

云存储应用技术

第一章：云存储概述

丁烨

dingye@dgut.edu.cn

网络空间安全学院

2019-09-05



東莞理工學院
DONGGUAN UNIVERSITY OF TECHNOLOGY

❖ 任课老师

❖ 丁焯，9A-305 室，dingye@dgut.edu.cn

❖ 课程网站

❖ <https://dingye.me/cloud.html>

❖ 课件、作业要求、各类通知、消息均在此网站发布

❖ 教材与参考书

❖ 以讲义为主

❖ 成绩

❖ 考勤 10%、平时作业 30%、大作业 60%

❖ 教学目标

- ❖ 掌握基本的云存储概念
- ❖ 掌握云存储技术相关的软件及工具库
- ❖ 了解作为云存储技术人员所具备的相关能力

❖ 教学方法

- ❖ 理论课：讲授知识点和云存储基本思想，提供相关阅读资料
- ❖ 实验课：通过项目实践真正体会和运用云存储技术相关的软件及工具库

❖ 课程安排

- ❖ 理论课：12 节，共 36 学时
- ❖ 实验课：6 节，共 18 学时，实验课为 UNIX 环境

❖ 理论课

1. 云存储概述
2. 存储技术基础
3. 虚拟化技术
4. 网络存储
5. 文件托管服务
6. 分布式云存储
7. 对象存储
8. 消息队列
9. 内容分发网络
10. 分布式数据库
11. 复习及期末大作业
12. 期末大作业点评

❖ 实验课

1. 虚拟化技术
2. 网络存储
3. 文件托管服务
4. 对象存储
5. 消息队列
6. 期末大作业

- ❖ 大作业
- ❖ 题目会在第 3-4 周公布
- ❖ 学生每 5-8 人一组，实现一套简易的云储存系统
- ❖ 需要准备口试 PPT
- ❖ 无论小组人数多少，组内成员均获得同样的成绩

大数据与数据科学

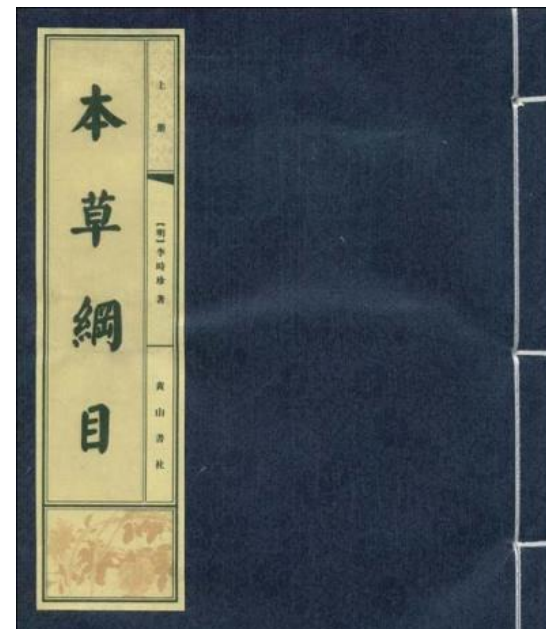
大数据架构

云存储概述



《神农本草经》

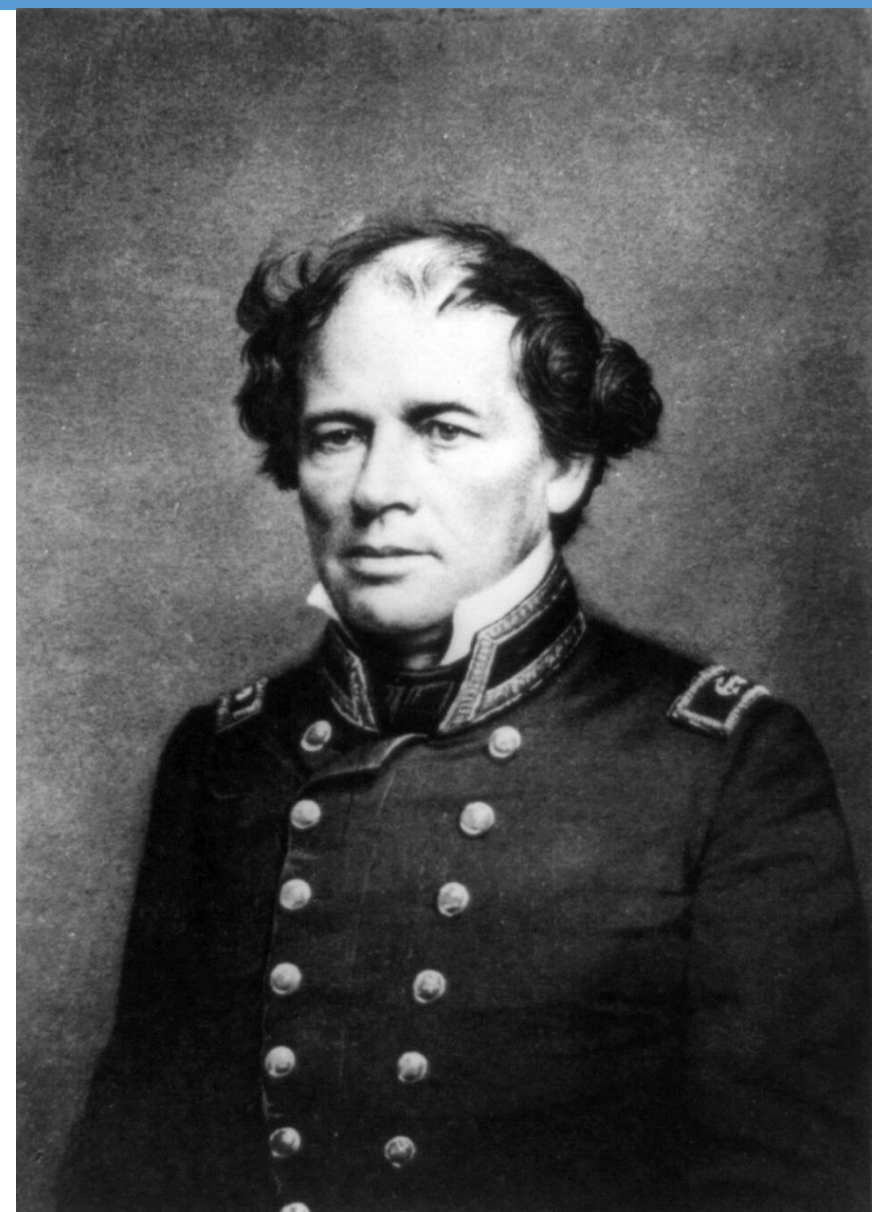
上古，先秦，秦汉时期
多位医家集结整理
上中下三卷，载药 365 种



《本草纲目》

明朝李时珍
历时 27 年编纂，1590 年出版
共 52 卷，载药 1892 种，方剂 11096 个

- ❖ 马修·方丹·莫里
- ❖ Matthew Fontaine Maury
- ❖ 1806 年出生于美国弗吉尼亚
- ❖ 1824 年刚刚达到入伍年龄便进入了美国海军学校
- ❖ 1839 年，已经晋升为海军上尉的莫里在一次事故中不幸腿部致残
- ❖ 不适合于服役远航的莫里在 1842 年被任命为主管海图和仪器库的负责人



