

前言

近年来，随着电子商务的蓬勃发展，以及人们对数字化阅读方式的日益追求，当当云阅读电子商务网站作为在线阅读与销售电子书籍的平台，扮演着越来越重要的角色。该网站背后有着庞大的电子书内容数据需要存储，其中文本数据占据主要部分。与此同时，网站已经按照书籍的类别进行了细致分类，并且包括了大量的图片数据与视频数据。

在技术支持方面，当当云阅读网站需求多方面的支持：一方面需要存储海量的书籍数据，另一方面需要支持用户流畅的在线阅读体验，同时还需要处理用户的注册信息与支付信息。因此，针对以上三个方面的需求，当当云阅读迫切需要一套考虑内存与外存的存储架构来解决这些问题。本文将探讨如何设计并实施这样一套存储架构，以满足当当云阅读网站的多方面需求。

设计原则

本系统整体采用微服务构架，对于这种大型系统在设计时需要遵照的原则：

- 高性能：提供快速的访问体验。
- 高可用：网站服务一直可以正常访问。
- 可伸缩：通过硬件增加/减少，提高/降低处理能力。
- 安全性：提供网站安全访问和数据加密、安全存储等策略。
- 扩展性：方便地通过新增/移除方式，增加/减少新的功能/模块。
- 敏捷性：按需应变，快速响应；

