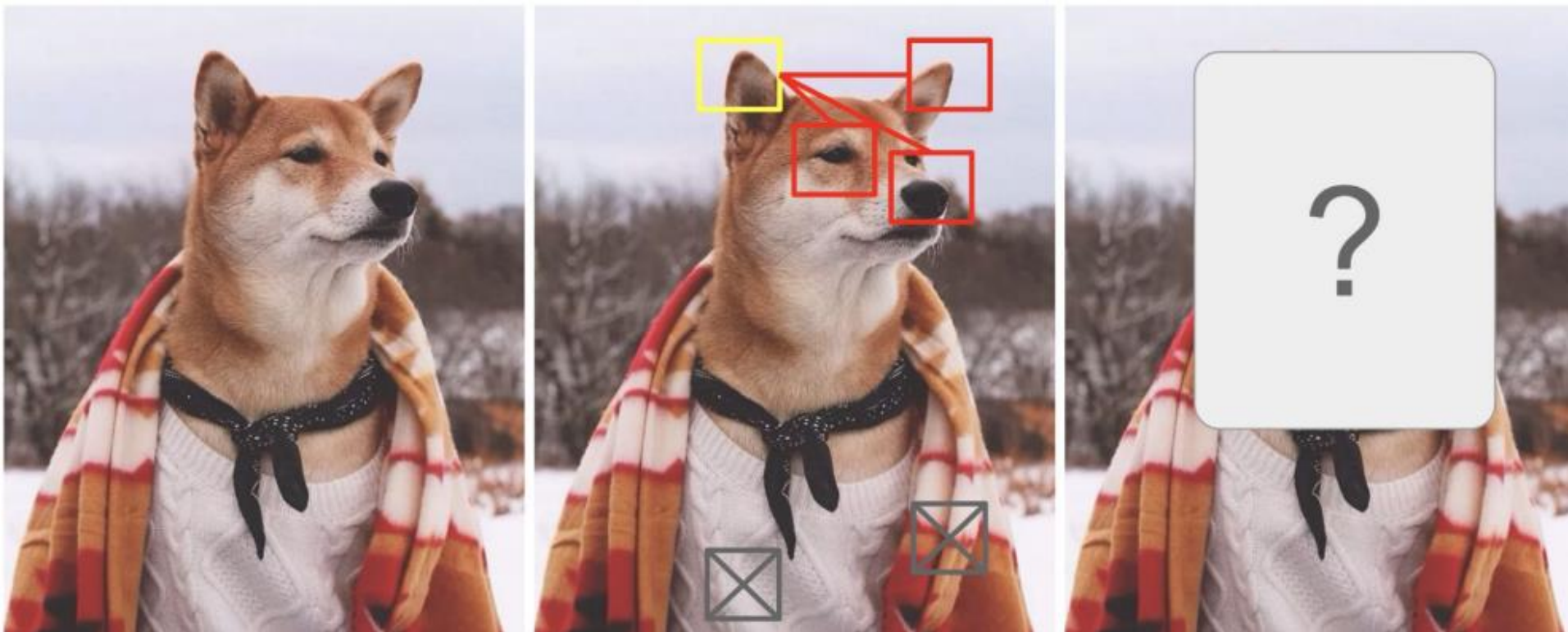


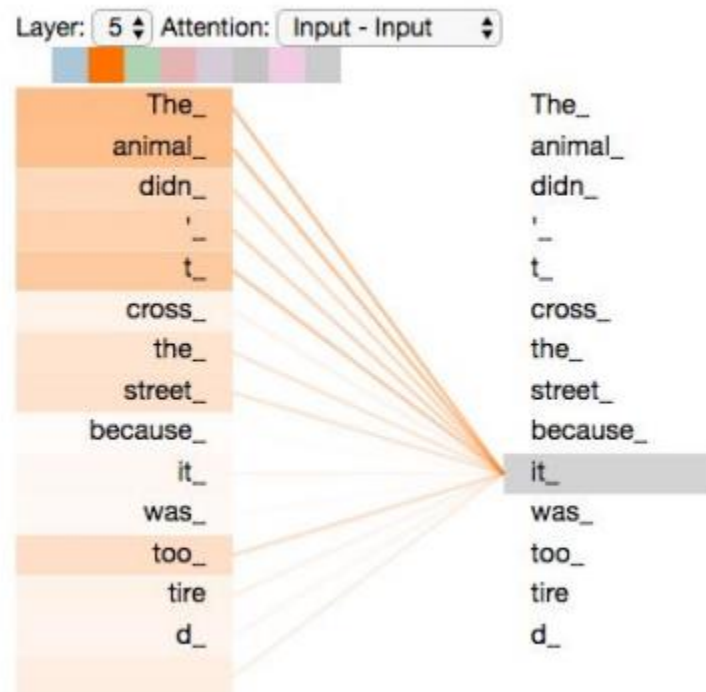
Transformers

▪Context



Context

- “The animal didn't cross the street because it was too tired”
- What is “it”?

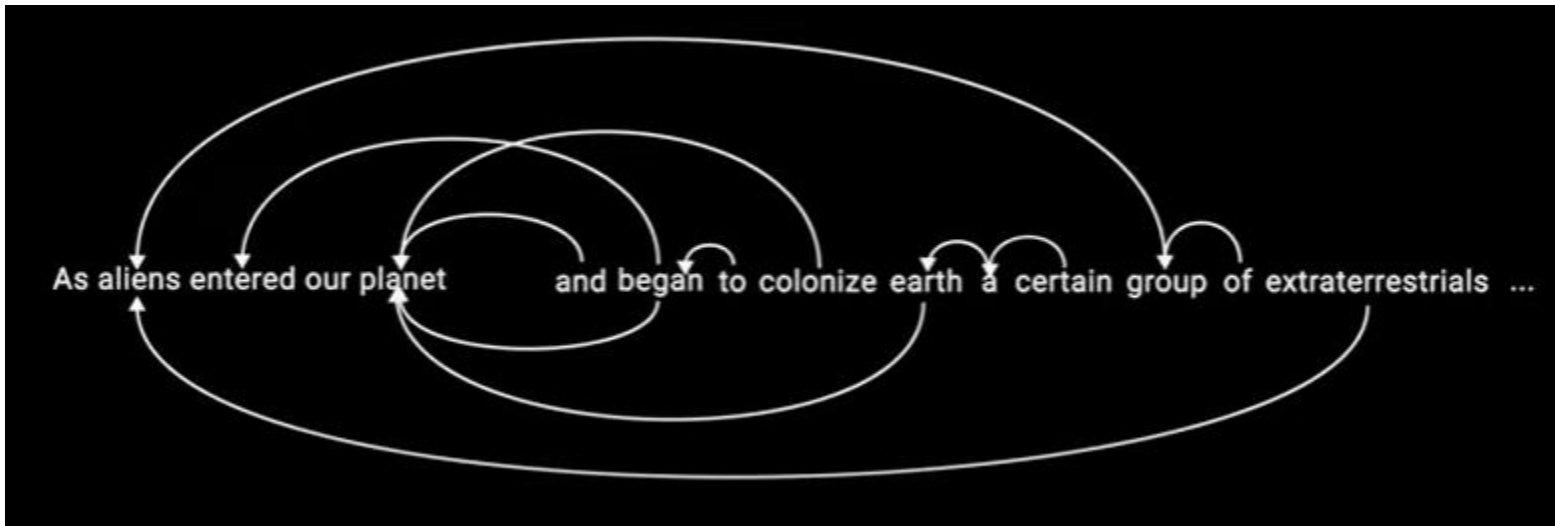


Context



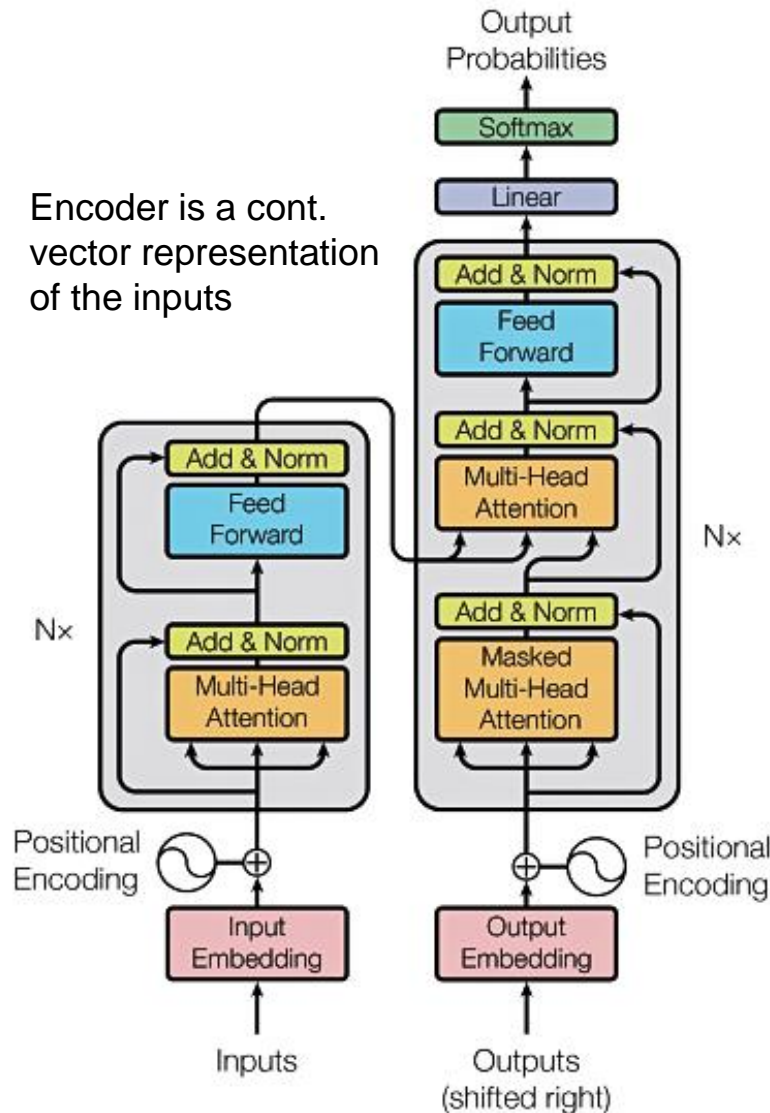
Attention

- Attention mechanism has an infinite reference window.

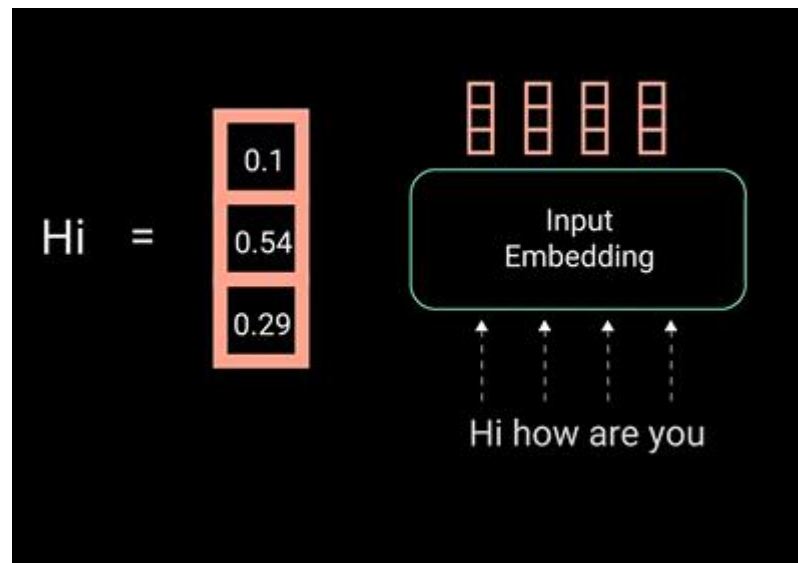


Attention is all you need

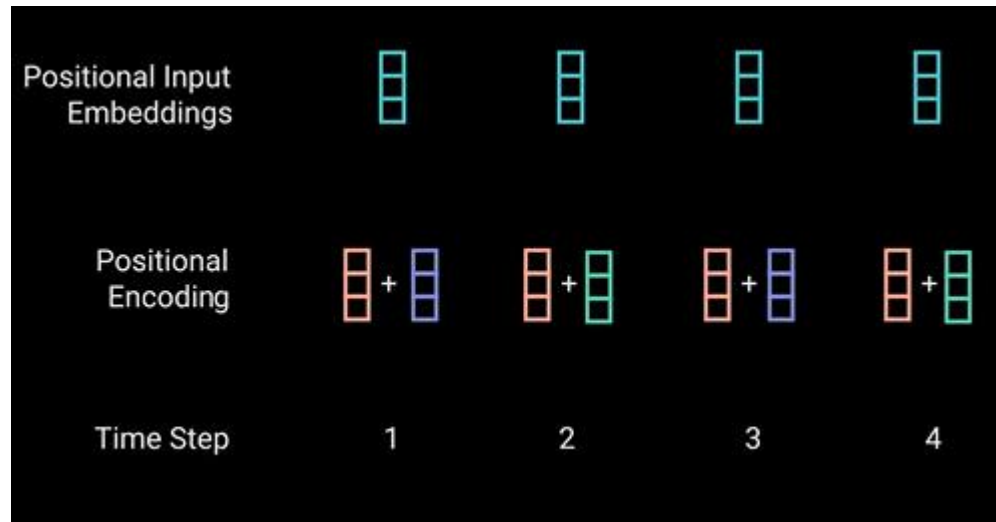
Encoder is a cont.
vector representation
of the inputs



