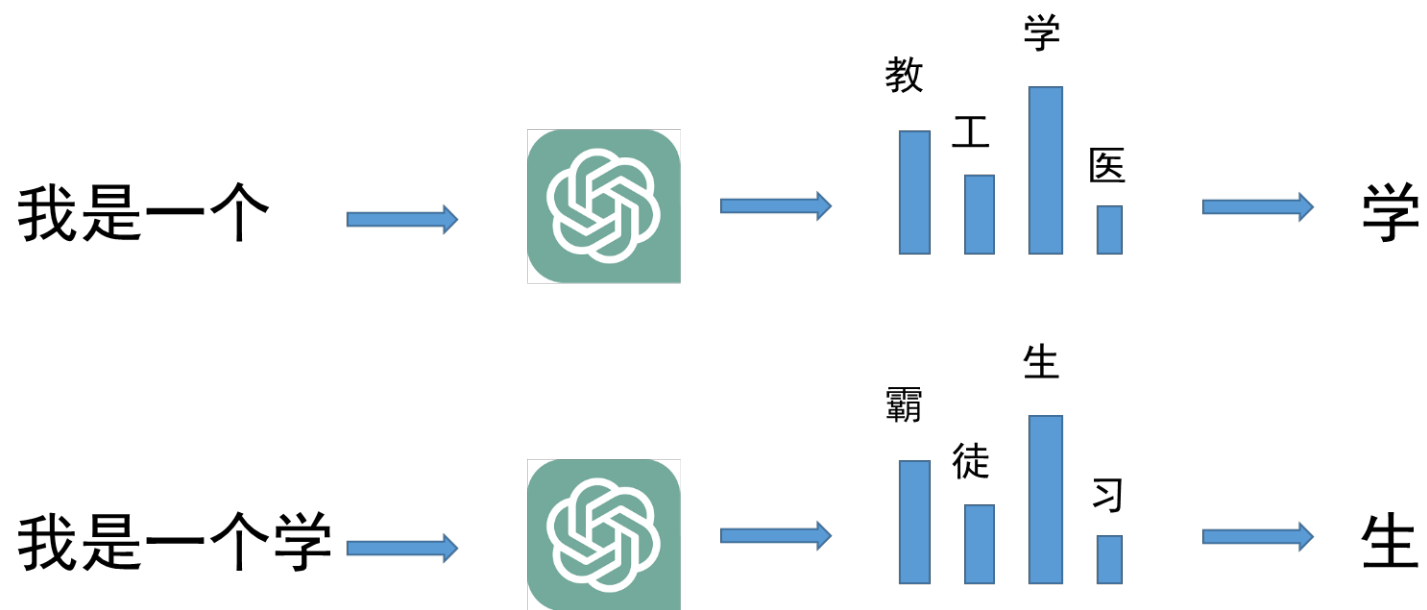


# Transformer——大模型的基石

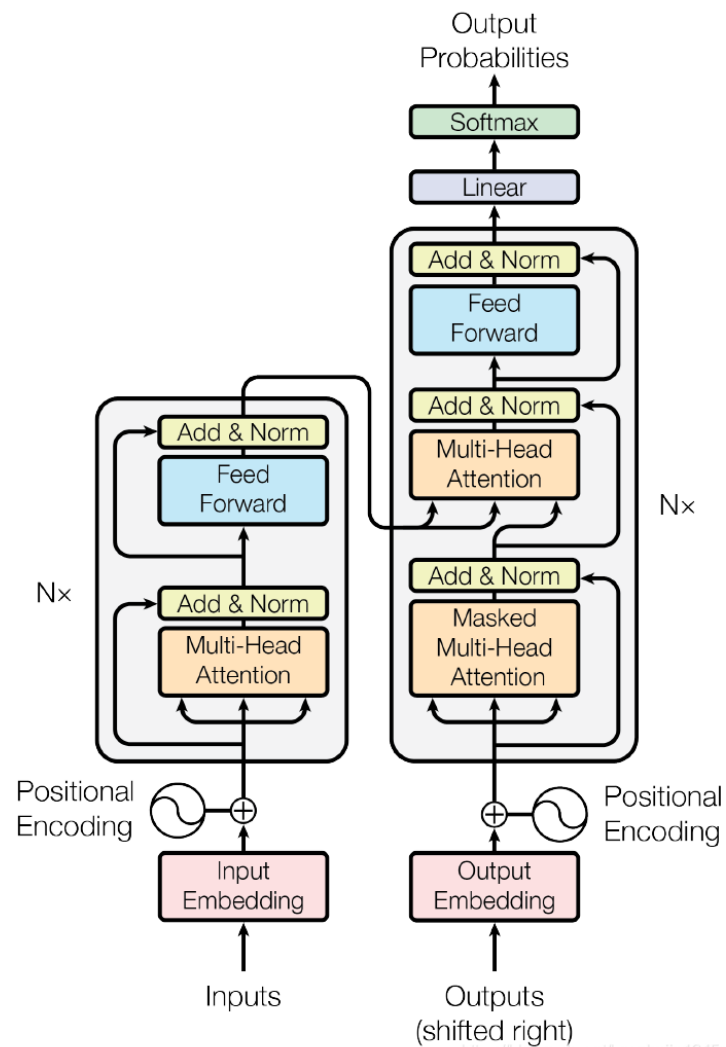
# GPT: 生成式预训练转换模型

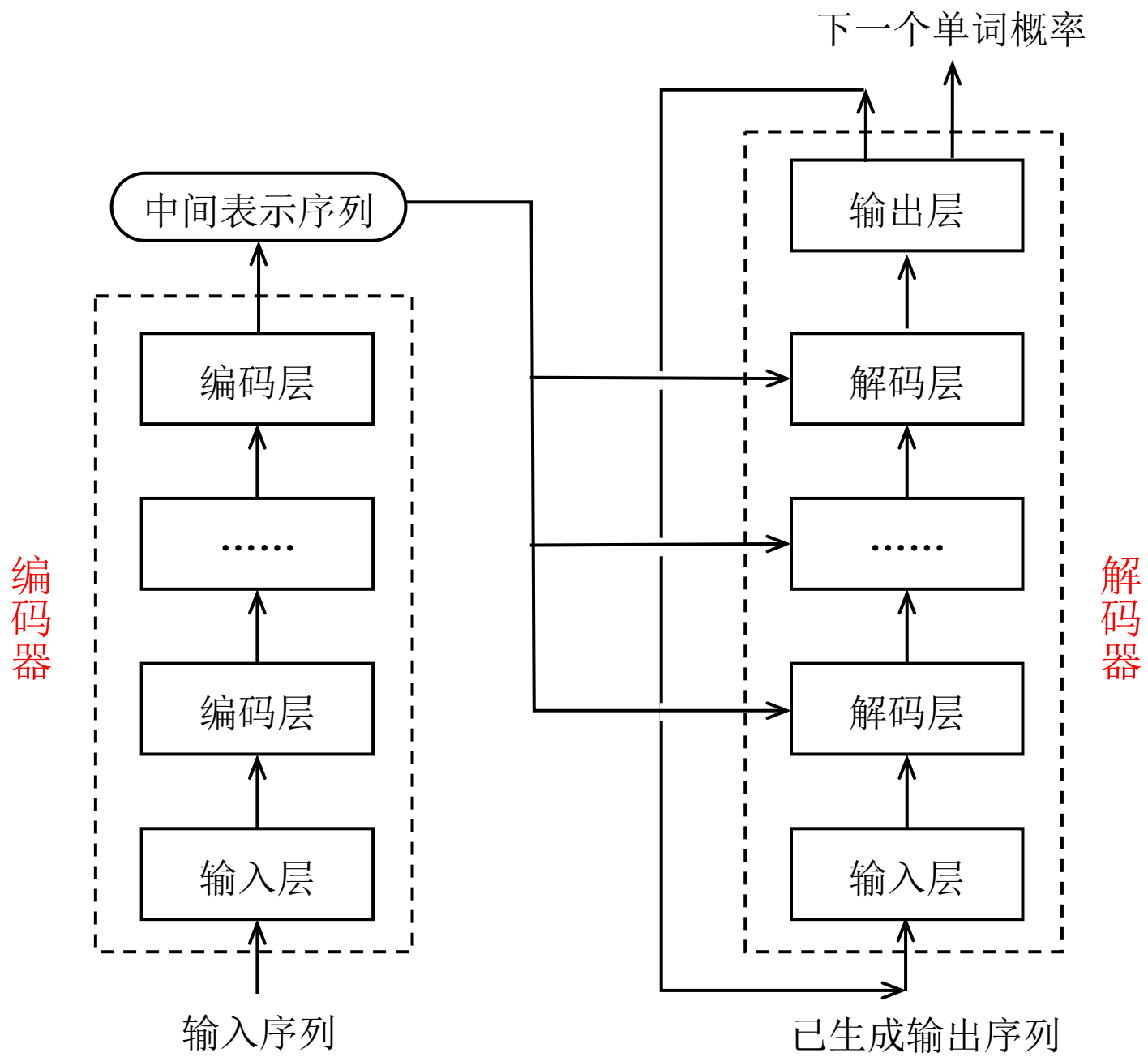
---

## ► 文字接龙

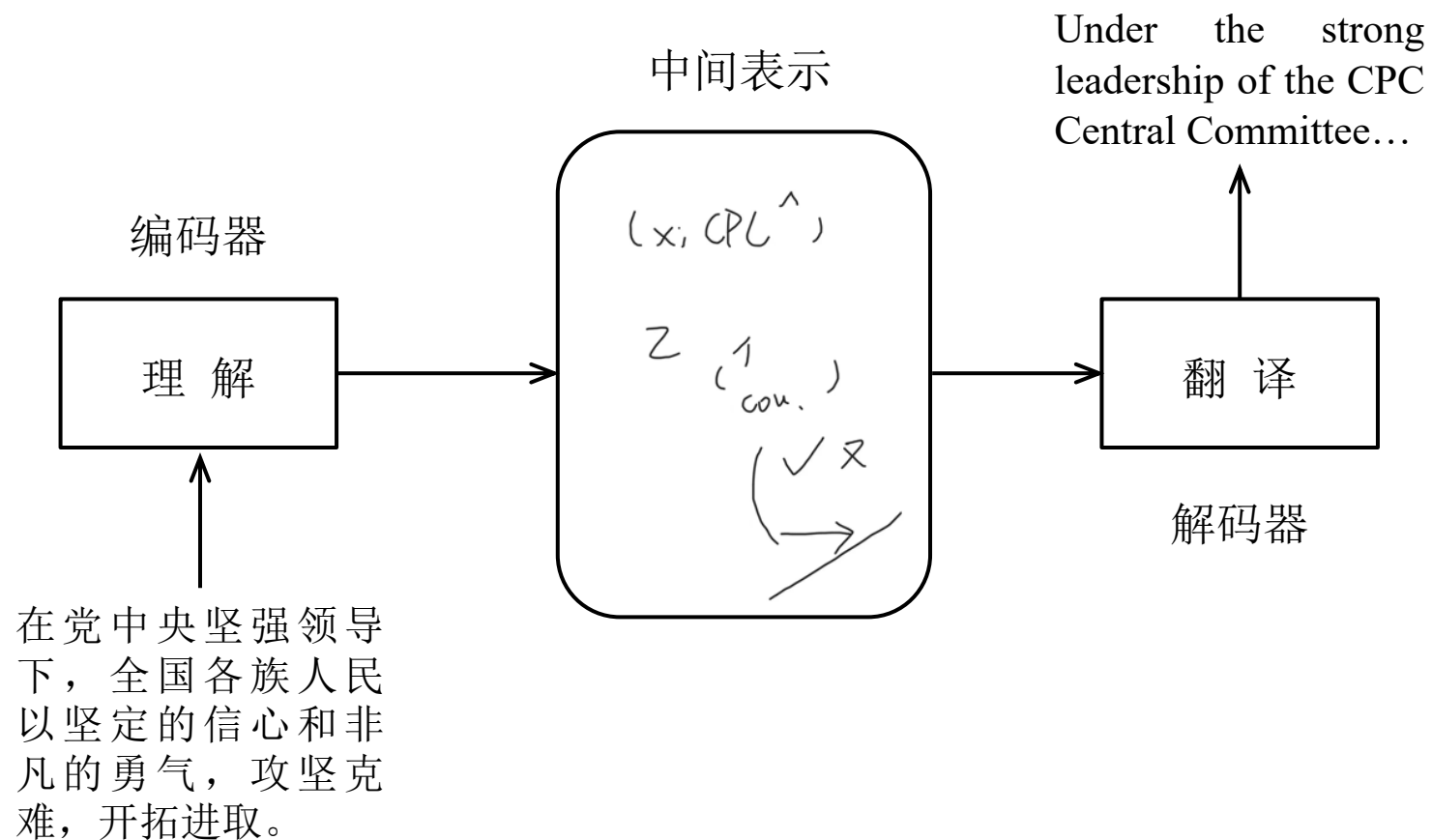


# Transformer: 转换模型

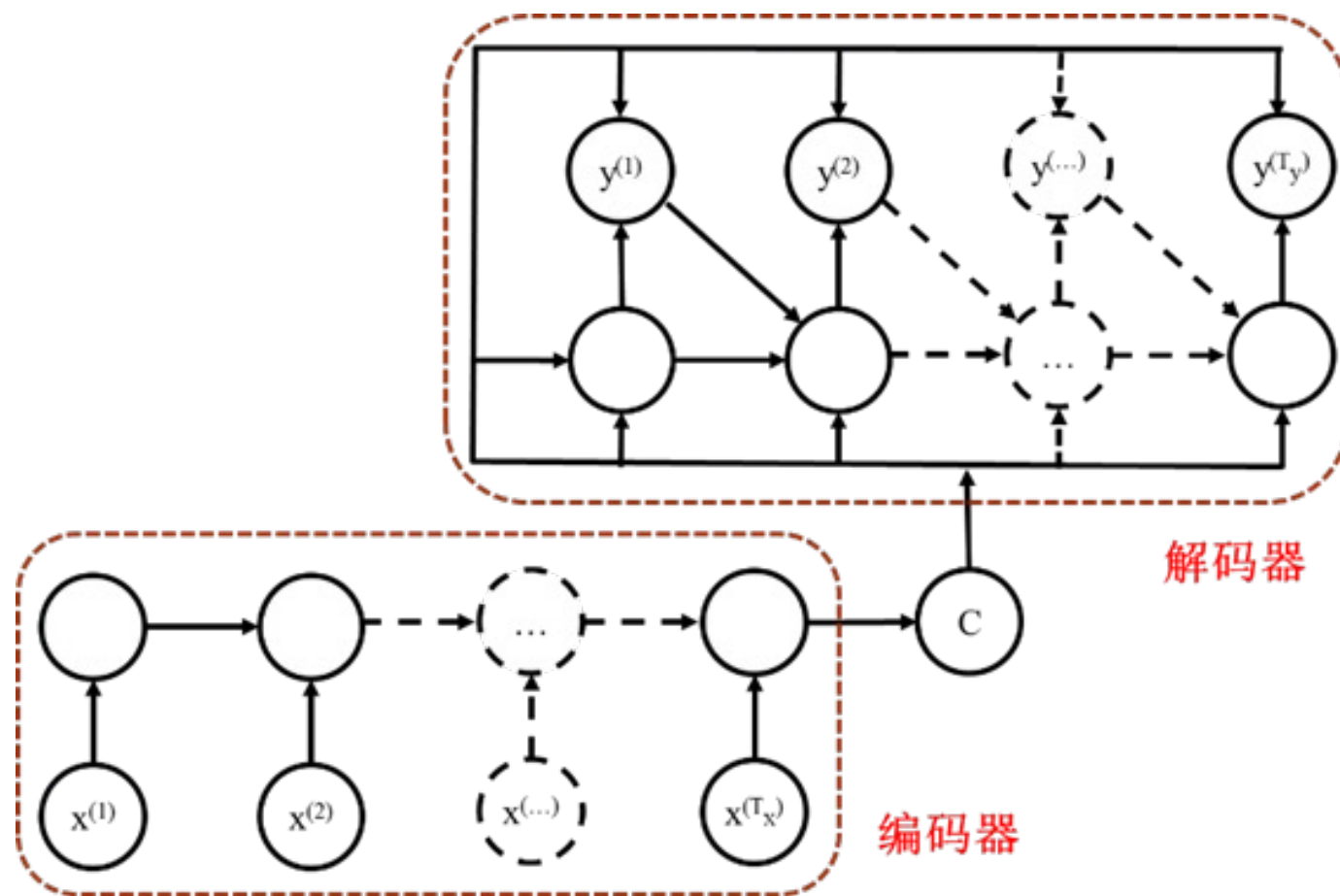




# 编码器-解码器模型



# 编码器-解码器模型



# 关键技术

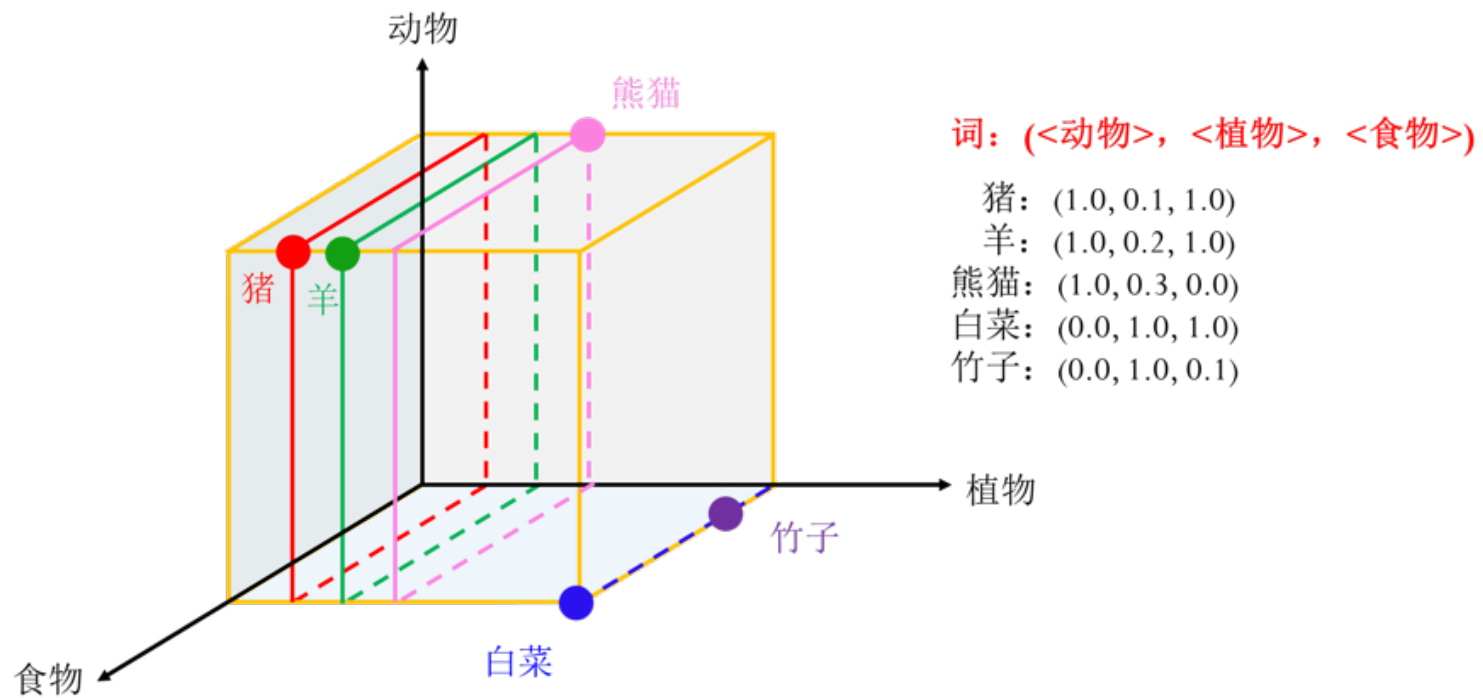
---

- ▶ 两个关键技术
  - ▶ 词向量的动态表示
  - ▶ 注意力机制



# 词向量：静态与动态

- ▶ 静态性
- ▶ 词的多义性





一个词可以用其上下文表示

---

“我吃了一个非常美味的苹果”



“我用苹果跟朋友联系”



# 注意力机制

---



# 键-值数据库查询

---

- ▶ 键-值对数据

- ▶  $\{(k_1, v_1), (k_2, v_2), \dots, (k_n, v_n)\}$

- ▶ 查询  $q$

- ▶ 例：

- ▶ (年龄  $i$ -收入  $i$ )

