

PR #25523 完整报告

sgl-project/sglang

[Diffusion] Default NVFP4 backend to FlashInfer TRTLLM

合并时间: 2026-05-25 18:14

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/25523>

执行摘要

- 一句话: 默认 NVFP4 后端切换为 FlashInfer TRTLLM, 支持 swizzled scale 布局
- 推荐动作: 该 PR 值得精读, 特别是 scale layout 的处理设计和默认后端的切换策略。建议重点关注 `_swizzled_nvfp4_scales_to_linear` 的实现以及 review 中暴露的潜在风险, 未来可能需统一后端间的 scale 处理路径。

功能与动机

PR 跟随 #25857 中 Wan2.2 NVFP4 checkpoint 的落地方案, 进一步将 diffusion ModelOpt NVFP4 的默认后端从 cudnn 改为 flashinfer_trtllm, 因为 trtllm 在 Blackwell GPU 上对常见形状 (如 FLUX.1-dev、Wan2.2) 具有显著性能优势 (见 PR 性能表格: FLUX.1-dev 生成时间降低约 34%)。此外, 官方 FLUX.2 NVFP4 导出的 weight scale 是 FlashInfer/CUTLASS-swizzled 布局, 而 SGLang 转换的 repos 保持线性布局, 需要新增布局追踪和转换逻辑以正确加载这类 checkpoint。

实现拆解

1. 变更默认后端选择逻辑

修改 `python/sglang/multimodal_gen/runtime/platforms/cuda.py` 中的 `get_modelopt_flashinfer_fp4_backend`, 将默认后端从 `cudnn` (Blackwell) 或 `auto` 统一改为 `trtllm`。同时重写 `get_modelopt_fp4_gemm_op`, 移除旧版的有条件优先逻辑, 直接优先尝试 `flashinfer.mm_fp4`, 失败则回退到 `sgl_kernel.cutlass_scaled_fp4_mm`。

2. 新增 swizzled scale layout 支持

在 `python/sglang/multimodal_gen/runtime/layers/quantization/modelopt_quant.py` 中新增 `_swizzled_nvfp4_scales_to_linear` 函数, 将 FlashInfer/CUTLASS 的 swizzled layout 转换回 row-major 线性布局。在 `ModelOptFp4Config` 中添加 `checkpoint_weight_scale_layout` 字段 (默认 `linear`), 并在 `process_weights_after_loading` 中根据该字段决定是否对 `weight_scale` 进行转换。

3. 配置合并与推断逻辑

修改 `transformer_load_utils.py` 中的 `_merge_modelopt_fp4_configs`, 当两个配置的 `checkpoint_weight_scale_layout` 不同时优先保留非 `linear` 的值。修改 `quantization_utils.py` 中的 `_build_nvfp4_config_from_safetensors`, 在检测到 packed QKV

时自动推断 layout 为 `swizzled`。

4. 测试与文档

新增测试 `test_flux2_swizzled_scale_checkpoint_flashinfer_trtllm_matches_cudnn`，验证从 `swizzled scale checkpoint` 加载后，`flashinfer_trtllm` 输出与 `cudnn` 后端一致。更新 `docs_new/docs/sglang-diffusion/quantization.mdx` 和 `docs/diffusion/quantization.md` 以反映默认后端变更；更新 `envs.py` 和对应文档注释。

