- ■ **Disk Data Model**

  a sequence of $r$ same sized **stripe units** ("scsi" interface) with indicies 0,1,$\cdots$,r-1. Stripe units may be (groups of) sectors, tracks, cylinders, or $\cdots$.

- ■ **User Data Model**

  a sequence of same sized data stripe units with indicies 0,1,2,$\cdots$ . Data stripe units typically created by file system, database management system, or other (sophisticated) user program.

- ■ **Reliability Groups a.k.a. Stripes**

  user data is partitioned into fixed sized groups with each group containing additional redundancy information.

- **Disk Array Data Layouts**

Data Layouts: map stripe unit indicies to disk sectors, tracks, or cylinders $\cdots$

| | |
|---|---|
| $n$ disks | $r$ stripe units/disk |
| $k$ stripe width | $b$ number of stripes |
| $g$ groups | $n = k \cdot g$ |

$m$ user data stripe units/stripe

$c = k - m$ redundant stripe units/stripe

Total number of data stripe units $b(k - c)$.

Total number of parity stripe units $b \cdot c$.

The total number of stripe units within for a completely filled/utilized disk array is $bk = nr$.

$$\mathbf{d}_0, \mathbf{d}_1, \mathbf{d}_2, \mathbf{d}_3, \cdots \qquad \text{user data}$$

$$\{\mathbf{d}_0, \mathbf{d}_1, \mathbf{d}_2, \mathbf{c}_0\} \qquad \text{stripes}$$

$$\{\mathbf{d}_3, \mathbf{d}_4, \mathbf{d}_5, \mathbf{c}_1\}$$

$$\{\mathbf{d}_6, \mathbf{d}_7, \mathbf{d}_8, \mathbf{c}_2\}$$

$$\ddots \qquad \mathbf{c}_i \text{ redundant data}$$

data stripe unit indicies $DSUI = \{0, 1, 2, \dots \mathbf{b(k\text{-}c)\text{-}1}\}$

*disk*: $\quad DSUI \mapsto \{0, 1, 2, \dots \mathbf{n\text{-}1}\}$

*offset*: $DSUI \mapsto \{0, 1, 2, \dots \mathbf{r\text{-}1}\}$

*checkDisk*: $\quad DSUI \times \{0, 1, \dots, \mathbf{c\text{-}1}\} \mapsto \{0, 1, 2, \dots, \mathbf{n\text{-}1}\}$

*checkOffset*: $DSUI \times \{0, 1, \dots, \mathbf{c\text{-}1}\} \mapsto \{0, 1, 2, \dots, \mathbf{r\text{-}1}\}$

- **Layout Taxonomy**

|  |  | typically |
|---|---|---|
| **Level 0** | **just a bunch of disks**<br>**JBOD, no redundancy** | $m = k = 1$ |
| **Level 1** | **mirroring** | $m = 1, k = 2$ |
| **Level 2** | **fine-grained interleaving**<br>**with ECC error correction** | $m = 10, k = 14$<br>$m = 20, k = 25$ |
| **Level 3** | **fine-grained interleaving**<br>**with dedicated parity disk** | $k = m+1$ |
| **Level 4** | **stripe unit interleaving;**<br>**dedicated parity disk** | $k = m+1$ |
| **Level 5** | **stripe-unit interleaving;**<br>**distributed parity** | $k = m+1$ |
| **Level 6** | **Level 5 with additional**<br>**redundant stripe-units;**<br>**typically one more** | $k = m+2$ |

- **Workloads**

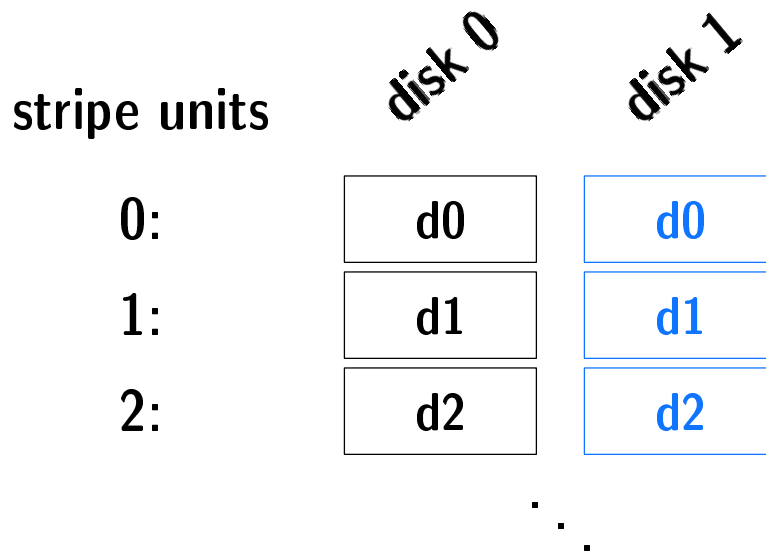| | |
|---|---|
| *Large Operations* | Parallel read or write accesses one stripe unit from each disk; high data transfer rates obtained. |
| *Small Operations* | Independent read or write accesses one data stripe unit from each disk; high numbers of i/o operations obtained. |