- ▶ 查询：中年人的收入

- ▶ 按照“年龄”与“中年人”的相似程度计算加权平均值



# 注意力机制

---

$$v = \sum_{i=1}^n \alpha(q, k_i) v_i$$

$$\sum_{i=1}^n \alpha(q, k_i) = 1$$

$$\alpha(q, k_i) = \text{softmax}(\text{sim}(q, k_i)) = \frac{e^{\text{sim}(q, k_i)}}{\sum_{j=1}^n e^{\text{sim}(q, k_j)}}$$

$\text{sim}(q, k_i)$ : 表示 $q$ 与 $k_i$ 的相似性

$\alpha(q, k_i)$ : 注意力



# 注意力机制

- ▶  $q, k_i, v_i$  均用维度为  $d$  的行向量表示

- ▶ 键矩阵:

$$K = \begin{bmatrix} k_1 \\ k_2 \\ \vdots \\ k_n \end{bmatrix}_{n \cdot d}$$

- 值矩阵:

$$V = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix}_{n \cdot d}$$

- 查询矩阵:

$$Q = \begin{bmatrix} q_1 \\ q_2 \\ \vdots \\ q_m \end{bmatrix}_{m \cdot d}$$

- ▶ 相似性:

$$\underset{\text{标量}}{\text{sim}(q_i, k_j)} = \frac{q_i k_j^T}{\sqrt{d}}$$

- ▶ 注意力:

$$\text{att}(q_i, K, V) = \text{softmax}\left(\frac{q_i K^T}{\sqrt{d}}\right) \cdot V$$

$d$  维行向量

对矩阵的每一行做softmax

$$\text{att}(Q, K, V) = \text{softmax}\left(\frac{Q K^T}{\sqrt{d}}\right) \cdot V$$

$m \cdot d$  矩阵

# 自注意力机制

---

- ▶ 长度为 $n$ 的输入序列 $x$ :  $x_1, x_2, \dots, x_n$
- ▶ 序列矩阵:

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}_{n \cdot d}$$

- ▶ 自注意力机制:

$$att(Q, K, V) = att(X, X, X) = softmax\left(\frac{XX^T}{\sqrt{d}}\right) \cdot X_{n \cdot d}$$

- ▶  $att(Q, K, V)$ 的第 $i$ 行为序列中第 $i$ 个单词的动态向量表示
  - ▶ 列数为向量的维度 $d$
  - ▶ 行数为序列长度 $n$



