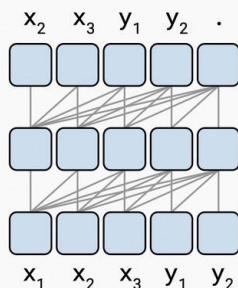


↑  
等价于 Transformer

## Language model

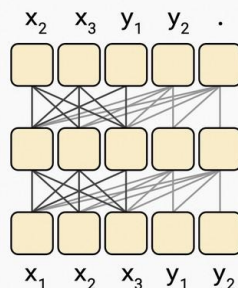


causal decoder  
(Next token prediction)

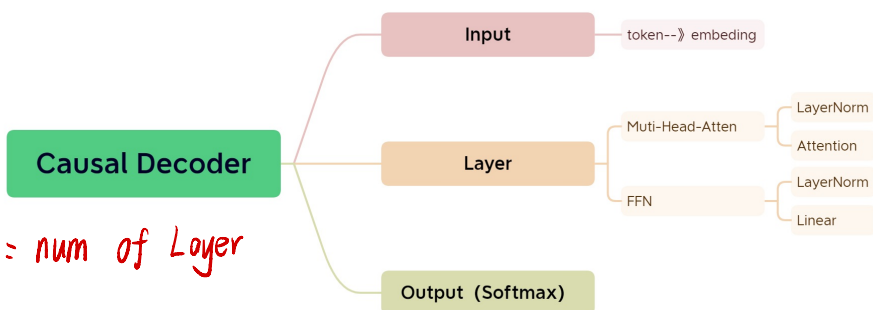
最常见的架构

LLam系列

## Prefix LM



chat GPT  
系列



$L$ : num of Layer

$H$ : hidden dim

$V$ : vocab size

Presented with xmind

对于 Attention:  $\text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V$

代码中有三个投影 `nn.Linear(hidden_dim, hidden_dim)`

单个线性层参数为:  $h^2 + h$

三个投影 + 输出 Layer  $\Rightarrow 4h^2 + 4h$

对于FFN中Linear  
维度先升后降  $\Rightarrow [b, s, h] \rightarrow [b, s, 4h] \rightarrow [b, s, h]$

权重  $\Rightarrow [h, 4h] \quad [4h, h]$

$$\Rightarrow \text{计算量} \Rightarrow 4h^2 + 4h + 4h^2 + h \Rightarrow 8h^2 + 5h$$

对于LayerNorm  $\Rightarrow 2h$

对于input 与 output 一般有 tie weight: 正常为  $vh$

$$\text{总参数为: } (4h^2 + 4h + 8h^2 + 5h + 2h \times 2) \cdot L + vh$$

$$\Rightarrow 12Lh^2 + 13Lh + vh$$

一般  $h$  比较大, 近似为  $12Lh^2$