

# 分布式系统中的一致性模型

叶家炜

Jun 13, 2026

在网络分区与延迟不可避免的现实世界里，我们如何在「正确性」与「可用性」之间做出取舍？

分布式系统与单机程序最显著的差异在于，跨节点的事务必须跨越不可靠的网络进行通信。网络分区、消息丢失与节点宕机不再是异常，而是常态。单机环境下事务的 ACID 属性在分布式场景下被重新诠释，核心问题转化为「如何在多个副本之间维持读写操作的可见性与顺序」。一致性模型正是对「读写可见性」进行形式化约束的数学工具，它并不等同于正确性，而是定义了在一定故障模型下系统能够提供的保证边界。本文将从读写语义出发，逐层剖析强一致性与弱一致性家族，并结合 CAP 与 PACELC 定理讨论工程实践中的权衡。

## 1 分布式系统中的读写语义

读写操作在分布式系统中可抽象为两个视角：客户端视角关注「我何时能看到哪些写入」，存储系统视角则关心「如何在多副本间同步状态以满足客户端预期」。在单机上，「最新写入」这一直觉天然成立；而在分布式环境下，跨地域的写操作可能在不同节点以不同顺序到达，导致后续读操作返回陈旧或相互矛盾的结果。举例而言，客户端 A 在北京写入值 X，客户端 B 在上海立即读取，可能得到旧值或默认值，因为网络延迟与分区使得「最新」这一概念失去了全局意义。正是由于这种失效，一致性模型才需要以形式化方式约束「可见性」与「顺序」。

## 2 强一致性家族

### 2.1 线性一致性

线性一致性要求存在一个全局时序，所有操作看起来像在单机上按该时序顺序执行，且必须遵守实时顺序。实时顺序指若操作 A 在全局时钟意义上先于操作 B 完成，则 A 必须排在 B 之前。线性一致性等价于原子寄存器模型，即任何读操作都能立即看到最近一次成功写入的值。典型实现包括多数派写入加租约读、Multi-Paxos 或 Raft 等共识算法。以 Raft 为例，领导者节点在收到写请求后，首先将日志条目复制到多数派节点，只有当多数派确认后才会向客户端返回成功；读操作则通过租约机制确保领导者身份在有效期内，从而避免陈旧读。代价是高延迟与可用性下降：在 CAP 视角下，系统更倾向于牺牲可用性以换取一致性。

### 2.2 顺序一致性

顺序一致性同样要求存在一个全局顺序，且所有进程看到的操作顺序一致，但不再强制遵守实时顺序。这意味着一个读操作可能「滞后」于另一进程已完成的写操作，只要全局顺序不被打破即可。早期的分布式共享内存系统曾采用此模型，它允许一定程度的「滞后」读，从而降低同步开销，但无法满足需要严格实时性的场景。

### 2.3 因果一致性

因果一致性利用 Lamport 的 happens-before 关系，而非全序。向量时钟或因果多播可用于追踪操作间的因果依赖。例如在社交网络时间线中，用户 A 回复用户 B 的评论必须在 B 的原始评论之后可见，但与 A 无关的其他用户操作则可并行执行。实现上，向量时钟为每个节点维护一个计数器数组，消息携带向量时钟，接收方通过比较向量判断因果顺序，从而在不引入全局时钟的前提下保证因果可见性。

## 3 弱一致性家族

### 3.1 最终一致性

最终一致性的非正式定义是：在无新写入且网络分区恢复后，所有副本将收敛到相同状态。Dynamo、Cassandra 与 Riak 等系统均采用此模型。冲突解决可通过 Last-Writer-Wins 时间戳或向量时钟加应用层合并函数完成。由于不提供即时保证，最终一致性允许在分区期间出现短暂不一致，从而换取更高的可用性与更低的延迟。

### 3.2 读己之写一致性

读己之写一致性是一种会话语义，要求在同一会话内，读操作必须能看到自己之前的写操作。实现方式包括粘性会话、单调读或向量时钟。例如在电商购物车场景中，用户更新购物车后立即刷新页面，必须看到最新修改，否则体验将受损。粘性会话通过将同一用户请求路由到同一副本实现，而向量时钟则在分布式环境下记录会话级因果关系。

### 3.3 单调读与单调写一致性

单调读禁止出现「时光倒流」现象，即若一次读返回版本  $V$ ，则后续读不能返回版本小于  $V$  的值。单调写则要求同一进程的写操作按其发起顺序被所有副本接受。在缺乏全局时钟的情况下，工程上可通过本地逻辑时钟或租约机制近似实现单调性，但仍需权衡同步成本与一致性强度。

### 3.4 带界限的陈旧读

带界限的陈旧读将「多久之前」的读限定在可接受范围内。PNUTS 系统中的时间线一致性即属此类：它允许读操作返回最多  $t$  秒前的状态，通过在副本间异步传播更新来降低跨地域延迟，同时通过界限  $t$  控制陈旧程度。

## 4 一致性模型与 CAP、PACELC 定理

CAP 定理指出，在网络分区存在的情况下，一致性与可用性不可兼得。PACELC 定理进一步将权衡扩展到正常运行场景：若系统在无分区时，则需在延迟与一致性之间做出选择。映射到实际系统，HBase、etcd 与 ZooKeeper 选择 CP 模式，在分区时牺牲可用性；Cassandra 与 Dynamo 选择 AP 模式，优先保证可用性；MongoDB 副本集则在正常情况下提供因果一致性，分区时退化为最终一致性。工程启示在于，不存在 universally 最好的模型，只有最适合具体负载与 SLA 的模型。

## 5 工业实践中的权衡

银行转账场景要求线性一致性，以确保「先扣款后转账」的原子性。电商库存扣减可采用线性一致性或顺序一致性配合幂等重试，既保证不超卖，又允许短暂重试。社交动态推送适合因果一致性加最终一致性，在保证评论顺序的前提下容忍短暂延迟。实时推荐系统则依赖读己之写加最终一致性，确保用户看到自己最新行为产生的推荐。指标采集场景下，最终一致性配合聚合窗口即可满足监控需求。

## 6 实现一致性模型的关键技术

共识算法如 Paxos、Raft 与 Viewstamped Replication 为强一致性提供基础。Quorum 系统通过  $W + R > N$  的数学约束确保读写操作交叠，从而实现线性一致性。租约与领导者机制在降低读延迟的同时维护一致性。冲突检测与解决可借助 CRDT、OT 或向量时钟。混合一致性则允许同一系统内部不同数据路径采用不同模型，例如 TiDB 中 TiKV 提供强一致性，而 TiFlash 提供最终一致性的列存副本。

## 7 形式化验证与测试

Jepsen 测试框架通过在现实故障下注入网络分区与节点宕机，验证系统是否满足线性一致性。TLA+ 与 PlusCal 可对算法进行模型检验，提前发现死锁或活锁。混沌工程实践揭示了 ElasticSearch、MongoDB 与 etcd 等系统在历史版本中暴露的一致性 Bug，证明形式化方法与混沌测试的互补价值。

## 8 未来趋势

弱一致性正成为默认选项，强一致性按需开启，例如 CockroachDB 支持 follower read 与严格可串行化。跨地域多活面临「全球一致性」难题，新硬件如 RDMA 与持久内存有望降低一致性开销。形式化方法在工业界的落地也将加速一致性模型的可靠实现。

## 9 结论

一致性模型并非非黑即白，而是一条连续光谱。理解「可见性」「顺序」与「实时性」三要素，可在 CAP 与 PACELC 约束下做出最优权衡。没有银弹，理解取舍即是分布式系统设计的艺术。