大数据与数据科学

每个时代都有大数据



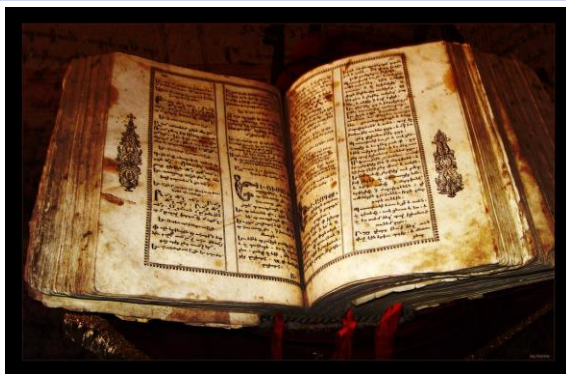
六分仪



经典的书籍、教材



指南针



大量快发霉的航海日志

前人视为垃圾

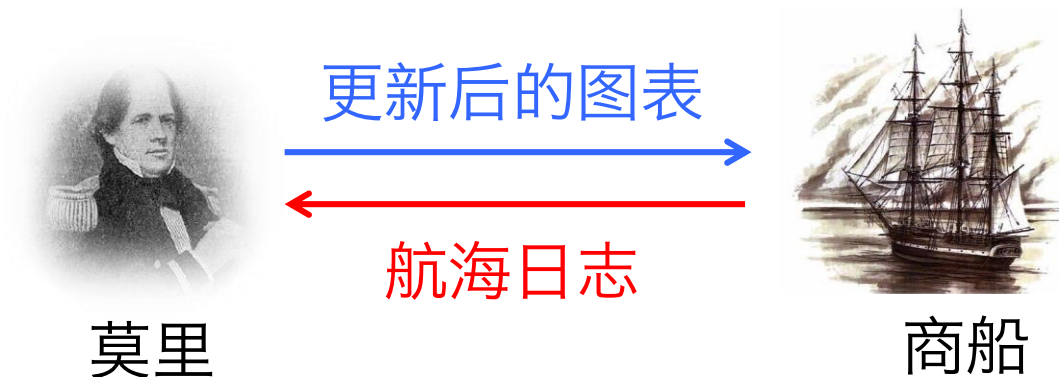
变废
为宝

莫里的目标

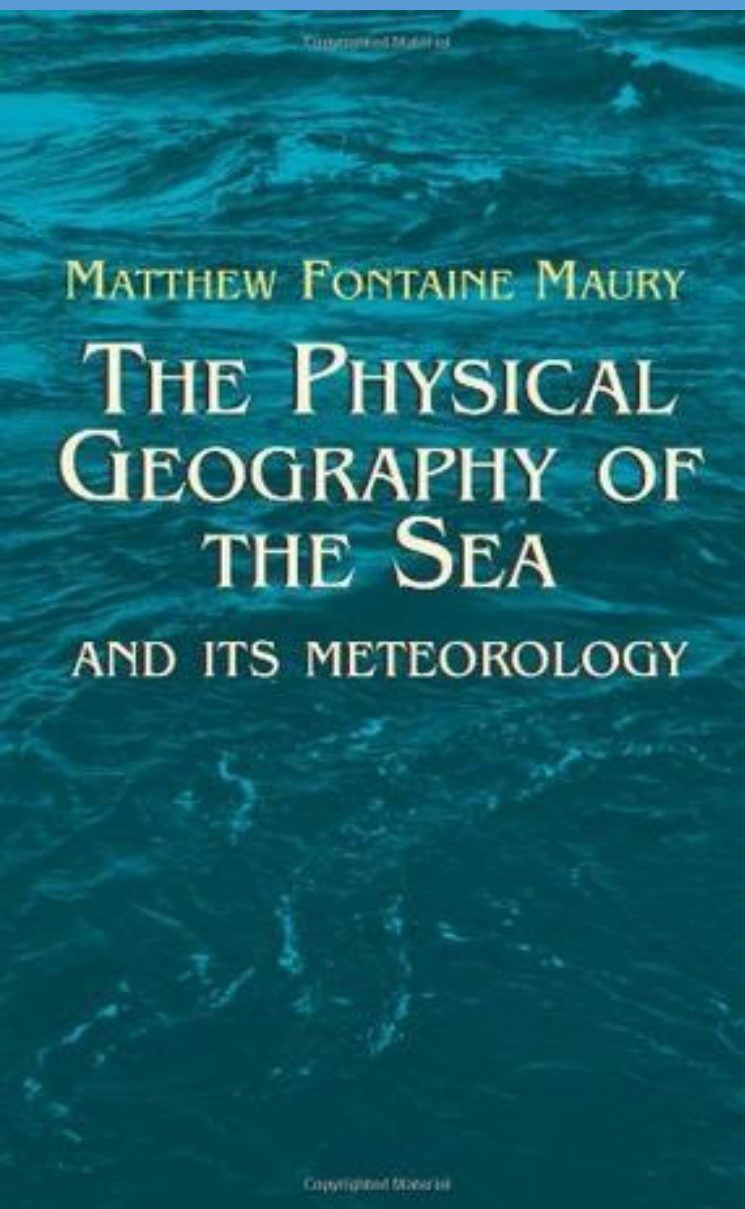


精确的图表

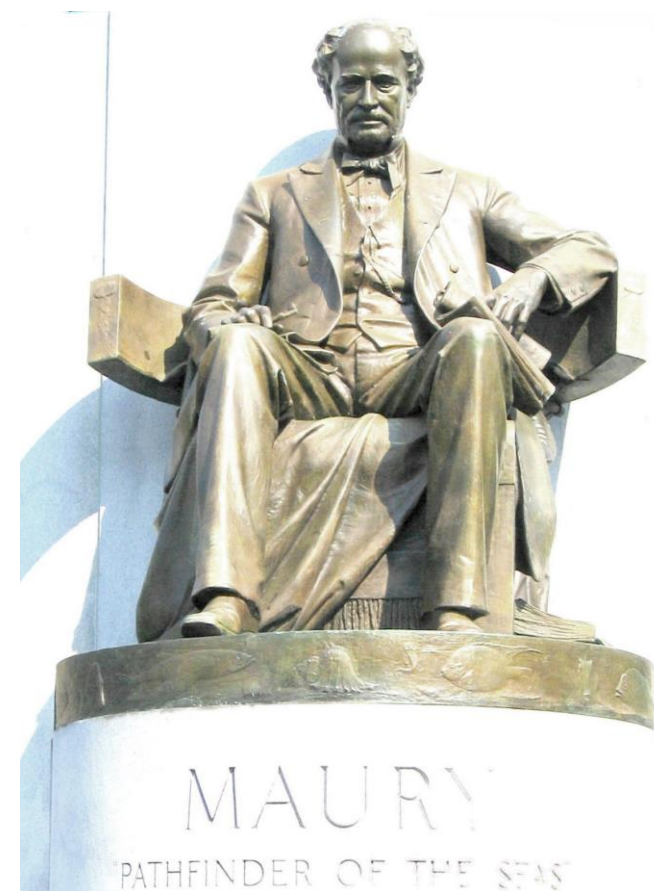
与商船交换信息，在自愿基础上以互利互惠的合作方式开创了国际气象界公开交换环境资料的传统



- ❖ 让商船定期向海里扔掷标有日期、位置、风向以及当时洋流情况的瓶子
- ❖ 寻回瓶子，记录信息，更新数据



- ❖ 1855 年,莫里出版权威著作《海洋物理地理学和气象学》, 被誉为海洋学的奠基人
- ❖ 当时, 他已经绘制了 120 万个数据点
- ❖ 四个国家授予了他爵士爵位, 包括梵蒂冈在内的其他八个国家还颁给了他金牌奖章
- ❖ 即使到今天,美国海军颁布的导航图上仍然有他的名字



信息技术令人类获取
数据的能力大幅提高



交通数据



金融数据



物联网数据



零售数据



社交网络数据



科学数据

无处不在的大数据

据麦肯锡全球学会预计，美国各行各业需要 14-19 万名拥有“深度分析”专长的工作者

每一天，有……

2940亿封电子邮件发送

平均每个地球人每天发送42封

1288个新移动应用可被下载

日均下载量已达3500万

10万3680小时视频

上传到YouTube

2亿5000万张照片

上传到Facebook

如果把它们都印出来，叠起来能有

80个埃菲尔铁塔那么高

每年上传的总数据为35ZB

即35,000,000PB

100TB数据上传到Facebook

如果用2TB的硬盘储存

每年Facebook要新购11吨硬盘

20亿小时电视与电影

在Netflix上观看

整个因特网的流量信息可以装满

1亿6800万张DVD光盘

需要67节载重50吨的车皮运送

2亿3000万条tweets

在Twitter上发布

据预测，到2020年

每年数据量增长60%

其中非结构化数据增长80%

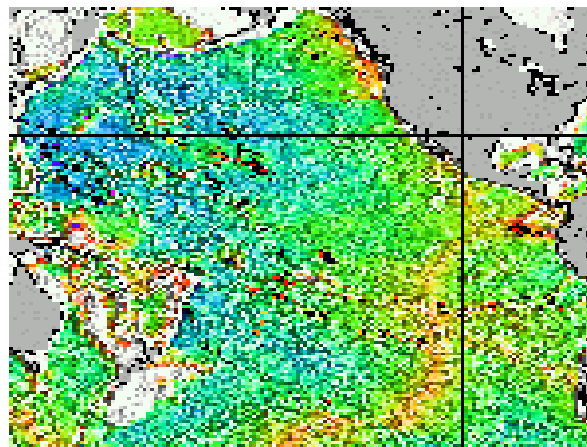
- ❖ 计算机基础存储单位：字节
- ❖ 1 KB = 1,024 字节
- ❖ 1 MB = 1,048,576 字节
- ❖ 1 GB = 1,073,741,824 字节
- ❖ 1 TB = 1,099,511,627,776 字节
- ❖ 1 PB = 1,125,899,906,842,624 字节
- ❖ 1 EB = 1,152,921,504,606,846,976 字节
- ❖ 1 ZB = 1,180,591,620,717,411,303,424 字节

