

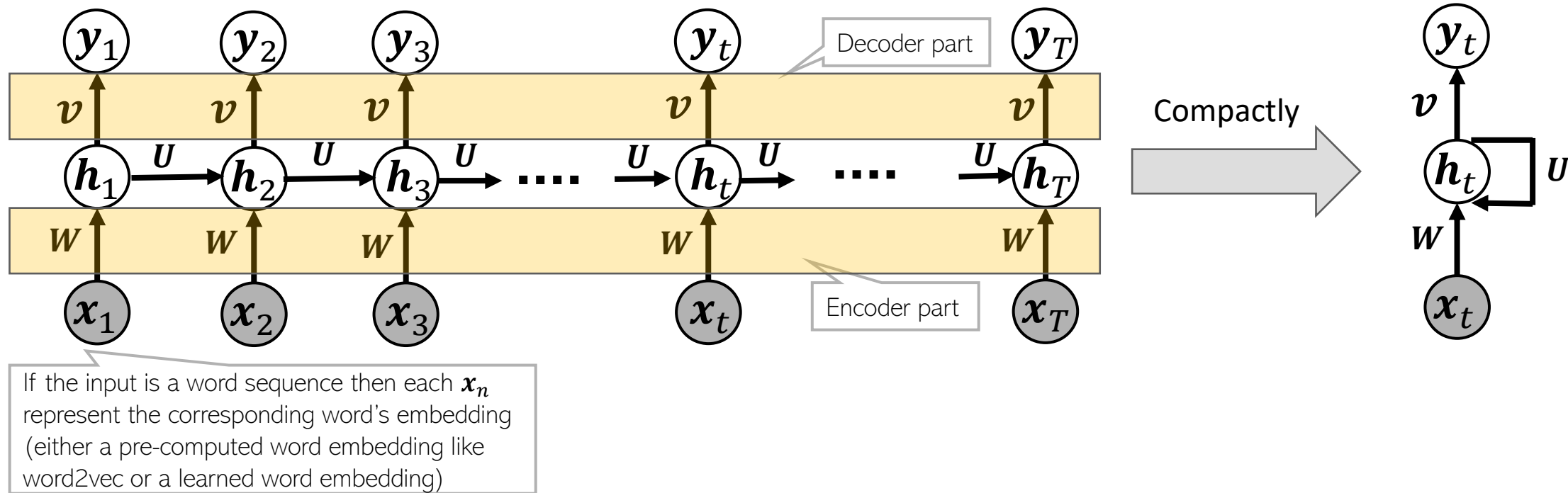
Deep Neural Networks: Attention Mechanism and Transformers

CS771: Introduction to Machine Learning

Piyush Rai

Recap: RNNs

- RNNs are used when each input or output or both are **sequences of tokens**

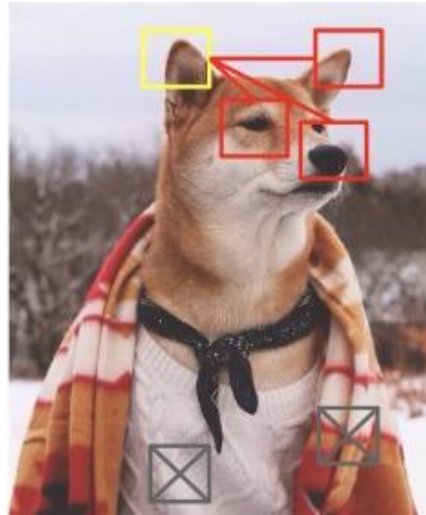


- Hidden state h_t is supposed to remember everything up to time $t - 1$. However, in practice, RNNs have difficulties remembering the distant past
 - Variants such as LSTM, GRU, etc mitigate this issue to some extent
- Slow processing is another major issue (e.g., can't compute h_t before computing h_{t-1})



Attention

- We use attention to “focus” on some part of interest in an input
 - Other nearby relevant parts help us focus
 - Other irrelevant parts do not contribute in the process

