# 物理存储系统 & 数据存储结构
# Physical Storage Systems & Data Storage Structures

李文根/Wengen Li

Email: lwengen@tongji.edu.cn
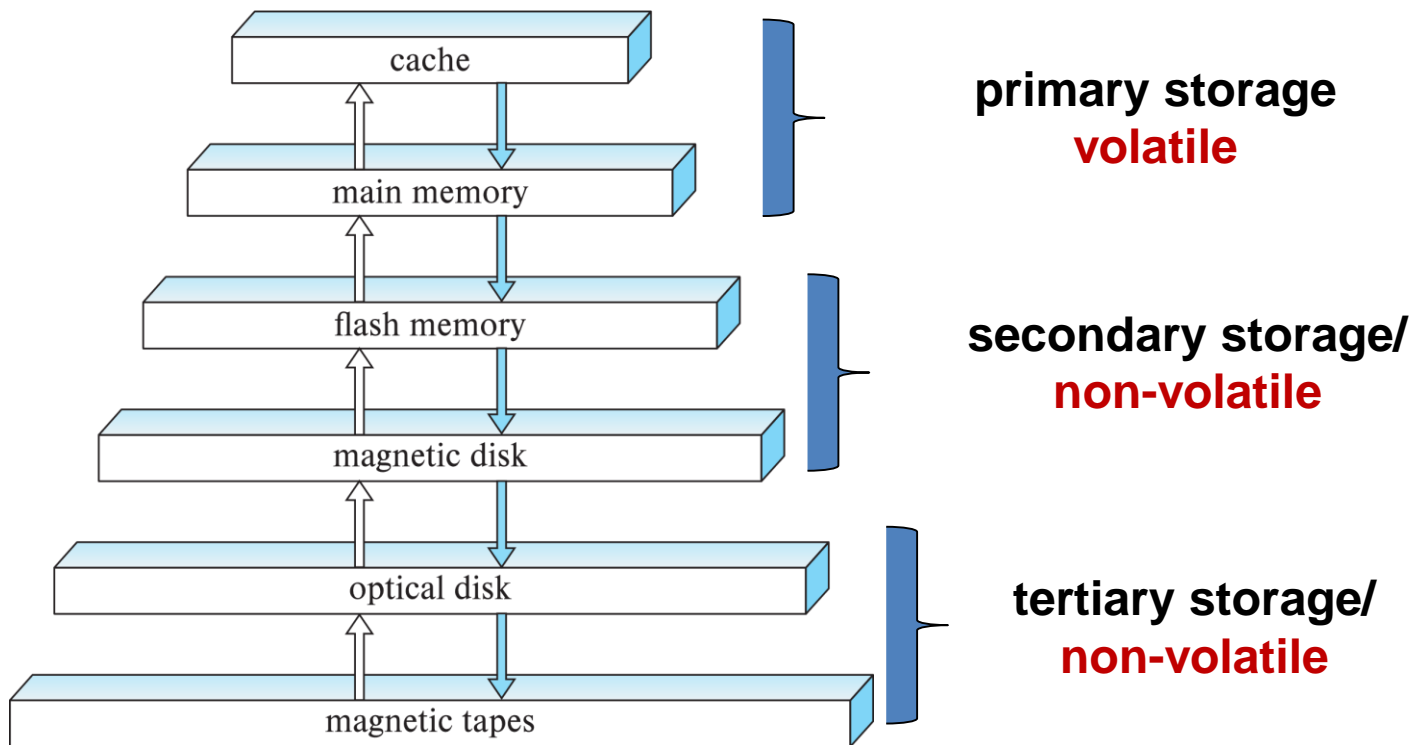
先进数据与机器智能系统实验室（ADMIS）

https://admis.tongji.edu.cn/main.htm

同济大学 电子与信息工程学院 计算机科学与技术系
TONGJI UNIVERSITY  College of Electronics and Information Engineering  Department of Computer Science and Technology

# 课程概要

# ▶ 目录

primary storage
volatile

secondary storage/
non-volatile

tertiary storage/
non-volatile

cache

main memory

flash memory

magnetic disk

optical disk

magnetic tapes

- **Speed**
  - the speed with which data can be accessed
- **Cost**
  - the cost for storing each unit of data
- **Reliability**
  - volatile storage (易失性存储)：lose contents when power is switched off
  - non-volatile storage (非易失性存储)：contents persist when power is switched off
    - secondary (二级) and tertiary（三级）storage
    - battery-backed up main-memory

- **Cache (高速缓存)**
  - the fastest and most costly form of storage
  - volatile

- **Main memory (主存/内存)**
  - fast access (about 100 nanoseconds, 1 nanosecond = $10^{-9}$ seconds)
  - generally too small (or too expensive) to store the entire database
    - capacities of up to tens of gigabytes widely used currently
    - capacities have gone up and per-byte cost has decreased steadily
  - **volatile** — contents of main memory will be lost if a power failure or system crash occurs

# ▶ 闪存

- **Flash memory (闪存)**
  - **USB**, **SSD**
  - Non-volatile, data survives power failure
  - Data can be written at a location only once, but the location can be erased and written again
    - support 10K – 1M write/erase cycles
  - Reads are roughly as fast as main memory, but writes are slow (few microseconds)
  - Cost per unit of storage is similar to main memory
  - Widely used in embedded devices such as digital cameras

# Magnetic disk (磁盘)

- Primary medium for the long-term storage of data
  - Typically stores the entire database
- Data must be moved from disk to main memory for access, and written back for storage
  - Much slower access than main memory
- Capacities range up to several TBs currently
  - Much larger capacity and lower cost/byte than main memory/flash memory
  - Growing constantly and rapidly with technology improvements
- Survives power failures and system crashes
  - disk failure can destroy data, but is rare

# 光学存储器

- **Optical storage (光学存储器)**
  - Non-volatile, data is read optically from a spinning disk using a laser
  - CD-ROM (640 MB) and DVD (4.7 to 17 GB) are the most popular forms
  - Write-once and read-many (WORM) optical disks used for archival storage
    - CD-R, DVD-R, DVD+R
  - Multiple write versions are also available
    - CD-RW, DVD-RW, DVD+RW, DVD-RAM
  - Reads and writes are slower than magnetic disk

# Optical storage (光学存储器)

- 2010年5月，日本东京大学化学教授大越慎一(Shin-ichi Ohkoshi)研究团队发现一种材料，可以用来制造更便宜的大容量超级光盘，可储存容量是目前**一般DVD的5千倍**

- 材料是一种透明的新型氧化钛，平常是能导电的黑色金属状态，在受到光的点击后会转变成棕色的半导体。在室温下受到光的照射，能够任意在金属和半导体之间转变，因而产生储存数据功能。这种材料所制成的新光碟，容量是蓝光光盘的1千倍，而蓝光光盘的容量则是一般DVD的5倍。一般一张DVD为4.7G，一张蓝光光盘为25G。最新的蓝光协会表示，新蓝光光盘容量可达128G。东京大学的最新光盘容量将达25000G，即25T

- 2012年富士胶片开发利用双光子吸收热量的新型光盘记录方式，可实现每层25GB的记录密度，与蓝光光盘相同，且该技术有可能实现多达20层的多层化。富士新技术单张盘片容量可达15TB

# ▶ 磁带存储器

- **Tape storage (磁带存储器)**
  - **Non-volatile**, used primarily for backup (to recover from disk failure), and for archival data
    - IBM、EMC、Dell…
  - **Sequential-access** – much slower than magnetic disk and flash memory (direct access)
  - **Very high capacity** (1TB to tens of TB tapes are available)
  - Can be **removed** from drive (磁带驱动器)
  - Storage cost is much **cheaper** than disk, but the drive is **expensive**

# ▶ 磁带存储器（续）

- 磁带是一种古老的数据存储方式。人类需要保存的很多数据都是冷数据，是不常用的数据，大部分是为了备份、保存，需求量在不断增长，磁带的优势就正是存储容量超大、成本低廉、保存时间长等。

- 富士LTO系列磁带，Ultrium 8系列磁带采用BaFe（钡铁氧体磁性材料），并有富士专利的纳米超薄涂层技术，磁带长度960m，宽度12.65mm，厚度5.6um。LTO Ultrium 8系列磁带容量高达30TB（未压缩时是12TB），是前代产品的两倍，而且速度可达750MB/s（未压缩时是360MB/s），性能及容量都要比HDD硬盘要有优势，适合长期保存重要数据。

- LTO 8系列磁带有两种类型，一种是可以重复擦写数据的，另外一种是WROM，写入一次，多次读取类型的，这种类型的可以防止数据被篡改或者意外删除，提高了安全性。

