第09章 注意力机制

欧新宇



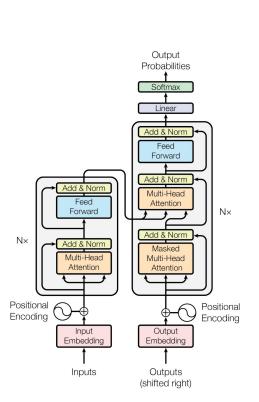




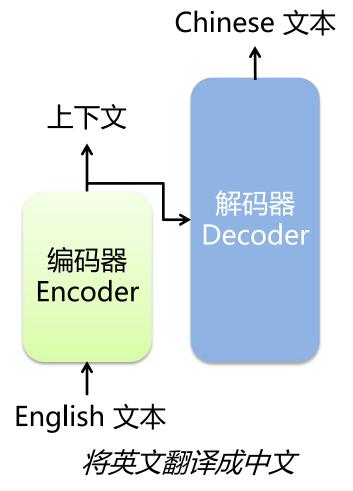
- Transformer 的架构概述
- Transformer 的模块
- Transformer 的训练和预测



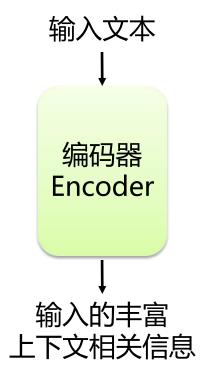




原始 Seq2Seq



编码器模型



Basis:

