

# PR #44212 完整报告

vllm-project/vllm

[Perf] Improve multimodal item handling from  $O(n)$  to  $O(\log n)$  per step

合并时间: 2026-06-03 19:00

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/44212>

## 执行摘要

- 一句话: 二分查找加速多模态特征遍历, 每步  $O(n) \rightarrow O(\log n)$
- 推荐动作: 此 PR 是典型的  $O(n) \rightarrow O(\log n)$  优化范例, 推荐精读。关键设计决策包括: 二分查找边界处理 (使用 `offset+length` 而不是 `offset`)、`encoder-decoder` 特殊处理、以及 `request_cached_ids` 的清理策略。这些细节值得在类似优化中参考。

## 功能与动机

Voxtral Realtime 在长转录会话中因处理大量多模态项 (最高 32K) 而显著变慢。PR body 指出 'vLLM is not very efficient in handling a large amount of multimodal items'。

## 实现拆解

1. 新增 `get_mm_features_in_window` 工具函数(`vllm/multimodal/utils.py`): 基于 `bisect` 在已按 `offset` 排序的 `mm_features` 列表中定位与给定 token 窗口重叠的特征范围, 返回 `(lo, hi)` 索引, 复杂度  $O(\log n)$ 。
2. 调度器改用二分查找(`vllm/v1/core/sched/scheduler.py`): 在 `_try_schedule_encoder_inputs` 中将原有的 `for i, mm_feature in enumerate(mm_features)` 线性扫描替换为先调用 `get_mm_features_in_window` 获得 `lo/hi`, 再只遍历该子范围。对于 `encoder-decoder` 模型, 由于所有输入 `offset=0`, 强制 `lo=0`。
3. 模型运行器同步优化(`vllm/v1/worker/gpu_model_runner.py`): 在 `_gather_mm_embeddings` 中做相同替换, 移除原有的 `break/continue` 线性判断。
4. `EncoderCacheManager` 添加 per-request 索引 (`vllm/v1/core/encoder_cache_manager.py`): 引入 `request_cached_ids` 字典 (`request_id \rightarrow set of input_id`), 在 `check_and_update_cache`、`allocate` 中同步记录, 将 `get_cached_input_ids` 从扫描所有 `mm_features` 的  $O(n)$  降为直接字典查询的  $O(1)$ 。在 `free_encoder_input` 中确保始终清理此字典 (即使 `mm_hash` 已被驱逐), 并修复因早期返回导致遗漏的 bug。
5. Voxtral 处理器微优化(`vllm/transformers_utils/processors/voxtral.py`): 将 `torch.tensor(audio)` 改为 `torch.from_numpy(audio)`, 避免 CPU 到 tensor 的数据拷贝。
6. 测试补充(`tests/v1/core/test_encoder_cache_manager.py`): 添加 `test_free_request_with_duplicate_mm_hashes` 测试, 覆盖重复 `mm_hash` 时

request\_cached\_ids 的完整清理。

关键文件：

- vllm/multimodal/utils.py (模块 多模态工具; 类别 source; 类型 core-logic; 符号 get\_mm\_features\_in\_window) : 新增 get\_mm\_features\_in\_window 函数, 是整个优化的核心工具, 被调度器和模型运行器复用。
- vllm/v1/core/sched/scheduler.py (模块 调度器; 类别 source; 类型 core-logic) : 核心调度路径中使用二分查找替换线性扫描, 是性能提升的关键之一。
- vllm/v1/core/encoder\_cache\_manager.py (模块 编码缓存管理; 类别 source; 类型 core-logic) : 引入 request\_cached\_ids 字典, 将 get\_cached\_input\_ids 从  $O(n)$  降为  $O(1)$ , 并修复清理逻辑。
- vllm/v1/worker/gpu\_model\_runner.py (模块 模型运行器; 类别 source; 类型 core-logic) : 模型运行器使用相同二分查找优化 \_gather\_mm\_embeddings, 避免线性扫描。
- tests/v1/core/test\_encoder\_cache\_manager.py (模块 测试; 类别 test; 类型 test-coverage; 符号 test\_free\_request\_with\_duplicate\_mm\_hashes) : 新增 test\_free\_request\_with\_duplicate\_mm\_hashes 测试, 覆盖重复 mm\_hash 场景下的 cleanup, 确保正确性。
- vllm/transformers\_utils/processors/voxtral.py (模块 Voxtral 处理器; 类别 source; 类型 core-logic) : 使用 torch.from\_numpy 替代 torch.tensor 实现零拷贝音频转换, 减少不必要的 CPU 到 tensor 拷贝。

关键符号: get\_mm\_features\_in\_window, \_try\_schedule\_encoder\_inputs, \_gather\_mm\_embeddings, check\_and\_update\_cache, allocate, get\_cached\_input\_ids, free\_encoder\_input, free