# 存储设备的层次结构

- **Primary storage (主存储器)**
  - Fastest media but volatile (cache, main memory)
- **Secondary storage (辅助存储器)**
  - Non-volatile, moderately fast access time
  - Also called on-line storage, e.g., flash memory, magnetic disks
- **Tertiary storage (三级存储器)**
  - Non-volatile, slow access time
  - Also called off-line storage, e.g., magnetic tape, optical storage

cache

main memory

flash memory

magnetic disk

optical disk

magnetic tapes

# ▶ 目录

- **物理存储介质概述**
- **磁盘**
- **闪存**
- **RAID**
- **三级存储**
- **存储访问**
- **文件组织**
- **文件中记录的组织**
- **数据字典存储**

spindle: 主轴

arm assembly: 磁盘臂组件

arm: 磁盘臂

cylinder: 柱面

platter: 盘片

track: 磁道

sector: 扇区

read-write head: 读写头



**Schematic diagram of magnetic disk**



**Photo of magnetic disk drive**

- **Read-write head (读写头)**
  - positioned very close to the platter surface (almost touching it)
  - reads or writes magnetically encoded information
- **Surface of platter is divided into circular tracks (磁道)**
  - about 50K-100K tracks per platter on typical hard disks
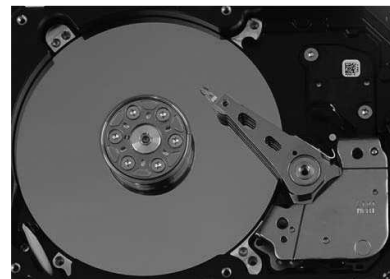- **Each track is divided into sectors (扇区)**
  - sector is the smallest unit of data that can be read or written, typical sector size: 512 bytes
  - typical sectors per track: 500-1000 (inner tracks), 1000-2000 (outer tracks)
- **To read/write a sector**
  - disk arm swings to position head on the right track
  - platter spins continually, and data is read/written as sector passes under head
- **Head-disk assemblies (磁头-磁盘装置)**
  - multiple disk platters on a single spindle (1 to 5 usually)
  - one head per platter, mounted on a common arm assembly.
- **Cylinder(柱面) $i$ consists of the $i$-th tracks of all the platters(盘片)**

- **Earlier generation disks were susceptible to head-crashes**
  - Surface of earlier generation disks had metal-oxide coatings (金属氧化物涂层) which would disintegrate on head crash and damage all data on disk
  - Current disks are less susceptible to such disastrous failures
- **Disk controller**(磁盘控制器) – interfaces between the computer system and the disk drive hardware
  - Accept high-level commands to read or write a sector
  - Initiate actions such as moving the disk arm to the right track and actually reading or writing the data
  - Compute and attach checksums(校验和) to each sector to verify that data is read back correctly.
  - Performs remapping of bad sectors(坏扇区重映射)

- **Multiple disks connected to a computer system through a disk controller**
  - Controller functionality (checksum, bad sector remapping) is often carried out by individual disks to reduce the load on controller



- **Disk interface standards families**
  - ATA (AT adaptor/attachment) range of standards (1994-2002,7 standards)
  - **SATA** (Serial ATA)、PATA (Parallel ATA)
  - SCSI (Small Computer System Interconnect) range of standards
  - Several variants of each standard (different speeds and capabilities)

18

- **Access time(访问时间)**
  - the time it takes from when a read or write request is issued to when data transfer begins, including
    - **Seek time(寻道时间)** – the time it takes to reposition the arm over the correct track.
      - average seek time is about 1/2 the worst case seek time
      - 4 to 10 milliseconds on typical disks
    - **Rotational latency(旋转延迟)** – the time it takes for the sector to be accessed to appear under the head
      - average latency is 1/2 of the worst case latency.
      - 4 to 11 milliseconds on typical disks (5400 to 15000 rpm)

# 磁盘的性能度量（续）

- **Data-transfer rate**(数据传输率)
  - The rate at which data can be retrieved from or stored to the disk
  - Max rate: 50 to 200 MB per second, lower for inner tracks
  - Multiple disks may share a controller, so the rate that controller can handle is also important
    - E.g., ATA-5: 66 MB/sec, SATA: 150 MB/sec, Ultra 320 SCSI: 320 MB/s
    - Fiber Channel (FC2Gb): 256 MB/s

# 磁盘的性能度量（续）

- **Mean time to failure** (平均故障时间, MTTF)
  - The average time the disk is expected to run continuously without any failure
  - Probability of failure of new disks is quite low, corresponding to a "theoretical MTTF" of 500,000 to 1,200,000 hours for a new disk
    - E.g., an MTTF of 1,200,000 hours for a new disk means that given 1000 new disks, one will fail every 1200 hours on the average
  - MTTF decreases as disk ages

# ▶ 磁盘块访问的优化

- **Block** – a contiguous sequence of sectors from a single track
  - data is transferred between disk and main memory in blocks
  - sizes range from 512 bytes to several kilobytes
    - small blocks: more transfers from disk
    - large blocks:  more space wasted due to partially filled blocks
    - typical block sizes range from **4 to 16 kilobytes**
- **Disk-arm-scheduling**(磁盘臂调度) algorithms order pending accesses to tracks so that disk arm movement is minimized
  - **elevator algorithm(电梯算法)** : move disk arm in one direction (from outer to inner tracks or vice versa), processing next request in that direction, till no more requests in that direction, then reverse direction and repeat

- **File organization** – optimize block access time by organizing the blocks according to the way of data access
  - E.g., store related information on the same or nearby cylinders
  - Files may get fragmented(碎片化) over time
    - data is inserted to/deleted from the file
    - free blocks on disk are scattered, and newly created file has its blocks scattered over the disk
    - sequential access to a fragmented file results in increased disk arm movement
  - To speed up file access, some systems have utilities to defragment the file system

- **Non-volatile write buffers** speed up disk writes by writing blocks to a non-volatile RAM buffer immediately
  - Non-volatile RAM:  battery backed up RAM or flash memory
    - Even if power fails, the data is safe and will be written to disk when power returns
  - Controller writes to disk when the disk has no other requests or the non-volatile RAM is full
  - Database operations can continue without waiting for data to be written to disk
  - Writes can be reordered to minimize disk arm movement

- **Log disk(日志磁盘)**
  - a disk devoted to writing a sequential log of block updates
  - used exactly like non-volatile RAM
    - write to log disk is very fast since seeking is not required
    - no need for special hardware (NV-RAM)
- **File systems typically reorder writes to disk to improve performance**
  - Journaling file systems(日志文件系统) write data in safe order to NV-RAM or log disk
  - Reordering without journaling: risk of corruption of file system data

# ▶ 目录

- **物理存储介质概述**
- **磁盘**
- **闪存**
- **RAID**
- **三级存储**
- **存储访问**
- **文件组织**
- **文件中记录的组织**
- **数据字典存储**

# 闪存（Flash Storage）

- **NOR flash vs NAND flash**
  - **NAND flash**
    - used widely for storage, cheaper than NOR flash
    - requires page-at-a-time read (page: 512 bytes to 4 KB)
      - 20 to 100 microseconds for a page read
      - Not much difference between sequential and random read
    - Page can only be written once
      - Must be erased to allow rewrite
- **Solid state disks (SSD，固态磁盘)**
  - Use standard block-oriented disk interfaces
  - Transfer rate of up to 500 MB/sec using SATA, and up to 3 GB/sec using NVMe PCIe

# ▶ 闪存（续）

- Erase happens in units of erase block
  - Takes 2 to 5 ms
  - Erase block typically 256 KB to 1 MB (128 to 256 pages)
- Remapping of logical page addresses to physical page addresses avoids waiting for erase
  - Flash translation table (转换表) tracks mapping, also stored in a label field of flash page
- After 100,000 to 1,000,000 erases, erase block becomes unreliable and cannot be used
  - wear leveling (损耗均衡)

# SSD的性能指标

- **Random reads/writes per second**
  - Typical 4 KB reads: 10,000 reads per second (10,000 IOPS)
  - Typical 4KB writes: 40,000 IOPS
  - SSDs support parallel reads
    - Typical 4KB reads:
      - 100,000 IOPS with 32 requests in parallel (QD-32) on SATA
      - 350,000 IOPS with QD-32 on NVMe PCIe
    - Typical 4KB writes:
      - 100,000 IOPS with QD-32, even higher on some models
- **Data transfer rate for sequential reads/writes**
  - 400 MB/sec for SATA3, 2 to 3 GB/sec using NVMe PCIe
- **Hybrid disks**: combine small amount of flash cache with larger magnetic disk

# 存储类存储器

- **Storage class memory (存储类存储器)**
  - Allow to read and write bytes or words
- 3D-XPoint memory technology pioneered by Intel
- Available as Intel Optane SSD
  - SSD interface shipped from 2017
    - Allows lower latency than flash SSDs
  - Non-volatile memory interface announced in 2018
    - Supports direct access to words, at speeds comparable to main-memory speeds

# ▶ 目录

- **物理存储介质概述**
- **磁盘**
- **闪存**
- **RAID**
- **三级存储**
- **存储访问**
- **文件组织**
- **文件中记录的组织**
- **数据字典存储**

