

明白了，你关心的是这一步：

```
Q = Q.view(B, L, H, D).permute(0, 2, 1, 3)
```

为什么要这么做，以及数学和逻辑上的原因。

---

## 1 背景

- 输入 `input` : (B, L, 512)
  - B = batch size
  - L = token 数量
  - 512 = embedding 维度
- 多头注意力：
  - 假设有 H = 8 个头
  - 每个头维度 D = 512 / 8 = 64

我们先把向量投影成 Q、K、V：

```
Q = Linear(input) # 仍然 (B, L, 512)
```

然后 reshape 成多头：

```
Q = Q.view(B, L, H, D) # (32, 128, 8, 64)
```

注意：现在头是第 2 个维度（从 0 开始计数）。

---

## 2 为什么 `permute(0, 2, 1, 3)`

最终计算注意力：

```
Attention(Q, K, V) = softmax(Q @ K^T / sqrt(D)) @ V
```

- Q 的 shape 需要是 (B, H, L, D)
- K 的 shape 需要是 (B, H, D, L) (最后两维要交换，用 `transpose(-2, -1)`)

为什么？

- Q @ K^T 的矩阵乘法要在 **head\_dim D** 上相乘

- token 数量 L 要在最后形成矩阵 (L, L) 的注意力分数

所以 permute :

```
Q.view(B,L,H,D).permute(0,2,1,3) → (B,H,L,D)
```

- 0 : batch 不动
- 2 : head 移到第 1 个维度
- 1 : token 数量移到第 2 个维度
- 3 : head\_dim 保持最后

这样每个头都可以独立计算注意力。

---

### ③ 数学上理解

1. reshape 是物理内存不变 · 只是把 512 维拆成 8 个 64 维
2. permute 是逻辑重排 · 把维度顺序调整成 (B, H, L, D) · 方便矩阵乘法操作：

```
[  
scores = Q @ K^T  
]
```

- 对每个 head 计算 : (L,D) @ (D,L) → (L,L)
  - batch 独立 · head 独立
- 

### 小结

1. view(B,L,H,D) → 拆头
2. permute(0,2,1,3) → 调整维度顺序 (B,H,L,D)
3. 这样就可以直接做 Q @ K^T · 每个 head 都独立计算注意力分数

换句话说：reshape 只是切块 · permute 才是把 token、head、embedding 放到正确维度 · 方便矩阵乘法。

所以矩阵的“形状变化”本质上是指索引、访问和乘法规则的变化 · 而内存是线性的。

完全理解 · 这种高维张量确实很难在脑海里“直观想象”。我给你一个可视化思路 · 帮你把这些维度变换在脑海中变得直观一些。我们以多头注意力为例。

---

### ① 把高维拆成“盒子 + 网格”

假设你有张量 Q 形状 (batch, num\_heads, seq\_len, head\_dim) 。可以这样分解想象：

1. **batch** → 想成“样本盒子” · 每个样本独立计算。
  2. **num\_heads** → 每个样本内部又有几个“小头盒子” · 每个头独立做注意力。
  3. **seq\_len** → 每个头盒子里有一行行的 token ( 序列位置 ) · 每行是一个 token 的特征。
-

4. **head\_dim** → 每行的“格子”· 表示 token 的特征维度。

形象比喻：每个 batch 是一个文件夹 · 每个文件夹里有 num\_heads 个表格 · 每个表格有 seq\_len 行 · 每行有 head\_dim 列。

## 2 点积的操作想象

当做  $Q @ K^T$  :

- $K$  的最后两维转置 ( $seq\_len, head\_dim \rightarrow head\_dim, seq\_len$ )
- 对每行 (一个 token 的向量) 做点积 → 得到一个长度为  $seq\_len$  的行向量
- 结果 : 每个 token 对所有 token 的注意力权重

在我们的表格比喻里 :

- 每行 (token) 和另一张表格 (同一个 head) 的所有行 (token) 做向量内积
- 最终每个表格变成“方阵” ( $seq\_len \times seq\_len$ ) · 记录每个 token 看其他 token 的注意力

## 3 为什么 batch 和 num\_heads 不变

- **batch** : 每个文件夹独立做计算
- **num\_heads** : 每个小表格独立做计算
- 底层实现 : PyTorch 用 批量矩阵乘法 · 一次处理所有 batch 和 head · 不需要循环

所以输出形状 ( $batch, num\_heads, seq\_len, seq\_len$ ) 很自然 : 每个 token 对每个 token 的注意力在每个 head 和每个 batch 中都有记录。

