

# PR #18016 完整报告

sgl-project/sglang

[Feature] Add SiMM as sglang HiCache Storage backend

合并时间: 2026-04-14 08:12

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/18016>

## 执行摘要

- 一句话: 添加 SiMM 作为 HiCache 分布式存储后端, 支持 RDMA 零拷贝缓存加速。
- 推荐动作: 建议精读此 PR, 重点关注存储后端集成模式, 如配置优先级设计、NUMA 感知优化以及 RDMA 集成。对于计划扩展分布式缓存的团队, 可参考 HiCacheSiMM 的接口实现和错误处理机制。

## 功能与动机

根据 PR body 描述, SiMM 是一个针对 AI 工作负载的分布式、高性能、弹性缓存加速层。集成 SiMM 作为 HiCache 后端是为了解决传统 GPU/CPU 缓存容量有限的问题, 通过 RDMA 网络实现跨服务器的分布式内存池, 提供近乎无限的缓存存储空间, 从而提升大模型多轮对话等场景的吞吐和延迟性能。基准测试显示在 DeepSeek R1 模型上, 使用 SiMM 后请求吞吐和令牌吞吐均有提升。

## 实现拆解

1. 新增后端实现类: 在 `python/sglang/srt/mem_cache/storage/simm/hicache_simm.py` 中创建 `HiCacheSiMM` 类, 继承自 `HiCacheStorage`, 实现 `set/get` 等存储接口。关键符号包括 `SiMMConfig` 配置类、`get_current_process_numa` 和 `get_numa_nic_mapping` 辅助函数, 支持从 `extra_config` 或环境变量加载配置。
2. 扩展命令行选项: 修改 `python/sglang/srt/server_args.py`, 在 `--hicache-storage-backend` 的 `choices` 中添加 `'simm'` 选项, 使用户可以通过 CLI 启用 SiMM 后端。
3. 注册后端到工厂: 更新 `python/sglang/srt/mem_cache/storage/backend_factory.py`, 在 `_create_builtin_backend` 函数中添加 `'simm'` 分支, 并调用 `register_backend` 注册 `HiCacheSiMM`, 使其能被存储工厂动态创建。
4. 调整缓存控制器: 修改 `python/sglang/srt/managers/cache_controller.py` 中的 `attach_storage_backend` 方法, 将 `'simm'` 添加到支持 `interface_v1` 的后端列表中, 确保与现有接口兼容。
5. 补充测试和文档: 新增 `python/sglang/srt/mem_cache/storage/simm/test_simm.py` 包含单操作和批量操作测试; 添加 `python/sglang/srt/mem_cache/storage/simm/README.md` 详细说明安装、部署和配置步骤。

关键文件:

- python/sclang/srt/mem\_cache/storage/simm/hicache\_simm.py (模块 存储后端; 类别 source; 类型 core-logic; 符号 SiMMConfig, from\_file, load\_from\_extra\_config, get\_current\_process\_numa) : 核心后端实现, 定义 HiCacheSiMM 类和 SiMMConfig, 处理与 SiMM 库的交互及 NUMA 优化。
- python/sclang/srt/mem\_cache/storage/simm/test\_simm.py (模块 存储后端; 类别 test; 类型 test-coverage; 符号 generate\_batch\_query\_keys, create\_mock\_host\_kv\_cache, MockHostKVCache, test\_single\_operation) : 单元测试文件, 验证 SiMM 后端的基本存储操作和批量接口。
- python/sclang/srt/server\_args.py (模块 服务器参数; 类别 source; 类型 configuration) : 命令行参数扩展, 添加 'simm' 到 HiCache 存储后端选项列表。
- python/sclang/srt/mem\_cache/storage/backend\_factory.py (模块 存储后端; 类别 source; 类型 dependency-wiring) : 存储后端工厂类更新, 注册并实例化 HiCacheSiMM 后端。
- python/sclang/srt/managers/cache\_controller.py (模块 缓存控制; 类别 source; 类型 entrypoint) : 缓存控制器调整, 将 'simm' 纳入支持 interface\_v1 的后端类型。
- python/sclang/srt/mem\_cache/storage/simm/README.md (模块 存储后端; 类别 docs ; 类型 documentation) : 文档说明, 详细介绍 SiMM 安装、部署和配置方法。

