

PR #25599 完整报告

sgl-project/sclang

[PD] Add conclude_state to fake KV backend

合并时间: 2026-05-18 19:56

原文链接: <http://prhub.com.cn/sgl-project/sclang/pull/25599>

执行摘要

- 一句话: 为 fake KV backend 添加 abort 与状态管理
- 推荐动作: 该 PR 逻辑清晰、改动量小, 属于常规维护级别的 bugfix/ 功能补强, 值得合并。但其重要性较低, 不需精读。

功能与动机

fake KV backend 用于 warmup 请求测试, 原有 poll 逻辑无法模拟传输中断场景, 需支持通过 abort 设置失败状态以测试 PD 框架的容错行为。

实现拆解

1. 新增 `conclude_state` 属性: 在 `FakeKVSender.__init__` 和 `FakeKVReceiver.__init__` 中添加 `self.conclude_state: Optional[KVPoll] = None`, 用于记录传输的最终状态 (Success 或 Failed)。
2. 修改 `poll` 方法: 两个类的 `poll` 入口都先检查 `conclude_state` 非空时直接返回该状态 (短路优先); `FakeKVSender.poll` 移除冗余 `else` 分支, 使用 `not self.has_sent` 替代 `self.has_sent is False`, 与 `FakeKVReceiver.poll` 风格一致; 每种成功路径都将 `conclude_state` 设置为 `KVPoll.Success` 后再返回。
3. 新增 `abort` 方法: 两个类都添加 `abort(self)` 方法, 将 `conclude_state` 设置为 `KVPoll.Failed`, 外部调用后可立即终止传输。
4. 代码风格统一: 同步调整 `FakeKVReceiver.poll` 也增加 `conclude_state` 优先检查逻辑, 保证行为对称。
5. 无测试、配置或部署配套变更: 仅修改了 `python/sclang/srt/disaggregation/fake/conn.py` 一个文件。

关键文件:

- `python/sclang/srt/disaggregation/fake/conn.py` (模块 KV 缓存; 类别 source; 类型 core-logic; 符号 abort): 唯一修改文件, 为 `FakeKVSender` 和 `FakeKVReceiver` 新增 `abort` 与 `conclude_state` 状态管理。

关键符号: `abort`, `poll`

关键源码片段

python/sglang/srt/disaggregation/fake/conn.py

唯一修改文件，为 FakeKVSender 和 FakeKVReceiver 新增 abort 与 conclude_state 状态管理。

```
# file: python/sglang/srt/disaggregation/fake/conn.py
# FakeKVSender 和 FakeKVReceiver 新增 abort 支持
class FakeKVSender(BaseKVSender):
    def __init__(self, mgr, bootstrap_addr, bootstrap_room, dest_tp_ranks, pp_rank):
        self.kv_mgr = mgr
        self.has_sent = False
        # 新增: 记录传输终结状态, 用于 poll 短路返回
        self.conclude_state: Optional[KVPoll] = None

    def poll(self) -> KVPoll:
        # 优先返回终结状态 (由 abort 或成功路径设置)
        if self.conclude_state is not None:
            return self.conclude_state
        if not self.has_sent:
            return KVPoll.WaitingForInput
        logger.debug("FakeKVSender poll success")
        self.conclude_state = KVPoll.Success
        return KVPoll.Success

    def send(self, kv_indices, state_indices=None):
        self.has_sent = True
        logger.debug(...)

    def abort(self):
        # 新增: 外部调用 abort 后 poll 将返回 Failed
        self.conclude_state = KVPoll.Failed

class FakeKVReceiver(BaseKVReceiver):
    def __init__(self, mgr, bootstrap_addr, bootstrap_room=None):
        self.bootstrap_done = False
        self.has_sent_metadata = False
        self.require_staging: bool = False
        self.conclude_state: Optional[KVPoll] = None # 新增

    def poll(self) -> KVPoll:
        if self.conclude_state is not None: # 新增: 先行检查终结状态
            return self.conclude_state
        # ... 原有逻辑 ...
        self.conclude_state = KVPoll.Success
        return KVPoll.Success

    def abort(self):
        # 新增: 与 FakeKVSender 行为对称
```

```
self.conclude_state = KVPoll.Failed

# 修改要点: 移除冗余 else 分支, 使用 not self.has_sent
# 之前:
# if self.has_sent is False:
# return KVPoll.WaitingForInput
# else:
# logger.debug(...)
# return KVPoll.Success

# 之后 (见 poll 方法) :
# if not self.has_sent:
# return KVPoll.WaitingForInput
# logger.debug(...)
# self.conclude_state = KVPoll.Success
# return KVPoll.Success
```

评论区精华

Review 中 [gemini-code-assist\[bot\]](#) 指出 `FakeKVSender.poll` 的 `else` 分支是冗余的 (因为 `if` 块已 `return`) , 建议将其改为 `not self.has_sent` 并扁平化, 与 `FakeKVReceiver.poll` 实现风格一致。作者采纳该建议并在同一 PR 中修正。

- 冗余 `else` 分支清理 (style): 作者采纳建议, 在后续 commit 中修正。

风险与影响

- 风险: 变更范围极小, 仅影响 fake KV backend (测试 / warmup 专用) , 不涉及生产路径。主要风险是 `conclude_state` 被设置后无法重置, 可能导致后续 `poll` 调用一直返回同一个状态; 但 fake backend 在本轮测试中通常只使用一次, 实际风险较低。
- 影响: 影响范围仅限于 PD 场景下的 fake KV backend 使用方, 使其能在测试中通过 `abort` 模拟传输失败, 提高测试覆盖的完整性。对用户无直接感知, 对系统无性能影响。
- 风险标记: 状态不可重置, 低影响路径

关联脉络

- PR #25638 [Move module-level helpers out of scheduler.py](#): 同属调度模块重构序列, 但本 PR 与重构无直接关联, 仅因近期大量重构 PR 被列在一起。