

PR #25983 完整报告

sgl-project/sglang

feat(model_runner): remove pool/backend refs from ForwardBatch via ForwardContext

合并时间: 2026-05-22 05:01

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/25983>

执行摘要

- 一句话: 引入 ForwardContext 分离 ForwardBatch 职责, 移除运行时引用
- 推荐动作: 此 PR 是重要的架构解耦重构, 值得所有涉及 attention 或 model runner 开发的工程师精读。关键设计决策包括:
 - 冻结数据类的选择 (防止意外修改, 鼓励 `dataclasses.replace`) 。
 - 注意力后端在 `__init__` 中缓存 pool 引用 vs. 实时通过 `get_forward_context()` 派生。
 - `hisparse_coordinator` 通过后端 `@property` 延迟读取 `model_runner`, 避免循环依赖。
 - 多 runner 场景下 ForwardContext 的 save/restore 策略。建议重点关注 `forward_context.py` 的设计和 `_forward_raw` 中的上下文发布模式。

功能与动机

PR body 指出: ForwardBatch 是 per-batch input data 的数据类, 却携带了四个运行时引用 (`req_to_token_pool`、`token_to_kv_pool`、`attn_backend`、`hisparse_coordinator`), 它们与批次数据无关。将它们分离后, ForwardBatch 可专注于批次状态, 注意力后端只负责自己的 pool 引用, PDmux/TBO/spec 也能统一通过 context manager 切换 active backend。

实现拆解

1. 新增 ForwardContext 核心模块: 创建 `python/sglang/srt/model_executor/forward_context.py`, 定义 ForwardContext 冻结数据类 (当前仅含 `attn_backend` 字段) 以及 `forward_context()` 上下文管理器。同时提供 `set_forward_context`、`get_forward_context`、`get_attn_backend`、`get_token_to_kv_pool`、`get_req_to_token_pool` 等全局函数。
2. ModelRunner 入口发布 ForwardContext: 在 `ModelRunner._forward_raw` 中, 每次前向开始时发布 `ForwardContext(attn_backend=self.attn_backend)`, 结束时恢复。PDmux 场景下, `forward_decode` 使用 `ForwardContext(self.decode_attn_backend)` 嵌套上下文以覆盖活跃后端; TBO 和 speculative draft 类似使用 `dataclasses.replace` 或新建 ForwardContext 实现覆盖。
3. Pattern A: 注意力后端缓存池引用: 约 21 个注意力后端 (FlashInfer、FlashAttention、Triton、DSV4/DSA、Hardware backends 等) 在 `__init__` 中从 `model_runner` 捕获 `req_to_token_pool` 和 `token_to_kv_pool`, 体内部从 `forward_batch.xxx` 改为 `self.xxx`。DSV4 系列的 `Compressor` 和 `C4Indexer` 显式接收 `attn_backend` 参数, 替代从 `forward_batch` 读取。

4. 删除 ForwardBatch 字段：从 ForwardBatch dataclass 和 init_new 中移除 req_to_token_pool、token_to_kv_pool、attn_backend、hisparse_coordinator。同时清理所有手动构造 ForwardBatch(...) 的调用点（CUDA graph runner、piecewise CUDA graph runner、CPU graph runner、two_batch_overlap、speculative 路径等），不再传入这些字段。
5. 配套修正与 bug 修复：TBO 的 filter_batch 不再拷贝已删除字段；修正 BreakableCudaGraphRunner._warmup 缺少 forward_context 包装导致的崩溃；修复 BaseSWAKVPool 缺少 invalidate_loc_cache 导致的 AttributeError；压缩器 API 同步更新，Compressor.forward 和 _get_states/_get_state_pool 参数从 forward_batch 改为 attn_backend。

关键文件：

- python/sglang/srt/model_executor/forward_context.py（模块 上下文层；类别 source；类型 data-contract；符号 ForwardContext, set_forward_context, has_forward_context, get_forward_context）：核心新增文件，定义了 ForwardContext 冻结数据类、forward_context 上下文管理器和全局访问函数（get_forward_context, get_attn_backend, get_token_to_kv_pool, get_req_to_token_pool）。所有模型层和注意力后端通过此模块获取运行时状态。
- python/sglang/srt/model_executor/model_runner.py（模块 模型执行；类别 source；类型 data-contract；符号 _do_forward）：入口文件，在 _forward_raw 中发布 ForwardContext，在 forward_decode 中处理 PDMux 覆写，是上下文生命周期的起始点。
- python/sglang/srt/model_executor/cuda_graph_runner.py（模块 CUDA 图捕获；类别 source；类型 data-contract；符号 run_once）：CUDA graph 捕获路径中需要在 capture_one_batch_size 内包装 ForwardContext，确保 TBO 等钩子能读取正确后端。
- python/sglang/srt/layers/attention/dsv4/compress_hip.py（模块 DSv4 压缩器；类别 source；类型 core-logic；符号 _get_states, _get_state_pool, forward）：修改 CompressorHip 的 _get_states 和 _get_state_pool 方法签名，由接受 ForwardBatch 改为接受 AttentionBackend，体现 Pattern A 的池引用迁移。
- python/sglang/srt/model_executor/cpu_graph_runner.py（模块 CPU 图；类别 source；类型 data-contract；符号 run_once）：CPU graph 捕获路径也需要 ForwardContext 包装，与 CudaGraphRunner 类似。

