

Redis - string系列 (54~69)

Redis 核心技术与业务驱动的系统设计深度解析

一、Redis 字符串类型的底层编码与细节 (深度理解)

Redis 的字符串类型 (string) 在表面上看似简单，但其高性能与高内存效率的实现，关键在于底层编码的智能选择。理解编码机制，才能真正掌握 Redis “快” 与 “省” 的本质。

1. 三种编码格式：原理、适用场景与设计哲学

(1) int 编码：为纯数值与极致性能而生

- 触发条件

字符串为纯整数、长度 ≤ 39 字节、可安全存入 64 位有符号整数范围。

- 底层实现

不存储字符串形式，直接将数值作为 `long long` 存入 Redis 对象。

- 优势

- a. 省内存：无需 SDS (Simple Dynamic String) 结构，无字符数组开销。
- b. 运算极快：`INCR`、`DECR` 等操作直接使用 CPU 整数运算，无需解析字符串。
- c. 确定性：整数运算无精度损失，结果绝对可靠。

- 关键意义

`int` 编码是 Redis 高性能计数器的基石，适用于点赞数、访问量、库存等高频计数场景。

(2) embstr 编码：为短字符串极致优化

- 触发条件

非纯整数字符串，长度 ≤ 44 字节。

- 核心思想

“短字符串非常多，让它们一次内存分配搞定。”

- **内存布局**

Redis 对象与 SDS 连续分配，仅需一次 `malloc`，释放也仅需一次 `free`。

- **性能优势**

- 减少内存碎片
- 降低 `malloc/free` 系统调用次数
- 提升 CPU Cache 命中率

- **典型场景**

用户名、状态标志、小型 JSON、配置项等短字符串存储。

(3) raw 编码：为可扩展字符串设计

- **触发条件**

字符串长度 > 44 字节。

- **底层实现**

使用 SDS (Simple Dynamic String)，支持动态扩容、预分配、O(1) 长度获取。

- **设计意义**

适用于频繁 `APPEND`、`SETRANGE` 或存储大文本、大 JSON 等场景。

- **本质区别**

`embstr` 追求最小化内存分配次数；`raw` 追求字符串的可扩展性与大容量存储。

2. 编码选择与限制：实际开发中的“坑”与应对

(1) 纯数字 ≠ 一定使用 int 编码

Redis 对 `int` 编码的判断极为严格：

- `"123"` → `int`
- `"000123"` → `embstr` (前导零被视为字符串)
- `"123.45"` → `embstr` (小数点被视为非整数)
- 超长整数字符串 → `raw`

后果：若未触发 `int` 编码，则无法使用 `INCR` / `DECR`，否则报错。

(2) 浮点数的处理与建议

`INCRBYFLOAT` 虽支持浮点运算，但底层仍以字符串存储，存在：

- 解析与格式化开销
- 浮点数精度误差（如 `0.1 + 0.2 ≠ 0.3`）
- 不适合高频计数或金融场景

实践建议：

- 高频统计 → 使用整数
- 金额存储 → 放大 100/1000 倍后以整数存储，显示时再还原

(3) 大厂为何不完全依赖 Redis 编码？

并非 Redis 不够强大，而是业务模型通常更复杂，需要：

- 版本号、状态位、校验字段等元信息
 - 统一的数据访问层，以屏蔽底层编码差异
 - 将 Redis 视为“高性能工具”，而非“业务模型载体”
-

二、Redis 作为缓存的核心场景与问题（真实世界视角）

1. 缓存穿透、缓存更新与缓存雪崩

(1) 缓存穿透：请求不存在的数据

- 问题本质

数据在 Redis 与 DB 中均不存在，但请求持续冲击系统。

- 解决方案

- a. 缓存空值（设置较短 TTL）
- b. 布隆过滤器（Bloom Filter）预判存在性

- c. 请求参数合法性校验

(2) 缓存更新策略：一致性与复杂度权衡

- 主动更新

一致性高，但实现复杂，高并发下易产生竞态条件。

- 过期淘汰

实现简单，接受短暂不一致，是 Redis 的典型使用方式。

- Redis 设计哲学

“追求极速，但不保证绝对实时一致。”