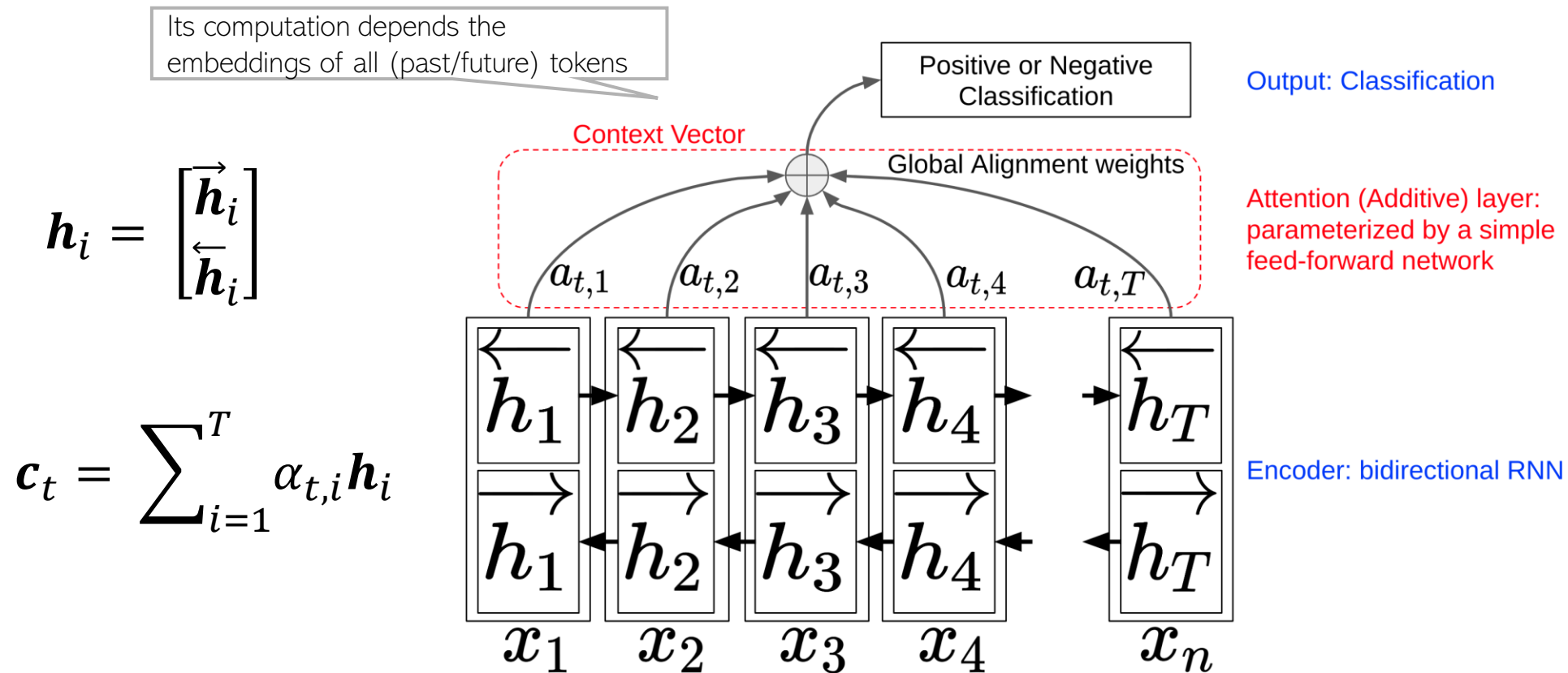


- In sequence modeling problems, we can use attention between input and output tokens (between encoder and decoder parts), as well as among the inputs only (only within the encoder part)



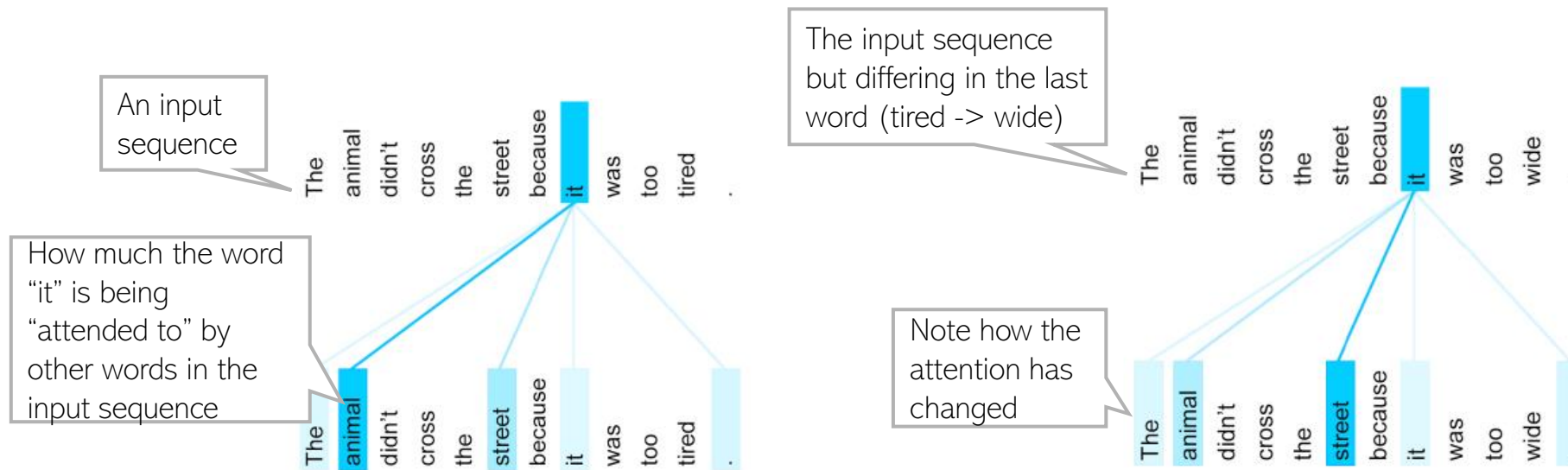
RNN with Attention

- RNNs have also been augmented with attention to help remember the distant past
- Attention mechanism for a bi-directional RNN encoder-decoder model



Self-Attention

- With self-attention, each token \mathbf{x}_n can “attend to” all other tokens of the same sequence when computing this token’s embedding \mathbf{h}_n



