



南湖架构 HuanCun Non-blocking Cache 设计与实现

王凯帆 蔺嘉炜 张林隽

2022/8/25

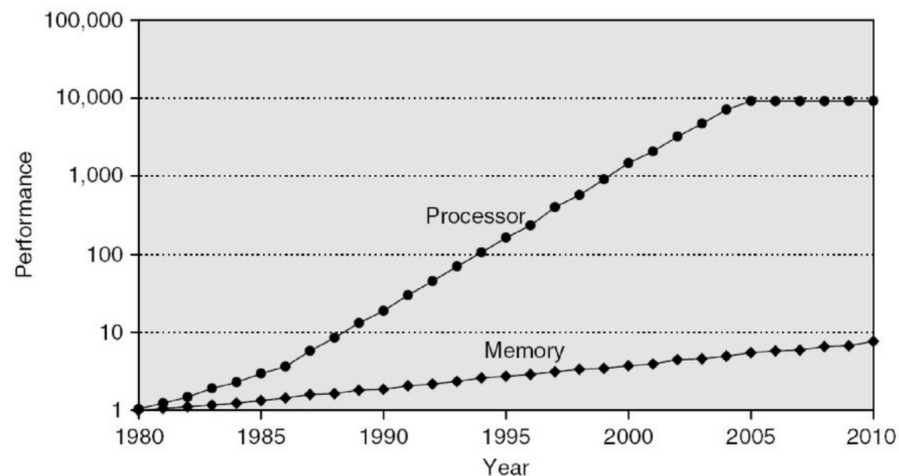
背景

• 内存墙矛盾

- 程序员想要低延迟地访问大量的内存
- 内存的性能增长速度远低于处理器

• 缓存技术的演进

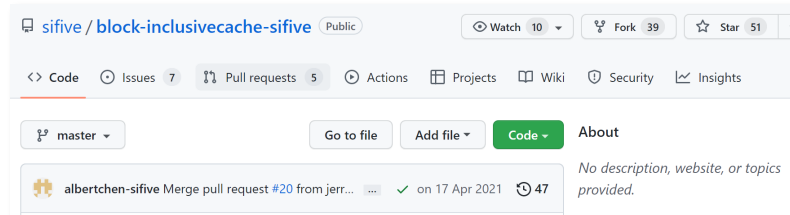
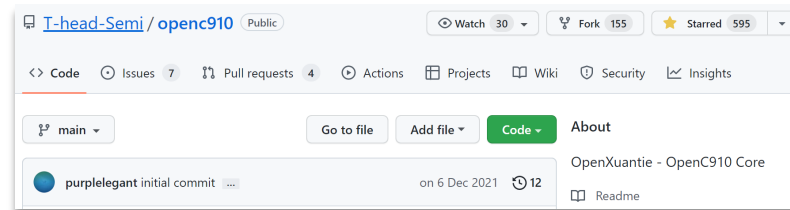
- 延迟：单级缓存 → 多级缓存
- 带宽：Blocking → Non-Blocking
- 替换算法
- 预取算法



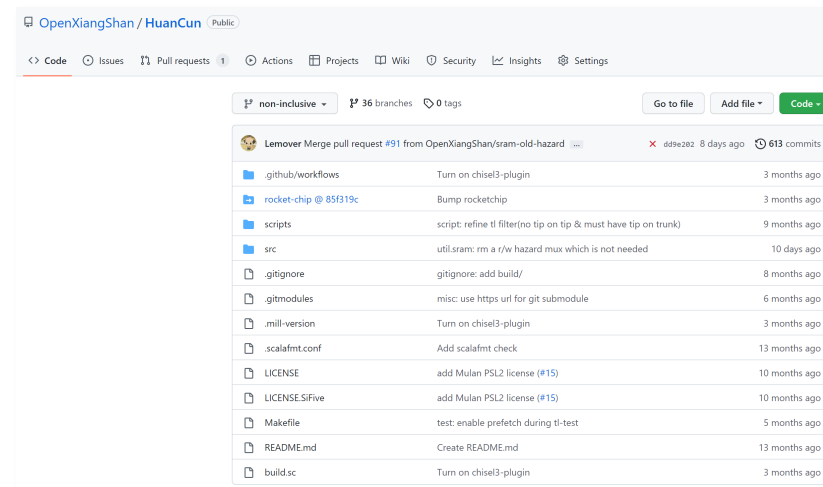
Processor mem accesses/sec vs DRAM accesses/sec