# 独立磁盘冗余阵列（RAID）

- **RAID: Redundant Arrays of Independent Disks**

  - Disk organization techniques that manage a large numbers of disks, providing a view of a single disk of

    - High capacity and high speed by using multiple disks in parallel

    - High reliability by storing data redundantly. Data can be recovered even if a disk fails

# ▶ RAID （续）

- The chance that some disk out of a set of N disks will fail is much higher than the chance that a specific single disk will fail
  - E.g., a system with 100 disks, each with MTTF of 100,000 hours (approx. 11 years), will have a system MTTF of 1000 hours (around 41 days)
  - Techniques for using redundancy to avoid data loss are critical with large numbers of disks

- **Originally a cost-effective alternative to large and expensive disks**
  - I in RAID originally stood for "inexpensive"
  - Now the "I" is interpreted as independent

- Redundancy – store extra information that can be used to rebuild information lost in a disk failure
  - Mirroring (镜像): Duplicate every disk, and each logical disk consists of two physical disks
  - Every write is carried out on both disks
    - Reads can take place from either disk
  - If one disk in a pair fails, data is still available in the other
    - Data loss would occur only if a disk fails, and its mirror disk also fails before the system is repaired
    - Probability of combined event is very small except for dependent failure modes such as fire or building collapse or electrical power surges
- Mean time to data loss (平均数据丢失时间) depends on mean time to failure (平均故障时间) and mean time to repair (平均修复时间)

# 通过并行提高性能

- Two main goals of **parallelism** in a disk system
  - Load balance multiple small accesses to increase throughput
  - Parallelize large accesses to reduce response time
- Bit-level striping (位级拆分): split the bits of each byte across multiple disks
  - In an array of eight disks, write bit $i$ of each byte to disk $i$
  - Each access can read data at eight times the rate of a single disk
- Block-level striping (块级拆分): with $n$ disks, block $i$ of a file goes to disk ($i$ mod $n$) + 1
  - Requests for different blocks can run in parallel if the blocks reside on different disks
  - A request for a long sequence of blocks can utilize all disks in parallel

- Schemes to provide redundancy at lower cost by using disk striping (磁盘拆分) combined with parity bits(奇偶校验位)
  - Different RAID levels, have different cost, performance and reliability characteristics
  - **RAID Level 0: Block striping, non-redundant. (无冗余拆分)**
    - ✓ **Used in high-performance applications where data lose is not critical**
  - **RAID Level 1: Mirrored disks with block striping (镜像磁盘)**
    - ✓ **Popular for applications such as storing log files in a database system**

(a) RAID 0: nonredundant striping

(b) RAID 1: mirrored disks

- **RAID Level 2**: Memory-Style Error-Correcting-Codes (ECC) with bit striping. (内存风格的纠错码, 奇偶校验位)
- **RAID Level 3**: Bit-Interleaved Parity (位交叉的奇偶校验)
  - a single parity bit is enough for error correction, not just detection, since we know which disk has failed
    - When writing data, corresponding parity bits must also be computed and written to a parity bit disk
    - To recover data in a damaged disk, compute XOR of bits from other disks (including parity bit disk)

(c) RAID 2: memory-style error-correcting codes

(d) RAID 3: bit-interleaved parity

- **RAID Level 3 (Cont.)**
  - Faster data transfer than with a single disk, but fewer I/Os per second since every disk has to participate in every I/O.
  - Subsumes Level 2 (provides all its benefits, at lower cost).
- **RAID Level 4**
  - Block-Interleaved Parity(块交叉的奇偶校验); uses block-level striping, and keeps a parity block on a separate disk for corresponding blocks from N other disks.
  - When writing data block, corresponding block of parity bits must also be computed and written to parity disk
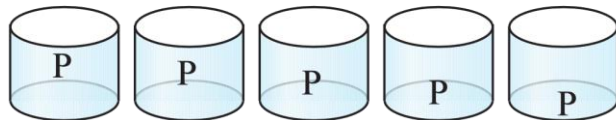  - To find value of a damaged block, compute XOR of bits from corresponding blocks (including parity block) from other disks.



(e) RAID 4: block-interleaved parity

- **RAID Level 4 (Cont.)**
  - Provides higher I/O rates for independent block reads than Level 3
    - block read goes to a single disk, so blocks stored on different disks can be read in parallel
  - Provides high transfer rates for reads of multiple blocks than no-striping
  - Before writing a block, parity data must be computed
    - Can be done by using old parity block, old value of current block and new value of current block (2 block reads + 2 block writes)
      - More efficient for writing large amounts of data sequentially
  - Parity block becomes a bottleneck for independent block writes since every block write also writes to parity disk

- **RAID Level 5**
  - Block-Interleaved Distributed Parity(块交叉的分布奇偶校验); partitions data and parity among all N + 1 disks, rather than storing data in N disks and parity in 1 disk.
  - E.g., with 5 disks, parity block for n-th set of blocks is stored on disk (n mod 5) + 1, with the data blocks stored on the other 4 disks.

RAID 5: block-interleaved distributed parity

| P0 | 0 | 1 | 2 | 3 |
|----|-----|-----|-----|-----|
| 4 | P1 | 5 | 6 | 7 |
| 8 | 9 | P2 | 10 | 11 |
| 12 | 13 | 14 | P3 | 15 |
| 16 | 17 | 18 | 19 | P4 |

- **RAID Level 5 (Cont.)**
  - Higher I/O rates than Level 4.
    - Block writes occur in parallel if the blocks and their parity blocks are on different disks.
  - Subsumes Level 4: provides same benefits, but avoids bottleneck of parity disk.
- **RAID Level 6**
  - P+Q Redundancy scheme; similar to Level 5, but stores extra redundant information to guard against multiple disk failures.
  - Better reliability than Level 5 at a higher cost; not used as widely.

(d) RAID 6: P + Q redundancy

# ▶ RAID级别的选择

- Factors in choosing RAID level
  - Monetary cost
  - Performance: number of I/O operations per second, and bandwidth during normal operation
  - Performance during failure
  - Performance during rebuild of failed disk
    - Including time taken to rebuild failed disk
- RAID 0 is used only when data safety is not important
  - E.g. data can be recovered quickly from other sources
- Level 2 and 4 never used since they are subsumed by 3 and 5
- Level 3 is not used anymore since bit-striping forces single block reads to access all disks, wasting disk arm movement, which block striping (level 5) avoids
- Level 6 is rarely used since levels 1 and 5 offer adequate safety for almost all applications
- So competition is between **1 and 5** only

# ▶ RAID级别的选择（续）

- Level 1 provides much better write performance than level 5
  - Level 5 requires at least 2 block reads and 2 block writes to write a single block, whereas Level 1 only requires 2 block writes
  - Level 1 preferred for high update environments such as log disks

- Level 1 had higher storage cost than level 5
  - disk drive capacities increasing rapidly (50%/year) whereas disk access times have decreased much less (x 3 in 10 years)
  - I/O requirements have increased greatly, e.g. for Web servers
  - When enough disks have been bought to satisfy required rate of I/O, they often have spare storage capacity
    - so there is often no extra monetary cost for Level 1!

- Level 5 is preferred for applications with low update rate, and large amounts of data

- Level 1 is preferred for all other applications

- **Software RAID**:  RAID implementations done entirely in software, with no special hardware support
- **Hardware RAID**:  RAID implementations with special hardware
  – Use non-volatile RAM to record writes that are being executed
  – Beware:  power failure during write can result in corrupted disk
    - E.g., failure after writing one block but before writing the second in a mirrored system
    - Such corrupted data must be detected when power is restored
      – Recovery from corruption is similar to recovery from failed disk
      – NV-RAM helps to efficiently detect potentially corrupted blocks
        » Otherwise all blocks of disk must be read and compared with mirror/parity block

- **Hot swapping(热交换):** replacement of disk while system is running, without power down
  - Supported by some hardware RAID systems
  - reduces time to recovery, and improves availability (可用性) greatly
- Many systems maintain spare disks which are kept online, and used as replacements for failed disks immediately on detection of failure
  - Reduces time to recovery greatly
- Many hardware RAID systems ensure that a single point of failure will not stop the functioning of the system by using
  - Redundant power supplies with battery backup
  - Multiple controllers and multiple interconnections to guard against controller/interconnection failures

# ▶ 目录

- **物理存储介质概述**
- **磁盘**
- **闪存**
- **RAID**
- **三级存储**
- **存储访问**
- **文件组织**
- **文件中记录的组织**
- **数据字典存储**

