# Chapter 10: Storage and File Structure Part 1

# Chapter 10: Storage and File Structure Part 1

- Overview of Physical Storage Media
- Magnetic Disks
- RAID
- Tertiary Storage
- Storage Access

# 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 batter-backed up main-memory.

# Physical Storage Media

- **Cache** – fastest and most costly form of storage; volatile; managed by the computer system hardware.

- **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
  - Widely used in embedded devices such as digital cameras, phones, and USB keys

# 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
    - Much slower access than main memory (more on this later)
  - **direct-access** – possible to read data on disk in any order, unlike magnetic tape
  - Capacities range up to roughly 1.5 TB as of 2009
    - Much larger capacity and cost/byte than main memory/flash memory
    - Growing constantly and rapidly with technology improvements (factor of 2 to 3 every 2 years)
  - Survives power failures and system crashes
    - disk failure can destroy data, but is rare

# 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

  - Blu-ray disks: 27 GB to 54 GB

  - 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
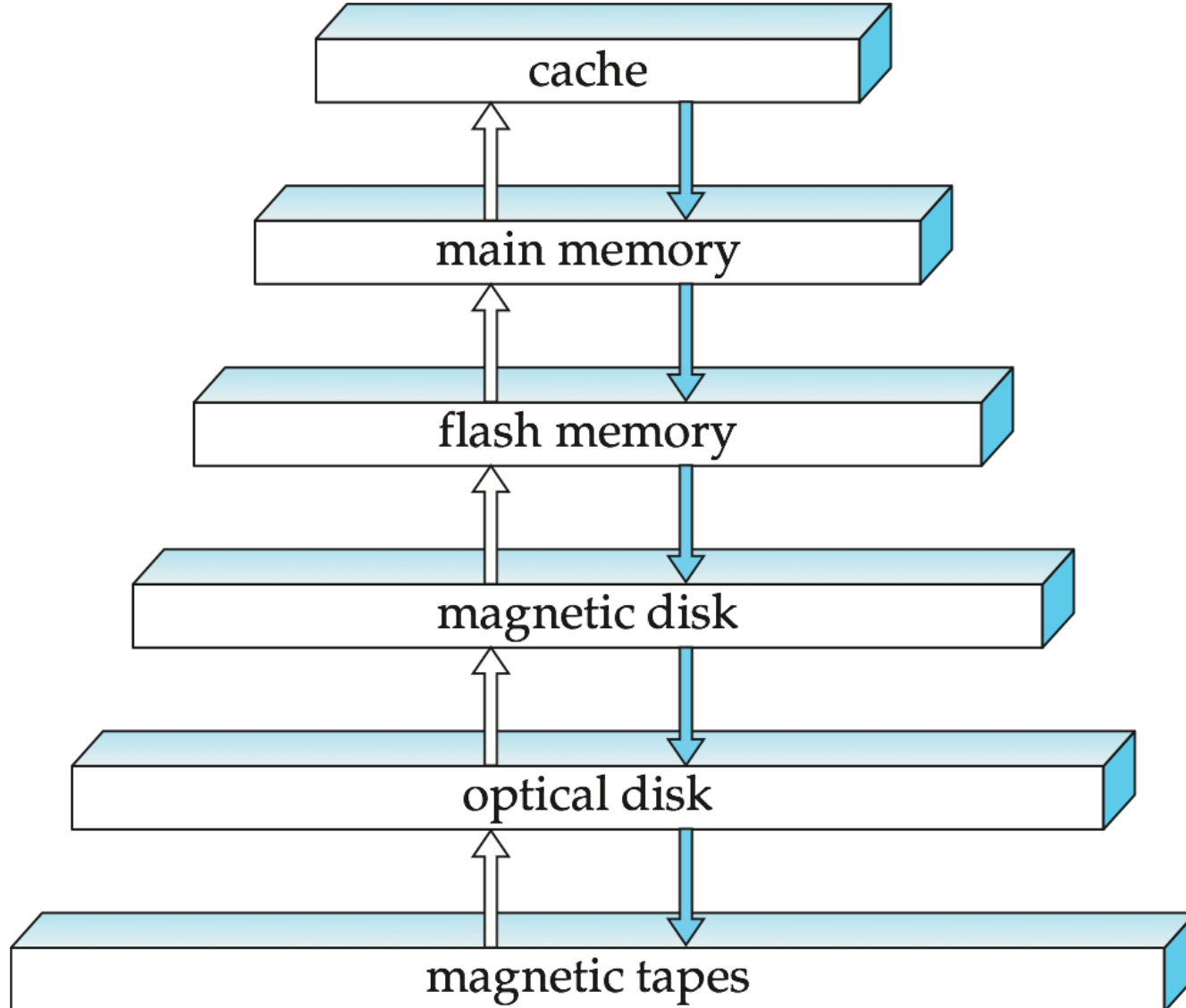
# Physical Storage Media (Cont.)