- Attention helps capture the context better and in a much more “global” manner
 - “Global”: Long ranges captures and in both directions (previous and ahead)



Self-Attention

6

- For an N length sequence, the attention scores for each token \mathbf{x}_n are computed using

- A **query vector** \mathbf{q}_n associated with that token
- N **key vectors** $\mathbf{K} = \{\mathbf{k}_1, \mathbf{k}_2, \dots, \mathbf{k}_N\}$ (one per token)
- N **value vectors** $\mathbf{V} = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_N\}$ associated with the key vectors

Provided in form of a D -dim embedding (e.g., word2vec)

$N \times D$ matrix of original embeddings from the input layer

Linear projection by $D \times d$ matrix

Learnable

Row n is \mathbf{q}_n

$N \times d$ matrix of "queries"

$$\mathbf{Q} = \mathbf{XW}_Q$$

$N \times d$ matrix of N "keys"

Assuming same size d as query

$$\mathbf{K} = \mathbf{XW}_K$$

$D \times d$ matrix. Learnable

$N \times K$ matrix of "value" vectors of the N keys

$$\mathbf{V} = \mathbf{XW}_V$$

$D \times K$ matrix. Learnable

- One way to compute the attention score is

How much token i attends to token n

$$\alpha_{n,i} = \frac{\exp(\mathbf{q}_n^\top \mathbf{k}_i)}{\sum_{j=1}^N \exp(\mathbf{q}_n^\top \mathbf{k}_j)}$$

Dot-product attention (query and key assumed d dimensional)

- Given attention scores, encoder's hidden state for \mathbf{x}_n is

\mathbf{Q} and \mathbf{K} are assumed $N \times d$

$N \times v$

$$\mathbf{h}_n = \sum_{i=1}^N \alpha_{n,i} \mathbf{v}_i$$

Attention-weighted sum of the value vectors of all the tokens in the sequence

Thus the encoding of \mathbf{x}_n depends on all the tokens in the sequence

$N \times v$

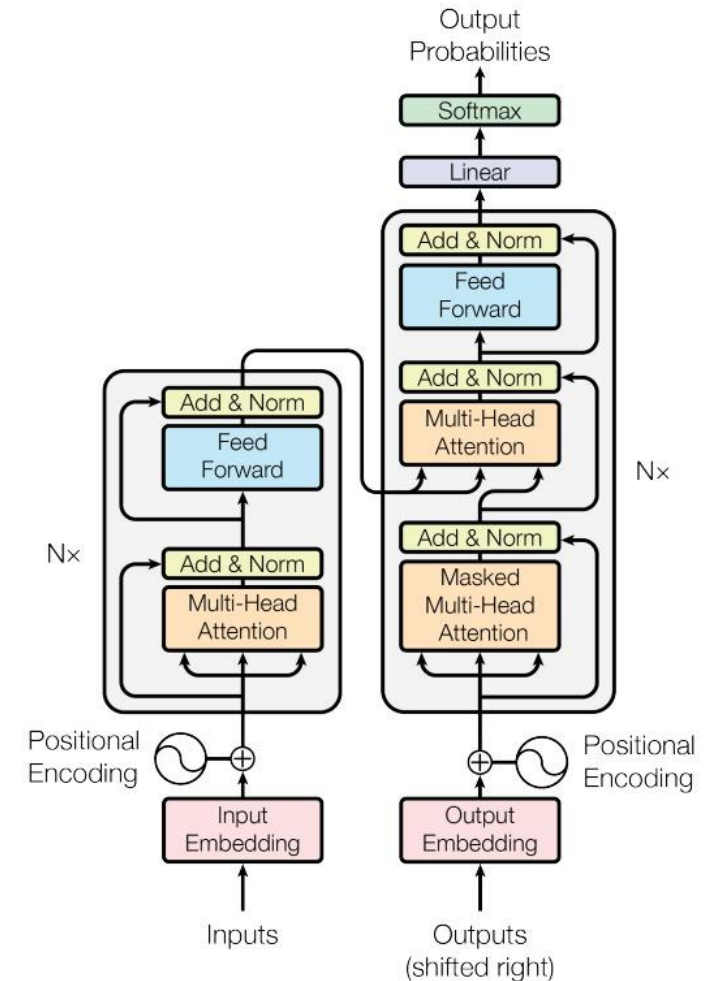
$$\mathbf{H} = \text{softmax} \left(\frac{\mathbf{QK}^\top}{\sqrt{d}} \right) \mathbf{V}$$

Dividing by \sqrt{d} ensures variance of the dot product is 1

"Scaled" dot-product attention

Transformers

- Transformers also use the idea of attention
 - “self-attention” over all input tokens
 - “self-attention” over each output token and previous tokens
 - “cross-attention” between output tokens and input tokens
- Transformer also compute embeddings of all tokens in parallel
- Transformers are based on the following key ideas*
 - “Self-attention” and “cross-attention” for computing the hidden states
 - Positional encoding
 - Residual connections
- Attention helps capture the context better and in a much more “global” manner in sequence data



Positional Encoding

- Transformers also need a “positional encoding” for each token of the input since they don’t process the tokens sequentially (unlike RNNs)
- Let $\mathbf{p}_i \in \mathbb{R}^d$ be the positional encoding for location i . One way to define it is