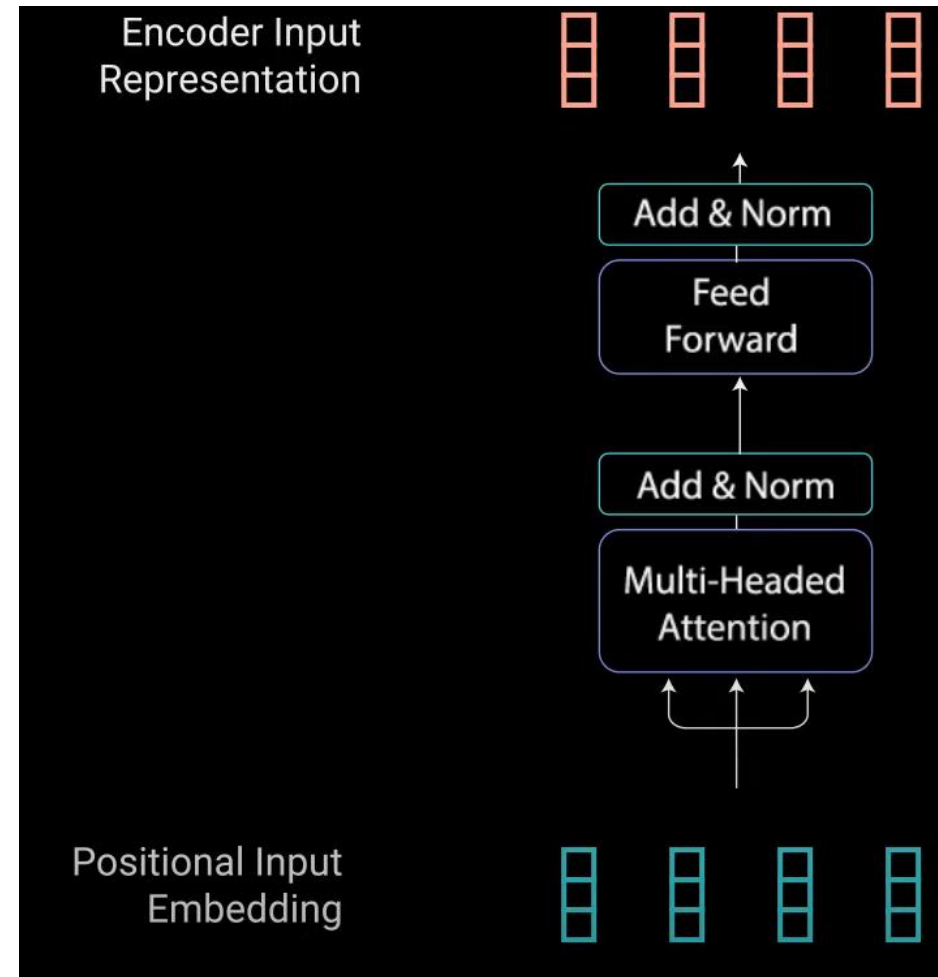
1. Input Embedding



2. Positional Encoding

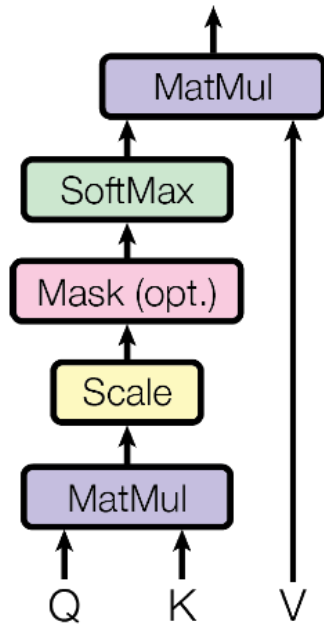


3-4 Encoder layer

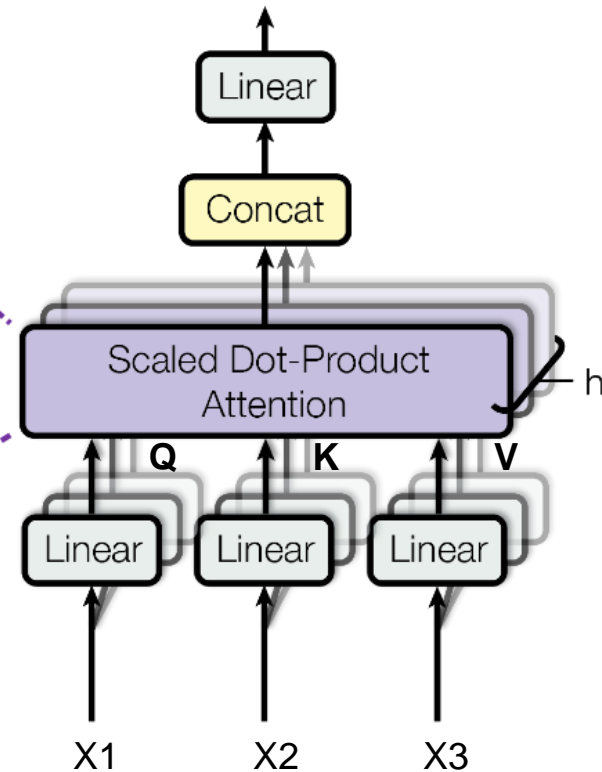


Multi-Head Attention

Scaled Dot-Product Attention

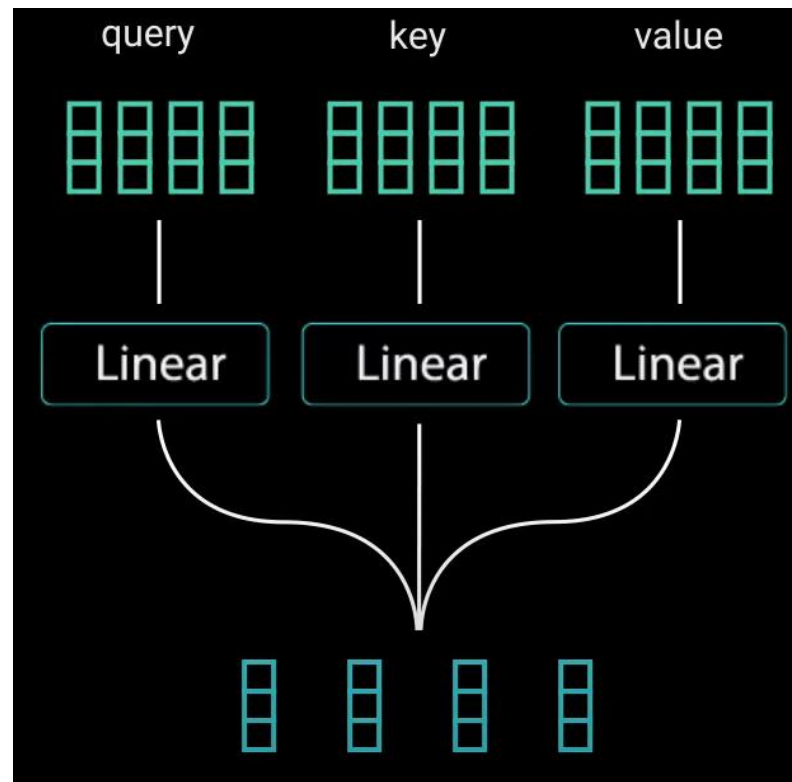


Multi-Head Attention



Self Attention

- Associates each individual word in the input to other words in the input.



Self attention details

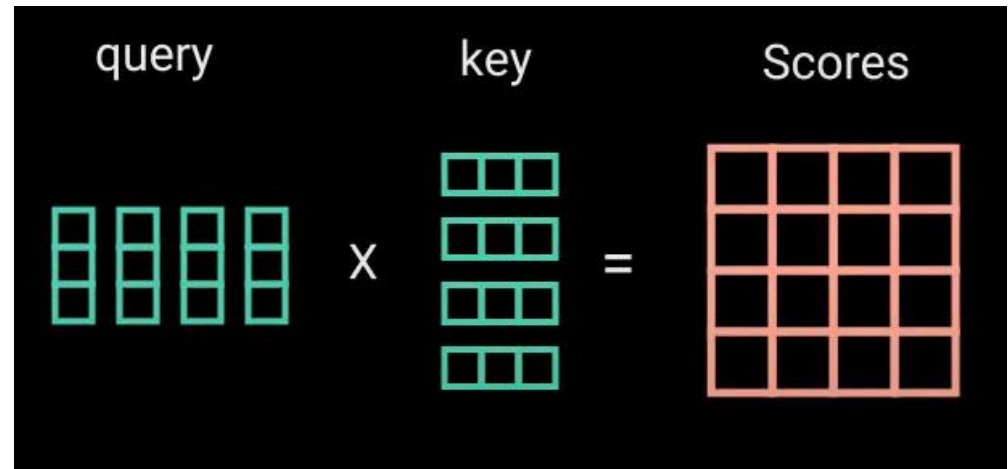
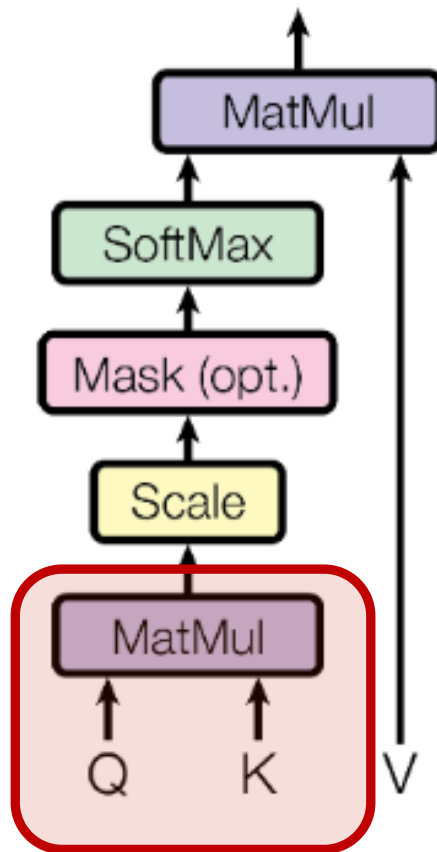
- Given an input sequence, X
- Project to Queries, Keys and Values using linear transforms.

