

下面我把“**8个Q头 + 2个KV头**”在这段代码中是如何体现、如何运算的逻辑，逐行拆开解释到最底层，让你彻底理解**GQA (Grouped Query Attention)** 是怎么跑的。

你会看到：

代码并没有直接写“**8个头 / 2个头**”，但它们“**隐藏在维度和 repeat_kv 中**”。

✓ 先给你答案（一句话）

8个Q头来自 n_local_heads, 2个KV头来自 n_local_kv_heads，而 repeat_kv(..., n_rep) 负责把2个KV头复制成8头以匹配Q。

1 关键参数（你必须先理解）

代码中的这几个参数定义了“多少个头”：

名字	含义
self.n_local_heads	Q 的头数，例如 8个head
self.n_local_kv_heads	KV 的真实头数，例如 2个KV head
self.n_rep = n_local_heads / n_local_kv_heads	KV 要复制几次：例如 $8 / 2 = 4$ 次

所以：

```
Q: 8 heads
K: 2 heads → 被 repeat 成 8 heads
V: 2 heads → 被 repeat 成 8 heads
```

这就是**GQA (Grouped Query Attention)** 的核心！

2 Q、K、V 初始 shape（重要）

```
xq = xq.view(bsz, seq_len, self.n_local_heads, self.head_dim)
xk = xk.view(bsz, seq_len, self.n_local_kv_heads, self.head_dim)
xv = xv.view(bsz, seq_len, self.n_local_kv_heads, self.head_dim)
```

如果：

- bsz = 32

- seq_len = 512
- n_local_heads = 8
- n_local_kv_heads = 2
- head_dim = 64

那么 Q、K、V reshape 后变成：

```
xq: (32, 512, 8, 64)  
xk: (32, 512, 2, 64)  
xv: (32, 512, 2, 64)
```

注意：

Q 有 8 个 head · 但 K/V 只有 2 个 head !

3 核心逻辑：repeat_kv —— 2 KV heads → 8 heads

这行代码至关重要：

```
xk = repeat_kv(xk, self.n_rep)  
xv = repeat_kv(xv, self.n_rep)
```

你已经知道：

```
n_rep = n_local_heads / n_local_kv_heads = 8 / 2 = 4
```

repeat_kv 是把 KV 的 head 复制 n_rep 次。

所以：

```
(32, 512, 2, 64)  
→ repeat 4 次  
→ (32, 512, 8, 64)
```

现在形状匹配了 Q：

```
Q: (32, 512, 8, 64)  
K: (32, 512, 8, 64)  
V: (32, 512, 8, 64)
```

4 transpose 成 Attention 所需格式

```
xq = xq.transpose(1, 2)    # → (32, 8, 512, 64)
xk = xk.transpose(1, 2)    # → (32, 8, 512, 64)
xv = xv.transpose(1, 2)    # → (32, 8, 512, 64)
```

最终全部变成：

```
(batch, heads, seq_len, head_dim)
```

5 注意力核心运算（此时 8 个 head 全部独立计算）

Flash Attention 或普通计算：

```
(batch=32, heads=8, Q_len=512, K_len=512)
```

点积过程：

```
xq @ xk.transpose(-2, -1)
→ shape: (32, 8, 512, 512)
```

所以：

- 每个 batch 的每个 head 都单独计算
- KV 头已经被复制，因此也有 8 个 head

✿ 最终你问的问题：8 个头 + 2 个 KV 头在哪里体现？

全部在下面三行：

① 这里体现 Q 有 8 个 head, KV 有 2 个 head

```
xq = xq.view(bsz, seq_len, self.n_local_heads, self.head_dim)
xk = xk.view(bsz, seq_len, self.n_local_kv_heads, self.head_dim)
```

```
xv = xv.view(bsz, seq_len, self.n_local_kv_heads, self.head_dim)
```

② 这里体现 KV 被扩展为 8 个 head

```
repeat_kv(xk, self.n_rep)  
repeat_kv(xv, self.n_rep)
```

③ 这里体现 最终 8 个 head 并行执行 Attention

```
xq.transpose(1, 2)  
xk.transpose(1, 2)  
xv.transpose(1, 2)  
# → (bsz, 8 heads, seq_len, head_dim)
```