- **Operations**  small reads, writes, read-modify-writes (r m w) & large reads, writes, read-modify-writes.

- **Relative Efficiency**  $\dfrac{\text{RAID operations / sec.}}{\text{single disk operations / sec.}}$

- **RAID Level 1**     **mirroring**

stripe units

disk 0     disk 1

| | disk 0 | disk 1 |
|---|---|---|
| 0: | d0 | d0 |
| 1: | d1 | d1 |
| 2: | d2 | d2 |

$disk(a) = 0$

$offset(a) = a$

$checkDisk(a) = 1$

$checkOffset(a) = a$

- ## RAID Level 1     mirroring

$g$ **stripes**      $n = 2 \cdot g$ **disks**

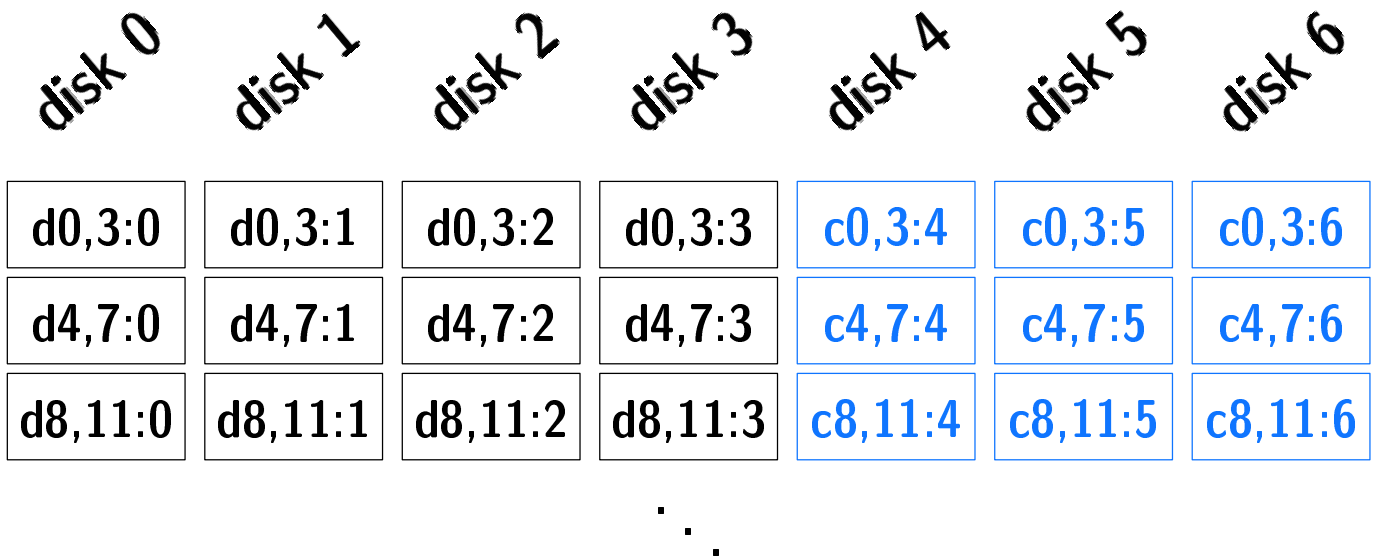$k = 2$, $m = 1$

| | | Relative Efficiency |
|---|---|:---:|
| small | read: | $2g$ |
| | write: | $g$ |
| | r m w: | $9g/8$ |
| large | read: | $2g/s$ |
| | write: | $g/s$ |
| | r m w: | $9g/8s$ |

$s$ slowdown:    Within large operations, all disks complete their individual tasks to finish the operation; $s$ is less than 2.

