

PR #43314 完整报告

vllm-project/vllm

[CI] Fix test_lora_with_spec_decode on V2 model runner

合并时间: 2026-05-22 14:24

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/43314>

执行摘要

- 一句话: 修复 V2 Runner 中 LoRA CUDA Graph 捕获遗漏问题
- 推荐动作: 建议精读, 特别是理解 CUDA Graph 捕获时 LoRA 内核被跳过的根本原因和修复方案。该 PR 展示了动态图捕获与 Python 层短路之间的微妙交互, 对理解 vLLM 的 LoRA 实现有重要参考价值。

功能与动机

CI 测试 `tests/v1/e2e/spec_decode/test_lora_with_spec_decode.py::test_batch_inference_correctness` 在 PR #43139 引入后失败率从 0% 升至 24% (76/100)。根本原因是 V2 Runner 在 CUDA Graph 捕获时使用 `LoRAMapping(0,0)` 导致 Python 层短路, LoRA 内核未被捕获, 回放时静默跳过 LoRA 计算。

实现拆解

1. 在 `model_runner.py` 的 dummy run 分支中, 保留原有的 `set_adapter_mapping(LoRAMapping(0,0))` 调用 (来自 #43139 的 workaround), 该调用设置所有 token 的 LoRA 索引为 -1。
2. 新增对每个唯一 `PunicaWrapper` 的 `kernel_meta` 属性修改: 遍历 `adapter_manager.punica_wrapper_mapping`, 对 `token_mapping_meta` 和 `prompt_mapping_meta` 设置 `no_lora_flag_cpu[0] = False` 和 `num_active_loras_cpu[0] = 1`。这确保 Python 层不会因全 -1 索引而短路, 从而使得 LoRA 内核被 CUDA Graph 捕获。
3. 捕获时内核仍不会执行有效计算: 由于 `active_lora_ids` 仍为全 -1, 每个 CTA 在 GPU 内核中通过 `if lora_id == -1: return` 提前退出。
4. 测试文件优化: 将 `prompts`、`lora_request`、`sampling_params` 变量声明提前, 对 LLM 实例使用 `try...finally` 确保资源释放, 改进匹配率断言消息。

关键文件:

- `vllm/v1/worker/gpu/model_runner.py` (模块 模型执行器; 类别 source; 类型 core-logic; 符号 `execute_model`): 核心修复文件, 在 dummy run 分支中新增对 `PunicaWrapper kernel_meta` 的修改, 确保 LoRA 内核被 CUDA Graph 捕获。
- `tests/v1/e2e/spec_decode/test_lora_with_spec_decode.py` (模块 测试; 类别 test; 类型 test-coverage; 符号 `test_batch_inference_correctness`): 测试文件, 修复间接依赖问题

并改进资源清理，确保 CI 稳定性。

关键符号：execute_model

关键源码片段

vllm/v1/worker/gpu/model_runner.py

核心修复文件，在 dummy run 分支中新增对 PunicaWrapper kernel_meta 的修改，确保 LoRA 内核被 CUDA Graph 捕获。

```
# vllm/v1/worker/gpu/model_runner.py (dummy run branch inside execute_model)
if self.lora_config:
    # 原有 workaround: 将所有 token 映射到 index 0, 经 convert_mapping 后变为全 -1
    adapter_manager = self.lora_manager._adapter_manager
    adapter_manager.set_adapter_mapping(
        LoRAMapping(
            index_mapping=(0,) * input_batch.num_tokens_after_padding,
            prompt_mapping=(0,) * input_batch.num_reqs,
            is_prefill=True,
        )
    )
    # 新增修复: 强制让 Python 层不短路, 使得 CUDA Graph 捕获到 LoRA 内核
    seen_wrappers: set[int] = set()
    for punica_wrapper in adapter_manager.punica_wrapper_mapping.values():
        if id(punica_wrapper) in seen_wrappers:
            continue
        seen_wrappers.add(id(punica_wrapper))
        # 对 token 和 prompt 映射元数据均进行覆写
        for kernel_meta in (
            punica_wrapper.token_mapping_meta, # 类型: KernelMeta
            punica_wrapper.prompt_mapping_meta, # 类型: KernelMeta
        ):
            # 设置 no_lora_flag=False 使得 prepare_tensors 不会短路
            kernel_meta.no_lora_flag_cpu[0] = False
            # 设置 num_active_loras=1 使得内核被调度, 但 active_lora_ids 仍为全 -1
            kernel_meta.num_active_loras_cpu[0] = 1
```

