[544] File Systems

Meenakshi Syamkumar

Learning Objectives

- compare the performance characteristics of different kinds of block devices (HDDs and SSDs)
- describe different kinds of file systems
- interpret the output of tools like "mount" and "df" to understand the structure of a mount namespace

Outline

Block Devices (overview, HDD, SSD)

File Systems

Demos

Block Devices

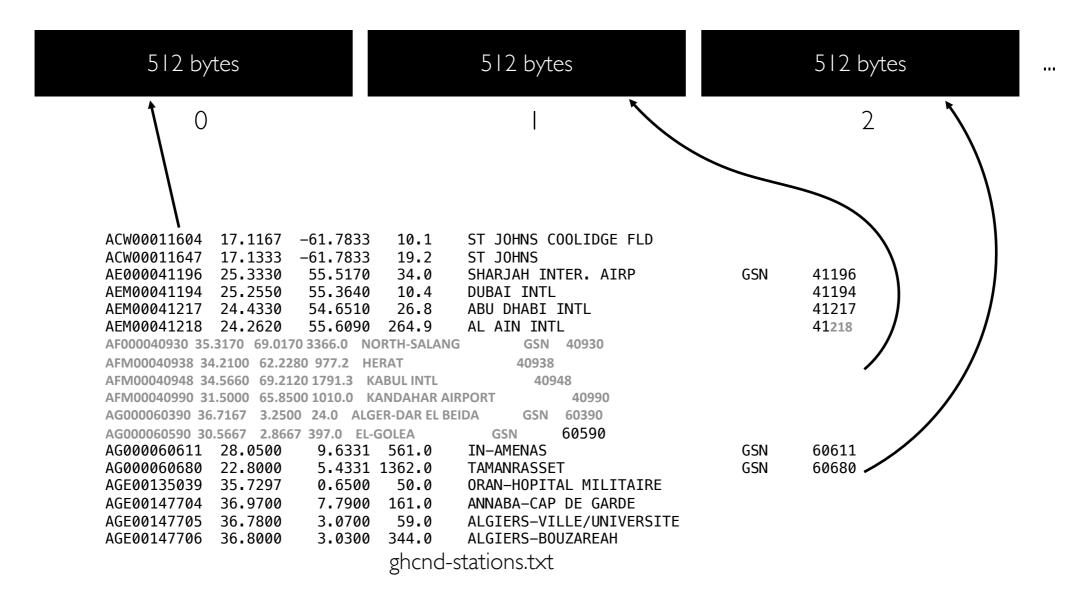
Memory is byte addressable



Block storage devices are accessed in units of blocks (512 bytes, few KBs, etc)



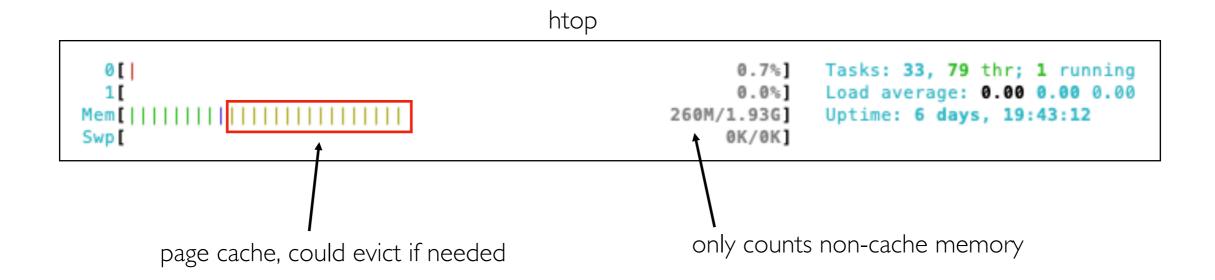
Optimizing Disk I/O with Memory: Caching and Buffering



We might want to process one line a time, but it would be wasteful to repeatedly read the same block from the device

