

PR #22911 完整报告

sgl-project/sglang

[perf] support return_routed_experts with overlap scheduling

合并时间: 2026-04-22 05:42

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/22911>

执行摘要

- 一句话: 支持返回路由专家数据时使用重叠调度, 提升 MoE 模型推理吞吐量。
- 推荐动作: 建议工程师精读此 PR, 重点关注 RoutedExpertsOutput 类的设计, 它通过分离复制和完成操作, 实现了高效的重叠调度模式, 这种模式在性能优化中值得借鉴。同时, 注意配置 `disable_overlap_schedule` 的影响, 以最大化性能收益。

功能与动机

PR body 中的性能对比图显示, 之前的实现没有支持重叠调度, 导致 `return_routed_experts` 功能存在性能瓶颈。作者通过速度测试展示了优化后的吞吐量提升, 从 172.58 tok/s 到 260.37 tok/s, 旨在减少数据传输延迟并提高推理效率。

实现拆解

1. 入口点改造: 在 `python/sglang/srt/layers/moe/routed_experts_capturer.py` 中新增 `RoutedExpertsOutput` 数据类, 提供 `copy_to_cpu()` 和 `finalize()` 方法, 用于异步 GPU 张量复制和最终写入宿主缓存。
2. 核心逻辑调整: 修改 `_RoutedExpertsCapturerReal.on_forward_end()` 方法, 添加 `no_copy_to_cpu` 参数, 当重叠调度启用时返回 `RoutedExpertsOutput` 对象而非直接同步数据; 同时重构 `_get_local_range()` 方法以统一本地范围计算。
3. 数据结构扩展: 在 `python/sglang/srt/managers/utils.py` 的 `GenerationBatchResult` 类中添加 `routed_experts_output` 字段, 并在 `copy_to_cpu()` 方法中调用其 `copy_to_cpu()` 以集成异步复制。
4. 调度器集成: 在 `python/sglang/srt/managers/scheduler_output_processor_mixin.py` 的 `process_batch_result_prefill()` 和 `process_batch_result_decode()` 中, 添加对 `routed_experts_output.finalize()` 的调用, 确保数据在复制完成后写入缓存。
5. 配套传递: 在 `python/sglang/srt/model_executor/model_runner.py` 和 `python/sglang/srt/managers/tp_worker.py` 等文件中, 将 `routed_experts_output` 作为输出的一部分传递, 以支持跨模块数据流。

关键文件:

- `python/sglang/srt/layers/moe/routed_experts_capturer.py` (模块 MoE 捕获器; 类别 `source`; 类型 `entrypoint`; 符号 `RoutedExpertsOutput`, `copy_to_cpu`, `finalize`, `on_forward_end`): 核心变更文件, 引入了 `RoutedExpertsOutput` 类并重构了

on_forward_end 方法，支持异步复制和重叠调度。

- python/sclang/srt/managers/utils.py (模块 调度工具; 类别 source; 类型 dependency-wiring) : 扩展 GenerationBatchResult 数据结构, 添加 routed_experts_output 字段并集成到 copy_to_cpu 方法中。
- python/sclang/srt/managers/scheduler_output_processor_mixin.py (模块 调度器; 类别 source; 类型 core-logic) : 调度器输出处理器的关键修改, 在结果处理中调用 routed_experts_output.finalize() 以确保数据最终化。
- python/sclang/srt/model_executor/model_runner.py (模块 模型运行器; 类别 source; 类型 data-contract) : 模型运行器的数据契约扩展, 将 routed_experts_output 作为 ModelRunnerOutput 的一部分返回。

关键符号: RoutedExpertsOutput.copy_to_cpu, RoutedExpertsOutput.finalize, _RoutedExpertsCapturerReal.on_forward_end, GenerationBatchResult.copy_to_cpu

