

PR #25190 完整报告

sgl-project/sglang

fix(nvfp4): make process_weights_after_loading hot-reload-safe via alias-when-same-shape

合并时间: 2026-05-14 07:57

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/25190>

执行摘要

本 PR 修复 NVFP4 量化模型在热重载 (update_weights_from_disk) 时因 PR #25107 删除源 scale 而导致崩溃或产生错误 token 的问题。核心方案是引入 `alias_or_bind_derived_param` 工具函数, 让派生张量与源 Parameter 在形状兼容时共享同一存储, 从而既保留 ~15 GiB/rank 的内存节省, 又确保热重载后重新派生正确。该 PR 修改了 2 个文件, 新增 59 行, 删除 19 行, 已在 4xGB300 节点上验证通过。

功能与动机

PR #25107 在 `process_weights_after_loading` 中删除了不再需要的源 scales (`input_scale`, `weight_scale_2` 等), 以回收约 15 GiB/rank 内存。然而这破坏了热重载流程: `update_weights_from_disk` 会从 safetensors 重新加载原始权重, 但源 scales 已被删除, 再次调用 `process_weights_after_loading` 时引发 `AttributeError`; 即使添加跳过标志, 新加载的原始权重也与缓存的重排布局不匹配, 导致 GEMM 输出垃圾 token。因此需要一个既能保持内存节省又能安全支持热重载的解决方案。

实现拆解

1. 新增通用工具函数: 在 `python/sglang/srt/layers/utils/common.py` 中添加 `alias_or_bind_derived_param`。该函数检查派生张量是否能广播到源 Parameter 的形状, 如果可以, 则将派生值就地写入源存储并将派生属性名注册为源参数的别名 (两个名字对应同一个 Parameter 对象); 否则回退到 `copy_or_rebind_param` 分配独立参数。
2. 修改 Linear 后处理: 在 `ModelOptFp4LinearMethod.process_weights_after_loading` 中, 将 TRTLLM 和 CUTLASS 分支对 `weight_scale_interleaved` 的处理从 `copy_or_rebind_param + del weight_scale` 改为 `alias_or_bind_derived_param(layer, 'weight_scale', 'weight_scale_interleaved', ...)`。移除 `del layer.input_scale`, `layer.weight_scale_2`。alpha/input_scale_inv 等标量参数保持 `copy_or_rebind_param` (因 fused-QKV 下游内核对标量形状有依赖)。
3. 修改 MoE 后处理: 在 `ModelOptNvFp4FusedMoEMethod.process_weights_after_loading` 中, 移除对 `w13_input_scale/w2_input_scale` 的删除; TRTLLM 分支将 `w13_blockscale_swizzled/w2_blockscale_swizzled` 别名到已原地重排的 `w13_weight_scale/w2_weight_scale` (释放占位符); 非 TRTLLM 分支将 `*_blockscale_swizzled` 从 `copy_or_rebind_param + del` 改为别名。标量保留独立参数。

4. 测试：复用已有测试 `test_update_weights_from_disk_blackwell.py`，不新增测试文件。验证了两次热重载后 decode logprobs 在 `atol=1e-4` 内一致，内存节省保持。

关键源码片段

以下是新增的 `alias_or_bind_derived_param` 完整实现，它是本 PR 的核心：

```
def alias_or_bind_derived_param(
    module: torch.nn.Module,
    source_name: str,
    derived_name: str,
    derived_value: torch.Tensor,
) -> None:
    """将派生张量绑定到派生属性名称。
```

当 `derived_value`` 可广播到源 Parameter 的形状（且 dtype 匹配）时，将其广播后写入源存储，并将 `derived_name`` 注册为源参数的别名。

两个属性名共享同一底层缓冲区，因此：

- `apply()` 可以通过 `derived_name`` 读取派生数据
- `update_weights_from_disk`` 可以持续填充 `source_name``（加载器会重新运行 `process_weights_after_loading`` 来原地重新派生）
- 峰值 GPU 内存为源参数的大小，而不是源+派生之和。

当形状不可广播时，回退到通过 `copy_or_rebind_param`` 分配独立的 Parameter。
"""

```
derived_value = derived_value.detach()
source = getattr(module, source_name, None)
if isinstance(source, Parameter) and source.data.dtype == derived_value.dtype:
    try:
        broadcast = torch.broadcast_to(derived_value, source.data.shape)
    except RuntimeError:
        broadcast = None
    if broadcast is not None:
        source.data.copy_(broadcast)
        source.requires_grad_(False)
        setattr(module, derived_name, source)
    return
copy_or_rebind_param(module, derived_name, derived_value)
```

评论区精华

本 PR 没有 review 评论，但 PR body 作者详细阐述了关键设计决策：为什么纯跳过标志不可接受（产生垃圾 token），为什么标量不别名（fused-QKV 内核假定标量形状），以及形状不兼容时为什么要回退（正确性优先）。这些论述直接体现了技术权衡。

风险与影响

- 风险：仅影响 NVFP4 量化路径；`alias_or_bind_derived_param`` 的广播假设可能在某些非标准模型上不成立，此时回退到独立分配会增加峰值内存，但正确性得到保证。测试覆盖了

主流场景 (Kimi-K2.5 标准维度) 。

- 影响: 修复了 NVFP4 模型热重载问题, 使 `/update_weights_from_disk` 恢复正常; 保留了 PR #25107 的 ~15 GiB/rank 内存节省; 无推理延迟影响 (后处理仅一次调用) 。

关联脉络

- PR #25107: 最初的 NVFP4 内存优化 (删除源 scales), 本 PR 修复了其热重载兼容性问题。
- PR #24925: 同样涉及 NVFP4 相关后端集成 (tokenspeed_mla), 但无直接代码冲突。
- PR #25001: MLALoRA 支持, 与本 PR 无直接关联, 但均涉及参数绑定模式。
 - 本 PR 引入的 `alias_or_bind_derived_param` 可能为未来其他量化方案提供通用参数共享模式。