

# PR #16793 完整报告

sgl-project/sglang

Add deterministic mode for XPU operations

合并时间: 2026-04-30 13:39

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/16793>

## 执行摘要

- 一句话: 为 XPU 操作添加确定性模式
- 推荐动作: 值得精读, 特别是 `batch_invariant_ops.py` 中通过 `get_dispatch_device_backend` 实现多后端调度的设计模式, 以及 `layernorm` 中条件分支的编写方式, 对后续扩展其他硬件有参考意义。PR 表明团队在为 Intel XPU 提供一阶支持, 建议关注后续 `intel_dev` 分支的集成情况。

## 功能与动机

在 XPU 设备上推理时, 需要确定性计算以保证结果可重现。本 PR 解决了此前 `batch_invariant_ops` 和 `RMSNorm` 仅支持 CUDA 的问题, 使 XPU 也能享受确定性模式。

## 实现拆解

1. 在 `python/sglang/srt/utils/common.py` 中新增 `get_dispatch_device_backend()` 函数, 通过检查 `is_cuda_alike()` 和 `is_xpu()` 决定返回 "CUDA" 或 "XPU"; 同时在 `set_random_seed` 中添加 XPU 的种子设置, 在 `get_device_core_count` 中添加 XPU 的 EU 核心数获取。
2. 在 `python/sglang/srt/batch_invariant_ops/batch_invariant_ops.py` 中导入新工具函数, 将原先硬编码为 "CUDA" 的 `dispatch` key 替换为 `get_dispatch_device_backend()` 的返回值; 同时将 `NUM_SMS` 的获取改为调用 `get_device_core_count()`, 使 Triton 内核调度适应不同设备。
3. 在 `python/sglang/srt/layers/layernorm.py` 的 `RMSNorm` 类的 `forward_xpu` 方法中, 新增分支: 当 `is_batch_invariant_mode_enabled()` 时, 根据条件选择 `forward_native` 或 `rms_norm_batch_invariant`, 确保在 XPU 上也能使用确定性批量不变算子。
4. 配置和测试: PR 未包含单独的测试文件, 但作者提及已有测试 `test/srt/test_deterministic.py` 中的 `TestTritonDeterministic::test_single` 在 XPU 上通过。

关键文件:

- `python/sglang/srt/utils/common.py` (模块 工具函数; 类别 `source`; 类型 `core-logic`; 符号 `get_dispatch_device_backend`, `set_random_seed`, `get_device_core_count`): 核心工具函数, 新增了设备后端分发逻辑并扩展了随机种子和设备核心数获取, 是跨后端支持的基础。

- python/sglang/srt/layers/layernorm.py (模块 模型层; 类别 source; 类型 core-logic; 符号 forward\_xpu) : RMSNorm 的 forward\_xpu 方法新增了确定性模式分支, 使 XPU 在启用 batch invariant 时能使用确定性算子。
- python/sglang/srt/batch\_invariant\_ops/batch\_invariant\_ops.py (模块 算子注册; 类别 infra; 类型 infrastructure; 符号 enable\_batch\_invariant\_mode) : 注册确定性算子的核心文件, 通过动态 dispatch key 支持 XPU。

关键符号: get\_dispatch\_device\_backend, set\_random\_seed, get\_device\_core\_count, forward\_xpu, enable\_batch\_invariant\_mode, is\_batch\_invariant\_mode\_enabled

## 关键源码片段

### python/sglang/srt/utils/common.py

核心工具函数, 新增了设备后端分发逻辑并扩展了随机种子和设备核心数获取, 是跨后端支持的基础。

```
# 新增函数: 返回当前加速器后端的 dispatch key
@lru_cache(maxsize=1)
def get_dispatch_device_backend():
    if is_cuda_alike():
        dispatch_key = "CUDA"
    elif is_xpu():
        dispatch_key = "XPU"
    else:
        raise RuntimeError("No supported accelerator (CUDA/XPU) available")
    return dispatch_key

# 随机种子设置中新增 XPU 支持
def set_random_seed(seed: int) -> None:
    random.seed(seed)
    np.random.seed(seed)
    torch.manual_seed(seed)
    if torch.cuda.is_available():
        torch.cuda.manual_seed_all(seed)
    # 添加 XPU 种子设置
    if torch.xpu.is_available():
        torch.xpu.manual_seed_all(seed)

# 设备核心数获取中新增 XPU 支持
def get_device_core_count(device_id: int = 0) -> int:
    if (hasattr(torch, "cuda") and torch.cuda.is_available()) or is_musa():
        return torch.cuda.get_device_properties(device_id).multi_processor_count
    elif hasattr(torch, "xpu") and torch.xpu.is_available():
        return torch.xpu.get_device_properties(device_id).gpu_eu_count
    return 0
```