- **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 multiple **petabytes** (1 petabyte = $10^{12}$ bytes)
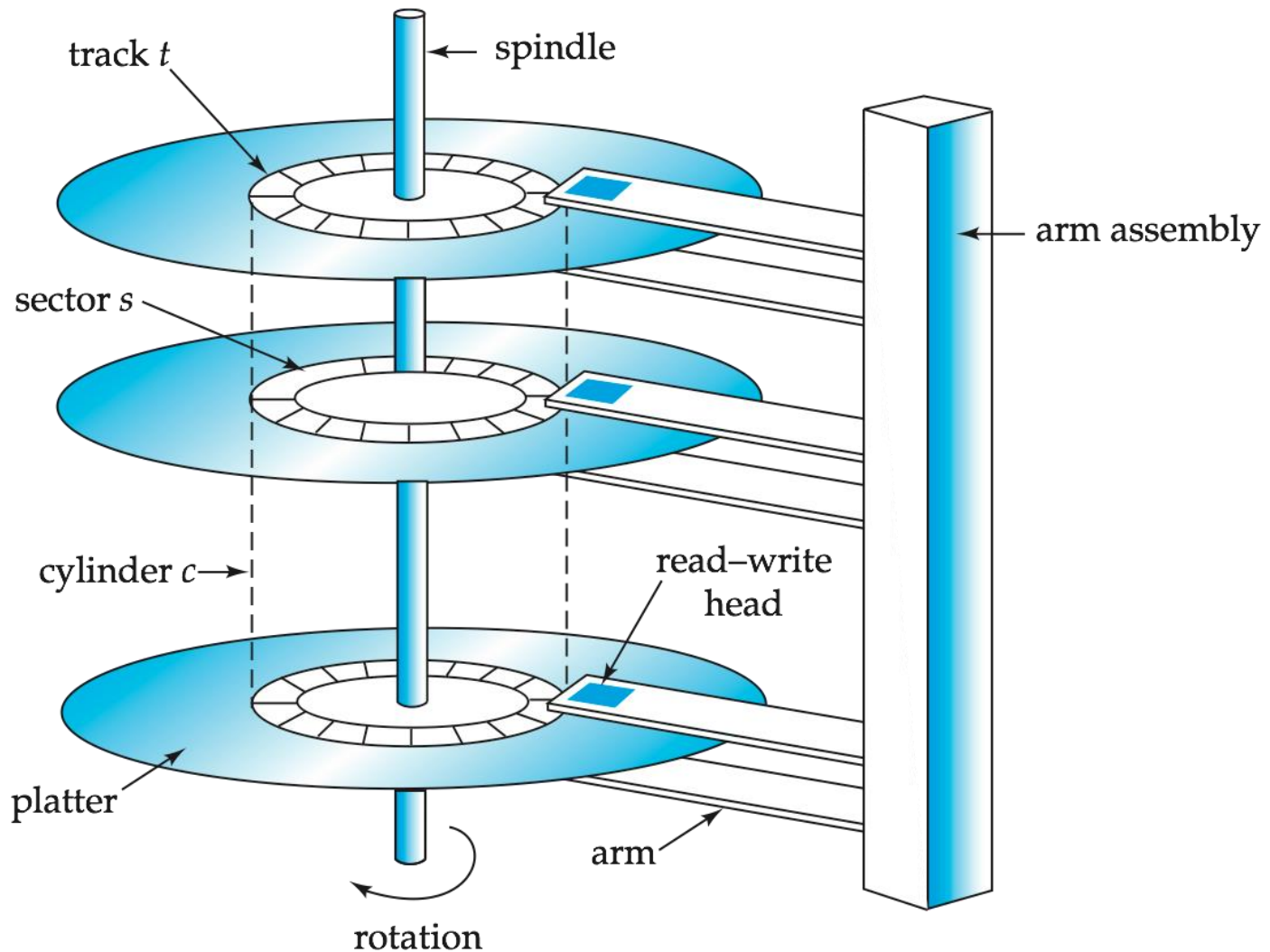
# Storage Hierarchy

