

PR #43768 完整报告

vllm-project/vllm

[BugFix] Fix hard-coded timeout for multi-API-server startup

合并时间: 2026-05-28 15:09

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/43768>

执行摘要

- 一句话: 修复多 API Server 启动硬编码超时
- 推荐动作: 此 PR 为简单的 Bugfix, 解决了实际部署中遇到的超时问题。推荐快速合并并 cherry-pick 到 v0.22.0, 因为问题已影响用户。同时关注 njhill 提出的启动序列重构, 以彻底避免此类问题。

功能与动机

在 PR #42585 中引入了 60 秒的硬编码超时, 但在实际环境中 (如 HF Hub 限流导致的重试) 该超时不足, 导致 API 服务器启动失败。用户反馈见 <https://github.com/vllm-project/vllm/pull/42585#issuecomment-4552626521>。变更旨在通过环境变量 `VLLM_ENGINE_READY_TIMEOUT_S` 提供可配置的超时。

实现拆解

1. 修改默认超时来源: 在 `vllm/v1/utils.py` 的 `gather_actual_addresses` 方法中, 将 `timeout` 参数的默认值从 60.0 改为 `envs.VLLM_ENGINE_READY_TIMEOUT_S`。该环境变量已在其他地方使用, 可提供更灵活的配置。
2. 无其他变更: 仅涉及一行改动, 未引入新配置或测试。

关键文件:

- `vllm/v1/utils.py` (模块启动流程; 类别 source; 类型 core-logic): 修改了 `gather_actual_addresses` 方法的默认超时值, 从硬编码 60 秒改为环境变量 `VLLM_ENGINE_READY_TIMEOUT_S`。

关键符号: `gather_actual_addresses`

关键源码片段

`vllm/v1/utils.py`

修改了 `gather_actual_addresses` 方法的默认超时值, 从硬编码 60 秒改为环境变量 `VLLM_ENGINE_READY_TIMEOUT_S`。

```
# vllm/v1/utils.py
# 在 MultiApiServer 类中, gather_actual_addresses 方法用于等待所有 API 服务器子进程
# 报告其绑定的 ZMQ 地址。原默认超时 60 秒在 HF Hub 限流重试时不够用。
```

```

# 现改为使用环境变量 VLLM_ENGINE_READY_TIMEOUT_S, 用户可灵活配置。
def gather_actual_addresses(
    self,
    timeout: float = envs.VLLM_ENGINE_READY_TIMEOUT_S, # 从硬编码 60.0 改为环境变量
) -> tuple[list[str], list[str]]:
    """Return (inputs, outputs) reported by each child, indexed by
    ``client_index``. Raises ``RuntimeError`` on timeout or premature
    child exit."""
    n = len(self._address_pipes)
    inputs: list[str | None] = [None] * n
    outputs: list[str | None] = [None] * n
    pending: dict[connection.Connection, int] = {
        pipe: i for i, pipe in enumerate(self._address_pipes)
    }
    sentinel_to_idx: dict[Any, int] = {
        proc.sentinel: i for i, proc in enumerate(self.processes)
    }

    deadline = time.monotonic() + timeout
    try:
        while pending:
            remaining = deadline - time.monotonic()
            if remaining <= 0:
                missing = [self.processes[i].name for i in pending.values()]
                raise RuntimeError(
                    f"Timed out after {timeout:.1f}s waiting for "
                    f"API server(s) to report bound ZMQ addresses: "
                    f"{missing}"
                )
            waitables: list[Any] = list(pending.keys()) + list(
                sentinel_to_idx.keys()
            )
            ready = connection.wait(waitables, timeout=remaining)
            # ... 后续处理逻辑不变

```

评论区精华

- njhill最初对超时原因表示疑惑，认为 API 服务器启动不应超过 60 秒。但随后调查发现，`engine_args.create_engine_config()` 中的 HF Hub 查找可能因限流重试而耗时更长。
- njhill认可修复的合理性，同时表示将尝试重构启动序列，使 HF Hub 查找与引擎启动重叠，从根本上减少延迟。
- ehfd请求将此修复 cherry-pick 到 v0.22.0 版本。
- 超时原因和重构方向 (design): 当前修复被接受，njhill 计划后续重构启动序列以从根本上优化。

风险与影响

- 风险：
 - 低风险：仅将默认超时从硬编码 60 秒改为环境变量值，若环境变量未设置则可能回退到默认值（需确认默认值，通常为 600 秒或 3600 秒），因此风险极低。
 - 兼容性：如果用户之前显式传递了 `timeout` 参数，则不受影响。仅影响使用默认超时的调用处。
- 影响：
 - 用户影响：解决了因 HF Hub 限流或网络慢导致的多 API 服务器启动超时问题，提升了部署稳定性。
 - 系统影响：无性能影响，仅启动阶段的超时行为可变。
 - 团队影响：小改动，易于理解和合并。
 - 风险标记：暂无

关联脉络

- PR #42585 前序引入硬编码超时的 PR: 本 PR 修复了 PR #42585 中引入的硬编码 60 秒超时问题。