$$Q = W^Q X$$

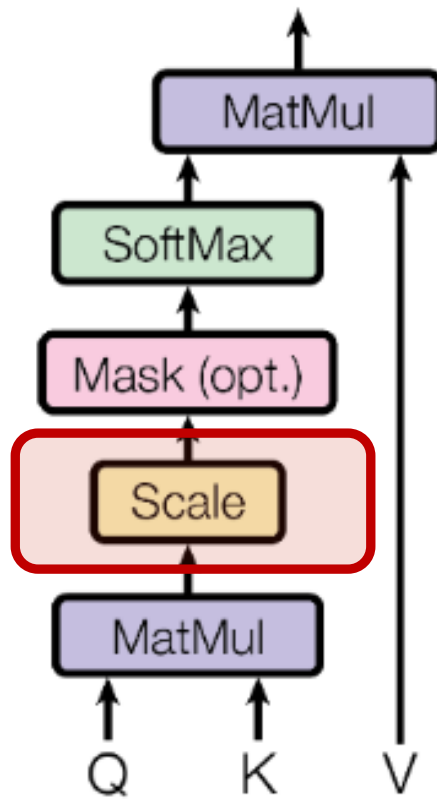
$$K = W^K X$$

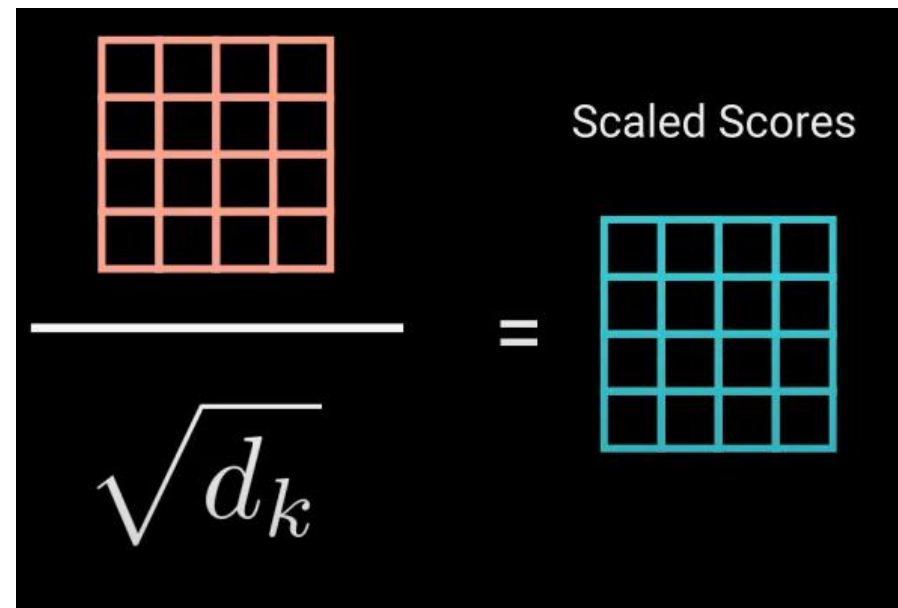
$$V = W^V X$$





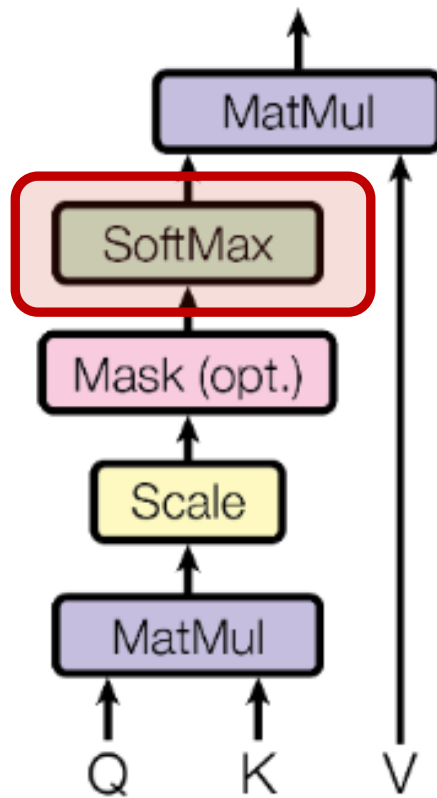
	Hi	how	are	you
Hi	98	27	10	12
how	27	89	31	67
are	10	31	91	54
you	12	67	54	92





$$\frac{\text{Grid of Scores}}{\sqrt{d_k}} = \text{Scaled Scores}$$

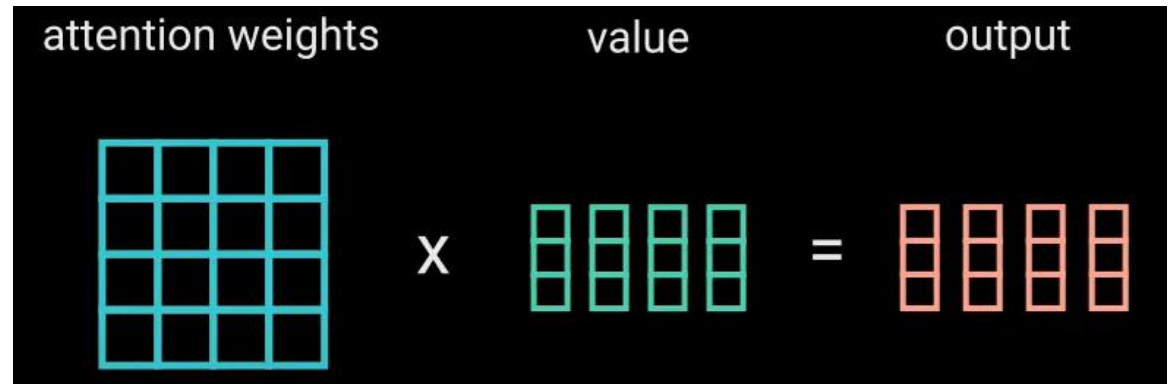
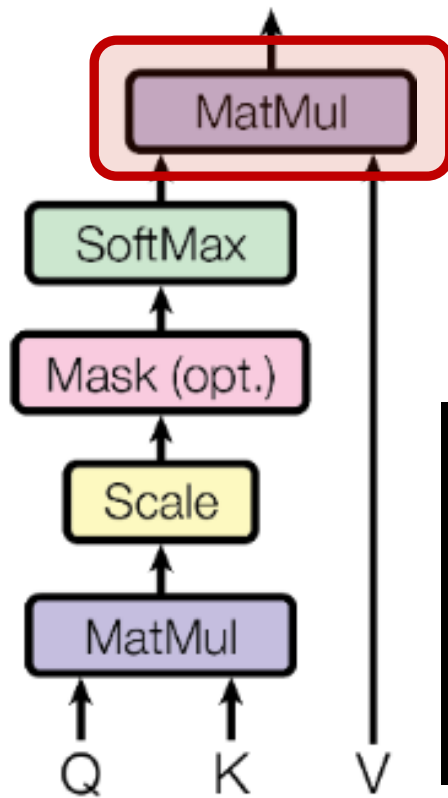
The equation shows a 4x4 grid of scores (represented by an orange grid) divided by the square root of the key dimension d_k (represented by the symbol $\sqrt{d_k}$) to produce the Scaled Scores (represented by a cyan grid).

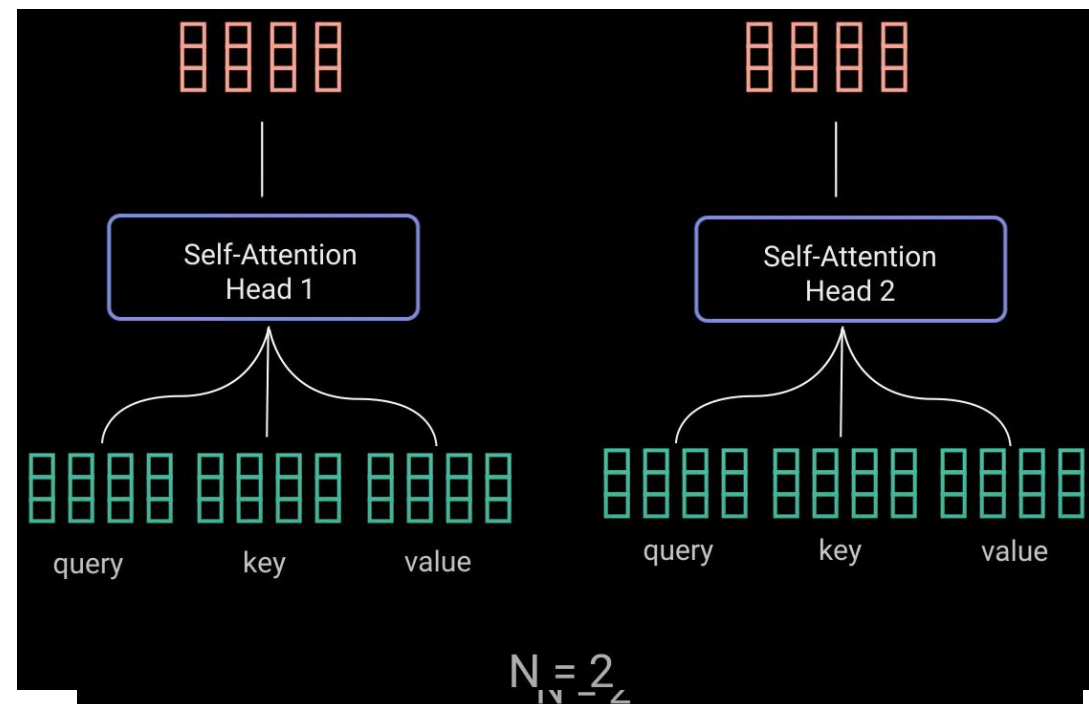
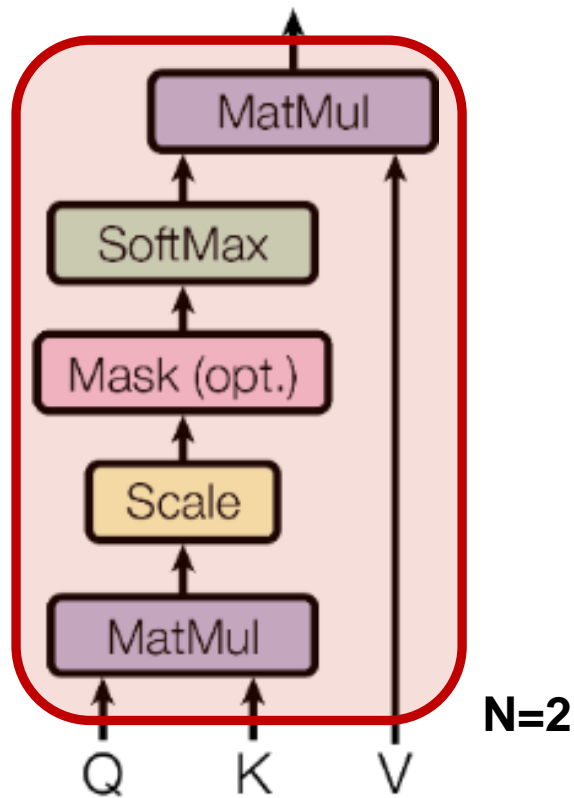


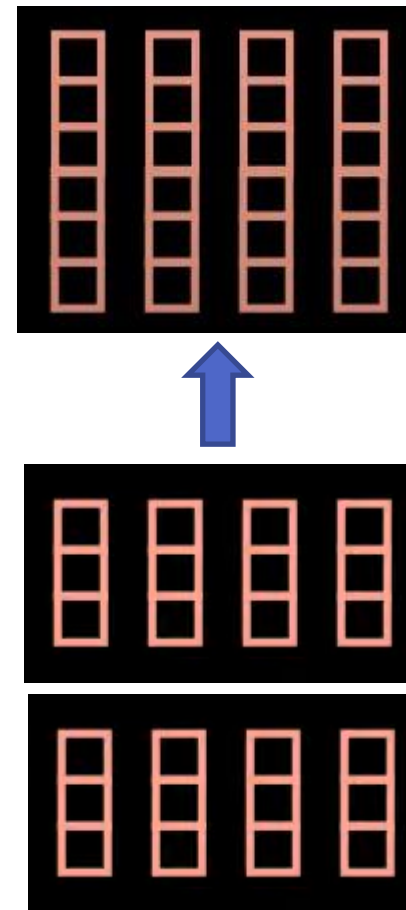
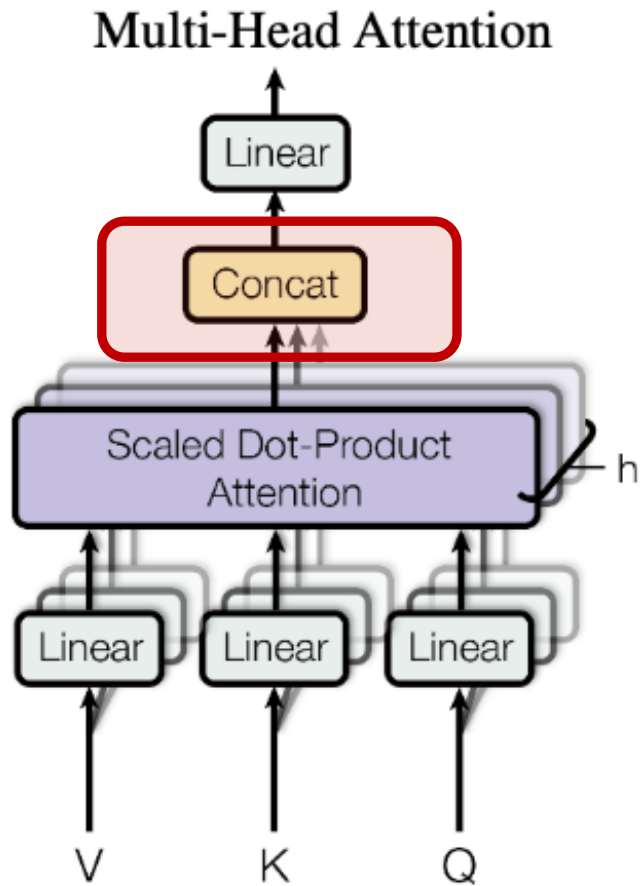
Softmax(