- the Linux page cache stores pages from files in RAM (usually 4KB pages, often larger than device blocks)
- Python (and other) programs might buffer chunks of data to avoid asking Linux too many times for small pieces of data

Optimizing Disk I/O with Memory: Caching and Buffering



We might want to process one line a time, but it would be wasteful to repeatedly read the same block from the device

- the Linux page cache stores pages from files in RAM (usually 4KB pages, often larger than device blocks)
- Python (and other) programs might buffer chunks of data to avoid asking Linux too many times for small pieces of data

Small Reads (<4KB): Performance

goal: collect all station IDs

```
ACW00011604
            117.1167 -61.7833
                                  10.1
                                          ST JOHNS COOLIDGE FLD
ACW00011647
            17.1333
                      -61.7833
                                  19.2
                                          ST JOHNS
            25.3330
                       55.5170
                                  34.0
                                                                           GSN
AE000041196
                                          SHARJAH INTER. AIRP
                                                                                   41196
AEM00041194 25.2550
                       55.3640
                                                                                   41194
AEM00041217 | 24.4330
                       54.6510
                                26.8
                                          ABU DHABI INTL
                                                                                   41217
AEM00041218 | 24.2620
                       55.6090 264.9
                                          AL AIN INTL
                                                                                   41218
AF000040930 35.8170 69.0170 3366.0 NORTH-SALANG
                                                 GSN 40930
                                                                                          ghcnd-stations.txt
AFM00040938 34 2100 62.2280 977.2 HERAT
                                                 40938
AFM00040948 34 5660 69.2120 1791.3 KABUL INTL
AFM00040990 31 5000 65.8500 1010.0 KANDAHAR AIRPORT
                                                       40990
AG000060390 36.7167 3.2500 24.0 ALGER-DAR EL BEIDA
                                                 GSN 60390
                                                      60590
AG000060590 30.5667 2.8667 397.0 EL-GOLEA
                                             GSN
AG000060611 28.0500
                         9.6331 561.0
                                           IN-AMENAS
                                                                           GSN
                                                                                   60611
AG000060680 22.8000
                         5.4331 1362.0
                                          TAMANRASSET
                                                                           GSN
                                                                                   60680
            35.7297
AGE00135039
                         0.6500
                                50.0
                                          ORAN-HOPITAL MILITAIRE
                         7.7900 161.0
AGE00147704
            36.9700
                                          ANNABA-CAP DE GARDE
AGE00147705
            36.7800
                         3.0700
                                59.0
                                          ALGIERS-VILLE/UNIVERSITE
AGE00147706 36.8000
                         3.0300 344.0
                                           ALGIERS-BOUZAREAH
```

```
start = time.time()
with open("ghcnd-stations.txt") as f:
    for line in f:
        stations.append(line[:11])
print(time.time() - start)
```

simple version that reads everything: 66 ms

format issue: no good way to ready one column without everything else

(similar to issues with bad cache line usage)

"optimized" version that only reads stations: 171 ms

Hard Disk Drives (HDDs)

Steps to read/write

- I. move head to correct track
- 2. wait for spinning disk to rotate until data is under head

these steps dominate unless transferring lots of data (few MBs)

3. transfer the data



Layout

- assign block numbers to platter locations so sequential (like 5,6,7,8, ...) reads/writes will be fast
- programmers should assume random accesses (like 2, 9, 5, 1, ...) will be slow

Capacity vs. I/O and Short Stroking

Storage resources

- I. capacity
- 2. I/O (input/output often more limited when using HDDs)



Short Stroking

- head moves over platter faster near outside track
- smaller block addrs correspond to outside tracks
- strategy: only use outside tracks
- pros: faster I/O
- cons: less space

