

# PR #26045 完整报告

sgl-project/sglang

Apply apply\_group\_norm\_silu to LTX-2 latent upsampler

合并时间: 2026-06-02 17:28

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/26045>

## 执行摘要

- 一句话: LTX-2 upsampler 融合 GroupNorm+SiLU kernel 加速 6.5-14x
- 推荐动作: 值得精读。PR 展示了清晰的 kernel 融合集成实践: 从定位热点、替换调用、编写多层次测试到添加基准测试, 每一步都有详细解释。对于希望在 SGLang 或其他推理框架中应用类似优化的读者有很好参考价值。

## 功能与动机

LTX-2 LatentUpsampler 位于每次生成的关键路径, 其内部使用手动的  $\text{silu}(\text{norm}(x))$  模式, 未利用已有的 Triton 融合内核。通过接入 `apply_group_norm_silu`, 可以在不改变输出的情况下获得显著性能提升, 符合 SGLang 扩散路线图中的内核融合方向。

## 实现拆解

1. 导入融合核: 在 `latent_upsampler.py` 顶部增加 `from sglang.jit_kernel.diffusion.group_norm_silu import apply_group_norm_silu` 导入。
2. 修改 `ResBlock.forward`: 将 `self.norm1(x); self.activation(x)` 替换为 `apply_group_norm_silu(x, self.norm1, self.activation)`, 并添加注释说明第二个 norm 站点因模式不符 ( $\text{silu}(\text{norm} + \text{residual})$ ) 暂未融合。
3. 修改 `LatentUpsampler.forward`: 在 2D 路径和 3D 路径的初始卷积后, 同样将 `initial_norm + initial_activation` 替换为融合调用。变更覆盖 `initial_norm`、`norm1` 两处, 每处每个前向调用触发一次融合。
4. 新增基准测试: 在 `bench_group_norm_silu.py` 中添加三个 LTX-2 实际 shape 用例 (`small`、`pre_720p`、`post_720p`), 其中大 case 因显存需求标记为 `LARGE_CASES`, 通过 `--cases large` 或 `--cases all-large` 单独启用。
5. 编写单元测试: 新建 `test_latent_upsampler_group_norm_silu.py`, 包含 19 个测试: CPU 上 `atol=0, rtol=0` 的零容忍数值对照、CUDA 上 `bf16/fp16` 的生产路径容忍对照, 以及通过 `patch` 验证融合调用次数的测试。

关键文件:

- `python/sglang/multimodal_gen/runtime/models/upsampler/latent_upsampler.py` (模块上采样器; 类别 `source`; 类型 `core-logic`; 符号 `ResBlock.forward`, `LatentUpsampler.forward`): 核心修改文件: 替换三个  $\text{silu}(\text{norm}(x))$  站点为融合调用, 是性能提升的直接来源。

- python/sglang/jit\_kernel/benchmark/diffusion/bench\_group\_norm\_silu.py (模块 JIT 内核基准; 类别 source; 类型 core-logic; 符号 CASES, LARGE\_CASES, parse\_cases) : 添加 LTX-2 实际 shape 的基准用例 (small, pre\_720p, post\_720p), 便于量化性能提升和验证大显存场景。
- python/sglang/multimodal\_gen/test/unit/test\_latent\_upsampler\_group\_norm\_silu.py (模块 测试验证; 类别 test; 类型 test-coverage; 符号 \_resblock\_eager\_reference, \_latent\_upsampler\_eager\_reference, test\_resblock\_forward\_parity, test\_latent\_upsampler\_forward\_parity) : 新测试文件: 覆盖 CPU 零容忍精度、CUDA 生产路径容忍、以及融合调用次数验证, 确保正确性。

关键符号: apply\_group\_norm\_silu, ResBlock.forward, LatentUpsampler.forward, native\_group\_norm\_silu, \_resblock\_eager\_reference, \_latent\_upsampler\_eager\_reference

## 关键源码片段

[python/sglang/multimodal\\_gen/runtime/models/upsampler/latent\\_upsampler.py](#)

