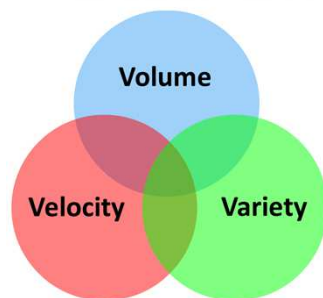


大数据的背景与趋势



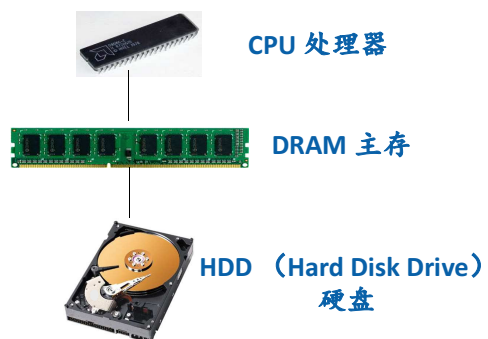
陈世敏

中科院计算所
计算机体系结构
国家重点实验室
©2015-2018 陈世敏

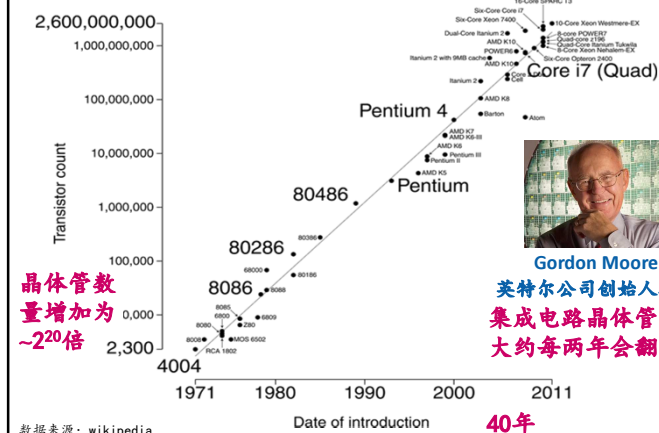
Outline

- 计算机硬件的发展
- 数据管理系统的发展
- 大数据的挑战
- 大数据管理系统

80年代的计算机系统



摩尔定律 (Moore's Law)

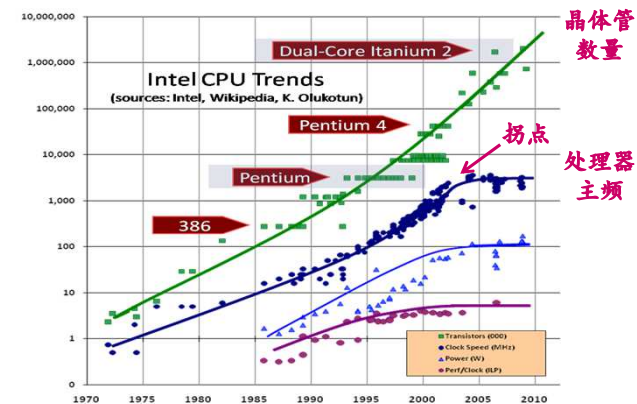


CPU体系结构的发展(2004年前)

- 提高串行程序效率

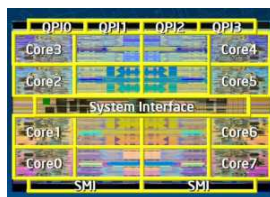
- 提高主频
- 流水线 (Pipeline)
- 超标量 (Super-Scalar)
- 乱序执行 (Out-of-order Execution)
- 向量指令 (SIMD/Vector Instructions)
- 多级高速缓存 (Multi-level Cache)

主频增加这一趋势于2004-2005年间结束



单核单线→多核多线→众核

- 功耗、散热等限制处理器主频的进一步增加
- 业界不得不转向使用多核
 - 从主频增加→核的增加
 - 双核、4核、6核、8核...、22核 (Intel Broadwell)
- 研究片上网络、高速缓存等



Nehalem EX

多种类型的处理器

- GPU

- 大量的并行单元
 - Nvidia Tesla P100有3584个并行计算单元
- 从专用到通用: OpenGL, CUDA, OpenCL
- 图形卡 (Dedicated Graphics Card)
 - 插在外部总线(例如PCIe)上
 - 有独立的显卡内存
- 集成在处理器芯片中 (Integrated Graphics)
 - 使用部分主存
 - 例如: Intel HD graphics, AMD APU

