

PR #26672 完整报告

sgl-project/sglang

[AMD] Work around HIP TPOT regression from Event.wait() in MTP seq lens resolution

合并时间: 2026-05-29 15:18

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/26672>

执行摘要

- 一句话: 修复 AMD MI355 上 MTP seq_lens 同步性能回退
- 推荐动作: 该 PR 为针对特定 AMD GPU 型号性能回退的临时 workaround, 代码量小且逻辑清晰, 适合快速合并以解除 AMD MI355 上的性能阻塞。建议在 AMD CI 中增加 MTP 性能测试以覆盖此场景, 并跟踪后续 AMD 驱动或 PyTorch 版本更新是否能移除该 workaround。

功能与动机

在 AMD MI355 上进行 MTP 测试时, 发现 `resolve_seq_lens_cpu` 路径中使用 `Event.wait()` 发布 `publish_ready` 会导致 TPOT 性能回退。该问题仅在 HIP 后端观察到, 其他后端应保持原有的同步行为。

实现拆解

1. 修改 `python/sglang/srt/managers/overlap_utils.py` 中 `resolve_seq_lens_cpu` 方法: 当 `publish_ready` 不为 `None` 时, 根据后端类型选择不同的同步方式。
2. 若为 HIP 后端 (`_is_hip` 为 `True`), 则调用 `self.publish_ready.synchronize()` 作为临时规避措施。
3. 其他后端继续使用 `self.publish_ready.wait()` 保持原有行为。
4. 该改动仅影响 MTP (Multi-Token Prediction) 推测解码路径, 不涉及其他功能。

关键文件:

- `python/sglang/srt/managers/overlap_utils.py` (模块 调度器; 类别 `source`; 类型 `core-logic`; 符号 `resolve_seq_lens_cpu`): 该文件实现了 MTP 推测解码中的序列长度同步逻辑。PR 在此处增加了一个 HIP 后端条件分支, 用 `synchronize()` 替代 `wait()` 以规避 AMD MI355 上的 TPOT 回退。

关键符号: `resolve_seq_lens_cpu`

关键源码片段

`python/sglang/srt/managers/overlap_utils.py`

该文件实现了 MTP 推测解码中的序列长度同步逻辑。PR 在此处增加了一个 HIP 后端条件分支, 用 `synchronize()` 替代 `wait()` 以规避 AMD MI355 上的 TPOT 回退。

```
def resolve_seq_lens_cpu(self, batch: ScheduleBatch) -> None:
    # ... 前面代码省略 ...
    fi = batch.spec_info.future_indices if batch.spec_info is not None else None
    if fi is None:
        return
    if self.publish_ready is not None:
        if _is_hip:
            # 临时 workaround: Event.wait() 在 AMD MI355 上导致 TPOT 回退
            self.publish_ready.synchronize()
        else:
            self.publish_ready.wait()
    batch.seq_lens = self.new_seq_lens_buf[fi]
    # ... 后续逻辑不变 ...
```

评论区精华

无实质性讨论。PR 由 gemini-code-assist 机器审核无反馈，并由 HaiShaw 批准。

- 暂无高价值评论线程

风险与影响

- 风险：
 - 回归风险：低。仅修改 HIP 后端的一个分支，其他后端行为不变。但缺少针对 HIP 后端的性能回归测试，无法保证该 workaround 在所有 AMD GPU 型号上均有效。
 - 功能性风险：synchronize() 是阻塞调用，会等待事件完成，可能略微增加 CPU 等待时间，但避免了更严重的 TPOT 回退。
 - 兼容性：无影响，因为改动仅限 HIP 后端。
- 影响：
 - 用户影响：AMD MI355 用户在使用 MTP 时将获得 TPOT 性能改善。其他 AMD GPU 或 NVIDIA GPU 用户无影响。
 - 系统影响：极小，仅影响 MTP 推测解码路径中的序列长度 CPU 同步环节。
 - 团队影响：低，代码改动小，但需在后续版本中确认是否需要更通用的修复。
 - 风险标记：缺少 HIP 性能测试覆盖

关联脉络

- PR #25083 fix(mooncake): honour MOONCAKE_PROTOCOL so EFA hardware can select efa transport: 同为 AMD 平台相关问题修复，涉及 HIP 后端同步机制。
- PR #26591 [AMD] Pin compressed-tensors<0.16.0 for srt_hip (fixes ROCm 7.2 nightly build): 同为 AMD HIP 后端的兼容性修复，体现 AMD 平台的持续适配。