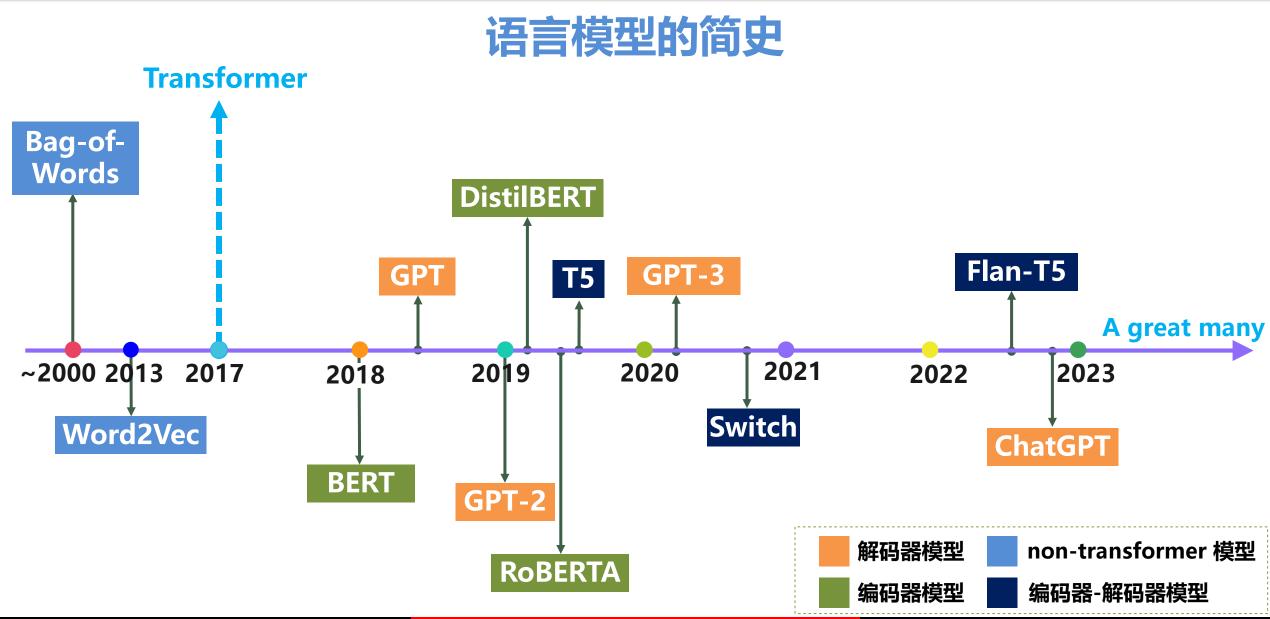
- BERT
- 大多数 embedding 模型

解码器模型

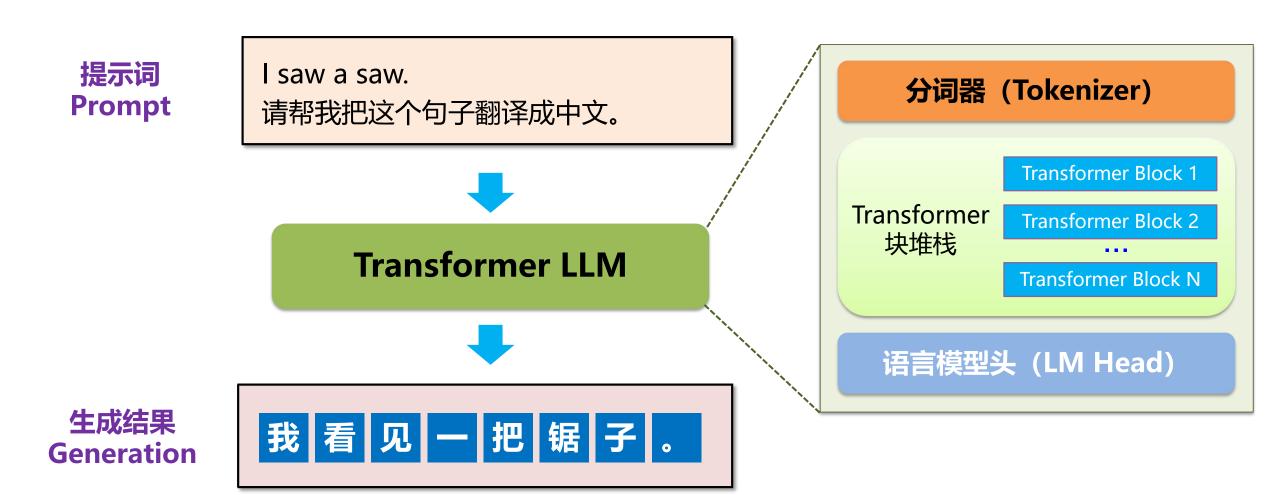


Basis:

