# CS343 - Operating Systems

## Storage Arrays & RAID Implementation

Dr. John Jose

Assistant Professor

Department of Computer Science & Engineering

Indian Institute of Technology Guwahati, Assam.

http://www.iitg.ac.in/johnjose/

# Disk Storage Systems Management

❖ Disk Organization & Structure

❖ Disk Attachment

❖ Disk Scheduling

❖ Disk Management

❖ Swap-Space Management

❖ Storage Arrays

❖ RAID Structure

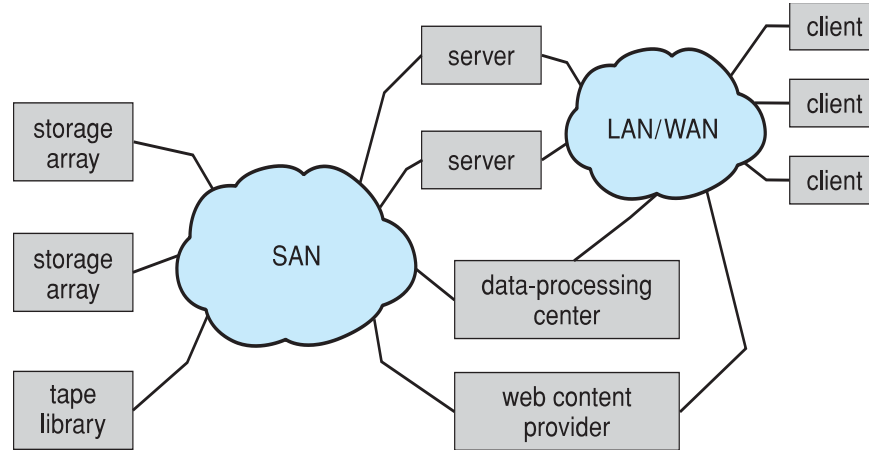❖ Stable-Storage Implementation

# Objectives

- ❖ To describe the physical structure of secondary storage devices and its effects on the uses of the devices

- ❖ To explain the performance characteristics of mass-storage devices

- ❖ To evaluate disk scheduling algorithms

- ❖ To discuss operating-system services provided for mass storage, including RAID

# Storage Array

❖ Attach multiple disks -  arrays of disks

❖ Storage array has controller that provides features to attached hosts

    ❖ A few to thousands of disks

    ❖ Ports to connect hosts to array

    ❖ Memory, controlling software

    ❖ Support RAID, hot spares, hot swap
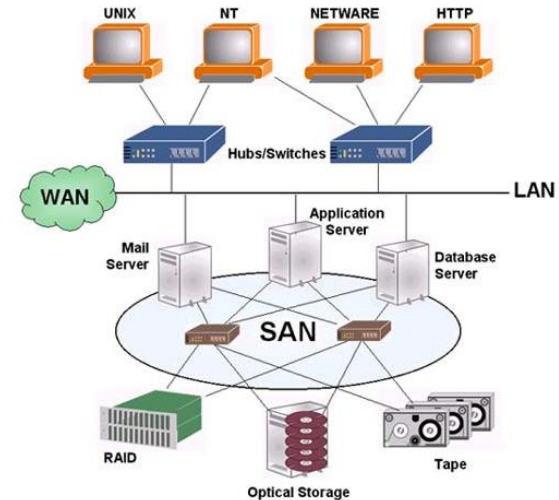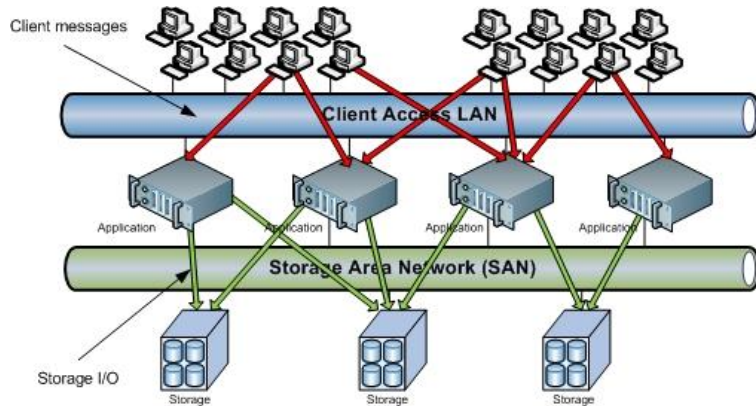
    ❖ Shared storage → more efficiency