# 多头注意力机制

---

- ▶ 分别对  $Q$ 、 $K$ 、 $V$  线性变换后再计算注意力

$$U_i = \text{att}(QW_q^{(i)}, KW_k^{(i)}, VW_v^{(i)}), i = 1, 2, \dots, h$$

$$U = [U_1; U_2; \dots; U_h]_{n \cdot (d \cdot h)}$$

- ▶ 多头注意力机制

$$\text{multi\_att}(Q, K, V) = U \cdot W_o$$

$$W_o: (d \cdot h) \cdot d$$

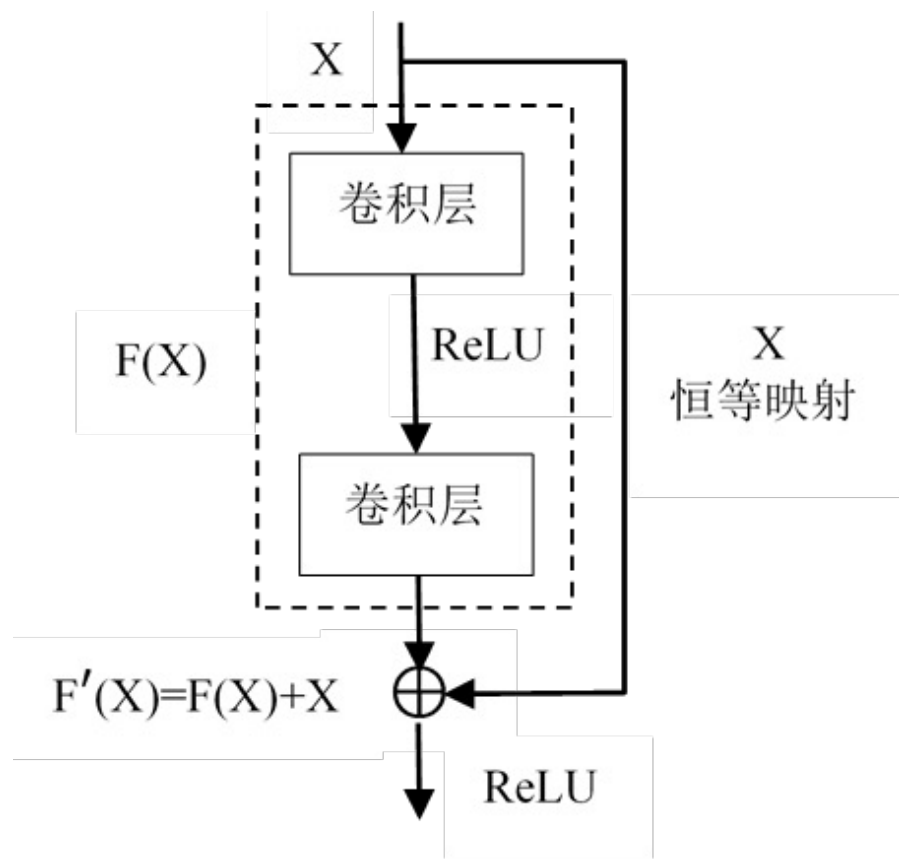
$$\text{multi\_att}(Q, K, V): n \cdot d$$

