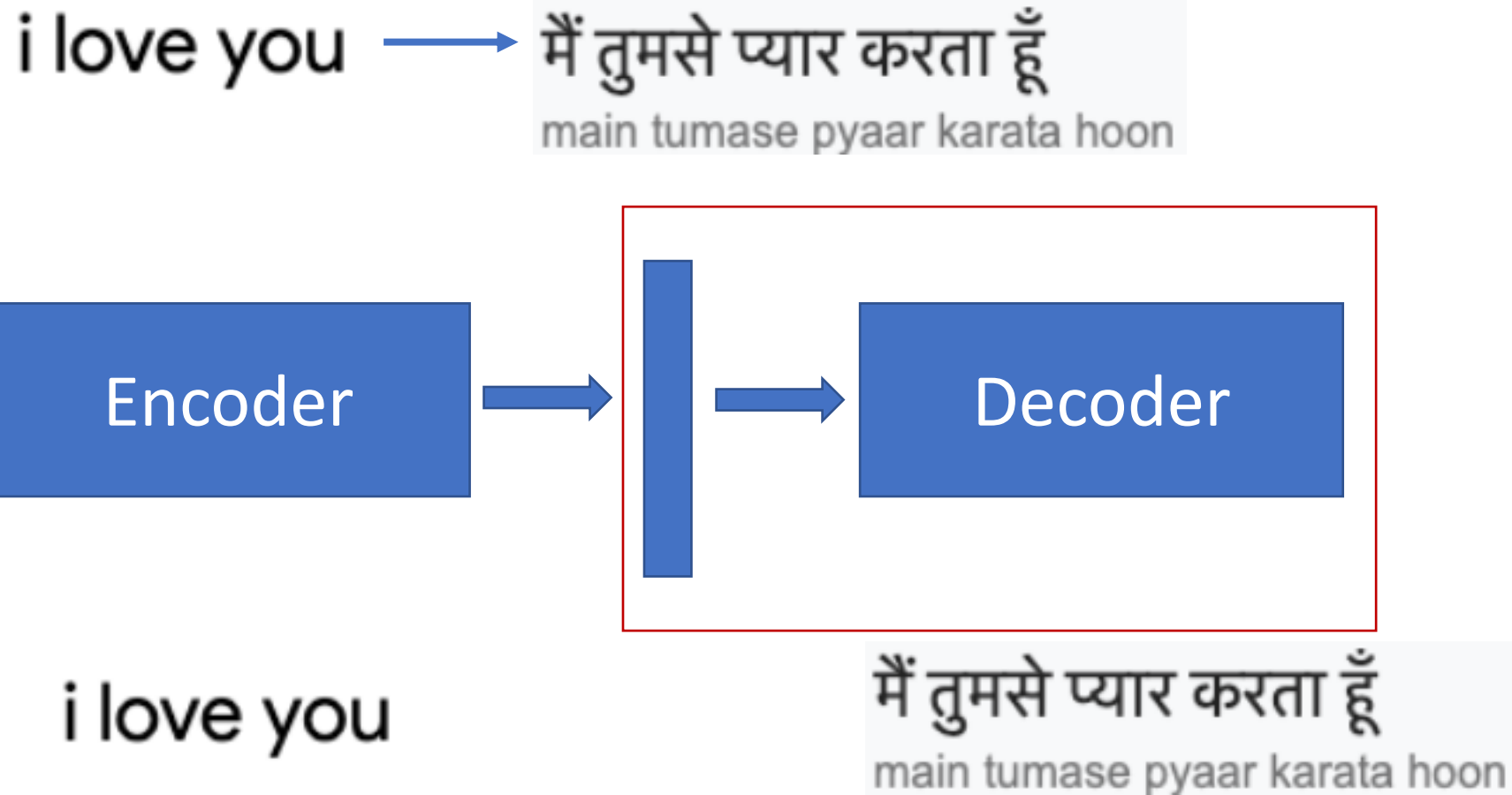


Sequence-to-Sequence (Seq2Seq) Model

i love you → मैं तुमसे प्यार करता हूँ
main tumase pyaar karata hoon

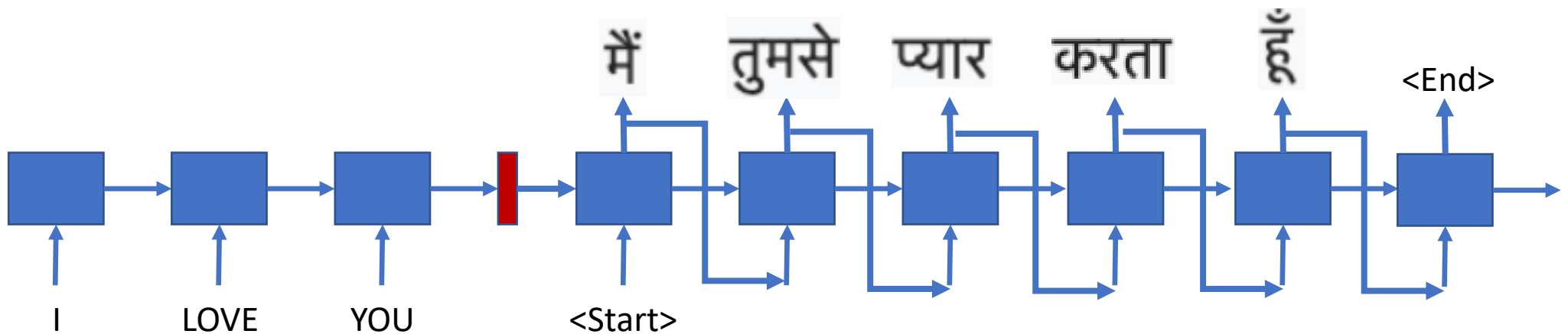


Sequence-to-Sequence (Seq2Seq) Model



Sequence-to-Sequence (Seq2Seq)

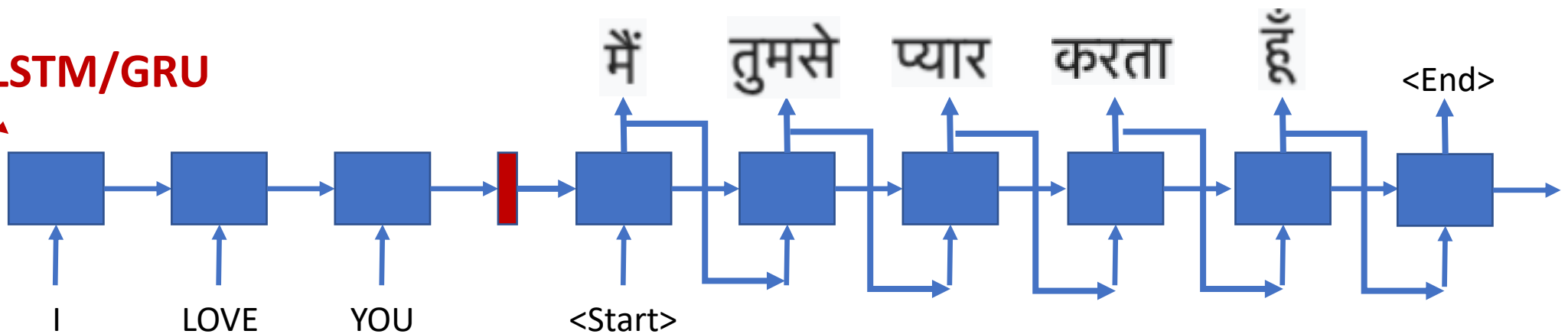
i love you → मैं तुमसे प्यार करता हूँ
main tumase pyaar karata hoon



Sequence-to-Sequence (Seq2Seq)

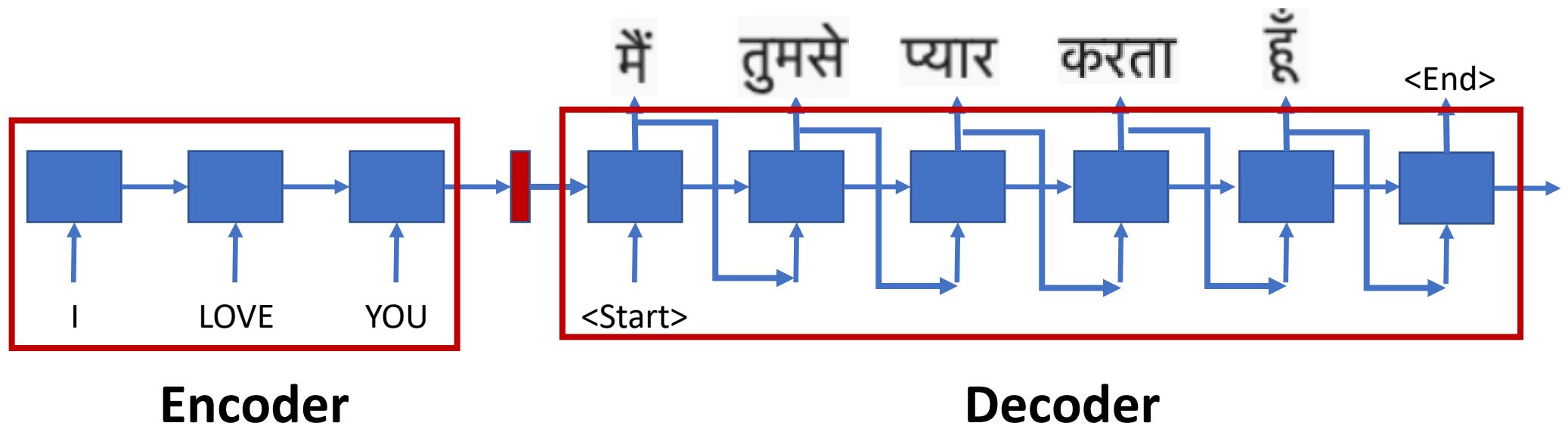
i love you → मैं तुमसे प्यार करता हूँ
main tumase pyaar karata hoon

RNN/LSTM/GRU



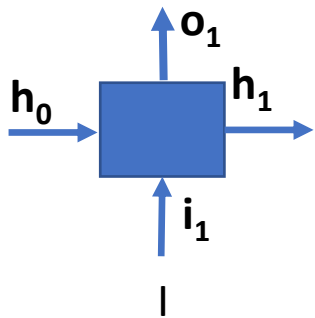
Sequence-to-Sequence (Seq2Seq)

i love you → मैं तुमसे प्यार करता हूँ
main tumase pyaar karata hoon



Sequence-to-Sequence (Seq2Seq)

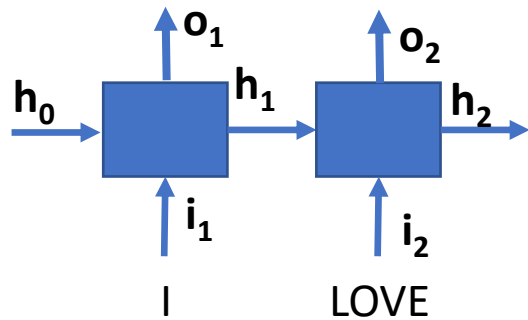
i love you → मैं तुमसे प्यार करता हूँ
main tumase pyaar karata hoon



Encoder

Sequence-to-Sequence (Seq2Seq)

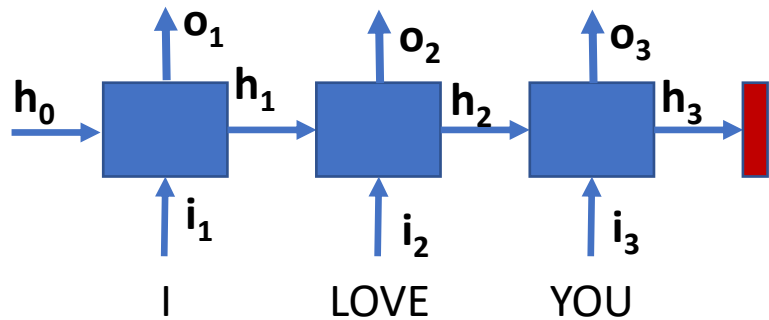
i love you → मैं तुमसे प्यार करता हूँ
main tumase pyaar karata hoon



Encoder

Sequence-to-Sequence (Seq2Seq)

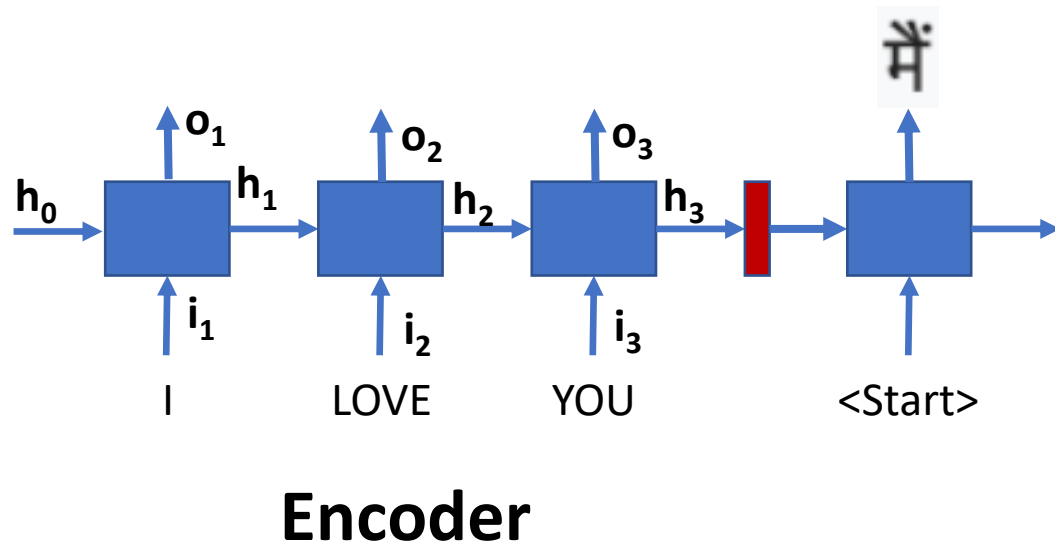
i love you → मैं तुमसे प्यार करता हूँ
main tumase pyaar karata hoon



Encoder

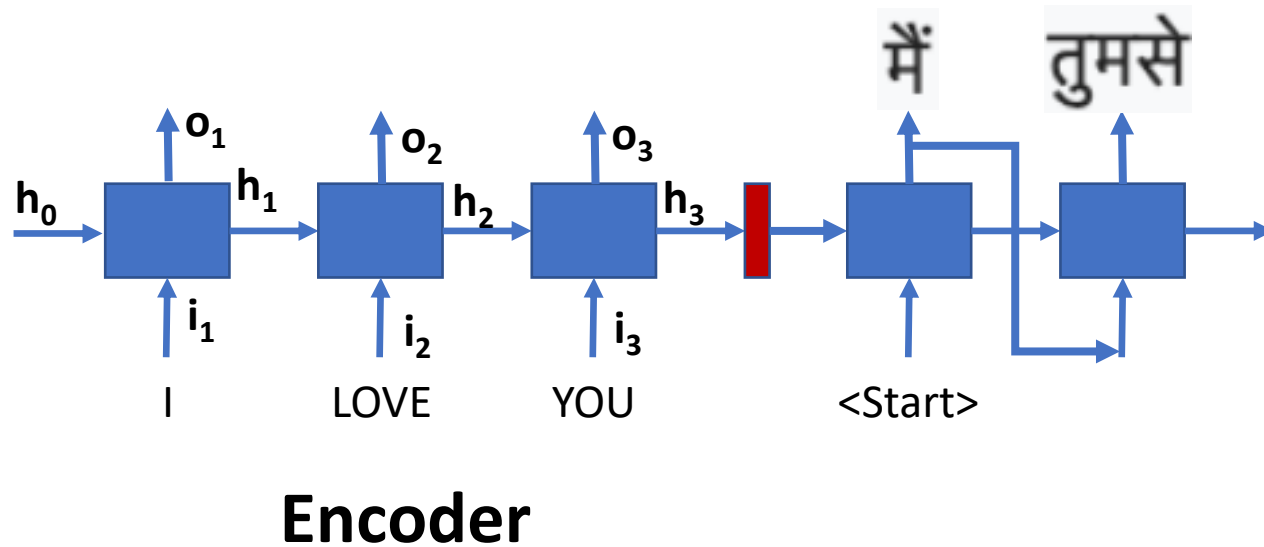
Sequence-to-Sequence (Seq2Seq)

i love you → मैं तुमसे प्यार करता हूँ
main tumase pyaar karata hoon



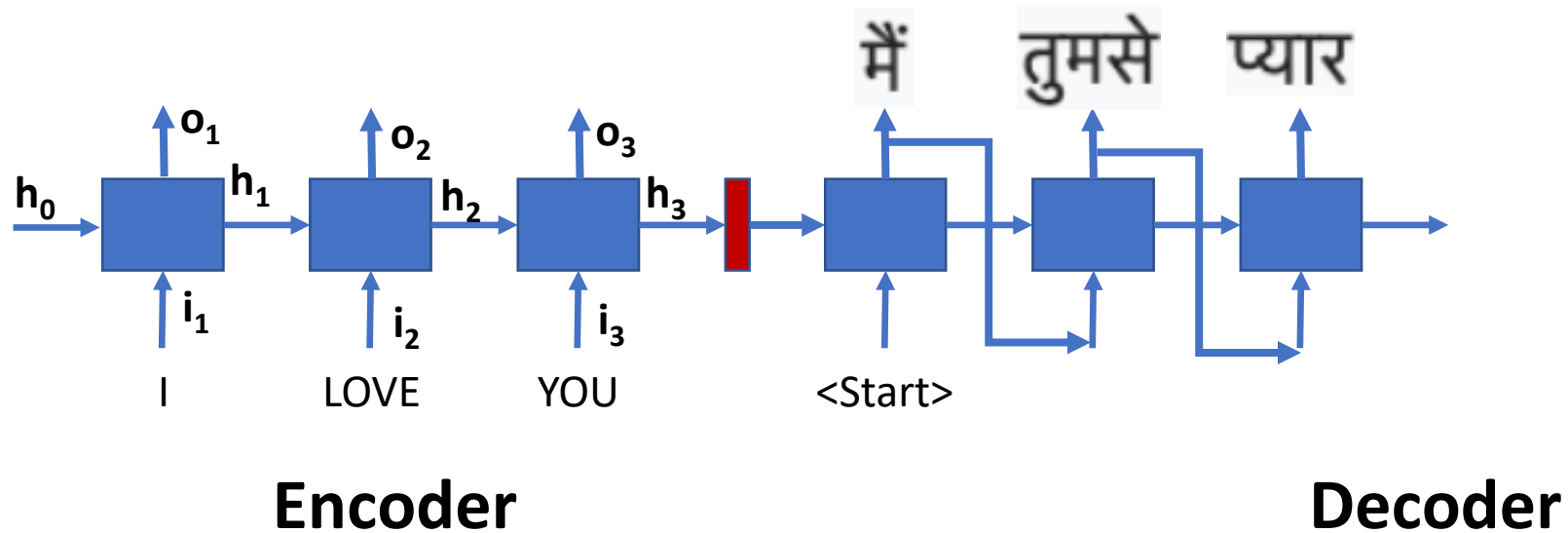
Sequence-to-Sequence (Seq2Seq)

i love you → मैं तुमसे प्यार करता हूँ
main tumase pyaar karata hoon



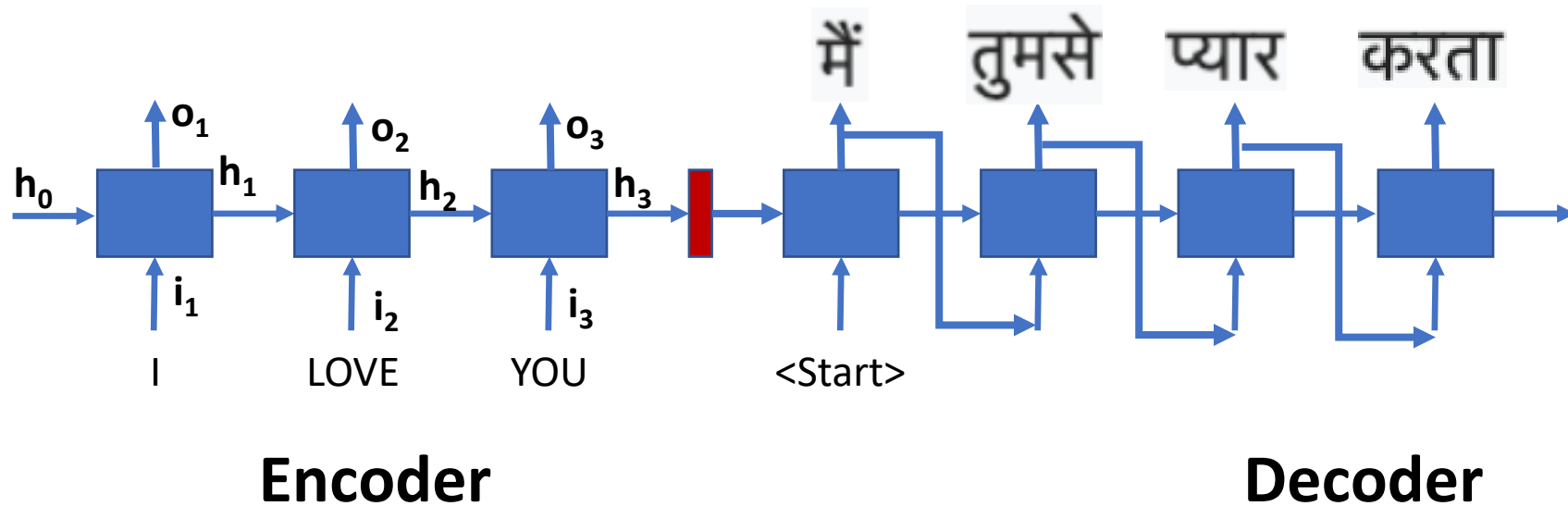
Sequence-to-Sequence (Seq2Seq)

i love you → मैं तुमसे प्यार करता हूँ
main tumase pyaar karata hoon



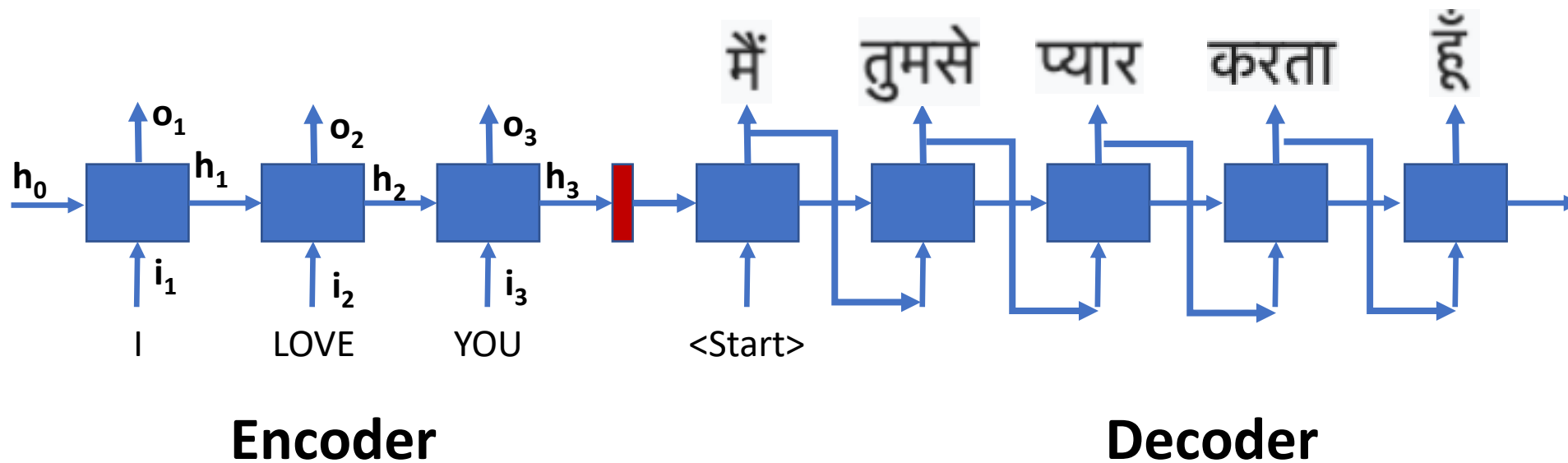
Sequence-to-Sequence (Seq2Seq)

i love you → मैं तुमसे प्यार करता हूँ
main tumase pyaar karata hoon



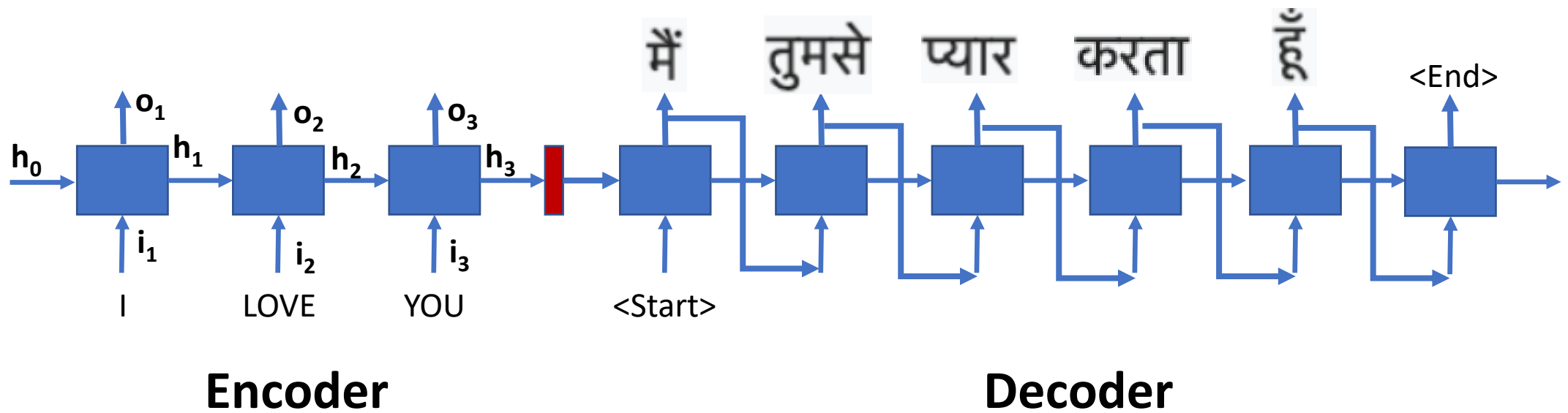
Sequence-to-Sequence (Seq2Seq)

i love you → मैं तुमसे प्यार करता हूँ
main tumase pyaar karata hoon



Sequence-to-Sequence (Seq2Seq)

i love you → मैं तुमसे प्यार करता हूँ
main tumase pyaar karata hoon



Attention

NEURAL MACHINE TRANSLATION BY JOINTLY LEARNING TO ALIGN AND TRANSLATE

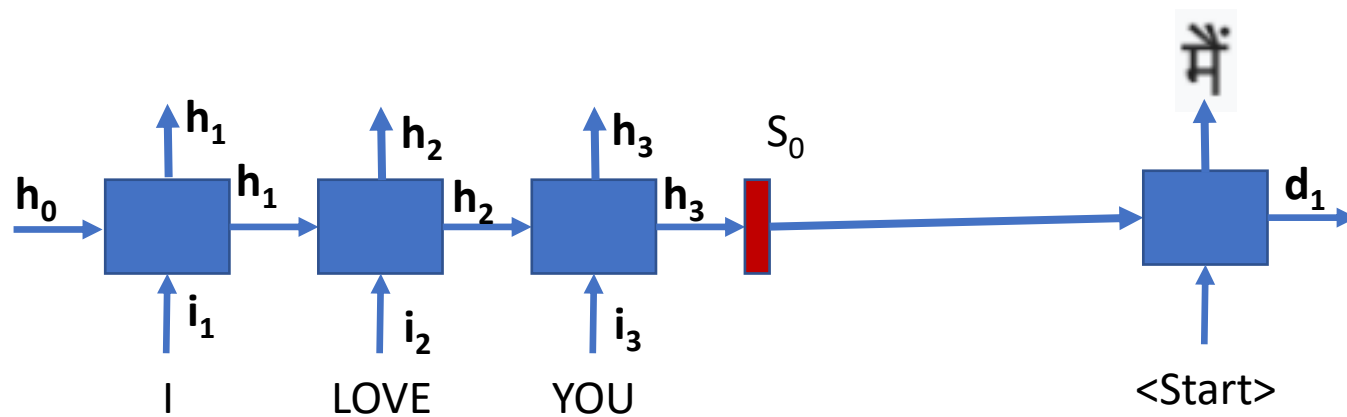
Dzmitry Bahdanau
Jacobs University Bremen, Germany

KyungHyun Cho **Yoshua Bengio***
Université de Montréal

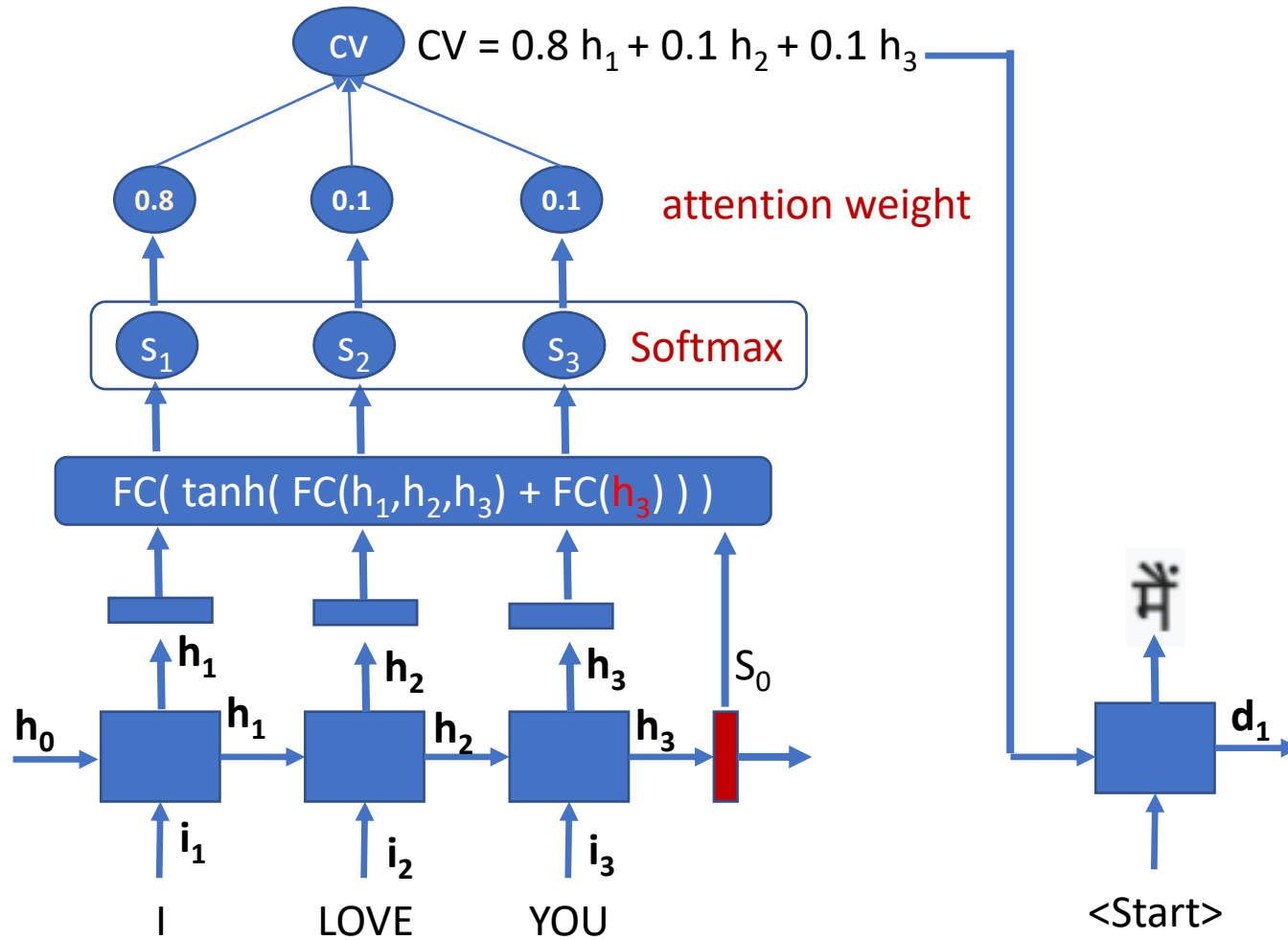
ABSTRACT

Neural machine translation is a recently proposed approach to machine translation. Unlike the traditional statistical machine translation, the neural machine translation aims at building a single neural network that can be jointly tuned to maximize the translation performance. The models proposed recently for neural machine translation often belong to a family of encoder–decoders and encode a source sentence into a fixed-length vector from which a decoder generates a translation. In this paper, we conjecture that the use of a fixed-length vector is a bottleneck in improving the performance of this basic encoder–decoder architecture, and propose to extend this by allowing a model to automatically (soft-)search for parts of a source sentence that are relevant to predicting a target word, without having to form these parts as a hard segment explicitly. With this new approach, we achieve a translation performance comparable to the existing state-of-the-art phrase-based system on the task of English-to-French translation. Furthermore, qualitative analysis reveals that the (soft-)alignments found by the model agree well with our intuition.

