

PR #41869 完整报告

vllm-project/vllm

PD disagg with NIXL Connector: GDN support (Qwen3.5)

合并时间: 2026-05-14 22:33

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/41869>

执行摘要

- 一句话: 为 NIXL PD 分离添加 GDN(准 Qwen3.5) 支持
- 推荐动作: 值得精读, 特别是 MambaConvSplitInfo 的泛化模式, 展示了如何在保持向后兼容的同时扩展数据结构。derive_mamba_conv_split 中的异构 TP 推理逻辑值得参考。

功能与动机

关联 Issue #41886, 要求为 NIXL 的 Prefill/Decode 分离添加 GDN 模型支持。GDN 模型 (如 Qwen3.5) 具有不同的 SSM 布局: conv 状态是 [Q, K, V] 而非 [x, B, C], 时间状态形状为 (num_v_heads, v_dim, k_dim) 而非 (num_heads, head_dim)。需要扩展传输层以支持这种布局。

实现拆解

1. 泛化 MambaConvSplitInfo 数据结构 (ssm_conv_transfer_utils.py): 将固定字段 x_local/b_local 替换为 local_proj_dims: tuple[int,int,int], 将 x_bytes/b_bytes 属性重构为统一的 proj_bytes 属性, 并相应调整 local_conv_offsets 和 remote_conv_offsets 方法, 使其适用于 Mamba2 和 GDN 的子投影布局。
2. 扩展 derive_mamba_conv_split 函数: 允许 mamba_type 为 GDN_ATTN 和 MAMBA2, 根据时间状态形状推断子投影维度 (对于 GDN, 使用 num_v_heads、v_dim、k_dim 重建 key_dim 和 value_dim)。
3. 简化 NIXL worker 中的异构 TP 处理 (worker.py): 移除 _build_mamba_remote 中重复的异构 TP 偏移计算, 统一通过 self._conv_decomp.remote_conv_offsets(local_offset, tp_ratio) 获取偏移, 并支持负 tp_ratio 情况 (P_TP > D_TP 时的反向缩放)。
4. 测试与 CI 配套: 新增参数化单元测试 test_derive_mamba_conv_split 覆盖 Mamba2 和 GDN 在多种 TP 下的子投影维度计算; 在集成测试脚本中加入 Qwen3.5 配置, 并在 test_accuracy.py 中添加其精度阈值 0.33; 调整 CI 超时从 20 到 25 分钟。

关键文件:

- vllm/distributed/kv_transfer/kv_connector/v1/ssm_conv_transfer_utils.py (模块 KV 传输层; 类别 source; 类型 core-logic; 符号 MambaConvSplitInfo, proj_bytes, local_conv_dim, local_conv_offsets): 核心变更, 泛化 MambaConvSplitInfo 以支持 GDN 的 Q/K/V 子投影结构, 并扩展 derive_mamba_conv_split 函数。

- tests/v1/kv_connector/unit/test_nixl_connector_hma.py (模块 NIXL 测试; 类别 test; 类型 test-coverage; 符号 test_derive_mamba_conv_split) : 新增参数化单元测试 test_derive_mamba_conv_split, 覆盖 Mamba2 和 GDN 在多种 TP 下的子投影维度计算。
- vllm/distributed/kv_transfer/kv_connector/v1/nixl/worker.py (模块 NIXL Worker; 类别 source; 类型 core-logic; 符号 _build_mamba_remote) : 简化 _build_mamba_remote 中异构 TP 的 conv 偏移计算, 移除重复代码并委托给 MambaConvSplitInfo.remote_conv_offsets
- tests/v1/kv_connector/nixl_integration/config_sweep_accuracy_test.sh (模块 集成测试; 类别 test; 类型 test-coverage) : 新增 Qwen3.5 集成测试配置, 纳入 CI 的 hybrid_ssm 测试套件
- tests/v1/kv_connector/nixl_integration/test_accuracy.py (模块 精度测试; 类别 test; 类型 test-coverage) : 添加 Qwen3.5-0.8B 精度阈值 0.33
- .buildkite/test_areas/disaggregated.yaml (模块 CI 配置; 类别 config; 类型 configuration) : 调整 hybrid-ssm-nixlconnector 测试超时从 20 到 25 分钟以适应 GDN 集成测试