关键符号：ForwardContext, set_forward_context, get_forward_context, get_attn_backend, get_token_to_kv_pool, get_req_to_token_pool, forward_context, _do_forward, run_once (cuda_graph_runner), run_once (cpu_graph_runner), run_once (piecewise_cuda_graph_runner), _get_states, _get_state_pool, invalidate_loc_cache

关键源码片段

python/sglang/srt/model_executor/model_runner.py

入口文件，在 _forward_raw 中发布 ForwardContext，在 forward_decode 中处理 PDMux 覆写，是上下文生命周期的起始点。

在 _forward_raw 结尾处：

```

# Publish attn_backend for the duration of this forward.
# Save/restore so callers can nest override scopes.
with forward_context(ForwardContext(attn_backend=self.attn_backend)):
    with torch.inference_mode(), run_ctx or empty_context():
        run_once()

# 在 forward_decode 中 PDmux 覆写:
if pdmux_override:
    with forward_context(ForwardContext(attn_backend=self.decode_attn_backend)):
        return _do_forward()
return _do_forward()

```

评论区精华

1. duplicate invalidate_loc_cache (correctness) : gemini-code-assist[bot] 指出 deepseek_v4_memory_pool.py 中 invalidate_loc_cache 方法被定义两次 (high priority)。作者在后续 commit 中删除重复定义, 并确认 BaseSWAKVPool 和 DeepSeekV4TokenToKVPool 已有该方法 (由 PR #25889 引入)。
 2. ForwardContext 创建方式 (design) : ch-wan 在自审中建议直接使用 ForwardContext(attn_backend=...) 而非 dataclasses.replace(get_forward_context(), ..), 以简化代码。最终代码采纳简化写法。
 3. 注释清理 (style) : ch-wan 在多个文件 (breakable_cuda_graph_runner.py, frozen_kv_mtp_worker.py) 中要求删除冗余注释, 保持代码简洁。
- duplicate invalidate_loc_cache in DeepSeekV4TokenToKVPool (correctness): 已修复: 移除了 DeepSeekV4TokenToKVPool 中重复的 invalidate_loc_cache 定义, 该方法已由基类 BaseSWAKVPool 提供 (PR #25889 添加)。
 - ForwardContext 使用 dataclasses.replace 还是直接构造 (design): 已采纳: 最终代码使用 ForwardContext(attn_backend=self.decode_attn_backend) 直接构建新上下文。
 - 清理多余注释 (style): 已清理: 删除了多个文件中的模板式注释。

风险与影响

- 风险:
 1. 多 runner pool 误写风险: 在 target + spec draft 多 runner 场景下, 若 ForwardContext 恢复顺序错误, 可能导致 target 前向写入 draft 的 KV pool。当前设计通过 _forward_raw 的 save/restore 机制和每步显式 forward_context 包装防御。
 2. 全局变量并发安全: ForwardContext._current 是模块级全局变量, 而非 contextvars.ContextVar。目前每个 worker 进程单线程运行前向, 无竞态; 但若未来引入工作线程共享进程, 需迁移到 context var。
 3. 回归风险: Pattern A 中注意力后端缓存的 pool 引用需与 model_runner 实际 pool 保持一致, 若后端初始化后 pool 被替换 (如 frozen-KV MTP 中短暂替换), 可能引用失效。当前 frozen-KV 通过封装上下文解决。
 4. 测试覆盖: 缺少单元测试验证 ForwardContext 嵌套覆盖后的行为, 可能导致深层次遗漏。- 影响: 对用户透明, 无 API 或配置变更。对系统: 降低了 ForwardBatch 与注意

力后端的耦合，使得添加新后端或切换后端（PDmux、TBO、frozen-KV）更加统一。对团队：未来若需新增 per-forward 控制字段（如调试标识、特殊模式），只需扩展 ForwardContext 冻结类，无需调整 ForwardBatch，减少了跨模块的修改范围。 - 风险标记：多 runner 场景 context 嵌套风险，全局变量非线程安全，后端 pool 引用一致性回归

关联脉络

- 暂无明显关联 PR