

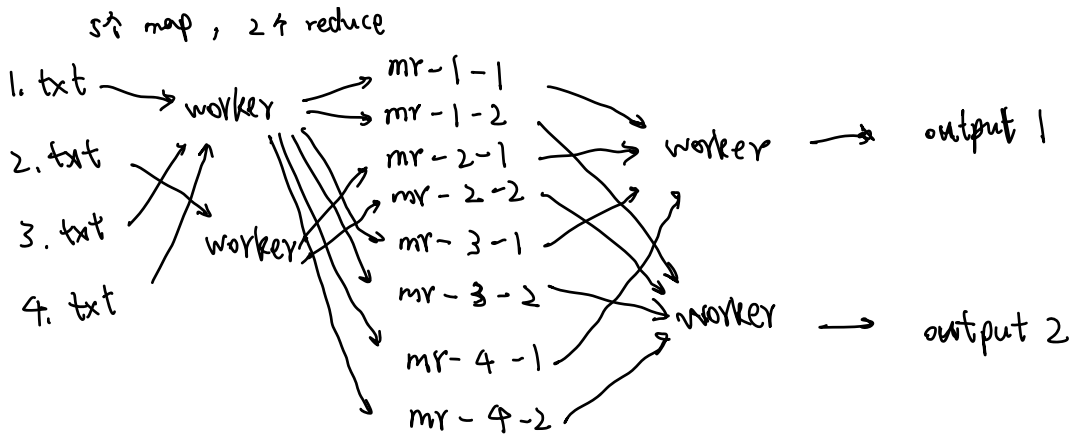
lab 1

map-reduce 流程

- ① split: 把 input files 分割为 m 个子片段
- ② fork: 创建一个 master 和多个 worker, 整个计算过程会执行 m 个 map task 和 r 个 reduce task, master 用于调度 workers
- ③ map: 当 worker 被分配到 map task, 它会去读指定的片段, 调用 `mapf()` 处理原生数据, 生成 $k-v$ 中间数据, 通过 `hash` 函数把中间数据归类到 r 个文件, 写入本地 disk
- ④ reduce: 当 worker 被分配到 reduce 任务, 它通过 RPC 读取中间文件。由于这些中间文件是无序的 (key 与 key 之间未排序), 当一个 reduce task 读完它对应的多个中间文件后, 需要统一对 key 排序。排序后对每个 key 调用 `reducef()` 进行聚合, 然后把结果输出到对应分区的 output file.

错误处理 (lab 中的容错机制比论文简单很多)

