# 3. SSD Controller

**Special Topics in Computer Systems:**
Modern Storage Systems
(IC820-01)

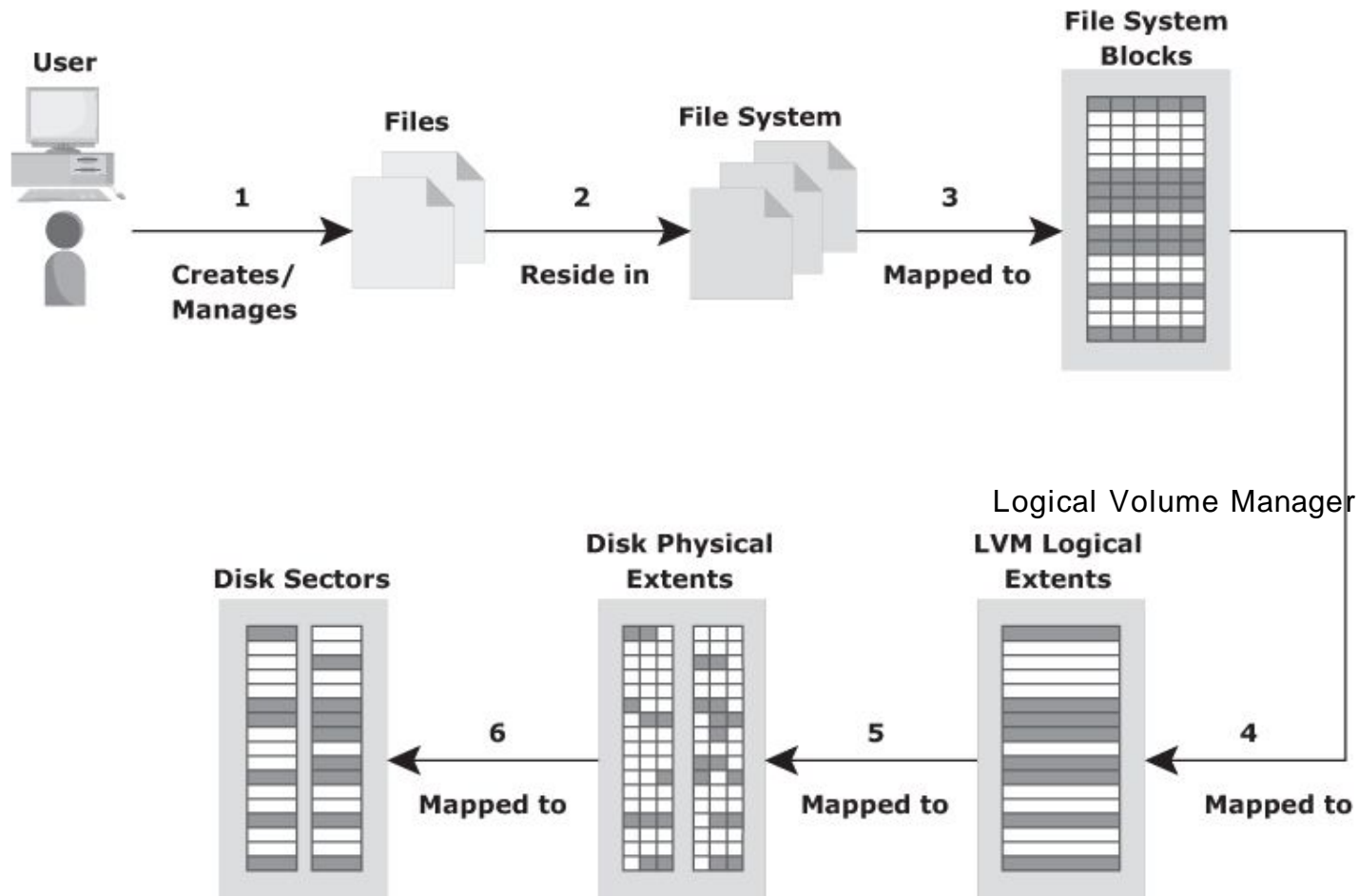**Instructor:**

Prof. Sungjin Lee (sungjin.lee@dgist.ac.kr)
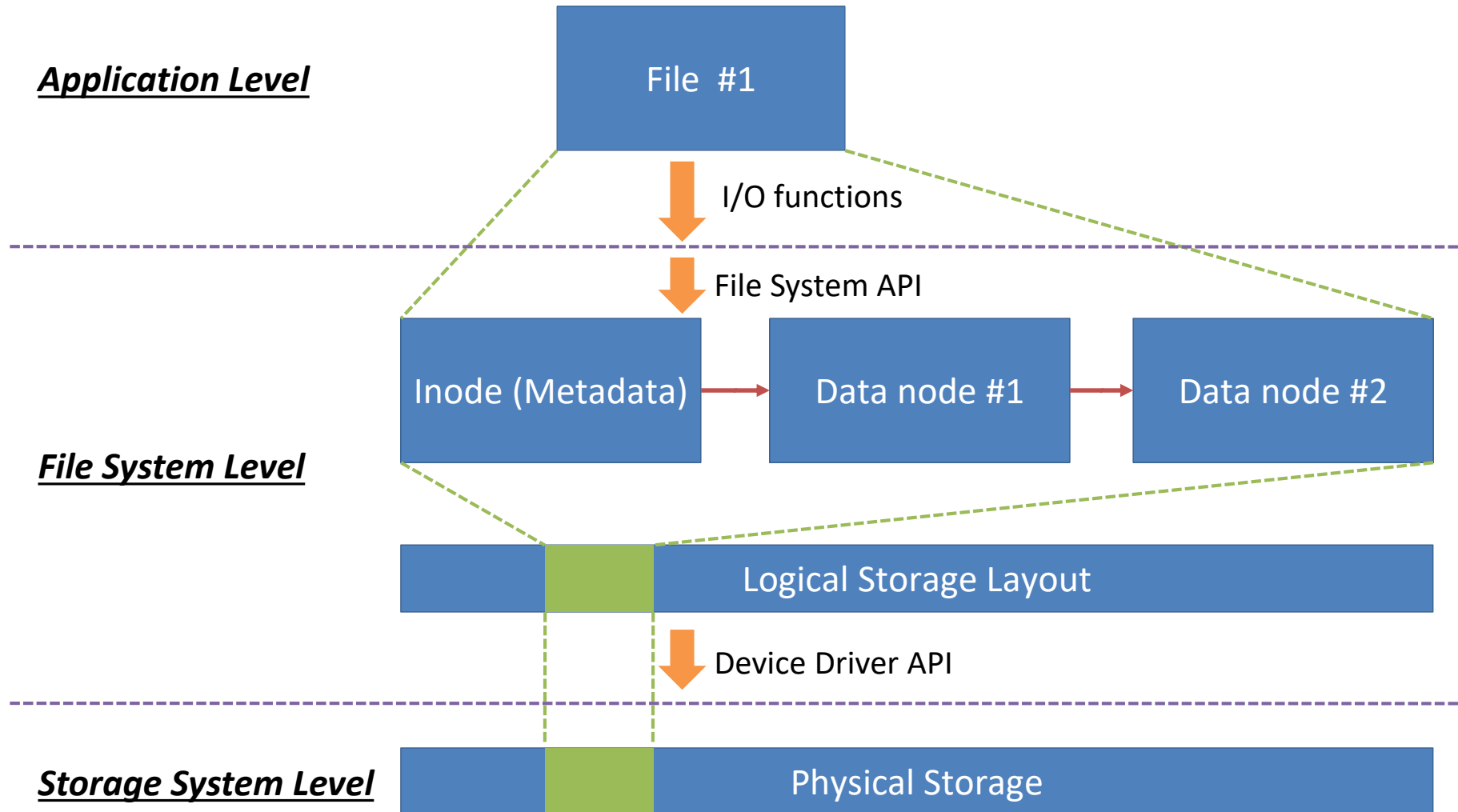
# Outline

- **Storage Abstraction & Protocols**
- **Flash Packages**
- **SSD Controller**
- **Flash Array & I/O Parallelism**
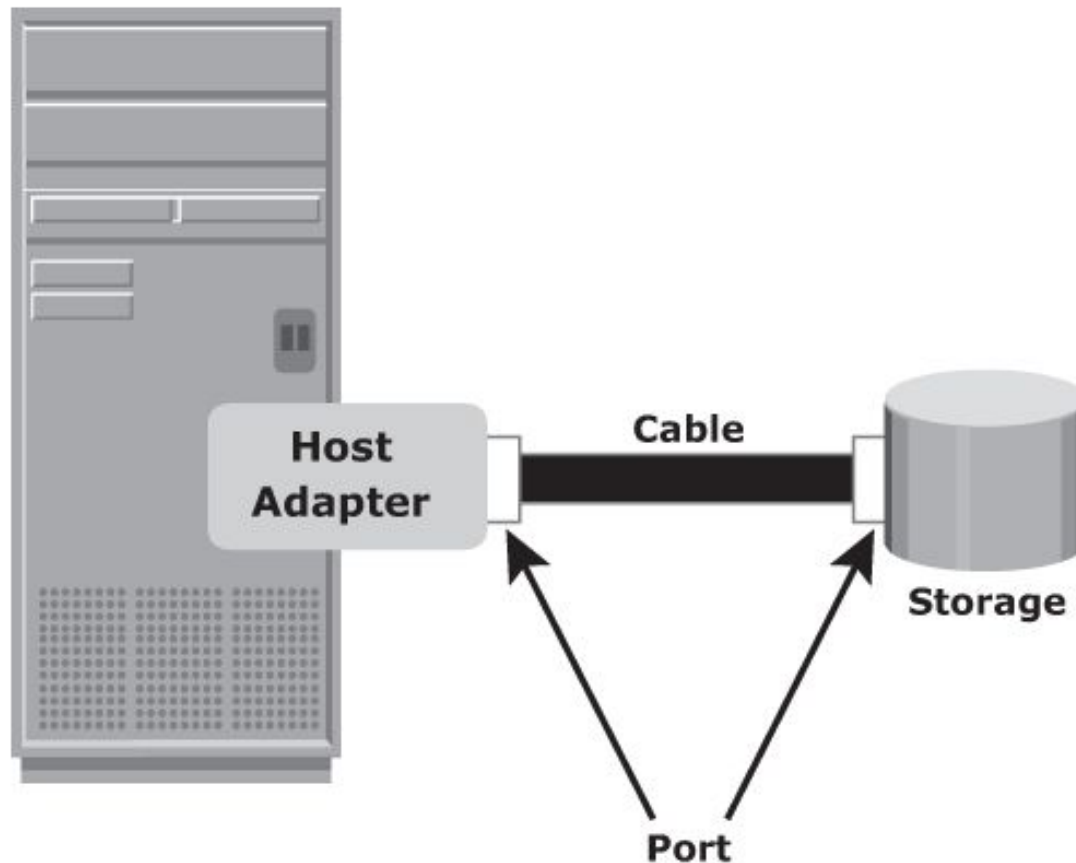
# File to Storage