- ▶  $\text{multi\_att}(Q, K, V)$  的第  $i$  行为序列中第  $i$  个单词的动态向量表示

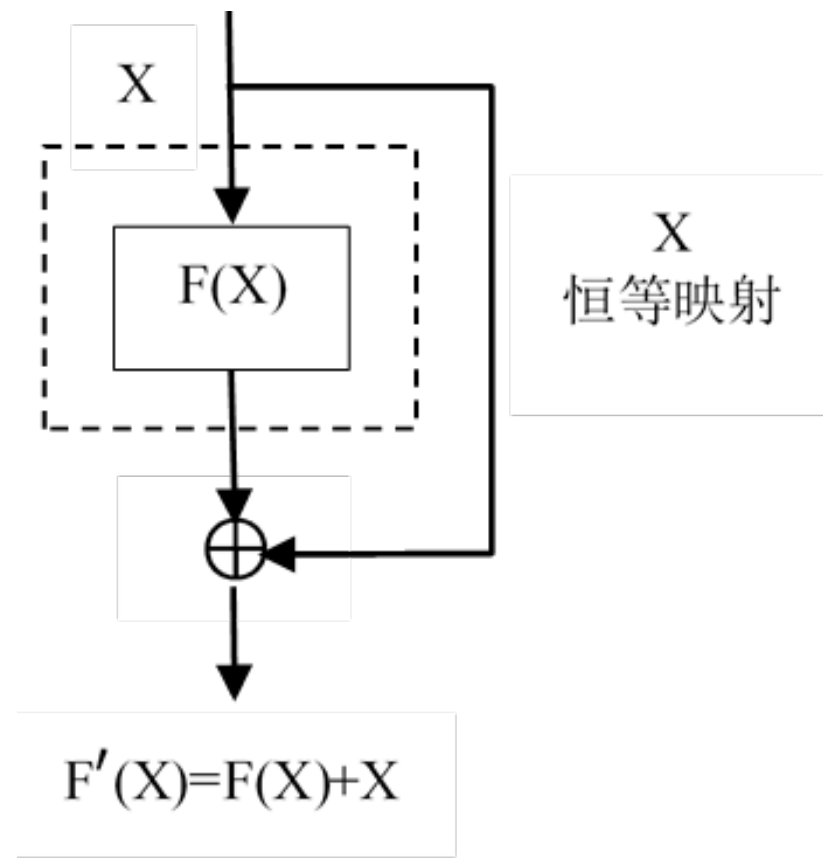


# 残差连接

残差网络

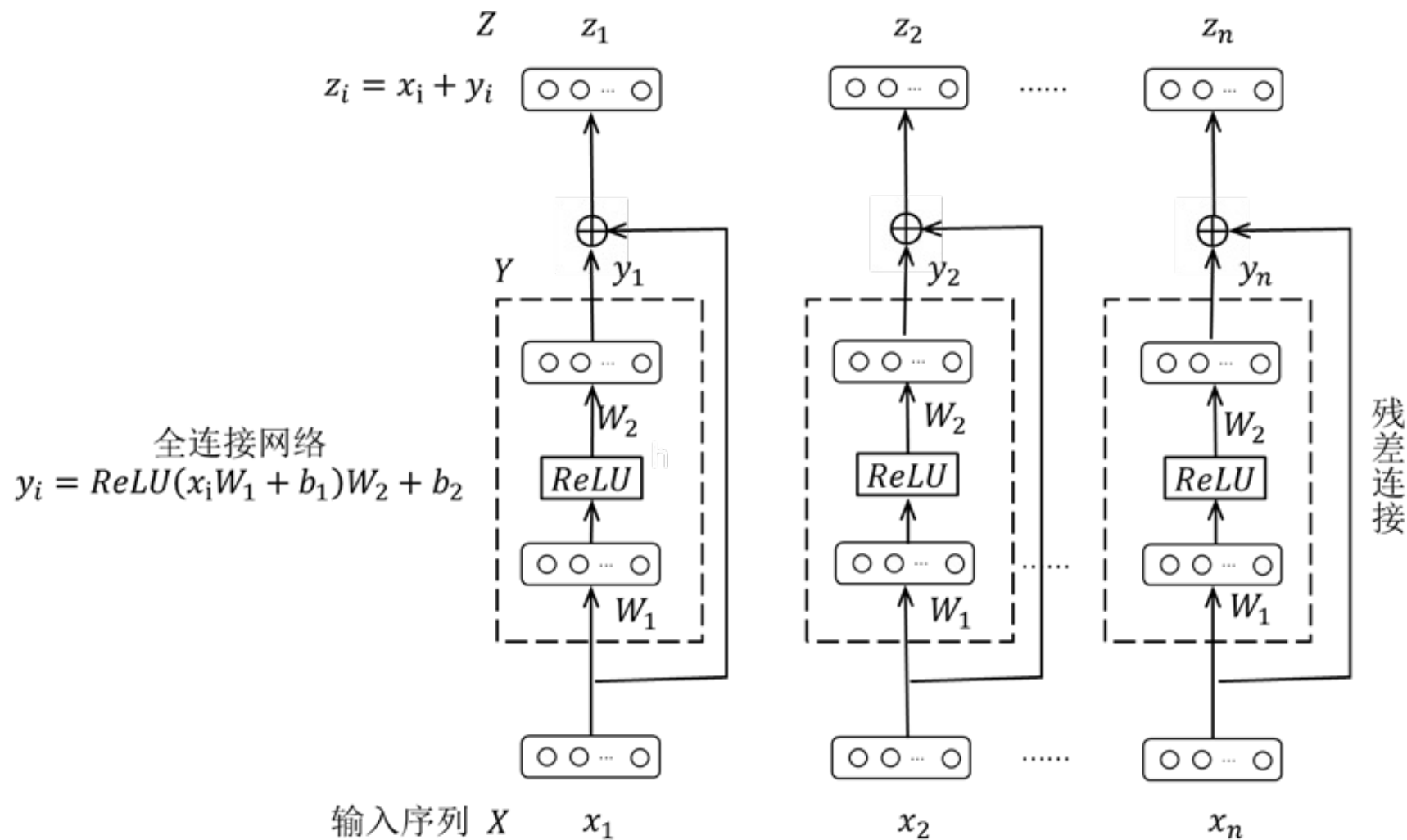


一般形式

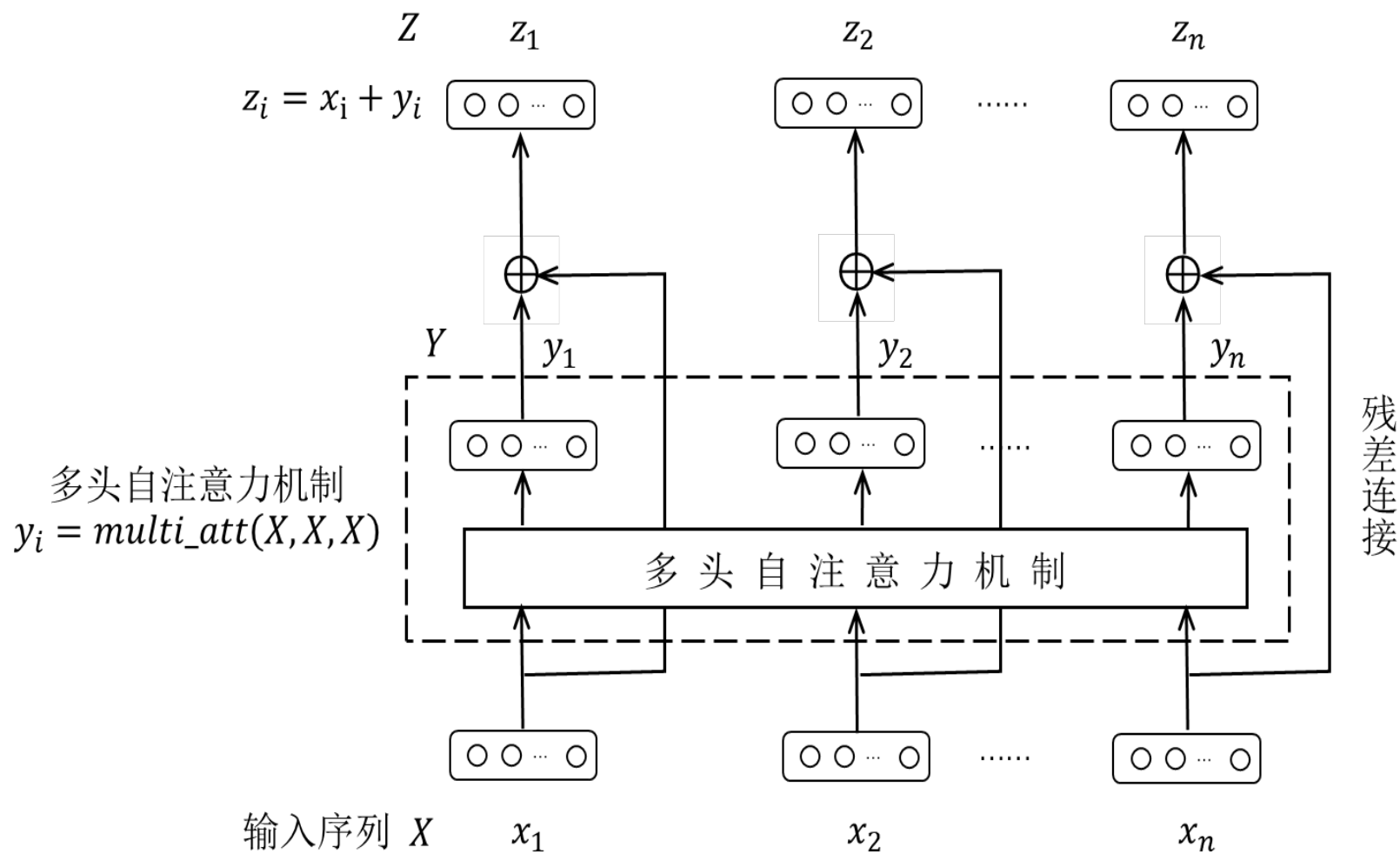




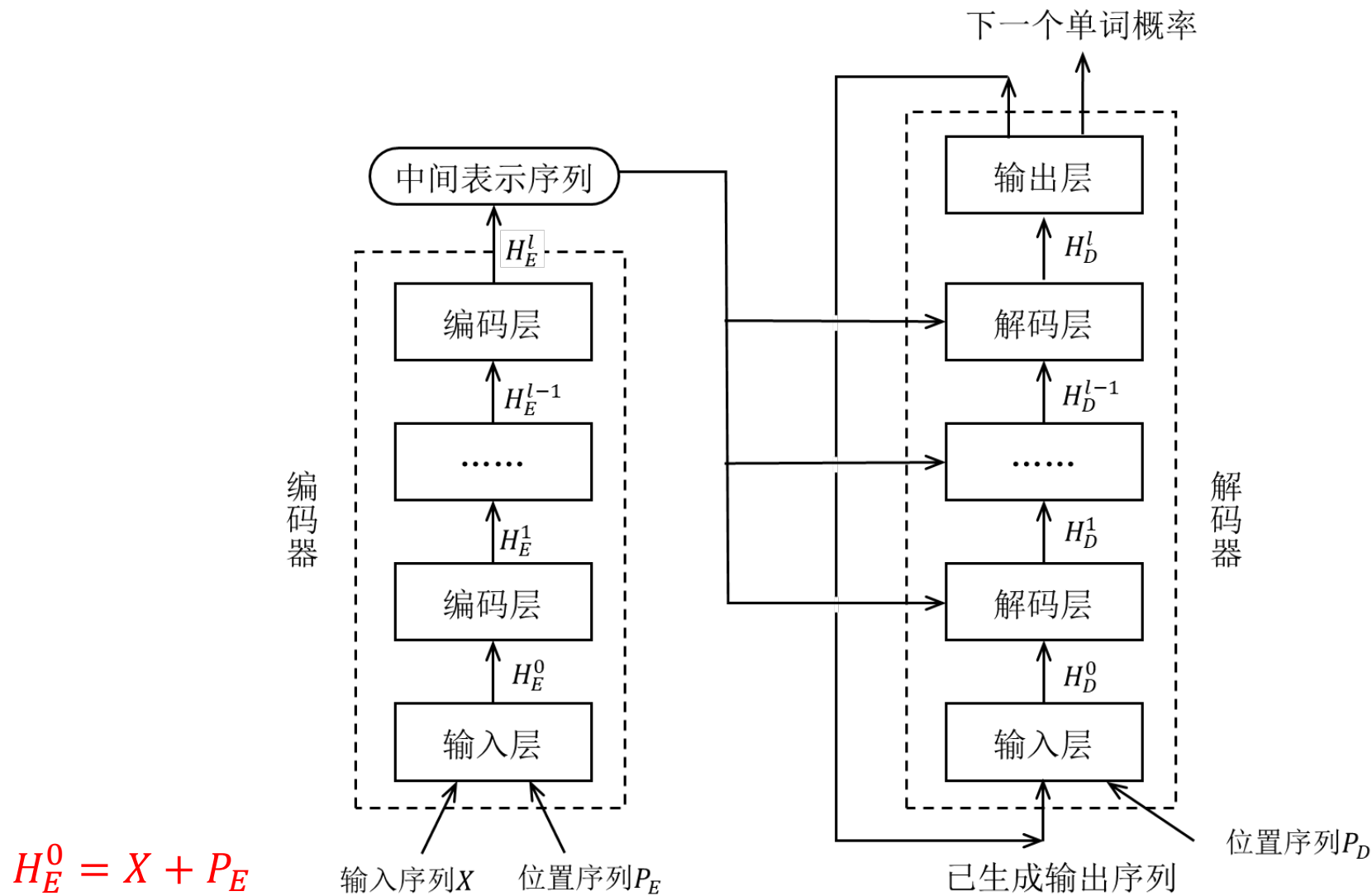