核心修改文件: 替换三个 silu(norm(x)) 站点为融合调用, 是性能提升的直接来源。

```
class ResBlock(torch.nn.Module):
    # ... __init__ unchanged ...

    def forward(self, x: torch.Tensor) -> torch.Tensor:
        residual = x
        x = self.conv1(x)
        # 融合 GroupNorm + SiLU, 替代原来的 self.norm1(x) + self.activation(x)
        # 第二个 norm (norm2) 之后是 silu(norm2 + residual), 当前 helper 不覆盖
        x = apply_group_norm_silu(x, self.norm1, self.activation)
        x = self.conv2(x)
        x = self.norm2(x)
        x = self.activation(x + residual)
        return x

class LatentUpsampler(torch.nn.Module):
    # ... __init__ unchanged ...

    def forward(self, latent: torch.Tensor) -> torch.Tensor:
        b, _, f, _, _ = latent.shape
        if self.dims == 2:
            x = rearrange(latent, "b c f h w -> (b f) c h w")
            x = self.initial_conv(x)
            # 合并 initial_norm 和 initial_activation 为一个融合调用
            x = apply_group_norm_silu(x, self.initial_norm, self.initial_activation)
            for block in self.res_blocks:
                x = block(x)
            # ... upsampler and post blocks ...
        else:
```

```
x = self.initial_conv(latent)
# 3D 路径同样融合
x = apply_group_norm_silu(x, self.initial_norm, self.initial_activation)
# ... remaining 3D path ...
```

## python/sclang/jit\_kernel/benchmark/diffusion/bench\_group\_norm\_silu.py

添加 LTX-2 实际 shape 的基准用例 (small, pre\_720p, post\_720p) , 便于量化性能提升和验证大显存场景。

```
CASES = [
    # ... 已有通用用例 ...
    # LTX-2 latent upsampler 的实际 shape (num_groups=32)
    Case("ltx2_upsampler_small", (1, 512, 8, 45, 80), 32),
    Case("ltx2_upsampler_pre_720p", (1, 512, 16, 90, 160), 32),
]

# 大显存用例: post_720p 约 940 MB, 原生路径中间变量约 5 GB
# 只在 H100/H200 等大显存 GPU 运行, 通过 --cases large 或 --cases all-large 启用
LARGE_CASES = [
    Case("ltx2_upsampler_post_720p", (1, 512, 16, 180, 320), 32),
]

CASE_BY_NAME = {case.name: case for case in CASES + LARGE_CASES}

def parse_cases(text: str) -> list[Case]:
    if text == "all":
        return CASES
    if text == "large":
        return LARGE_CASES
    if text == "all-large":
        return CASES + LARGE_CASES
    # ... 原有按名称解析逻辑 ...
```

## 评论区精华

Reviewer **BBuf** 指出测试文件过于复杂, 建议简化。作者随后大幅精简测试代码 (从 462 行减至 306 行), 去除大量 AI 生成注释, 最终获得 **BBuf** 批准并合并。讨论反映了社区对测试可维护性的重视。

- 测试文件复杂度问题 (testing): 作者简化了测试文件 (移除大量 AI 生成注释, 减少行数), **BBuf** 最终批准合并。

## 风险与影响

- 风险:
  1. 精度回归风险: 融合核可能存在数值差异。CPU 测试使用 float32 atol=0, rtol=0 确保完全一致; CUDA 测试设置了合理的容忍阈值 (bf16: 7e-2, fp16: 3e-3), 并通过

patch 验证实际调用了 Triton 内核。风险较低。

2. 大显存案例 OOM: post\_720p 基准测试约 940 MB 张量, 加上 native 对比路径可能占用 5 GB, 在 24 GB GPU 上可能 OOM, 已通过 LARGE\_CASES 隔离, 不影响常规 CI。

3. 未覆盖的 residual 路径: ResBlock 中第二个 norm 站点 (silu(norm2(x) + residual)) 未融合, 因为现有 helper 不支持。此局限已明确注释, 不构成 bug, 但后续可扩展。

• 影响:

- 用户侧: 所有 LTX-2 生成受益——每个 upsampler 调用融合 9 个站点 (默认配置), 端到端延迟降低。高分辨率 (如 720p) 加速更显著 (约 14x kernel 加速)。
- 系统侧: 新增 19 个单元测试和 3 组基准测试, 覆盖 CPU 和 CUDA 路径, 增强了扩散模块的测试基础。基准测试标记为独立运行, 不影响 CI 资源。
- 团队侧: 展示了将 Triton 融合内核应用于生产模块的操作流程, 包括性能验证、精度测试和显存保护策略。该模式可推广到其他 diffusion 模块。
- 风险标记: 精度回归风险 (已全面测试), 大内存基准案例需 H100-class GPU, 未融合的 residual 加法路径 (有意为之, 可后续扩展)

## 关联脉络

- 暂无明显关联 PR