现有开源工作

- 平头哥 openC910 - L2
 - 采用 ACE 总线协议
 - Verilog 编写
- Sifive inclusive-cache - LLC
 - 功能有裁剪：不支持向下的一致性维护
 - 时序较差，不满足流片需求

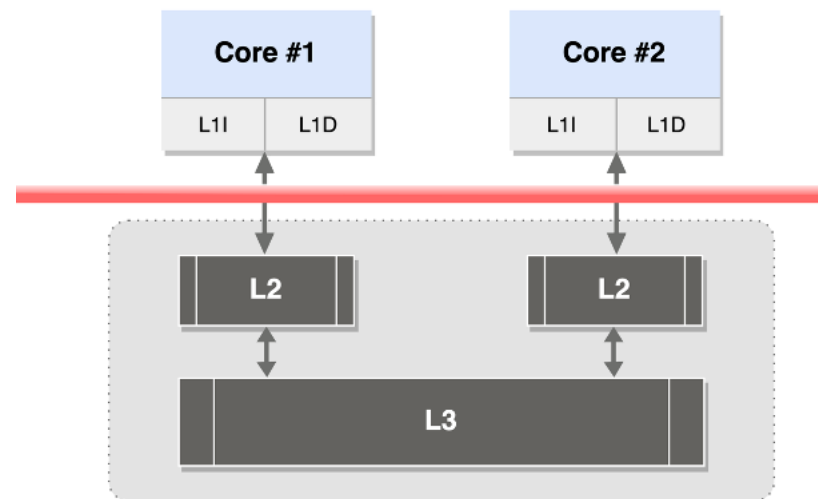


- HuanCun 目标
 - 提升整体缓存**性能**，缓解香山处理器的访存瓶颈
 - 搭建一个开源的高性能缓存研究**真实**平台

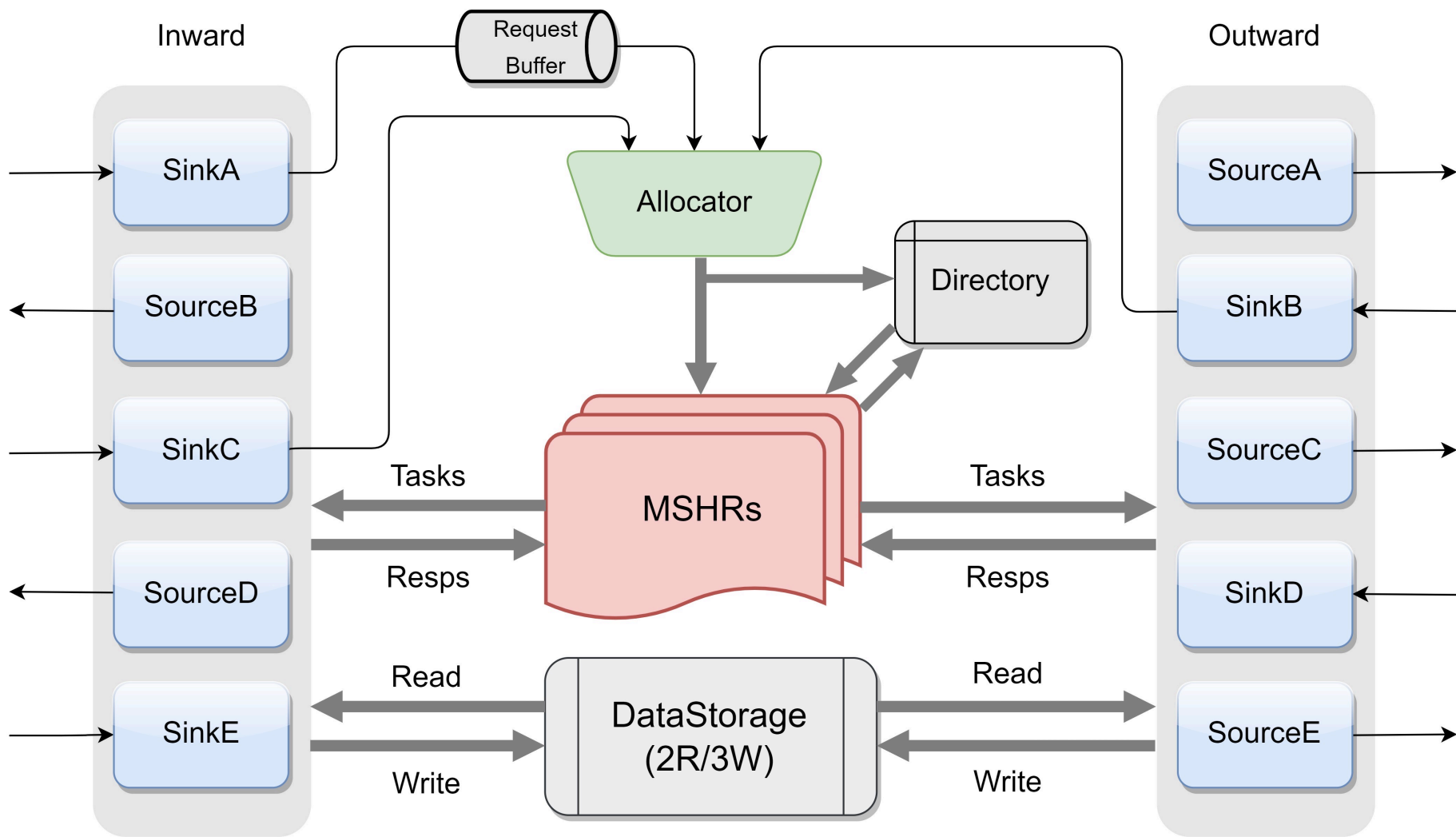


缓存结构设计空间

- 整体框架设计维度
 - 一致性：**Directory Based**
 - 包含关系：**Optional Inclusive/Non-inclusive**
 - 总线：**TileLink**
- 替换算法：**Optional PLRU/Random/SRRIP**
- 预取算法：**Best-Offset Predictor**
- **支持 I/D Cache 一致性维护**
- 细节处理
 - 一致性控制策略
 - 延迟优化
 - 并发度提升

