# Database Applications (15-415)

# DBMS Internals- Part I
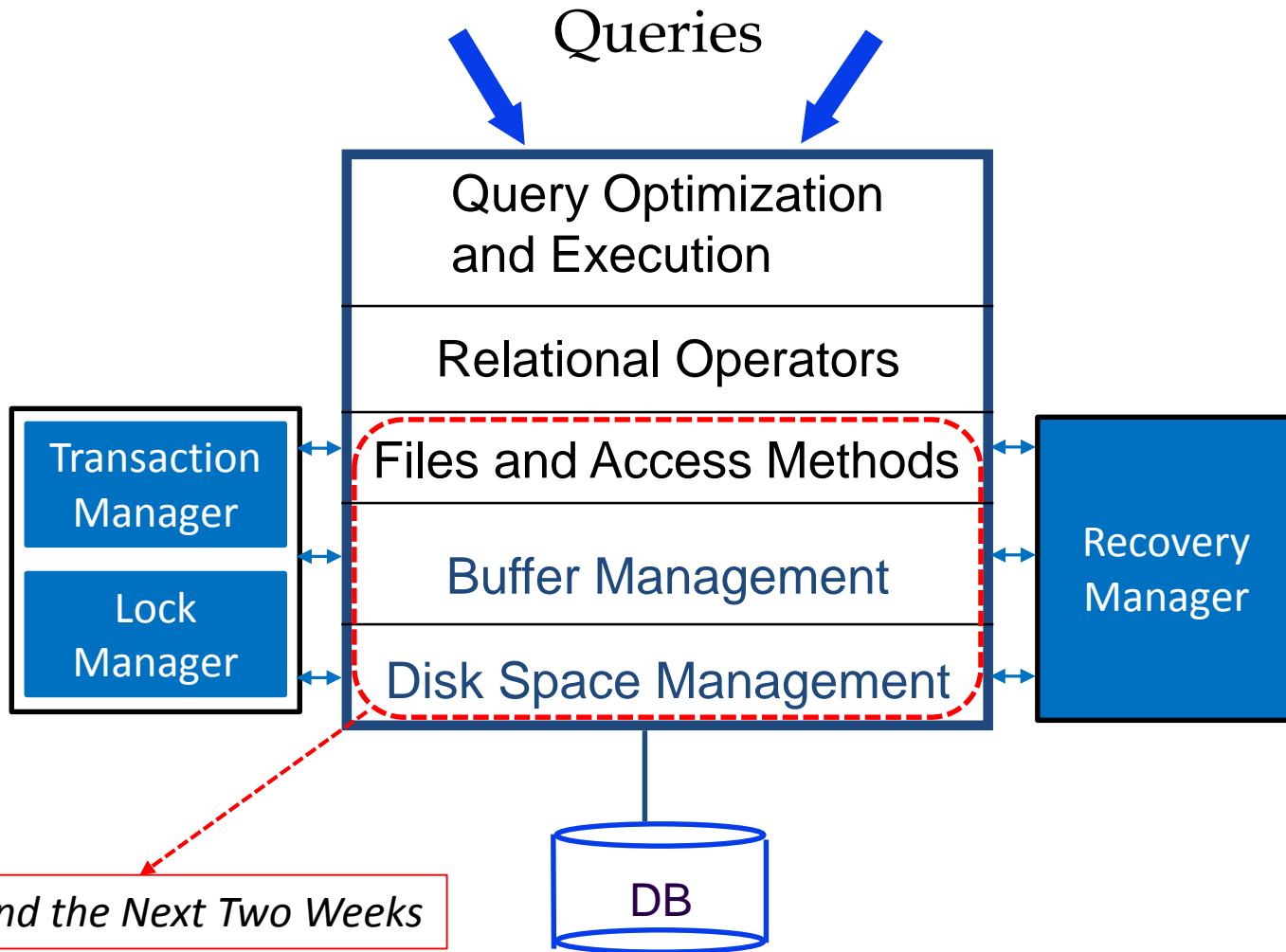# Lecture 10, February 15, 2015
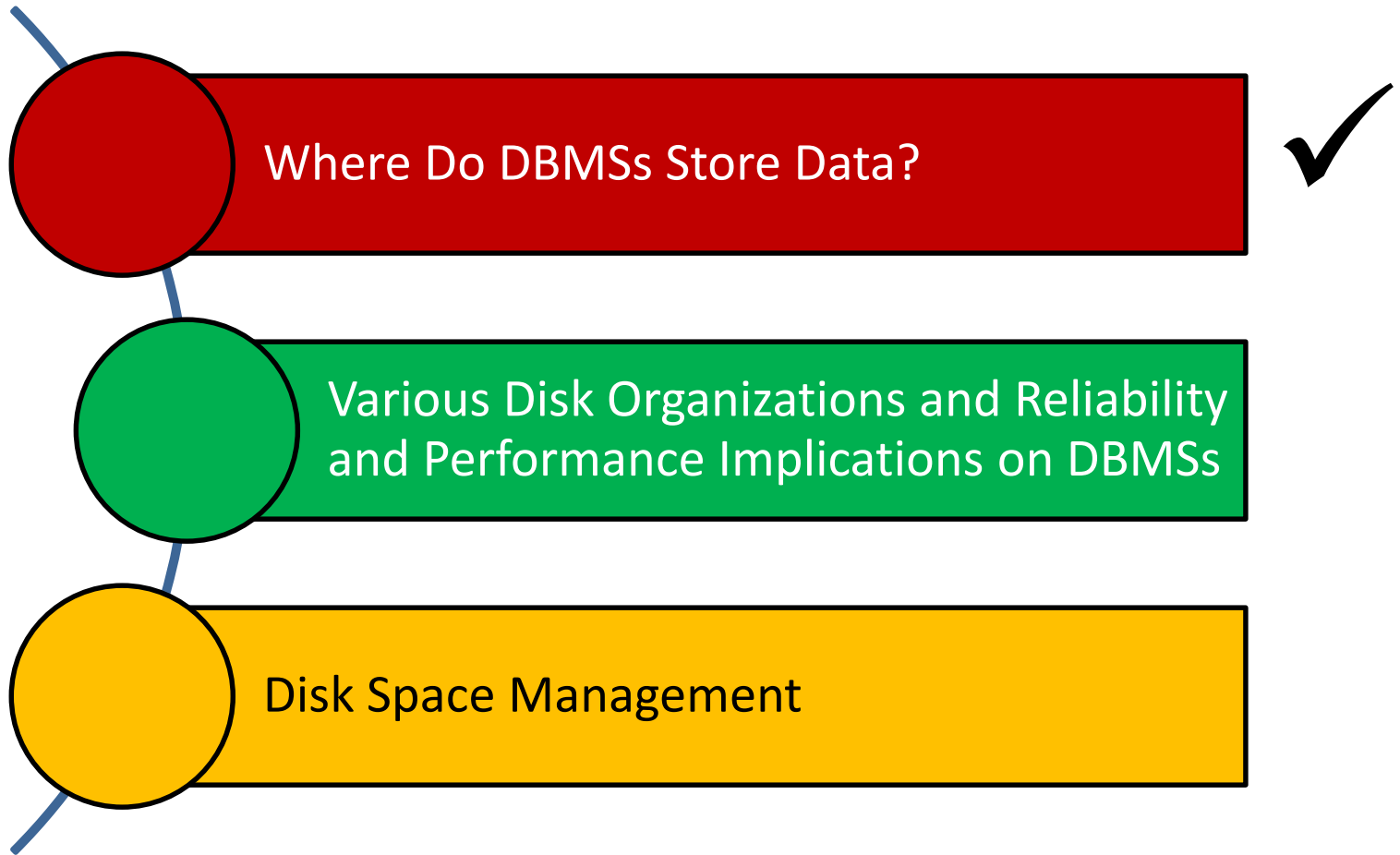
## Mohammad Hammoud

