

PR #25614 完整报告

sgl-project/sglang

Move profiler controls to SchedulerProfilerManager

合并时间: 2026-05-18 18:33

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/25614>

执行摘要

- 一句话: 将 profiler 方法从 Mixin 迁移至独立 Manager 组件
- 推荐动作: 值得精读, 尤其对于参与调度器重构的开发者。可以了解 Mixin 方法向组件迁移的标准步骤: 创建目标类、复制方法、更新调用点和测试、删除旧代码。此 PR 是典型的安全重构 (behaviour-preserving refactoring), 值得学习。

功能与动机

PR body 说明这是一个机械剪切粘贴操作, 作为 migrate-profiler-mixin 机制移动。目的是将 profile 控制逻辑从 SchedulerProfilerMixin 中分离出来, 归入独立的 SchedulerProfilerManager 组件, 减少调度器类的多重继承复杂度。该 PR 是调度器大规模重构 (将 Mixin 功能逐步迁移到 scheduler_components 目录下的独立组件) 的一部分。

实现拆解

1. 在 profiler_manager.py 中新增 import: 添加 logging、os、time、torch 等, 以及从被删除文件中复制来的 NPU 平台补丁和模块 logger。
2. 复制方法: 将 SchedulerProfilerMixin 中的 6 个 @staticmethod 方法 (_init_profile, _start_profile, _merge_profile_traces, _stop_profile, _profile_batch_predicate, _profile) 原样粘贴到 SchedulerProfilerManager 类中, 去除 @staticmethod 装饰器, 简化 self 类型注解。
3. 删除旧文件: 删除 scheduler_profiler_mixin.py 及其内的 SchedulerProfilerMixin 类。
4. 更新调度器类: 在 scheduler.py 中删除对 SchedulerProfilerMixin 的导入和继承, 将调用点从 self._profile(self.profiler_manager, req) 改为 self.profiler_manager._profile(req), 并调整 _profile_batch_predicate 调用。
5. 更新测试: 在 test_profile_merger.py 中将导入 SchedulerProfilerMixin 改为导入 SchedulerProfilerManager, 并更新签名检查的目标。
6. 未解决的问题: review 中建议在 __post_init__ 中初始化 rpd_profile_path 和 profile_prefix, 但该建议未被采纳, PR 已合并。

关键文件:

- python/sglang/srt/managers/scheduler_profiler_mixin.py (模块 调度器; 类别 source; 类型 deletion; 符号 SchedulerProfilerMixin, _init_profile, _start_profile, _merge_profile_traces): 源文件, 被完整删除, 包含所有 profile 控制方法。

- `python/sglang/srt/managers/scheduler_components/profiler_manager.py` (模块 分析器; 类别 `source`; 类型 `dependency-wiring`; 符号 `_init_profile`, `_start_profile`, `_merge_profile_traces`, `_stop_profile`): 目标文件, 接收所有 `profile` 方法, 添加大量 `import` 和模块级代码。
- `python/sglang/srt/managers/scheduler.py` (模块 调度器; 类别 `source`; 类型 `dependency-wiring`): 调度器主文件, 更新导入和调用方式。
- `test/registered/unit/utils/test_profile_merger.py` (模块 测试; 类别 `test`; 类型 `test-coverage`): 测试文件, 更新导入路径。

关键符号: `_init_profile`, `_start_profile`, `_merge_profile_traces`, `_stop_profile`, `_profile_batch_predicate`, `_profile`

