

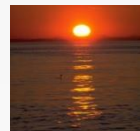


# Chapter 11: Storage and File Structure

**Database System Concepts, 5th Ed.**

©Silberschatz, Korth and Sudarshan

See [www.db-book.com](http://www.db-book.com) for conditions on re-use





# Classification of Physical Storage Media

- ❑ Speed with which data can be accessed
- ❑ Cost per unit of data
- ❑ Reliability
  - ❑ data loss on power failure or system crash
  - ❑ physical failure of the storage device
- ❑ Can differentiate storage into:
  - ❑ **volatile storage**: loses contents when power is switched off
  - ❑ **non-volatile storage**:
    - ▶ Contents persist even when power is switched off.
    - ▶ Includes secondary and tertiary storage, as well as battery-backed up main-memory.





# Physical Storage Media

- **Cache** – fastest and most costly form of storage; volatile; managed by the computer system hardware
  - (Note: “Cache” is pronounced as “cash”)
- **Main memory:**
  - fast access (10s to 100s of nanoseconds; 1 nanosecond =  $10^{-9}$  seconds)
  - generally too small (or too expensive) to store the entire database
    - ▶ capacities of up to a few Gigabytes widely used currently
    - ▶ Capacities have gone up and per-byte costs have decreased steadily and rapidly (roughly factor of 2 every 2 to 3 years)
  - **Volatile** — contents of main memory are usually lost if a power failure or system crash occurs.





# Physical Storage Media (Cont.)

## □ Flash memory

- Data survives power failure
- Data can be written at a location only once, but location can be erased and written to again
  - ▶ Can support only a limited number (10K – 1M) of write/erase cycles.
  - ▶ Erasing of memory has to be done to an entire bank of memory
- Reads are roughly as fast as main memory
- But writes are slow (few microseconds), erase is slower





# Physical Storage Media (Cont.)

## □ Flash memory

### □ NOR Flash

- ▶ Fast reads, very slow erase, lower capacity
- ▶ Used to store program code in many embedded devices

### □ NAND Flash

- ▶ Page-at-a-time read/write, multi-page erase
- ▶ High capacity (several GB)
- ▶ Widely used as data storage mechanism in portable devices





# Physical Storage Media (Cont.)

## □ Magnetic-disk

- Data is stored on spinning disk, and read/written magnetically
- Primary medium for the long-term storage of data; typically stores entire database.
- Data must be moved from disk to main memory for access, and written back for storage
- **direct-access** – possible to read data on disk in any order, unlike magnetic tape
- Survives power failures and system crashes
  - ▶ disk failure can destroy data: is rare but does happen





# Physical Storage Media (Cont.)

## □ Optical storage

- non-volatile, data is read optically from a spinning disk using a laser
- CD-ROM (640 MB) and DVD (4.7 to 17 GB) most popular forms
- Write-one, read-many (WORM) optical disks used for archival storage (CD-R, DVD-R, DVD+R)
- Multiple write versions also available (CD-RW, DVD-RW, DVD+RW, and DVD-RAM)
- Reads and writes are slower than with magnetic disk
- **Juke-box** systems, with large numbers of removable disks, a few drives, and a mechanism for automatic loading/unloading of disks available for storing large volumes of data