关键文件：

- `python/sglang/multimodal_gen/runtime/layers/quantization/modelopt_quant.py`（模块 量化层；类别 `source`；类型 `data-contract`；符号 `_swizzled_nvfp4_scales_to_linear`, `ModelOptFp4Config.checkpoint_weight_scale_layout`, `ModelOptFp4LinearMethod.process_weights_after_loading`）：核心变更：新增 `swizzled scale` 转换函数、配置字段和 `weights` 加载逻辑
- `python/sglang/multimodal_gen/runtime/platforms/cuda.py`（模块 平台层；类别 `source`；类型 `dependency-wiring`；符号 `get_modelopt_flashinfer_fp4_backend`, `get_modelopt_fp4_gemm_op`）：后端选择逻辑核心：更改默认后端为 `trtllm` 并简化 GEMM op 获取流程
- `python/sglang/jit_kernel/tests/diffusion/test_diffusion_nvfp4_scaled_mm.py`（模块 NVFP4 测试；类别 `test`；类型 `test-coverage`；符号 `test_flux2_swizzled_scale_checkpoint_flashinfer_trtllm_matches_cudnn`）：新增端到端测试，验证 `swizzled scale checkpoint` 在 `trtllm` 与 `cudnn` 后端输出一致
- `python/sglang/multimodal_gen/runtime/loader/transformer_load_utils.py`（模块 模型加载器；类别 `source`；类型 `core-logic`；符号 `_merge_modelopt_fp4_configs`）：合并两处配置源时保留非 `linear` 的 `scale layout`
- `python/sglang/multimodal_gen/runtime/utils/quantization_utils.py`（模块 量化工具；类别 `source`；类型 `core-logic`；符号 `_build_nvfp4_config_from_safetensors`）：从 `safetensors` 推断 `scale layout` 并设置 `checkpoint_weight_scale_layout`
- `python/sglang/multimodal_gen/envs.py`（模块 环境变量；类别 `source`；类型 `documentation`）：更新环境变量注释以反映默认行为
- `docs_new/docs/sglang-diffusion/quantization.mdx`（模块 文档（新）；类别 `other`；类型 `documentation`）：更新文档以匹配新的默认后端

关键符号：`_swizzled_nvfp4_scales_to_linear`, `get_modelopt_flashinfer_fp4_backend`, `get_modelopt_fp4_gemm_op`, `ModelOptFp4LinearMethod.process_weights_after_loading`, `_merge_modelopt_fp4_configs`, `_build_nvfp4_config_from_safetensors`, `test_flux2_swizzled_scale_checkpoint_flashinfer_trtllm_matches_cudnn`

关键源码片段

```
python/sglang/multimodal_gen/runtime/layers/quantization/modelopt_quant.py
```

核心变更：新增 `swizzled scale` 转换函数、配置字段和 `weights` 加载逻辑

```

# python/sglang/multimodal_gen/runtime/layers/quantization/modelopt_quant.py

# 新增函数: 将 FlashInfer/CUTLASS 的 swizzled FP4 scales 转换回 row-major 线性布局
# 用于加载官方 FLUX.2 等以 swizzled 格式存储的 checkpoint
def _swizzled_nvfp4_scales_to_linear(scales: torch.Tensor) -> torch.Tensor:
    scale_ndim = scales.ndim
    if scale_ndim == 2:
        scales = scales.unsqueeze(0) # 增加 batch 维度使处理统一
    assert scales.ndim == 3
    B, M, K = scales.shape
    M_padded = round_up(M, 128) # 对齐到 128 的倍数
    K_padded = round_up(K, 4) # 对齐到 4 的倍数
    if M != M_padded or K != K_padded:
        padded = torch.zeros((B, M_padded, K_padded), dtype=scales.dtype, device=scales.device)
        padded[:, :M, :K] = scales
        scales = padded
    # 反 swizzle: reshape + permute 还原线性布局
    linear = scales.reshape(B, M_padded // 128, K_padded // 4, 32, 4, 4)
    linear = linear.permute(0, 1, 4, 3, 2, 5).contiguous()
    linear = linear.reshape(B, M_padded, K_padded)[: , :M, :K]
    return linear.squeeze(0) if scale_ndim == 2 else linear

# 在 process_weights_after_loading 中, 先根据配置转换 scale, 再根据后端处理
scales = layer.weight_scale
if getattr(self.quant_config, 'checkpoint_weight_scale_layout', 'linear') == 'swizzled':
    scales = _swizzled_nvfp4_scales_to_linear(scales)
_, flashinfer_backend = _get_fp4_gemm_op()
if flashinfer_backend == 'trtllm':
    # trtllm 后端需要额外 padding 和 shuffle
    weight, _ = pad_nvfp4_weight(w_swapped, n_alignment=128, k_alignment=0)
    if scales.shape[0] != weight.shape[0]:
        pad_n = weight.shape[0] - scales.shape[0]
        scales = torch.nn.functional.pad(scales, (0, 0, 0, pad_n))
    # ... 应用 shuffle 并写入 layer.weight_scale_interleaved

```

python/sglang/multimodal_gen/runtime/platforms/cuda.py

后端选择逻辑核心: 更改默认后端为 trtllm 并简化 GEMM op 获取流程

```

# python/sglang/multimodal_gen/runtime/platforms/cuda.py

@classmethod
@lru_cache(maxsize=1)
def get_modelopt_flashinfer_fp4_backend(cls) -> str:
    backend = envs.SGLANG_DIFFUSION_FLASHINFER_FP4_GEMM_BACKEND
    default_backend = 'trtllm' # PR 核心变更: 默认后端统一为 trtllm
    if backend is None:
        return default_backend
    backend = backend.lower()
    # 支持多个别名