Here C denotes the maximum possible length of a sequence

$$p_{i,2j} = \sin\left(\frac{i}{C^{2j/d}}\right), \quad p_{i,2j+1} = \cos\left(\frac{i}{C^{2j/d}}\right)$$

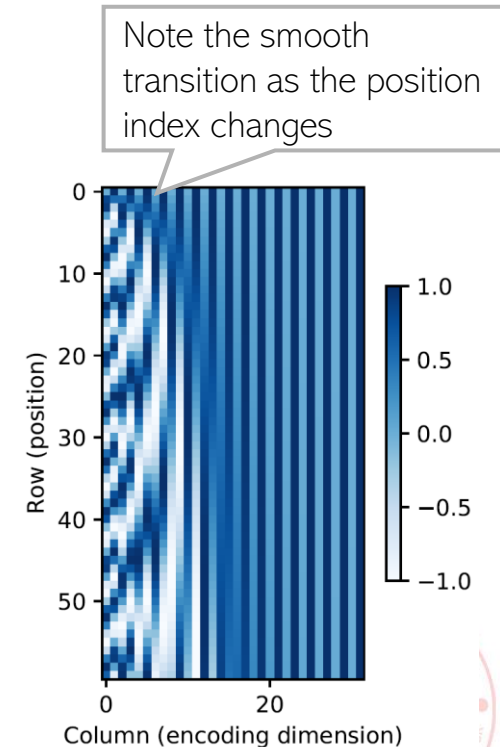
Positional encoding vector for location i assuming $d = 4$

$$\mathbf{p}_i = \left[\sin\left(\frac{i}{C^{0/4}}\right), \cos\left(\frac{i}{C^{0/4}}\right), \sin\left(\frac{i}{C^{2/4}}\right), \cos\left(\frac{i}{C^{2/4}}\right) \right]$$

- Given the positional encoding, we add them to the token embedding

$$\hat{\mathbf{x}}_i = \mathbf{x}_i + \mathbf{p}_i$$

- The above positional encoding is pre-defined but can also be learned



Zooming into the encoder and the decoder..

FF operation is applied "position-wise" (for each token separately) but all FF blocks in the layer have same weights

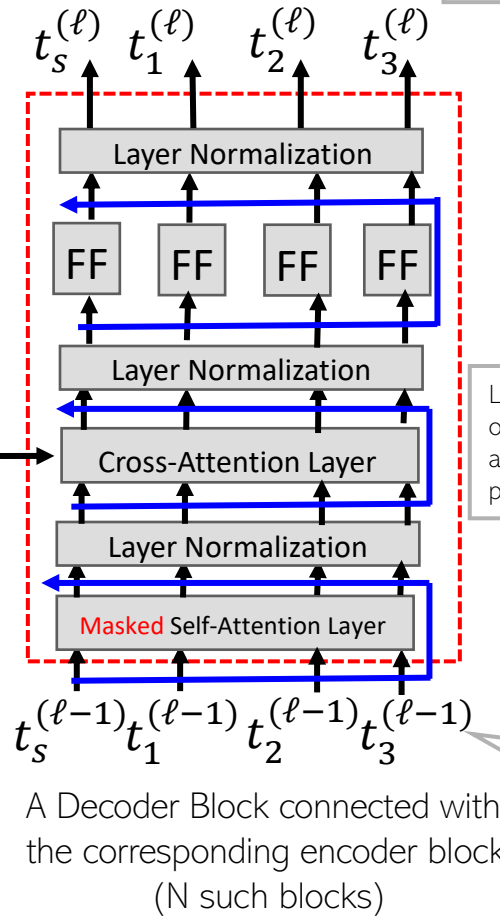
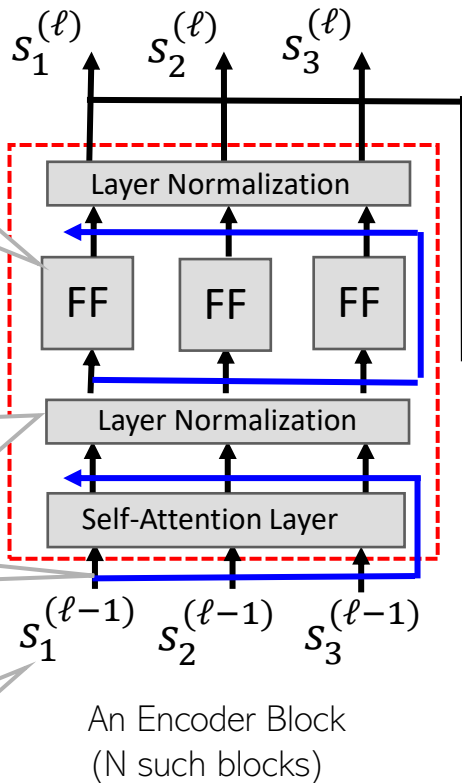
Each FF (feed-forward) is usually a linear layer + ReLU nonlinearity + another linear layer

Layer normalization (batch normalization is difficult since difference input sequences can be of different lengths)

Blue arrows are the residual connections

Input ("source") token embeddings

Fixed in the first layer (obtained from an [input embedding table](#)) and [learned for subsequent layers](#)