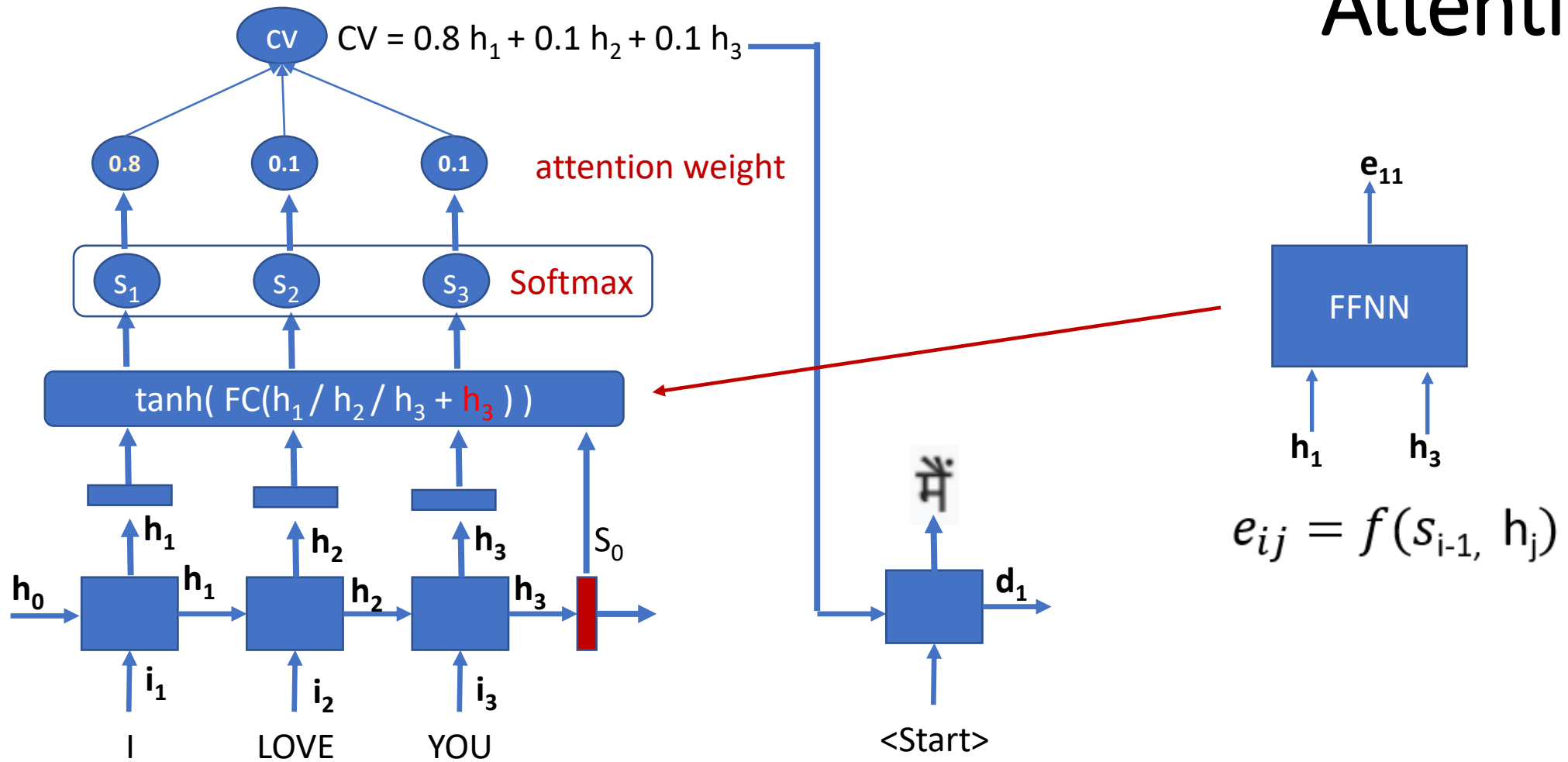
Attention



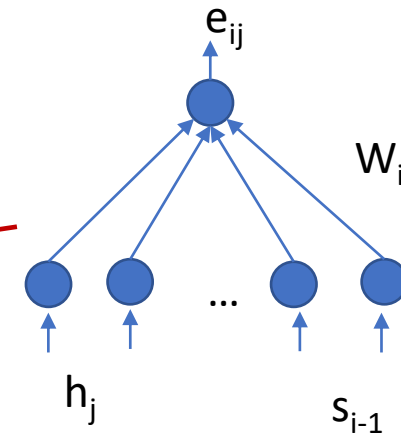
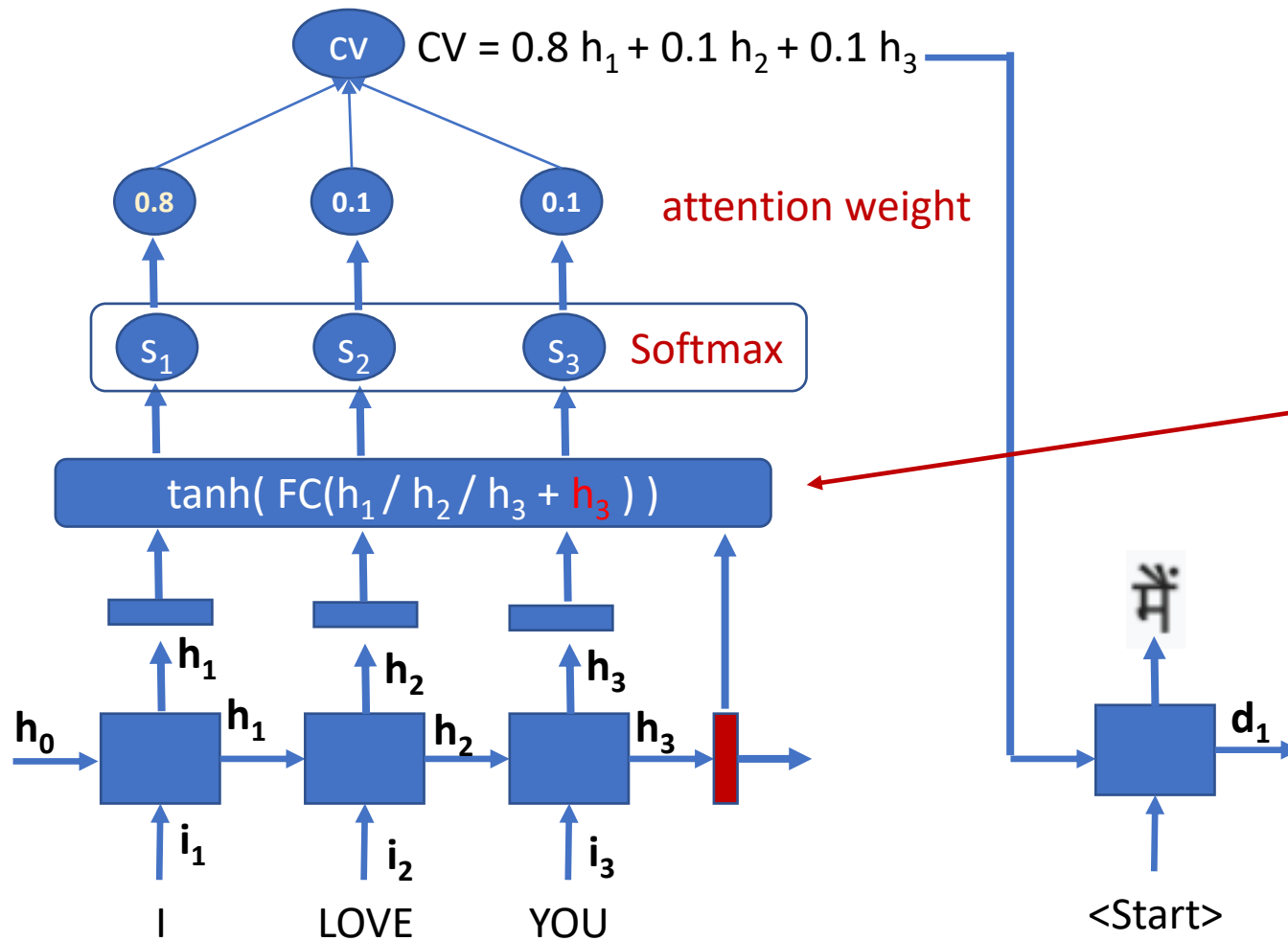
Attention



Attention



Attention

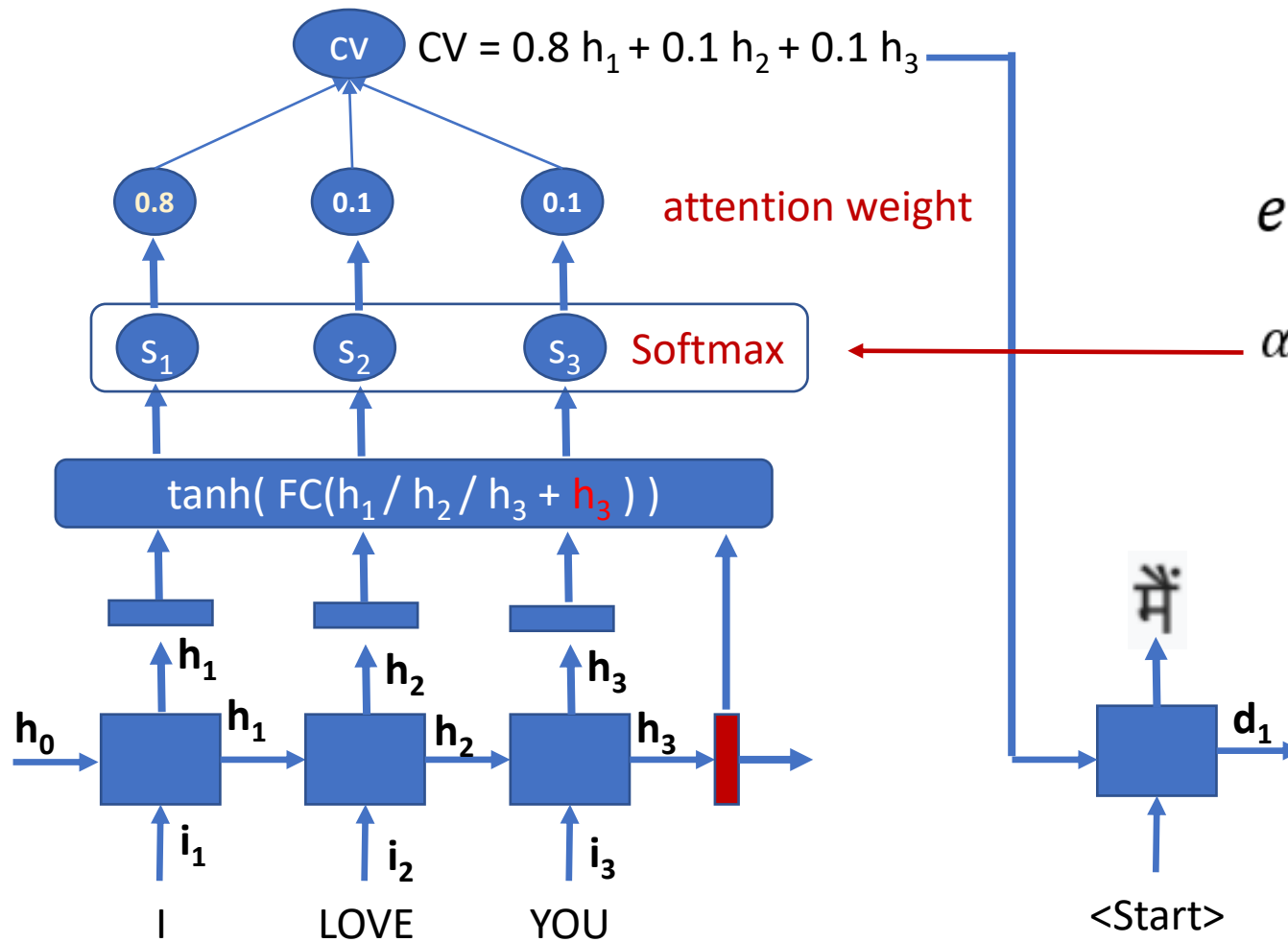


$$e_{ij} = f(s_{i-1}, h_j)$$

$$e_{ij} = \tanh(W^T x)$$

x = concatenation of s_{i-1} and h_j

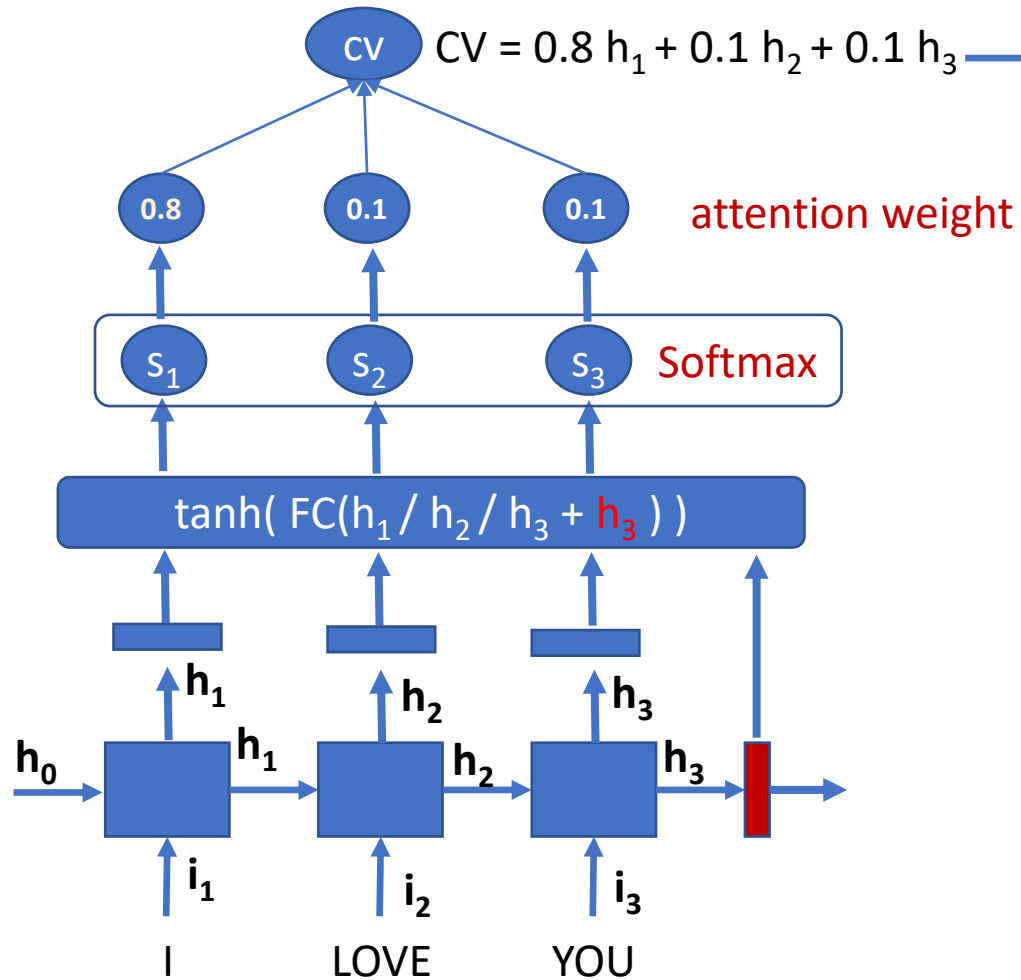
Attention



$$e_{ij} = f(s_{i-1}, h_j)$$

$$\alpha_{ij} = \text{softmax}(e_{ij}) = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})}$$

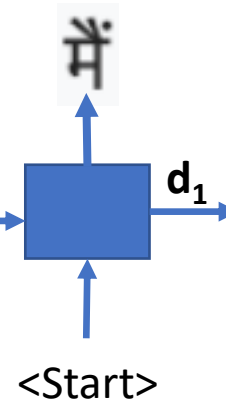
Attention



$$e_{ij} = f(s_{i-1}, h_j)$$

$$\alpha_{ij} = \text{softmax}(e_{ij}) = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})}$$

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j$$



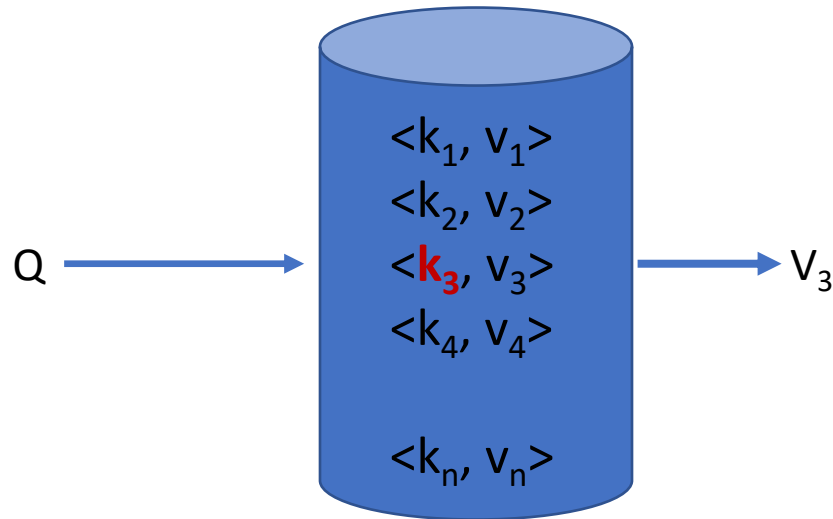
Attention – general concept

Q : Query

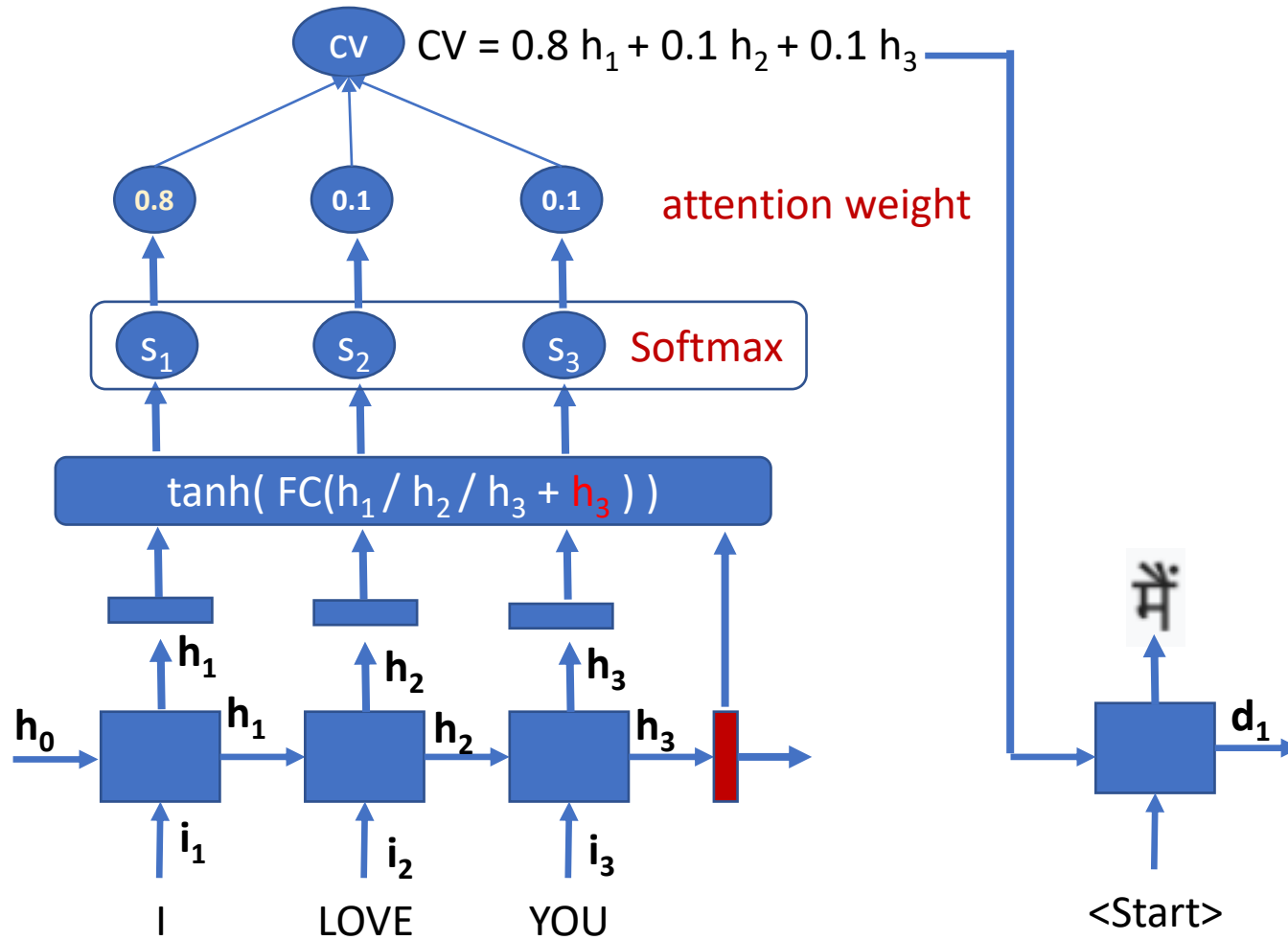
K : Key

V : Value

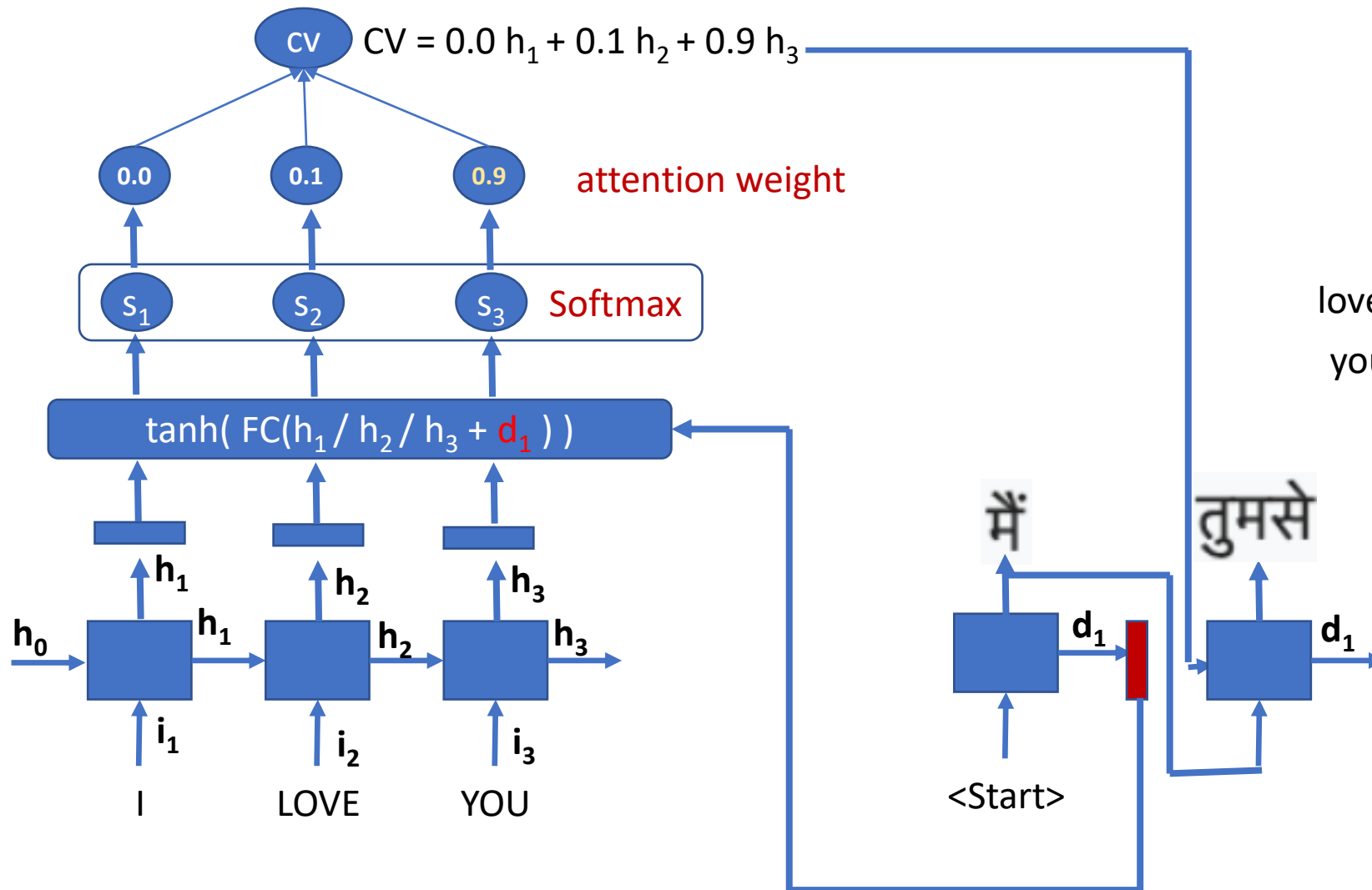
$$A(q, K, V) = \sum_i \frac{\exp(e_{qk_i})}{\sum_j \exp(e_{qk_j})} v_i$$



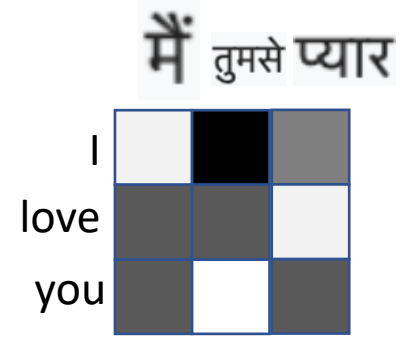
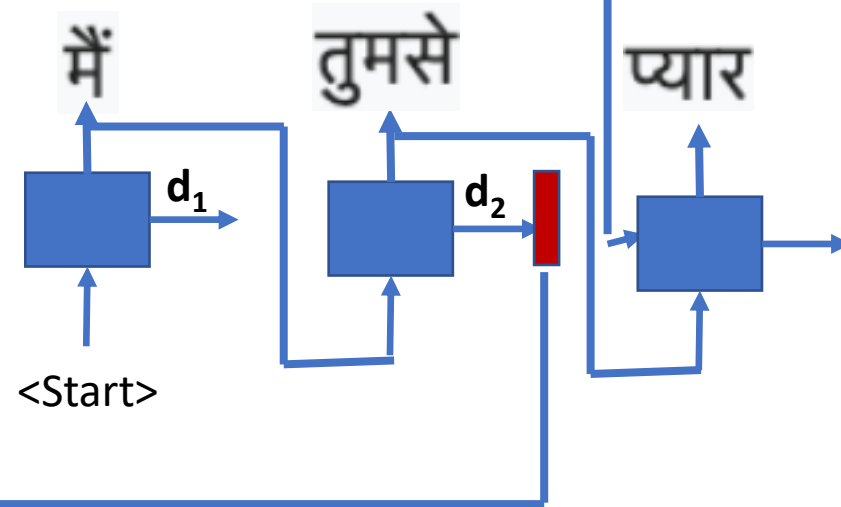
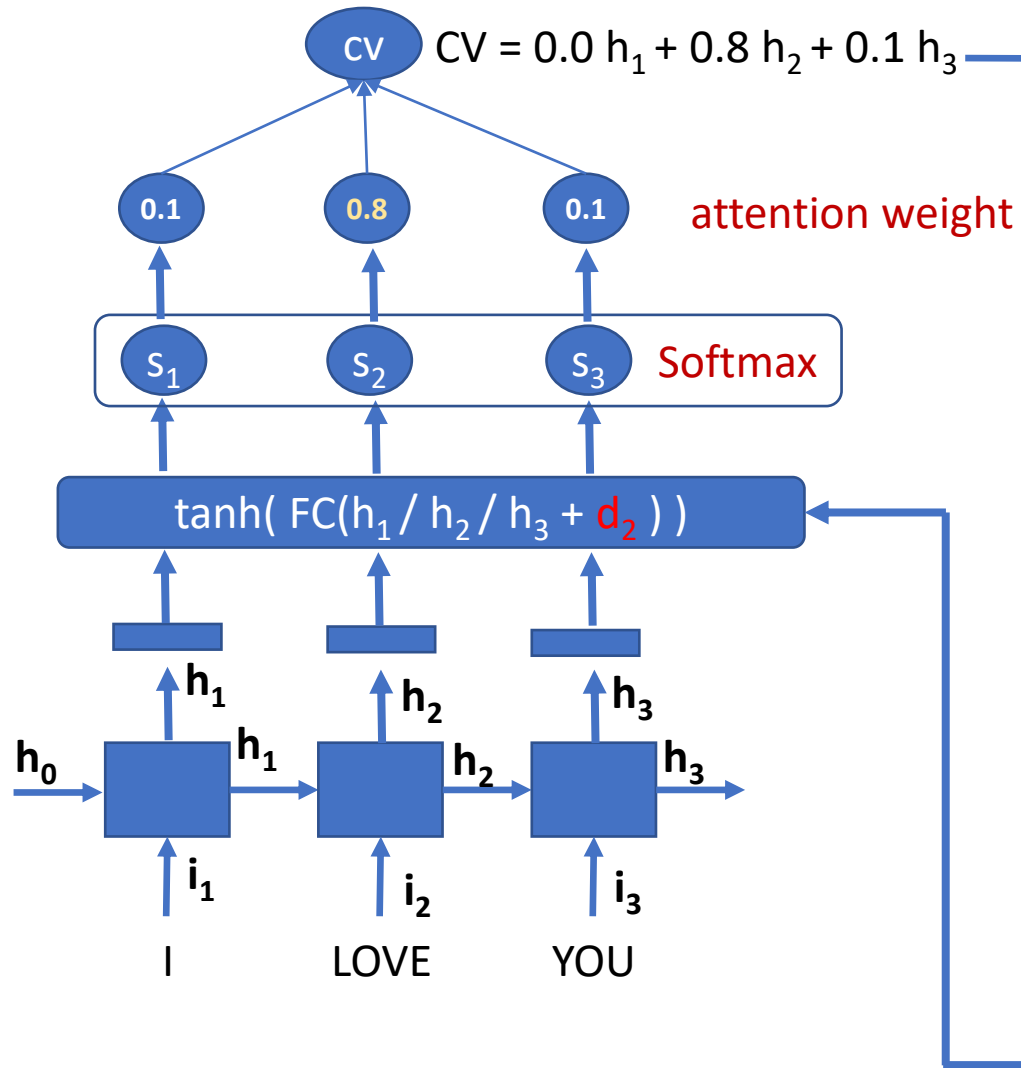
Attention



Attention



Attention



Attention Types

Computing $e_1, \dots, e_n \in \mathbb{R}$ from $k_1, \dots, k_n \in \mathbb{R}^{d_1}$ and $q \in \mathbb{R}^{d_2}$

- Basic dot-product attention:

$$e_i = q^T k_i \in \mathbb{R}$$

where $d_1 = d_2$

- Multiplicative attention:

$$e_i = q^T W k_i \in \mathbb{R}$$

where $W \in \mathbb{R}^{d_2 \times d_1}$ is a weight matrix

- Additive attention:

$$e_i = W_3^T \tanh(W_1 k_i + W_2 q) \in \mathbb{R}$$

where $W_1 \in \mathbb{R}^{d_3 \times d_1}$, $W_2 \in \mathbb{R}^{d_3 \times d_2}$ are weight matrices and $W_3 \in \mathbb{R}^{d_3}$ is a weight vector.

