

PR #22782 完整报告

sgl-project/sglang

[HiCache]Fix CP support for hybrid model

合并时间: 2026-04-15 23:50

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/22782>

执行摘要

- 一句话: 为混合模型 (如 KV+Mamba) 的 HiCache 添加上下文并行 (CP) 支持。
- 推荐动作: 该 PR 值得快速浏览, 以了解 HiCache 如何扩展支持混合模型的上下文并行。重点关注 HybridCacheController 的参数扩展设计, 以及 review 中关于参数来源不一致的讨论, 这反映了代码库中潜在的配置管理问题。

功能与动机

PR 标题和提交信息表明, 此变更旨在为混合模型 (hybrid model) 添加 CP (上下文并行) 支持。PR 正文未明确说明动机, 但从代码变更和 review 讨论推断, 这是为了扩展 HiCache 功能, 使其在混合架构 (如结合 KV 缓存和 Mamba 线性注意力) 中支持上下文并行, 确保缓存控制器能正确接收并处理 CP 相关的配置参数 (attn_cp_rank 和 attn_cp_size)。

实现拆解

1. 扩展 HybridCacheController 的构造函数: 在 python/sglang/srt/mem_cache/hybrid_cache/hybrid_cache_controller.py 中, 为 HybridCacheController.__init__ 方法新增两个参数 attn_cp_rank (默认 0) 和 attn_cp_size (默认 1), 并将它们传递给父类 BaseHiCacheController.__init__, 以便缓存控制器能感知 CP 配置。
2. 更新混合堆栈构造函数: 在 python/sglang/srt/mem_cache/hybrid_cache/hybrid_pool_assembler.py 中, 修改 build_nsa_hybrid_stack 和 build_mamba_hybrid_stack 函数, 在创建 HybridCacheController 实例时, 从 params 对象传入 attn_cp_rank 和 attn_cp_size 参数, 确保 CP 配置在两种混合缓存堆栈中均被正确设置。
3. 未包含测试或配置配套改动: 本次变更仅涉及源码文件, 没有添加或修改测试文件、配置文件或部署脚本, 表明这是一个核心逻辑的增量扩展。

关键文件:

- python/sglang/srt/mem_cache/hybrid_cache/hybrid_cache_controller.py (模块 混合缓存; 类别 source; 类型 entrypoint; 符号 HybridCacheController.init): 这是 HiCache 混合缓存控制器的入口点, 新增 CP 参数扩展了其配置能力, 直接影响缓存初始化。
- 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 build_nsa_hybrid_stack, build_mamba_hybrid_stack): 包含两个混合缓存堆栈构造函数, 负责创建 HybridCacheController 实例, 新增 CP 参数传递确保 CP 支持生效。

关键符号: HybridCacheController.init, build_nsa_hybrid_stack,
build_mamba_hybrid_stack

关键源码片段

[python/sglang/srt/mem_cache/hybrid_cache/hybrid_cache_controller.py](#)

这是 HiCache 混合缓存控制器的入口点, 新增 CP 参数扩展了其配置能力, 直接影响缓存初始化。

```
class HybridCacheController(BaseHiCacheController):
    def __init__(
        self,
        token_to_kv_pool_allocator: BaseTokenToKVPoolAllocator,
        mem_pool_host: Any,
        page_size: int,
        tp_group: torch.distributed.ProcessGroup,
        load_cache_event: threading.Event,
        write_policy: str = "write_through_selective",
        io_backend: str = "",
        storage_backend: Optional[str] = None,
        prefetch_threshold: int = 256,
        model_name: Optional[str] = None,
        storage_backend_extra_config: Optional[dict] = None,
        pp_rank: int = 0,
        pp_size: int = 1,
        attn_cp_rank: int = 0, # 新增: 上下文并行排名, 默认0
        attn_cp_size: int = 1, # 新增: 上下文并行大小, 默认1
        transfer_layer_num: Optional[int] = None,
        enable_storage_metrics: bool = False,
    ):
        startup_storage_backend = storage_backend
        super().__init__(
            token_to_kv_pool_allocator=token_to_kv_pool_allocator,
            mem_pool_host=mem_pool_host,
            page_size=page_size,
            tp_group=tp_group,
            load_cache_event=load_cache_event,
            write_policy=write_policy,
            io_backend=io_backend,
            storage_backend=None,
            prefetch_threshold=prefetch_threshold,
            model_name=model_name,
            storage_backend_extra_config=storage_backend_extra_config,
            pp_rank=pp_rank,
            pp_size=pp_size,
            attn_cp_rank=attn_cp_rank, # 将CP参数传递给父类BaseHiCacheController
            attn_cp_size=attn_cp_size,
            enable_storage_metrics=enable_storage_metrics,
        )
```

```
# ... 其余初始化逻辑保持不变
```

