

PR #25516 完整报告

sgl-project/sclang

refactor: remove ModelWorkerBatch indirection

合并时间: 2026-05-18 09:05

原文链接: <http://prhub.com.cn/sgl-project/sclang/pull/25516>

执行摘要

- 一句话: 移除 ModelWorkerBatch 中间层, 简化批量数据流
- 推荐动作: 值得精读, 尤其关注 `_overlap_forward_isolation` 的上下文管理器设计、一次性覆盖模式以及跨流张量保活策略。可作为架构重构的参考案例。

功能与动机

消除长期存在的冗余中间层 ModelWorkerBatch, 简化数据流和降低维护成本。参考 `schedule_batch.py` 中原有的 TODO 注释: "ModelWorkerBatch seems a bit redundant and we consider removing it in the future."

实现拆解

1. 移除 ModelWorkerBatch 类: 删除 `schedule_batch.py` 中的 ModelWorkerBatch 数据类和 `get_model_worker_batch` 方法, 清理所有引用。
2. 迁移一次性覆盖字段: 将 `seq_lens_cpu_cache`、`capture_hidden_mode`、`return_hidden_states_before_norm` 三个字段移到 ScheduleBatch, 由 `ForwardBatch.init_new` 消费后自动重置。
3. 引入快照恢复机制: 在 `scheduler.py` 中添加 `_overlap_forward_isolation` 上下文管理器, 在 Spec V2 重叠前向期间对 ScheduleBatch 进行完整快照, 并在完成后恢复, 避免 V2 的原地修改泄漏到下一次调度。同时替换 `sampling_info` 为前向专用副本, 防止多次 `init_new` 重复累积惩罚。
4. 添加流记录防御: 在 `spec_utils.py` 中新增 `record_stream_for_v2_verify` 和 `record_stream_each`, 并在 Spec V2 工作器中调用, 确保跨流张量 (如 `input_ids`、`out_cache_loc`) 在验证阶段不被释放。
5. 全面适配调用方: 修改所有 Spec 工作器 (`eagle_worker_v2.py`、`multi_layer_eagle_worker_v2.py`、`eagle_worker.py`、`multi_layer_eagle_worker.py`、`dflash_worker.py`) 和硬件后端 (`mlx/tp_worker.py`) 的方法签名, 将 ModelWorkerBatch 替换为 ScheduleBatch。
6. 调整工具函数: `overlap_utils.py` 中的 `resolve_future` 和 `tp_worker.py` 中的 `forward_batch_embedding` 等适配新参数类型。

关键文件:

- `python/sglang/srt/managers/scheduler.py` (模块 调度器; 类别 source; 类型 core-logic; 符号 `record_batch_in_overlap`, `_overlap_forward_isolation`) : 核心调度器, 实现 batch 快照恢复和跨流引用的 ring buffer 机制
- `python/sglang/srt/managers/schedule_batch.py` (模块 批处理; 类别 source; 类型 core-logic; 符号 `get_model_worker_batch`, `ModelWorkerBatch`) : 删除 `ModelWorkerBatch` 数据结构和转换方法, 新增一次性覆盖字段
- `python/sglang/srt/model_executor/forward_batch_info.py` (模块 前向批信息; 类别 source; 类型 data-contract; 符号 `_init_ngram_embedding_info`, `_compute_mrope_positions`) : `ForwardBatch.init_new` 直接消费 `ScheduleBatch`, 消费一次性覆盖字段
- `python/sglang/srt/speculative/eagle_worker_v2.py` (模块 推测解码; 类别 source; 类型 core-logic; 符号 `draft`, `forward_batch_generation`, `verify`) : Spec V2 工作器, 方法签名从 `ModelWorkerBatch` 改为 `ScheduleBatch`
- `python/sglang/srt/speculative/multi_layer_eagle_worker_v2.py` (模块 推测解码; 类别 source; 类型 core-logic; 符号 `draft`, `forward_batch_generation`) : 多层 Eagle V2 工作器, 适配 `ScheduleBatch` 并新增 `capture_hidden_mode` 传递
- `python/sglang/srt/hardware_backend/mlx/tp_worker.py` (模块 MLX 硬件; 类别 source; 类型 core-logic; 符号 `async_forward_batch_generation_mlx`, `_async_extend_batch`) : MLX 硬件后端, 适配 `ScheduleBatch` 签名变动
- `python/sglang/srt/speculative/spec_utils.py` (模块 推测工具; 类别 source; 类型 core-logic; 符号 `record_stream_each`, `record_stream_for_v2_verify`) : 新增 `record_stream_for_v2_verify` 和 `record_stream_each` 工具函数
- `python/sglang/srt/managers/tp_worker.py` (模块 工作进程; 类别 source; 类型 core-logic; 符号 `forward_batch_embedding`) : `TpModelWorker` 核心方法适配 `ScheduleBatch` 签名
- `python/sglang/srt/managers/overlap_utils.py` (模块 重叠执行; 类别 source; 类型 core-logic; 符号 `resolve_future`) : 调整 `resolve_future` 等工具函数参数类型
- `python/sglang/srt/speculative/dflash_worker.py` (模块 推测解码; 类别 source; 类型 dependency-wiring) : 适配 `ScheduleBatch` 签名变化
- `python/sglang/srt/speculative/multi_layer_eagle_worker.py` (模块 推测解码; 类别 source; 类型 core-logic) : 多层 Eagle V1 工作器适配 `ScheduleBatch`
- `python/sglang/srt/speculative/eagle_worker.py` (模块 推测解码; 类别 source; 类型 core-logic) : Eagle V1 工作器适配 `ScheduleBatch`

