# RAID

- Redundant Array of Independent Disks

- A set of physical disk drives viewed by the OS as a single logical drive

- Data are distributed across the physical drives. May improve performance.

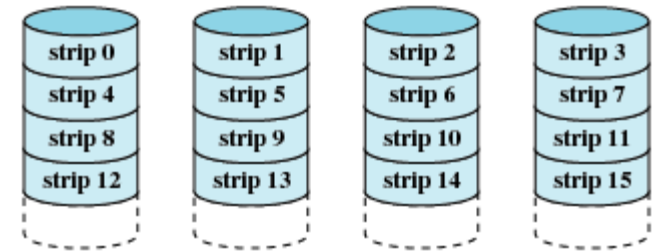- Redundant disk stores parity information. Recoverability, reliability.

# RAID Levels

**Table 6.3** RAID Levels

| Category | Level | Description | Disks Required | Data Availability | Large I/O Data Transfer Capacity | Small I/O Request Rate |
|---|---|---|---|---|---|---|
| Striping | 0 | Nonredundant | $N$ | Lower than single disk | Very high | Very high for both read and write |
| Mirroring | 1 | Mirrored | $2N$ | Higher than RAID 2, 3, 4, or 5; lower than RAID 6 | Higher than single disk for read; similar to single disk for write | Up to twice that of a single disk for read; similar to single disk for write |
| Parallel access | 2 | Redundant via Hamming code | $N + m$ | Much higher than single disk; comparable to RAID 3, 4, or 5 | Highest of all listed alternatives | Approximately twice that of a single disk |
| | 3 | Bit-interleaved parity | $N + 1$ | Much higher than single disk; comparable to RAID 2, 4, or 5 | Highest of all listed alternatives | Approximately twice that of a single disk |
| Independent access | 4 | Block-interleaved parity | $N + 1$ | Much higher than single disk; comparable to RAID 2, 3, or 5 | Similar to RAID 0 for read; significantly lower than single disk for write | Similar to RAID 0 for read; significantly lower than single disk for write |
| | 5 | Block-interleaved distributed parity | $N + 1$ | Much higher than single disk; comparable to RAID 2, 3, or 4 | Similar to RAID 0 for read; lower than single disk for write | Similar to RAID 0 for read; generally lower than single disk for write |
| | 6 | Block-interleaved dual distributed parity | $N + 2$ | Highest of all listed alternatives | Similar to RAID 0 for read; lower than RAID 5 for write | Similar to RAID 0 for read; significantly lower than RAID 5 for write |

*Note:* $N$ = number of data disks; $m$ proportional to $\log N$
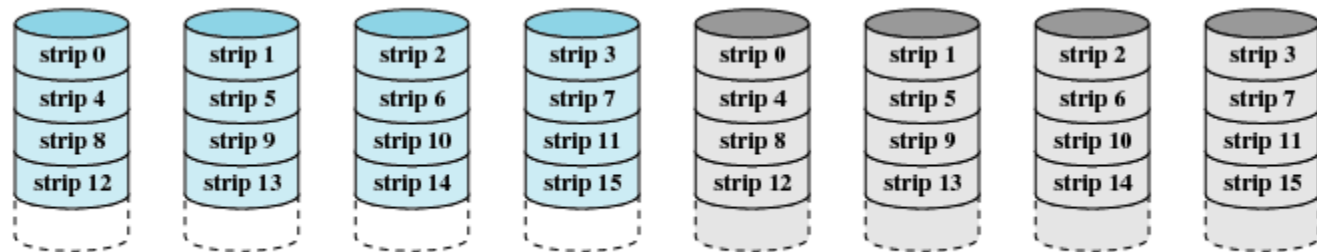
# RAID 0 (Non-redundant)

- The logical disk is divided into strips, mapped round robin to consecutive physical disks

- Improve performance in disk read/write

- Not fault tolerant

| strip 0 | strip 1 | strip 2 | strip 3 |
|---------|---------|---------|---------|
| strip 4 | strip 5 | strip 6 | strip 7 |
| strip 8 | strip 9 | strip 10 | strip 11 |
| strip 12 | strip 13 | strip 14 | strip 15 |

(a) RAID 0 (non-redundant)