# Today…

- **Last Session:**
  - SQL- Part III & Quiz I

- **Today's Session:**
  - DBMS Internals- Part I
    - Background on Disks and Disk Arrays
    - Disk Space Management
    - Buffer Management

- **Announcements:**
  - Quiz I grades are out
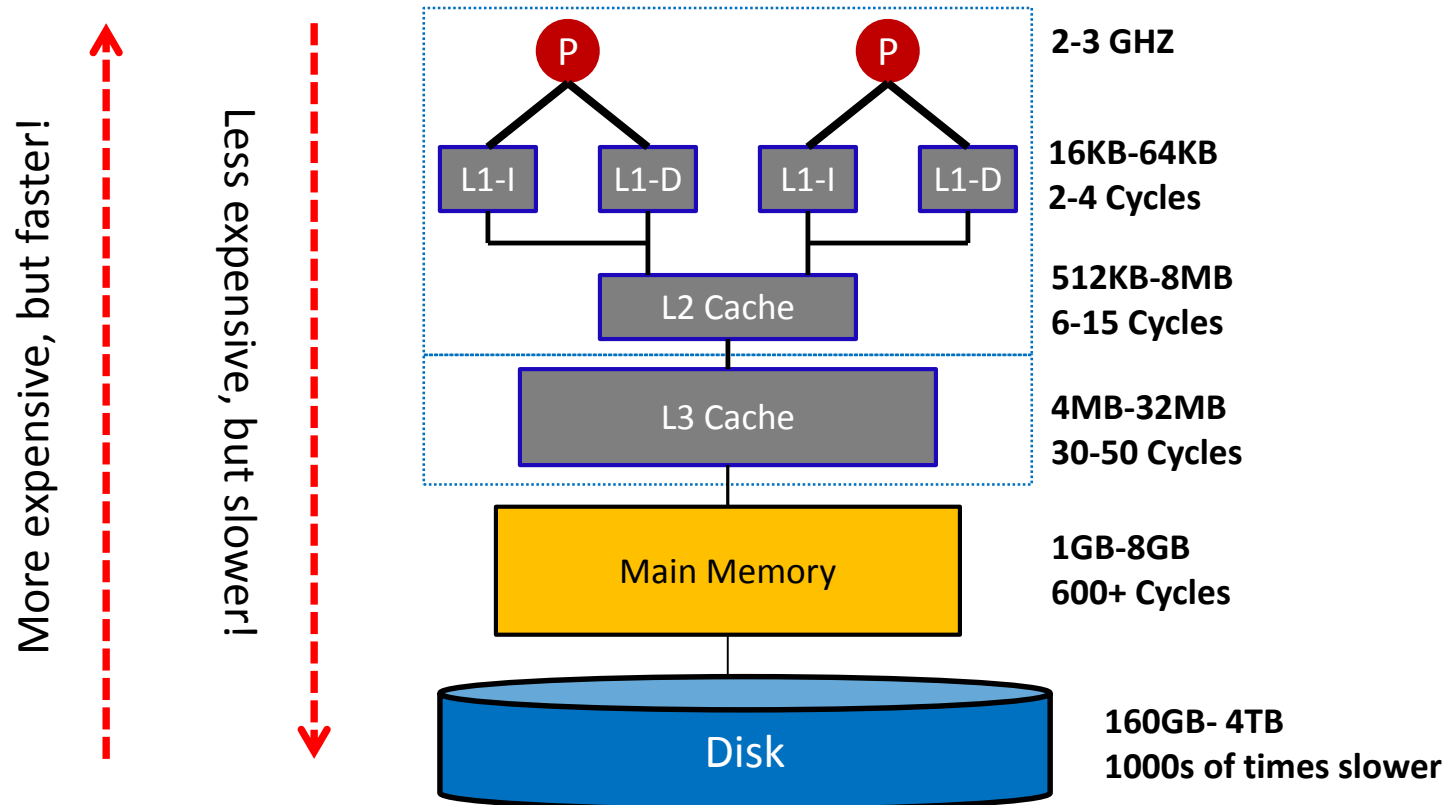  - Project 1 is due on Tuesday, Feb 17 by midnight

# DBMS Layers



Queries

Query Optimization and Execution

Relational Operators

Files and Access Methods

Buffer Management

Disk Space Management

Transaction Manager

Lock Manager

Recovery Manager

DB

*Today and the Next Two Weeks*

# Outline



Where Do DBMSs Store Data? ✔

Various Disk Organizations and Reliability and Performance Implications on DBMSs

Disk Space Management

# The Memory Hierarchy

- Storage devices play an important role in database systems

- How systems arrange storage?

# Where to Store Data?

- Where do DBMSs store information?
  - DBMSs store large amount of data (what about Big Data?)– as of now, we assume *centralized* DBMSs

  - Typically, buying enough memory to store all data is prohibitively expensive (let alone that memories are *volatile*)

  - Thus, databases are usually stored on disks (or tapes for backups)

# But, Is Memory Gone?

- Data must be brought into memory to be processed!
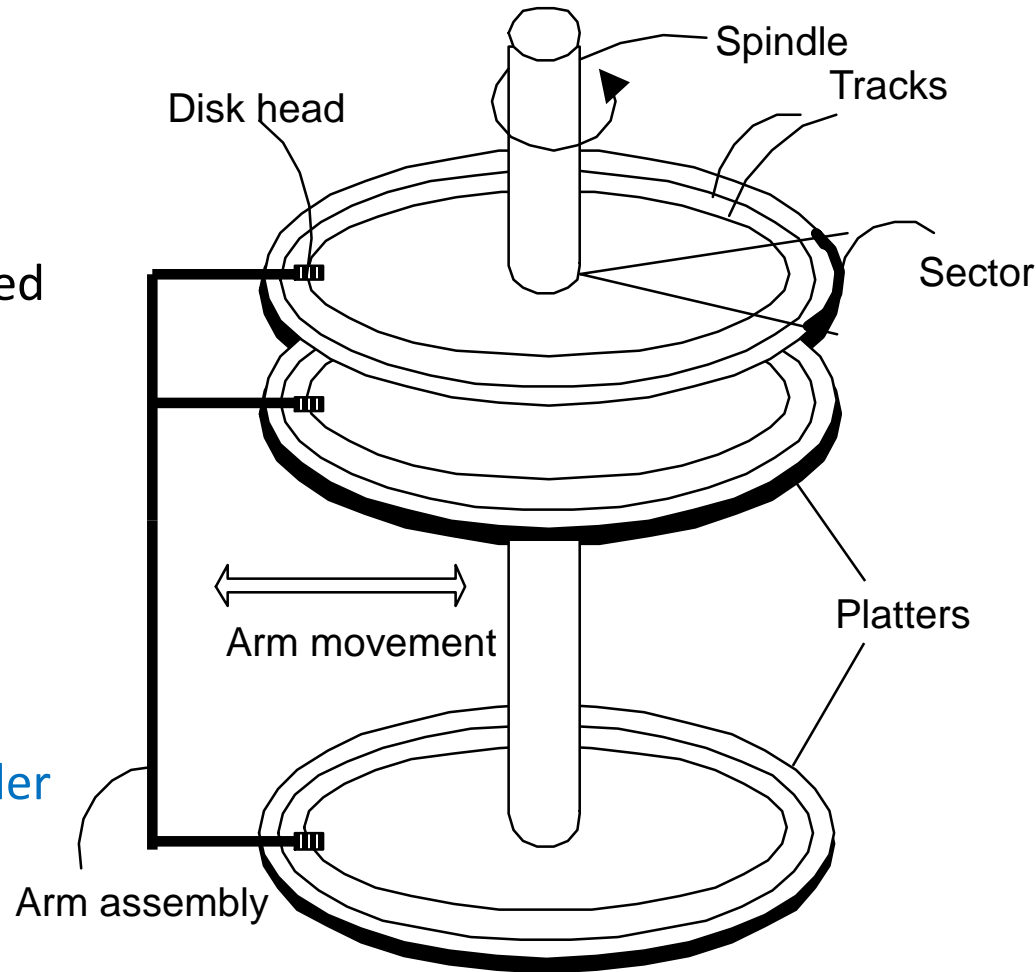  - READ: transfer data from disk to main memory (RAM)

