

# PR #25418 完整报告

sgl-project/sglang

integrate flash\_mla\_sparse\_fwd

合并时间: 2026-06-03 16:09

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/25418>

## 执行摘要

- 一句话: 集成 flash\_mla\_sparse\_fwd 加速 DSv4 预填并修复长序列 chunk prefill 错误
- 推荐动作: 此 PR 核心价值显著, 性能改进已用 benchmark 验证。建议仔细审查 `_forward_prefill_sparse` 的缓存管理逻辑, 确保跨层一致性; 并考虑将特性默认开启以获取更多反馈。值得关注的设计决策包括: 全量反量化 vs 选择性反量化、int64 索引转换、以及阈值硬编码的后续优化。

## 功能与动机

PR 描述指出 flash\_mla\_with\_kvcache 加载逻辑复杂导致性能瓶颈, 切换到 flash\_mla\_sparse\_fwd 可获 1.35x 内核加速和整体 1.1x 提升; 同时 chunk prefill 仅支持  $\leq 8192$ , 在 32768 时失败 (关联 issue #25484), 需要路由到稀疏内核绕过。

## 实现拆解

1. 新增 `dequant_k_cache.py`: 编写 Triton 内核 `_dequantize_k_cache_paged_kernel`, 将分页 KV 缓存中 fp8 K nope + bf16 rope 反量化为 bf16 平面工作空间 (每 token 512 元素), 支持直接喂给 flash\_mla\_sparse\_fwd。
2. 新增 `sparse_prefill_utils.py`: 从 vllm 适配 `combine_topk_swa_indices` 内核, 将 query 的 topk 压缩索引和 SWA 位置索引合并为单行索引; 配套 `build_swa_token_ids` 内核构建 SWA 令牌 ID; 定义 `SparsePrefillChunkCache` 在首个稀疏层缓存跨层不变的元数据 (如 dequant 后的工作空间和索引), 后续层直接复用。
3. 修改 `deepseek_v4_backend.py`: 在 `DSV4AttnMetadata` 中添加 `c4_sparse_raw_indices` 字段; `init_flashmla_related` 接受 `is_prefill` 参数, 在预填时额外分配原始索引; 新增 `_forward_prefill_sparse` 方法, 调用 `dequantize_k_cache_paged` 和 `combine_topk_swa_indices` 准备输入, 然后调用 `flash_mla_sparse_fwd`; 主 `forward` 方法根据 `forward_mode.is_extend` 和序列长度路由到稀疏或常规路径。
4. 修改 `metadata.py`: 添加 `_LARGE_INDEXER_QUERY_THRESHOLD = 11673`, 当 `c4_seq_lens.numel()` 超过该阈值时强制使用 JIT 索引器元数据, 避免 `deep_gemm` 在大 batch 时共享内存不足。
5. 修改 `indexer.py`: 在 C4 索引器路径中, 若 `core_metadata.c4_sparse_raw_indices` 非空则直接使用原始索引, 避免重新计算。

6. 修改 `deepseek_v4.py` 和 `deepseek_v4_nextn.py`: 在预填前调用 `init_flashmla_related(is_prefill=True)` 以分配稀疏预填所需的额外中间张量。
7. 新增 `environ.py` 配置: 添加 `SGLANG_OPT_FLASHMLA_SPARSE_PREFILL` 环境变量控制开关, 默认关闭, 需显式启用。

关键文件:

- `python/sglang/srt/layers/attention/dsv4/sparse_prefill_utils.py` (模块 稀疏预填; 类别 source; 类型 core-logic; 符号 `combined_topk_width`, `combine_topk_swa_indices`, `build_swa_token_ids`, `_build_swa_token_ids_kernel`): 新增的核心文件: 实现稀疏预填索引组合器, 包括 `combine_topk_swa_indices` Triton 内核和 `SparsePrefillChunkCache` 缓存类, 是 `flash_mla_sparse_fwd` 的前置数据准备。
- `python/sglang/srt/layers/attention/dsv4/dequant_k_cache.py` (模块 反量化; 类别 source; 类型 core-logic; 符号 `dequantize_k_cache_paged`, `_dequantize_k_cache_paged_kernel`, `dequantize_k_cache_paged_ref`): 新增文件: 提供分页 KV cache 的反量化 Triton 内核, 将 `fp8 nope + bf16 rope` 输出为 `bf16` 平面张量, 供 `flash_mla_sparse_fwd` 使用。
- `python/sglang/srt/layers/attention/deepseek_v4_backend.py` (模块 注意力后端; 类别 source; 类型 core-logic; 符号 `init_flashmla_related`, `_forward_prefill_sparse`): 注意力后端的核修改: 集成 `_forward_prefill_sparse` 方法, 并在 `DSV4AttnMetadata` 中添加字段以支持稀疏预填。
- `python/sglang/srt/layers/attention/dsv4/metadata.py` (模块 元数据; 类别 source; 类型 core-logic): 添加了大查询阈值 `_LARGE_INDEXER_QUERY_THRESHOLD` 和自动切换 JIT indexer 的逻辑, 避免 `deep_gemm` 在大 batch 时崩溃。
- `python/sglang/srt/layers/attention/dsv4/indexer.py` (模块 索引器; 类别 source; 类型 core-logic): 修改 C4 索引器, 在稀疏预填模式下直接使用 `core_metadata` 中的原始索引。
- `python/sglang/srt/models/deepseek_v4.py` (模块 模型; 类别 source; 类型 data-contract): 模型前向函数中调用 `init_flashmla_related(is_prefill=True)` 以分配稀疏预填所需中间张量。
- `python/sglang/srt/models/deepseek_v4_nextn.py` (模块 NextN 模型; 类别 source; 类型 data-contract): 与 `deepseek_v4.py` 相同, 在 nextn 模型中同步修改调用。
- `python/sglang/srt/envIRON.py` (模块 环境变量; 类别 source; 类型 core-logic): 添加 `SGLANG_OPT_FLASHMLA_SPARSE_PREFILL` 环境变量控制开关。

关键符号: `combined_topk_width`, `combine_topk_swa_indices`, `build_swa_token_ids`, `SparsePrefillChunkCache.build`, `dequantize_k_cache_paged`, `_dequantize_k_cache_paged_kernel`, `init_flashmla_related`, `_forward_prefill_sparse`

## 关键源码片段

`python/sglang/srt/layers/attention/dsv4/sparse_prefill_utils.py`

新增的核心文件: 实现稀疏预填索引组合器, 包括 `combine_topk_swa_indices` Triton 内核和 `SparsePrefillChunkCache` 缓存类, 是 `flash_mla_sparse_fwd` 的前置数据准备。

```
"""Per-query sparse-index combiner for the FlashMLA sparse prefill path.
```

Adapts vllm's ``combine\_topk\_swa\_indices`` to sglang's flat-workspace layout.

Reference:

[https://github.com/vllm-project/vllm/blob/124fac10cb0ea83aee2ffeabac0b413d6b759b26/vllm/models/deepseek\\_v4/common/ops/cache\\_utils.py#L476](https://github.com/vllm-project/vllm/blob/124fac10cb0ea83aee2ffeabac0b413d6b759b26/vllm/models/deepseek_v4/common/ops/cache_utils.py#L476)

For each query token in a prefill chunk, emits one row of combined indices into the chunk's bf16 KV workspace:

```
[ topk indices into compressed cache (rebased) ]
[ swa positional indices (rebased)           ]
[ -1 padding up to a multiple of 128       ]
```

The workspace is a single flat ``(total\_workspace\_tokens, 512)`` tensor formed by concatenating, per request, that request's compressed-region gather followed by all requests' SWA-region gathers.