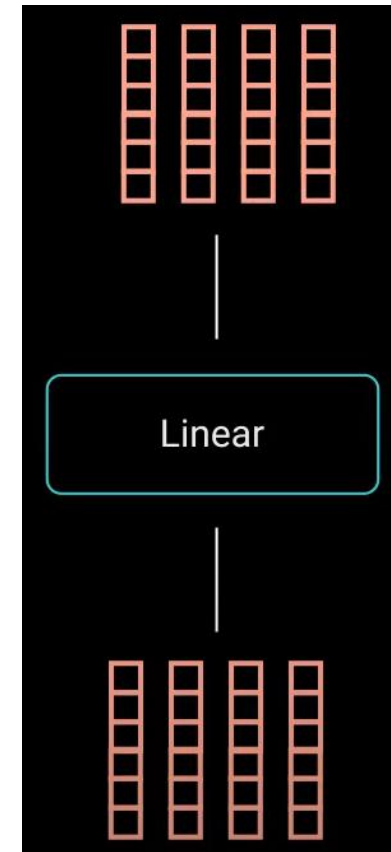
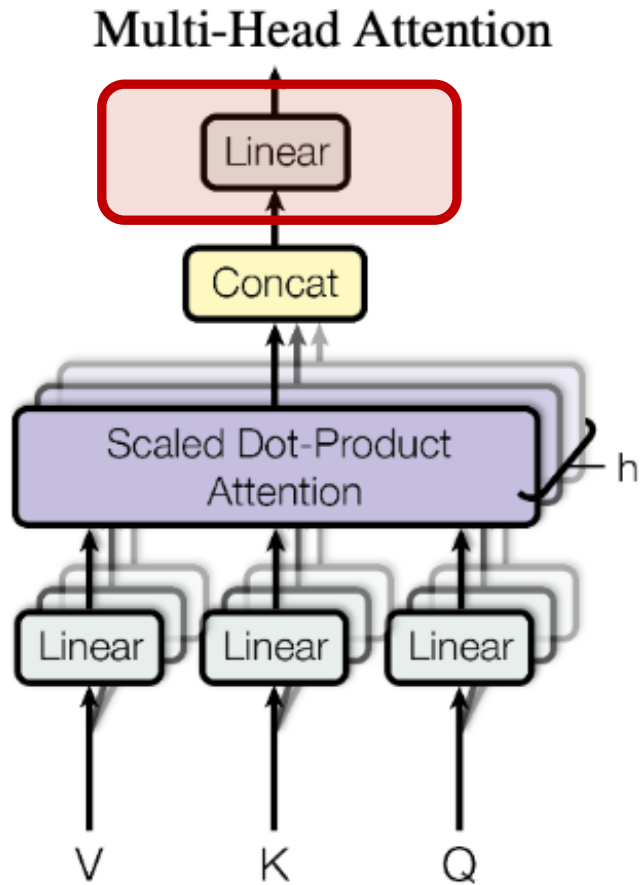
) =

	Hi	how	are	you
Hi	0.7	0.1	0.1	0.1
how	0.1	0.6	0.2	0.1
are	0.1	0.3	0.6	0.1
you	0.1	0.3	0.3	0.3

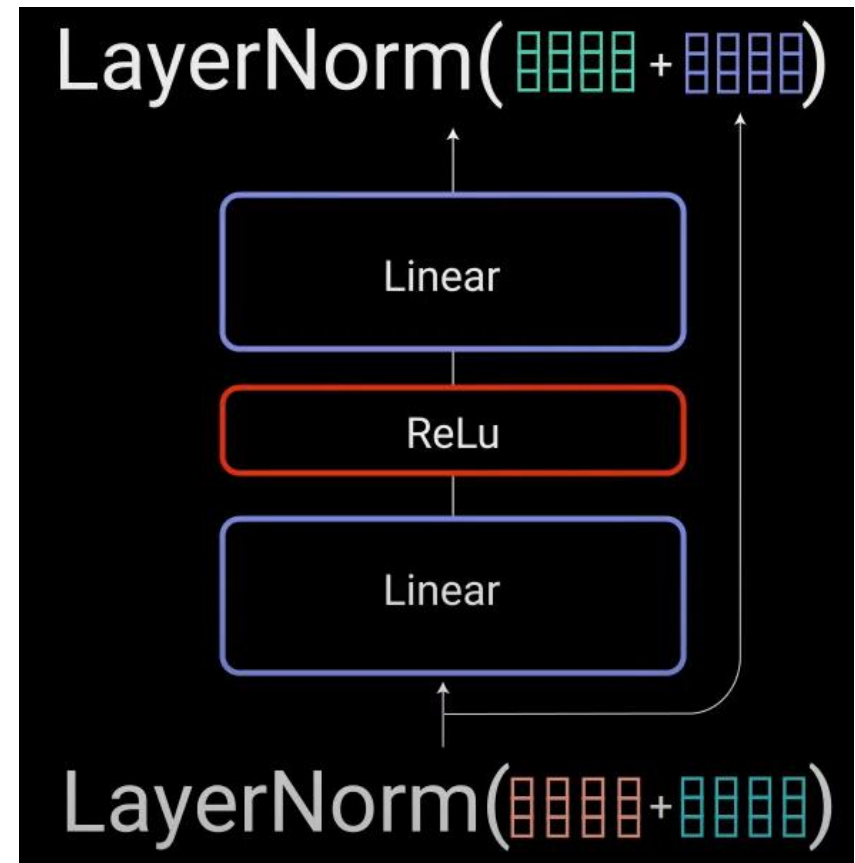
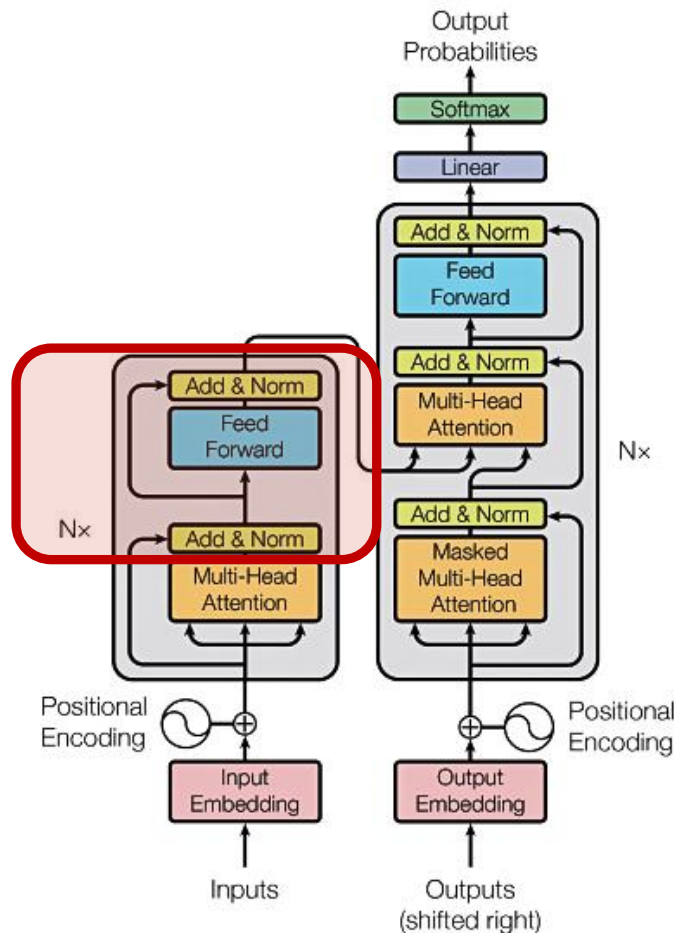