# Physical Storage Media (Cont.)

## □ Tape storage

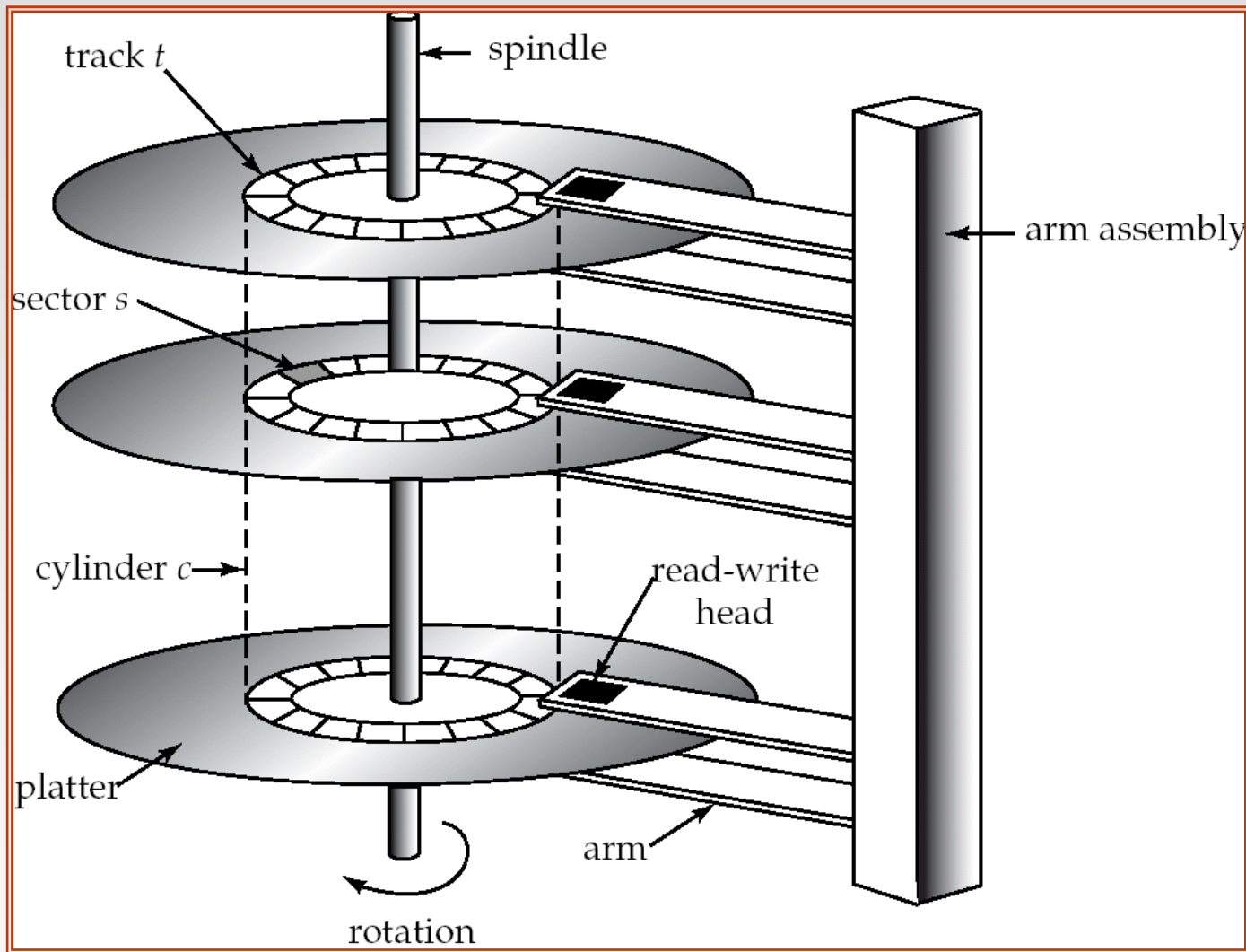
- non-volatile, used primarily for backup (to recover from disk failure), and for archival data
- **sequential-access** – much slower than disk
- very high capacity (40 to 300 GB tapes available)
- tape can be removed from drive  $\Rightarrow$  storage costs much cheaper than disk, but drives are expensive
- Tape jukeboxes available for storing massive amounts of data
  - ▶ hundreds of terabytes (1 terabyte =  $10^9$  bytes) to even a petabyte (1 petabyte =  $10^{12}$  bytes)







