

PR #41627 完整报告

vllm-project/vllm

[EC Connector] Non blocking EC Connector lookup

合并时间: 2026-06-02 16:48

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/41627>

执行摘要

- 一句话: 为 EC 连接器引入非阻塞查询机制, 允许调度器延迟等待多模态编码缓存预取请求。
- 推荐动作: 建议精读该 PR 的核心设计: `ensure_cache_available` 接口定义与调度器集成点, 特别是如何通过非阻塞延迟提升吞吐。对于连接器实现者, 应关注后续 PR #42998 以了解完整用法。代码质量和测试覆盖良好, 值得 merge。

功能与动机

当使用 Encoder Cache (EC) 连接器时, 多模态请求可能依赖远程或 CPU 预取的编码缓存, 若缓存尚未就绪 (如正在异步传输), 调度器若仍然尝试调度将导致等待或失败。通过非阻塞延迟机制, 调度器可跳过未就绪请求, 让其他请求继续, 提升整体吞吐和资源利用率。PR body 明确说明: 'When a waiting request has multimodal features whose encoder cache is still being staged (e.g., remote → CPU prefetch in progress), the scheduler now skips the request for the current step and re-queues it in `step_skipped_waiting`, allowing other requests to proceed.'

实现拆解

1. 在 `vllm/distributed/ec_transfer/ec_connector/base.py` 中添加 `ensure_cache_available(self, request, num_computed_tokens)` 方法, 作为可选接口 (默认返回 `True`)。该方法接收请求对象和已计算 token 数, 返回 `True` 表示就绪, `False` 表示仍在传输中。子类可覆盖以触发异步预取并返回状态。
2. 在 `vllm/v1/core/sched/scheduler.py` 的 `schedule` 方法中, 在计算 `num_computed_tokens` 后, 插入检查: 若启用了 `ec_connector` 且请求包含 `mm_features` 且 `ensure_cache_available` 返回 `False`, 则将该请求从当前队列弹出并重新放入 `step_skipped_waiting` (等待后续步骤), 随后 `continue` 处理下一个请求。该检查位于 `KVConnector` 处理后, 确保已结合外部缓存信息计算 `num_computed_tokens`。
3. 在 `tests/v1/core/test_scheduler.py` 中添加两个测试函数:
 - `test_ec_connector_ensure_cache_available_defers_request`: 验证当 `ensure_cache_available` 返回 `False` 时, 延迟请求不会分配任何 KV 缓存或编码缓存, 但后续文本请求可正常调度; 当返回 `True` 且 `has_cache_item` 为 `True` 时, 延迟请求被正常调度。

- `test_ec_connector_pending_prefetch_only_checks_future_mm_features`: 验证仅超过当前 `num_computed_tokens` 的 MM 特征会被传递, 避免对已缓存的旧特征触发不必要的预取。

4. 配套调整: 移除之前可能存在的 `has_pending_prefetch` 抽象方法, 改为非抽象方法以保持向后兼容; 同时将参数从 `mm_hashes` 改为更灵活的 `request` 和 `num_computed_tokens`, 使基类可自行计算需检查的特征。

关键文件:

- `vllm/distributed/ec_transfer/ec_connector/base.py` (模块 连接器; 类别 `source`; 类型 `core-logic`; 符号 `ensure_cache_available`): 定义新的 `ensure_cache_available` 接口, 作为可选方法 (默认返回 `True`), 允许子类实现非阻塞缓存就绪检查。
- `vllm/v1/core/sched/scheduler.py` (模块 调度器; 类别 `source`; 类型 `core-logic`; 符号 `schedule`): 在调度循环中集成延迟检查, 当确保缓存不可用时跳过请求并重新排队。
- `tests/v1/core/test_scheduler.py` (模块 测试; 类别 `test`; 类型 `test-coverage`; 符号 `test_ec_connector_ensure_cache_available_defers_request`, `test_ec_connector_pending_prefetch_only_checks_future_mm_features`): 新增两个测试函数, 覆盖延迟请求、正常调度、以及未来特征过滤。

关键符号: `ensure_cache_available`, `test_ec_connector_ensure_cache_available_defers_request`, `test_ec_connector_pending_prefetch_only_checks_future_mm_features`