- 主流LLMs模型
- GPT、Claude、文心等

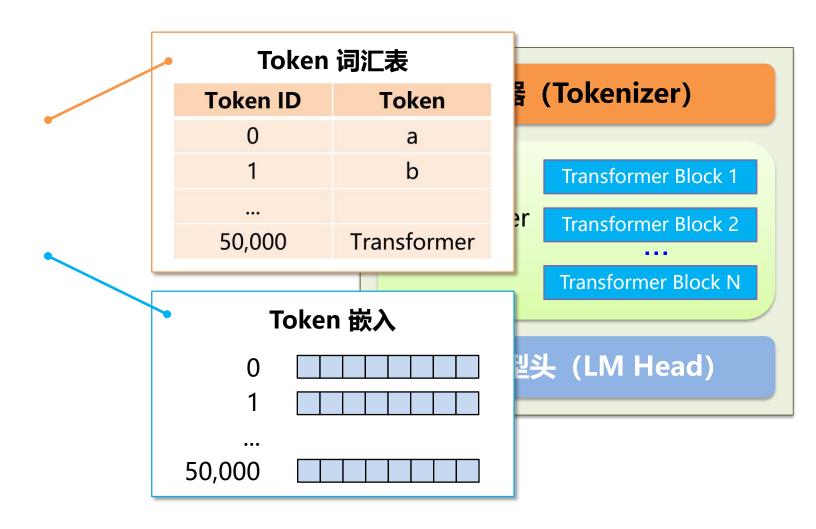


Transformer 架构的主要组件



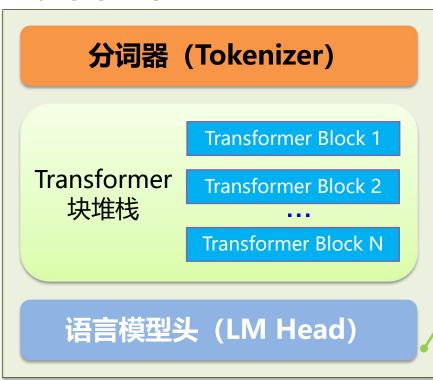
Transformer 架构的主要组件

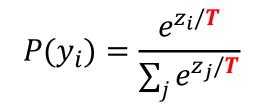
Transformer LLM

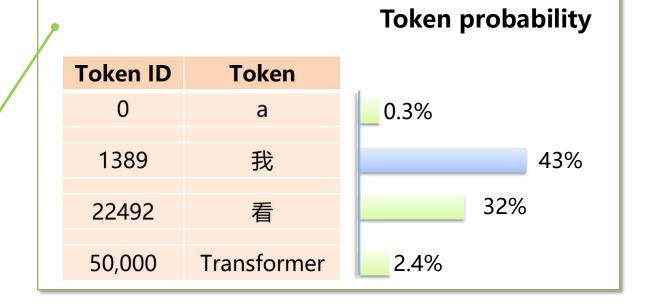


Transformer 架构的主要组件

Transformer LLM







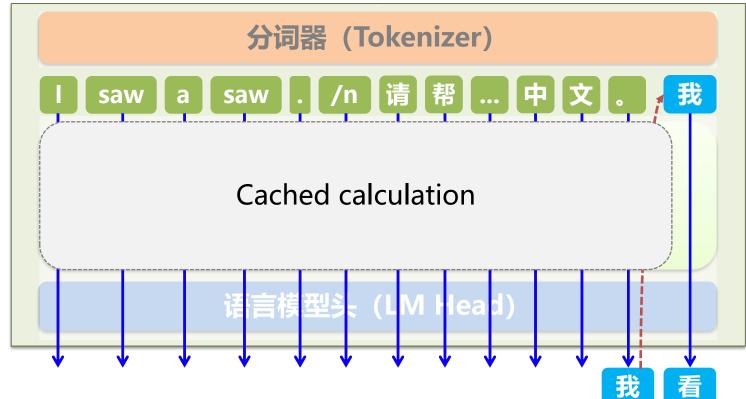
Transformer 架构的主要组件

I saw a saw.

