

PR #26295 完整报告

sgl-project/sglang

Refactor HiCache stack dispatch into strategies

合并时间: 2026-05-26 00:06

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/26295>

执行摘要

- 一句话: 重构 HiCache 堆栈分发为策略模式
- 推荐动作:

功能与动机

原有的 `attach_hybrid_pool_to_unified_cache` 函数通过 `if/elif` 链判断不同的 KV 池类型并执行对应的构建逻辑, 导致代码冗长且难以扩展。为了支持更多 KV 池组合 (如 DeepSeekV4、Mamba、SWA 等), 需要一种更灵活、可插拔的方式。PR 描述中提到 'Replace the if/elif chain in `attach_hybrid_pool_to_unified_cache` with a `_StackStrategy` registry', 明确表达了重构动机。

实现拆解

1. 定义策略基类和结果数据结构: 在 `hybrid_pool_assembler.py` 中新增 `StackStrategy` 抽象基类和 `StackBuildResult` 数据类。`StackStrategy` 声明 `matches(kvcache, components)` 和 `build(...)` 方法, 子类通过 `matches` 判断是否适用于当前 KV 池和组件集合, `build` 执行具体的堆栈构建逻辑并返回 `StackBuildResult`。
2. 实现具体策略类: 为每种支持的 KV 池组合创建策略类:
 - `_DeepSeekV4Strategy`: 用于 DeepSeekV4 (需要 FULL + SWA 组件)。
 - `_MambaStrategy`: 用于 Mamba 池 (需要 FULL + MAMBA 组件)。
 - `_SwaStrategy`: 用于 SWA 池 (需要 FULL + SWA 组件)。
 - `_DsaStrategy`: 用于 DSA 池 (需要 FULL 组件)。
 - `_PlainKvStrategy`: 作为 fallback, 处理普通 KV 池 (需要 FULL 组件), 并显式排除 DeepSeekV4、SWA、Mamba 等特殊情况, 避免错误匹配。
3. 引入策略注册和选择机制: 维护一个全局列表 `_STRATEGIES`, 通过 `_select_strategy(kvcache, components)` 遍历列表并返回第一个 `matches` 返回 True 的策略。新增 `register_stack_strategy(strategy)` 函数, 允许下游 fork 将自定义策略添加到列表头部, 实现优先匹配。
4. 通用应用函数: 新增 `_apply_stack_result(cache, result, kvcache, params, server_args)` 函数, 将 `StackBuildResult` 中构建好的 `HostPoolGroup`、`HybridCacheController` 以及各个组件的 `host pool` 和 `sidecar` 附着到 `UnifiedRadixCache` 对象上。该函数通过 `_COMPONENT_HOST_ATTR` 映射自动将 `host`

pool 设置到对应的组件属性，减少了每个策略中的重复代码。

5. 重构入口函数：保留 `attach_hybrid_pool_to_unified_cache` 的对外签名，但内部实现简化为：调用 `_select_strategy` 获取策略，调用策略的 `build` 方法，然后调用 `_apply_stack_result`。同时将原来内联的本地导入 (`DeepSeekV4TokenToKVPool` 等) 分散到各策略类中，使代码更加模块化。
6. 新增单元测试：创建 `test_unified_radix_hicache_dispatch.py`，使用 Mock 对象模拟各种 KV 池和组件组合，测试策略选择正确性、注册顺序、fallback 行为、未知组合异常抛出以及自定义策略注册功能。

关键文件：

- `python/sglang/srt/mem_cache/hybrid_cache/hybrid_pool_assembler.py` (模块 混合缓存；类别 `source`；类型 `core-logic`；符号 `_COMPONENT_HOST_ATTR`, `StackBuildResult`, `StackStrategy`, `matches`)：核心重构文件：将条件分支替换为策略模式，新增策略基类、具体策略、注册选择机制和通用应用函数。
- `test/registered/unit/mem_cache/test_unified_radix_hicache_dispatch.py` (模块 单元测试；类别 `test`；类型 `test-coverage`；符号 `_mock_kvcache`, `TestUnifiedRadixHiCacheDispatch`, `test_strategy_registry_ordering`, `test_deepseek_v4_full_swa`)：新增单测，全面验证策略注册顺序、选择逻辑、fallback 和自定义策略扩展，是保证重构正确性的关键。

