

PR #24978 完整报告

sgl-project/sglang

[MUSA]: Add flashinfer sampling backend

合并时间: 2026-05-15 11:23

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/24978>

执行摘要

- 一句话: 为 MUSA 添加 FlashInfer 采样后端
- 推荐动作: 值得阅读, 特别是对 MUSA 后端的适配方式。设计决策中采用了与 CUDA 后端类似的接口封装, 便于未来统一。

功能与动机

Add FlashInfer sampling backend support for MUSA to enable optimized sampling operations on MUSA devices.

实现拆解

1. C++ 接口声明: 在 `sgl-kernel/include/sgl_kernel_musa_ops.h` 中添加了 `min_p_sampling_from_probs` 和 `top_p_sampling_from_probs` 函数声明, 供 Torch 绑定使用。
2. Torch 算子注册: 在 `sgl-kernel/csrc/common_extension_musa.cc` 中通过 `TORCH_LIBRARY_EXPAND` 注册上述算子, 关联 MUSA 后端实现。
3. Python 封装: 在 `sgl-kernel/python/sgl_kernel/musa.py` 中新增内部函数和导出函数 (如 `top_k_renorm_probs`、`top_p_sampling_from_probs` 等), 统一处理类型转换和参数传递。
4. 采样器集成: 在 `python/sglang/srt/layers/sampler.py` 中添加 `if is_musa():` 条件分支, 从 `sgl_kernel` 导入所需采样函数, 使得 MUSA 设备运行时自动使用新后端。
5. 包入口和构建: 在 `sgl-kernel/python/sgl_kernel/__init__.py` 导出新函数, 并在 `sgl-kernel/setup_musa.py` 中将 FlashInfer 的 `sampling.cu` 加入编译列表。
6. 依赖更新: 更新 `python/pyproject_other.toml`、`sgl-kernel/pyproject_musa.toml` 和 `3rdparty/amd/wheel/sglang/pyproject.toml` 中的 `torchada` 版本至 0.1.56。

关键文件:

- `sgl-kernel/python/sgl_kernel/musa.py` (模块 采样封装; 类别 source; 类型 dependency-wiring; 符号 `_top_k_renorm_probs_internal`, `top_k_renorm_probs`, `_top_p_renorm_probs_internal`, `top_p_renorm_probs`): 添加了所有 FlashInfer 采样函数的 Python 封装, 是 MUSA 采样后端的核心接口
- `python/sglang/srt/layers/sampler.py` (模块 采样器集成; 类别 source; 类型 dependency-wiring): 添加了 MUSA 分支的采样函数导入, 使采样器在 MUSA 设备上使

用新后端

- `sgl-kernel/include/sgl_kernel_musa_ops.h` (模块 采样接口; 类别 `source`; 类型 `core-logic`): 声明了新的采样函数接口, 供 Torch 绑定使用
- `sgl-kernel/csrc/common_extension_musa.cc` (模块 扩展注册; 类别 `source`; 类型 `core-logic`): 注册了新的采样算子到 Torch 库, 连接 C++ 实现和 Python 调用
- `sgl-kernel/python/sgl_kernel/__init__.py` (模块 包入口; 类别 `source`; 类型 `core-logic`): 导出新采样函数到 `sgl_kernel` 包级别

关键符号: `top_k_renorm_probs`, `top_p_renorm_probs`, `top_p_sampling_from_probs`, `top_k_top_p_sampling_from_probs`, `min_p_sampling_from_probs`

关键源码片段

`sgl-kernel/python/sgl_kernel/musa.py`

添加了所有 FlashInfer 采样函数的 Python 封装, 是 MUSA 采样后端的核心接口

```
def _top_p_sampling_from_probs_internal(
    probs: torch.Tensor,
    indices: Optional[torch.Tensor],
    maybe_top_p_arr: Optional[torch.Tensor],
    top_p_val: float,
    deterministic: bool,
    generator: Optional[torch.Generator],
) -> torch.Tensor:
    # 获取设备并转换概率为 float
    device = probs.device
    probs = probs.float()
    # 类型转换: top_p 数组转为 float (若提供)
    maybe_top_p_arr = maybe_top_p_arr.float() if maybe_top_p_arr is not None else None
    # 预分配输出张量 ( int32 )
    samples = torch.empty(probs.size(0), dtype=torch.int32, device=device)
    # 调用底层 MUSA 算子
    torch.ops.sgl_kernel.top_p_sampling_from_probs.default(
        probs, samples, indices, maybe_top_p_arr, top_p_val, deterministic, generator,
    )
    return samples

def top_p_sampling_from_probs(
    probs: torch.Tensor,
    top_p: Union[torch.Tensor, float],
    indices: Optional[torch.Tensor] = None,
    deterministic: bool = True,
    generator: Optional[torch.Generator] = None,
    check_nan: bool = False,
) -> torch.Tensor:
    # 可选的 NaN 检查
    if check_nan and torch.any(torch.isnan(probs)):
```

```
raise ValueError("Input probs contains NaN.")
# 将标量 top_p 参数转换为 (tensor, val) 统一格式
return _top_p_sampling_from_probs_internal(
    probs, indices, *_to_tensor_scalar_tuple(top_p), deterministic, generator
)
```

评论区精华

Review 中 `gemini-code-assist` 指出了三个关键错误：在 `musa.py` 中错误地将 `probs.device` 作为上下文管理器使用（`with probs.device as device`），这会导致运行时 `AttributeError`。建议改为直接赋值变量后使用。作者已确认修复。此外 reviewer `yeahdongcn` 要求 rebase 并参考 `flashinfer` 官方实现的方式（使用 `probs.device` 而不是上下文管理器），并最终批准了 PR。

- `torch.device` 上下文管理器错误 (correctness): 作者已修复 (fixed.)
- 建议参考 `flashinfer` 官方实现使用 `probs.device` (design): 作者已采纳并修复, `yeahdongcn` 最终批准。

风险与影响

- 风险：
 - 技术风险：新后端可能在某些 MUSA 设备上不稳定，缺少充分的单元测试（没有测试文件变更）。
 - 性能风险：新后端可能引入性能回归，但预期是优化。
 - 兼容性风险：`torchada` 版本更新可能影响其他依赖；但只在 MUSA 下生效，不影响其他后端。
- 影响：
 - 用户：MUSA 用户可使用 `FlashInfer` 采样，获得性能提升。
 - 系统：增加了条件导入，非 MUSA 环境无变化。
 - 团队：需要维护 MUSA 特有的采样代码，但核心逻辑来自 `FlashInfer`，降低维护成本。
 - 风险标记：MUSA 新后端，`torch.device` 错误已修复，依赖版本变更，缺少单元测试

关联脉络

- 暂无明显关联 PR