# Magnetic Hard Disk Mechanism



**NOTE:** Diagram is schematic, and simplifies the structure of actual disk drives





# Magnetic Disks

- **Read-write head**
  - Positioned very close to the platter surface (almost touching it)
  - Reads or writes magnetically encoded information.
- Surface of platter divided into circular **tracks**
  - Over 50K-100K tracks per platter on typical hard disks
- Each track is divided into **sectors**.
  - Sector size typically 512 bytes
  - Typical sectors per track: 500 (on inner tracks) to 1000 (on outer tracks)
- To read/write a sector
  - disk arm swings to position head on right track
  - platter spins continually; data is read/written as sector passes under head





# Magnetic Disks (Cont.)

- Head-disk assemblies
  - multiple disk platters on a single spindle (1 to 5 usually)
  - one head per platter, mounted on a common arm.
- **Cylinder**  $i$  consists of  $i^{\text{th}}$  track of all the platters
- Earlier generation disks were susceptible to “head-crashes” leading to loss of all data on disk
  - Current generation disks are less susceptible to such disastrous failures, but individual sectors may get corrupted





# Performance Measures of Disks

- **Access time** – the time it takes from when a read or write request is issued to when data transfer begins. Consists of:
  - **Seek time** – time it takes to reposition the arm over the correct track.
    - ▶ Average seek time is  $1/2$  the worst case seek time.
      - Would be  $1/3$  if all tracks had the same number of sectors, and we ignore the time to start and stop arm movement
    - ▶ 4 to 10 milliseconds on typical disks
  - **Rotational latency** – 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 r.p.m.)





# Performance Measures (Cont.)

- **Data-transfer rate** – the rate at which data can be retrieved from or stored to the disk.
  - 25 to 100 MB per second max rate, lower for inner tracks
  - Multiple disks may share a controller, so 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





# Performance Measures (Cont.)

- **Mean time to failure (MTTF)** – the average time the disk is expected to run continuously without any failure.
  - Typically 3 to 5 years
  - 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 relatively new disks, on an average one will fail every 1200 hours
  - MTTF decreases as disk ages





# 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, and
  - ▶ **high reliability** by storing data redundantly, so that data can be recovered even if a disk fails
- 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 (approx. 41 days)





# Improvement of Reliability via Redundancy

- **Redundancy** – store extra information that can be used to rebuild information lost in a disk failure
- E.g., **Mirroring** (or **shadowing**)
  - Duplicate every disk. 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 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







# Improvement of Reliability via Redundancy

- Mean time to data loss depends on mean time to failure, and mean time to repair
  - E.g. MTTF of 100,000 hours, mean time to repair of 10 hours gives mean time to data loss of  $500 \cdot 10^6$  hours (or 57,000 years) for a mirrored pair of disks (ignoring dependent failure modes)





# Improvement in Performance via Parallelism

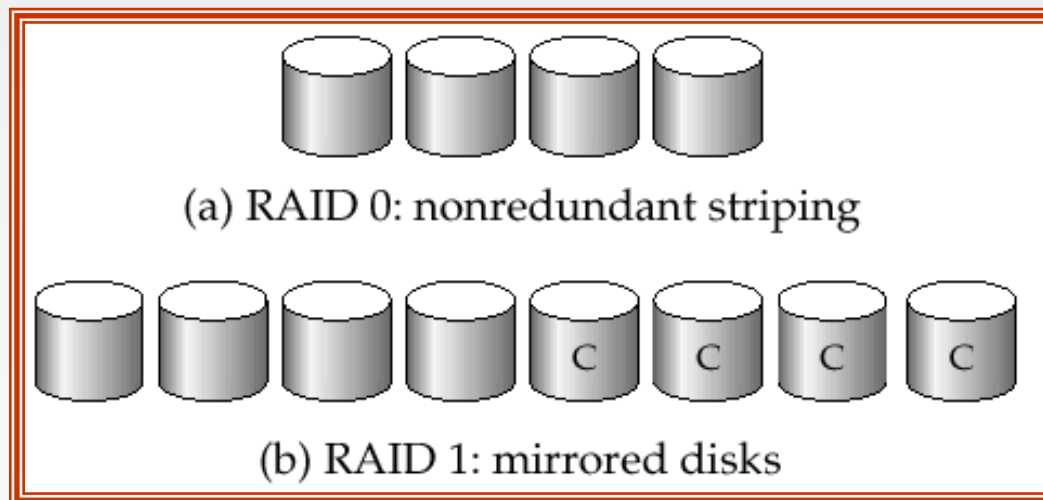
- Two main goals of parallelism in a disk system:
  1. Load balance multiple small accesses to increase throughput
  2. Parallelize large accesses to reduce response time.
- Improve transfer rate by striping data across multiple disks.
- **Bit-level striping** – split the bits of each byte across multiple disks
  - But seek/access time worse than for a single disk
    - ▶ Bit level striping is not used much any more
- **Block-level striping** – with  $n$  disks, block  $i$  of a file goes to disk  $(i \bmod n) + 1$ 
  - 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