## 关键源码片段

### vllm/multimodal/utils.py

新增 get\_mm\_features\_in\_window 函数, 是整个优化的核心工具, 被调度器和模型运行器复用。

```
import bisect
from .inputs import MultiModalFeatureSpec

def get_mm_features_in_window(
    mm_features: list[MultiModalFeatureSpec],
    start: int,
    end: int,
) -> tuple[int, int]:
    """Return (lo, hi) indices for features overlapping [start, end).

    Assumes mm_features are sorted by offset and non-overlapping, so
    offset + length is also sorted.
    """
    # bisect_left on start+1 using end offset (offset+length) to find first
    # feature whose end >= start, i.e., overlapping on the left.
```

```

lo = bisect.bisect_left(
    mm_features,
    start + 1,
    key=lambda f: f.mm_position.offset + f.mm_position.length,
)
# bisect_left on end using start offset to find first feature whose
# offset >= end, i.e., beyond the window.
hi = bisect.bisect_left(
    mm_features,
    end,
    key=lambda f: f.mm_position.offset,
)
return lo, hi

```

### vllm/v1/core/sched/scheduler.py

核心调度路径中使用二分查找替换线性扫描，是性能提升的关键之一。

```

def _try_schedule_encoder_inputs(self, request, num_new_tokens, ...):
    # ... earlier setup ...
    # Use bisect to narrow iteration from O(n) to O(log n)
    lo, hi = get_mm_features_in_window(
        mm_features,
        start=num_computed_tokens,
        end=num_computed_tokens + num_new_tokens + shift_computed_tokens,
    )
    # For encoder-decoder, all inputs sit at start_pos=0, so lo=0 always.
    if self.is_encoder_decoder:
        lo = 0

    for i in range(lo, hi):
        mm_feature = mm_features[i]
        start_pos = mm_feature.mm_position.offset
        num_encoder_tokens = mm_feature.mm_position.length
        # ... rest of scheduling logic ...

```

### vllm/v1/core/encoder\_cache\_manager.py

引入 `request_cached_ids` 字典，将 `get_cached_input_ids` 从  $O(n)$  降为  $O(1)$ ，并修复清理逻辑。

```

class EncoderCacheManager:
    def __init__(self, cache_size: int):
        # ... existing code ...
        self.cached: dict[str, set[str]] = {}
        # Per-request cache: request_id -> set of input_ids cached
        self.request_cached_ids: dict[str, set[int]] = {}
        # ... rest ...

    def get_cached_input_ids(self, request: Request) -> set[int]:
        """Get all cached multimodal input IDs for a request.

```

O(1) lookup using per-request index, vs O(n) scanning all mm\_features.

```
"""
```

```
return self.request_cached_ids.get(request.request_id, set())
```

```
def free_encoder_input(self, request: Request, input_id: int) -> None:
```

```
    """Free the request's reference to the encoder input."""
```

```
    req_id = request.request_id
```

```
    mm_hash = request.mm_features[input_id].identifier
```

```
    # Always clean up request_cached_ids, even if the mm_hash was
```

```
    # already evicted (e.g., by can_allocate).
```

```
    if req_id in self.request_cached_ids:
```

```
        self.request_cached_ids[req_id].discard(input_id)
```

```
        if not self.request_cached_ids[req_id]:
```

```
            del self.request_cached_ids[req_id]
```

```
    # If mm_hash not in cache or no references, early return.
```

```
    if not self.cached.get(mm_hash, None):
```

```
        return
```

```
    # ... existing refcount logic ...
```

## 评论区精华

Review 中讨论了几个关键点：

- Mrv2 适用性：NickLucche 询问是否同样适用于 Mrv2（未直接回应，但可后续跟进）。
- free\_encoder\_input cleanup: ywang96 指出存在早期返回导致 request\_cached\_ids 未清理的问题，andylo2 随后修复并补充测试。
- defaultdict vs dict: njhill 建议使用 defaultdict，andylo2 解释使用普通 dict 更安全，避免在 get 时自动创建空条目，最终保留 dict。
  - 是否同样适用于 Mrv2 (question): 未直接回应，但 PR 作者可能认为类似需求可后续跟进。
  - free\_encoder\_input 中 early return 导致 request\_cached\_ids 未清理 (correctness): andylo2 确认并修复，将 cleanup 移到 early return 之前，并添加测试覆盖。
  - 使用 defaultdict 简化 request\_cached\_ids (style): 保留普通 dict，使用 setdefault。

## 风险与影响

- 风险：
  1. 排序假设风险：二分查找假设 mm\_features 已按 offset 排序且不重叠，若运行时出现异常排序可能导致越界或遗漏（暂无 runtime 验证）。
  2. encoder-decoder 特殊处理：强制将 lo 置为 0 的逻辑需与其他路径保持一致，避免误判。
  3. per-request cache 状态管理：request\_cached\_ids 引入额外状态，若清理不及时或遗漏可能导致内存泄漏（测试已覆盖重复 hash 场景）。

4. torch.from\_numpy 连续性: 要求输入 numpy 数组连续, voxtral 中 audio 已由 pad 保证, 但未来改动可能引入隐患。 - 影响: 直接影响 Voxtral Realtime 等长序列多模态模型, profile 显示 scheduler.schedule() 和 GpuModelRunner.\_preprocess() 耗时显著下降。对普通单 / 少多模态输入的模型几乎无影响 (二分查找开销极小)。系统整体吞吐在长转录场景预计提升明显。团队可借鉴此模式优化其他  $O(n)$  遍历热点。 - 风险标记: 二分查找假设排序不重叠, encoder-decoder 特殊处理, per-request cache 清理, torch.from\_numpy 要求连续内存

## 关联脉络

- 暂无明显关联 PR