🔨 最后总结（最清楚的一句话）

Q 有 8 个头 · K/V 只有 2 个头 · Attention 计算必须 head 对齐 · 所以 K/V 被 repeat 4 次 → $2 \times 4 = 8$ · 最终 Q/K/V 都是 8 个 head → 多头并行。

为什么？

你问的核心问题是：

为什么 Q 有 8 个头 · 而 K/V 只用 2 个头 · 然后再复制成 8 个头？为什么不直接让 K/V 也有 8 个头？

这个问题非常关键 · 它背后是整个 LLaMA、Qwen、Gemma 等模型为什么选择 **GQA (Grouped Query Attention)** 的核心动机。

下面我会用 直觉解释 + 数学解释 + 工程/硬件解释 三个角度让你完全明白——不是“形象解释” · 而是“原理解释”。

█████ 最核心的一句话：

因为 K/V 不需要像 Q 一样太多头 · 少头的 K/V 在保持模型性能几乎不变的情况下 · 大幅降低显存和计算量。

然后再复制 (repeat) 使得 Q 能用更多注意力头来表达更丰富的语义。

这就是为什么 Q 头多 · K/V 头少 · 最后再复制来匹配。

1. 原因一：K/V 的作用和 Q 不对称

在 Attention 中：

$$\text{Attention}(Q, K, V) = \text{softmax}(QK^T / \sqrt{d}) * V$$

从公式可以看出：

 Q : 控制“我去看哪里”

代表的是查询能力。

更多的 Q heads → 更细粒度的注意力 → 更丰富的语义表达。

 K/V : 控制“我提供什么信息”

K/V 的功能是：

- K : 提供地址
- V : 提供内容

这部分不需要像 Q 那么多头。

研究结果：**K/V** 头数减少对性能影响很小。

因此：

- Q : 保持多头 (例如 32)
- K/V : 只需要少头 (例如 4)

2. 原因二：K/V 是跨 token 缓存的 (KV Cache)

在推理/生成时：

- Q 每次重新计算 (每生成一个 token)
- K/V 会一直缓存下来 (KV Cache)

如果模型有：

```
32 Q heads
32 K heads
32 V heads
seq=2048
```

KV Cache 大小为：

```
batch * seq_len * heads * head_dim * 2
```

非常巨大。

现在 GQA 用：

```
32 Q heads  
4 KV heads
```

KV Cache 直接减少 **8 倍！！！**

- ❖ 推理成本最重的是 KV Cache
- ❖ 减少 KV heads 成本显著下降
- ❖ 而模型性能几乎不变

这就是 大型语言模型普遍选择 **GQA** 的根本原因。

■ 3. 原因三：计算量减少

注意力计算复杂度：

```
QKT : seq_len2 * heads
```

如果 K/V 只有 1/4 的 head：

```
K/V heads = Q heads / 4
```

计算复杂度减少 4 倍。

■ 4. 原因四：复制 KV (repeat_kv) 很便宜

repeat_kv 做的事情：

```
K(2 heads) → repeat → K(8 heads)
```

这是简单的 **tensor expand (广播)**，没有重新计算，只是视图变换，是 O(1) 操作。

但 Attention 需要 K/V 头数 = Q 头数，所以必须做复制。

这就是：

Q：保持丰富的视角

K/V：降低成本 + 被复制给 Q 使用

5. 原因五：多头注意力的本质是分组线性变换

本质上：

- Q：需要多个“方向”查询 → 多头越多越好
- K/V：只提供信息 → 不需要太多方向

这是模型结构本身的数学属性。

⌚ 最终总结（评价最高的那种一段话）

注意力中 Q 的功能是“决定看哪里”，因此需要更多的 head 来捕获多维语义空间。

而 K/V 的作用只是“提供位置和内容”，不需要这么多 head。

减少 K/V 的 head 数量可以成倍降低推理显存和计算量，而不会明显降低性能。

为了让 Q 和 K/V head 数一致进行 Attention，模型把少数 K/V head 重复复制成 Q 的数量。

这就是 GQA 的原理。