

PR #25862 完整报告

sgl-project/sglang

Address overlap future token map by request-pool index

合并时间: 2026-05-21 07:34

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/25862>

执行摘要

- 一句话: 将 FutureMap 索引从环形缓冲区切换为请求池索引
- 推荐动作: 此 PR 是消除一类重叠调度竞态的重要修复, 设计思路值得借鉴——通过将临时槽位映射为具有语义的唯一标识符来避免错误。建议阅读 `overlap_utils.py` 的变更, 理解如何通过统一索引根除整类问题。虽然缺少测试, 但 CI 已全部通过, 可以合入。合入后应关注重叠调度场景的稳定性。

功能与动机

原环形分配器在槽位重用时会引发分块预填充膨胀、索引环绕别名、脏读以及跨流竞态。PR body 明确指出 'The ring allocator slot id is a temporary counter — slot K denotes different requests at different times, driving sizing miscalculations, wrap-around aliasing, slot-reuse staleness, and a cross-stream indices-tensor lifetime race.' 必须将槽位与请求生命周期绑定。

实现拆解

实现分步拆解

1. 重构 FutureMap 类(`overlap_utils.py`): 构造函数改为接收 `ReqToTokenPool`, 用 `req_to_token_pool.req_to_token.shape[0]` 设定缓冲区大小, 移除 `future_ct`、`future_limit`、`future_buffer_len`、`is_empty_slice` 等环形分配器状态。
2. 移除 `alloc_future_indices` 方法: 不再由 FutureMap 分配槽位, 改为调用者直接使用 `batch.req_pool_indices` 构造 `FutureIndices`。 `FutureIndices` 数据类简化, 删除 `interval` 字段。
3. 更新调度器接线(`scheduler.py`): `init_overlap` 中创建 FutureMap 时传递 `self.req_to_token_pool`; `run_batch` 中直接 `FutureIndices(indices=batch.req_pool_indices)` 替代原来的 `self.future_map.alloc_future_indices(bs)`。
4. 调整接口签名(`spec_info.py`): `create_future_map` 方法简化参数列表, 只接受 `device` 和 `req_to_token_pool`。
5. 统一预构建代码路径(`decode_schedule_batch_mixin.py`): `process_prebuilt` 中同样使用 `FutureIndices(indices=self.req_pool_indices)` 替代 `future_map.alloc_future_indices(...)`。
6. 增强类型安全性: `store_to_map` 中显式将 `next_token_ids` 转换为 `torch.int64` 并展平为 1D, 避免形状或类型不匹配。

7. 保留对推测解码的支持：推测解码相关的缓冲区（topk_p_buf 等）同样基于 req_pool_size，索引方式一致。

注意：本次改动未包含专门的新测试，依赖现有 CI 覆盖。

关键文件：

- python/sglang/srt/managers/overlap_utils.py（模块 调度器；类别 source；类型 core-logic；符号 alloc_future_indices, is_empty_slice）：核心逻辑变更，重构 FutureMap 索引方式为请求池索引，移除环形分配器状态
- python/sglang/srt/managers/scheduler.py（模块 调度器；类别 source；类型 dependency-wiring）：适配新 FutureMap 创建方式与 run_batch 中使用 FutureIndices
- python/sglang/srt/speculative/spec_info.py（模块 推测解码；类别 source；类型 core-logic）：create_future_map 接口简化，移除计算 buffer 大小的参数
- python/sglang/srt/disaggregation/decode_schedule_batch_mixin.py（模块 分离推理；类别 source；类型 dependency-wiring）：预构建分支使用新 FutureIndices 构造方式

关键符号：FutureMap.init, FutureMap.store_to_map, FutureMap.resolve_future, SpeculativeAlgorithm.create_future_map, Scheduler.init_overlap, Scheduler.run_batch