- ① 当 worker 失效 (master 发现任务超时), master 把 task 状态改回为 ready, 重新分配给可用 worker 即可



输入文件

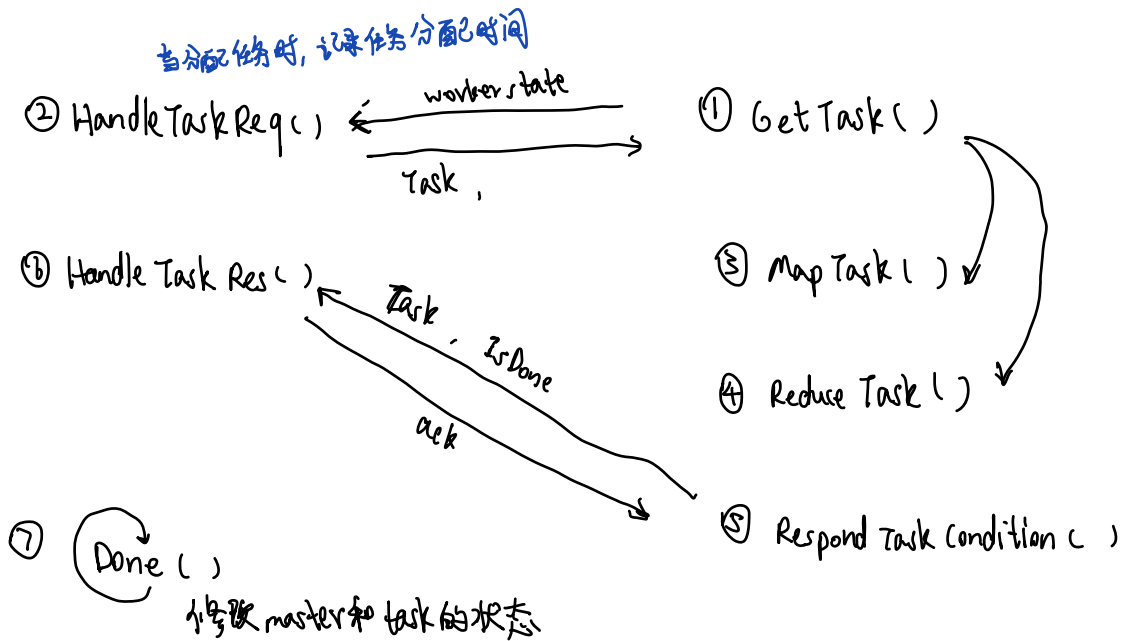
中间文件

输出文件

用 $\text{ihash}(\text{key}) \% N_{\text{Reduce}}$ 划分

master

worker



1. task 超时 → t.rdy

2. map task 全部完成 → m. MAP-FZN

3. reduce task 全部完成 → m. RED-FZN

master 状态: INIT, MAP-FZN, REDUCE-FZN

Task 状态: Task Rdy, Task FZN, Task Running

Lab 2

① check Election Loop () 检查是否需要选举

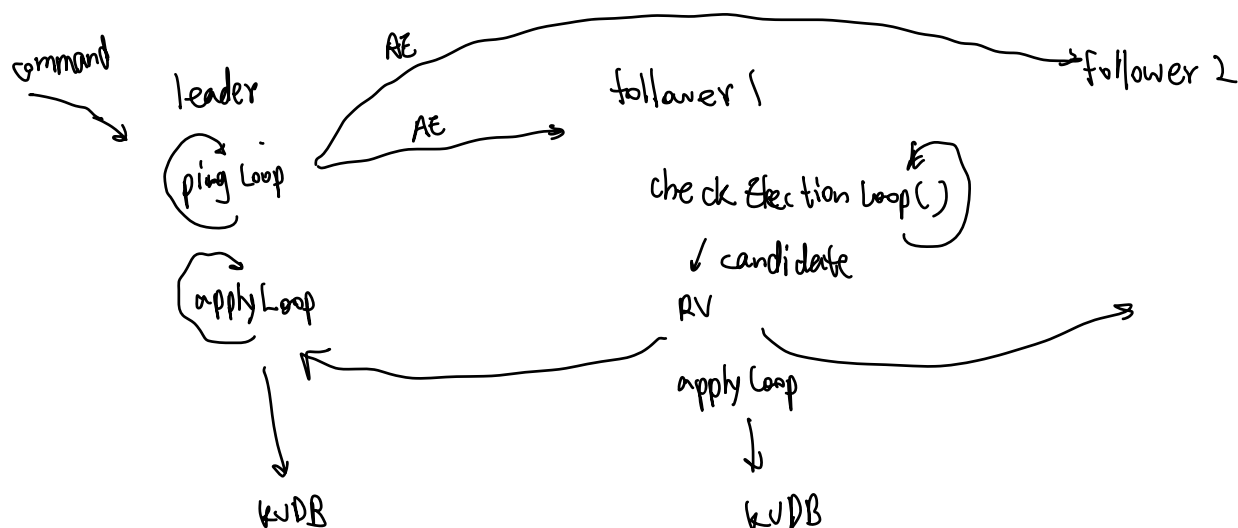
当 server 不为 leader 时，不断循环检查 heart beat 是否超时，一旦超过随机超时时间就调用 send Request Vote 发起选举

② ping Loop () leader 定期发送 AE

一旦 server 成为 leader 就会调用 ping Loop () 直到发现自己不再是 leader

③ apply Loop ()

每个 server 都通过这个循环检查是否有新的 committed 日志，把新日志通过 chan ApplyMsg 应用到状态机 (KV DB)