■ **In most cases, users manage their data in the form of files**

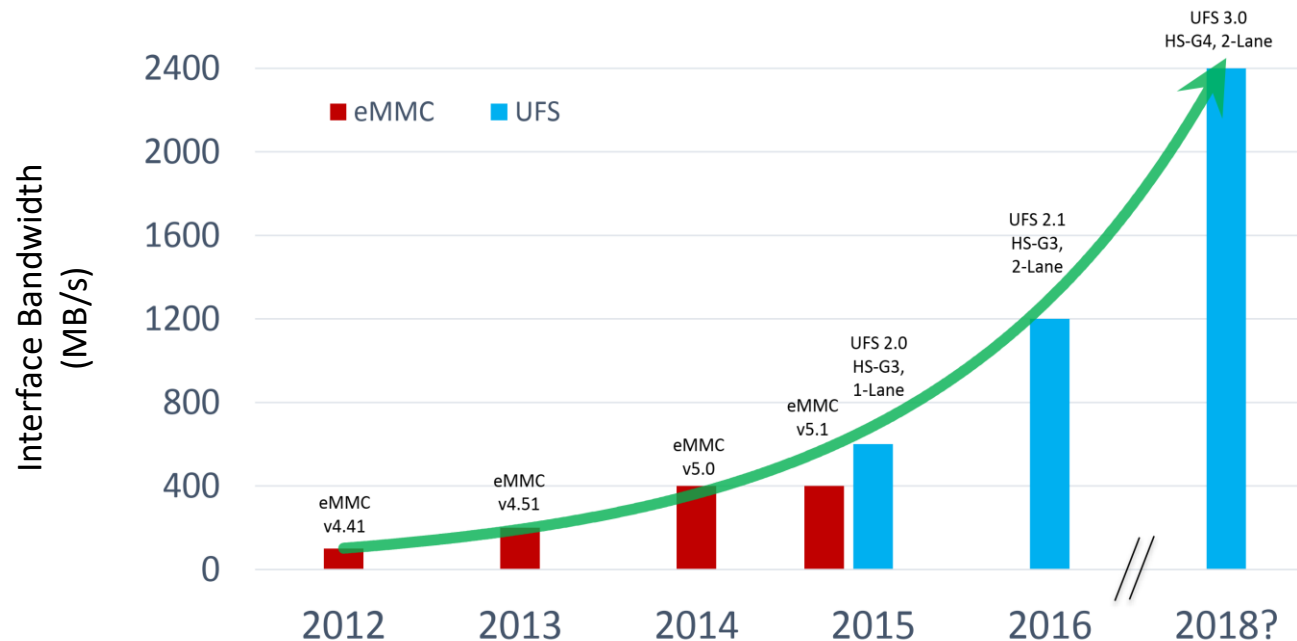# File to Storage (Cont.)

# Storage Interface & Protocols



*More detailed introduction to storage interface & protocols will be given as a separate lecture*

# Storage Interface & Protocols (Cont.)

- **There are various types of storage interfaces depending on target applications' requirements**

- **Case #1: mobile storage interface**
  - High throughput and low latency (everybody wants high-speed storage)
  - *Low active power* (e.g., Max 1.62W) and *limited form factor*

# Storage Interface & Protocols (Cont.)

- **There are various types of storage interfaces depending on target applications' requirements**

- **Case #2: desktop/server storage interface**
  - *High performance* is the primary goal
  - Servers require *many* storage devices and high *reliability*

| Interface | Mnemonic Meaning | Transfer Speed | Characteristics |
|-----------|------------------|----------------|------------------|
| **SATA** | Serial ATA | Up to 2GB/s | Low cost |
| **SAS** | Serial Attached SCSI | Up to 3GB/s        HDD 150 | Supports multiple ports Error detection/correction |
| **FC** | Fibre Channel | Up to 16GB/s | Predominately SCSI commands and features |
| **NVMe** | Nonvolatile memory express over Pcie | 4 GB/s per lane | Up to 32 lanes High command queue support |
| **NVDIMM** | Nonvolatile memory on memory channel | Up to 1 GB/s over 64-bit bus | Very slow latency No interrupt Deterministic |

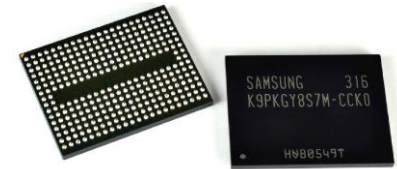* T. Coughlin et.al., "Digital Storage and Memory Technology (Part1)," IEEE Technical Trend Paper, Nov. 2017

# Outline

- **Storage Abstraction & Protocols**
- **Flash Packages**
- **SSD Controller**
- **Flash Array & I/O Parallelism**

Review:
# NAND Flash Chip

- **A group of blocks compose a NAND flash die**
- **Two or four dies are then packed as a single NAND chips**

<NAND flash chip specification>

| | |
|---|---|
| Page Read to SRAM | 50 $\mu$s |
| Page Program (Write) from SRAM | 500 $\mu$s |
| Block Erasure | 4 $m$s |
| Serial Access to SRAM | 100 $\mu$s (per 4 KB) |
| Page Size | 16 KB |
| Block Size | 4 MB |
| Die Size | 8 Gb |
| Dies per Package | 1, 2, or 4 |

- **Key properties**
  - Read is 10x faster than write: 50 $\mu$s vs 500 $\mu$s
  - Erasure is slowest: 4 $m$s
  - Throughput is not so high: 36.4 MB/s for reads (Note: recent SSDs offer 5 GB/s)

Review:
# How to Improve Performance?

Read throughput (36.4 MB/s) = 16 KB / (50 + 400) $\mu s$

■ **Strategy 1: Reduce the latency of three I/O primitives**

▪ The latency of three I/O primitives get longer as NAND cells scale down

| | **SLC (2D)** | **MLC (2D)** | **TLC (2D)** | **MLC (3D)** |
|---|---|---|---|---|
| Page Program | 25 $\mu s$ | 50 $\mu s$ | 100 $\mu s$ | 50 $\mu s$ |
| Page Read | 250 $\mu s$ | 500 $\mu s$ | 3 $m$s | 500 $\mu s$ |

▪ 3D NAND improves I/O latency with larger cells

■ **Strategy 2: Improve the bandwidth of I/O interface**

| | **Interface** | **Throughput (MB/s)** |
|---|---|---|
| Before 2006 | No Standard (SDR) | 40 MB/s |
| 2006 | NV-SDR | 50 MB/s |
| 2008 | NV-DDR | 200 MB/s |
| 2011 | NV-DDR2 | 533 MB/s |
| 2014 | NV-DDR3 | 800 MB/s |

Review:
# How to Improve Performance? (Cont.)
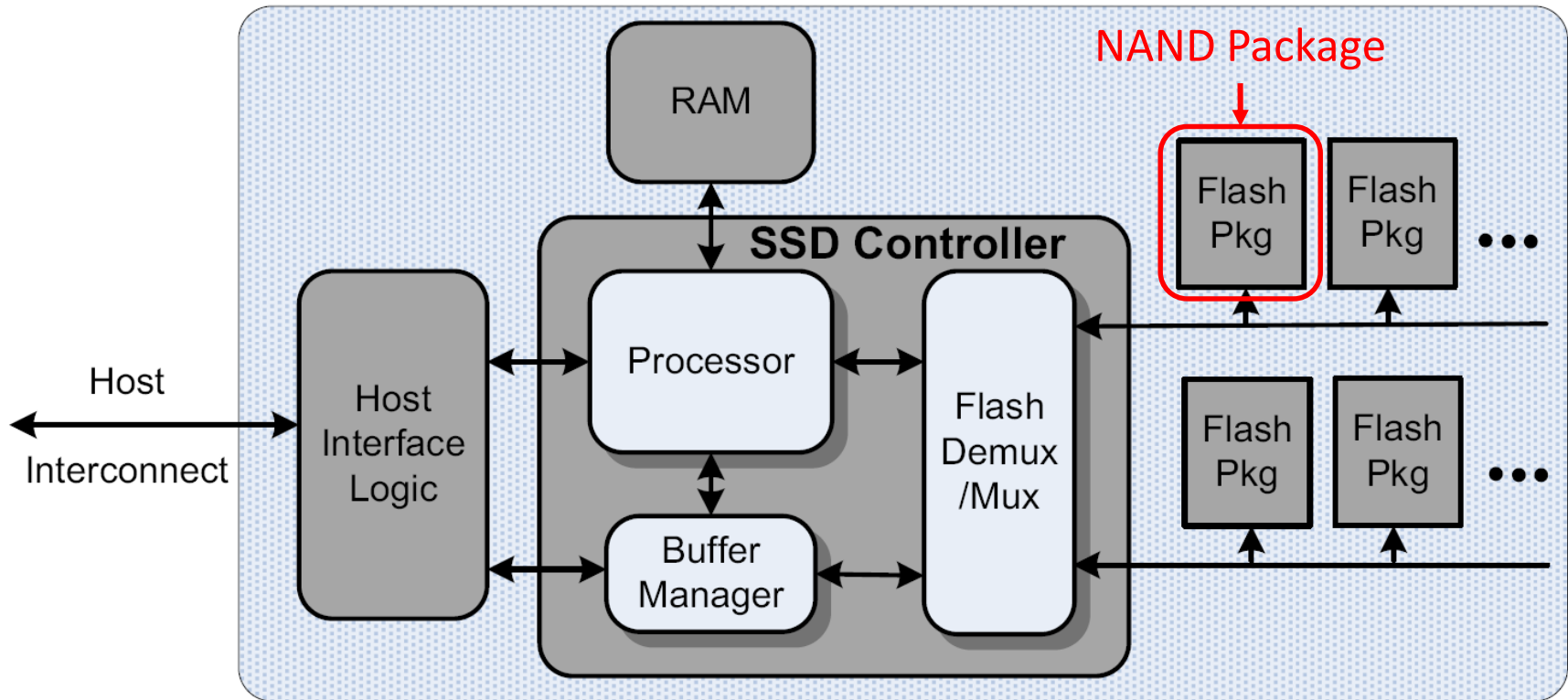
- **Strategy 3: Aggregate many NAND chips and access them in parallel**
  - Example: 36.4 MB/s x 64-128 NAND dies = **2.3-4.7 GB/s**
  - This is what an SSD controller does!
    - To achieve optimal performance, sophisticated algorithms are needed!