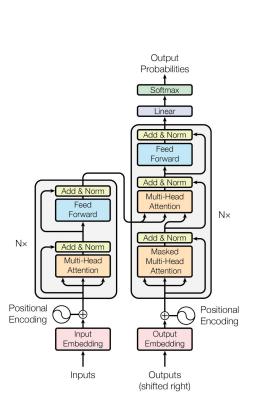
请帮我把这个句子翻译成中文。



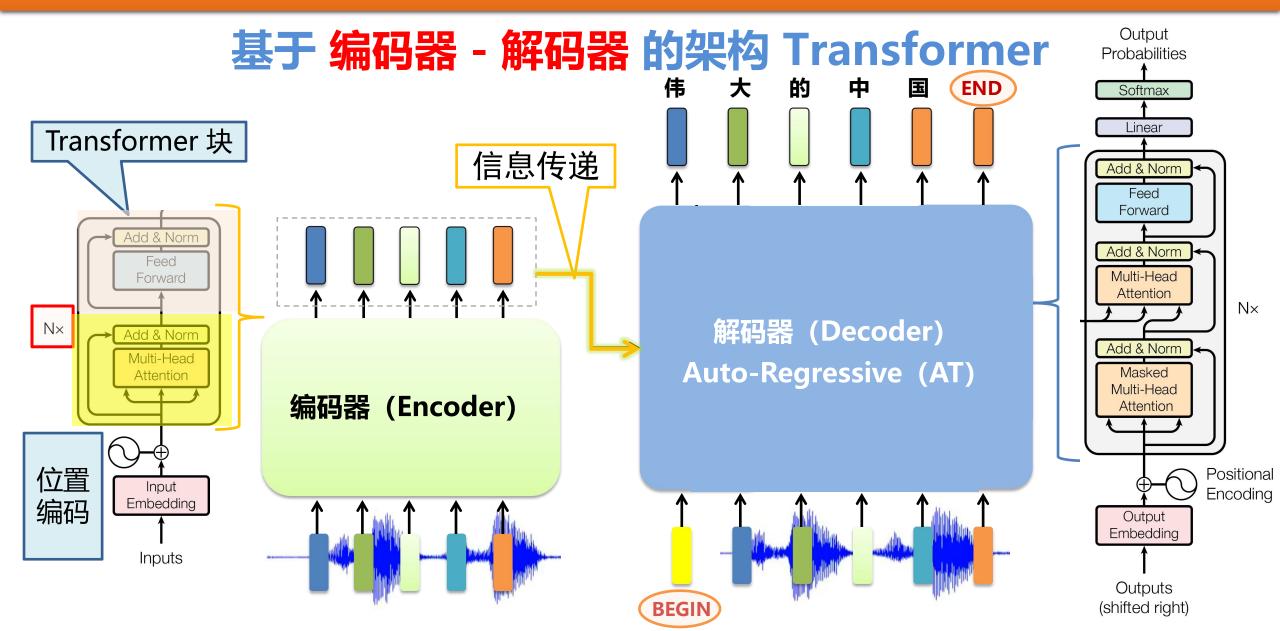
Transformer LLM



10/27



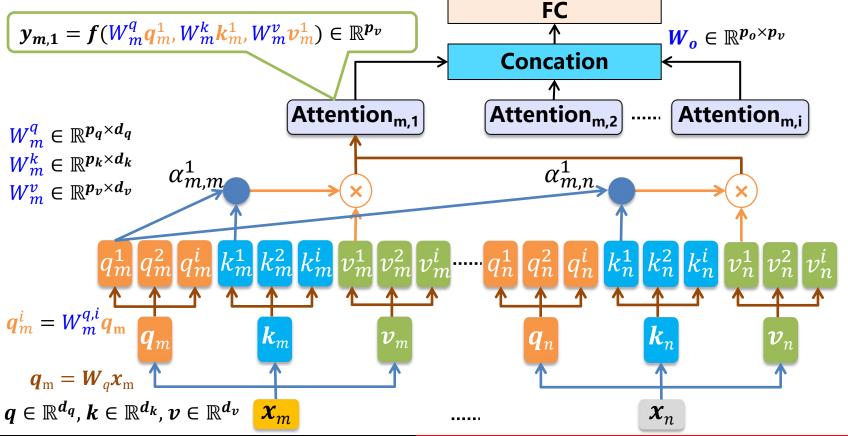
Ahish Vaswani, Noam Shazeer, Niki Parmar, et.al. Attention is all you need. NIPS2017.

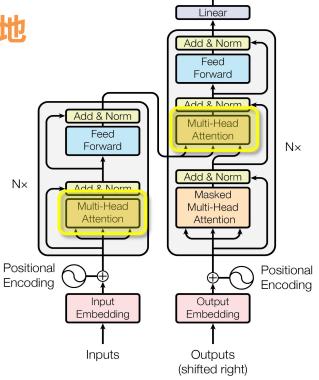


多头注意力

● 多头注意力:对于同一组 query, key, value,通过多组线性变换并行地

捕捉不同的特征。短距离关系(句子)和长距离关系(段落)。





Output

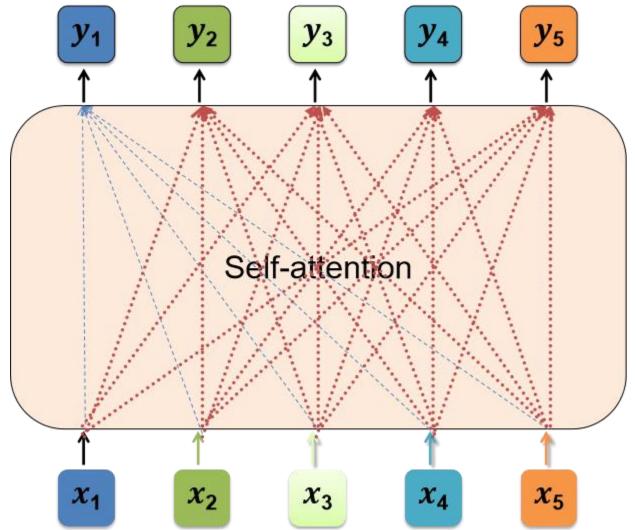
Probabilities

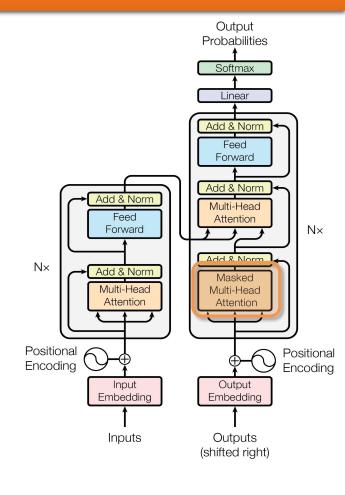
Softmax

● 多头注意力的输出

$$egin{aligned} oldsymbol{y_m} &= oldsymbol{W_o} egin{bmatrix} oldsymbol{y_{m,1}} \ oldsymbol{y_{m,i}} \ oldsymbol{y_{m,i}} \end{bmatrix} \in \mathbb{R}^{p_o} \end{aligned}$$