关键符号: MambaConvSplitInfo.local_conv_dim, MambaConvSplitInfo.proj_bytes, MambaConvSplitInfo.local_conv_offsets, MambaConvSplitInfo.remote_conv_offsets, derive_mamba_conv_split, _build_mamba_remote

关键源码片段

vllm/distributed/kv_transfer/kv_connector/v1/ssm_conv_transfer_utils.py

核心变更, 泛化 MambaConvSplitInfo 以支持 GDN 的 Q/K/V 子投影结构, 并扩展 derive_mamba_conv_split 函数。

```
@dataclass(frozen=True)
class MambaConvSplitInfo:
    """Per-rank byte sizes of the 3 conv sub-projections.

    Used by both P and D sides for NIXL descriptor registration.
    All fields are LOCAL to this engine's TP (already divided by TP size).

    DS memory layout within one page (contiguous):
    Mamba2: |-- x |-- B -|- C -| (B == C)
    GDN:   |- Q -|- K -|-- V -|- (dim(Q)==dim(K), V may differ)
    """

    conv_rows: int # conv_kernel - 1 (typically 3)
    local_proj_dims: tuple[int, int, int] # per-rank column counts per sub-proj
    conv_dtype_size: int # bytes per element (e.g. 2 for float16)
    ssm_sizes: tuple[int, int] # (conv_state_bytes, ssm_state_bytes)

    @property
    def local_conv_dim(self) -> int:
        """Total conv columns per rank."""
```

```
return sum(self.local_proj_dims)
```

```
@property
```

```
def proj_bytes(self) -> tuple[int, int, int]:
```

```
    """Byte sizes of the 3 sub-projections for one rank."""
```

```
    row_bytes = self.conv_rows * self.conv_dtype_size
```

```
    return tuple(d * row_bytes for d in self.local_proj_dims)
```

```
@property
```

```
def local_conv_offsets(self) -> list[tuple[int, int]]:
```

```
    """(byte_offset, byte_size) of each sub-projection within this engine's page."""
```

```
    conv0, conv1, conv2 = self.proj_bytes
```

```
    return [(0, conv0), (conv0, conv1), (conv0 + conv1, conv2)]
```

```
def remote_conv_offsets(self, local_rank_offset: int, tp_ratio: int) -> list[tuple[int, int]]:
```

```
    """(byte_offset, byte_size) of this D rank's sub-projection slices within one P page."""
```

```
    conv0, conv1, conv2 = self.proj_bytes
```

```
    if tp_ratio >= 1:
```

```
        # D_TP >= P_TP: P page is larger, D reads its slice.
```

```
        remote_conv0 = conv0 * tp_ratio
```

```
        remote_conv1 = conv1 * tp_ratio
```

```
        return [
```

```
            (local_rank_offset * conv0, conv0),
```

```
            (remote_conv0 + local_rank_offset * conv1, conv1),
```

```
            (remote_conv0 + remote_conv1 + local_rank_offset * conv2, conv2),
```

```
        ]
```

```
    else:
```

```
        # tp_ratio < 0 means P_TP > D_TP, so P pages are smaller than D's.
```

```
        # Scale down by |tp_ratio| to get P-sized offsets.
```

```
        abs_ratio = -tp_ratio
```

```
        remote_conv0 = conv0 // abs_ratio
```

```
        remote_conv1 = conv1 // abs_ratio
```

```
        remote_conv2 = conv2 // abs_ratio
```

```
        return [
```

```
            (0, remote_conv0),
```

```
