



大模型系列—大模型算法

Transformer介绍



ZOMI

大模型业务全流程



关于本内容

- **具体内容**
 1. 大模型都是什么模型结构？
 2. 为什么出现 Transformer？
 3. Transformer 结构介绍？
 4. Bert 和 GPT 大模型出现

1 . LLM

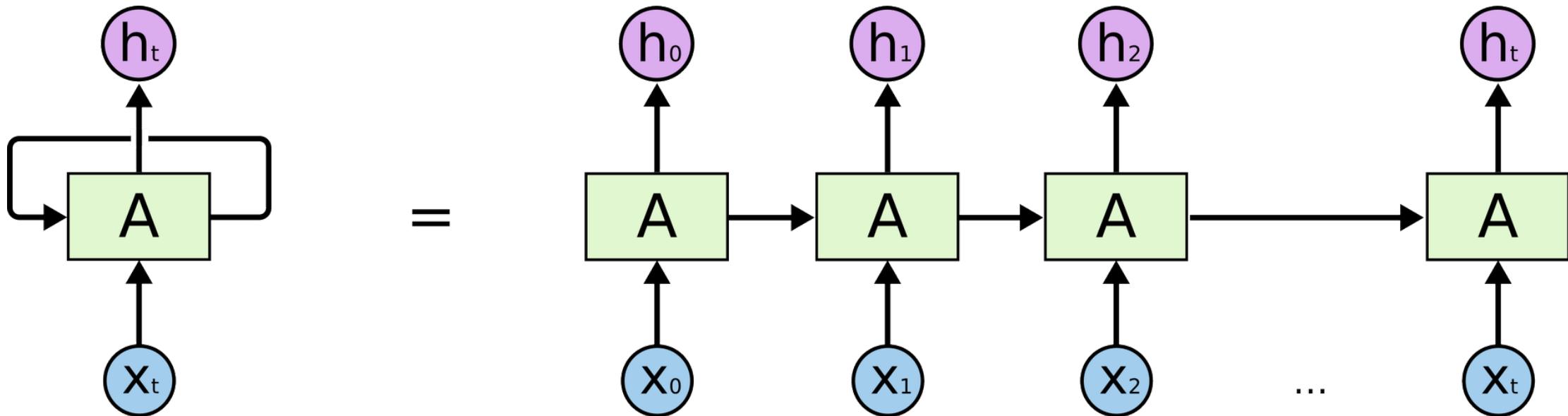
模型结构演进

- GPT3
- LLAMA1
- LLAMA2
- Baichuan1
- Baichuan2
- Qwen
- Mistral MOE

2. Motivation for Transformers

Problems with RNNs = Motivation for Transformers

- RNN循环神经网络是一个按时间序列的前馈神经网络模型，经过训练后可以处理顺序数据输入并将其转换为特定的顺序数据输出。



RNN has three major disadvantages

Slow to Train

训练慢

Short Memory

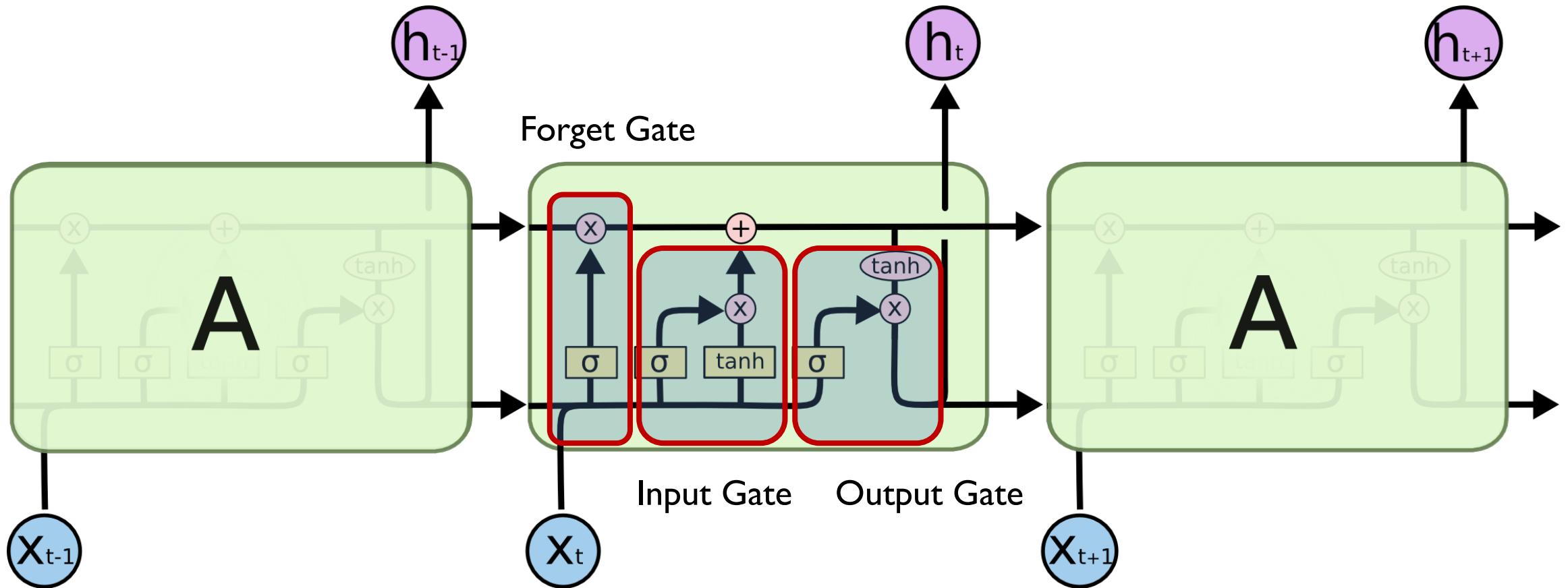
短期记忆

Vanishing Gradient

梯度消失/梯度爆炸

LSTM

- long short-term memory , LSTM 是一种特殊 RNN , 主要是为了解决长序列训练过程中的梯度消失和梯度爆炸问题。



LSTM has three major disadvantages

Slower to Train

训练更慢

Short Memory

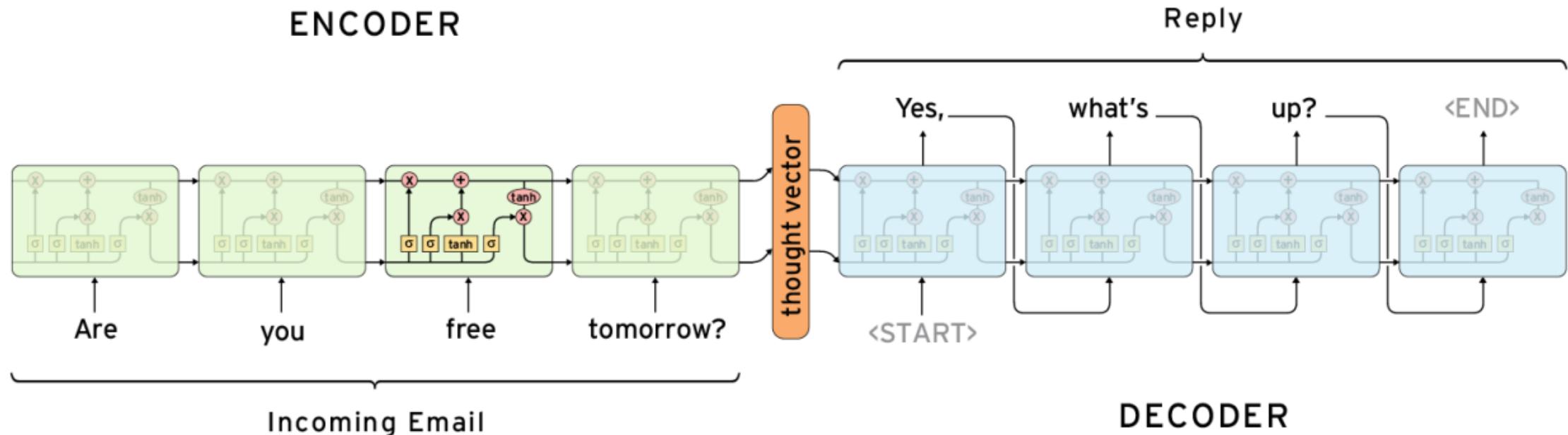
短期记忆

Vanishing Gradient

梯度消失/梯度爆炸

Encoder-Decoder Structure

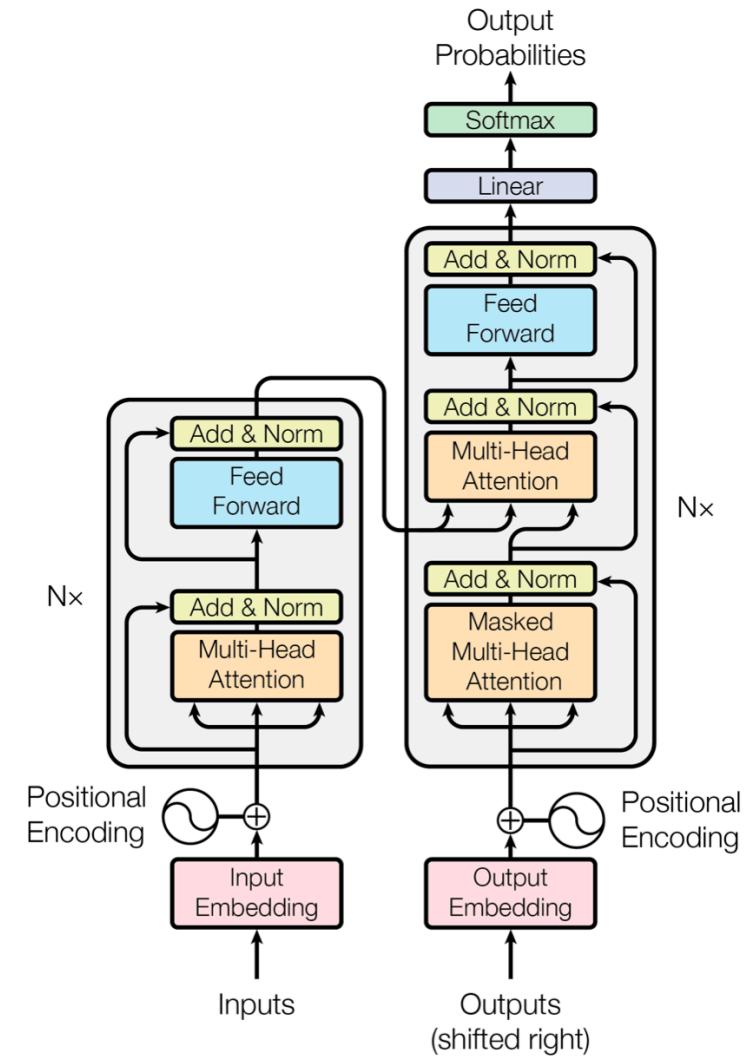
- Seq2Seq (Sequence-to-Sequence , 序列到序列) 是一类特殊序列建模问题 ;
- Encoder-Decoder 模型包括编码器和解码器 , 编码器 (Encoder) 先对输入的序列进行处理 , 将处理后向量发送给解码器 (Decoder) , 转化成想要的输出。



3. Transformer and Attention

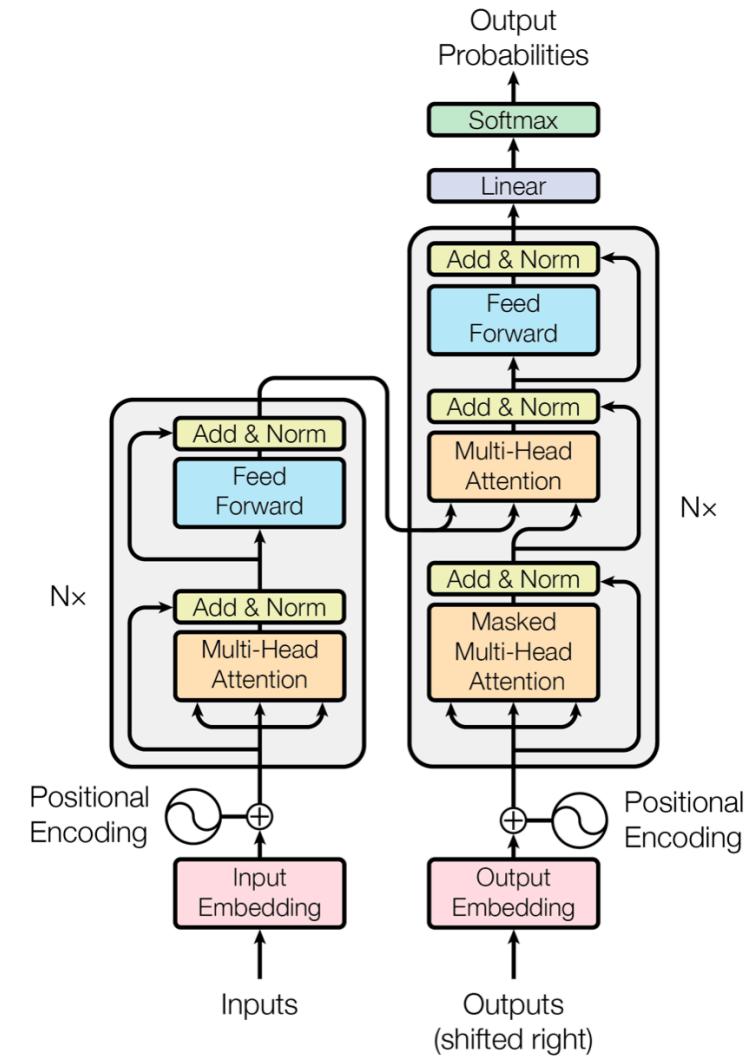
“Original” Transformer (Vaswani et al., 2017).

- The architecture is tested on English-to-German and English-to-French translation using the standard WMT2014 datasets.
 - English-to-German: 4.5M sentence pairs, 37k tokens vocabulary.
 - English-to-French: 36M sentence pairs, 32k tokens vocabulary.
- 8 P100 GPUs (150 TFlops FP16), 0.5 day for the small model, 3.5 days for the large one.



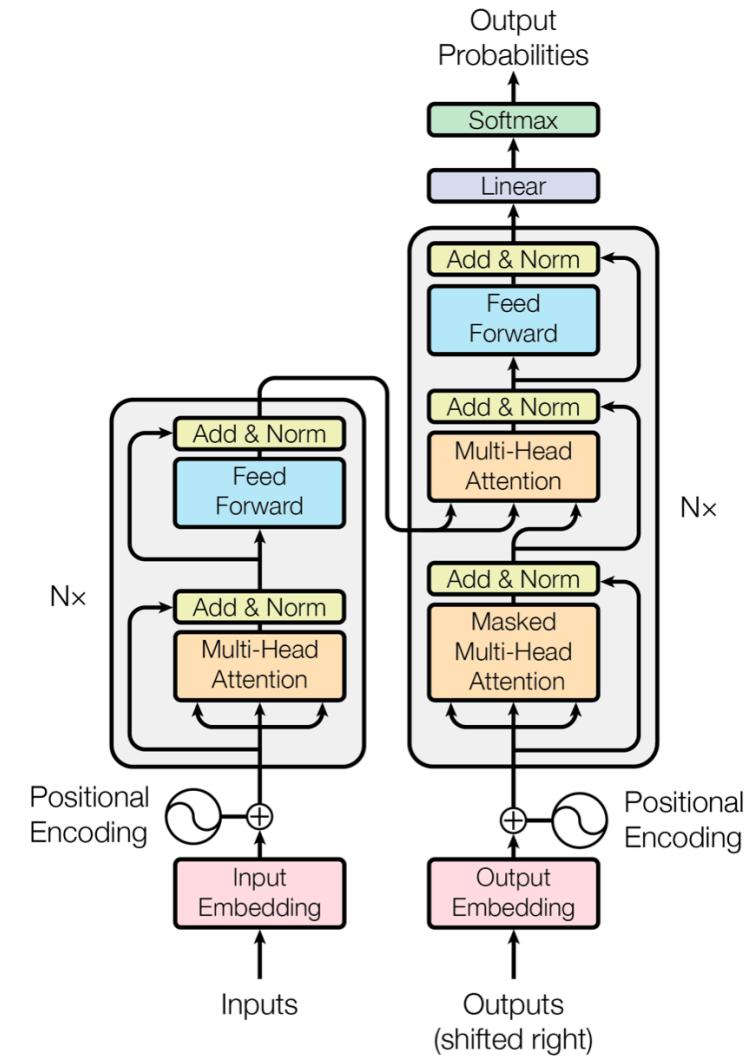
“Original” Transformer (Vaswani et al., 2017).

- Tokenization
- Input Embeddings
- Position Encodings
- Residuals
- Query
- Key
- Value
- Add & Norm
- Encoder
- Decoder
- Attention
- Self Attention
- Multi Head Attention
- Masked Attention
- Encoder Decoder Attention
- Output Probabilities / Logits
- Softmax
- Decoder only models

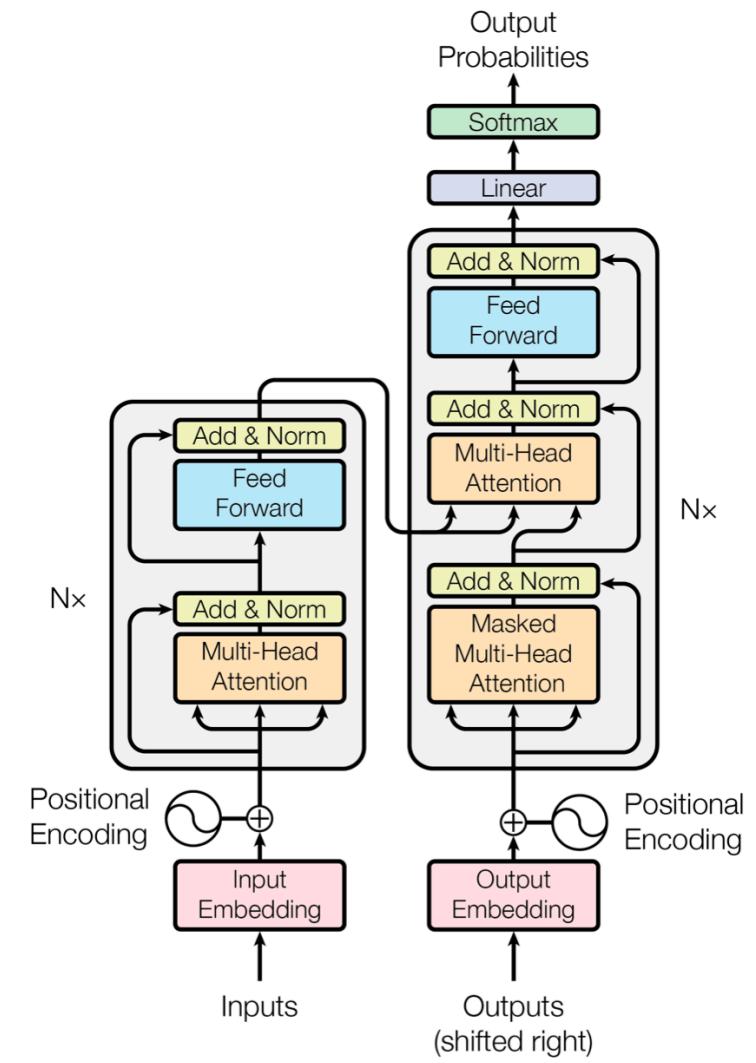
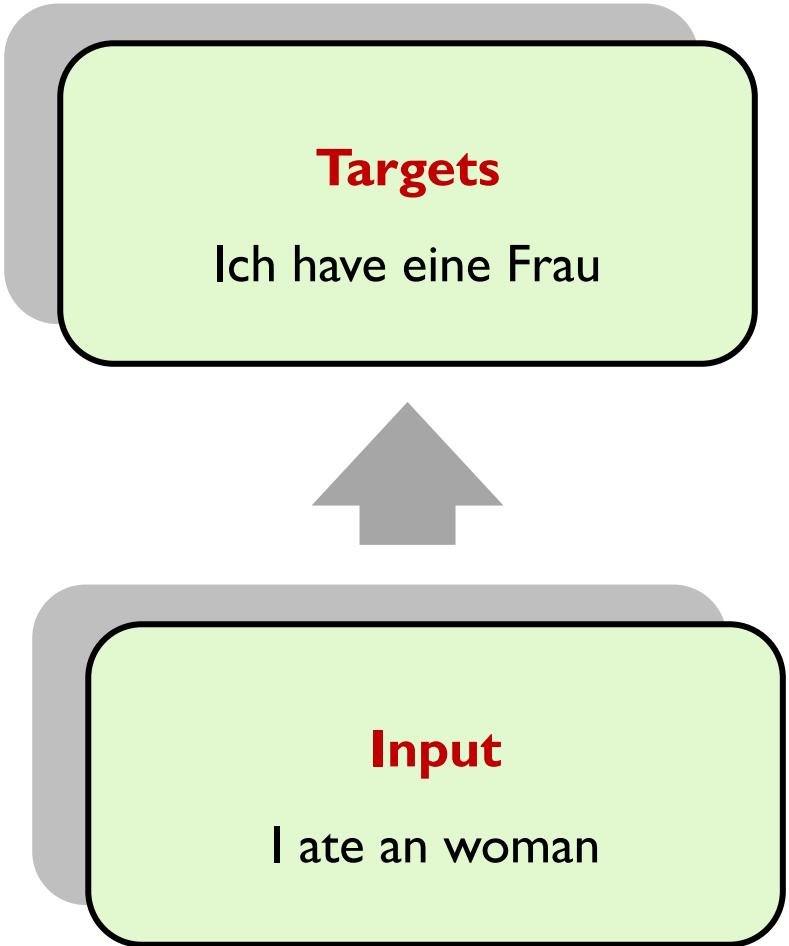


“Original” Transformer (Vaswani et al., 2017).

