

PR #26805 完整报告

sgl-project/sglang

Add the KV-canary verify JIT kernel and reference implementation

合并时间: 2026-05-31 09:52

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/26805>

执行摘要

- 一句话: 新增 KV-canary 验证 JIT 内核与参考实现
- 推荐动作: 推荐精读此 PR。理由: 1) 展示了 JIT 核心理念——用 Python 描述 CUDA 内核并通过 TVM FFI 调用, 这在项目中具有代表性; 2) 链式哈希验证的设计模式值得学习; 3) Review 中关于 CUDA 安全性的讨论对于编写正确内核有借鉴意义。此外, 建议关注其与其系列 PR (如 #26806、#26807) 的关联。

功能与动机

KV-canary 旨在对 KV 缓存进行轻量级、确定性验证, 及时发现内存损坏、计算错误等问题。该 PR 实现验证内核, 作为整个验证流程的消费端: 先有写内核写入 canary 数据, 再有验证内核读取并检查完整性。此 PR 还随附了无关独立用途的共享内核基础代码 (常量、公共头文件) 以及内核无关的测试支架。

实现拆解

使用 Markdown 按以下步骤拆解实现过程:

1. 常量与数据结构定义(`consts.py`, `csrc/kv_canary/consts.cuh`) - 定义 canary 缓冲区字段偏移、故障原因 (`FailReason`)、链锚点常量等。 - 提供 `splitmix64` 哈希函数用于链式校验。
2. CUDA 验证内核实现(`csrc/kv_canary/canary_verify.cuh`) - 实现 `CanaryVerifyKernel`, 每个线程处理一个验证条目。 - 从 `verify_plan` 中读取槽索引、期望位置、前一槽索引, 从 `canary_buf` 中加载存储值, 执行链式哈希匹配、位置匹配、令牌匹配 (可选), 若不一致则向违规环写入报告。 - 注册健康计数器 (`slot_run_counter`、`kernel_run_counter`)。
3. Python 绑定与数据结构(`verify.py`) - 定义 `CanaryLaunchTag` 枚举区分启动类型。 - `VerifyOrWriteContext` 数据类持有共享上下文张量。 - `VerifyPlan` 数据类持有验证计划张量 (`verify_slot_indices`、`verify_expected_tokens` 等)。 - `launch_canary_verify_kernel` 函数利用 JIT 内核机制 (`load_jit`) 加载编译 CUDA 内核并启动。
4. 参考实现(`verify_ref.py`) - `launch_canary_verify_kernel_torch_reference` 函数使用纯 PyTorch 操作逐条目执行与 CUDA 内核相同的验证逻辑, 用于输出对比测试。
5. 基准测试与单元测试 - `bench_verify.py` 使用 Triton 基准框架对 `verify` 内核进行多维度性能测试。 - `test_const_sync.py` 验证 Python 常量与 CUDA 头文件常量的一致性。 - `test_utils.py` 测试基准用例生成逻辑。

关键文件:

- `python/sglang/jit_kernel/kv_canary/verify.py` (模块 验证内核; 类别 `source`; 类型 `core-logic`; 符号 `CanaryLaunchTag`, `VerifyOrWriteContext`, `VerifyPlan`, `launch_canary_verify_kernel`) : 核心 Python 绑定文件, 定义验证内核的启动入口、数据结构 `VerifyPlan` 和 `VerifyOrWriteContext`, 以及枚举 `CanaryLaunchTag`; 是 JIT 内核与上层管理的桥梁。
- `python/sglang/jit_kernel/csrc/kv_canary/canary_verify.cuh` (模块 验证内核; 类别 `other`; 类型 `core-logic`; 符号 `CanaryVerifyKernel`, `run`) : CUDA 验证内核实现, 是功能的核心: 每个线程处理一个验证条目, 执行链式哈希校验和数据比对, 并写入违规环。
- `python/sglang/jit_kernel/kv_canary/verify_ref.py` (模块 参考实现; 类别 `source`; 类型 `core-logic`; 符号 `launch_canary_verify_kernel_torch_reference`, `_to_signed_int64`, `compute_slot_hash`) : 纯 PyTorch 参考实现, 用于验证 CUDA 内核正确性; 包含链式哈希计算和故障报告逻辑。
- `python/sglang/jit_kernel/kv_canary/consts.py` (模块 常量定义; 类别 `source`; 类型 `core-logic`; 符号 `FailReason`, `splitmix64`, `splitmix64_mix3`) : 定义 `canary` 数据布局、故障原因枚举、链锚点及 `splitmix64` 哈希函数, 是内核和参考实现的基础。
- `python/sglang/jit_kernel/tests/kv_canary/test_const_sync.py` (模块 同步测试; 类别 `test`; 类型 `test-coverage`; 符号 `test_int_consts_sync`, `test_enums_sync`) : 自动验证 Python 和 CUDA 头文件中的常量定义是否一致, 确保跨语言同步, 是正确性的重要保证。