持久化

对于 Raft 节点，需要持久的状态分别是 **current Term**, **vote For**, **log []**

因此，每当3个变量发生变更就调用一次 **persist ()**

触发时机可以归为3个：

- ① 投票时
- ② AE 追加新日志
- ③ 升 term

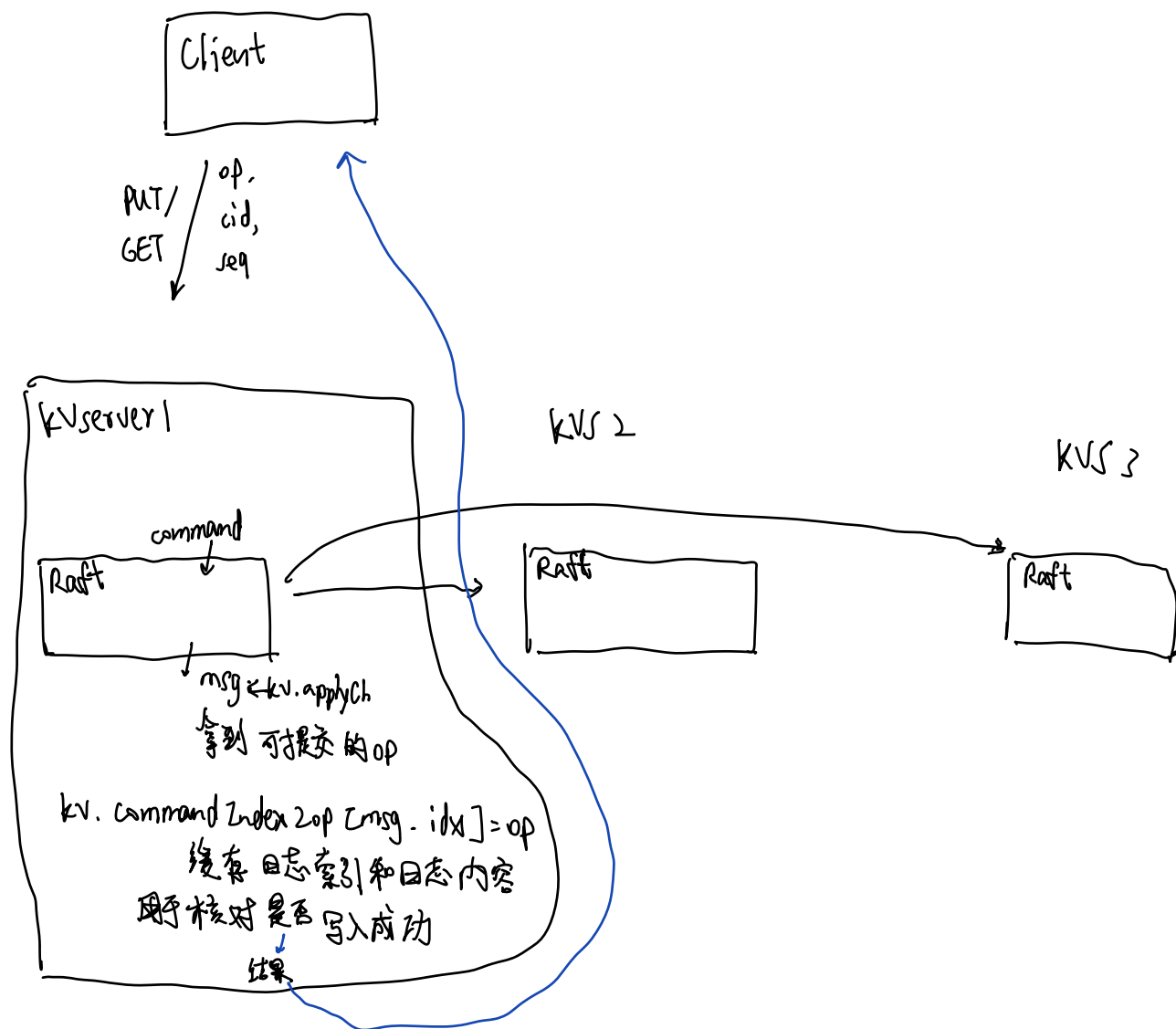
应该重置 timeout 的 3 个时机

① 收到当前 leader 的 AE

② server 开启一次选举

③ 向一个 peer 投票

lab 3



日志压缩与快照 (注意 lab 的持久化都由 Raft 模块完成)