"""

```
from dataclasses import dataclass, field
from typing import Optional
```

```
import torch
import triton
import triton.language as tl
```

```
from sglang.srt.layers.attention.dsv4.dequant_k_cache import DIM_NOPE, DIM_ROPE
from sglang.srt.utils import ceil_align
```

```
# FlashMLA sparse prefill asserts ``params.topk % B_TOPK == 0``.
# B_TOPK is 64 for h_q=64 and 128 for h_q=128; pad to 128 to satisfy both.
SPARSE_PREFILL_TOPK_ALIGNMENT = 128
# Bf16 workspace per-token width, matching ``dequantize_k_cache_paged``'s
# output: 448 fp8 nope (dequanted) + 64 bf16 rope = 512.
WORKSPACE_DIM = DIM_NOPE + DIM_ROPE
```

```
def combined_topk_width(topk: int, window_size: int) -> int:
    """Width of the padded combined_indices last dim that
    ``combine_topk_swa_indices`` would produce for these args."""
    return ceil_align(topk + window_size, SPARSE_PREFILL_TOPK_ALIGNMENT)
```

```
def combine_topk_swa_indices(
    topk_indices: torch.Tensor,
    query_start_loc: torch.Tensor,
    seq_lens: torch.Tensor,
    gather_lens: torch.Tensor,
    compressed_base: torch.Tensor,
    swa_base: torch.Tensor,
    window_size: int,
```

```

compress_ratio: int,
topk: int,
out_indices: Optional[torch.Tensor] = None,
out_lens: Optional[torch.Tensor] = None,
) -> tuple[torch.Tensor, torch.Tensor]:
    """Combine topk + SWA indices into a single ``flash_mla_sparse_fwd`` row.

Args:
    topk_indices: (num_tokens, K) int32, per-query indices into compressed cache
        (request-local).
    query_start_loc: (num_reqs+1,) int32, cumulative query lengths.
    seq_lens: (num_reqs,) int32, full sequence lengths.
    gather_lens: (num_reqs,) int32, trailing tokens dequanted into SWA region.
    compressed_base: (num_reqs,) int32, flat offset for compressed region.
    swa_base: (num_reqs,) int32, flat offset for SWA region.
    window_size: SWA window size.
    compress_ratio: compress ratio (>=1 even when topk==0).
    topk: configured topk; 0 for SWA-only layers.
    out_indices, out_lens: optional preallocated buffers.
Returns:
    combined_indices: (num_tokens, padded_topk_swa) int32.
    combined_lens: (num_tokens,) int32.
    """
    assert topk_indices.dtype == torch.int32
    # ... kernel launch logic

```

### python/sglang/srt/layers/attention/dsv4/dequant\_k\_cache.py

新增文件：提供分页 KV cache 的反量化 Triton 内核，将 fp8 nope + bf16 rope 输出为 bf16 平面张量，供 flash\_mla\_sparse\_fwd 使用。

```

from typing import Optional

import torch
import triton
import triton.language as tl

from sglang.srt.layers.quantization.fp8_kernel import is_fp8_fnuz

fp8_dtype = torch.float8_e4m3fnuz if is_fp8_fnuz() else torch.float8_e4m3fn

# v4 KV cache layout (see dsv4.index_buf_accessor._set_k_and_s_triton_kernel):
# per-token: 448 fp8 nope + 64 bf16 rope (= 576 contiguous bytes) +
# 7 ue8m0 scales padded to 8 bytes.
# per-page: [token 0..P-1 nope+rope (P*576 bytes)] [token 0..P-1 scale (P*8 bytes)]
# padded up to a multiple of 576.
DIM_NOPE = 448
DIM_ROPE = 64
TILE_SIZE = 64 # one nope scale tile = 64 fp8 values
NUM_SCALE_TILES = DIM_NOPE // TILE_SIZE # 7

```

```
NOPE_ROPE_BYTES = DIM_NOPE + DIM_ROPE * 2 # 576
PADDED_SCALE_PER_TOKEN = NUM_SCALE_TILES + 1 # 8
```

```
def dequantize_k_cache_paged(
    quant_k_cache: torch.Tensor,
    page_table_1_flattened: torch.Tensor,
    page_size: int,
    out: Optional[torch.Tensor] = None,
) -> torch.Tensor:
    """Dequantize the DeepSeek v4 paged KV cache for a list of token IDs.
```

Args:

```
    quant_k_cache: (num_pages, bytes_per_page_padded) uint8.
    page_table_1_flattened: (num_tokens,) int — token IDs into the cache.
    page_size: number of tokens per page.
    out: optional (num_tokens, 1, DIM_NOPE + DIM_ROPE) bf16 destination.
```