# 光学存储器

- Compact disk-read only memory (CD-ROM)
  - Removable disks, 640 MB per disk
  - Seek time about 100 msec (optical read head is heavier and slower)
  - Higher latency (3000 RPM) and lower data-transfer rates (3-6 MB/s) compared to magnetic disks
- Digital Video Disk (DVD)
  - DVD-5  holds 4.7 GB , and DVD-9 holds 8.5 GB
  - DVD-10 and DVD-18 are double sided formats with capacities of 9.4 GB and 17 GB
  - Slow seek time, for same reasons as CD-ROM
- Record once versions (CD-R and DVD-R) are popular
  - Data can only be written once, and cannot be erased
  - High capacity and long lifetime; used for archival storage
- Multi-write versions (CD-RW, DVD-RW, DVD+RW and DVD-RAM) also available

# 磁带（Magnetic Tape）

- Hold large volumes of data and provide high transfer rates
  - Few GB for DAT (Digital Audio Tape) format
  - 10-40 GB with DLT (Digital Linear Tape) format
  - 100 GB+ with Ultrium format
  - 330 GB with Ampex helical scan format
  - Transfer rates from few to 10s of MB/s
- **Currently the cheapest storage medium**
  - Tapes are cheap, but cost of drives is very high

- Very slow access time in comparison to magnetic disks and optical disks
  - Limited to sequential access.
  - Some formats provide faster seek (10s of seconds) at cost of lower capacity
- Used mainly for backup, for storage of infrequently used information, and as an off-line medium for transferring information from one system to another.
- Tape jukeboxes(自动磁带机) used for very large capacity storage
  - terabyte ($10^{12}$ bytes) to petabye ($10^{15}$ bytes)

- **物理存储介质概述**
- **磁盘**
- **闪存**
- **RAID**
- **三级存储**
- **存储访问**
- **文件组织**
- **文件中记录的组织**
- **数据字典存储**

- **Storage of database file**
  - Partitioned into fixed-length storage units called **blocks**. Blocks are units of both storage allocation and data transfer
  - Database system seeks to minimize the number of block transfers between the disk and memory. We can reduce the number of disk accesses by keeping as many blocks as possible in main memory
- **Buffer**
  - The portion of main memory available to store copies of disk blocks
  - **Buffer manager:** responsible for allocating buffer space in main memory

# 缓冲区管理器（Buffer Manager）

- Programs call the buffer manager when they need a block from disk
  - If the block is already in the buffer, buffer manager returns the address of the block in main memory
  - If the block is not in the buffer, the buffer manager
    - allocates space in the buffer for the block
      - Replacing (throwing out) some other block, if required, to make space for the new block.
      - Replaced block is written back to disk only if it was modified since the most recent time that it was written to/fetched from the disk
    - reads the block from the disk to the buffer, and returns the address of the block in main memory to requester

- **LRU (Lest Recently Used，最近最少使用)**
  - Most operating systems replace the block **least recently used**
  - Idea behind LRU – use past patterns of block references as a predictor of future references
  - LRU can be a bad strategy for certain access patterns involving repeated scans of data
    - E.g., computing the join of two relations r and s by a nested loops

      for each tuple $t_r$ of r do
        for each tuple $t_s$ of s do
          if the tuples $t_r$ and $t_s$ match

            …

# 缓冲区替换策略（续）

- **Pinned block(被钉住的块)**
  - memory block that is not allowed to be written back to disk
- **Toss-immediate (立即丢弃) strategy**
  - free the space occupied by a block as soon as the final tuple of that block has been processed
- **Most recently used (MRU) (最近最常使用) strategy**
  - system must pin the block currently being processed. After the final tuple of that block has been processed, the block is unpinned, and becomes the most recently used block

- **物理存储介质概述**
- **磁盘**
- **闪存**
- **RAID**
- **三级存储**
- **存储访问**
- **文件组织**
- **文件中记录的组织**
- **数据字典存储**

# ▶ 文件组织

- The database is stored as a collection of files. Each file is a sequence of records. A record is a sequence of fields

- **One approach**
  - Assume that the record size is fixed
  - Each file has records of one particular type only
  - Different files are used for different relations

  **Note**: this case is easy to implement. We will consider variable length records later

# 定长记录

- **Simple approach**
  - Store record $i$ starting from byte $n \cdot (i - 1)$, where $n$ is the size of each record
  - Record access is simple but records may cross blocks
    - Modification: don't allow records to cross block boundaries
- **Alternative methods for deleting record $i$**
  - move records $i + 1, \ldots, n$ to $i, \ldots, n - 1$
  - move record $n$ to $i$
  - do not move records, but link all free records on a free list

- Delete record 3 and move all records

| | | | | |
|---|---|---|---|---|
| record 0 | 10101 | Srinivasan | Comp. Sci. | 65000 |
| record 1 | 12121 | Wu | Finance | 90000 |
| record 2 | 15151 | Mozart | Music | 40000 |
| record 3 | 22222 | Einstein | Physics | 95000 |
| record 4 | 32343 | El Said | History | 60000 |
| record 5 | 33456 | Gold | Physics | 87000 |
| record 6 | 45565 | Katz | Comp. Sci. | 75000 |
| record 7 | 58583 | Califieri | History | 62000 |
| record 8 | 76543 | Singh | Finance | 80000 |
| record 9 | 76766 | Crick | Biology | 72000 |
| record 10 | 83821 | Brandt | Comp. Sci. | 92000 |
| record 11 | 98345 | Kim | Elec. Eng. | 80000 |

| | | | | |
|---|---|---|---|---|
| record 0 | 10101 | Srinivasan | Comp. Sci. | 65000 |
| record 1 | 12121 | Wu | Finance | 90000 |
| record 2 | 15151 | Mozart | Music | 40000 |
| record 4 | 32343 | El Said | History | 60000 |
| record 5 | 33456 | Gold | Physics | 87000 |
| record 6 | 45565 | Katz | Comp. Sci. | 75000 |
| record 7 | 58583 | Califieri | History | 62000 |
| record 8 | 76543 | Singh | Finance | 80000 |
| record 9 | 76766 | Crick | Biology | 72000 |
| record 10 | 83821 | Brandt | Comp. Sci. | 92000 |
| record 11 | 98345 | Kim | Elec. Eng. | 80000 |

- Delete record 3 and move the final record

| record 0 | 10101 | Srinivasan | Comp. Sci. | 65000 |
|----------|-------|------------|------------|-------|
| record 1 | 12121 | Wu | Finance | 90000 |
| record 2 | 15151 | Mozart | Music | 40000 |
| record 3 | 22222 | Einstein | Physics | 95000 |
| record 4 | 32343 | El Said | History | 60000 |
| record 5 | 33456 | Gold | Physics | 87000 |
| record 6 | 45565 | Katz | Comp. Sci. | 75000 |
| record 7 | 58583 | Califieri | History | 62000 |
| record 8 | 76543 | Singh | Finance | 80000 |
| record 9 | 76766 | Crick | Biology | 72000 |
| record 10 | 83821 | Brandt | Comp. Sci. | 92000 |
| record 11 | 98345 | Kim | Elec. Eng. | 80000 |

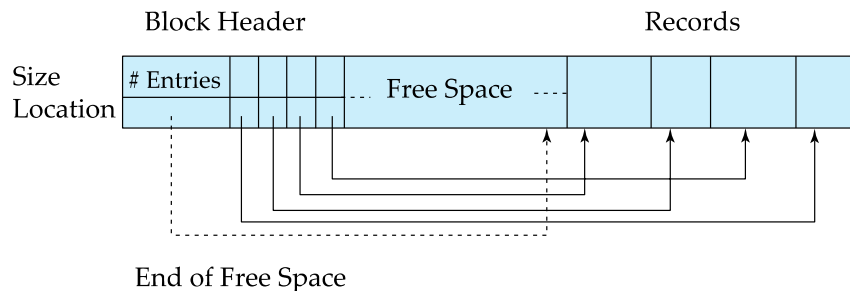| record 0 | 10101 | Srinivasan | Comp. Sci. | 65000 |
|----------|-------|------------|------------|-------|
| record 1 | 12121 | Wu | Finance | 90000 |
| record 2 | 15151 | Mozart | Music | 40000 |
| record 11 | 98345 | Kim | Elec. Eng. | 80000 |
| record 4 | 32343 | El Said | History | 60000 |
| record 5 | 33456 | Gold | Physics | 87000 |
| record 6 | 45565 | Katz | Comp. Sci. | 75000 |
| record 7 | 58583 | Califieri | History | 62000 |
| record 8 | 76543 | Singh | Finance | 80000 |
| record 9 | 76766 | Crick | Biology | 72000 |
| record 10 | 83821 | Brandt | Comp. Sci. | 92000 |

# 空闲列表 (Free List)

- Store the address of the first deleted record in the file header, and use the first record to store the address of the second deleted record, and so on

- These stored addresses are pointers since they "point" to the location of a record

| | | | | |
|---|---|---|---|---|
| header | | | | |
| record 0 | 10101 | Srinivasan | Comp. Sci. | 65000 |
| record 1 | | | | |
| record 2 | 15151 | Mozart | Music | 40000 |
| record 3 | 22222 | Einstein | Physics | 95000 |
| record 4 | | | | |
| record 5 | 33456 | Gold | Physics | 87000 |
| record 6 | | | | |
| record 7 | 58583 | Califieri | History | 62000 |
| record 8 | 76543 | Singh | Finance | 80000 |
| record 9 | 76766 | Crick | Biology | 72000 |
| record 10 | 83821 | Brandt | Comp. Sci. | 92000 |
| record 11 | 98345 | Kim | Elec. Eng. | 80000 |