所有 KVServer 在收到来自 Raft 的 command 后, 会检查 raft 的 log 长度。超过阈值则生成快照 (包括 kvDB, cid2seq, lastIncludedIndex, lastIncludedTerm, currentTerm, VoteFor, logC, 其中后5个属于 RaftState)

这样我们就可以清空原来的 logC, 新 log 的真正 index 为 i + lastIncludedIndex

处理流程

① 每个 Raft 节点检查自己的日志大小，当超过阈值，^{kvServer} 需要将 kvDB 和 cid2seq 序列化并发送到 Raft 库。

② Raft 收到上述调用会生成快照 (上述两部分 + Raft State)，落盘。
这两步每个节点独立完成，快照都是写入状态机的 committed 日志，不违反一致性。

③ 在 lab 中，leader 一旦发现 follower 打算接收的下一条日志 nextIndex[i] 比 log[i] 中的日志都小，说明 follower 的快照落后太多，把准备发送的 AE 改为 Install Snapshot

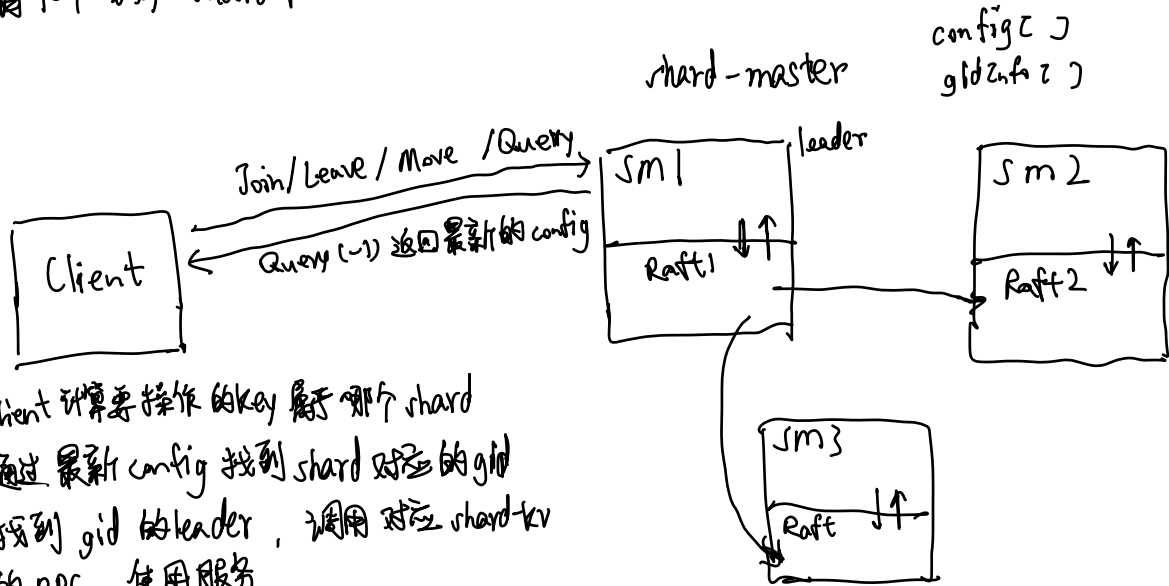
④ follower 用 leader 的 snapshot 覆盖自己的 snapshot 并同步 lastIncludedIndex 和 lastIncludedTerm 使数据与 leader 一致