关键符号：`_select_strategy`, `_apply_stack_result`, `register_stack_strategy`, `StackStrategy.matches`, `StackStrategy.build`, `_DeepSeekV4Strategy.matches`, `_DeepSeekV4Strategy.build`, `_MambaStrategy.matches`, `_MambaStrategy.build`, `_SwaStrategy.matches`, `_SwaStrategy.build`, `_DsaStrategy.matches`, `_DsaStrategy.build`, `_PlainKvStrategy.matches`, `_PlainKvStrategy.build`, `attach_hybrid_pool_to_unified_cache`

关键源码片段

`python/sglang/srt/mem_cache/hybrid_cache/hybrid_pool_assembler.py`

核心重构文件：将条件分支替换为策略模式，新增策略基类、具体策略、注册选择机制和通用应用函数。

```
# --- 策略基类与结果数据类 ---
```

```
@dataclass
```

```
class StackBuildResult:
```

```
    host_pool_group: HostPoolGroup
    cache_controller: HybridCacheController
    component_host_pools: dict[ComponentType, Any]
    sidecars: list[SidecarPoolSpec] = field(default_factory=list)
    register_req_to_token_counter: bool = False
    transfer_layer_num: int = 0
    pools_desc: str = ""
```

```
class StackStrategy:
```

```
    """策略基类，用于判断是否匹配给定 KV 池与组件集合，并执行构建。"""
```

```

def matches(self, kvcache: Any, components: set[ComponentType]) -> bool:
    raise NotImplementedError

def build(self, *, cache, kvcache, params, server_args, load_cache_event,
          attn_cp_group=None, attn_tp_group=None) -> StackBuildResult:
    raise NotImplementedError

# --- 具体策略示例: DeepSeekV4 ---
class _DeepSeekV4Strategy(StackStrategy):
    def matches(self, kvcache, components):
        from sglang.srt.mem_cache.deepseek_v4_memory_pool import DeepSeekV4TokenToKVPool
        return isinstance(kvcache, DeepSeekV4TokenToKVPool) and components == {
            ComponentType.FULL, ComponentType.SWA}

    def build(self, *, cache, kvcache, params, server_args, load_cache_event,
              attn_cp_group=None, attn_tp_group=None):
        # 具体构建逻辑, 最终返回 StackBuildResult
        ...

# --- 入口函数: 保持接口不变, 内部使用策略 ---
def attach_hybrid_pool_to_unified_cache(cache, params, server_args, *,
                                       load_cache_event, attn_cp_group=None,
                                       attn_tp_group=None) -> None:
    kvcache = params.token_to_kv_pool_allocator.get_kvcache()
    components = set(cache.components.keys())
    strategy = _select_strategy(kvcache, components) # 遍历 _STRATEGIES 匹配
    result = strategy.build(...)
    _apply_stack_result(cache, result, kvcache, params, server_args)

```

test/registered/unit/mem_cache/test_unified_radix_hicache_dispatch.py

新增单测, 全面验证策略注册顺序、选择逻辑、fallback 和自定义策略扩展, 是保证重构正确性的关键。

```

class TestUnifiedRadixHiCacheDispatch(unittest.TestCase):
    # 验证策略注册顺序: _DeepSeekV4Strategy 必须在 _SwaStrategy 之前
    def test_strategy_registry_ordering(self):
        order = [type(s) for s in _STRATEGIES]
        self.assertLess(order.index(_DeepSeekV4Strategy),
                        order.index(_SwaStrategy))
        self.assertEqual(order[-1], _PlainKvStrategy)

    # 验证 DeepSeekV4 使用 FULL+SWA 时选择 _DeepSeekV4Strategy
    def test_deepseek_v4_full_swa(self):
        from sglang.srt.mem_cache.deepseek_v4_memory_pool import DeepSeekV4TokenToKVPool
        kvcache = _mock_kvcache(DeepSeekV4TokenToKVPool)
        strategy = _select_strategy(kvcache, {FULL, SWA})
        self.assertIsInstance(strategy, _DeepSeekV4Strategy)

    # 未知组合应抛出 AssertionError