Transformer

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukaszkaizer@google.com

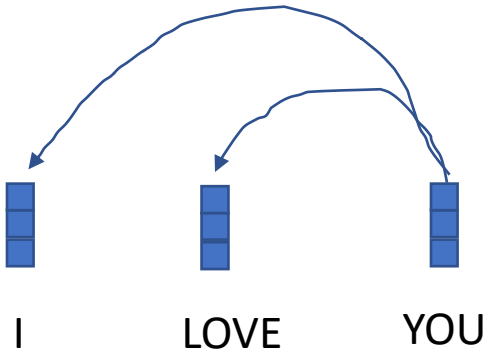
Illia Polosukhin* ‡
illia.polosukhin@gmail.com

Abstract

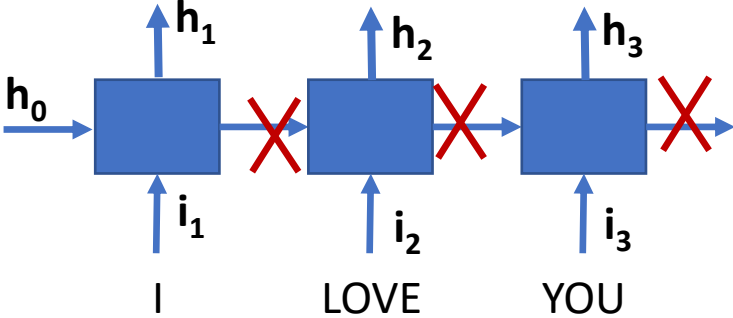
The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.0 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature.

1 Introduction

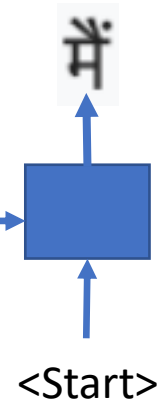
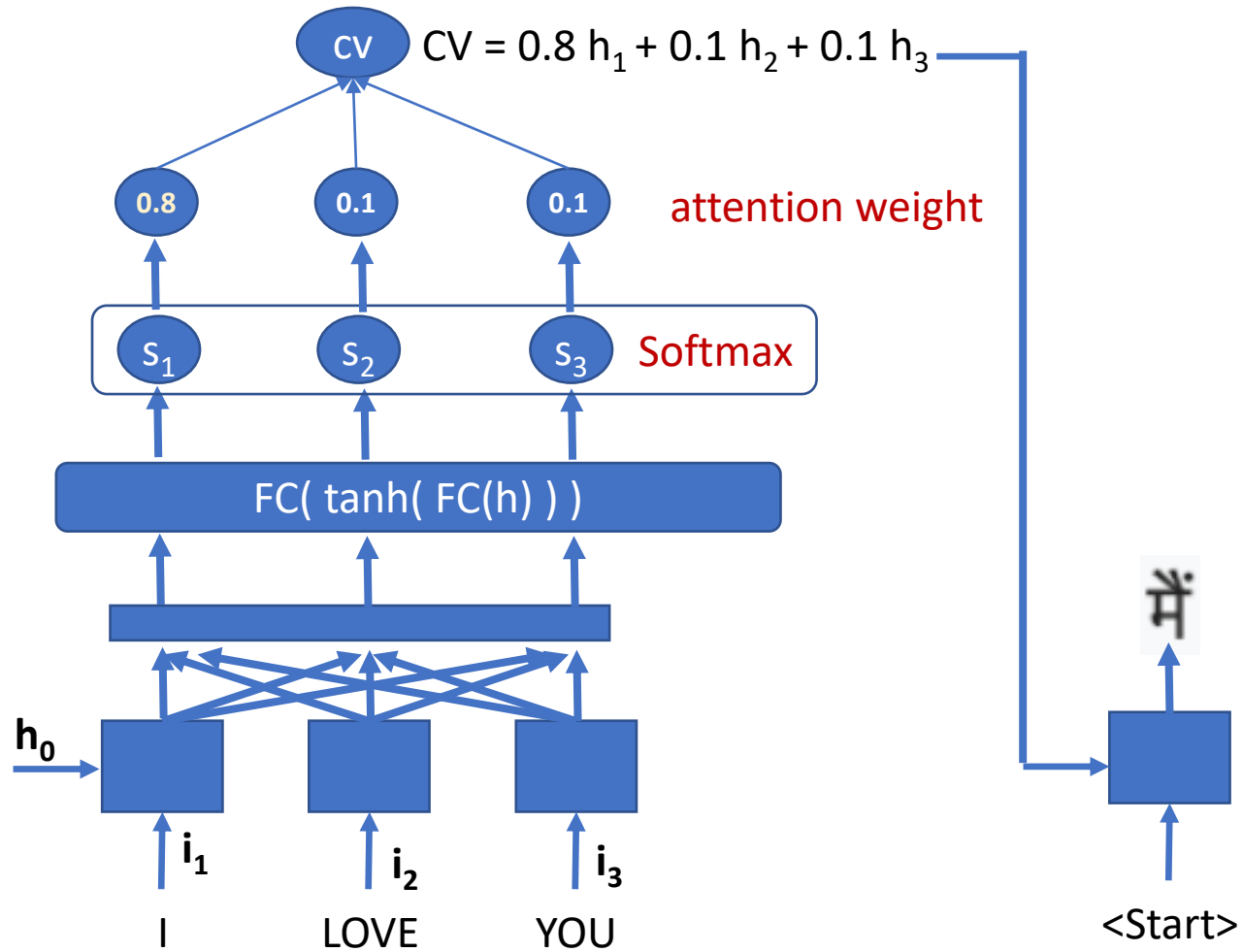
Self Attention



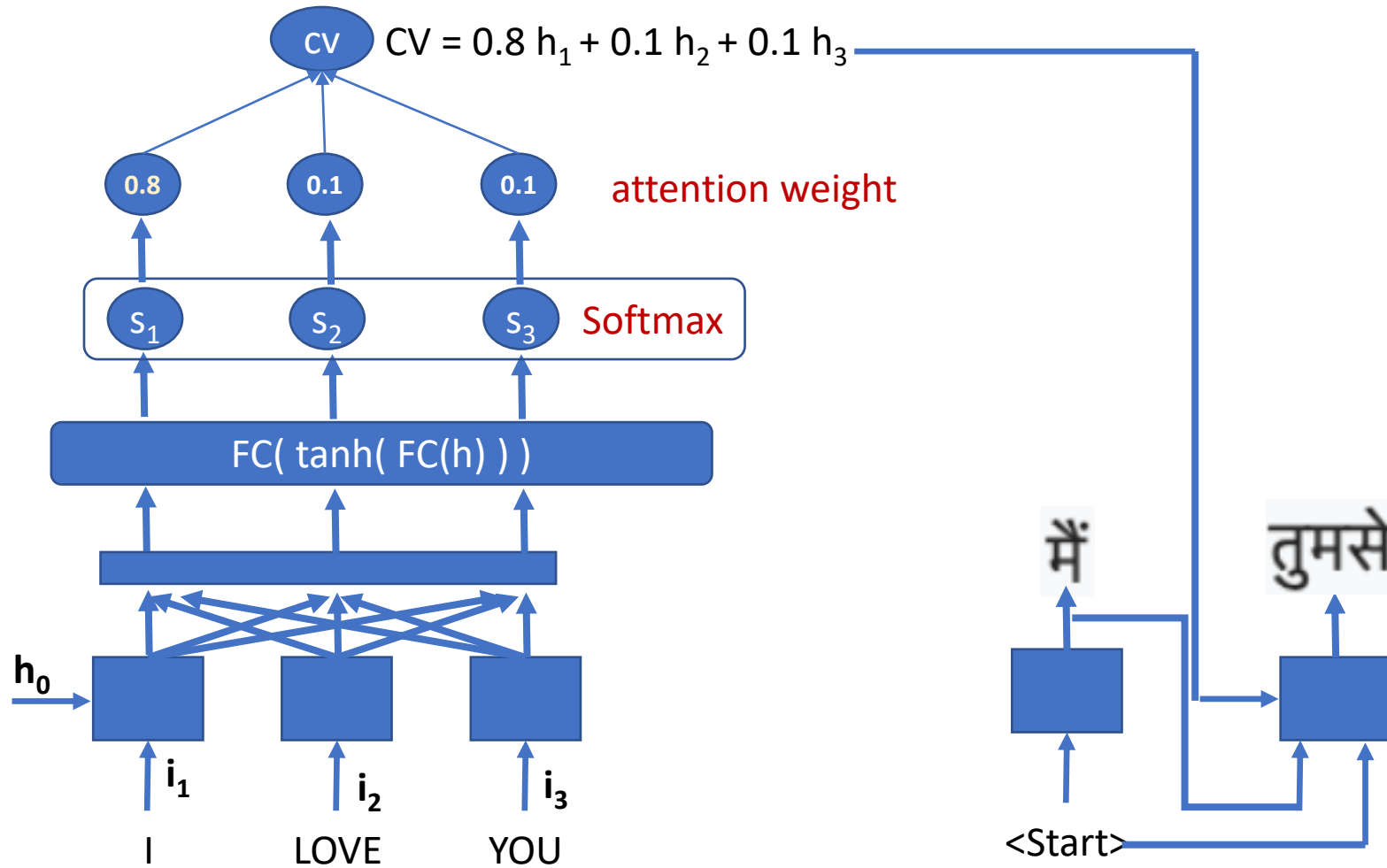
Transformer



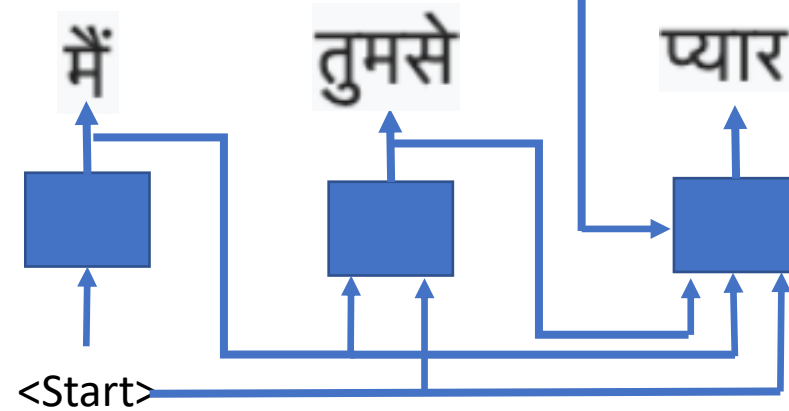
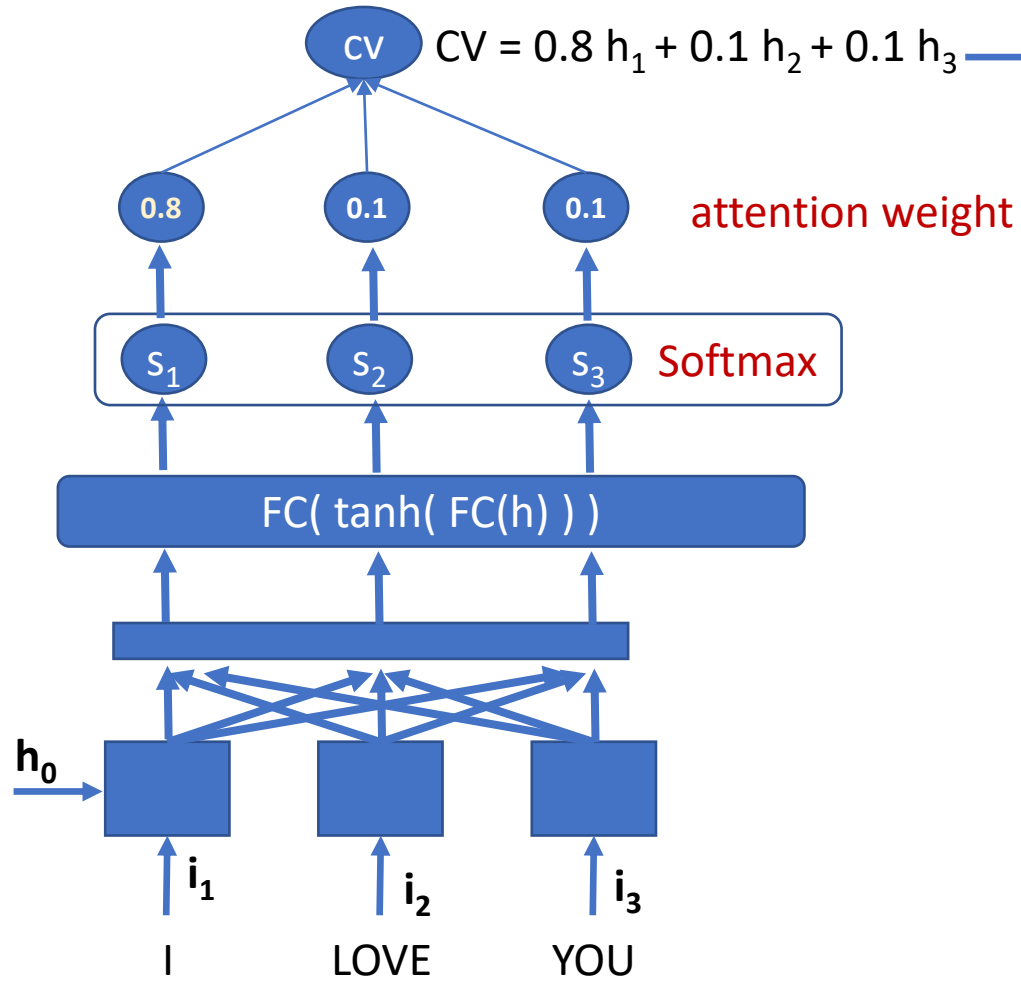
Attention

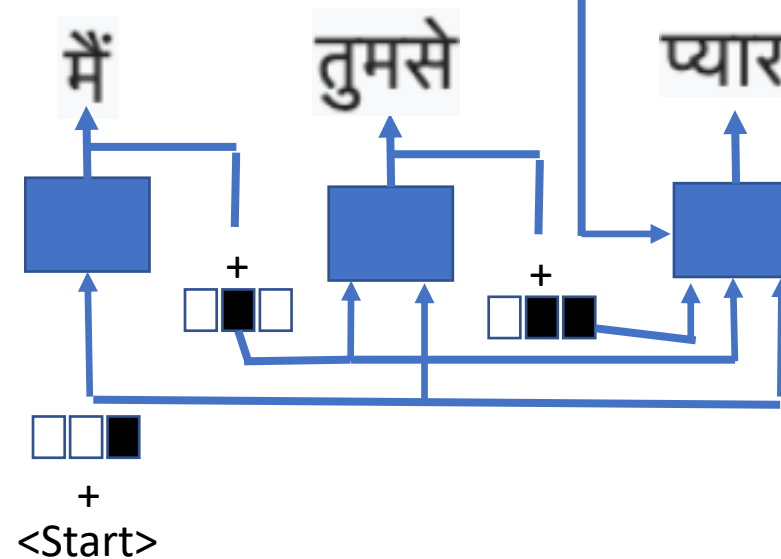
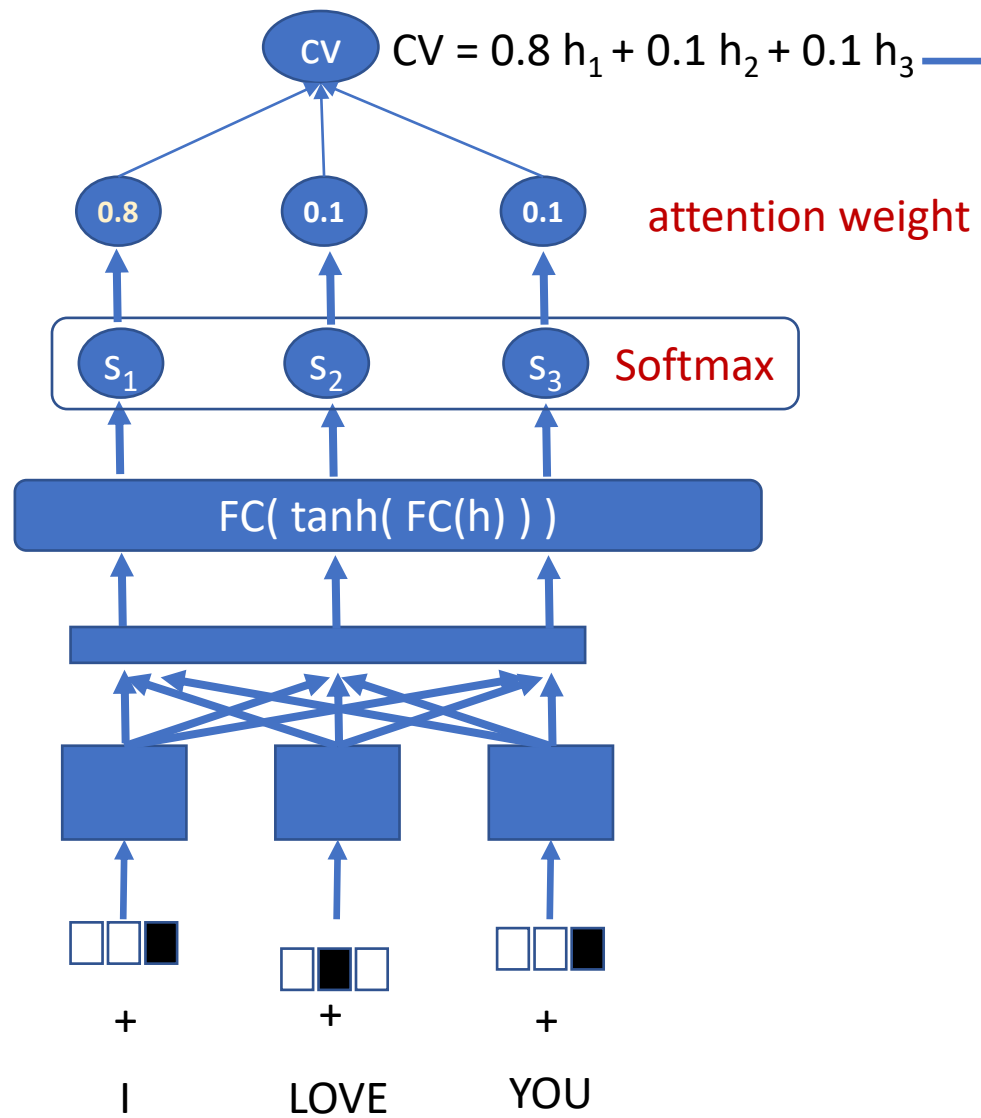


Attention

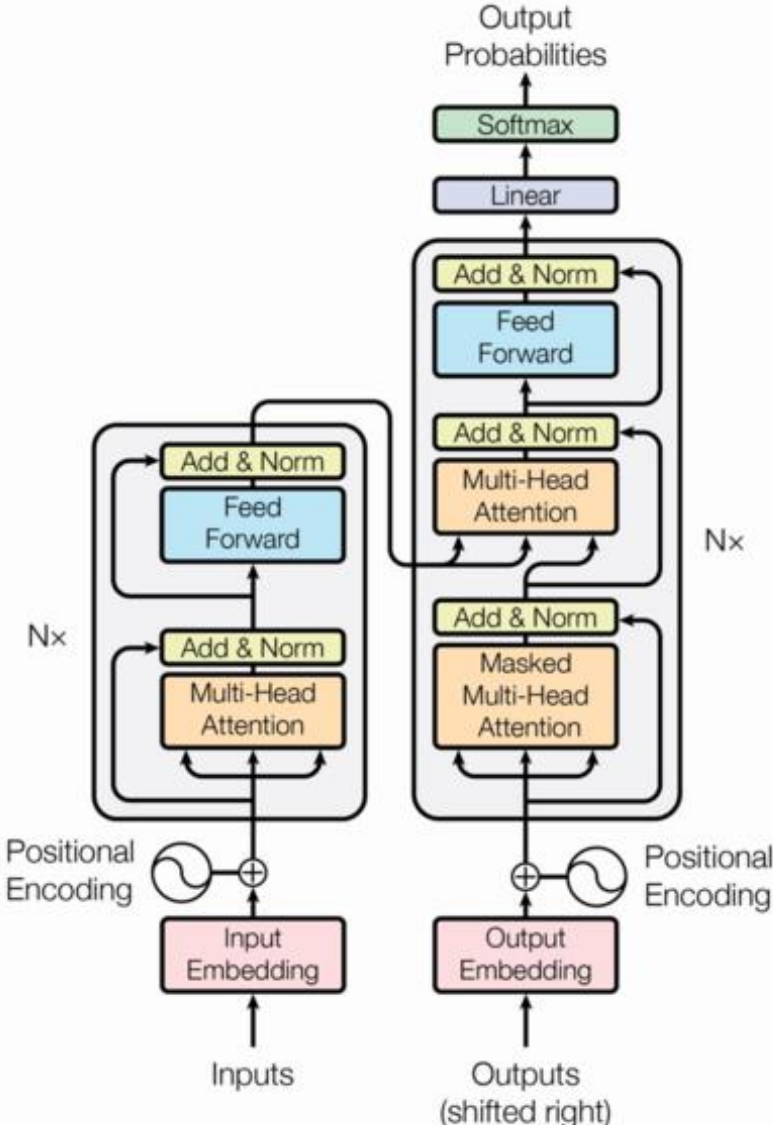


Attention

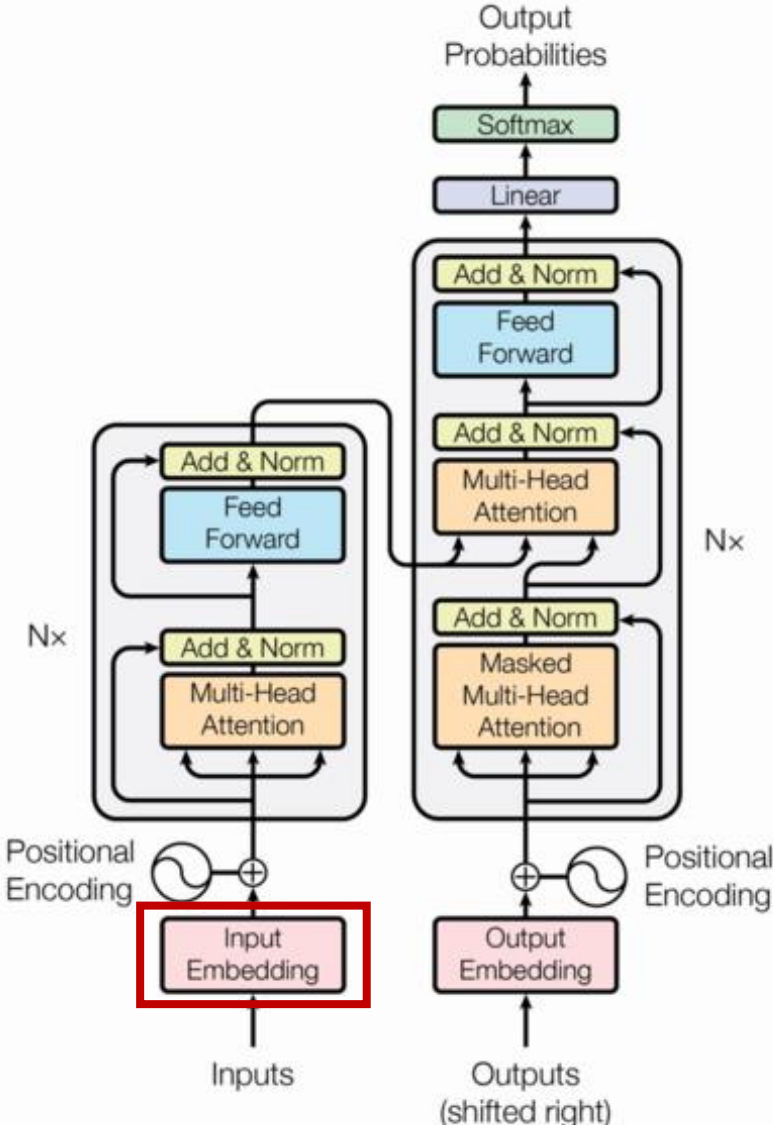




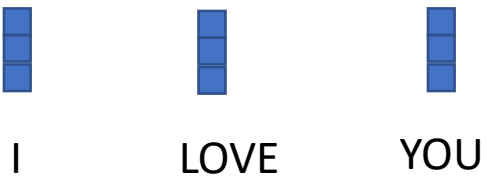
Transformer: Attention Is All You Need



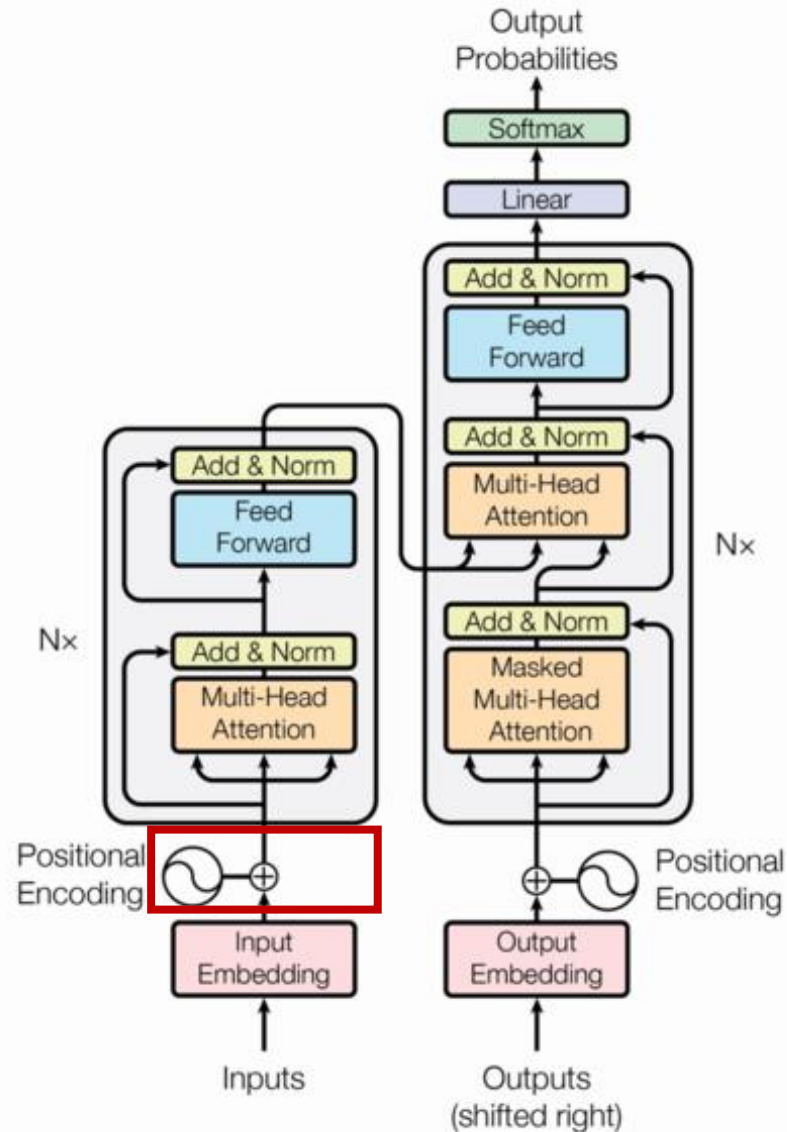
Transformer: Attention Is All You Need



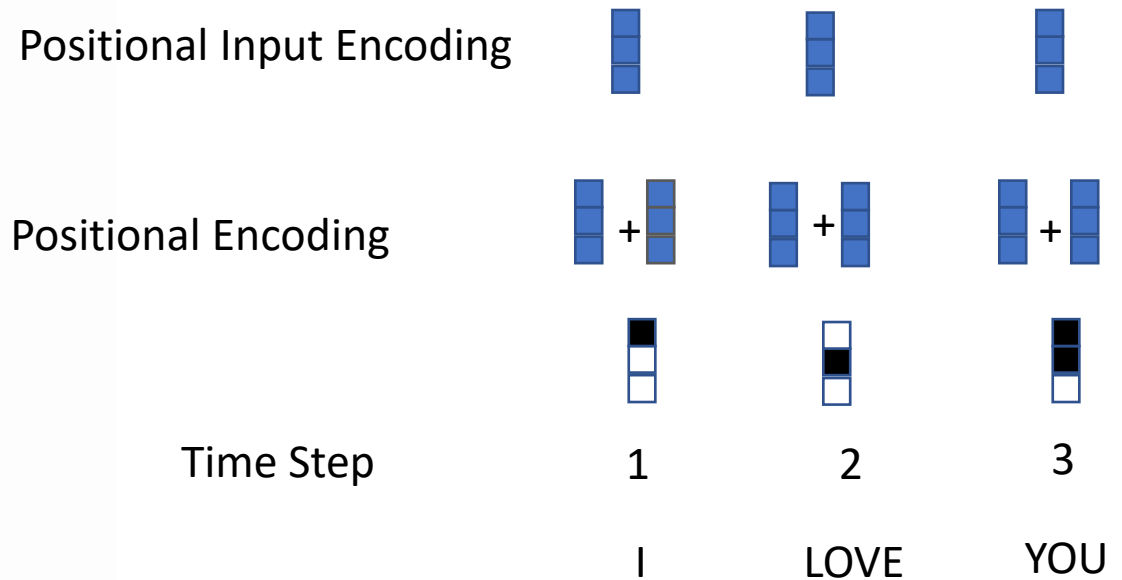
Word Embedding layer



Transformer: Attention Is All You Need



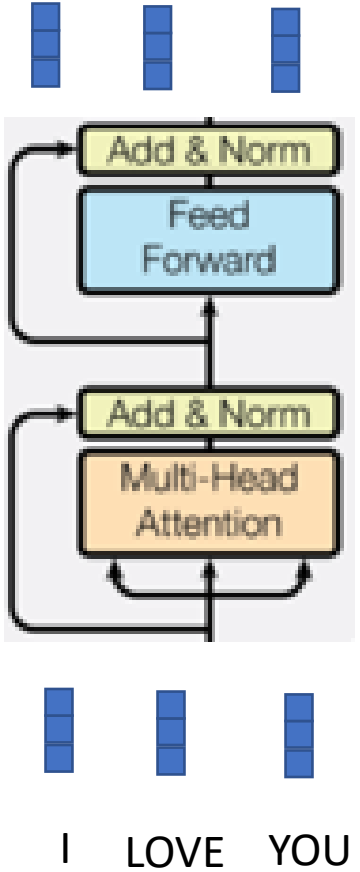
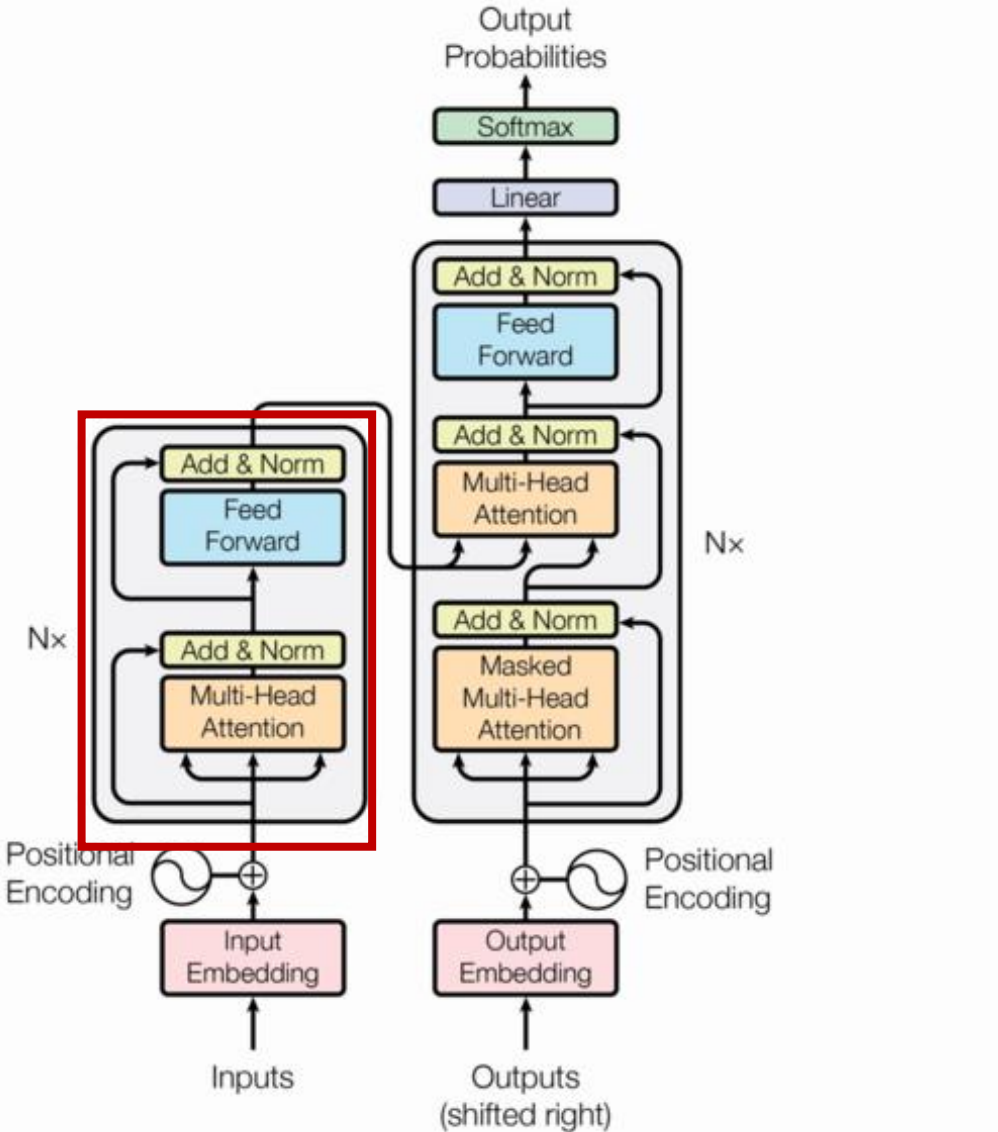
Positional Encoding