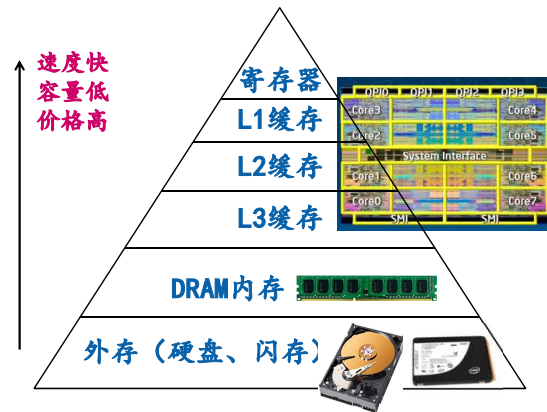


- Xeon Phi

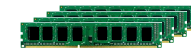
- ARM

- Dark Silicon(暗硅)→应用加速器

存储层次结构



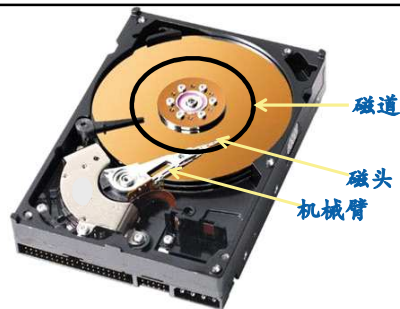
内存



- 容量 → 符合摩尔定律，指数级增加
- 带宽：有一定的办法增加
- 访问速度：比指令执行慢100倍
 - 访存墙问题

硬盘访问

- 硬盘容量：
呈指数级增长



- 硬盘的性能

- 访问速度：受限于机械臂的移动，盘片的转速
 - 大约每次访问需要10ms
- 带宽：受限于盘片的转速
 - 传输速度大约为100MB/s
- 顺序访问比随机访问好很多

闪存(Flash)与固态硬盘(SSD: Solid State Drive)

- 闪存

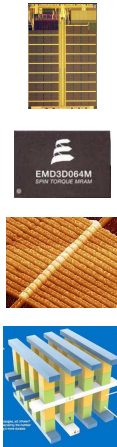
- 发明于1980年，与DRAM技术有一定的相似性
- 最早用于取代ROM作BIOS存储
- 后来用于数字电子设备：相机、手机、mp3、U盘、microSD卡等，大量生产，价格降低

- 固态硬盘

- 2009年开始出现以闪存为存储介质的固态硬盘
- 优点：没有机械装置，随机读性能比硬盘高100倍，顺序读或顺序写性能好于硬盘
- 缺点：随机写性能差，重写次数有限制(例如，5000次)，超过即报废

新型的非易失存储技术

- 集成电路特征尺寸已经接近极限
 - 当前的特征尺寸是14纳米
 - DRAM每个比特依靠存储电荷来区别0/1
 - 特征尺寸变小→存储电荷数变少→稳定性变差
- 业界在研发新型的存储技术来替代DRAM
 - Phase Change Memory, STT-RAM, Memristor
 - 3D-XPoint
- 目标是完美存储技术
 - 特征尺寸可以进一步减小
 - 与DRAM的读写速度相似, 支持的读写次数相似
 - 非易失: 不需要定时刷新, 节能, 可靠
- 对计算机系统的各方面都会产生深远的影响



体系结构和硬件技术的巨大发展



计算机系统的发展



小结

- 计算机硬件的发展
- 数据管理系统的发展
- 大数据的挑战
- 大数据管理系统

关系型数据库



- E.F. Codd于1970年提出了数据管理的关系模型，并因此于1981年获得图灵奖



- Jim Gray参与了第一个关系型数据库原型系统（System R）的实现，并由于对数据库和事务处理的多项贡献获得了1998年的图灵奖



- Michael Stonebraker主持了另一个早期关系型数据库原型系统（Ingress）的实现，并实现了一系列的系统（Postgres, C-Store, H-Store, 等），2015年获得图灵奖



- 三大数据库产品Oracle, Microsoft SQL Server, IBM DB2的最初实现都是在1970末到1980年代

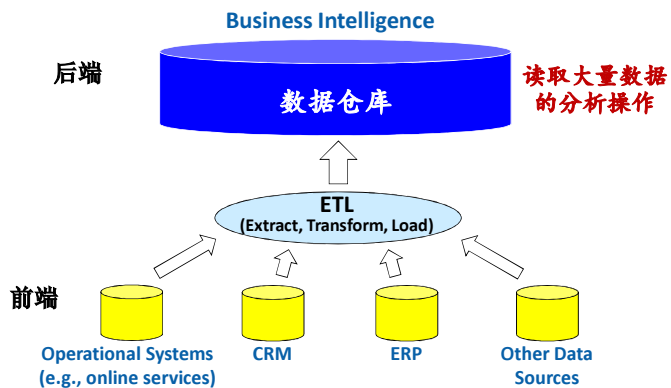
事务处理 (Transaction Processing)



事务处理系统

- 早期的数据库系统主要针对事务处理应用
- 典型例子：银行业务，订票，购物等
- 大量并发用户，少量随机读写操作

数据仓库 (Data Warehouse) 1990年代出现



多种发展

2000年代出现

- 数据流处理
- 地理信息系统
- 多媒体数据库
- 用于Web的后端
-

2010年代?