(3) 缓存雪崩：大量 Key 同时失效

- 诱因

相同 TTL 设置、Redis 重启、热点数据集中过期。

- 解决方案

- a. 随机化过期时间（基础策略）
- b. 多级缓存（如 Redis + 本地缓存）
- c. 热点数据永不过期 + 异步刷新

2. 视频网站缓存案例：典型最终一致性设计

- 为什么不缓存视频文件？

Redis 内存成本高，视频为冷数据，更适合 CDN 或对象存储。

- 缓存什么？

播放量、点赞数、评论数等高频更新的小数据。

- 写入策略

- a. 播放 → `Redis.INCR`
- b. 定时/异步 → 同步至 MySQL

这是典型的“最终一致性”架构，兼顾性能与数据持久化。

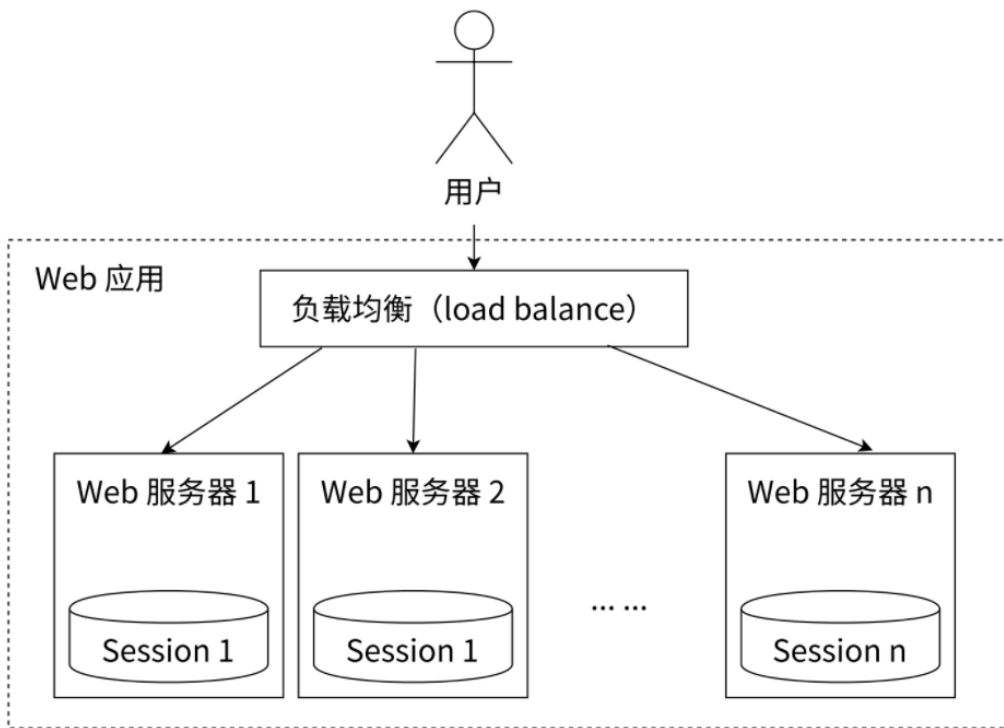
三、Session 共享问题与 Redis（分布式系统必修）

1. Session 问题的本质

负载均衡环境下，Session 绑定到单机导致状态丢失。