Solid State Drives (SSDs) - Flash

Reading and writing

- no moving parts
- inherently parallel

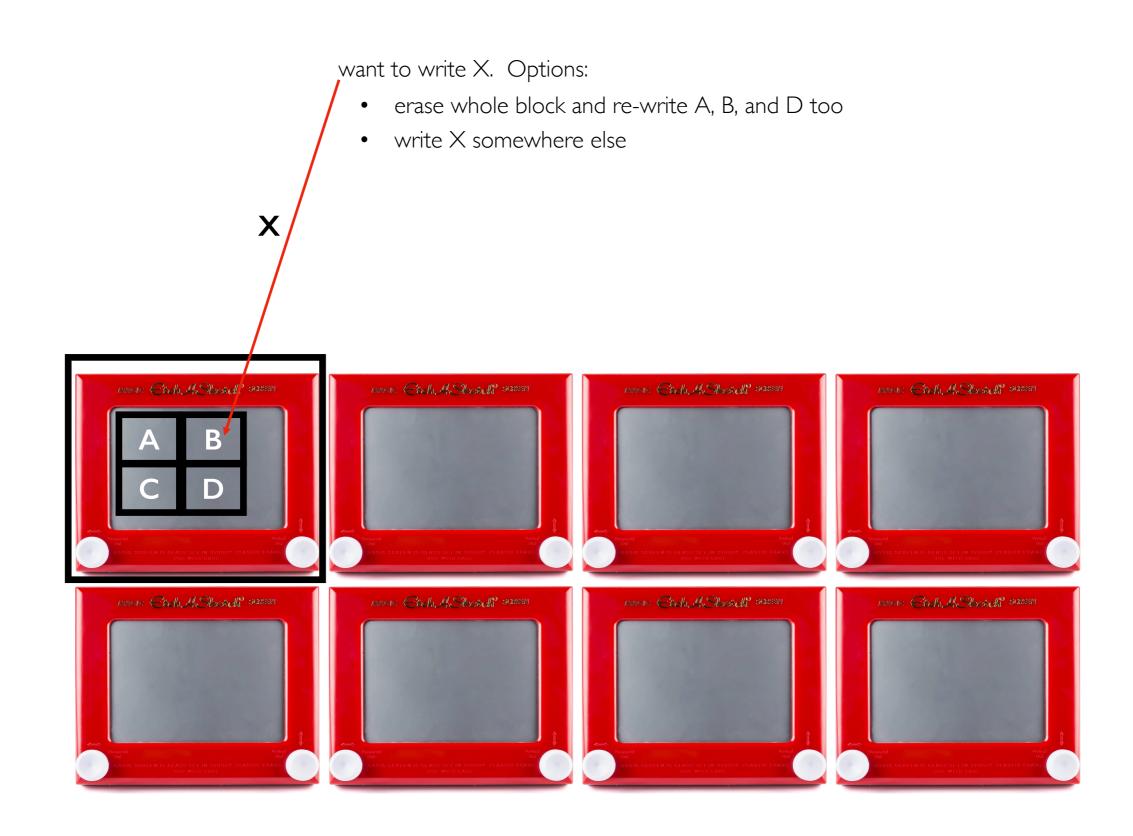
SSD internals:

- "block" and "page" have different meanings in this context
- "page" => unit that we can read or write (couple KBs)
- pages cannot be individually re-written
- "block" => unit that is erased together (maybe 100s of KBs)





