

PR #23633 完整报告

sgl-project/sglang

[MUSA] Use MUSA-optimized operators in piecewise CUDA graph

合并时间: 2026-05-12 08:55

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/23633>

执行摘要

- 一句话: MUSA 设备启用 piecewise CUDA graph 并优化算子
- 推荐动作: 该 PR 值得精读, 特别关注 MUSA 后端的 FX 补丁和 fake kernel 注册模式。对于跨平台算子层的开发者来说, 这是一种通用的解决 torch.compile 兼容性的方法。PR 中的 review 讨论也提供了良好的代码健壮性实践 (使用 `_replace`、`getattr`、条件注册)。建议后续增加 MUSA 硬件上的 CI 测试。

功能与动机

MUSA devices previously could not use piecewise CUDA graph due to multiple incompatibilities: torch.compile on MUSA fails to serialize custom device types in FX-generated code, and several operators (SiluAndMul, RMSNorm, FP8 quantization) lacked fake kernel registrations needed for torch.compile tracing. As a workaround, these operators fell back to native PyTorch implementations when piecewise CUDA graph was enabled, resulting in suboptimal performance.

实现拆解

1. 新增 FX 代码生成补丁: 创建 `python/sglang/srt/hardware_backend/musa/utils/patch_torch.py`, 实现 `patch_fx_custom_device()` 函数, 通过正则替换将 FX 生成的 `device(type='musa', index=N)` 转换为 `torch.device('musa:N')`, 解决 torch.compile 序列化问题。同时通过 `_replace` 方法安全修改 `PythonCode` 对象。
2. 注册 Fake Kernel: 在 `activation.py` 中注册 `aten::_fused_swiglu_forward` 的 fake kernel, 在 `fp8_kernel.py` 中注册 `sgl_kernel::sgl_per_token_group_quant_8bit_v2` 的 fake kernel, 均使用 `register_fake_if_exists` 条件注册以避免算子不存在时出错。这使得 torch.compile 能在 MUSA 设备上正确跟踪这些算子。
3. 移除 Piecewise CUDA Graph Guard: 删除 `activation.py` 中 `SiluAndMul.forward_musa` 和 `layernorm.py` 中 `RMSNorm.forward_musa` 内的回退检查 (`if not disable_piecewise_cuda_graph: return self.forward_native(x)`), 使 MUSA 优化路径 (如 `nn.SwishGLU`) 在 PCG 启用时也被使用。
4. 集成补丁调用: 在 `piecewise_cuda_graph_runner.py` 的 `set_torch_compile_config()` 中检测 MUSA 设备后调用 `patch_fx_custom_device()`, 确保 PCG 初始化时应用补丁。

5. 修复 MUSA 兼容性: 在 `sgl-kernel/python/sgl_kernel/utils.py` 中将 `is_arch_support_pdl()` 的条件从仅检查 `torch.version.hip` 扩展为同时检查 `musa`。同时根据 review 建议将补丁代码从 `utils/patch_torch.py` 迁移到专门的 MUSA 后端目录, 保持组织清晰。

测试配套: 本次变更未包含直接测试文件, 可能需要在 MUSA 硬件上验证 PCG 功能与精度。

关键文件:

- `python/sglang/srt/hardware_backend/musa/utils/patch_torch.py` (模块 MUSA 后端; 类别 source; 类型 dependency-wiring; 符号 `_replace_device_repr`, `patch_fx_custom_device`, `patched`): 核心新增文件, 实现 FX 代码生成补丁, 修复 `torch.compile` 在 MUSA 上的 device 序列化问题, 是本次 PR 的关键使能技术。
- `python/sglang/srt/layers/activation.py` (模块 激活层; 类别 source; 类型 core-logic; 符号 `_`): 移除 `SiluAndMul` 中的 PCG guard, 并注册 fake kernel 使 `torch.compile` 能跟踪 SwishGLU 算子, 直接影响激活层性能。
- `python/sglang/srt/layers/quantization/fp8_kernel.py` (模块 量化层; 类别 source; 类型 core-logic; 符号 `_`): 注册 `sgl_per_token_group_quant_8bit_v2` 的 fake kernel, 使 FP8 量化在 `torch.compile` 下可跟踪。
- `python/sglang/srt/model_executor/piecewise_cuda_graph_runner.py` (模块 计算图运行; 类别 source; 类型 data-contract): 在 `set_torch_compile_config` 中集成 MUSA 补丁调用, 是 PCG 初始化的关键入口。
- `python/sglang/srt/utils/patch_torch.py` (模块 工具层; 类别 source; 类型 dependency-wiring): 增加了 `is_musa` 导入和 `_is_musa` 检查, 为后续 MUSA 条件逻辑提供基础。
- `python/sglang/srt/layers/layernorm.py` (模块 归一化层; 类别 source; 类型 core-logic): 移除 `RMSNorm.forward_musa` 中的 PCG guard, 使 MUSA 优化路径始终生效。
- `sgl-kernel/python/sgl_kernel/utils.py` (模块 sgl 内核工具; 类别 source; 类型 core-logic): 修复 `is_arch_support_pdl` 支持 MUSA, 避免 MUSA 设备上功能异常。

关键符号: `_replace_device_repr`, `patch_fx_custom_device`, `patched`, `SiluAndMul.forward_musa`, `RMSNorm.forward_musa`, `set_torch_compile_config`, `is_arch_support_pdl`

关键源码片段

`python/sglang/srt/hardware_backend/musa/utils/patch_torch.py`

核心新增文件, 实现 FX 代码生成补丁, 修复 `torch.compile` 在 MUSA 上的 device 序列化问题, 是本次 PR 的关键使能技术。

```
import re
from dataclasses import replace as _dataclass_replace

import torch
import torch.fx.graph as fx_graph

# 正则匹配 device(type='xxx', index=N) 格式
```

```
_DEVICE_REPR_RE = re.compile(r"\bdevice\(type='([\^']+)'(?:,\s*index=(\d+))?\)")
```

```
def _replace_device_repr(m: re.Match) -> str:
    """将匹配到的 device 表达式替换为 torch.device 调用"""
    dev_type = m.group(1)
    dev_index = m.group(2)
    if dev_index is not None:
        # 例如 device(type='musa', index=0) -> torch.device('musa:0')
        return f"torch.device('{dev_type}:{dev_index}')"
    return f"torch.device('{dev_type}')"

def patch_fx_custom_device() -> None:
    """
    Fix FX codegen serialization for non-standard devices (e.g., torch_musa).
    通过包装 _gen_python_code, 在生成的代码字符串中替换 device 表示,
    并确保 torch 模块在 graph globals 中可用。
    """
    original = fx_graph.CodeGen._gen_python_code

    def patched(self, nodes, root_module, namespace, **kwargs):
        result = original(self, nodes, root_module, namespace, **kwargs)
        new_src = _DEVICE_REPR_RE.sub(_replace_device_repr, result.src)
        if new_src is result.src:
            # 没有需要替换的, 返回原始结果
            return result
        # 保证 torch 被添加到 globals 中
        result.globals.setdefault("torch", torch)
        # 使用 _replace 安全地更新 src, 保留其他字段 (如 _lineno_map)
        if hasattr(result, "_replace"):
            return result._replace(src=new_src)
        return _dataclass_replace(result, src=new_src)

    fx_graph.CodeGen._gen_python_code = patched
```

python/slang/srt/layers/activation.py

移除 SiluAndMul 中的 PCG guard, 并注册 fake kernel 使 torch.compile 能跟踪 SwishGLU 算子, 直接影响激活层性能。

```
elif _is_musa:
    # 导入条件 fake kernel 注册工具
    from slang.srt.utils.patch_torch import register_fake_if_exists

    # 注册 aten::_fused_swiglu_forward 的 fake kernel, 使 torch.compile 能跟踪 nn.SwishGLU
    @register_fake_if_exists("aten::_fused_swiglu_forward")
    def _(x):
        # 仅需返回正确的输出形状和 dtype, 不需要实际计算
        d = x.shape[-1] // 2
```

```

        output_shape = x.shape[:-1] + (d,)
        return torch.empty(output_shape, dtype=x.dtype, device=x.device)

# ... 省略无关代码 ...

class SiluAndMul(MultiPlatformOp):
    # ... 其他方法 ...

    def forward_musa(self, x: torch.Tensor) -> torch.Tensor:
        # 移除原来对 PCG 的 guard, 现在总是使用 MUSA 优化路径
        if not hasattr(self, "_musa_swish_glu"):
            # nn.SwishGLU 在 MUSA 上性能优于 silu_and_mul, 故优先使用
            self._musa_swish_glu = nn.SwishGLU()
        return self._musa_swish_glu(x)

