

PR #42244 完整报告

vllm-project/vllm

Avoid silent weights corruption when loading Nemotron Nano VL with reusable-buffer loaders like runai distributed streaming

合并时间: 2026-05-11 20:03

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/42244>

执行摘要

- 一句话: 修复 Nemotron Nano VL 权重加载损坏
- 推荐动作: 建议精读该 PR, 特别是生成器耗尽和多模态权重克隆的设计, 可作为多模态模型权重加载的参考模式。

功能与动机

修复 issue #41749: 当使用 `runai_streamer` 等可重用缓冲区加载器时, Nemotron Nano VL 模型权重会出现静默损坏。原代码将所有检查点张量分区为三个列表后再加载, 但可重用缓冲区加载器会在迭代间覆写底层缓冲区, 导致列表中持有过期引用。

实现拆解

1. 将 LLM 权重改为生成器流式加载: 在 `vllm/model_executor/models/nano_nemotron_vl.py` 的 `load_weights` 中, 将原 `llm_weights` 列表替换为生成器 `llm_weights_gen()`, 每个张量在参数加载前即被消费, 避免持有跨迭代的引用。
2. 多模态权重添加 `detach().clone()`: 将 `mlp1`、`vision_model`、`sound_encoder` 的权重张量在收集时通过 `w.detach().clone()` 复制, 使其独立于底层流式缓冲区。
3. 显式耗尽生成器: 在 `self.language_model.load_weights(llm_weights_iter)` 之后, 添加 `for _ in llm_weights_iter: pass` 以确保即使 LLM 加载器提前停止, 多模态权重也能被完整收集。
4. 更新测试: 在 `tests/models/multimodal/test_nano_nemotron_vl.py` 中新增 `_FakeTensor` 类以模拟 `detach().clone()` 行为, 并修改现有测试用例使用该哨兵类。

关键文件:

- `vllm/model_executor/models/nano_nemotron_vl.py` (模块 模型加载; 类别 `source`; 类型 `data-contract`; 符号 `llm_weights_gen`): 核心修复文件, 重构 `load_weights` 方法以支持可重用缓冲区加载器。
- `tests/models/multimodal/test_nano_nemotron_vl.py` (模块 测试; 类别 `test`; 类型 `test-coverage`; 符号 `_FakeTensor`, `detach`, `clone`): 新增 `_FakeTensor` 类验证 `detach().clone()` 行为, 确保测试覆盖新逻辑。

关键符号: `NemotronH_Nano_VL_V2.load_weights`, `llm_weights_gen`

关键源码片段

vllm/model_executor/models/nano_nemotron_vl.py

核心修复文件，重构 load_weights 方法以支持可重用缓冲区加载器。

```
def load_weights(self, weights: Iterable[tuple[str, torch.Tensor]]):
    mm_config = self.model_config.multimodal_config
    load_multimodal_weights = not all(
        mm_config.get_limit_per_prompt(modality) == 0
        for modality in ("image", "video", "audio")
    )
    adapter_dict = dict(self.mlp1.named_parameters())

def is_llm(name: str) -> bool:
    return name.startswith("language_model")

def is_adapter_weights(weight: tuple[str, torch.Tensor]):
    return weight[0].startswith("mlp1")

def is_vision_weights(name: str) -> bool:
    return name.startswith("vision_model.radio_model.")

def is_sound_weights(name: str) -> bool:
    return name.startswith("sound")

# LLM 权重通过生成器惰性流式加载，每个张量在迭代器前进前即被复制到参数中
# 从而避免可重用缓冲区流式加载器导致的过期引用损坏。
# 多模态组件（mlp1, vision, sound）在添加时立即 detach+clone，
# 使其独立于流式加载器的可重用缓冲区，然后在 LLM 加载完成后处理。
adapter_weights: list[tuple[str, torch.Tensor]] = []
vision_weights: list[tuple[str, torch.Tensor]] = []
sound_weights: list[tuple[str, torch.Tensor]] = []

def llm_weights_gen():
    for name, w in weights:
        if is_llm(name):
            # 去掉 'language_model.' 前缀
            yield ".".join(name.split(".")[1:]), w
        elif is_adapter_weights((name, w)):
            if not load_multimodal_weights:
                continue
            trimmed_name = ".".join(name.split(".")[1:])
            adapter_weights.append((trimmed_name, w.detach().clone()))
        elif is_vision_weights(name):
            if not load_multimodal_weights:
                continue
            # 转换: vision_model.radio_model.* → radio_model.*
            hf_key = name[len("vision_model.") :]
            vision_weights.append((hf_key, w.detach().clone()))
```

```

elif is_sound_weights(name):
    if not load_multimodal_weights:
        continue
    assert self.sound_encoder is not None
    sound_weights.append((name, w.detach().clone()))

# 完全耗尽生成器，确保即使 LLM 加载器提前停止迭代，
# 所有多模态张量也已被缓冲
llm_weights_iter = llm_weights_gen()
self.language_model.load_weights(llm_weights_iter)
for _ in llm_weights_iter:
    pass

if load_multimodal_weights:
    for trimmed_name, w in adapter_weights:
        param = adapter_dict[trimmed_name]
        with torch.no_grad():
            default_weight_loader(param, w)
    self.vision_model.load_weights(vision_weights)
    if self.sound_encoder is not None and len(sound_weights) > 0:
        self.sound_encoder.load_weights(sound_weights)

```

评论区精华

gemini-code-assist[bot] 指出：`llm_weights_gen` 依赖于 `self.language_model.load_weights` 完全消耗生成器；若加载器提前停止，多模态权重将只被部分收集。作者 noa-neria 随后增加了显式耗尽循环修复此问题。

- 生成器耗尽的安全性问题 (correctness): 作者增加显式耗尽循环 `for _ in llm_weights_iter: pass`，确保生成器被完整遍历。

风险与影响

- 风险：修改集中在权重加载路径，对推理性能和正确性无直接风险。但若其他模型也存在类似的可重用缓冲区加载器问题，此 fix 仅针对 Nemotron Nano VL，其他模型可能需要单独处理。
- 影响：用户：使用 runai distributed streaming 或其他可重用缓冲区加载器加载 Nemotron Nano VL 的用户将不再遇到静默权重损坏。系统：仅影响该模型的权重加载流程，未改动公共 API。团队：需注意类似模式可能在多模态模型间传播。
- 风险标记：核心加载路径变更，依赖加载器行为假设，仅修复单个模型

关联脉络

- PR #41383 [Nixl][PD] Lease renewal TTL KV blocks on P: 同样涉及分布式权重 / 数据传输的可重用缓冲区模式，提供了上下文背景。