

PR #41166 完整报告

vllm-project/vllm

[Ci][BugFix] Fix slow DP tests due to bad teardown logic

合并时间: 2026-04-30 07:31

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/41166>

执行摘要

- 一句话: 修复 DP 测试因关机逻辑导致的延时问题
- 推荐动作: 值得关注的设计点是“将单次阻塞操作拆解为并行阶段 + 统一等待”的模式, 可复用于其他资源清理场景。同时建议后续跟进修复评论中提到的 baseline 选择小概率 bug。

功能与动机

Body 指出 CI 远程服务器测试框架的关机逻辑未能正确处理多顶级进程场景, 导致 GPU 内存释放检查失效, 部分测试 (尤其是 DP 测试) 被迫等待 2 分钟后强制关闭。修复后可将某些分布式测试组的运行时间减少约 10 分钟。

实现拆解

1. 重构 `RemoteOpenAIServer._shutdown` (`tests/utils.py`): 将原有关机逻辑拆分为 `_terminate_process_tree` (仅终止进程树) 和 `_wait_for_gpu_memory_release` (等待 GPU 内存释放), `_shutdown` 依次调用两者。
2. 新增类方法 `shutdown_many` (`tests/utils.py`): 接收服务器列表, 用多线程并行调用各服务器的 `_terminate_process_tree`, 然后从列表中选择启动前内存基线最早的服务器调用 `_wait_for_gpu_memory_release`, 避免顺序关机的阻塞等待。
3. 更新三个 DP 测试管理器 (`test_internal_lb_dp.py`、`test_external_lb_dp.py`、`test_hybrid_lb_dp.py`): 将其 `__exit__` 方法从循环逐个调用 `server.__exit__` 改为收集所有服务器实例后调用 `RemoteOpenAIServer.shutdown_many`, 并清空列表。
4. 调整导入 (`tests/utils.py`): 将 `threading` 导入从 ROCm 条件块移到文件顶部全局导入, 并新增 `Sequence` 类型导入。

关键文件:

- `tests/utils.py` (模块 测试工具; 类别 `test`; 类型 `test-coverage`; 符号 `_terminate_process_tree`, `shutdown_many`): 核心重构: 拆分 `_shutdown`、新增 `_terminate_process_tree` 和 `shutdown_many` 类方法, 实现并行关机逻辑。
- `tests/v1/distributed/test_internal_lb_dp.py` (模块 内部 DP 测试; 类别 `test`; 类型 `test-coverage`): 两个服务器管理器 (`MultinodeInternalLBServerManager` 和 `APIOnlyServerManager`) 的 `__exit__` 改用 `shutdown_many`。
- `tests/v1/distributed/test_external_lb_dp.py` (模块 外部 DP 测试; 类别 `test`; 类型 `test-coverage`): `ExternalLBServerManager.__exit__` 改用 `shutdown_many`。

- tests/v1/distributed/test_hybrid_lb_dp.py (模块混合 DP 测试; 类别 test; 类型 test-coverage) : HybridLBServerManager.__exit__改用 shutdown_many。