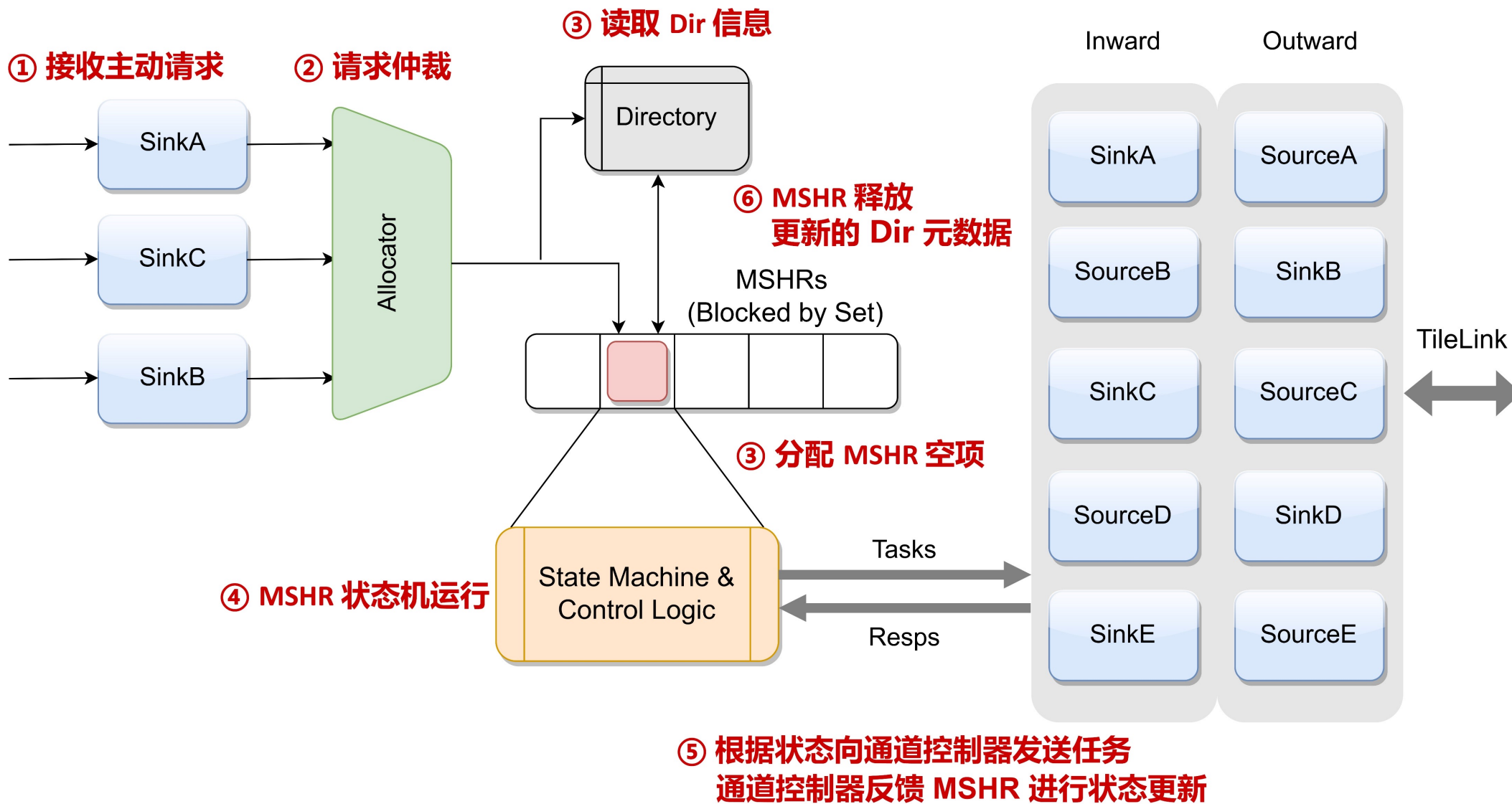


香山 HuanCun 硬件结构





Transaction 处理逻辑



挑战——死锁问题

- ① 请求独占 MSHR & 流水级
- ② 缓存资源容量有限
- ③ MSHR 会等待请求完成才释放资源
- ④ 复杂的多级缓存总线请求依赖易造成循环等待

• 示例场景

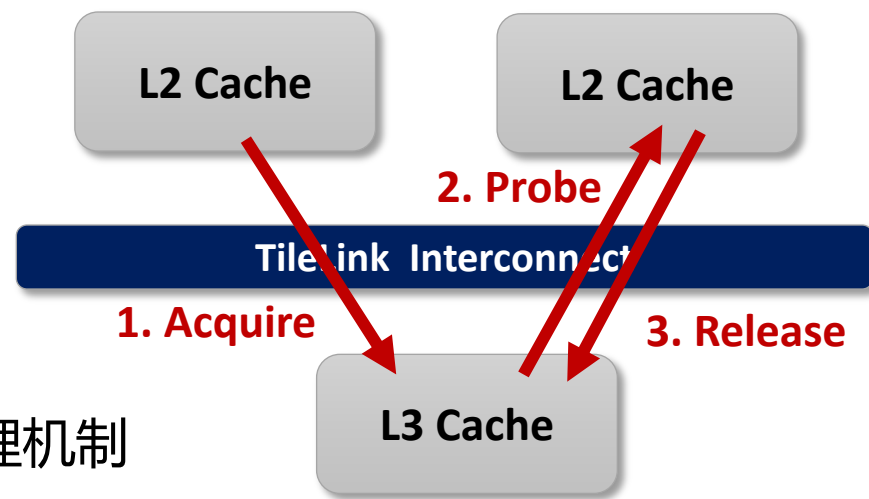
- Release 与 Probe 的资源互等

• 解决方法：