Returns:

```
(num_tokens, 1, DIM_NOPE + DIM_ROPE) bfloat16.
```

```
"""
```

```
assert quant_k_cache.is_contiguous()
assert page_table_1_flattened.dtype in (torch.int32, torch.int64)
```

```
# The buffer's dtype is whatever the pool exposes (often bf16);
# reinterpret to byte-space first.
quant_k_cache_u8 = quant_k_cache.view(torch.uint8)
num_tokens = page_table_1_flattened.shape[0]
bytes_per_page = quant_k_cache_u8.shape[-1]
s_offset_bytes = page_size * NOPE_ROPE_BYTES
```

```
# Three typed views over the same underlying bytes.
buf_fp8 = quant_k_cache_u8.view(fp8_dtype).reshape(-1)
buf_bf16 = quant_k_cache_u8.view(torch.bfloat16).reshape(-1)
buf_uint8 = quant_k_cache_u8.reshape(-1)
```

if out is None:

```
    out = torch.empty(
        (num_tokens, 1, DIM_NOPE + DIM_ROPE),
        dtype=torch.bfloat16,
        device=quant_k_cache.device,
    )
```

else:

```
    assert out.shape == (num_tokens, 1, DIM_NOPE + DIM_ROPE)
    assert out.dtype == torch.bfloat16
```

```
_dequantize_k_cache_paged_kernel[(num_tokens,)]( # launch one program per token
    out,
    buf_fp8,
```

```

    buf_bf16,
    buf_uint8,
    page_table_1_flattened,
    out.stride(0),
    BYTES_PER_PAGE=bytes_per_page,
    PAGE_SIZE=page_size,
    DIM_NOPE=DIM_NOPE,
    DIM_ROPE=DIM_ROPE,
    TILE_SIZE=TILE_SIZE,
    NUM_SCALE_TILES=NUM_SCALE_TILES,
    NOPE_ROPE_BYTES=NOPE_ROPE_BYTES,
    PADDED_SCALE_PER_TOKEN=PADDED_SCALE_PER_TOKEN,
    S_OFFSET_BYTES=s_offset_bytes,
)
return out

```

## python/sglang/srt/layers/attention/deepseek\_v4\_backend.py

注意力后端的核心修改：集成 `_forward_prefill_sparse` 方法，并在 `DSV4AttnMetadata` 中添加字段以支持稀疏预填。

# 在 `DSV4AttnMetadata` 类中新增字段

```

class DSV4AttnMetadata:
    # ... existing fields ...
    c4_sparse_raw_indices: Optional[torch.Tensor] = field(init=False, default=None)

```

# 修改 `init_flashmla_related` 以支持预填

```

def init_flashmla_related(self, is_prefill: bool = False):
    assert self.c4_sparse_topk in (512, 1024)
    self.c4_sparse_page_indices = ... # 原有计算
    if is_prefill:
        # 稀疏预填需要额外的原始索引张量
        self.c4_sparse_raw_indices = torch.empty_like(self.c4_sparse_page_indices)
    self.c1_flashmla_metadata = _create_flashmla_metadata()
    self.c4_flashmla_metadata = _create_flashmla_metadata()
    self.c128_flashmla_metadata = _create_flashmla_metadata()

```

class `DSV4Metadata`:

```

    core_attn_metadata: DSV4AttnMetadata
    # 新增稀疏预填缓存（懒初始化）
    sparse_prefill_cache: Optional[SparsePrefillChunkCache] = None

```

```

def copy_(self, other: DSV4Metadata):
    # 每次 copy 时清空缓存，确保 cuda-graph 重播时重新构建
    self.sparse_prefill_cache = None

```

# 新增 `_forward_prefill_sparse` 方法（核心稀疏预填逻辑）

```

def _forward_prefill_sparse(self, q, layer_id, compress_ratio, forward_batch,
                            token_to_kv_pool, core_attn_metadata, attn_sink) -> torch.Tensor:
    """Unified prefill via flash_mla_sparse_fwd kernel."""