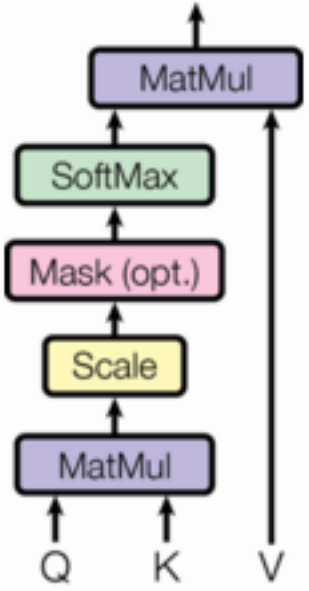
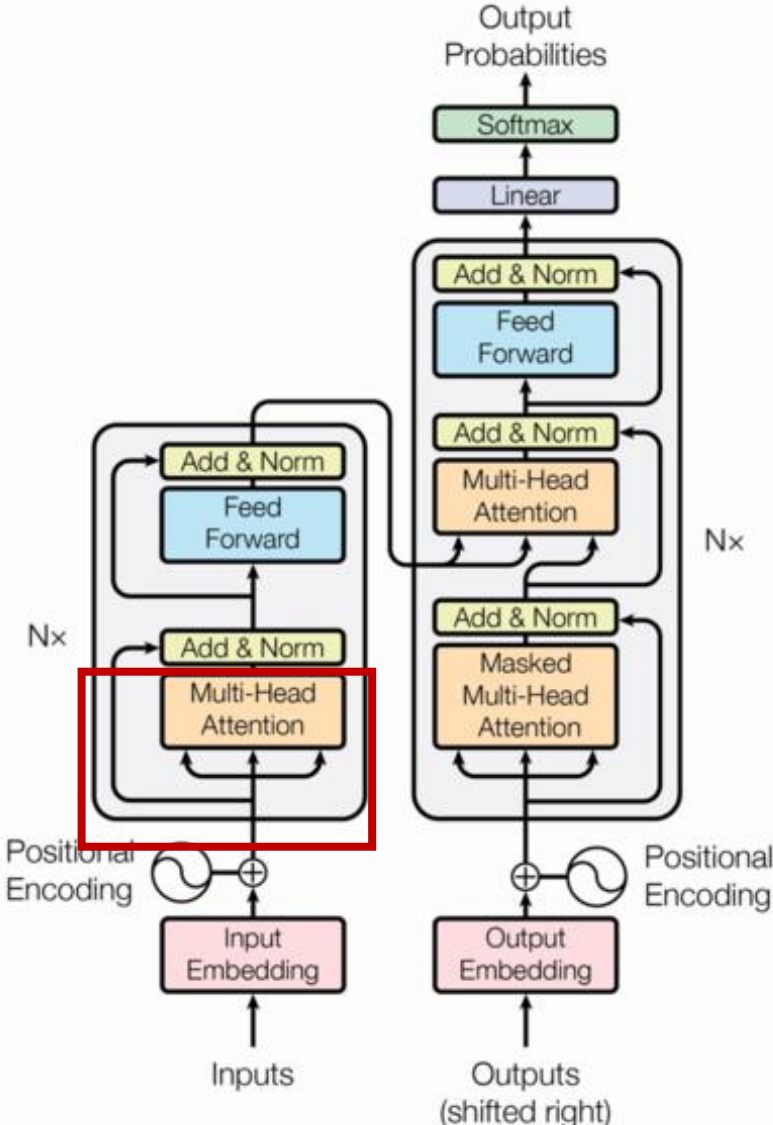
$$PE(pos, 2i + 1) = \cos\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$

$$PE(pos, 2i) = \sin\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$

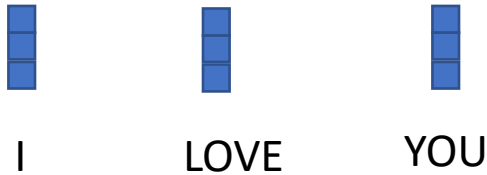
Transformer: Attention Is All You Need



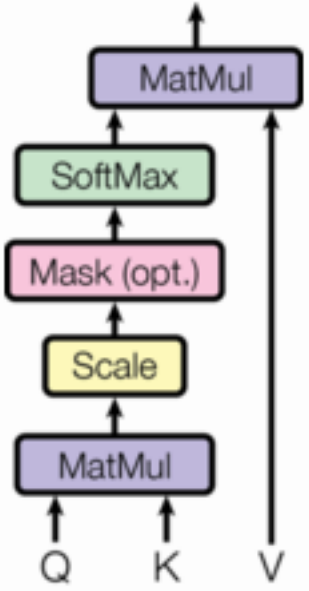
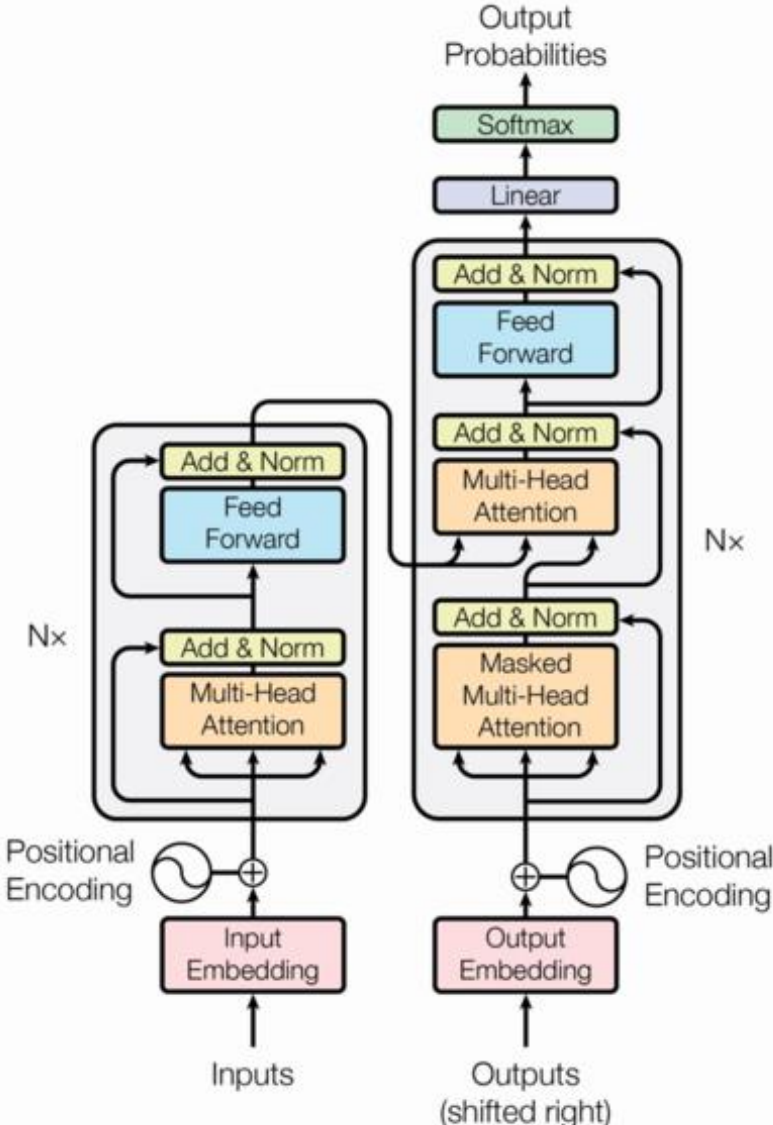
Transformer: Attention Is All You Need



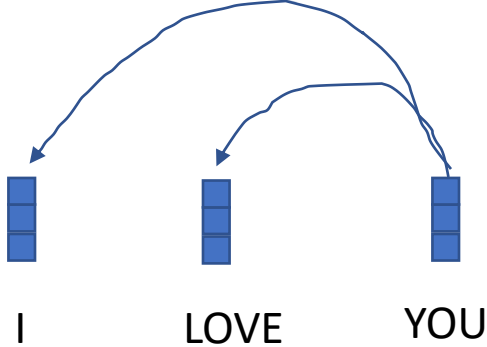
Self Attention



Transformer: Attention Is All You Need

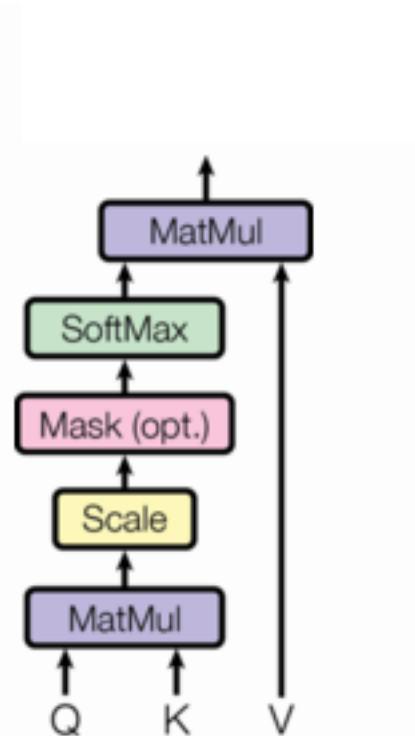
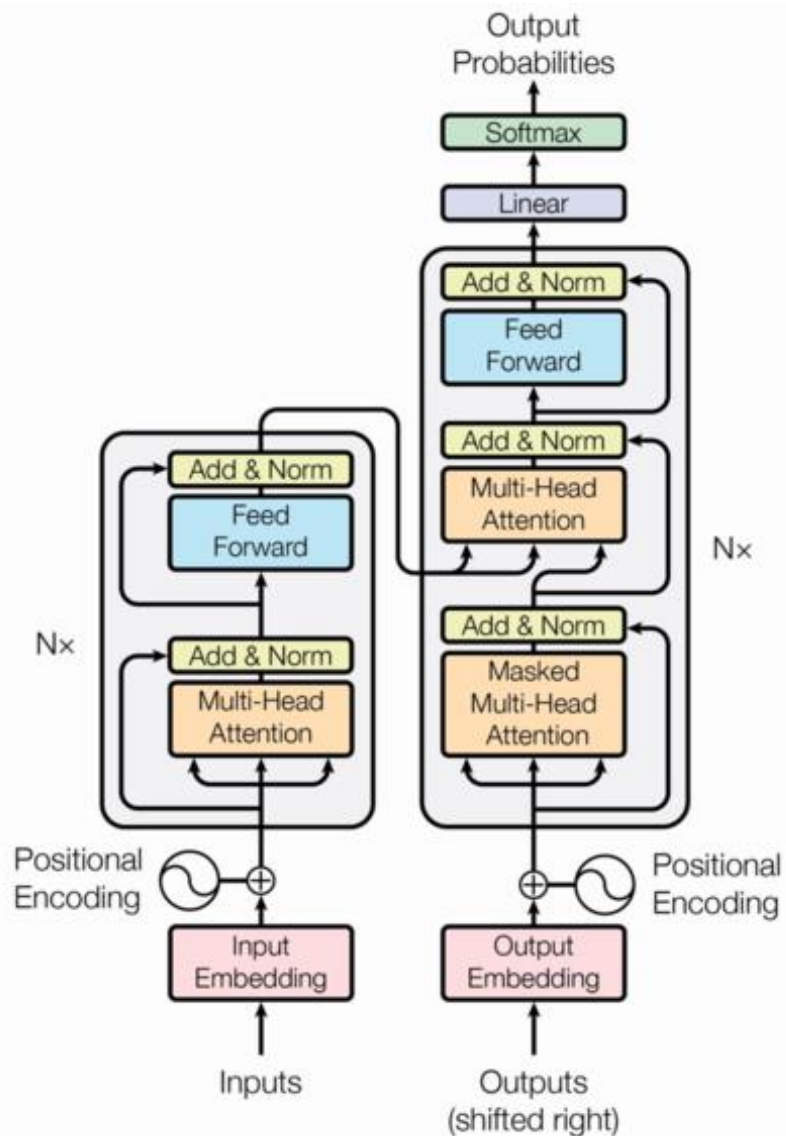


Self Attention

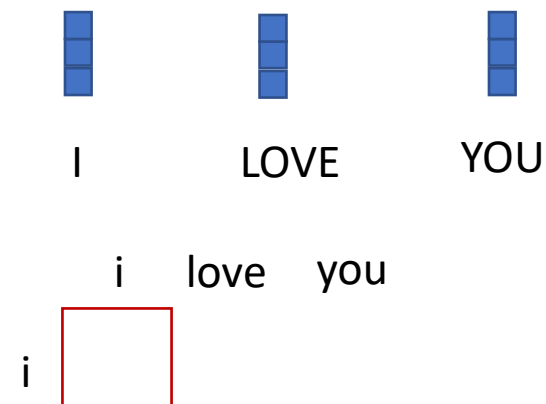


Softmax(Q . K)

Transformer: Attention Is All You Need

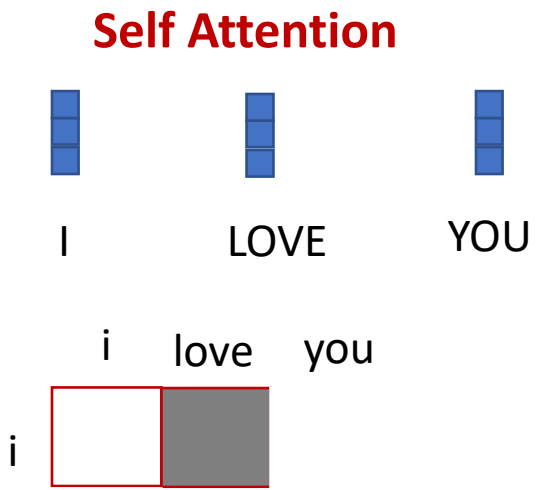
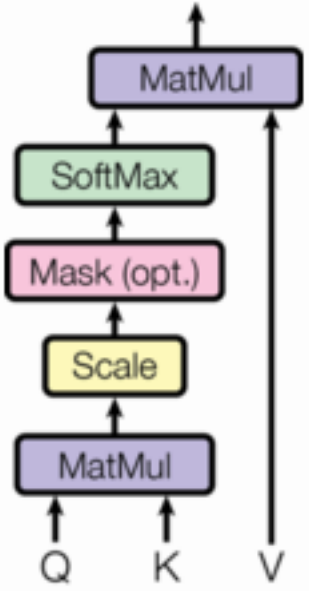
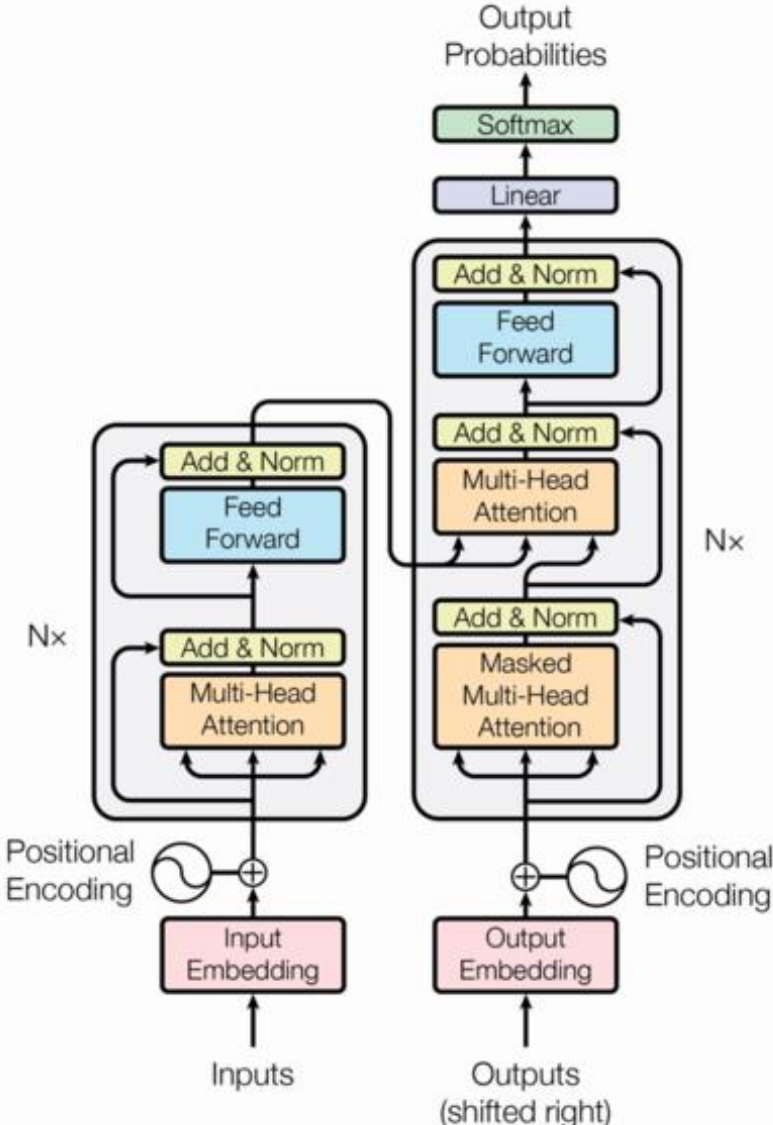


Self Attention



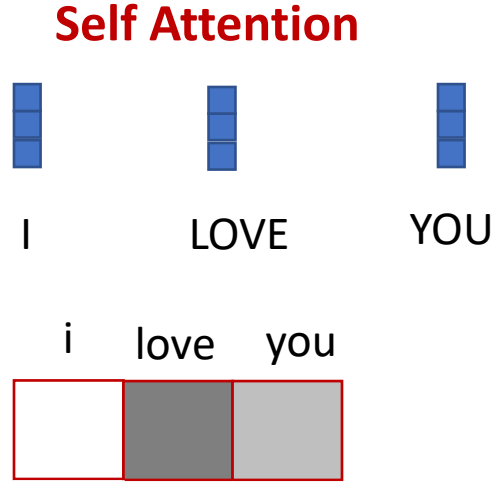
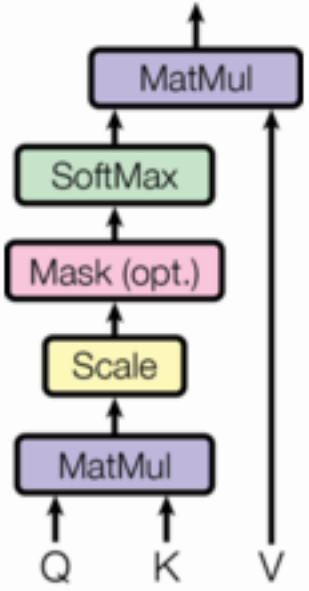
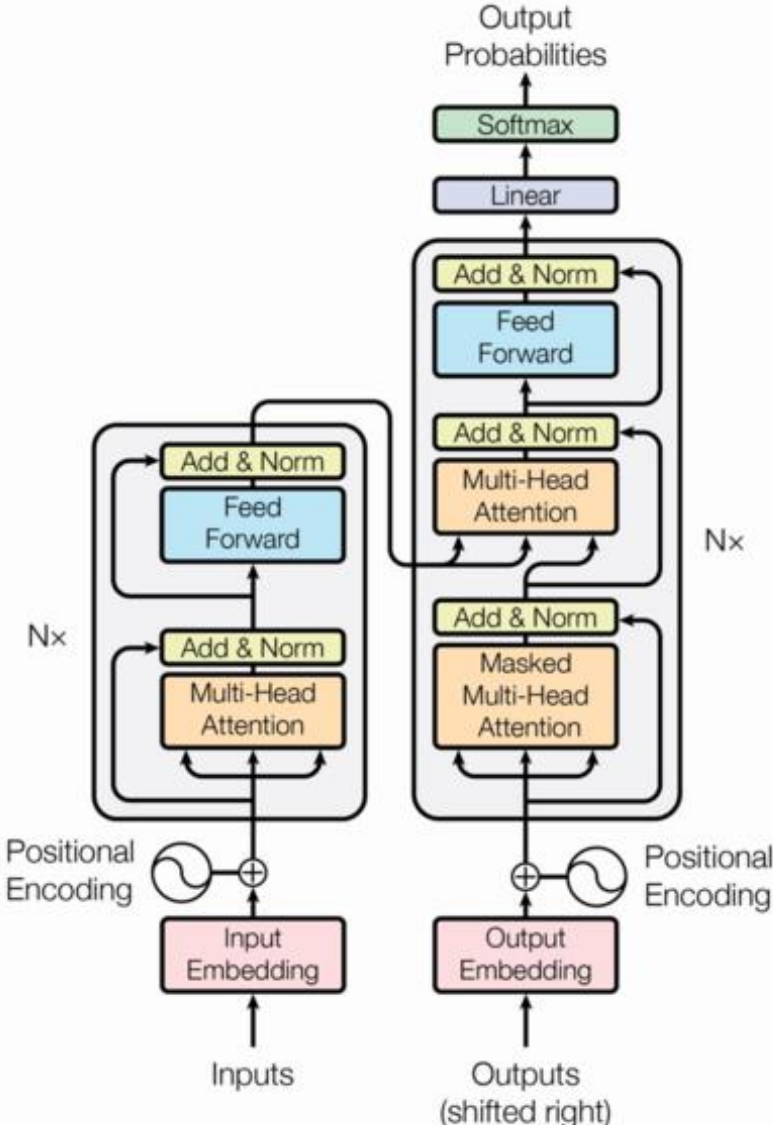
Softmax(I . I)

Transformer: Attention Is All You Need



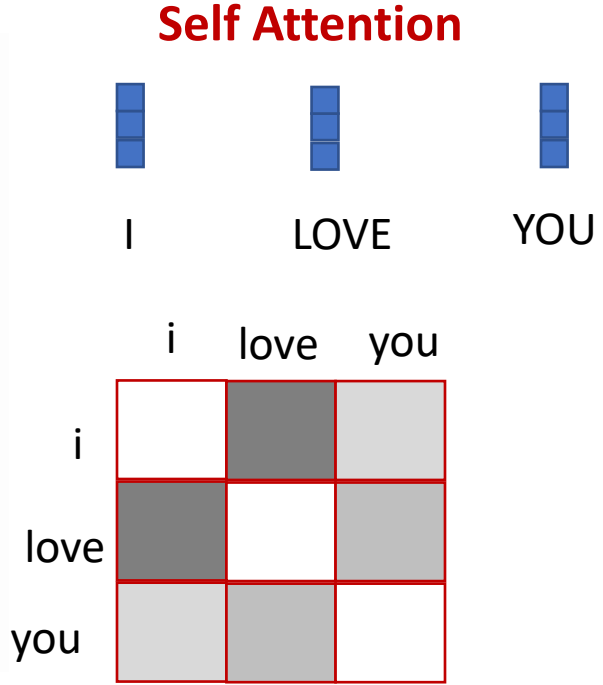
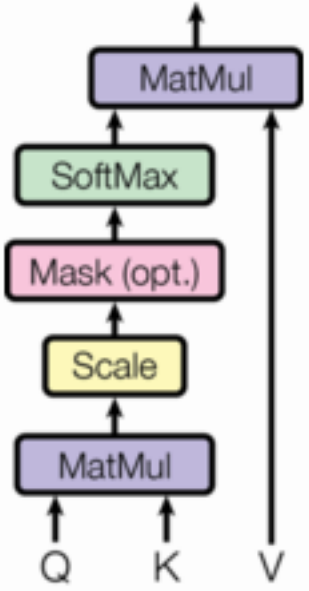
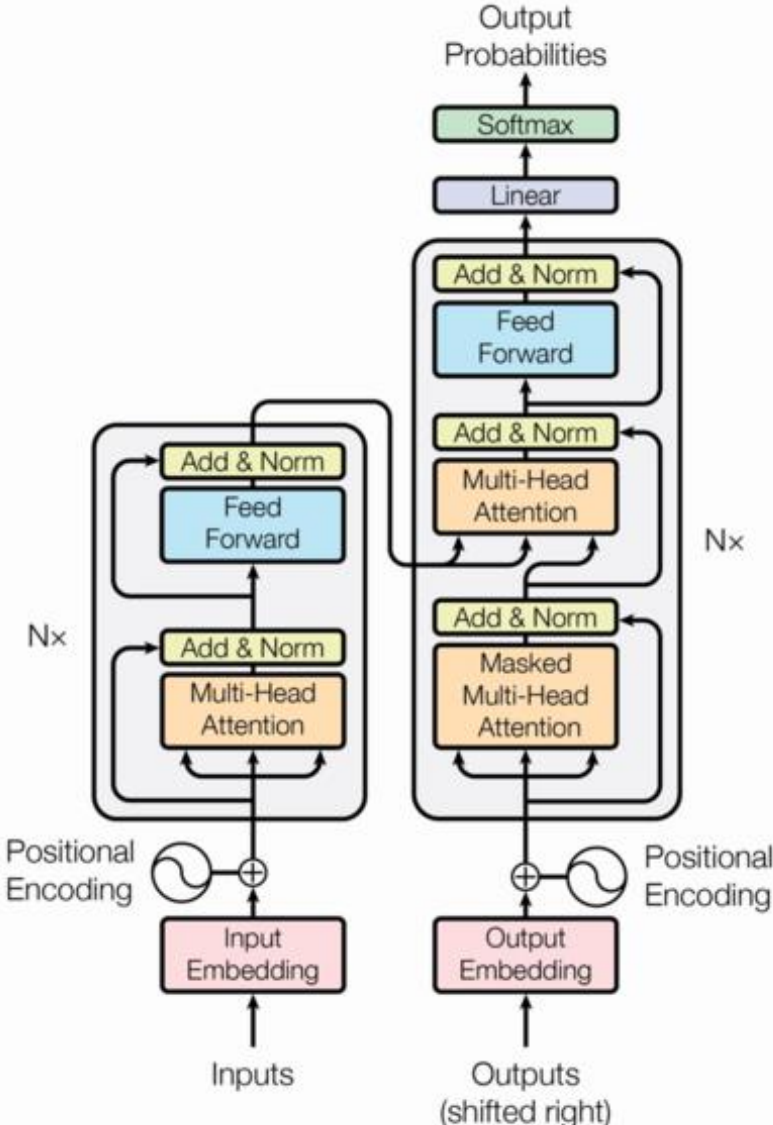
Softmax(I . LOVE)

Transformer: Attention Is All You Need

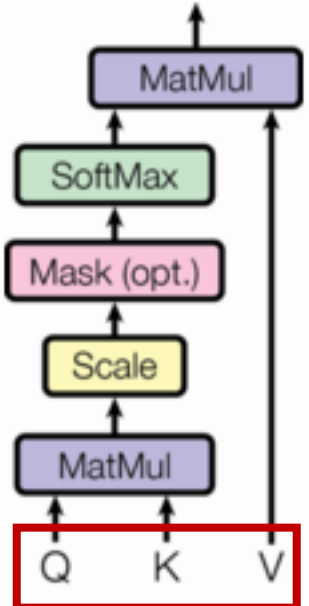
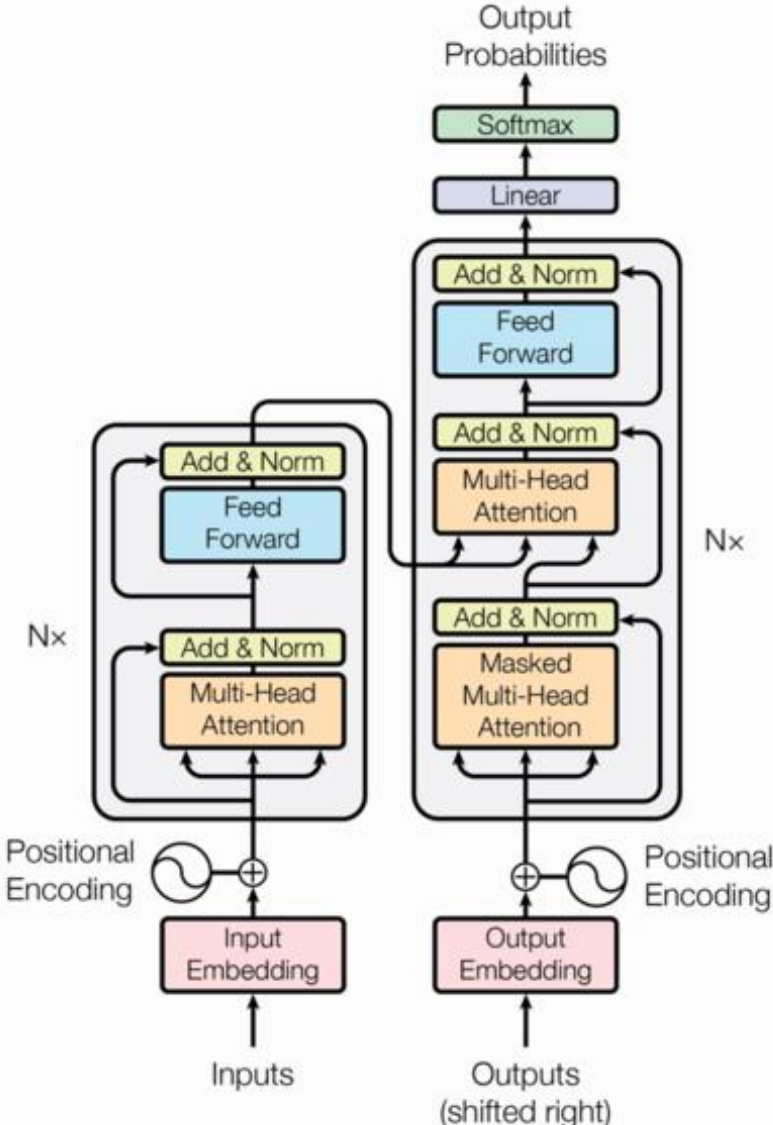