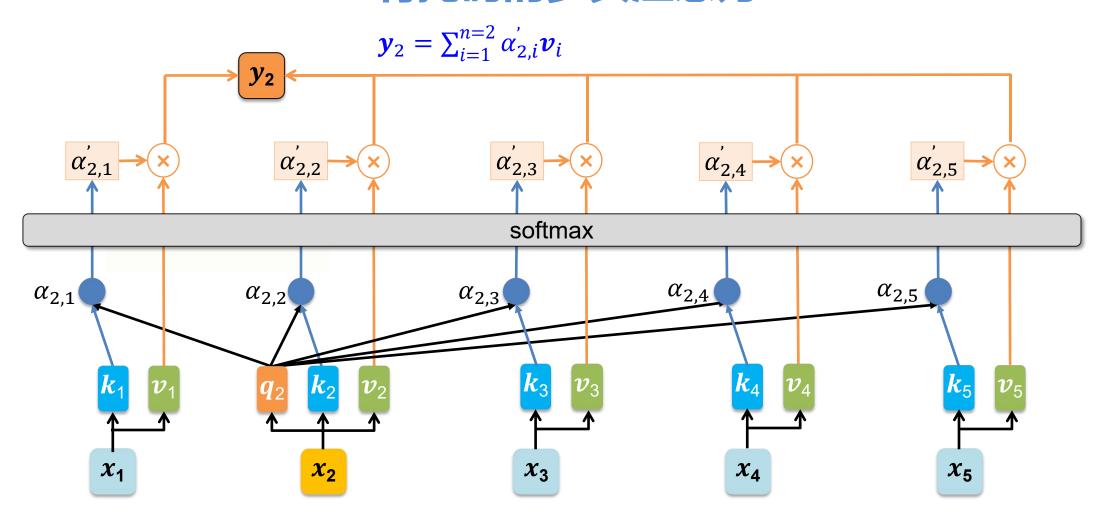
有掩码的多头注意力





欧新宇 | ouxinyu@alumni.hust.edu.cn

有掩码的多头注意力



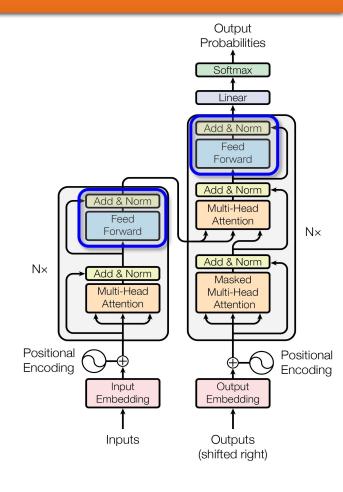
基于位置的前馈网络

● 卷积神经网络

- ✓ 输入 (卷积层) 形状: (b, h, w, c)
- ✓ **输出 (全连接层) 形状:** (b, h×w×c)

• Transformer

- ✓ 编码器: 将输入形状由 (b, n, d) 变换为 (bn, d)
- ✓ 解码器: 将输出形状由 (bn, d) 变回 (b, n, d)
- ✓ 等价于两个核窗口为1的一维卷积



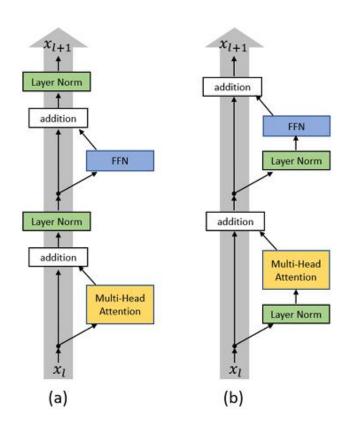
Self-attention

Output 层归一化 Probabilities Softmax Linear Add & Norm Feed Forward Add & Norm norm norm Multi-Head Forward $N \times$ Add & Norm N× Layer Norm Add & Norm x + yMulti-Head Multi-Head Attention len Positional Positional Encoding Encoding x_2 均值: μ 标准差: σ Output \boldsymbol{x} Embedding Embedding Inputs Outputs (shifted right) FC