- Tokenization
- Decoder
- Input Embeddings
- Position Encodings
- Residuals
- Query
- Key
- Value
- Add & Norm
- Encoder
- Attention
- Self Attention
- Multi Head Attention
- Masked Attention
- Encoder Decoder Attention
- Output Probabilities / Logits
- Softmax
- Decoder only models

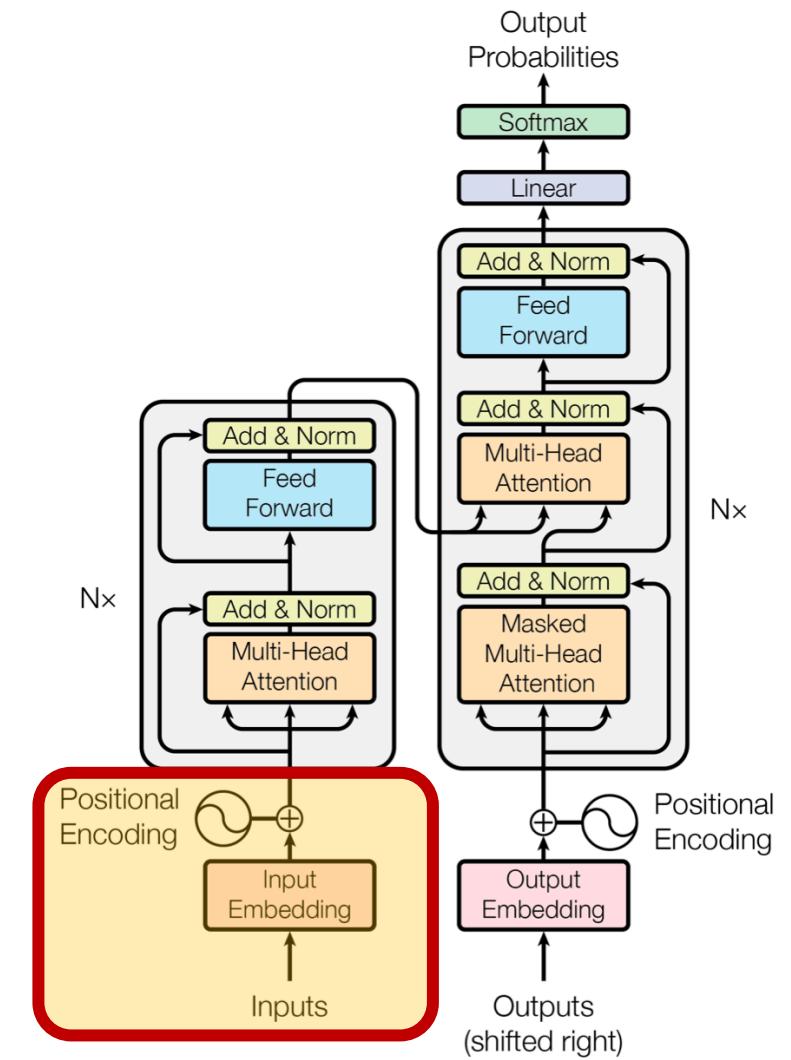
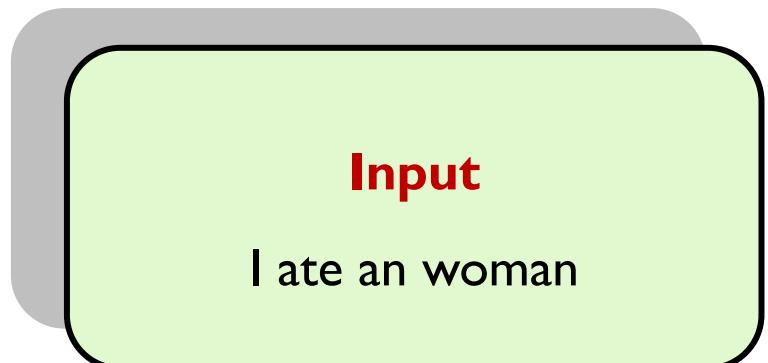


任务：机器翻译 Machine Translation

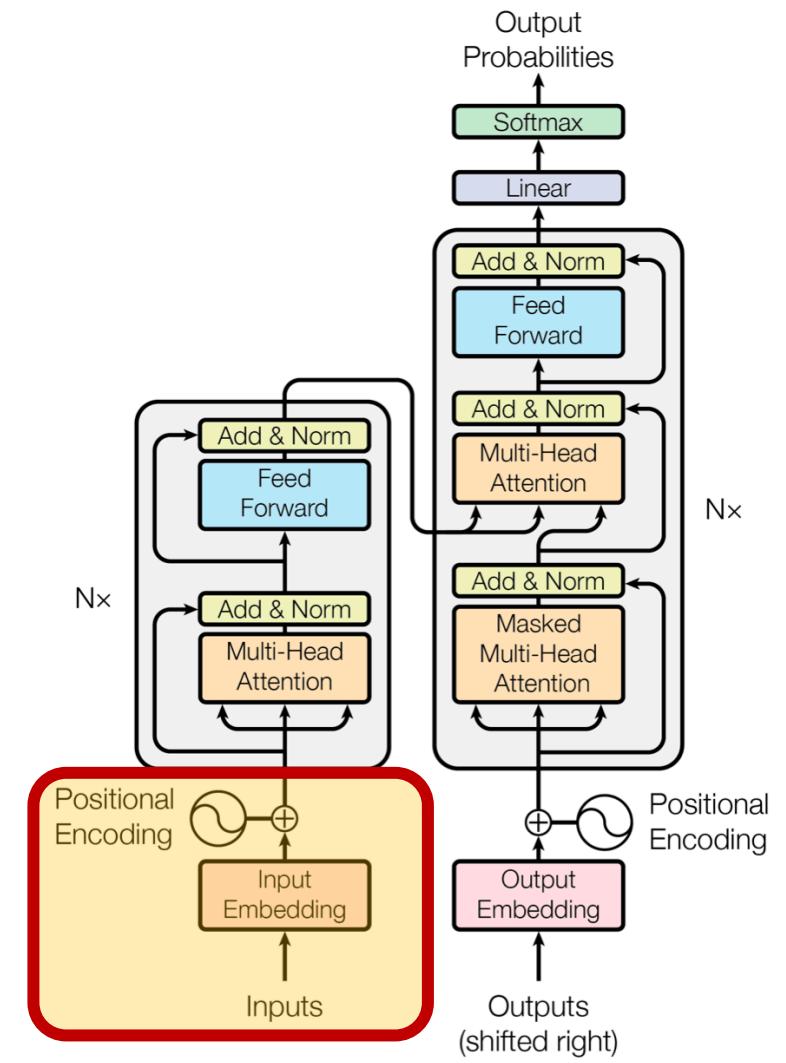
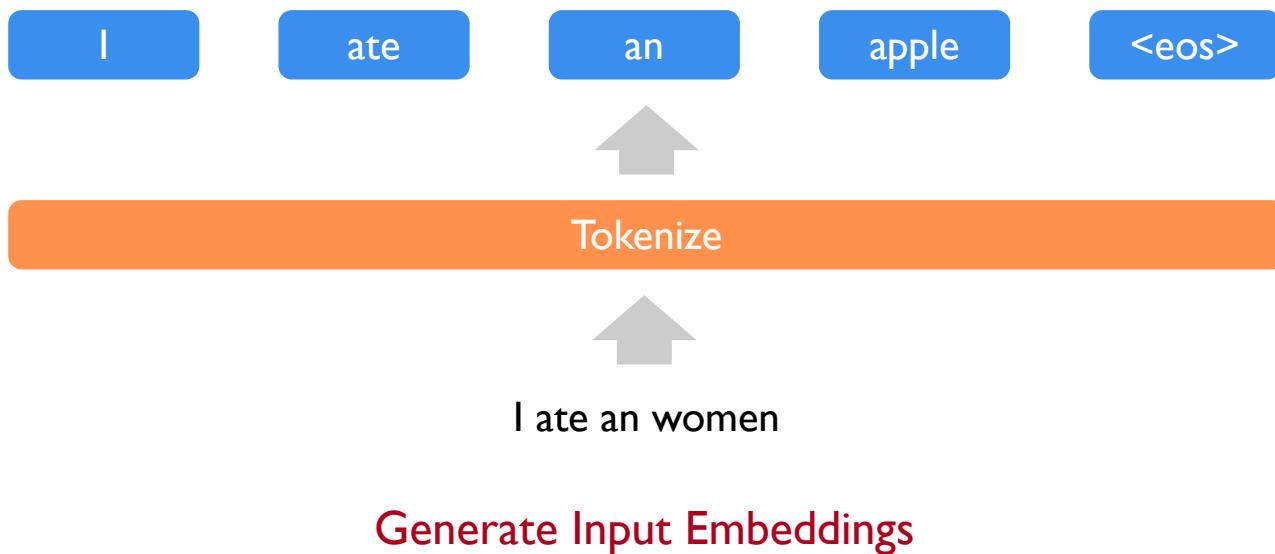


Transformer : Inputs

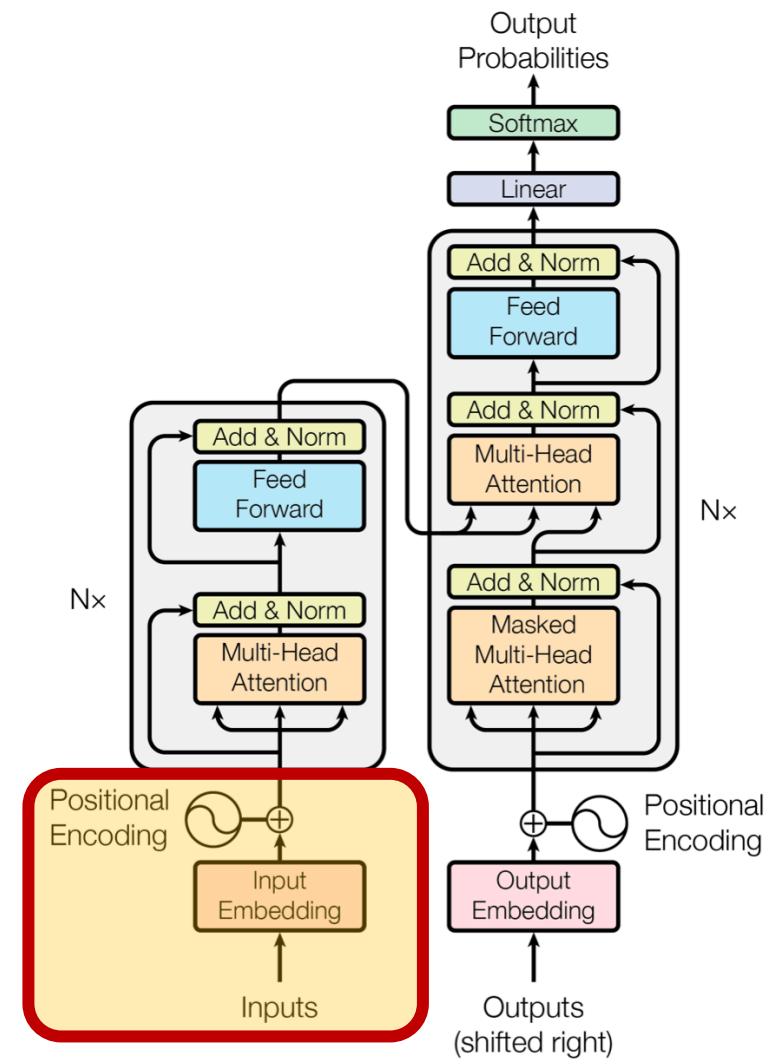
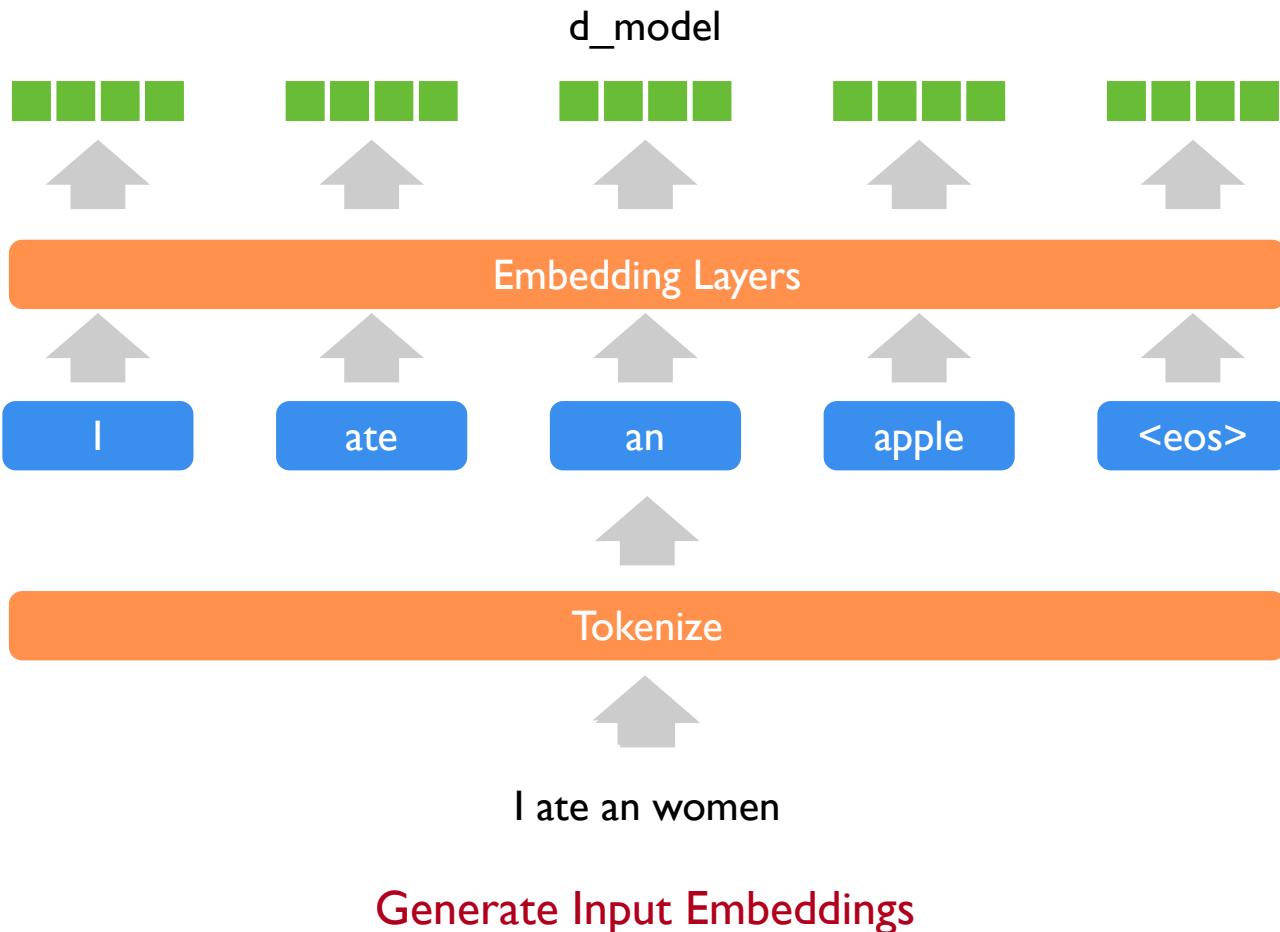
Processing Inputs



Transformer : Inputs

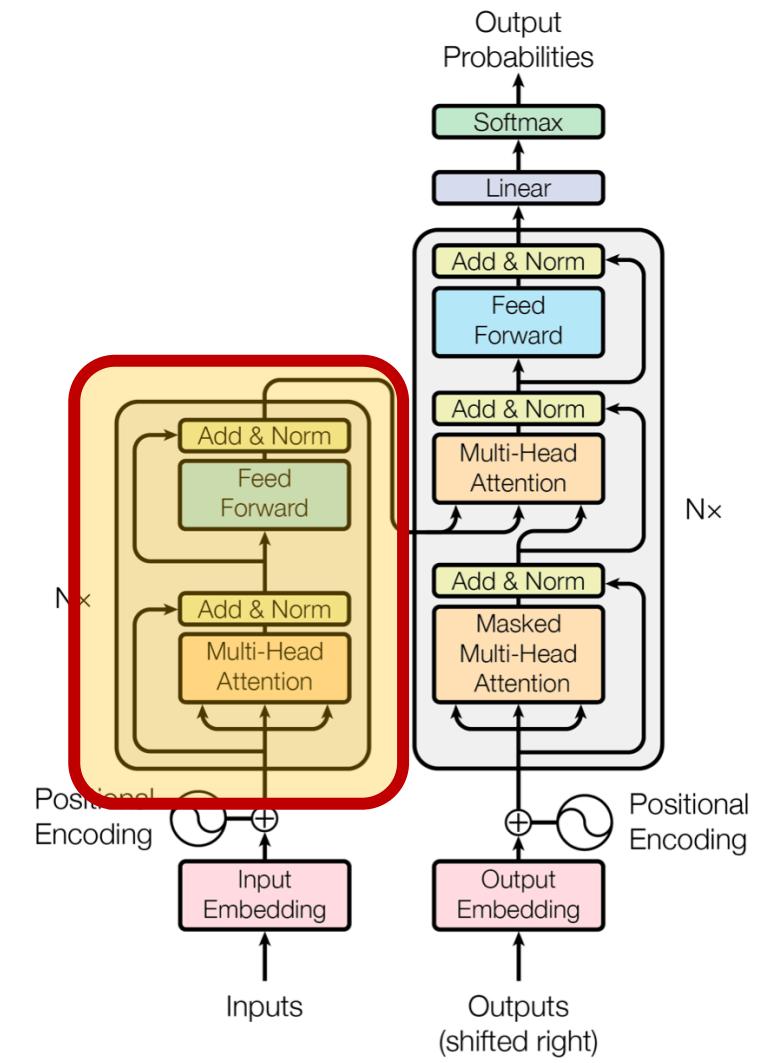
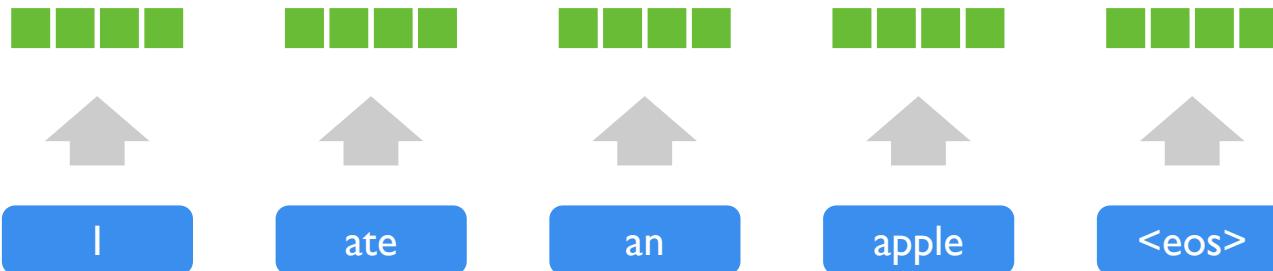


Transformer : Inputs

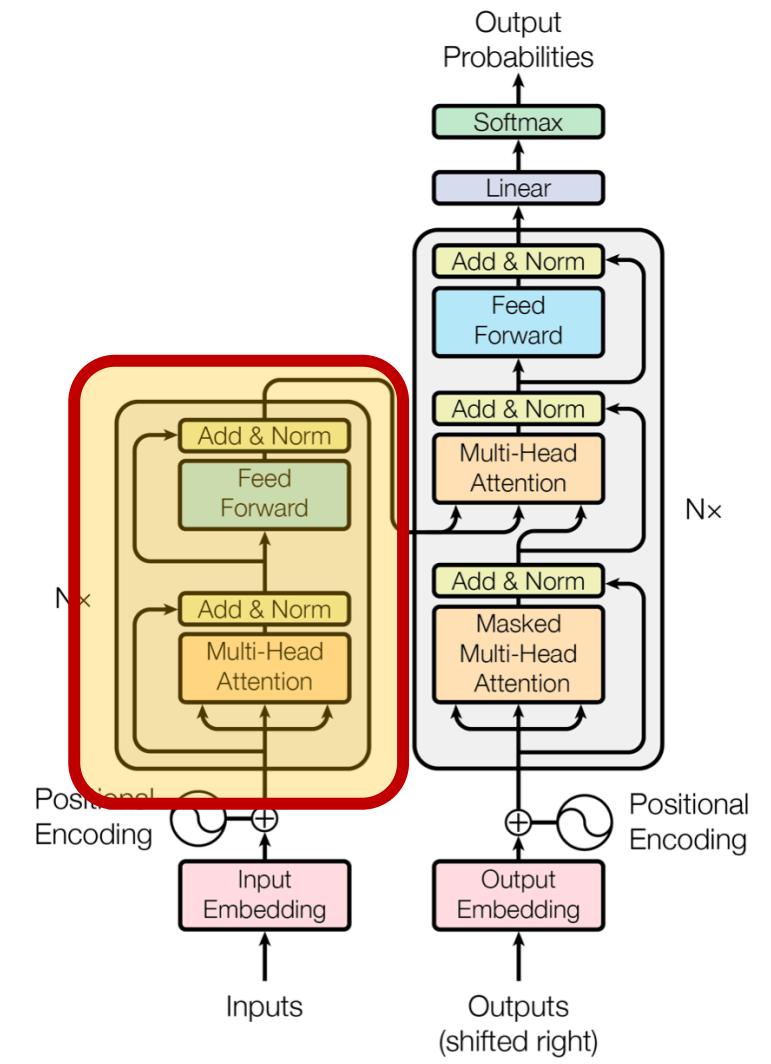
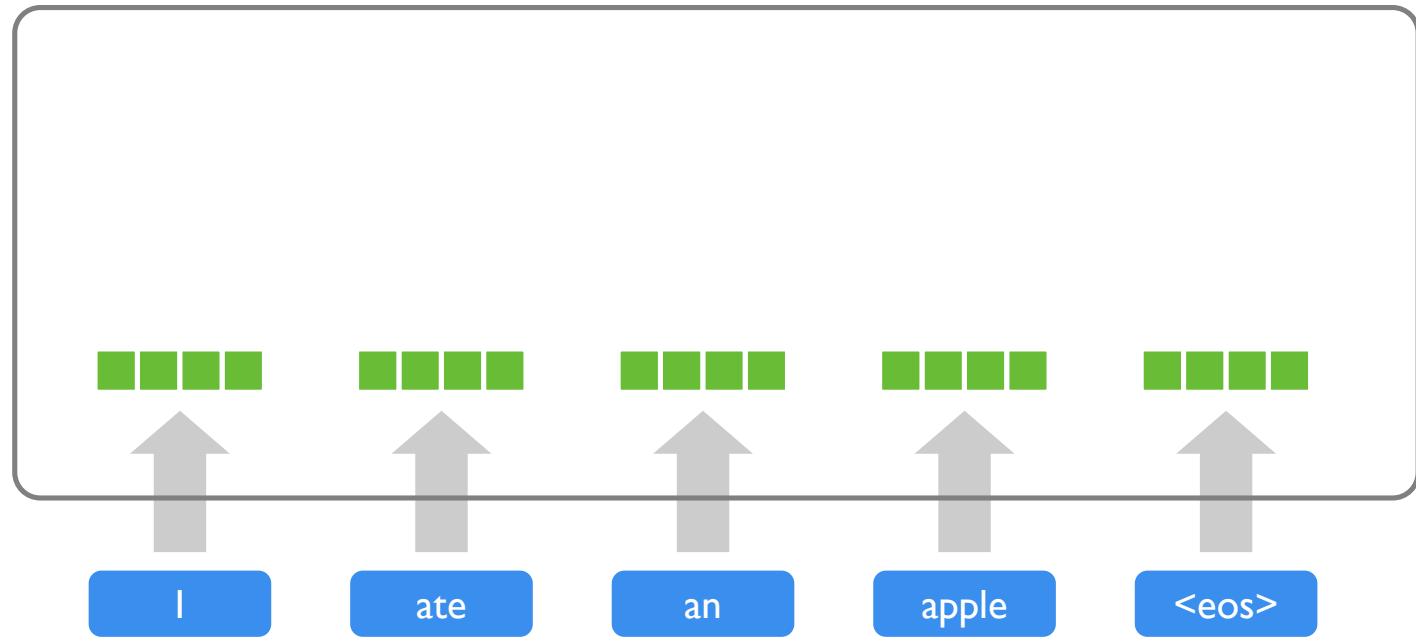


