

PR #39949 完整报告

vllm-project/vllm

[Spec Decode] Support hybrid attention models in extract_hidden_states

合并时间: 2026-05-14 01:45

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/39949>

执行摘要

- 一句话: 支持混合注意力模型的隐藏状态提取
- 推荐动作: 值得精读的设计决策: HMA 支持检查的泛化机制 (通过 SupportsHMA 接口和 supports_hma 函数) 使 KV connector 能够声明兼容性; HiddenStateCacheSpec 的隔离分组方式使纯缓存层不干扰正常缓存分配。建议关注后续对 max_memory_usage_bytes 的 CP 修复。

功能与动机

Hidden-state extraction breaks on hybrid-attention models (e.g. Qwen3.5) because kv_transfer_config force-disables HMA and unify_hybrid_kv_cache_specs cannot fold MambaSpec into a uniform type. This PR fixes that by introducing HiddenStateCacheSpec and making HMA disable conditional on connector capability.

实现拆解

1. 引入 HiddenStateCacheSpec 数据类型 (vllm/v1/kv_cache_interface.py), 表示仅用于存储的缓存层 (非 K/V 对), 供 CacheOnlyAttentionLayer 使用。
2. 修改 KV 缓存分组逻辑 (vllm/v1/core/kv_cache_utils.py): 在 get_kv_cache_groups 函数中, 先过滤出 HiddenStateCacheSpec 层, 避免它们影响统一化的页大小计算, 然后以对齐后的页大小重新添加为独立的缓存组。
3. 启用 HMA 自动适配 (vllm/config/vllm.py): 对于 kv_transfer_config 不为 None 的情况, 通过 KVConnectorFactory 获取 connector 类并检查是否继承 SupportsHMA (使用新的 supports_hma 函数)。仅当 connector 不支持 HMA 时才强制禁用混合 KV 缓存管理器。特别处理 MultiConnector, 需要所有子连接器都支持 HMA。
4. 调整 ExampleHiddenStatesConnector (vllm/distributed/kv_transfer/kv_connector/v1/example_hidden_states_connector.py): 使其继承 SupportsHMA 接口, 简化 ReqMeta 数据结构 (移除 slot_mapping 等字段), 改为从 attn_metadata 直接获取 slot_mapping, 并更新 extract_from_kv_cache 函数以直接使用块索引导出。
5. 修改 ExtractHiddenStatesProposer (vllm/v1/spec_decode/extract_hidden_states.py): 记录其所属的 kv_cache_group_id, 以便在 common_attn_metadata 中选择正确的分组。同时修改 validate_same_kv_cache_group 方法, 实际查找层所在的组。

6. 调整 `extract_hidden_states` 模型 (`vllm/model_executor/models/extract_hidden_states.py`) : 改为使用 `HiddenStateCacheSpec` 替代 `MLAAttentionSpec`, 并修正 `basic_cache` 函数的方向 (缓存视图改为维护序列索引) 。

7. 补充测试: 新增 `test_extract_hidden_states_qwen35_hybrid_smoke` 冒烟测试 (`tests/v1/kv_connector/extract_hidden_states_integration/test_extraction.py`) , 确保混合模型端到端工作; 同时修复现有测试 (`force fork` 以避免 `spawn` 模式下的模型注册问题)

。

关键文件:

- `vllm/distributed/kv_transfer/kv_connector/v1/example_hidden_states_connector.py` (模块 连接器; 类别 `source`; 类型 `core-logic`; 符号 `ExampleHiddenStatesConnector`, `request_finished_all_groups`) : 核心连接器改造: 继承 `SupportsHMA`, 简化 `ReqMeta`, 使用 `attn_metadata slot_mapping`
- `vllm/v1/core/kv_cache_utils.py` (模块 KV 缓存; 类别 `source`; 类型 `dependency-wiring`) : KV 缓存分组核心修改: 过滤 `HiddenStateCacheSpec` 层并单独分组, 避免影响统一化
- `vllm/config/vllm.py` (模块 配置; 类别 `source`; 类型 `dependency-wiring`) : HMA 自动禁用逻辑重构, 根据 `connector` 支持情况决定是否保留 HMA
- `vllm/v1/spec_decode/extract_hidden_states.py` (模块 推测解码; 类别 `source`; 类型 `core-logic`) : 记录 KV 缓存组 ID, 验证提取层属于同一组
- `vllm/model_executor/models/extract_hidden_states.py` (模块 模型执行; 类别 `source`; 类型 `data-contract`) : 模型定义改用 `HiddenStateCacheSpec`, 修正缓存方向
- `tests/v1/kv_connector/extract_hidden_states_integration/test_extraction.py` (模块 测试; 类别 `test`; 类型 `test-coverage`; 符号 `test_extract_hidden_states_qwen35_hybrid_smoke`) : 新增 Qwen3.5 混合模型冒烟测试