关键符号: `launch_canary_verify_kernel`, `launch_canary_verify_kernel_torch_reference`, `VerifyPlan.allocate`, `splitmix64`, `splitmix64_mix3`, `ComputeSlotHash`, `CanaryVerifyKernel.run`

关键源码片段

`python/sglang/jit_kernel/kv_canary/verify.py`

核心 Python 绑定文件, 定义验证内核的启动入口、数据结构 `VerifyPlan` 和 `VerifyOrWriteContext`, 以及枚举 `CanaryLaunchTag`; 是 JIT 内核与上层管理的桥梁。

```
# python/sglang/jit_kernel/kv_canary/verify.py
```

```
class CanaryLaunchTag(IntEnum):
```

```
    """Unique tag per (head | tail) × (K | V) × (FULL | SWA) launch."""
```

```
    HEAD_K_FULL = 0
```

```
    HEAD_V_FULL = 1
```

```
    # ... 其余 6 个 tag 类似
```

```
@dataclass(frozen=True, slots=True, kw_only=True)
```

```
class VerifyPlan:
```

```
    """Flat verify entries consumed by launch_canary_verify_kernel."""
```

```
    verify_slot_indices: torch.Tensor # shape [verify_capacity], int64
```

```
    verify_expected_tokens: torch.Tensor # shape [verify_capacity], int64, -1 表示跳过
```

```
    verify_expected_positions: torch.Tensor
```

```
    verify_prev_slot_indices: torch.Tensor # -1 表示链头
```

```
    verify_num_valid: torch.Tensor # shape [1], int32, 实际有效条目数
```

```
enable: torch.Tensor # shape [1], int32, 1= 运行
```

```
@classmethod
```

```
def allocate(cls, *, verify_capacity: int, device: torch.device) -> "VerifyPlan":
```

```
    # 分配 CUDA graph 捕获尺寸的固定缓冲区, torch.zeros 填充
```

```
    return cls(
```

```
        verify_slot_indices=torch.zeros(verify_capacity, dtype=torch.int64, device=device),
```

```
        verify_expected_tokens=torch.zeros(verify_capacity, dtype=torch.int64, device=device),
```

```
        # ... 其他字段类似
```

```
        verify_num_valid=torch.zeros(1, dtype=torch.int32, device=device),
```

```
        enable=torch.ones(1, dtype=torch.int32, device=device),
```

```
    )
```

```
def launch_canary_verify_kernel(*, context: VerifyOrWriteContext, plan: VerifyPlan) -> None:
```

```
    """编译并启动 CUDA verify 内核。"""
```

```
    _assert_contiguous(context.canary_buf, "canary_buf")
```

```
    # 构建编译参数, 加载 TVM 模块
```

```
    args = make_cpp_args(USE_AGGREGATE_VIOLATIONS="true", ZERO_INIT_OKAY="true")
```

```
    module = _jit_canary_verify_module(args) # cache_once 装饰
```

```
    # 准备内核启动参数 (通过 Python 函数设置 grid/block)
```

```
    module["canary::CanaryVerifyKernel<{}>::run".format(", ".join(args))](
```

```
        # 传递所有张量指针和标量维度
```

```
    )
```

python/sglang/jit_kernel/csrc/kv_canary/canary_verify.cuh

CUDA 验证内核实现, 是功能的核心: 每个线程处理一个验证条目, 执行链式哈希校验和数据比对, 并写入违规环。

```
// csrc/kv_canary/canary_verify.cuh
```

```
namespace canary {
```

```
template <bool USE_AGGREGATE_VIOLATIONS, bool ZERO_INIT_OKAY>
```

```
struct CanaryVerifyKernel {
```

```
    static __global__ void run(const Params p) {
```

```
        int tid = blockIdx.x * blockDim.x + threadIdx.x;
```

```
        if (tid == 0) {
```

```
            atomicAdd(reinterpret_cast<long long*>(p.kernel_run_counter), 1LL);
```

```
        }
```

```
        int32_t active = min(*p.verify_num_valid, p.verify_capacity);
```

```
        if (active <= 0) return; // 避免负数 active 转无符号越界
```

```
        if (tid >= static_cast<uint32_t>(active)) return;
```

```
        // 加载验证条目
```

```
        int64_t slot_idx = p.verify_slot_indices[tid];
```

```
        if (slot_idx == TOKEN_TO_KV_SLOT_PADDING) return; // 跳过填充槽
```

```
        // ... 加载期望值与存储值, 执行比较
```

```
        // 若不一致, 原子写入 violation_ring
```

```

    }
};

} // namespace canary

