

PR #23331 完整报告

sgl-project/sglang

[BugFix] Resolve adaptive speculative decoding conflicts for Qwen3.5 (hybrid GDN)

合并时间: 2026-05-20 06:09

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/23331>

执行摘要

- 一句话: 修复 Qwen3.5 混合 GDN 模型上自适应推测解码的三个冲突问题
- 推荐动作: 建议所有涉及推测解码、混合注意力路由的开发者精读。重点关注:
 - `_is_full_attn` 的 `isinstance` 优先设计, 它提供比 ID 列表更健壮的分发。
 - `effective_max_speculative_num_draft_tokens` 的解耦方式, 可作为未来自适应参数分配的模板。
 - Review 讨论中关于“保留 API 兼容 vs 全面重构”的权衡决策。
 - 重构后的 `adaptive_spec_params.py` 模块划分 (候选步骤解析与配置加载分离) 提升了可测试性。

功能与动机

来自 Issue #23330 和 PR body。用户报告在 Qwen3.5 (hybrid GDN) 模型上启用 adaptive speculative decoding 时遇到 `AttributeError` 和 `CUDA illegal memory access`。PR 详细分析了三个并发问题的根本原因: MTP 构造函数中 `hf_config` 被缓存共享导致 `full_attention_layer_ids` 错误、GDN kernel 的 stride 与物理分配不一致、Mamba 中间池分配步数小于运行时实际最大步数。

实现拆解

1. 隔离 MTP 配置对象: 在 `qwen3_5_mtp.py` 和 `qwen3_next_mtp.py` 构造函数开始时添加 `copy.deepcopy(config)`, 确保 draft 模型的 `hf_config` 修改不污染 target 模型共享的配置。同时在 `hybrid_linear_attn_backend.py` 中重写 `_is_full_attn` 方法, 优先通过 `isinstance(layer, RadixLinearAttention)` 和 `isinstance(layer, RadixAttention)` 路由, 避免依赖可能被污染的 `full_attn_layers` 列表。
2. 修正 GDN 内核步长计算: 在 `fused_sigmoid_gating_recurrent.py` 的 `fused_sigmoid_gating_delta_rule_update` 函数中, 从 `intermediate_states_buffer.shape[1]` 直接推导实际步长, 替代原来依赖调用者传入的 `cache_steps` 参数。这样无论自适应模式如何切换候选步数, 内核访问的始终是物理分配的 stride。
3. 扩大 Mamba 中间池分配: 在 `server_args.py` 中新增 `effective_max_speculative_num_draft_tokens()` 方法, 当启用自适应推测时通过 `resolve_candidate_steps_from_config` 计算 `max(candidate_steps)+1`; 否则返回

`speculative_num_draft_tokens`。在 `model_runner_kv_cache_mixin.py` 的 `_init_pools` 中替换原有直接使用 `speculative_num_draft_tokens` 的地方，确保池分配能容纳所有候选步数。

4. 重构候选步骤解析并更新测试：将候选步骤的默认值 (1,3,7) 和去重、排序逻辑从 `AdaptiveSpeculativeParams.__init__` 内联提取为独立模块级函数 `_resolve_candidate_steps` 和 `resolve_candidate_steps_from_config`；构造函数改为接受 `cfg_path`，由模块函数统一加载配置。测试文件 `test_adaptive_spec_params.py` 适配新签名，新增 `test_params_loads_config_path` 验证文件加载路径，并用 `_make_params_from_config` 辅助方法减少重复。