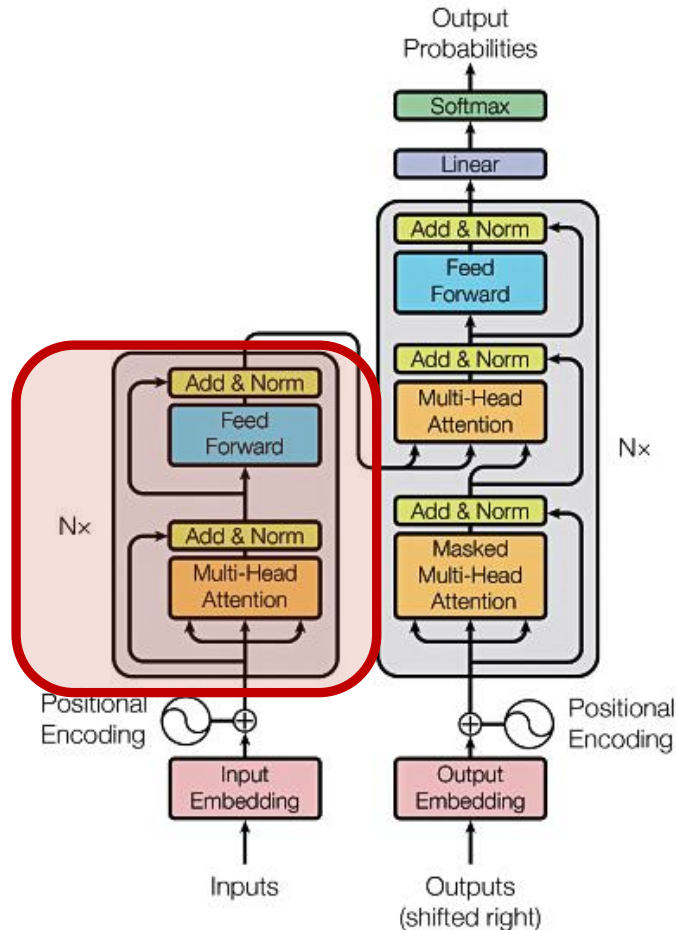




Residual Connection, Layer Normalization and Pointwise FW NW



Transformer Encoder



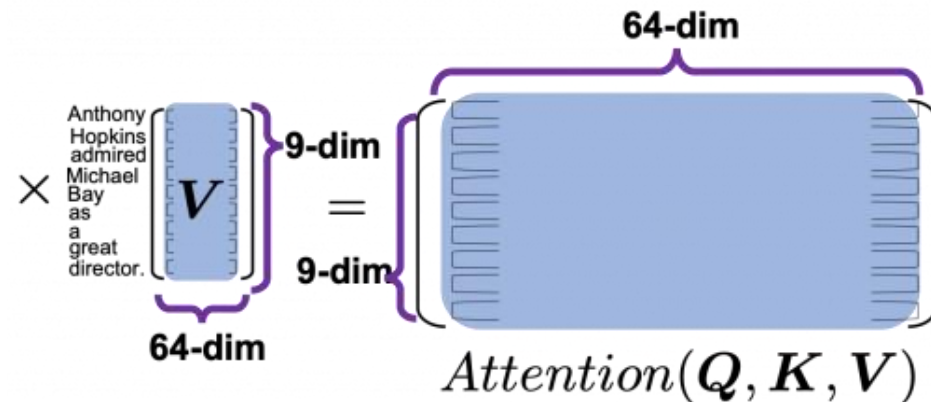
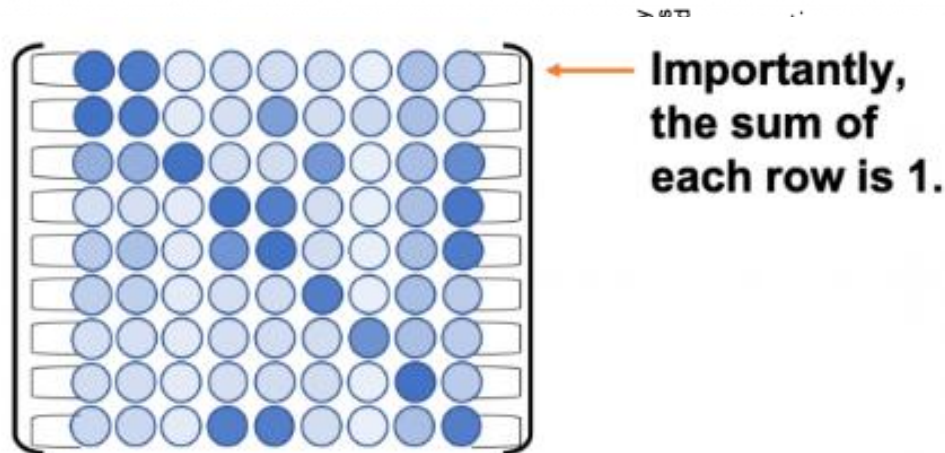
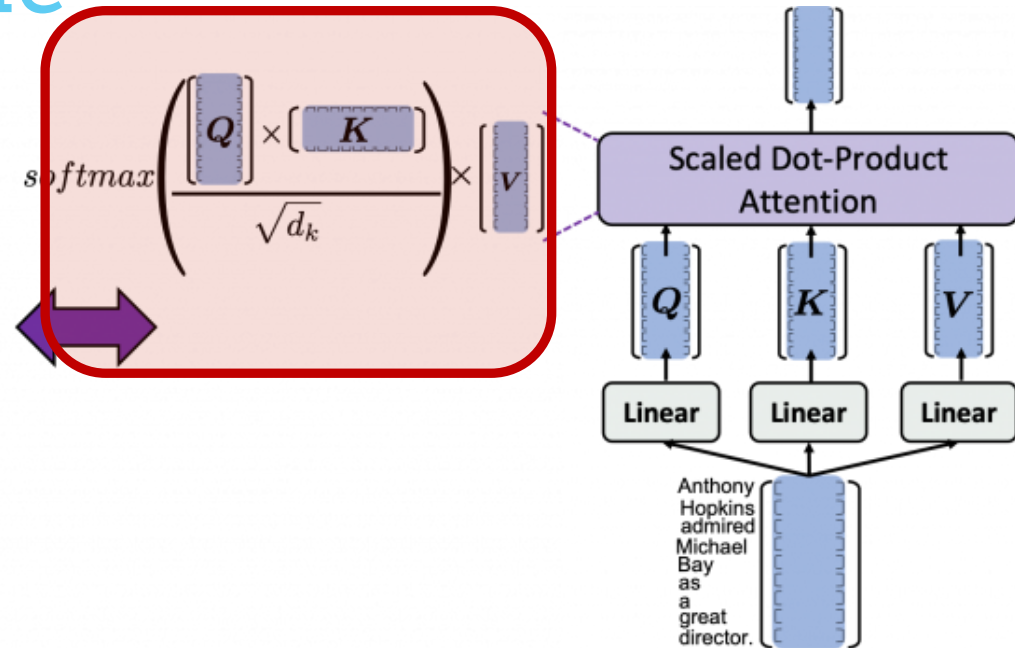
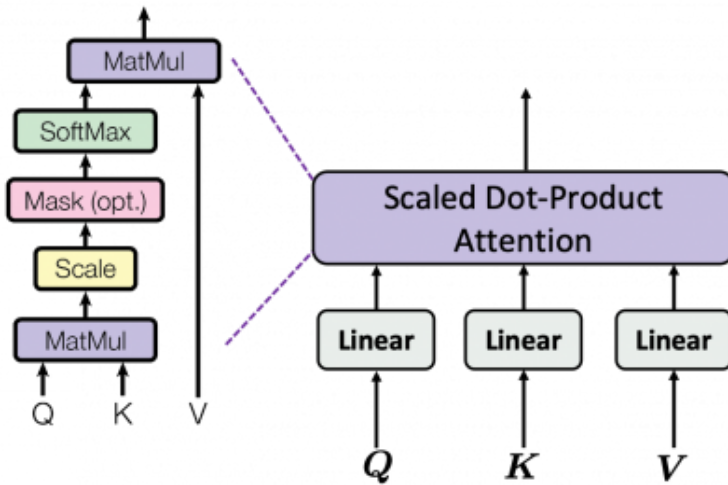
Transformer Encoder

N Times

Transformer Encoder

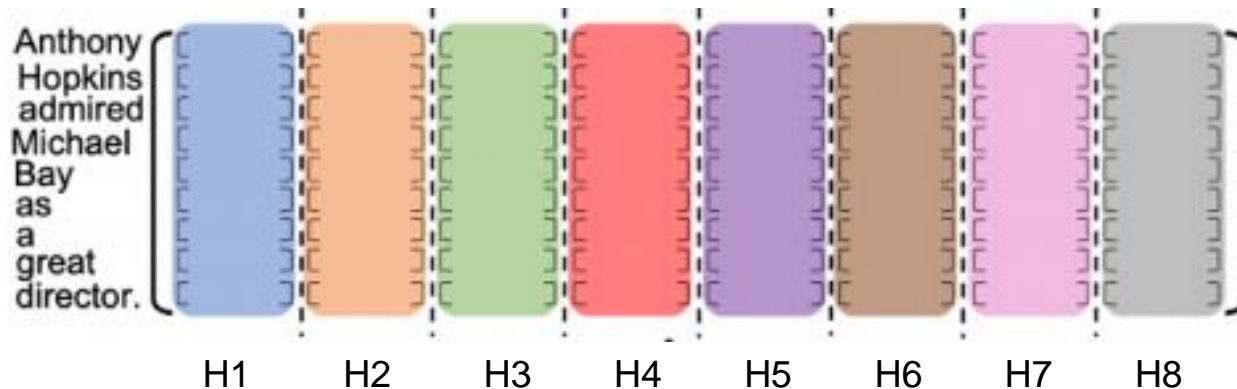
Transformer Encoder

Encoder Example

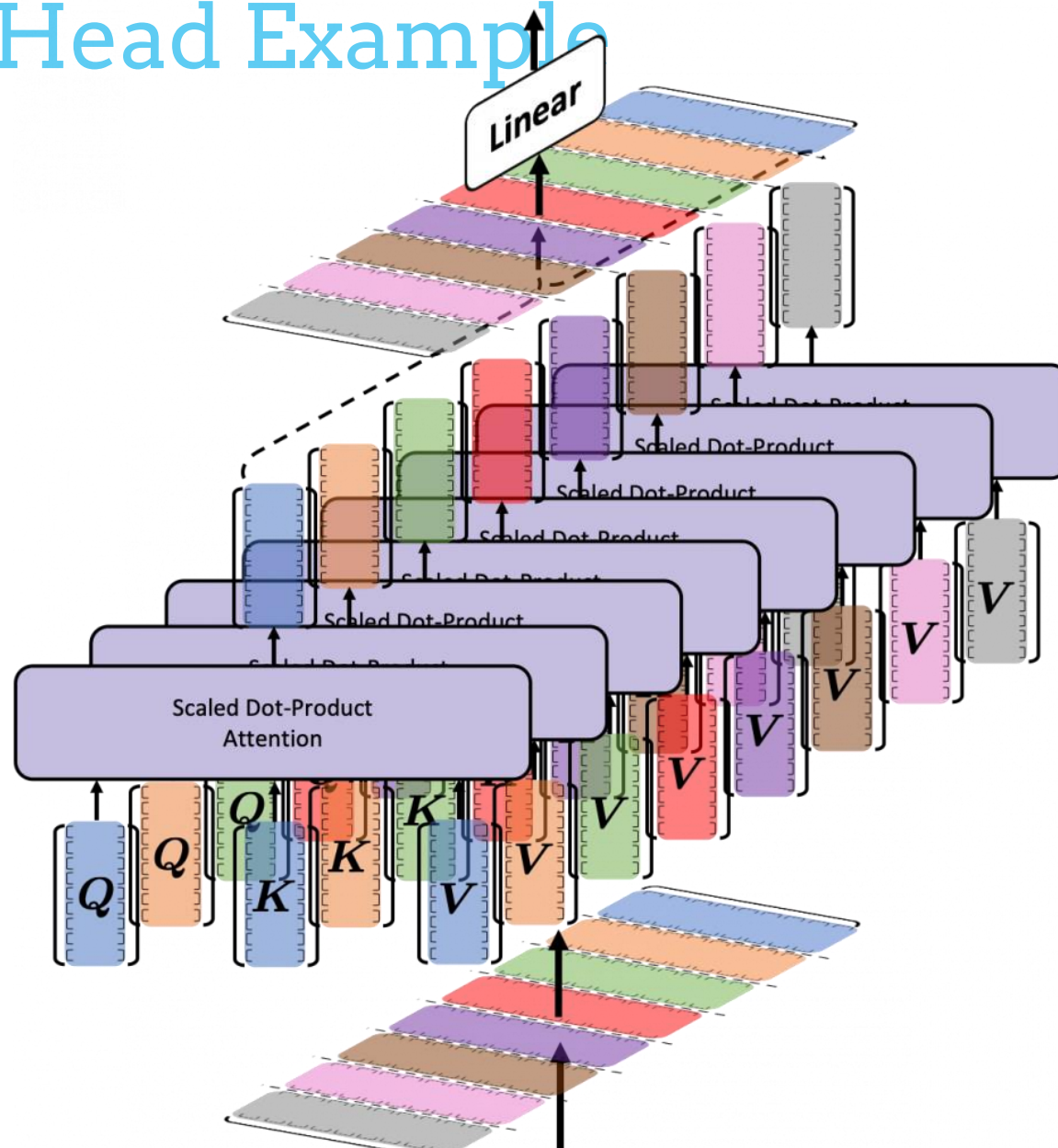


Multi Head Example

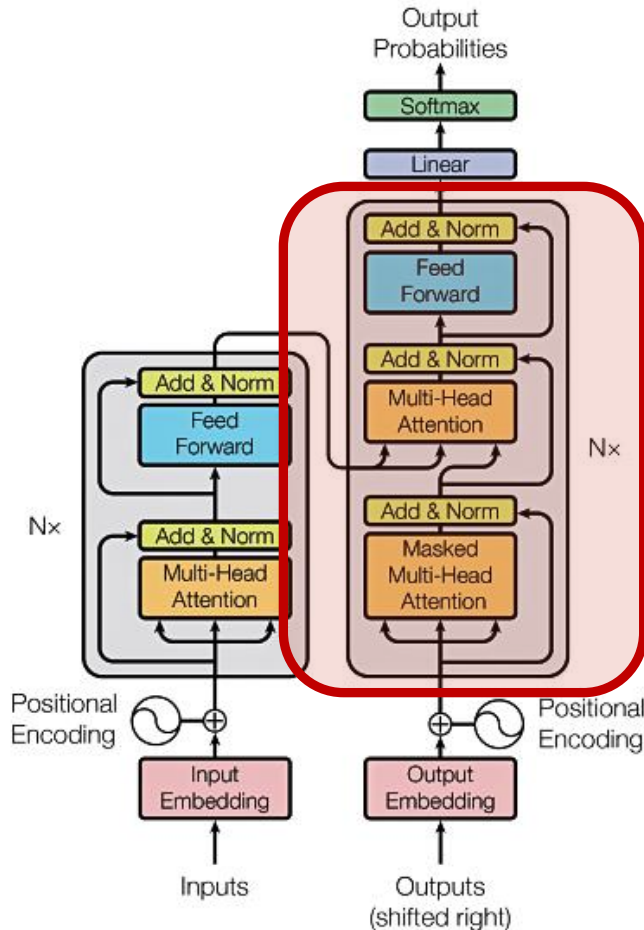
- No of tokens=9
- Token encoded dimension =512
- Split each token into 8 equal parts of 64 dimension