Softmax(I . YOU)

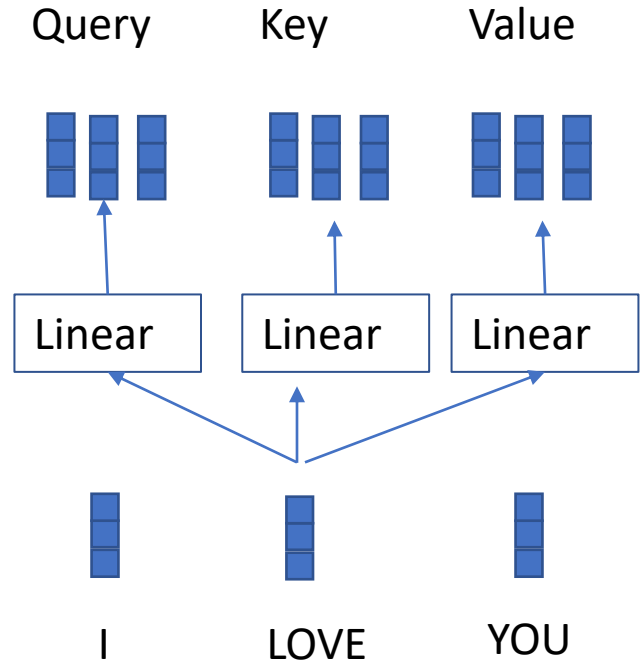
Transformer: Attention Is All You Need



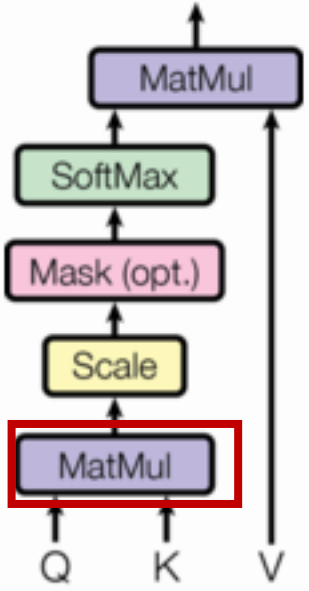
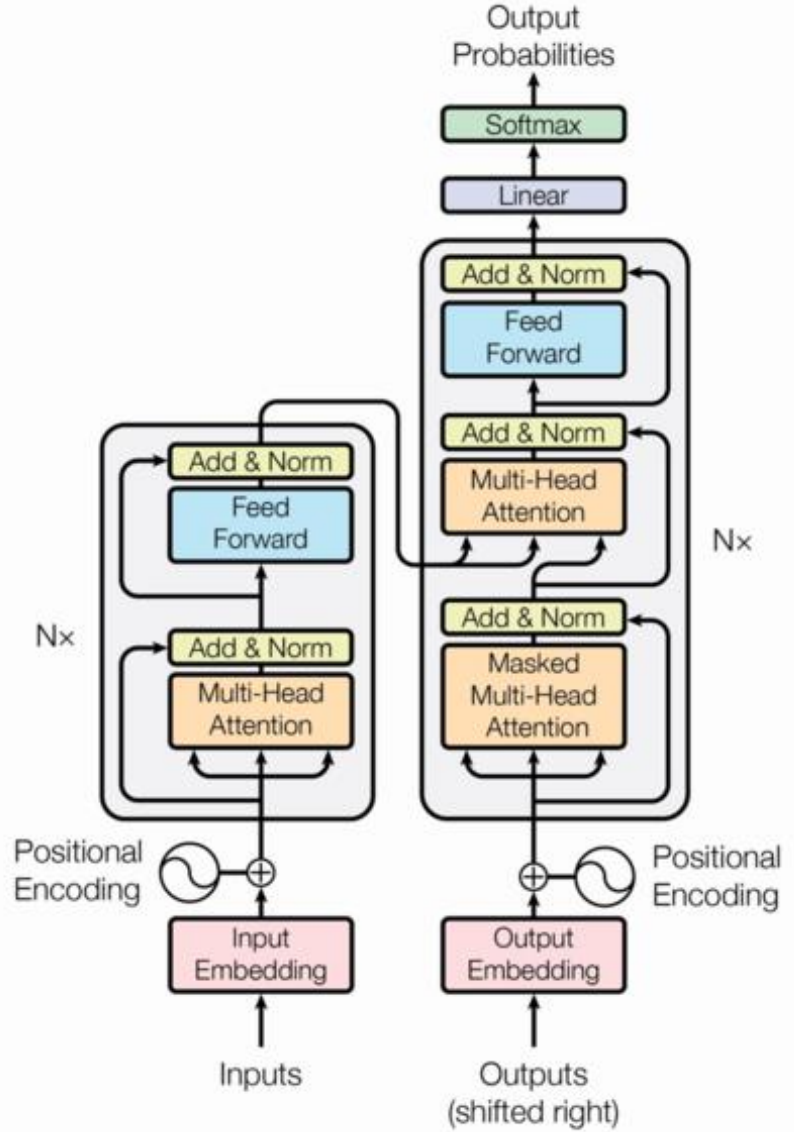
Transformer: Attention Is All You Need



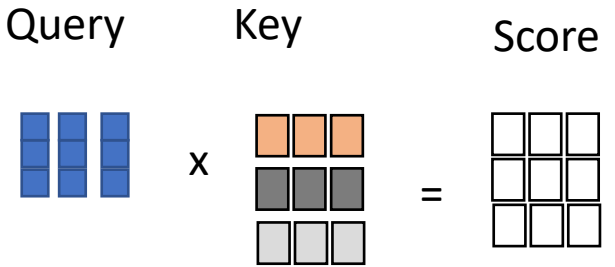
Self Attention



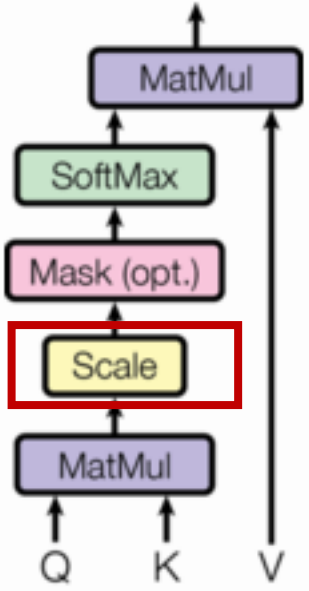
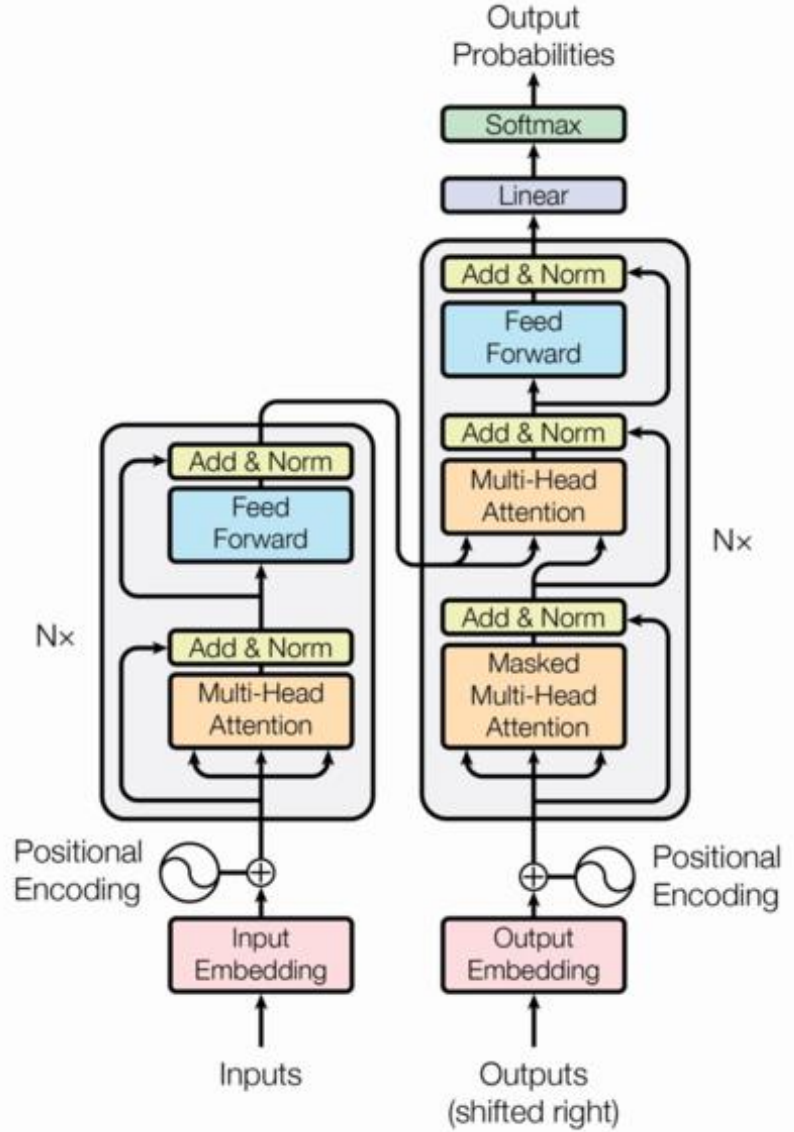
Transformer: Attention Is All You Need



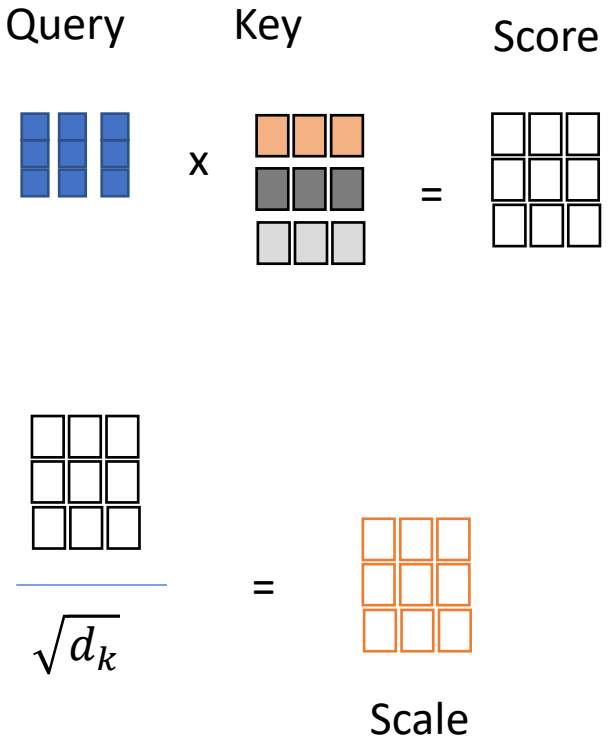
Self Attention



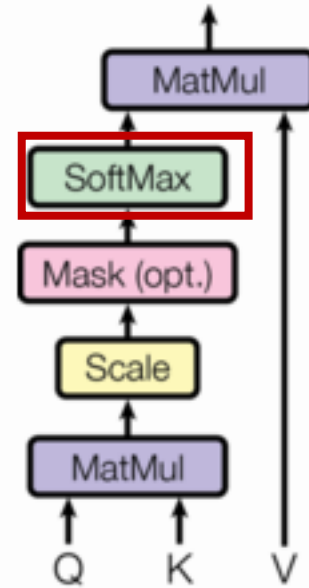
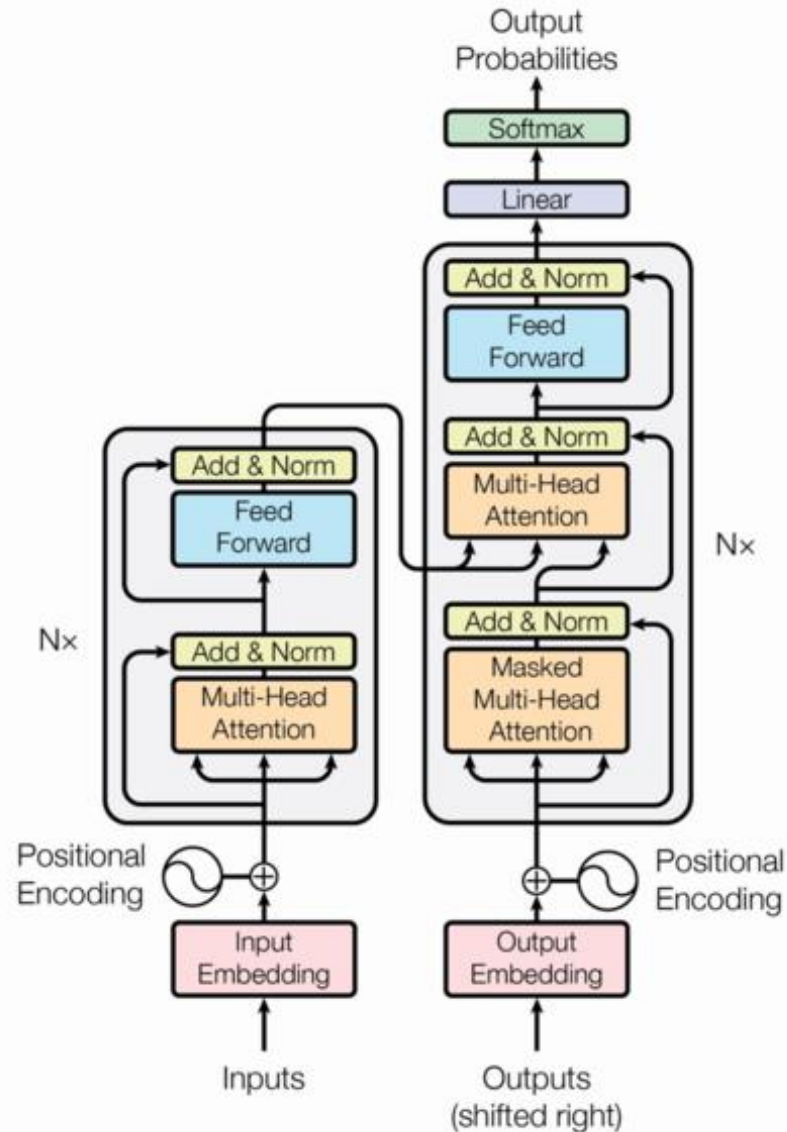
Transformer: Attention Is All You Need



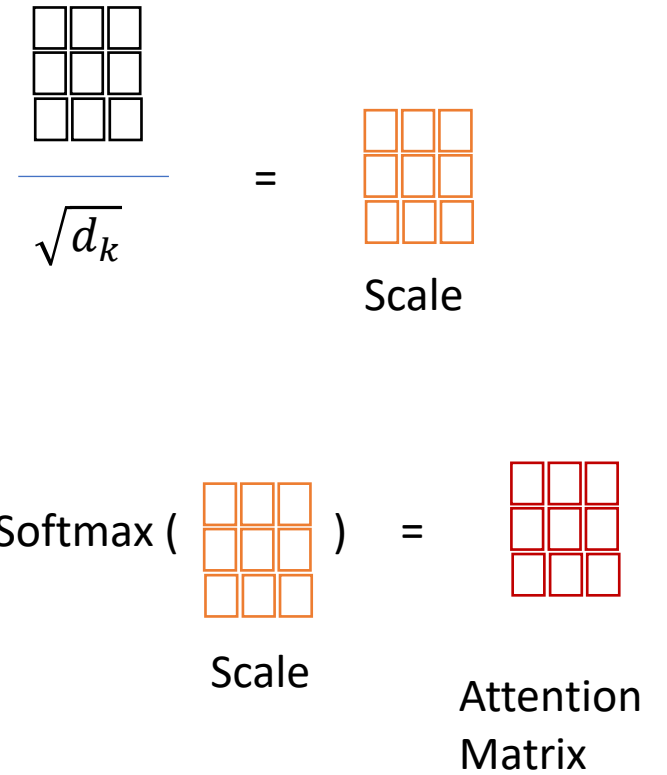
Self Attention



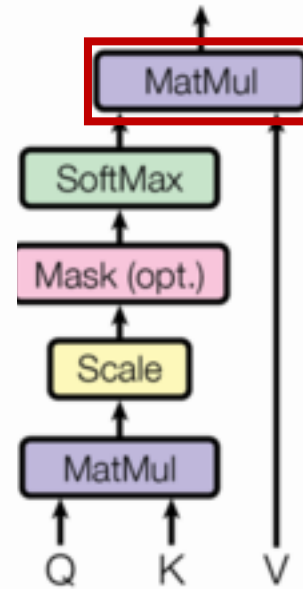
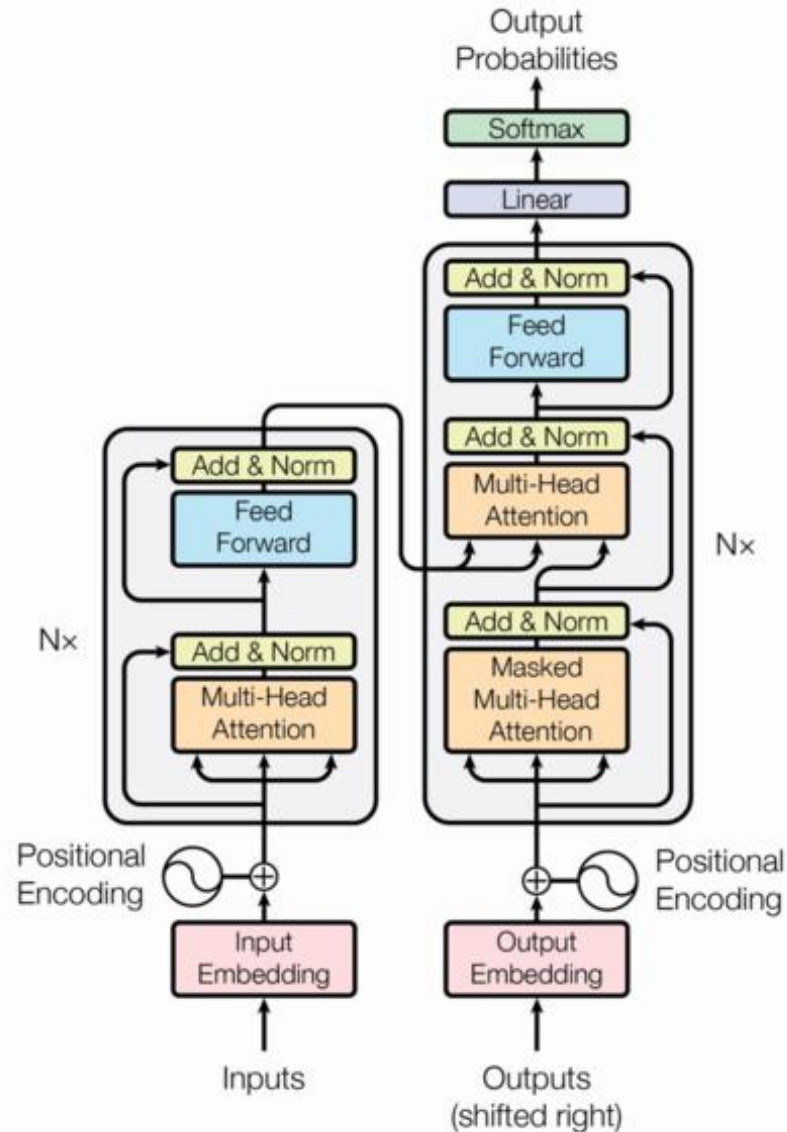
Transformer: Attention Is All You Need



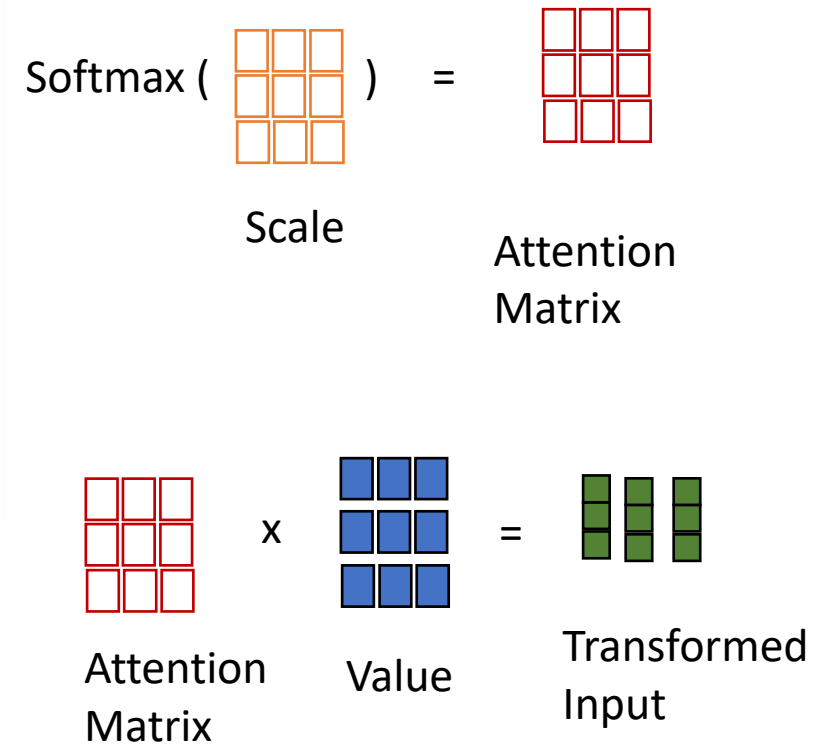
Self Attention



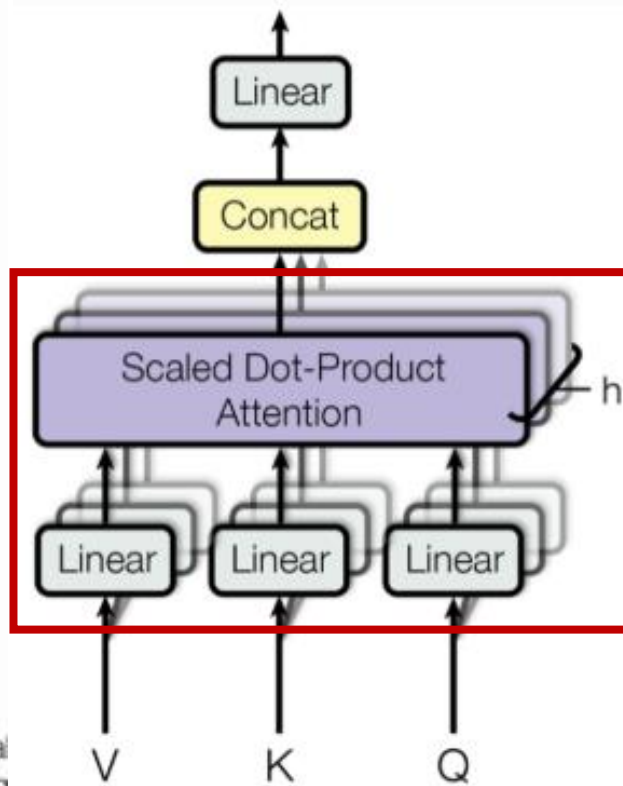
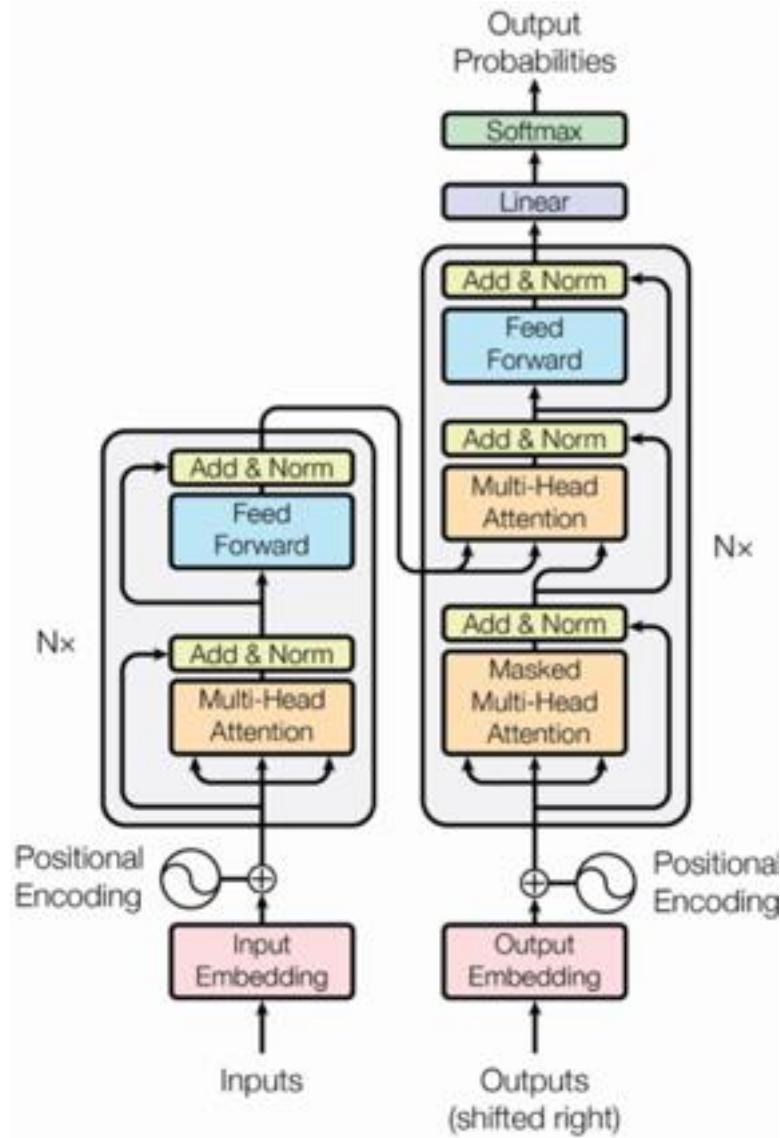
Transformer: Attention Is All You Need



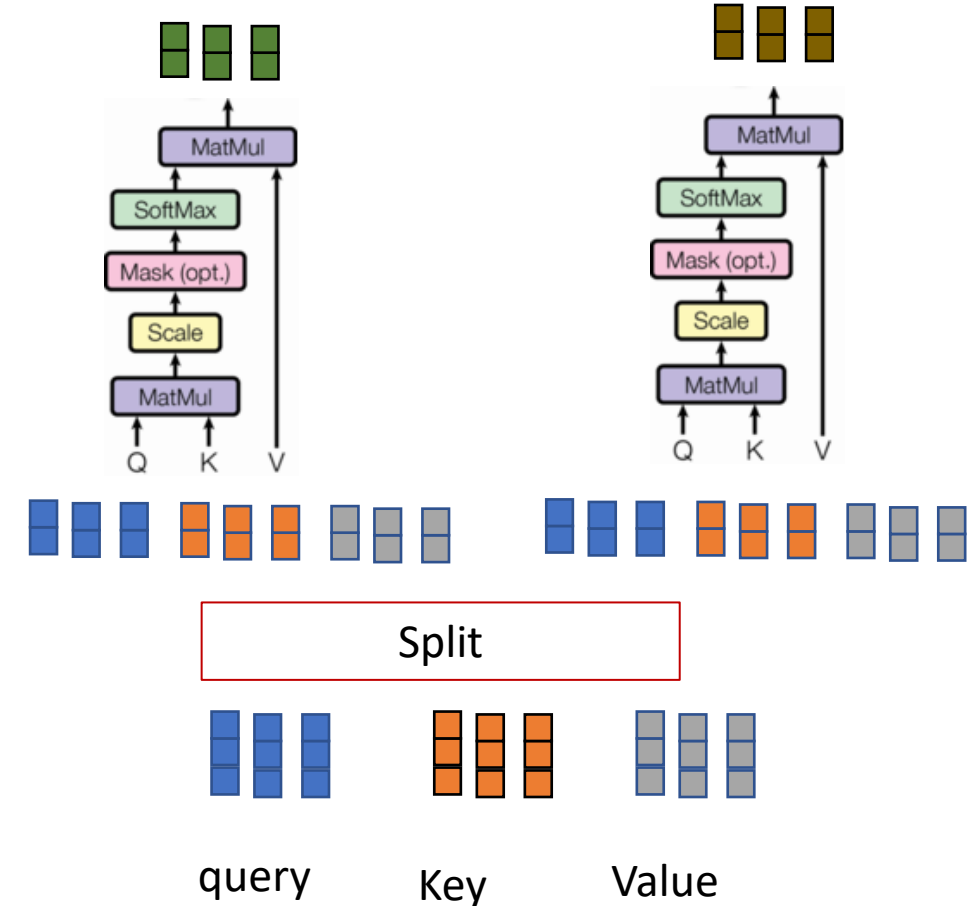
Self Attention



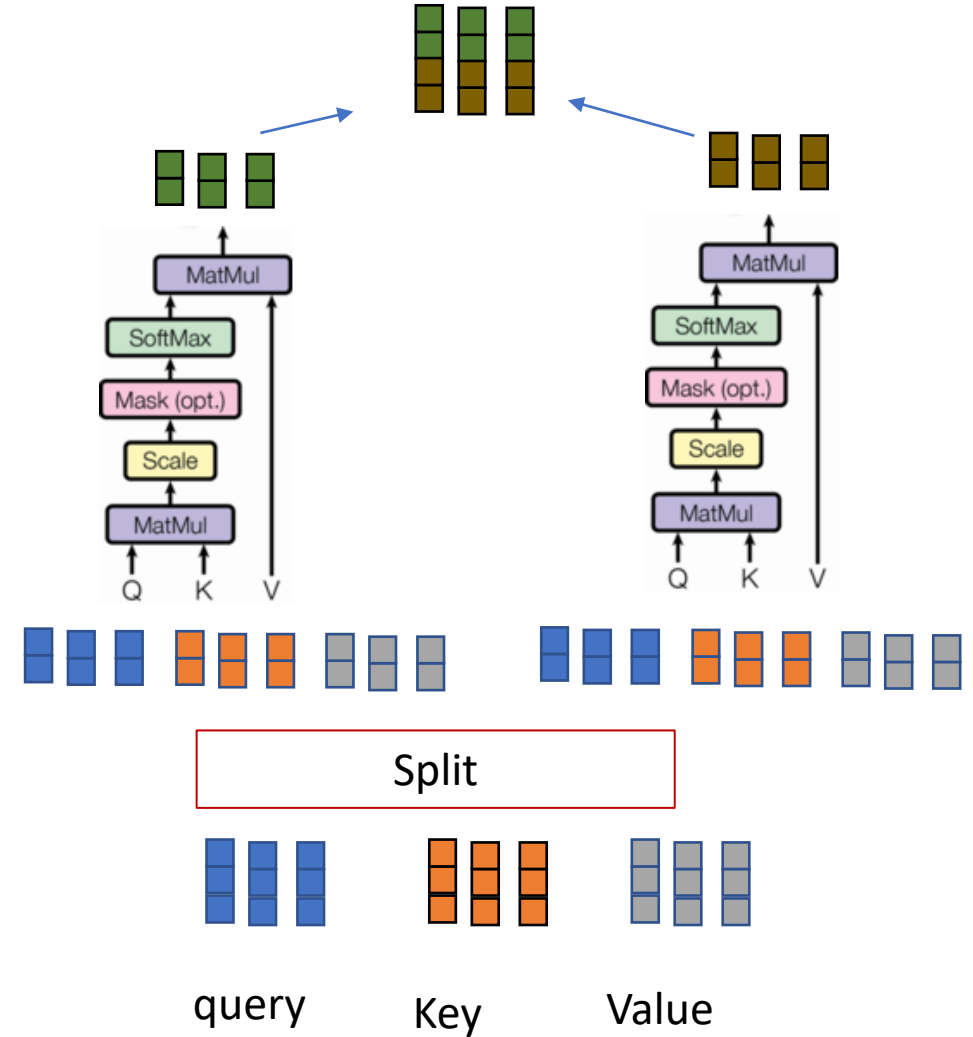
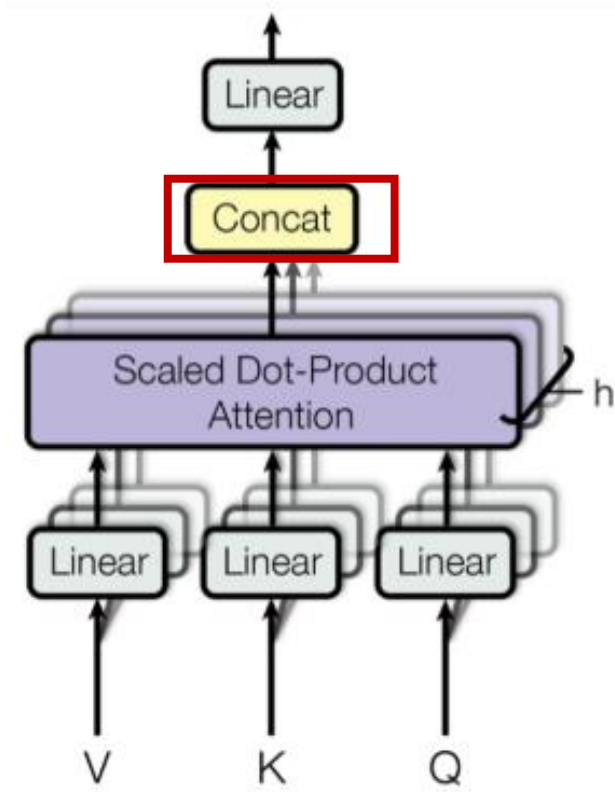
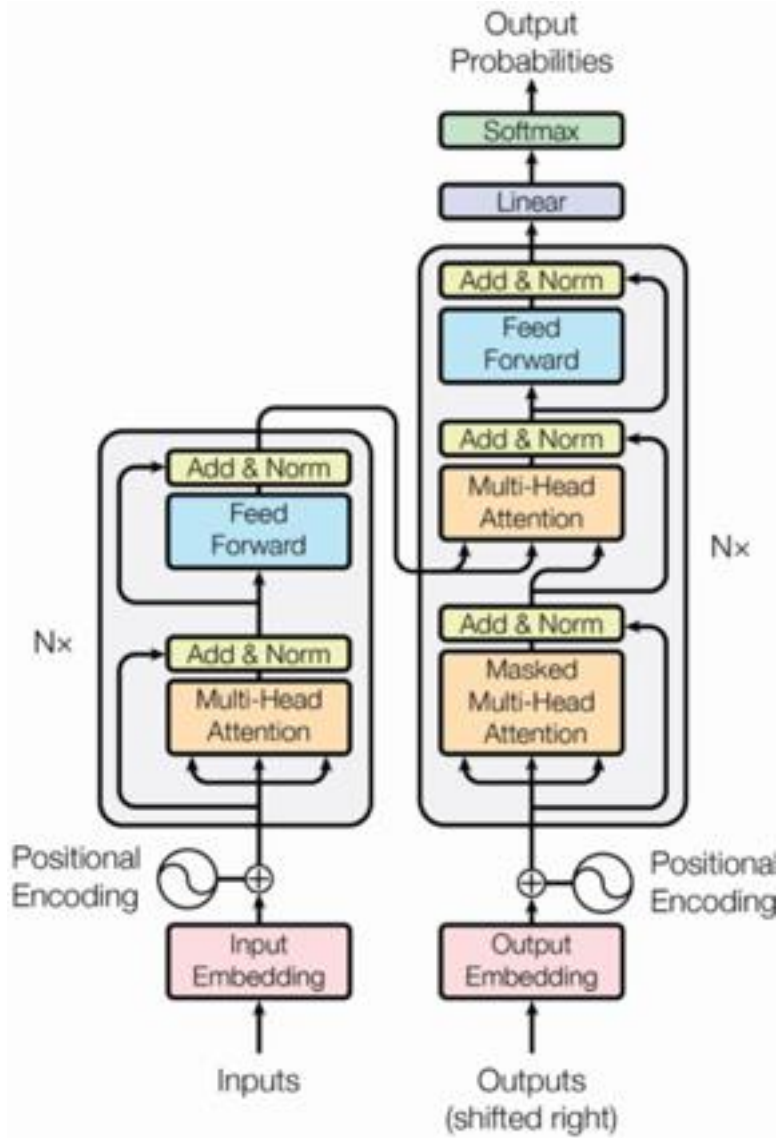
Transformer: Attention Is All You Need