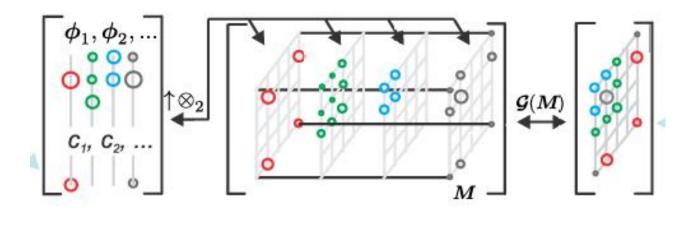
Transformer

shortcut

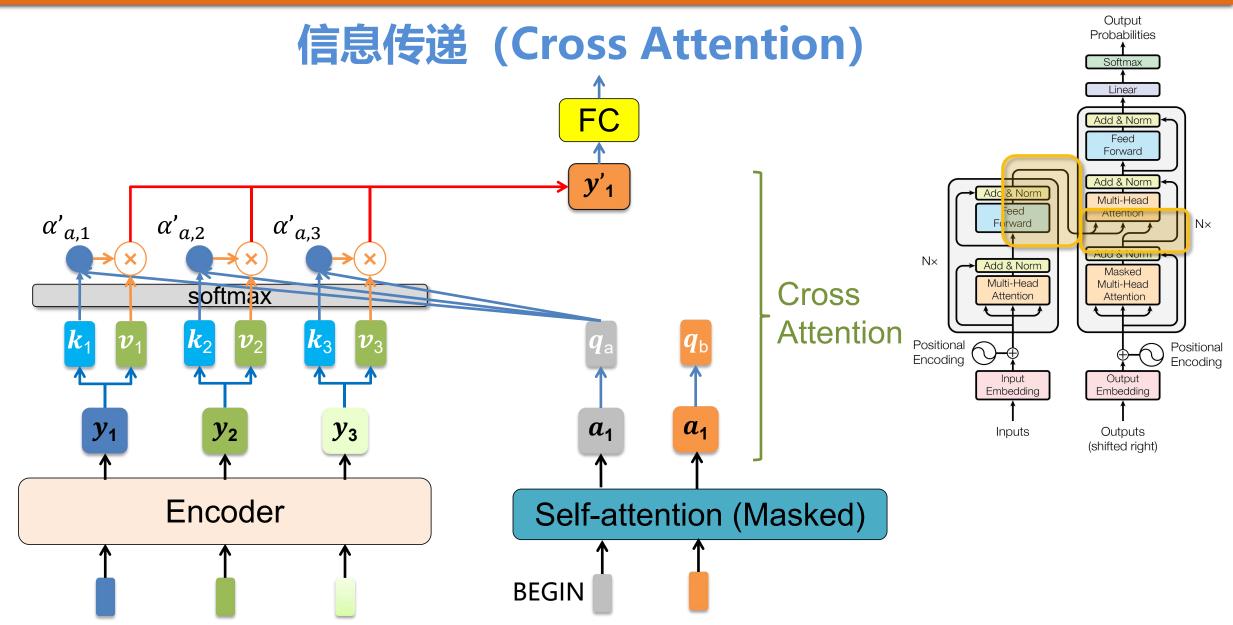
Tips1: Add & Norm



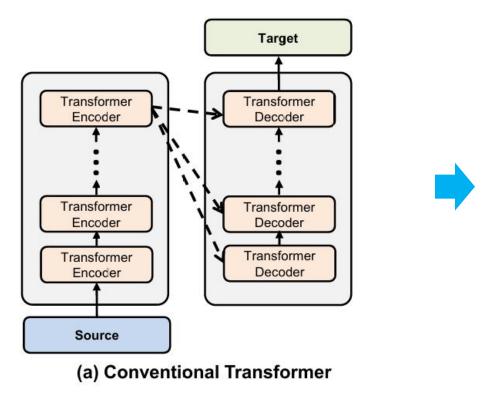
 On Layer Normalization in the Transformer Architecture



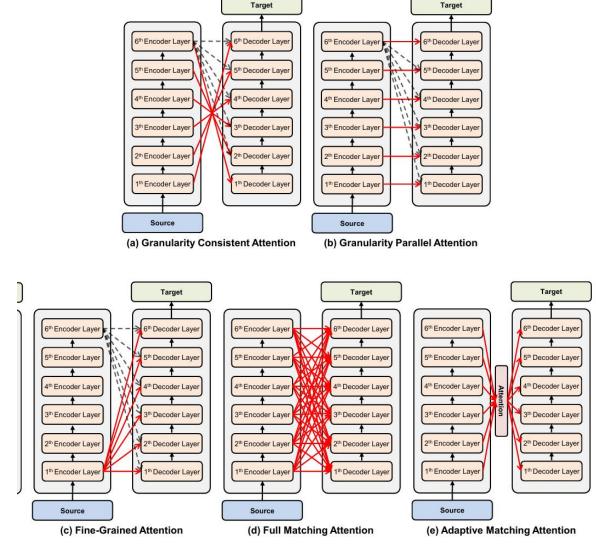
- PowerNorm: Rethinking Batch normalization in Transformers
- A Deeper Look at Power Normalizations