SATA Interface

SSD Controller

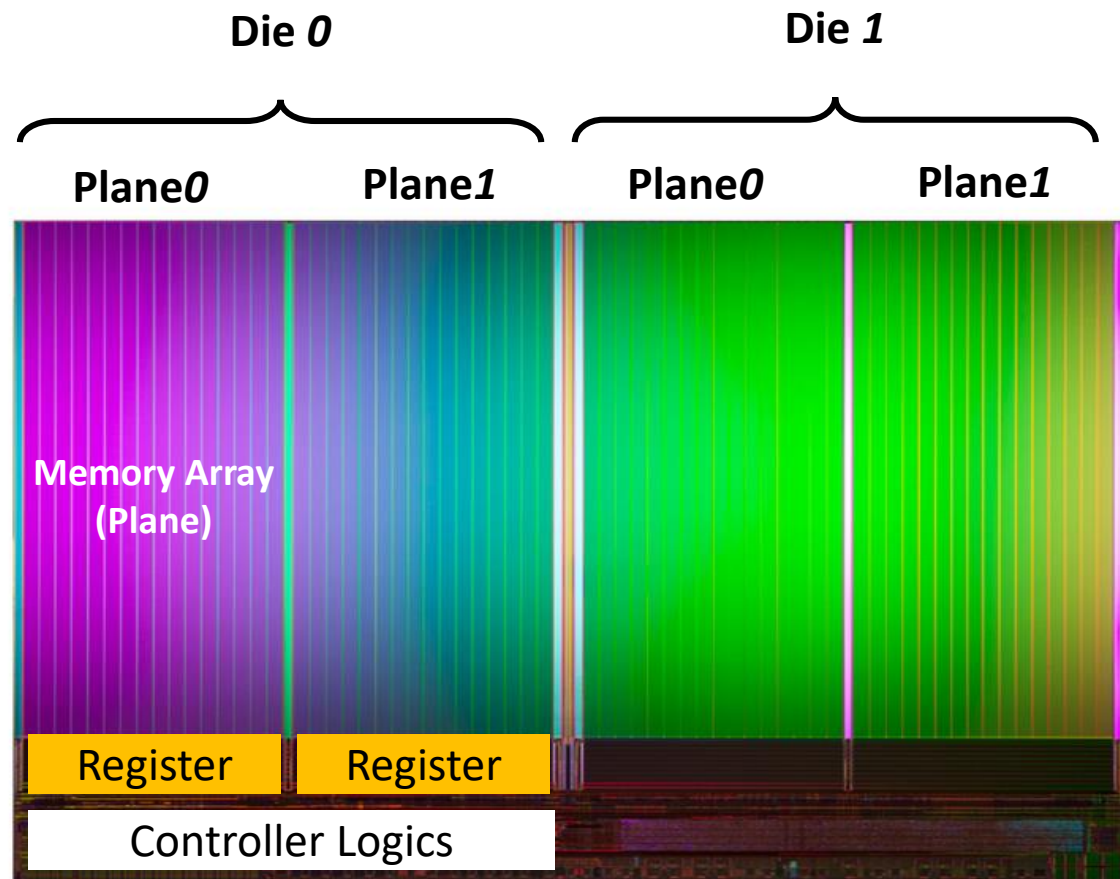NAND chip

# SSD Overview

- **Composed of three main modules**
  - (1) A host interface, (2) an SSD Controller, and (3) an array of NAND packages (or chips)
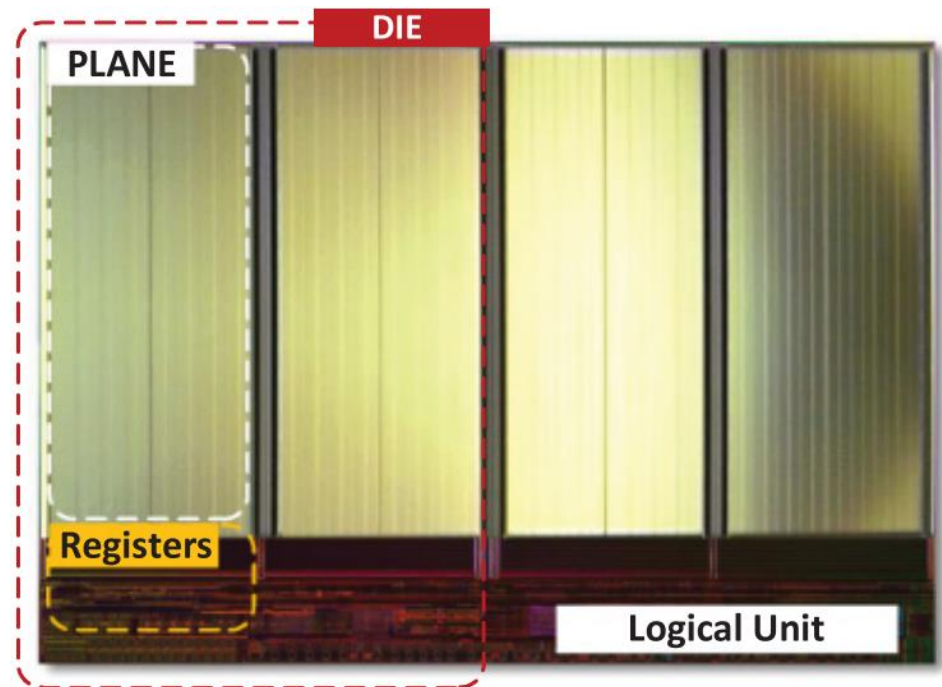
# NAND Package

- **Employing cache and data registers**
- **Multiple planes**
- **Multiple dies**



Die *0* — Die *1*

Plane*0* — Plane*1* — Plane*0* — Plane*1*

Memory Array (Plane)

Register — Register

Controller Logics

# NAND Package (Cont.)

- **Flash package employs multiple <span style="color:red">dies</span>, each consisting of one or more <span style="color:red">planes</span>**

  - Planes *share* multiple peripherals, which enable all the target pages connected to their wordlines/bitlines

- **Each plane employs a set of registers to cache/buffer the data brought by flash interface**

die          control logic
Plane                .

              .
read/write/erase       handle
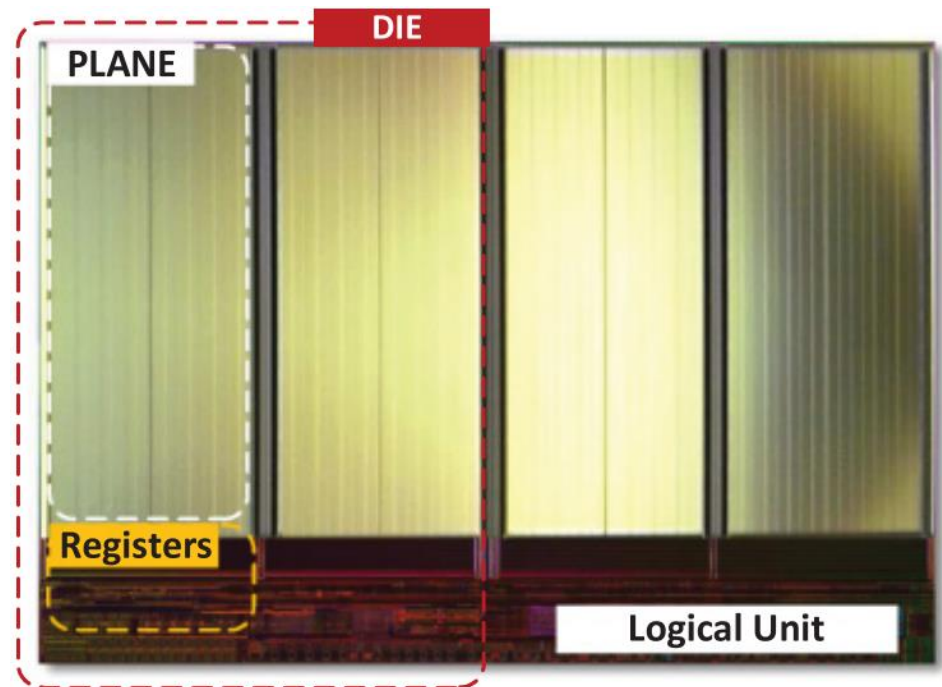                  die      plane

         .