# Storage Area Network

❖ Common in large storage environments

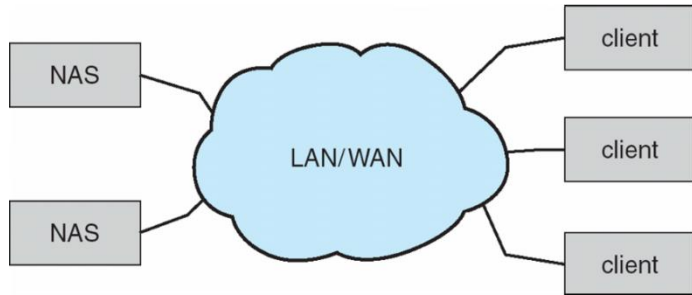❖ Multiple hosts attached to multiple storage arrays - flexible

# Storage Area Network

❖ SAN is one or more storage arrays

    ❖ Connected to one or more Fibre Channel switches

❖ Hosts also attach to the switches

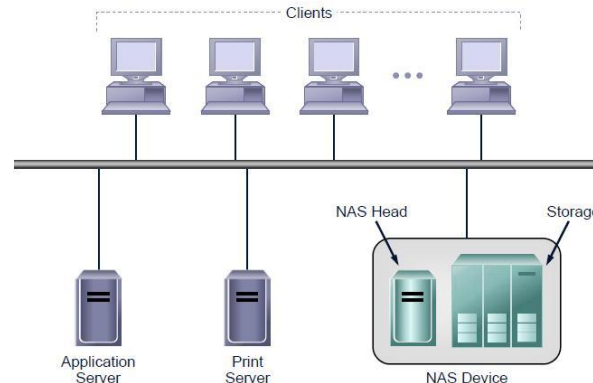❖ Easy to add or remove storage, add new host and allocate it storage

❖ Network-attached storage (**NAS**) is storage made available over a network rather than over a local connection (such as a bus)

  ❖ Remotely attaching to file systems

❖ Implemented via remote procedure calls (RPCs) between host and storage over typically standard computer network protocols.
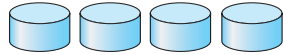
# RAID Structure

- ❖ RAID – redundant array of inexpensive disks

- ❖ Multiple disk drives provides reliability via **redundancy**

- ❖ Use of multiple disks working cooperatively

- ❖ Increases the **mean time to failure**

- ❖ Frequently combined with **NVRAM** to improve write performance

- ❖ Disk **striping** (**RAID 0**) uses a group of disks as one storage unit

- ❖ RAID is arranged into six different levels

# RAID Structure

❖ RAID schemes improve performance and improve the reliability of the storage system by storing redundant data

  ❖ **Mirroring** or **shadowing** (**RAID 1**) keeps duplicate of each disk

  ❖ Striped mirrors (**RAID 1+0**) or mirrored stripes (**RAID 0+1**) provides high performance and high reliability

  ❖ **Block interleaved parity** (**RAID 4, 5, 6**) uses much less redundancy

❖ RAID within a storage array can still fail if the array fails, so automatic **replication** of the data between arrays is common

❖ Frequently, a small number of **hot-spare** disks are left unallocated, automatically replacing a failed disk and having data rebuilt onto them
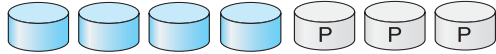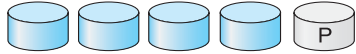
# RAID Levels



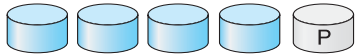(a) RAID 0: non-redundant striping.