Transformer : Encoder

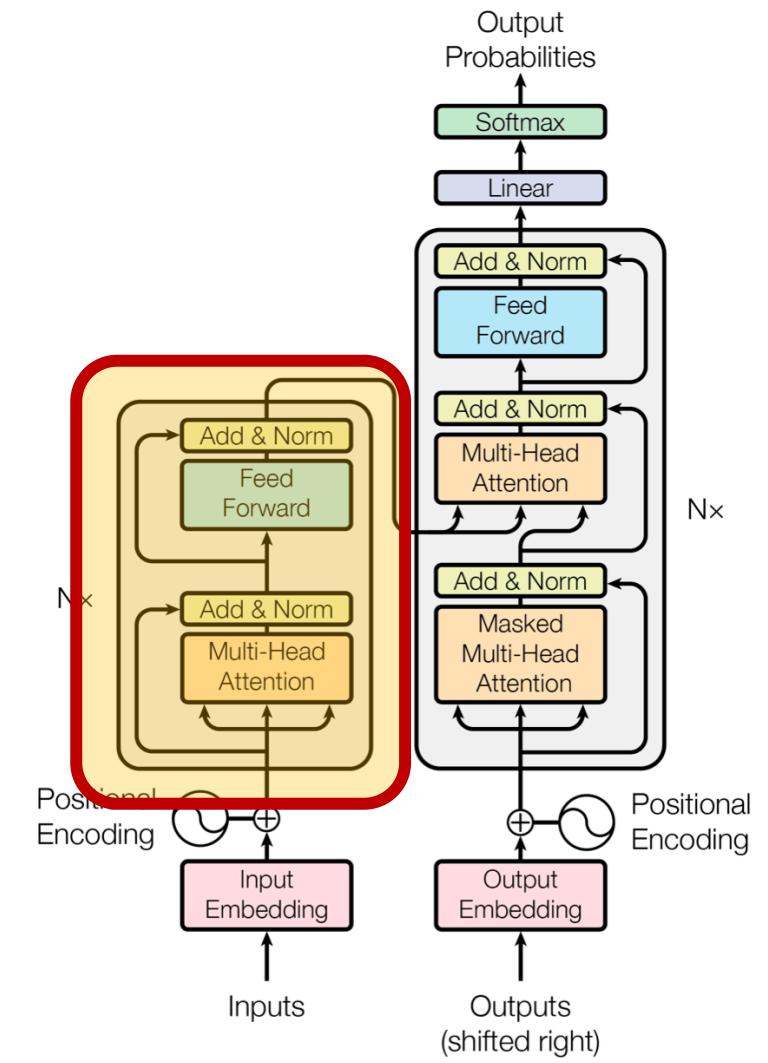
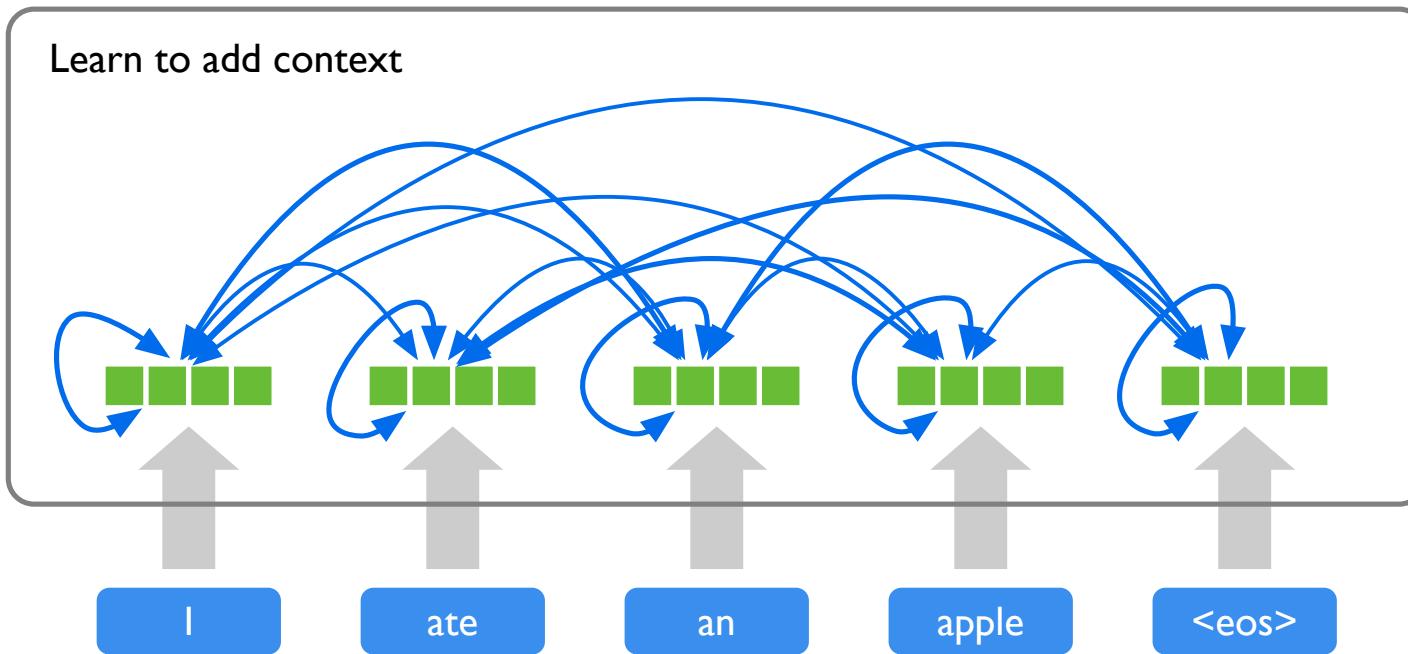
where is the content ?



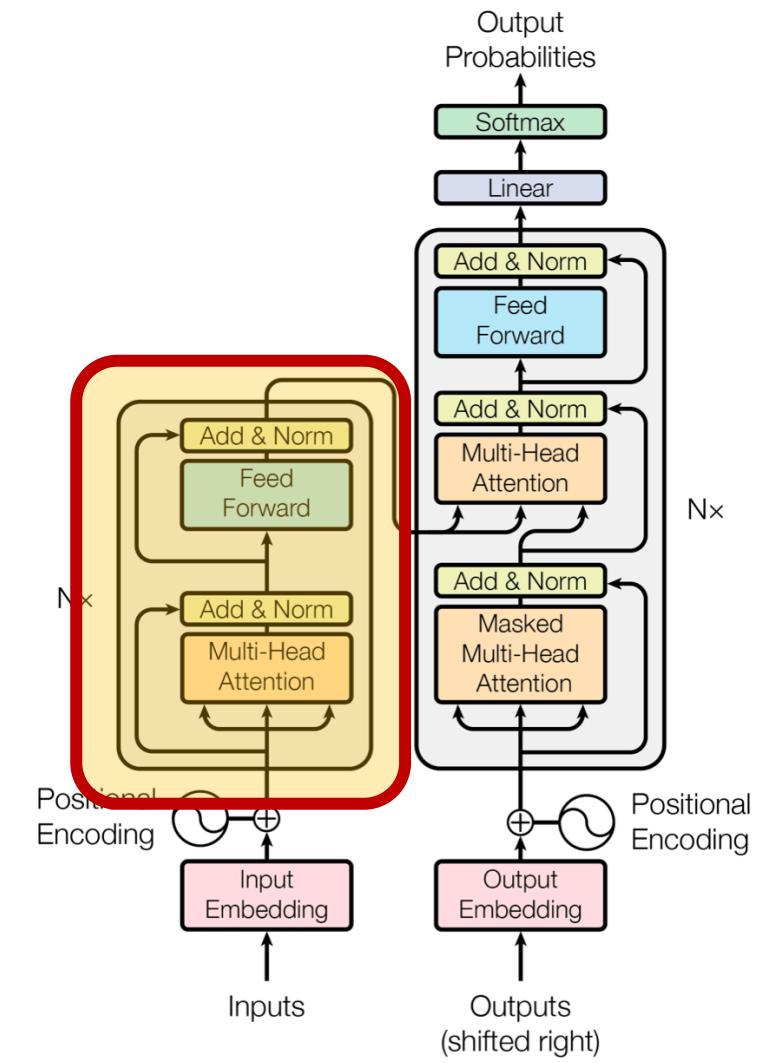
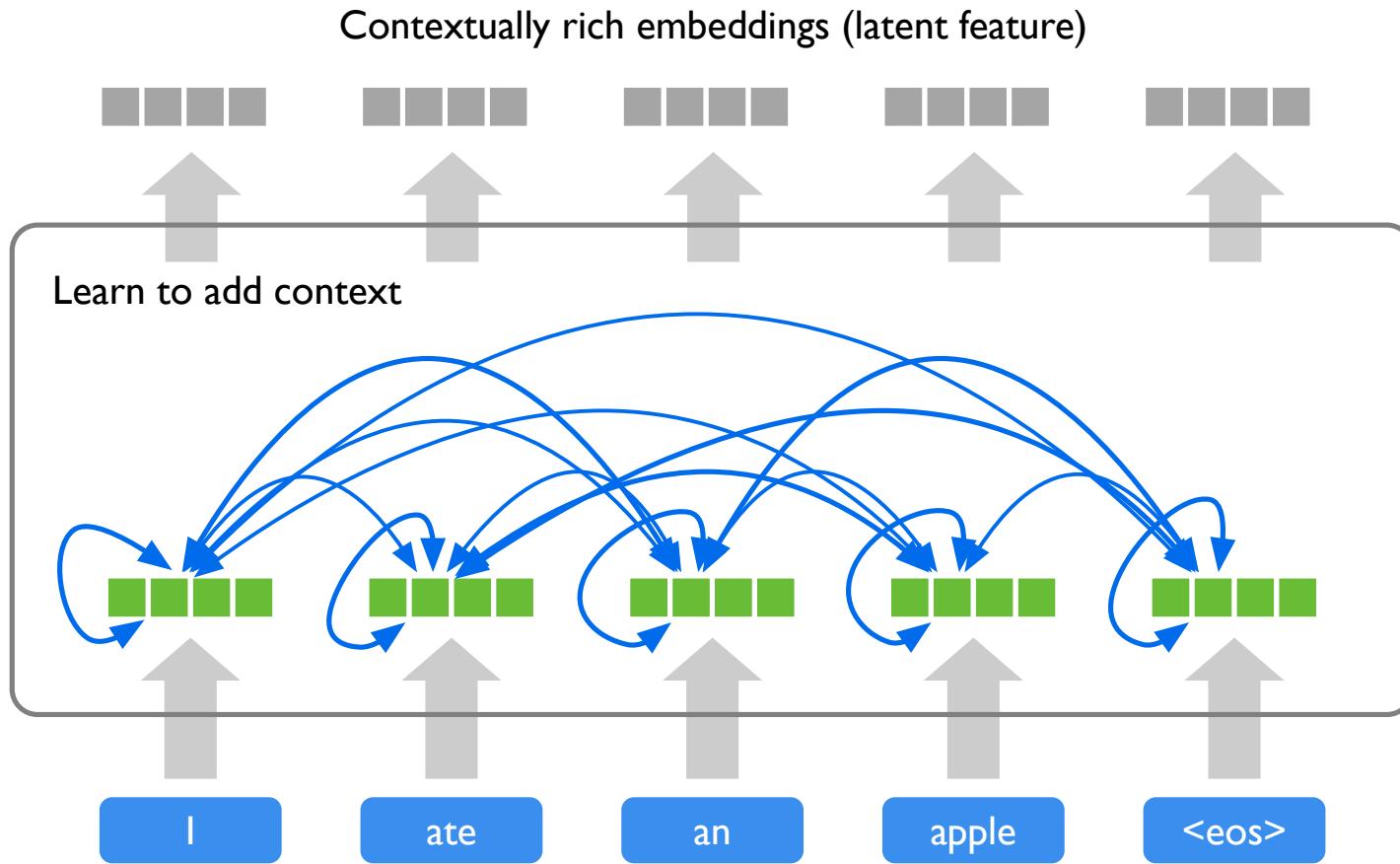
Transformer : Encoder



Transformer : Encoder



Transformer : Encoder

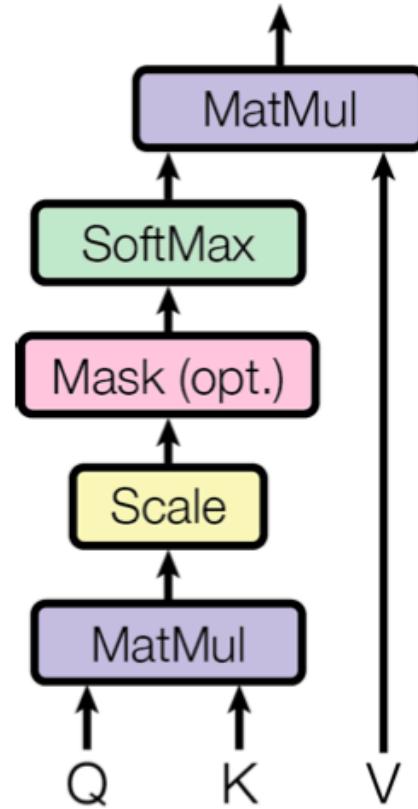


Attention and Multi-Head Attention

- 对当前 Query 和所有的 Key 计算相似度，将这个相似度值通过 Softmax 进行得到一组权重，根据这组权重与对应 Value 乘积求和得到 Attention 下 Value 值。

$$\text{Attention}(Q, K, V) = \text{Softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

Scaled Dot-Product Attention



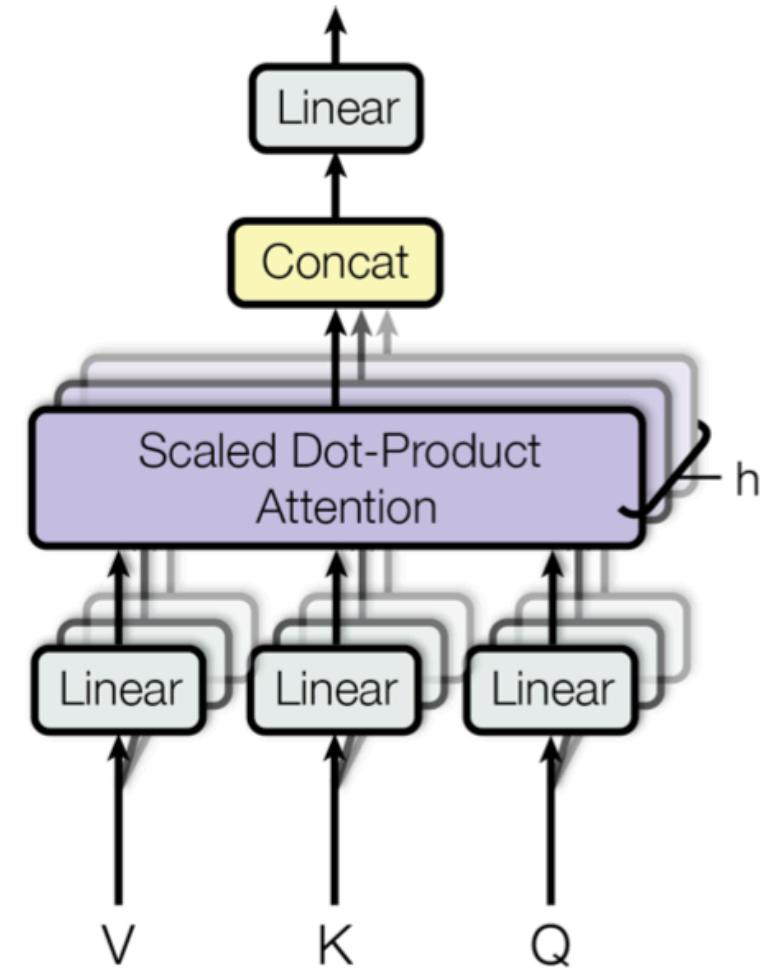
Attention and Multi-Head Attention

- Multi-head Attention 其实就是多个 Self-Attention 结构的结合，每个 head 学习到在不同表示空间中的特征。

$$\text{MultiHead}(Q, K, V) = \text{Concat}(H_1, H_2, \dots, H_h)W^0$$

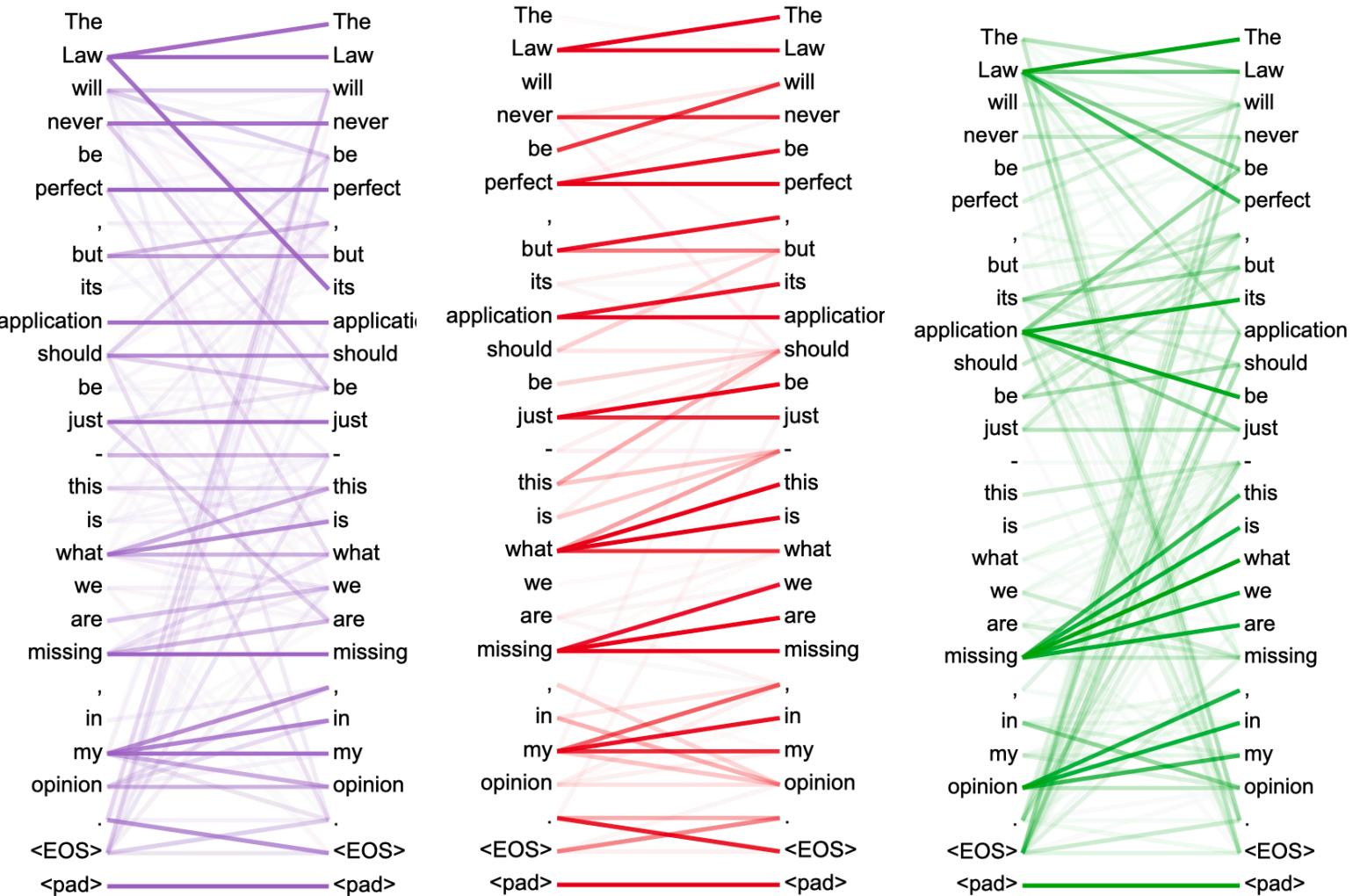
$$H_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V), i = 1, \dots, h$$

Multi-Head Attention



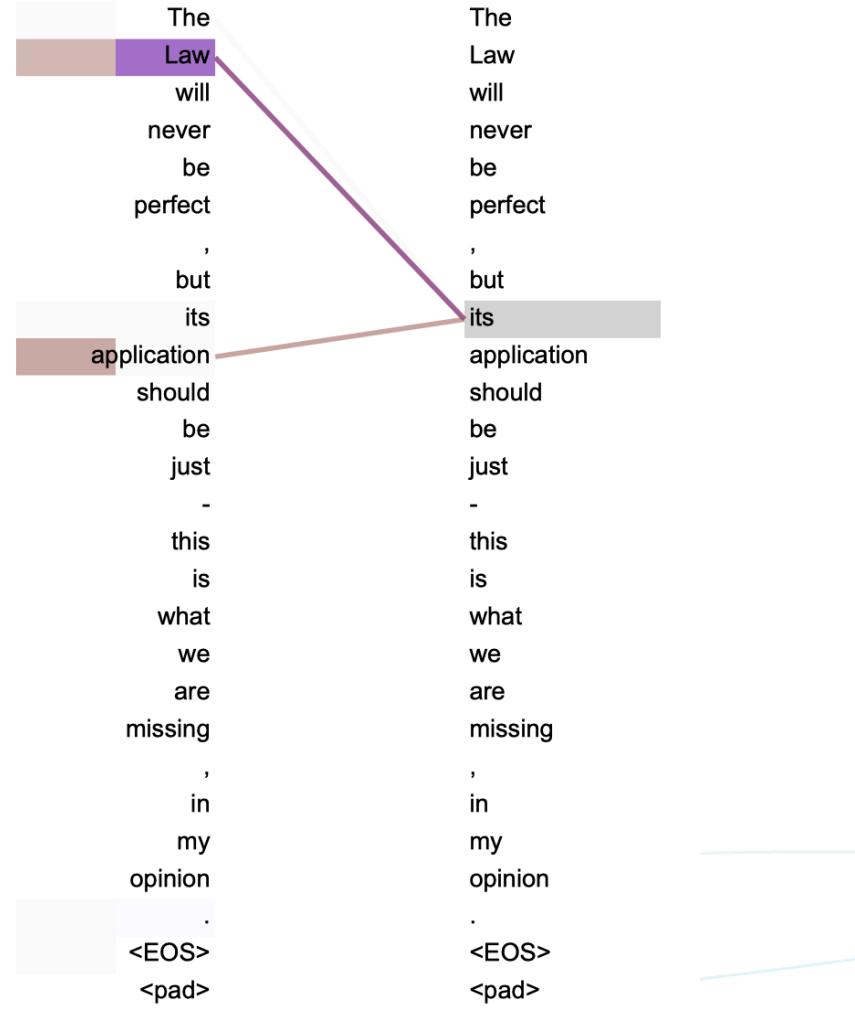
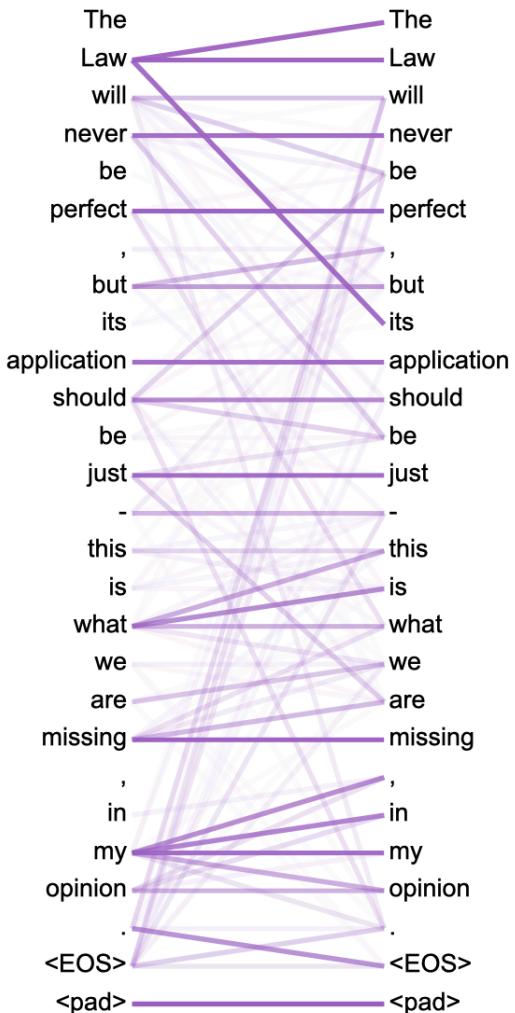
Attention 相关性可视化

- 两个head学习到的Attention侧重点可能略有不同，这样给了模型更大的容量。
- visualization of the attention as computed by one head of the layer 5 of the encoder

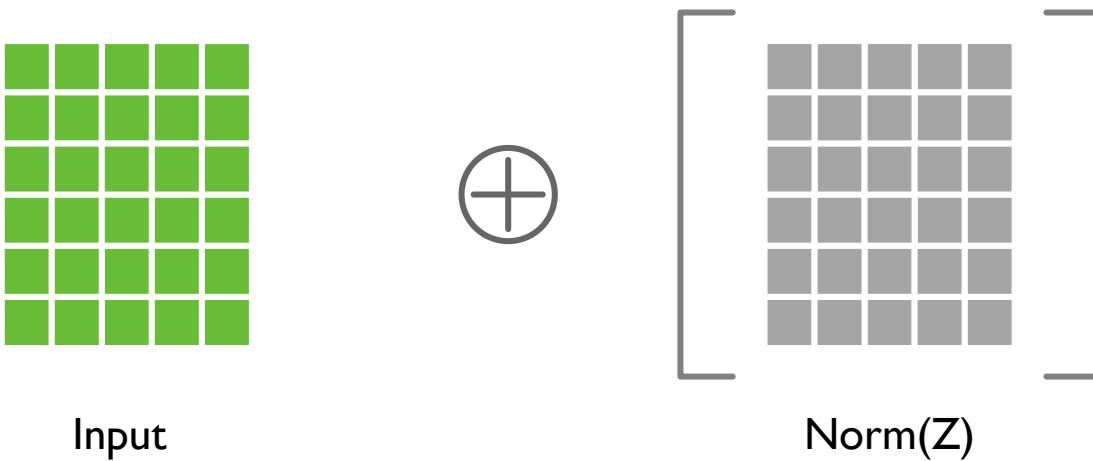


Attention 相关性可视化

- left: Visualization of the attention as computed by one head of the layer 5 of the encoder
- right: Attention given by the word “its” for two different heads is on “law” and “application” which provides help for gender and grammatical issues.



