# 转换器中的多层感知机

# Multilayer Perceptron in Transformer

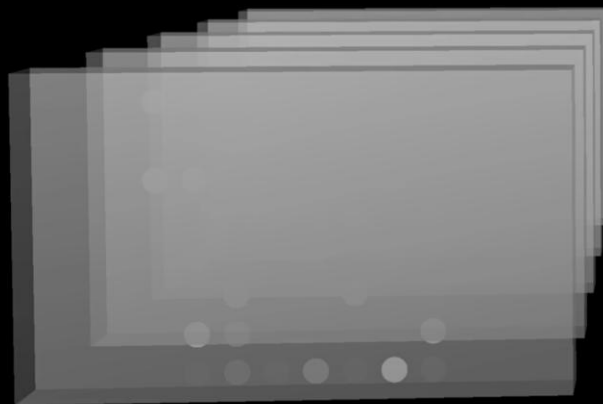# 1 GPT-3中的MLP

- 主要存储世界知识（facts）
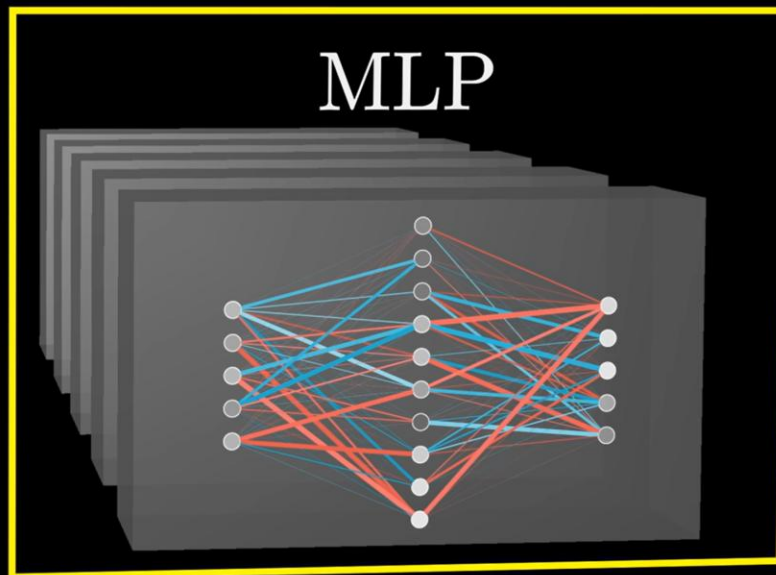
# 1 GPT-3中的MLP

- 2层全连接神经网络



MLP

$$\mathrm{FFN}\bigl(\vec{E}\bigr) = \boldsymbol{W}^D \operatorname{ReLu}\bigl(\boldsymbol{W}^U \vec{E} + \boldsymbol{b}_1\bigr) + \boldsymbol{b}_2$$
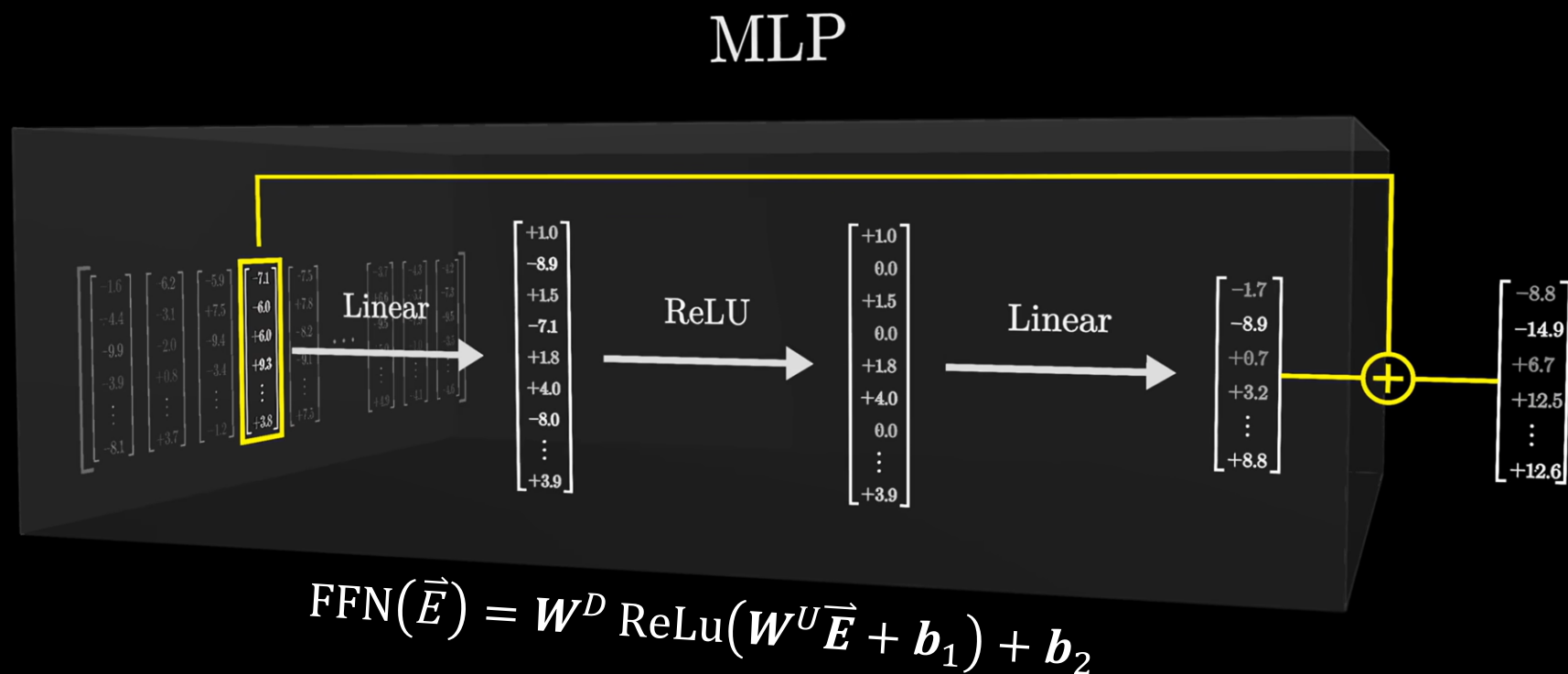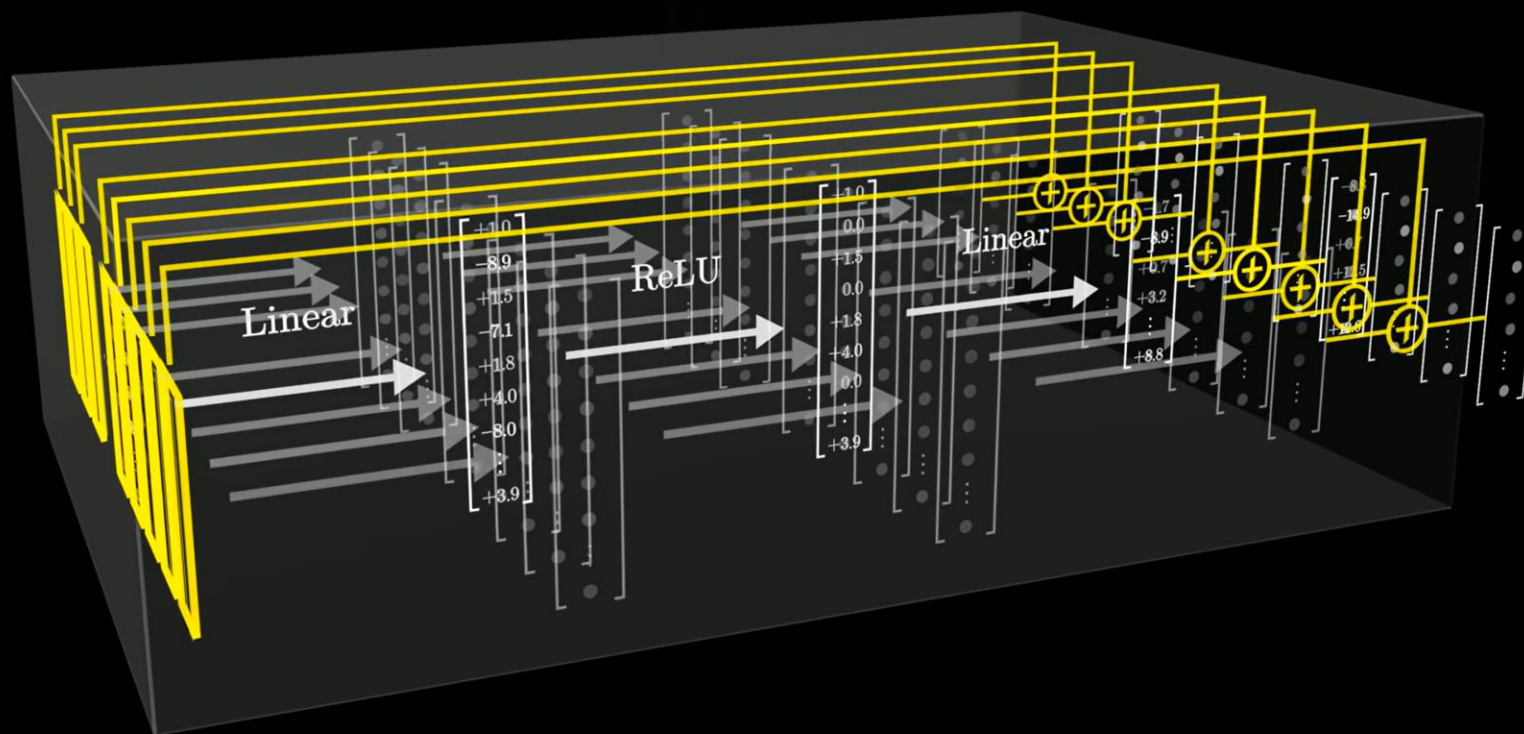
多层感知机（MLP）是最典型的全连接前向神经网络；"前向神经网络（FFN）"强调信息单向传播，"全连接神经网络"强调层间连接方式，三者在多数教学语境下常指同一类模型，但严格来说 FFN 的概念更广。

