

# PR #21415 完整报告

sgl-project/sglang

[diffusion] fix: fix qwen-image with nunchaku

合并时间: 2026-03-26 16:31

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/21415>

## 执行摘要

本 PR 修复了 Qwen-Image 扩散模型与 nunchaku 量化的兼容性问题，通过替换线性层为 ReplicatedLinear 并调整前向传播逻辑，确保模型在量化配置下正常运行。影响范围限于使用该量化方案的用户，但存在权重加载潜在风险，需工程师关注实现细节。

## 功能与动机

此变更旨在解决 Qwen-Image 模型在使用 nunchaku 量化时的兼容性问题，避免因权重加载失败导致推理错误。动机来源于 Issue 评论中 gemini-code-assist[bot] 的总结: 'This pull request addresses a compatibility issue with the Qwen-Image model when used with nunchaku quantization.' 确保模型能顺利集成量化配置。

## 实现拆解

主要改动集中在 `python/sglang/multimodal_gen/runtime/models/dits/qwen_image.py` 文件的 `QwenImageBlock` 类中:

- 初始化部分: 将 `to_add_out`、`to_out`、`img_mod` 和 `txt_mod` 中的 `nn.Linear` 替换为 `ReplicatedLinear`，并添加 `quant_config` 和 `prefix` 参数。例如: `python self.to_add_out = ReplicatedLinear(self.inner_dim, self.dim, bias=out_bias, quant_config=quant_config, prefix=f"{prefix}.to_add_out", )`
- 前向传播部分: 调整 `forward` 方法中 `img_mod_params` 和 `txt_mod_params` 的获取，从单返回值改为解包元组: `python img_mod_params, _ = self.img_mod[1](temb_img_silu) # 原为 img_mod_params = self.img_mod[1](temb_img_silu) txt_mod_params, _ = self.txt_mod[1](temb_txt_silu)`

## 评论区精华

review 中仅有的实质性讨论来自 gemini-code-assist[bot]，重点指出 `prefix` 参数设置问题:

```
The prefix for ReplicatedLinear appears to be incorrect. Since this layer is at index 1 within the img_mod nn.Sequential module, its name in the model hierarchy will be ..img_mod.1. To ensure correct weight loading, the prefix should be updated to f'{prefix}.img_mod.1'.
```

此建议未在合并前被采纳，讨论戛然而止，可能留下未修复隐患。

## 风险与影响

- 技术风险: prefix 参数错误可能导致模型权重加载异常, 影响输出准确性; forward 方法调整不当可能引发运行时错误; 缺乏专门测试, 无法验证修复效果。
- 影响分析: 直接影响使用 Qwen-Image 与 nunchaku 量化的用户, 修复后应能正常推理, 但风险需监控。对系统其他部分无波及, 属于局部优化。

## 关联脉络

从近期 PR 看, 本 PR 与以下变更相关:

- PR 21348: 修复 MxInt4 MoE 量化问题, 共享 quant 标签, 体现团队对量化兼容性的持续关注。
- PR 21387: 优化 diffusion 模型的 Triton 内核, 同属 diffusion 模块演进, 显示该领域的技术迭代趋势。整体上, 这些 PR 共同推动扩散模型在多模态和量化方向的完善。