Most likely output token at step m

"auto-regressive" generation

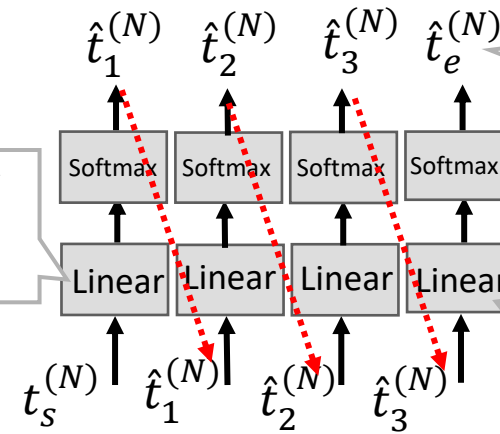
Output $\hat{t}_{m-1}^{(N)}$ from the previous step ($m-1$) in the sequence

$$\hat{t}_m^{(N)} = \operatorname{argmax}_{i=1,\dots,V} \operatorname{softmax}(\mathbf{W}t_m^{(N)})$$

t_e is a special end of sequence (EOS) token

With weight matrix \mathbf{W} of size $V \times D$ where V is vocab size and D is the dimensionality of the last decoder block embeddings $t_m^{(N)}$

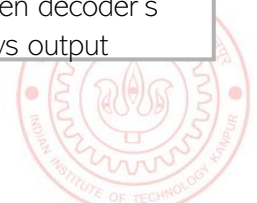
Like FF, linear operation is also applied position-wise



Output ("target") token embeddings

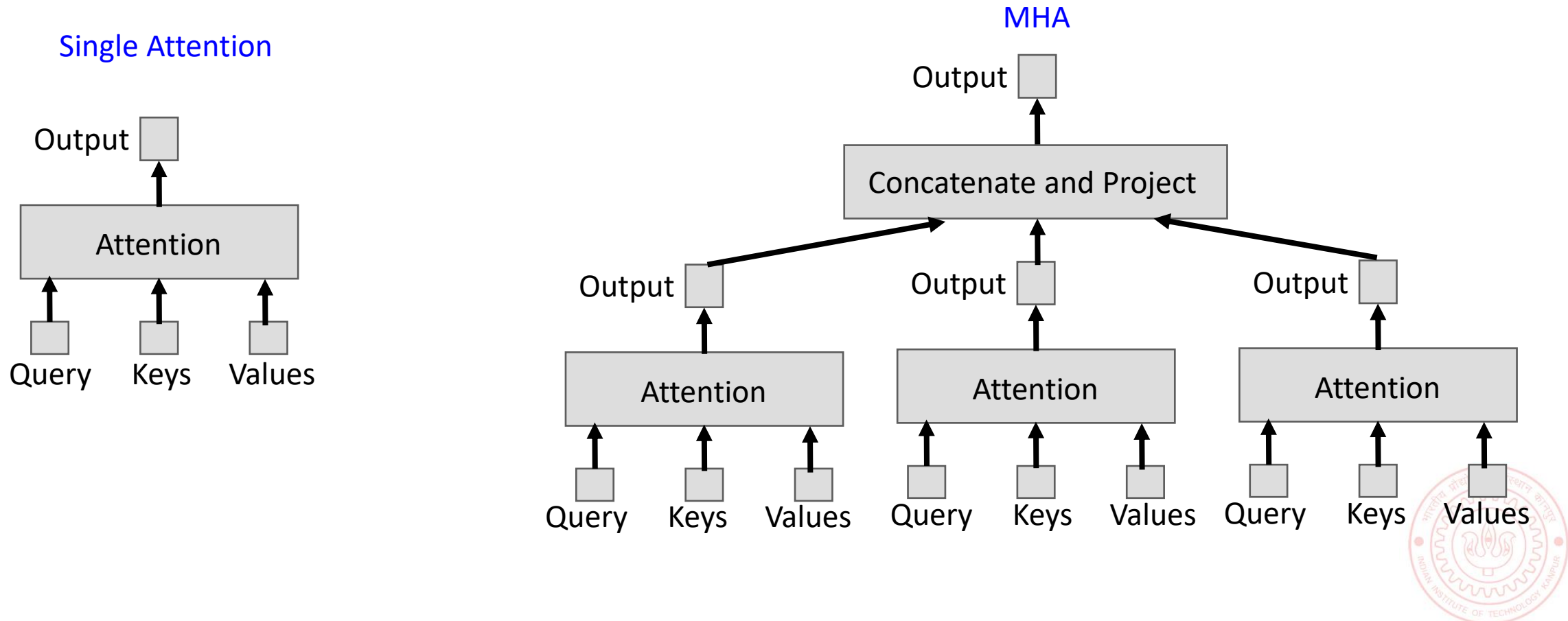
t_s is a special start of sequence (SOS) token

