

PR #25631 完整报告

sgl-project/sglang

Move idle-metrics logging to SchedulerMetricsReporter

合并时间: 2026-05-18 18:42

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/25631>

执行摘要

- 一句话: 将空闲指标日志移至 MetricsReporter 组件
- 推荐动作: 作为重构序列中的一环, 推荐阅读以理解团队如何将大混入类拆解为细粒度组件。SchedulerMetricsReporter 中的 `_maybe_log_idle_metrics` 方法展示了组件如何通过 `self.scheduler` 反向引用调度器状态, 是一种常见的组件交互模式。

功能与动机

该 PR 是调度器重构链的一部分, 旨在逐步将 Scheduler 类的混合职责拆解为独立的组件类。`_maybe_log_idle_metrics` 原本位于 SchedulerRuntimeCheckerMixin 中, 该混入类混合了运行时检查与指标记录双重职责。将其移至 SchedulerMetricsReporter 使得指标收集逻辑更加内聚, 并让 SchedulerRuntimeCheckerMixin 可以最终被移除 (后续 PR)。

实现拆解

1. 在 metrics_reporter.py 中新增 `_maybe_log_idle_metrics` 方法: 复制原方法主体, 将所有 `self.X` 状态访问替换为 `self.scheduler.X` (如 `self.scheduler.running_batch`、`self.scheduler.waiting_queue` 等), 因为 SchedulerMetricsReporter 持有指向调度器的反向引用。
2. 在 scheduler_runtime_checker_mixin.py 中删除整个文件和类: 由于方法已移出, 该文件不再需要, 彻底删除。
3. 在 scheduler.py 中进行三处配套调整: 删除导入 SchedulerRuntimeCheckerMixin; 从 Scheduler 类的基类列表中移除 SchedulerRuntimeCheckerMixin; 在 `on_idle()` 中将 `self._maybe_log_idle_metrics()` 改为 `self.metrics_reporter._maybe_log_idle_metrics()`。

关键文件:

- python/sglang/srt/managers/scheduler_runtime_checker_mixin.py (模块 调度器; 类别 source; 类型 deletion; 符号 SchedulerRuntimeCheckerMixin, `_maybe_log_idle_metrics`): 该文件被完全删除, 因为 `_maybe_log_idle_metrics` 已移出, 整个混入类不再有用途。
- python/sglang/srt/managers/scheduler_components/metrics_reporter.py (模块 调度器; 类别 source; 类型 core-logic; 符号 `_maybe_log_idle_metrics`): 新增 `_maybe_log_idle_metrics` 方法, 核心迁移目标文件。

- python/sglang/srt/managers/scheduler.py (模块 调度器; 类别 source; 类型 dependency-wiring) : 进行导入、基类列表和调用点的配套更新。

关键符号: `_maybe_log_idle_metrics`

关键源码片段

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

新增 `_maybe_log_idle_metrics` 方法, 核心迁移目标文件。

```
# 此方法从 SchedulerRuntimeCheckerMixin 原样迁移至 SchedulerMetricsReporter ,
# 唯一变化是所有 self.X 状态引用改为 self.scheduler.X (因为
# SchedulerMetricsReporter 通过 self.scheduler 持有调度器引用), 池统计
# 方法调用改为 self.scheduler.pool_stats_observer.X() 形式。
```

```
def _maybe_log_idle_metrics(self):
```

```
    """Collect and log metrics every 30 seconds during idle."""
```

```
    if (
```

```
        not self.current_scheduler_metrics_enabled
```

```
        or time.perf_counter() <= self.metrics_collector.last_log_time + 30
```

```
):
```

```
    return
```

```
self.scheduler.pool_stats_observer.get_pool_stats().update_scheduler_stats(
    self.stats
```

```
)
```

```
self.stats.num_streaming_sessions = (
```

```
    self.scheduler.pool_stats_observer.streaming_session_count()
```

```
)
```

```
self.stats.streaming_session_held_tokens = (
```

```
    self.scheduler.pool_stats_observer.session_held_tokens()
```

```
)
```

```
priority_enabled = self.scheduler.enable_priority_scheduling
```

```
self.stats.num_running_reqs = QueueCount.from_reqs(
```

```
    self.scheduler.running_batch.reqs, priority_enabled
```

```
)
```

```
self.stats.gen_throughput = 0
```

```
self.stats.num_queue_reqs = QueueCount.from_reqs(
```

```
    self.scheduler.waiting_queue, priority_enabled
```

```
)
```

```
self.stats.num_grammar_queue_reqs = len(self.scheduler.grammar_manager)
```

```
if self.scheduler.disaggregation_mode == DisaggregationMode.PREFILL:
```

```
    self.stats.num_prefill_bootstrap_queue_reqs = QueueCount.from_reqs(
```

```
        self.scheduler.disagg_prefill_bootstrap_queue.queue, priority_enabled
```

```
)
```

```
    self.stats.num_prefill_inflight_queue_reqs = QueueCount.from_reqs(
```

```
        self.scheduler.disagg_prefill_inflight_queue, priority_enabled
```

```
)
```

```
if self.scheduler.disaggregation_mode == DisaggregationMode.DECODE:
    self.stats.num_decode_prealloc_queue_reqs = QueueCount.from_reqs(
        self.scheduler.disagg_decode_prealloc_queue.queue, priority_enabled
    )
    self.stats.num_decode_transfer_queue_reqs = QueueCount.from_reqs(
        self.scheduler.disagg_decode_transfer_queue.queue, priority_enabled
    )

self.metrics_collector.log_stats(self.stats)
```

评论区精华

该 PR 无 review 评论 (review 数为 0)，提交信息明确描述为纯重定位，无行为变化。

- 暂无高价值评论线程

风险与影响

- 风险：由于是纯代码移动，几乎不存在回归风险。主要风险在于对调度器状态的间接访问（`self.scheduler.X`）可能存在遗漏或不一致的属性名，但通过逐一检查原方法与迁移后方法的每一行可确保一致性。未涉及性能关键路径，且已通过 CI（虽然 `pr-test` 显示缺少 `run-ci` 标签，但作者认为无风险直接合入）。
- 影响：对最终用户无影响，因为无行为变更。对开发团队而言，这是调度器架构重构的一步，有利于后续维护和代码可读性。`SchedulerRuntimeCheckerMixin` 文件被删除，相关依赖的开发者需注意该混入类已不再存在。
- 风险标记：低风险纯重构

关联脉络

- PR #25635 Stand up SchedulerOutputStreamer; migrate output-streaming state to it: 同一重构系列：将输出流状态从混入移动到独立组件，与本 PR 的模式完全相同。
- PR #25636 Carve out SchedulerBatchResultProcessor for batch-result state: 同样将批量结果状态从混入抽离到独立组件，反映一致的重构方向。