# Storage Hierarchy (Cont.)

- **primary storage:** Fastest media but volatile (cache, main memory).

- **secondary storage:** next level in hierarchy, non-volatile, moderately fast access time
  - also called **on-line storage**
  - E.g. flash memory, magnetic disks

- **tertiary storage:** lowest level in hierarchy, non-volatile, slow access time
  - also called **off-line storage**
  - E.g. magnetic tape, optical storage

# 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.**
  - A sector is the smallest unit of data that can be read or written.
  - Sector size typically 512 bytes
  - Typical sectors per track: 500 to 1000 (on inner tracks) to 1000 to 2000 (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
- Head-disk assemblies
  - multiple disk platters on a single spindle (1 to 5 usually)
  - one head per platter, mounted on a common arm.
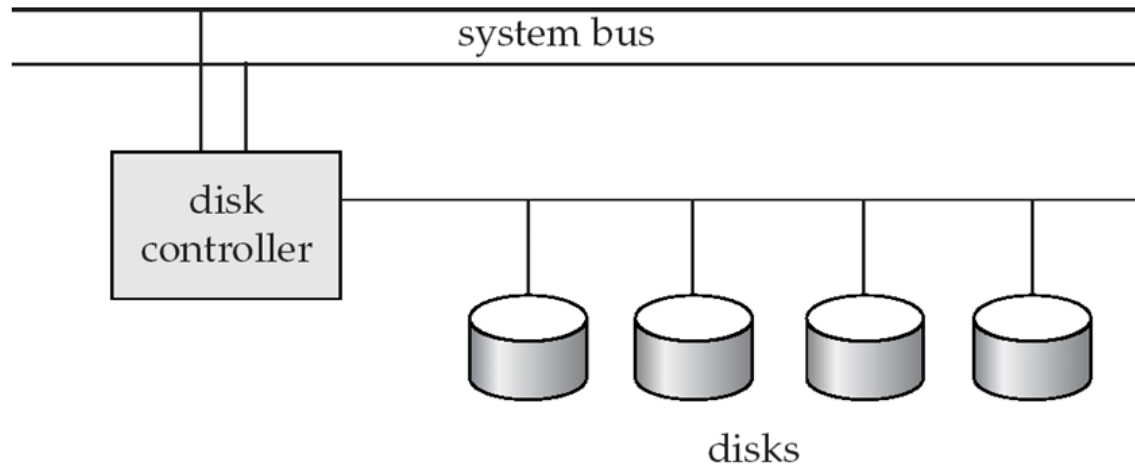- **Cylinder** $i$ consists of $i^{th}$ track of all the platters

# Magnetic Disks (Cont.)

- Earlier generation disks were susceptible to head-crashes
  - Surface of earlier generation disks had metal-oxide coatings which would disintegrate on head crash and damage all data on disk
  - Current generation disks are less susceptible to such disastrous failures, although individual sectors may get corrupted
- **Disk controller** – interfaces between the computer system and the disk drive hardware.
  - accepts high-level commands to read or write a sector
  - initiates actions such as moving the disk arm to the right track and actually reading or writing the data
  - Computes and attaches **checksums** to each sector to verify that data is read back correctly
    - If data is corrupted, with very high probability stored checksum won't match recomputed checksum
  - Ensures successful writing by reading back sector after writing it
  - Performs remapping of bad sectors