■ **RAID Level 2**     **fine-grain interleaving with ECC**

| disk 0 | disk 1 | disk 2 | disk 3 | disk 4 | disk 5 | disk 6 |
|--------|--------|--------|--------|--------|--------|--------|
| d0,3:0 | d0,3:1 | d0,3:2 | d0,3:3 | c0,3:4 | c0,3:5 | c0,3:6 |
| d4,7:0 | d4,7:1 | d4,7:2 | d4,7:3 | c4,7:4 | c4,7:5 | c4,7:6 |
| d8,11:0 | d8,11:1 | d8,11:2 | d8,11:3 | c8,11:4 | c8,11:5 | c8,11:6 |

$disk(a) = \{\, 0, 1, ..., m\text{-}1 \,\}$

$offset(a) = a/m$

$checkDisk(a) = \{\, m, m\text{+}1, ..., k\text{-}1 \,\}$

$checkOffset(a) = a/m$

- **RAID Level 2**      **fine-grain interleaving with ECC**

$$g \ \text{stripes} \qquad n = k \cdot g \ \text{disks}$$

$$k = m + c \ \text{stripe width}$$

This level attempts to provide good performance with less redundant data.

|  |  | Relative Efficiency |
|---|---|---|
| small | read: | $g/s$ |
|  | write: | $g/2s$ |
|  | r m w: | $g/s$ |
| large | read: | $g \cdot m/s$ |
|  | write: | $g \cdot m/s$ |
|  | r m w: | $g \cdot m/s$ |

Level 2 redundancy not needed; disk internal ecc guarantees one bit error in $10^{14}$.

- **RAID Level 3**     fine-grained interleaving, dedicated parity disk

|  | disk 0 | disk 1 | disk 2 | disk 3 | disk 4 |
|---|---|---|---|---|---|
| 0: | d0,3:0 | d0,3:1 | d0,3:2 | d0,3:3 | c0,3 |
| 1: | d4,7:0 | d4,7:1 | d4,7:2 | d4,7:3 | c4,7 |
| 2: | d8,11:0 | d8,11:1 | d8,11:2 | d8,11:3 | c8,11 |

$$disk(a) = \{ 0, 1, ..., \text{m-1} \}$$

$$offset(a) = \text{a/m}$$

$$checkDisk(a) = \text{k-1} = \text{m}$$

$$checkOffset(a) = \text{a/m}$$

- **RAID Level 3** — fine-grained interleaving, dedicated parity disk

$g$ stripes $\qquad n = k \cdot g$ disks

$k = m + 1$ stripe width

This level provides good performance with less redundant data. Level 2 and Level 3 performances are identical with one redundant disk.

|  |  | Relative Efficiency |
|---|---|:---:|
| small | read: | $g/s$ |
|  | write: | $g/2s$ |
|  | r m w: | $g/s$ |
| large | read: | $g \cdot m / s$ |
|  | write: | $g \cdot m / s$ |
|  | r m w: | $g \cdot m / s$ |

# RAID XI

Un-interleave data to improve small operations.

- **RAID Level 4**  **stripe-unit interleaving, dedicated parity disk**

|  | disk 0 | disk 1 | disk 2 | disk 3 | disk 4 |
|---|---|---|---|---|---|
| 0: | d0 | d1 | d2 | d3 | c0,3 |
| 1: | d4 | d5 | d6 | d7 | c4,7 |
| 2: | d8 | d9 | d10 | d11 | c8,11 |

$\ddots$

$disk(a) = a\%m$

$offset(a) = a/m$

$checkDisk(a) = k\text{-}1 = m$

$checkOffset(a) = a/m$

- **RAID Level 4**     stripe-unit interleaving, dedicated parity disk

$$g \ \text{stripes} \qquad n = k \cdot g \ \text{disks}$$

$$k = m + 1 \ \text{stripe width}$$

Level 4 provides better small operation performance than Level 3.

|         |        | Relative Efficiency |
|---------|--------|---------------------|
| small   | read:  | $g \cdot m$         |
|         | write: | $g / 2$             |
|         | r m w: | $g$                 |
| large   | read:  | $g \cdot m / s$     |
|         | write: | $g \cdot m / s$     |
|         | r m w: | $g \cdot m / s$     |

- **RAID Level 5**      stripe-unit interleaving,
                        distributed parity stripe units

Level 4 with distributed parity stripe units.