Multi Head Example

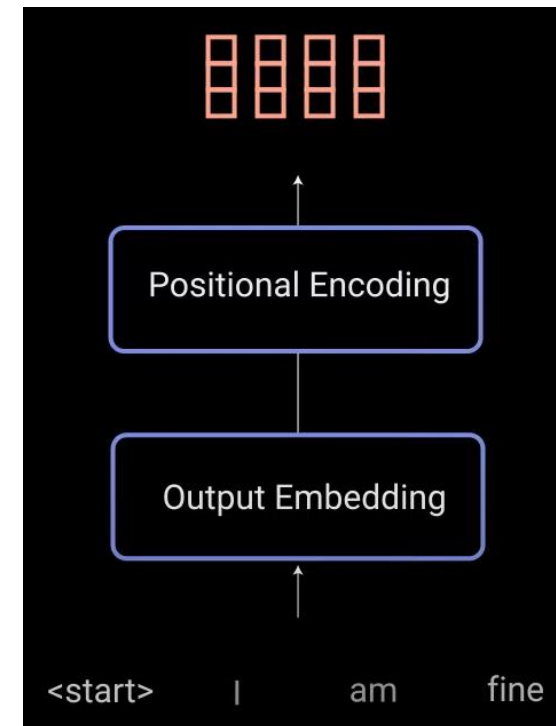
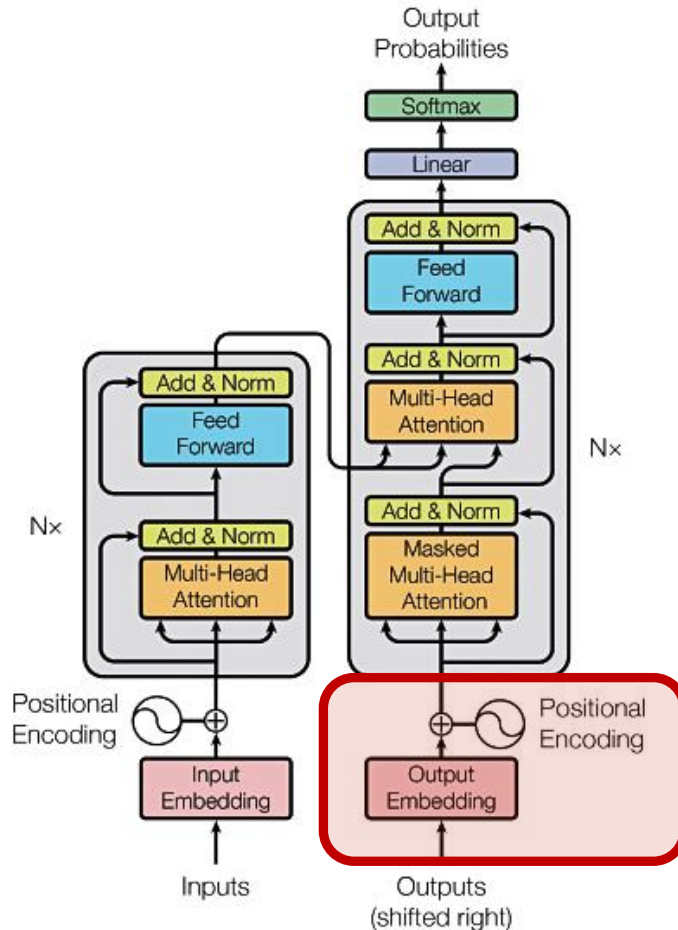


Transformer Decoder

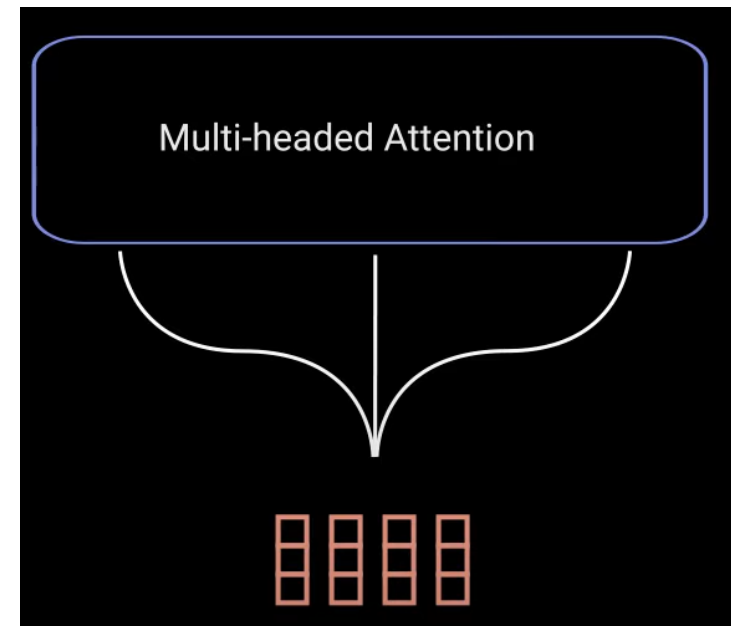
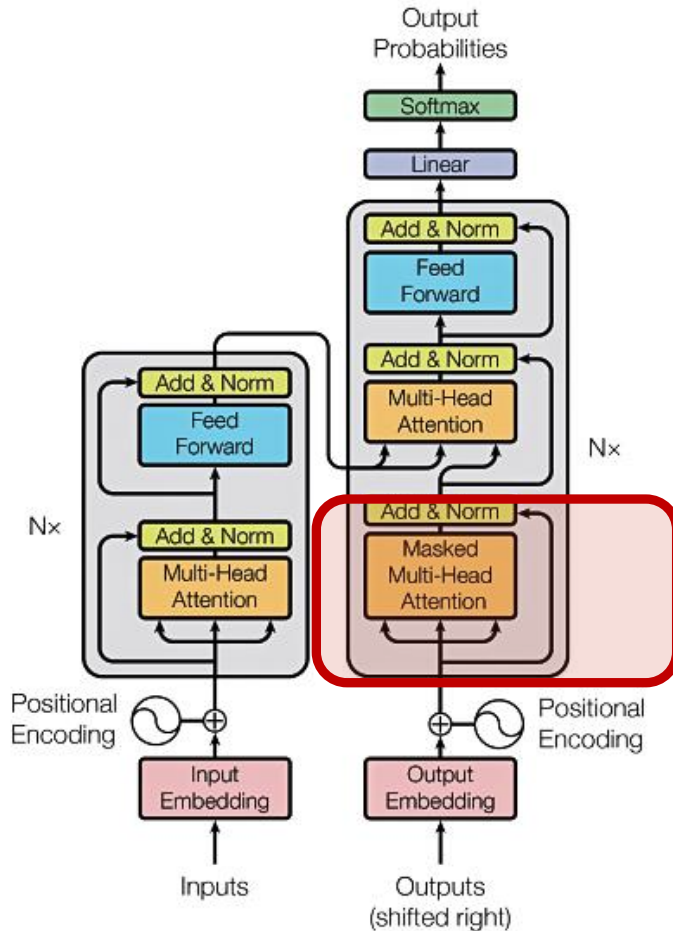


- Decoder is to generate text
- It is auto regressive – Takes previous outputs as inputs
- Stops generation when `<end>` token is generated.

Output Embedding and positional Encoding



Decoder MH Attention1



Look ahead Mask

- Prevents decoder from looking at future tokens

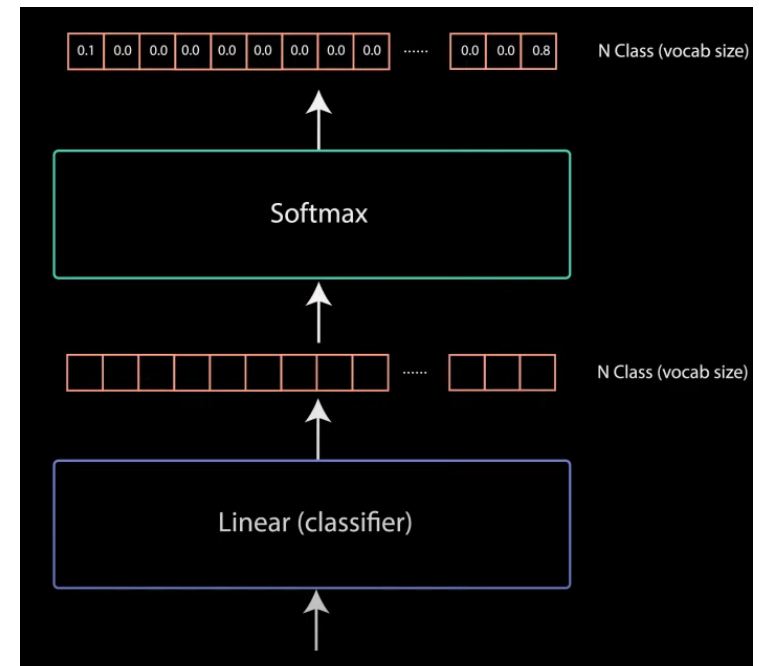
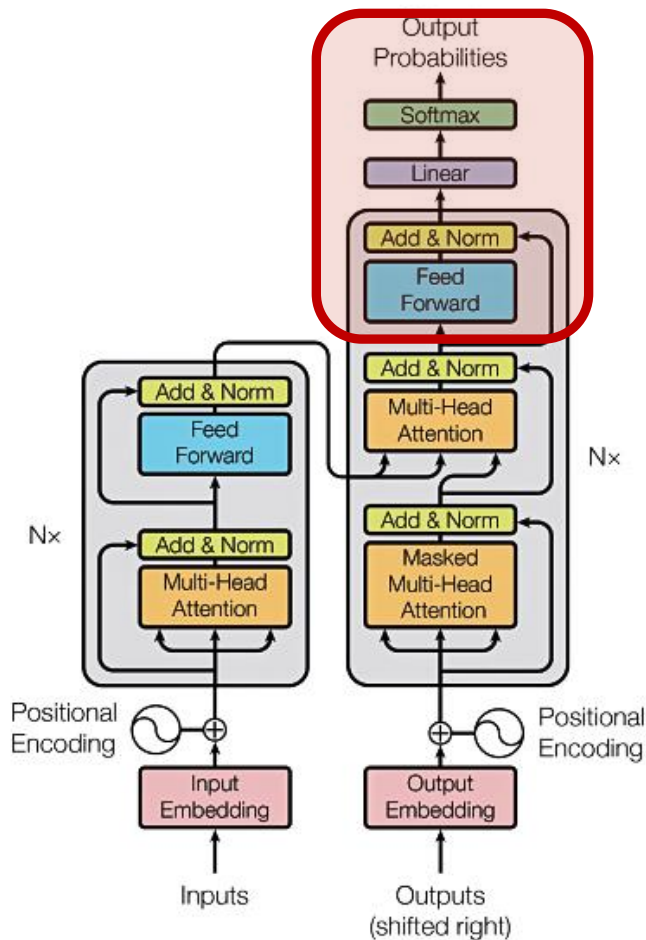
Scaled Scores Look-Ahead Mask Masked Scores

$$\begin{bmatrix} 0.7 & 0.1 & 0.1 & 0.1 \\ 0.1 & 0.6 & 0.2 & 0.1 \\ 0.1 & 0.3 & 0.6 & 0.1 \\ 0.1 & 0.3 & 0.3 & 0.3 \end{bmatrix} + \begin{bmatrix} 0 & -\text{inf} & -\text{inf} & -\text{inf} \\ 0 & 0 & -\text{inf} & -\text{inf} \\ 0 & 0 & 0 & -\text{inf} \\ 0 & 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 0.7 & -\text{inf} & -\text{inf} & -\text{inf} \\ 0.1 & 0.6 & -\text{inf} & -\text{inf} \\ 0.1 & 0.3 & 0.6 & -\text{inf} \\ 0.1 & 0.3 & 0.3 & 0.3 \end{bmatrix}$$

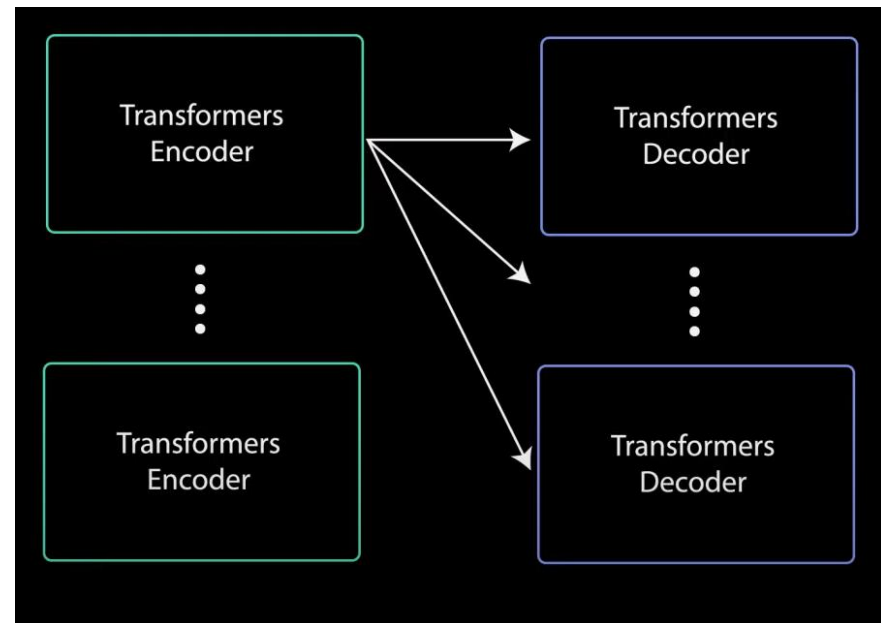
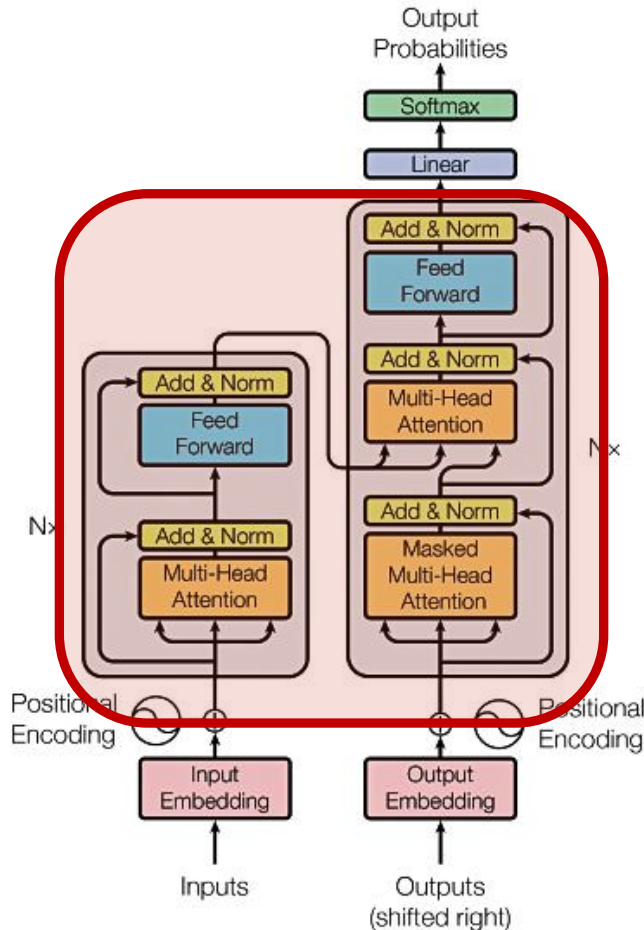
Softmax($\begin{bmatrix} 0.7 & -\text{inf} & -\text{inf} & -\text{inf} \\ 0.1 & 0.6 & -\text{inf} & -\text{inf} \\ 0.1 & 0.3 & 0.6 & -\text{inf} \\ 0.1 & 0.3 & 0.3 & 0.3 \end{bmatrix}$) =