- **Variable-length records**
  - Storage of multiple record types in a file
  - Record types that allow variable lengths for one or more fields
  - Record types that allow repeating fields, e.g., array and multiset (used in some old data models)

Block Header — Records

Size / Location | # Entries | Free Space

End of Free Space

- **Slotted page** (分槽的页) header contains
  - number of record entries
  - end of free space in the block
  - location and size of each record
- Records can be moved around within a page to keep them contiguous with no empty space between them; entry in the header must be updated
- Pointers should not point directly to record — instead they should point to the entry for the record in header

# ▶ 目录

- **物理存储介质概述**
- **磁盘**
- **闪存**
- **RAID**
- **三级存储**
- **存储访问**
- **文件组织**
- **文件中记录的组织**
- **数据字典存储**

# ▶ 文件中记录的组织

- **Sequential** (顺序)
  - store records in sequential order, based on the value of the search key of each record
- **Heap** (堆)
  - a record can be placed anywhere in the file where there is space
- **Hashing** (散列)
  - a hash function computed on some attribute of each record
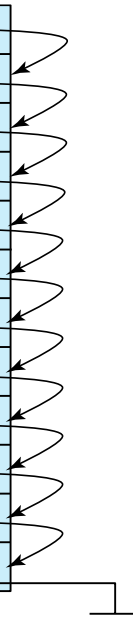  - the result specifies in which block of the file the record should be placed

- Suitable for applications that require sequential processing of the entire file

- The records in the file are ordered by a search-key

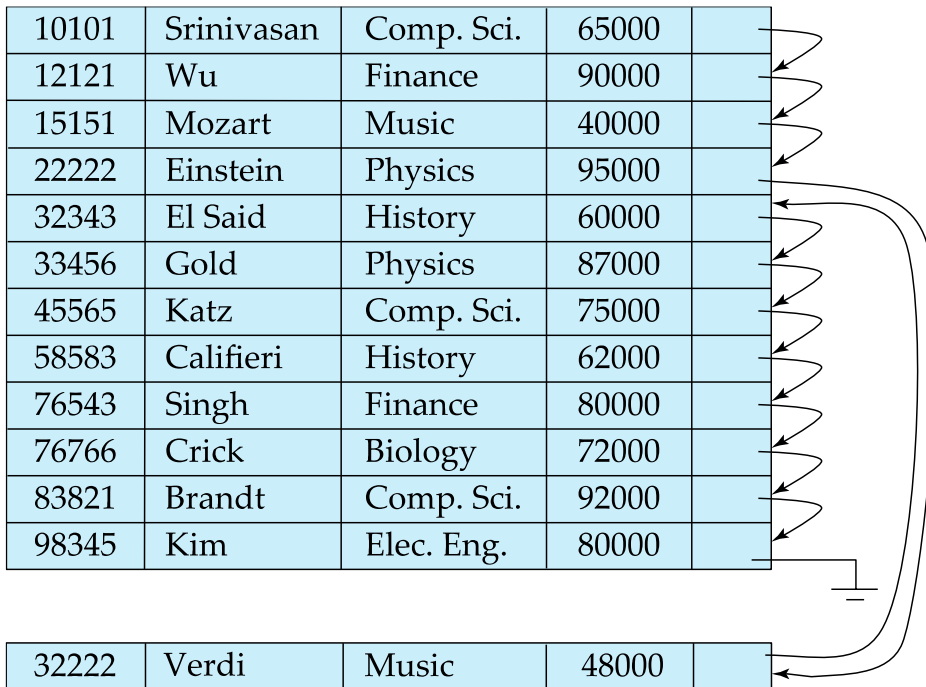| | | | | |
|---|---|---|---|---|
| 10101 | Srinivasan | Comp. Sci. | 65000 | |
| 12121 | Wu | Finance | 90000 | |
| 15151 | Mozart | Music | 40000 | |
| 22222 | Einstein | Physics | 95000 | |
| 32343 | El Said | History | 60000 | |
| 33456 | Gold | Physics | 87000 | |
| 45565 | Katz | Comp. Sci. | 75000 | |
| 58583 | Califieri | History | 62000 | |
| 76543 | Singh | Finance | 80000 | |
| 76766 | Crick | Biology | 72000 | |
| 83821 | Brandt | Comp. Sci. | 92000 | |
| 98345 | Kim | Elec. Eng. | 80000 | |

- Deletion – use pointer chains
- Insertion –locate the position where the record is to be inserted
  - if there is free space insert there
  - if no free space, insert the record in an overflow block
  - In either case, pointer chain must be updated
- Need to reorganize the file from time to time to restore sequential order

| 10101 | Srinivasan | Comp. Sci. | 65000 | |
| 12121 | Wu | Finance | 90000 | |
| 15151 | Mozart | Music | 40000 | |
| 22222 | Einstein | Physics | 95000 | |
| 32343 | El Said | History | 60000 | |
| 33456 | Gold | Physics | 87000 | |
| 45565 | Katz | Comp. Sci. | 75000 | |
| 58583 | Califieri | History | 62000 | |
| 76543 | Singh | Finance | 80000 | |
| 76766 | Crick | Biology | 72000 | |
| 83821 | Brandt | Comp. Sci. | 92000 | |
| 98345 | Kim | Elec. Eng. | 80000 | |

| 32222 | Verdi | Music | 48000 | |

- Store several relations in one file using a multi-table clustering file organization

*department*

| dept_name | building | budget |
|-----------|----------|--------|
| Comp. Sci. | Taylor | 100000 |
| Physics | Watson | 70000 |

*instructor*

| ID | name | dept_name | salary |
|-----|------|-----------|--------|
| 10101 | Srinivasan | Comp. Sci. | 65000 |
| 33456 | Gold | Physics | 87000 |
| 45565 | Katz | Comp. Sci. | 75000 |
| 83821 | Brandt | Comp. Sci. | 92000 |

Multi-table clustering of *department* and *instructor*

| Comp. Sci. | Taylor | 100000 | |
|-----------|--------|--------|------|
| 10101 | Srinivasan | Comp. Sci. | 65000 |
| 45565 | Katz | Comp. Sci. | 75000 |
| 83821 | Brandt | Comp. Sci. | 92000 |
| Physics | Watson | 70000 | |
| 33456 | Gold | Physics | 87000 |

# ▶ 多表聚簇文件组织（续）

- Good for queries involving *department* ⋈ *instructor*, and for queries involving one single department and its instructors
- Bad for queries involving only *department*
- Results in variable size records
- Can add pointer chains to link records of a particular relation

- **物理存储介质概述**
- **磁盘**
- **闪存**
- **RAID**
- **三级存储**
- **存储访问**
- **文件组织**
- **文件中记录的组织**
- <span style="color:red">**数据字典存储**</span>

- **Data dictionary** (also called system catalog) stores metadata, i.e., data about data, such as
- **Information about relations**
  - names of relations
  - names and types of attributes of each relation
  - names and definitions of views
  - integrity constraints
- **User and accounting information, including passwords**
- **Statistical and descriptive data**
  - the number of tuples in each relation
- **Physical file organization information**
  - How relation is stored (sequential/hash/…)
  - Physical location of relation
- **Information about indices (Chapter 14)**

- **Catalog structure**
  - relational representation on disk
  - specialized data structures designed for efficient access
- A possible catalog representation:

*relation_metadata* = (*relation_name, number_of_attributes, storage_organization, location*)
*attribute_metadata* = (*attribute_name, relation_name, domain_type, position, length*)
*user_metadata* = (*user_name, encrypted_password, group*)
*index_metadata* = (*index_name, relation_name, index_type, index_attributes*)
*view_metadata* = (*view_name, definition*)

# 数据库及存储技术（补充）

# ▶ 内存数据库

- **内存数据库**
  - 将数据放在内存中直接操作的数据库。相对于磁盘，内存的数据读写速度要高出几个数量级，将数据保存在内存中相比从磁盘上访问能够极大地提高应用的性能同时，内存数据库抛弃了磁盘数据管理的传统方式，基于全部数据都在内存中重新设计了体系结构，并且在数据缓存、快速算法、并行操作方面也进行了相应的改进，所以数据处理速度比传统数据库的数据处理速度要快很多，一般都在10倍以上。内存数据库的最大特点是其"主拷贝"或"工作版本"常驻内存，即活动事务只与实时内存数据库的内存拷贝打交道。

