

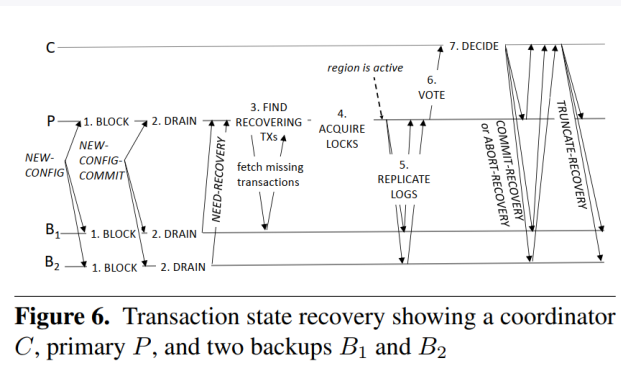
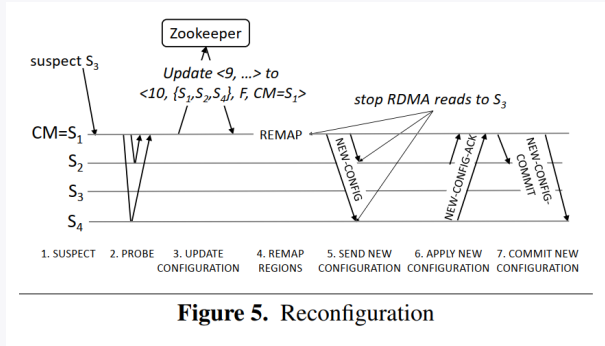
FaRM

优缺点

- 优点
  - 软硬一体化
  - 同时优化CAP
- 缺点
  - backup 为了保证一致性，不提供读，降低吞吐量
  - 集群只是在数据中心中，没有跨区域
  - 应用程序需要处理部分事务一场

可用性

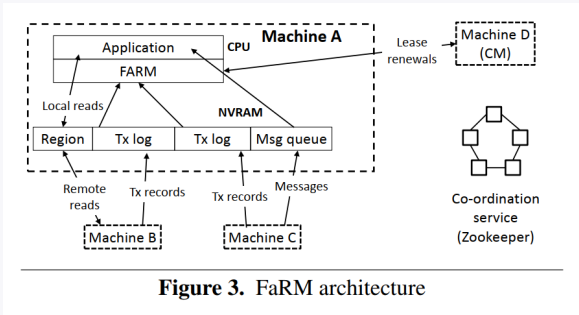
- 错误检测
  - 使用极短租约 (5ms)
- Mechine Down
  - 重新配置
    - 使用 zookeeper 序列号保证一致
- 事务状态恢复
  - RDMA 不被处理
  - 选择要恢复的事务
  - 当接收到NEW-CONFIGCOMMIT消息时，所有机器处理其日志中的所有记录
  - 恢复数据 & 恢复分配器状态



目标

- 在 CAP 三个方面都不妥协
- 高可用的支持分布式事务的高性能内存存储系统

架构



- Cluster 集群
  - 对象 包含 64 bit 版本和数据内容，版本用于并发控制 相当于数据
  - 区域 (region) 多个对象聚合在一起存储，成为 Region 每个 Regions 2G
- 架构组成
  - 集群中的每一台机器
    - 主从复制 1 个 Primary 多个 backup，对象写入读取都在 primary，backup 只做热备份
    - Tx log 事务 log，用于事务恢复
    - Msg Queue 消息队列，用于 RPC
  - Machine
    - Machine 之间可以互相通信 (分布式事务)
  - CM (ConfigManager) 管理 leases, regions 在 machine 间的分布，检测失败，故障恢复等管理工作
  - Zookeeper 协调服务器

软件支持

- 读写
  - 原子读
    - 本地读写 or 远程单边 RDMA
    - RDMA Fetch-and-Add (增加操作)
    - RDMA Compare-and-Swap (比较并交换操作)
  - 原子写
    - 一致性协议
    - 大 Value 支持 RPC
- 事务
  - 乐观并发控制
  - 版本验证，数据是否被改变
  - 验证有效性
  - 无锁读
- 架构设计
  - 使用主从备份取代复制状态机
  - 降低 CPU 瓶颈

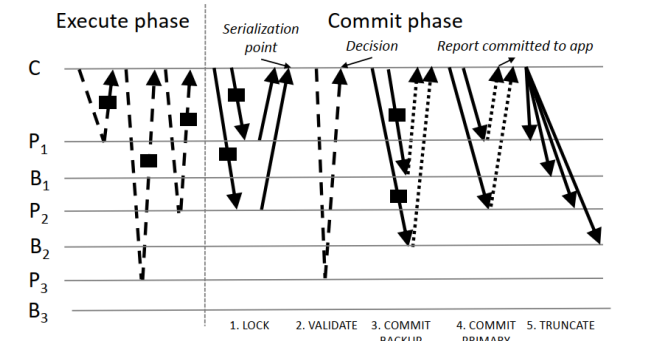


Figure 4. FaRM commit protocol with a coordinator C, primaries on  $P_1, P_2, P_3$ , and backups on  $B_1, B_2, B_3$ .  $P_1$  and  $P_2$  are read and written.  $P_3$  is only read. We use dashed lines for RDMA reads, solid ones for RDMA writes, dotted ones for hardware acks, and rectangles for object data.

一致性事务

- 保证一致性
  - Lock (加锁阶段)
    - 因为使用版本控制，所以只针对写请求加锁
    - 协调器将 lock 记录写入对应写入对象的机器的日志中
    - 其中包括了 primary 所有写入对象的版本和新值
    - 检查锁 锁失败的情况
      - 旧版本：事物更旧
      - Lock：已经被上锁
    - 发送消息报告上锁是否成功
  - Validate (验证阶段)
    - 只针对读请求，只发送给 primary
    - coordinator 从对象的 primaries 读取读验证
    - 如果任何对象被改变，则验证失败同时事务回滚
  - Commit backups
    - 只针对写请求，只发送给 backup，用来保证一致性
    - coordinator 写一个 COMMITBACKUP 记录给非易失性日志，写给每一个 备份服务器，并且等待备份服务器的 ack 回复
  - Commit primaries
    - 只针对写请求，只发送给 primary
    - 在收到所有 backup 的 ack 应答后，协调器给每一个 primary 写 COMMITPERIMARY 记录到日志中
    - Primary 通过在相应位置更新对象来处理这些记录
      - 增加版本
      - 解锁
  - Truncate
    - 内存优化
    - 直到主从服务器被 truncate 之前，都在他们的日志里保存这个记录。协调器会在接受所有的 primary acks 之后惰性的对 primaries 和 backup 进行 truncate logs
  - 原子读写 + 事务乐观并发控制 + 两阶段提交

通信支持

- RDMA
  - RDMA (Remote Direct Memory Access) 是一种网络通信技术，允许计算机系统直接在远程主机的内存中读取和写入数据，而无需涉及中央处理器 (CPU) 的干预
  - 处理不同 mechine 之间的数据不使用 RPC
  - 数据包较小时：RDMA 吞吐大于 RPC
  - 使用单边 RDMA

硬件支持

- 消除存储和网络瓶颈
  - InfiniBand是一种高性能网络技术，常用于构建高性能计算和数据中心环境。FDR表示 InfiniBand的一代规范，提供了每秒传输速率为 56 Gbps的数据传输能力
  - 分布式 UPS (Uninterruptible Power Supply) 是一种分布式的不间断电源系统
  - DRAM + SSD + UPS
  - NVRAM: non-volatile RAM