- 更改请求分配策略
- 设计若干非同构 MSHR，引入嵌套处理和跨层处理机制



死锁产生的4个条件



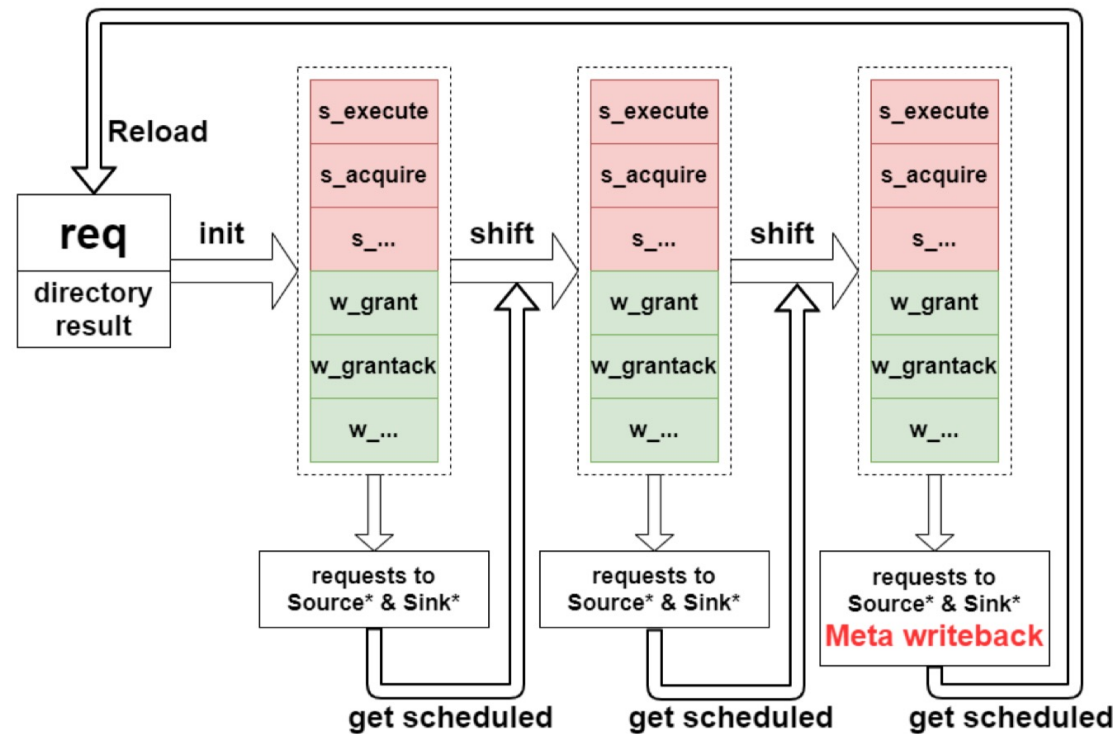
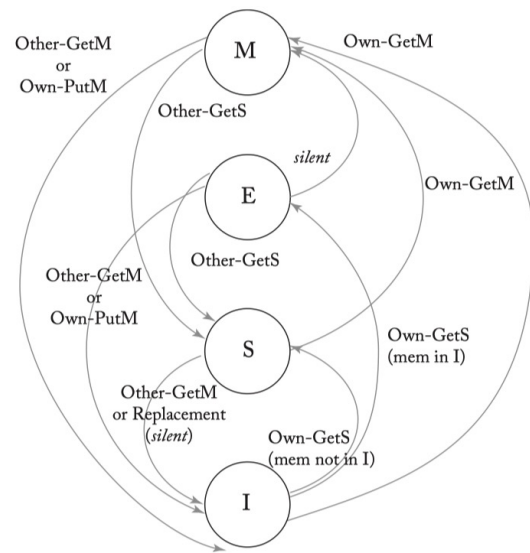
挑战二——状态爆炸问题

• 状态空间：

- 多样的总线请求类型
- 多样的总线响应情况
- Meta 数据的多重组合
 - Client Meta + Self Meta
- 各种嵌套处理情形
- 维护 Non-Inclusive 的复杂策略