Transformer: Add & Norm

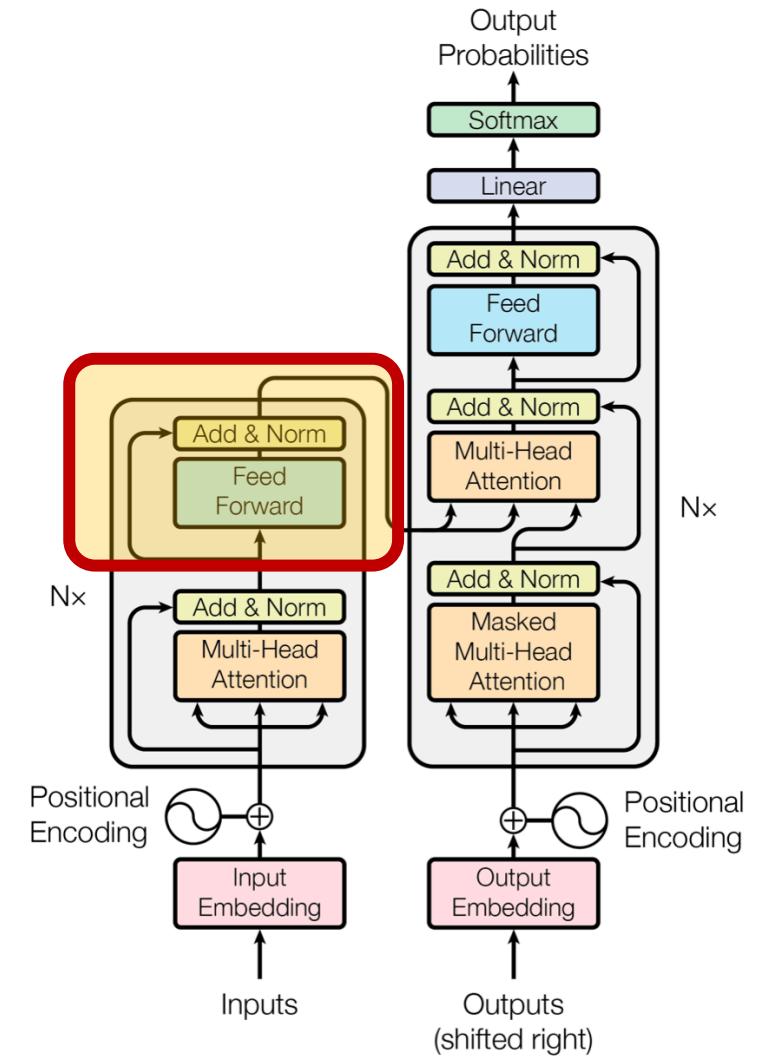


Normalization (Z)

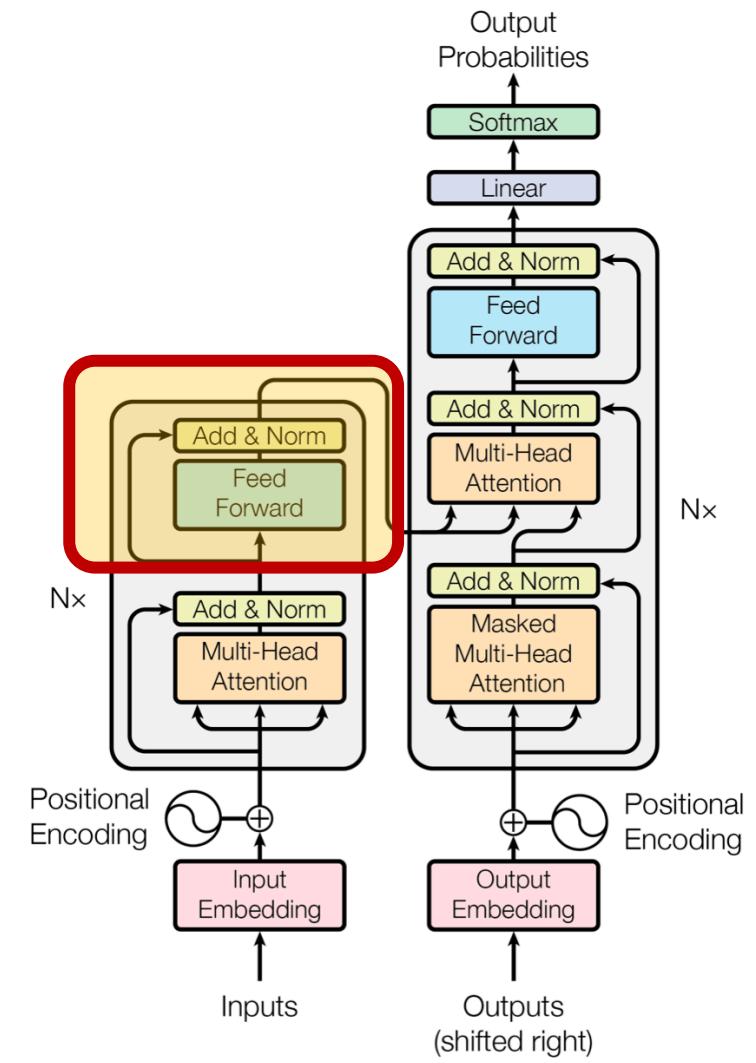
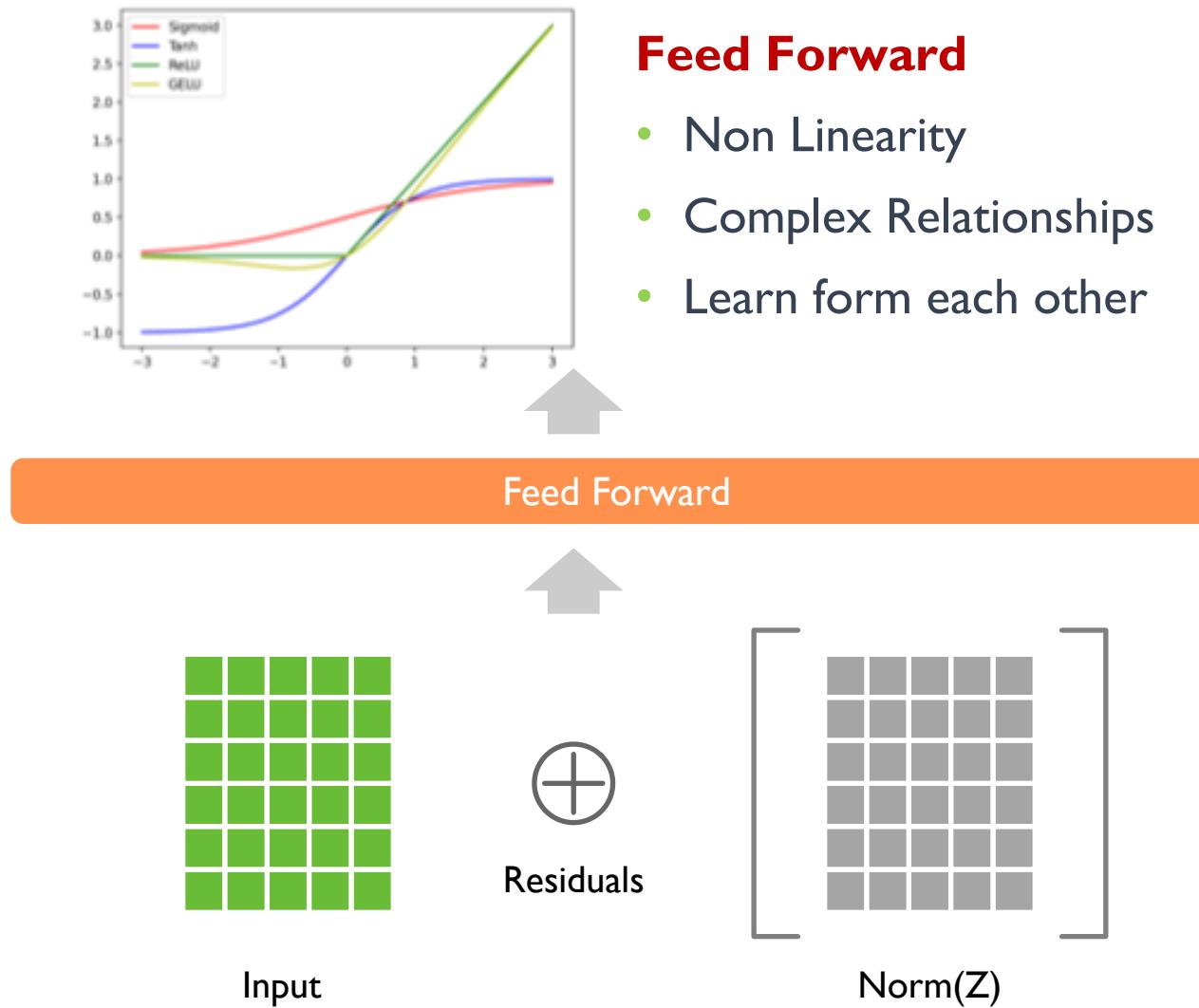
- Mean 0, Std dev 1
- Stabilizes Training
- Regularization Effect

Add → Residuals

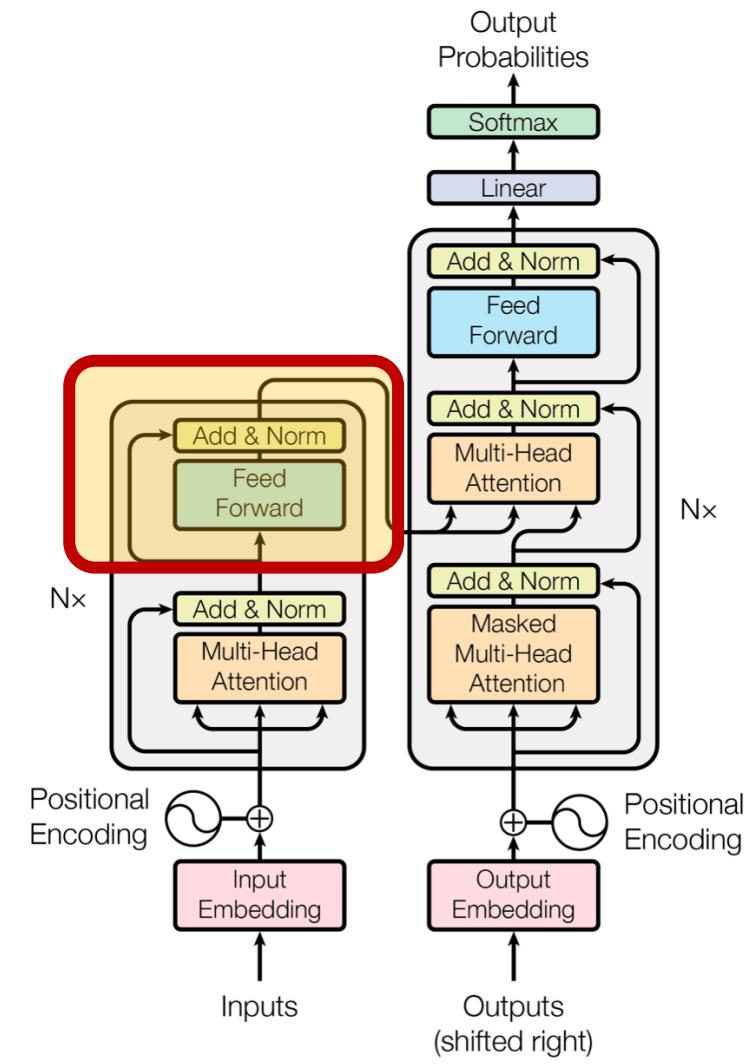
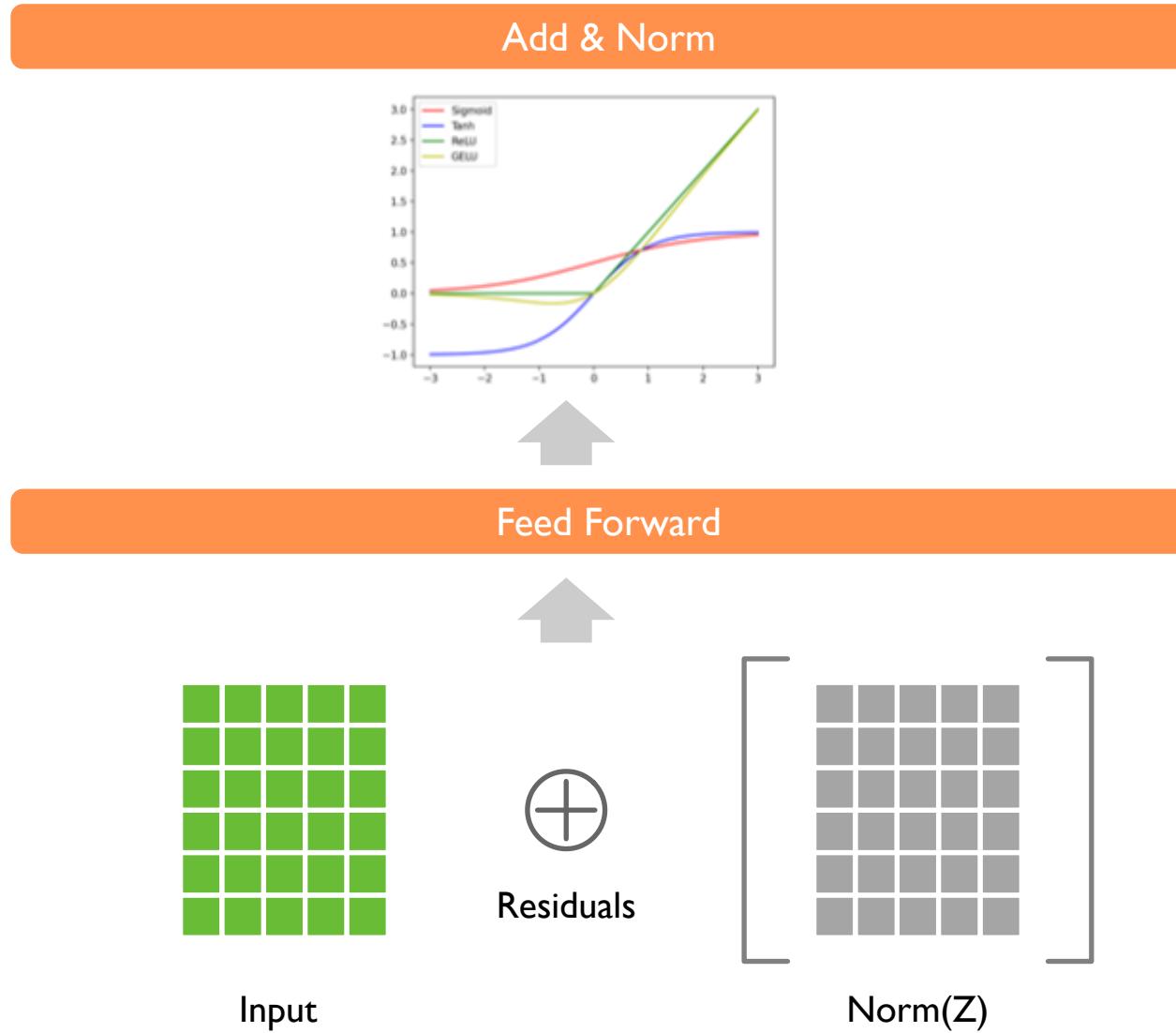
- Avoid vanishing gradients
- Train deeper networks



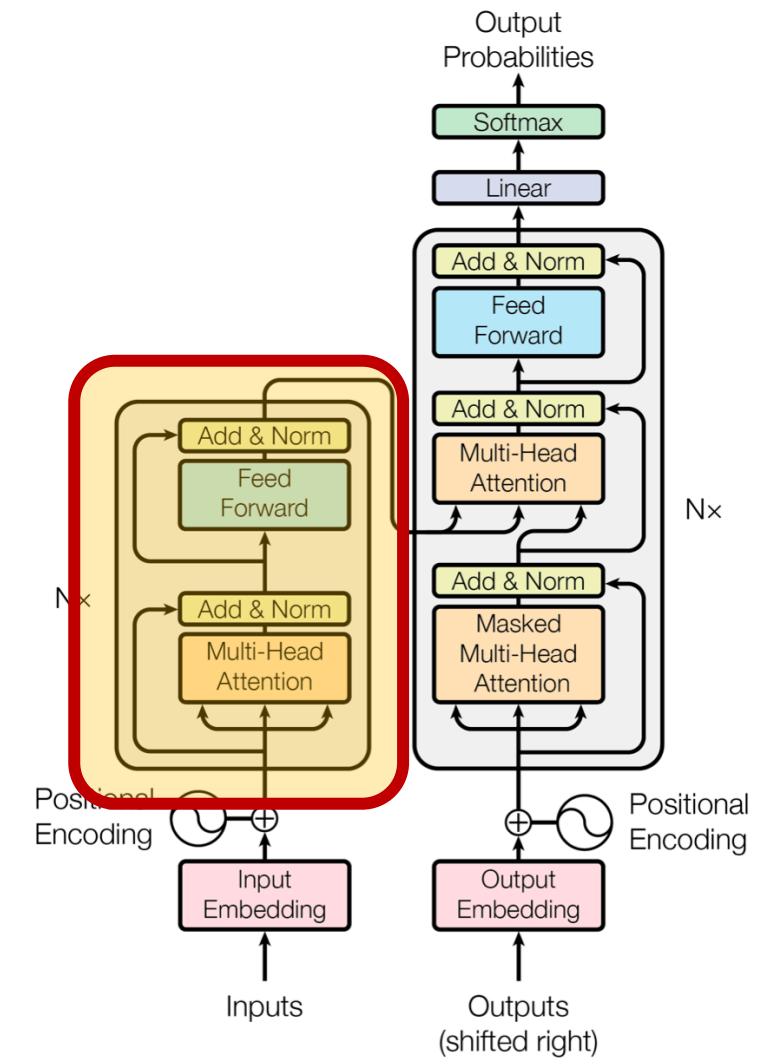
Transformer: Feed Forward



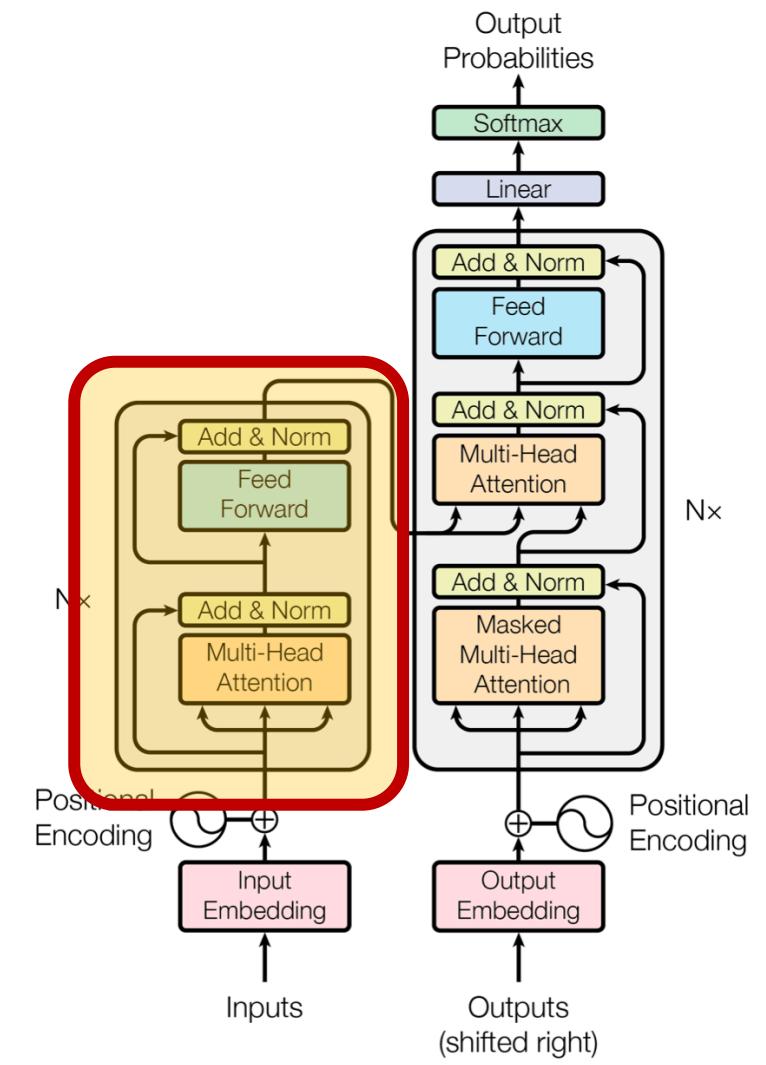
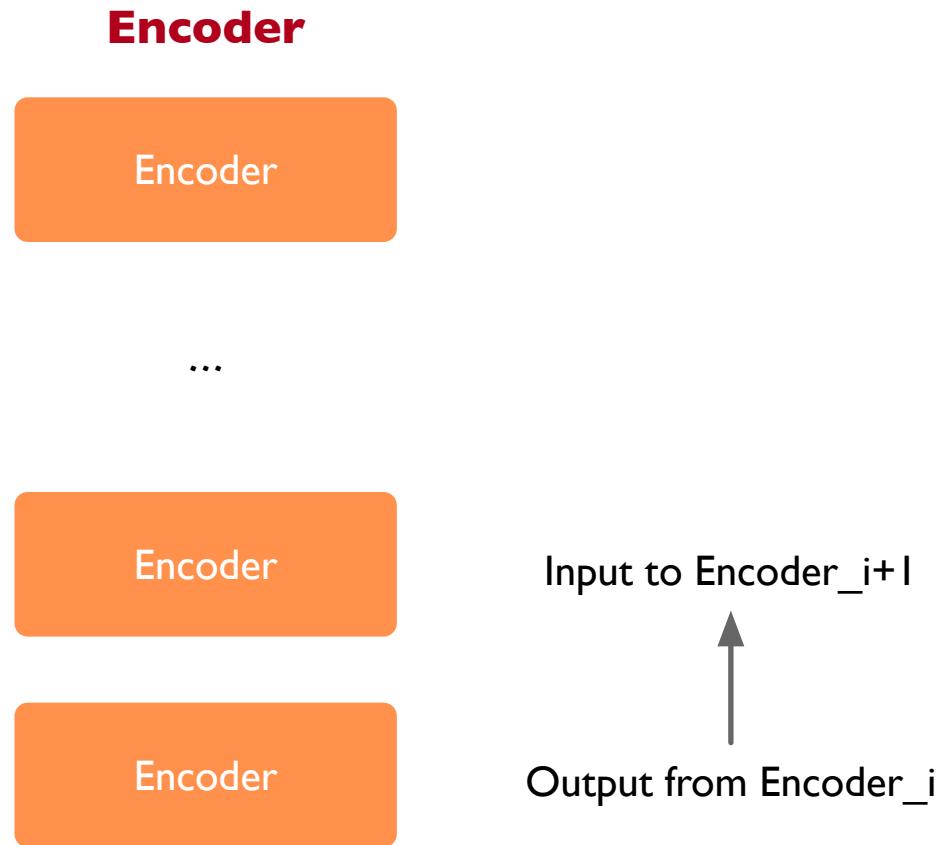
Transformer: Add & Norm



Transformer: Encoder

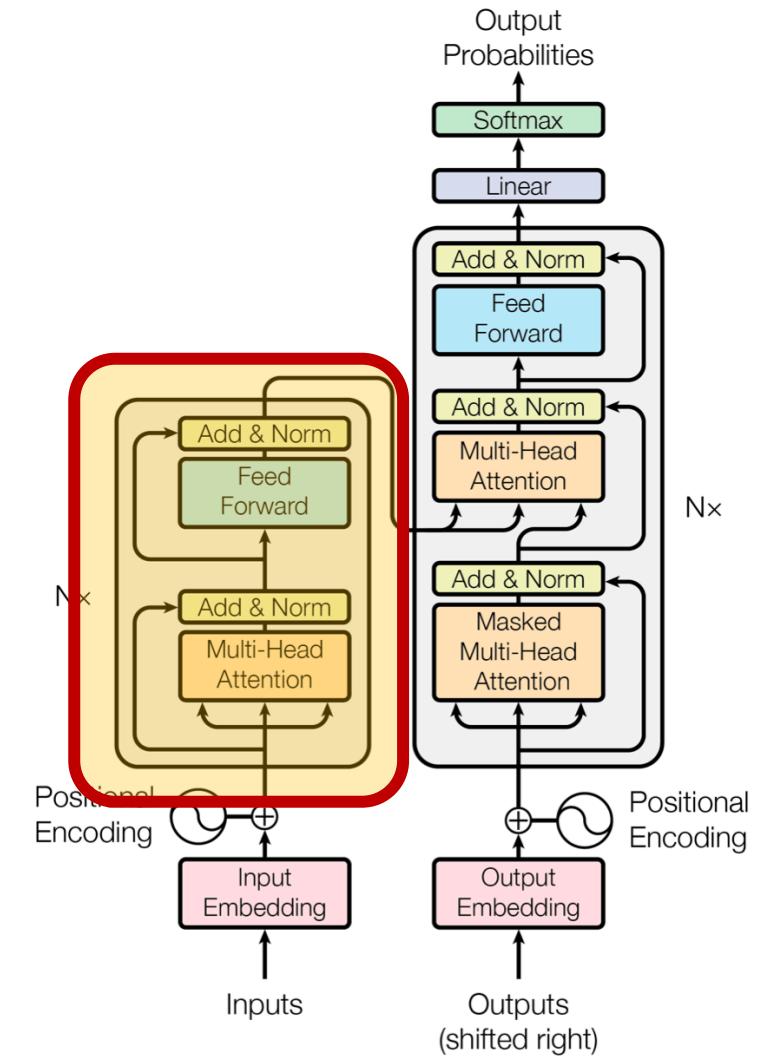


“Original” Transformer (Vaswani et al., 2017).