14

# Physical Plane and Die

- **Each plane defines a different set of block addresses, and the pages/blocks on different planes operate in parallel**
  - The operation type for the requests across different plane should be same
  - The page and plane addresses should be identical
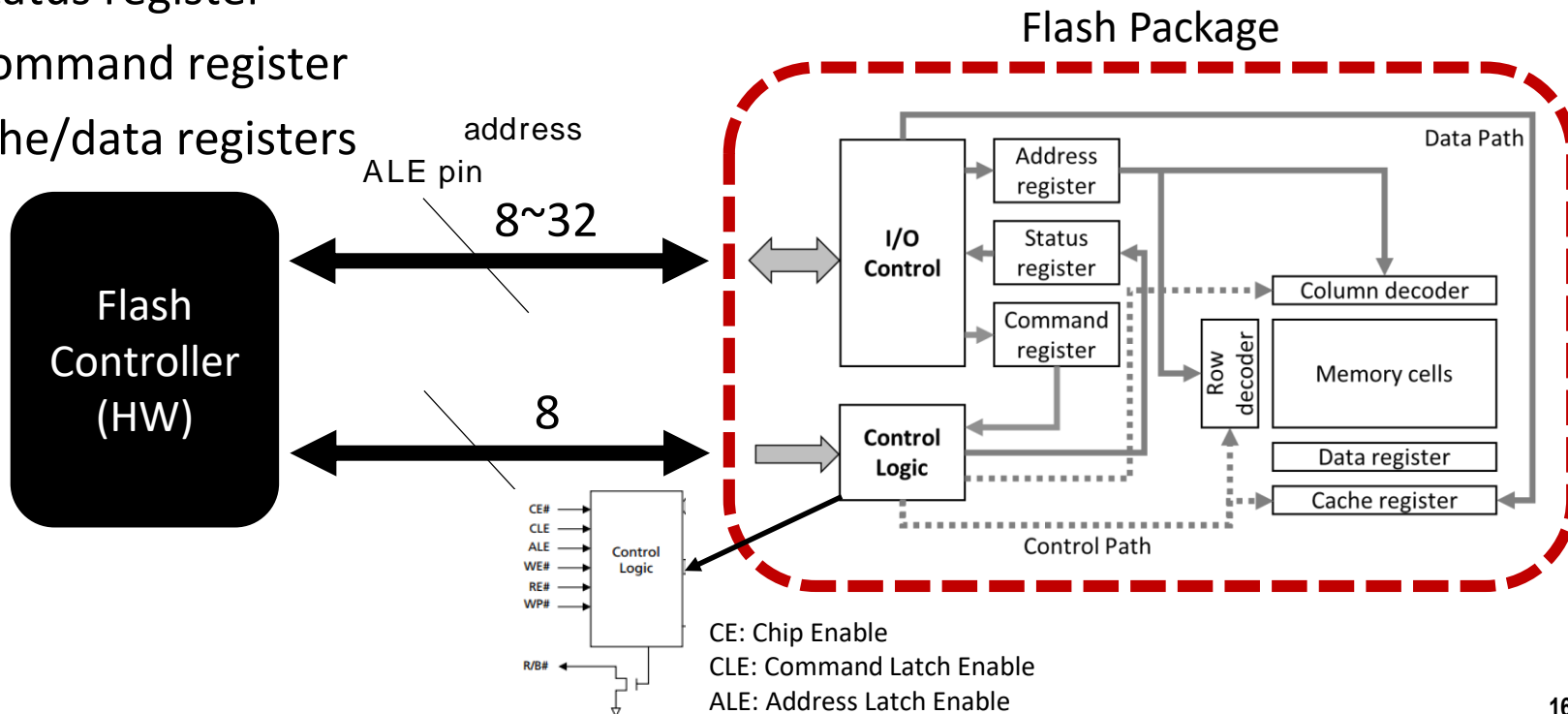- **All individual dies can be simultaneously activated, but share data-path (described later)**
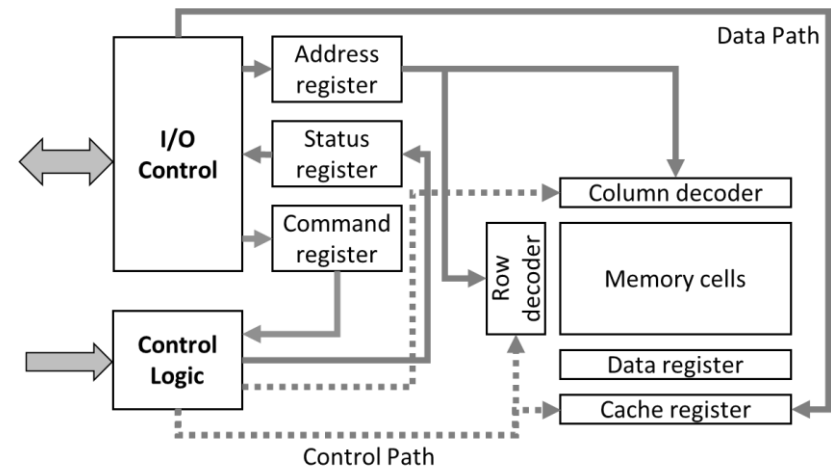
die

data- path(bus)            .

# Flash Microarchitecture

- **All data, commands and addresses are <span style="color:red">multiplexed</span> onto same I/O pins and received by I/O control circuit**

- **Each component of flash transactions is latched by**
  - An address register
  - A status register
  - A command register
  - Cache/data registers



Flash Package

address

ALE pin

8~32

8

Flash Controller (HW)

Flash Package

Data Path

I/O Control — Address register, Status register, Command register

Control Logic

Row decoder

Column decoder

Memory cells

Data register

Cache register

Control Path

CE#
CLE
ALE
WE#
RE#
WP#

Control Logic

R/B#

CE: Chip Enable
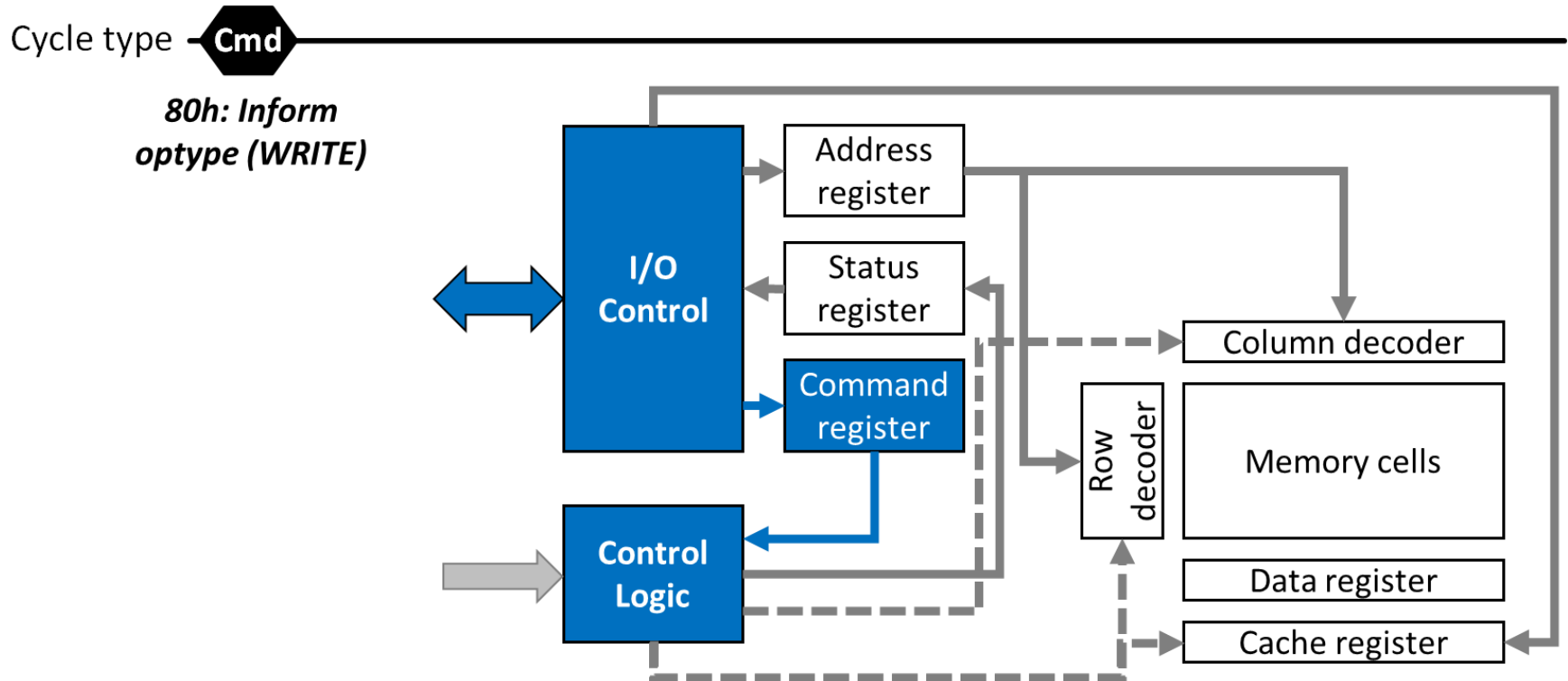CLE: Command Latch Enable
ALE: Address Latch Enable

16

# Flash Interface and Protocol

- **To issue a read/write request through the multiplexed I/O pins, a flash controller should obey the protocol that flash interface defines for all memory transactions**
  - Basic operations (e.g., reads and writes)
  - Cache mode operations
  - Multi-plane mode operations
  - Copyback, etc.

# Page Write

- The first command is needed to indicate that this is a page write

- The column address (C-ADDR) and row address (R-ADDR) are shipped in a serial order