(b) RAID 1: mirrored disks.

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

(d) RAID 3: bit-interleaved parity.

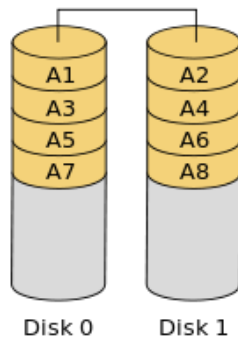(e) RAID 4: block-interleaved parity.
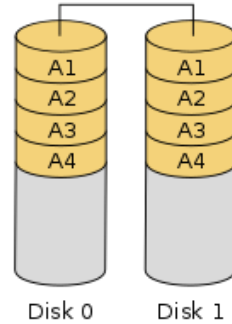
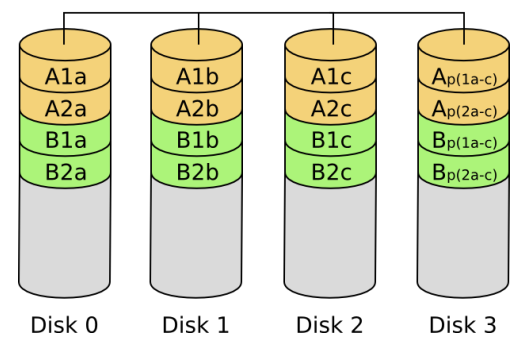(f) RAID 5: block-interleaved distributed parity.
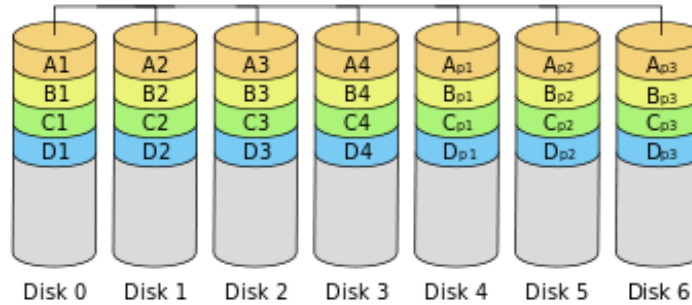
(g) RAID 6: P + Q redundancy.

## RAID 0

| Disk 0 | Disk 1 |
|--------|--------|
| A1 | A2 |
| A3 | A4 |
| A5 | A6 |
| A7 | A8 |

## RAID 1

| Disk 0 | Disk 1 |
|--------|--------|
| A1 | A1 |
| A2 | A2 |
| A3 | A3 |
| A4 | A4 |

## RAID 3

| Disk 0 | Disk 1 | Disk 2 | Disk 3 |
|--------|--------|--------|--------|
| A1a | A1b | A1c | $A_{p(1a-c)}$ |
| A2a | A2b | A2c | $A_{p(2a-c)}$ |
| B1a | B1b | B1c | $B_{p(1a-c)}$ |
| B2a | B2b | B2c | $B_{p(2a-c)}$ |

## RAID 2

| Disk 0 | Disk 1 | Disk 2 | Disk 3 | Disk 4 | Disk 5 | Disk 6 |
|--------|--------|--------|--------|--------|--------|--------|
| A1 | A2 | A3 | A4 | $A_{p1}$ | $A_{p2}$ | $A_{p3}$ |
| B1 | B2 | B3 | B4 | $B_{p1}$ | $B_{p2}$ | $B_{p3}$ |
| C1 | C2 | C3 | C4 | $C_{p1}$ | $C_{p2}$ | $C_{p3}$ |
| D1 | D2 | D3 | D4 | $D_{p1}$ | $D_{p2}$ | $D_{p3}$ |

## RAID 4

| Disk 0 | Disk 1 | Disk 2 | Disk 3 |
|--------|--------|--------|--------|
| A1 | A2 | A3 | $A_p$ |
| B1 | B2 | B3 | $B_p$ |
| C1 | C2 | C3 | $C_p$ |
| D1 | D2 | D3 | $D_p$ |

# RAID Levels



(a) RAID 0: non-redundant striping.