Note the one position (towards right) shift between decoder's input vs output

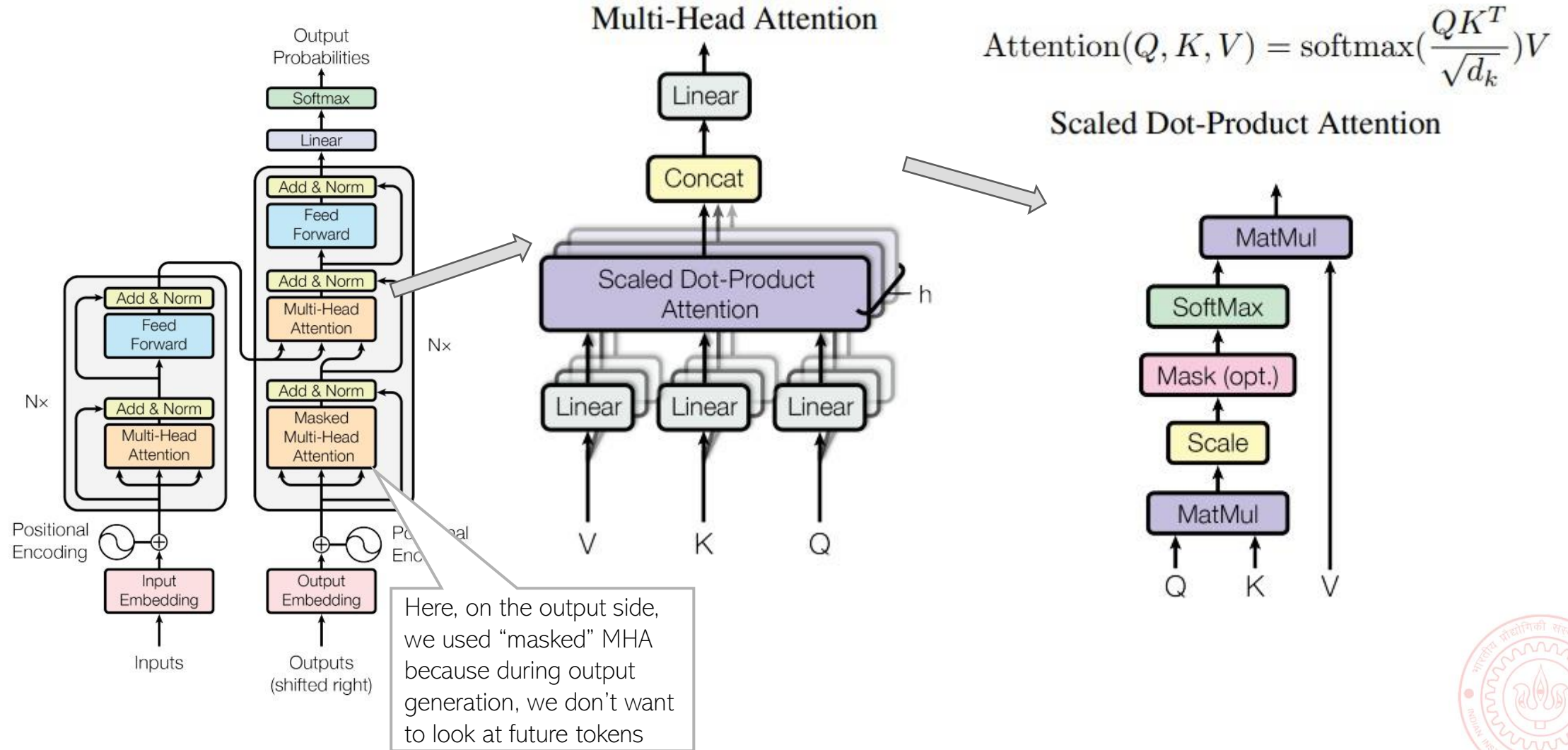


Multi-head Attention (MHA)

- A single attention function can capture only one notion of similarity
- Transformers therefore use multi-head attention (MHA)



(Masked) Multi-head Attention (MHA)



Computing Attention Efficiently

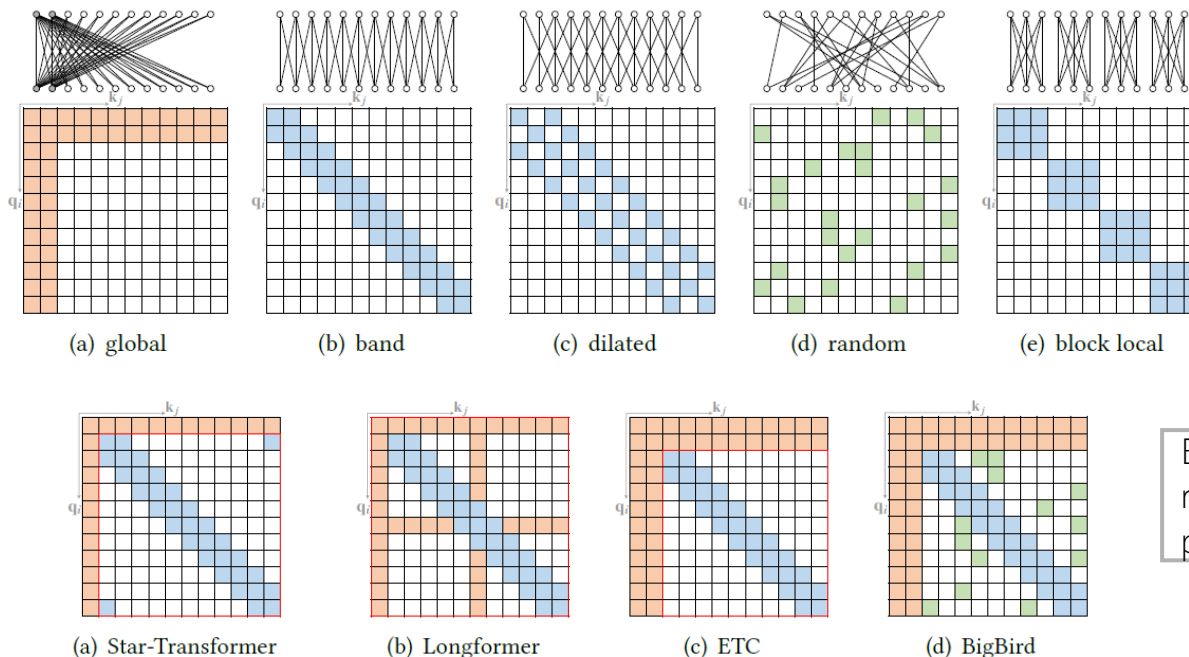
- The standard attention mechanism is inefficient for large sequences