# 1 GPT-3中的MLP

- 对每个经注意力机制处理后的嵌入向量同时处理

# 2 第一层

• 类比："你问我猜"

# 2 第一层

- 类比："你问我猜"

# 2 第一层

- 类比："你问我猜"

# 2 第一层

- 类比："你问我猜"

# 2 第一层

- 类比："你问我猜"

# 2 第一层

- 类比："你问我猜"



是某种金属吗？

MLP

# 2 第一层

- 类比："你问我猜"

# 2 第一层

- 维度与参数数量

# 2 第一层

- 激活函数：ReLU

# 3 第二层

- "下投影"矩阵

# 3 第二层

- 视作列向量组合

$$\begin{bmatrix} \vec{C}_0 & \vec{C}_1 & \vec{C}_2 & \vec{C}_3 & \vec{C}_4 & \cdots & \vec{C}_m \end{bmatrix} \begin{bmatrix} n_0 \\ n_1 \\ n_2 \\ n_3 \\ n_4 \\ \vdots \\ n_m \end{bmatrix} + \begin{bmatrix} \vec{B} \end{bmatrix} = \begin{bmatrix} -1.7 \\ -8.9 \\ +0.7 \\ +3.2 \\ \vdots \\ +8.8 \end{bmatrix}$$

MLP

# 3 第二层

- 表示为对"你问我猜"的答案的加权组合

$$n_0\vec{C}_0 + n_1\vec{C}_1 + n_2\vec{C}_2 + n_3\vec{C}_3 + n_4\vec{C}_4 + \cdots + n_m\vec{C}_m$$

# 3 第二层

- 表示为对"你问我猜"的答案的加权组合

$$\boxed{n_0 \vec{\mathbf{C}}_0} + n_1 \vec{\mathbf{C}}_1 + n_2 \vec{\mathbf{C}}_2 + n_3 \vec{\mathbf{C}}_3 + n_4 \vec{\mathbf{C}}_4 + \cdots + n_m \vec{\mathbf{C}}_m$$



篮球 $\overrightarrow{\text{Basketball}}$

芝加哥公牛 $\overrightarrow{\text{Chicago Bulls}}$

23号球员 $\overrightarrow{\text{Number 23}}$

1963年出生 $\overrightarrow{\text{Born 1963}}$

迈克尔·乔丹

# 3 第二层

- MLP完整结构

# 3 参数数量分析

- "上投影"矩阵参数数量：

$$4 \times 12{,}288 \times 12{,}288 = 603{,}979{,}776$$



| Up-projection | 49,152      12,288  n_neurons * d_embed = 603,979,776  per layer |
| Down-projection | |
| Unembedding | 50,257      12,288  n_vocab * d_embed = 617,558,016 |

