

# PR #41035 完整报告

vllm-project/vllm

[Model Runner V2] Apply synthetic mode to probabilistic rejection sampler

合并时间: 2026-05-13 04:37

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/41035>

## 执行摘要

- 一句话: 将合成拒绝采样融合到统一内核
- 推荐动作: 该 PR 展示了如何将两个独立代码路径合并而不损失性能匹配。值得关注的设计决策: 故意保留 LSE 计算以对齐运行时间。建议推测解码相关开发者精读内核分支。

## 功能与动机

MRV2 当前使用独立的合成拒绝采样路径, 这并不理想, 因为合成接受率应紧密近似实际拒绝采样。将合成模式集成到概率拒绝采样内核中, 能使合成路径的行为与标准路径更一致, 便于维护和性能匹配。

## 实现拆解

步骤 1: 重命名核心符号和文件

- 将 `probabilistic_rejection_sampler_utils.py` 重命名为 `rejection_sampler_utils.py`。
- 将 `_probabilistic_rejection_kernel` 重命名为 `_rejection_kernel`, `probabilistic_rejection_sample` 重命名为 `rejection_sample`。
- 删除 `probabilistic` 前缀, 因为现在只有 `standard` 和 `synthetic` 模式。

步骤 2: 添加合成模式分支

- 在 `_rejection_kernel` 中新增 `SYNTHETIC_MODE` 编译常量和 `synthetic_conditional_rates_ptr` 指针参数。
- 在贪心采样路径 (`temp == 0`) 中, 合成模式使用随机数与条件率比较决定接受, 而非比较 `target argmax`。
- 在随机采样路径中, 同样用条件率代替概率比测试, 但保持 LSE 计算以匹配性能。

步骤 3: 删除独立合成实现

- 移除 `synthetic_rejection_sampler_utils.py` 及其内核函数 `_synthetic_rejection_sample_kernel` 和入口 `synthetic_rejection_sample`。

步骤 4: 简化入口调用

- 在 `RejectionSampler.__call__` 中去掉 `if rejection_sample_method == "standard"` 与 `elif rejection_sample_method == "synthetic"` 的分支, 统一调用 `rejection_sample`。
- 通过 `synthetic_conditional_rates` 参数 (None 时为标准模式) 决定是否启用合成模式。

## 步骤 5: 测试与 CI 更新

- 重命名测试文件 `test_probabilistic_rejection_sampler_utils.py` → `test_rejection_sampler_utils.py`。
- 新增 `test_synthetic_rejection_sample` 参数化测试, 验证在不同温度和无条件率下接受率统计一致性。
- 更新 CI 配置 `.buildkite/test_areas/model_runner_v2.yaml` 中的文件路径。

### 关键文件:

- `vllm/v1/worker/gpu/spec_decode/rejection_sampler_utils.py` (模块 推测解码; 类别 source; 类型 rename-or-move; 符号 `_probabilistic_rejection_kernel`, `_rejection_kernel`, `probabilistic_rejection_sample`, `rejection_sample`): 核心变更文件: 重命名并整合合成模式到内核, 新增 `SYNTHETIC_MODE` 分支和 `synthetic_conditional_rates_ptr` 参数。
- `vllm/v1/worker/gpu/spec_decode/synthetic_rejection_sampler_utils.py` (模块 推测解码; 类别 source; 类型 deletion; 符号 `_synthetic_rejection_sample_kernel`, `synthetic_rejection_sample`): 删除文件: 合成拒绝采样的独立实现被合并到统一内核中。
- `vllm/v1/worker/gpu/spec_decode/rejection_sampler.py` (模块 推测解码; 类别 source; 类型 dependency-wiring): 入口文件: 统一调用 `rejection_sample`, 删除标准 / 合成分支, 简化流程。
- `tests/v1/spec_decode/test_rejection_sampler_utils.py` (模块 单元测试; 类别 test; 类型 rename-or-move; 符号 `test_synthetic_rejection_sample`): 测试文件: 重命名并新增合成拒绝采样测试用例, 验证接受率统计一致性。
- `.buildkite/test_areas/model_runner_v2.yaml` (模块 CI 配置; 类别 config; 类型 configuration): CI 配置: 更新测试文件路径以反映重命名。