$O(T^2)$ storage and computation cost for a T length sequence

$$\mathbf{H} = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}}\right) \mathbf{V}$$

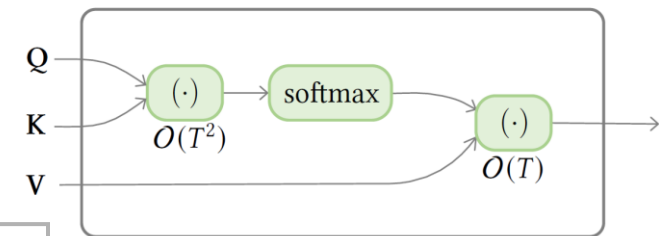
- Many ways to make it more efficient, e.g.,

Sparse Attention



Linearized Attention

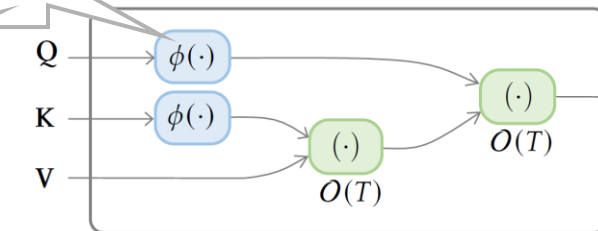
$$\exp(\mathbf{Q}\mathbf{K}^\top) \approx \phi(\mathbf{Q})^\top \phi(\mathbf{K})$$



(a) standard self-attention

A nonlinear projection based on kernels

E.g., kernel random features projection

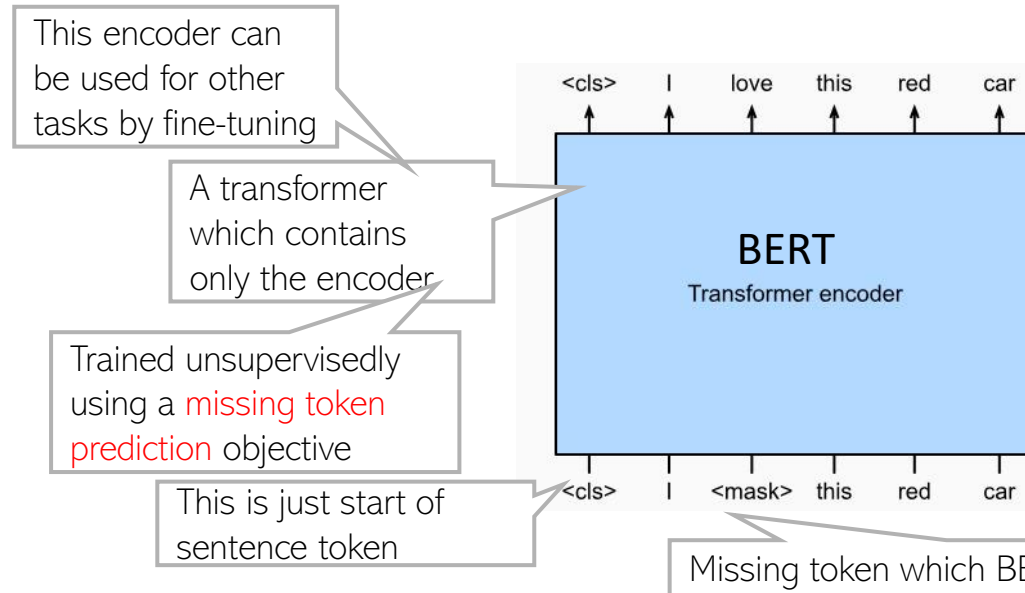
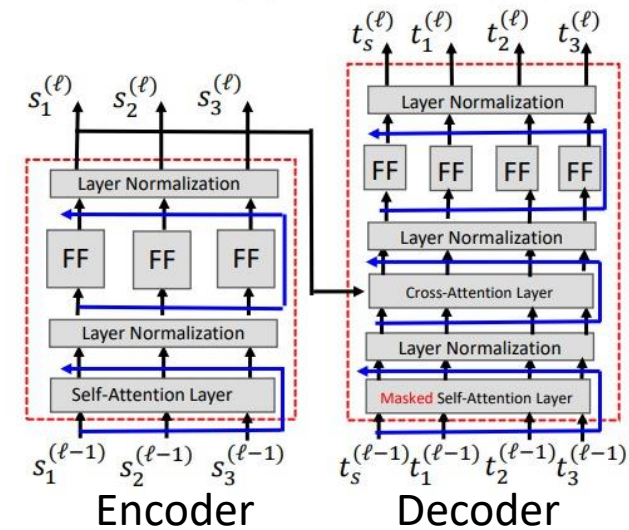