Solid State Drives (SSDs) - Flash



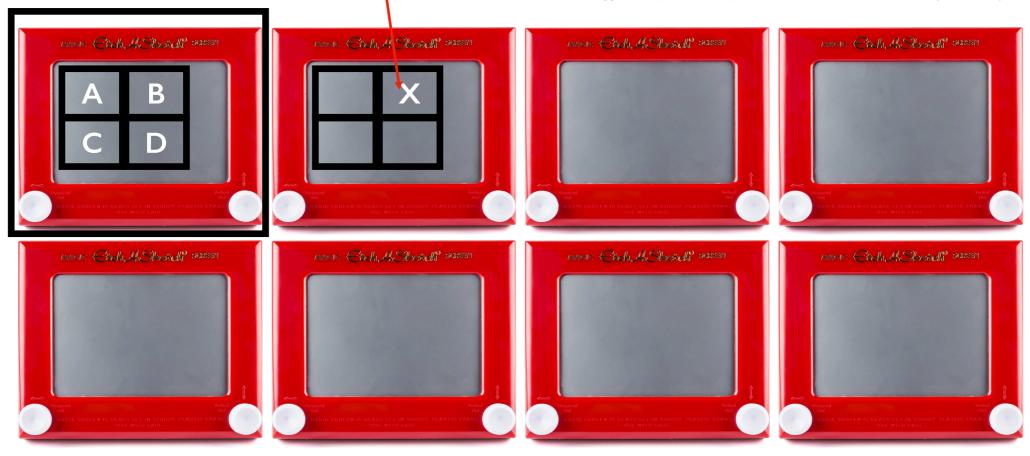
Solid State Drives (SSDs) - Flash

want to write X. Options:

- erase whole block and re-write A, B, and D too
- write X somewhere else

disadvantages

- need extra bookkeeping (in SSD) to know where data is
- need to eventually move things around to reclaim the space wasted by B
- **strategy**: sequentially write whole blocks (when possible)



HDDs vs. SSDs

Metrics

- capacity: how many bytes can we store?
- latency: how long does it take to start transferring data
- IOPS (I/O operations, of some max size, per second): how many small/random transfers can we do per second
- throughput: how many bytes can we transfer per second

Metric: Relative to HDDS, SSDS are:

capacity worse

latency much better (no moving parts)

random IOPS even better -- low latency AND in parallel

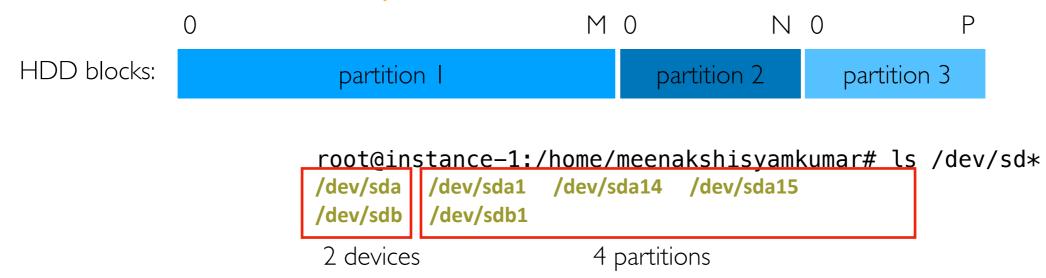
throughput (sequential) little better

throughput (random writes) better (but block erase is a concern)

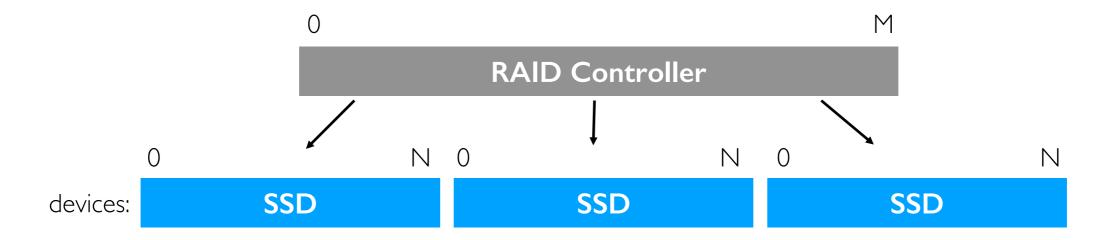
throughput (random reads) much better

Partitions and RAID

Block devices can be divided into partitions:



RAID controllers (Redudant Array of Inexpensive Disks) can make multiple devices appear as one:



Many configs use redundancy (e.g., same data on >1 disk) to avoid data loss when one device dies.