640 KB 足够所有人用了！
——比尔·盖茨，1981

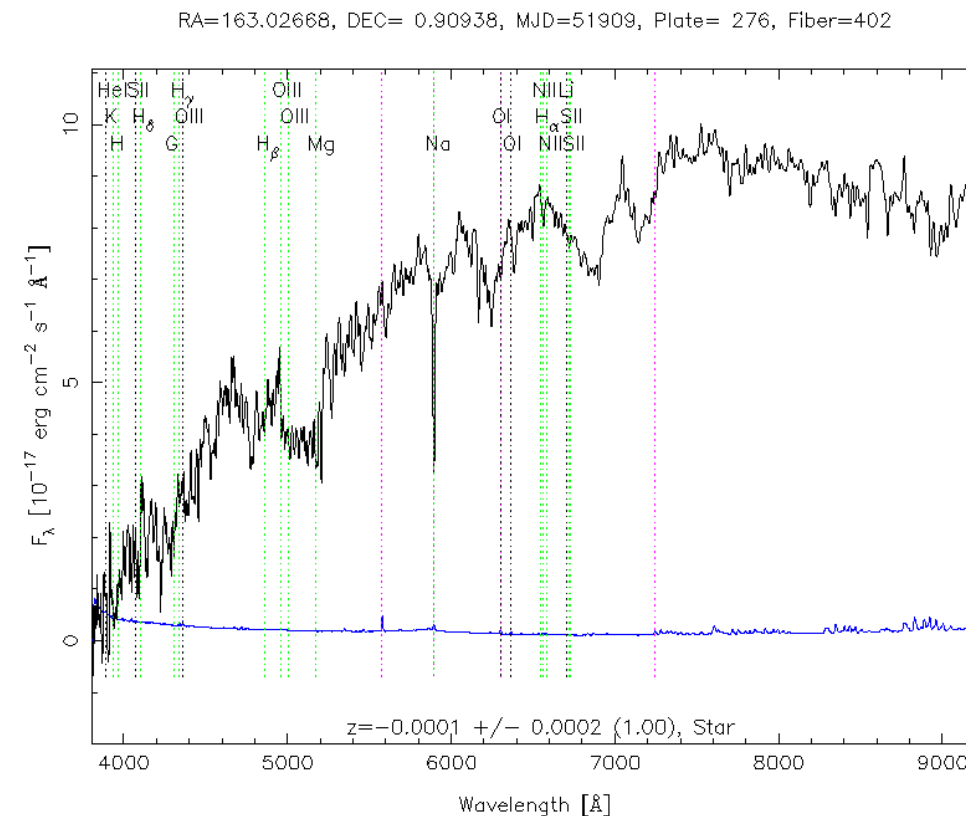


数据科学研究

- ❖ 绘制宇宙：斯隆数字化巡天
- ❖ 绘制地图是一个逐步提高人类知识的主要活动过程
- ❖ 斯隆数字化巡天绘制出人类历史上最大的天体图，覆盖全天四分之一星空
- ❖ 它还能对成百万的天体进行定位和光度绝对定标，并记录 100,000 个类星体（目前所知的最远的天体）的距离
- ❖ 一晚上的观测会产生多至 200 GB 的数据
- ❖ 重要事件：计算机科学家 Jim Gray 和天文学家 Alex Szlay 的合作



- ❖ 在斯隆巡天提供的浩瀚数据的基础上，科学家们已经有了许多重大的进展
- ❖ 比如小行星带的组成
- ❖ 比如褐矮星的发现
- ❖ 比如星系晕的来源
- ❖ 在人类探索宇宙的道路上，斯隆巡天提供的大数据是我们强有力的武器



❖ 斯隆之后

❖ **LSST**, 96 年提出

❖ 2014 年八月开建, 2022 年建成

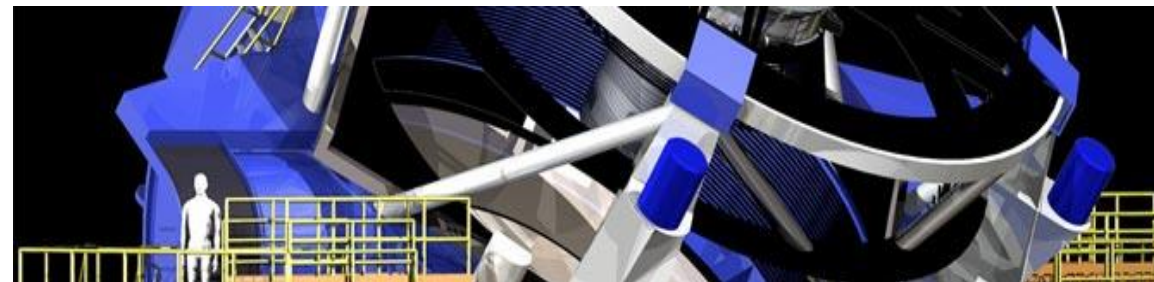
❖ 8.4 米长, 30 亿像素相机

❖ 每晚 30 TB 原始数据

❖ **GWAC**, 正在进行时

❖ 72 个 (36 个一组, 分两组)

❖ 每晚 20 TB 原始数据



- ❖ 谷歌工程师 vs. 疾病控制与预防中心
- ❖ 2009 年出现了一种新的流感病毒。这种甲型 H1N1 流感结合了导致禽流感 and 猪流感的病毒的特点，在短短几周之内迅速传播开来，有的评论家甚至警告说，可能会出现类似于 1918 年在西班牙爆发的夺走了数千万人性命的大规模流感
- ❖ 在研发出疫苗之前，公共卫生专家能做的只是减慢它传播的速度。但要做到这一点，他们必须先知道这种流感出现在哪里
- ❖ 美国，和所有其他国家一样，都要求医生在发现新型流感病例时告知疾病控制与预防中心 (CDC)
- ❖ 但由于人们可能患病多日后才会去医院，同时这个信息传达回疾控中心也需要时间，通告新流感病例时往往会有一两周的延迟
- ❖ 对于一种飞速传播的疾病，信息滞后两周的后果将是致命的。这种滞后会导致公共卫生机构在疫情爆发的关键时期无所适从