关键源码片段

`python/sglang/srt/managers/scheduler_components/profiler_manager.py`

目标文件, 接收所有 `profile` 方法, 添加大量 `import` 和模块级代码。

```
def _init_profile(
    self,
    output_dir: Optional[str],
    start_step: Optional[int],
    num_steps: Optional[int],
    activities: Optional[List[str]],
    with_stack: Optional[bool],
    record_shapes: Optional[bool],
    profile_by_stage: bool,
    profile_id: str,
    merge_profiles: bool = False,
    profile_prefix: str = "",
    profile_stages: Optional[List[str]] = None,
) -> ProfileReqOutput:
    # 如果启用了 V2 profile, 委托给 ProfileManager 处理
    if envs.SGLANG_PROFILE_V2.get():
        return self._profile_manager.configure(
            output_dir=output_dir,
            start_step=start_step,
            num_steps=num_steps,
            activities=activities,
            with_stack=with_stack,
            record_shapes=record_shapes,
            profile_by_stage=profile_by_stage,
            profile_id=profile_id,
            merge_profiles=merge_profiles,
            profile_prefix=profile_prefix,
            profile_stages=profile_stages,
        )

    # 如果 profile 正在运行, 拒绝新的初始化请求
```

```

if self.profile_in_progress:
    return ProfileReqOutput(
        success=False,
        message="Profiling is already in progress. Call /stop_profile first.",
    )

# 存储配置参数
self.profile_by_stage = profile_by_stage
self.merge_profiles = merge_profiles

if output_dir is None:
    output_dir = os.getenv("SGLANG_TORCH_PROFILER_DIR", "/tmp")
if activities is None:
    activities = ["CPU", "GPU"]

self.torch_profiler_output_dir = Path(output_dir).expanduser()
self.torch_profiler_with_stack = with_stack
self.torch_profiler_record_shapes = record_shapes
self.profiler_activities = activities
self.profile_id = profile_id
self.profile_prefix = profile_prefix

# 根据 start_step 计算开始前向计数
if start_step:
    self.profiler_start_forward_ct = max(start_step, self.get_forward_ct() + 1)

# 根据 num_steps 设置目标前向计数或分阶段计数
if num_steps:
    if self.profile_by_stage:
        self.profiler_prefill_ct = 0
        self.profiler_decode_ct = 0
        self.profiler_target_prefill_ct = num_steps
        self.profiler_target_decode_ct = num_steps
    elif start_step:
        self.profiler_target_forward_ct = (
            self.profiler_start_forward_ct + num_steps
        )
    else:
        self.profiler_target_forward_ct = self.get_forward_ct() + num_steps
    # 当达到目标前向计数时，调用者会收到通知
else:
    self.profiler_target_forward_ct = None

return ProfileReqOutput(success=True, message="Succeeded")

```

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

调度器主文件，更新导入和调用方式。

导入部分：删除 SchedulerProfilerMixin 的导入

```

# from sglang.srt.managers.scheduler_profiler_mixin import SchedulerProfilerMixin # 被删除

# 类定义：从继承列表中移除 SchedulerProfilerMixin
class Scheduler(
    SchedulerOutputProcessorMixin,
    SchedulerUpdateWeightsMixin,
    # SchedulerProfilerMixin, # 已移除
    SchedulerMetricsMixin,
    ...
):
    ...

    # 调用点：从 self._profile(self.profiler_manager, req) 改为 self.profiler_manager._profile(req)
    (ProfileReq,
     lambda req: self.profiler_manager._profile(req),
    ),

    # 另一处调用
    # self._profile_batch_predicate(self.profiler_manager, batch) 改为
    self.profiler_manager._profile_batch_predicate(batch)

```

评论区精华

仅有一条 review 评论：gemini-code-assist[bot] 建议在 `__post_init__` 中显式初始化 `rpd_profile_path` 和 `profile_prefix`，以避免动态属性创建。该评论标记为 medium 优先级，但未在 PR 合并前解决，作者可能认为当前代码已足够（属性在对应方法中动态赋值）。

- 初始化 `rpd_profile_path` 和 `profile_prefix (style)`：未采纳，PR 已合并。

风险与影响

- 风险：风险极低，因为方法体完全一致（机械复制）。主要风险在于调用点变更：`scheduler.py` 中两处调用从 `self._xxx(profiler_manager, ...)` 改为 `self.profiler_manager._xxx(...)`，需要确保 `profiler_manager` 属性在调用时已正确初始化。此外，review 指出的 `rpd_profile_path` 和 `profile_prefix` 未在 `__post_init__` 中初始化，在方法中首次赋值前若被意外访问会引发 `AttributeError`，但目前没有这样的调用路径。
- 影响：影响范围限于调度器模块和共用测试。对用户无直接感知，对系统功能无变化。对团队来说，后续调度器组件化重构可以继续将其他 Mixin 方法迁移到对应组件，该 PR 提供了可复用的模式。
- 风险标记：属性未初始化风险，调用路径变更

关联脉络

- PR #25638 Move module-level helpers out of scheduler.py: 同属调度器组件化重构系列，也修改了 `scheduler.py` 文件。
- PR #25637 Move batch-result processing to SchedulerBatchResultProcessor: 同样是 Mixin 方法向独立组件迁移的类似操作。