2. Redis 作为 Session 存储的优势

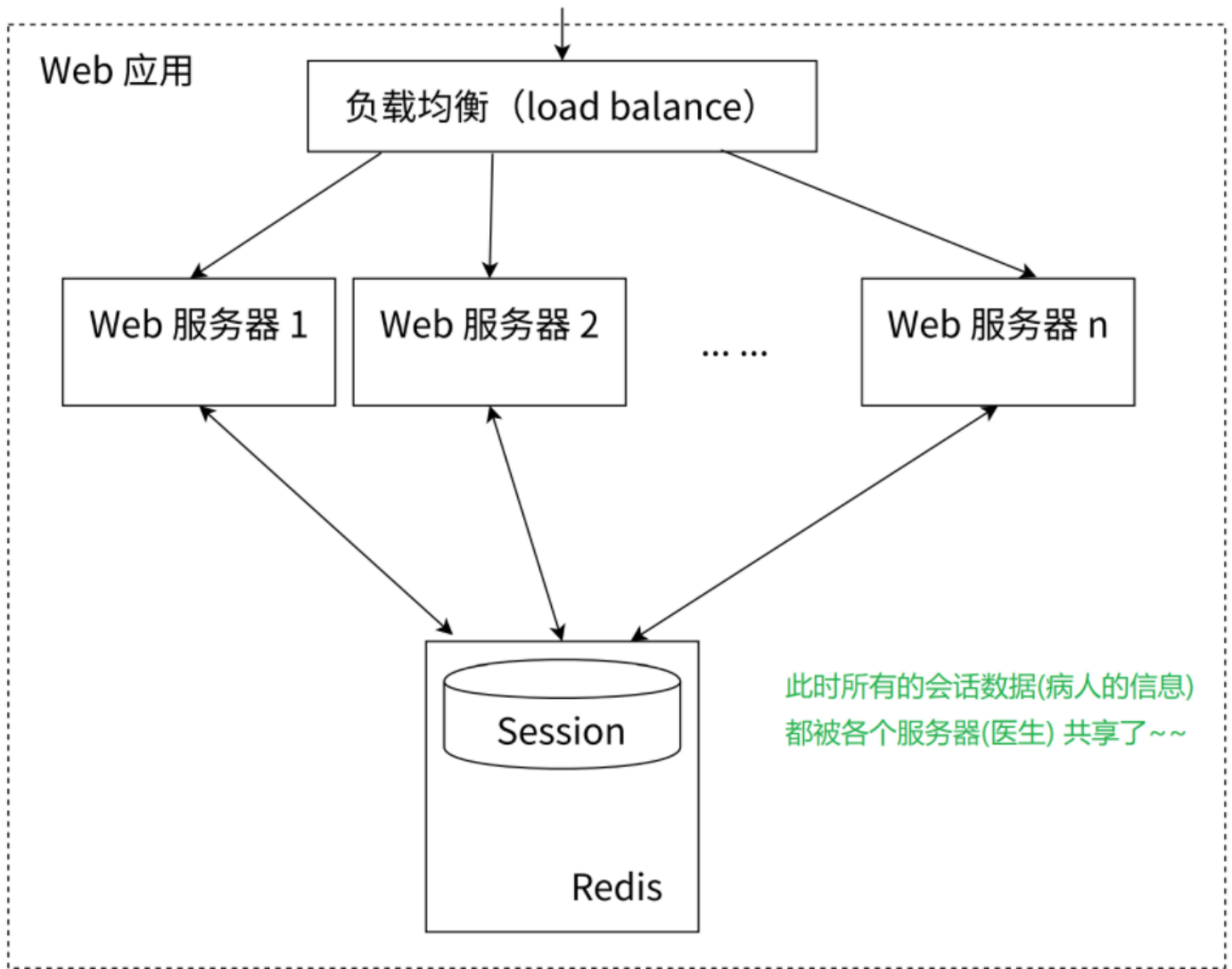
- 高性能读写
- 原生支持 TTL（自动过期）
- 支持多客户端并发访问
- 充当“状态中心”，解耦服务与状态存储



如果每个应用服务器, 维护自己的会话数据, 此时彼此之间不共享, 用户请求访问到不同的服务器上, 就可能会出现一些不能正确处理的情况了~~

3. 类比：“医院病历系统”

- 病历 (Session) 不能存放在单个医生 (服务节点) 手中
- Redis 扮演“病历中心”角色, 统一管理、高效分发
- 服务节点无状态化, 便于水平扩展



四、Redis 字符串命令的深度解析与设计哲学

1. SET 系列命令：原子性即生命线

(1) SETNX：分布式锁的原子原语

- 语义：“如果不存在则设置”，类比 CPU 的 CAS 操作。
- 为什么不用 `EXISTS` + `SET` ?

