

PR #26496 完整报告

sgl-project/sglang

Changes for SM120 perf and usability for NVFP4

合并时间: 2026-06-05 06:29

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/26496>

执行摘要

- 一句话: SM120 NVFP4 性能与可用性优化
- 推荐动作: 值得精读, 该 PR 展示了针对特定硬件 (SM120) 进行系统性性能优化的典型方法: 从后端选择、autotune 触发、kernel 配置到量化修复, 覆盖了整个推理链路。设计权衡 (如后端切换原因、配置一致性处理) 有参考价值。建议重点关注 `_should_run_flashinfer_autotune` 和 `try_get_optimal_moe_config` 的变更逻辑。

功能与动机

根据 issue #19637 (SM120 Performance Optimization Plan), 社区对 SM120 上 NVFP4 模型的性能与功能完善有迫切需求。该 PR 旨在修复已知问题、优化后端选择策略并调整内核配置, 以提升推理吞吐和稳定性。

实现拆解

1. 后端选择策略调整: 在 `python/sglang/srt/layers/quantization/fp4_utils.py` 的 `initialize_fp4_gemm_config` 中, 移除了 SM120 上优先使用 `flashinfer_cudnn` 的逻辑, 改为回退到 `flashinfer_cutlass`, 解决 NaN 问题。
2. MoE 自动后端选择: 在 `python/sglang/srt/server_args.py` 的 `_handle_moe_kernel_config` 中, 当 `quantization=modelopt_fp4` 且设备为 SM120 时, 将 `moe_runner_backend` 设为 `flashinfer_cutlass`, 覆盖默认的 `flashinfer_trtllm` (后者仅支持 SM100)。
3. 扩展 FlashInfer autotune 覆盖范围: 在 `python/sglang/srt/model_executor/model_runner.py` 的 `_should_run_flashinfer_autotune` 中, 新增 `fp4_gemm_needs_autotune` 分支, 使 NVFP4 GEMM 在 FlashInfer CUTLASS/CuteDSL 后端上也能触发 autotune。
4. MoE 配置一致性放宽: 在 `python/sglang/srt/layers/moe/moe_runner/triton_utils/fused_moe_triton_config.py` 的 `try_get_optimal_moe_config` 中, 将 `down_moe` 的 `BLOCK_SIZE_M` 硬断言改为警告并自动覆盖为 `up` 配置的值, 避免因配置不匹配导致崩溃。
5. 新增 SM120 特定 MoE 调优配置: 为 `NVIDIA_RTX_PRO_6000_Blackwell_Server_Edition` 新增 `up/down` 两个 JSON 配置文件, 通过精细的 `BLOCK_SIZE_M/N/K`、`GROUP_SIZE_M`、`num_warps`、`num_stages` 参数提升 Triton MoE kernel 在 SM120 上的执行效率。

6. 禁用 DeepGEMM 避免误用：在 `python/sglang/srt/layers/deep_gemm_wrapper/configurer.py` 中，将 `DEEPGEMM_BLACKWELL` 门限从 `is_blackwell_supported` 收窄为 `is_sm100_supported`，防止在 SM120 上产生错误警告。
7. AWQ 跳过层修复：在 `python/sglang/srt/layers/quantization/awq/awq.py` 中添加条件，当层属于 `modules_to_not_convert` 时跳过 MoE 量化，修复了之前 AWQ 量化可能错误应用于不应转换层的问题。

关键文件：

- `python/sglang/srt/model_executor/model_runner.py`（模块 模型运行器；类别 source；类型 core-logic；符号 `_should_run_flashinfer_autotune`）：核心调度路径，扩展 autotune 判断逻辑以包含 FP4 GEMM，确保 NVFP4 模型也能触发 FlashInfer autotune。
- `python/sglang/srt/layers/moe/moe_runner/triton_utils/fused_moe_triton_config.py`（模块 MoE 配置；类别 source；类型 core-logic；符号 `try_get_optimal_moe_config`）：调整 `down_moe` 配置一致性处理方式，将硬断言改为 `warning + override`，提升鲁棒性。
- `python/sglang/srt/server_args.py`（模块 服务器参数；类别 source；类型 core-logic；符号 `_handle_moe_kernel_config`）：控制 MoE 后端的自动选择，在 SM120 上为 `modelopt_fp4` 选择 `flashinfer_cutlass`。
- `python/sglang/srt/layers/quantization/fp4_utils.py`（模块 量化工具；类别 source；类型 core-logic；符号 `initialize_fp4_gemm_config`）：调整 NVFP4 GEMM 后端自动选择，移除 `flashinfer_cudnn` 特例，回落至 `flashinfer_cutlass`。
- `python/sglang/srt/layers/deep_gemm_wrapper/configurer.py`（模块 DeepGEMM 配置；类别 source；类型 core-logic）：防止 DeepGEMM 在 SM120 上误用，收窄启用门限。

关键符号：`_should_run_flashinfer_autotune`, `try_get_optimal_moe_config`,
`initialize_fp4_gemm_config`, `_handle_moe_kernel_config`