    **I/O Time**

  - WRITE: transfer data from RAM to disk

  - I/O time dominates the time taken for database operations!

  - To minimize I/O time, it is necessary to store and locate data *strategically*

# Magnetic Disks

- Data is stored in disk blocks

- Blocks are arranged in concentric rings called tracks

- Each track is divided into arcs called sectors (whose size is fixed)

- The block size is a multiple of sector size

- The set of all tracks with the same diameter is called cylinder

- To read/write data, the arm assembly is moved in or out to position a head on a desired track

# Accessing a Disk Block

- ## What is I/O time?

  - ### The time to move the disk heads to the track on which a desired block is located

  - ### The waiting time for the desired block to rotate under the disk head

  - ### The time to actually read or write the data in the block once the head is positioned

# Accessing a Disk Block

- What is I/O time?
  - Seek Time
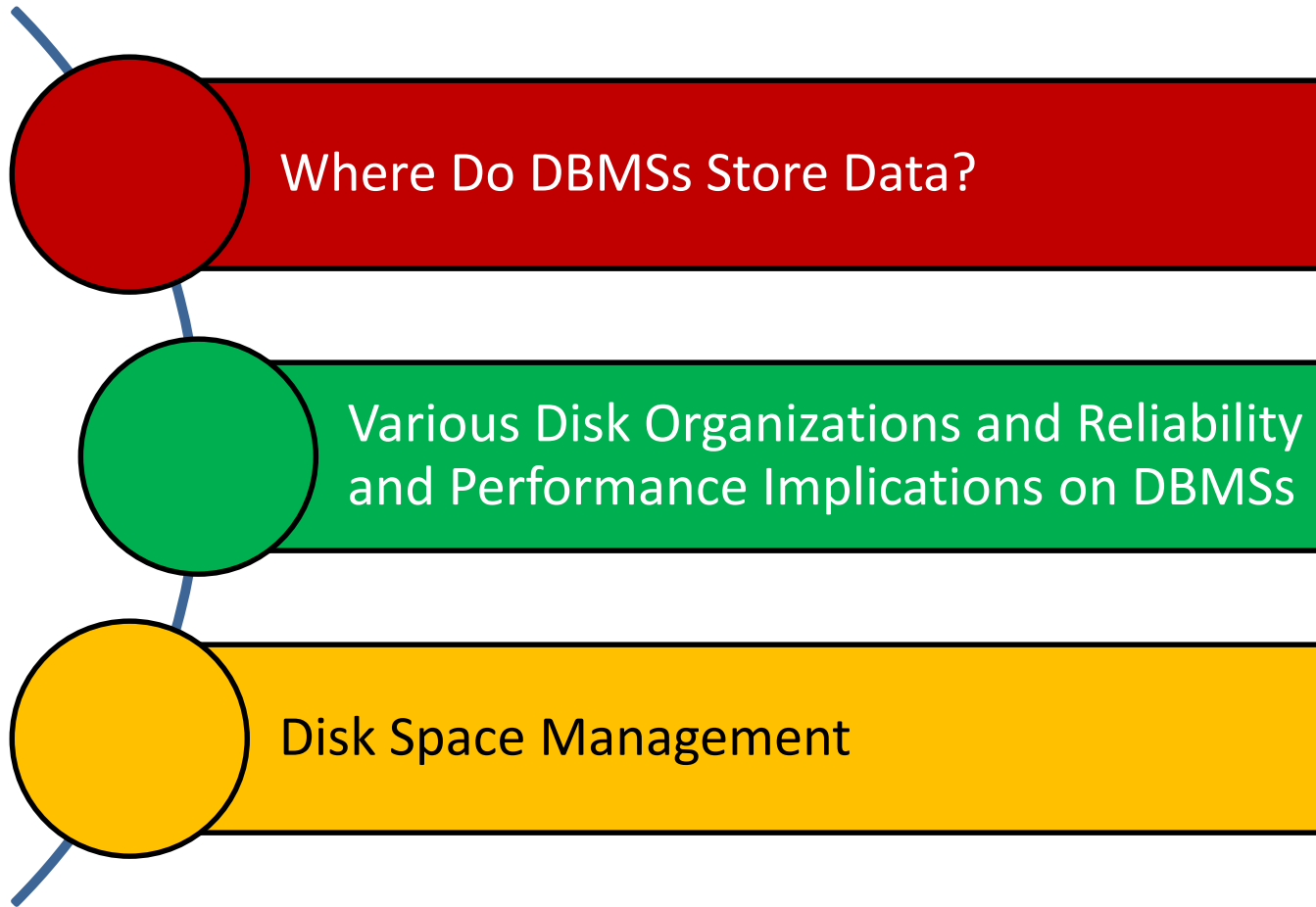  - Rotational Time
  - Transfer Time
- I/O time = seek time + rotational time + transfer time

# Implications on DBMSs

- Seek time and rotational delay dominate!

- Key to lower I/O cost: reduce seek/rotation delays!

- How to minimize seek and rotational delays?
  - Blocks on same track, followed by
  - Blocks on same cylinder, followed by
  - Blocks on adjacent cylinder
  - Hence, _sequential_ arrangement of blocks of a file is a big win!