# Disk Subsystem



- Multiple disks connected to a computer system through a controller
  - Controllers functionality (checksum, bad sector remapping) often carried out by individual disks; reduces load on controller
- Disk interface standards families
  - ATA (AT adaptor) range of standards
  - SATA (Serial ATA)
  - SCSI (Small Computer System Interconnect) range of standards
  - SAS (Serial Attached SCSI)
  - Several variants of each standard (different speeds and capabilities)

# Disk Subsystem

- Disks usually connected directly to computer system

- In **Storage Area Networks (SAN)**, a large number of disks are connected by a high-speed network to a number of servers

- In **Network Attached Storage (NAS)** networked storage provides a file system interface using networked file system protocol, instead of providing a disk system interface

# 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.)

■ **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. SATA: 150 MB/sec, SATA-II 3Gb (300 MB/sec)
  - ‣ Ultra 320 SCSI: 320 MB/s, SAS (3 to 6 Gb/sec)
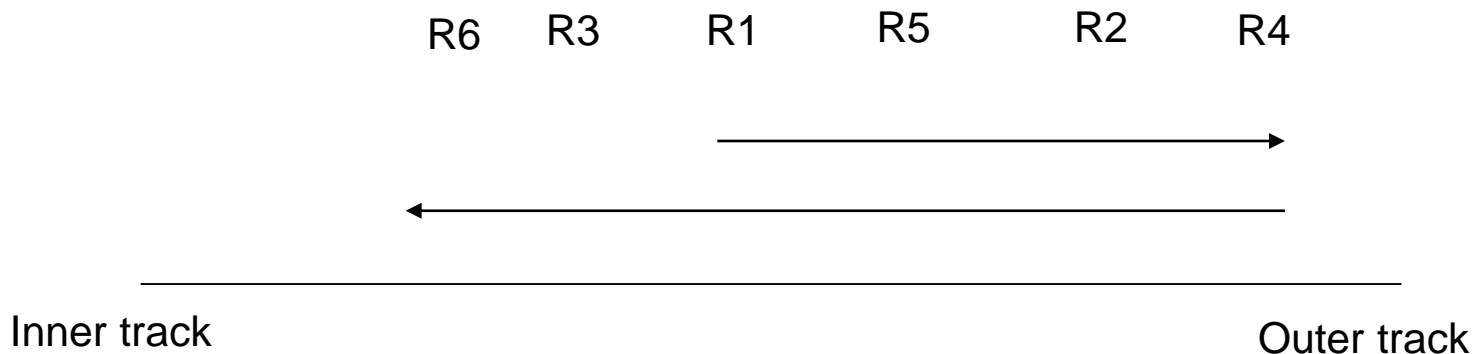  - ‣ Fiber Channel (FC2Gb or 4Gb): 256 to 512 MB/s

# 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

# Optimization of Disk-Block Access

- **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
        - Smaller blocks: more transfers from disk
        - Larger blocks:  more space wasted due to partially filled blocks
        - Typical block sizes today range from 4 to 16 kilobytes
- **Disk-arm-scheduling** algorithms order pending accesses to tracks so that disk arm movement is minimized
    - **elevator algorithm**:

R6     R3     R1     R5     R2     R4

Inner track                                                                    Outer track

# Optimization of Disk Block Access (Cont.)

- **File organization** – optimize block access time by organizing the blocks to correspond to how data will be accessed
    - E.g. Store related information on the same or nearby cylinders.
    - Files may get **fragmented** over time
        - E.g. if data is inserted to/deleted from the file
        - Or 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
    - Some systems have utilities to defragment the file system, in order to speed up file access