关键符号: `_overlap_forward_isolation`, `record_batch_in_overlap`, `ForwardBatch.init_new`, `record_stream_for_v2_verify`, `record_stream_each`, `draft`, `verify`, `forward_batch_generation`

关键源码片段

[python/sglang/srt/managers/scheduler.py](#)

核心调度器, 实现 batch 快照恢复和跨流引用的 ring buffer 机制

```
# python/sglang/srt/managers/scheduler.py
```

```
@contextmanager
```

```
def _overlap_forward_isolation(self, batch: ScheduleBatch):
```

```
    """
```

```
    使 ScheduleBatch 在一次 overlap forward 中具有事务性：
```

1. 快照 V2 字段，以便 forward 后恢复。
2. 替换 sampling_info 为前向专用副本，避免多次 init_new 重复累积惩罚。
3. 将 (batch, snapshot) 固定到 batch_record_buf 中 2 个迭代周期，
 确保 GPU 张量在 forward stream 完成前不被释放。

```
    """
```

```
# 1. 快照：仅对 spec V2 完整保存所有 dataclass 字段
```

```
snapshot_v2_full = batch.is_spec_v2
```

```
sched_snapshot = (
```

```
    {f.name: getattr(batch, f.name) for f in dataclasses.fields(batch)}
```

```
    if snapshot_v2_full
```

```
    else None
```

```
)
```

```
sched_sampling_info = batch.sampling_info
```

```
# 2. 替换 sampling_info 为前向副本 (orchestrator=None, 共享已累计惩罚缓冲区)
```

```
if sched_sampling_info is not None:
```

```
    batch.sampling_info = sched_sampling_info.copy_for_forward()
```

```
# 3. 将 (batch, snapshot) 固定到 ring buffer, 确保张量存活
```

```
# 注意：必须在 sampling_info 替换之后执行，以固定前向副本
```

```
self.record_batch_in_overlap(batch)
```

```
try:
```

```
    yield
```

```
finally:
```

```
    # 恢复快照
```

```
    if sched_snapshot:
```

```
        for k, v in sched_snapshot.items():
```

```
            setattr(batch, k, v)
```

```
    # 恢复 sampling_info
```

```
    batch.sampling_info = sched_sampling_info
```

python/sglang/srt/managers/schedule_batch.py

删除 ModelWorkerBatch 数据结构和转换方法，新增一次性覆盖字段

```
# python/sglang/srt/managers/schedule_batch.py
```

```
# 数据流注释更新：
```

```
# ScheduleBatch -> ForwardBatch
```

```
# ForwardBatch 由 ScheduleBatch 通过 ForwardBatch.init_new 直接构造。
```

```
@dataclass
```

```
class ScheduleBatch(...):
```

```
# ... 原有字段 ...

# 全新: 一次性前向覆盖, init_new 消费后自动重置为默认值
seq_lens_cpu_cache: torch.Tensor = None
capture_hidden_mode: Optional[CaptureHiddenMode] = None
return_hidden_states_before_norm: bool = False
```

```
# ... 其他字段保持不变 ...
```