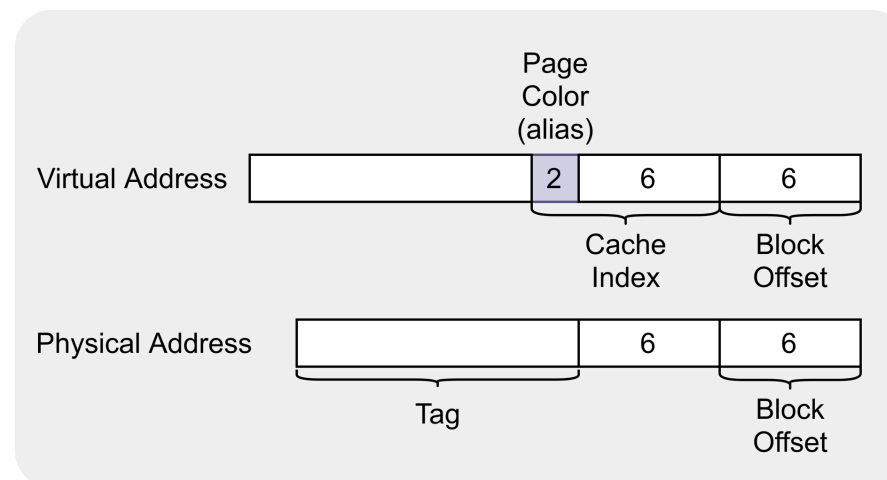
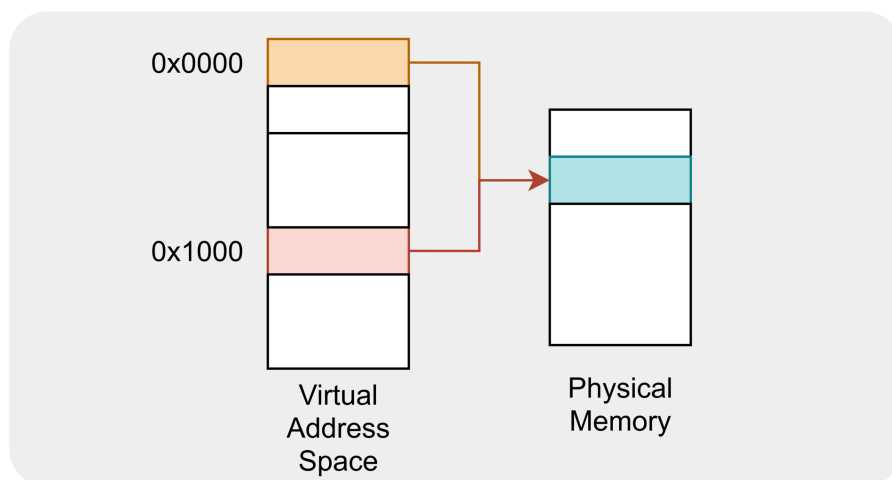
• 解决方法：