```

```

# 1. 获取或构建稀疏预填缓存
metadata = forward_batch.metadata
if metadata.sparse_prefill_cache is None:
    # 首次调用时构建缓存: 反量化所有页、构建索引
    workspace = dequantize_k_cache_paged(...) # 全量反量化
    indices, lengths = combine_topk_swa_indices(...)
    metadata.sparse_prefill_cache = SparsePrefillChunkCache(
        workspace=workspace, combined_indices=indices, combined_lens=lengths)
# 2. 从缓存获取数据并调用 flash_mla_sparse_fwd
cache = metadata.sparse_prefill_cache
out = flash_mla_sparse_fwd(q, cache.workspace, cache.combined_indices, ...)
return out

```

## 评论区精华

Review 中有几个关键讨论点:

- 反量化范围(DarkSharpness): 询问 `_forward_prefill_sparse` 是否会反量化所有 c4 cache 还是仅选定的。Fridge003 回答是全量, 因为选定的每层不同但稀疏预填缓存只计算一次, 需要作者确认。
- 硬编码阈值(Fridge003): 对 `metadata.py` 中 `_LARGE_INDEXER_QUERY_THRESHOLD=11673` 的硬编码提出质疑, 建议避免硬编码。该疑问未解决但未阻止合并。
- `numel` 使用(Fridge003): 在 `metadata.py` 中 `c4_seq_lens.numel() > threshold` 判断, 疑问 `numel` 应为 batch size 而非 query 数? 该点未解决。
- `int32` 溢出(Fridge003): 在 `dequant_k_cache.py` 中质疑 `tl.int32` 是否导致索引溢出。后续多次 commit 将索引张量改为 `int64` 修复。
- 添加参考实现与自测试(Fridge003): 要求在 `dequant_k_cache.py` 中添加 torch 参考实现和 `__main__` 自测试。已在后续 commit 中完成。
- 反量化范围 (correctness): 未完全解决, 全量反量化被接受, 但长期应优化为仅反量化需要的页。
- 硬编码阈值 (design): 未修改, 硬编码保留, 后续动态计算可能改进。
- `int32` 索引溢出 (correctness): 在后续 commit 中改为 `int64` 修复。
- 添加参考实现与自测试 (testing): 已在后续 commit 中添加了 `dequantize_k_cache_paged_ref` 和 `self-test`。
- 索引改为 `int64` (correctness): 已在后续 commit 中全部改为 `int64`。

## 风险与影响

- 风险:
  1. 全量反量化开销: 每次预填都会反量化所有 c4 缓存页 (而非仅选中部分), 在长上下文时可能增加内存带宽压力, 但 PR 声称 1.1x 加速已覆盖此成本。
  2. 稀疏预填缓存生命周期: `SparsePrefillChunkCache` 在首个稀疏层构建后跨层复用, 但若批次内请求的 KV cache 被修改 (如新的压缩令牌), 缓存可能失效。当前设计假设首次构建后的元数据在整个 chunk 内不变, 需通过测试验证。

3. int64 偏移一致性: 索引计算改为 int64 以修复 IMA, 但若空间偏移仍使用 int32 可能引入隐藏 bug。
4. 硬编码阈值: `_LARGE_INDEXER_QUERY_THRESHOLD` 的选取可能不适用于所有硬件配置, 如 SM90 的 H20 需验证。
5. 环境变量默认关闭: 特性默认未启用, 可能降低发现问题的几率。
6. 缺乏测试配套: 未直接包含端到端测试, 仅 self-test 在 `__main__` 中, 回归风险较高。
  - 影响: 对 DeepSeek V4 用户: 启用环境变量后预填性能提升约 1.1x, 长序列 (>8192) chunk prefill 不再崩溃。对系统: 新增两个 Triton 内核文件和约 800 行代码, 增加了 kernel 缓存占用; SparsePrefillChunkCache 占用额外显存。对开发团队: 需要维护两条预填路径; 后续 flash\_mla\_sparse\_fwd 若能全面优于现有路径可移除旧代码。
  - 风险标记: 全量反量化带宽开销, 稀疏预填缓存生命周期, int64 偏移一致性, 硬编码阈值适用性, 环境变量默认关闭, 缺少端到端测试

## 关联脉络

- PR #25502 Route long input seq to sparse prefill kernel: PR 描述中提到 cherry-pick 自该 PR, 用于路由长序列到稀疏预填内核以修复 chunk prefill bug。
- PR #25484 DeepSeek-V4-Pro on 8 \* H20-3e: DeepGEMM kernel exceeds shared memory limit: 动机中提到解决 chunk prefill 在 32768 失败的问题, 该 issue 记录了相关错误。