



系统通过心跳机制来确定某个节点是否失效。超时时间不能设置的太长, 因为这样系统恢复的时间就很长(因为失效要被修复)。也不能设置得太短, 因为这会引发很多不必要的切换, 给系统带来负担。

DDIA读书笔记

指标: 如何评价一个数据系统

- 可靠性: 即使发生故障, 系统也能正常运行的能力
- 可扩展性: 负载增加时, 系统能有效的保持其性能的能力
- 可维护性: 简单, 易操作, 可演化

数据模型: 如何理解数据世界

- 关系型数据库 (RDB)
 - 支持一对多, 多对一, 多对多 (尽管实现起来比较麻烦)
- 关系模型
 - 强事务
 - 由于应用广泛, 目前有很多查询优化器
 - 采用 schema on write, 比较笨重
 - 适合一对多的场景, 便于一次性将所有数据读出
- 文档模型
 - 采用 schema on read, 比较灵活
 - 在读的时候才知道存储格式
 - write的时候直接存储
- 图模型: 适合多对多场景并能记录路径关系

数据该如何存储

- 事务型数据库(OLTP)
 - 基于日志结构的存储引擎
 - 只支持追加的数据文件。每行为一个KV对
 - 如果更新某个key的值, 旧版本的值不会被覆盖, 而是新增一行记录。查询时某个key的值以最新的为准。
 - 优点: 充分利用了磁盘寻道特性, 比随机写性能更优; 读取时一次寻址, 即可将value从磁盘读取到内存
 - 性能分析: 写: O(1), 每次在文件的末尾追加一行; 读: O(n), 从后往前遍历文件中的每一行
 - 基于SSTable的LSM-Tree
 - TiDB 底层索引结构
 - SSTable (sorted string table) 是一种存储模式
 - LSM-Tree (Log Structured Merge-Tree) 是基于该存储模式在内存中的数据结构。
 - 哈希索引
 - 将key加载到内存, 并建立key与文件中存储位置的映射 (hash map)。
 - 查询时, 通过key找到value在磁盘文件中的偏移量, 定位到磁盘文件指定位置读取数据。
 - 随着时间的推移, 磁盘文件会越来越大。可以通过分段控制文件的大小, 将一个文件拆分成多个文件。每个文件在内存中维护一个hash map。对于一次查询, 会先检查最新的hash map。如果该hash map没有这个key, 则去第二新的hash map去找, 以此类推。
 - 崩溃恢复: 如果数据库重启, 内存中的hash map将丢失。如果重新扫描每个段文件并建立相应的索引将比较耗时。一个可行的方案是, 将每个段的hash map按照存在磁盘上。恢复时只需将快照加载到内存中即可。
 - 优点: 追加、分段、合并都是顺序写, 相比随机写要快的多, 且不易碎片化
 - 并发和崩溃恢复时, 更简单。因此每次写入, 都是先在磁盘上追加, 然后更新内存, 即先持久化。内存丢失的数据, 可以完全依赖磁盘恢复现场。
 - 缺点: 所有的key都必须放入内存: 当key很多时, 内存不够用
 - 使用场景: key不多, 且value更新频繁, 如记录用户粉丝数, 点赞数等
 - 分析型数据库(OLAP)
 - 又称数据仓库
 - 列式存储

数据该长什么样

- 普通文本: CSV, JSON, XML, 二进制
- 特定协议: Thrift, Protocol Buffer
- 特定格式: Avro

数据写到一半, 系统挂了怎么办

- ACID
 - 隔离性: 并发进行多个事务时不应该互相干扰。例如某个事务多次写入, 另一个事务观察到的应该是其全部完成或者一个都没有完成的结果, 而不应该看到中间的部分结果。
 - 读未提交 (read uncommitted): 解决脏读
 - 读已提交 (read committed): InnoDB默认级别
 - 可重复读 (repeatable read): 可重复读一般基于快照来实现。因此又有快照隔离的说法。由于多个正在运行的事务可能会在不同的时间点查看数据库状态, 所以数据库保留了对象多个不同的提交版本。并基于版本来控制事务读取时的可见性。该技术称为多版本并发控制(MVCC)并未解决幻读问题。比如新增数据, 但此时原来的事务读取到的还是原来的数据。通过加行锁或间隙锁, 可以避免幻读问题
 - 串行化: 实现可串行化隔离级别
- 事务
 - 2PL (两阶段锁)