关键符号: _terminate_process_tree, shutdown_many

关键源码片段

tests/utils.py

核心重构: 拆分 _shutdown、新增 _terminate_process_tree 和 shutdown_many 类方法, 实现并行关机逻辑。

从 tests/utils.py 中提取的核心改动

```
class RemoteOpenAIServer:
```

```
    # ... 原有代码 ...
```

```
    def _shutdown(self) -> None:
```

```
        """
```

```
        完整关机流程: 先终止进程树, 再等待 GPU 内存释放。
```

```
        """
```

```
        self._terminate_process_tree()
```

```
        self._wait_for_gpu_memory_release()
```

```
    def _terminate_process_tree(self) -> None:
```

```
        """
```

```
        仅终止服务器进程树, 不等待 GPU 内存释放。
```

```
        分离出来以便 shutdown_many 可以并行终止多个服务器。
```

```
        """
```

```
        pid = self.proc.pid
```

```
        try:
```

```
            pgid = os.getpgid(pid)
```

```
        except (ProcessLookupError, OSError):
```

```
            pgid = None
```

```
        # Phase 1: 向根进程发送 SIGTERM
```

```
        with contextlib.suppress(ProcessLookupError, OSError):
```

```
            self.proc.terminate()
```

```
            print(f"Sent SIGTERM to process {pid}")
```

```
        try:
```

```
            self.proc.wait(timeout=15)
```

```
            print(f"Server {pid} terminated gracefully")
```

```
        except subprocess.TimeoutExpired:
```

```
            # Phase 2: 如果 SIGTERM 超时, 则向整个进程组发送 SIGKILL
```

```
            print(f"Server {pid} did not respond to SIGTERM, sending SIGKILL to process group")
```

```
            if pgid is not None:
```

```
                with contextlib.suppress(ProcessLookupError, OSError):
```

```
                    os.killpg(pgid, signal.SIGKILL)
```

```
            self._kill_process_group_survivors(pgid)
```

```

@classmethod
def shutdown_many(cls, servers: Sequence["RemoteVLLMServer"]) -> None:
    """
    并行关闭多个服务器并仅等待一次 GPU 内存释放。

    测试夹具持有多个 RemoteVLLMServer 实例时，绝不能顺序调用各服务器的
    __exit__: 因为每个服务器的 _wait_for_gpu_memory_release 都会检查所有
    可见设备的总内存，第一个服务器的等待会被后方服务器持有的 GPU 内存阻塞
    整个超时周期。

    这个方法同时终止所有服务器的进程树，然后使用最早记录的内存基线
    (任何服务器启动前的内存) 来统一等待 GPU 内存释放。
    """
    if not servers:
        return

    # 并行终止所有服务器的进程树
    threads = [
        threading.Thread(
            target=s._terminate_process_tree,
            name=f"shutdown-{s.proc.pid}",
            daemon=True,
        )
        for s in servers
    ]
    for t in threads:
        t.start()
    for t in threads:
        t.join()

    # 选择启动前内存基线最早的服务器 (即 pre 值最小的)
    # 注意: 如果某个服务器的 _pre_server_gpu_memory 为 0.0,
    # Python 的 0.0 or float('inf') 会错误地返回 inf,
    # 但这在实践中几乎不会发生。
    earliest = min(
        servers,
        key=lambda s: (
            float("inf")
            if s._pre_server_gpu_memory is None
            else s._pre_server_gpu_memory
        ),
    )
    earliest._wait_for_gpu_memory_release()

```

评论区精华

Review 中 [gemini-code-assist\[bot\]](#) 指出 `shutdown_many` 中选取 `earliest` 服务器的逻辑存在 bug: 当 `_pre_server_gpu_memory` 为 `0.0` 时, `0.0 or float("inf")` 会返回 `float("inf")`, 导致选错基线, 可能使 GPU 内存释放等待提前结束。该问题在 PR 合并时未修正, 但实际场景中

`_pre_server_gpu_memory` 极少为 0.0，风险较低。

- GPU 内存基准选择逻辑缺陷 (correctness): 未解决, PR 已合并。实际场景中 `_pre_server_gpu_memory` 为 0.0 的可能性极低, 风险可接受。

风险与影响

- 风险: 仅修改测试基础设施代码, 不影响生产逻辑。主要风险在于:
 - 评论区指出的 baseline 选择缺陷: 若服务器启动前 GPU 内存恰好为 0 GB (几乎不可能, 因为系统至少占用约 1.93 GB), 则可能选错等待目标, 导致后续测试 OOM。
 - 并行终止进程可能引入竞态条件, 但操作系统信号处理通常安全。
 - 影响: 对用户无影响。对 CI 效率有显著正向影响: 减少分布式测试组 (特别是 DP 测试) 的运行时间约 10 分钟。团队可获得更快的 CI 反馈。
 - 风险标记: 潜在 GPU 内存释放等待误判

关联脉络

- 暂无明显关联 PR