### python/sglang/srt/layers/layernorm.py

RMSNorm 的 forward\_xpu 方法新增了确定性模式分支，使 XPU 在启用 batch invariant 时能使用确定性算子。

```
def forward_xpu(
    self,
    x: torch.Tensor,
    residual: Optional[torch.Tensor] = None,
    post_residual_addition: Optional[torch.Tensor] = None,
) -> Union[torch.Tensor, Tuple[torch.Tensor, torch.Tensor]]:
    if self.variance_size_override is not None:
        return self.forward_native(x, residual, post_residual_addition)
    # 新增：当启用 batch invariant 模式时，使用确定性计算
    if is_batch_invariant_mode_enabled():
        # 若存在 residual 或为 FSDP 训练，回退到 native 实现
        if (
            residual is not None
            or get_global_server_args().rl_on_policy_target == "fsdp"
        ):
            return self.forward_native(x, residual, post_residual_addition)
        # 否则使用 batch invariant 的 rmsnorm，保证结果确定性
        return rms_norm_batch_invariant(
            x,
            self.weight.data,
            self.variance_epsilon,
        )
    # 原有 XPU 优化路径……
    if residual is not None:
        # ...
```

[python/sclang/srt/batch\\_invariant\\_ops/batch\\_invariant\\_ops.py](#)

注册确定性算子的核心文件，通过动态 dispatch key 支持 XPU。

```
# 导入新工具函数
from sclang.srt.utils.common import (
    calc_diff,
    get_bool_env_var,
    get_device_core_count,
    get_dispatch_device_backend,
)

def enable_batch_invariant_mode(enable_bmm: bool = True):
    # ...
    dispatch_key = get_dispatch_device_backend()
    # 注册时使用动态 dispatch_key 替代硬编码的 "CUDA"
    _batch_invariant_LIB.impl("aten::mm", mm_batch_invariant, dispatch_key)
    _batch_invariant_LIB.impl("aten::addmm", addmm_batch_invariant, dispatch_key)
    _batch_invariant_LIB.impl("aten::_log_softmax", _log_softmax_batch_invariant, dispatch_key)
    _batch_invariant_LIB.impl("aten::mean.dim", mean_batch_invariant, dispatch_key)
    if enable_bmm:
```

```
_batch_invariant_LIB.impl("aten::bmm", bmm_batch_invariant, dispatch_key)
```

## 评论区精华

- gemini-code-assist[bot]: 指出在 layernorm.py 的 forward\_xpu 中, 对 self.forward\_native 的调用缺少 \*\*kwargs 参数, 建议修复以提高与 forward\_cuda 的一致性。该问题在后续提交中已修正。
- mingfeima: 要求添加测试覆盖多个后端。作者 jthakurH 回应已有测试 test\_deterministic, 并在此 PR 后能在 XPU 通过。此外, mingfeima 建议与 issue #15299 对齐, 并建议先提交到 intel\_dev 分支。作者随后创建了 PR #17915。
- polisettyvarma: 多次催促 review, 最终 mingfeima 批准。
  - forward\_xpu 中 forward\_native 调用缺少 \*\*kwargs (correctness): 该问题在后续提交中已修复, 最终代码包含了 post\_residual\_addition 参数。
  - 需要测试覆盖多个后端 (testing): mingfeima 接受了该回应, 最终 PR 被批准。但本次 PR 未附带新测试文件。

## 风险与影响

- 风险:
  - 兼容性风险: 新增的 get\_dispatch\_device\_backend 可能影响未来设备后端新增时的调试, 但当前逻辑清晰, 风险较低。
  - 性能风险: 在 layernorm 中添加了分支判断 is\_batch\_invariant\_mode\_enabled(), 会增加少量运行时开销, 但仅在模式启用时生效, 对默认行为无影响。
  - 测试覆盖: 虽然作者声称有测试, 但本次 PR 未附带新测试文件, 且评论中 mingfeima 也要求更多测试覆盖, 现有测试可能不足以覆盖全部 XPU 场景。
  - 回归风险: batch\_invariant\_ops 中将 dispatch key 从固定 "CUDA" 改为动态, 若其他后端 (CUDA) 行为异常, 可能需要回退。
- 影响:
  - 用户: XPU 用户现在可以启用 --enable-batch-invariant 以确保推理结果完全确定, 有利于调试和精度验证。
  - 系统: 无侵入性变更, 现有 CUDA 路径行为保持不变。增加了对 XPU 的显式支持, 但未改变其他后端。
  - 团队: 维护者需要留意多后端统一抽象, 未来新增后端时需在后端 get\_dispatch\_device\_backend 中添加对应分支。
  - 风险标记: 缺少 XPU 全面测试覆盖, CUDA 路径可能受影响 (回归风险), 性能微增 (分支判断)

## 关联脉络

- PR #17915 Add deterministic mode for XPU operations: 同一功能线的分支提交, 旨在先并入 intel\_dev 分支后再合入 main, 展示了团队的 upstream 策略。