- Data should be then transferred, and the second command initiates a write

- Once the write is performed, users need to send the last command to check up the status of target
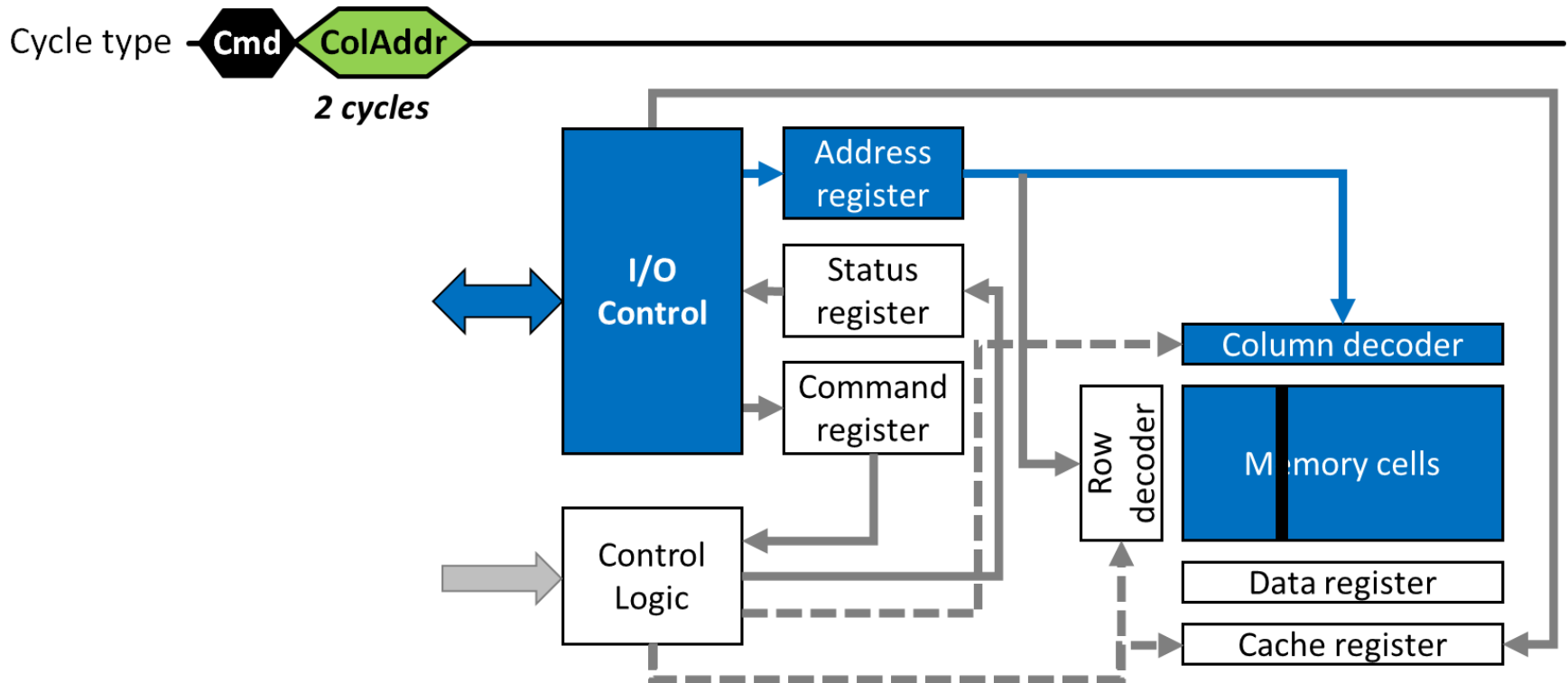


18

# Page Write

- The first command is needed to indicate that this is a page write
- The column address (C-ADDR) and row address (R-ADDR) are shipped in a serial order
- Data should be then transferred, and the second command initiates a write
- Once the write is performed, users need to send the last command to check up the status of target
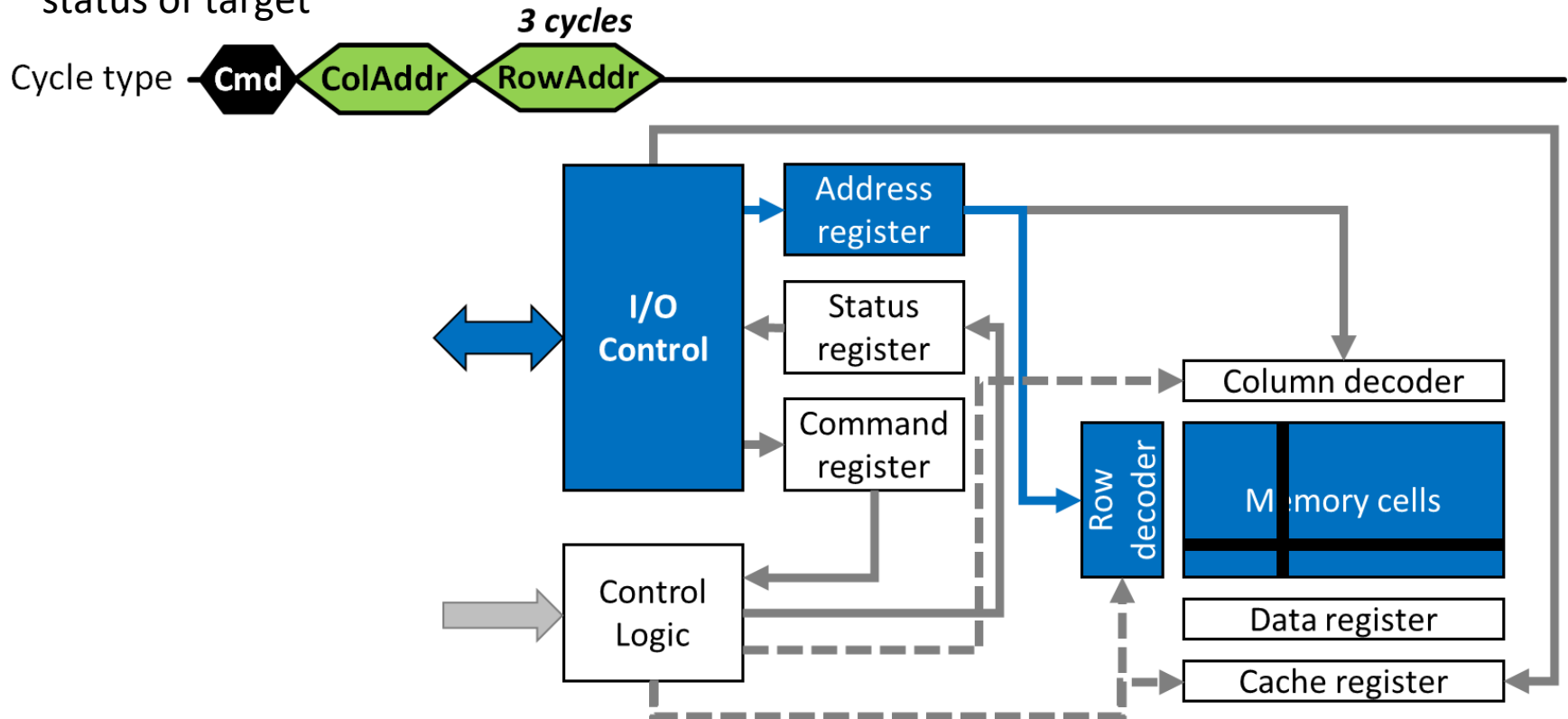
# Page Write

- The first command is needed to indicate that this is a page write
- The column address (C-ADDR) and row address (R-ADDR) are shipped in a serial order
- Data should be then transferred, and the second command initiates a write
- Once the write is performed, users need to send the last command to check up the status of target

# Page Write

- The first command is needed to indicate that this is a page write
- The column address (C-ADDR) and row address (R-ADDR) are shipped in a serial order
- Data should be then transferred, and the second command initiates a write
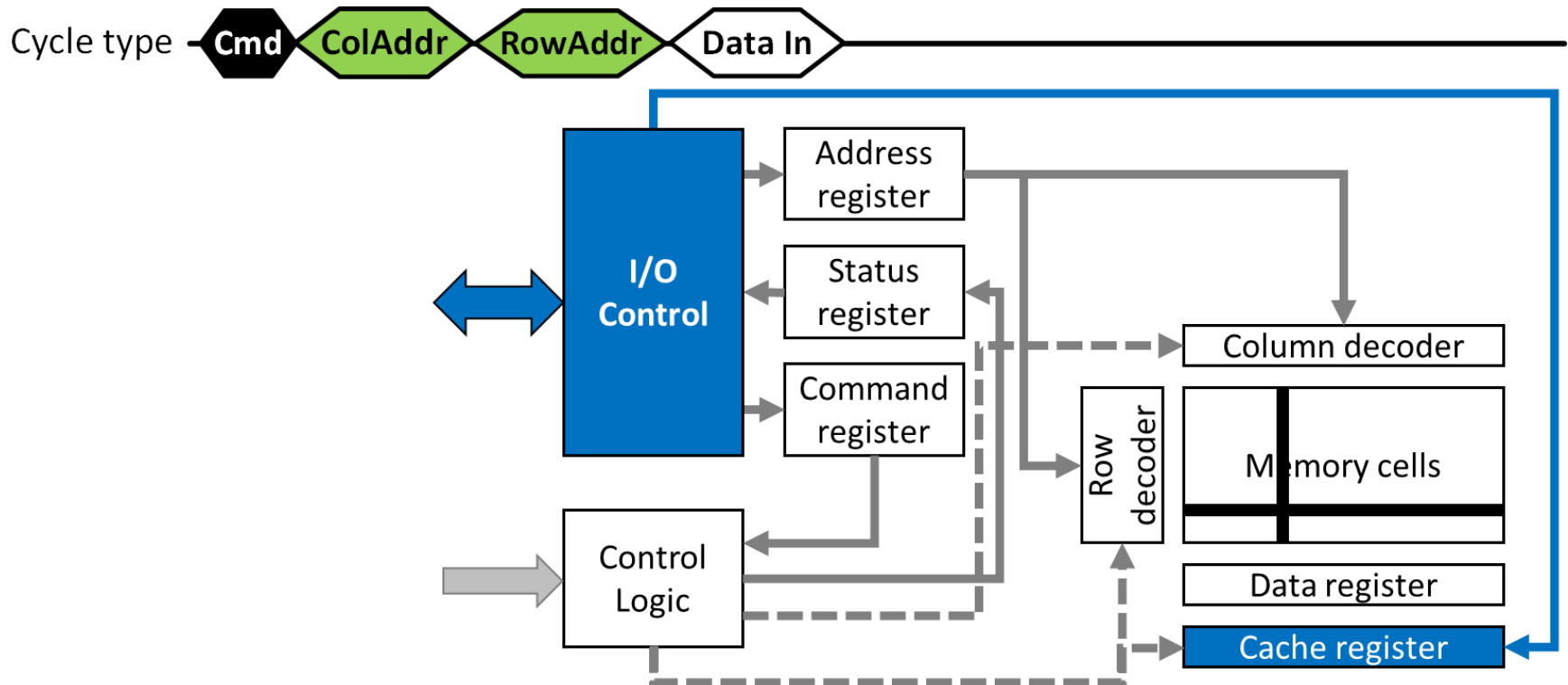- Once the write is performed, users need to send the last command to check up the status of target

# Page Write

■ The first command is needed to indicate that this is a page write

■ The column address (C-ADDR) and row address (R-ADDR) are shipped in a serial order