⑤ 通知 kvServer 拉取新的 kvDB 和 cid2seq，使状态机也与 leader 一致

③. ④. ⑤ 为 leader 与 follower 交互实现

lab 4

有 10 个分片 shard 1 - shard 10

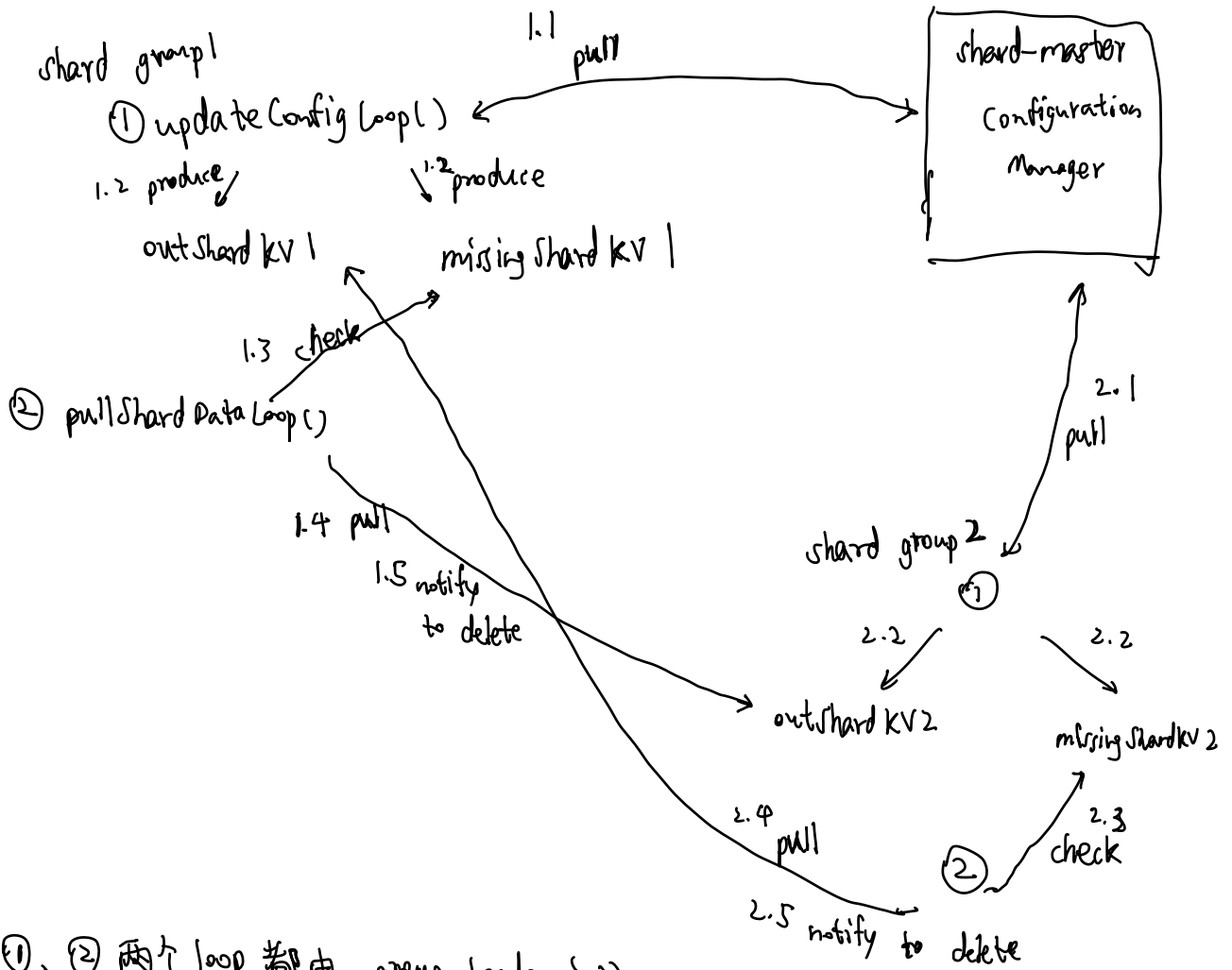


- ① Client 计算要操作的 key 属于哪个 shard
- ② 通过最新 config 找到 shard 对应的 gid
- ③ 找到 gid 的 leader, 调用对应 shard-kv 的 RPC, 使用服务

GET/Append/PUT

shard-kv

- ① 通过 update(ConfigLoop) 不断从 sm 拉取新配置, 根据 shard 变化, 记录要迁出的数据与迁入数据: outShardKV, missingShardKV



①、②两个loop都由 group leader 完成
 通过 Raft 把 config 和 进入数据以日志
 方式同步到 group 中的 follower