|     | disk 0 | disk 1 | disk 2 | disk 3 | disk 4 |
|-----|--------|--------|--------|--------|--------|
| 0:  | d0     | d1     | d2     | d3     | p0,3   |
| 1:  | d4     | d5     | d6     | p4,7   | d7     |
| 2:  | d8     | d9     | p8,11  | d10    | d11    |
| 3:  | d12    | p12,15 | d13    | d14    | d15    |
| 4:  | p16,19 | d16    | d17    | d18    | d19    |

$\cdots$

$$P = m \text{ - } Q\%k \quad R = a\%m \quad Q = a/m$$

$$disk(a) = \begin{cases} R & \text{if } R < P \\ R + 1 & \text{otherwise} \end{cases}$$

$$offset(a) = Q$$

$$checkDisk(a) = P$$

$$checkOffset(a) = Q$$

- **RAID Level 5**     stripe-unit interleaving,
                      distributed parity stripe units

$g$ stripes     $n = k \cdot g$ disks

$k = m + 1$ stripe width

Level 5 provides better than Level 4 small operation performance.

|        |        | Relative Efficiency |
|--------|--------|---------------------|
| small  | read:  | $g \cdot k$         |
|        | write: | $g \cdot k / 4$     |
|        | rmw:   | $g \cdot k / 2$     |
| large  | read:  | $g \cdot m / s$     |
|        | write: | $g \cdot m / s$     |
|        | rmw:   | $g \cdot m / s$     |

- **RAID Level 5      Left Symmetric layout**

|  | disk 0 | disk 1 | disk 2 | disk 3 | disk 4 |
|---|---|---|---|---|---|
| 0: | d0 | d1 | d2 | d3 | p0,3 |
| 1: | d5 | d6 | d7 | p4,7 | d4 |
| 2: | d10 | d11 | p8,11 | d8 | d9 |
| 3: | d15 | p12,15 | d12 | d13 | d14 |
| 4: | p16,19 | d16 | d17 | d18 | d19 |

$$P = m - Q\%k \quad Q = a/m$$

$$disk(a) = a\%k$$

$$offset(a) = Q$$

$$checkDisk(a) = P$$

$$checkOffset(a) = Q$$

- **RAID Level 6**    stripe-unit interleaving, distributed check units

Level 5 with additional check stripe units.

|  | disk 0 | disk 1 | disk 2 | disk 3 | disk 4 | disk 5 |
|---|---|---|---|---|---|---|
| 0: | d0 | d1 | d2 | d3 | c0,3:0 | c0,3:1 |
| 1: | d4 | d5 | d6 | c4,7:0 | c4,7:1 | d7 |
| 2: | d8 | d9 | c8,11:0 | c8,11:1 | d10 | d11 |
| 3: | d12 | c12,15:0 | c12,15:1 | d13 | d14 | d15 |
| 4: | c16,19:0 | c16,19:1 | d16 | d17 | d18 | d19 |

$$P = m - Q\%k \quad R = a\%m \quad Q = a/m$$

$$disk(a) = \begin{cases} R & \text{if } R < P \\ R + c & \text{otherwise} \end{cases}$$

$$offset(a) = Q$$

$$checkDisk(a) = \{P, P+1, ..., P+(c-1)\}$$

$$checkOffset(a) = Q$$

- **RAID Level 6**      stripe-unit interleaving,
  distributed check stripe units

$g$   stripes

$n = k \cdot g$   disks

$k = m{+}2$   typical stripe width

Level 6 provides better than Level 4 small operation performance as well as better reliability.

|  |  | Relative Efficiency |
|---|---|---|
| small | read: | $g \cdot k$ |
|  | write: | $g \cdot k / 4$ |
|  | r m w: | $g \cdot k / 2$ |
| large | read: | $g \cdot m / s$ |
|  | write: | $g \cdot m / s$ |
|  | r m w: | $g \cdot m / s$ |

- ## RAID Level 6     "Left Symmetric" layout

|  | disk 0 | disk 1 | disk 2 | disk 3 | disk 4 | disk 5 |
|---|---|---|---|---|---|---|
| 0: | d0 | d1 | d2 | d3 | c0,3:0 | c0,3:1 |
| 1: | d6 | d7 | c4,7:0 | c4,7:1 | d4 | d5 |
| 2: | c8,11:0 | c8,11:1 | d8 | d9 | d10 | d11 |

$\cdots$

**assume** $m = c\lambda$ ;    $P = m - (Q\%k)\lambda$    $Q = a/m$

$disk(a) = a\%k$

$offset(a) = Q$

$checkDisk(a) = \{\, P, P+1, ..., P+(c\text{-}1)\,\}$

$checkOffset(a) = Q$