tests/v1/e2e/spec_decode/test_lora_with_spec_decode.py

测试文件，修复间接依赖问题并改进资源清理，确保 CI 稳定性。

```
# tests/v1/e2e/spec_decode/test_lora_with_spec_decode.py (部分)
prompts = [LORA_TEST_PROMPT_MAP[lora_path]] * 100
lora_request = LoRARequest("adapter", 1, lora_path)
sampling_params = SamplingParams(
    temperature=0.0, top_p=1.0, top_k=-1, seed=SEED, max_tokens=128
)

# 无 speculative decoding 的参考 LLM
ref_llm = LLM(model=model_name, ..., enable_lora=True, ...)
```

```

try:
    ref_outputs = ref_llm.generate(prompts, sampling_params, lora_request=lora_request)
finally:
    del ref_llm
    torch.accelerator.empty_cache()
    cleanup_dist_env_and_memory()

# 带 speculative decoding 的测试 LLM
lora_spec_llm = LLM(model=model_name, ..., speculative_config=..., enable_lora=True, ...)
try:
    lora_spec_outputs = lora_spec_llm.generate(prompts, sampling_params, lora_request=lora_request)
    matches = 0
    for ref_output, spec_output in zip(ref_outputs, lora_spec_outputs):
        if ref_output.outputs[0].text == spec_output.outputs[0].text:
            matches += 1
        else:
            print(f"ref_output: {ref_output.outputs[0].text}")
            print(f"spec_output: {spec_output.outputs[0].text}")
    # Heuristic threshold: 90%
    threshold = int(0.90 * len(ref_outputs))
    print(f"match ratio: {matches}/{len(ref_outputs)}")
    assert matches > threshold, (
        f"match ratio {matches}/{len(ref_outputs)} <= {threshold}"
    )
finally:
    del lora_spec_llm
    torch.accelerator.empty_cache()
    cleanup_dist_env_and_memory()

```

评论区精华

gemini-code-assist[bot] 提出测试中断言 `matches > threshold` 与注释 "at least 90%" 不一致，建议改为 `matches >= threshold`。该问题未在可见讨论中明确解决，但实际提交代码仍使用 `>`，可能存在微小 flaky 风险。

- 测试断言阈值逻辑与注释不一致 (correctness): 未采纳（或未回应），实际提交仍使用 `>`。可能存在微小 flaky 风险。
- 修复代码方案确认 (design): 修订按 suggestion 实现，保持原有 `set_adapter_mapping` 调用不变，新增覆写逻辑。

风险与影响

- 风险:
 1. 回归风险: 改动仅影响 V2 Model Runner 的 dummy run 分支 (CUDA Graph 捕获路径)，对正常推理路径无影响。

2. 性能风险: 在 dummy run 中遍历所有 PunicaWrapper 并设置元数据, 但数量通常很小 (<10), 开销可忽略。
3. 兼容性: 仅当 lora_config 存在且 execute_model 进入 dummy run 分支时触发, 不影响无 LoRA 或非 V2 Runner 场景。

- 影响:

1. 测试稳定性: 修复 V2 Runner + LoRA + Spec Decode 的集成测试, 使匹配率从 76/100 恢复至 100/100。
2. 用户: 无感知, 但确保 CUDA Graph 优化与 LoRA 兼容。
3. 系统: 避免静默错误 (推理结果错误但无报错)。 - 风险标记: 核心路径变更, 测试断言边界条件未对齐

关联脉络

- PR #43139 [Model Runner V2] Fix lora Triton Error [CUDA]: device-side assert triggered: 本 PR 引入的 workaround (LoRAMapping(0,0)) 是此 PR 的修复根因, 导致 CUDA Graph 捕获时 LoRA 内核被跳过。