

PR #42080 完整报告

vllm-project/vllm

[feat] Add FP8 per-tensor Q scale support to Triton attention backend

合并时间: 2026-05-20 00:02

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/42080>

执行摘要

- 一句话: 修复 Triton 注意力 FP8 查询缩放错误
- 推荐动作: 值得精读, 特别是 Triton 内核中通过 `constexpr` 实现编译时降级的设计模式。对于需要扩展量化支持的开发者有参考价值。

功能与动机

Triton 注意力后端错误地允许 FP8 查询用于每张量 KV 缓存模式, 但内核忽略了缩放因子, 产生错误结果。本 PR 通过在内核中添加适当的 Q 缩放支持来修复此问题。

实现拆解

1. 内核层(`vllm/v1/attention/ops/triton_unified_attention.py`): 在内核 `kernel_unified_attention` 中新增 `q_scale` 指针参数和 `Q_IS_FP8 constexpr`, 推导 `USE_FP8_Q_DESCALE` 控制是否折叠缩放。在循环前初始化 `score_scale` 和 `value_scale`: 当 `USE_FP8_Q_DESCALE` 为真时, 将 `q_scale` 和 `k_scale` 折叠到 `score_scale`, 并将 `v_scale` 设为 `value_scale`。在点积和 epilogue 中分别使用这些缩放值。
2. 启动器更新: 在 `unified_attention wrapper` 中, 根据 `q.dtype` 计算 `Q_IS_FP8` 并传递给内核; 将 `q_descale` 指针传递给内核的 `q_scale` 参数。
3. 后端层(`vllm/v1/attention/backends/triton_attn.py`): 在 `TritonAttentionImpl.forward` 中, 当查询 `dtype` 为 FP8 且 KV 缓存模式为 `FP8_PER_TENSOR` 时, 将 `layer._q_scale` 作为 `q_descale` 传递给 `unified_attention`。移除了原来断言 `_q_scale_float == 1.0` 的代码。
4. 测试配套(`tests/kernels/attention/test_triton_unified_attention.py`): 修改 FP8 参数化测试, 使用非 1 缩放 (`q_scale=0.75, k_scale=0.5, v_scale=0.25`) 并显式设置 `kv_quant_mode=FP8_PER_TENSOR`。新增 `test_triton_unified_attn_bf16_query_fp8_kv` 测试 bf16 查询与 FP8 KV 缓存的解量化路径。

关键文件:

- `tests/kernels/attention/test_triton_unified_attention.py` (模块 测试; 类别 `test`; 类型 `test-coverage`; 符号 `test_triton_unified_attn_bf16_query_fp8_kv`): 添加了新的 bf16 Q + FP8 KV 测试, 并使用非 1 缩放测试 FP8 路径, 确保缩放正确应用。
- `vllm/v1/attention/backends/triton_attn.py` (模块 注意力后端; 类别 `source`; 类型 `core-logic`): 后端层修改, 传递 `q_descale` 到内核, 并根据条件移除断言, 使 FP8 查询缩放正确驱动。

- vllm/v1/attention/ops/triton_unified_attention.py (模块 注意力操作; 类别 infra; 类型 infrastructure) : 内核修改, 添加 Q 缩放指针和 USE_FP8_Q_DESCALE constexpr, 实现缩放折叠逻辑。