# 全连接网络+残差连接



# 多头注意力机制+残差连接

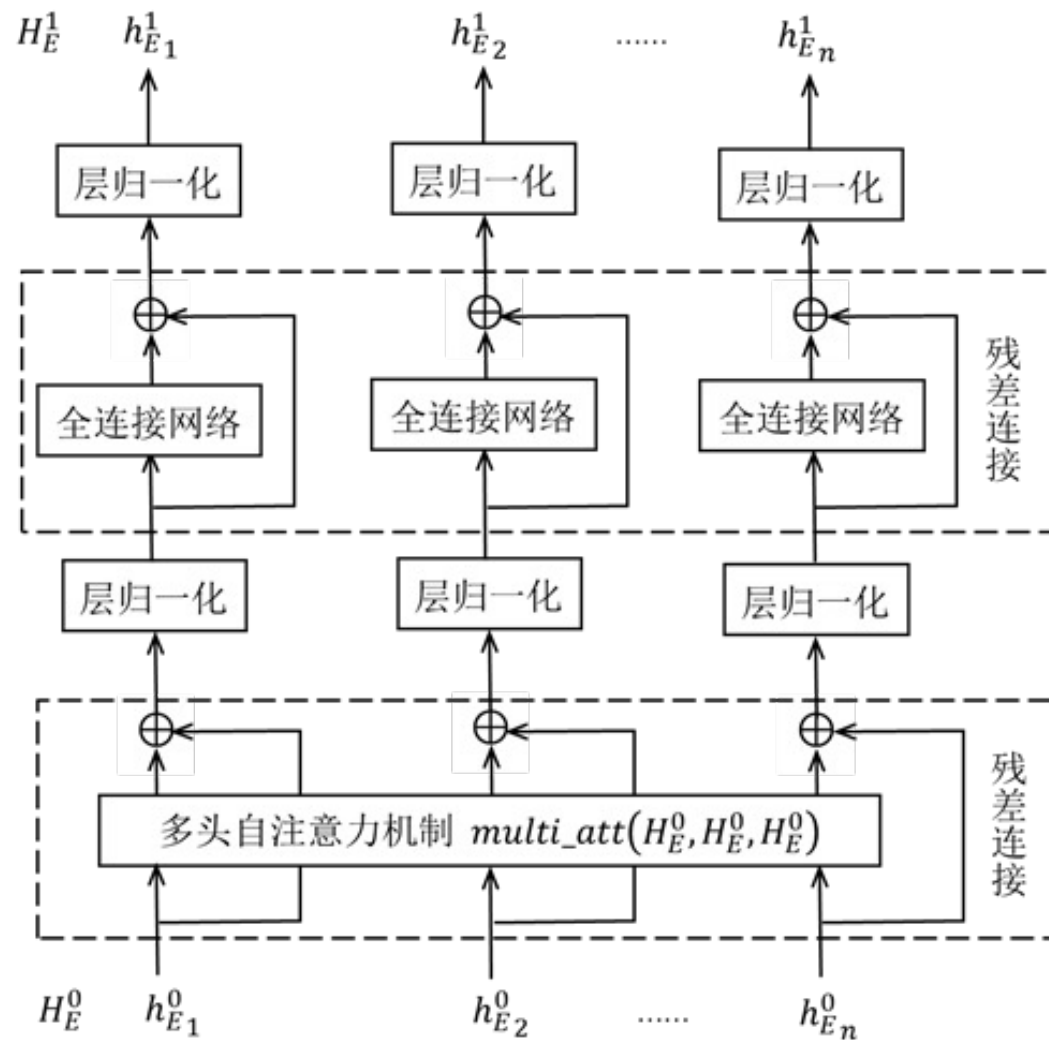


# Transformer

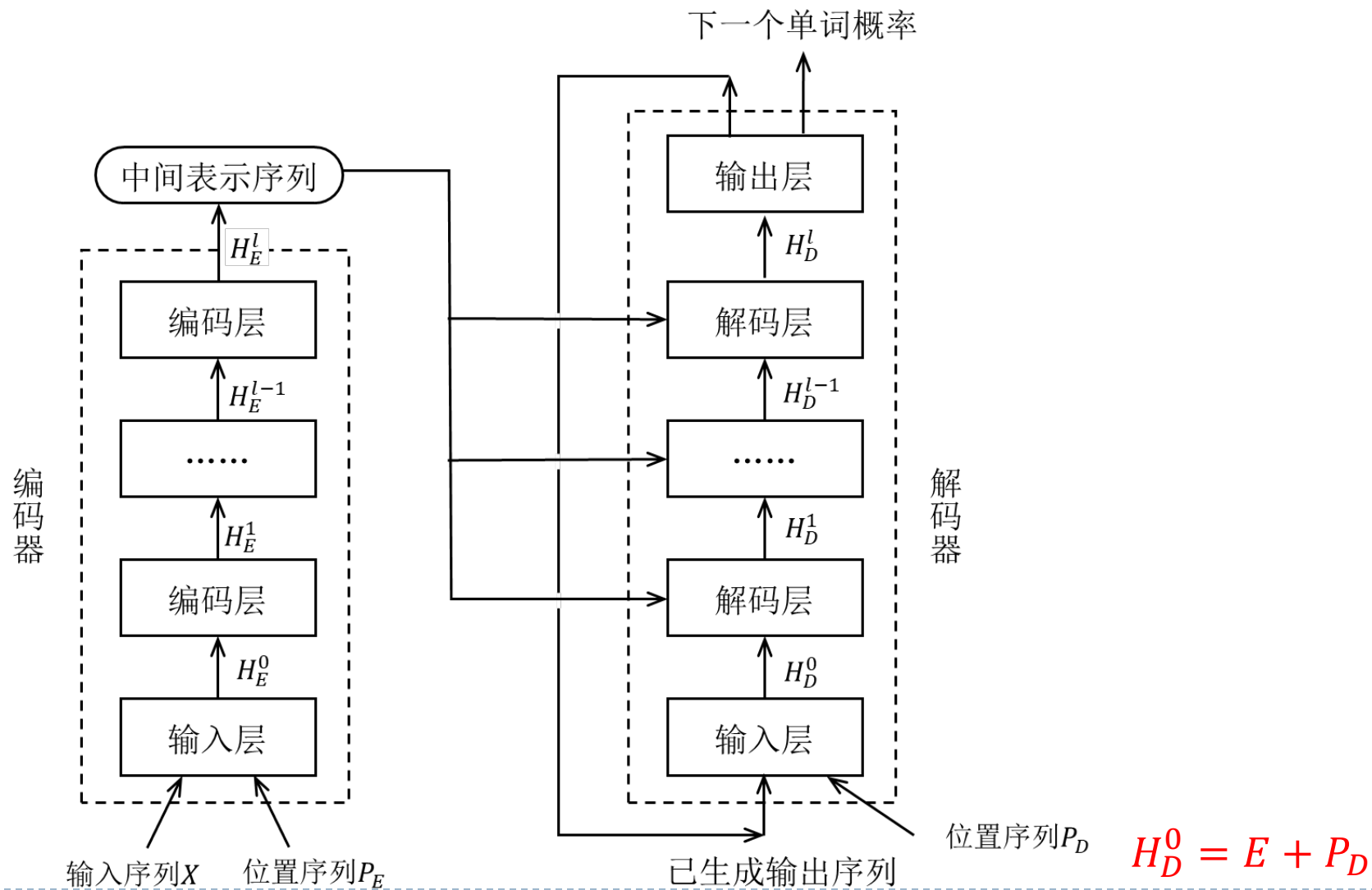


# 编码层

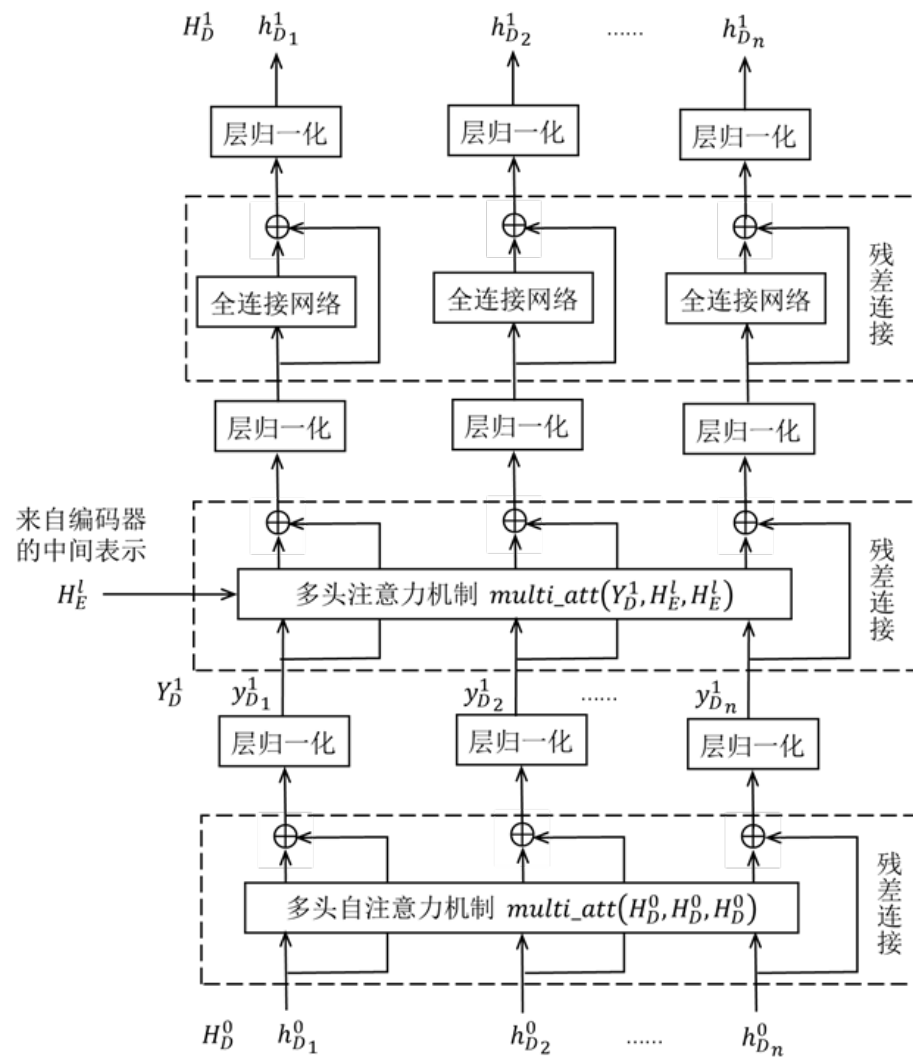
$$H_E^0 = X + P_E$$



# Transformer

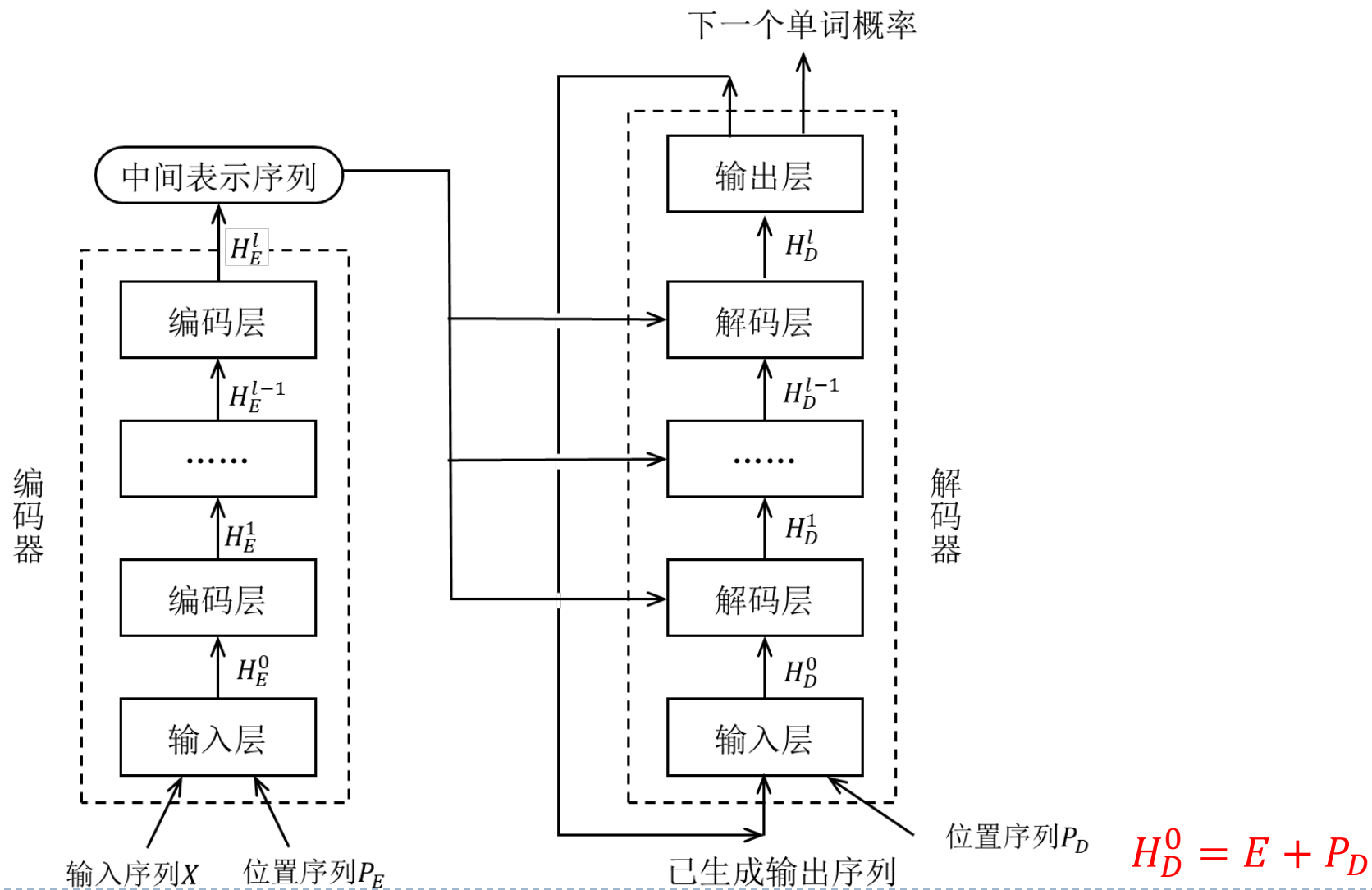


# 解码层



$$H_D^0 = E + P_D$$