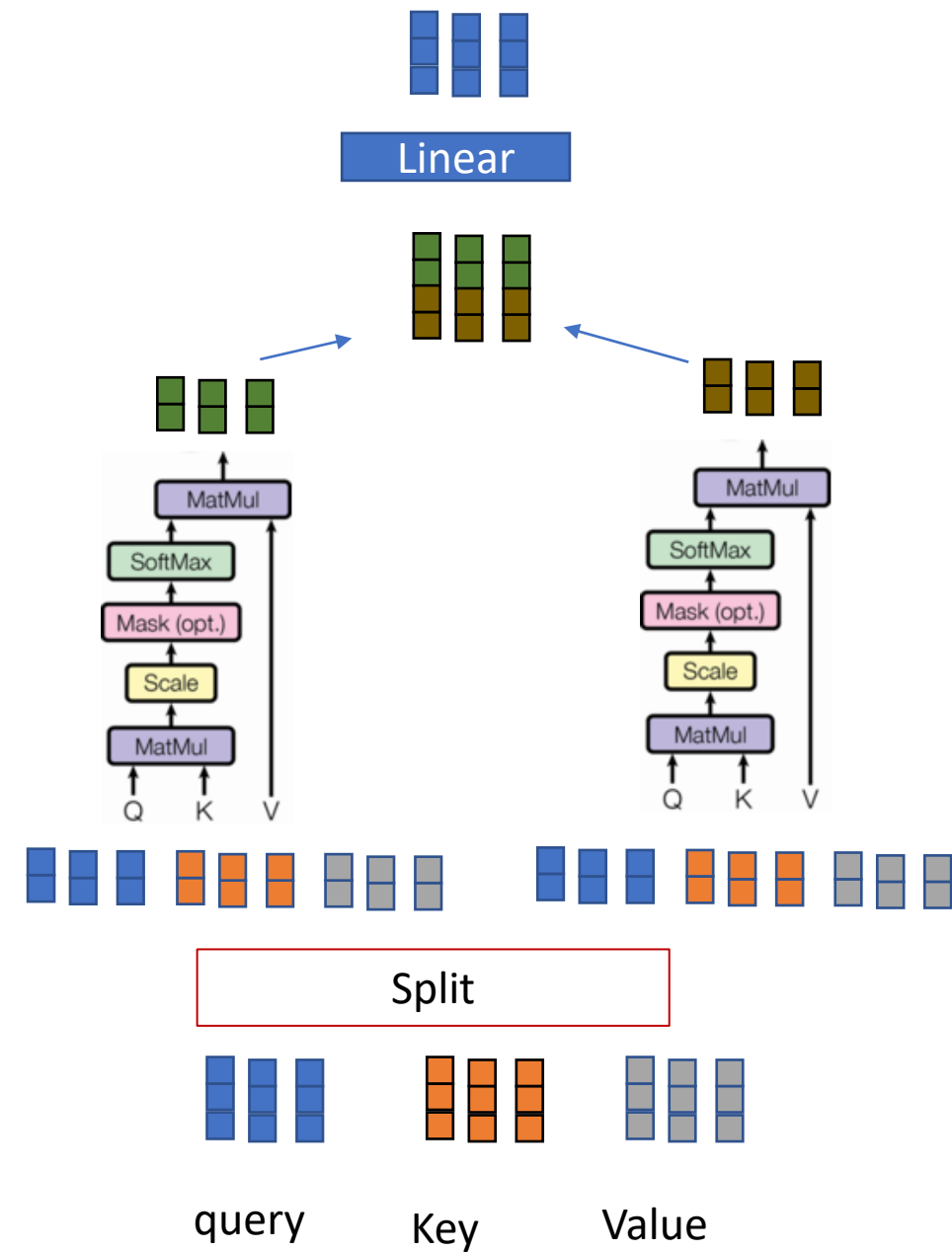
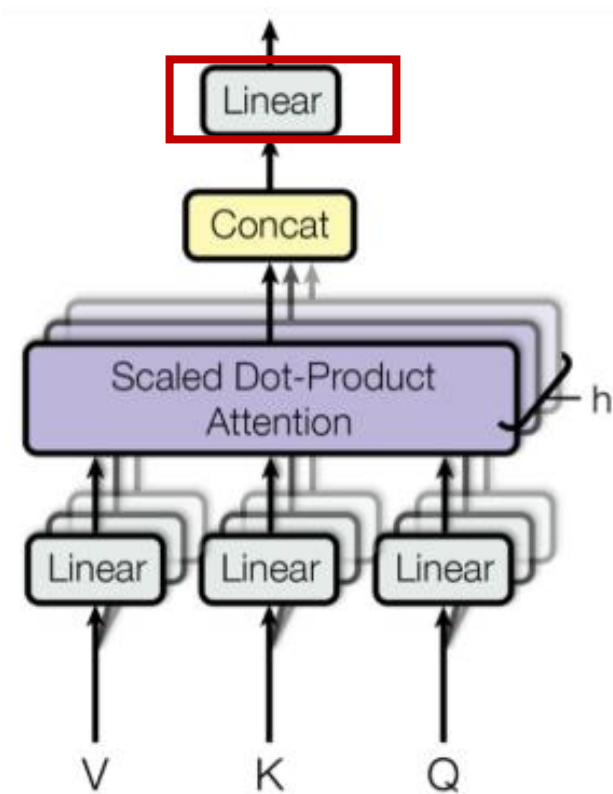
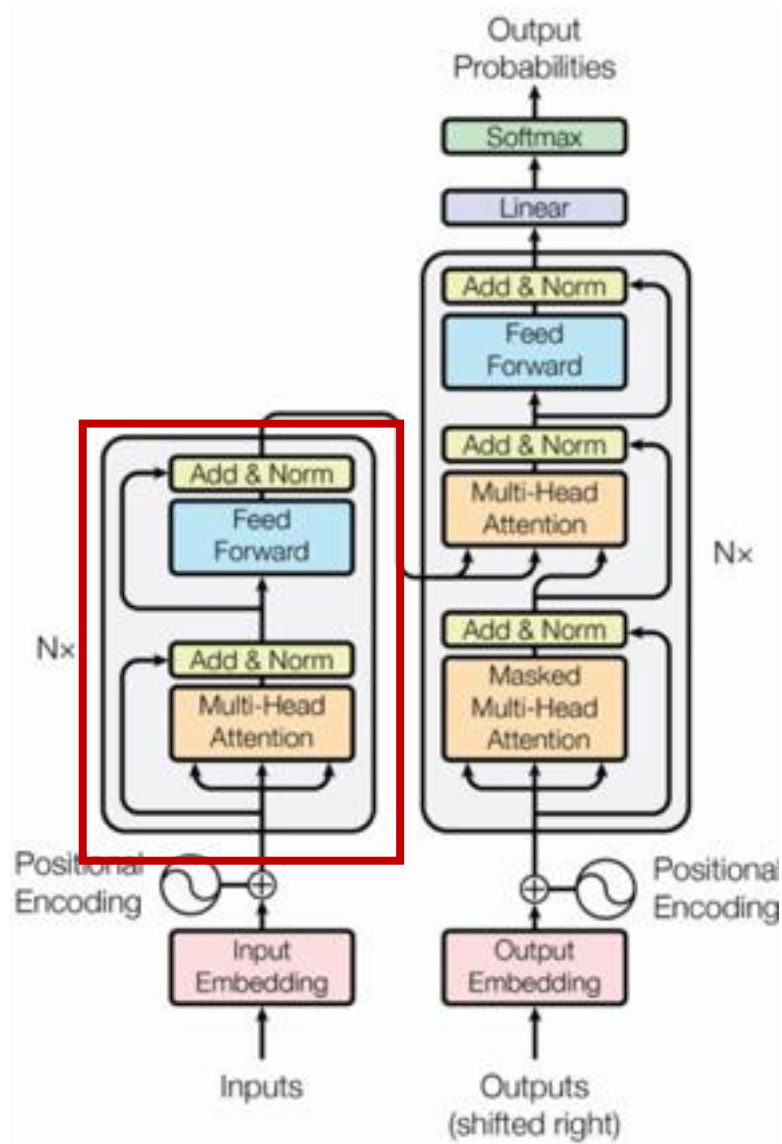


Multi Head Attention

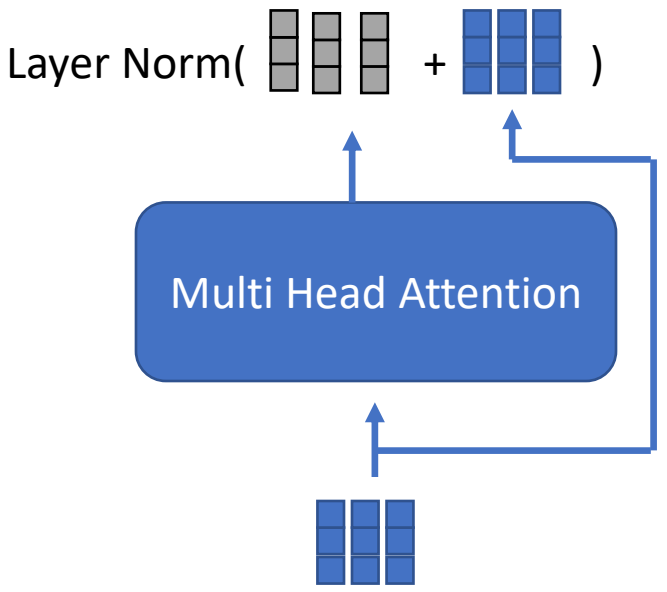
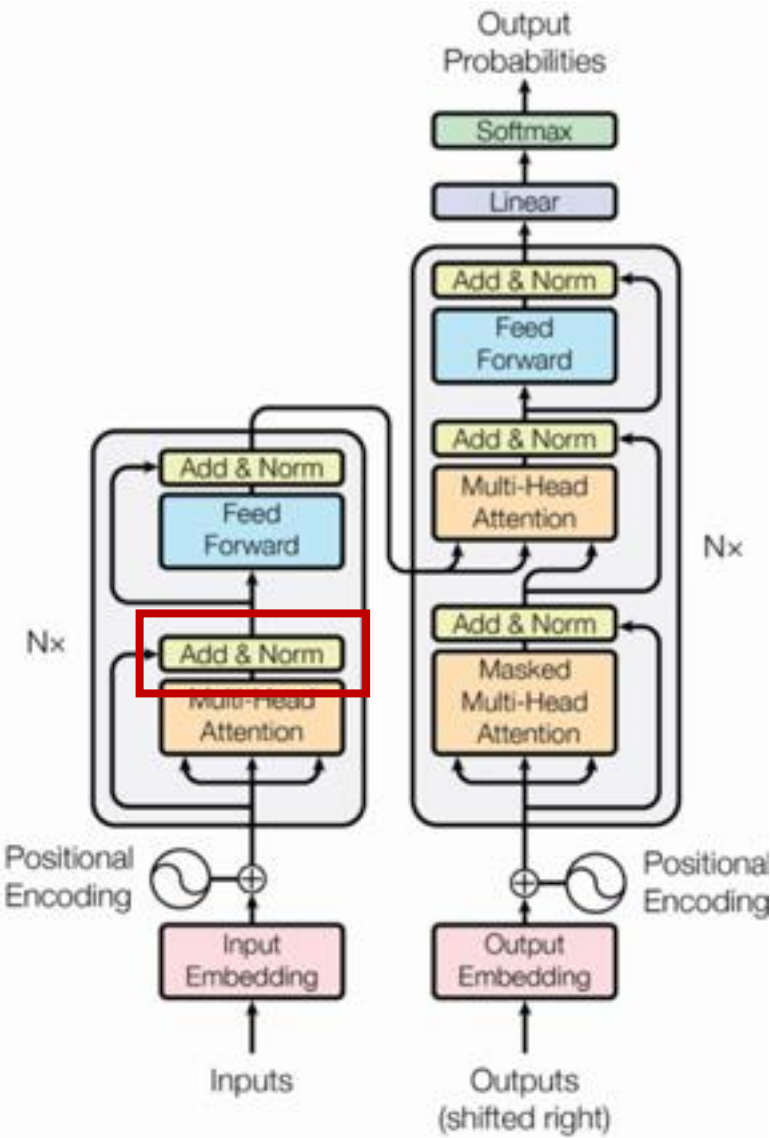


Transformer: Attention Is All You Need

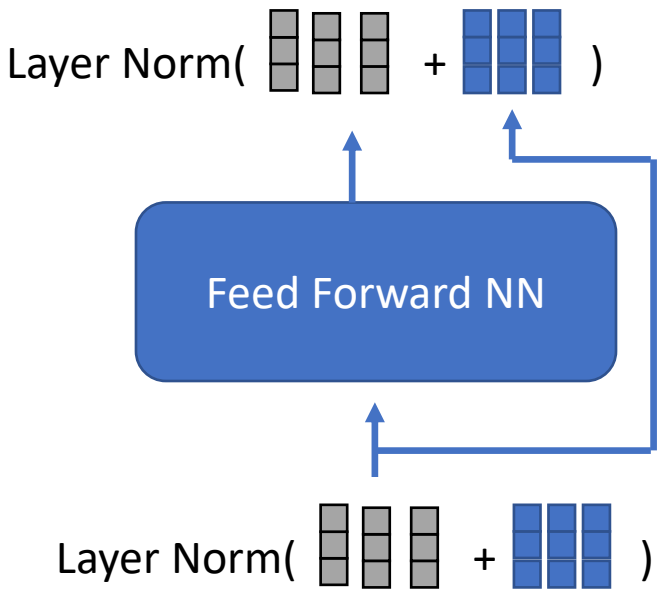
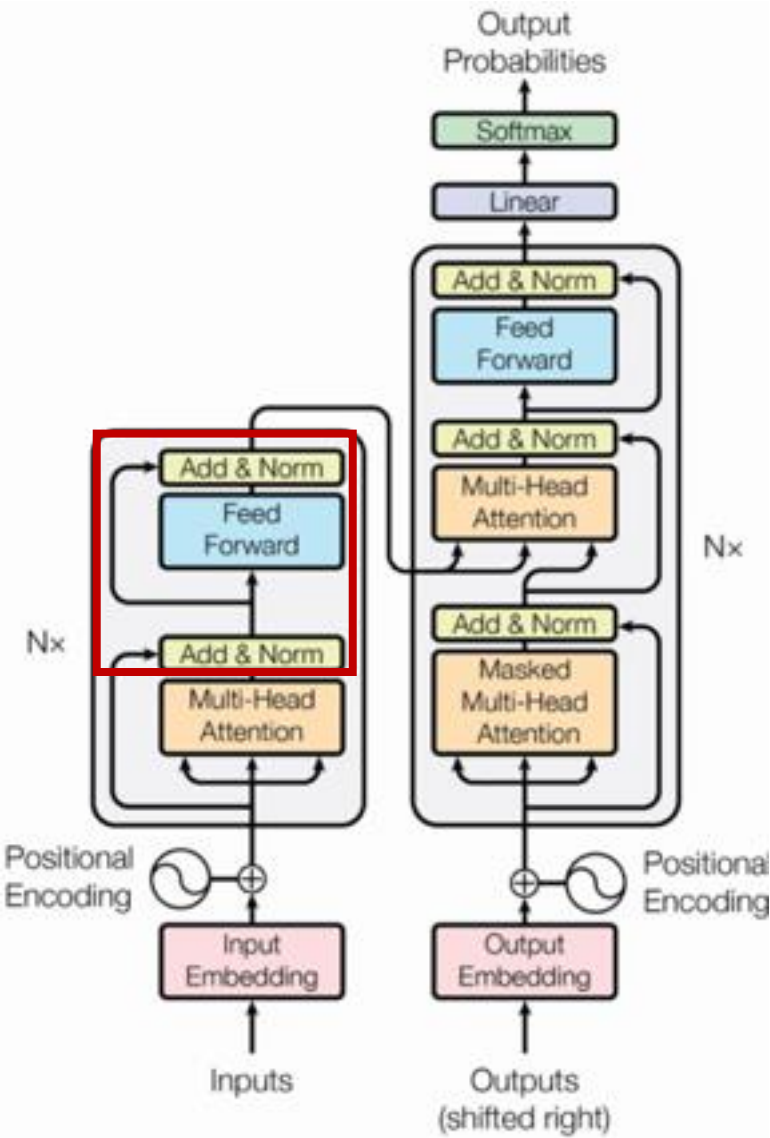


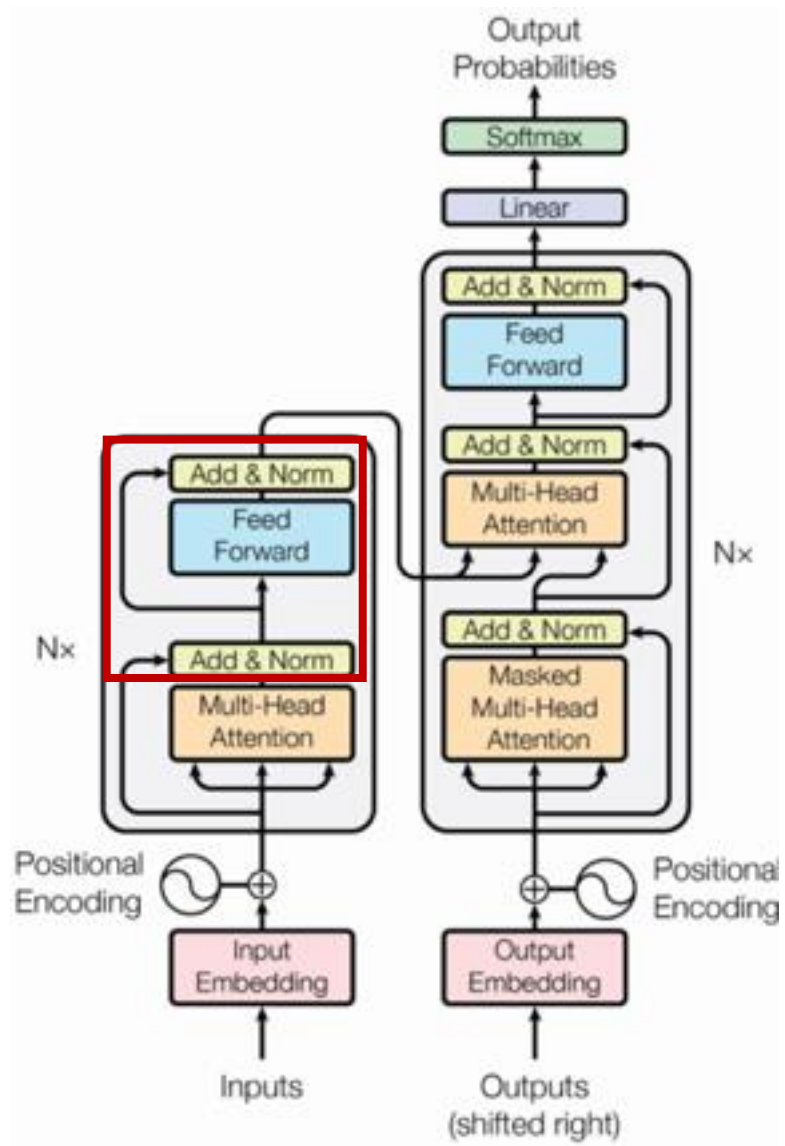


Transformer: Attention Is All You Need

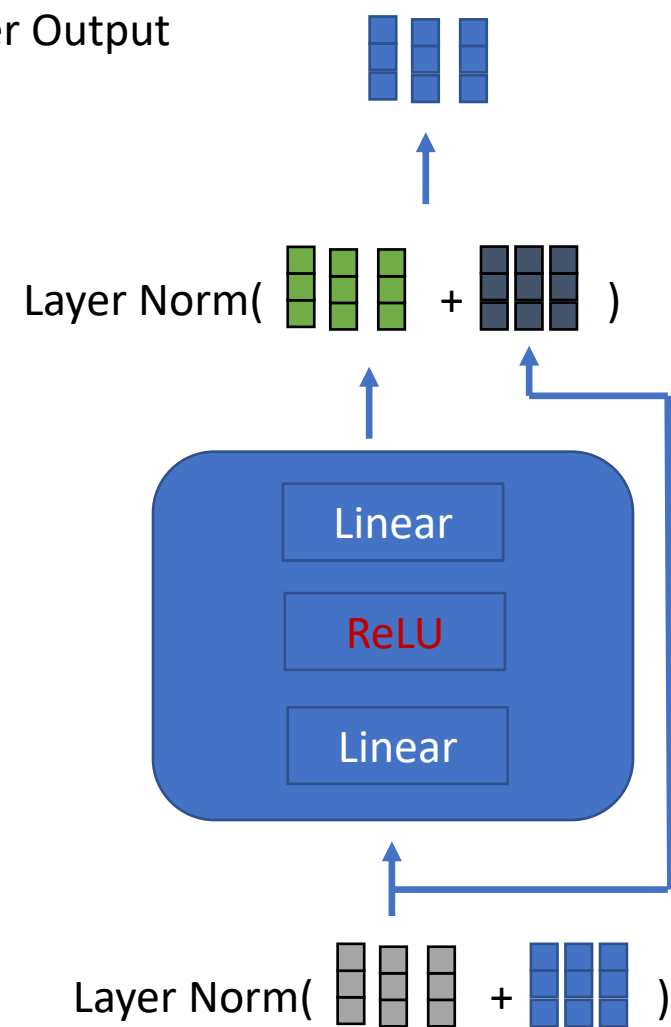


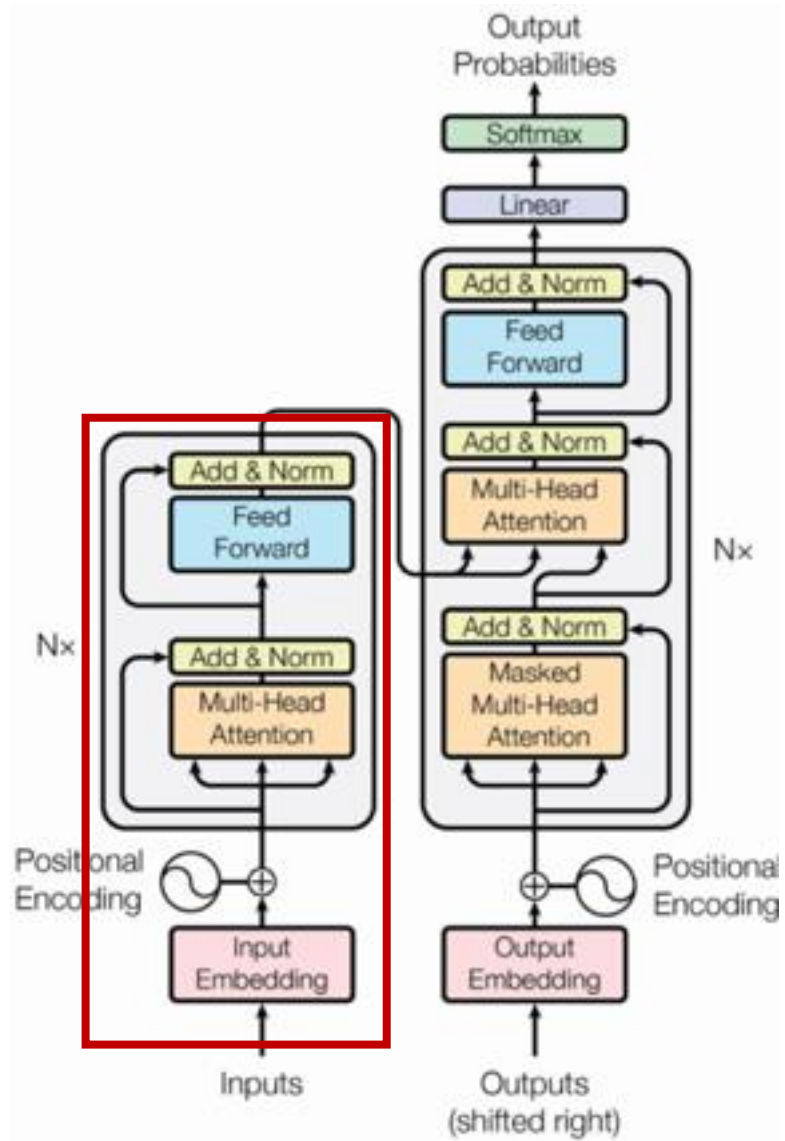
Transformer: Attention Is All You Need



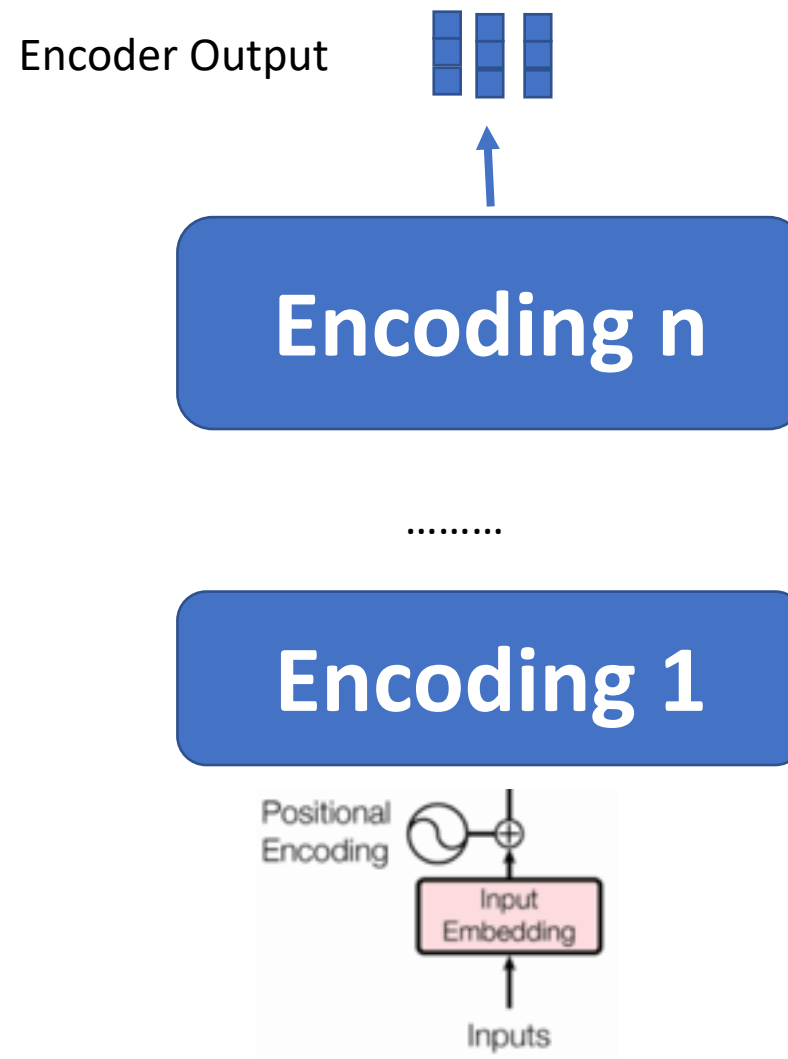
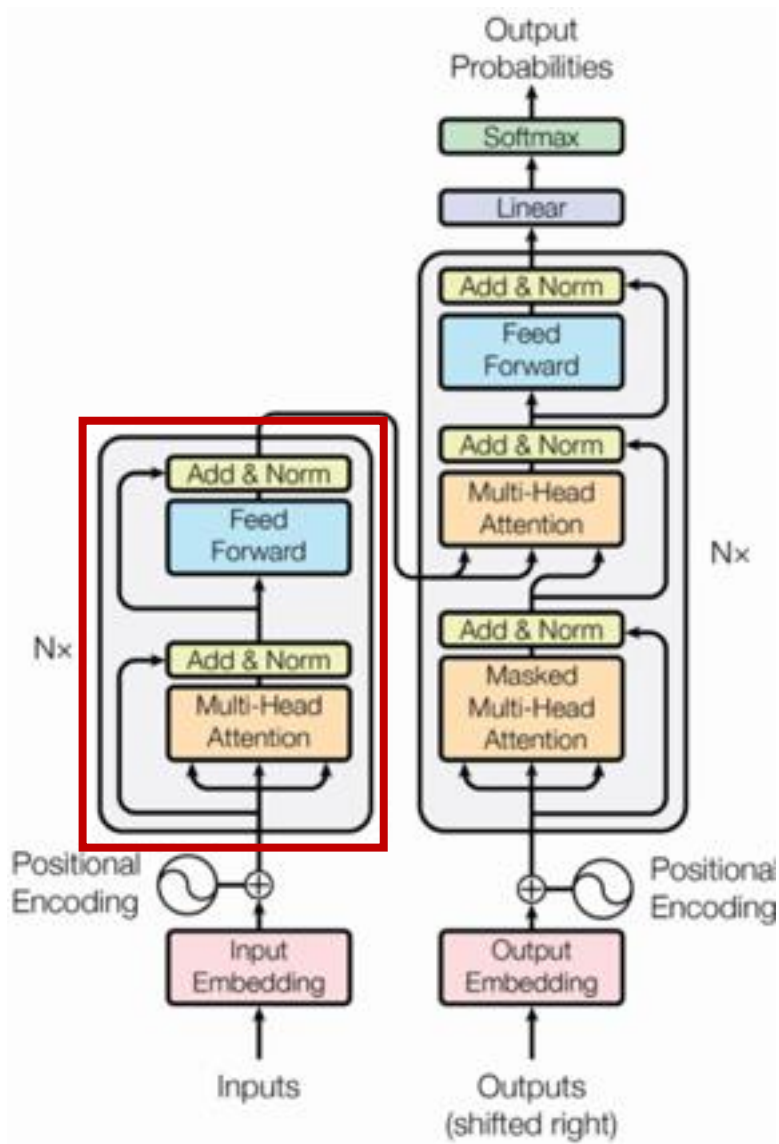