关键符号: SiMMConfig.from\_file, SiMMConfig.load\_from\_extra\_config, get\_current\_process\_numa, get\_numa\_nic\_mapping, HiCacheSiMM.init, HiCacheSiMM.warmup, generate\_batch\_query\_keys, create\_mock\_host\_kv\_cache, test\_single\_operation, test\_batch\_operation

## 关键源码片段

### python/sclang/srt/mem\_cache/storage/simm/hicache\_simm.py

核心后端实现, 定义 HiCacheSiMM 类和 SiMMConfig, 处理与 SiMM 库的交互及 NUMA 优化。

```
import json
import logging
import os
from dataclasses import dataclass
from typing import Dict, List

from sclang.srt.mem_cache.hicache_storage import HiCacheStorage, HiCacheStorageConfig
from sclang.srt.mem_cache.memory_pool_host import HostKVCache

# 尝试导入第三方SiMM库, 若失败则抛出友好错误提示
try:
    from simm.kv import Store
except ImportError as e:
    raise ImportError(
        "请按照 https://github.com/scitix/SiMM 安装SiMM以使用SGLang的SiMM后端。"
    ) from e
```

```
logger = logging.getLogger(__name__)
```

```
@dataclass
```

```
class SiMMConfig:
```

```
    """SiMM后端配置类，支持从extra_config字典或环境变量加载。"""
```

```
    manager_address: str # SiMM管理服务地址
```

```
    clnt_threadpool_size: int = 10 # 客户端线程池大小
```

```
    enable_profile: bool = False # 是否启用性能分析
```

```
@staticmethod
```

```
def load_from_extra_config(extra_config: dict) -> "SiMMConfig":
```

```
    """从extra_config字典加载配置，优先使用此方式。"""
```

```
    if "manager_address" not in extra_config:
```

```
        raise ValueError("extra_config中必须包含manager_address字段")
```

```
    return SiMMConfig(
```

```
        manager_address=extra_config.get("manager_address"),
```

```
        clnt_threadpool_size=extra_config.get("clnt_threadpool_size", 10),
```

```
        enable_profile=extra_config.get("enable_profile", False),
```

```
    )
```

```
def get_current_process_numa() -> int:
```

```
    """获取当前进程所在的NUMA节点编号，失败时返回-1。"""
```

```
    try:
```

```
        with open("/proc/self/stat", "r") as f:
```

```
            fields = f.read().split()
```

```
            if len(fields) < 39:
```

```
                return -1
```

```
            current_cpu = int(fields[38]) # 第39字段为处理器编号
```

```
            # 通过sysfs解析NUMA节点
```

```
            numa_path = f"/sys/devices/system/cpu/cpu{current_cpu}/node0"
```

```
            if os.path.exists(numa_path) and os.path.islink(numa_path):
```

```
                import re
```

```
                link_target = os.readlink(numa_path)
```

```
                match = re.search(r"node(\d+)$", link_target)
```

```
                if match:
```

```
                    return int(match.group(1))
```

```
            return -1
```

```
    except Exception:
```

```
        return -1
```

```
class HiCacheSiMM(HiCacheStorage):
```

```
    """HiCache的SiMM存储后端实现。"""
```

```
    def __init__(
```

```
        self, storage_config: HiCacheStorageConfig = None, mem_pool: HostKVCache = None
```

```
    ):
```

```
        super().__init__(storage_config, mem_pool)
```

```
        extra_config = getattr(storage_config, "extra_config", None) if storage_config else None
```

```
        # 优先从extra_config加载配置，确保与CLI参数一致
```

```
        if extra_config is not None and extra_config.get("manager_address"):
```

```
self.config = SiMMConfig.load_from_extra_config(extra_config)
else:
    # 降级到环境变量（已不推荐，此处为向后兼容）
    self.config = SiMMConfig.from_file()
# 初始化SiMM存储客户端，启用NUMA感知的RDMA设备选择
self.store = Store(self.config.manager_address, numa_aware=True)
logger.info(f"SiMM后端初始化完成，管理地址: {self.config.manager_address}")
```