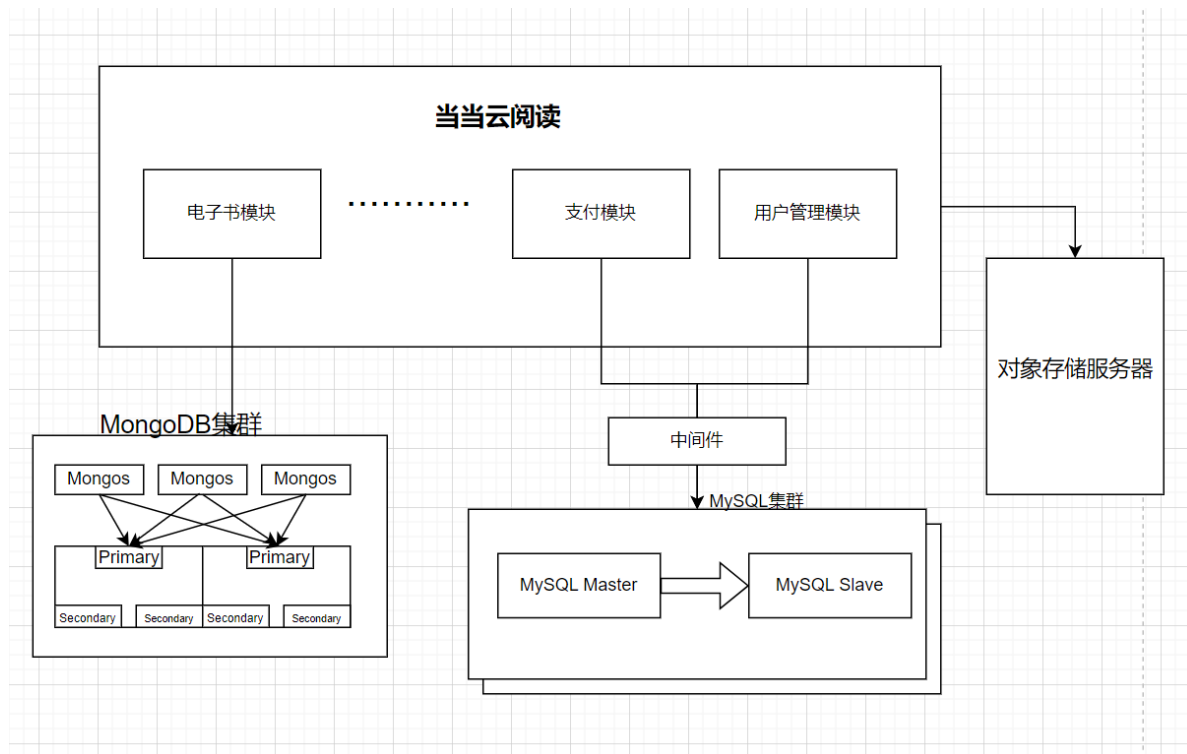
数据分析及可选技术选型

- 文本数据的存储：由于电子书的内容以**文本数据**为主：
 1. 可以使用分布式文件系统（如HDFS）进行存储，它可以提供高可靠性的数据存储，并支持海量数据的存储和处理。电子书的文本数据、图片数据和视频数据都可以存储在HDFS中。对于书籍的分类信息，可以在数据层面进行建模，例如使用标签系统或者目录结构来表示。
 2. 采用NoSQL非关系型数据库MongoDB。
 - 电子书内容主要是文本数据,适合存储在文档型数据库MongoDB中。MongoDB支持Embedding/Nesting数据格式,可以很好地对应电子书分类的树形结构。
 - 同时MongoDB性能优异,支持水平扩展,能很好支撑大量读写访问。
- 图片和视频数据的存储：是非结构化数据，需要处理和存储大规模的二进制文件图片和视频由于所占用的存储空间比较大，因此一般的存储方式是，**在数据库中只保存图片和视频的ID或URL，实际的图片和视频则以文件的方式单独存储。**
 1. 可以选择具有高可扩展性和高可用性的云服务提供商，如阿里云OSS、腾讯云COS等
 2. 使用分布式文件系统HDFS。
 - 图片和视频数据体积较大,采用HDFS可以实现高容量和高可靠性存储。

- HDFS支持数据备份和灾难恢复,保证数据安全。
- 3. 对于图片和视频这类静态资源,可以使用CDN(内容分发网络)来存储和分发。CDN可以提高资源加载速度,减轻主服务器的负担。
- 用户信息的存储: 这部分数据有**明确的字段和关系**,关系型数据库更高效。用户的注册信息和支付信息等敏感数据需要保证安全性和隐私性,需要采用加密算法对敏感数据进行加密保护。同时,可以使用访问控制和身份验证机制来限制对用户信息的访问权限。
 - MySQL, Oracle等主流关系型数据库性能可靠,经过长期优化,可以很好支持高并发访问。

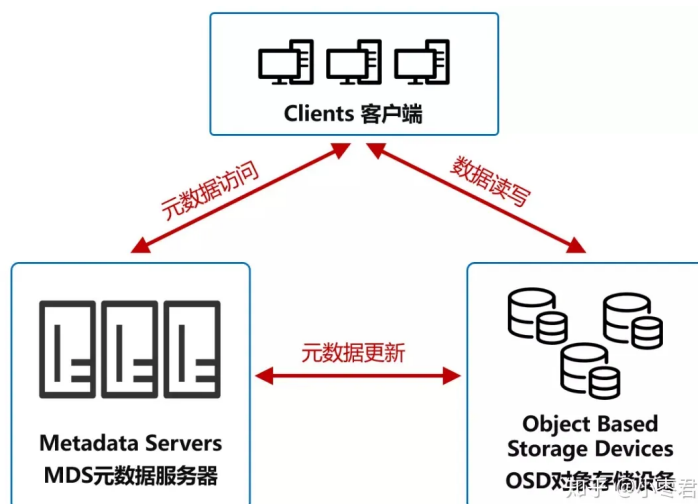
当当云阅读技术选型

初步架构如下:

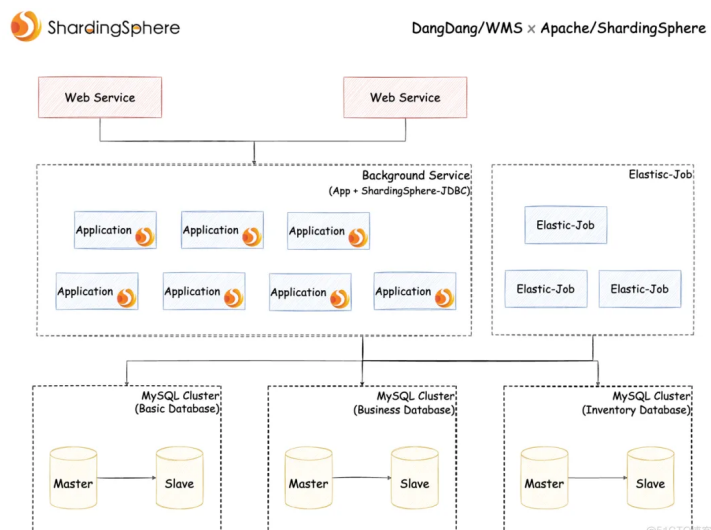


由于MongoDB 本身就是一个基于分布式文件存储的数据库,并且其作为一个面向文档存储的数据库,操作起来比较简单和容易。所以本系统的文本数据采用MongoDB来存储。