两条命令非原子，并发下锁失效。

- 现代最佳实践：

```
SET lock_key value NX EX 10
```

兼具“不存在才设置”与“自动过期”，原子完成。

(2) SETEX: 原子设值与过期

- **经典陷阱:** SET + EXPIRE 非原子，中间崩溃导致永不过期。
- **SETEX 价值:** 原子保证“设值”与“设过期”同时成功或失败。
- **适用场景:** 验证码、临时 Token、会话状态等。

2. 数值操作命令: INCR / DECR 系列

(1) INCR 为何高效?

- **原子性来源:** Redis 单线程执行命令，无竞态条件。
- **对值的要求:** 必须为整数或不存在（隐式初始化为 0）。
- **设计理念:** 让“计数器”成为一等公民，简化开发。

(2) INCRBYFLOAT 的隐藏成本

- 底层仍为字符串存储，存在浮点精度误差。
- **金融场景严禁使用**，应使用整数放大法。

3. 字符串操作命令与 SDS 底层

(1) SDS: Redis 字符串的真实身份

- 记录长度 (O(1) 获取)
- 二进制安全 (可存任意数据)
- 支持动态扩容

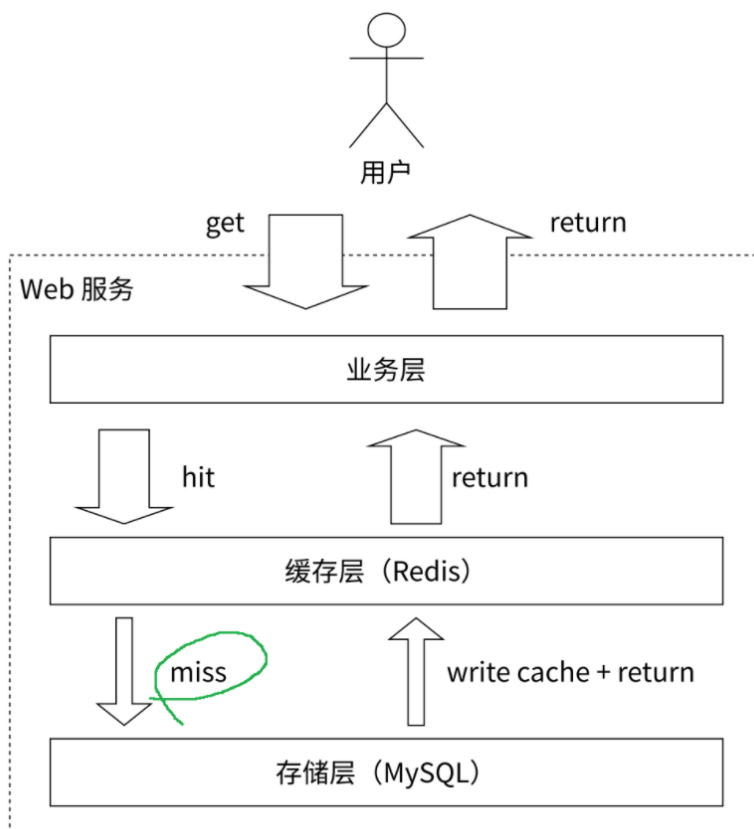
(2) APPEND 与 SETRANGE 的注意事项

- APPEND 触发扩容时有一定成本，不建议用于无限日志。
- SETRANGE 的 offset 单位为字节，对 UTF-8 中文易破坏编码，仅适用于二进制或固定结构数据。

(3) STRLEN 返回字节长度

- 对中文字符串返回字节数（如“你好”返回6），非字符数。

五、string类型的应用场景



整体的思路: 应用服务器访问数据的时候，先查询Redis.如果Redis上数据存在了，就直接从Redis取数据交给应用服务器.不继续访问数据库了.