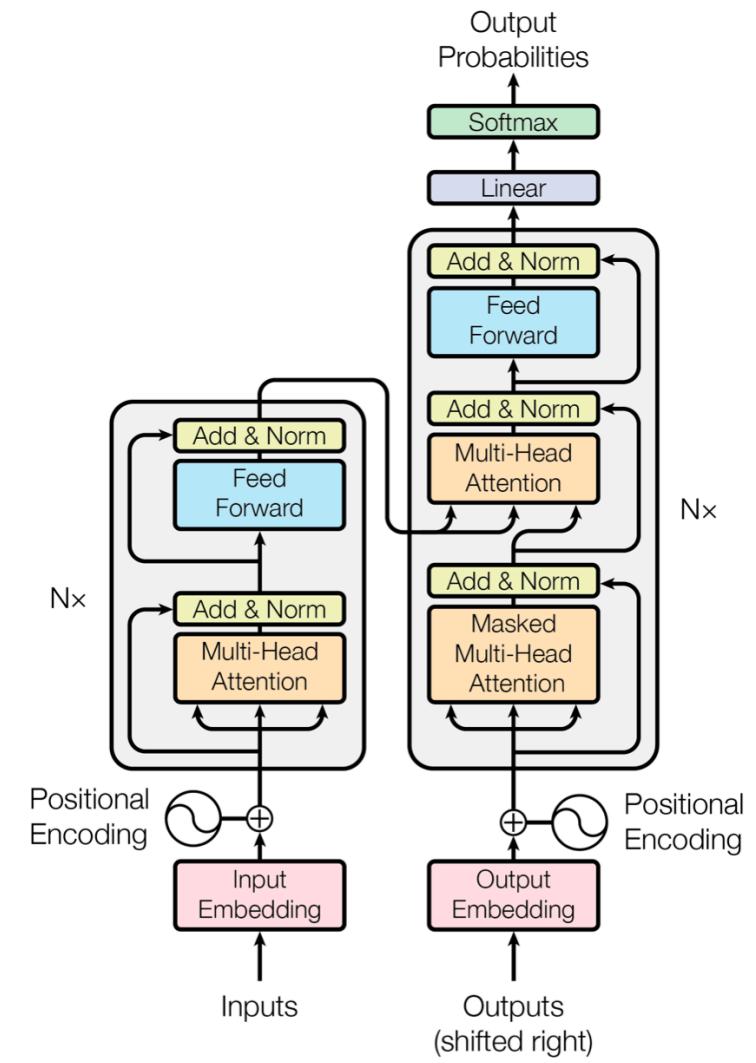
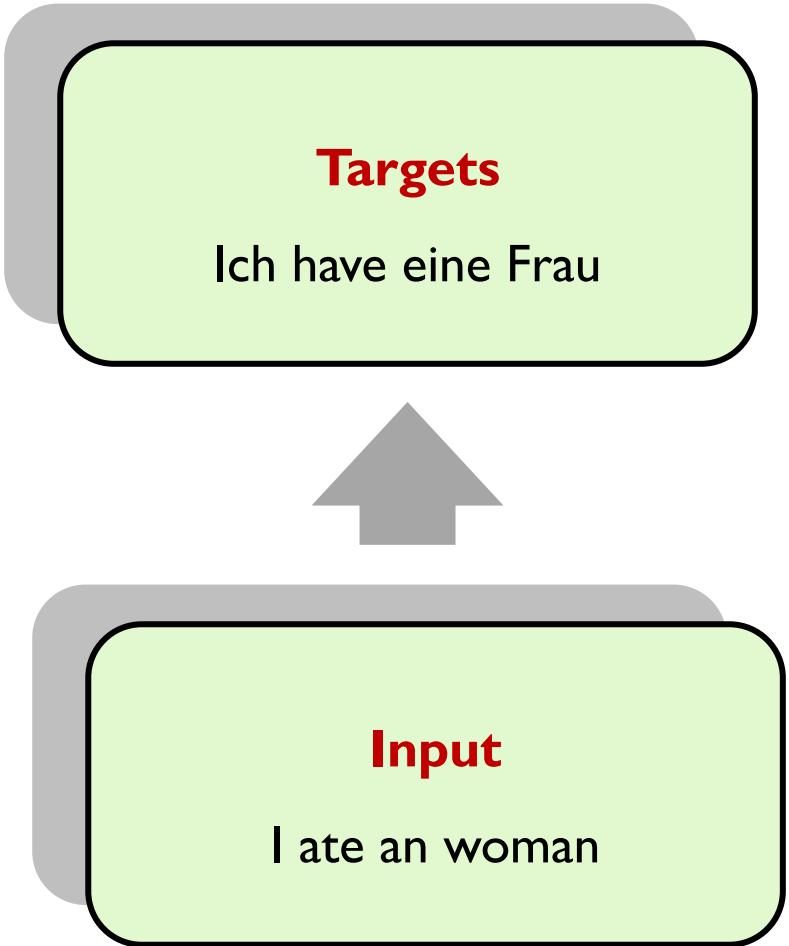


“Original” Transformer (Vaswani et al., 2017).

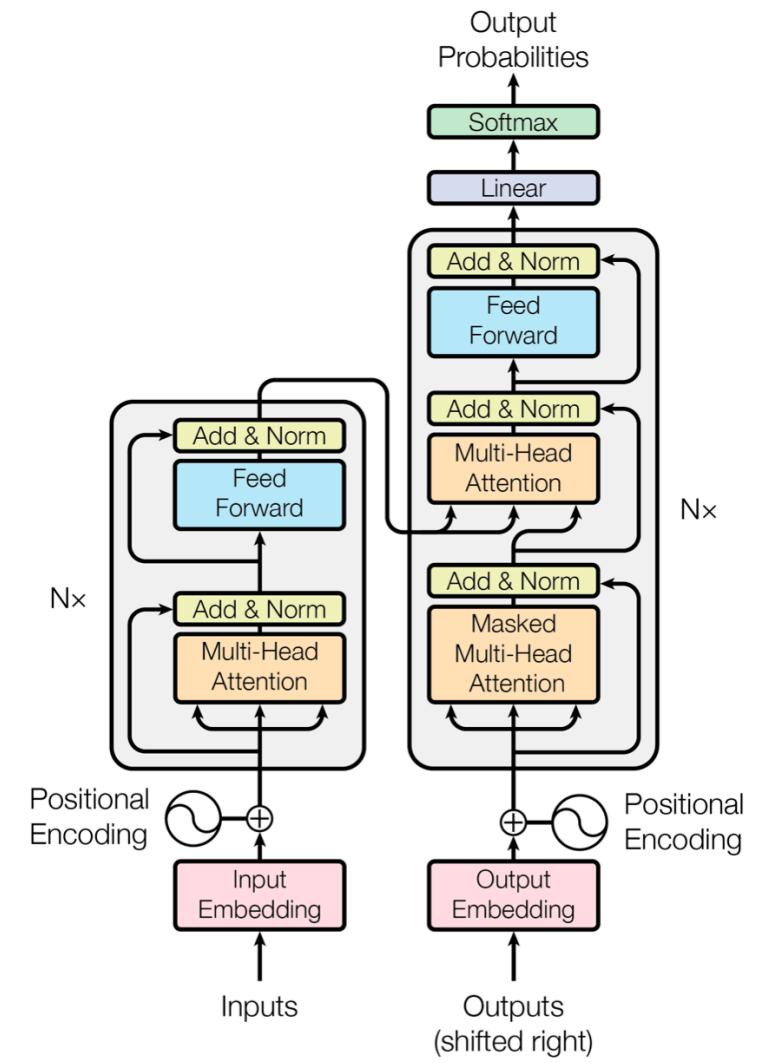
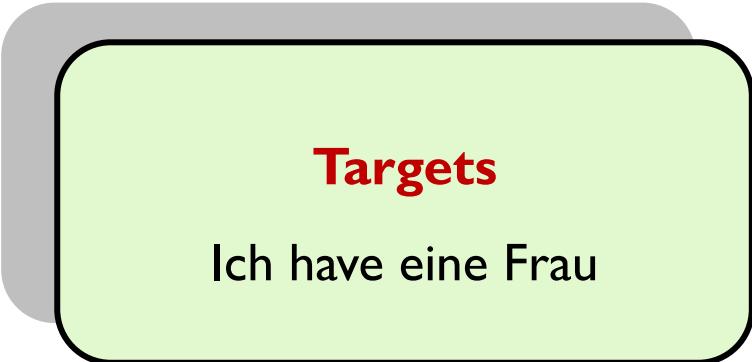
- ✓ Tokenization
- ✓ Input Embeddings
- ✓ Position Encodings
- ✓ Residuals
- ✓ Query
- ✓ Key
- ✓ Value
- ✓ Add & Norm
- ✓ Encoder
- ✓ Decoder
- ✓ Attention
- ✓ Self Attention
- ✓ Multi Head Attention
 - Masked Attention
 - Encoder Decoder Attention
 - Output Probabilities / Logits
 - Softmax
 - Decoder only models



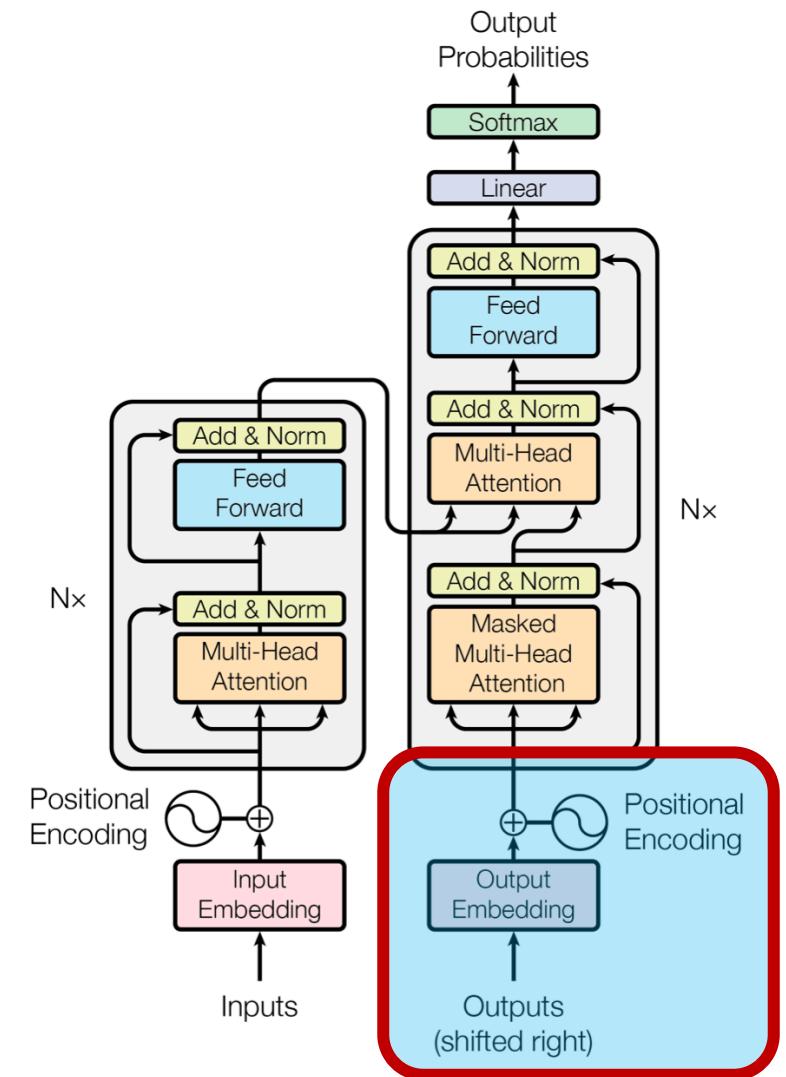
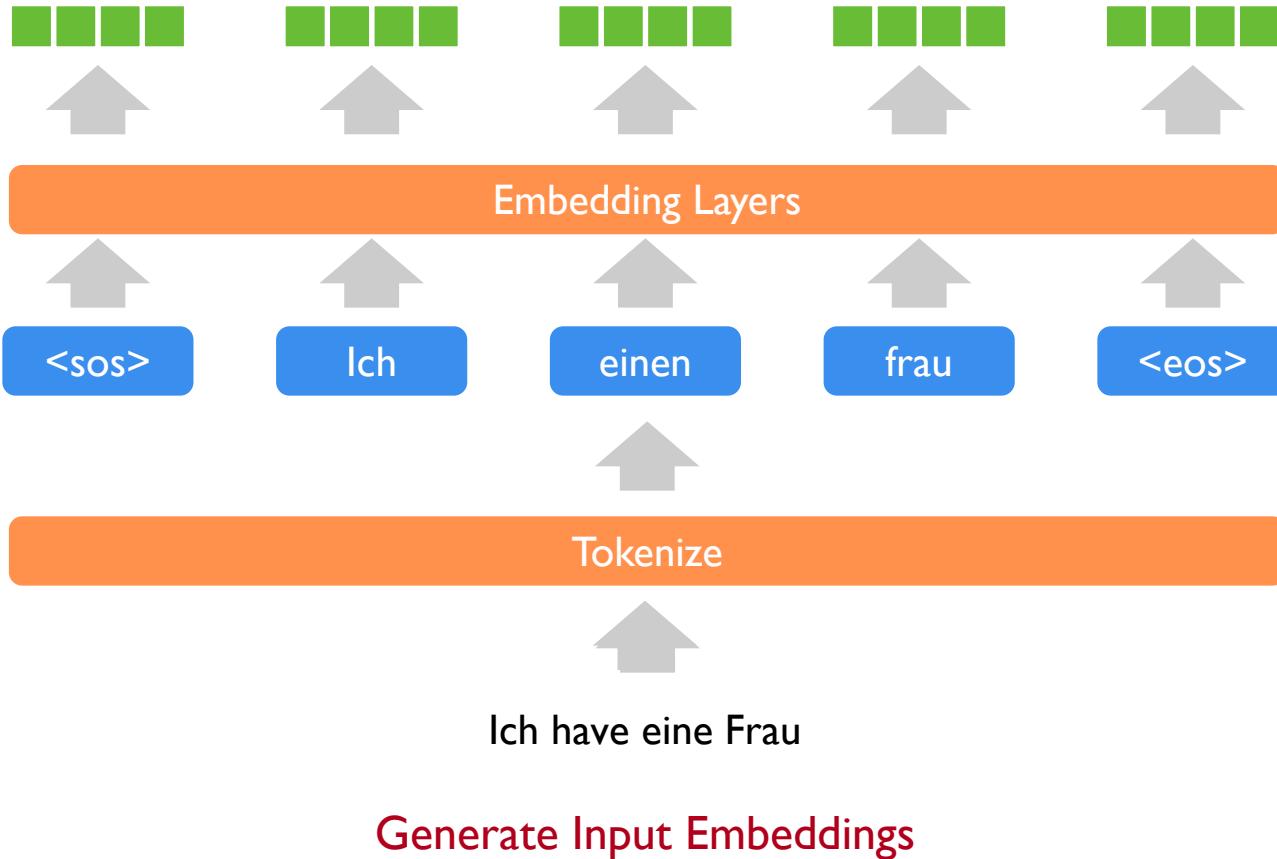
任务：机器翻译 Machine Translation



Transformer: Target

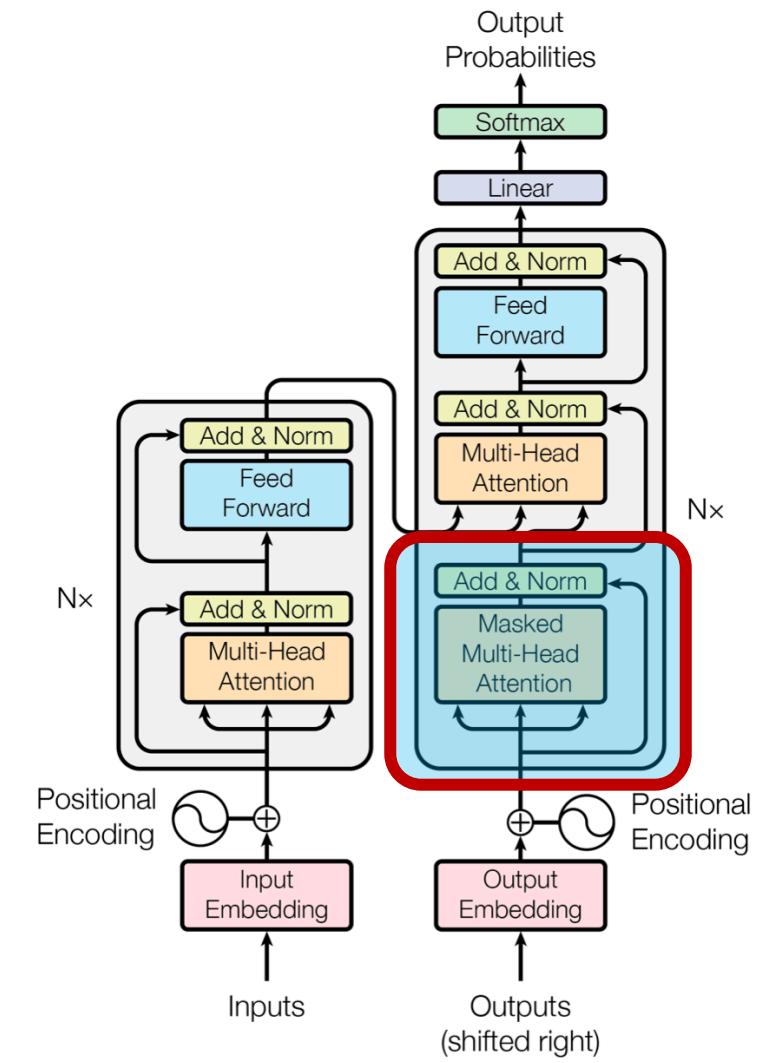
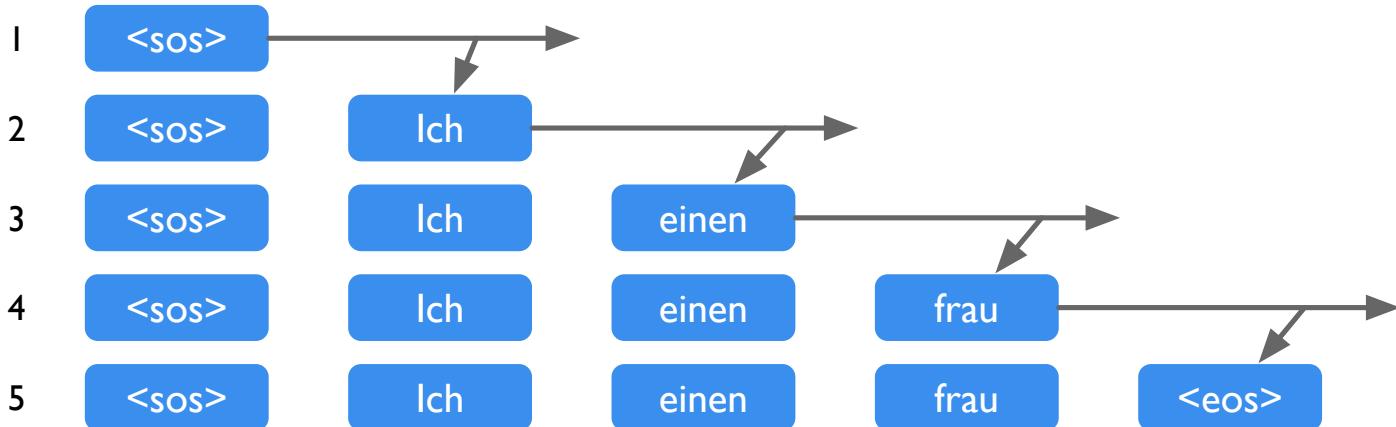