```

```

backend = {
    'flashinfer_cudnn': 'cudnn',
    'flashinfer_cutlass': 'cutlass',
    'flashinfer_trtllm': 'trtllm',
    'trtllm': 'trtllm',
    'cudnn': 'cudnn',
    'auto': 'auto',
}.get(backend, backend)
if backend not in {'auto', 'cudnn', 'cutlass', 'trtllm'}:
    logger.warning('Unsupported backend %r, falling back to %r', backend, default_backend)
    return default_backend
return backend

```

```

@classmethod
@lru_cache(maxsize=1)
def get_modelopt_fp4_gemm_op(cls) -> tuple[Callable | None, str | None]:
    requested_backend = envs.SGLANG_DIFFUSION_FLASHINFER_FP4_GEMM_BACKEND
    try:
        from flashinfer import mm_fp4 as flashinfer_mm_fp4
        return flashinfer_mm_fp4, cls.get_modelopt_flashinfer_fp4_backend()
    except ImportError:
        logger.warning(
            'flashinfer.mm_fp4 unavailable, falling back to cutlass (requested: %r)',
            requested_backend or 'flashinfer_trtllm (default)',
        )
    try:
        from sgl_kernel import cutlass_scaled_fp4_mm as cutlass_fp4_gemm
        return cutlass_fp4_gemm, None
    except ImportError:
        return None, None

```

评论区精华

review 评论: scale un-swizzling 逻辑应移至后端通用位置

[gemini-code-assist\[bot\]](#) 指出, 当前 swizzled scale 的转换在 `process_weights_after_loading` 中只在进入 `trtllm` 后端分支前执行, 但转换后的 `scales` 仅用于 `trtllm` 路径; 对于其他后端 (`cudnn`、`cutlass`、`sgl_kernel` 回退), `layer.weight_scale` 仍保留 swizzled 布局, 可能导致错误。建议将转换逻辑提前到所有后端共享部分, 并更新 `layer.weight_scale` 参数。该评论未得到显式回复或关闭, PR 仍然合并。

- swizzled scale 转换逻辑仅对 `trtllm` 后端生效, 其他后端可能出错 (correctness): PR 合并时该问题未被解决; 当前实现仅转换局部 `scales` 变量, 仅 `trtllm` 路径使用, 其他后端仍使用原始 `layer.weight_scale`。

风险与影响

- 风险:

1. 后向兼容性风险：已有用户依赖 cudnn 后端时未设置环境变量，切换默认值后可能面临不同的数值行为或性能变化（但性能应更好）。若用户希望保留原后端需显式设置 `SGLANG_DIFFUSION_FLASHINFER_FP4_GEMM_BACKEND=cudnn`。
 2. swizzled scale 覆盖不全：根据 review 评论，非 trtllm 后端（如 cudnn 或 cutlass）遇到 swizzled layout 的检查点时，`process_weights_after_loading` 中转换后的 `scales` 未写回 `layer.weight_scale`，可能导致这些后端加载失败或产生错误结果。需要确认其他后端是否也期望 linear 布局。
 3. 依赖 flashinfer：新默认路径依赖 `flashinfer.mm_fp4`，若用户环境未安装 flashinfer 或版本不兼容，会自动回退到 cutlass。回退路径的数值一致性已通过测试验证，但性能可能下降。
 - 影响：用户影响：使用 diffusion NVFP4 的用户（FLUX、Wan2.2 等）将默认获得 `flashinfer_trtllm` 后端的性能提升。需关注兼容 checkpoints 的 layout 标记，若加载官方 FLUX.2 导出可自动识别 swizzled layout。
 - 系统影响：无直接系统级变更。
 - 团队影响：维护成本增加，需跟踪不同后端对 scale layout 的处理一致性。
- 风险标记：后向兼容性风险，swizzled scale 未覆盖非 trtllm 后端，依赖 flashinfer 导致回退可能

关联脉络

- PR #25857 [diffusion] NVFP4 checkpoint support for Wan2.2: 本 PR 紧跟 #25857 的 checkpoint 落地方案，在此基础上更改默认后端并增加 swizzled scale 处理。