	<start>	I	am	fine
<start>	1	0	0	0
I	0.37	0.62	0	0
am	0.26	0.31	0.43	0
fine	0.21	0.26	0.26	0.26

Linear Classifier



Transformer – Encoder & Decoder



Transformers, GPT-2, and BERT

1. A transformer uses an **encoder stack** to model input, and uses **decoder stack** to model output (using input information from encoder side)
2. If we do not have input, we just want to model the “next word”, we can get rid of the encoder side of a transformer and output “next word” one by one. This gives us **GPT**
3. If we are only interested in training a language model for the input for some other tasks, then we do not need the decoder of the transformer, that gives us **BERT**

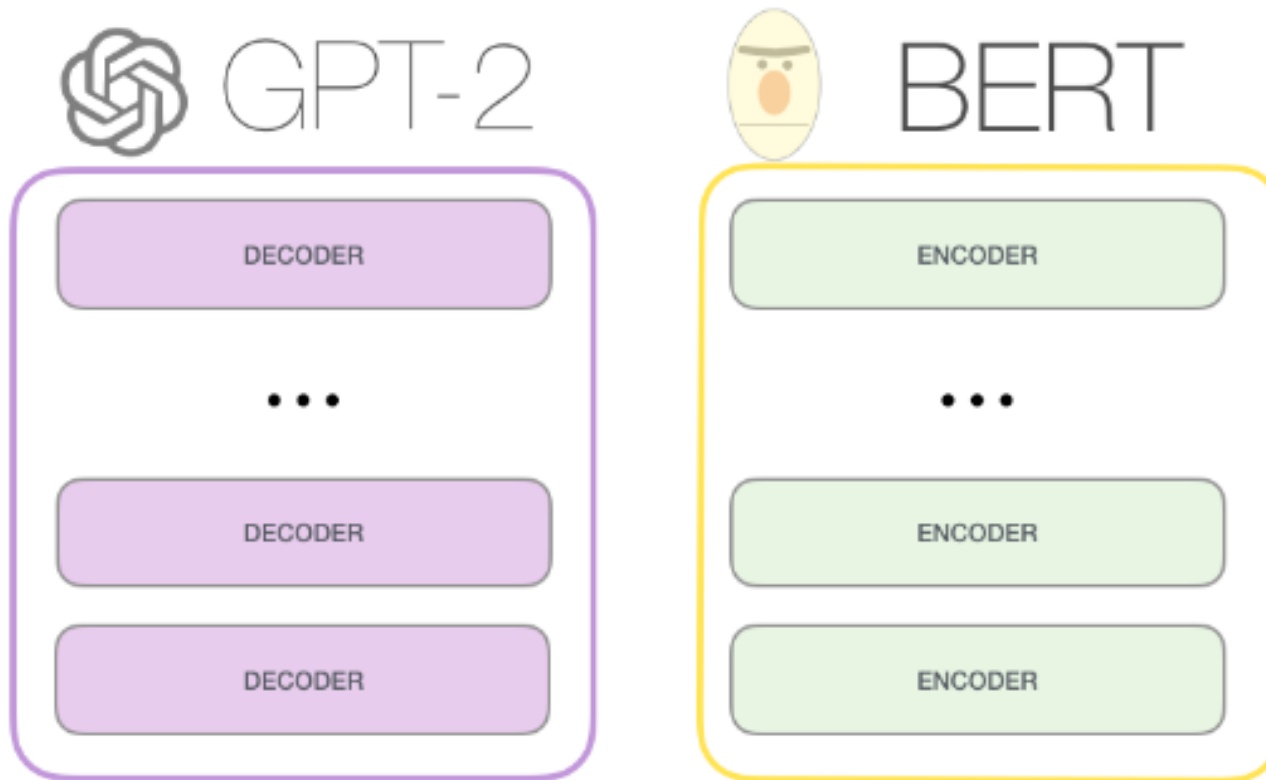
Training a Transformer

- Transformers typically use semi-supervised learning with
 - Unsupervised pretraining over a very large dataset of general text
 - Followed by supervised **fine-tuning** over a focused data set of inputs and outputs for a particular task
- Tasks for pretraining and fine-tuning commonly include:
 - language modeling
 - next-sentence prediction (aka completion)
 - question answering
 - reading comprehension
 - sentiment analysis
 - paraphrasing

Pre-trained Models

- Since training a model requires huge datasets of text and significant computation, researchers often use common pretrained models
- Examples (~ December 2021) include
 - ▶ Google's [BERT](#) model
 - ▶ Huggingface's various [Transformer models](#)
 - ▶ OpenAI's and [GPT-3 models](#)

GPT-2 & BERT

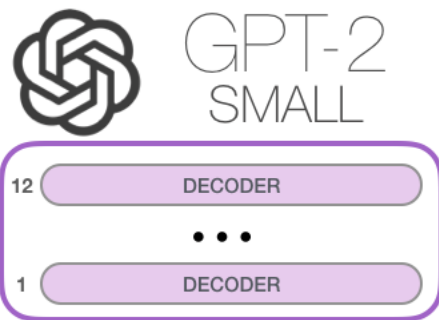


GPT - Variants

GPT released June 2018

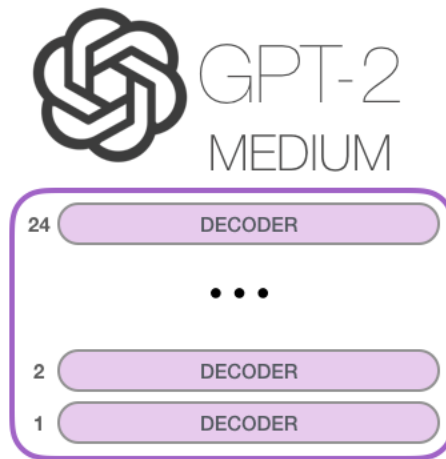
GPT-2 released Nov. 2019 with 1.5B parameters

GPT-3 released in 2020 with 175B parameters



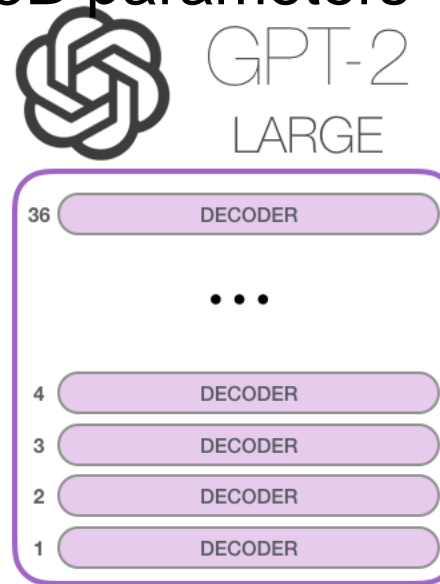
Model Dimensionality: 768

117M parameters



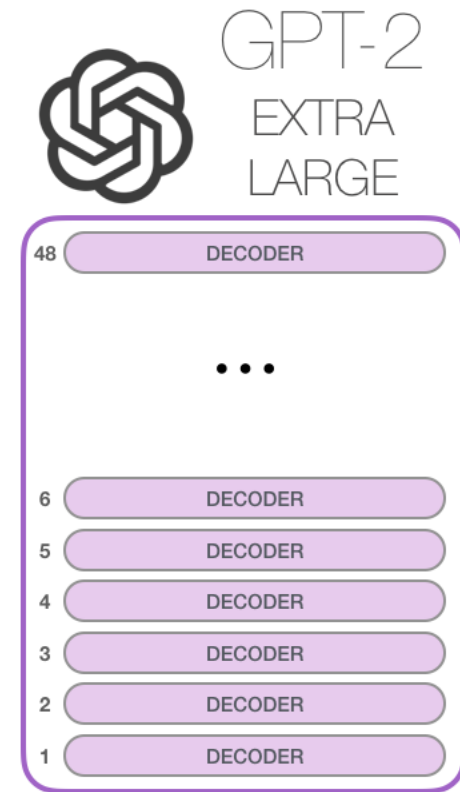
Model Dimensionality: 1024

345M



Model Dimensionality: 1280

762M



Model Dimensionality: 1600

1542M

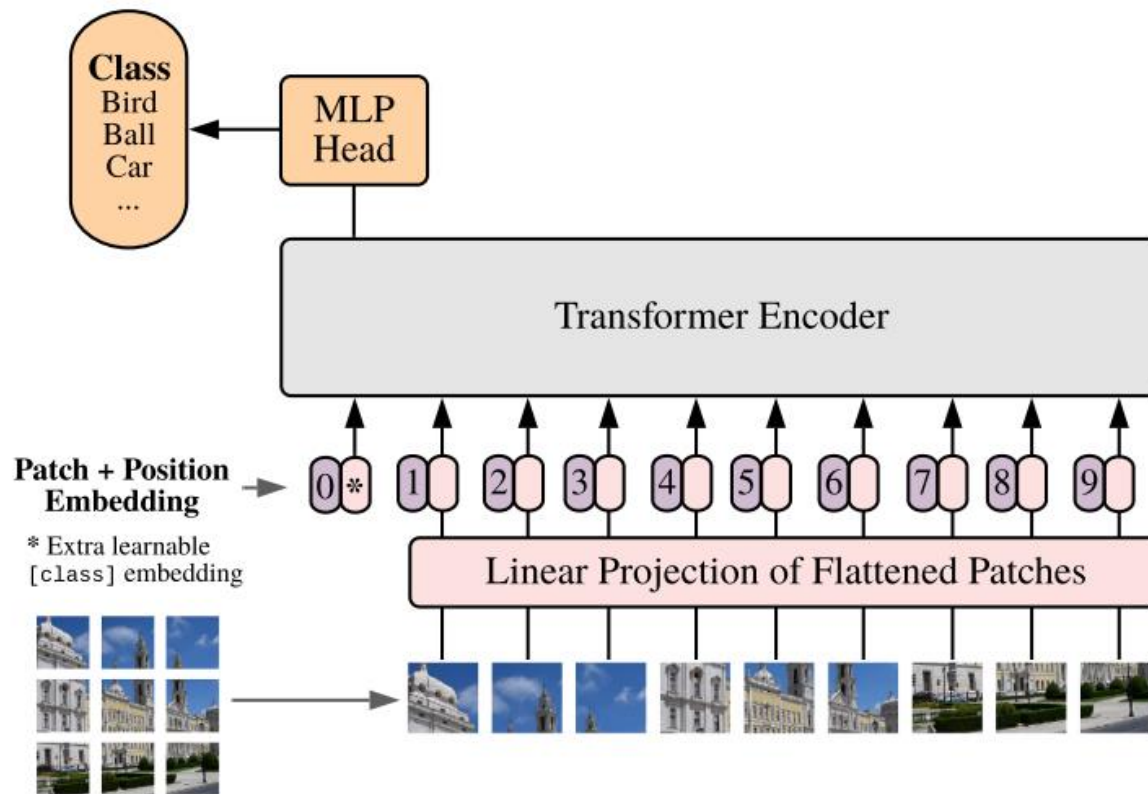


ViT – Vision Transformers

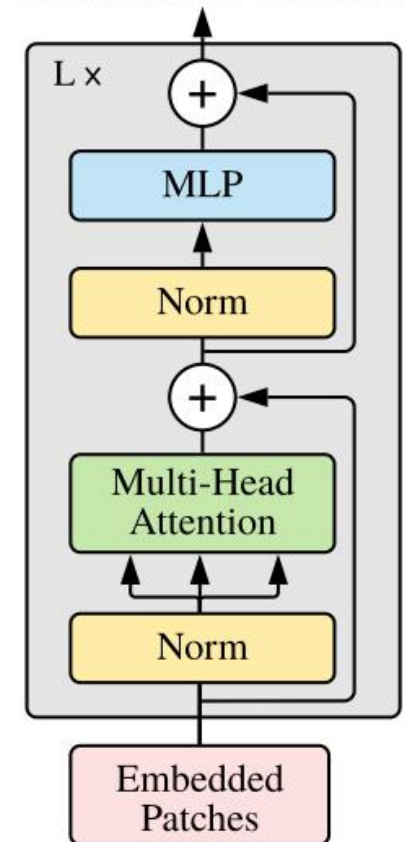


- Divide an input image into 196 (14×14) small images of size (16×16)
- Treat it as embedding in NLP
- Use it as an input for traditional transformer encoder (like in BERT)
- Use 12 transformer layers (Norm, Multi-head attention, etc.)
- Take the last output, use it as input for Dense Layer with 1000 classes