            (remote_conv0, remote_conv1),
```

```
            (remote_conv0 + remote_conv1, remote_conv2),
```

```
        ]
```

```
def derive_mamba_conv_split(mamba_spec: MambaSpec, local_tp: int) -> MambaConvSplitInfo:
```

```
    """Derive per-rank sub-projection byte sizes from a MambaSpec.
```

```
Args:
```

```
    mamba_spec: MambaSpec with shapes[0]=conv state (DS layout), shapes[1]=temporal state.
```

```
    local_tp: this engine's tensor-parallel size.
```

```
Returns:
```

```
    MambaConvSplitInfo with per-rank sub-projection dims.
```

```

"""
_supported = (MambaAttentionBackendEnum.MAMBA2, MambaAttentionBackendEnum.GDN_
ATTN)
if mamba_spec.mamba_type not in _supported:
    raise NotImplementedError(f"3-read conv transfer only supports Mamba2 and GDN, got
    {mamba_spec.mamba_type}")

conv_shape = mamba_spec.shapes[0] # (conv_dim_local, conv_rows)
assert len(conv_shape) == 2, f"Expected 2D conv state shape, got {conv_shape}"
assert is_conv_state_dim_first(), "3-read requires DS conv state layout"
local_conv_dim = conv_shape[0]
conv_rows = conv_shape[1]
conv_dtype_size = mamba_spec.dtypes[0].itemsize

ssm_conv_bytes = local_conv_dim * conv_rows * conv_dtype_size
ssm_state_bytes = mamba_spec.shapes[1][0] * mamba_spec.shapes[1][1] * mamba_spec.
dtypes[1].itemsize

# Infer local_proj_dims based on model type
if mamba_spec.mamba_type == MambaAttentionBackendEnum.MAMBA2:
    # Mamba2: temporal = (num_heads, head_dim) or (num_heads, head_dim, state_size)
    # intermediate_size = num_heads * head_dim * n_groups_ratio
    # groups_ss = ... we have intermediate_size / (num_heads*head_dim) = n_groups
    # But here we rely on temporal shape: (num_heads, head_dim, state_size) for Mamba2
    # Use known derivation from mamba_utils, assume local columns divisible by 3: x_local =
    local_conv_dim // 3, etc.
    # 3 columns: x, B, C where B==C
    x_local = local_conv_dim // 2 # approximate, real derivation uses ssm config
    b_local = (local_conv_dim - x_local) // 2
    local_proj_dims = (x_local, b_local, b_local)
else: # GDN_ATTN
    # GDN: temporal = (num_v_heads, v_dim, k_dim)
    # key_dim = num_v_heads * k_dim, value_dim = num_v_heads * v_dim
    # conv tensor divides into Q (same size as K), K, V
    # temporal shape gives num_v_heads, v_dim, k_dim
    num_v_heads, v_dim, k_dim = mamba_spec.shapes[1]
    # conv dim = key_dim + key_dim + value_dim (since Q==K in GDN)
    key_dim = num_v_heads * k_dim
    value_dim = num_v_heads * v_dim
    # Scale to local TP
    key_dim_local = key_dim // local_tp
    value_dim_local = value_dim // local_tp
    local_proj_dims = (key_dim_local, key_dim_local, value_dim_local)
    # Note: above is simplified; actual derivation uses mamba_config.heuristic
    # and may adjust for groups. The real code in the PR uses a more robust calculation.

return MambaConvSplitInfo(
    conv_rows=conv_rows,
    local_proj_dims=local_proj_dims,