关键符号: `_rejection_kernel`, `rejection_sample`, `_synthetic_rejection_sample_kernel` (deleted), `synthetic_rejection_sample` (deleted)

## 关键源码片段

### `vllm/v1/worker/gpu/spec_decode/rejection_sampler_utils.py`

核心变更文件: 重命名并整合合成模式到内核, 新增 `SYNTHETIC_MODE` 分支和 `synthetic_conditional_rates_ptr` 参数。

```
@triton.jit
def _rejection_kernel(
    # ... 参数列表 (略去大多数字段) ...
    synthetic_conditional_rates_ptr, # [num_speculative_steps] 合成条件接受率
    vocab_num_blocks,
    PADDED_VOCAB_NUM_BLOCKS: tl.constexpr,
    HAS_DRAFT_LOGITS: tl.constexpr,
    SYNTHETIC_MODE: tl.constexpr, # 编译常量: 启用合成模式
):
    req_idx = tl.program_id(0)
    req_state_idx = tl.load(idx_mapping_ptr + req_idx)
```

```

start_idx = tl.load(cu_num_logits_ptr + req_idx)
end_idx = tl.load(cu_num_logits_ptr + req_idx + 1)
num_tokens = end_idx - start_idx
seed = tl.load(seed_ptr + req_state_idx)
temp = tl.load(temp_ptr + req_state_idx).to(tl.float32)

rejected_step = 0
accepted = True
for i in range(num_tokens - 1):
    if accepted:
        logit_idx = start_idx + i
        draft_sampled = tl.load(draft_sampled_ptr + logit_idx + 1).to(tl.int64)
        if temp == 0.0:
            # 贪心采样: 加载 target 各块的 argmax
            target_blocks = tl.arange(0, PADDED_VOCAB_NUM_BLOCKS)
            target_blocks_mask = target_blocks < vocab_num_blocks
            target_local_max = tl.load(
                target_local_max_ptr + logit_idx * target_local_max_stride + target_blocks,
                mask=target_blocks_mask, other=float('-inf'))
            max_target_block_idx = tl.argmax(target_local_max, axis=0)
            target_argmax = tl.load(
                target_local_argmax_ptr + logit_idx * target_local_argmax_stride + max_target_
                block_idx
            ).to(tl.int64)

            if SYNTHETIC_MODE:
                # 合成模式: 用随机数和条件率决定是否接受 draft token
                pos = tl.load(pos_ptr + logit_idx)
                u = tl_rand64(seed, pos, includes_zero=False)
                rate = tl.load(synthetic_conditional_rates_ptr + i)
                accepted &= u < rate
            else:
                # 标准模式: 仅接受与 target argmax 匹配的 draft token
                accepted &= target_argmax == draft_sampled
            tl.store(sampled_ptr + req_idx * sampled_stride + i,
                    draft_sampled if accepted else target_argmax)
        else:
            # 随机采样 (temperature > 0)
            target_logit = tl.load(
                target_logits_ptr + logit_idx * target_logits_stride + draft_sampled
            ).to(tl.float32)
            target_lse = _compute_global_lse(
                target_local_max_ptr, target_local_max_stride,
                target_local_sumexp_ptr, target_local_sumexp_stride,
                logit_idx, vocab_num_blocks, PADDED_VOCAB_NUM_BLOCKS)
            target_log_prob = target_logit - target_lse
            pos = tl.load(pos_ptr + logit_idx)
            u = tl_rand64(seed, pos, includes_zero=False)
            if HAS_DRAFT_LOGITS:

```

```

draft_logit = tl.load(
    draft_logits_ptr + req_state_idx * draft_logits_stride_0 + i * draft_logits_stride_
    1 + draft_sampled
).to(tl.float32)
draft_lse = _compute_global_lse(
    draft_local_max_ptr, draft_local_max_stride,
    draft_local_sumexp_ptr, draft_local_sumexp_stride,
    logit_idx, vocab_num_blocks, PADDED_VOCAB_NUM_BLOCKS)
draft_log_prob = draft_logit - draft_lse
else:
    # 无草稿 logits: 假设 one-hot 分布, log probability = 0
    draft_log_prob = 0.0

if SYNTHETIC_MODE:
    # 合成模式: 仅基于条件率, 忽略概率比
    rate = tl.load(synthetic_conditional_rates_ptr + i)
    accepted &= u < rate
else:
    # 标准模式: 概率比测试  $p(x) > u * q(x)$ 
    accepted &= target_log_prob > tl.log(u) + draft_log_prob
tl.store(sampled_ptr + req_idx * sampled_stride + i, draft_sampled)
# 更新 rejected_step 等 (省略)

```

## 评论区精华

### 讨论 1: `_resample_kernel` 是否需要适配

- gemini-code-assist 指出 `_resample_kernel` 未处理 `SYNTHETIC_MODE`, 应在拒绝时直接采样 `target logits` 而非残差分布。
- TheEpicDolphin 回应这是故意设计, 保留残差逻辑以匹配标准拒绝采样的性能。

### 讨论 2: 是否应跳过 LSE 计算

- gemini-code-assist 建议在合成模式下跳过昂贵的 `target_lse` 和 `draft_lse` 计算。
- TheEpicDolphin 解释仍保留这些计算是为了精确匹配标准拒绝采样的运行时间, 确保合成模式能真实反映标准模式性能。

### 讨论 3: 参数类型提示

- gemini-code-assist 指出 `synthetic_conditional_rates` 应为 `Optional[torch.Tensor]`。
- TheEpicDolphin 已修复。
- `_resample_kernel` 未适配合成模式 (correctness): TheEpicDolphin 回应这是故意设计, 保留残差逻辑以匹配标准拒绝采样的性能。
- LSE 计算开销在合成模式下是否必要 (performance): TheEpicDolphin 解释保留计算是为了匹配标准拒绝采样的运行时间。
- 参数类型提示应为 `Optional` (style): TheEpicDolphin 回复已修复。

## 风险与影响

- 风险：性能风险：合成模式下仍计算 LSE 和 log-probability，可能引入不必要的开销。但作者有意为之以确保性能匹配，实际影响有限。正确性风险：如果 `synthetic_conditional_rat` 计算错误或传递错误，可能导致接受率与预期不符。测试覆盖了多种参数组合。兼容性风险：移除了 `rejection_sample_method` 配置键（但 PR 中实际仍保留？从 code 看 `rejection_sampler.py` 中不再有分支，但 `SpeculativeConfig` 可能仍使用该字段？本 PR 未展示配置部分，需确认。文件变更中未见配置删除，可能在后端处理。但需注意如果外部代码依赖 `rejection_sample_method` 字段，可能中断。回归风险：核心路径变更可能影响标准拒绝采样。现有的 `test_stochastic_rejection_sample` 和 `test_greedy_rejection_sample` 应能捕获问题。
- 影响：用户影响：使用 V2 模型运行器和合成拒绝采样的用户将自动受益于统一路径，无需手动选择模式。标准拒绝采样用户无感知。系统影响：移除了一个独立 Triton 内核，减少代码量和二进制大小。团队影响：降低了推理采样模块的复杂度，便于未来维护和优化。
- 风险标记：核心路径变更，性能匹配开销，配置键兼容性，重采样路径同步

## 关联脉络

- PR #40662 Add per-position acceptance rates: 前置依赖，添加每个位置接受率，使得合成拒绝采样成为可能。