Vision Transformer (ViT)



Transformer Encoder

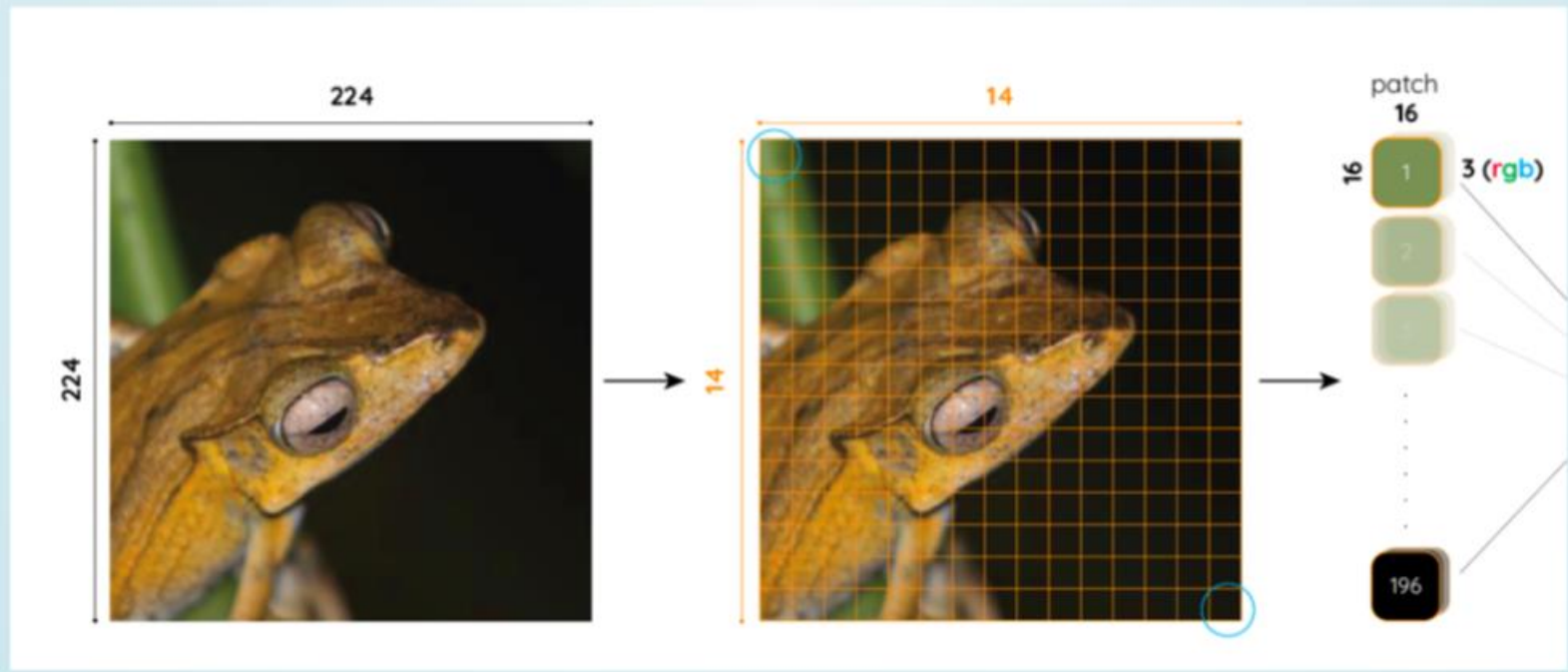


Patch Embedding

Image of size (3, 224, 224)

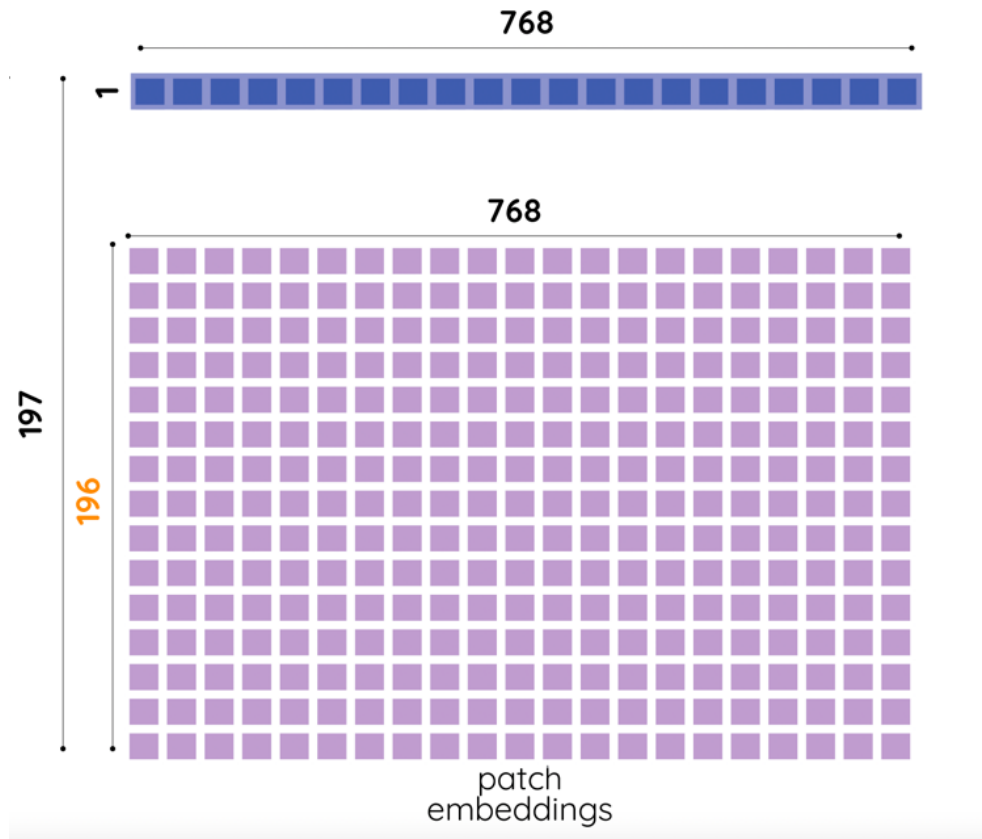
Divided into 196 (14×14) patches of size $3 \times 16 \times 16$

$16 \times 14 = 224$ (original image size)



Each patch is converted into a vector of size $3 \times 16 \times 16 = 768$

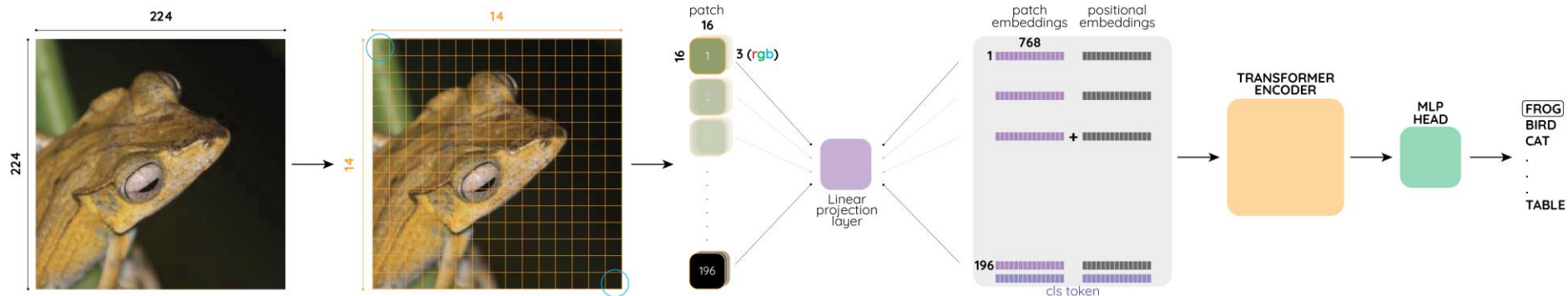
[CLS] Token



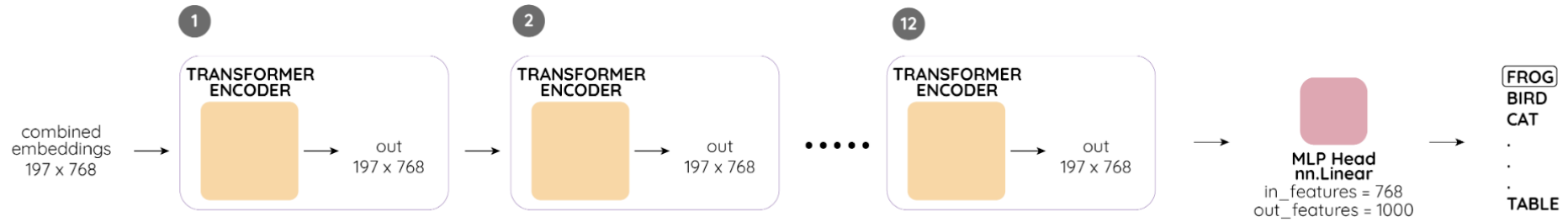
[CLS] token is a vector of size $(1,768)$

The final patch matrix has size $(197,768)$, 196 from patches and 1 [CLS] token

ViT - Summary



ViT Layers



How Good is ViT?

- Worse than Resnet when trained just on ImageNet
- Performance improved when pre-trained on big (and I mean it) dataset
- Pretrained outperforms much bigger CNNs

ViT – in Numbers

Model	Layers	Hidden size D	MLP size	Heads	Params
ViT-Base	12	768	3072	12	86M
ViT-Large	24	1024	4096	16	307M
ViT-Huge	32	1280	5120	16	632M

Table 1: Details of Vision Transformer model variants.

	Ours-JFT (ViT-H/14)	Ours-JFT (ViT-L/16)	Ours-I21K (ViT-L/16)	BiT-L (ResNet152x4)	Noisy Student (EfficientNet-L2)
ImageNet	88.55 ± 0.04	87.76 ± 0.03	85.30 ± 0.02	87.54 ± 0.02	88.4/88.5*
ImageNet Real	90.72 ± 0.05	90.54 ± 0.03	88.62 ± 0.05	90.54	90.55
CIFAR-10	99.50 ± 0.06	99.42 ± 0.03	99.15 ± 0.03	99.37 ± 0.06	—
CIFAR-100	94.55 ± 0.04	93.90 ± 0.05	93.25 ± 0.05	93.51 ± 0.08	—
Oxford-IIIT Pets	97.56 ± 0.03	97.32 ± 0.11	94.67 ± 0.15	96.62 ± 0.23	—
Oxford Flowers-102	99.68 ± 0.02	99.74 ± 0.00	99.61 ± 0.02	99.63 ± 0.03	—
VTAB (19 tasks)	77.63 ± 0.23	76.28 ± 0.46	72.72 ± 0.21	76.29 ± 1.70	—
TPUv3-core-days	2.5k	0.68k	0.23k	9.9k	12.3k

Critics

- Better results - only with more data
- The cost of training from scratch is ridiculously high (30k\$)
- Is it really that different from Convolutions?