3

# RAID 1 (Mirrored)

- Each disk is mirrored by another disk

- Good performance if the hardware supports concurrent read/write to the mirrored pair
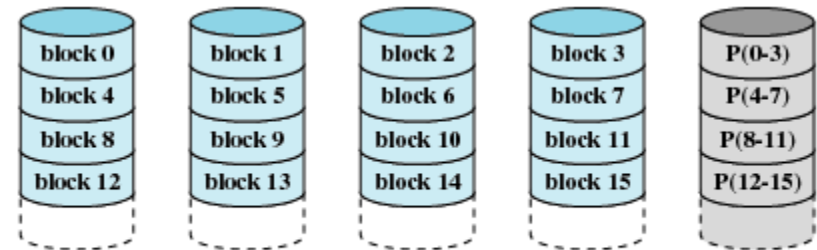
- Reliable, but expensive

| strip 0 | strip 1 | strip 2 | strip 3 | strip 0 | strip 1 | strip 2 | strip 3 |
|---------|---------|---------|---------|---------|---------|---------|---------|
| strip 4 | strip 5 | strip 6 | strip 7 | strip 4 | strip 5 | strip 6 | strip 7 |
| strip 8 | strip 9 | strip 10 | strip 11 | strip 8 | strip 9 | strip 10 | strip 11 |
| strip 12 | strip 13 | strip 14 | strip 15 | strip 12 | strip 13 | strip 14 | strip 15 |

(b) RAID 1 (mirrored)

4

# Parity strip

- Computed and updated at write, verified at read

- Every write results in two read and two write of strips

- A corrupted strip can be recovered

To compute the parity strip...
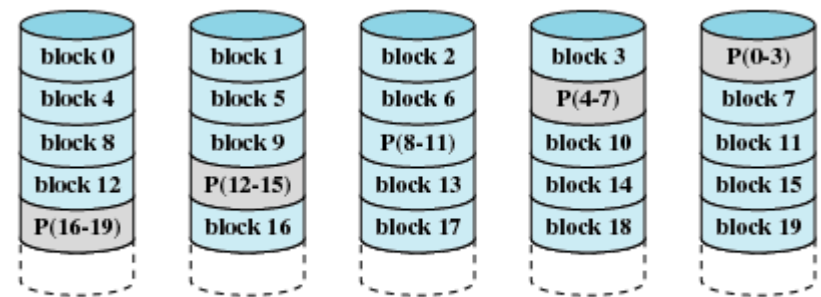$P(0\text{-}3) := b0 \oplus b1 \oplus b2 \oplus b3$

To recover the block 0...
$b0 = P(0\text{-}3) \oplus b1 \oplus b2 \oplus b3$

| block 0 | block 1 | block 2 | block 3 | P(0-3) |
|---------|---------|---------|---------|--------|
| block 4 | block 5 | block 6 | block 7 | P(4-7) |
| block 8 | block 9 | block 10 | block 11 | P(8-11) |
| block 12 | block 13 | block 14 | block 15 | P(12-15) |

**(e) RAID 4 (block-level parity)**

5

# RAID 5 (Block-level distributed parity)

- Having all parity strips on one disk may make it a bottleneck.  Instead, we can distribute the parity strips among the disks

- If a single disk fails, the system can regenerate the data lost

- Reliable. Good performance with special hardware

| block 0 | block 1 | block 2 | block 3 | P(0-3) |
|---------|---------|---------|---------|---------|
| block 4 | block 5 | block 6 | P(4-7) | block 7 |
| block 8 | block 9 | P(8-11) | block 10 | block 11 |
| block 12 | P(12-15) | block 13 | block 14 | block 15 |
| P(16-19) | block 16 | block 17 | block 18 | block 19 |

(f) RAID 5 (block-level distributed parity)