(b) RAID 1: mirrored disks.

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

(d) RAID 3: bit-interleaved parity.

(e) RAID 4: block-interleaved parity.

(f) RAID 5: block-interleaved distributed parity.
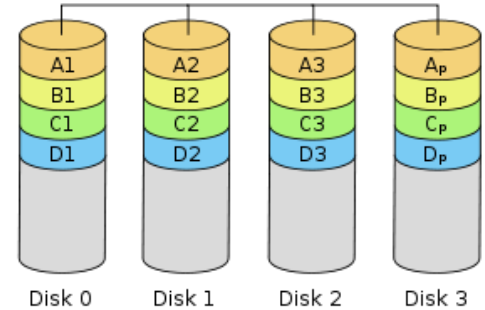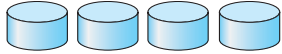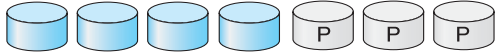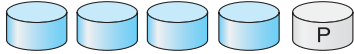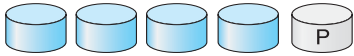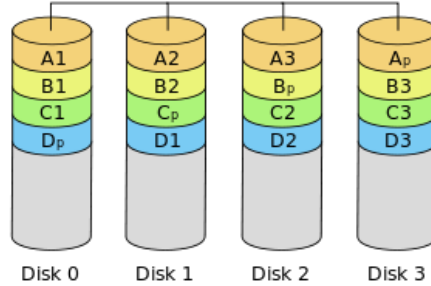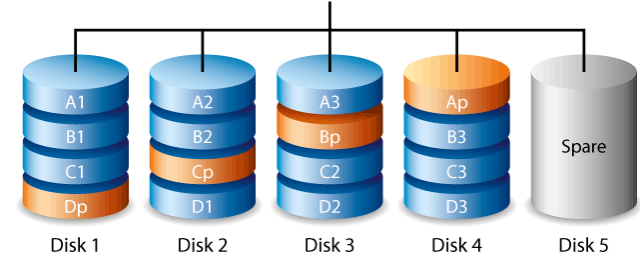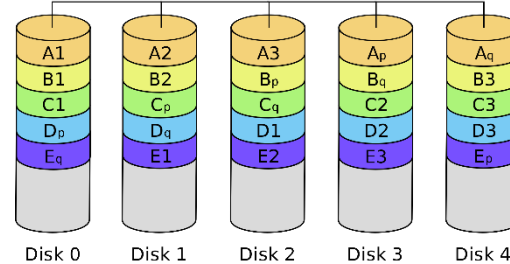
(g) RAID 6: P + Q redundancy.

RAID 5
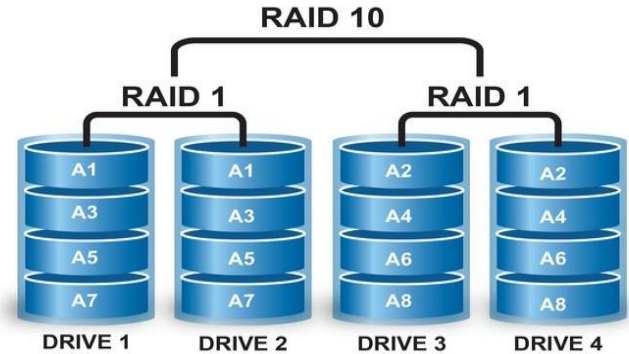
Disk 0    Disk 1    Disk 2    Disk 3

RAID 5+Spare

Disk 1    Disk 2    Disk 3    Disk 4    Disk 5

RAID 6

Disk 0    Disk 1    Disk 2    Disk 3    Disk 4

# RAID Levels



a) RAID 0 + 1 with a single disk failure.

b) RAID 1 + 0 with a single disk failure.