# 3 参数数量分析

- "下投影"矩阵参数数量：



"Down projection"

$$\begin{bmatrix} +0.5 & +8.4 & -4.7 & -8.6 & +4.7 & +5.4 & +8.1 & \cdots & -9.6 \\ -5.3 & +2.3 & +8.9 & +8.9 & +1.1 & +8.2 & +2.8 & \cdots & -0.3 \\ +2.1 & +1.0 & +8.4 & +8.3 & -2.1 & +9.2 & -6.5 & \cdots & -7.2 \\ +0.1 & -9.5 & +8.9 & +6.5 & -9.6 & -6.4 & -3.3 & \cdots & +6.1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ -4.2 & -0.2 & +2.0 & -9.6 & +1.9 & -1.3 & +6.1 & \cdots & +7.8 \end{bmatrix} \begin{bmatrix} +1.0 \\ +0.0 \\ +1.5 \\ +0.0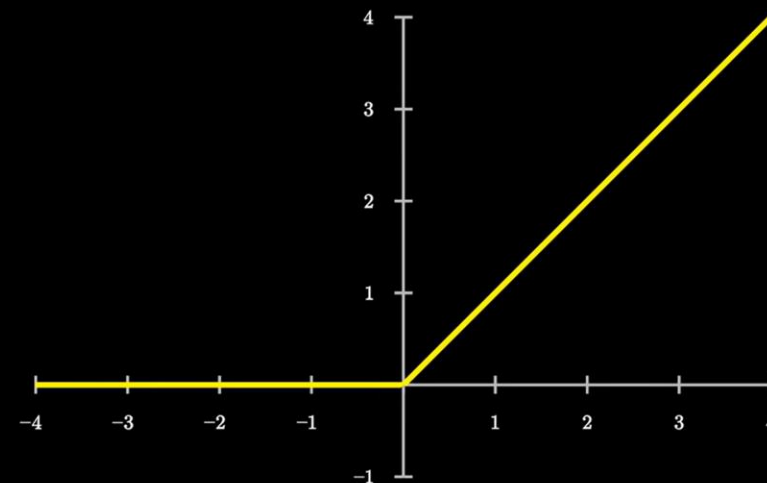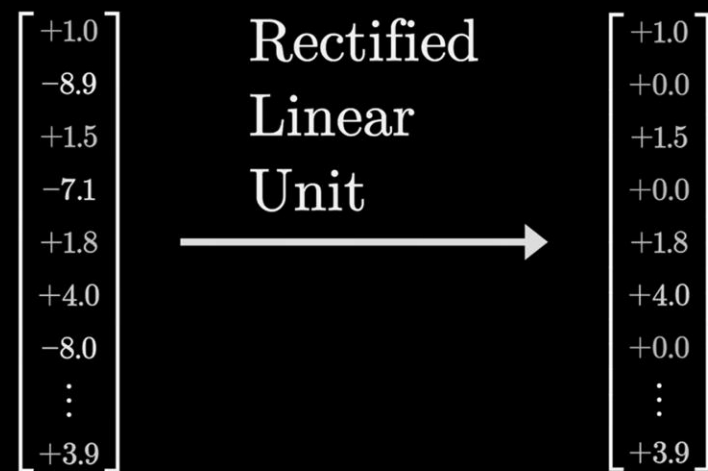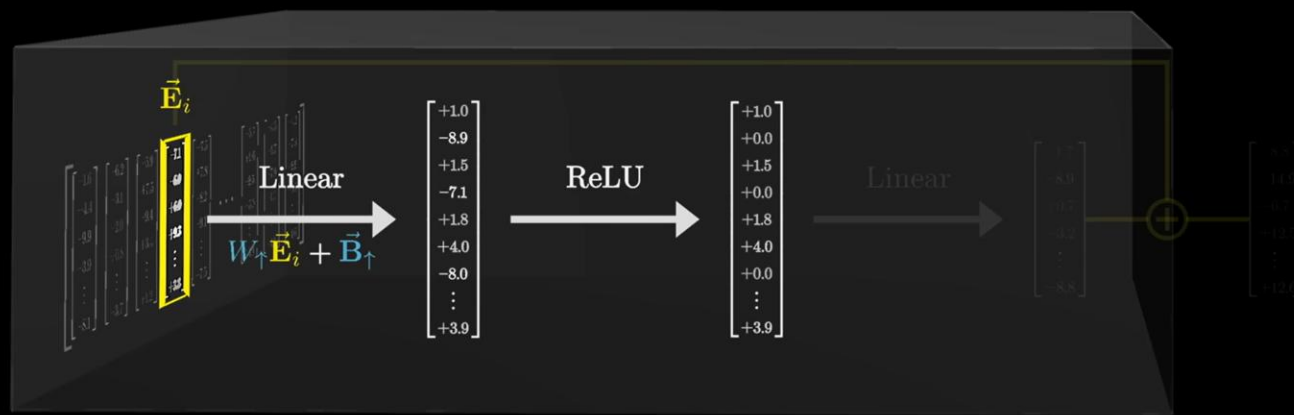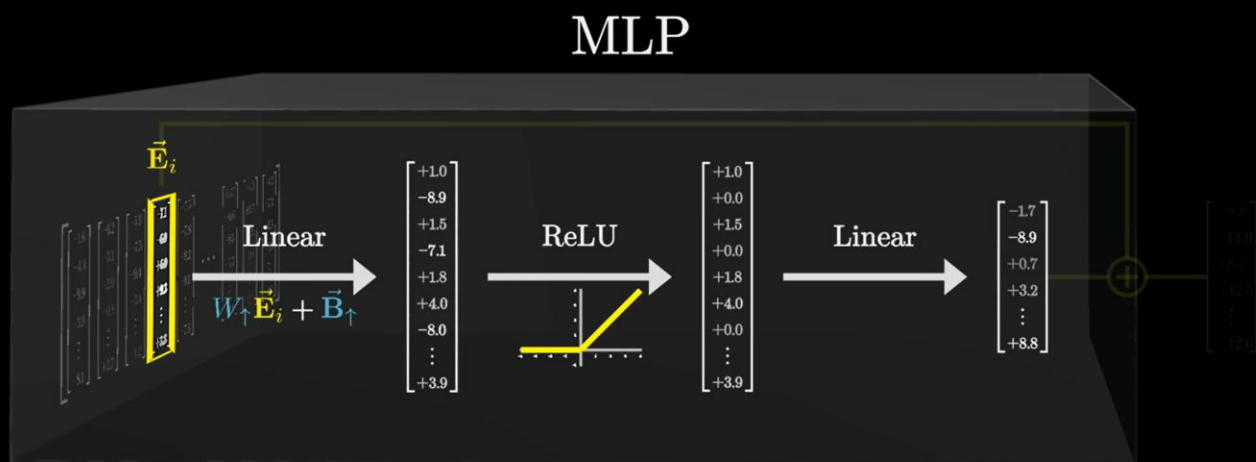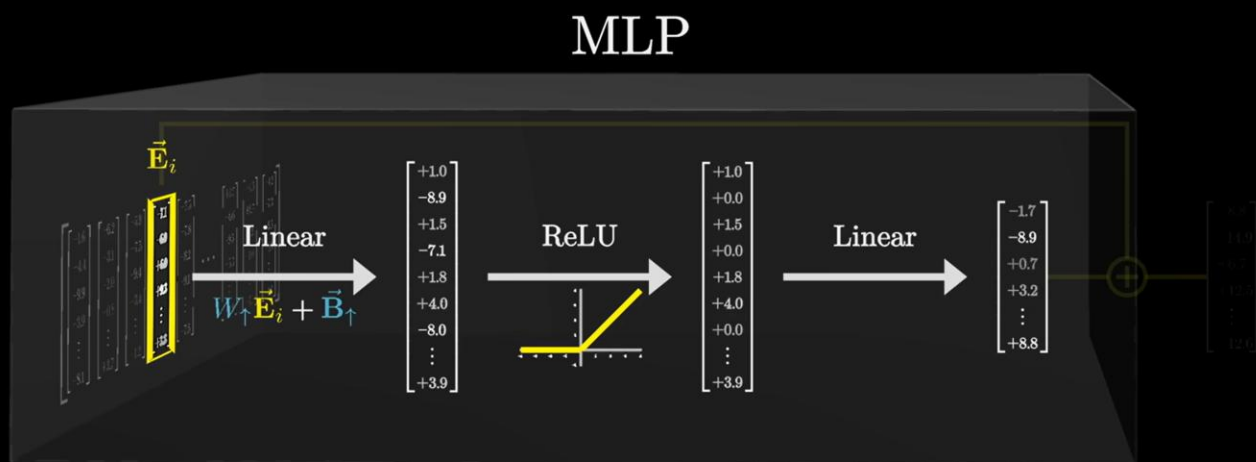 \\ +1.8 \\ +4.0 \\ +0.0 \\ \vdots \\ +3.9 \end{bmatrix} + \begin{bmatrix} +1.5 \\ -6.3 \\ +5.7 \\ +2.2 \\ \vdots \\ -1.6 \end{bmatrix} = \begin{bmatrix} -1.7 \\ -8.9 \\ +0.7 \\ +3.2 \\ \vdots \\ +8.8 \end{bmatrix} \Bigg\} \ 12{,}288$$

