

# PR #42187 完整报告

vllm-project/vllm

[ModelRunnerV2] Avoid pipeline parallel bubbles

合并时间: 2026-06-03 05:02

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/42187>

## 执行摘要

- 一句话: 重构 PP 调度避免气泡, 吞吐最高提升 3 倍
- 推荐动作: 此 PR 是 ModelRunnerV2 的重要性能里程碑, 强烈建议精读。重点关注 PPHandler 的延迟消费与专用通信器设计, 这是一种解耦同步通信的通用模式。同时关注 CPU mirror 与 GPU 状态一致性维护策略。建议在后续 PR 中增加更多极端场景测试 (如高并发出错恢复)。

## 功能与动机

原有 ModelRunnerV2 的 PP 实现存在调度气泡: decode tokens 未按 `pp_size` 步调度, 广播与 p2p 通信在同一流水线上串行化, 导致 GPU 空闲等待。PR body 中给出了详细的性能数据对比, 并引用了动画展示数据流。此外, 混合负载下经常出现的 IMA 内部内存错误也被此 PR 修复。

## 实现拆解

1. 引入 PPHandler: 在 `vllm/v1/worker/gpu/pp_utils.py` 中新建 PPHandler 类。它在非末级 rank 上维护一个 FIFO 队列, 按 `pp_size` 步延迟消费广播结果。使用专用 NCCL 组和 CUDA 流执行广播, 与主流的 p2p 隐藏状态通信并行。
2. 修改 `model_runner.py`: 新增 `update_pp_decode_requests()` 方法, 非末级 rank 每步从 PPHandler 获取先前步骤的采样输出并更新状态; 将原 `postprocess` 拆分为 `postprocess_sampled`, 支持过滤的 `idx_mapping` 和可选的 `query_start_loc`; 在 `add_requests` 中传递 `max_tokens` 以便 PP 控制流使用; 在 `_remove_request` 中通知 PPHandler 释放索引。
3. 调整 InputBatch 数据契约: 在 `input_batch.py` 中添加 `num_computed_tokens_np`、`prefill_len_np`、`num_computed_prefill_tokens_np`、`max_seq_len_np` 等 CPU 镜像字段, 用于 PP 控制流决策。修改 `_post_update_kernel` 支持负索引跳过 (filtered rows), 并将 pool kernel 分离为 `_post_update_num_computed_tokens_kernel`。
4. 修改调度器: 在 `async_scheduler.py` 和 `scheduler.py` 中增加 PP throttle 条件, 当 PPHandler 活跃时限制每一步调度的 token/req 数量, 确保流水线节奏对齐。在 `vllm/v1/core/sched/async_scheduler.py` 中新增 `self.pp_throttle` 布尔标志。
5. 配套改动: 在 `parallel_state.py` 中新增 `make_sibling_device_group` 方法以创建专用广播组; 在 `buffer_utils.py` 中新增 `set_default_max_concurrency` 控制 UVA 缓冲区池大小;

在 `states.py` 中修改 `remove_request` 返回释放的索引。

6. 新增测试: `tests/v1/distributed/test_pp_dp_v2.py` 覆盖 low/mid concurrency 和 abort mid-decode 场景, 使用 DP=2, PP=2 验证组合稳定性。

关键文件:

- `vllm/v1/worker/gpu/pp_utils.py` (模块 PP 工具; 类别 source; 类型 core-logic; 符号 PendingRecv, compute\_need\_sampled\_mask, PPHandler, init) : 核心实现: 引入 PPHandler 类、PendingRecv 数据结构和 compute\_need\_sampled\_mask, 定义了整个延迟广播机制。
- `vllm/v1/worker/gpu/model_runner.py` (模块 模型运行器; 类别 source; 类型 core-logic; 符号 update\_pp\_decode\_requests, postprocess\_sampled, postprocess\_pool, postprocess\_num\_computed\_tokens) : 主入口: 集成 PPHandler, 修改 execute\_model 和 postprocess 流程, 添加 update\_pp\_decode\_requests 方法。
- `vllm/v1/worker/gpu/input_batch.py` (模块 输入批处理; 类别 source; 类型 core-logic; 符号 \_post\_update\_pool\_kernel, \_post\_update\_num\_computed\_tokens\_kernel, post\_update\_pool, post\_update\_num\_computed\_tokens) : 数据契约: 增加 CPU 镜像字段支持 PP 控制流, 修改 post\_update kernel 支持过滤和分离。
- `tests/v1/distributed/test_pp_dp_v2.py` (模块 集成测试; 类别 test; 类型 test-coverage; 符号 \_gpu\_skip\_reason, \_engine\_args, \_generate, test\_pp\_dp\_v2\_low\_concurrency) : 新增集成测试: 覆盖低中高并发和请求取消场景, 验证 PP+DP 组合功能。
- `vllm/v1/worker/gpu/states.py` (模块 请求状态; 类别 source; 类型 core-logic; 符号 remove\_request, is\_prefilling) : 修改 remove\_request 返回释放索引, 供 PPHandler 更新 generation 计数器。
- `vllm/distributed/parallel_state.py` (模块 分布式状态; 类别 source; 类型 core-logic; 符号 make\_sibling\_device\_group) : 新增 make\_sibling\_device\_group 以创建专用 NCCL 广播组, 避免与 p2p 通信串行。
- `vllm/v1/worker/gpu/buffer_utils.py` (模块 缓冲区工具; 类别 source; 类型 core-logic; 符号 set\_default\_max\_concurrency) : 新增 set\_default\_max\_concurrency 用于 UVA 缓冲区池大小计算, 需在 PPHandler 构造前调用。
- `vllm/v1/core/sched/async_scheduler.py` (模块 异步调度器; 类别 source; 类型 core-logic) : 添加 PP throttle 条件, 限制步骤并发以匹配流水线节奏。

关键符号: PPHandler.init, PPHandler.on\_req\_idx\_freed, PPHandler.get\_prev\_sampled\_outputs, PPHandler.receive, compute\_need\_sampled\_mask, GPUModelRunner.update\_pp\_decode\_requests, GPUModelRunner.postprocess\_sampled, GPUModelRunner.\_remove\_request, RequestState.remove\_request, parallel\_state.make\_sibling\_device\_group, set\_default\_max\_concurrency

关键源码片段

`vllm/v1/worker/gpu/model_runner.py`

主入口：集成 PPHandler，修改 execute\_model 和 postprocess 流程，添加 update\_pp\_decode\_requests 方法。

```
def update_pp_decode_requests(self):
    # For non-last PP ranks, update decode requests with sampler output from
    # the prior step in which they were scheduled (pp_size steps ago).
    if self.pp_handler is None:
        return
    outputs = self.pp_handler.get_prev_sampled_outputs()
    if outputs is None:
        return
    sampled_tokens = outputs["sampled_tokens"]
    num_sampled = outputs["num_sampled"]
    num_rejected = outputs["num_rejected"]
    idx_mapping = outputs["idx_mapping"]

    # This postprocess call uses a filtered idx_mapping; query_start_loc is
    # omitted because the per-request computed token increment was already
    # applied in immediate postprocess; here we only adjust for rejected tokens.
    self.postprocess_sampled(
        idx_mapping=idx_mapping,
        sampled_tokens=sampled_tokens,
        num_sampled=num_sampled,
        num_rejected=num_rejected,
        query_start_loc=None,
    )

def postprocess_sampled(
    self,
    idx_mapping: torch.Tensor, # [num_reqs]; May include -1 for masked entries
    sampled_tokens: torch.Tensor,
    num_sampled: torch.Tensor,
    num_rejected: torch.Tensor,
    query_start_loc: torch.Tensor | None = None,
) -> None:
    # Write output tokens to all_token_ids
    self.req_states.write_sampled_tokens(
        idx_mapping, sampled_tokens, self.max_num_reqs
    )

    # Update staging for num_computed_tokens (delta for rejected only if query_start_loc is
    # None)
    post_update(
        idx_mapping,
        self.req_states.num_computed_tokens.gpu(),
        self.req_states.num_computed_prefill_tokens.gpu(),
        self.req_states.prefill_len.gpu(),
        num_sampled,
        num_rejected,
```

```

        query_start_loc,
        self.req_states.all_token_ids.gpu(),
        self.max_num_reqs,
    )

    # Apply any model-specific state update (e.g. Mamba)
    self.model_state.postprocess_state(idx_mapping, num_sampled)

