

# PR #21599 完整报告

sgl-project/sglang

[SPEC][1/N] feat: add adaptive speculative\_num\_steps for EAGLE topk=1

合并时间: 2026-04-21 05:25

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/21599>

## 执行摘要

- 一句话: 为 EAGLE 推测解码添加自适应步数调整, 根据接受长度动态切换运行时状态。
- 推荐动作: 建议技术管理者和工程师精读此 PR, 重点关注: 1) 自适应策略的 EMA 设计和滞后阈值如何平衡响应速度与稳定性; 2) 运行时状态切换机制如何实现零开销原子操作, 避免 CUDA 图重捕获; 3) CUDA 图兼容性检查的风险及潜在解决方案。对于实施类似动态调整的系统具有较高参考价值。

## 功能与动机

PR body 指出, 静态 speculative\_num\_steps 无法适应变化的工作负载: 步数过小导致草稿模型能力未充分利用, 步数过大则生成过多候选令牌被浪费。自适应机制旨在根据实际接受率动态调整步数, 以最大化吞吐量, 并提供了基准数据展示静态步数下的吞吐量瓶颈。

## 实现拆解

1. 新增自适应策略模块 (adaptive\_spec\_params.py): 定义 AdaptiveSpeculativeParams 类, 使用 EMA (默认 alpha=0.2) 跟踪批次平均接受长度, 根据候选步数列表 (默认 [1,3,7]) 和滞后阈值决策何时上下调整步数。关键方法 update() 在满足预热批次和更新间隔后触发 \_recompute\_params() 进行步数切换。
2. 新增运行时状态管理模块 (adaptive\_runtime\_state.py): 引入 SpecRuntimeState 数据类型封装各阶段 (draft、verify、extend) 的注意力后端和 CUDA 图运行器; 定义 AdaptiveSpecWorker 协议要求 worker 实现 build\_adaptive\_runtime\_state() 和 apply\_runtime\_state(); AdaptiveController 负责初始化时预构建所有候选步数的状态池, 并在 on\_verify\_complete() 中根据策略决策原子切换状态。
3. 修改 EAGLEWorker 核心逻辑 (eagle\_worker.py): 集成 AdaptiveController, 在初始化时根据 server\_args.speculative\_adaptive 标志创建控制器并注册初始状态; 实现协议方法以支持状态构建和应用, 通过 \_override\_worker\_state() 临时覆盖服务器参数捕获 CUDA 图。
4. 测试和基准配套: 新增单元测试 test\_adaptive\_spec\_params.py 验证策略逻辑 (如预热、间隔、滞后阈值), 新增端到端测试 test\_adaptive\_speculative.py 启动真实服务器测试状态切换和 GSM8K 准确性; 新增基准脚本 bench\_adaptive\_speculative.py 对比自适应与静态服务器性能。

5. 文档和配置更新：新增文档 `adaptive_speculative_decoding.md` 说明使用方法和参数，更新 `server_args.py` 添加 `speculative_adaptive` 和 `speculative_adaptive_config` 命令行选项。

关键文件：

- `python/sglang/srt/speculative/adaptive_spec_params.py`（模块 推测解码；类别 source；类型 core-logic；符号 `load_adaptive_config`, `AdaptiveSpeculativeParams`, `init`, `update`）：定义了自适应策略核心类 `AdaptiveSpeculativeParams`，负责 EMA 跟踪接受长度和步数决策逻辑。
- `python/sglang/srt/speculative/adaptive_runtime_state.py`（模块 推测解码；类别 source；类型 core-logic；符号 `SpecRuntimeState`, `AdaptiveSpecWorker`, `build_adaptive_runtime_state`, `apply_runtime_state`）：定义了运行时状态管理核心，包括 `SpecRuntimeState` 数据类、`AdaptiveSpecWorker` 协议和 `AdaptiveController` 控制器，负责状态池和原子切换。
- `python/sglang/srt/speculative/eagle_worker.py`（模块 推测解码；类别 source；类型 core-logic；符号 `apply_runtime_state`, `build_adaptive_runtime_state`, `_override_worker_state`, `init`）：修改 `EAGLEWorker` 以支持自适应协议，集成 `AdaptiveController` 并实现状态构建和应用方法，是功能落地的核心。
- `test/registered/unit/spec/test_adaptive_spec_params.py`（模块 单元测试；类别 test；类型 test-coverage；符号 `TestAdaptiveSpeculativeParams`, `test_initial_steps_snap_to_nearest_candidate_preferring_larger_step`, `test_update_respects_warmup_and_interval`, `test_empty_batches_do_not_consume_warmup_or_shift_steps`）：新增单元测试，验证 `AdaptiveSpeculativeParams` 的策略逻辑，覆盖预热、间隔、滞后阈值等关键行为。
- `benchmark/bench_adaptive_speculative.py`（模块 基准测试；类别 source；类型 dependency-wiring；符号 `build_phase_plan`, `send_request`, `run_phase`, `summarize_phases`）：新增基准测试脚本，对比自适应与静态服务器的吞吐量、延迟和接受长度，用于性能验证。

