

PR #27458 完整报告

sgl-project/sglang

[spec] Consolidate the per-decode KV alloc reserve into one helper

合并时间: 2026-06-07 05:58

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/27458>

执行摘要

- 一句话: 统一 spec decode KV 分配预留计算并移动 helper
- 推荐动作: 建议精读, 因为展示了如何通过集中化计算消除重复逻辑并解决导入循环, 是一个教科书级的纯重构案例。特别关注 `pr_fix_toggle.py` 中 revert target 的迁移方式。

功能与动机

PR body 指出 `2 * get_alloc_len_per_decode` 在四个站点重复书写, 容易不一致; `managers/utils.py` 本不属于分配逻辑, 且存在导入循环 (`managers/utils.py` 顶导入 `schedule_batch`, 导致 `schedule_batch` 和 `model_runner` 必须通过函数内局部导入来使用这些 helper)。

实现拆解

1. 在 `mem_cache/common.py` 中新增 `get_alloc_len_per_decode` (从 `utils.py` 搬入)、`get_alloc_reserve_per_decode` (`2 * get_alloc_len_per_decode`) 和 `get_req_to_token_extra_context_len` (封装原 `_init_pools` 中的行宽头寸计算逻辑)。
2. 从 `managers/utils.py` 中删除原 `get_alloc_len_per_decode` 并移除不再需要的 `ServerArgs` 导入。
3. 在 `model_runner_kv_cache_mixin.py` 的 `_init_pools` 中, 用 `get_req_to_token_extra_context_len` 替换内联的 `extra_max_context_len` 计算和函数内导入。
4. 在 `eagle_info_v2.py` 的 `prepare_for_decode` 中, 改用 `get_alloc_reserve_per_decode` 并更新断言消息。
5. 在 `schedule_batch.py` 的 `_new_tokens_required_next_decode_spec_v2` 中, 改用 `get_alloc_reserve_per_decode` 并移除函数内导入。
6. 更新 `pr_fix_toggle.py` 中 `_PR_REVERT_YAML_26972` 的 target 从 `ModelRunnerKVCacheMixin._init_pools` 改为 `mem_cache.common.get_req_to_token_extra_context_len`, 使其仍能通过 patch 准确回退 #26972 的改动。

关键文件:

- `python/sglang/srt/mem_cache/common.py` (模块 缓存层; 类别 source; 类型 core-logic ; 符号 `get_alloc_len_per_decode`, `get_alloc_reserve_per_decode`, `get_req_to_token_extra_context_len`): 核心文件, 新增了三个 helper 函数并作为分配逻辑

辑的新归属地。

- python/sglang/srt/managers/utils.py (模块 调度器; 类别 source; 类型 core-logic; 符号 get_alloc_len_per_decode) : 被移除原有函数, 简化了该模块的职责。
- python/sglang/srt/model_executor/model_runner_kv_cache_mixin.py (模块 模型执行器; 类别 source; 类型 data-contract) : 简化了 _init_pools 中行宽头寸计算, 消除了函数内导入。
- python/sglang/srt/speculative/eagle_info_v2.py (模块 推测解码; 类别 source; 类型 dependency-wiring) : 使用统一 helper 并更新断言消息。
- python/sglang/srt/managers/schedule_batch.py (模块 调度器; 类别 source; 类型 dependency-wiring) : 使用统一 helper 并移除函数内导入。
- python/sglang/srt/debug_utils/pr_fix_toggle.py (模块 调试工具; 类别 source; 类型 dependency-wiring) : 调整 #26972 回退 patch 的目标, 使其匹配重构后的代码位置。

关键符号: get_alloc_len_per_decode, get_alloc_reserve_per_decode,
get_req_to_token_extra_context_len

关键源码片段

python/sglang/srt/mem_cache/common.py

核心文件, 新增了三个 helper 函数并作为分配逻辑的新归属地。

```
def get_alloc_len_per_decode(server_args: Optional[ServerArgs] = None) -> int:
    """单个 decode 步骤每个 request 的 KV 分配长度。
    根据 speculative 算法和 page_size 计算最坏情况下的 token 数。
    """
    if server_args is None:
        server_args = get_global_server_args()

    if server_args.speculative_algorithm is None:
        return 1

    spec_steps = server_args.speculative_num_steps or 1
    spec_topk = server_args.speculative_eagle_topk or 1
    spec_tokens = server_args.max_speculative_num_draft_tokens
    page_size = server_args.page_size

    if page_size == 1 or spec_topk == 1:
        return max(spec_steps * spec_topk, spec_tokens)
    else:
        # page_size > 1 + topk > 1 (spec v2 tree): worst-case page-aligned tree
        num_new_pages_per_topk = (
            (page_size - 1) + spec_steps + page_size - 1
        ) // page_size
        return max(num_new_pages_per_topk * page_size * spec_topk, spec_tokens)

def get_alloc_reserve_per_decode(server_args: Optional[ServerArgs] = None) -> int:
```

```
"""每个 decode step 预留的 KV 总长度 (double buffer) 。
对应原先各处手写的 `2 * get_alloc_len_per_decode`。
"""
return 2 * get_alloc_len_per_decode(server_args)
```

```
def get_req_to_token_extra_context_len(server_args: ServerArgs) -> int:
    """req_to_token 行宽超出 model context length 的头寸。
    用于容纳 decode 过度分配 (kv_committed_len + get_alloc_reserve_per_decode) ，
    尤其 spec v2 page>1 topk>1 的洞形 draft 足迹可能超出默认 num_draft_tokens 头寸。
    """
    extra = 4 + (server_args.max_speculative_num_draft_tokens or 0)
    if (
        server_args.speculative_algorithm is not None
        and server_args.page_size > 1
        and (server_args.speculative_eagle_topk or 1) > 1
    ):
        extra = max(extra, get_alloc_reserve_per_decode(server_args))
    return extra
```

评论区精华

只有 [gemini-code-assist\[bot\]](#) 的自动评议，没有实质性人类讨论。

- 无人工审查讨论 (other): 无需要处理的讨论。

风险与影响

- 风险：行为保持重构，核心风险在于 (1) `pr_fix_toggle.py` 中反向开关的目标从 `_init_pools` 改为 `get_req_to_token_extra_context_len`，若 `patch-apply` 逻辑有偏差可能导致调试时无法正确回退 #26972 修复；(2) `get_req_to_token_extra_context_len` 的参数类型从 `Optional[ServerArgs]` 改为 `ServerArgs` (非可选)，调用方需确保传入非空 `server_args`。
- 影响：影响范围限于 `speculative decoding` 相关模块。对用户无感知，系统行为不变。对开发者来说，未来若需调整 `decode KV` 预留量只需修改一个入口，降低了维护成本。
- 风险标记：revert toggle 目标变更，参数类型收紧

关联脉络

- PR #26972 Spec v2 tree drafting (topk>1) with page_size>1: 本 PR 统一了 #26972 引入的 `'2 * get_alloc_len_per_decode'` 逻辑，并更新了对应的 revert toggle 目标。
- PR #27364 [perf] reduce radix cache match overhead by changing the match algorithm: 无关性能优化，但同一时间段内对 `radix cache` 进行了修改，与本 PR 在 `mem_cache` 模块有间接关联。