# RAID Levels

- ❑ RAID organizations, or RAID levels, have differing cost, performance and reliability characteristics
- ❑ **RAID Level 0: Block striping; non-redundant.**
  - ❑ Used in high-performance applications where data lost is not critical.
- ❑ **RAID Level 1: Mirrored disks** with block striping
  - ❑ Offers best write performance.
  - ❑ Popular for applications such as storing log files in a database system.





# RAID Levels (Cont.)

- ❑ **RAID Level 2: Memory-Style Error-Correcting-Codes (ECC)** with bit striping.
- ❑ **RAID Level 3: Bit-Interleaved Parity**
  - ❑ a single parity bit is enough for error correction, not just detection
    - ▶ 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



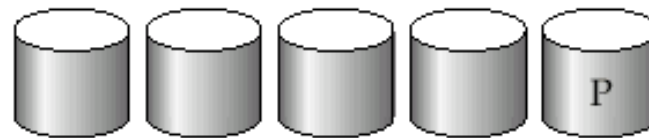
(d) RAID 3: bit-interleaved parity





# RAID Levels (Cont.)

- 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.
- **RAID Level 4: Block-Interleaved Parity**; 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 Levels (Cont.)

## □ 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
- 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)
  - ▶ Or by recomputing the parity value using the new values of blocks corresponding to the parity block
    - 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





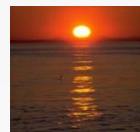
# RAID Levels (Cont.)

- **RAID Level 5: Block-Interleaved Distributed Parity**; 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  $(n \bmod 5) + 1$ , with the data blocks stored on the other 4 disks.



(f) 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



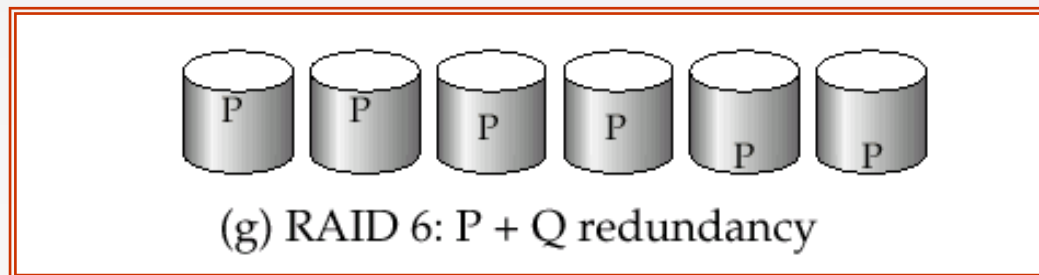


# RAID Levels (Cont.)

## □ 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.







# Choice of RAID Level

- ❑ 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 since bit-striping forces single block reads to access all disks, wasting disk arm movement
- ❑ Level 6 is rarely used since levels 1 and 5 offer adequate safety for most applications
- ❑ So competition is mainly between 1 and 5





# Choice of RAID Level (Cont.)

- 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