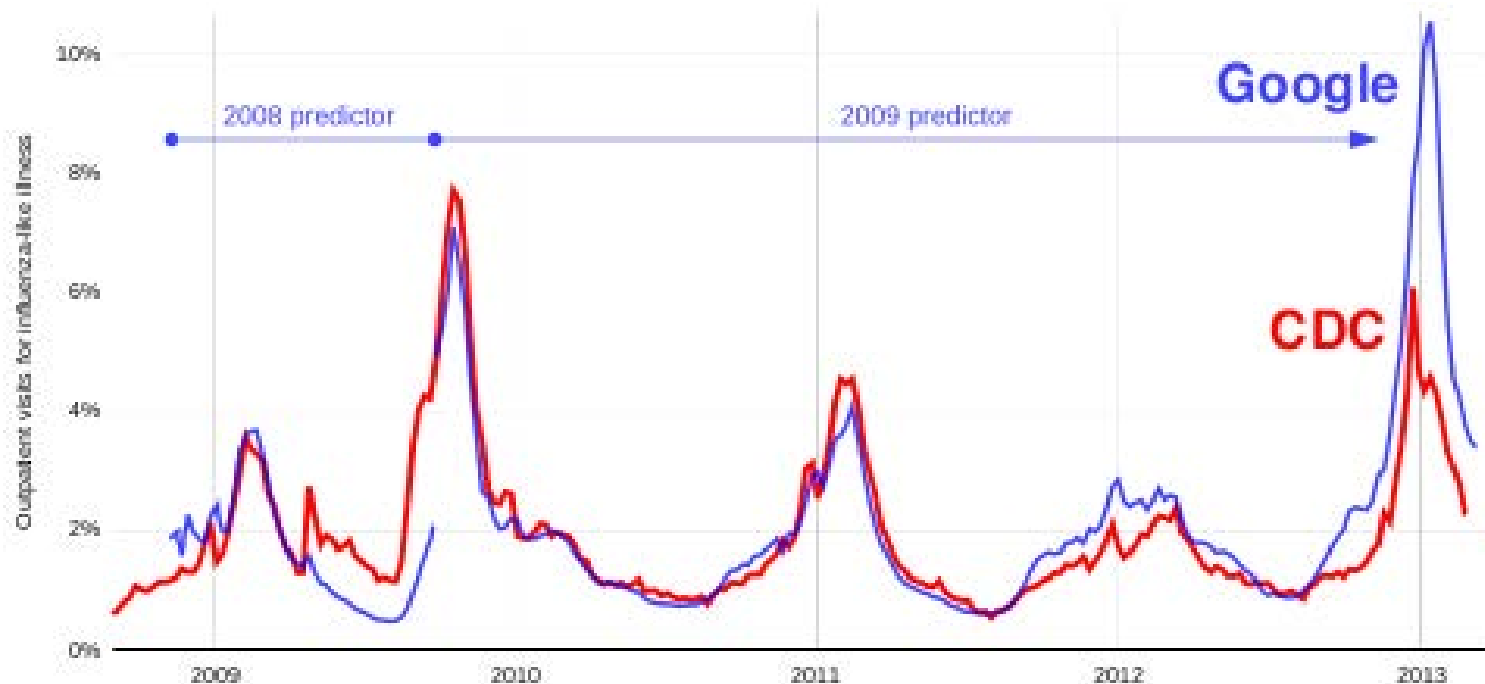
保存多年的
搜索记录

5000 万条
最频繁检索词条

45 条与流感
相关的词条组合

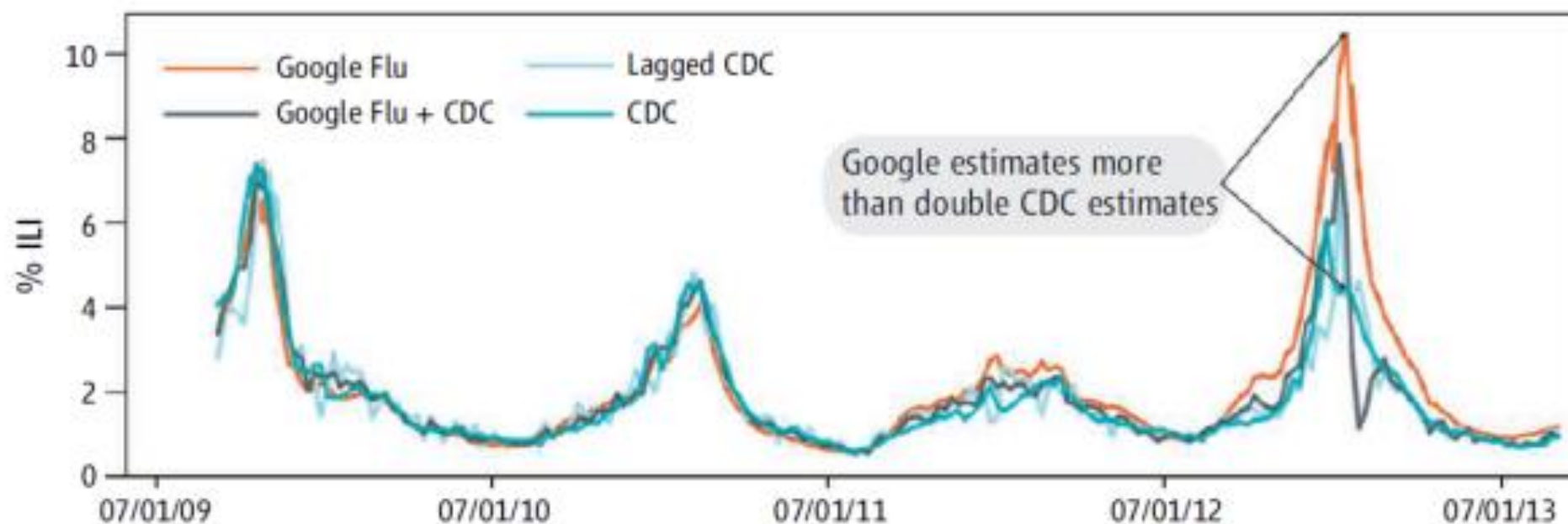
数学模型

每天收到 30 亿条
的搜索指令



预测无延时，相关性高达 97%

- ❖ 2009 年，谷歌对预测系统进行了更新，在数据库中加入疾病控制与预测中心在 H1N1 时期的数据
- ❖ 同时，谷歌将与流感高度相关的检索词条的组合扩展到了 160 条
- ❖ 其他公司也曾试图确定这些相关的词条，但是他们缺乏像谷歌一样庞大的数据资源、处理能力和统计技术



Google 对用户搜索、邮件等数据进行分析
实现**在线广告**的精准投放

美国总统奥巴马的竞选团队依据选民的微博
实时分析选民喜好

美国疾病控制和预防中心依据网民搜索
分析全球范围内流感等**病疫的传播状况**

欧洲粒子物理研究中心（CERN）分析大型强子对撞机（LHC）记录的
超过 100 PB 的数据，“几乎”**发现了上帝粒子**

Amazon 利用 Kindle 数据分析用户阅读习惯
向用户推荐其可能购买的书籍

对城市交通数据进行深度分析与挖掘
对城市道路规划提出建议

对冲基金依据购物网站的顾客评论
分析企业产品销售状况

.....

数据

知识

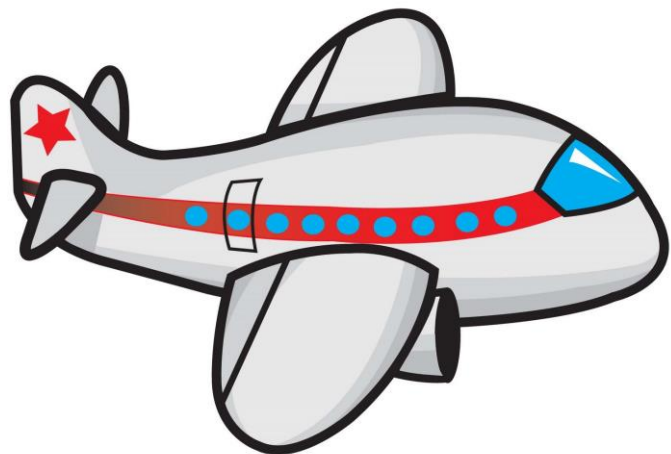
价值

大数据存储、管理、处理、挖掘、应用

数据是**未来的新石油**

数据不仅能帮助科学研究
还能够产生商业价值

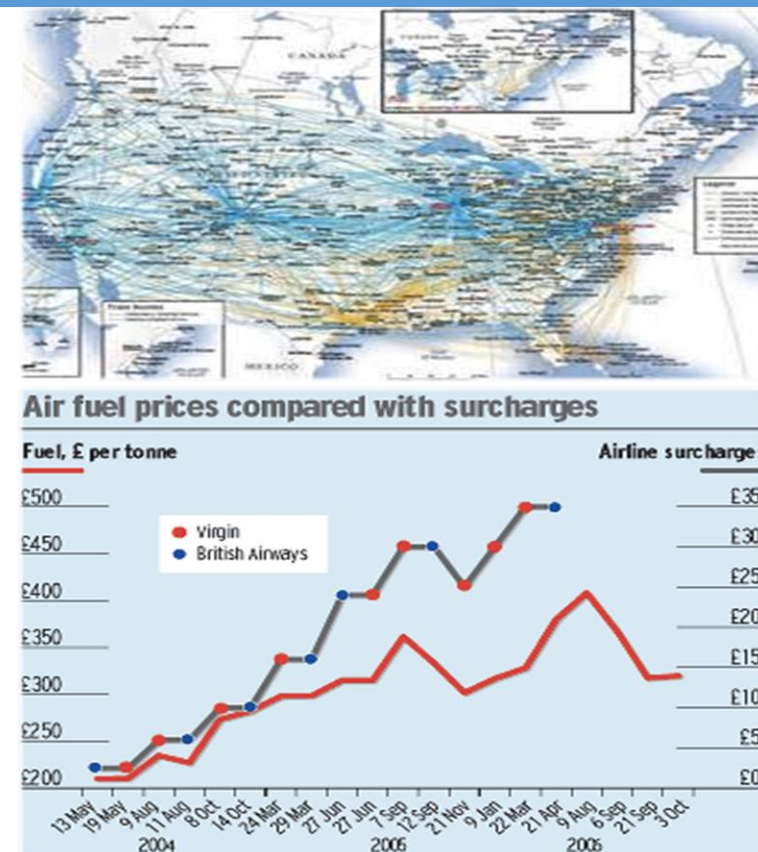
奥伦·埃齐奥尼准备乘坐从西雅图到洛杉矶的飞机去参加弟弟的婚礼。他知道飞机票越早预订越便宜，于是提前几个月就在网上预订了机票



