

PR #32553 完整报告

vllm-project/vllm

[P/D] Prefill compute optimizations with bi-directional KV cache transfers between P and D nodes

合并时间: 2026-04-30 18:14

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/32553>

执行摘要

- 一句话: 支持 D→P 双向 KV 传输以消除冗余预填充计算
- 推荐动作: 建议精读调度器变更和示例代理设计, 重点关注阈值调优和 HMA 兼容性。此 PR 体现了在现有框架上演进新功能的设计模式: 通过配置门控最小化风险。

功能与动机

避免 Prefill 节点上的冗余计算, 提升预填充吞吐量; 支持多轮对话场景中重用上一轮 Decode 节点的 KV 缓存; 支持缓存驱逐恢复——当 Prefill 节点本地缓存被驱逐时, 仍可从 Decode 节点拉取。PR body 明确指出: "To avoid redundant compute cycles on Prefill nodes and hence improve prefill throughput."

实现拆解

1. 调度器配置扩展: 在 `vllm/distributed/kv_transfer/kv_connector/v1/nixl/scheduler.py` 的 `__init__` 中添加 `kv_recompute_threshold` (默认 64) 和 `is_bidirectional_kv_xfer_enabled` (默认关闭) 两个配置项。
2. 远程 Token 匹配逻辑: 修改 `get_num_new_matched_tokens`, 当请求标记了 `do_remote_decode` 并携带 `remote_block_ids` 时, 计算可以远程拉取的 token 数并与阈值比较, 决定是否拉取。
3. 请求状态更新: 调整 `update_state_after_alloc`, 使双向传输模式下 P 节点请求也能加入 `_reqs_in_batch`; 引入 `_remote_blocks_processed` 标志防止重复拉取。
4. 请求完成时返回块信息: 在 `request_finished` 中增加 D 节点返回 `remote_block_ids` 和 `remote_num_tokens` 的逻辑, 供代理缓存并在下一轮传递给 P 节点; 同时扩展了对 `FINISHED_STOPPED` 状态的支持。
5. 代理示例与测试: 新增 `examples/online_serving/disaggregated_serving/disagg_proxy_multiturn.py` 示例代理, 管理 `conversation_id` 到 KV 参数的映射; 新增 `tests/v1/kv_connector/unit/test_bidirectional_kv_transfer.py` 单元测试套件, 覆盖多轮生命周期、阈值跳过、部分远程覆盖、异常路径等场景。

关键文件:

- `vllm/distributed/kv_transfer/kv_connector/v1/nixl/scheduler.py` (模块 KV 传输调度器; 类别 `source`; 类型 `core-logic`; 符号 `init`, `get_num_new_matched_tokens`,

update_state_after_alloc, request_finished) : 核心逻辑变更: 添加了双向 KV 传输的配置、远程拉取 token 匹配、请求状态更新和完成时参数返回。

- examples/online_serving/disaggregated_serving/disagg_proxy_multiturn.py (模块 代理示例; 类别 source; 类型 dependency-wiring; 符号 CachedKVEntry, ConversationKVCache, get, put) : 新增代理示例, 展示了如何通过 conversation_id 映射来管理多轮对话的 KV 缓存, 供部署参考。
- tests/v1/kv_connector/unit/test_bidirectional_kv_transfer.py (模块 单元测试; 类别 test ; 类型 test-coverage; 符号 _make_p_node_turn2_request, _make_connector_with_fake_worker, _make_p_node_recv_metadata, _do_load_kv) : 新增单元测试套件, 覆盖核心场景和边缘条件, 保证代码质量。

关键符号: init, get_num_new_matched_tokens, update_state_after_alloc, request_finished

关键源码片段

vllm/distributed/kv_transfer/kv_connector/v1/nixl/scheduler.py

核心逻辑变更: 添加了双向 KV 传输的配置、远程拉取 token 匹配、请求状态更新和完成时参数返回。

```
# vllm/distributed/kv_transfer/kv_connector/v1/nixl/scheduler.py
```

```
def __init__(self, vllm_config, ...):
```

```
    # ... 已有代码 ...
```

```
    # Threshold to decide whether to compute kv cache locally
    # or pull from a remote node: minimum number of remote
    # tokens to amortize the xfer latencies
```

```
    self.kv_recompute_threshold: int = int(
        vllm_config.kv_transfer_config.get_from_extra_config(
            "kv_recompute_threshold", 64 # default 64 tokens
        )
    )
```

```
    # Bi-directional KV transfer feature supports KV block
    # transfers from D node to P node
```

```
    self.is_bidirectional_kv_xfer_enabled = (
        vllm_config.kv_transfer_config.get_from_extra_config(
            "bidirectional_kv_xfer", False # opt-in
        )
    )
```

```
def get_num_new_matched_tokens(self, request, num_computed_tokens):
```

```
    # ... 已有 do_remote_prefill / do_remote_decode 逻辑 ...
```

```
    # NEW: P node pulling from D node (bidirectional)
    if (params is not None
```

```

and params.get("do_remote_decode")
and params.get("remote_block_ids")
and all(p in params for p in (
    "remote_engine_id", "remote_request_id",
    "remote_host", "remote_port"))
):
remote_num_tokens = params.get("remote_num_tokens") or 0
count = (min(remote_num_tokens, request.num_prompt_tokens)
        - num_computed_tokens)
if count > 0:
    # Check threshold: skip if too small
    if (self.kv_recompute_threshold > 0
        and count < self.kv_recompute_threshold):
        logger.debug("Skipping remote pull for %s: "
                    "%d remote tokens < threshold %d",
                    request.request_id, count,
                    self.kv_recompute_threshold)
        return 0, False
    return count, True # pull KV from D

return 0, False

