

# Persistence: Solid-State Storage Devices

## CS 537: Introduction to Operating Systems

Louis Oliphant & Tej Chajed

University of Wisconsin - Madison

Spring 2024

## Administrivia

- **Discussion Sessions Cancelled This Week (Instead Additional Office Hours)**
- Project 7 due April 30th @ 11:59pm
- Final Exam:
  - Lec 1 - May 8th, 12:25-2:25 (Biochem 1125)
  - Lec 2 - May 6th, 2:45-4:45 (Sterling Hall 1310)
  - McBurney: May 6th, 2:40-6:50 (Nancy Nicholas Hall 1135)
  - If you can't take the exam for a *legitimate reason* at your designated time, please fill out the [alternate exam form](#) to take the exam with the other lecture. Legitimate Reasons include:
    - Another exam at the same time, Religious conflict, University Sanctioned conflict, Scheduled Medical conflict, Civic Duty (e.g. jury duty), Military Service, Family Caregiving Responsibility, Family Emergency, Serious Illness, 3 or more exams scheduled during a 24 hour period

# Review FSCK, Journaling & Log-Structured File Systems

- FSCK
  - fsck attempts to scan and correct inconsistencies found in the file system.
  - build **used data blocks** from inode table, checks inodes and directory entries for consistency
- Data Journaling and Metadata (or ordered) Journaling
  - Understand protocol of what gets written where and what waits occur to insure consistency
- Log-structured File System
  - Layout on disk – checkpoint region, segments (data, inodes, imap, segment summary),
  - Memory caching – imap and buffered writes
  - Garbage Collection – block liveness, which blocks to clean
  - Crash Recovery – multiple CRs, roll forward

## Quiz 21 FSCK & LFS

<https://tinyurl.com/cs537-sp24-q21>



# Solid-State Storage Devices

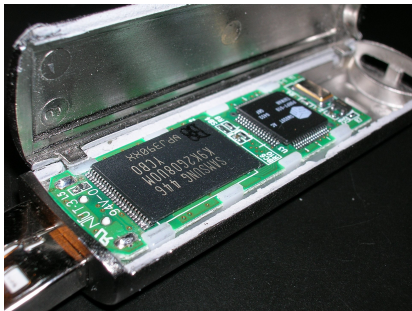
- Physical Storage System
  - SLC, MLC, TLC
  - Banks, Blocks, and Pages
- Flash-based Operations
  - Read (a page), Erase (a block), Program (a page)
- Flash Translation Layer (FTL)
- Log-Structured FTL
- Garbage Collection
- Mapping Tables
- SSD Performance and Cost

# NAND FLASH

Single Level Cell (SLC) = 1 bit per cell  
(faster, more reliable)

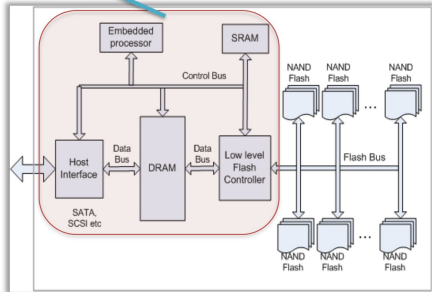
Multi Level Cell (MLC) = 2 bits per cell  
(slower, less reliable)

Triple Level Cell (TLC) = 4 bits per cell  
(even more so)



# SSD STRUCTURE

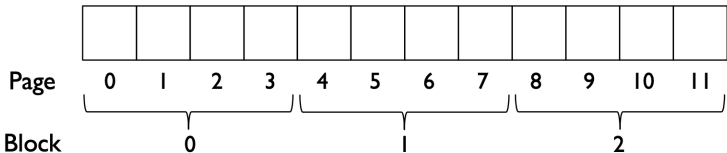
Flash Translation Layer  
(Proprietary firmware)



Simplified block diagram of an SSD

# SSD PROPERTIES

Page ~ 4KB,  
Block ~ 128 KB  
or 256 KB



Read

Write

Failures: Block likely to fail after a certain number of erases  
(~10000 for MLC flash, ~100,000 for SLC flash)



# SSD OPERATIONS

Read a page: Retrieve contents of entire page (e.g., 4 KB)

- Cost: 25—75 microseconds
- Independent of page number, prior request offsets

Erase a block: Resets each page in the block to all 1s

- Cost: 1.5 to 4.5 milliseconds
- Much more expensive than reading!
- Allows each page to be written