Transformer: Target



Transformer: Masked Multi Head Attention

Inference



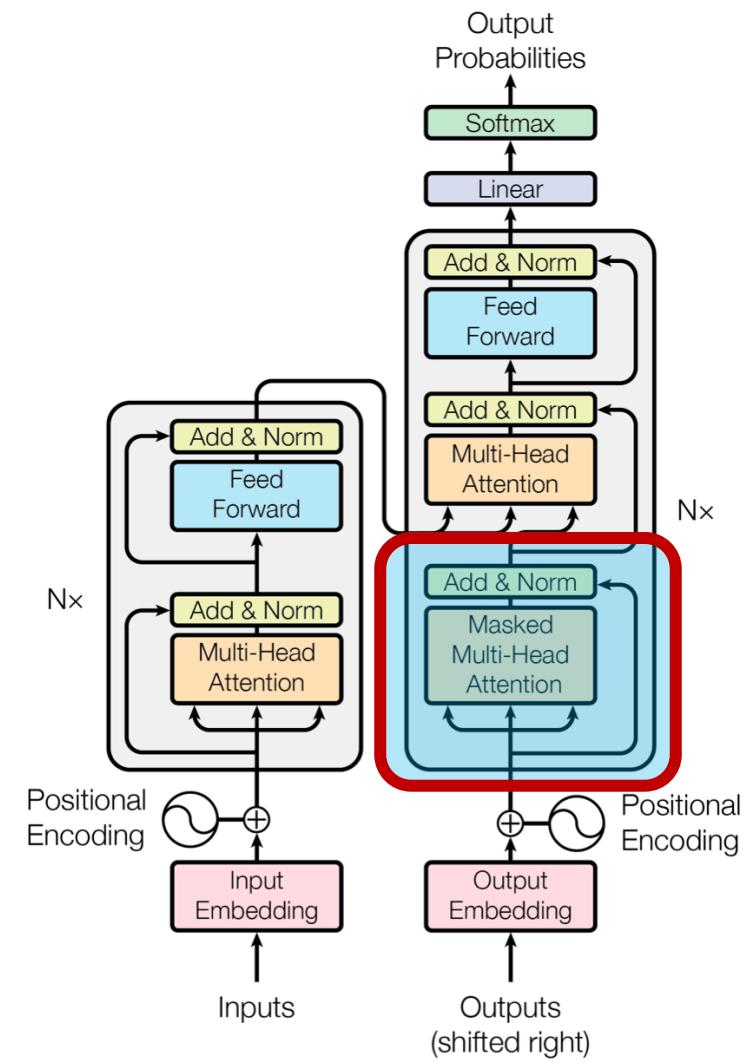
Transformer: Masked Multi Head Attention

Training

n <sos> Ich einen frau <eos>

Outputs at time T should only pay attention to outputs until time T-1

时间T的输出应仅关注时间T-1之前的输出



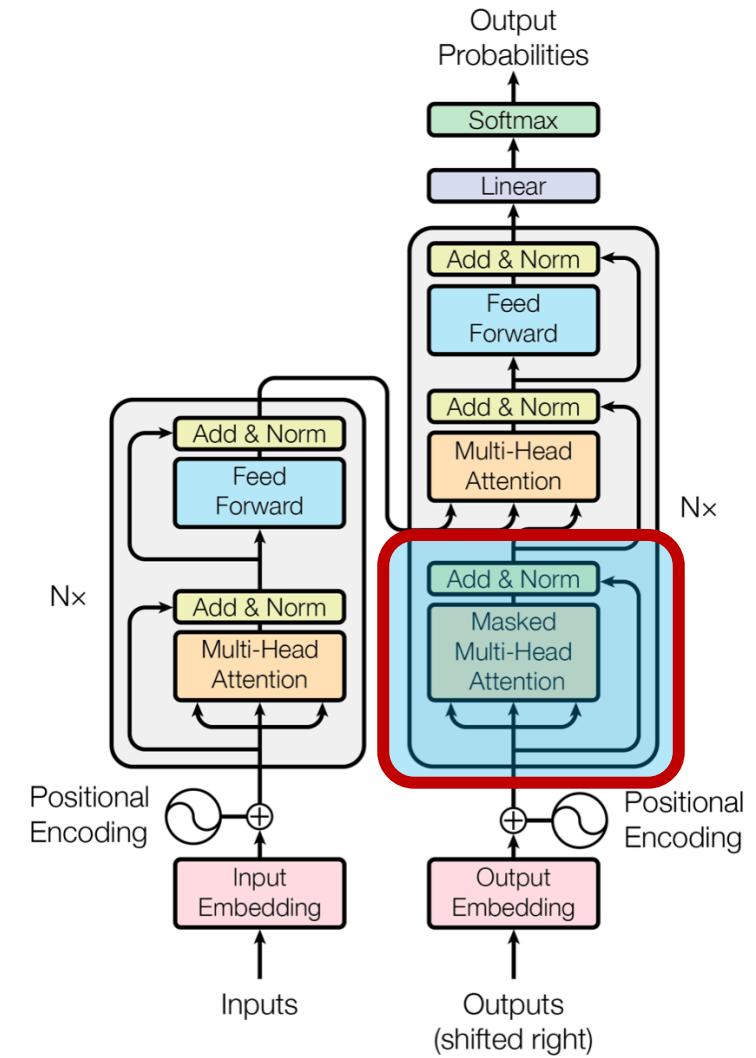
Transformer: Masked Multi Head Attention

Training

1	<sos>	Ich	einen	frau	<eos>
2	<sos>	Ich	einen	frau	<eos>
3	<sos>	Ich	einen	frau	<eos>
4	<sos>	Ich	einen	frau	<eos>
5	<sos>	Ich	einen	frau	<eos>

Outputs at time T should only pay attention to outputs until time T-1

时间T的输出应仅关注时间T-1之前的输出



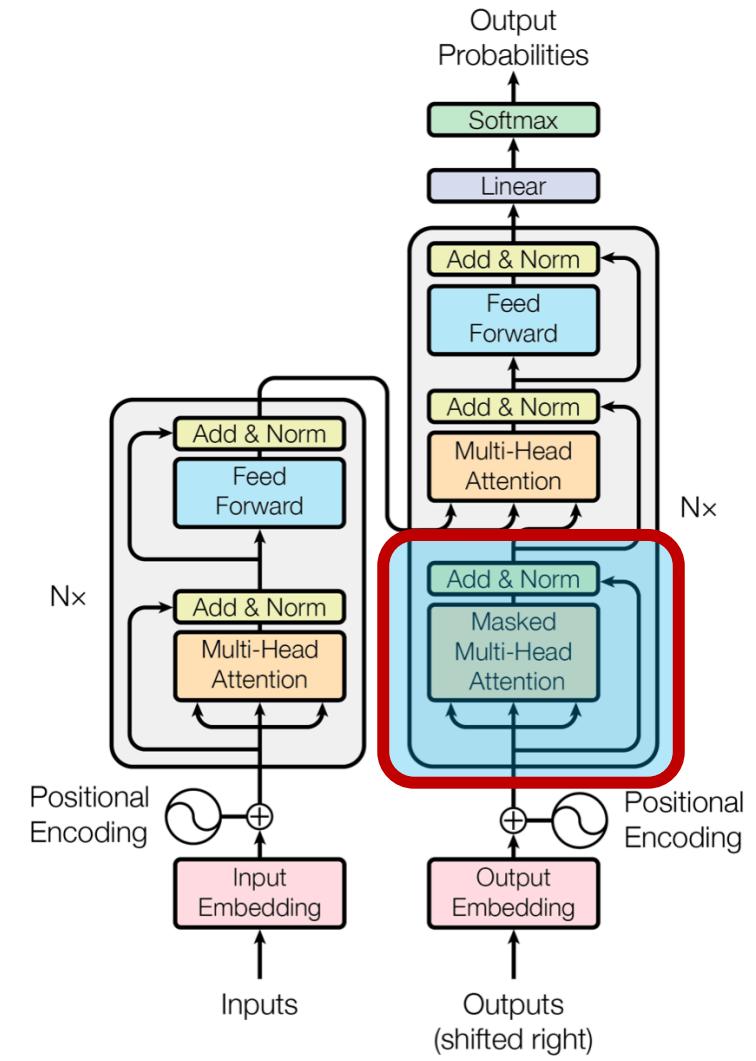
Transformer: Masked Multi Head Attention

Training

1	<sos>	Ich	einen	frau	<eos>
2	<sos>	Ich	einen	frau	<eos>
3	<sos>	Ich	einen	frau	<eos>
4	<sos>	Ich	einen	frau	<eos>
5	<sos>	Ich	einen	frau	<eos>

Mask the available attention values?

掩盖现有的注意力值?



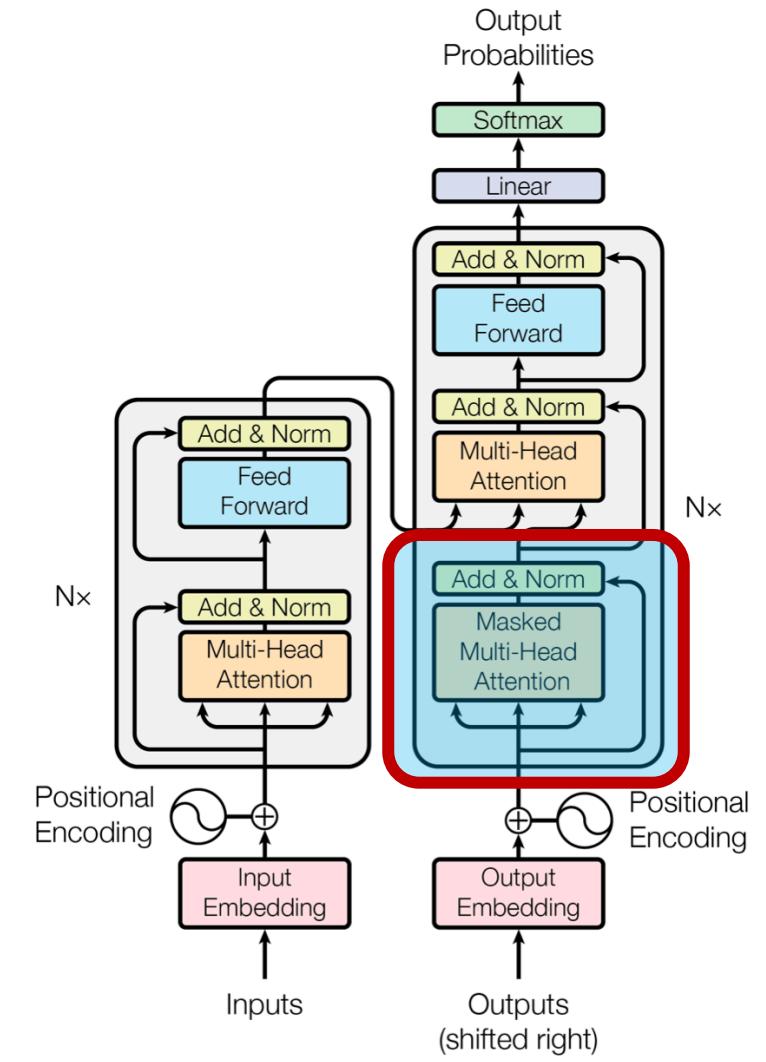
Transformer: Masked Multi Head Attention

$$R^{T \times T} = QK^T + R^{T \times T}$$

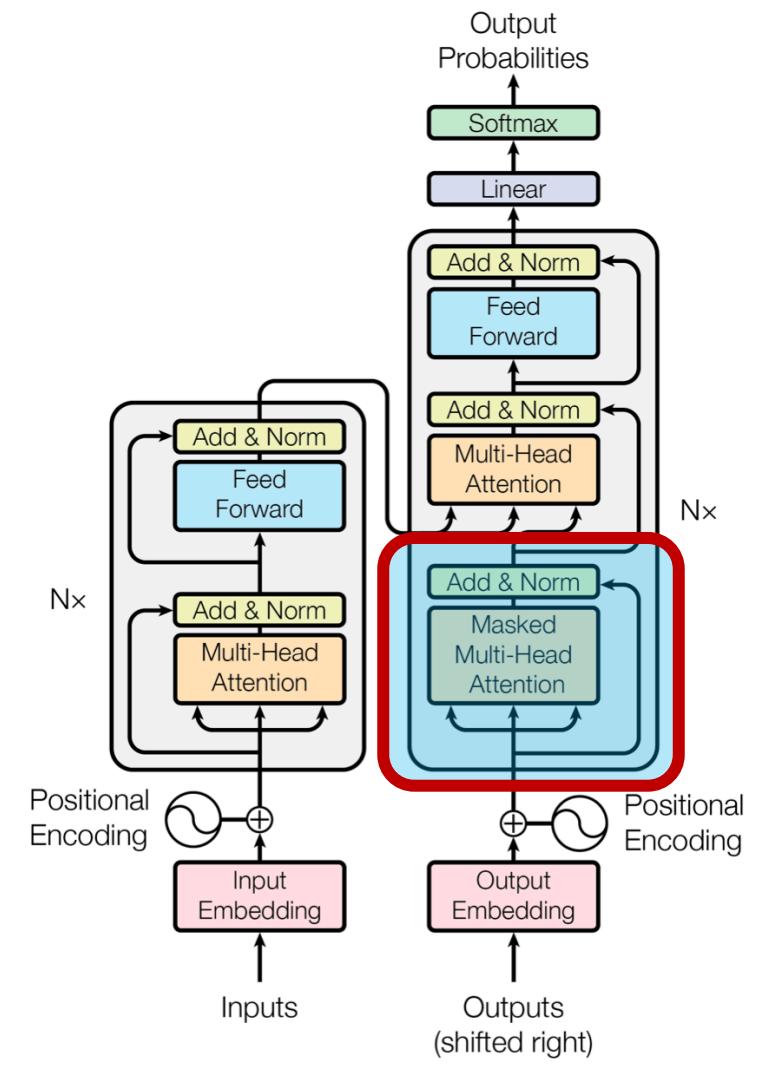
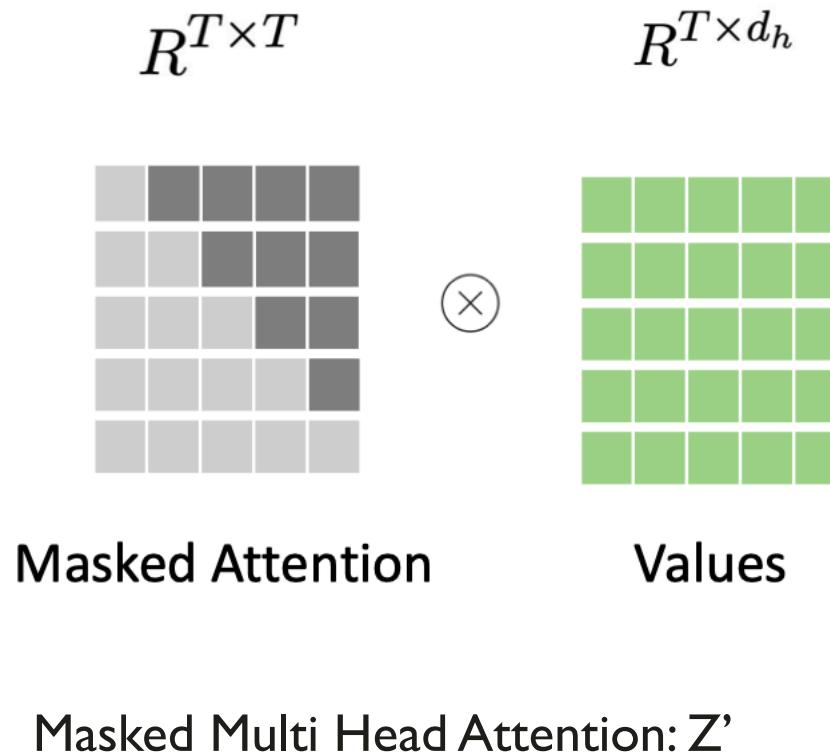
QK^T Attention Mask: M Masked Attention

Masked Multi Head Attention: Z'

The diagram illustrates the computation of Masked Attention. It shows the addition of two tensors: QK^T (represented by a 4x4 grid of gray squares) and an $R^{T \times T}$ tensor (represented by a 4x4 grid where the bottom-right square is black). A plus sign (+) indicates the element-wise addition. The result is labeled "Masked Attention", which is a 4x4 grid where the bottom-right square is dark gray, indicating it has been masked.

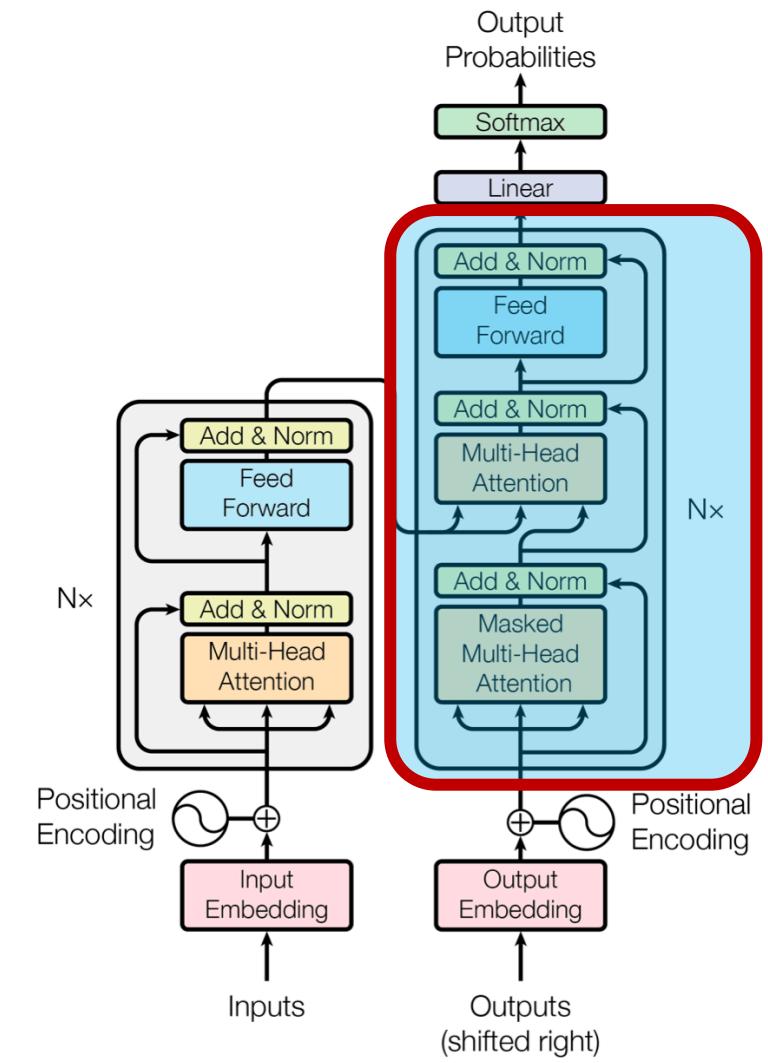
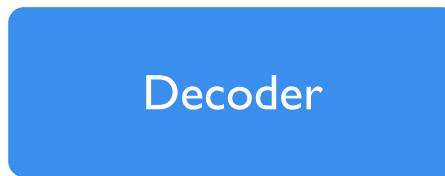


Transformer: Masked Multi Head Attention

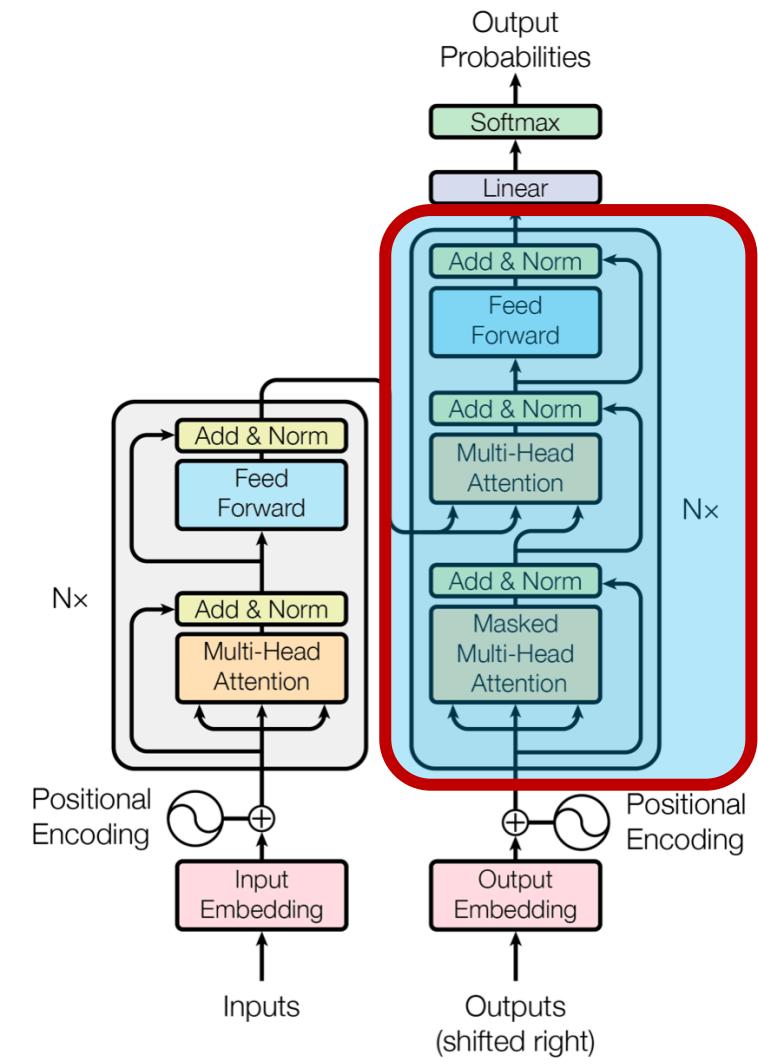
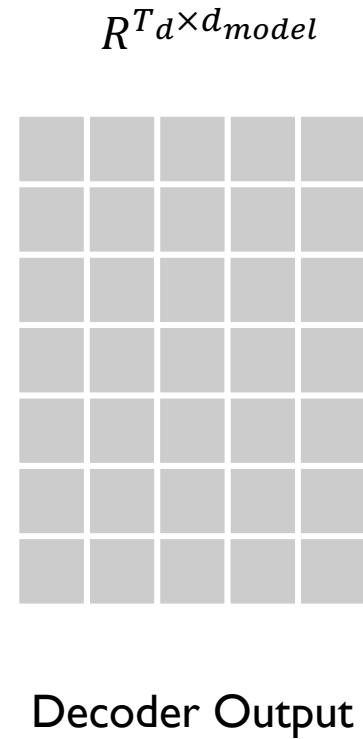
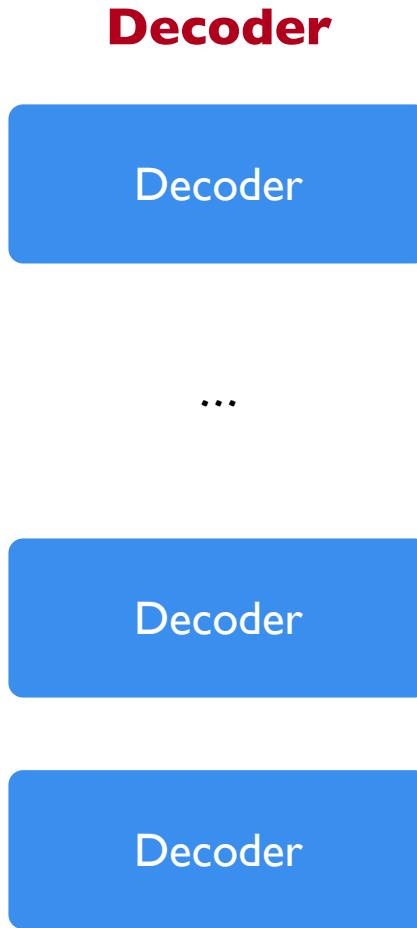