python/sclang/srt/mem_cache/hybrid_cache/hybrid_pool_assembler.py

包含两个混合缓存堆栈构造函数，负责创建 HybridCacheController 实例，新增 CP 参数传递确保 CP 支持生效。

```
def build_nsa_hybrid_stack(
    radix_cache: "HiRadixCache",
    params: "CacheInitParams",
    server_args: "ServerArgs",
    *,
    layer_num: int,
    prefetch_threshold: int,
    enable_storage_metrics: bool = False,
    extra_config: Optional[dict] = None,
):
    # ... 前置逻辑：构建host_pool_group等
    cache_controller = HybridCacheController(
        params.token_to_kv_pool_allocator,
        host_pool_group,
        radix_cache.page_size,
        radix_cache.tp_group,
        load_cache_event=load_cache_event,
        write_policy=server_args.hicache_write_policy,
        io_backend=server_args.hicache_io_backend,
        storage_backend=server_args.hicache_storage_backend,
        prefetch_threshold=prefetch_threshold,
        model_name=server_args.served_model_name,
        storage_backend_extra_config=extra_config,
        pp_rank=radix_cache.pp_rank, # 注意：pp参数来自radix_cache，而非params
        pp_size=radix_cache.pp_size,
        attn_cp_rank=params.attn_cp_rank, # 新增：从params传入CP排名
        attn_cp_size=params.attn_cp_size, # 新增：从params传入CP大小
        transfer_layer_num=layer_num,
        enable_storage_metrics=enable_storage_metrics,
    )
    # ... 后续设置和日志

def build_mamba_hybrid_stack(
    mamba_cache: "HiMambaRadixCache",
    params: "CacheInitParams",
    server_args: "ServerArgs",
    *,
    transfer_layer_num: int,
    prefetch_threshold: int,
    enable_storage_metrics: bool = False,
    extra_config: Optional[dict] = None,
):
    # ... 前置逻辑
```

```
cache_controller = HybridCacheController(
    params.token_to_kv_pool_allocator,
    host_pool_group,
    params.page_size,
    params.tp_cache_group,
    load_cache_event=load_cache_event,
    write_policy=server_args.hicache_write_policy,
    io_backend=server_args.hicache_io_backend,
    storage_backend=server_args.hicache_storage_backend,
    prefetch_threshold=prefetch_threshold,
    model_name=server_args.served_model_name,
    storage_backend_extra_config=extra_config,
    pp_rank=params.pp_rank, # pp参数来自params, 与上一个函数来源不同
    pp_size=params.pp_size,
    attn_cp_rank=params.attn_cp_rank, # 新增: 同样从params传入CP参数
    attn_cp_size=params.attn_cp_size,
    transfer_layer_num=transfer_layer_num,
    enable_storage_metrics=enable_storage_metrics,
)
# ... 后续设置和日志
```

评论区精华

review 中只有一条来自 `gemini-code-assist[bot]` 的评论，指出 `hybrid_pool_assembler.py` 中 `build_nsa_hybrid_stack` 函数使用 `radix_cache.pp_rank` 和 `radix_cache.pp_size`，而 `build_mamba_hybrid_stack` 使用 `params.pp_rank` 和 `params.pp_size`，存在不一致性。建议统一使用 `params` 以确保参数来源一致和最新。但此评论未被采纳或回应，PR 最终以当前状态合并，表明团队可能认为现有实现是可接受的，或计划后续处理。

- 参数来源不一致性 (design): 评论未被采纳，PR 以当前状态合并，表明团队可能接受现有差异或计划后续处理。

风险与影响

- 风险：1. 兼容性风险：新增参数 `attn_cp_rank` 和 `attn_cp_size` 带有默认值（0 和 1），向后兼容，不会破坏现有调用。2. 一致性风险：review 中指出的 `pp_rank` 和 `pp_size` 参数来源不一致（`radix_cache` vs `params`）可能导致在不同混合堆栈中流水线并行配置不一致，但风险较低，因为两者可能指向相同值。3. 测试覆盖不足：变更未附带单元测试或集成测试，无法验证 CP 支持在混合模型场景下的正确性，存在潜在回归风险。4. 性能影响：新增参数传递是轻量级的，预计对性能无影响。
- 影响：1. 对用户的影响：用户在使用混合模型（如 KV+Mamba）并启用上下文并行时，HiCache 将能正确初始化，支持更高效的分布式推理。2. 对系统的影响：扩展了 HiCache 控制器的配置能力，使其能处理 CP 参数，提升系统在混合架构下的灵活性和可扩展性。3. 对团队的影响：为后续开发混合模型 CP 功能奠定了基础，但需注意 review 中未解决的一致性疑虑。
- 风险标记：参数来源不一致，缺少测试覆盖

关联脉络

- PR #22767 [HiCache] Fix memory host free logic when share_indices_with_anchor enabled: 同属 HiCache 模块的 bugfix, 涉及内存管理, 但本 PR 是功能扩展, 关联性较弱。
- PR #22882 [HiSparse][BugFix]: Fix the memory leak issue during health checks.: 同属缓存相关模块 (HiSparse) 的修复, 但本 PR 针对 HybridCache, 关联性一般。