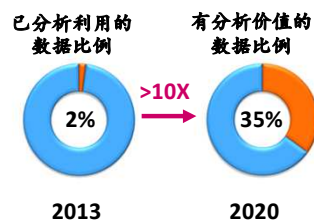
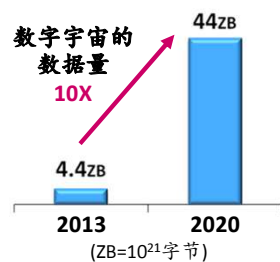
大数据

Outline

- 计算机硬件的发展
- 数据管理系统的发展
- 大数据的挑战
- 大数据管理系统

大数据分析的重要性

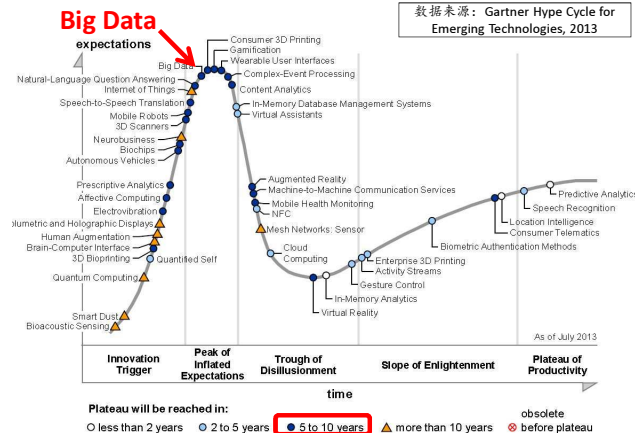
数据来源: EMC Digital Universe with Research & Analysis by IDC, 2014



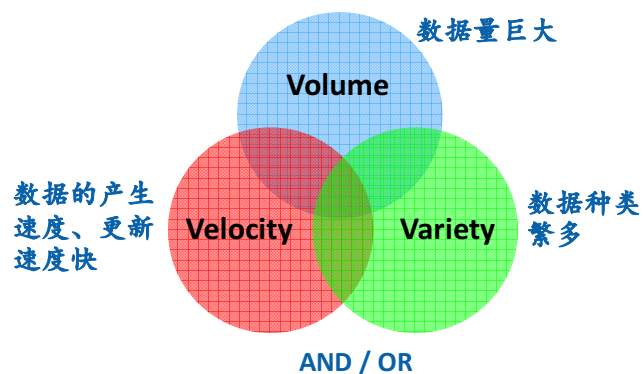
大数据分析将有**超过100倍**的巨大增长空间

大数据分析的重要性

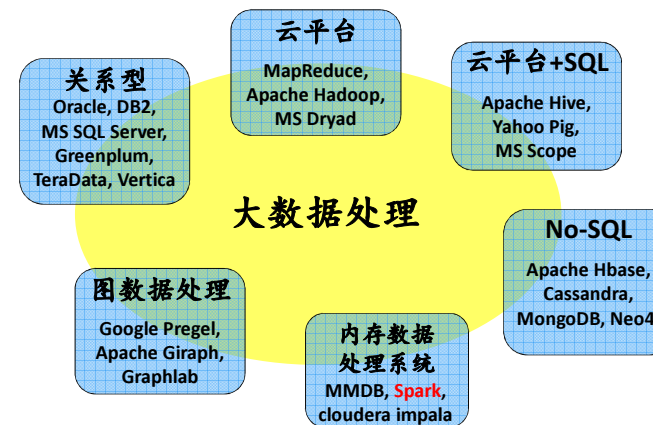
数据来源: Gartner Hype Cycle for Emerging Technologies, 2013



大数据的概念 (Big Data) 也是三个重要挑战



大数据管理系统



小结

- 计算机体系结构和硬件技术的巨大发展
 - 摩尔定律、存储系统的发展
 - 新的计算设备、云计算的出现
- 数据管理系统的发展
 - 1970-1980: 关系型系统出现, 事务处理
 - 1990: 数据仓库开始流行
 - 2000: 数据流等多种发展
 - 2010: 大数据
- 大数据的挑战
 - Volume, Velocity, Variety
 - 一个方面超出了传统处理能力就是大数据