关键文件：

- `python/sglang/srt/speculative/adaptive_spec_params.py`（模块 推测解码配置；类别 source；类型 core-logic；符号 `_resolve_candidate_steps`, `resolve_candidate_steps_from_config`）：核心配置模块，重构候选步骤解析逻辑，提取 `_resolve_candidate_steps` 和 `resolve_candidate_steps_from_config` 函数，构造函数改为接受 `cfg_path`，是三个 bug 修复中配置加载统一化的关键。
- `python/sglang/srt/layers/attention/hybrid_linear_attn_backend.py`（模块 注意力后端；类别 source；类型 core-logic；符号 `_is_full_attn`）：混合注意力后端，重写 `_is_full_attn` 方法，使用 `isinstance` 优先路由，避免 config 污染导致的路由错误。
- `python/sglang/srt/server_args.py`（模块 服务参数；类别 source；类型 core-logic；符号 `effective_max_speculative_num_draft_tokens`）：新增 `effective_max_speculative_num_draft_tokens` 方法，为池分配提供安全上限，是修复 Mamba buffer 溢出的关键。
- `test/registered/unit/spec/test_adaptive_spec_params.py`（模块 测试；类别 test；类型 test-coverage；符号 `_make_params_from_config`, `test_params_loads_config_path`）：测试适配新构造函数签名，新增 `cfg_path` 加载验证，确保重构后正确性。
- `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）：使用 `effective_max_speculative_num_draft_tokens` 替代原始值初始化池，确保分配安全。
- `python/sglang/srt/layers/attention/fla/fused_sigmoid_gating_recurrent.py`（模块 FLA 内核；类别 source；类型 core-logic）：修复 stride 计算，从 `intermediate_states_buffer.shape[1]` 推导，避免自适应模式下 stride 错误。
- `python/sglang/srt/models/qwen3_5_mtp.py`（模块 模型定义；类别 source；类型 data-contract）：添加 `copy.deepcopy(config)` 隔离 MTP 配置污染。
- `python/sglang/srt/models/qwen3_next_mtp.py`（模块 模型定义；类别 source；类型 data-contract）：同样添加 `copy.deepcopy(config)` 隔离配置污染。
- `python/sglang/srt/speculative/adaptive_runtime_state.py`（模块 推测解码；类别 source；类型 dependency-wiring）：调整导入以匹配重构后的 `adaptive_spec_params` 接口。

关键符号：`_resolve_candidate_steps`, `resolve_candidate_steps_from_config`, `_is_full_attn`, `effective_max_speculative_num_draft_tokens`, `fused_sigmoid_gating_delta_rule_update`

关键源码片段

python/sglang/srt/speculative/adaptive_spec_params.py

核心配置模块，重构候选步骤解析逻辑，提取 `_resolve_candidate_steps` 和 `resolve_candidate_steps_from_config` 函数，构造函数改为接受 `cfg_path`，是三个 bug 修复中配置加载统一化的关键。

```
# _resolve_candidate_steps: 独立函数，负责解析候选步骤，保证 initial_steps 一定在候选集中
# 避免 AdaptiveController.register() 预构建的运行时状态泄漏
def _resolve_candidate_steps(initial_steps: int, cfg: dict[str, object]) -> list[int]:
    """Return sorted, deduplicated candidate steps; inserts *initial_steps* when missing."""
    raw = cfg.get("candidate_steps") or (1, 3, 7) # 默认候选集 [1,3,7]
    candidates: set[int] = set(raw)

    if initial_steps not in candidates:
        log_info_on_rank0(
            logger,
            f"Adding initial speculative_num_steps={initial_steps} to "
            f"candidate_steps={sorted(candidates)} so the pre-built "
            f"runtime state is reused.",
        )
        candidates.add(initial_steps)

    return sorted(candidates)

def resolve_candidate_steps_from_config(
    initial_steps: int, cfg_path: str | None
) -> list[int]:
    """Load adaptive config and resolve candidate steps."""
    cfg = load_adaptive_config(cfg_path)
    return _resolve_candidate_steps(initial_steps, cfg)

class AdaptiveSpeculativeParams:
    def __init__(
        self,
        initial_steps: int,
        cfg_path: str | None = None, # 参数从 dict 改为路径字符串
    ):
        cfg = load_adaptive_config(cfg_path)
        self.candidate_steps = _resolve_candidate_steps(initial_steps, cfg)
        # 剩余初始化 (ema_alpha、warmup_batches 等) 不变
        self.ema_alpha = cfg.get("ema_alpha", 0.2)
        self.update_interval = cfg.get("update_interval", 5)
        self.warmup_batches = cfg.get("warmup_batches", 10)
        # ...
```

python/sglang/srt/layers/attention/hybrid_linear_attn_backend.py

混合注意力后端，重写 `_is_full_attn` 方法，使用 `isinstance` 优先路由，避免 config 污染导致的路由错误。

```
class HybridLinearAttnBackend(AttentionBackend):
    def _is_full_attn(
        self, layer: Optional[RadixAttention], layer_id: Optional[int] = None
    ) -> bool:
        # 优先根据运行时类型分发：即使 full_attn_layers 被配置污染误导，
        # 实际 layer 对象的类型仍能给出正确判断。
        if isinstance(layer, RadixLinearAttention):
            return False
        if isinstance(layer, RadixAttention):
            return True

        # Fallback: 当 layer 对象不可用时，使用 layer_id 回退
        if layer is not None:
            layer_id = layer.layer_id
        assert layer_id is not None, "either layer or layer_id must be provided"
        return layer_id in self.full_attn_layers

# 调用点统一使用 _is_full_attn
def forward_decode(self, ..., layer, ...):
    if self._is_full_attn(layer, kwargs.get("layer_id")):
        return self.full_attn_backend.forward_decode(...)
    else:
        return self.linear_attn_backend.forward_decode(...)
```