More on that later…

# Outline

**Where Do DBMSs Store Data?**

**Various Disk Organizations and Reliability and Performance Implications on DBMSs** ✓

Disk Space Management

# Many Disks vs. One Disk

- Although disks provide cheap, non-volatile storage for DBMSs, they are usually bottlenecks for DBMSs
  - Reliability
  - Performance

- How about adopting multiple disks?
  1. More data can be held as opposed to one disk
  2. Data can be stored redundantly; hence, if one disk fails, data can be found on another
  3. Data can be accessed concurrently

# Many Disks vs. One Disk

- Although disks provide cheap, non-volatile storage for DBMSs, they are usually bottlenecks for DBMSs
  - Reliability
  - Performance

- How about adopting multiple disks?
  1. More data can be held as compared to one disk    **Capacity!**
  2. Data can be stored redundantly; hence, if one disk fails, data can be found on another    **Reliability!**
  3. Data can be accessed concurrently    **Performance!**
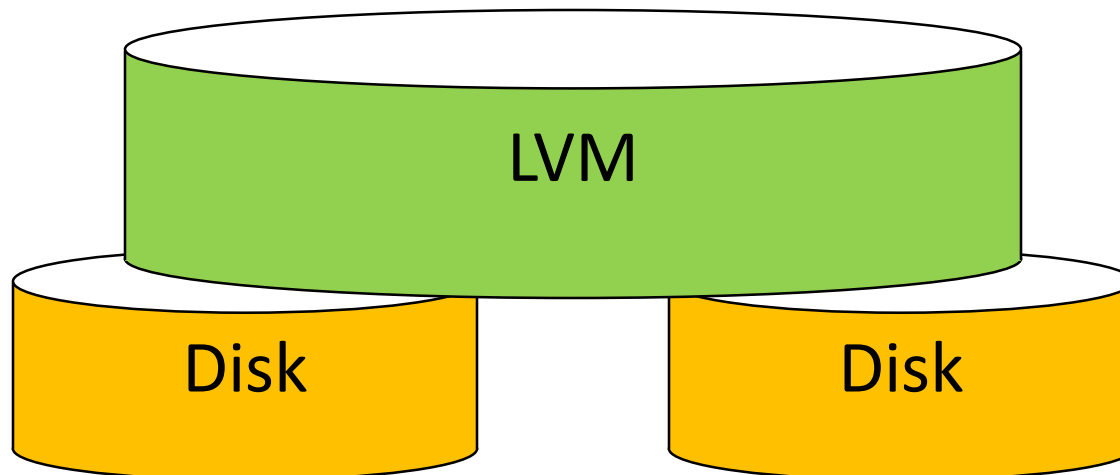
# Multiple Disks

Discussions on:

Reliability

Performance

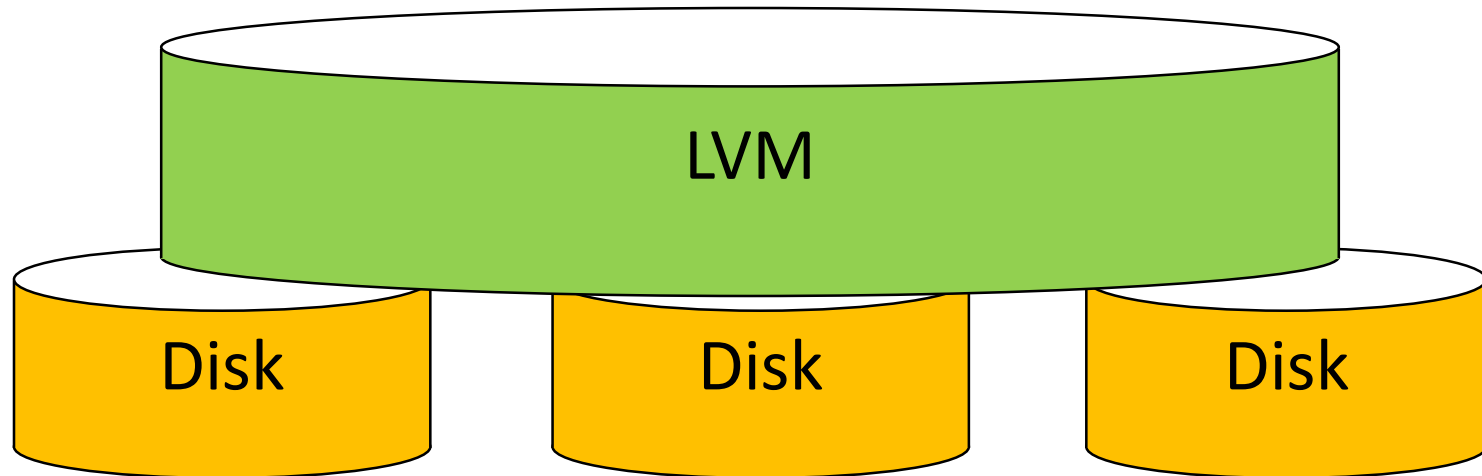Reliability + Performance

# Logical Volume Managers (LVMs)

- But, disk addresses used within a file system are assumed to refer to one particular disk (or sub-disk)

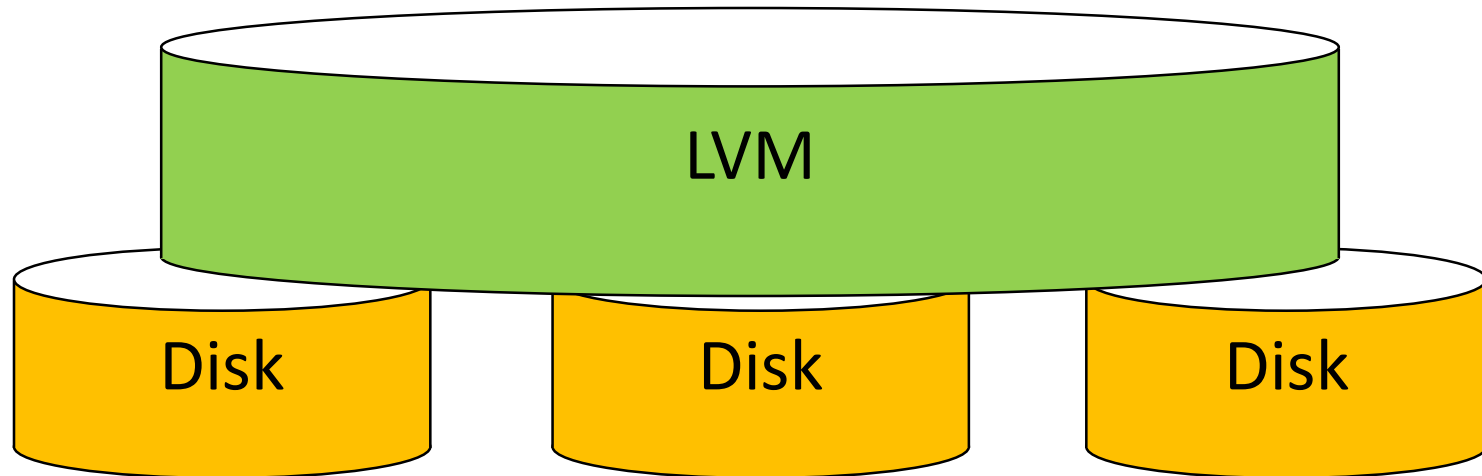- What about providing an abstraction that makes a number of disks *appear* as one disk?

# Logical Volume Managers (LVMs)



- **What can LVMs do?**
  - Spanning:
    - LVM transparently maps a *larger* address space to <u>different</u> disks
  - Mirroring:
    - Each disk can hold a separate, identical copy of data
    - LVM directs writes to the same block address on each disk
    - LVM directs a read to *any* disk (e.g., to the less busy one)

# Logical Volume Managers (LVMs)



- ## What can LVMs do?