在飞机上，埃齐奥尼好奇地问邻座乘客花了多少钱购买机票。当得知那个人的机票比他买得更晚，但是票价却比他便宜得多时，他感到非常气愤



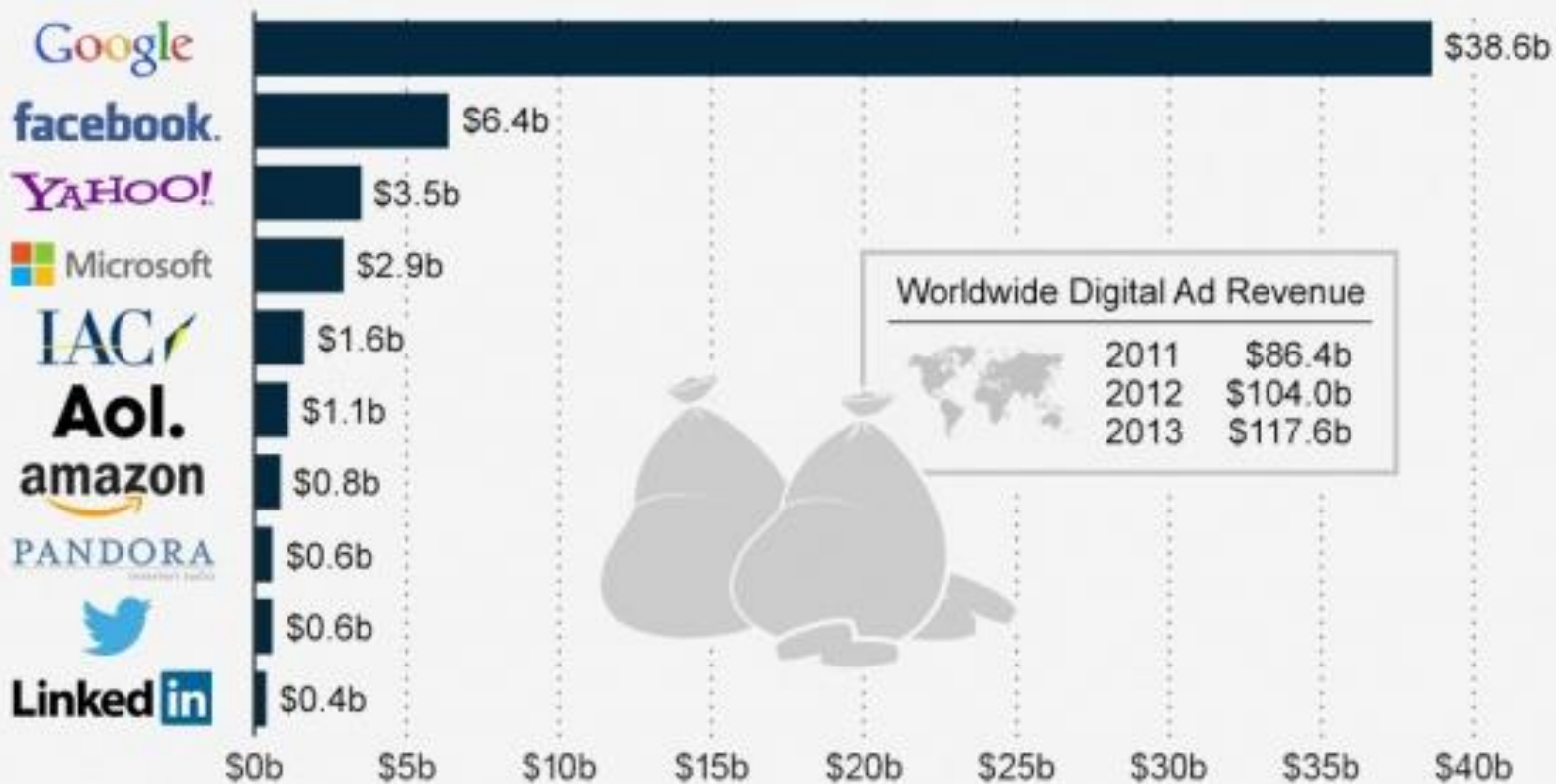
- ❖ 在风投的支持下，埃齐奥尼创立了一家名为 **Farecast** 的科技创业公司，专门预测机票价格的走势以及增降幅度，帮助消费者抓住最佳购买时机
- ❖ 为了提高预测的准确性,埃齐奥尼找到了一个机票预订数据库
- ❖ 到 2012 年为止，Farecast 系统用了**近十亿条**价格记录来帮助预测美国国内航班的票价
- ❖ Farecast 票价预测的准确度已经高达 **75%**，平均每张机票可节省 **50 美元**
- ❖ 埃齐奥尼计划将这项技术应用到其他领域，比如宾馆预订、二手车购买等
- ❖ 但是在他实现计划之前，微软公司以 **1.1 亿美元**的价格收购了 Farecast 公司。而后，这个系统被并入必应搜索引擎



- ❖ 谷歌的金库：精准广告投放
- ❖ 2013 年，谷歌从在线广告获得 386 亿美元的收益
- ❖ 基本相当于全球在线广告收入的三分之一
- ❖ 这 386 亿美元的收益占到谷歌 2013 年总营收的 65%

Google to Rake in 33% of Online Ad Revenues This Year

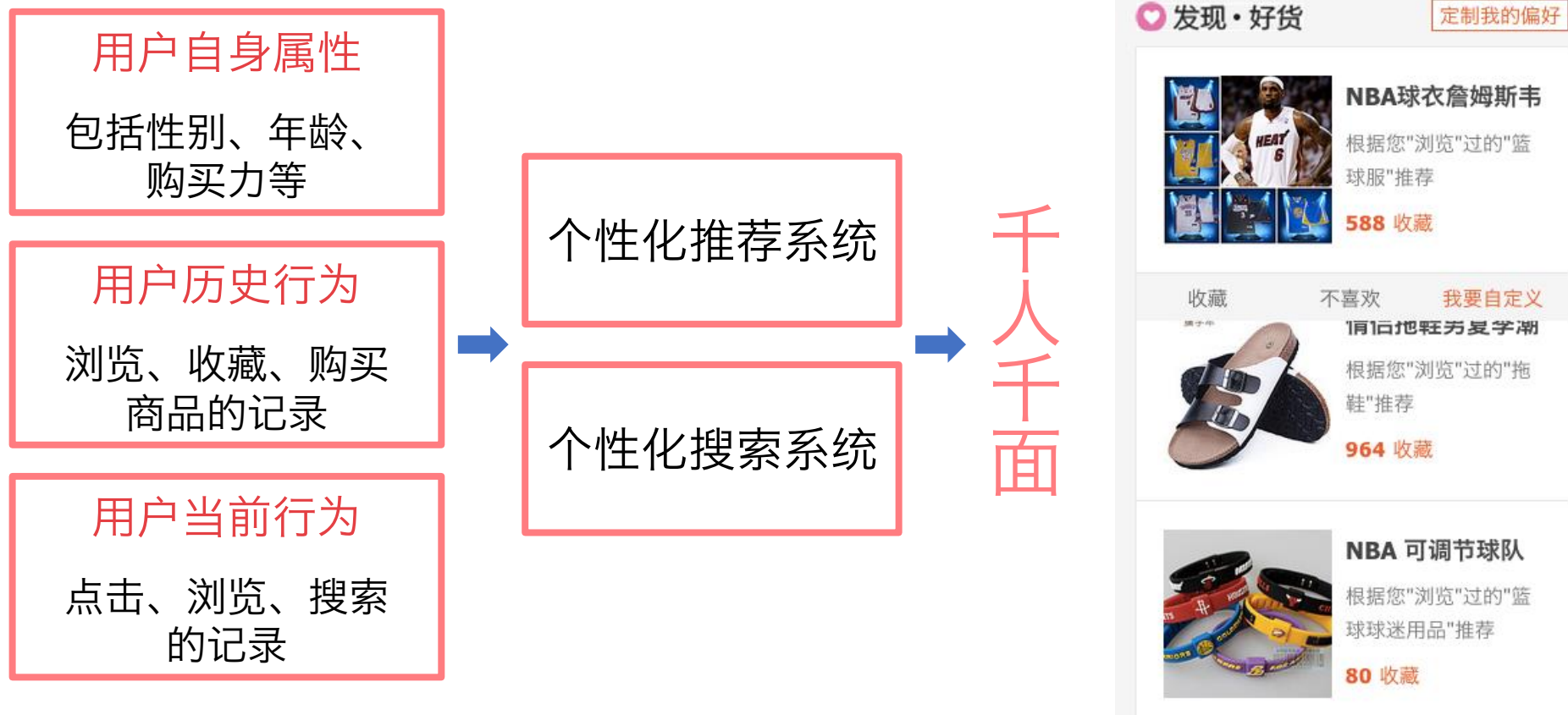
Worldwide digital advertising revenue forecast for the 10 largest ad publishers in 2013