- 使用分立的 Flag 寄存器联合表示状态
- 构造合理抽象，尽可能将状态机解耦



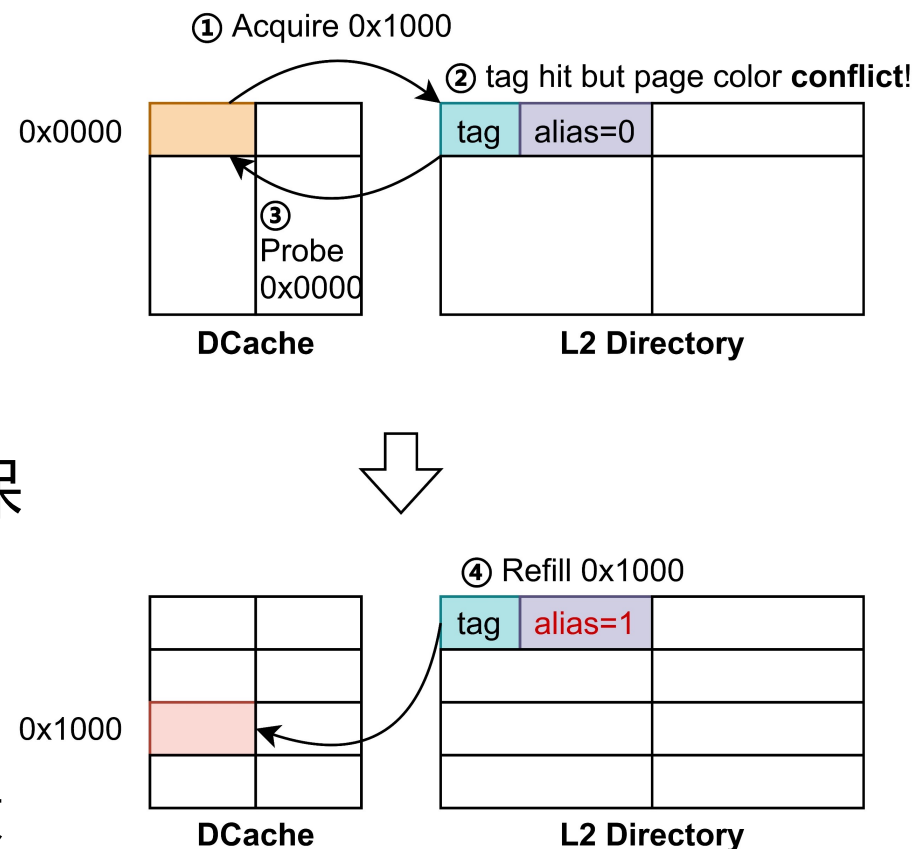
挑战三—— Cache Alias 问题

- 当 VIPT Cache 索引和块偏移所占 bit 数超过了页偏移 (4K, 12bit) 时出现
- 两个虚页映射到同一个物理页时，这两个虚页可能会被索引到 Cache 不同 Set 中，最终结果为：一个物理块在 Cache 中有多份拷贝
- 当处理器写该物理块时，只会索引到一份拷贝，于是导致 Cache 一致性的丧失



Cache Alias 问题的硬件处理

- **解决方案**：利用 L2 Cache，保证一个物理块在上层的一个 VIPT cache 中最多只占有一项