关键源码片段

`vllm/v1/core/sched/scheduler.py`

在调度循环中集成延迟检查, 当确保缓存不可用时跳过请求并重新排队。

```
# 检查 EC Connector 是否报告缓存未就绪 (仍在预取中)
if (
    self.ec_connector is not None
    and request.mm_features
    and not self.ec_connector.ensure_cache_available(
        request, num_computed_tokens
    )
):
    # 移除当前请求并放入 step_skipped_waiting, 下次再试
    request_queue.pop_request()
    step_skipped_waiting.prepend_request(request)
    continue
```

`tests/v1/core/test_scheduler.py`

新增两个测试函数, 覆盖延迟请求、正常调度、以及未来特征过滤。

```
# --- Step 1: ensure_cache_available 返回 False → 请求应被延迟 ---
scheduler.ec_connector.ensure_cache_available = Mock(return_value=False)

scheduler.add_request(request_deferred)
scheduler.add_request(request_behind)
```

```

output = scheduler.schedule()

# 确保 ensure_cache_available 被正确调用
scheduler.ec_connector.ensure_cache_available.assert_called_once_with(
    request_deferred, 0
)
# 延迟请求不得被调度
assert request_deferred.request_id not in output.num_scheduled_tokens
# 后续文本请求仍应被调度
assert request_behind.request_id in output.num_scheduled_tokens
assert output.num_scheduled_tokens[request_behind.request_id] == 20

# --- Step 2: 预取完成, 缓存可用 → 请求恢复调度 ---
scheduler.ec_connector.ensure_cache_available = Mock(return_value=True)
scheduler.ec_connector.has_cache_item = Mock(return_value=True)

output = scheduler.schedule()

# 现在延迟请求应被正常调度
assert request_deferred.request_id in output.num_scheduled_tokens
assert output.num_scheduled_tokens[request_deferred.request_id] == NUM_TOKENS

```

评论区精华

- API 命名争议: orozery 建议将 `has_pending_prefetch` 改为 `ensure_cache_available`, 以匹配 `has_cache_item` 命名体系, 且避免提及具体实现细节 (如 CPU prefetch)。
- 方法可选性: orozery 将方法从 `abstractmethod` 改为非抽象, 默认返回 `True`, 避免破坏现有连接器实现。
- 解耦决策: fakeOfan 质疑该 PR 作为独立部分尚未被任何连接器调用, orozery 回应应将核心变更与具体连接器解耦, 核心变更是对调度器的低风险扩展, 而具体连接器实现 (#42998) 会在后续完善。
- 语义澄清: fakeOfan 指出 `_future_mm_hashes` 中可能存在语义混淆, 需区分编码缓存与 KV 缓存; 最终通过传入 `request` 和 `num_computed_tokens` 在基类实现过滤。
- API 命名: `has_pending_prefetch` → `ensure_cache_available` (design): 采纳建议, 改为 `ensure_cache_available`。
- 方法是否作为 `abstractmethod` (design): 改为非抽象方法, 默认返回 `True`。
- 核心变更是否应独立提交 (design): 核心变更独立合入, 具体连接器实现在后续 PR。
- 确保只检查未来 mm features (correctness): 通过 `filter` 实现并添加单元测试验证。
- 语义混淆: 编码缓存 vs KV 缓存 (correctness): 最终通过传入 `request` 和 `num_computed_tokens` 在基类实现过滤, 避免混淆。

风险与影响

- 风险:

1. 回归风险：调度器核心路径新增分支，若 `ec_connector` 实现返回不当值可能影响所有多模态请求。但默认返回 `True` 保证无影响。
 2. 性能风险：每步调度增加一次 `ensure_cache_available` 调用，但方法体内通常为轻量状态检查。
 3. 兼容性：现有连接器未实现该方法，默认行为保持不变；后续连接器可能需覆盖此方法。
 4. 测试覆盖：新增测试覆盖延迟路径和未来特征过滤，但仍需集成测试与真实连接器验证。
- 影响：
 - 用户：无直接影响，需配合具体连接器实现（如 #42998）才能体现。
 - 系统：为调度器引入非阻塞延迟能力，有望提升多模态请求场景的整体吞吐（避免因等待单个请求而阻塞后续请求）。
 - 团队：该 PR 为连接器框架层变更，后续连接器实现者需要理解并可能覆盖 `ensure_cache_available`；评审者需关注调度器逻辑正确性。
 - 风险标记：核心路径变更，依赖后续实现，缺少集成测试

关联脉络

- PR #42998 WIP: Specific EC connector implementation (mentioned in review): 在 review 中被视为该核心变更的后续具体实现，直接利用 `ensure_cache_available` 机制。