## 4 直观小例子 (数字)

假设 :

- $batch=2, num\_heads=2, seq\_len=3, head\_dim=4$
- $Q.shape = (2,2,3,4)$
- $K.shape = (2,2,3,4)$

每个 batch、每个 head 内的  $3 \times 4$  矩阵点积转置 :

$Q[0,0] \ (3 \times 4) @ K[0,0].T \ (4 \times 3) \rightarrow scores[0,0] \ (3 \times 3)$

- 3行3列 → 每个 token 对 3 个 token 的注意力
- 总输出 :  $(2,2,3,3) \rightarrow 2 \text{个样本} \cdot 2 \text{个头} \cdot \text{每个头} 3 \times 3 \text{矩阵}$

你可以把它想成一个 多维立体表格 : batch 和 head 是外层索引 · 内部是  $seq\_len \times seq\_len$  的注意力方阵。

明白了，你想直观理解 (`batch`, `seq_len`, `num_heads`, `head_dim`) 这个顺序的张量维度。我给你一个形象化思路，用“盒子、表格、行列”去想，你脑海里就能“看到”它：

---

## 1 将维度拆开理解

张量形状：

```
(batch, seq_len, num_heads, head_dim)
```

- **batch** → 样本盒子
  - 想成你一次喂模型的 `batch_size` 个句子，每个句子独立处理。
- **seq\_len** → 每个句子的 token 列表
  - 想成表格里的行，每行是一条 token 的信息。
- **num\_heads** → 每个 token 内部的多头注意力分支
  - 想成每行有 `num_heads` 个小格子，每个小格子独立计算注意力。
- **head\_dim** → token 在每个头的特征维度
  - 小格子里真正的数字特征，比如 64 维的向量。

形象比喻：

- `batch` → 文件夹
- `seq_len` → 表格里的行 (token)
- `num_heads` → 每行里的子表格 (每个头)
- `head_dim` → 子表格里的格子 (向量维度)

## 2 可视化一个小例子

假设：`batch=2, seq_len=3, num_heads=4, head_dim=5`

- 你有 2 个句子 (`batch=2`) → 2 个大文件夹
- 每个句子 3 个 token → 每个文件夹里 3 行表格
- 每个 token 有 4 个头 → 每行里有 4 个小格子
- 每个头有 5 维特征 → 小格子里有 5 个数字

如果你要取第一个句子的第二个 token 的第三个 head 的向量：

```
x[0, 1, 2, :] # 取 5 个特征
```

### 3 为什么有时候维度顺序会换

- 有些实现为了方便矩阵乘法，把 `num_heads` 提前，变成 `(batch, num_heads, seq_len, head_dim)`
- 这样可以直接做 `batch matmul`，方便一次处理所有头
- 本质上只是 `reshapetranspose`，数据本身不变，只是视角不同

你可以把它想成：把表格里每行的小格子“旋转”到最前面，让矩阵乘法可以一次处理整列。

为什么batch和num\_heads不变？

好，你问得非常关键：

“**batch** 和 **num\_heads** 为什么在矩阵乘法中不会消失？

不是比喻，我要底层、数学、源码、硬件级别的解释。”

下面我会从数学 → 张量代数 → PyTorch 源码 → GPU 实现

完整给你解释，确保你从本质上理解。



## 结论提前说

**batch** 和 **num\_heads** 不“参与”矩阵乘法，所以不会消失，它们只是被当成“外层循环维度/并行维度”。

真正进入矩阵乘法的是内部两个维度：

```
Q: (... , seq_len, head_dim)  
KT: (... , head_dim, seq_len)
```

参与乘法的是：

```
seq_len × head_dim * head_dim × seq_len
```

外面的维度 (`batch, num_heads`)

被广播成 GPU 的并行任务 **grid/block**，不参与乘法，不会抵消。

## 第一层：数学本质（张量乘法规则）

矩阵乘法只定义在：

```
A: (... , m, k)  
B: (... , k, n)  
C = A @ B → (... , m, n)
```

注意：

- 能被“抵消”的只有内部维度 **k**
- 其他维度（批次、组、通道等）不会参与乘法，而是做 逐元素批处理（**batching**）

这在数学上叫：

\*\*“broadcasted batch matrix multiplication”