- 具体做法：
 - L1 向下 Acquire 时携带别名位
 - 当 Client Meta Miss 时，L2 将该别名位信息保存在 Client Meta-data 中
 - 当 Client Meta Hit 时，检查总线上的别名位信息与 Meta-data 是否一致
 - 如果不一致，认为发生了 Alias 问题，在 Grant 数据之前先把之前的 L1 数据块 Invalidate 掉



性能优化

• 提升并发度

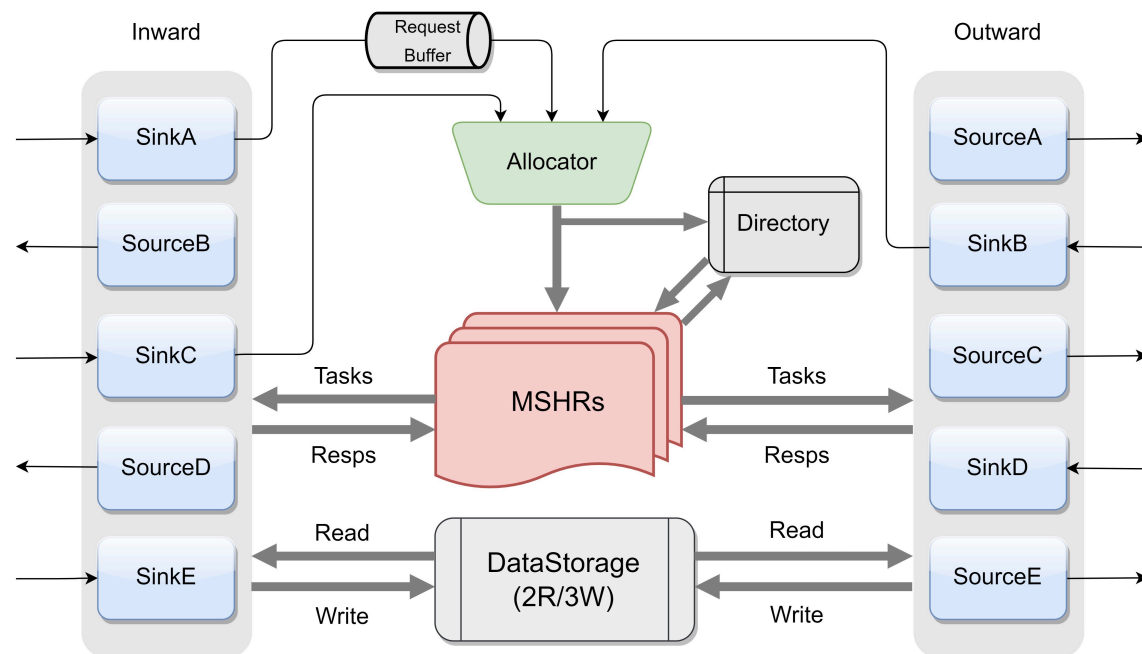
- 增加请求队列 Request Buffer
- 优化 MSHRs 向通道控制器的仲裁