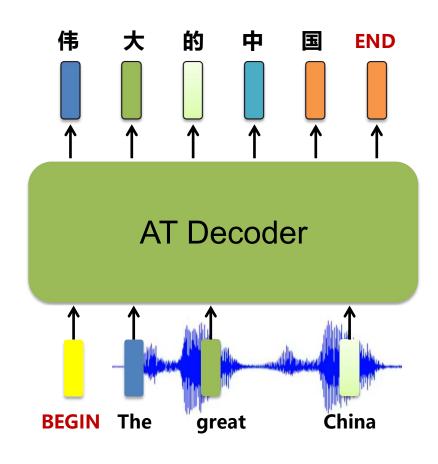
Tips2: 交叉注意力

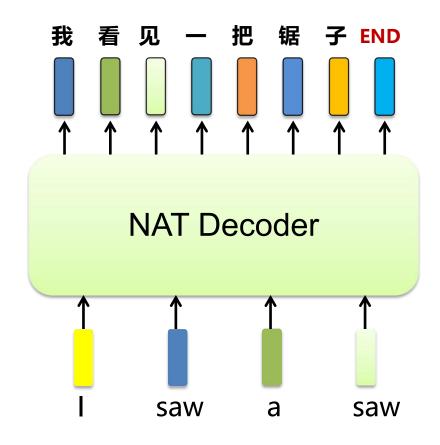


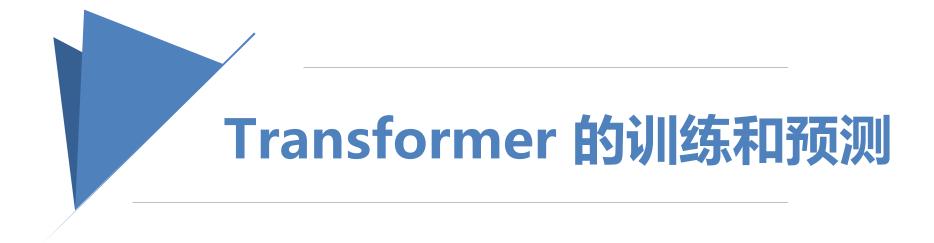
■ Rethinking and Improving Natural Language Generation with Layer-Wise Multi-View Decoding