# Optimization of Disk Block Access (Cont.)

- **Nonvolatile 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 then writes to disk whenever the disk has no other requests or request has been pending for some time
  - Database operations that require data to be safely stored before continuing can continue without waiting for data to be written to disk
  - *Writes can be reordered to minimize disk arm movement*
- **Log disk** – a disk devoted to writing a sequential log of block updates
  - Used exactly like nonvolatile RAM
    - Write to log disk is very fast since no seeks are required
    - No need for special hardware (NV-RAM)
- File systems typically reorder writes to disk to improve performance
  - **Journaling file systems** write data in safe order to NV-RAM or log disk
  - Reordering without journaling: risk of corruption of file system data

# Flash Storage

- NOR flash vs NAND flash
- NAND flash
  - used widely for storage, since it is much cheaper than NOR flash
  - requires page-at-a-time read (page: 512 bytes to 4 KB)
  - transfer rate around 20 MB/sec
  - **solid state disks**: use multiple flash storage devices to provide higher transfer rate of 100 to 200 MB/sec
  - erase is very slow (1 to 2 millisecs)
    - erase block contains multiple pages
    - **remapp**ing of logical page addresses to physical page addresses avoids waiting for erase
      - **translation table** tracks mapping
        - also stored in a label field of flash page
      - remapping carried out by **flash translation layer**
    - after 100,000 to 1,000,000 erases, erase block becomes unreliable and cannot be used
      - **wear leveling**

# 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, 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)
  - Techniques for using redundancy to avoid data loss are critical with large numbers of disks
- Originally a cost-effective alternative to large, expensive disks
  - I in RAID originally stood for ``inexpensive''
  - Today RAIDs are used for their higher reliability and bandwidth.
    - The "I" is interpreted as independent

# 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

- **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 $500*10^6$ 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
    - 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.
    - 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 ($i$ 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

# RAID Levels

- Schemes to provide redundancy at lower cost by using disk striping combined with parity bits

  - Different 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 loss 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.
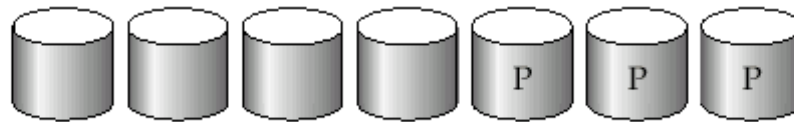
(a) RAID 0: nonredundant striping

(b) RAID 1: mirrored disks

# 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, 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
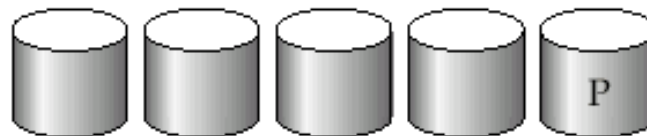
(d) RAID 3: bit-interleaved parity

# RAID 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.

  - Subsumes Level 2 (provides all its benefits, at lower cost).

- **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
  - 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)
    - ▸ 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.



(g) RAID 6: P + Q redundancy

# 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 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 most applications

# 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

# Hardware Issues

- **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 detected potentially corrupted blocks
        - Otherwise all blocks of disk must be read and compared with mirror/parity block

# Hardware Issues (Cont.)

- **Latent failures**: data successfully written earlier gets damaged

    - can result in data loss even if only one disk fails

- **Data scrubbing:**

    - continually scan for latent failures, and recover from copy/parity

- **Hot swapping**: replacement of disk while system is running, without power down

    - Supported by some hardware RAID systems,

    - reduces time to recovery, and improves availability 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

# Optical Disks

- **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

  - Blu-ray DVD: 27 GB (54 GB for double sided disk)

  - 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

# Magnetic Tapes

- 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, and 330 GB with Ampex helical scan format
  - Transfer rates from few to 10s of MB/s
- Tapes are cheap, but cost of drives is very high
- Very slow access time in comparison to magnetic and optical disks
  - limited to sequential access.
  - Some formats (Accelis) 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 used for very large capacity storage
  - Multiple petabyes ($10^{15}$ bytes)