- ❖ 2013年11月12日凌晨，阿里公布了淘宝天猫“双十一”购物狂欢节全天的销售额：
- ❖ 支付宝全天成交金额为 **571 亿**
- ❖ 比去年的 350 亿增长 **63%**





阿里巴巴
保存多年的
各类数据宝藏



	A	B
1	user_id	10944750
2	性别	女
3	年龄	23
4	收货地址	广东省广州市
5

用户资料数据

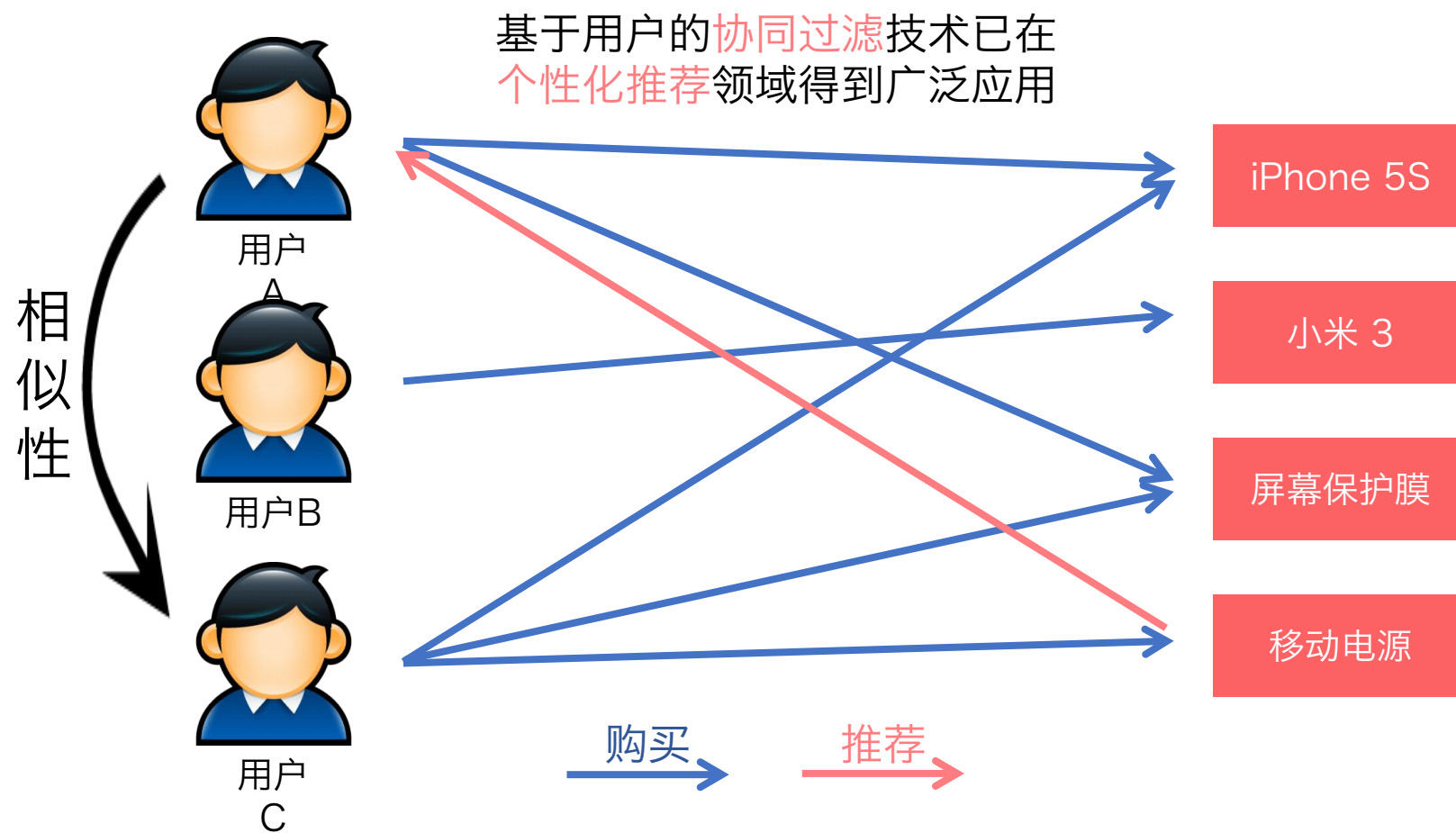
	A	B
1	brand_id	13451
2	商品类型	短裙
3	品牌	ZARA
4	价格	299
5	历史销量	3485
6

商品资料数据

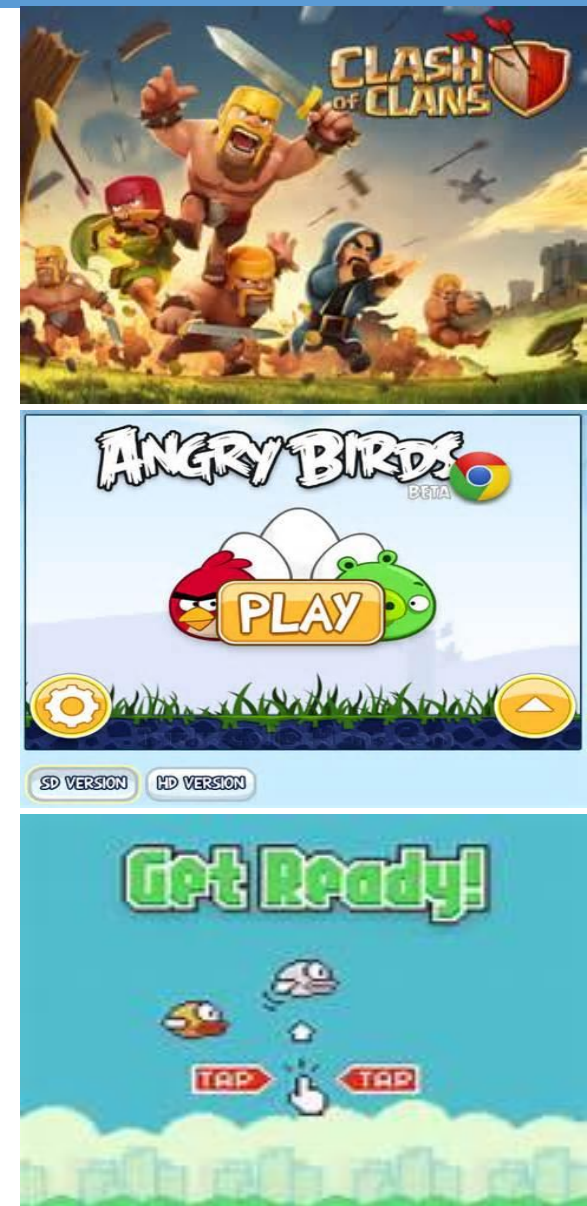
	A	B	C	D
1	user_id	Time	type	brand_id
2	10944750	6月4日	0	13451