■ Data should be then transferred, and the second command initiates a write

■ Once the write is performed, users need to send the last command to check up the status of target
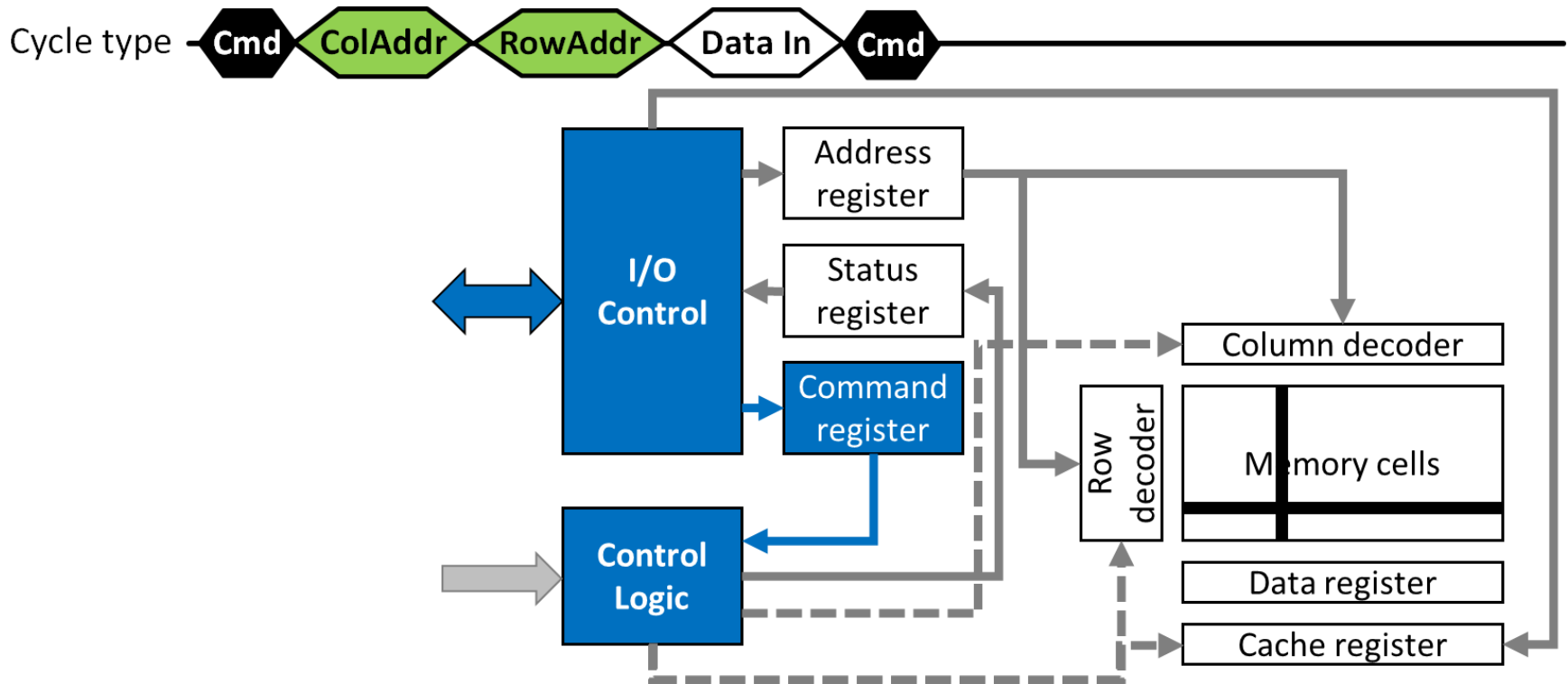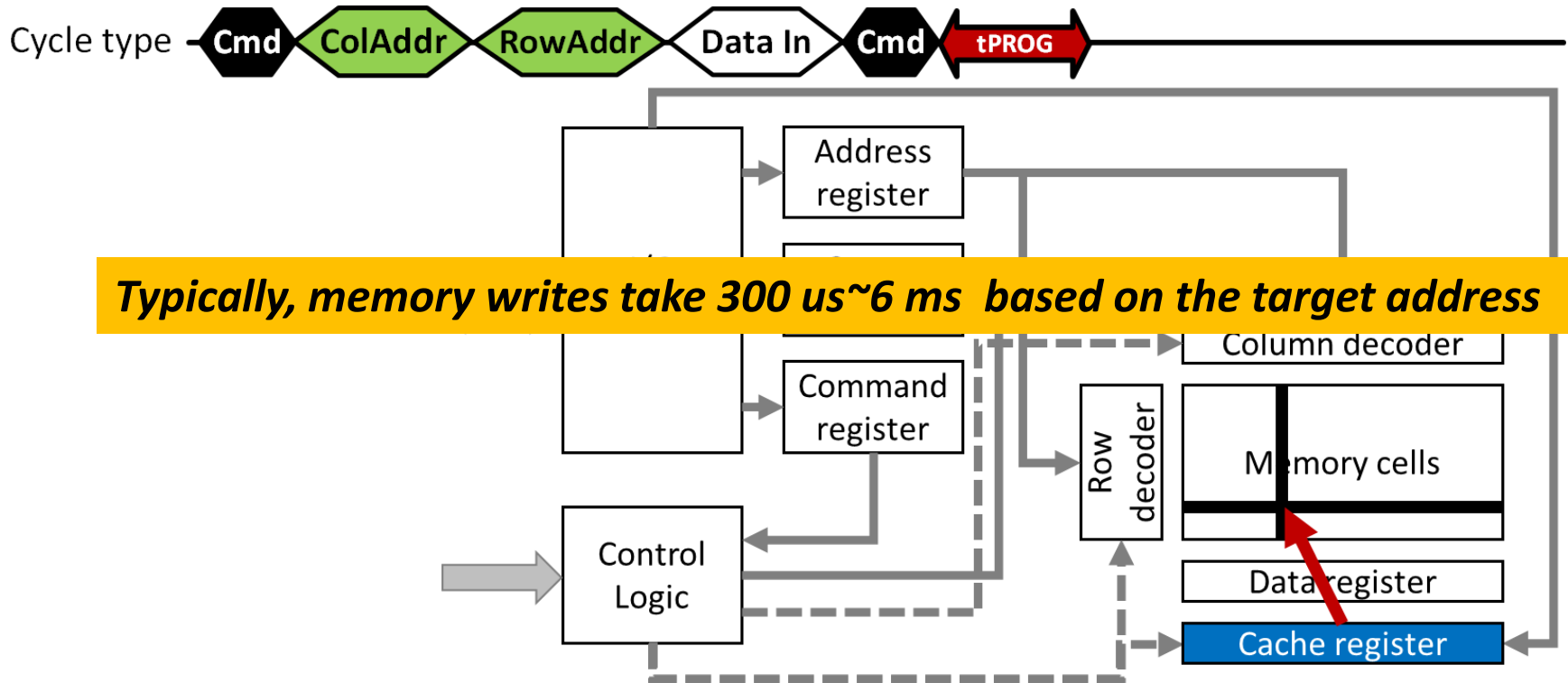
# Page Write

- The first command is needed to indicate that this is a page write

- The column address (C-ADDR) and row address (R-ADDR) are shipped in a serial order

- Data should be then transferred, and the second command initiates a write

- Once the write is performed, users need to send the last command to check up the status of target



*Typically, memory writes take 300 us~6 ms  based on the target address*

# Waveform: Page Read

**Figure 18:    PAGE READ Operation**

| | |
|---|---|
| CLE | |
| CE# | |
| WE# | |
| ALE | |
| R/B# | $t_R$ |
| RE# | |
| I/Ox | 00h — Address (5 Cycles) — 30h — Data Output (Serial Access) |

read
cmd

read
from nand

///// Don't Care

**24**

# Waveform: Page Write

**Figure 23: PROGRAM and READ STATUS Operation**



R/B#

$t_{PROG}$

I/Ox: 80h — Address (5 cycles) — D$_{IN}$ — 10h — 70h — Status

I/O 0 = 0  PROGRAM successful
I/O 0 = 1  PROGRAM error

# Waveform: Block Erasure



Figure 28: BLOCK ERASE Operation

# Outline

- **Storage Abstraction & Protocols**
- **Flash Packages**
- **SSD Controller**
- **Flash Array & I/O Parallelism**

# SSD Overview

# Simplest Controller Architecture

- **ARM CPU deals with host commands, managing NAND packages directly**



- Problems
  - Excessive firmware overhead
  - System bus bottleneck
  - Flash bus bottleneck

Processor Core

Data/Control

Commands / addresses / status

System bus

Host Interface

Flash Controller

Flash bus (40 MB/sec)

Flash Pkg.

Flash Pkg.

Flash Pkg.

Flash Pkg.

...

Processor Core          busy
(Host Interface, Flash Controller)

System bus          bottleneck
Flash bus          bottleneck

Bandwidth = 500 MB/sec
        (SATA 3.0)

# High-level Flash Memory Controller

- **Offload some management duty to high-level controller**
- **Use high-bandwidth bus for fast data transfer**



Processor Core

Control

System bus

Data

Host Interface

25
Processor
Controller

Bandwidth = 500 MB/sec
   (SATA 3.0)

High-level
Flash Memory
Controller

Low-level
Flash Memory
Controller

Flash
Pkg.    Flash
Pkg.    Flash
Pkg.    Flash
Pkg.    ...