python/sclang/srt/model_executor/forward_batch_info.py

ForwardBatch.init_new 直接消费 ScheduleBatch, 消费一次性覆盖字段

```
# python/sclang/srt/model_executor/forward_batch_info.py

@classmethod
def init_new(
    cls,
    batch: ScheduleBatch, # 现在直接接收 ScheduleBatch
    model_runner: ModelRunner,
):
    # 消费一次性覆盖字段并重置
    capture_hidden_mode = batch.capture_hidden_mode
    batch.capture_hidden_mode = None
    seq_lens_cpu_cache = batch.seq_lens_cpu_cache
    batch.seq_lens_cpu_cache = None
    return_hidden_states_before_norm = batch.return_hidden_states_before_norm
    batch.return_hidden_states_before_norm = False

    # 若未覆盖, 则从 batch 的 spec_info 等推导默认值
    if capture_hidden_mode is None:
        if batch.return_hidden_states:
            capture_hidden_mode = CaptureHiddenMode.FULL
        elif batch.spec_info is not None:
            capture_hidden_mode = getattr(
                batch.spec_info, "capture_hidden_mode", CaptureHiddenMode.NULL
            )
        else:
            capture_hidden_mode = CaptureHiddenMode.NULL

    # ... 后续构造 ForwardBatch 的逻辑 ...
```

评论区精华

1. 调度器快照的完整性: Review 指出 `dataclasses.fields(batch)` 只能捕获定义为 `dataclass` 字段的属性, 可能遗漏动态添加的临界张量。作者后续通过 `attr_snapshot` 使用 `dataclasses.fields` 明确只覆盖声明字段, 未处理动态属性 (视为已知限制)。
2. `record_stream` 冗余: Review 指出 `eagle_worker_v2.py` 和 `multi_layer_eagle_worker_v2.py` 中手动的 `record_stream` 调用可能冗余, 因为

`_record_sb_tensors_on_stream` 已覆盖相同字段。建议将调用移动到 `rebind` 之后以集中处理。作者后续提交将相关逻辑提取到 `spec_utils.py`, 但未删除手动调用。

- 调度器快照可能遗漏动态属性 (`correctness`): 作者未回应, 但该设计为已知限制, 适用于当前所有字段均为 `dataclass` 字段的假设。
- `eagle_worker_v2` 中 `record_stream` 冗余 (`performance`): 作者后续提交将辅助函数移至 `spec_utils`, 但未删除手动调用, 保留了双重记录。
- `multi_layer_eagle_worker_v2` 中 `record_stream` 冗余 (`performance`): 同上, 最终保留双重调用以防 `rebind` 后张量不同。

风险与影响

- 风险:
 1. 快照遗漏风险: `_overlap_forward_isolation` 仅快照 `dataclass` 字段, 若未来因需求在运行时添加非字段属性, 可能导致 GPU 张量提前释放。
 2. `record_stream` 双重记录: 手动的 `record_stream` 与 `_record_sb_tensors_on_stream` 可能重复记录, 带来微小性能开销但无功能影响。
 3. 状态管理复杂度: 快照恢复和一次性覆盖机制增加了 `ScheduleBatch` 的隐式状态转换, 可能引入难以调试的临时错误。- 影响: 用户侧: 无直接影响, 功能等价。系统侧: 减少一次数据拷贝和转换, 轻微提升性能; 降低后续开发的心智负担。团队侧: 消除长期 TODO, 简化代码维护; 新的状态管理模式需要团队成员理解并遵循一次性覆盖契约。
- 风险标记: 快照可能遗漏动态属性, `record_stream` 双重记录带来微小开销, 新增状态管理增加复杂性

关联脉络

- PR #22822 [Refactor] Refactor DeepEP dispatcher: 同为重构 MoE 相关分布式组件, 减少间接层, 思路相似
- PR #23760 [MoE] Unify DeepEPMoE+MoriEPMoE through AITER MoeRunner pre/post-permute: 消除冗余 MoE 实现, 与本 PR 消除 `ModelWorkerBatch` 的简化目标一致