Outline

Block Devices (overview, HDD, SSD)

File Systems

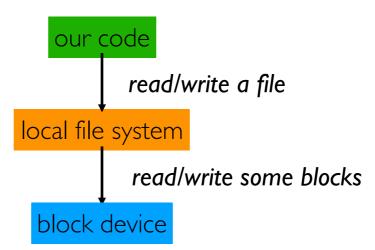
Demos

File Systems

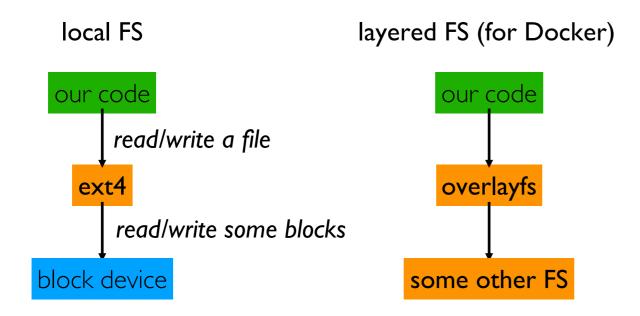
Difficult: writing code to store data in blocks

Easier: writing code to store data in files

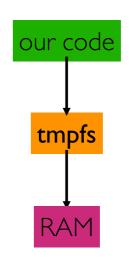
Files systems *abstract* storage for us. We write to data **blocks** without thinking about it by writing data to **files** in a local file system.



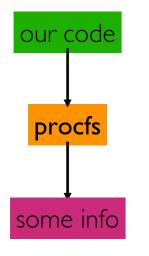
Types of File System (FS)



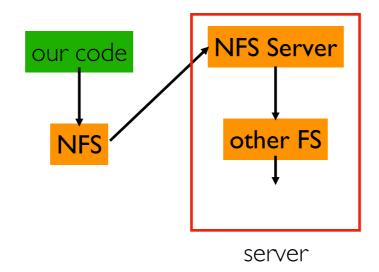
in-memory FS (Temp Files)



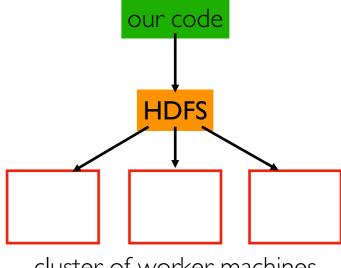




network FS

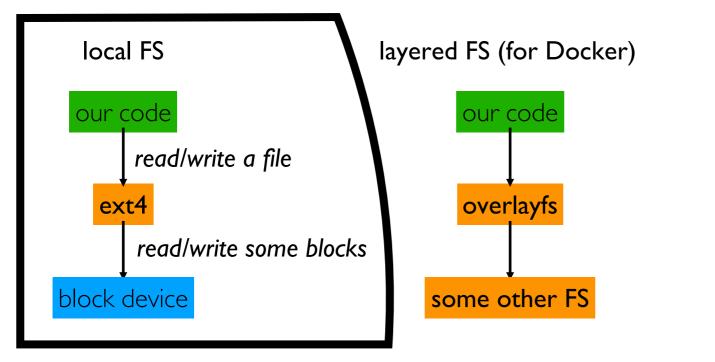


distributed FS

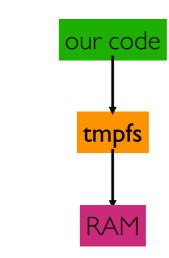


cluster of worker machines

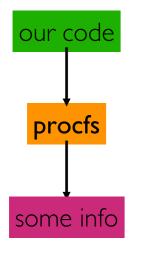
Types of File System (FS)



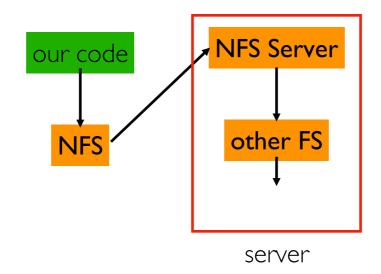
in-memory FS (Temp Files)



