

# 分布式系统一致性模型

杨子凡

Jun 30, 2026

分布式系统把数据复制到多个节点以求高可用和高吞吐，但节点宕机、网络分区与消息乱序却让「同一份数据在同一时刻看到相同值」变成奢望。某支付系统在网络抖动时出现两次扣款，购物车在跨地域同步时丢失商品——这些故障本质都是「一致性」被牺牲的代价。为什么强一致在分布式环境下如此昂贵？本文先用 CAP 定理剖析根本矛盾，再沿着「强到弱」的光谱逐一拆解模型，最后结合真实系统给出工程落地的权衡思路。

## 1 分布式一致性问题的根源

在单机上，读写寄存器天然满足顺序；一旦跨机，消息延迟、丢失或重排便打破了全局时钟。节点 A 在时刻  $t_1$  写入  $x=1$ ，节点 B 在  $t_2$  收到更新之前读到  $x=0$ ，这种「先写后读」被违背的现象，正是分布式一致性要解决的核心。CAP 定理指出，在网络可能分区的前提下，一致性与可用性只能二选一。证明思路可简单概括为：若分区发生且两侧节点仍需对外服务，则至少一侧无法获得最新写入，违背一致性；若强制要求一致性，则至少一侧必须拒绝服务，违背可用性。代价不仅体现在延迟与吞吐，更在于协议复杂度与运维心智负担。

## 2 一致性模型光谱

### 2.1 强一致性：Linearizability

Linearizability 的直观含义是：系统对外呈现一个原子对象，所有操作都可以被排成一个全序，且这个全序必须尊重真实时间。形式化地说，对于任意两个操作  $o_1$  与  $o_2$ ，若  $o_1$  的响应时刻早于  $o_2$  的调用时刻，则在全序中  $o_1$  一定排在  $o_2$  之前。实现该性质最经典的算法是 Multi-Paxos 或 Raft。以 Raft 为例，日志条目必须先被多数派确认，再由 Leader 应用到状态机。etcd、Consul、CockroachDB 均采用此协议，保证任意时刻最多只有一个 Leader，且读请求可线性化地从 Leader 或加租约的 Follower 返回。

### 2.2 强一致性的近亲：Sequential Consistency

Sequential Consistency 去掉了「真实时间」约束，只要求存在一个全序，使每个线程内部的顺序得以保留。换言之，跨线程的操作只要最终顺序一致即可，不必严格遵守物理钟。典型实现仍可复用 Paxos，但允许读请求在任意副本执行，只要保证线程内因果即可。两者差异在 Jepsen 测试中极易暴露：Linearizability 禁止「新写入被旧读取覆盖」，Sequential Consistency 则可能允许。

## 2.3 因果一致性

因果一致性仅维护「因果序」。若操作 a 因果先决於 b，则所有副本都必须以 a 先於 b 的顺序观察。实现方式通常借助向量时钟：每个副本维护一个长度等于节点数的数组，写操作递增自身计数，读操作携带向量一并返回；当收到消息时，只有其向量在逐分量意义上小于等于本地向量，才允许交付。COPS 系统利用该机制在数据中心间提供低延迟的因果广播，同时允许跨数据中心异步复制。

## 2.4 最终一致性

最终一致性仅承诺「若不再有新写入，所有副本将收敛到相同值」。Dynamo、Cassandra、Riak 均采用可配置的 NWR 策略：N 为副本总数，W 为写成功所需确认数，R 为读成功所需响应数。满足  $W+R>N$  即可提供「Quorum」级的一致性；若放宽到  $W=1$ 、 $R=1$ ，则吞吐最高但可能读到旧值。冲突消解可使用 Last-Writer-Win，也可使用向量时钟合并或 CRDT。CRDT 的核心思想是让并行操作满足交换律、结合律与幂等律，从而无需协调即可安全合并。

## 2.5 最终一致性的细化子模型

在最终一致性框架下，还可进一步约束读写行为。Read-your-writes 要求一旦客户端写入成功，后续读必须能看到该值；Monotonic Reads 要求若已读到版本 v，则不会再读到版本小于 v 的值；Consistent Prefix 保证读到的序列是全局更新前缀的子序列；Bounded Staleness 则用时间或版本号量化「过时」上限。这些子模型可在不牺牲全部可用性的前提下，显著降低业务层面的异常概率。