| | | |
|---|---|---|
| Up-projection | 49,152      12,288 n_neurons * d_embed = 603,979,776 | per layer |
| Down-projection | 12,288      49,152 d_embed * n_neurons = 603,979,776 | per layer |
| Unembedding | 50,257      12,288 n_vocab * d_embed = 617,558,016 | |

# 3 参数数量分析

- Bias部分参数可忽略

$$4 \times 12,288 \times 12,288 = 603,979,776$$

$$4 \times 12,288$$

$$\frac{4 \times 12,288}{603,979,776} \approx 0.00008$$

# 3 参数数量分析

- × 96层网络



96 Layers

# 3 参数数量分析

- GPT-3参数总量

Total weights: 175,181,291,520

Organized into 27,938 matrices

GPT-3

| | | |
|---|---|---|
| Embedding | $d\_embed * n\_vocab = 617,558,016$ <br> 12,288 $\quad$ 50,257 | |
| Key | $d\_query * d\_embed * n\_heads * n\_layers = 14,495,514,624$ <br> 128 $\quad$ 12,288 $\quad$ 96 $\quad$ 96 | |
| Query | $d\_query * d\_embed * n\_heads * n\_layers = 14,495,514,624$ <br> 128 $\quad$ 12,288 $\quad$ 96 $\quad$ 96 | |
| Value | $d\_value * d\_embed * n\_heads * n\_layers = 14,495,514,624$ <br> 128 $\quad$ 12,288 $\quad$ 96 $\quad$ 96 | |
| Output | $d\_embed * d\_value * n\_heads * n\_layers = 14,495,514,624$ <br> 12,288 $\quad$ 128 $\quad$ 96 $\quad$ 96 | |
| Up-projection | $n\_neurons * d\_embed * n\_layers = 57,982,058,496$ <br> 49,152 $\quad$ 12,288 $\quad$ 96 | |
| Down-projection | $d\_embed * n\_neurons * n\_layers = 57,982,058,496$ <br> 12,288 $\quad$ 49,152 $\quad$ 96 | |
| Unembedding | $n\_vocab * d\_embed = 617,558,016$ <br> 50,257 $\quad$ 12,288 | |

# 3 参数数量分析

- GPT-3参数总量

Total weights: 175,181,291,520
Organized into 27,938 matrices

**GPT-3**

| | | | | | |
|---|---|---|---|---|---|
| **Embedding** | 12,288<br>d_embed | 50,257<br>* n_vocab | | | = 617,558,016 |
| **Key** | 128<br>d_query | 12,288<br>* d_embed | 96<br>* n_heads | 96<br>* n_layers | = 14,495,514,624 |
| **Query** | 128<br>d_query | 12,288<br>* d_embed | 96<br>* n_heads | 96<br>* n_layers | = 14,495,514,624 |
| **Value** | 128<br>d_value | 12,288<br>* d_embed | 96<br>* n_heads | 96<br>* n_layers | = 14,495,514,624 |
| **Output** | 12,288<br>d_embed | 128<br>* d_value | 96<br>* n_heads | 96<br>* n_layers | = 14,495,514,624 |
| **Up-projection** | 49,152<br>n_neurons | 12,288<br>* d_embed | 96<br>* n_layers | | = 57,982,058,496 |
| **Down-projection** | 12,288<br>d_embed | 49,152<br>* n_neurons | 96<br>* n_layers | | = 57,982,058,496 |
| **Unembedding** | 50,257<br>n_vocab | 12,288<br>* d_embed | | | = 617,558,016 |

# 3 参数数量分析

- GPT-3参数总量



$\times 96$

Attention
604M Parameters

MLP
1.2B Parameters

Layer Norm
49K Parameters