  - ### Spanning:
    - LVM transparently maps a *larger* address space to <u>different</u> disks

  - ### Mirroring:
    - Each disk can hold a separate, identical copy of data
    - LVM directs writes to the same block address on each disk
    - LVM directs a read to *any* disk (e.g., to the less busy one)

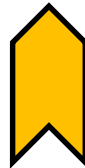  **Mainly Provides Redundancy!**

# Multiple Disks

Discussions on:

Reliability

Performance

Reliability + Performance

# Data Striping

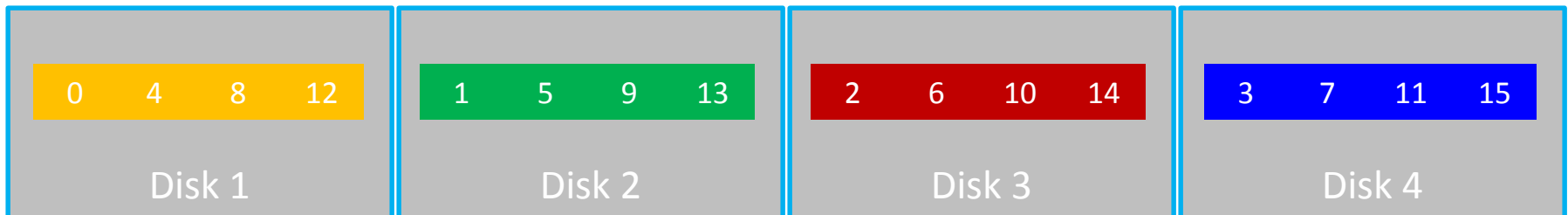- To achieve parallel accesses, we can use a technique called data striping

Logical File

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

Stripe Length = # of disks

Striping Unit

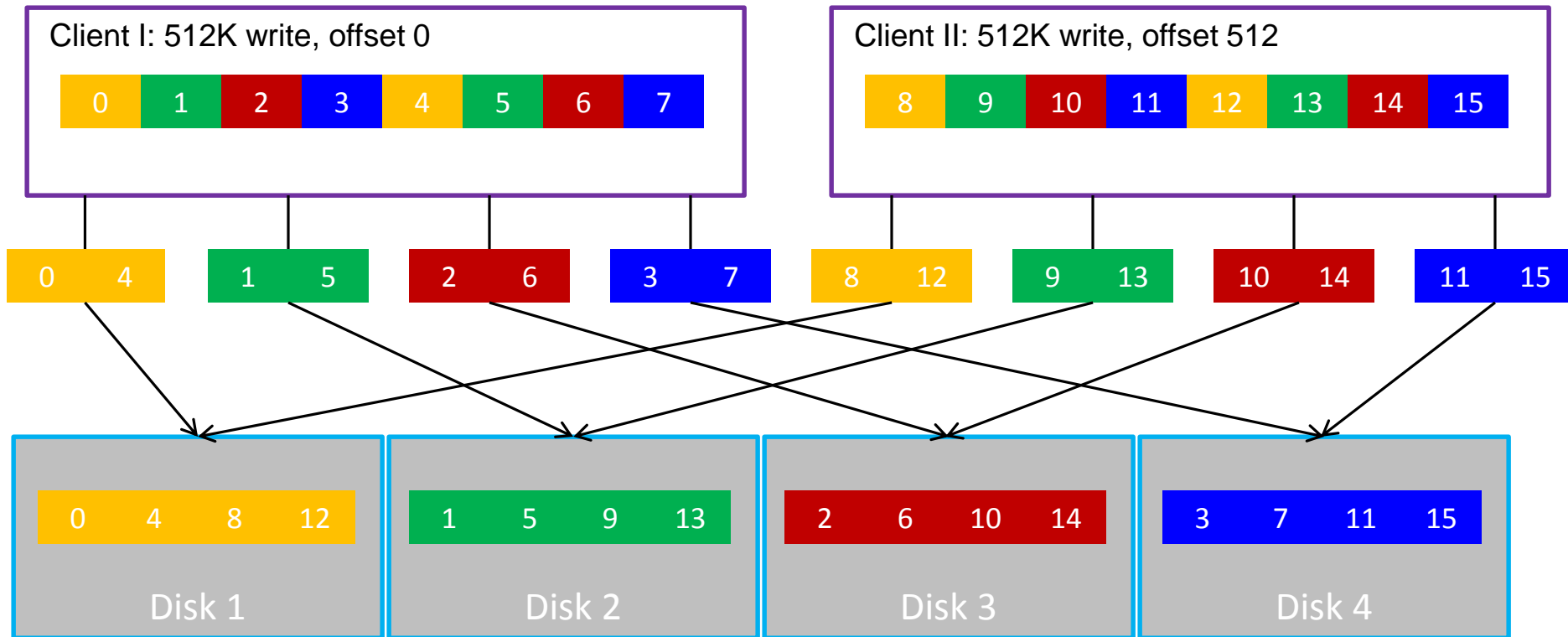| Disk 1 | Disk 2 | Disk 3 | Disk 4 |
|--------|--------|--------|--------|
| 0  4  8  12 | 1  5  9  13 | 2  6  10  14 | 3  7  11  15 |

# Data Striping

- To achieve parallel accesses, we can use a technique called data striping

# Data Striping

|  | Disk 1 | Disk 2 | Disk 3 | Disk 4 |
|---|---|---|---|---|
| **Stripe 1** | Unit 1 | Unit 2 | Unit 3 | Unit 4 |
| **Stripe 2** | Unit 5 | Unit 6 | Unit 7 | Unit 8 |
| **Stripe 3** | Unit 9 | Unit 10 | Unit 11 | Unit 12 |
| **Stripe 4** | Unit 13 | Unit 14 | Unit 15 | Unit 16 |
| **Stripe 5** | Unit 17 | Unit 18 | Unit 19 | Unit 20 |

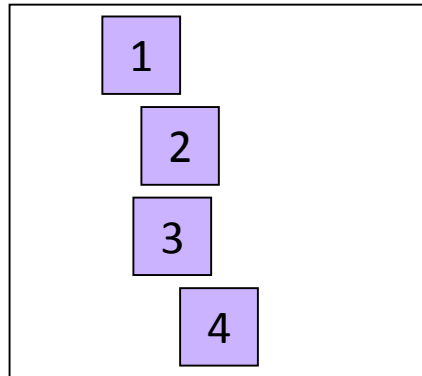Each stripe is written across all disks *at once*

Typically, a unit is either:
- A bit ➜ **Bit Interleaving**
- A byte ➜ **Byte Interleaving**
- A block ➜ **Block Interleaving**
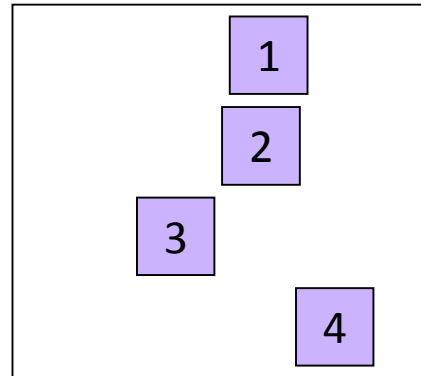
# Striping Unit Values: Tradeoffs