Transformer: Masked Multi Head Attention

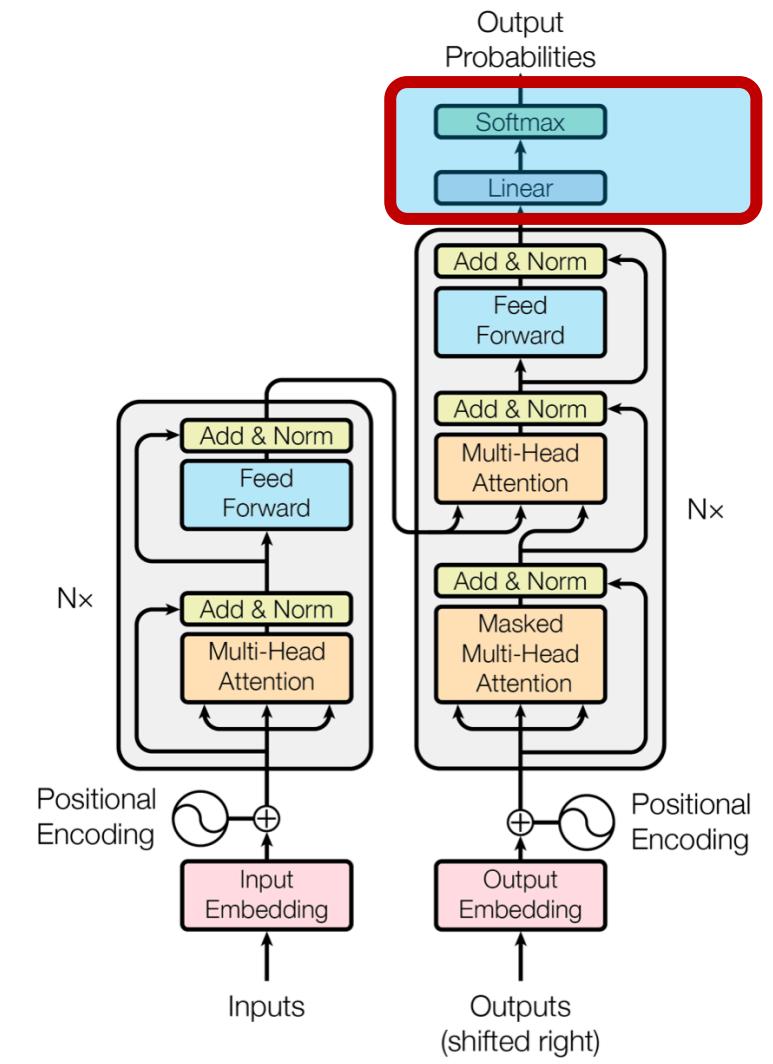
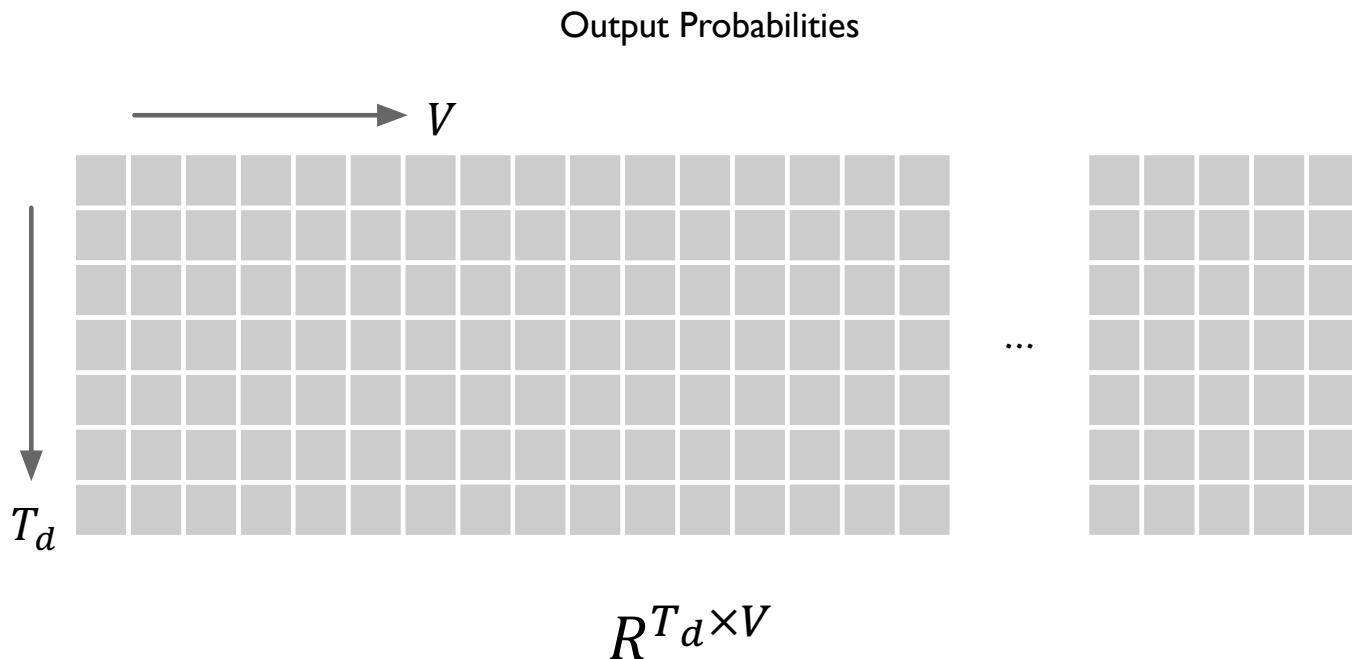
Decoder



Transformer: Masked Multi Head Attention

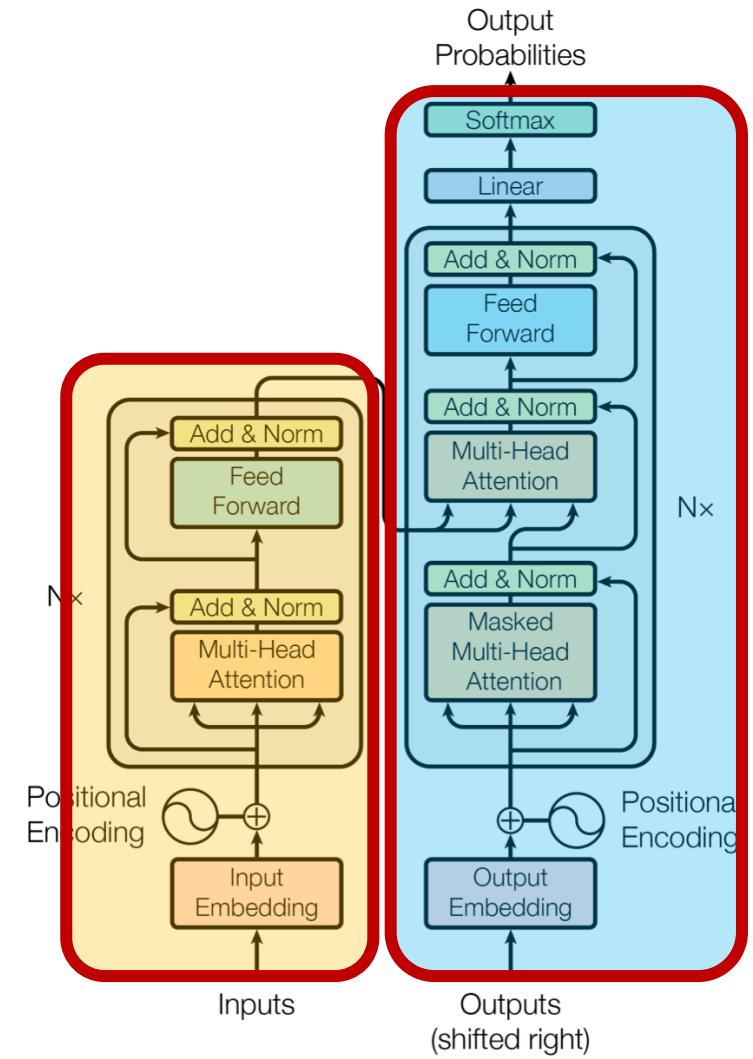


Transformer: Masked Multi Head Attention



“Original” Transformer (Vaswani et al., 2017).

- ✓ Tokenization
- ✓ Input Embeddings
- ✓ Position Encodings
- ✓ Residuals
- ✓ Query
- ✓ Key
- ✓ Value
- ✓ Add & Norm
- ✓ Encoder
- ✓ Decoder
- ✓ Attention
- ✓ Self Attention
- ✓ Multi Head Attention
- ✓ Masked Attention
- ✓ Encoder Decoder Attention
- ✓ Output Probabilities / Logits
- ✓ Softmax
- ✓ Decoder only models



Transformer 最终效果

Table 2: The Transformer achieves better BLEU scores than previous state-of-the-art models on the English-to-German and English-to-French newstest2014 tests at a fraction of the training cost.

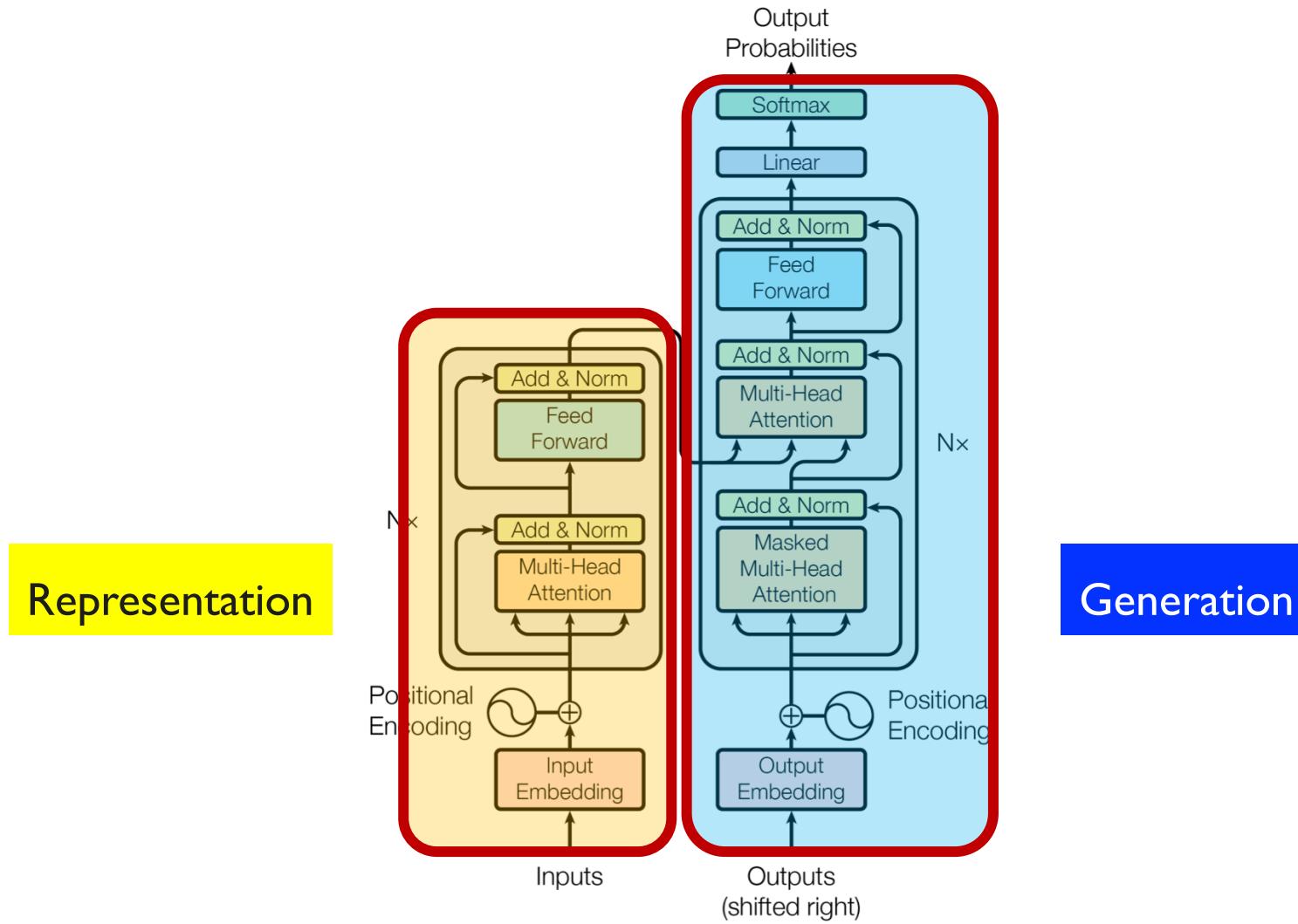
Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [18]	23.75			
Deep-Att + PosUnk [39]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [38]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [9]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [32]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [39]		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble [38]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [9]	26.36	41.29	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1	$3.3 \cdot 10^{18}$	
Transformer (big)	28.4	41.8	$2.3 \cdot 10^{19}$	

Transformer 优点

1. 模型在处理序列输入时，可以对整个序列输入进行并行计算，突破 RNN 不能并行计算限制
2. 相比 CNN，计算两个位置之间的关联所需的操作次数不随距离增长
3. Attention 机制更具可解释性，各注意头 attention head 可以学习不同的关联特征
4. 模型的规模可以随着 Encoder 和 Decoder 的堆叠而不断增加，不需要考虑特征减少

4. BERT vs GPT

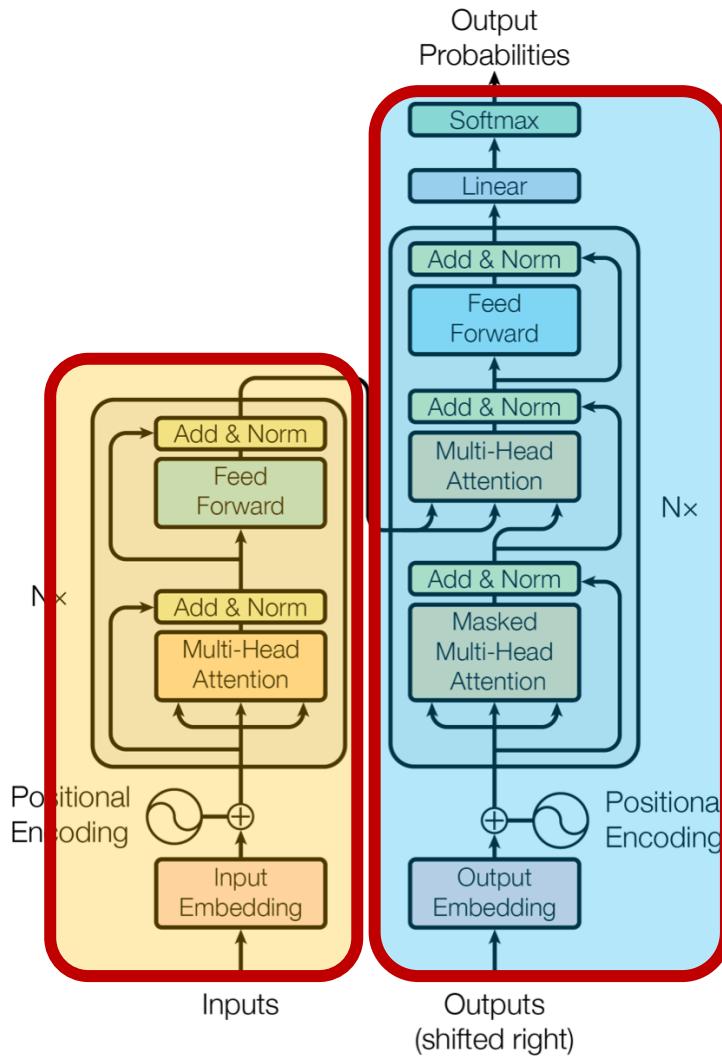
从 Transformers 到大模型



从 Transformers 到大模型

- **Input** – input tokens
- **Output** – hidden states
- Model can see all timesteps
- Does not usually output tokens, so no inherent auto-regresss

Representation



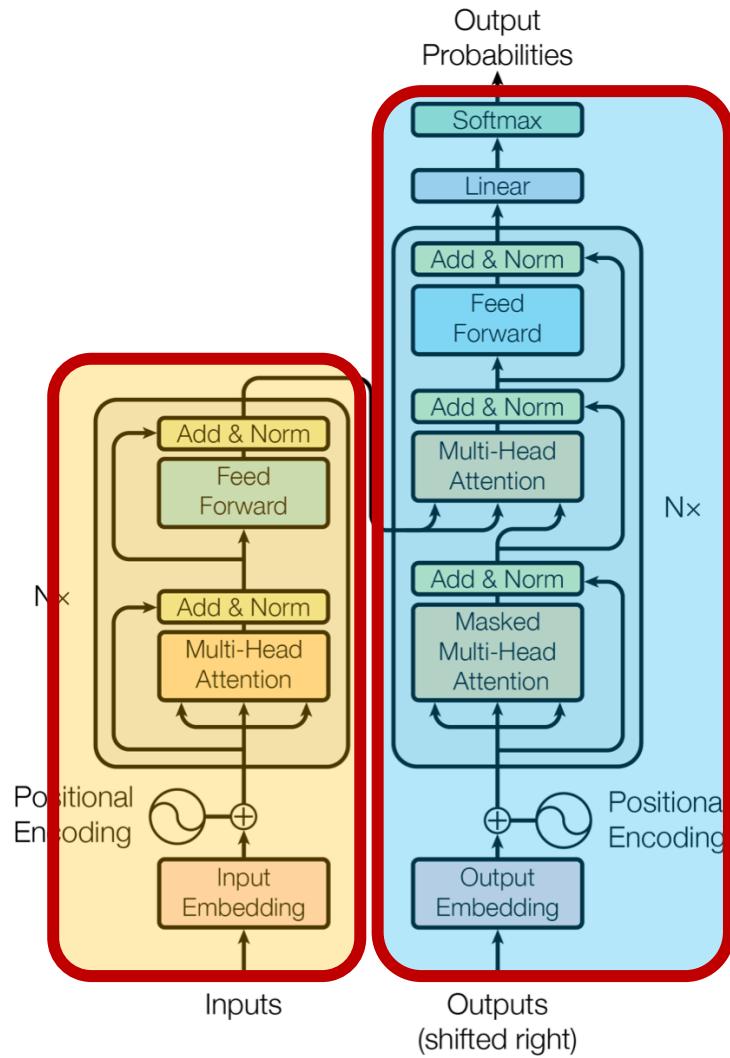
- **Input** – output tokens and hidden states*
- **Output** – output tokens
- Model can only see previous timesteps
- Model is auto-regressive with previous timesteps' outputs