- **定义**
  - 设有数据库系统DBS，DB为DBS中的数据库，DBM(t)为在时刻t，DB在内存的数据集，DBM(t)属于DB。TS为DBS中所有可能的事务构成的集合。AT(t)为在时刻t处于活动状态的事务集，AT(t)属于TS。Dt(T)为事务T在时刻t所操作的数据集，
  - Dt(T)属于DB。若在任意时刻t，均有: 任意T属于AT(t) Dt(T)属于DBM(t)成立，则称DBS为一个内存数据库系统，简称MMDBS;DB为一个内存数据库，简称MMDB
  - 常见的例子有MySQL的MEMORY存储引擎、eXtremeDB、TT、FastDB、SQLite、Microsoft SQL Server Compact等

- 传统的数据库系统是关系型数据库，开发这种数据库的目的，是处理永久、稳定的数据。关系数据库强调维护数据的完整性、一致性，但很难顾及有关数据及其处理的定时限制，不能满足工业生产管理实时应用的需要，因为实时事务要求系统能较准确地预报事务的运行时间。对磁盘数据库而言，由于磁盘存取、内外存的数据传递、缓冲区管理、排队等待及锁的延迟等使得事务实际平均执行时间与估算的最好情况执行时间相差很大，如果将整个数据库或其主要的"工作"部分放入内存，使每个事务在执行过程中没有I/O，则为系统较准确估算和安排事务的运行时间，使之具有较好的动态可预报性提供了有力的支持，同时也为实现事务的定时限制打下了基础。这就是内存数据库出现的主要原因。

- 内存数据库所处理的数据通常是"短暂"的，即有一定的有效时间，过时则有新的数据产生，而当前的决策推导变成无效。所以，实际应用中采用内存数据库来处理实时性强的业务逻辑处理数据。而传统数据库旨在处理永久、稳定的数据，其性能目标是高的系统吞吐量和低的代价，处理数据的实时性就要考虑得相对少一些。实际应用中利用传统数据库这一特性存放相对实时性要求不高的数据。

- 在实际应用中这两种数据库常常结合使用，而不是以内存数据库替代传统数据库。

- 闪存（Flash Memory）是一种长寿命的非易失性（在断电情况下仍能保持所存储的数据信息）的存储器，数据删除不是以单个的字节为单位而是以固定的区块为单位，区块大小一般为256KB到20MB。闪存是电子可擦除只读存储器（EEPROM）的变种，EEPROM与闪存不同的是，它能在字节水平上进行删除和重写而不是整个芯片擦写，这样闪存就比EEPROM的更新速度快。由于其断电时仍能保存数据，闪存通常被用来保存设置信息，如在电脑的BIOS（基本输入输出程序）、PDA（个人数字助理）、数码相机中保存资料等。另一方面，闪存不像RAM（随机存取存储器）一样以字节为单位改写数据，因此不能取代RAM。

- 闪存卡（Flash Card）是利用闪存（Flash Memory）技术达到存储电子信息的存储器，一般应用在数码相机，掌上电脑，MP3等小型数码产品中作为存储介质，所以样子小巧，有如一张卡片，所以称之为闪存卡。根据不同的生产厂商和不同的应用，闪存卡大概有SmartMedia（SM卡）、Compact Flash（CF卡）、MultiMediaCard（MMC卡）、Secure Digital（SD卡）、Memory Stick（记忆棒）、XD-Picture Card（XD卡）和微硬盘（MICRODRIVE）这些闪存卡虽然外观、规格不同，但是技术原理都是相同的。
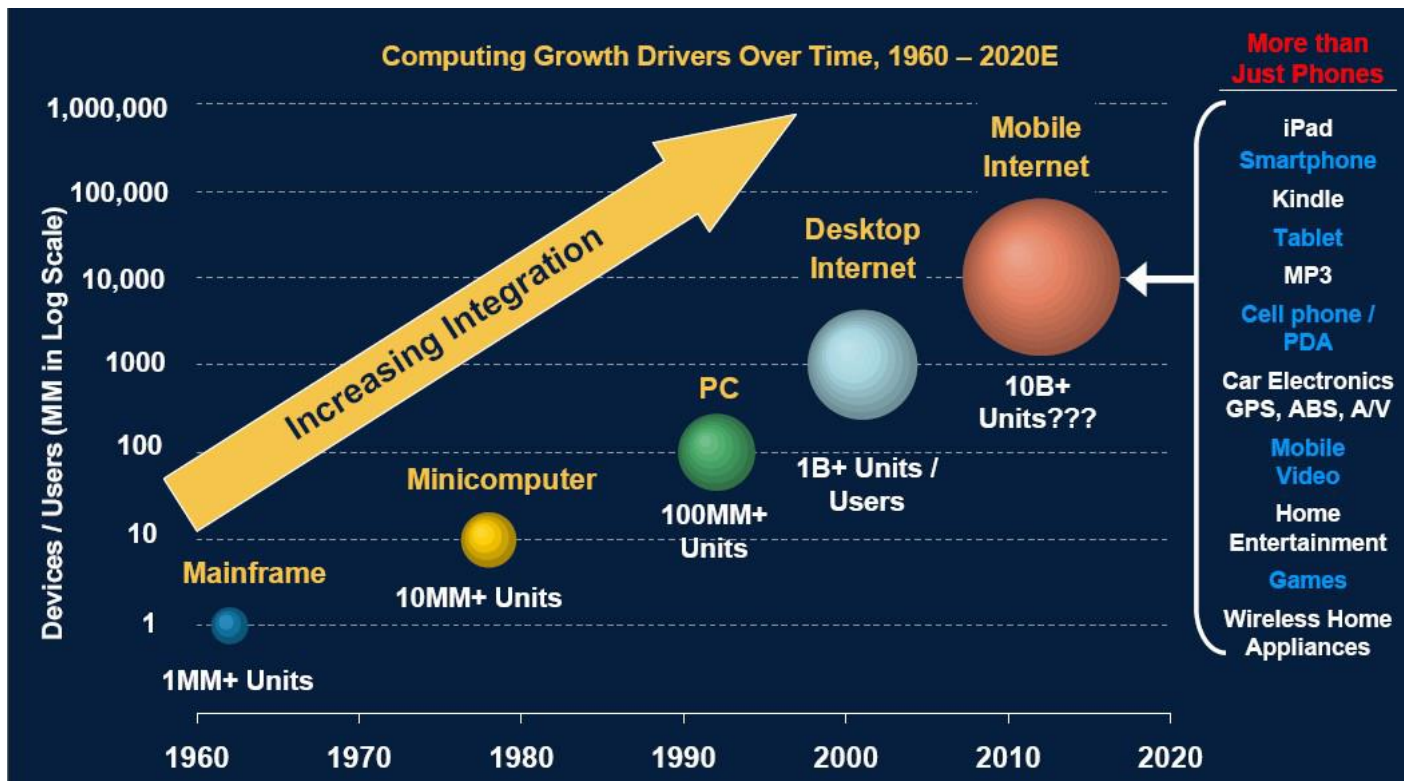
# 大数据及云存储

# The Future is Full of Opportunity

- Designing a next Internet – GENI
  - Driving advances in all fields of science and engineering
- Wreckless driving
- Personalized education
- Predictive, preventive, personalized medicine
- Quantum computing (量子计算)
- Personalized health monitoring => quality of life
- Data-intensive supercomputing
- Neurobotics (神经学机器人)
- Synthetic biology (合成生物学)
- The algorithmic lens => Cyber-enabled Discovery and Innovation

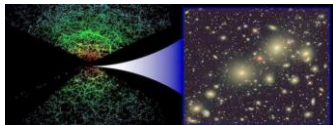Computing Growth Drivers Over Time, 1960 – 2020E

# ▶ 大数据的困难

- ## 容量大

美国国会图书馆存档信息量：约 **80TB**

| | | |
|---|---|---|
| 科学计算 |  | 新墨西哥州的天文望远镜**每年**产生**80TB**的图像信息 |
| 生物信息 |  | 第一个中国人的全基因组图谱，**1177亿**碱基对 |
| 电子商务 | 淘宝网 Taobao.com | **每月**交易**21亿**笔，产生**300TB**交易日志信息 |
| 网络生活 | facebook® | **7亿**用户、**400亿**张照片，总容量超过**1500TB** |

**Kilo**
**Mega**
**Giga**
**Tera**
**Peta**
**Exa**
**Zeta**
**Yotta**

**79**

# ▶ 大数据的困难

- **数据中心—看上去，规模庞大**


脸谱(Facebook) 数据中心


苹果(Apple) 数据中心


微软(Microsoft)数据中心


谷歌(Google) 数据中心

# 大数据的困难

- **数据中心—走进去，结构复杂**

# ▶ 新技术

- **针对存储**
  - **石英玻璃数据存储技术**：可保存数据数亿年，电脑硬盘驱动器的数据保存时间就只有10年，使用较频繁的闪存盘数据保存时间至多5年，而磁带的数据保存时限则达15至30年。
  - **双写入存储技术**：提升至10TB/平方英寸，目前硬盘中所使用的垂直记录技术的存储密度极限大约是每平方英寸数百GB。 两种新写入方式名为"热辅助磁记录技术"(Thermally-Assisted Magnetic Recording：TAR)，和 "位式记录技术" (bit-patterned recording：BPR)

