

PR #35736 完整报告

vllm-project/vllm

[Bugfix] Fix Ray compiled-DAG SHM channel stalls by detaching zero-copy `np.ndarray` logprobs buffers

合并时间: 2026-04-16 23:49

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/35736>

执行摘要

- 一句话: 修复 Ray compiled DAG 在 logprobs 请求下 SHM 通道阻塞, 通过复制只读 numpy 数组脱离零拷贝缓冲区。
- 推荐动作: 该 PR 值得精读, 特别是 `detach_zero_copy_from_model_runner_output` 函数的实现, 展示了如何安全地处理 Ray SHM 中的零拷贝对象。关注注释中关于 `prompt_logprobs_dict` 和 `cu_num_generated_tokens` 的设计决策, 以及如何平衡性能与稳定性。对于涉及 Ray 或分布式执行的开发者, 此修复提供了重要的技术洞察。

功能与动机

Issue #35319 报告了在 $PP > 1$ 的多节点推理中, 带有 logprobs 的请求初始成功但后续崩溃, 出现 `RayChannelTimeoutError: Timed out acquiring the read lock.`。Ray 通道 API 明确警告, 如果零拷贝反序列化对象 (如 `np.ndarray`) 仍在作用域内, 后续 `read()` 调用可能阻塞。vLLM 的 `ModelRunnerOutput.logprobs` 可能包含由 Ray SHM 支持的 numpy 数组, 这些数组通常标记为只读, 保持其引用会阻塞通道并导致超时。因此, 需要复制这些数组以脱离 Ray 缓冲区。

实现拆解

1. 新增核心函数: 在 `vllm/v1/executor/ray_utils.py` 中, 添加 `detach_zero_copy_from_model_runner_output` 函数及其辅助函数 `_copy_if_readonly`。该函数检查 `ModelRunnerOutput.logprobs` 中的 numpy 数组是否为只读, 如果是则复制它们, 从而脱离 Ray SHM 缓冲区。注释解释了为何跳过 `prompt_logprobs_dict` (由 PyTorch 张量支持) 和 `cu_num_generated_tokens` (普通 Python 列表)。
2. 修改 `FutureWrapper`: 在同一文件中, 更新 `FutureWrapper.result` 方法, 在 `ray.get()` 后调用 `detach_zero_copy_from_model_runner_output`, 确保单工作器和多工作器场景下输出都已脱离缓冲区。
3. 集成到执行器: 在 `vllm/v1/executor/ray_executor.py` 中, 导入新函数并在阻塞模式下获取输出后立即调用它, 覆盖无连接器和有连接器的情况。
4. 添加测试配套: 新增 `tests/v1/executor/test_ray_utils.py` 文件, 包含 `_make_readonly` 辅助函数和 `test_detach_zero_copy_from_model_runner_output_copies_only_numpy_views` 测试, 验证只读 numpy 数组被正确复制, 而其他部分保持不变。

关键文件:

- `vllm/v1/executor/ray_utils.py` (模块 执行器; 类别 `source`; 类型 `core-logic`; 符号 `detach_zero_copy_from_model_runner_output, _copy_if_readonly`): 核心修复逻辑文件, 新增了脱离零拷贝缓冲区的函数并修改了 `FutureWrapper`, 直接解决通道阻塞问题。
- `tests/v1/executor/test_ray_utils.py` (模块 测试; 类别 `test`; 类型 `test-coverage`; 符号 `_make_readonly, test_detach_zero_copy_from_model_runner_output_copies_only_numpy_views`): 新增测试文件, 验证脱离零拷贝函数的正确性, 确保只复制只读 `numpy` 数组且其他部分不变。
- `vllm/v1/executor/ray_executor.py` (模块 执行器; 类别 `source`; 类型 `entrypoint`): 集成修复到 `Ray` 执行器, 在阻塞模式下获取输出后立即调用脱离函数, 确保实际执行中缓冲区的及时释放。

关键符号: `detach_zero_copy_from_model_runner_output, _copy_if_readonly, FutureWrapper.result, test_detach_zero_copy_from_model_runner_output_copies_only_numpy_views`