# 3 工程落地：协议与权衡

## 3.1 单对象与多对象事务

线性一致的读写寄存器仅能保证单键原子性。若需跨键事务，传统做法是两阶段提交（2PC）：第一阶段询问所有参与者是否可提交，第二阶段统一提交或回滚。Percolator 在 Bigtable 之上引入主次提交记录，把事务状态持久化到独立行，从而把锁范围缩小到单行。Spanner 更进一步，借助 TrueTime API 把时间戳误差控制在 7 ms 以内，再用 2PC+Multi-Paxos 实现全球线性一致的分布式事务。

## 3.2 共识与复制的分野

共识算法如 Raft、Paxos 追求强一致，代价是写入必须等多数派确认。另一条路径是「Gossip+Anti-Entropy」：节点周期性交换 Merkle 树或版本向量，异步修复差异。Cassandra 便采用此策略，在牺牲实时一致性的同时换取极高的写入吞吐。

## 3.3 冲突检测与物理时钟

向量时钟可精确刻画因果，但无法处理物理时间需求。Spanner 的 TrueTime 通过原子钟与 GPS 提供区间时间戳，事务提交需等待区间结束以消除时钟偏差。CockroachDB 则采用混合逻辑时钟（HLC），在事件驱动下

融合物理时钟与逻辑计数，兼顾精度与可用性。

### 3.4 量化一致性-延迟权衡

Probabilistically Bounded Staleness (PBS) 用概率模型预测「读到旧值的概率」随延迟的变化曲线。监控指标通常包括 Replication Lag (主从延迟毫秒数) 和 Clock Skew (节点间时钟偏差)，两者共同决定了实际系统偏离 Linearizability 的程度。

## 4 真实系统案例

Google Spanner 把 TrueTime 与 2PC 结合，实现全球范围的线性一致读写，典型延迟在几十毫秒量级，适合金融级交易。Amazon Dynamo 选择 Sloppy Quorum 与向量时钟，允许临时写入仅被  $W$  个节点确认，读时最多尝试  $R$  个节点，极致追求可用性，适合购物车、会话存储。TiDB 把 Raft 复制单元下放到每个 Region，同时在 TiKV 之上构建 Percolator 风格的事务层，既保证单 Region 强一致，又支持跨 Region 分布式事务。Kafka 的 In-Sync Replica (ISR) 机制仅保证「至少一次」投递：当 Leader 故障且 ISR 为空时，可能出现数据丢失，因此不提供 Linearizability。

## 5 选择一致性模型的决策路径

若业务要求「转账必须原子且立即可见」，则必须选择 Linearizability，对应 Raft/Paxos 方案；若仅需「好友动态按因果顺序出现」，则向量时钟或因果广播即可；若允许「购物车偶尔少商品，刷新后恢复」，则最终一致性加上可接受的冲突消解策略最为经济。成本检查清单包括：P99 延迟 SLA、数据规模、团队对运维复杂度的承受力，以及是否需要多机房容灾。

## 6 未来趋势与开放问题

Jepsen 与 Porcupine 等工具已能自动注入网络分区并验证 Linearizability 违规；RDMA 与持久内存的普及正在把一次网络 RTT 从百微秒降至数微秒，这将显著改变现有协议的延迟瓶颈；跨云、异地多活场景对「有界延迟一致性」提出新需求，即允许在毫秒级时间窗口内出现不一致，但必须可量化、可证明。形式化方法在工业界普及仍需降低心智负担与工具链门槛。

一致性模型没有银弹，强一致带来正确性，却付出延迟与复杂度；弱一致换取吞吐，却把异常处理交给业务。建议读者用 Maelstrom 教学框架或 Jepsen 在 minikube 上运行 etcd，观察 Linearizability 被违反时的日志，从而建立直观感受。理解权衡的本质，才能在业务需求与技术代价之间做出最适合的选择。

## 7 附录与延伸阅读

必读论文包括 Lamport 的《Time, Clocks, and the Ordering of Events in a Distributed System》、Gilbert 与 Lynch 的《Brewer's Conjecture and the Feasibility of Consistent, Available, Partition-Tolerant Web Services》，以及 Viotti 与 Vukolić 的综述《Consistency in Non-Transactional Distributed Storage Systems》。书籍推荐《Designing Data-Intensive Applications》第 5、9 章。开源实现可参考 etcd、TiKV、FoundationDB；教学工具可尝试 Maelstrom。相关博客与播客包括 The Morning

---

Paper 与 Distributed Systems Podcast。