而对于图片和视频数据,选择对象存储服务(如下图[1]),它提供可扩展和耐久的存储解决方案,允许添加新数据而不影响现有数据的性能。Amazon S3、Azure Blob Storage和阿里云OSS以及腾讯云COS是常用的对象存储服务,这些厂商的对象存储服务以及相当成熟,综合考虑成本可以选择直接使用云服务厂商提供的服务。

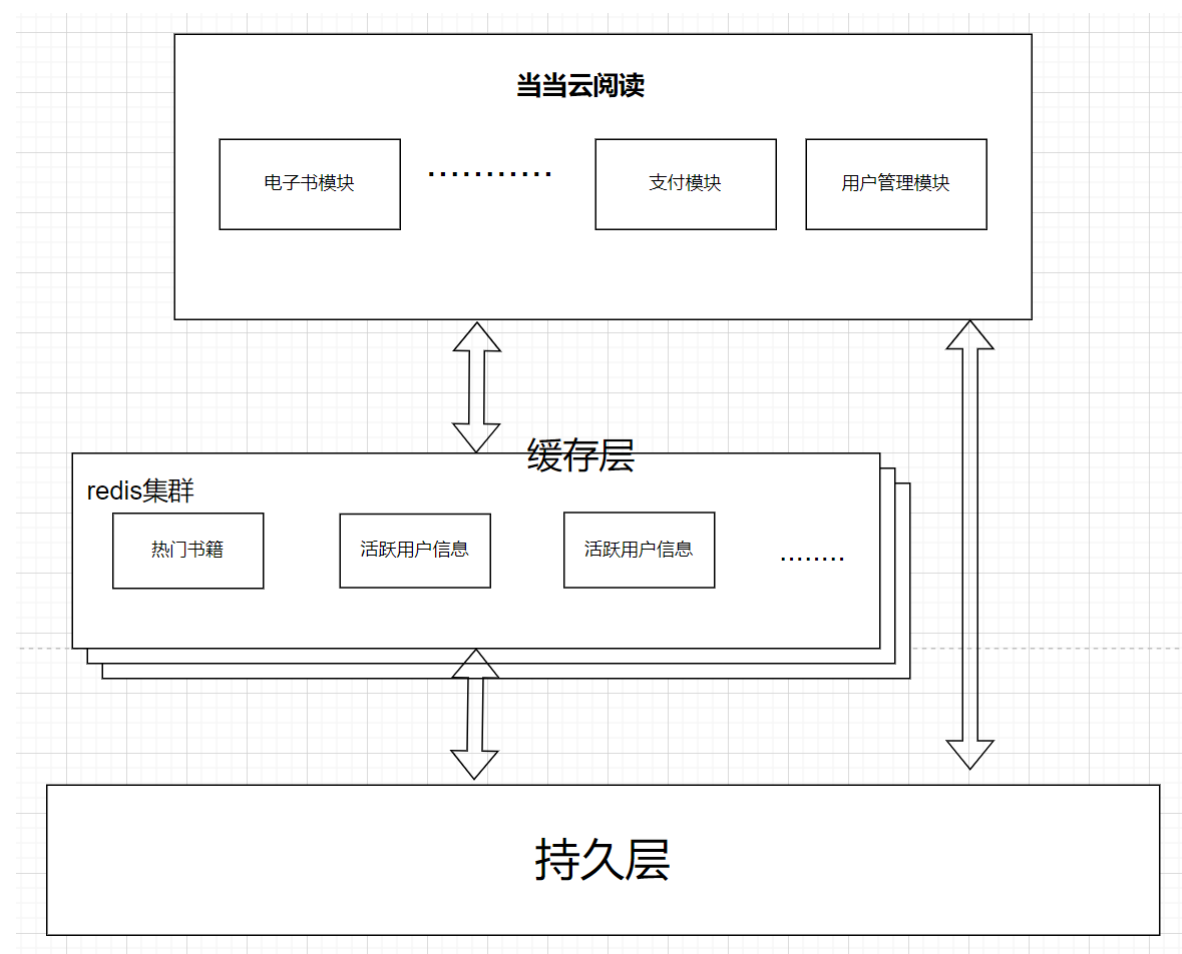


对于用户信息的存储、支付信息等结构化数据。本系统选择开源的MySQL，因为其开源的特性我们在使用的过程中可以针对自身业务定制化开发中间件。如当当网的WMS系统（如下图[2]）中自研分库分表产品Sharding-JDBC。