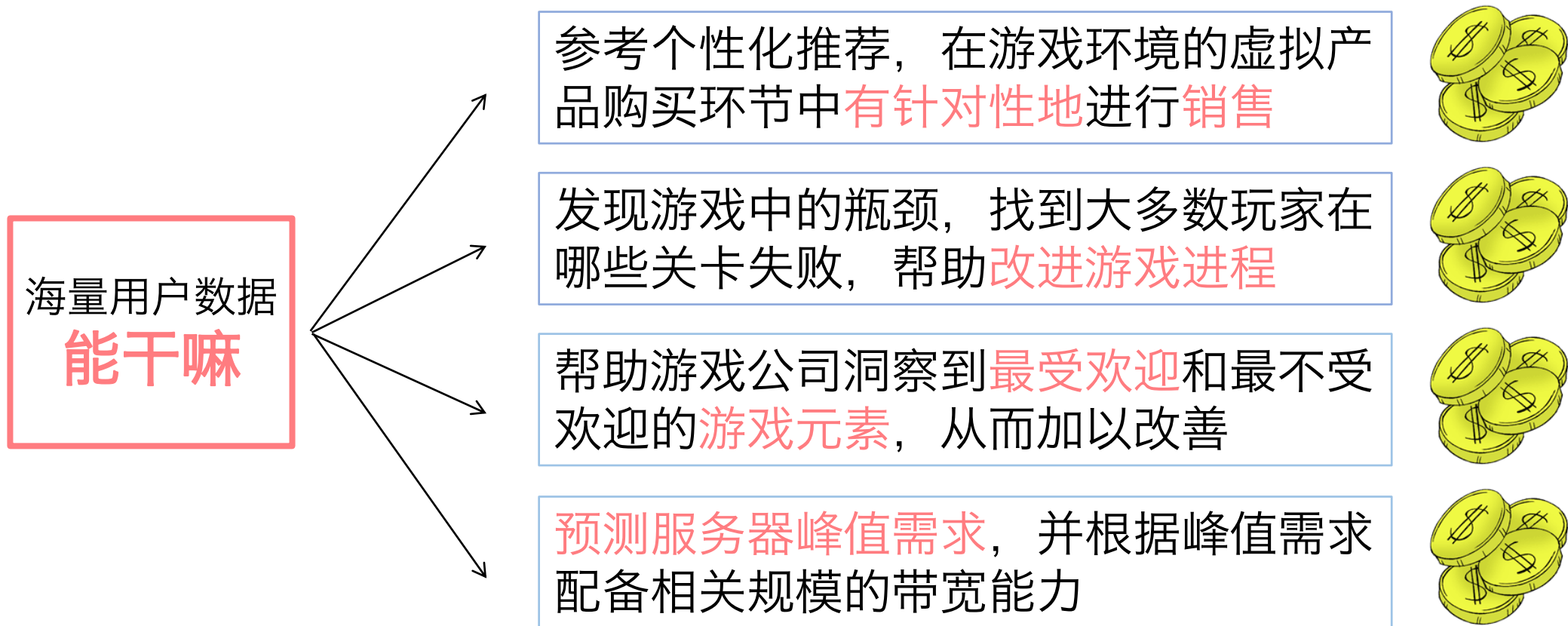
字段	字段说明	提取说明
user_id	用户标记	抽样&字段加密
Time	行为时间	精度到天级别&隐藏年份
action_type	用户对品牌的行为类型	包括点击、购买、加入购物车、收藏4种行为 (点击: 0 购买: 1 收藏: 2 购物车: 3)
brand_id	品牌数字ID	抽样&字段加密

用户行为数据



- ❖ 游戏产业的数据收集
- ❖ 2013 年，仅在美国市场，游戏行业就创造了超过 200 亿美元的收入，而且正在迅速增长中
- ❖ 无论是通过微博连接的在线社交游戏，还是用 Xbox 玩的离线游戏，当玩家开始玩游戏时，都会产生大量不同格式的数据信息
- ❖ 他们在游戏过程中所进行的一切操作都会创造出的海量数据流，包括玩家之间是如何相互配合的、他们玩多久、与谁在什么时间玩了游戏、玩家在虚拟产品上的花费如何、游戏过程在与谁聊天了等等





20Q: 一个简单的大数据应用

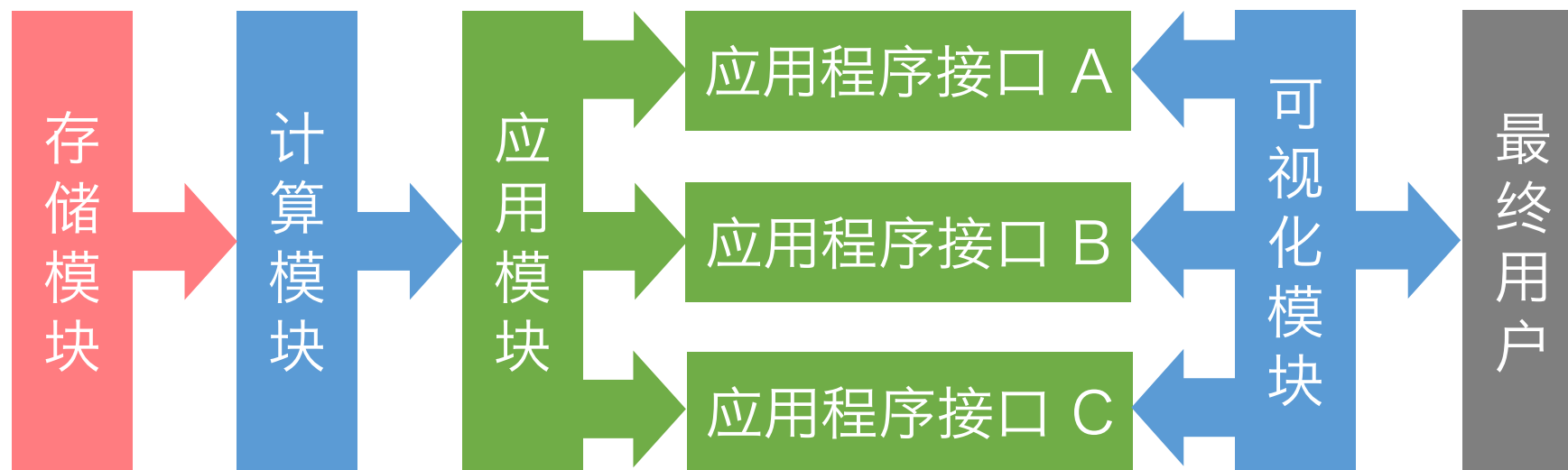
<http://20q.net/>

大数据为我们带来了
新的思维模式和新的发展机会

大数据与数据科学

大数据架构

云存储概述



- ❖ **存储模块**：存储大规模海量数据
- ❖ 文件存储、对象存储、消息队列、数据库等
- ❖ 本课程的主要内容

- ❖ **计算模块**: 对大规模海量数据进行计算和分析
- ❖ **分布式计算引擎**: Spark、Flink、Hadoop 等
- ❖ **机器学习引擎**: TensorFlow、PyTorch 等

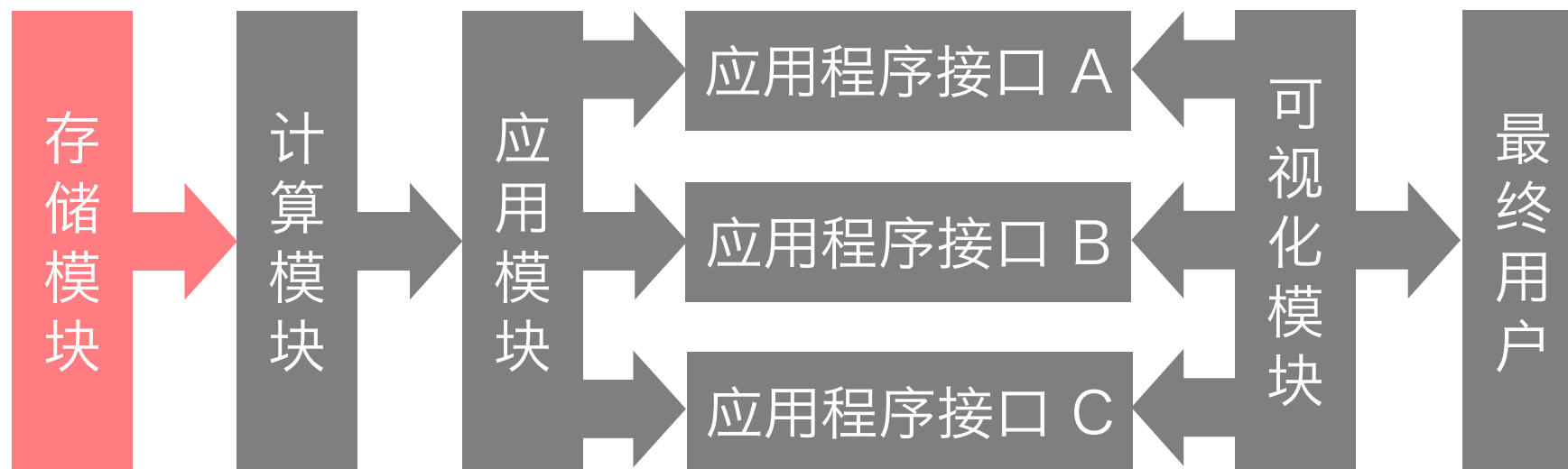


- ❖ **应用模块**：控制计算引擎工作，与最终用户进行抽象沟通
- ❖ 应用程序接口（API, Application Programming Interface）
- ❖ 五大最主流互联网应用程序框架：
- ❖ <https://github.com/topics/framework>