- Small striping unit values
  - Higher parallelism (**+**)
  - Smaller amount of data to transfer (**+**)
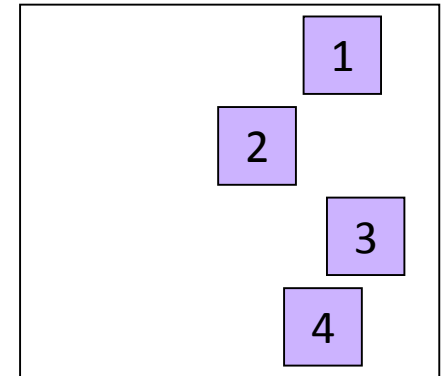  - Increased seek and rotational delays (**-**)

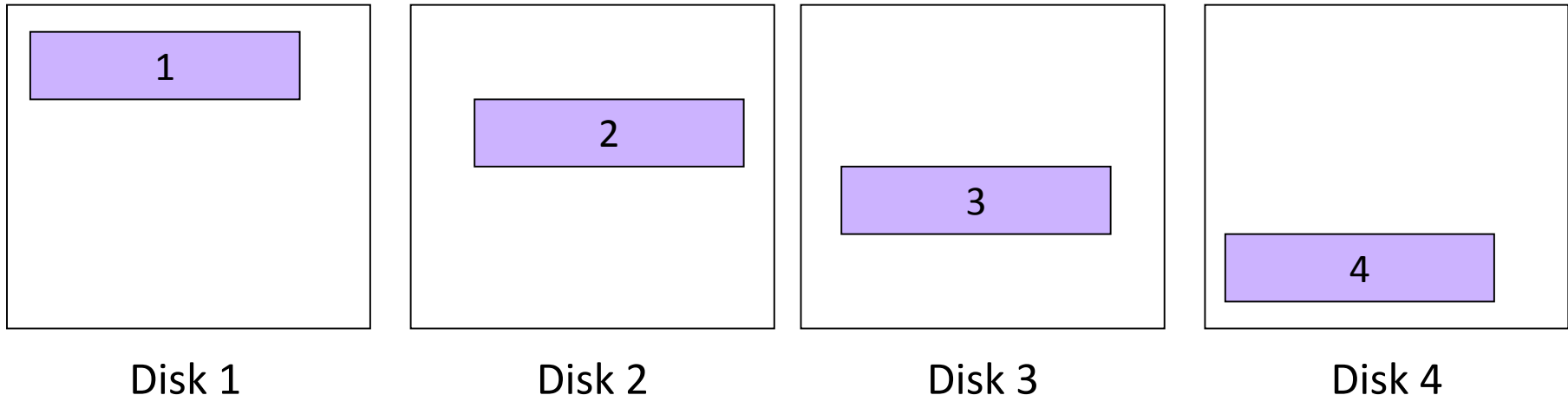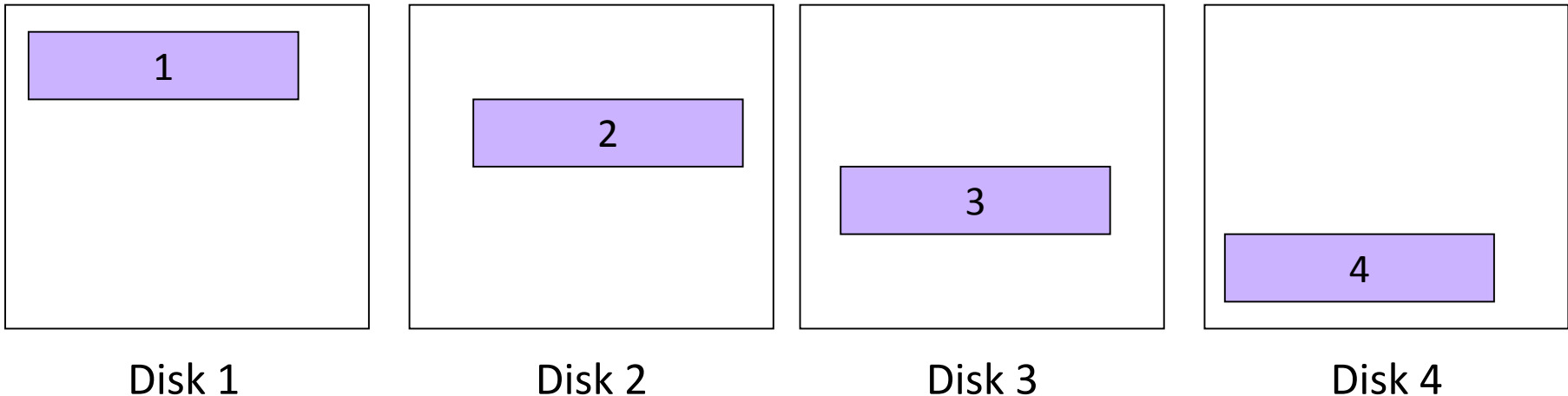| Disk 1 | Disk 2 | Disk 3 | Disk 4 |
|--------|--------|--------|--------|
| 1 2 3 4 | 1 2 3 4 | 1 2 3 4 | 1 2 3 4 |

# Striping Unit Values: Tradeoffs

- Large striping unit values
  - Lower parallelism (-)
  - Larger amount of data to transfer (-)
  - Decreased seek and rotational delays (+)
  - A request can be handled completely on a separate disk! (- or +)
  - But, multiple requests could be satisfied at once! (+)

| Disk 1 | Disk 2 | Disk 3 | Disk 4 |
|--------|--------|--------|--------|
| 1 | 2 | 3 | 4 |

# Striping Unit Values: Tradeoffs

- Large striping unit values
  - Lower parallelism
  - Larger amount of data to transfer
  - Decreased seek and rotational delays
  - A request can be handled completely on a separate disk!
  - **Number of requests = *Concurrency Factor***

| | | | |
|---|---|---|---|
| 1 | | | |
| | 2 | | |
| | | 3 | |
| | | | 4 |
| Disk 1 | Disk 2 | Disk 3 | Disk 4 |

# Multiple Disks

Discussions on:

Reliability

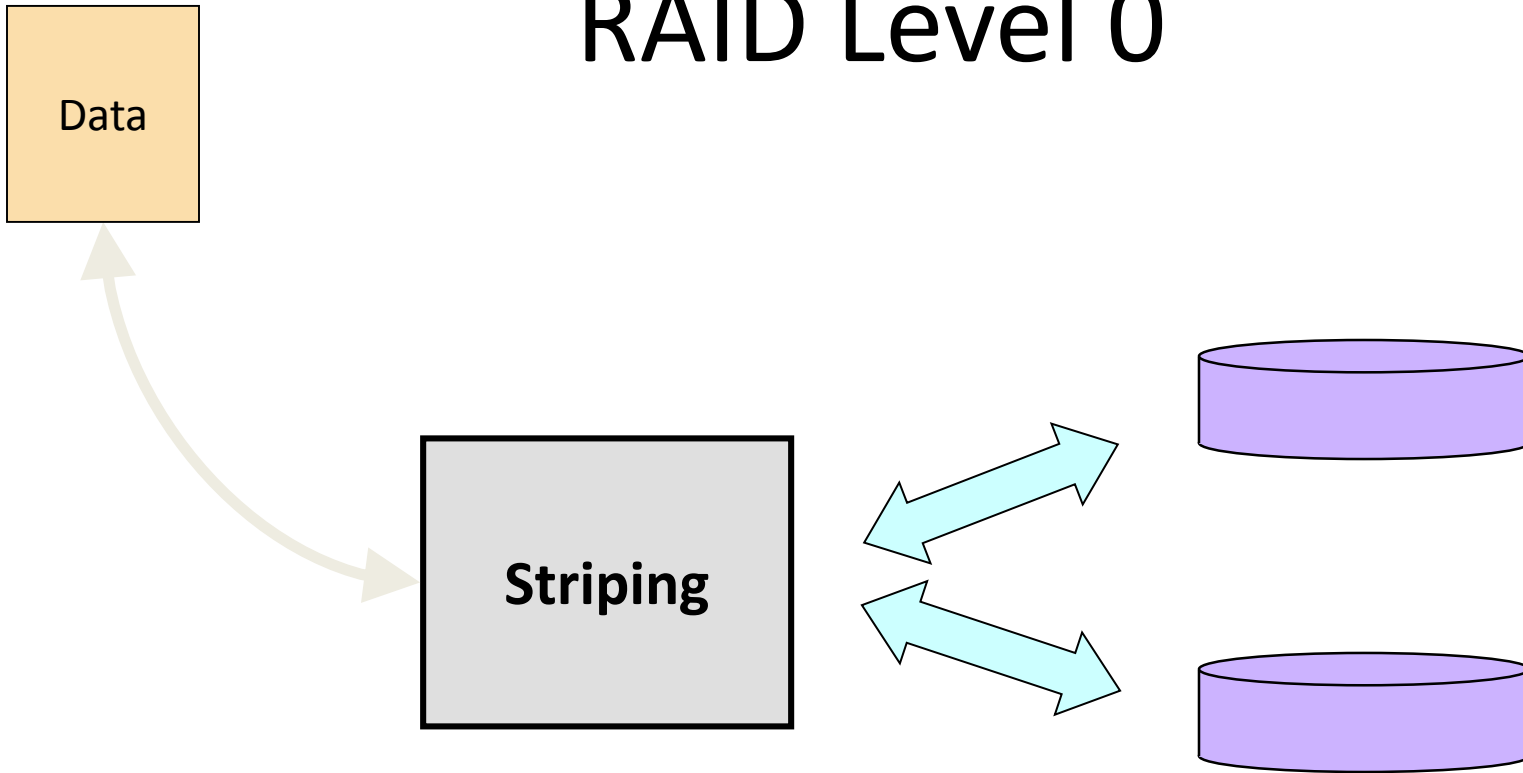Performance

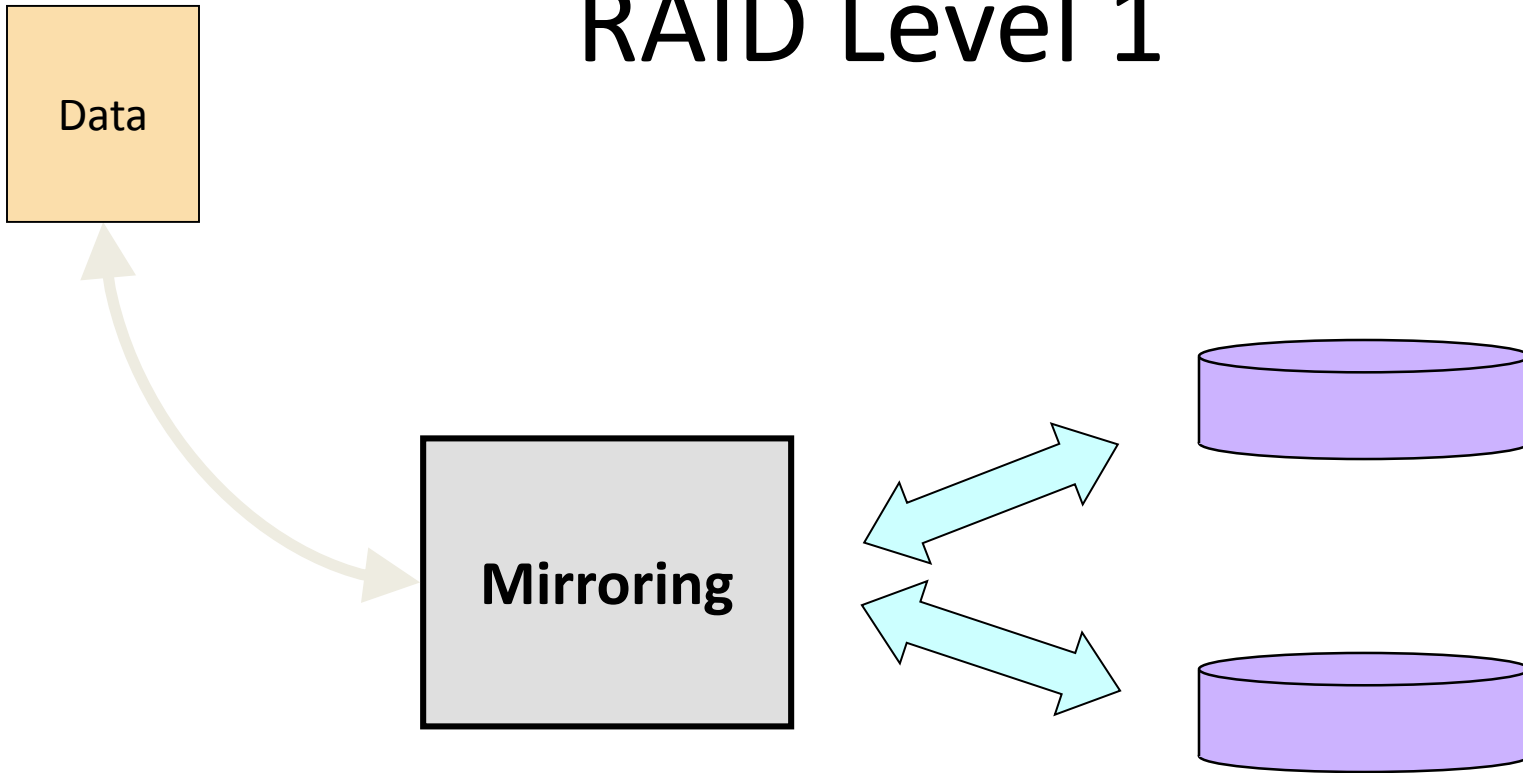Reliability + Performance

# Redundant Arrays of Independent Disks

- A system depending on **_N_** disks is much more likely to fail than one depending on one disk
  - If the probability of one disk to fail is **_f_**
  - Then, the probability of N disks to fail is $(1-(1-f)^N)$

- How would we combine reliability with performance?
  - Redundant Arrays of Inexpensive Disks (**RAID**) *combines* mirroring and striping
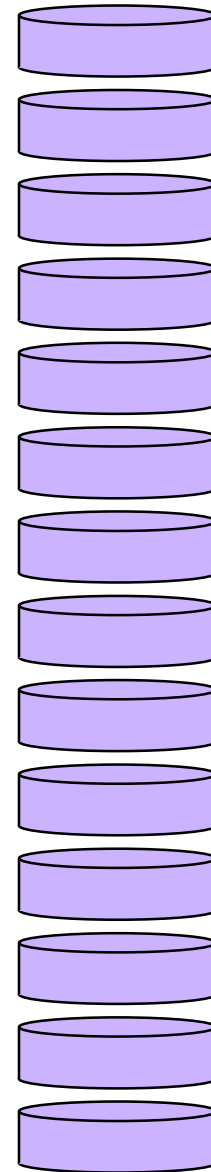
Nowadays, Independent!