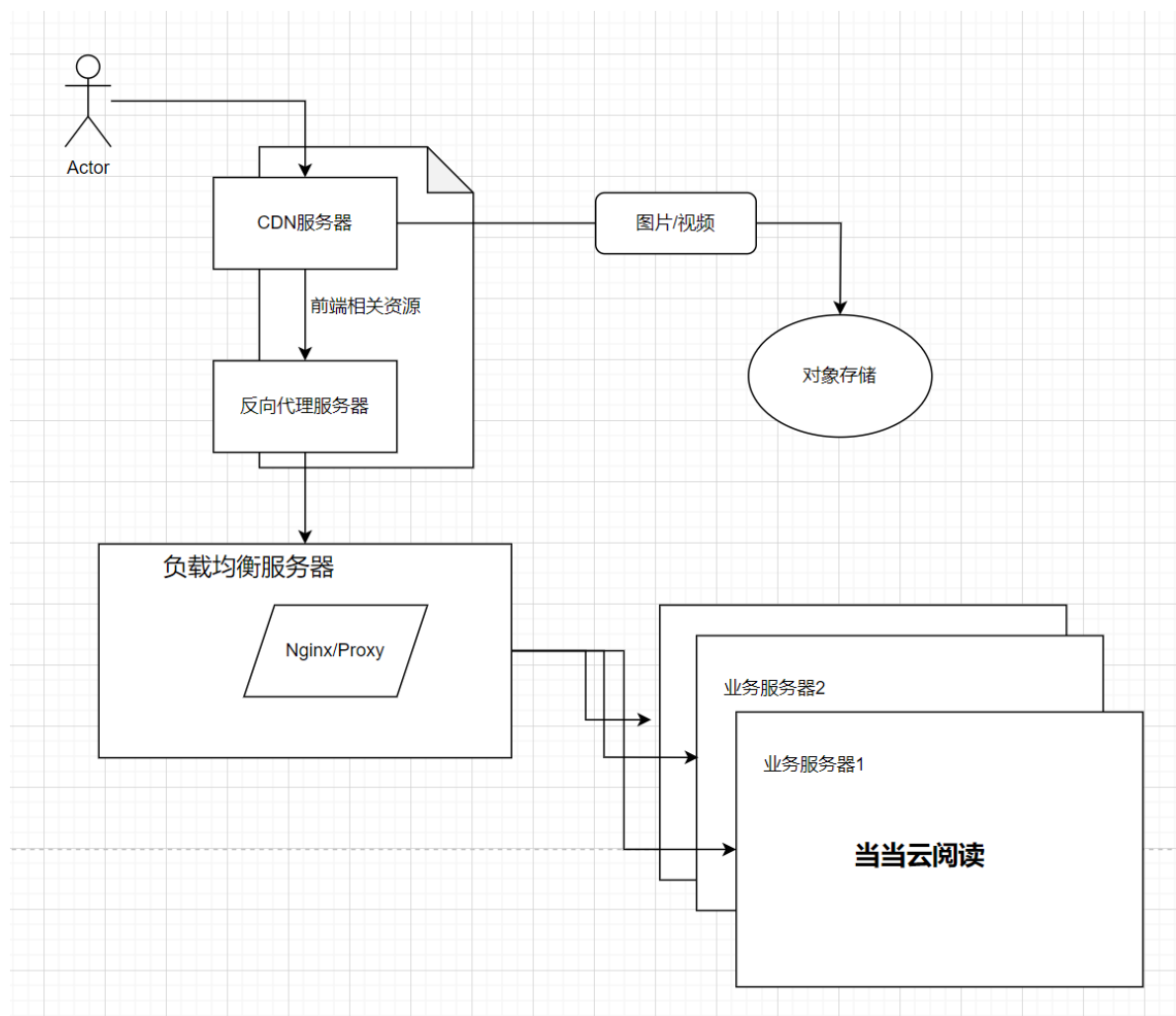
性能优化

a. 为了保证高性能，可以使用Redis或Memcached来实现缓存层。使用缓存主要源于热点数据的存在，大部分网站访问都遵循28原则（即80%的访问请求，最终落在20%的数据上），所以我们可以对热点数据进行缓存，减少这些数据的访问路径，提高用户体验。这些缓存层可以将频繁访问的数据存储在内存中，减少数据库查询的数量，提高页面加载时间。