关键源码片段

python/sglang/srt/managers/overlap_utils.py

核心逻辑变更，重构 FutureMap 索引方式为请求池索引，移除环形分配器状态

```
# python/sglang/srt/managers/overlap_utils.py
from typing import TYPE_CHECKING
import torch
from sglang.srt.speculative.spec_utils import spec_need_hidden_states

if TYPE_CHECKING:
    from sglang.srt.mem_cache.memory_pool import ReqToTokenPool
    from sglang.srt.speculative.eagle_info import EagleDraftInput

@dataclass
class FutureIndices:
    indices: torch.Tensor # 直接使用 req_pool_indices，不再需要 interval

class FutureMap:
    def __init__(
        self,
        device: torch.device,
        spec_algo: SpeculativeAlgorithm,
        req_to_token_pool: ReqToTokenPool,
    ):
        self.device = device
        self.spec_algo = spec_algo
        # 缓冲区大小等于请求池的行数，slot 0 用作 padding row（与 KV cache pool 共享）
        self.req_pool_size = req_to_token_pool.req_to_token.shape[0]
```

```

if self.spec_algo.is_none():
    self.buf_initialized = True
    self.token_ids_buf = torch.empty(
        (self.req_pool_size,), dtype=torch.int64, device=self.device,
    )
else:
    self.buf_initialized = False

def store_to_map(
    self, future_indices: FutureIndices, batch_result: GenerationBatchResult
):
    if self.spec_algo.is_none():
        indices = future_indices.indices
        if indices.shape[0] == 0:
            # DP attention idle rank: 无需存储
            return
        # 类型安全: 展平并转 int64
        self.token_ids_buf[indices] = batch_result.next_token_ids.view(-1).to(torch.int64)
    else:
        # 推测解码分支类似, 索引方式一致
        ...

```

python/sglang/srt/managers/scheduler.py

适配新 FutureMap 创建方式与 run_batch 中使用 FutureIndices

```

# python/sglang/srt/managers/scheduler.py
from sglang.srt.managers.overlap_utils import FutureIndices

def init_overlap(self):
    ...
    self.future_map = self.spec_algorithm.create_future_map(
        self.device,
        self.req_to_token_pool, # 之前传递 max_running_requests/chunked_prefill_size/context_len
    )
    ...

def run_batch(self, batch, ...):
    ...
    with self._overlap_forward_isolation(batch):
        # 旧方式: future_indices = self.future_map.alloc_future_indices(len(batch.seq_lens))
        future_indices = FutureIndices(indices=batch.req_pool_indices)
        with self.forward_stream_ctx:
            self.forward_stream.wait_stream(self.schedule_stream)
    ...

```

评论区精华

仅有一条来自 `gemini-code-assist[bot]` 的评论, 关注 `store_to_map` 中 `next_token_ids` 的赋值安全性: 建议显式 `view(-1).to(torch.int64)` 以避免形状 (如 `(bs,1)`) 或类型 (如 `int32`) 不

符。作者后续提交中已采纳该建议（见提交 dbb530a2、2105d23）。

- store_to_map 中 next_token_ids 的类型安全 (correctness): 作者在后续提交中采纳建议, 添加了类型转换 (dbb530a2, 2105d23)。

风险与影响

- 风险:

1. 核心路径变更: overlap_utils.py 和 scheduler.py 是调度核心, 任何错误会导致生成异常或崩溃。环形分配器的移除改变了错误模式, 新的 req_pool_idx 索引依赖请求池的正确维护。
2. 缺少测试覆盖: 没有新增测试文件, 只有现有 CI 验证。可能遗漏边界情况, 如 DP attention 空闲 rank 的空 indices 处理、推测解码与重叠调度结合的场景。
3. 向后兼容性: 接口变化 (create_future_map 参数减少) 影响所有调用方, 但已同步更新。
4. 性能影响: 移除环形分配器后, 缓冲区分配更简单, Advanced indexing 可能比之前的 slice 稍慢, 但差异极微。- 影响: 直接影响所有启用重叠调度 (--enable-overlap) 的功能, 包括非推测和推测解码 (EAGLE、EAGLE3 等) 以及分块预填充。修复了多个潜在的静默错误, 提高了正确性和稳定性。对用户而言, 重叠调度相关的场景会变得可靠。团队需要确保没有遗漏其他使用 FutureMap 的地方。- 风险标记: 核心路径变更, 缺少测试覆盖, 环形分配器移除

关联脉络

- PR #25819 disagg prebuilt: drop dead prepare_for_extend shift: 同时修改了 decode_schedule_batch_mixin.py, 与本 PR 在 process_prebuilt 方法中的 future_indices 构造有直接交互。