关键符号: `extract_from_kv_cache`, `ExampleHiddenStatesConnector.save_kv_layer`, `ExampleHiddenStatesConnector.build_connector_meta`, `get_kv_cache_groups`, `ExtractHiddenStatesProposer.validate_same_kv_cache_group`, `VllmConfig.post_init`, `basic_cache`, `CacheOnlyAttentionBackend.get_kv_cache_spec`

关键源码片段

`vllm/distributed/kv_transfer/kv_connector/v1/example_hidden_states_connector.py`

核心连接器改造: 继承 `SupportsHMA`, 简化 `ReqMeta`, 使用 `attn_metadata slot_mapping`

```
# 从 vllm/distributed/.../example_hidden_states_connector.py 提取
# 展示 extract_from_kv_cache 和使用 attn_metadata slot_mapping 的 save_kv_layer
```

```
def extract_from_kv_cache(
    kv_cache: torch.Tensor,
    slot_mapping: torch.Tensor,
    num_tokens: int,
) -> torch.Tensor:
    """Extract data from KV cache."""
```

```

block_size = kv_cache.shape[1] # 已知 block_size 为首位
return kv_cache[slot_mapping // block_size, slot_mapping % block_size][:num_tokens]

```

```

class ExampleHiddenStatesConnector(KVConnectorBase_V1, SupportsHMA):
    # 继承 SupportsHMA 以便 kv_transfer_config 保留混合 KV 缓存管理器

```

```

def save_kv_layer(
    self,
    kv_layer: torch.Tensor,
    attn_metadata: AttentionMetadata,
    connector_metadata: ExampleHiddenStatesConnectorMetadata,
) -> None:
    # 从 attn_metadata 获取 slot_mapping, 替代之前从 ReqMeta 获取
    slot_mapping = attn_metadata.slot_mapping
    offset = 0
    for request in connector_metadata.requests:
        num_tokens = request.token_ids.shape[0]
        req_slot_mapping = slot_mapping[offset: offset + num_tokens]
        offset += num_tokens
        hidden_states = extract_from_kv_cache(
            kv_layer, req_slot_mapping, num_tokens
        )
        tensors = {"hidden_states": hidden_states.detach().cpu(), ...}
        safetensors.torch.save(tensors, request.filename)

```

vllm/v1/core/kv_cache_utils.py

KV 缓存分组核心修改: 过滤 HiddenStateCacheSpec 层并单独分组, 避免影响统一化

```

# 从 vllm/v1/core/kv_cache_utils.py 的 get_kv_cache_groups 函数末尾提取
# 过滤 HiddenStateCacheSpec 层, 统一页大小后重新添加

```

```

hidden_specs = {
    k: v for k, v in kv_cache_spec.items() if isinstance(v, HiddenStateCacheSpec)
}

```

```

filtered_spec = {
    k: v
    for k, v in kv_cache_spec.items()
    if not isinstance(v, HiddenStateCacheSpec)
}

```

```

# 统一过滤后的 spec 的页大小

```

```

filtered_spec = unify_kv_cache_spec_page_size(filtered_spec)
groups = _get_kv_cache_groups_uniform_page_size(filtered_spec)

```

```

# 将隐藏状态层重新添加为独立组, 页对齐到公共页大小

```

```

if hidden_specs:
    common_page = get_uniform_page_size([g.kv_cache_spec for g in groups])
    for name, spec in hidden_specs.items():
        per_token = spec.num_kv_heads * spec.head_size * get_dtype_size(spec.dtype)