如果Redis上数据不存在，再读取MySQL.把读到的结果，返回给应用服务器同时，把这个数据也写入到Redis中

Redis这样的缓存，经常用来存储“热点”数据（“热点数据”即为高频被使用的数据，定义方式需结合业务场景）

上述策略有一个明显的问题：随着时间的推移，肯定是会有越来越多的key在Redis上访问不到，从而从MySQL中读取到Redis中，此时Redis中的数据就会越来越多？？

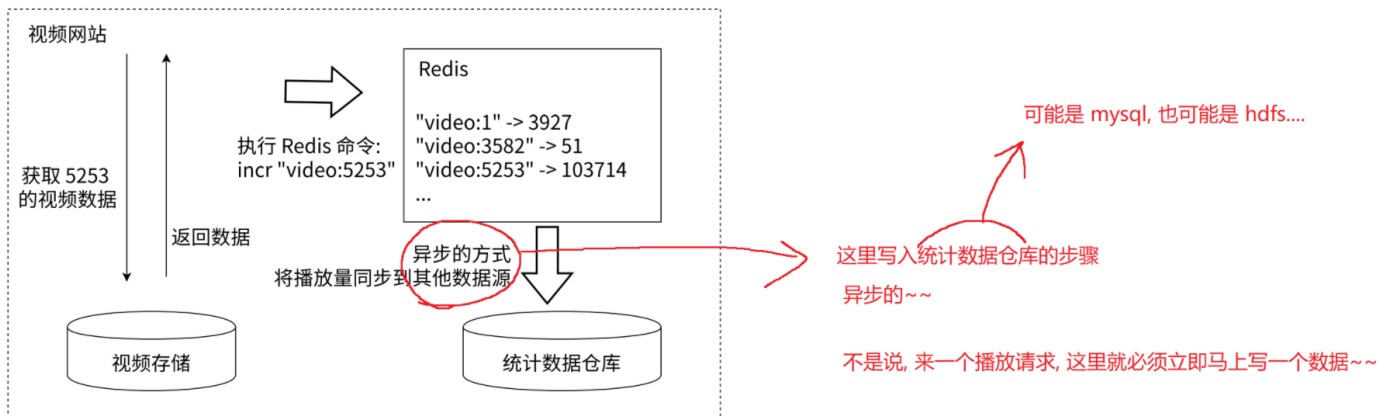
- 把数据写给Redis的同时，给这个key设置一个过期时间
- Redis也在内存不足时提供了淘汰策略

企业为什么总喜欢收集用户数据

统计-->明确用户需求-->迭代产品

Redis并不擅长统计数据

想想统计数据，MySQL打个表，limit限制多少位就可以拿到前多少位的数据，而Redis就很麻烦



六、验证码机制、业务驱动的技术选型与系统设计思维

1. 验证码：系统防护设施，而非用户功能

(1) 为什么不自研验证码？

- 安全对抗成本高（脚本、自动化攻击、模型迭代）
- 投入产出比低，专业平台（如极验、阿里云验证码）提供行为识别、风控策略等整套安全能力。

(2) 验证码规则设计的业务逻辑

- **5分钟有效期**：平衡操作时长与安全风险
- **60秒间隔**：防脚本刷请求，不影响正常用户
- **5次上限**：容错与防爆破的平衡
- **本质目标**：最小化用户体验损失，最大化系统安全。

(3) 验证码的真实防护目标

- 防刷接口、防撞库、防短信轰炸、防流量冲击
- 是“请求过滤器”，为系统入口设置第一道闸门。

(4) 校验必须后端完成

- 前端校验仅为体验优化，所有安全判断必须后端调用第三方接口复核。

- 原则：前端不可信、客户端不可信、请求参数不可信。

2. 业务驱动技术选型：为什么“业务优先”