RAID 0 1

RAID 1

RAID 0 — RAID 0

| | | | |
|---|---|---|---|
| A1 | A2 | A1 | A2 |
| A3 | A4 | A3 | A4 |
| A5 | A6 | A5 | A6 |
| A7 | A8 | A7 | A8 |
| Disk 1 | Disk 2 | Disk 3 | Disk 4 |

RAID 10

RAID 1 — RAID 1

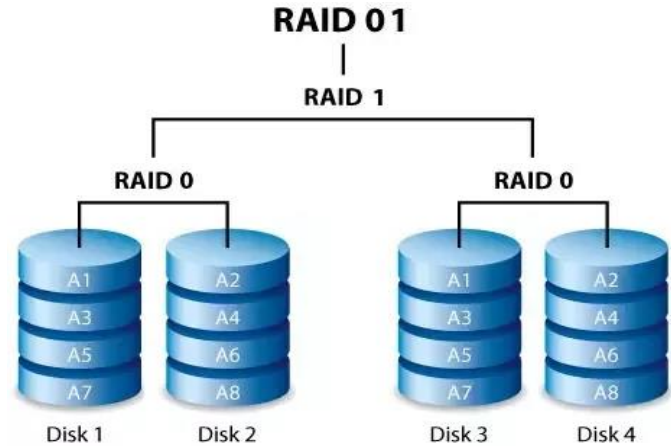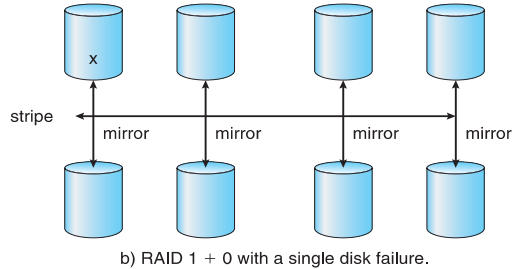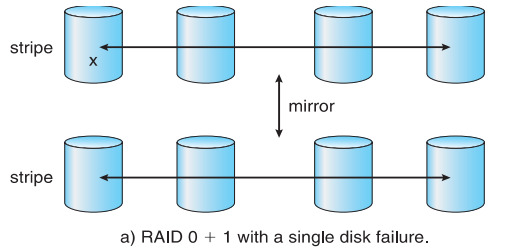| | | | |
|---|---|---|---|
| A1 | A1 | A2 | A2 |
| A3 | A3 | A4 | A4 |
| A5 | A5 | A6 | A6 |
| A7 | A7 | A8 | A8 |
| DRIVE 1 | DRIVE 2 | DRIVE 3 | DRIVE 4 |

# Stable-Storage Implementation

❖ Write-ahead log scheme requires stable storage

❖ Stable storage means data is never lost (due to failure, etc)

❖ To implement stable storage:

  ❖ Replicate information on more than one nonvolatile storage media with independent failure modes

  ❖ Update information in a controlled manner to ensure that we can recover the stable data after any failure during data transfer or recovery

# Stable-Storage Implementation

❖ Disk write has 1 of 3 outcomes

  ❖ **Successful completion -** The data were written correctly on disk

  ❖ **Partial failure -** A failure occurred in the midst of transfer, so only some of the sectors were written with the new data, and the sector being written during the failure may have been corrupted

  ❖ **Total failure -** The failure occurred before the disk write started, so the previous data values on the disk remain intact

❖ If failure occurs during block write, recovery procedure restores block to consistent state

  ❖ System maintains 2 physical blocks per logical block

    ❖ Write to 1st physical, When successful, write to 2nd physical

    ❖ Declare complete only after second write completes successfully

**johnjose@iitg.ac.in**
**http://www.iitg.ac.in/johnjose/**