```

```

new_bs = max(common_page // per_token, 1)
aligned = replace(spec, block_size=new_bs, page_size_padded=common_page)
groups.append(KVCacheGroupSpec([name], aligned))

```

```
return groups
```

vllm/config/vllm.py

HMA 自动禁用逻辑重构，根据 connector 支持情况决定是否保留 HMA

从 vllm/config/vllm.py 的 __post_init__ 中提取，HMA 自动禁用检查

```

if self.scheduler_config.disable_hybrid_kv_cache_manager is None:
    if self.kv_transfer_config is not None:
        from vllm.config.kv_transfer import KVTransferConfig
        from vllm.distributed.kv_transfer.kv_connector.factory import KVConnectorFactory
        from vllm.distributed.kv_transfer.kv_connector.v1.base import supports_hma

        connector_cls = KVConnectorFactory.get_connector_class(
            self.kv_transfer_config
        )
        all_support_hma = supports_hma(connector_cls)
        # MultiConnector 继承 SupportsHMA，但需确保所有子连接器都支持
        if all_support_hma and connector_cls.__name__ == "MultiConnector":
            sub_ktcs = self.kv_transfer_config.kv_connector_extra_config.get(
                "connectors", []
            )
            all_support_hma = all(
                supports_hma(
                    KVConnectorFactory.get_connector_class(
                        KVTransferConfig(**sub)
                    )
                )
                for sub in sub_ktcs
            )
        if not all_support_hma:
            need_disable_hybrid_kv_cache_manager = True
            logger.warning(
                "Turning off hybrid kv cache manager because "
                "connector %s does not subclass `SupportsHMA`."
                "This will reduce performance on models with "
                "sliding window or Mamba attention.",
                connector_cls.__name__,
            )

```

评论区精华

- benchislett 担忧手动特殊化 (gpu_model_runner.py 中的手动步进视图) 会偏离标准实践，增加维护成本，最终 APPROVED 但有保留。

- gemini-code-assist 指出 HiddenStateCacheSpec.max_memory_usage_bytes 未考虑上下文并行 (DCP/PCP) ，会导致内存高估，配置检查可能失败。该问题未在本 PR 中修复。
- NickLucche 认为 MultiConnector HMA 支持检查的代码“过深”，但设计方案确保了子连接器都支持 HMA 时的正确行为。
- 手动特殊化 vs 通用填充 (design): 作者认为当前方案是最小改动，且已与 main 上的 strided view 功能对齐。benchislett 最终 APPROVED，但保留了对脆弱性的担忧。
- HiddenStateCacheSpec 内存计算未考虑上下文并行 (correctness): 该问题未在本 PR 中修复，可能在后续 PR 中处理。
- MultiConnector HMA 支持检查 (design): 当前实现已合并，是确保正确性的必要折中。

风险与影响

- 风险:
 - 手动特殊化: gpu_model_runner.py 中新增的手动步进视图代码与现有原生填充方案不一致，随 HMA 演进可能产生维护债务。
 - HiddenStateCacheSpec 内存计算: max_memory_usage_bytes 未考虑上下文并行，在 CP 配置下可能导致初始化失败（需后续修复）。
 - chunked prefill 兼容性: issue 评论指出在 chunked prefill 模式下隐藏状态保存过早，导致形状不匹配，但此问题在本 PR 范围外。
 - 回归风险: 修改了 HMA 自动禁用逻辑，可能影响现有 KV connector 用户（但已通过检查向后兼容性）。
 - 影响: 对用户: 现在可在 Qwen3.5 等混合注意力模型上使用 extract_hidden_states 推测解码，扩展了功能适用场景。对系统: 涉及配置、KV 缓存管理、分布层 connector 和模型执行的联动改动，影响面中等。对团队: 新增手动特殊化路径需要后续维护；内存计算问题需修复。
 - 风险标记: 手动特殊化维护成本，内存计算未考虑上下文并行，chunked prefill 兼容性未验证

关联脉络

- 暂无明显关联 PR