关键源码片段

`vllm/v1/executor/ray_utils.py`

核心修复逻辑文件, 新增了脱离零拷贝缓冲区的函数并修改了 `FutureWrapper`, 直接解决通道阻塞问题。

```
import numpy as np # 新增导入, 用于处理numpy数组
```

```
def detach_zero_copy_from_model_runner_output(output: "ModelRunnerOutput") -> None:
    """从ModelRunnerOutput中脱离Ray SHM通道的零拷贝缓冲区。
```

```
    Ray编译DAG的SHM通道可能返回零拷贝对象 (如`np.ndarray`),
    这些对象由Ray的共享内存对象存储支持。Ray通道文档明确警告,
    如果此类对象仍在作用域内, 后续读取可能阻塞。
```

```
    vLLM可以在`ModelRunnerOutput.logprobs`中返回numpy支持的logprobs。
    如果这些数组由Ray SHM支持 (通常是只读的), 在调度器迭代中保持它们
    会导致通道阻塞并最终触发`RAY_CGRAPH_get_timeout`。
```

```
    复制只读numpy数组, 使返回的输出不再引用Ray的共享内存缓冲区。
```

```
    我们故意不处理`prompt_logprobs_dict`: 这些条目是`LogprobsTensors`,
    由PyTorch拥有的CPU张量支持, 而不是从Ray通道解码的NumPy视图。
```

```
    """
```

```
    if output.logprobs is None:
        return # 如果logprobs为None, 无需处理
```

```
    token_ids, logprobs, ranks, cu_num_generated_tokens = output.logprobs
```

```
    def _copy_if_readonly(arr):
        # 辅助函数: 仅当arr是只读numpy数组时才复制, 避免不必要的开销
```

```

if isinstance(arr, np.ndarray) and not arr.flags.writeable:
    return arr.copy()
return arr

# `cu_num_generated_tokens` 已经是普通的Python列表（或None），
# 所以它永远不会别名Ray SHM缓冲区，可以直接重用。
token_ids_c = _copy_if_readonly(token_ids)
logprobs_c = _copy_if_readonly(logprobs)
ranks_c = _copy_if_readonly(ranks)
if token_ids_c is token_ids and logprobs_c is logprobs and ranks_c is ranks:
    return # 如果没有数组被复制，则提前返回，避免不必要的对象重建

output.logprobs = type(output.logprobs)(
    token_ids_c, logprobs_c, ranks_c, cu_num_generated_tokens
) # 重建logprobs元组，仅替换被复制的数组

```

评论区精华

- `prompt_logprobs_dict` 处理: kouroshHakha 指出 `ModelRunnerOutput` 中的 `prompt_logprobs_dict` 未被处理，作者在函数文档中添加注释说明其由 PyTorch 张量支持，不受 Ray SHM 影响。
- `cu_num_tokens` 跳过: 讨论中建议为 `cu_num_generated_tokens` 的跳过添加注释，作者已实现，解释其为普通 Python 列表，无需复制。
- 错误处理特异性: gemini-code-assist[bot] 建议将 `except Exception:` 改为更具体的 `except ImportError:` 以避免掩盖其他问题，但最终代码中未显示此变更，可能已在提交中调整。所有讨论点均被解决，无未决疑虑。
 - `prompt_logprobs_dict` 未处理的说明 (design): 通过添加注释澄清设计决策，无需代码变更。
 - `cu_num_tokens` 跳过注释 (documentation): 注释已添加，增强了代码可读性。
 - 错误处理特异性建议 (correctness): 建议被考虑，但最终实现可能已优化；无进一步讨论。

风险与影响

- 风险: - 内存开销增加: 复制只读 numpy 数组会引入额外内存分配，但仅限于 `logprobs` 中的数组，且通常较小，影响可控。
- 性能影响: 如果误判可写数组为只读并进行复制，可能导致不必要的性能下降，但通过 `flags.writeable` 检查可最小化。
- Ray 依赖: 修复依赖于 Ray 通道行为，若 Ray API 变更可能失效，但当前基于 Ray 文档警告，风险较低。
- 测试覆盖: 新增测试确保了核心逻辑正确，但未覆盖所有边缘情况（如不同数据类型），但现有测试足以验证主要场景。
- 影响: - 用户影响: 使用 Ray compiled DAG 进行多节点流水线并行且开启 `logprobs` 的用户将不再遇到通道超时崩溃，提升了系统稳定性和可用性。

- 系统影响：修复针对特定崩溃场景，不会影响其他功能；内存开销轻微，但避免了更严重的阻塞问题。
- 团队影响：为 Ray 集成中的零拷贝问题提供了解决方案，可作为类似问题的参考；代码注释增强了可维护性。
- 风险标记：内存开销增加，Ray 集成依赖，核心路径变更

关联脉络

- PR #39990 Fix #33773: Replace unconditional pandas import with PlaceholderModule: 同为 bugfix PR，涉及依赖管理和导入优化，展示了 vLLM 中对可选依赖的处理模式，可作参考。
- PR #40011 [Bugfix] Fix LLM priority normalization for single-string prompts: 同为 v1 标签的 bugfix，修复前端逻辑错误，体现了项目对稳定性的持续改进。