- ❖ 可视化模块：绘制容易理解的图表展现数据分析结果
- ❖ D3.js: <https://d3js.org/>





大数据与数据科学

大数据架构

云存储概述

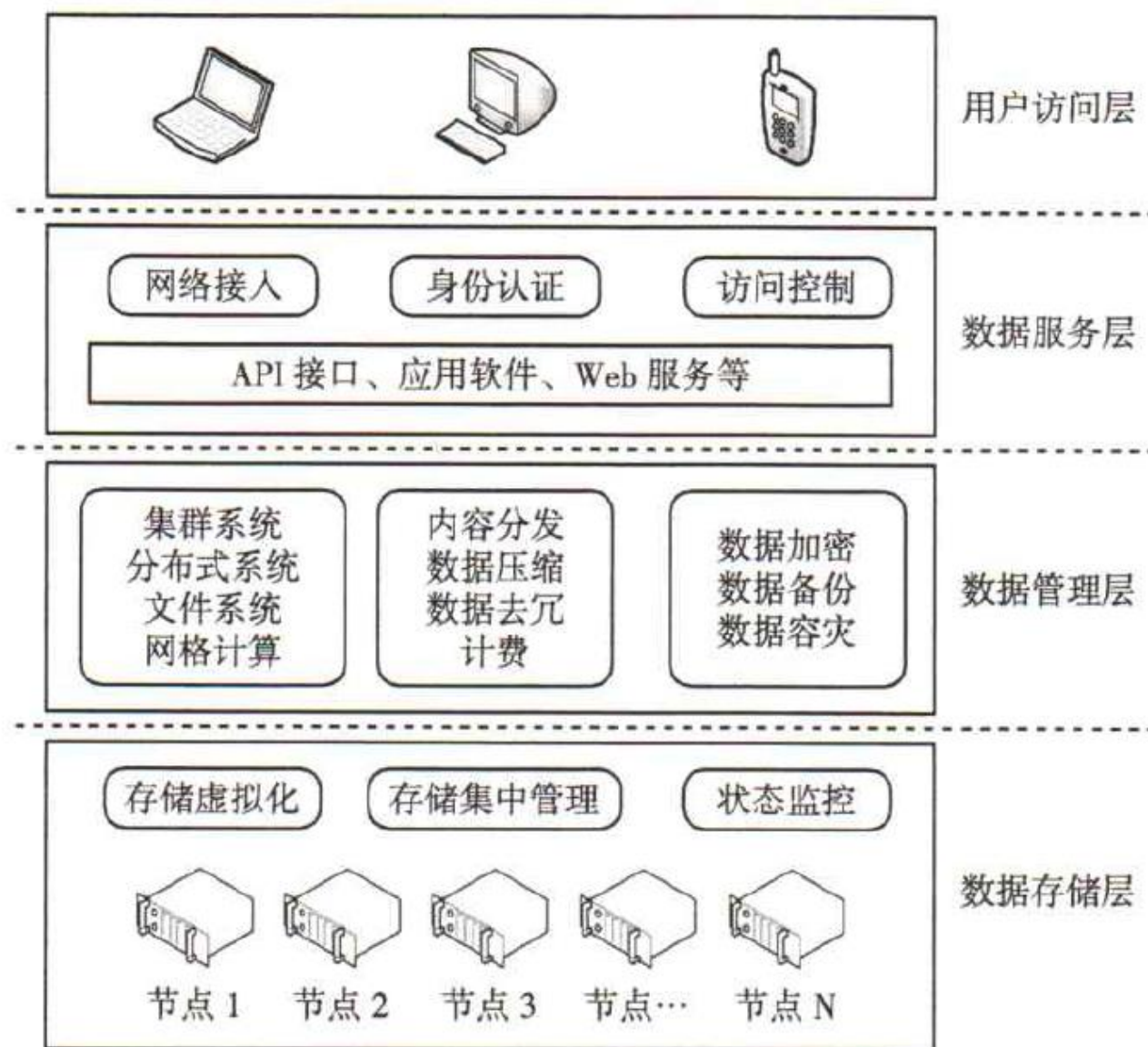
- ❖ 传统的数据存储
- ❖ 将数据存储于磁盘上，本地访问

- ❖ 网络存储
- ❖ 将数据存储于远程服务器上，通过数据通讯协议远程访问

- ❖ 云存储
- ❖ 更高级的网络存储
- ❖ 使用了更复杂的认证机制
- ❖ 通过存储集群可以存储大规模的海量数据
- ❖ 通过数据压缩和去冗余节约存储成本
- ❖ 具备多种数据管理、状态监控的功能

云存储概述

云存储的系统架构



❖ 基础网络存储

❖ 集中式网络传输协议:

- ❖ File Transfer Protocol, FTP
- ❖ SSH File Transfer Protocol, SFTP
- ❖ Server Message Block, SMB

❖ 分布式网络传输协议:

- ❖ BitTorrent
- ❖ eDonkey2000
- ❖ Gnutella

```
magnet:?xt=urn:btih:d21f8d9d004d99c2463acb4b3325495  
bade38693
```

```
ed2k://|file|/%E9%BB%91%E9%95%9C.Black.Mirror.S01E01  
.Chi_Eng.WEB-HR.AAC.1024X576.x264-  
YYeTs%E4%BA%BA%E4%BA%BA%E5%BD%B1%E8%A7%86.mkv|45445  
8192|4f4bd021833d78b384d43191b7c790eb|h=d4t7ysuqmqc  
xrc32ma6m5tk5vj6mbuc6|/
```

❖ 文件托管服务

amazon drive

Yandex Disk

Dropbox

SugarSync



百度网盘
让美好永远陪伴



Google Drive



❖ 分布式云存储



❖ 对象存储



❖ 消息队列

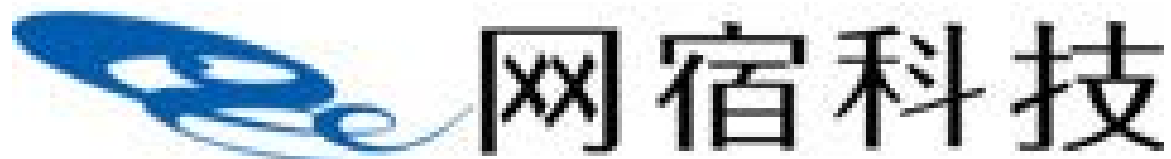


Apache RocketMQ

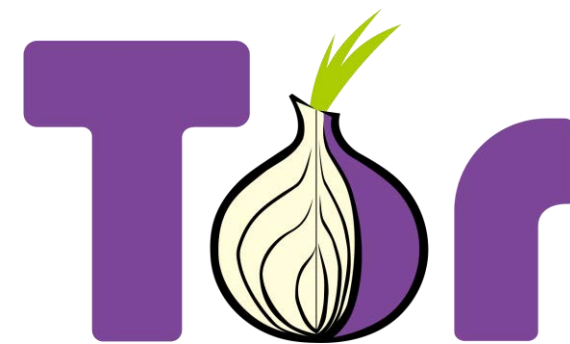


TARANTOOL

❖ 内容分发网络



❖ 分布式数据库



- ❖ 三大云服务平台
- ❖ Google Cloud: <https://cloud.google.com>
- ❖ Amazon Web Services: <https://aws.amazon.com>
- ❖ 阿里云: <https://www.aliyun.com>



Google Cloud



❖ 课外阅读

- ❖ 《云存储技术——分析与实践》，刘洋著，经济管理出版社
- ❖ <http://product.dangdang.com/24247525.html>
- ❖ 《Ahead in the Cloud》，Stephen Orban (GM of AWS)
- ❖ <https://www.amazon.com/Ahead-Cloud-Practices-Navigating-Enterprise/dp/1981924310/>
- ❖ 《Cloud Computing: Concepts, Technology & Architecture》，Thomas Erl
- ❖ <https://www.amazon.com/Cloud-Computing-Concepts-Technology-Architecture/dp/0133387526/>

Thanks!