• 降低延迟

- 通道控制器流水线优化
- 增加 Refill Buffer

• 一致性策略优化

- Refill Through , Filter Clean Block 等

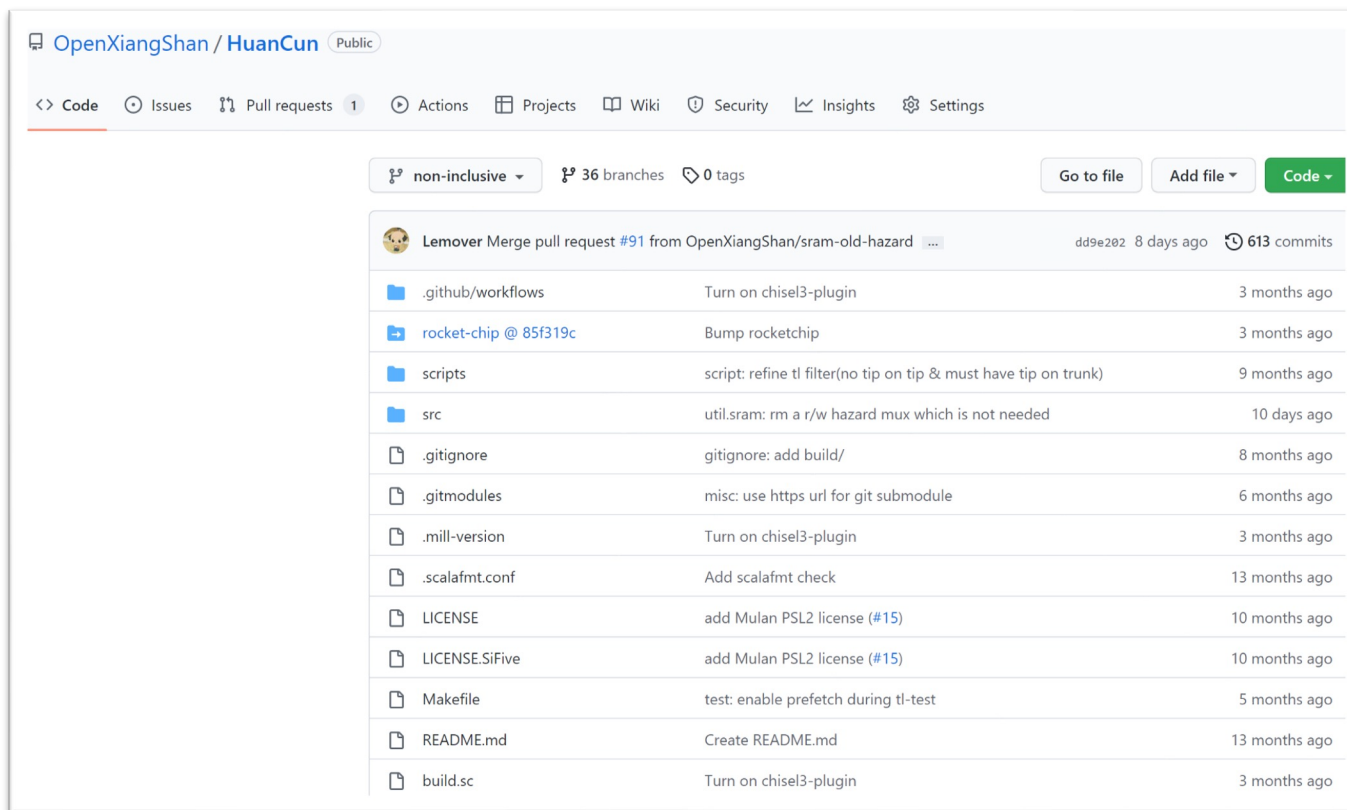


开源情况

- OpenXiangShan 项目组旗下子项目
- 源码、文档均在 Github 等主流代码托管平台上开放



<https://github.com/OpenXiangShan/HuanCun>



谢谢！
请各位批评指教！