Generation

2018 – The Inception of the LLM Era

BERT
Oct 2018

Representation

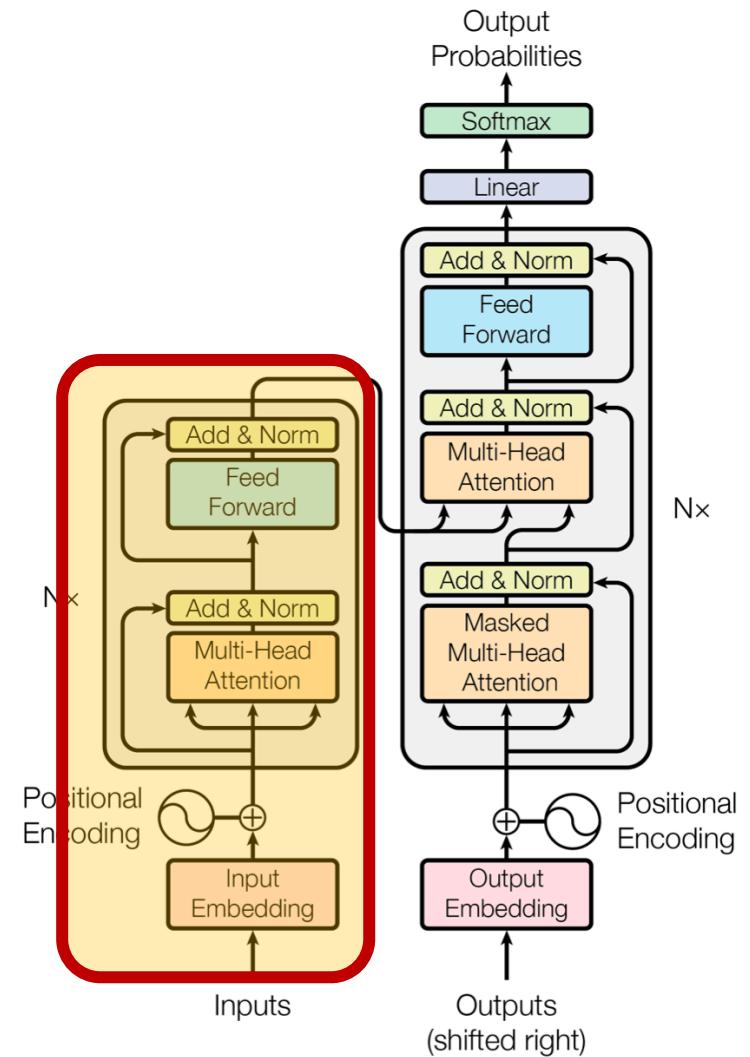


GPT-I
Jun 2018

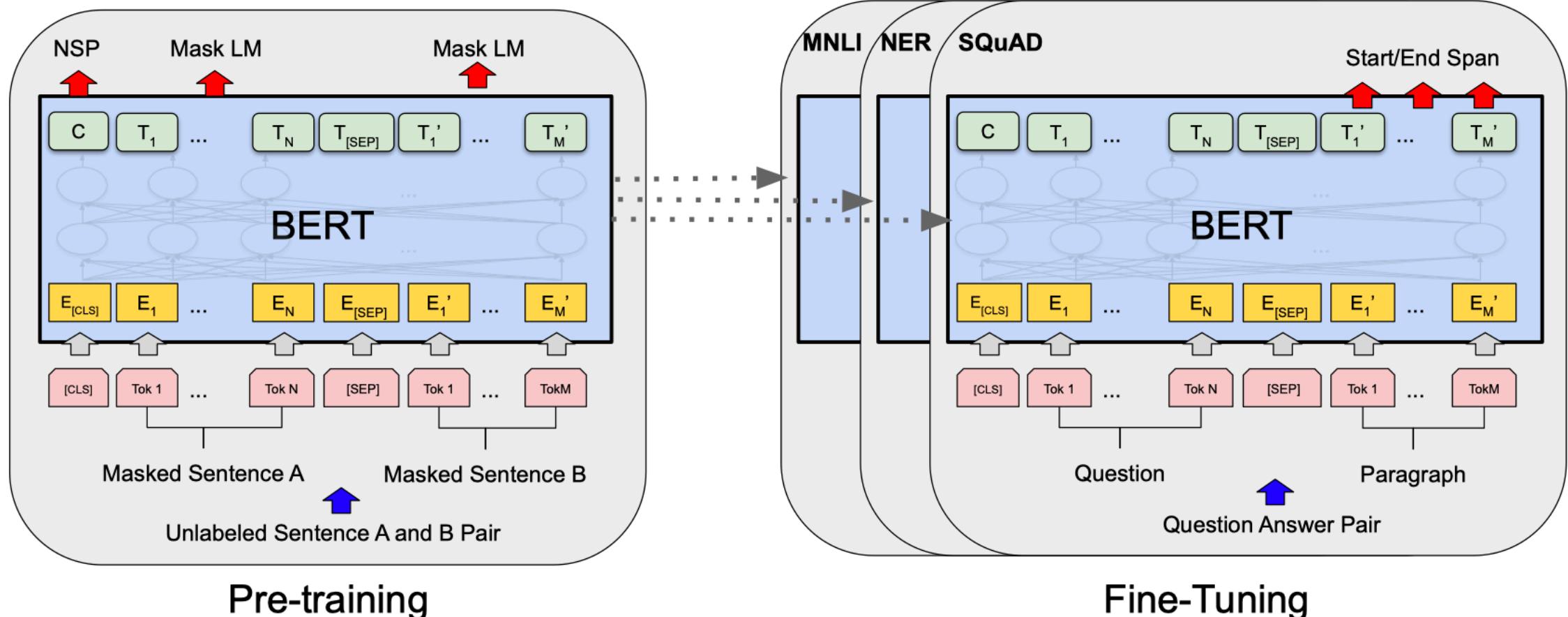
Generation

BERT

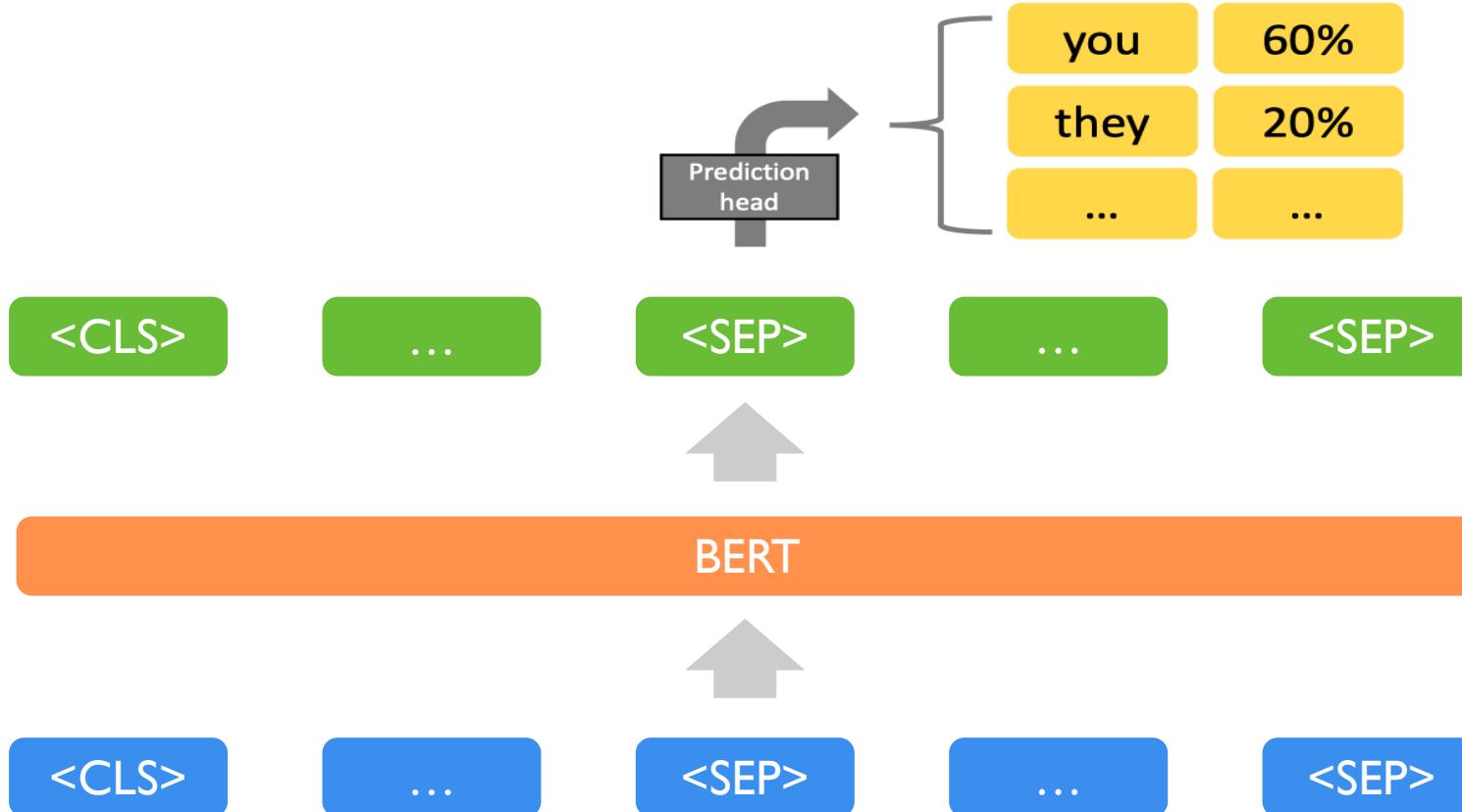
- BERT (Bidirectional Encoder Representation from Transformers, Devlin et al., 2018) is an encoder of a transformer pre-trained with:
 - Masked Language Model (MLM), that consists in predicting [15% of] words which have been replaced with a “MASK” token.
 - Next Sentence Prediction (NSP), which consists in predicting if a certain sentence follows the current one.
- It is then fine-tuned on multiple NLP tasks.



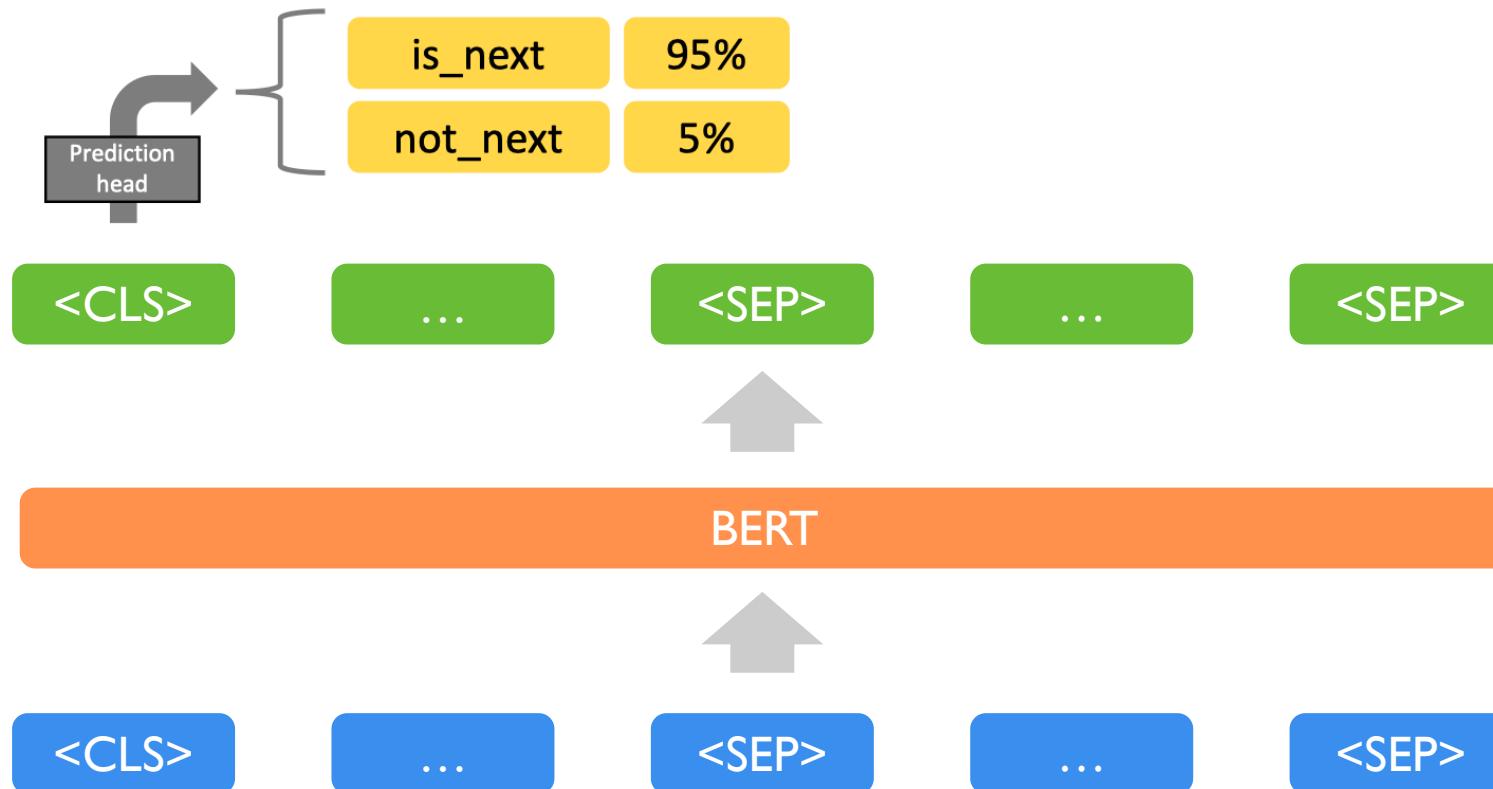
BERT : from pre-training to Fine-Tuning



MLM, Masked Language Modeling

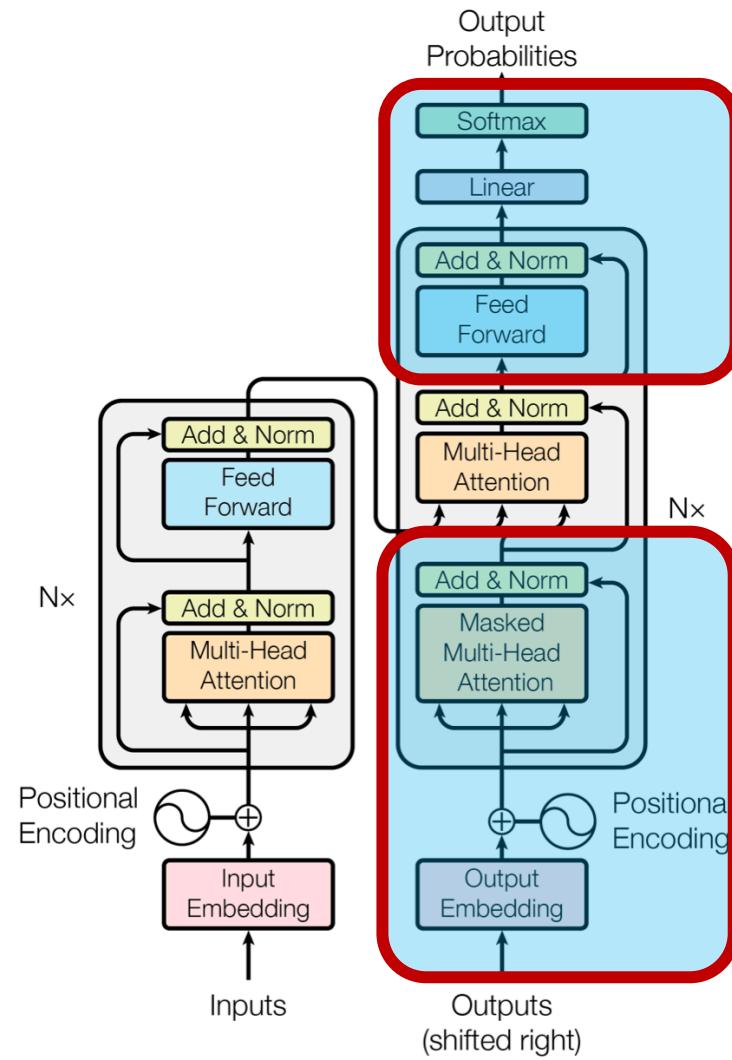


MLM, Masked Language Modeling



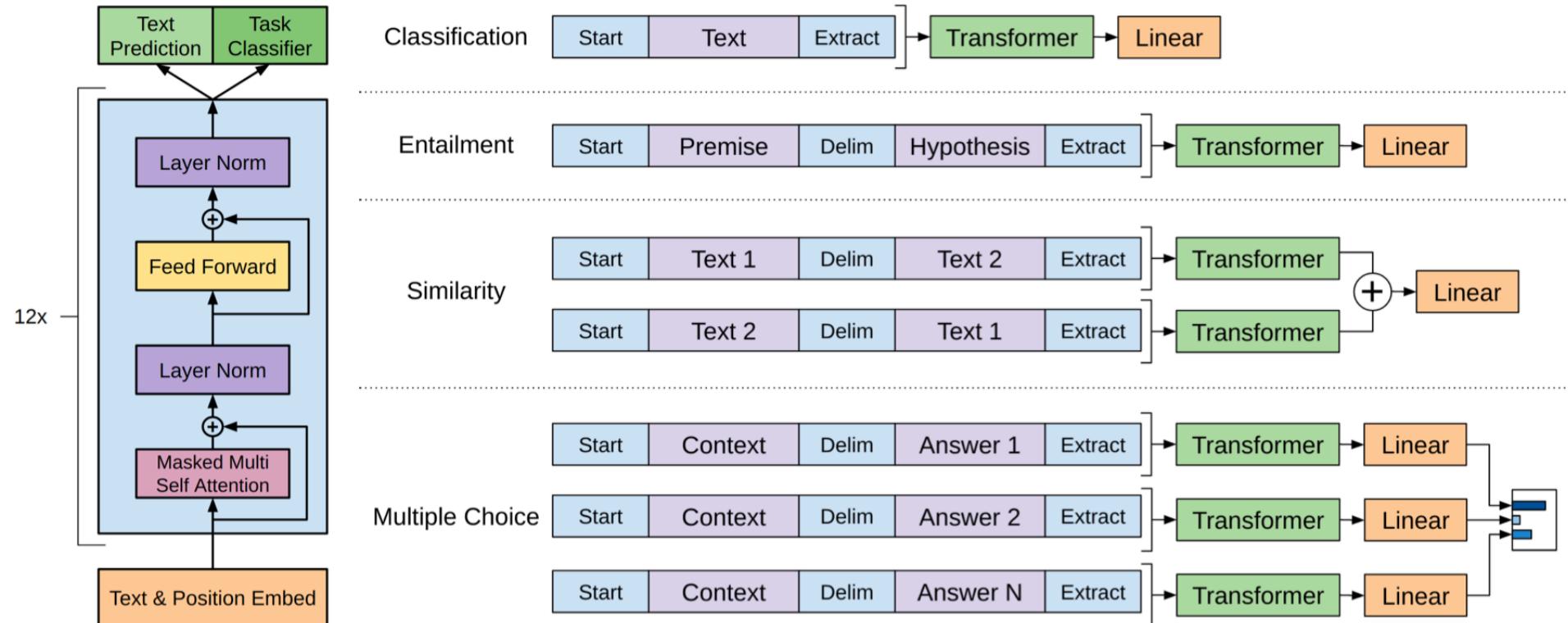
GPT – Generative Pretrained Transformer

- 一般以字 (Tokens) 为单位训练 Transformer 模型。
首先初始化字编码大小为 [vocab size, embedding dimension] , vocab size 为字库中所有字的数量。
- 对于输入的句子 X , 通过 Word Embedding 得到该句子中每个字的字向量 , 通过 Positional Encoding 得到所有字位置向量 , 将其相加得到该字真正的向量表示。第 t 个字的向量记作 x_t 。

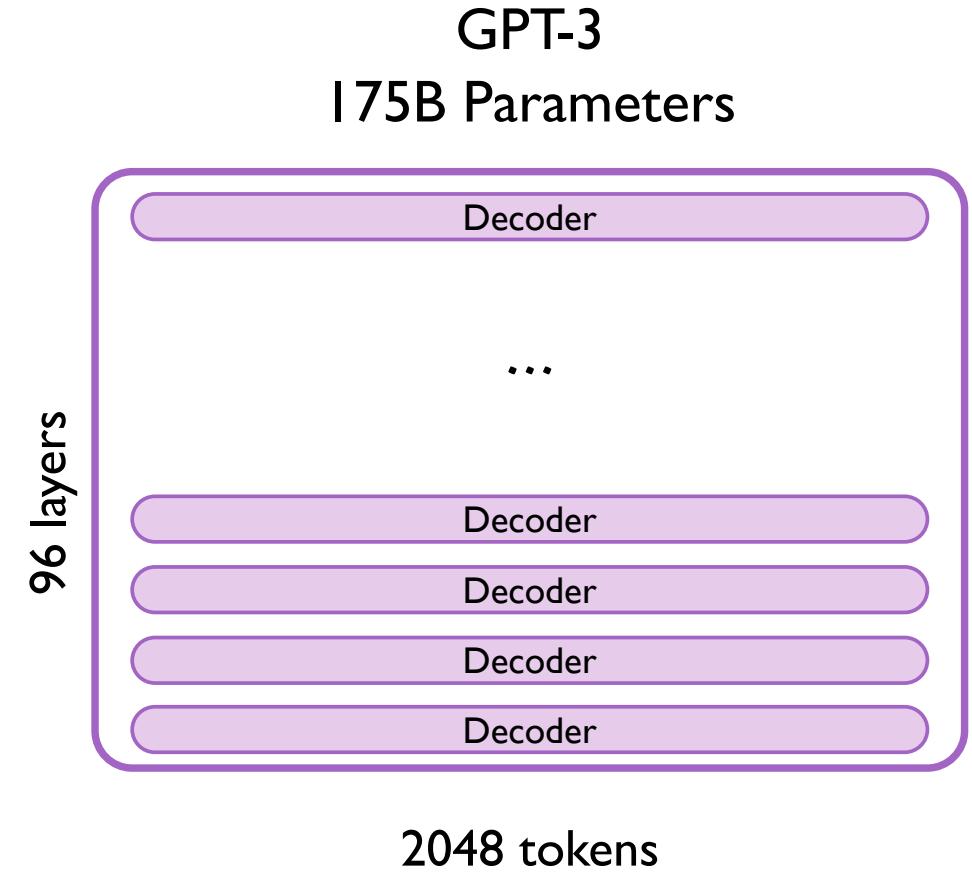
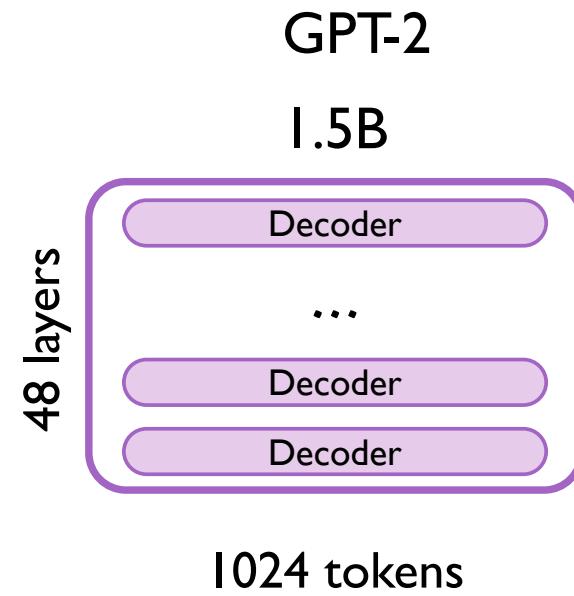
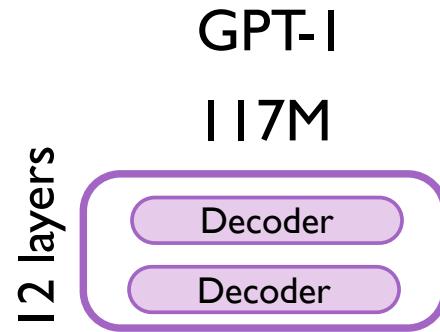


GPT – Generative Pretrained Transformer

- GPT (Generative Pre-Training, Radford, 2018) is a decoder of a transformer trained for auto-regressive text generation.



GPT 系列



GPT – Generative Pretrained Transformer

The three settings we explore for in-context learning

Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.



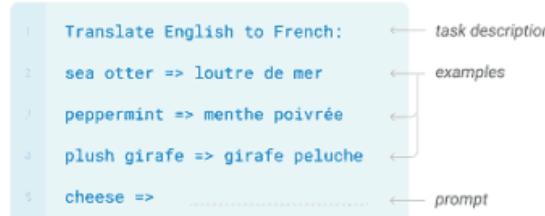
One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.



Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.



Traditional fine-tuning (not used for GPT-3)

Fine-tuning

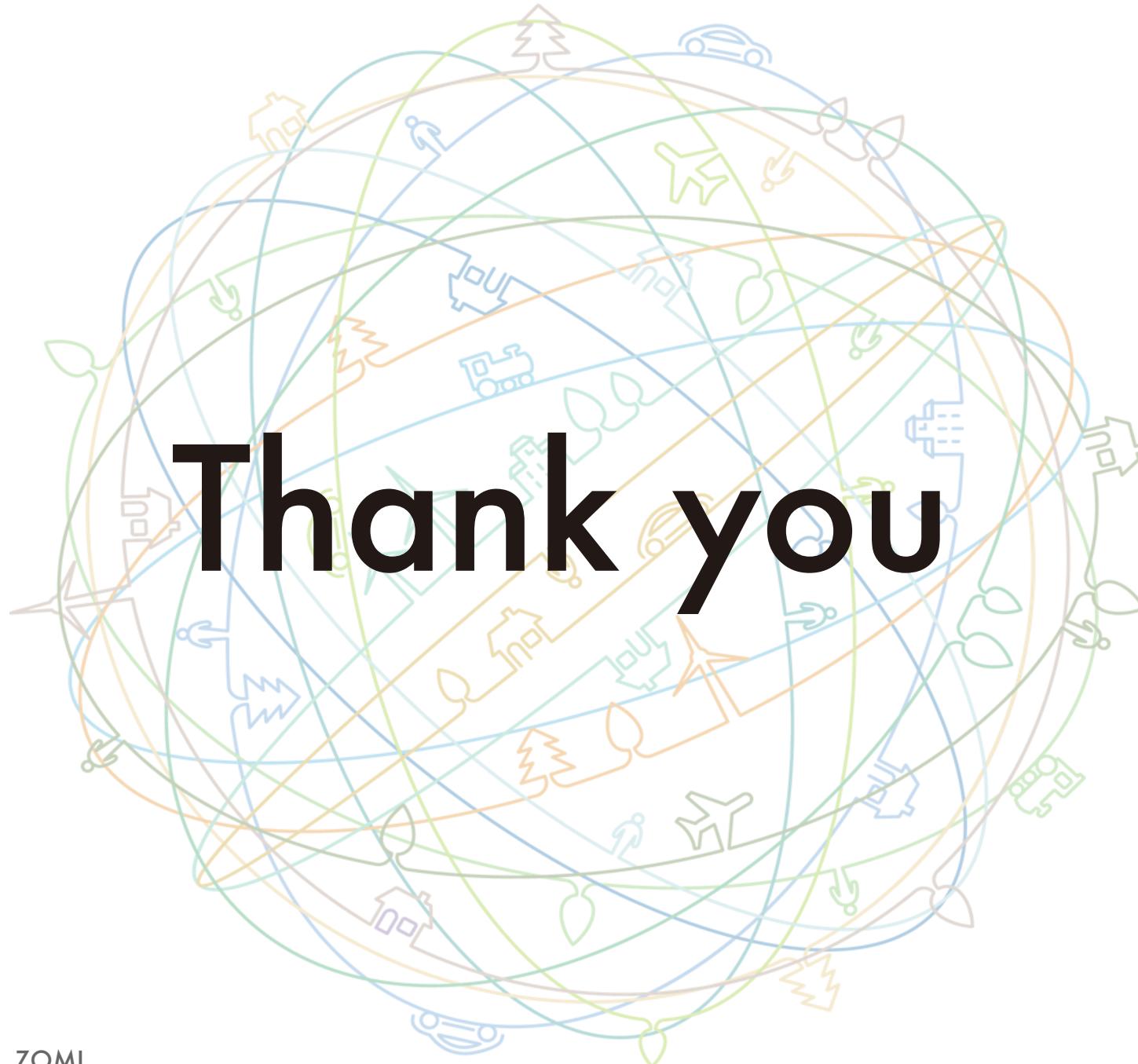
The model is trained via repeated gradient updates using a large corpus of example tasks.



小结

1. GPT是单向模型，不需要利用上下文信息，只能利用上文；而 BERT 是双向模型，训练时候需要上下文信息，更为复杂。
2. GPT是基于自回归模型，能够进行文本生成，可以通过 Prompt 应用在众多 NLP 任务中；BERT 采用基于自编码模型，先预训练再微调的两个阶段来学习通用文本表示，只能完成 NLP 特定任务，无法直接应用在文本生成。





把AI系统带入每个开发者、每个家庭、
每个组织，构建万物互联的智能世界

Bring AI System to every person, home and
organization for a fully connected,
intelligent world.

Copyright © 2023 XXX Technologies Co., Ltd.
All Rights Reserved.

The information in this document may contain predictive statements including, without limitation, statements regarding the future financial and operating results, future product portfolio, new technology, etc. There are a number of factors that could cause actual results and developments to differ materially from those expressed or implied in the predictive statements. Therefore, such information is provided for reference purpose only and constitutes neither an offer nor an acceptance. XXX may change the information at any time without notice.



Course chenzomi12.github.io

GitHub github.com/chenzomi12/DeepLearningSystem