

PR #41285 完整报告

vllm-project/vllm

[Model Runner v2] Fix v2 compile counter `num_gpu_runner_capture_triggers` and `num_cudagraph_captured`

合并时间: 2026-05-01 06:20

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/41285>

执行摘要

- 一句话: 修复 V2 模型运行器 CUDA Graph 计数器缺失
- 推荐动作: 值得精读, 尤其关注计数器放置位置的设计讨论。该 PR 展示了在多文件架构下如何正确维护跨模块计数器, 以及处理 review 中不同设计意见的决策过程。

功能与动机

作为 PR #39337 (V2 模型运行器) 的一部分, 运行 `VLLM_USE_V2_MODEL_RUNNER=1` `pytest tests/compile/test_config.py::test_use_cudagraphs[FULL_DECODE_ONLY-1]` 时, 测试失败, 原因是 `num_gpu_runner_capture_triggers` 计数器在 V2 模型运行器中从未被递增, 导致断言 `num_gpu_runner_capture_triggers` 期望为 1 但实际为 0。

实现拆解

1. 在 `vllm/v1/worker/gpu/model_runner.py` 中导入 `compilation_counter` 并在 `capture_model()` 方法中递增 `num_gpu_runner_capture_triggers`: 该计数器记录模型运行器调用 CUDA Graph 捕获的顶层次数, 放在 `needs_capture()` 检查通过之后、实际捕获之前, 以匹配 V1 语义。
2. 在 `vllm/v1/worker/gpu/cudagraph_utils.py` 中导入 `compilation_counter` 并在 `capture()` 方法中递增 `num_cudagraph_captured`: 在每次成功捕获一个 CUDA Graph 后 (即 `self.graphs[desc] = graph` 之后) 递增, 确保每个实际捕获的图形都被计数。
3. 无测试文件变更: 当前 PR 本身是修复测试失败, 并未新增测试文件, 但修复使已有的 `test_use_cudagraphs` 测试通过。

关键文件:

- `vllm/v1/worker/gpu/model_runner.py` (模块 模型运行器; 类别 source; 类型 data-contract): 添加 `compilation_counter` 导入并在 `capture_model()` 中递增 `num_gpu_runner_capture_triggers`, 这是计数器修复的核心入口。
- `vllm/v1/worker/gpu/cudagraph_utils.py` (模块 CUDA Graph 工具; 类别 source; 类型 dependency-wiring): 添加 `compilation_counter` 导入并在 `capture()` 方法中递增 `num_cudagraph_captured`, 确保每个实际捕获的 CUDA Graph 都被计数。

关键符号: `capture_model`, `capture`

关键源码片段

vllm/v1/worker/gpu/model_runner.py

添加 `compilation_counter` 导入并在 `capture_model()` 中递增 `num_gpu_runner_capture_triggers`，这是计数器修复的核心入口。

```
# vllm/v1/worker/gpu/model_runner.py
from vllm.compilation.counter import compilation_counter

@torch.inference_mode()
def capture_model(self) -> int:
    assert self.cudagraph_manager is not None
    if not self.cudagraph_manager.needs_capture():
        logger.warning(
            "Skipping CUDA graph capture. To turn on CUDA graph capture, "
            "ensure `cudagraph_mode` was not manually set to `NONE`"
        )
    return 0

# 记录顶层的 GPU 模型运行器捕获触发次数,
# 放在 needs_capture 检查之后、实际捕获之前,
# 以匹配 V1 语义
compilation_counter.num_gpu_runner_capture_triggers += 1

start_time = time.perf_counter()
gc.collect()
torch.accelerator.empty_cache()
start_free_gpu_memory = torch.cuda.mem_get_info()[0]
# ... 后续捕获逻辑
```

vllm/v1/worker/gpu/cudagraph_utils.py

添加 `compilation_counter` 导入并在 `capture()` 方法中递增 `num_cudagraph_captured`，确保每个实际捕获的 CUDA Graph 都被计数。

```
# vllm/v1/worker/gpu/cudagraph_utils.py
from vllm.compilation.counter import compilation_counter

# 在 capture() 方法内部，在成功捕获每个 CUDA Graph 后:
    self.graphs[desc] = graph
    # 每个 graph 捕获成功后递增计数器
    compilation_counter.num_cudagraph_captured += 1
    self._graphs_captured = True
    return captured_attn_states
```

评论区精华

Reviewer njhill 提出是否可以将 `compilation_counter.num_gpu_runner_capture_triggers` 的递增移至 `cudagraph_manager.capture()` 内部，以保持代码在单一文件中。作者

yewentao256 回应称，保持其在 `GPUModelRunner.capture_model()` 中更准确，因为其意图是计数顶层的模型运行器捕获尝试，且 `capture` 可能在其他上下文中被调用，移动会导致计数不匹配。该讨论达成共识，未变更设计。

- 计数器放置位置：`capture_model()` vs `cuda_graph_manager.capture()` (design):
yewentao256 解释保持其在 `capture_model()` 中更准确，因为 `capture` 可能在其他上下文中被调用，移动会导致计数不匹配。njhill 接受该论点并批准 PR。

风险与影响

- 风险：风险极低。变更仅添加了两行计数器递增代码，且计数器是一个已有模块（`vllm.compilation.counter`），被设计为线程安全的。不会影响 CUDA Graph 捕获逻辑或性能。潜在风险是如果未来 `capture` 方法在新上下文被调用而忘记调整计数器，但当前设计符合预期语义。
- 影响：直接影响仅限于 V2 模型运行器（`VLLM_USE_V2_MODEL_RUNNER=1`）下的编译计数器统计，进而影响依赖这些计数器的测试（如 `test_use_cuda_graphs`）和可能的外部监控。对用户无影响，对开发者而言修复了 V2 模型运行器下的测试失败，确保计数器准确反映 CUDA Graph 捕获状态。
- 风险标记：暂无

关联脉络

- PR #39337 V2 model runner: 本 PR 是 PR #39337（V2 模型运行器）的一部分，修复了该 PR 引入的编译计数器未正确递增的问题。