

PR #21123 完整报告

sgl-project/sglang

[VLM] reduce CPU peak memory in multimodal tensor hashing

合并时间: 2026-03-28 11:09

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/21123>

PR 分析报告: 减少多模态张量哈希的 CPU 峰值内存

执行摘要

本 PR 通过零拷贝优化改进了多模态张量哈希的 CPU 内存使用, 消除中间分配, 在基准测试中平均 TTFT 降低约 15.8%, 内存使用大幅下降, 准确性无损失, 是一个值得关注的性能优化案例。

功能与动机

为什么做: 在服务多模态模型时, 张量哈希过程存在内存瓶颈, 原实现使用 `.float()`、`.tobytes()` 和 `torch.concat()` 导致大量中间内存分配。PR body 中明确指出: “消除所有中间内存分配”, 目标是通过零拷贝方式减少 CPU 峰值内存, 提升服务性能。基准测试显示, 优化后平均 TTFT 从 1,580.10 ms 降至 1,330.12 ms, 节省约 250 ms。

实现拆解

做了什么: 主要修改文件 `python/sglang/srt/managers/mm_utils.py` 中的两个函数:

- `tensor_hash`: 重构 CPU 路径, 移除张量拼接和显式类型转换, 改用以下代码增量哈希:

```
python hasher = hashlib.sha256() for t in tensors: t = t.detach().contiguous()
hasher.update(memoryview(t.view(torch.uint8).numpy()))
```
- `hash_feature`: 对 NumPy 数组优化, 用 `memoryview(arr)` 替代 `arr.tobytes()`。GPU 路径保持不变, 确保兼容性。

评论区精华

讨论了什么: review 中主要有两个讨论线程:

1. 代码简化建议: `gemini-code-assist[bot]` 提议统一单张量和列表处理逻辑以减少重复, 但作者未采纳, 保持优化实现。

“This function contains duplicated logic... You can simplify this...”

2. 准确性验证: `mickqian` 询问准确性结果, 作者回复准确性测试显示无问题 (`ocrbench_scorer accuracy` 保持 0.899), 并补充性能数据。

风险与影响

风险:

- 兼容性风险：依赖 memoryview 和 view(torch.uint8)，需确保与 BFloat16 等数据类型兼容；PR 中已通过 view 转换处理。
- 回归风险：哈希逻辑变更可能影响缓存，但准确性测试验证无差异。影响：
- 性能提升：TTFT 减少 15.8%，内存使用从 MB 级降至接近零，提升服务响应速度和可扩展性。
- 代码维护：变更集中，但未采纳简化建议可能增加未来维护成本。

关联脉络

与历史 PR 的关系：本 PR 与近期 PR 如 #19915（多模态工具 bugfix）和 #21418（多模态传输优化）关联，共同构成多模态性能改进的演进脉络。这些 PR 显示团队持续优化多模态模块的内存和性能，本 PR 是这一趋势的具体体现。