(b) linearized self-attention



Popular Transformer Variants: BERT and GPT

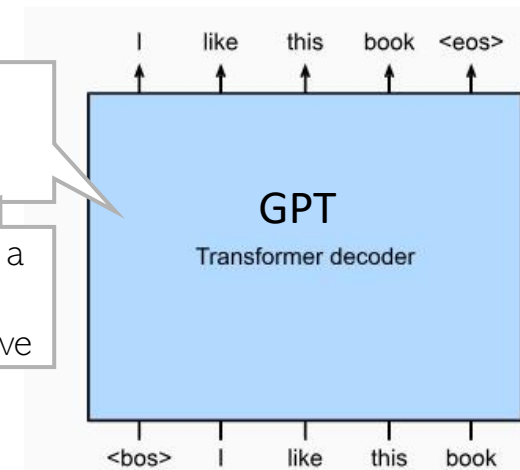
- The standard transformer architecture is an encoder-decoder model
- Some models use just the encoder or the decoder of the transformer
- **BERT** (Bidirectional Encoder Representations from Transformers)
 - Basic BERT can be learned to encoder token sequences
- **GPT** (Generative Pretrained Transformer)
 - Basic GPT can be used to generate token sequences similar to its training data



Also, no cross-attention since there is no encoder

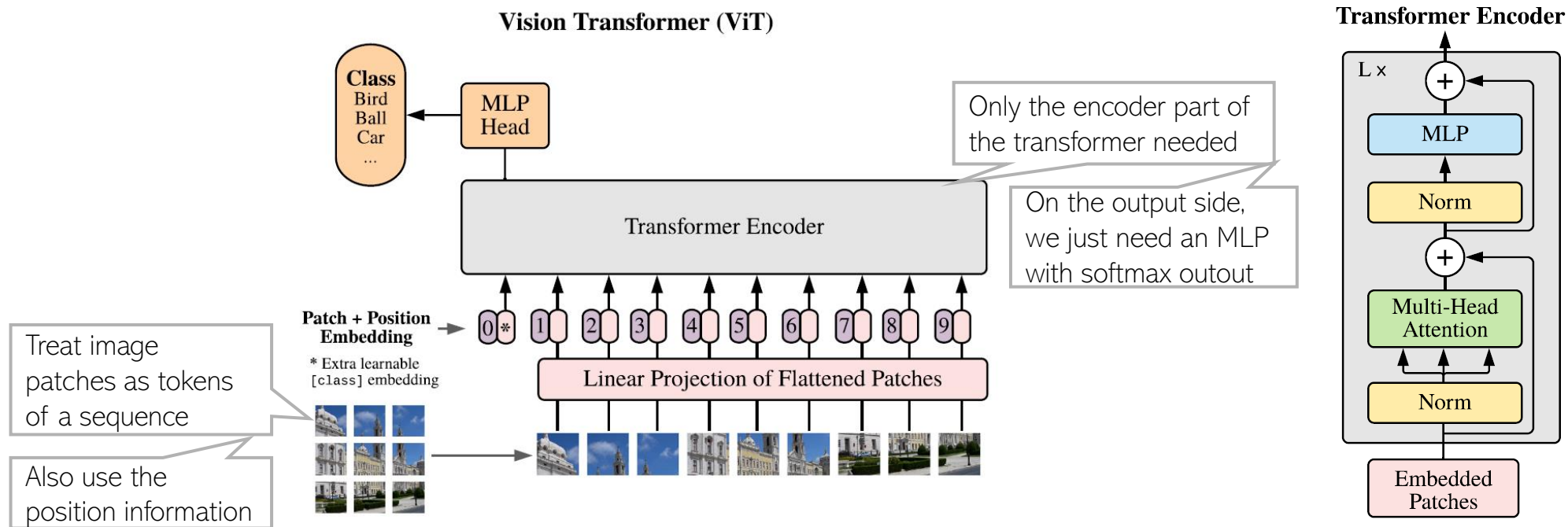
A transformer which contains only the decoder

Pre-trained using a **next token prediction** objective



Transformers for Images: ViT

- Transformers can be used for images as well[#]. For image classification, it looks like this



- Early work showed ViT can outperform CNNs given very large amount of training data
- However, recent work^{*} has shown that good old CNNs still rule! ViT and CNN perform comparably at scale, i.e., when both given large amount of compute and training data

[#] An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale (Dosovitskiy et al, 2020)

^{*}ConvNets Match Vision Transformers at Scale (Smith et al, 2023)