（广播批矩阵乘法）\*\*

所以 **batch** 和 **num\_heads** 的角色是：

- 它们不是矩阵维度，只是“多份矩阵”
- PyTorch、GPU 会针对它们生成多个矩阵乘法任务

就像：

有100个学校 → 你要做100次矩阵乘法  
每个学校8个班级 → 每个班级做一次矩阵乘法

**batch**、**num\_heads** 只是索引维度。

## ■ 第二层：PyTorch 的实现（源码级）

在 PyTorch 中，矩阵乘法是调用 **aten::matmul** 或 **bmm**：

当输入形状是：

(batch, num\_heads, seq\_len, head\_dim)

PyTorch 会判断：

- 最后两个维度是矩阵
- 前面的维度是 **batch** 维度

源码（`aten/src/ATen/native/LinearAlgebra.cpp`）中逻辑：

```
// pseudo-code
batch_dims = input.sizes()[:-2]
for each batch_index in batch_dims:
    do matrix_multiply( A[batch_index], B[batch_index] )
```

这里明确规定：

**batch** 维度用于循环调度，而不是矩阵乘法参与的维度。

## ■ 第三层：GPU 底层（CUDA / cuBLAS）

真正的矩阵乘法是调用 cuBLAS 的：

```
cublasSgemmBatched
```

注意函数名里的 **Batched**！

cuBLAS 会把 batch 和 head 维度当成：

- ⌚ 一个要并行处理的矩阵列表 ( List of Matrices)
- ⌚ 每个元素是一个 ( $\text{seq\_len} \times \text{head\_dim}$ ) 的小矩阵

GPU 会：

- 用一个 block 处理一个 batch-head 对应的矩阵乘法
- grid 负责调度所有 ( $\text{batch} * \text{num\_heads}$ ) 个乘法

所以：

**batch** 和 **num\_heads** 是并行维度，不是参与乘法的维度。

## ■ 第四层：为什么 **seq\_len** 和 **head\_dim** 会“抵消”，而 **batch/head** 不会？

因为矩阵乘法的定义要求：

```
(m × k) · (k × n) → (m × n)
```

$k$  ( $\text{head\_dim}$ ) 在中间，是求和维度，必然消失。

而 **batch/head**：

```
(batch, num_heads, m, k)
(batch, num_heads, k, n)
```

它们不是乘法维度，不满足矩阵乘法规则，所以不会被“抵消”。

数学上等价于：

```
for b in batch:  
    for h in num_heads:  
        C[b,h] = Q[b,h] @ K[b,h].T
```

它们只是在“有很多个独立的矩阵需要乘”，  
不会参与乘法本身。

## ■ 第五层：矩阵是怎么存储的？（内存角度）

PyTorch 按 row-major 存储：

```
batch 最外层 → num_heads → seq_len → head_dim
```

也就是说：

- batch 是最外层指针
- head 是 batch 里的子指针
- 真实的矩阵乘法操作的是内层的 2D 区域

GPU 做乘法时只取：

```
tensor[b][h] 的最后2维矩阵
```

它不会读取 batch/head 维度的内容进行乘法，  
因为它们根本不是矩阵的行列。

## ⌚ 最终直观总结（完全非比喻）

### ■ 为什么 batch 不消失？

因为它是矩阵列表的外层索引，不在乘法定义里。

### ■ 为什么 num\_heads 不消失？

因为它也是列表索引，每个 head 的矩阵独立计算。

### ■ 为什么 seq\_len 和 head\_dim 会参与乘法？

因为它们是实际矩阵的行列维度。

### ■ 为什么 head\_dim 会被消掉？

因为它是求和维度 ( summation dimension) 。

### 1. 先看最简单的矩阵乘法

矩阵乘法：

A: ( $m \times k$ )

B: ( $k \times n$ )

$$C = A @ B \rightarrow (m \times n)$$

重点来了：

$k$  是共同维度  $k$  这个维度在乘法后消失因为它被用来“求和”了所以  $k$  就叫求和维度 ( summation dimension)

### 2. 什么叫“被求和掉”？

拿两个矩阵例子：

$$A = [ a1 \ a2 \ a3 ] \text{ shape} = (1 \times 3)$$

$$B = [ b1$$

$$b2$$

$$b3 ] \text{ shape} = (3 \times 1)$$

结果是：

$$C = a1b1 + a2b2 + a3*b3 \text{ shape} = (1 \times 1)$$

你看到没？第一个维度 3 ( $k = 3$ ) 在结果中 没有出现因为它被用于求和所以叫求和维度。