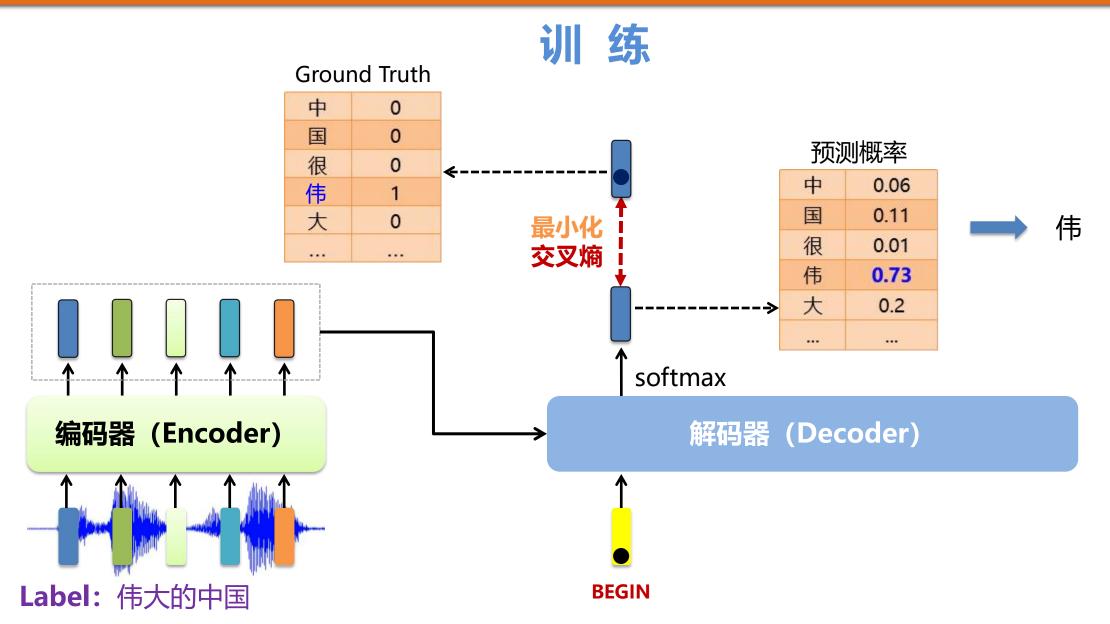
Tips3: 自回归解码器和非自回归解码器





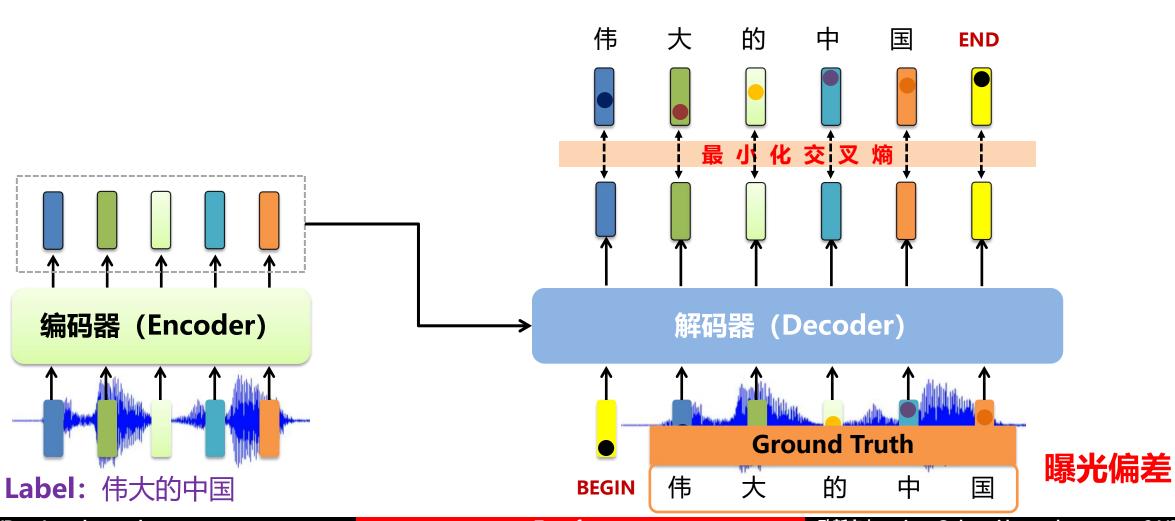


Transformer 的训练和预测



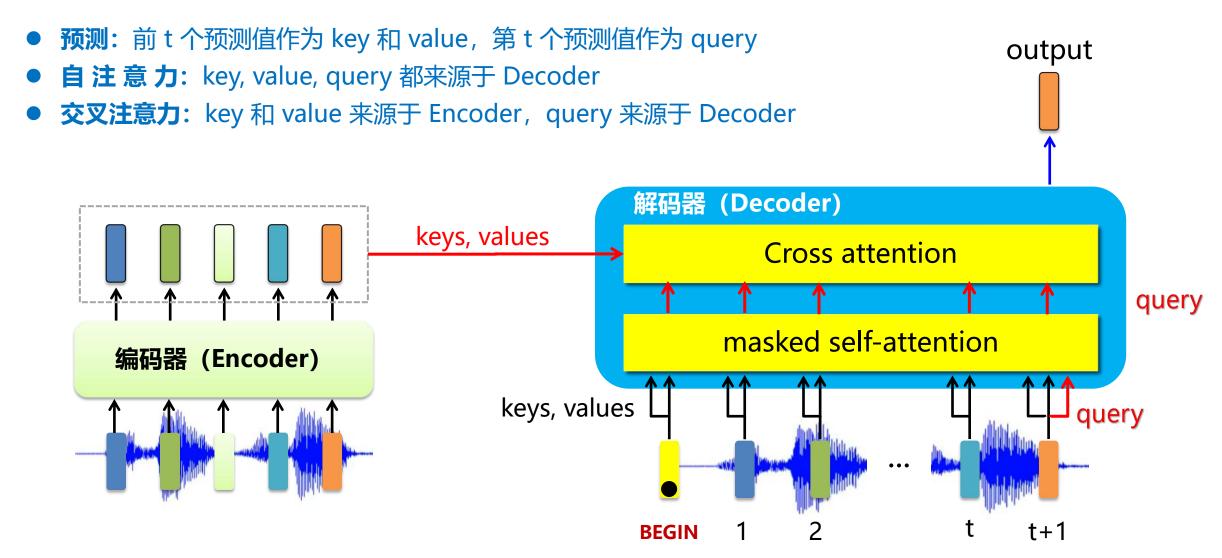
Transformer 的训练和预测

训练



Transformer 的训练和预测

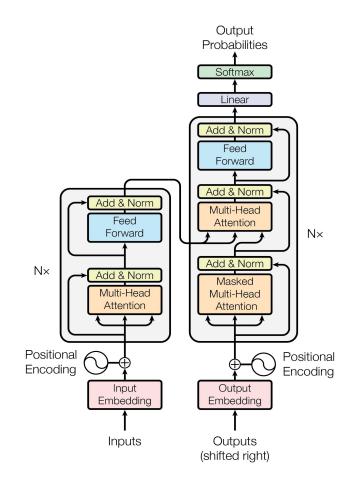
预 测



Transformer

小 结

- Transformer 是一个纯注意力的 编码器-解码器 模型
- 编码器 和 解码器 都由多个 transformer 块构成,主要包括 多头自注意力、基于位置的前馈网络、层归一化组成
- 编码器使用 self-attention 捕捉全局依赖;解码器采用 masked self-attention 来屏蔽未来信息,并实现历史序 列信息的获取
- 编码器和解码器之间使用 cross-attention 进行连接
- Positional Encoding 通过硬植入的方式将相对位置信息 融合到每一个 token 嵌入向量中。



欧老师的联系方式