关键源码片段

python/sclang/srt/layers/moe/routed_experts_capturer.py

核心变更文件, 引入了 RoutedExpertsOutput 类并重构了 on_forward_end 方法, 支持异步复制和重叠调度。

```
import dataclasses
import torch
from typing import Optional

@dataclasses.dataclass
class RoutedExpertsOutput:
    """Holds GPU tensors captured during forward for overlap scheduling.
    Call copy_to_cpu() inside forward stream (before copy_done.record()),
    then finalize() after copy_done.synchronize().
    """

    out_cache_loc: torch.Tensor # 输出缓存位置张量, 用于索引宿主缓存
    routed_experts: torch.Tensor # 路由专家张量, 需要异步复制到 CPU
    host_cache: "_RoutedExpertsHostCache" # 宿主缓存引用, 用于最终写入数据

    def copy_to_cpu(self):
        # 非阻塞复制张量到 CPU, 支持重叠调度中的异步操作
        self.out_cache_loc = self.out_cache_loc.to("cpu", non_blocking=True)
        self.routed_experts = self.routed_experts.to("cpu", non_blocking=True)

    def finalize(self):
        # 在复制完成后, 将数据写入宿主缓存, 确保数据一致性
        self.host_cache.buffer[self.out_cache_loc] = self.routed_experts
```

python/sclang/srt/managers/utils.py

扩展 GenerationBatchResult 数据结构, 添加 routed_experts_output 字段并集成到 copy_to_cpu 方法中。

```

from sglang.srt.layers.moe.routed_experts_capturer import RoutedExpertsOutput

@dataclasses.dataclass
class GenerationBatchResult:
    # ... 其他字段 ...
    # Routed experts: pending async D2H for overlap scheduling
    routed_experts_output: Optional[RoutedExpertsOutput] = None

    def copy_to_cpu(self, return_logprob: bool):
        # ... 复制其他张量 ...
        if self.routed_experts_output is not None:
            self.routed_experts_output.copy_to_cpu() # 异步复制路由专家数据
        if (x := self.expert_distribution_metrics) is not None:
            x.copy_to_cpu()
        self.copy_done.record() # 记录复制完成事件

```

python/sglang/srt/managers/scheduler_output_processor_mixin.py

调度器输出处理器的关键修改，在结果处理中调用 `routed_experts_output.finalize()` 以确保数据最终化。

```

def process_batch_result_prefill(
    self: Scheduler,
    batch: ScheduleBatch,
    result: Union[GenerationBatchResult, EmbeddingBatchResult],
):
    if self.is_generation:
        if result.copy_done is not None:
            result.copy_done.synchronize() # 等待复制事件完成
        if result.routed_experts_output is not None:
            result.routed_experts_output.finalize() # 最终化路由专家数据到缓存
            result.routed_experts_output = None # 清空引用，防止内存泄漏
        # ... 其他处理逻辑 ...

```

评论区精华

Review 中没有具体评论，只有两位审核者（zyzshishui 和 hnyls2002）的批准，表明变更被接受且无公开争议或疑虑。

- 暂无高价值评论线程

风险与影响

- 风险：主要风险在于异步复制可能引入数据一致性风险：如果 `copy_to_cpu()` 和 `finalize()` 调用时机不当（例如在 GPU 操作未完成时），可能导致数据损坏或丢失。此外，依赖 `disable_overlap_schedule` 配置（在 `model_runner.py` 中通过 `no_copy_to_cpu` 控制），如果配置错误可能回退到同步模式，影响性能优化效果。
- 影响：对用户而言，使用 `--enable-return-routed-experts` 的 MoE 模型推理吞吐量显著提升（约 51%），改善用户体验。系统层面，优化了调度和数据传输开销，提高了 GPU 利用率

，并支持更高效的重叠操作模式。团队需确保测试覆盖新异步路径，避免回归，但 PR 已通过 CI 测试 (run-ci 标签)，表明初步验证通过。

- 风险标记：异步复制风险，配置依赖

关联脉络

- PR #22933 [CPU] expand the interface of shared_expert without scaling factor: 同为 MoE 相关优化，扩展了 CPU 平台的 shared_expert 接口，与本 PR 的 GPU 性能优化形成互补。