```

## 评论区精华

主要讨论围绕三个方面：

### 1. 正确性问题 (via gemini-code-assist)

- `max_tokens` 可能为 `None` 导致 `add_request` 类型错误 (`sampling_params.max_tokens` 可为 `None`) 。
- CPU mirror `num_computed_tokens_np` 在 speculative decoding 拒词后未纠正，与 GPU 状态漂移，可能使 PPHandler 的 `compute_need_sampled_mask` 过早停止广播。
- `not_finishing` 条件 `np.maximum(old_computed, prefill_len) + 1 < max_seq_len` 存在 off-by-one，导致每请求最后一个 token 不被广播。以上问题在后续迭代中已修复（或确认无影响）。

### 2. 接口设计权衡 (via WoosukKwon)

- `postprocess_state` 为何从 `input_batch` 改为 `idx_mapping`? njhill 解释：延迟后 `input_batch` 可能已变化，只需传入过滤后的索引。
- dummy run 下 `intermediate_tensors` 为何需要特殊切片? njhill 指出：避免自复制冲突。
- Mamba 为何需要特殊处理? njhill 说明其 `postprocess_state` hook 使用 `num_sampled` 更新。

### 3. 辅助讨论

- `max_tokens` 默认值 0 vs 1 的合理性。
- `query_start_loc` 何时为 `None`。

最终 WoosukKwon 批准 (LGTM) 。

- `max_tokens` 可能为 `None` 导致类型错误 (correctness): 后续提交中已处理或确认当前使用 0 作为 N/A 默认值，无实际影响。
- CPU mirror 漂移导致 PP 控制流错误 (correctness): 该问题在 PR 迭代中通过调整 `postprocess` 逻辑（立即 + 延迟两阶段更新）解决。
- `not_finishing` 条件 off-by-one 导致最后 token 不广播 (correctness): 修正提交中已调整或确认实际逻辑经过验证无遗漏。
- `postprocess_state` 接口为何从 `input_batch` 改为 `idx_mapping` (design): WoosukKwon 接受解释。
- dummy run 下 `intermediate_tensors` 特殊处理 (design): WoosukKwon 理解。

## 风险与影响

- 风险:

1. 调度节奏失误: PPHandler 队列长度依赖 `pp_size`, 若计算错误或与调度器不同步, 可能导致广播乱序或死锁 (`pp_utils.py`) 。
  2. CPU mirror 一致性问题: `num_computed_tokens_np` 作为控制流依据, 若更新逻辑缺失 (如拒词未减), 会错误判断是否需要广播 (`input_batch.py`, `pp_utils.py`) 。
  3. 后端兼容性: `make_sibling_device_group` 依赖 NCCL 能力, 可能在某些分布式后端 (如 Gloo) 上失败或退化 (`parallel_state.py`) 。
  4. 残余 -1 索引处理: `post_update kernel` 和 `Mamba scatter kernel` 需要正确跳过 -1 行, 遗漏会导致 GPU 访存错误。
  5. 整体复杂度: PPHandler 与调度器的交替状态机增加了调试难度, 非 PP 路径也需要验证未受影响。
    - 影响: 用户: 使用 Pipeline Parallelism (特别是 large model 如 MiniMax-M2.7) 的用户获得显著吞吐提升 (最高 3.17x) 和 TTFT 降低 (最高 3.15x)。无需更改模型名称或 API, 仅需启用 `VLLM_USE_V2_MODEL_RUNNER=1` 环境变量。
    - 系统: 引入额外 NCCL 通信组和 CUDA 流, 增加少量显存和通信开销, 但收益远大于开销。PP 模式下调度器增加一个 `throttle` 条件, 可能略微增加调度延迟, 但整体 throughput 提升。团队: 维护成本上升, PPHandler 和专用通信组的实现需要与调度器和模型 runner 的其他部分保持同步。
- 风险标记: 调度逻辑核心变更, CPU/GPU 状态一致性, 专用 NCCL 通信组, 代码复杂度增加

## 关联脉络

- PR #44277 [ModelRunnerV2] Some other PP fix (推断): PR body 声明依赖此 PR 或类似, 确保基础功能就绪。