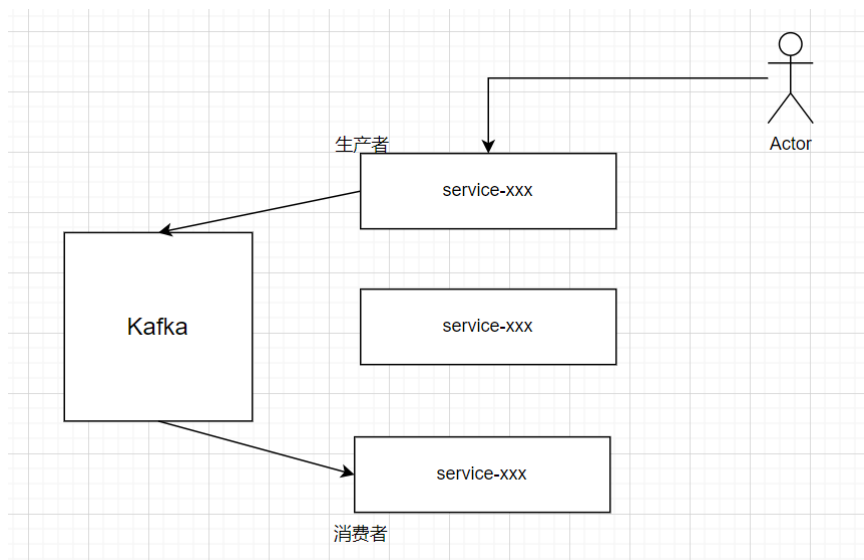
- 数据缓存：为了支持大量用户的实时阅读，应使用内存缓存（如Redis或Memcached）来存储频繁访问的数据，比如热门电子书的内容。这可以大幅降低数据库的读取压力，提高响应速度。

b. 为了处理大量用户，利用CDN（Content Delivery Network）和反向代理部署缓存前端相关资源进一步提高网站性能，同时使用负载均衡器将传入的流量分配到多个服务器上。这将确保没有一个单独的服务器会被请求压垮，并提高系统的整体性能。



- 后端服务采用微服务架构,每个功能组件作为一个独立的服务。
- 前端使用CDN加速页面响应。
- 负载均衡：使用负载均衡器分配用户请求到不同的服务器，可以防止任何单一服务器的过载，确保流畅的用户体验。

c. 为了支持实时更新，可以使用RabbitMQ或Apache Kafka这样的消息队列。这些队列可以处理大量的消息，并提供低延迟的更新，确保用户获得最新的信息。



- 采用Kafka消息中间件,连接后端各个微服务,实现异构数据交互与流程运作。
- 针对实时性要求较高的比如在线阅读功能,可以使用消息队列或发布-订阅模式来实现异步处理和解耦。通过将用户的阅读请求发送到消息队列中,然后由后台系统异步处理并返回结果,可以提高系统的并发能力和响应速度。

其他问题

灵活扩展性: 确保整个存储架构可以根据需求水平和垂直扩展。使用云服务可以方便地根据需要增减资源。

安全和备份: 实施严格的安全措施,包括网络隔离、加密传输和存储,以及定期备份。

高可用性: 通过在不同的地理位置部署多个数据中心,实现数据的高可用性和灾难恢复。

参考引用

[4918字, 详解商品系统的存储架构设计-腾讯云开发者社区-腾讯云 \(tencent.com\)](#)

[阿里电商架构演变之路-阿里云开发者社区 \(aliyun.com\)](#)

[什么是缓存架构, 什么又是后端分布式多级缓存架构, 全文解析 - 知乎 \(zhihu.com\)](#)

[什么是CDN? 它解决了什么难题? 5分钟让你明明白白! -腾讯云开发者社区-腾讯云 \(tencent.com\)](#)

图[1] [对象存储, 为什么那么火? - 知乎 \(zhihu.com\)](#)

图[2] [当当网的IT基础架构 当当网基本介绍mob6454cc7a88c0的技术博客51CTO博客](#)