关键符号：`AdaptiveSpeculativeParams.update`, `AdaptiveController.on_verify_complete`, `EAGLEWorker.build_adaptive_runtime_state`, `EAGLEWorker.apply_runtime_state`

## 关键源码片段

### `python/sglang/srt/speculative/adaptive_runtime_state.py`

定义了运行时状态管理核心，包括 `SpecRuntimeState` 数据类、`AdaptiveSpecWorker` 协议和 `AdaptiveController` 控制器，负责状态池和原子切换。

```
@dataclass
class SpecRuntimeState:
    """A complete set of runtime resources bound to a specific speculative decoding
    configuration.

    Each decode round runs three stages — draft, verify, extend — and every
    stage has shape-dependent resources (attention backends and CUDA graphs)
```

that must match the current configuration. Switching adaptive steps means swapping the entire state atomically.

```
"""
speculative_num_steps: int # 当前步数配置
speculative_num_draft_tokens: int # 对应草稿令牌数 (steps+1)
draft_attn_backend: "AttentionBackend | None" # 草稿阶段注意力后端
cuda_graph_runner: "EAGLEDraftCudaGraphRunner | None" # 草稿阶段 CUDA 图运行器
target_attn_backend: "AttentionBackend" # 验证阶段注意力后端
target_graph_runner: "CudaGraphRunner | CPUGraphRunner | None" # 验证阶段图运行器
draft_extend_attn_backend: "AttentionBackend | None" # 扩展阶段注意力后端
cuda_graph_runner_for_draft_extend: "EAGLEDraftExtendCudaGraphRunner | None" #
扩展阶段 CUDA 图运行器
```

```
class AdaptiveController:
```

```
    """Facade that owns adaptive decision-making and runtime state switching.
```

```
    Works with any worker that implements ``AdaptiveSpecWorker`` protocol.
```

```
    """
```

```
    def __init__(self, worker: AdaptiveSpecWorker, config_path: str | None = None):
```

```
        self.worker = worker
```

```
        cfg = load_adaptive_config(config_path) # 加载配置文件
```

```
        self.params = AdaptiveSpeculativeParams(
```

```
            initial_steps=worker.speculative_num_steps, config=cfg
```

```
        )
```

```
        self._states: dict[int, SpecRuntimeState] = {} # 状态池, 键为步数
```

```
    def init_states(self) -> None:
```

```
        """Build and register runtime states for all candidate steps."""
```

```
        for steps in self.params.candidate_steps:
```

```
            if steps in self._states:
```

```
                continue
```

```
            # 委托 worker 构建对应步数的运行时状态
```

```
            state = self.worker.build_adaptive_runtime_state(
```

```
                speculative_num_steps=steps,
```

```
                speculative_num_draft_tokens=steps + 1,
```

```
            )
```

```
            self._states[steps] = state
```

```
        self._activate(self.params.current_steps) # 激活初始状态
```

```
    def on_verify_complete(self, accept_lengths: list[int]) -> None:
```

```
        """Feed verify results; switch runtime state if EMA warrants it."""
```

```
        if self.params.update(accept_lengths): # 如果策略决策步数变化
```

```
            self._activate(self.params.current_steps) # 激活新步数对应状态
```

## [python/sglang/srt/speculative/eagle\\_worker.py](#)

修改 EAGLEWorker 以支持自适应协议, 集成 AdaptiveController 并实现状态构建和应用方法, 是功能落地的核心。

```
def __init__(self, server_args: ServerArgs, gpu_id: int, tp_rank: int, dp_rank: Optional[int],
```

```

        moe_ep_rank: int, attn_cp_rank: int, moe_dp_rank: int, nccl_port: int,
        target_worker: TpModelWorker):
# ... 原有初始化代码 ...
self.speculative_num_steps = server_args.speculative_num_steps
self.speculative_num_draft_tokens = server_args.speculative_num_draft_tokens

# Adaptive speculative
self.adaptive_controller: Optional[AdaptiveController] = None
if server_args.speculative_adaptive: # 检查是否启用自适应
    self.adaptive_controller = AdaptiveController(
        self, config_path=server_args.speculative_adaptive_config
    )
# ... 后续初始化 ...
if self.adaptive_controller is not None:
    # 注册初始运行时状态, 以便控制器管理
    self.adaptive_controller.register(
        SpecRuntimeState(
            speculative_num_steps=self.speculative_num_steps,
            speculative_num_draft_tokens=self.speculative_num_draft_tokens,
            draft_attn_backend=self.draft_attn_backend,
            cuda_graph_runner=self.cuda_graph_runner,
            target_attn_backend=self.target_worker.model_runner.attn_backend,
            target_graph_runner=self.target_worker.model_runner.graph_runner,
            draft_extend_attn_backend=self.draft_extend_attn_backend,
            cuda_graph_runner_for_draft_extend=self.cuda_graph_runner_for_draft_extend,
        )
    )
    self.adaptive_controller.init_states() # 预构建所有候选步数的状态

def apply_runtime_state(self, state: SpecRuntimeState) -> None:
    """Apply a pre-built runtime state to this worker."""
    if self.speculative_num_steps == state.speculative_num_steps:
        return # 步数未变, 无需切换
    logger.info(f"Switch adaptive runtime state: steps {self.speculative_num_steps} -> {state.speculative_num_steps}")
    self.speculative_num_steps = state.speculative_num_steps
    self.speculative_num_draft_tokens = state.speculative_num_draft_tokens
    # 原子更新各阶段资源引用
    self.draft_attn_backend = state.draft_attn_backend
    self.cuda_graph_runner = state.cuda_graph_runner
    self.target_worker.model_runner.attn_backend = state.target_attn_backend
    self.target_worker.model_runner.graph_runner = state.target_graph_runner
    self.draft_extend_attn_backend = state.draft_extend_attn_backend
    self.cuda_graph_runner_for_draft_extend = state.cuda_graph_runner_for_draft_extend