关键符号: kernel_unified_attention, unified_attention, TritonAttentionImpl.forward

关键源码片段

tests/kernels/attention/test_triton_unified_attention.py

添加了新的 bf16 Q + FP8 KV 测试, 并使用非 1 缩放测试 FP8 路径, 确保缩放正确应用。

```
# 在 test_triton_unified_attn 函数中, 当 q_dtype 不是 None 时:
if q_dtype is not None:
    # 使用非 1 缩放, 显式测试 FP8 解量化路径
    q_scale = torch.tensor(0.75, dtype=torch.float32)
    k_scale = torch.tensor(0.5, dtype=torch.float32)
    v_scale = torch.tensor(0.25, dtype=torch.float32)
    q_descale = q_scale
    scale_shape = (num_seqs, num_kv_heads)
    k_descale = torch.full(scale_shape, k_scale.item(), dtype=torch.float32)
    v_descale = torch.full(scale_shape, v_scale.item(), dtype=torch.float32)
    # 先除以缩放再量化, 模拟真实量化
    maybe_quantized_query = (query / q_scale).to(q_dtype)
    maybe_quantized_key_cache = (key_cache / k_scale).to(q_dtype)
    maybe_quantized_value_cache = (value_cache / v_scale).to(q_dtype)
    kv_quant_mode = KVQuantMode.FP8_PER_TENSOR

# 调用 unified_attention 时传递 kv_quant_mode
unified_attention(
    q=maybe_quantized_query,
    k=maybe_quantized_key_cache,
    v=maybe_quantized_value_cache,
    out=output,
    q_descale=q_descale,
    k_descale=k_descale,
    v_descale=v_descale,
    kv_quant_mode=kv_quant_mode,
    # 其他参数 ... # 注意: 这里需要包含所有必选参数
)
```

vllm/v1/attention/backends/triton_attn.py

后端层修改, 传递 q_descale 到内核, 并根据条件移除断言, 使 FP8 查询缩放正确驱动。

```
# FP8 per-tensor / auto path (original flow).
else:
    key_cache, value_cache = kv_cache.unbind(1)
    if (
        is_quantized_kv_cache(self.kv_cache_dtype)
        and key_cache.dtype != self.fp8_dtype
    )
```

```

):
    key_cache = key_cache.view(self.fp8_dtype)
    value_cache = value_cache.view(self.fp8_dtype)
    descale_shape = (
        attn_metadata.query_start_loc.shape[0] - 1,
        key_cache.shape[2],
    )
    # 仅当 query 为 FP8 且 KV 模式为每张量时才传递 q_descale
    q_descale = (
        layer._q_scale
        if (
            self._kv_quant_mode == KVQuantMode.FP8_PER_TENSOR
            and query.dtype == self.fp8_dtype
        )
        else None
    )
    k_descale = layer._k_scale.expand(descale_shape)
    v_descale = layer._v_scale.expand(descale_shape)
    k_scale_cache = None
    v_scale_cache = None

```

vllm/v1/attention/ops/triton_unified_attention.py

内核修改，添加 Q 缩放指针和 USE_FP8_Q_DESCALE constexpr，实现缩放折叠逻辑。

```

# in kernel_unified_attention, before main loop:
score_scale = scale
value_scale = 1.0
if USE_FP8_Q_DESCALE:
    # Fold q_scale and k_scale into score_scale, apply v_scale to output
    score_scale = scale * tl.load(q_scale) * tl.load(k_scale)
    value_scale = tl.load(v_scale)

# in dot product (inside loop):
# S += score_scale * tl.dot(Q, K) # previously scale * tl.dot(...)

# in epilogue (after loop):
if USE_FP8_Q_DESCALE:
    acc *= value_scale

```

评论区精华

gemini-code-assist[bot] 提出将 `q_descale` 传递耦合到 `kv_quant_mode` 可能限制其他模式，建议始终根据 `query.dtype` 传递。DomBrown 回应解释 FP8 查询仅发生在 FP8_PER_TENSOR 模式下，当前实现是安全的。最终 mgoin 批准合并。

- FP8 查询去缩放与 KV 量化模式的耦合 (correctness): 当前实现被接受，因为 FP8 查询仅与 FP8 每张量 KV 缓存一起使用。

风险与影响

- 风险：内核修改通过 constexpr 编译分支，不影响非 FP8 路径。测试覆盖了所有相关组合（FP8 Q+FP8 KV, bf16 Q+FP8 KV），但未覆盖 FP8 Q 与非 FP8 KV 的组合（这种组合在实际中不会出现）。风险较低。
- 影响：对使用 FP8 每张量 KV 缓存和 FP8 查询的用户，本 PR 修复了精度错误。对使用 bf16/fp16 查询的用户无影响。代码库增加约 30 行，维护成本低。
- 风险标记：核心 Kernel 变更，新增编译分支，未覆盖 FP8 Q 与非 FP8 KV 场景

关联脉络

- 暂无明显关联 PR