```

python/sglang/jit_kernel/kv_canary/verify_ref.py

纯 PyTorch 参考实现，用于验证 CUDA 内核正确性；包含链式哈希计算和故障报告逻辑。

```

# python/sglang/jit_kernel/kv_canary/verify_ref.py

def launch_canary_verify_kernel_torch_reference(
    *, context: VerifyOrWriteContext, plan: VerifyPlan, check_verify_expected_token: bool
) -> None:
    # 将数据拷贝到 CPU，逐条处理
    num_valid = int(plan.verify_num_valid[0].item())
    active = min(num_valid, plan.verify_slot_indices.shape[0])
    if active <= 0:
        return
    slot_indices = plan.verify_slot_indices[:active].tolist()
    expected_positions = plan.verify_expected_positions[:active].tolist()
    # ... 省略具体循环
    for k in range(active):
        s = slot_indices[k]
        if s == consts.TOKEN_TO_KV_SLOT_PADDING:
            continue
        # 加载存储值
        stored_token = buf_i64[s, consts.CANARY_FIELD_TOKEN].item()
        # 计算预期链哈希
        prev_slot = prev_slot_indices[k]
        if prev_slot != consts.TOKEN_TO_KV_SLOT_PADDING:
            expected_chain_hash = compute_slot_hash(buf_i64, prev_slot)
        else:
            expected_chain_hash = stored_chain_hash
        # 比较并记录违规
        fail_reason = consts.FailReason(0)
        if stored_chain_hash != expected_chain_hash:
            fail_reason |= consts.FailReason.VERIFY_CHAIN_HASH_MISMATCH
        if fail_reason:
            # 将违规写入 violation_ring 并原子更新索引
            ...

```

评论区精华

Review 中核心讨论包括：

- 内存越界风险：若 `verify_num_valid` 为负数（未初始化或错误），`active` 为负数后转为 `uint32_t` 会变成极大值，导致越界访问。建议添加 `active <= 0` 提前返回。
- 未初始化内存：`VerifyPlan.allocate` 中使用 `torch.empty` 创建 `verify_num_valid`，可能在 `plan` 内核执行前被 `verify` 内核读到垃圾值。建议改用 `torch.zeros`。

- 对齐要求: `slot_stride_bytes` 需为 8 的倍数, 否则 `canary_load_field` 的 `int64_t*` 转换导致未对齐访问。建议添加运行时对齐检查。
- 指针类型转换: 对 `int64_t*` 的 `atomicAdd` 应使用 `long long*` 而非 `unsigned long long*`。
- 模板参数格式化: 在 `f-string` 中直接写入元组会生成无效 C++, 应使用 `' , '.join(args)`。
- 标量提取简化: 将 `new_empty().copy_().item()` 简化为 `[0].item()`。

这些建议涉及正确性、性能和代码风格, 作者权衡后可能采纳了部分, 但 PR 已合并, 未确认是否全部修正。

- 负数 `active` 导致内存越界 (`correctness`): 建议合理, 但 PR 合并前未确认是否采纳。
- 未初始化 `verify_num_valid` (`correctness`): 建议合理, 可能当时已私下修复, 但 PR 中未体现更改。
- `slot_stride_bytes` 未对齐风险 (`performance`): 建议中等优先级, 可能部分内核已保证对齐但仍可提高健壮性。
- 指针类型转换与原子操作 (`style`): 代码风格建议, 不影响正确性但更清晰。
- 模板参数格式化与标量提取简化 (`style`): 两者均为样式改进, 不影响功能。

风险与影响

- 风险: 技术风险包括:
 - 未初始化 `verify_num_valid` 导致越界: `verify.py` 中使用 `torch.empty` 而非 `torch.zeros`, 若 `plan` 内核未写入就执行 `verify` 内核, 可能读到垃圾大数导致 GPU 越界, 引发挂起或错误结果。
 - 对齐失败: 若 `slot_stride_bytes` 非 8 倍数, CUDA 内核中 `int64_t*` 存取会未对齐, 导致性能下降或 `segfault`。
 - 整数符号转化: 负数 `active` 转为无符号后变成巨大值, 边界检查失效。
 - 参考实现未以相同方式处理: 参考实现同样可能含有缺陷, 但至少确保主机端与设备端一致。
 - 性能开销: 每次推理都启动多个 `verify` 内核, 可能增加 `latency`, 但 `benchmark` 可量化。目前 `benchmark` 仅含基础性能数据, 未与线上实际负载结合评估。
- 影响: 影响分析:
 - 对用户: 若启用 `KV-canary`, 该功能会正常收集验证结果并报告违规, 增强对 `KV` 缓存正确性的信心。默认不启用, 对现有用户无影响。
 - 对系统: 引入新的 CUDA 内核和 Python 组件, 增加二进制体积和启动时间。运行时会根据配置决定是否调用。
 - 对团队: 代码设计模块化, 易于扩展新的验证模式或适配不同注意力机制。测试和基准覆盖较好, 降低了回归风险。
- 风险标记: 未初始化内存风险, 负数 `active` 越界, 对齐未检查

关联脉络

- PR #26806 Add the `KV-canary` write JIT kernel and reference implementation: 写内核是对应的上游步骤, 先写 `canary` 数据再有 `verify` 内核验证, 两者构成完整的写入 - 验证链。

- PR #26807 Add the KV-canary plan JIT kernels: 计划内核负责编排验证条目，输出 verify 内核所需的 verify_plan，是 verify 的直接上游。