- Problems
  - ~~Excessive firmware overhead~~
  - ~~System bus bottleneck~~
  - Flash bus bottleneck

High-level requests
e.g., read 24 sectors from page 1 of block 7 of chip 1

Low-level requests
e.g., read 4 sectors from page 1 of block 7 of chip1
      read 4 sectors from page 2 of block 7 of chip1
      read 4 sectors from page 3 of block 7 of chip1

Flash bus (< 40 MB/sec)

# Bus-level Interleaving

■ **Add more buses for a group of flash packages and send commands in the interleaved manner**

- Problems
  - ~~Excessive firmware overhead~~
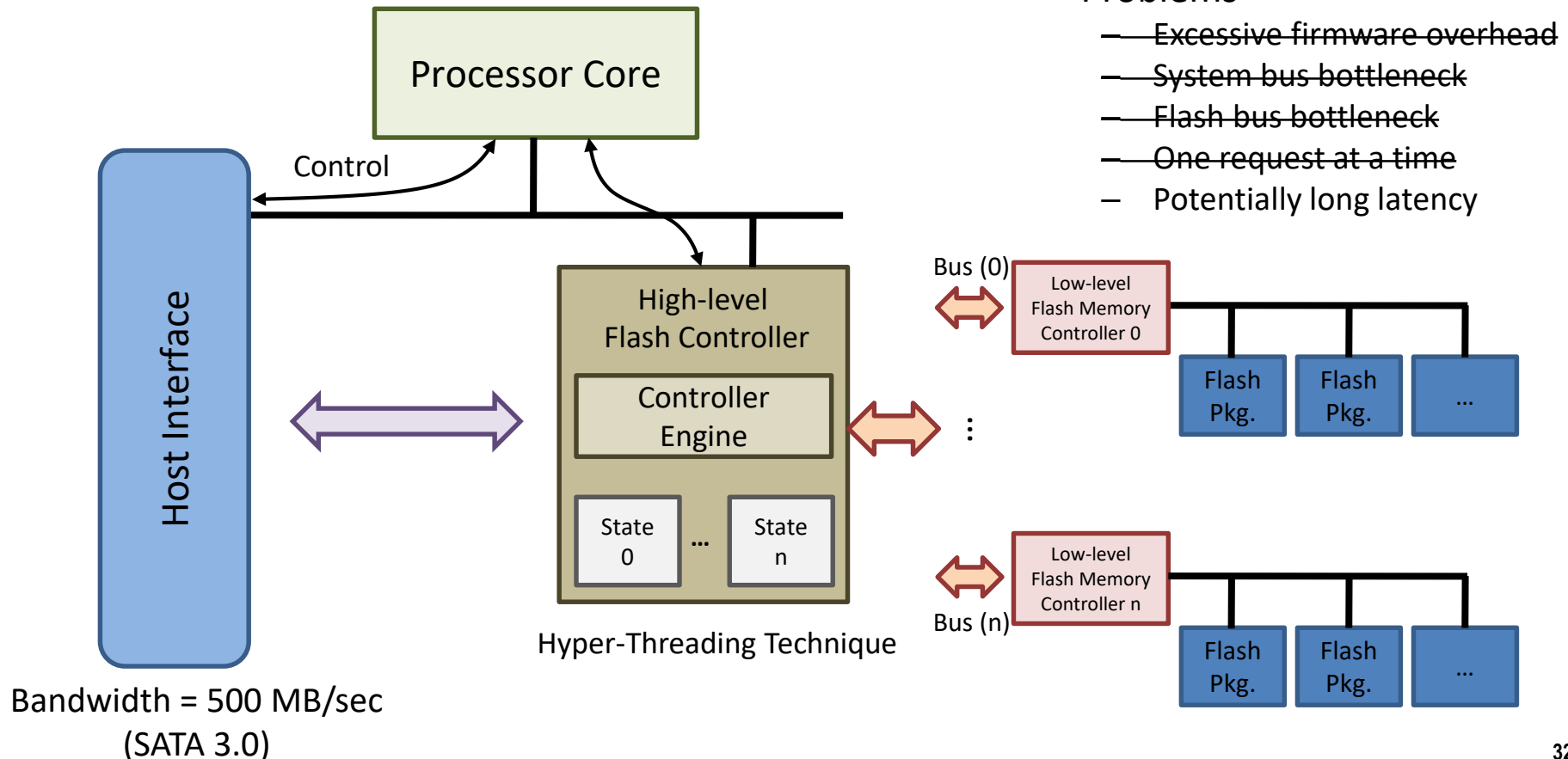  - ~~System bus bottleneck~~
  - ~~Flash bus bottleneck~~
  - One request at a time

Processor Core

Control

System bus

Host Interface

High-level
Flash
Memory
Controller

Bus (0)

Low-level
Flash Memory
Controller 0

Flash Pkg.　Flash Pkg.　...

Bus (n)

Low-level
Flash Memory
Controller n
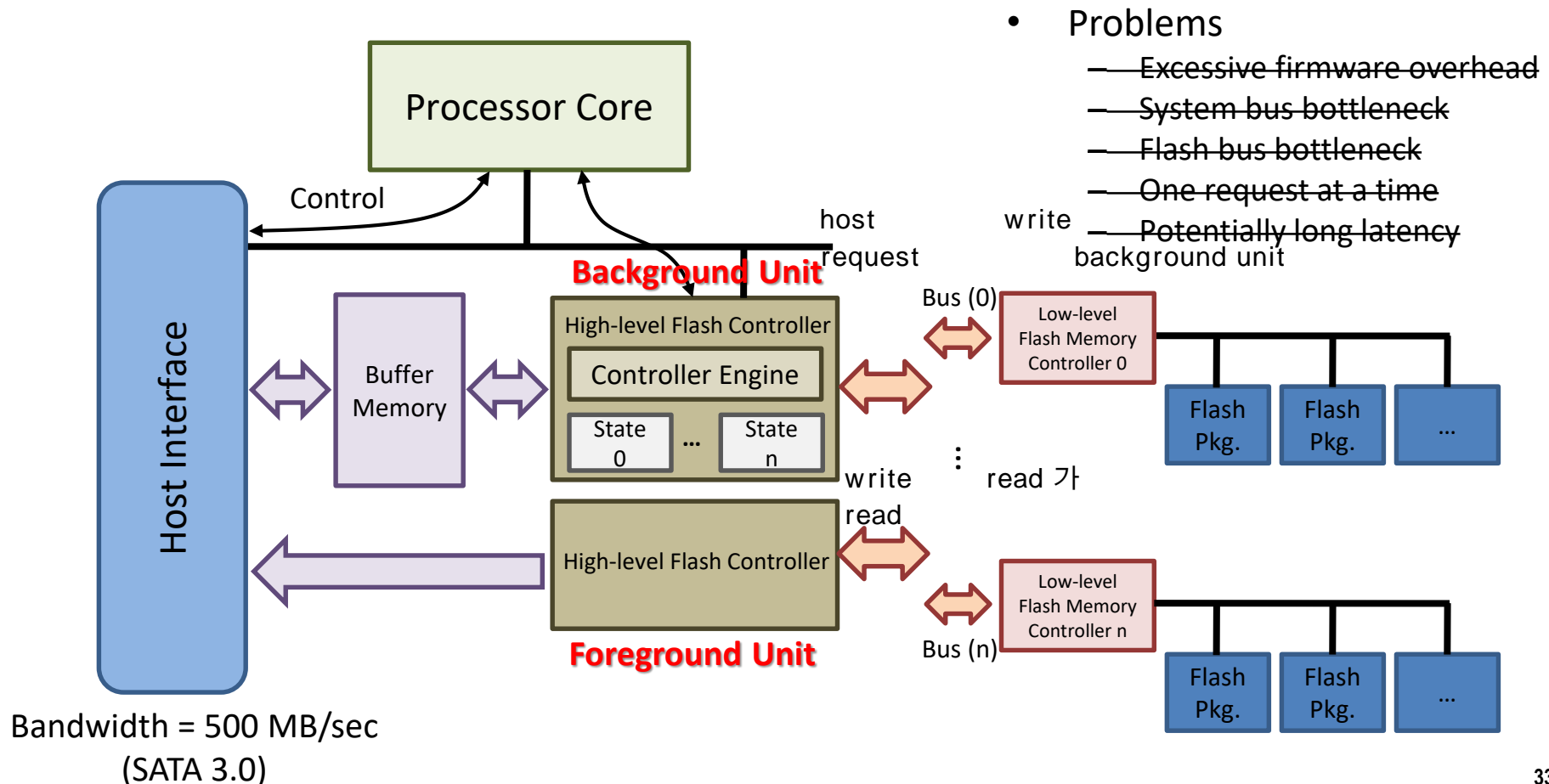
Flash Pkg.　Flash Pkg.　...

Bandwidth = 500 MB/sec
(SATA 3.0)

31

# Multiple Controller Units

■ **Enhance the high-level controller so that it is able to cope with multiple buses at the same time**
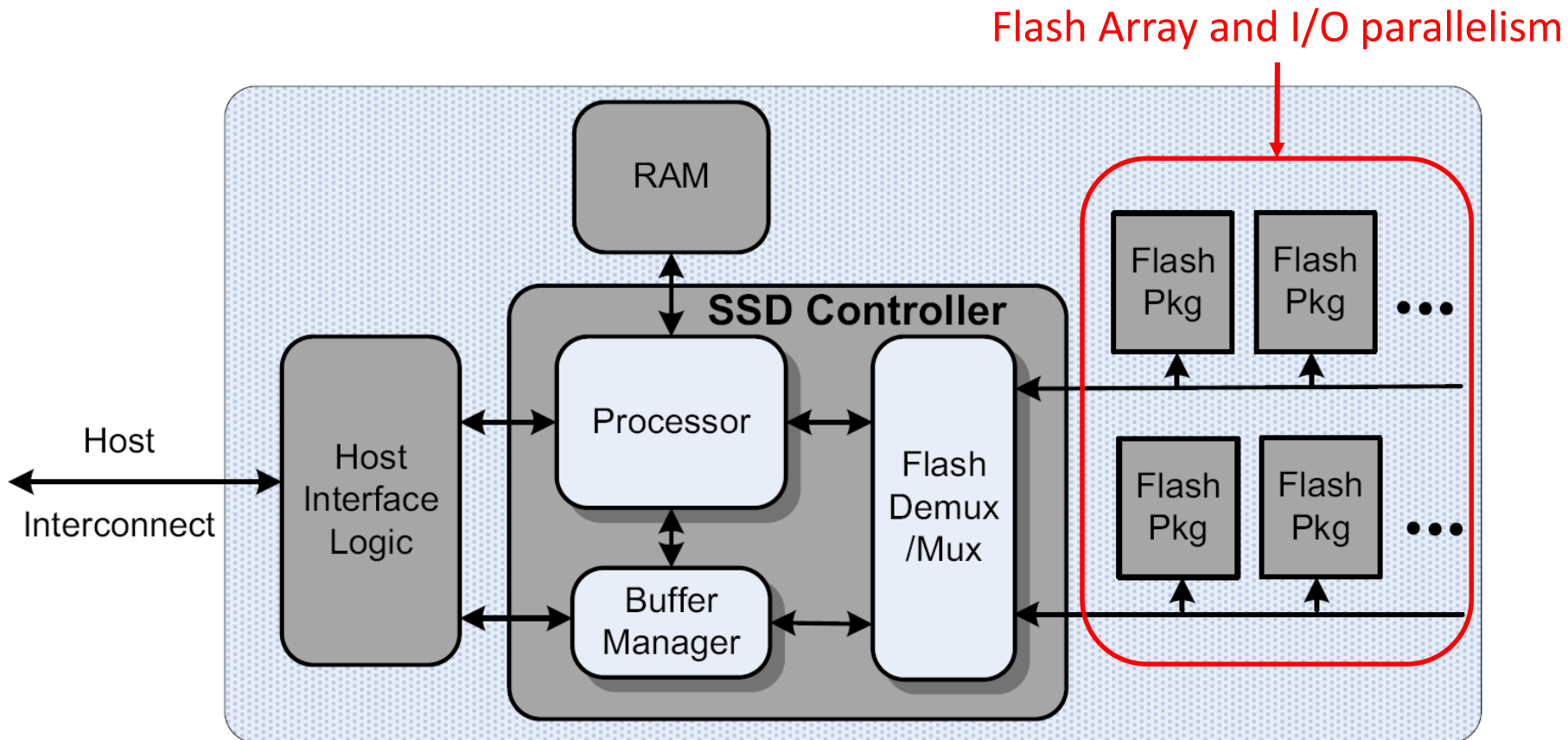