6

Table 6.4 RAID Comparison

| Level | Advantages | Disadvantages | Applications |
|---|---|---|---|
| 0 | I/O performance is greatly improved by spreading the I/O load across many channels and drives<br>No parity calculation overhead is involved<br>Very simple design<br>Easy to implement | The failure of just one drive will result in all data in an array being lost | Video production and editing<br>Image Editing<br>Pre-press applications<br>Any application requiring high bandwidth |
| 1 | 100% redundancy of data means no rebuild is necessary in case of a disk failure, just a copy to the replacement disk<br>Under certain circumstances, RAID 1 can sustain multiple simultaneous drive failures<br>Simplest RAID storage subsystem design | Highest disk overhead of all RAID types (100%) – inefficient | Accounting<br>Payroll<br>Financial<br>Any application requiring very high availability |
| 2 | Extremely high data transfer rates possible<br>The higher the data transfer rate required, the better the ratio of data disks to ECC disks<br>Relatively simple controller design compared to RAID levels 3, 4, & 5 | Very high ratio of ECC disks to data disks with smaller word sizes – inefficient<br>Entry level cost very high – requires very high transfer rate requirement to justify | No commercial implementations exist/not commercially viable |
| 3 | Very high read data transfer rate<br>Very high write data transfer rate<br>Disk failure has an insignificant impact on throughput<br>Low ratio of ECC (parity) disks to data disks means high efficiency | Transaction rate equal to that of a single disk drive at best (if spindles are synchronized)<br>Controller design is fairly complex | Video production and live streaming<br>Image editing<br>Video editing<br>Prepress applications<br>Any application requiring high throughput |
| 4 | Very high Read data transaction rate<br>Low ratio of ECC (parity) disks to data disks means high efficiency | Quite complex controller design<br>Worst write transaction rate and Write aggregate transfer rate<br>Difficult and inefficient data rebuild in the event of disk failure | No commercial implementations exist/not commercially viable |
| 5 | Highest Read data transaction rate<br>Low ratio of ECC (parity) disks to data disks means high efficiency<br>Good aggregate transfer rate | Most complex controller design<br>Difficult to rebuild in the event of a disk failure (as compared to RAID level 1) | File and application servers<br>Database servers<br>Web, e-mail, and news servers<br>Intranet servers<br>Most versatile RAID level |
| 6 | Provides for an extremely high data fault tolerance and can sustain multiple simultaneous drive failures | More complex controller design<br>Controller overhead to compute parity addresses is extremely high | Perfect solution for mission critical applications |

7

# Block-oriented disk

- Disk is block-oriented. One sector is read/written at a time.
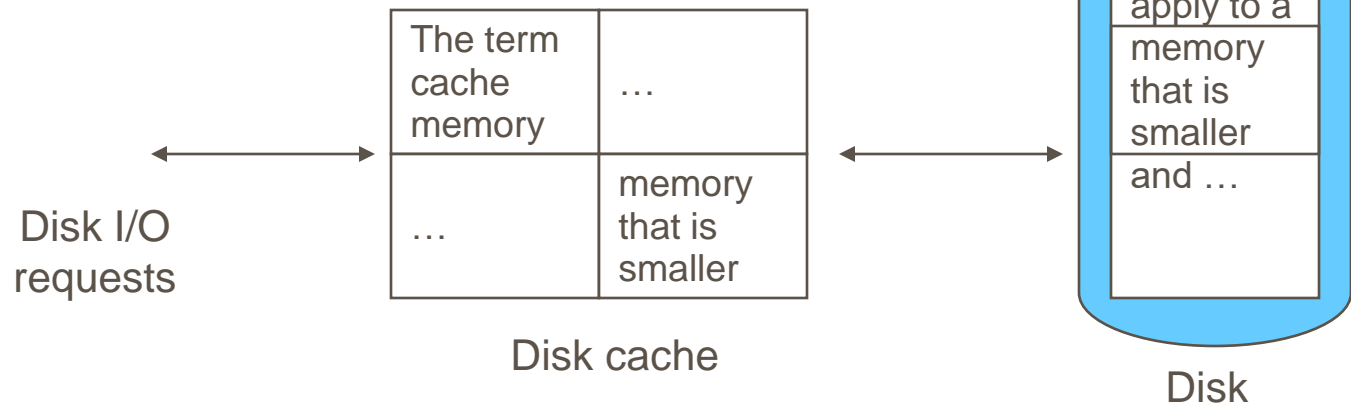
- In PC, a sector is 512 byte

The term cache memory is usually used to apply to a memory that is smaller and …

```
while (!feof(F)) {
  // read one char
  fscanf(F, "%c", &c);
  …
}
```

# Disk Cache

- Buffer in main memory for disk sectors
- Contains a copy of some of the sectors

Disk I/O requests

| The term cache memory | … |
| --- | --- |
| … | memory that is smaller |

Disk cache

The term cache memory is usually used to apply to a memory that is smaller and …

Disk

# Disk Cache, Hit and Miss

- When an I/O request is made for a particular sector, the OS checks whether the sector is in the disk cache.
  - If so, (cache hit), the request is satisfied via the cache.
  - If not (cache miss), the requested sector is read into the disk cache from the disk.