

PR #40648 完整报告

vllm-project/vllm

[Model Runner v2] Fix block table IMA issue

合并时间: 2026-04-29 23:30

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/40648>

执行摘要

- 一句话: 修复 v2 模型运行器中块表在 CuMem 唤醒后的非法内存访问
- 推荐动作: 建议精读该 PR, 特别是 `init_block_table_layout_tensors` 的设计和 `post_kv_cache_wake_up` 的抽象。关注 CUDA Graph 兼容性风险, 并评估在唤醒后是否需要重新捕获图。

功能与动机

在运行 `VLLM_USE_V2_MODEL_RUNNER=1` 的测试时出现 `CUDA illegal memory access`, 堆栈显示在 `flashinfer` 的 `build_attn_metadata` 中访问 `paged_kv_indptr_prefill_gpu[0]` 时出错。根本原因是 CuMem 唤醒后 block table 的 `data_ptr` 和 `stride` 变为悬空, 导致后续 `gather` 操作访问非法地址。

实现拆解

1. 在 `BlockTables` 类 (`vllm/v1/worker/gpu/block_table.py`) 中将原本在 `__init__` 中分散的指针 / 步幅张量创建逻辑提取为一个独立的 `init_block_table_layout_tensors` 方法, 并在构造函数末尾调用。
2. 在两个模型运行器文件 (`vllm/v1/worker/gpu/model_runner.py` 和 `vllm/v1/worker/gpu_model_runner.py`) 中均新增 `post_kv_cache_wake_up` 方法。前者委托给 `self.block_tables.init_block_table_layout_tensors()`, 后者委托给 `self.init_fp8_kv_scales()`。
3. 在 `GPUWorker.wake_up` (`vllm/v1/worker/gpu_worker.py`) 中, 将原有的仅针对 FP8 量化 KV 缓存的特殊处理替换为统一的逻辑: 当 `tags` 为 `None` 或包含 `'kv_cache'` 时, 直接调用 `self.model_runner.post_kv_cache_wake_up()`, 同时移除了不再需要的 `is_quantized_kv_cache` 导入。
4. 无测试文件变更, 但原测试用例可以通过该修复。

关键文件:

- `vllm/v1/worker/gpu/block_table.py` (模块 块表; 类别 `source`; 类型 `core-logic`; 符号 `init_block_table_layout_tensors`): 核心变更: 提取 `init_block_table_layout_tensors` 方法, 将指针 / 步幅张量的创建与 CuMem 唤醒后重建解耦
- `vllm/v1/worker/gpu_worker.py` (模块 GPU Worker; 类别 `source`; 类型 `dependency-wiring`): 修改 `wake_up` 方法, 用统一的 `post_kv_cache_wake_up` 调用替换

原有条件逻辑，并移除 `is_quantized_kv_cache` 导入

- `vllm/v1/worker/gpu/model_runner.py` (模块 模型运行器; 类别 source; 类型 data-contract; 符号 `post_kv_cache_wake_up`) : 新增 `post_kv_cache_wake_up` 方法, 委托给 `block_tables.init_block_table_layout_tensors()`
- `vllm/v1/worker/gpu_model_runner.py` (模块 模型运行器; 类别 source; 类型 data-contract; 符号 `post_kv_cache_wake_up`) : 新增 `post_kv_cache_wake_up` 方法, 委托给 `init_fp8_kv_scales()`