```

## 评论区精华

- 设计简化建议: gemini-code-assist 建议简化 AdaptiveSpeculativeParams.\_\_init\_\_ 中的 current\_steps 初始化逻辑, 作者可能已采纳以提高可读性。

- CUDA 图兼容性风险: chatgpt-codex-connector 警告在 cuda\_graph\_runner.py 和 eagle\_draft\_extend\_cuda\_graph\_runner.py 中新增的令牌数检查可能导致 DP 注意力模式下各 rank 执行路径分歧, 引发同步问题或性能回归; 此问题在 review 中未明确解决。
- 测试覆盖要求: Qiaolin-Yu 要求添加单元测试和端到端测试, 作者后续提交中添加了 test\_adaptive\_spec\_params.py 和 test\_adaptive\_speculative.py, 满足要求并获得批准。
- 类型注解改进: Qiaolin-Yu 建议在 adaptive\_runtime\_state.py 中使用具体类型而非 object, 作者回复“done”表示已修复。
- 功能扩展讨论: issue 评论中 liquanfeng 询问是否支持 EAGLE3, 作者回复已添加支持。
  - 简化自适应参数初始化逻辑 (design): 作者可能已采纳, 但 review 中未显式确认修改。
  - CUDA 图令牌数检查引入同步风险 (correctness): 此风险在 review 中未明确解决, 可能需后续验证或修复。
  - 要求添加单元测试和端到端测试 (testing): 作者后续提交中添加了 test\_adaptive\_spec\_params.py 和 test\_adaptive\_speculative.py, 满足要求并获得批准。

## 风险与影响

- 风险:
  - CUDA 图同步风险: cuda\_graph\_runner.py 中新增的 is\_num\_tokens\_supported 检查在 DP 注意力模式下可能因各 rank 本地令牌数不同导致 CUDA 图与 eager 执行路径分歧, 引发集体操作不同步或崩溃 (chatgpt-codex-connector 指出)。
  - 自适应策略振荡: AdaptiveSpeculativeParams 的 EMA 参数和滞后阈值配置不当可能导致步数频繁切换, 增加开销并降低吞吐量; 默认参数虽经过调优, 但需用户理解配置。
  - 兼容性限制: 目前仅支持 EAGLE topk=1, PR body 提到 EAGLEWorkerV2 支持为后续 TODO, 且未覆盖其他推测算法 (如 N-gram), 限制了使用范围。
  - 新增复杂性: 运行时状态池管理和原子切换增加了代码复杂度, 可能引入隐蔽 bug, 尤其是在多线程或分布式环境下。
  - 配置错误静默失败: chatgpt-codex-connector 指出 speculative\_adaptive\_config 标志未自动启用自适应行为, 若用户仅提供配置路径而遗漏 --speculative-adaptive, 系统将静默回退到静态解码。
- 影响:
  - 用户影响: 为用户提供了动态优化推测解码吞吐量的能力, 尤其是在接受率变化的工作负载下可提升性能; 新增命令行选项和配置文件增加了使用灵活性, 但需学习新参数。
  - 系统影响: 在核心推测解码路径中引入状态切换逻辑, 可能轻微增加运行时开销, 但设计为零开销引用交换; 预构建多个 CUDA 图可能增加内存占用, 但避免了在线重捕获。
  - 团队影响: 为推测解码模块增加了重要功能, 需要维护新组件和测试; 与近期历史 PR (如 PR 22908、22832) 同属推测解码改进线, 展现了该领域的持续投入。
  - 风险标记: CUDA 图同步风险, 自适应策略振荡, topk=1 限制, 配置静默失败

## 关联脉络

- PR #22908 [AMD] Resolve Qwen3.5 MTP (speculative decoding) radix cache conflict.: 同属推测解码功能线，涉及设备感知和缓存冲突修复，本 PR 的自适应机制可能与此类底层优化交互。
- PR #22832 [sgl] fix incorrect behavior in cuda graph draft extend: 修复 CUDA 图推测解码扩展逻辑，本 PR 新增自适应状态切换同样依赖 CUDA 图运行器，需关注兼容性。
- PR #22088 [sgl] add support for weight update function in spedec: 扩展 EAGLE 推测解码工作者功能，本 PR 在此基础上新增自适应参数调整，属于同一模块的持续演进。