```

tests/v1/kv_connector/unit/test_bidirectional_kv_transfer.py

新增单元测试套件，覆盖核心场景和边缘条件，保证代码质量。

```

# tests/v1/kv_connector/unit/test_bidirectional_kv_transfer.py

def _make_p_node_turn2_request(
    request_id, block_size, num_tokens,
    num_remote_blocks=3, remote_num_tokens=None):
    """Create a P-node Turn 2 request with remote_block_ids from D."""
    request = create_request(
        request_id=request_id,
        block_size=block_size,
        num_tokens=num_tokens,
        do_remote_decode=True,
    )
    if remote_num_tokens is None:
        remote_num_tokens = num_remote_blocks * block_size
    request.kv_transfer_params["remote_block_ids"] = [list(range(num_remote_blocks))]
    request.kv_transfer_params["remote_num_tokens"] = remote_num_tokens
    request.kv_transfer_params["remote_engine_id"] = "decode-engine"
    request.kv_transfer_params["remote_request_id"] = f"decode-{request_id}"
    request.kv_transfer_params["remote_host"] = "decode-host"
    request.kv_transfer_params["remote_port"] = 5678
    return request

def _make_connector_with_fake_worker(
    hand_shake_latency=0, cycles_before_done=0, do_handshake=True):

```

```
"""Create a NixlConnector with a FakeNixlConnectorWorker."""
vllm_config = create_vllm_config()
kv_cache_config = make_kv_cache_config(block_size=16, num_blocks=2)
connector = NixlConnector(
    vllm_config, KVConnectorRole.WORKER, kv_cache_config)
connector.connector_worker = FakeNixlConnectorWorker(
    vllm_config, connector.engine_id,
    hand_shake_latency=hand_shake_latency,
    kv_cache_config=kv_cache_config)
# ... handshake omitted for brevity ...
return connector, connector.connector_worker
```

评论区精华

1. 默认阈值讨论: NickLucche 认为 0 不合适, 希望跑 ablation 确定值。最终作者改为 64。
 2. 变量命名: NickLucche 建议 `remote_pull_threshold` 改进, 作者提议 `kv_recompute_threshold`, 被接受。
 3. 逻辑统一: NickLucche 建议将 P 拉取逻辑与已有 `do_remote_prefill` 分支统一, 作者在最终版本中调整了结构。
 4. FINISHED_STOPPED 扩展: 作者解释新状态是为了支持正常完成的请求也返回 KV 参数, 适用于 P→D 传输。
 5. 代理超时与块释放: markmc 指出代理超时可能影响 D 节点块释放时机, 作者解释沿用 P→D 设计。
- 默认阈值设置 (performance): 作者将默认值改为 64, 并添加了日志。
 - 变量命名 (design): 作者提出 `kv_recompute_threshold` 并被接受。
 - 统一代码分支 (design): 作者调整了代码结构, 合并了条件判断。
 - FINISHED_STOPPED 支持 (correctness): 作者解释: 正常完成的请求也需要返回 KV 参数, 适用于 P→D 常规传输。

风险与影响

- 风险:
 1. 阈值敏感: `kv_recompute_threshold` 默认 64 可能不适用于所有模型和网络延迟, 需实际调优。
 2. 块释放延迟: D 节点持有块等待 P 拉取, 若 P 异常或代理超时可能导致 OOM。
 3. HMA 兼容性: 新逻辑在 HMA (混合模型注意力) 模式下的 token 计数可能不准, 已在 review 中被指出并部分修复。
 4. 缺少集成测试: 单元测试覆盖充分, 但缺少端到端集成测试和性能基准来验证实际收益。
- 影响:
 - 用户: 通过 `kv_transfer_config` 可启用双向传输, 提升多轮对话性能; 默认禁用, 无影响。
 - 系统: 修改了调度器核心路径但通过配置开关隔离; 新增代理示例和测试, 降低了集成门槛。

- 团队：新测试套件便于后续回归；设计决策记录在 RFC #32733 中。
- 风险标记：默认阈值需调优，块释放延迟可能导致 OOM, HMA 兼容性未充分测试

关联脉络

- 暂无明显关联 PR