- **针对传输**
  - **量子网络瞬间通信不延迟:** 一个粒子可以传递有限的信息，而亿万个粒子联手就形成了量子网络，而这种信息传递没有任何的延迟。通过量子网络相互连接的量子计算机和量子服务器将应用量子纠缠实现无缝通讯
  - **光纤：** 一根光纤可供2.1亿对人同时通话，一根普通单模光纤中C波段168路，每路103Gb/s的超大容量超密集波分复用传输2240公里，传输总容量达到17.32Tb/s，相当于2.1亿对人在一根光纤上同时通话
  - **HMC技术：** 传输率可达1TB/s，Hybrid Memory Cube技术相比于当今主流的DDR3在能源效率上有至少7倍以上的优势。Hybrid Memory Cube技术使用堆叠技术将内存芯片压缩成一个紧凑的"立方体"，并配有新的高速传输接口。新的数据传输接口传输率可达到1TB/s。
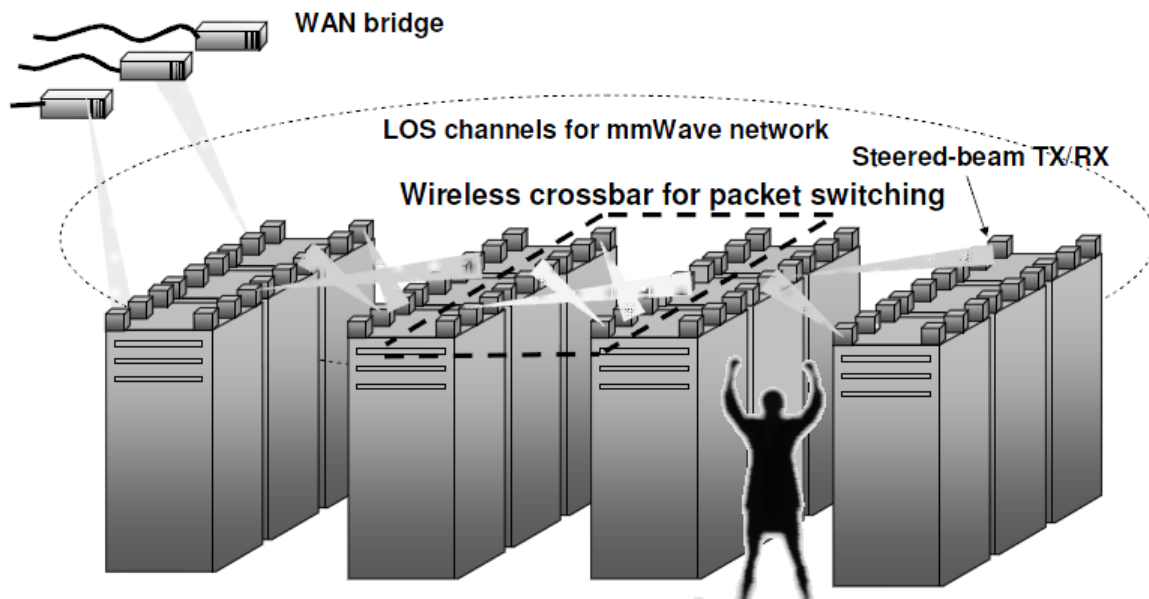
- **针对Hadoop**

  – 交互式数据分析系统Dremel

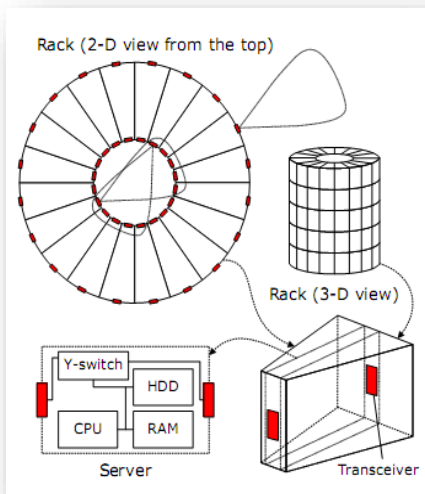    • 可以在几秒钟内处理PB级别的数据，并能轻松应对即时查询

  – 谷歌的图数据计算框架Pregel，微软的Trinity

    • 针对大规模图算法（如图遍历（BFS）、PageRank，最短路径（SSSP）等）。处理一个有着几十亿节点、上万亿条边的图，只需数分钟即可完成，其执行时间随着图的大小呈线性增长。

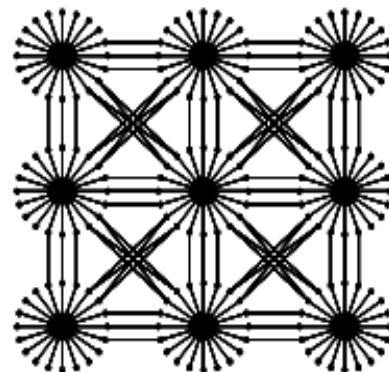# ▶ 新技术

- **针对数据中心**



WAN bridge

LOS channels for mmWave network

Steered-beam TX/RX

Wireless crossbar for packet switching

- **针对数据中心**
  - 扇形机架 + 柱形机箱
  - 机箱内: 冗余链接、强连通
  - 机箱间: 2维类Mesh拓扑

- **针对数据中心**
  - 微型数据中心
  - 台式数据中心

# ▶ 云存储（Cloud Storage）

- 云存储是云计算的存储部分，即虚拟化的、易于扩展的存储资源池。用户通过云计算使用存储资源池，但不是所有的云计算的存储部分都是可以分离的。

- 云存储意味着存储可以作为一种服务，通过网络提供给用户。用户可通过若干种方式使用存储，按使用（时间、空间或两者结合）付费：

  - 通过互联网开放接口（如REST），使得第三方网站可以通过云存储提供的服务为用户提供完整的Web服务；

  - 用户直接使用存储相关的在线服务，比如网络硬盘，在线存储，在线备份，或在线归档等服务；

  - 用户传送文件、或者服务商发布内容时的缓冲

- 云存储在云计算 (cloud computing)概念上延伸和发展出来的概念
  - 云计算是分布式处理(Distributed Computing)、并行处理(Parallel Computing)和网格计算(Grid Computing)的发展，是透过网络将庞大的计算处理程序自动分拆成无数个较小的子程序，再交由多部服务器所组成的庞大系统经计算分析之后将处理结果回传给用户。通过云计算技术，网络服务提供者可以在数秒之内，处理数以千万计甚至亿计的信息，达到和"超级计算机"同样强大的网络服务
  - 云存储的概念与云计算类似，它是指通过集群应用、网格技术或分布式文件系统等功能，将网络中大量各种不同类型的存储设备通过应用软件集合起来协同工作，共同对外提供数据存储和业务访问功能的一个系统
  - 云存储对使用者来讲，不是指某一个具体设备，而是指一个由许许多多个存储设备和服务器所构成的集合体。使用者使用云存储，并不是使用某一个存储设备，而是使用整个云存储系统带来的一种数据访问服务。所以严格来讲，云存储不是存储，而是一种服务！

- 与传统存储设备相比，云存储不仅仅是一个硬件，而是一个网络设备、存储设备、服务器、应用软件、公用访问接口、接入网、和客户端程序等多个部分组成的复杂系统。各部分以存储设备为核心，通过应用软件来对外提供数据存储和业务访问服务

- 云存储系统的结构模型由 4层组成
  - 存储层是云存储最基础的部分
  - 基础管理层是云存储最核心部分，也是最难实现的部分
  - 应用接口层是云存储最灵活多变的部分
  - 访问层

# ▶ 云存储（续）

- 云存储系统的结构模型由 4层组成：存储层、基础管理层、应用接口层、访问层
- **存储层**
  - 存储层是云存储最基础的部分。存储设备可以是FC光纤通道存储设备，可以是NAS和 iSCSI等IP存储设备，也可以是 SCSI或SAS等 DAS存储设备。云存储中的存储设备往往数量庞大且分布多不同地域，彼此之间通过广域网、互联网或者 FC光纤通道网络连接在一起
  - 存储设备之上是一个统一存储设备管理系统，可以实现存储设备的逻辑虚拟化管理、多链路冗余管理，以及硬件设备的状态监控和故障维护
- **基础管理层**
  - 基础管理层是云存储最核心的部分，也是云存储中最难以实现的部分。基础管理层通过集群、分布式文件系统和网格计算等技术，实现云存储中多个存储设备之间的协同工作，使多个的存储设备可以对外提供同一种服务，并提供更大更强更好的数据访问性能
  - CDN内容分发系统、数据加密技术保证云存储中的数据不会被未授权的用户所访问，同时，通过各种数据备份和容灾技术和措施可以保证云存储中的数据不会丢失，保证云存储自身的安全和稳定