```

评论区精华

- 使用 `_replace` 方法: `gemini-code-assist` 指出直接构造 `PythonCode` 对象可能因版本差异导致 `TypeError`, 建议使用 `NamedTuple` 的 `_replace` 方法以保留其他字段。开发者最终采用了该方案。
- 安全访问 `torch.version.musa`: `gemini-code-assist` 建议使用 `getattr(torch.version, 'musa', None)` 替代直接访问, 避免 `AttributeError`。开发者接受了该建议, 在 `sgl-kernel/utils.py` 中已修正。
- 条件注册 Fake Kernel: `gemini-code-assist` 建议使用 `register_fake_if_exists` 避免在算子未定义时出错。开发者采纳, 在 `fp8_kernel.py` 和 `activation.py` 中均改用条件注册。
- 代码迁移: `yeahdongcn` 建议将 `patch_torch.py` 中的补丁逻辑移到 `hardware_backend/musa/utils/` 下, 开发者通过独立 `commit` 完成迁移。
- TODO 实现: `yeahdongcn` 询问是否应实现 MUSA 专用的 `silu_and_mul` kernel 以替代 `nn.SwishGLU`, 但未在本次 PR 中解决。
 - 使用 `_replace` 替代直接构造 `PythonCode` (design): 开发者采纳, 最终代码使用 `result._replace(src=new_src)` 并回退到 `_dataclass_replace`。
 - 安全访问 `torch.version.musa` (correctness): 开发者采纳, 在 `sgl-kernel/utils.py` 中修改为 `getattr(torch.version, 'musa', None)`。
 - 条件注册 Fake Kernel (correctness): 开发者采纳, 在 `fp8_kernel.py` 和 `activation.py` 中使用条件注册。
 - 代码迁移到 MUSA 后端目录 (design): 开发者通过独立 `commit` 完成迁移。
 - 待办: 实现 MUSA 专用 `silu_and_mul` kernel (question): 未在本次 PR 中解决, 作为已知限制遗留。

风险与影响

- 风险:
 1. 正则替换覆盖不全面: `patch_fx_custom_device` 使用的正则可能无法处理所有 `device` 表示形式 (如复合表达式), 存在遗漏场景。

2. Fake Kernel 语义不足：注册的 fake kernel 仅返回空张量或简单形状，可能无法完全匹配真实算子的 shape 推导，导致图捕获或编译错误。
 3. 回退移除风险：移除 forward_native 回退后，若 nn.SwishGLU 或 RMSNorm 的 MUSA 实现存在边界 bug，在没有备用路径的情况下可能导致推理失败。
 4. MUSA 环境差异：不同 MUSA 设备或驱动版本可能影响算子行为，缺乏 CI 覆盖。
 5. 缺少测试覆盖：PR 未包含单元测试或集成测试，无法有效验证变更的正确性。 - 影响：
用户影响：MUSA 设备用户现在可以启用 --piecewise-cuda-graph 并享受优化的算子（如 nn.SwishGLU），预期提升推理吞吐和延迟。非 MUSA 用户无影响。系统影响：修改了多个核心层（激活、归一化、量化、PCG 运行器），但改动范围较小且通过条件分支隔离。团队影响：需要 MUSA 维护者验证功能正确性，并考虑后续添加专用 kernel 和测试。
- 风险标记：核心路径变更，缺少测试覆盖，多平台兼容性风险，正则替换可能遗漏，Fake Kernel 语义可能不完整

关联脉络

- 暂无明显关联 PR