

PR #25303 完整报告

sgl-project/sglang

[Spec]: Make Triton standalone spec test deterministic

合并时间: 2026-05-19 08:36

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/25303>

执行摘要

- 一句话: 使 Triton 独立推测解码测试确定化
- 推荐动作: 值得精读: PR body 中对非确定性原因的排查思路 (逐项检查掩码、位置、KV 槽、接受逻辑) 和实验对比 (纯确定性 vs 仅目标验证统一注意力) 展示了系统性的根因分析。配置级修复优先于内核修改的设计决策值得学习。

功能与动机

关联 Issue #22101 指出独立推测解码的 `test_gsm8k` 测试在贪心温度 =0 下准确率波动大 (0.635-0.864), 疑似验证 / 接受路径存在非确定性。PR 旨在稳定 Triton 后端的 CI 测试, 避免因数值漂移导致误判。

实现拆解

1. 新增 `enable_deterministic_inference` 属性: 在 `StandaloneServerBase` (文件 `python/sglang/test/server_fixtures/standalone_fixture.py`) 中新增类属性, 默认 `False`。此属性仅在 Triton 注意力后端测试中使用。
2. 条件追加服务端标志: 在 `get_server_args()` 方法中, 当 `enable_deterministic_inference` 为 `True` 时, 向服务启动参数添加 `--enable-deterministic-inference` 标志。此标志启用 SGLang 已有的确定性推理模式。
3. Triton 测试类启用: 在 `test/registered/spec/test_spec_standalone_extra.py` 中, 为 `TestStandaloneSpeculativeDecodingTriton` 类设置 `enable_deterministic_inference = True`。FA3 和 FlashInfer 后端保持不变, 因为它们不受该确定性问题的显著影响。

关键文件:

- `python/sglang/test/server_fixtures/standalone_fixture.py` (模块 测试夹具; 类别 `test`; 类型 `test-coverage`): 核心变更: 在测试夹具基类中新增 `enable_deterministic_inference` 属性, 并在服务启动参数中有条件追加 `--enable-deterministic-inference` 标志。此文件是测试基础设施的关键扩展点。
- `test/registered/spec/test_spec_standalone_extra.py` (模块 测试用例; 类别 `test`; 类型 `test-coverage`): 具体测试类变更: 为 `TestStandaloneSpeculativeDecodingTriton` 类设置 `enable_deterministic_inference = True`, 将标志引入实际测试。FA3 和 FlashInfer 保持默认 `False`。

关键符号: 未识别

关键源码片段

python/sglang/test/server_fixtures/standalone_fixture.py

核心变更：在测试夹具基类中新增 `enable_deterministic_inference` 属性，并在服务启动参数中有条件追加 `--enable-deterministic-inference` 标志。此文件是测试基础设施的关键扩展点。

新增类属性，默认 False，仅 Triton 后端测试启用

```
class StandaloneServerBase:
```

```
...
```

```
    enable_deterministic_inference: bool = False # 新增，默认关闭
```

```
    @classmethod
```

```
    def get_server_args(cls):
```

```
        assert cls.attention_backend, f"{cls.__name__} must set `attention_backend`"
```

```
        args = [
```

```
            "--trust-remote-code",
```

```
            "--cuda-graph-max-bs", "8",
```

```
            "--speculative-algorithm", "STANDALONE",
```

```
            "--speculative-draft-model-path", DEFAULT_DRAFT_MODEL_STANDALONE,
```

```
            "--speculative-num-steps", str(cls.speculative_num_steps),
```

```
            "--speculative-eagle-topk", str(cls.speculative_eagle_topk),
```

```
            "--speculative-num-draft-tokens", str(cls.speculative_num_draft_tokens),
```

```
            "--mem-fraction-static", 0.7,
```

```
            "--attention-backend", cls.attention_backend,
```

```
        ]
```

```
        # 仅当启用确定性推理时，追加标志
```

```
        if cls.enable_deterministic_inference:
```

```
            args.append("--enable-deterministic-inference")
```

```
        return args
```

test/registered/spec/test_spec_standalone_extra.py

具体测试类变更：为 `TestStandaloneSpeculativeDecodingTriton` 类设置

`enable_deterministic_inference = True`，将标志引入实际测试。FA3 和 FlashInfer 保持默认 False。

```
class TestStandaloneSpeculativeDecodingTriton(StandaloneServerBase, CustomTestCase):
```

```
    attention_backend = "triton"
```

```
    speculative_eagle_topk = 2
```

```
    speculative_num_draft_tokens = 7
```

```
    enable_spec_v2 = False
```

```
    enable_deterministic_inference = True # 新增，仅 Triton 后端启用确定性模式
```

```
class TestStandaloneSpeculativeDecodingFlashinfer(StandaloneServerBase, CustomTestCase):
```

```
    attention_backend = "flashinfer"
```

```
    speculative_eagle_topk = 2
```

```
    speculative_num_draft_tokens = 7
```

```
    enable_spec_v2 = False
```

```
    # enable_deterministic_inference 默认 False，不启用
```

评论区精华

仅有一个批准评论，无直接讨论。但 PR body 中提供了详尽的实验数据和分析，包括 100 次运行的准确性分布和性能对比，证明了确定性模式的有效性和不引入性能损失的权衡。

- 暂无高价值评论线程

风险与影响

- 风险：低风险：此变更仅影响测试配置，不修改任何产品代码或推理内核。确定性模式是 SGLang 已支持的选项，启用后准确率虽略下降（均值 0.6923 vs 0.7139），但方差大幅降低，能可靠通过 0.69 阈值。FA3 和 FlashInfer 后端不变，Triton 后端测试的更保守基准是工程设计上的合理选择。
- 影响：影响范围有限：只改变 CI 中 Triton 后端独立推测解码测试的配置，使其运行在确定性模式下。其他用户工作流不受影响。社区和 CI 验证将获得更稳定的准确率分数，减少失败重试和误报。
- 风险标记：测试覆盖变更，无产品代码修改

关联脉络

- PR #22101 [Bug] Standalone spec decoding gsm8k accuracy is flaky with high variance: 直接关联的 Issue，报告了相同的测试不稳定性问题，提供了问题背景和影响范围。
- PR #22100 [Workaround] Relax accuracy threshold for standalone spec decoding: 此前为缓解不稳定问题而放宽了准确率阈值，PR 中提及作为临时方案。