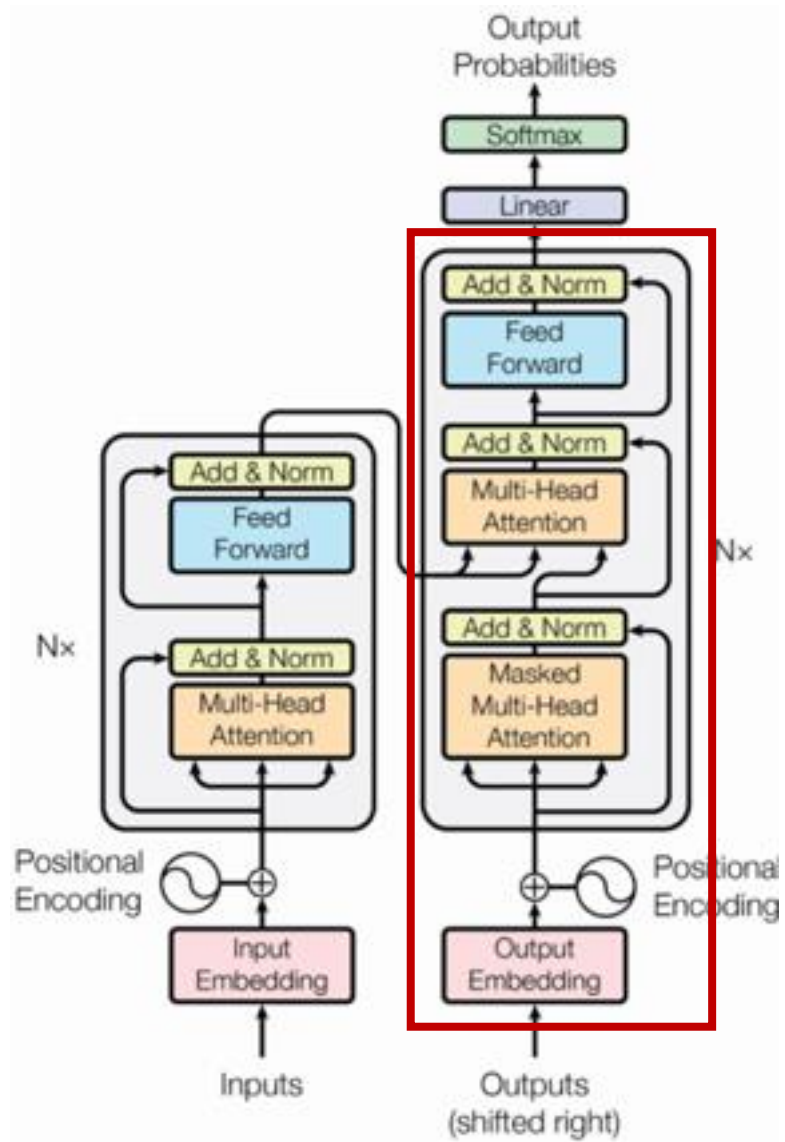
Encoder Output



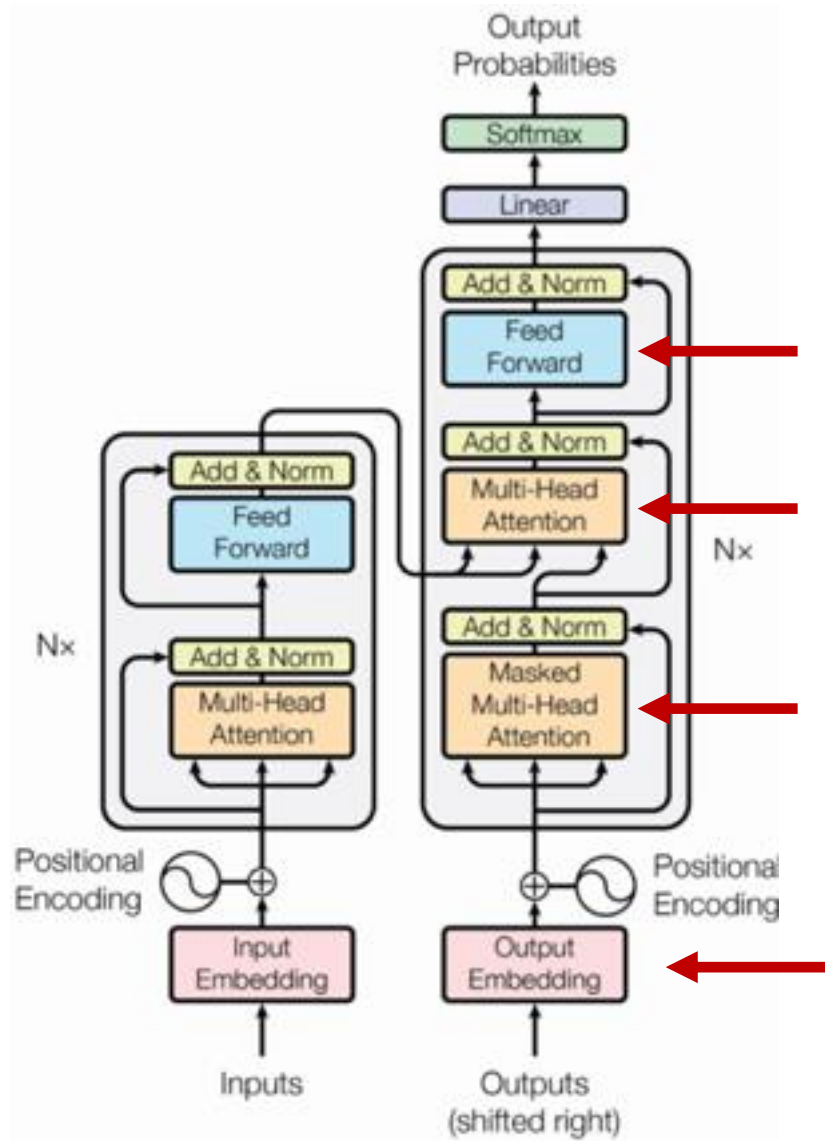


Encoding

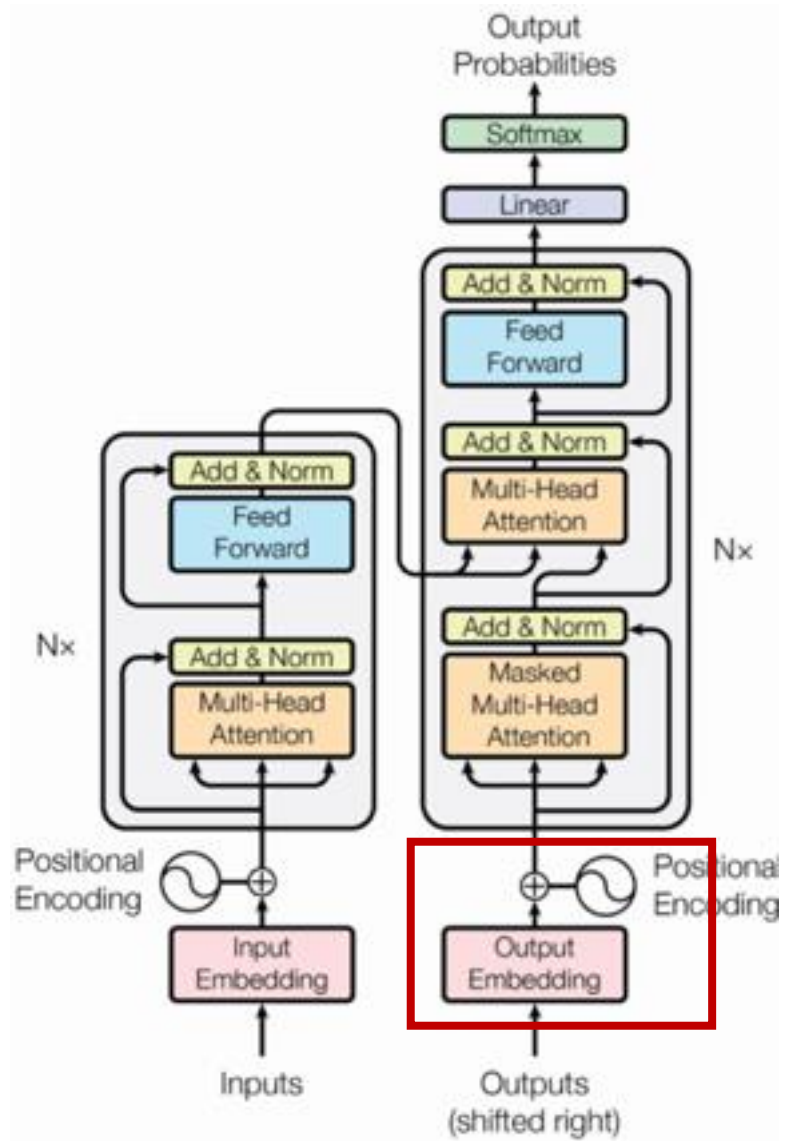




Decoding



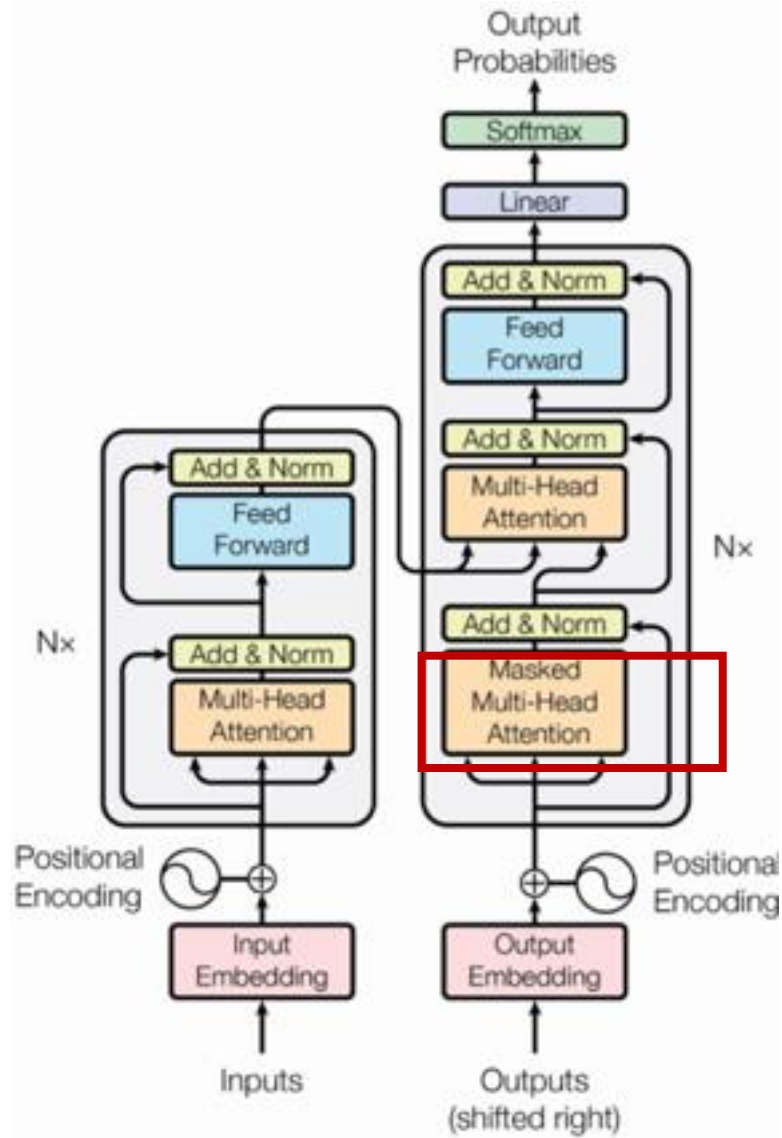
Decoding

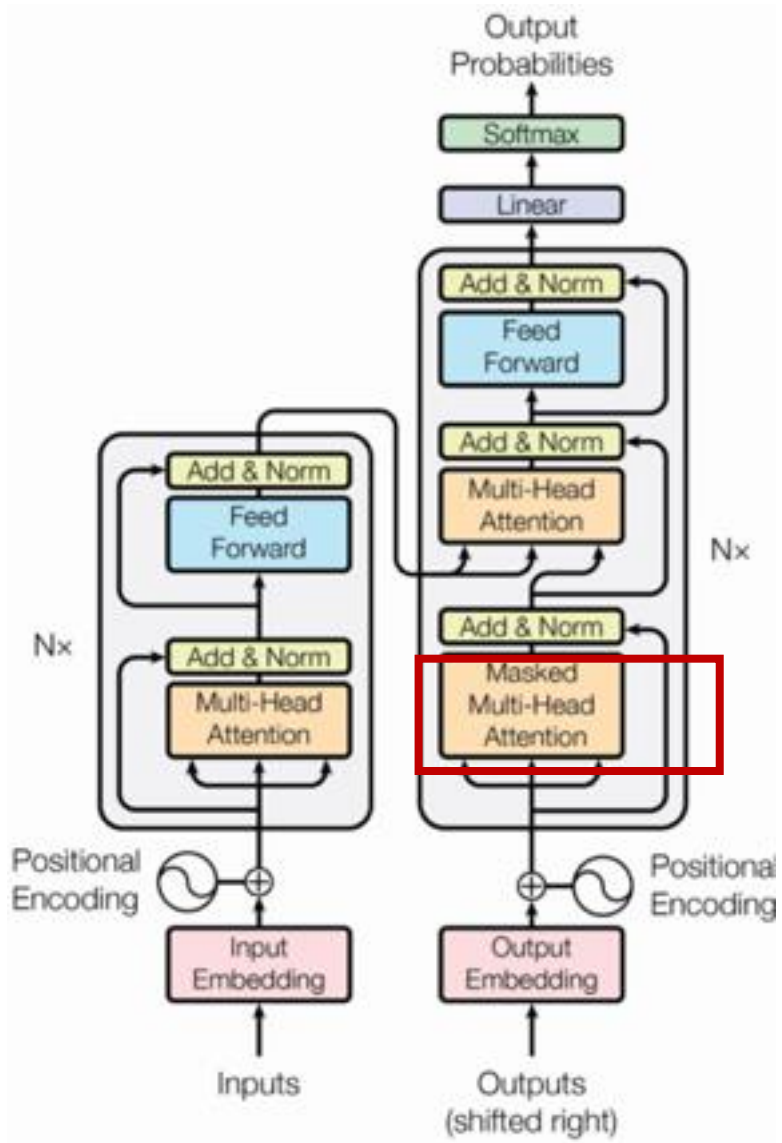


**Word Embedding
+
Positional Encoding**

Masked Attention

i love you → मैं तुमसे प्यार करता हूँ





Masked Attention

Query

Key

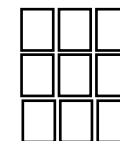
Score



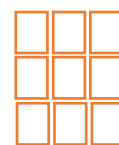
x



=

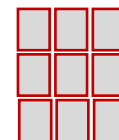


Softmax (

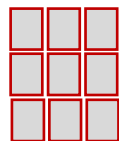


)

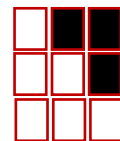
Scale



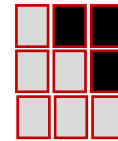
Attention
Matrix



x

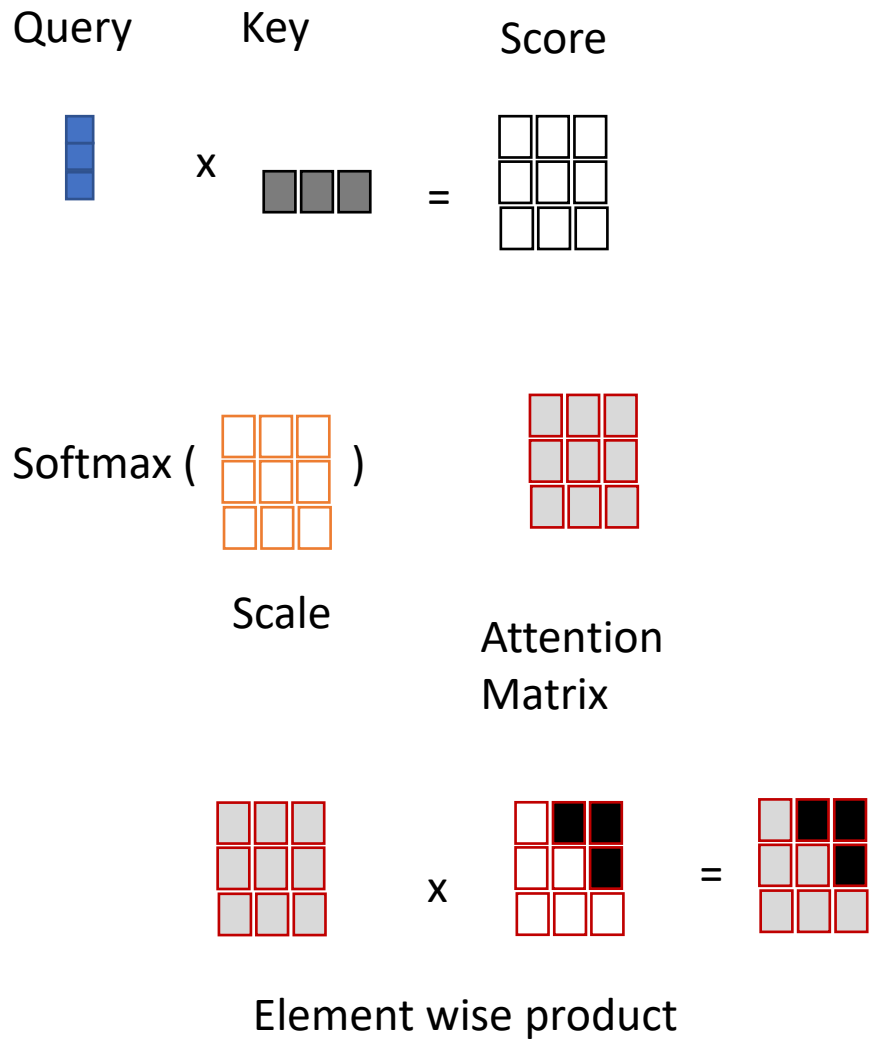
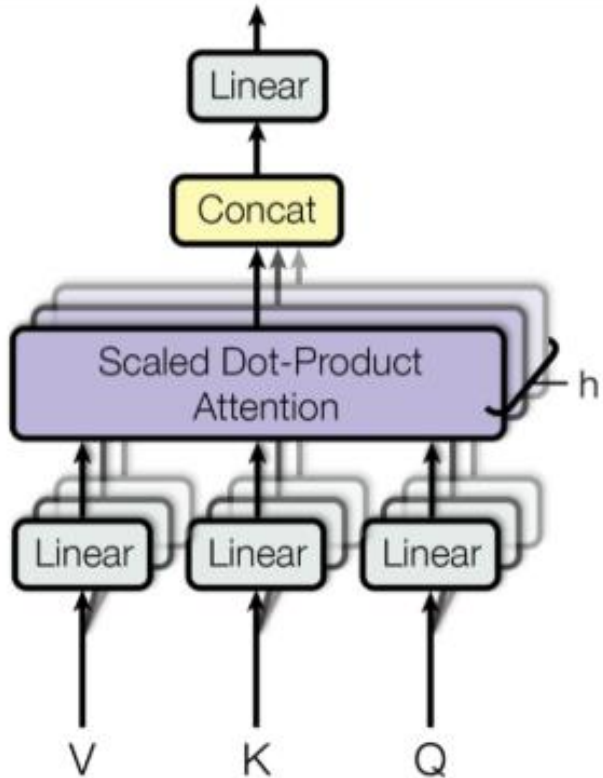
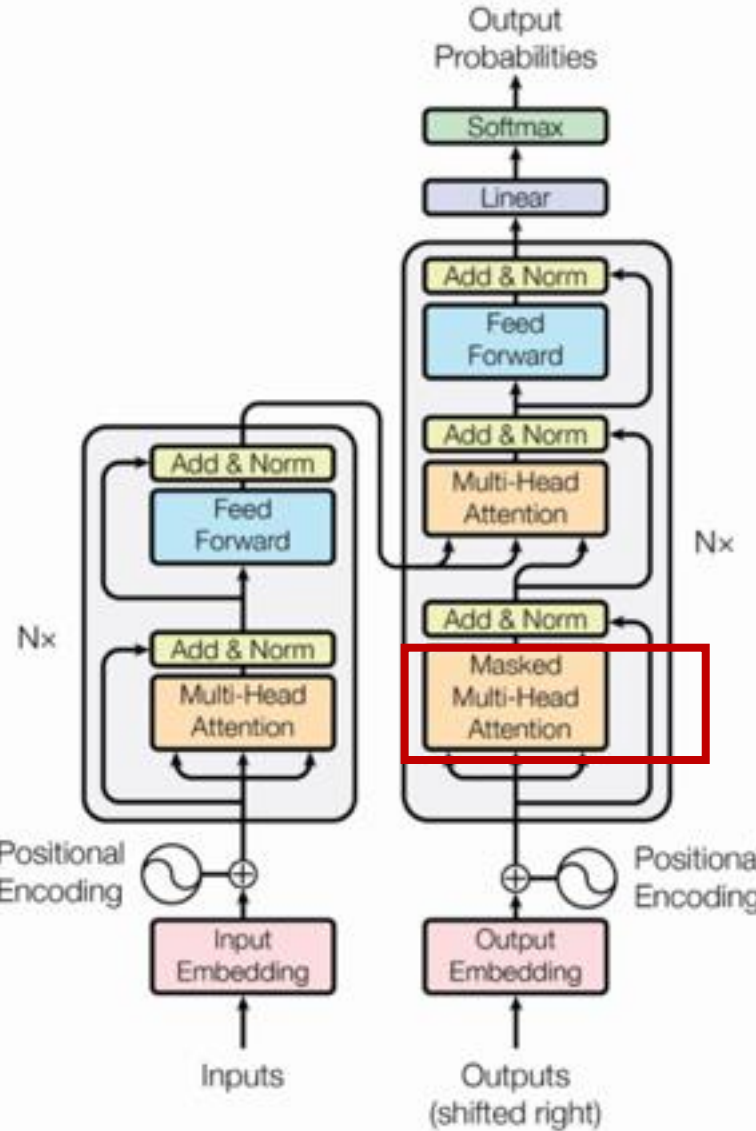


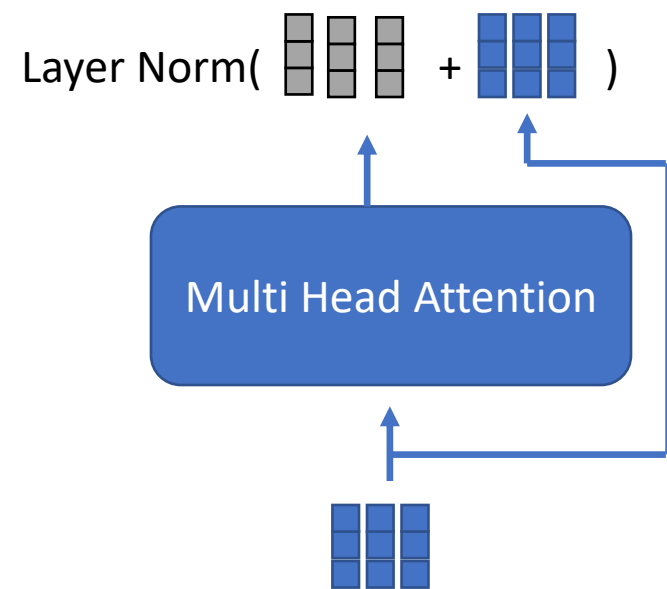
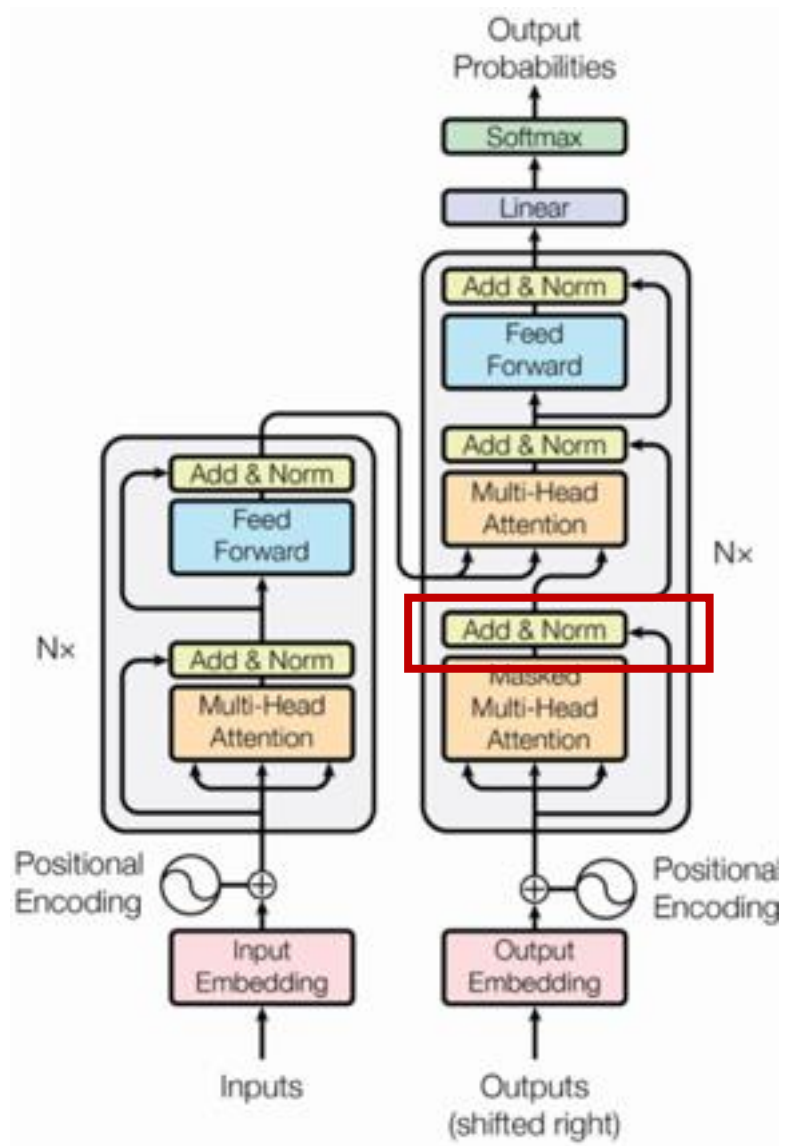
=