Program (i.e., write) a page: Change selected 1s to 0s

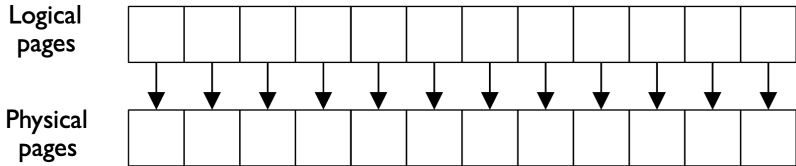
- Cost is 200 to 400 microseconds
- Faster than erasing a block, but slower than reading a page

# FLASH TRANSLATION LAYER

1. Translate reads/writes to logical blocks into reads/erases/programs
2. Reduce write amplification (extra copying needed to deal with block-level erases)
3. Implement wear leveling (distribute writes equally to all blocks)

Typically implemented in hardware in the SSD, but in software for some SSDs

# FTL: DIRECT MAPPING



Cons?



# FTL: LOG-STRUCTURED ADVANTAGES

Avoids expensive read-modify-write behavior

Better wear levelling: writes get spread across pages,  
even if there is spatial locality in writes at logical level

Challenges? Garbage!

# GARBAGE COLLECTION

Table: 100 → 0 101 → 1 2000 → 2 2001 → 3 Memory

---

Block:	0				1				2				
Page:	00	01	02	03	04	05	06	07	08	09	10	11	Flash Chip
Content:	a1	a2	b1	b2									
State:	V	V	V	V	i	i	i	i	i	i	i	i	

Table: 100 → 4 101 → 5 2000 → 2 2001 → 3 Memory

---

Block:	0				1				2				
Page:	00	01	02	03	04	05	06	07	08	09	10	11	Flash Chip
Content:	a1	a2	b1	b2	c1	c2							
State:	V	V	V	V	V	V	E	E	i	i	i	i	

# GARBAGE COLLECTION

Steps:

Read all pages in physical block

Write out the alive entries to the end of the log

Erase block (freeing it for later use)

Table: 100 → 4 101 → 5 2000 → 2 2001 → 3 Memory

Block:	0				1				2				
Page:	00	01	02	03	04	05	06	07	08	09	10	11	Flash Chip
Content:	a1	a2	b1	b2	c1	c2							
State:	V	V	V	V	V	V	E	E	i	i	i	i	

Table: 100 → 4 101 → 5 2000 → 6 2001 → 7 Memory

Block:	0				1				2				
Page:	00	01	02	03	04	05	06	07	08	09	10	11	Flash Chip
Content:					c1	c2	b1	b2					
State:	E	E	E	E	V	V	V	V	i	i	i	i	

# OVERHEADS

Garbage collection requires extra read+write traffic

Overprovisioning makes GC less painful

- SSD exposes logical space that is smaller than the physical space
- By keeping extra, “hidden” pages around, the SSD tries to defer GC to a background task (thus removing GC from critical path of a write)

Occasionally shuffle live (i.e., non-garbage) blocks that never get overwritten

- Enforces wear levelling

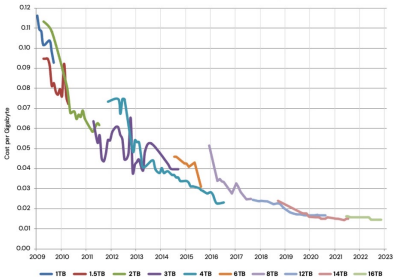


# OVERALL PERFORMANCE

Device	Random		Sequential	
	Reads (MB/s)	Writes (MB/s)	Reads (MB/s)	Writes (MB/s)
Samsung 840 Pro SSD	103	287	421	384
Seagate 600 SSD	84	252	424	374
Intel SSD 335 SSD	39	222	344	354
Seagate Savvio 15K.3 HDD	2	2	223	223

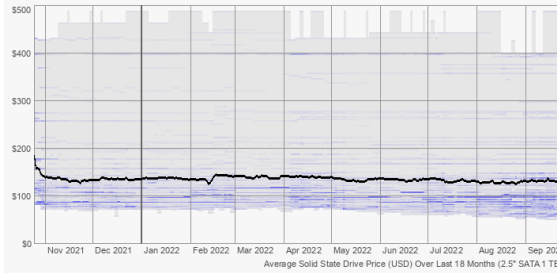
# COST?

**Backblaze Average Cost per Gigabyte by Drive Size Over Time**  
Drive sales grouped by drive size and month to compute average cost per month



Backblaze

~1.5 cents / GB



1TB ~ \$150 on average  
~15 cents / GB

Next Time – Distributed Systems