## 评论区精华

- 环境变量配置：reviewer 建议将所有环境变量移至 `--hicache-storage-backend-extra-config` 以保持一致性，作者回应已移除环境变量并支持 `extra-config` 优先加载。
- 日志和错误信息：reviewer 指出导入错误日志中误用 `'vLLM'`，应改为 `'SGLang'`，并建议添加更详细的调试信息；作者在后续提交中更新了日志内容。
- 代码风格优化：reviewer 建议使用 `TYPE_CHECKING` 包装导入以减少运行时开销，并将 `for-loop` 改为列表推导式以提升可读性。
- 文档完整性：reviewer 发现文档中 SiMM 官方链接缺失，作者回应将在 SiMM 开源后更新链接，并已添加安装说明到 README 中。
- Docker 文件位置：reviewer 询问是否应将 Dockerfile 移至 `3rdparty` 目录，作者最终删除 Dockerfile 并将安装指令整合到 README 中。
- 环境变量配置优化 (design)：采纳建议，配置优先级改为 `extra-config` 优先，减少环境变量依赖。
- 日志错误信息完善 (documentation)：作者更新日志内容，确保错误信息准确并提供解决指引。
- 代码风格和性能优化 (style)：部分优化被采纳，后续提交中调整了代码结构以提升可维护性。
- Docker 文件位置争议 (infra)：移除独立 Dockerfile，简化部署流程，依赖文档说明。

## 风险与影响

- 风险：- 依赖风险：SiMM 为第三方库，其稳定性和兼容性直接影响 SGLang 运行；若 SiMM 服务未正确部署或网络异常，可能导致缓存失败或性能下降。具体体现在 `hicache_simm.py` 的 `import` 块和 `Store` 初始化中。
- 配置复杂性：支持多种配置源（`extra-config`、环境变量、文件），配置错误可能引发运行时异常，如 `SiMMConfig.load_from_extra_config` 缺少 `manager_address` 时的 `ValueError`。
- 性能回归风险：RDMA 网络依赖和 NUMA 绑定可能在某些硬件环境下引入额外延迟，需确保 `get_numa_nic_mapping` 正确识别网卡，否则回退逻辑可能影响吞吐。
- 测试覆盖不足：当前测试集中在基础操作，缺少分布式场景和故障恢复的集成测试，可能掩盖生产环境问题。
- 影响：- 用户影响：用户可通过 `--hicache-storage-backend simm` 启用分布式缓存，提升大容量 KV 场景的扩展性，但需额外部署 SiMM 集群并配置 RDMA 网络，增加了使用复杂度。
- 系统影响：扩展了 HiCache 存储后端生态，支持更多弹性缓存方案；对现有存储接口无破坏性变更，兼容性良好。

- 团队影响：为存储团队引入新的维护点，需关注 SiMM 版本更新和问题排查；性能优化团队可基于此后端进一步探索 RDMA 优化机会。
- 风险标记：依赖第三方库，配置复杂性，性能波动，测试覆盖不足

## 关联脉络

- PR #22767 [HiCache] Fix memory host free logic when share\_indices\_with\_anchor enabled: 同样涉及 HiCache 内存管理，关注缓存释放逻辑，可参考内存泄漏修复模式。
- PR #22790 Refactor streaming session abort handling: 涉及流式会话和缓存交互，重构了中止处理，与本 PR 的存储后端扩展有协同测试价值。