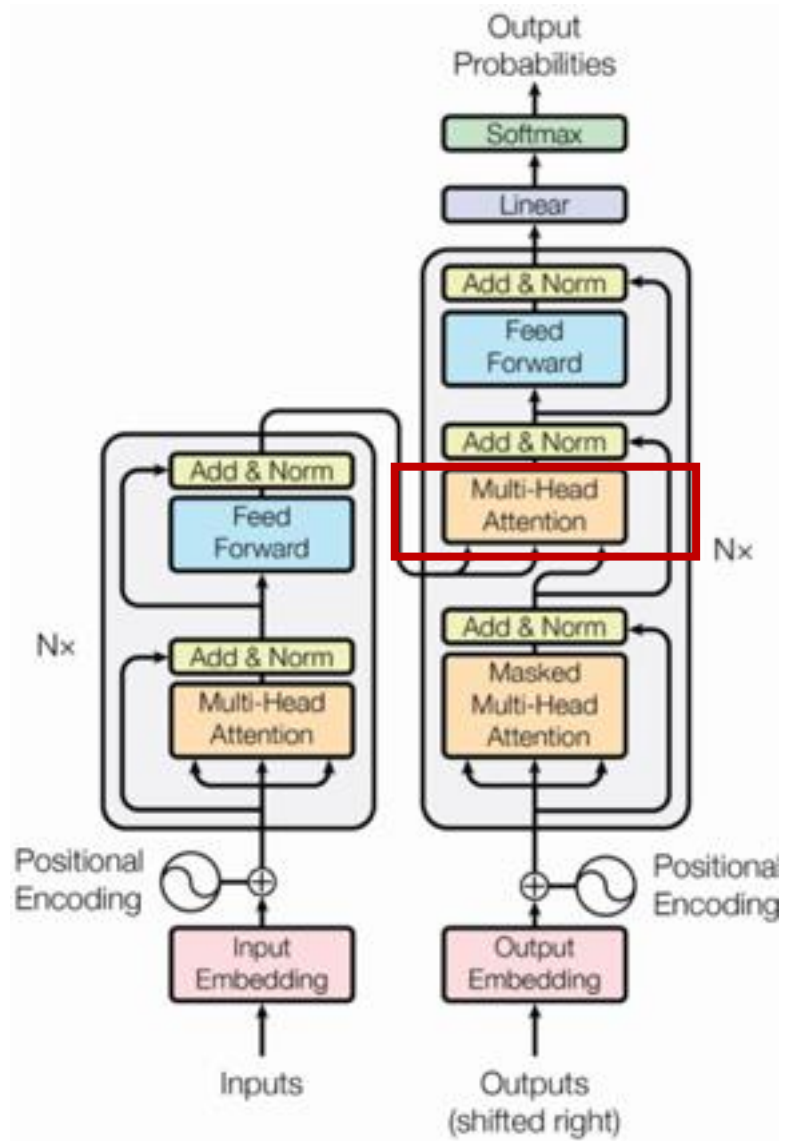


Element wise product

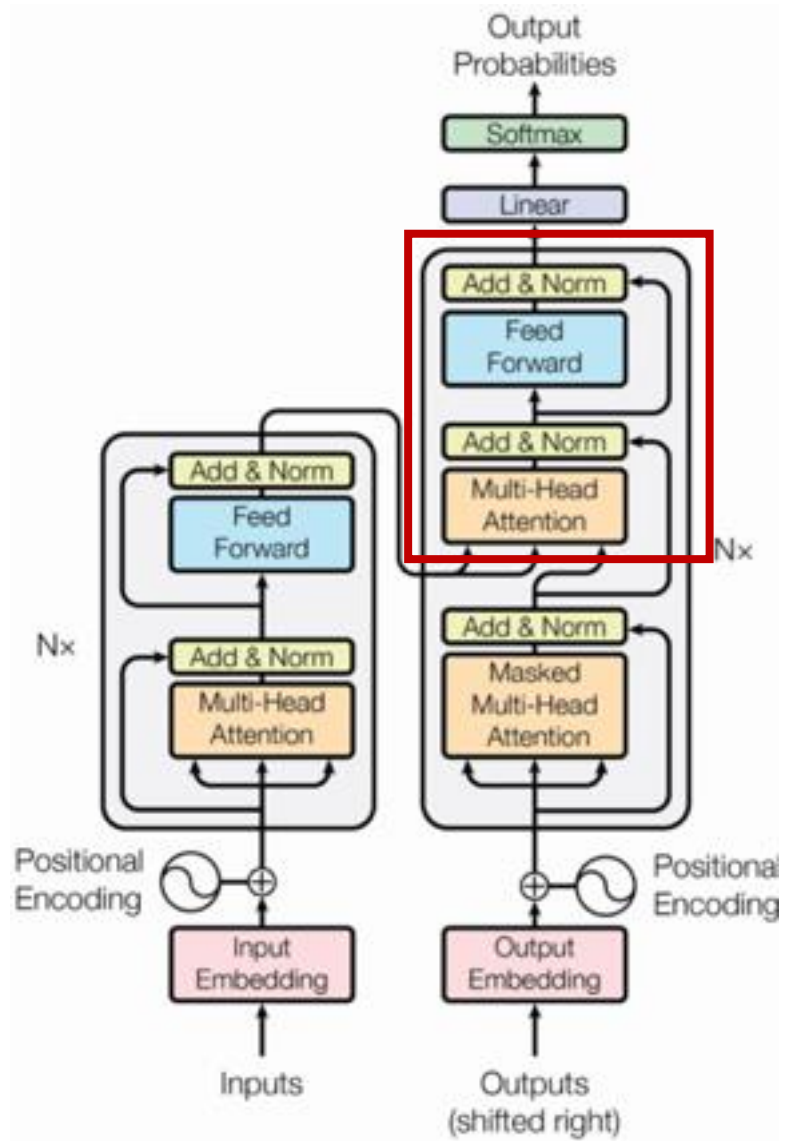
Masked Attention



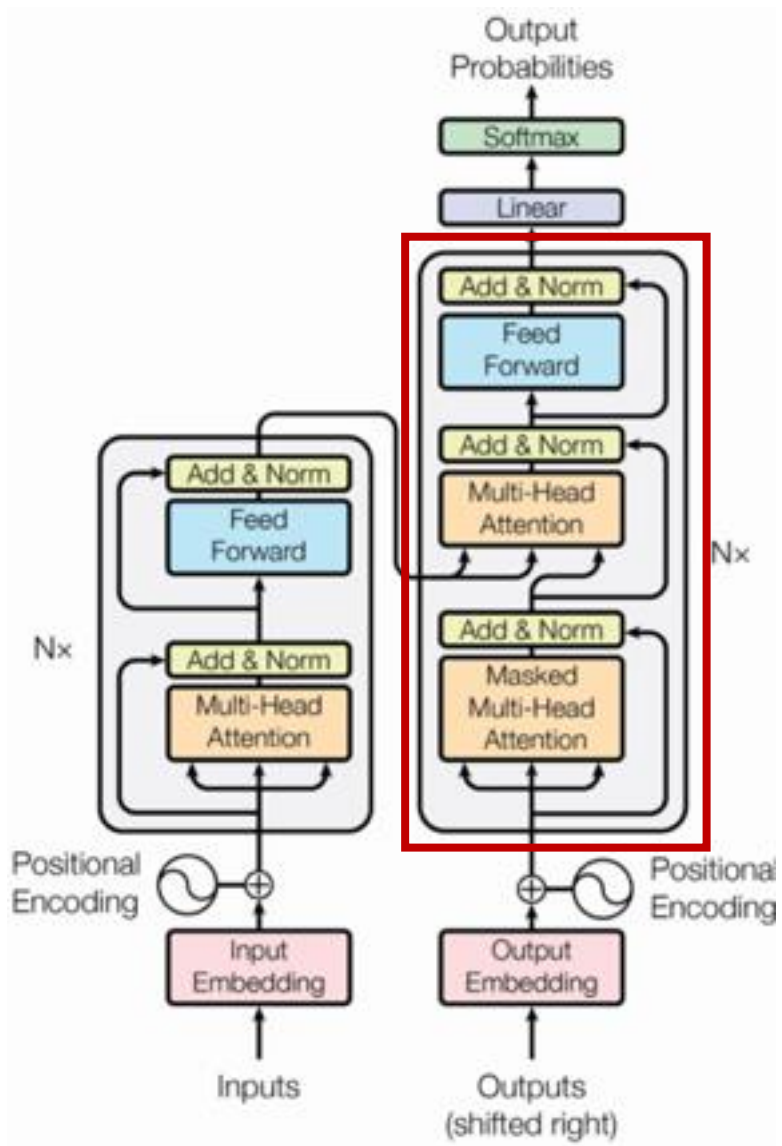




- K and V pairs from the encoder
- Q from the decoder



Other component of the decoder are same as that of the encoder

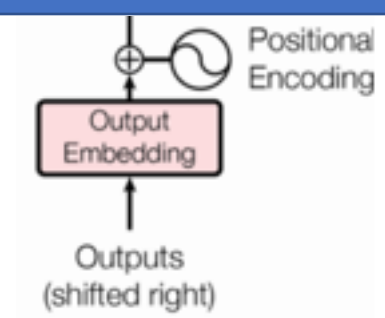


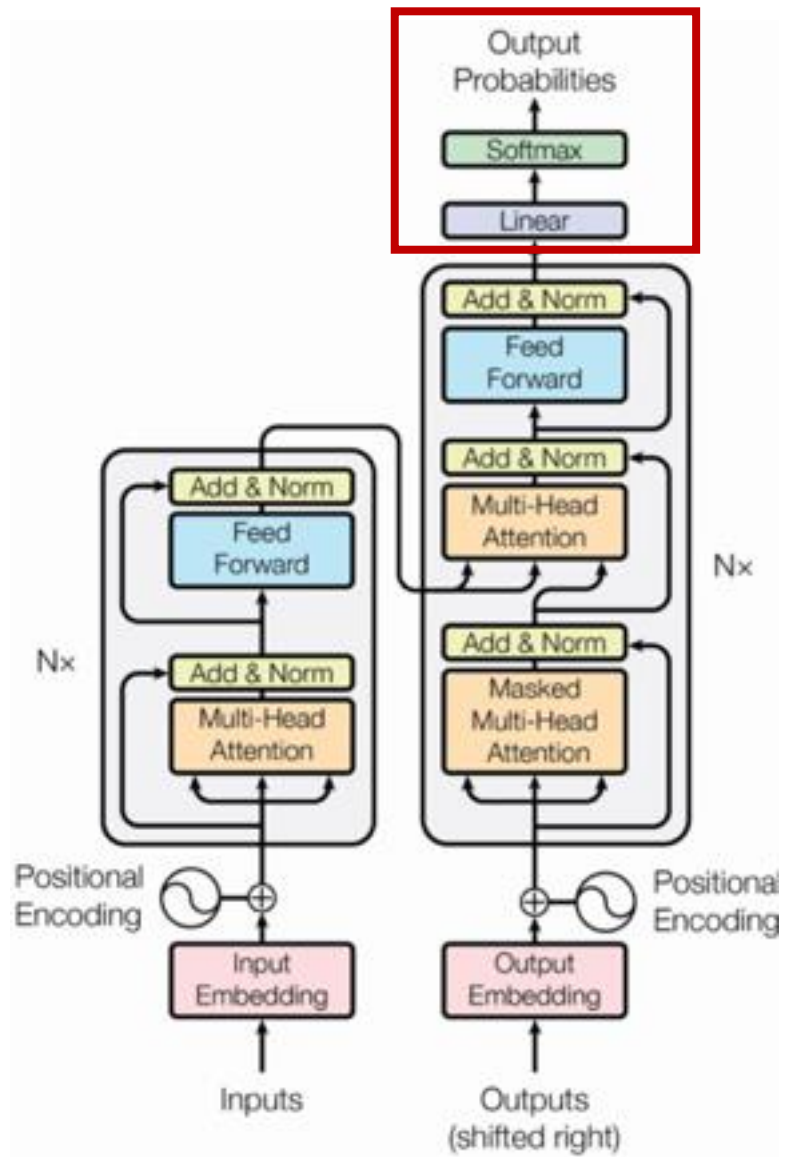
Decoder Output 

Decoding n

.....

Decoding 1



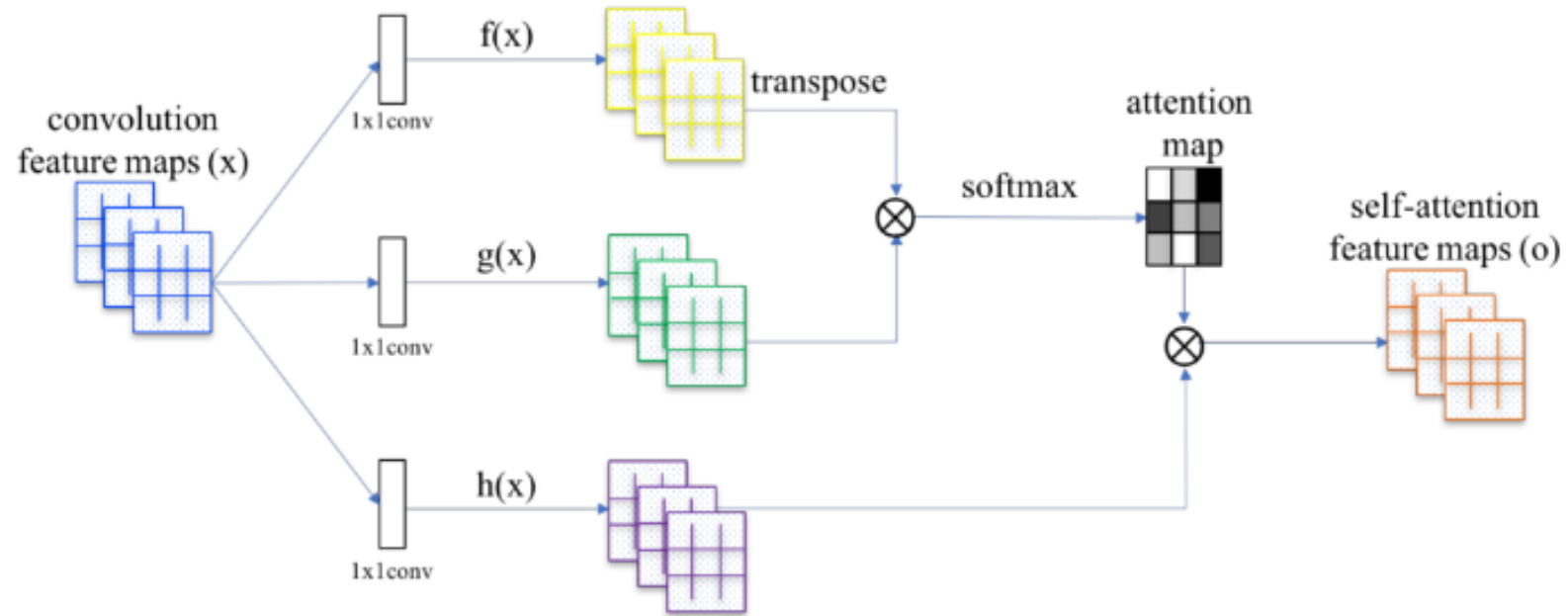


Classifier

Attention in CNN

- **Selp Attention**
- **Squeeze and Excitation (SE)**
- **Convolutional Block Attention Module (CBAM)**

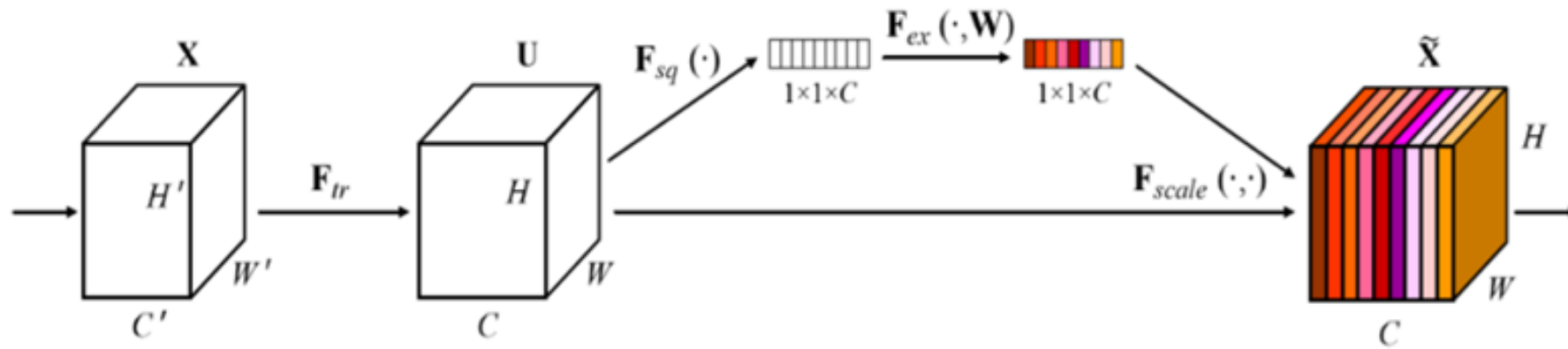
Self Attention Model in CNN



Squeeze and Excitation

Attention mechanism

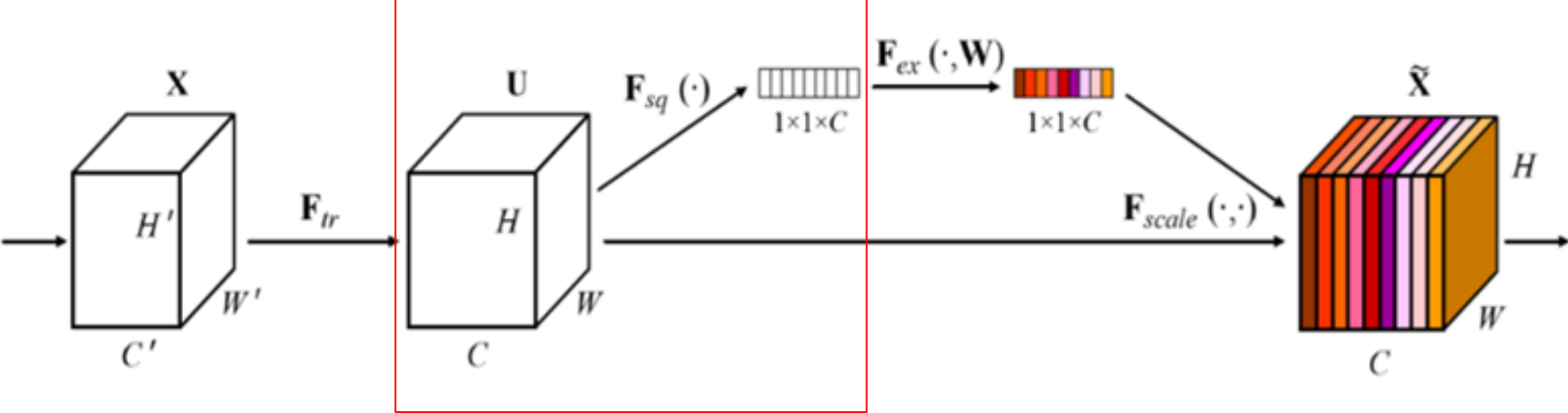
Squeeze and Excitation Model



It is a channel attention

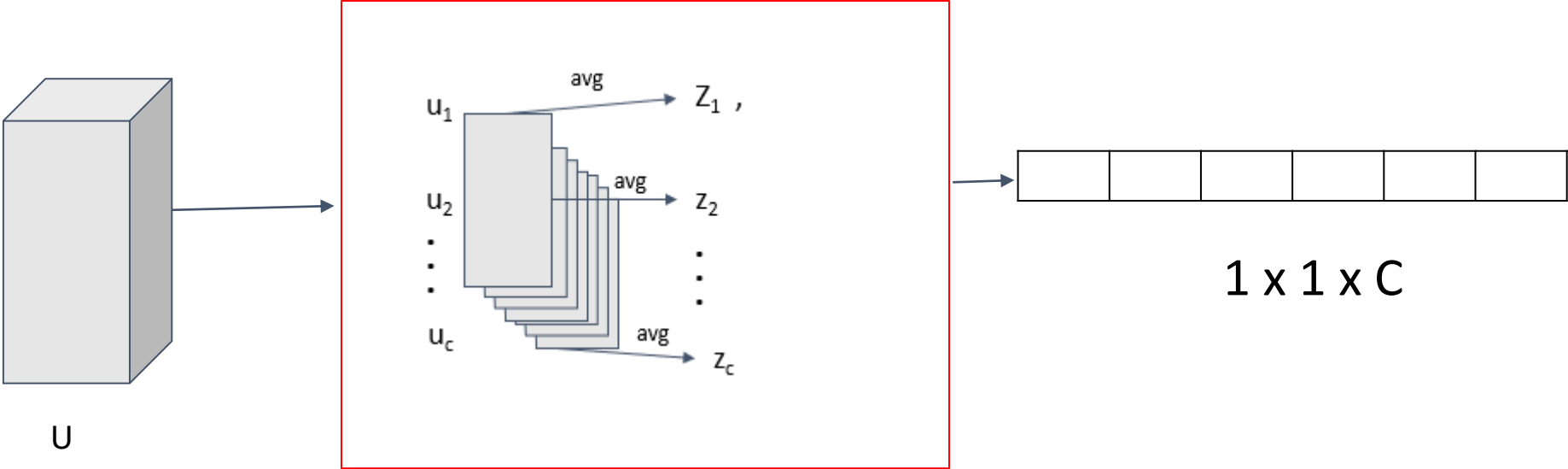
- **Squeeze operation:** Apply global average pooling to aggregate feature maps across their spatial dimensions $H \times W$ to produce a channel descriptor. It produces a $1 \times 1 \times C$ vector.
- **Excitation operation:** It applies fully-connected layers with one hidden layer using ReLU activation function on the Squeeze vector, and produces a channel weighted vector

Squeeze operation



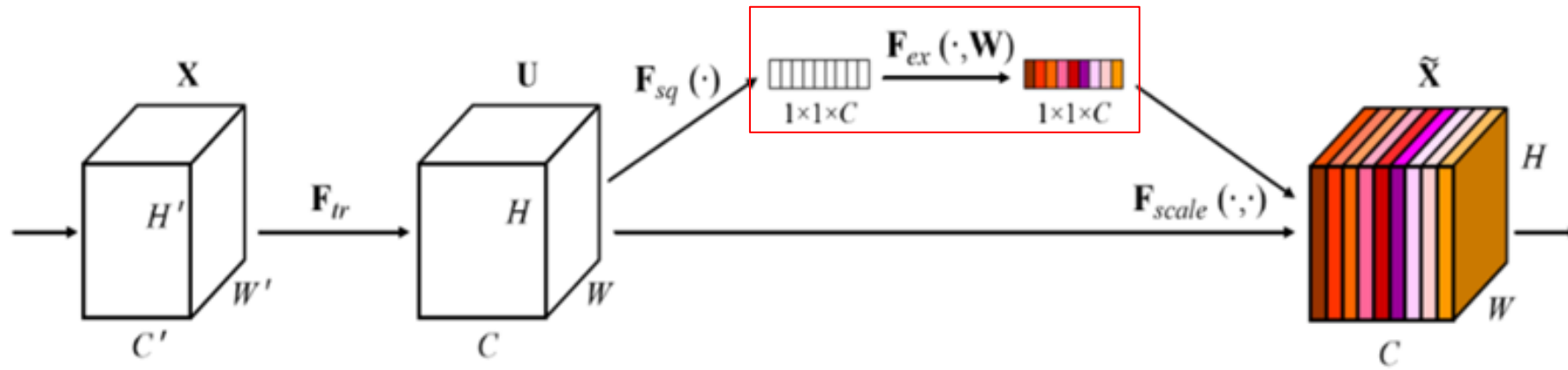
$$Z_c = F_{sq}(u_c) = \frac{1}{H \times W} \sum_{i=1}^H \sum_{j=1}^W u_c(i, j)$$

Squeeze operation



$$Z_1 = F_{sq}(u_1) = \frac{1}{H \times W} \sum_{i=1}^H \sum_{j=1}^W u_1(i, j)$$

Excitation operation



$$s = F_{ex}(z, W)$$

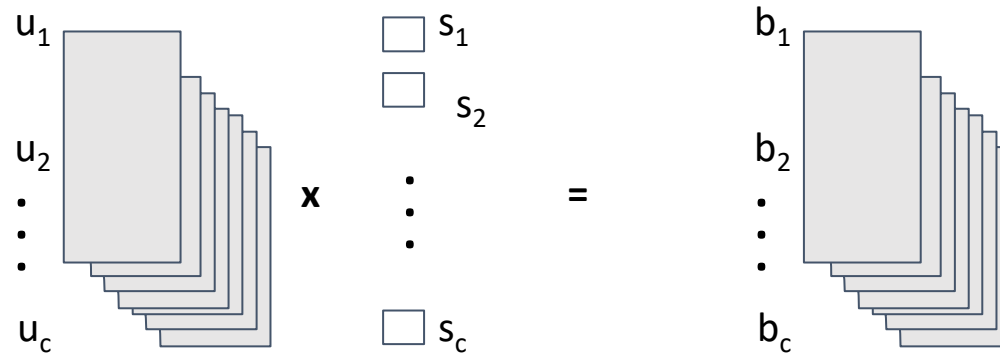
$$s = \text{sigmoid}(w_2(\text{Relu}(z, w_1)))$$

Output feature map

The final output feature map B is obtained by rescaling U with the activations s :

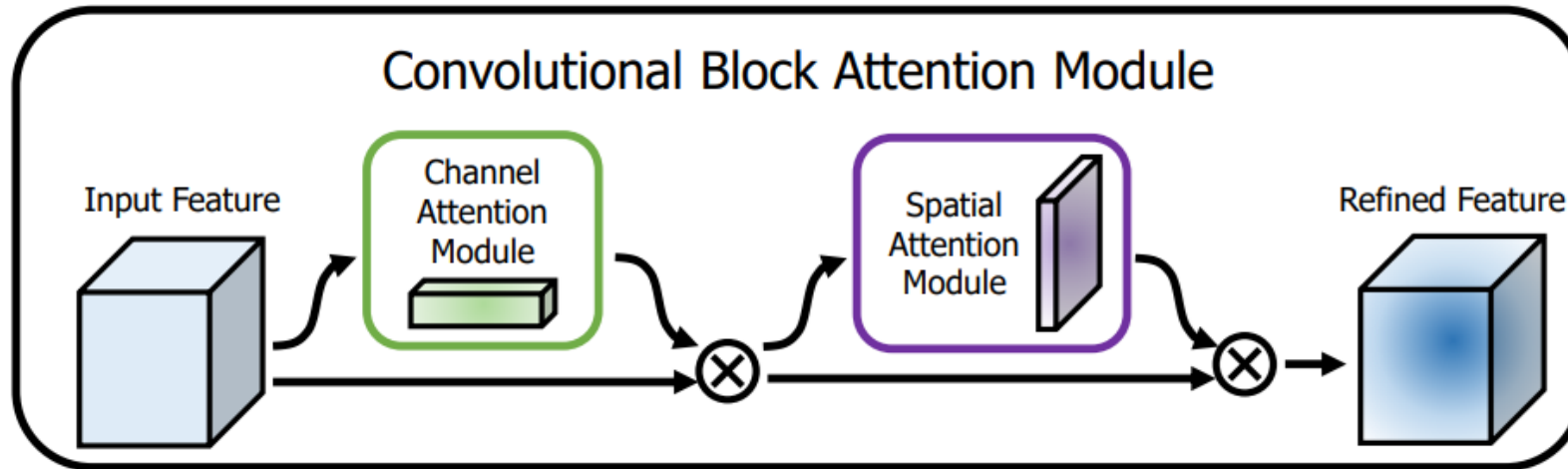
$$b_c = F_{scale}(u_c, s_c)$$

$F_{scale}(u_c, s_c)$ refers to channel-wise multiplication between the scalar s_c and the feature map u_c



Convolutional Block Attention Module (CBAM)

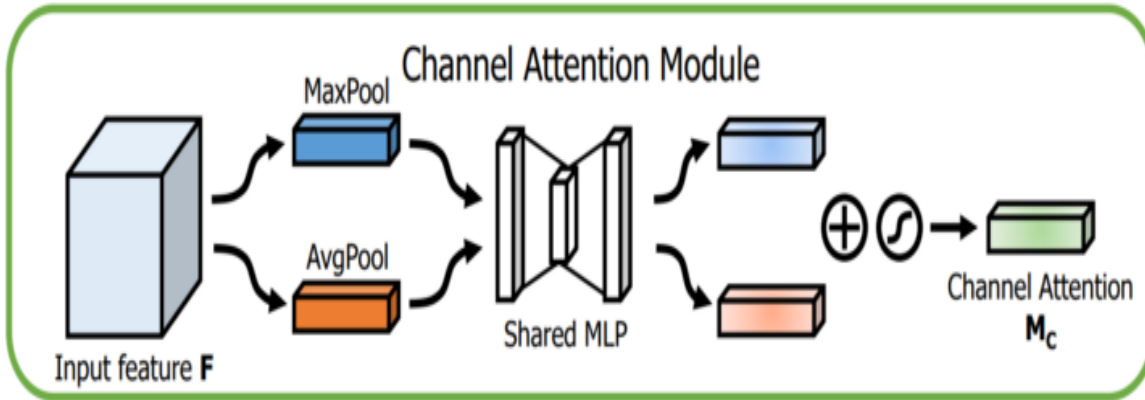
(Provide attention on both channel and spatial)



It has two components

- **Channel Attention module (CAM):** Which channel is important
- **Spatial Attention module (SAM):** Which region in the feature map is important

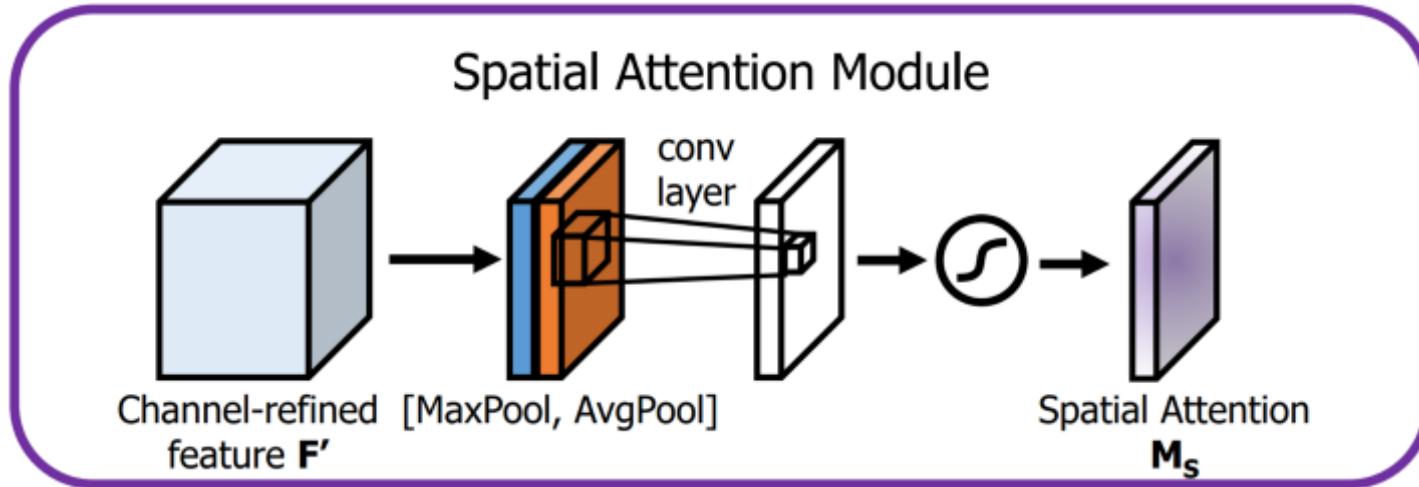
Channel Attention module (CAM)



It resembles Squeeze and Excite (SE) model.

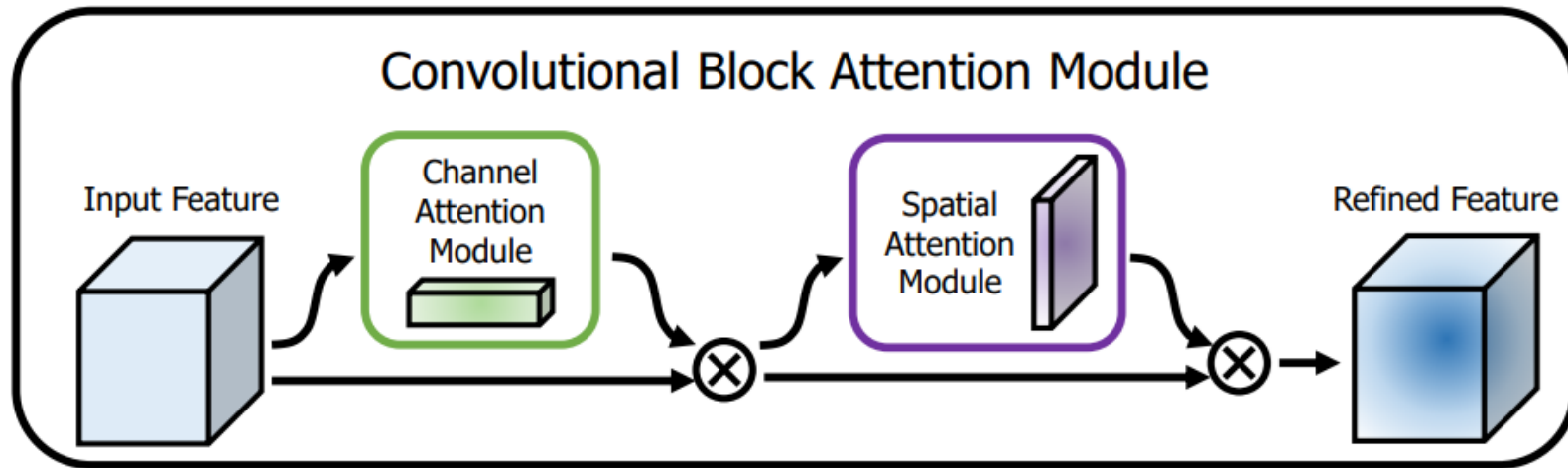
$$\begin{aligned}\mathbf{M}_c(\mathbf{F}) &= \sigma(MLP(AvgPool(\mathbf{F})) + MLP(MaxPool(\mathbf{F}))) \\ &= \sigma(\mathbf{W}_1(\mathbf{W}_0(\mathbf{F}_{avg}^c)) + \mathbf{W}_1(\mathbf{W}_0(\mathbf{F}_{max}^c))),\end{aligned}$$

Spatial Attention module (SAM)



- Apply Global Average Pool and Global Max Pooling across channels.
- Concatenate and pass it through a small convolutional block of 7x7 kernel size.

$$\begin{aligned}\mathbf{M}_s(\mathbf{F}) &= \sigma(f^{7 \times 7}([AvgPool(\mathbf{F}); MaxPool(\mathbf{F})])) \\ &= \sigma(f^{7 \times 7}([\mathbf{F}_{avg}^s; \mathbf{F}_{max}^s])),\end{aligned}$$



$$\mathbf{F}' = \mathbf{M}_c(\mathbf{F}) \otimes \mathbf{F},$$
$$\mathbf{F}'' = \mathbf{M}_s(\mathbf{F}') \otimes \mathbf{F}',$$