python/sglang/srt/server_args.py

新增 `effective_max_speculative_num_draft_tokens` 方法，为池分配提供安全上限，是修复 Mamba buffer 溢出的关键。

```
class ServerArgs:
    def effective_max_speculative_num_draft_tokens(self) -> Optional[int]:
        """Return the maximum draft-token count runtime speculative decoding may use."""
        if self.speculative_num_draft_tokens is None:
            return None
        if not self.speculative_adaptive:
            return self.speculative_num_draft_tokens

        from sglang.srt.speculative.adaptive_spec_params import (
            resolve_candidate_steps_from_config,
        )

        candidate_steps = resolve_candidate_steps_from_config(
            initial_steps=self.speculative_num_steps,
            cfg_path=self.speculative_adaptive_config,
        )
        # 注意：当前自适应推测仅支持 topk=1，因此每个状态需要 steps + 1 个槽位
        # 若未来支持 topk>1，此计算需调整。
```

```
return max(candidate_steps) + 1
```

评论区精华

Review 中最有价值的讨论集中在 `fused_sigmoid_gating_recurrent.py` 的 `stride` 修复方式。Reviewer `alphabetc1` 提出：

Can we change the input `cache_steps` at the caller? Ignoring this parameter and introducing a new `cache_stride_steps` instead feels a bit tricky.

作者 `EanWang211123` 回应：

Yeah I considered this approach, but I found that this function is called in too many places. Changing all of the call sites might require a relatively large amount of work. I think it might be more suitable to open a separate PR to refactor this function.

最终采用在函数内部通过 `shape[1]` 推导 `stride` 的方案，保留了 `cache_steps` 参数仅用于 API 兼容，未改动现有调用者。该权衡被 reviewer 接受，后续 approve。

- GDN `stride` 计算方式的选择 (design): 保留当前方案，在函数内部通过 `shape[1]` 派生 `stride`，`cache_steps` 仅用于 API 兼容。

风险与影响

- 风险：
 - 回归风险（中等）：三个 bug 修复涉及注意力路由、内存分配、配置文件加载等核心路径。`_is_full_attn` 的 `isinstance` 检查依赖新增的 `RadixLinearAttention` 导入，若未来线性注意力类层次调整可能需同步更新。
 - API 兼容性风险（低）：`AdaptiveSpeculativeParams` 构造函数参数从 `config: dict` 改为 `cfg_path: str`，若存在其他直接构造该类的代码（非 `server_args` 路径）会导致 break。测试已适配，内部使用也统一，但建议确认外部是否有调用。
 - 性能风险（低）：`deepcopy config` 增加微小开销，不影响推理性能。
 - 假设约束（中）：`effective_max_speculative_num_draft_tokens` 计算假设 `topk=1`（注释明确），若后续支持 `topk>1` 需重新计算。
- 影响：
 - 用户影响：修复 `Qwen3.5` 和 `Qwen3-Next` 混合 GDN 模型上自适应推测解码的崩溃和准确性问题，`GSM8k` 测试显示准确率保持 0.900，吞吐提升 6.6% (741→790 tok/s)，延迟降低 12.8%。所有使用自适应推测的用户间接受益于更安全的内存分配。
 - 系统影响：`model_runner_kv_cache_mixin.py` 中池分配逻辑变更，影响整个推测解码过程的 `req_to_token_pool` 和 `hybrid_req_to_token_pool` 初始化；`server_args.py` 新增方法可能影响序列化或导出。
 - 团队影响：需要回归测试覆盖非自适应和自适应场景下不同模型（包括纯 Mamba、纯注意力和混合模型）的 CPU/GPU 验收。
 - 风险标记：核心路径变更，配置类接口变更，`topk=1` 假设

关联脉络

- PR #21599 Support adaptive speculative decoding: 本 PR 直接修复自适应推测解码在混合 GDN 模型上的三个 bug, 是 PR #21599 的后续修复。