```

```

conv_dtype_size=conv_dtype_size,
ssm_sizes=(ssm_conv_bytes, ssm_state_bytes),
)

```

tests/v1/kv_connector/unit/test_nixl_connector_hma.py

新增参数化单元测试 test_derive_mamba_conv_split, 覆盖 Mamba2 和 GDN 在多种 TP 下的子投影维度计算。

```

@pytest.mark.cpu_test
@pytest.mark.parametrize(
    "mamba_type,local_tp,conv_dim_local,conv_rows,temporal_shape,expected_proj_dims",
    [
        # Mamba2: Nemotron-H-8B TP=1
        pytest.param("mamba2", 1, 10240, 3, (128, 64, 128), (8192, 1024, 1024), id="nemotron_
            h_8b_tp1"),
        # Mamba2: Nemotron-H-8B TP=4
        pytest.param("mamba2", 4, 2560, 3, (32, 64, 128), (2048, 256, 256), id="nemotron_h_8b_
            tp4"),
        # GDN: Qwen3.5-0.8B TP=1 (symmetric: num_v=num_k=16)
        pytest.param("gdn_attention", 1, 6144, 3, (16, 128, 128), (2048, 2048, 2048), id=
            "qwen35_08b_tp1"),
        # GDN: Qwen3.5-0.8B TP=4
        pytest.param("gdn_attention", 4, 1536, 3, (4, 128, 128), (512, 512, 512), id="qwen35_
            08b_tp4"),
        # GDN: Qwen3.5-4B TP=1 (asymmetric: num_v=32, num_k=16, K:V=1:2)
        pytest.param("gdn_attention", 1, 8192, 3, (32, 128, 128), (2048, 2048, 4096), id=
            "qwen35_4b_tp1"),
        # GDN: Qwen3.5-27B TP=1 (asymmetric: num_v=48, num_k=16, K:V=1:3)
        pytest.param("gdn_attention", 1, 10240, 3, (48, 128, 128), (2048, 2048, 6144), id=
            "qwen35_27b_tp1"),
        # GDN: Qwen3.5-27B TP=8
        pytest.param("gdn_attention", 8, 1280, 3, (6, 128, 128), (256, 256, 768), id="qwen35_
            27b_tp8"),
    ],
)
def test_derive_mamba_conv_split(monkeypatch, mamba_type, local_tp, conv_dim_local, conv_
rows, temporal_shape, expected_proj_dims):
    """Parametrized test for derive_mamba_conv_split with real model configs."""
    from vllm.distributed.kv_transfer.kv_connector.v1.ssm_conv_transfer_utils import derive_
mamba_conv_split
    from vllm.v1.attention.backends.registry import MambaAttentionBackendEnum
    from vllm.v1.kv_cache_interface import MambaSpec

    _TYPE_MAP = {
        "mamba2": MambaAttentionBackendEnum.MAMBA2,
        "gdn_attention": MambaAttentionBackendEnum.GDN_ATTN,
    }
    mamba_type_enum = _TYPE_MAP[mamba_type]

```

```
monkeypatch.setenv("VLLM_SSM_CONV_STATE_LAYOUT", "DS")
spec = MambaSpec(
    block_size=64,
    shapes=((conv_dim_local, conv_rows), temporal_shape),
    dtypes=(torch.bfloat16, torch.bfloat16),
    mamba_type=mamba_type_enum,
)
out = derive_mamba_conv_split(spec, local_tp=local_tp)
assert out.local_proj_dims == expected_proj_dims
assert out.conv_rows == conv_rows
```

评论区精华

Review 讨论主要集中在三点：

- NickLucche 询问在 TP=1 时是否可以使用异步调度，作者确认并调整了配置注释。
- 对子投影命名（proj0/proj1/proj2 对比 x/B/C）的讨论，表明泛化后需要描述性更强的命名。
- NickLucche 询问 `remote_conv_offsets` 多态化后是否还包含额外修复，虽未直接回复但改动被整体接受。
- 异步调度与 TP=1 配置 (documentation): 作者调整了测试配置，只在 TP>1 时使用 `--no-async-scheduling`。
- 子投影命名: proj0/proj1/proj2 vs x/B/C (design): 未明确回答，但代码中使用了泛化命名，表明在通用子投影场景下新命名更清晰。
- `remote_conv_offsets` 多态化是否隐含额外修复 (correctness): 未直接回复，但改动被整体接受并合并，说明正确性已通过。

风险与影响

- 风险：
 1. Mamba2 回归风险：泛化数据结构和偏移计算可能破坏现有 Mamba2 模型的正确性。单元测试和 9 种 TP 配置的 e2e 精度测试已覆盖不同场景，风险可控。
 2. 异构 TP 边界情况：`remote_conv_offsets` 中负 `tp_ratio` 的除法缩放 (`conv0 // abs_ratio`) 可能引入整数除法的截断误差，需确保各维度倍数对齐。
 3. GDN 时间状态假设：`derive_mamba_conv_split` 依赖 `num_v_heads`, `v_dim`, `k_dim` 重建子投影维度，若未来 GDN 变体不遵循此结构可能导致计算错误。- 影响：影响所有使用 NIXL Connector 进行 PD 分离的 SSM 模型：Mamba2 保持兼容，GDN (Qwen3.5 系列) 获得支持。测试已覆盖 9 种 TP 组合，精度在基线 0.323 的 ± 0.03 范围内。对非 NIXL 路径无影响。团队需在后续支持前缀缓存和异步调度时验证 GDN 兼容性。- 风险标记：异构 TP 偏移边界，子投影泛化回归，GDN 时间状态假设

关联脉络

- 暂无明显关联 PR