关键源码片段

`python/sglang/srt/model_executor/model_runner.py`

核心调度路径，扩展 autotune 判断逻辑以包含 FP4 GEMM，确保 NVFP4 模型也能触发 FlashInfer autotune。

```
def _should_run_flashinfer_autotune(self) -> bool:
    """Check if flashinfer autotune should be run."""
    if self.server_args.disable_flashinfer_autotune:
        return False

    # CuteDSL v1 (cutedsdl runner + deepep a2a) bypasses MoeRunner and must not
    # be autotuned -- its _dummy_run would dispatch more tokens per rank than
    # SGLANG_DEEPEP_NUM_MAX_DISPATCH_TOKENS_PER_RANK, tripping a DeepEP assert.
    if (
        self.server_args.moe_runner_backend == "flashinfer_cutedsdl"
        and self.server_args.moe_a2a_backend == "deepep"
    ):
        return False
```

```

backend_str = self.server_args.moe_runner_backend

# 判断 MoE runner 是否需要 autotune
moe_needs_autotune = backend_str in [
    "flashinfer_trtllm",
    "flashinfer_trtllm_routed",
    "flashinfer_mxfp4",
    "flashinfer_cutedsl",
    "flashinfer_cutlass",
]

from sglang.srt.layers.quantization.fp4_utils import get_fp4_gemm_runner_backend

model_uses_fp4 = self.model_config.quantization in (
    "modelopt_fp4",
    "modelopt_mixed",
)
# 如果模型使用 NVFP4 且后端是 CUTLASS / CuteDSL, 也需要 autotune
fp4_gemm_needs_autotune = model_uses_fp4 and (
    get_fp4_gemm_runner_backend().is_flashinfer_cutlass()
    or get_fp4_gemm_runner_backend().is_flashinfer_cutedsl()
)

if not (moe_needs_autotune or fp4_gemm_needs_autotune):
    return False

major, _ = torch.cuda.get_device_capability()
if major < 9:
    return False

if self.spec_algorithm.is_speculative():
    return not self.is_draft_worker

return True

```

python/sglang/srt/layers/moe/moe_runner/triton_utils/fused_moe_triton_config.py

调整 down_moe 配置一致性处理方式, 将硬断言改为 warning + override, 提升鲁棒性。

```

def try_get_optimal_moe_config(...):
    # ... 前面的代码获取 config 和 down_config ...
    if return_down_config:
        if (
            down_config is not None
            and config["BLOCK_SIZE_M"] != down_config["BLOCK_SIZE_M"]
        ):
            # 两个 kernel 共享同一个 moe_align_block_size 排序, 因此
            # down 配置必须使用 up 配置的 BLOCK_SIZE_M。
            logger.warning_once(

```

```
"down_moe config BLOCK_SIZE_M=%d does not match up config "  
"BLOCK_SIZE_M=%d at M=%d; overriding down BLOCK_SIZE_M to match.",  
down_config["BLOCK_SIZE_M"],  
config["BLOCK_SIZE_M"],  
M,  
)  
down_config["BLOCK_SIZE_M"] = config["BLOCK_SIZE_M"]  
return config, (down_config, max_block_m)  
return config
```

评论区精华

PR 未产生 Review 讨论，仅由 Fridge003 审批通过。PR body 中作者提供了性能对比数据，展示了约 17% TPS 提升，充分验证了变更的有效性。

- 暂无高价值评论线程

风险与影响

- 风险：

1. 默认后端切换回归风险：将 SM120 上 NVFP4 GEMM 后端从 flashinfer_cudnn 替换为 flashinfer_cutlass，可能在新场景下出现数值或性能回退。虽然主要原因 (NaN) 已修复，但仍需关注覆盖不足的情况。
2. SM120 特定逻辑影响：新增的许多条件分支（如 is_sm120_supported()）只针对 Blackwell 设备，不会影响其他架构，但增加了代码路径复杂性。
3. 缺少测试覆盖：本次变更新增了多个条件分支和配置，但未发现配套新增的自动化测试。特别是 autotune 触发条件的变化和 MoE 配置覆盖逻辑的变更，若无测试可能遗漏回归。
4. MoE 配置覆盖的潜在副作用：强制将 down_moe 的 BLOCK_SIZE_M 覆盖为 up 配置的值，虽然避免崩溃，但可能会略降低 down 部分的性能，需要后续验证。- 影响：用户视角：SM120 (Blackwell) 上使用 NVFP4 量化的模型（如 Qwen3.6-27B-NVFP4）将获得约 17% 的端到端 TPS 提升。AWQ 量化修复使得部分模型不再错误量化应跳过的层。

系统视角：MoE 后端自动选择逻辑更精细，autotune 覆盖更全面，但后端切换可能引入新的兼容性边界。DeepGEMM 不会在 SM120 上误用。

团队视角：此次改动涉及多文件协作（量化、调度、MoE 配置），后续维护者需要理解 SM120 专用逻辑。

- 风险标记：默认后端切换回归风险，SM120 特定逻辑影响，缺少配套测试覆盖

关联脉络

- PR #25239 [FlashInfer v0.6.12] Support FlashInfer 4over6 NVFP4: 同属 NVFP4 功能线，该 PR 提供了 FlashInfer NVFP4 支持，本 PR 在此基础上优化 SM120 性能。
- PR #23979 Enable DeepGEMM PDL on by default: 同样涉及 DeepGEMM 和 SM100/SM120 的启用策略，本 PR 进一步收窄了 DeepGEMM 在 SM120 的启用条件。