(1) 业务的本质

业务是公司为用户持续解决问题的能力，技术只是实现手段。

(2) 技术必须服务于业务

- 无“万能技术”，如 Redis 不适于冷数据，Kafka 不适于强一致。
- 技术选型取决于业务规模、并发模型、用户行为。

(3) 12306 的启示：业务规则优于技术硬扛

- 春运抢票是极端业务峰值，非单纯技术问题。
- 分时段放票是通过业务规则打散请求，将“瞬时洪峰”转为“可控流量”。
- **核心思想**：能用业务规则解决的，不依赖技术硬扩容。

(4) 系统设计优先级

1. **业务规则优化**：限流、验证码、分批处理
2. **系统架构优化**：缓存、异步、削峰填谷
3. **技术升级**：分布式、扩容、高可用

逆序操作易导致成本失控与系统复杂。

3. 教育类业务的模块与技术对应

(1) 模块分类与业务目标

- **内容交付**（直播、录播、训练营）：追求高质量与稳定体验
- **能力训练**（题库、刷题、答疑）：追求高频交互与即时反馈
- **评估管理**（模考、学习报告）：追求数据准确与并发可控

(2) 技术贴合业务痛点

- **直播课**：低延迟与抗抖动优于纯粹带宽
- **题库**：缓存命中率优于复杂 SQL
- **模考**：高并发与防作弊优于秒级一致性

技术选型的本质：为解决“最痛的点”投入最多的资源。

六、核心设计思想总结与升华

1. Redis 的定位与使用原则

- **编码是自动的，但后果是可知的**：开发者应知晓数据编码可能带来的行为差异。
- **Redis ≠ 数据库**：擅长高速、简单、高并发；不擅长复杂查询、强一致事务、冷数据存储。
- **分布式系统中的角色**：作为“中枢神经”，承担缓存、Session、计数器、轻量队列等核心功能。

2. 业务与技术的正确关系

- **业务决定问题形态**：技术是解决方案，而非目标。
- **架构决定系统上限**：良好的架构能使技术发挥最大价值。
- **12306 的终极启示**：系统扛不住，往往是规则不合理，而非机器不足。

3. 给开发者的实践原则

1. **先问业务**：用户在做什么？业务目标是什么？
2. **再问系统**：系统为何扛不住？瓶颈在哪里？
3. **最后选技术**：何种技术能最优雅地解决问题？

4. 系统设计的升维思考

- **验证码**：是系统防护设施，是入口限流器。
- **Redis 命令**：是并发模式的原子抽象，是“命令即模型”的设计哲学体现。

- **单线程模型**：是无锁、无切换、极致简洁的性能典范。
-

七、扩展建议与后续方向

若你希望进一步深化，可考虑以下方向：

1. 撰写系列博客

- 《Redis 字符串编码深度解析：从 int 到 SDS》
- 《业务驱动的系统设计：从验证码到分布式锁》
- 《12306 背后的架构思维：业务规则如何塑造技术方案》

2. 模拟面试与追问训练

- 以“Redis 字符串编码”为起点，展开到底层内存、并发模型、实际场景的连环追问。
- 以“验证码系统设计”为题，考察安全、用户体验、系统防护的全链路思维。

3. 实战项目设计

- 设计一个具备防刷、限流、缓存、Session 管理的小型高并发系统。
- 实现一个基于 Redis 的分布式计数器 + 异步持久化方案。

总结：

本文从 Redis 的底层编码与命令设计，延伸到缓存实战、Session 共享，再升华至业务驱动的系统设计思维，旨在帮助开发者不仅“会用”技术，更“懂为什么这样设计”，从而在真实业务中做出合理的技术选型与架构决策。技术为业务服务，思想为设计引路，这才是高级工程师的核心竞争力。