# RAID Level 0

Data

**Striping**

# RAID Level 1

Data

**Mirroring**

# RAID Level 2

Data

**Bit Interleaving; ECC**

Data bits

Check bits

# RAID Level 3

Data

Bit Interleaving;
Parity

Data bits

Parity bits

# RAID Level 4

Data

**Block Interleaving;
Parity**

Data blocks

Parity blocks

# RAID Level 5

Data

Block Interleaving; Parity

Data and parity blocks

# RAID 4 vs. RAID 5

- **What if we have a lot of small writes?**
  - RAID 5 is the best

- **What if we have mostly large writes?**
  - Multiples of stripes
  - Either is fine

- **What if we want to expand the number of disks?**
  - RAID 4: add a disk and re-compute parity
  - RAID 5: add a disk, re-compute parity, and shuffle data blocks among all disks to reestablish the check-block pattern (*expensive!*)
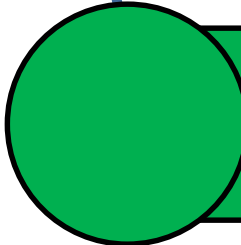
# Beyond Disks: Flash

- Flash memory is a relatively new technology providing the functionality needed to hold file systems and DBMSs
  - It is writable
  - It is readable
  - Writing is slower than reading
  - It is non-volatile
  - Faster than disks, but slower than DRAMs
  - Unlike disks, it provides random access
  - Limited lifetime
  - More expensive than disks

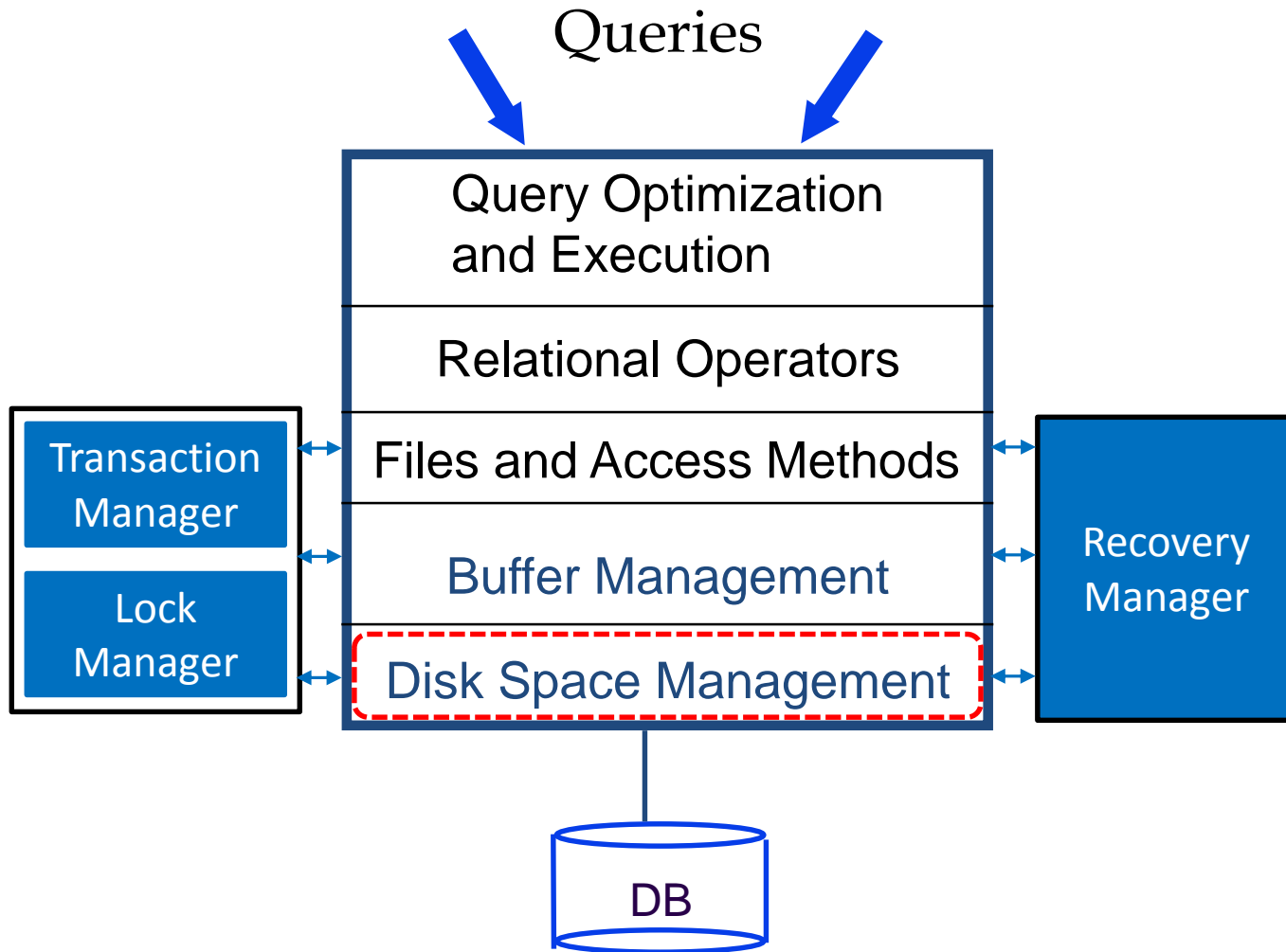# Outline

Where Do DBMSs Store Data?

Various Disk Organizations and Reliability and Performance Implications on DBMSs

Disk Space Management ✓

# DBMS Layers

# Disk Space Management

- DBMSs disk space managers
  - Support the concept of a page as a unit of data
    - Page size is usually chosen to be equal to the block size so that reading or writing a page can be done in 1 disk I/O

  - Allocate/de-allocate pages as a *contiguous* sequence of blocks on disks

  - Abstracts hardware (and possibly OS) details from higher DBMS levels

# What to Keep Track of?

- The DBMS disk space manager keeps track of:
    - Which disk blocks are in use
    - Which pages are on which disk blocks

- Blocks can be initially allocated contiguously, but allocating and de-allocating blocks usually create *"holes"*

- Hence, a mechanism to keep track of *free blocks* is needed
    - A list of free blocks can be maintained (*storage could be an issue*)
    - Alternatively, a bitmap with one bit per each disk block can be maintained (*more storage efficient and faster in identifying contiguous free areas!*)

# OS File Systems vs. DBMS Disk Space Managers

- Operating Systems already employ disk space managers using *their* "file" abstraction
  - "Read byte *i* of file *f*" ➔ "read block *m* of track *t* of cylinder *c* of disk *d*"

- DBMSs disk space managers usually pursue their own disk management without relying on OS file systems
  - Enables portability
  - Can address larger amounts of data
  - Allows *spanning* and *mirroring*

# Next Class

Queries

Query Optimization and Execution

Relational Operators

Files and Access Methods

Buffer Management

Disk Space Management

Transaction Manager

Lock Manager

Recovery Manager

Buffer Management and Parts of Files and Access Methods

DB