# ▶ 云存储（续）

- **应用接口层**
  - 应用接口层是云存储最灵活多变的部分。不同的云存储运营单位可根据实际业务类型，开发不同的应用服务接口，提供不同的应用服务。比如视频监控应用平台、IPTV和视频点播应用平台、网络硬盘应用平台，远程数据备份应用平台等

- **访问层**
  - 任何一个授权用户都可以通过标准的公用应用接口来登录云存储系统，享受云存储服务。云存储运营单位不同，云存储提供的访问类型和访问手段也不同

- 云存储系统是多设备、多应用、多服务协同工作集合体，其技术前提：
  - 云存储系统是多设备、多应用、多服务协同工作的集合体
  - WEB2.0技术
  - 应用存储的发展
  - 集群技术、网格技术和分布式文件系统
  - CDN内容分发、P2P技术、数据压缩技术、重复数据删除技术、数据加密技术
  - 存储虚拟化技术、存储网络化管理技术

# ▶ 云存储（续）

- **云存储得到众多厂商的关注和支持**

  - Amazon推出的Elastic Compute Cloud（EC2：弹性计算云）云存储产品，旨在为用户提供互联网服务形式同时提供更强的存储和计算功能

  - 内容分发网络服务提供商CDNetworks和业界著名的云存储平台服务商 Nirvanix发布了一项合作，并宣布结成战略伙伴关系，以提供业界的云存储和内容传送服务集成平台

  - 微软推出提供网络移动硬盘服务的Windows Live SkyDrive

  - EMC加入道里可信基础架构项目，致力于云计算环境下关于信任和可靠度保证的全球研究协作

  - IBM将云计算标准作为全球备份中心的3亿美元扩展方案的一部分

- 分布式、云化、闪存化、智能等存储不断发展，随着人工智能、物联网、区块链、人体增强 (Human Augmentation) 等技术快速发展，存储也将迎来新的形态，如：
  - **边缘存储**
  - **长期存储**
  - **生物存储/基因存储**
  - **区块链存储**
  - **量子存储**

- 加州大学伯克利分校的研究团队在2015年发表 "The Cloud is Not Enough: Saving IoT from the Cloud" 文章，指出物联网与互联网的七个不同之处：隐私和安全；可伸缩性；交互模型；延迟；带宽；可用性；持久性管理。边缘存储应该具有高可用性、超低延迟、高安全性、高隐私性、弱一致性、功耗低、空间小的特点，边缘存储仍处于发展早期，还有很多需研究探索

| | Web | IoT |
|---|---|---|
| Privacy & Security | Open for access | Personal sensitive data |
| Scalability | Power-law | Billlion devices + updates |
| Interaction Model | Human | Machine |
| Latency | Variable | Reactive |
| Bandwidth | Downstream | Upstream |
| Availability | None | Requirement |
| Durability Management | Cloud controls | Users control |

# ▶ 长期存储

- **蓝光光盘**
  - 2019中国数据与存储峰会，"Long Data"挑战国际上研究课题How to preserve information for 100 years？基于蓝光技术可提供解决办法，蓝光光盘具超过50年寿命，极低功耗，高安全性，较高密度，较低成本，对环境要求低。总体而言，光存储技术尤其适合大量冷数据的长期安全存储。
  - 针对于需要保存长达数十年的行业或场景，基于蓝光光盘的存储技术是一个选择，如电子档案、医疗影像、金融影像、备份归档等。一个标准机柜可保存1.22万片光盘，如果采用500GB光盘，裸容量可高达6PB。市场上一种M-Disc的产品，号称千年光盘

- **玻璃光盘**
  - 正在研究的玻璃光盘容量是蓝光光盘的几千倍，能承受超过几百度的高温，寿命能高达上万年甚至百亿年。如微软的玻璃光盘项目ProjectSilica

- **全息光存储**
  - 蓝光之后下一代变革性光存储技术包括两种：第一种是同轴多维全息光存储技术，列入国家重点研发计划。第二种是2014年得了诺贝尔奖的突破光的衍射极限项目，澳大利亚科学家把这个技术用到光上，把光斑从300纳米理论上可以减少到九个纳米，容量上得到巨大提高，至少可达每盘15TB，理想上可实现PB级。
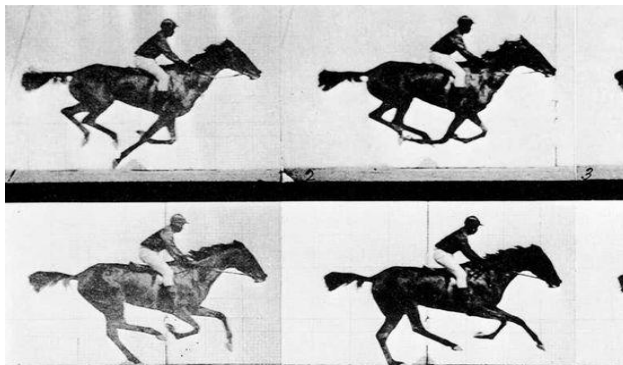
# ▶ 生物存储或基因存储

- **生物存储或基因存储**
  - 人类历史上，已经有过历经千百年的存储介质，如青铜器铭文、竹简、纸书、石碑等。不过，这类手段存储密度极低。《自私的基因》一书观点：包括人在内的动植物不过是DNA 繁衍的躯壳，你也可以看成是DNA 的存储器。
  - 可以反向让DNA 片段，即基因成为人类存放信息的存储器。基因存储将0 、1 数据通过一定的编码方法转换成DNA 中的A 、T 、C 、G 四种碱基，通过合成含有这些碱基序列的DNA 即可实现数据信息存储。

# ▶ 生物存储或基因存储

- **生物存储或基因存储**
  - 2017年7月，哈佛大学医学院利用CRISPR DNA编辑技术，将赛马视频录入大肠杆菌的基因组中，并以超过90%的准确率读取出来。
  - 2019年7月，美国布朗大学研究人员用基于生物小分子的存储系统累计存储了超过10万比特的数字图像信息，从中获得图像的准确率可达98% 以上。
  - 2020年2月，华尔街日报报道，波士顿初创公司Catalog Technologies 正在开发一种独特的方式来存储大量数据，将维基百科中14GB 数据存储在DNA 分子中，像"在试管中滴几滴水"。用分子生物学仪器"打印"合成分子序列，以DNA形式存储和表示数字信息，使用DNA测序仪和专有软件的来读取信息，软件可以将分子翻译成文本、图片甚至视频。

# ▶ 区块链存储

- 如何在激励数据提供方分享数据的同时，又能保护隐私？区块链和存储的结合
    - 区块链解决了数据确权、激励分享、数据资产交易和流转等问题。但区块链要发展，区块链基础设施要先行，目前区块链基础设施还处于非常早期，使得区块链应用数据的存储，具有很大的问题。无论是公有链还是联盟链，数据都是存放在中心化的存储上，前者可能选择公有云存储，后者选择类似NAS存储，都存在安全隐患、隐私泄露的可能性。实际上，去中心化的应用DApp（Distributed App）需要的是端到端的去中心化的基础设施，包括去中心化方式组织的存储。



**区块链存储：一种新的共享模式，存储空间来自多个中心**

第一类是**公链存储**：也即非中心化的存储+公链，例如IPFS+Filecoin、STORJ +Storj等；通常是跨越全球的存储池 + Token激励机制。
第二类是**许可链存储**：也即非中心化的存储+许可链，利用了区块链的特征，如分布式、不可篡改、可追踪、加密安全性等。核心是，**单一个体没有机会控制整个存储**。此种存储主要用于私有链或者联盟链；

| 本地存储时代 1957- | 云存储时代 2007- | 区块链存储时代 2017- | 中心化存储（如公有云存储），数据容易泄露。2017年6月21日，美国共和党全国委员会放在AWS S3的91.1 TB数据发生泄露，包含超过1.98亿名美国选民的敏感个人资料，例如姓名、出生日期、住址、电话号码以及选民注册细节信息；甚至还包括政治团体采用的先进情绪分析来预测个人选民如何处理热门问题，如枪支所有权、干细胞研究和堕胎权，宗教信仰和种族等信息。 |

云存储
成本高、管理难、安全性低　成本较低、管理较难、安全性低　多中心、安全性高、低成本

**原因：所有权和运营权相分离，可以保障数据的隐私，单一个体没有机会操控。**
备注：公有云存储所有权和使用权相分离，但所有权和运营权合为一体，存在隐患。

**100**

# ▶ 区块链存储（续）

- **区块链和存储的结合**
  - 2018年7月，Gartner技术成熟度曲线中区块链存储(Distributed Storage in Blockchain)列入科技诞生促动期(Technology Trigger)，预计2023~2028年进入到成熟应用的技术阶段，将有大量主流用户开始接纳



Hype Cycle for Blockchain Technologies, 2018

**101**