```

```
def test_unknown_combo_raises(self):
    from sglang.srt.mem_cache.deepseek_v4_memory_pool import DeepSeekV4TokenToKVPool
    from sglang.srt.mem_cache.swa_memory_pool import SWAKVPool
    for cls in (SWAKVPool, DeepSeekV4TokenToKVPool):
        kvcache = _mock_kvcache(cls)
        with self.assertRaises(AssertionError) as cm:
            _select_strategy(kvcache, {FULL})
        self.assertIn("No matching HiCache strategy", str(cm.exception))
```

评论区精华

- `gemini-code-assist[bot]` 提出 `_PlainKvStrategy.matches` 应显式排除 `DeepSeekV4TokenToKVPool`，避免当 `DeepSeekV4` 缓存仅配置 `FULL` 组件时被错误匹配（应抛出 `AssertionError`）。同时建议在 `test_unknown_combo_raises` 中增加对 `DeepSeekV4TokenToKVPool` 的测试。这些建议均在最终代码中得到采纳。
- `hzh0425` 表达了对重构方案的赞赏，并批准 PR：“That's amazing! I've thought about it for a long time, but couldn't come up with a suitable refactoring approach. This is so cool!”
- `_PlainKvStrategy` 应显式排除 `DeepSeekV4TokenToKVPool` (correctness): 作者已采纳并在最终代码中实现排除逻辑，测试用例 `test_unknown_combo_raises` 中也包含了 `DeepSeekV4TokenToKVPool`，验证其会引发 `AssertionError`。
- 测试用例应覆盖 `DeepSeekV4` 错误配置场景 (testing): 已采纳，测试文件中已有该场景。

风险与影响

- 风险：
 - 策略顺序敏感性：`_select_strategy` 顺序遍历 `_STRATEGIES` 列表，策略注册顺序决定了匹配优先级。若新加入的策略匹配条件过于宽泛，可能拦截本应由其他策略处理的场景。需要开发者理解 `isinstance` 的继承关系（如 `DeepSeekV4TokenToKVPool` 继承自 `BaseSWAKVPool`），并确保更具体的策略排在前面。当前通过 `_DeepSeekV4Strategy` 在 `_SwaStrategy` 前注册来规避。
 - 回归风险：重构涉及对核心路径 `attach_hybrid_pool_to_unified_cache` 的重写，虽然单元测试覆盖了各种组合，但端到端集成测试（如 `radix cache KL` 测试）已通过（15/15），降低了回归概率。但若未来新增 KV 池类型而忘记注册策略，将导致 `AssertionError`，可能影响线上部署。
 - 与下游 fork 的兼容性：新增的 `register_stack_strategy` API 为下游提供了扩展点，但若下游已直接使用内部导入的旧符号（如 `build_anchor_sidecar_stack`），需要适配新的策略化 API。
- 影响：
 - 对用户：无外部 API 变更，功能完全兼容，用户无感知。
 - 对系统：重构后策略匹配和调度路径几乎没有性能变化（单次 `isinstance` 对比），内存开销增加几个策略对象，可忽略。

- 对团队：代码可读性和可扩展性显著提升。后续添加新 KV 池组合（如 FlashAttention3 支持）只需新增策略子类并注册，无需修改核心函数。
- 对测试：新增的单元测试可直接在 CPU 环境运行（注册为 base-a-test-cpu），CI 覆盖更全面。
- 风险标记：策略顺序敏感，回归风险，下游兼容性

关联脉络

- PR #26249 [hisparse]: update user guide: 同样涉及 HiCache (HiSparse) 模块，虽然侧重点不同，但同属 HiCache 功能线，未来可能一起演进。