

PR #22989 完整报告

sgl-project/sglang

[Ray] Bind scheduler actors to GPU-local NUMA node

合并时间: 2026-04-17 05:52

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/22989>

执行摘要

- 一句话: 为 Ray 调度器 Actor 添加 GPU 本地 NUMA 绑定, 提升多 GPU 场景性能。
- 推荐动作: 该 PR 值得精读, 重点关注 NUMA 绑定在 Ray Actor 中的实现方式, 以及如何复用现有工具函数确保与 V1/V2 路径的互补性。设计决策展示了在分布式环境中处理进程绑定的优雅方案。

功能与动机

PR body 指出, Ray Actor 由 Ray 的 raylet 生成, 而非通过 multiprocessing.spawn, 因此 SGLANG_NUMA_BIND_V2 使用的 numactl 子进程包装路径 (通过 numa_utils 中的 configure_subprocess) 从未应用于它们。在默认 SGLANG_NUMA_BIND_V2=True 时, 调度器 Actor 完全未绑定, 导致性能损失。

实现拆解

1. 导入调整: 在 python/sglang/srt/ray/scheduler_actor.py 中, 从 sglang.srt.environ 导入 envs, 从 sglang.srt.managers.scheduler 导入 configure_scheduler_process (替换原 configure_scheduler), 并从 sglang.srt.utils.numa_utils 导入 get_numa_node_if_available 和 numa_bind_to_node。
2. NUMA 绑定逻辑: 在 SchedulerActor.__init__ 中, 调用 configure_scheduler_process 后, 检查 envs.SGLANG_NUMA_BIND_V2.get(), 若为 True, 则通过 get_numa_node_if_available 获取 NUMA 节点并调用 numa_bind_to_node 进行进程内绑定, 确保在调度器构造前完成。
3. 函数调用修正: 将 configure_scheduler 调用改为 configure_scheduler_process, 并传递 actual_gpu_id 参数。
4. 测试与配置: PR body 包含测试计划, 展示了性能改进数据 (如吞吐量提升 6.9%), 但未包含直接测试文件变更。

关键文件:

- python/sglang/srt/ray/scheduler_actor.py (模块 Ray 调度器; 类别 source; 类型 core-logic; 符号 SchedulerActor.init): 这是唯一修改的文件, 实现了 Ray 调度器 Actor 的 NUMA 绑定逻辑, 直接影响 Ray 部署下的性能。

关键符号: SchedulerActor.init, configure_scheduler_process, get_numa_node_if_available, numa_bind_to_node

关键源码片段

python/sglang/srt/ray/scheduler_actor.py

这是唯一修改的文件，实现了 Ray 调度器 Actor 的 NUMA 绑定逻辑，直接影响 Ray 部署下的性能。

```
from sglang.srt.environ import envs
from sglang.srt.managers.scheduler import Scheduler, configure_scheduler_process
from sglang.srt.utils.numa_utils import (
    get_numa_node_if_available,
    numa_bind_to_node,
)

# ... 在 configure_scheduler_process 调用后 ...

# Ray actors can't use the numactl subprocess-wrapping approach
# (SGLANG_NUMA_BIND_V2's normal path), so bind in-process via libnuma.
# The V1 path inside configure_scheduler_process already handles
# SGLANG_NUMA_BIND_V2=False.
if envs.SGLANG_NUMA_BIND_V2.get():
    numa_node = get_numa_node_if_available(server_args, actual_gpu_id)
    if numa_node is not None:
        numa_bind_to_node(numa_node) # 执行进程内NUMA绑定
        logger.info(
            f"[TP{tp_rank}] Bound to NUMA node {numa_node} for GPU {actual_gpu_id}"
        )

# 创建调度器，此时权重分配和NCCL初始化将在绑定的NUMA节点上执行
self.scheduler = Scheduler(...)
```

评论区精华

Review 中仅有一名审核者 (Qiaolin-Yu) 批准，无具体评论，表明变更被快速接受，可能因为逻辑清晰且基于现有 NUMA 工具。

- 暂无高价值评论线程

风险与影响

- 风险：风险较低：
- 回归风险：修改涉及核心调度器初始化路径，若 NUMA 绑定逻辑错误（如节点获取失败或绑定异常），可能导致进程启动失败或性能下降。但使用了现有辅助函数，且 V1 路径（SGLANG_NUMA_BIND_V2=False）保持不变，降低了风险。
- 兼容性风险：依赖 libnuma 库，需确保部署环境已安装；但 NUMA 绑定本就是可选功能，不影响基础功能。
- 测试覆盖不足：PR 未添加单元测试，依赖现有集成测试验证，可能掩盖边缘情况。
- 影响：影响范围：

- 用户影响：对使用 Ray 部署且启用 NUMA 绑定的用户，可显著提升性能（PR body 显示吞吐量提升达 6.9%），改善端到端延迟和 TTFT。
- 系统影响：确保调度器 Actor 的 NUMA 绑定生效，优化 GPU 内存访问和 NCCL 通信，提升多 GPU 系统资源利用率。
- 团队影响：强化了 Ray 与 NUMA 绑定的集成，为后续性能调优提供基础。
- 风险标记：核心路径变更，缺少测试覆盖

关联脉络

- PR #22994 use envs in server_args: 同样涉及环境变量 (envs) 的使用，本 PR 导入 envs 来读取 SGLANG_NUMA_BIND_V2，体现了环境变量处理的统一趋势。
- PR #22926 [misc] Configure logging before ServerArgs.post_init: 都涉及进程初始化配置的调整，本 PR 在调度器构造前进行 NUMA 绑定，类似地优化了初始化顺序。