# Transformer

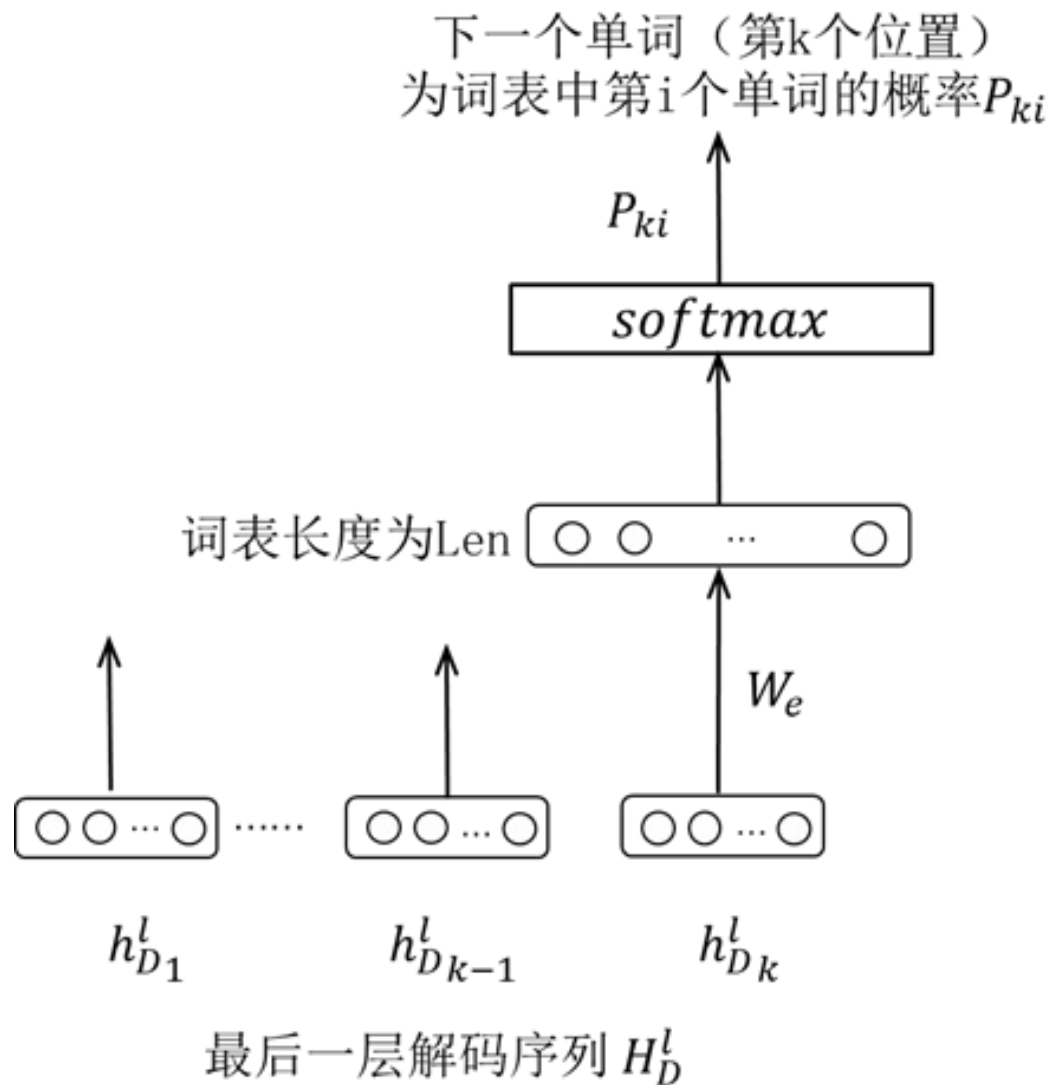


# 输出层

$W_e$ : 第i行为词表中第i个单词的词向量

$$P_k = \text{softmax}(h_{D_k}^l \cdot W_e^T)$$

$p_k$ 的维度为词表长度，其第i个元素 $p_{ki}$ 为第k个位置为词汇表中第i个单词的概率





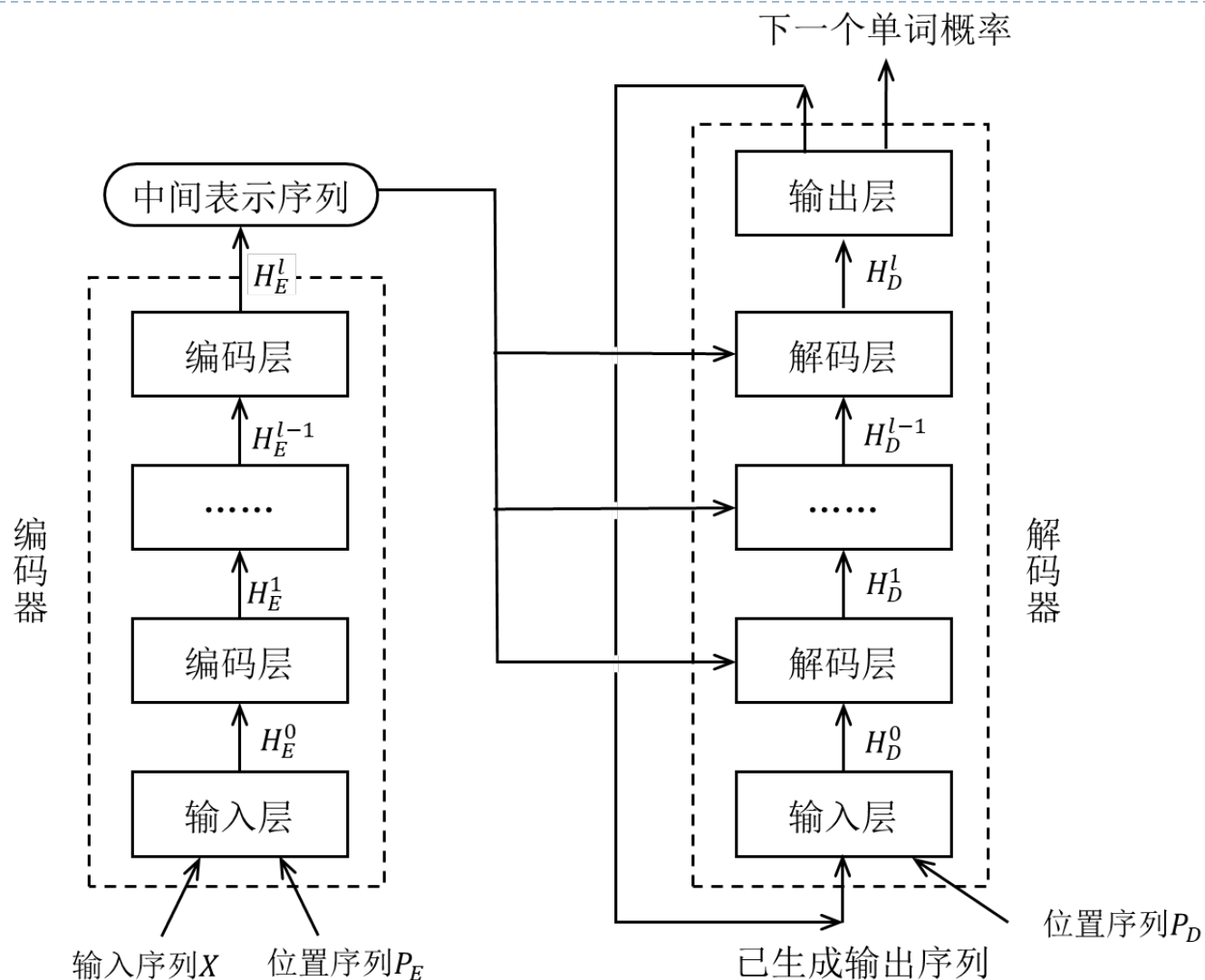
# 输出序列

## ► 编码器

- 运行一次

## ► 解码器

- 自回归
- 运行 $m+1$ 次
- $m$ 为输出序列的长度
- 结束符



# Transformer的训练

---

## ▶ 最大化似然函数

$$\max_{\theta} \prod_{w \in C} p(w = k | context(w), \theta)$$

## ▶ 损失函数

$$L(\theta) = - \sum_{w \in C} \ln(p(w = k | context(w), \theta))$$



# Transformer的应用

---

- ▶ GPT

- ▶ 只使用解码器
- ▶ 自回归

- ▶ BERT

- ▶ 只用编码器
- ▶ 双向（利用前后词预测中间词）