关键符号: `init_block_table_layout_tensors`, `post_kv_cache_wake_up`

关键源码片段

`vllm/v1/worker/gpu/block_table.py`

核心变更: 提取 `init_block_table_layout_tensors` 方法, 将指针 / 步幅张量的创建与 CuMem 唤醒后重建解耦

```
# vllm/v1/worker/gpu/block_table.py (head 版本)
```

```
class BlockTables:
    def __init__(self, ...):
        # ... 创建 block_tables 和 num_blocks ...

        # 在 head 版本中, 原本内联的指针创建被移到单独的方法
        self.init_block_table_layout_tensors()

    def init_block_table_layout_tensors(self) -> None:
        # 在初始化及 CuMem kv_cache 唤醒后重新生成指针 / 步幅张量。
        # 唤醒后, 底层 GPU 张量被 CuMem 重新分配, 缓存的 data_ptr 值
        # 会变成悬空指针; 同时 block_sizes_tensor 因存储于 kv_cache
        # 池中而内容未定义, 需重新创建。
        self.block_table_ptrs = self._make_ptr_tensor(
            [b.gpu for b in self.block_tables]
        )
        self.block_table_strides = torch.tensor(
            [b.gpu.stride(0) for b in self.block_tables],
            dtype=torch.int64,
            device=self.device,
        )
        self.block_sizes_tensor = torch.tensor(
            self.block_sizes, dtype=torch.int32, device=self.device
        )
        self.input_block_table_ptrs = self._make_ptr_tensor(
            self.input_block_tables
        )
```

`vllm/v1/worker/gpu_worker.py`

修改 `wake_up` 方法, 用统一的 `post_kv_cache_wake_up` 调用替换原有条件逻辑, 并移除 `is_quantized_kv_cache` 导入

```
# vllm/v1/worker/gpu_worker.py (head 版本)

def wake_up(self, tags: list[str] | None = None) -> None:
    from vllm.device_allocator.cumem import CuMemAllocator
    allocator = CuMemAllocator.get_instance()
    allocator.wake_up(tags)
    # 恢复 level 2 sleep 时保存的缓冲区
    if len(self._sleep_saved_buffers):
        model = self.model_runner.model
        for name, buffer in model.named_buffers():
            if name in self._sleep_saved_buffers:
                buffer.data.copy_(self._sleep_saved_buffers[name].data)
        self._sleep_saved_buffers = {}
    # 如果 KV 缓存被唤醒, 则统一调用 post_kv_cache_wake_up
    if tags is None or 'kv_cache' in tags:
        self.model_runner.post_kv_cache_wake_up()
```

评论区精华

核心讨论来自 gemini-code-assist 的自动 review, 它指出 `init_block_table_layout_tensors` 通过重新赋值创建新张量对象可能破坏 CUDA 图缓存, 因为新张量的虚拟地址可能不同于原图捕获时使用的地址, 建议使用 `.copy_()` 原地更新。但作者和合并者未采纳, PR 以重新分配方式合并。这可能是基于唤醒后 CUDA 图会被重新捕获的假设。

- Block table 元数据重新分配与 CUDA Graph 兼容性 (performance): 作者未直接回复, 但 PR 已合并, 当前实现仍采用重新分配方式。可能假设唤醒后 CUDA 图会重新捕获, 或该场景下旧图无需保留。

风险与影响

- 风险: CUDA Graph 兼容性风险: 重新分配 block table 元数据张量可能导致现有 CUDA 图引用旧地址而触发非法内存访问。如果唤醒后不重新捕获图, 可能出现不正确行为。此外, `wake_up` 方法中移除了 `is_quantized_kv_cache` 判断, 所有包含 'kv_cache' 标签的唤醒都会调用 `post_kv_cache_wake_up`, 对非 FP8 场景引入微额外开销但功能正确。
- 影响: 直接影响启用 `VLLM_USE_V2_MODEL_RUNNER=1` 并使用 CuMem 睡眠模式的用户。修复后可使相关测试通过, 并确保模型在 CuMem 唤醒后 block table 元数据正确重建。不影响默认模型运行器路径。对于 FP8 KV 缓存用户, 原有重置行为被保留且更加清晰。
- 风险标记: CUDA Graph 兼容性风险, `wake_up` 通用路径变更

关联脉络

- PR #39337 Model Runner v2 support base PR: PR body 指出该 PR 是 <https://github.com/vllm-project/vllm/pull/39337> 的一部分