- Problems
  - ~~Excessive firmware overhead~~
  - ~~System bus bottleneck~~
  - ~~Flash bus bottleneck~~
  - ~~One request at a time~~
  - Potentially long latency

Processor Core

Control

Host Interface

High-level Flash Controller

Controller Engine

State 0 ... State n

Hyper-Threading Technique

Bus (0)

Low-level Flash Memory Controller 0

Flash Pkg. | Flash Pkg. | ...

Bus (n)

Low-level Flash Memory Controller n

Flash Pkg. | Flash Pkg. | ...

Bandwidth = 500 MB/sec (SATA 3.0)

# Buffering and Prioritized Handling

- **Manage reads and writes with different priority**
- **Buffer writes in DRAM temporary and flush them later**

- Problems
  - ~~Excessive firmware overhead~~
  - ~~System bus bottleneck~~
  - ~~Flash bus bottleneck~~
  - ~~One request at a time~~
  - ~~Potentially long latency~~

Processor Core

Control

host request

write

background unit

**Background Unit**

Bus (0)

High-level Flash Controller

Controller Engine

State 0 ... State n

Low-level Flash Memory Controller 0

Flash Pkg.  Flash Pkg.  ...

Host Interface

Buffer Memory

write read

read

**Foreground Unit**

High-level Flash Controller

read

Low-level Flash Memory Controller n

Bus (n)

Flash Pkg.  Flash Pkg.  ...

Bandwidth = 500 MB/sec
(SATA 3.0)

33

# Outline

- **Storage Abstraction & Protocols**
- **Flash Packages**
- **SSD Controller**
- **Flash Array & I/O Parallelism**

# SSD Overview



Flash Array and I/O parallelism

# NAND Chip Parameters

| | |
|---|---|
| Page Read to Register | $25\mu s$ |
| Page Program (Write) from Register | $200\mu s$ |
| Block Erase | 1.5ms |
| Serial Access to Register (Data bus) | $100\mu s$ |
| Die Size | 2 GB |
| Block Size | 256 KB |
| Page Size | 4 KB |
| Data Register | 4 KB |
| Planes per die | 4 |
| Dies per package (2GB/4GB/8GB) | 1,2 or 4 |
| Program/Erase Cycles | 100 K |

(*Note: very old parameters)  2007

- **One page read**
  - Read a page from the media into the data register (25 us)
  - Shift 4 KB data out over the data bus (100 us)
- **One page write**
  - Shift 4 KB data into the data register (100 us)
  - Write a page from the data register to a flash page (200 us)

# The Serial Bus is a Primary Bottleneck

- ## Two page reads

Time to send control signals ( < 1 us)

Time to move 4 KB page data to on-chip register (25 us)

Time to move 4 KB on-chip data to flash page (200 us)

Time to move 4 KB page data to the controller through the serial bus (100 us)

Request 1 (4 KB = 1 page)    Request 2 (4 KB = 1 page)

Total time to read 8 KB page = 250 us

- ## Solution – Use parallelism of a flash array

  - (1) *Plane-level Pipelining*, (2) *Interleaving*, and (3) *Stripping*

# Flash Array

- **SSD handles requests on flash packages in parallel**
  - To achieve bandwidths or I/O rates greater than a single-chip



⟨Fully connected flash array architecture⟩

# Plane-level Pipelining: Read

- **Two page reads with pipelining**

bus .

Request 1 (4 KB = 1 page)

control logic
plane

Plane 0

| read | bus |

(Plane-pair)

Request 2 (4 KB = 1 page)

Plane 1

| read | | bus |

Total time to read 8 KB  (2 pages)
= 225 us

- The maximum read bandwidth improves to 20%
  - One page read requires about 100 usec

# Plane-level Pipelining: Write

- **Four page writes with pipelining**



- The maximum write bandwidths improve to 67%
  - One page write requires about 100 usec

# Interleaving Technique

■ **Exploiting inter-request parallelism using multiple flash packages**
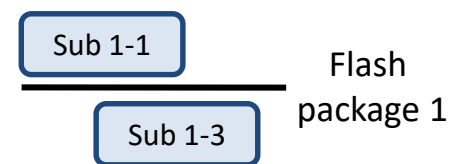


Request 2 — Flash package 2

Request 2   Request 1

Request 1 — Flash package 1

# Stripping Technique

- **Exploiting intra-request parallelism by spreading out a request across multiple packages**
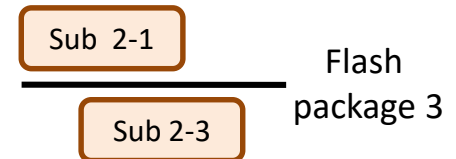
# Putting All Together
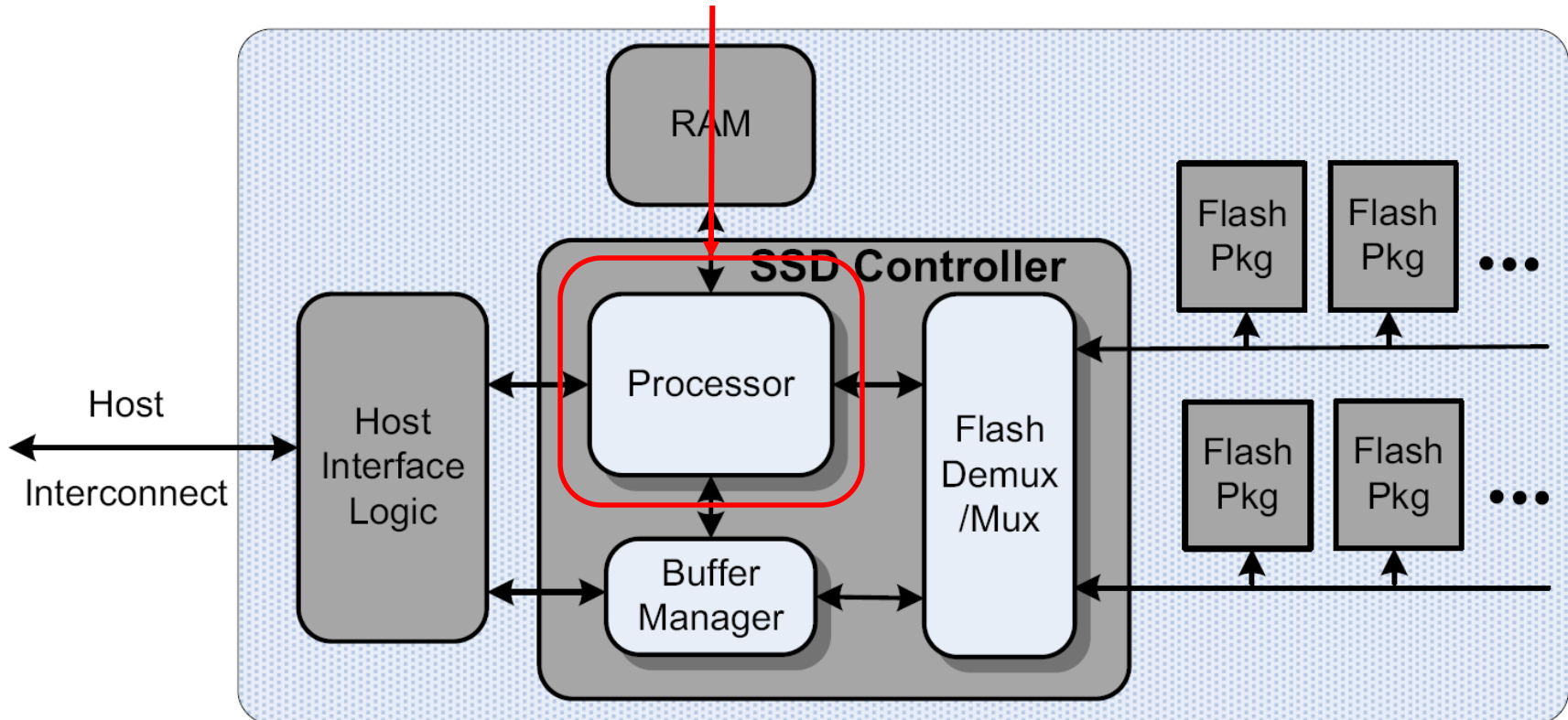
# SSD Firmware Algorithms



SSD Firmware Algorithms

# SSD Firmware Algorithms (Cont.)

- **Some form of mapping between logical address and physical flash location is required in NAND flash**
  - The SSD controller receives a list of LBAs with commands, which must be mapped to specific pages of blocks in flash packages

- **Major concerns in designing a logical block map (address mapping) algorithm**
  - Mapping table size
  - Garbage collection
  - Load balancing
  - Parallel access

  *Decide SSD performance & reliability*
  *(See in detail next week)*

*End of Chapter 3*