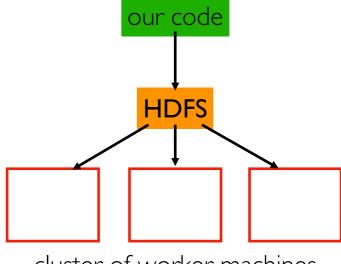




network FS

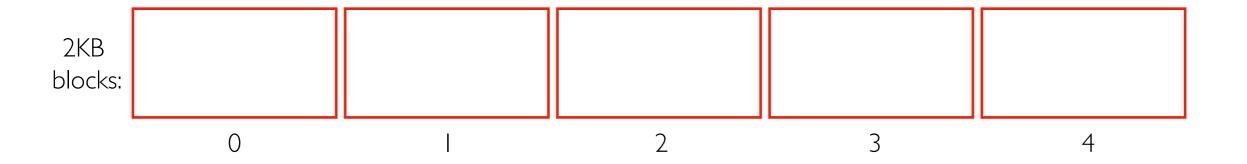


distributed FS



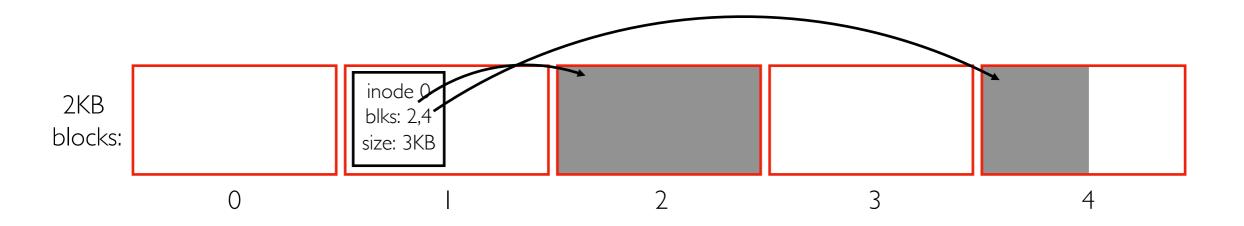
cluster of worker machines

Local File Systems



How does a local FS use blocks?

Local File Systems

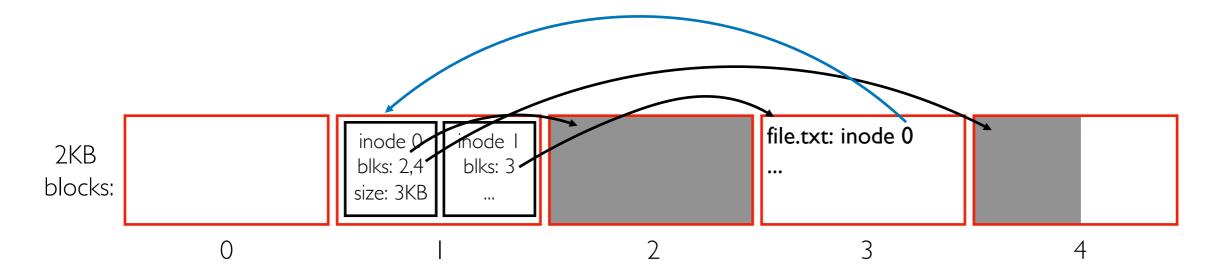


How does a local FS use blocks? Many possibilites. One example...

Files

- some metadata, like size, block locations
- each is represented by an "inode" structure (above file is fragmented)
- file extensions (like .txt) don't mean anything to the file system (just for documentation)

Local File Systems



How does a local FS use blocks? Many possibilites. One example...

Files

- some metadata, like size, block locations
- each is represented by an "inode" structure (above file is fragmented)
- file extensions (like .txt) don't mean anything to the file system (just for documentation)

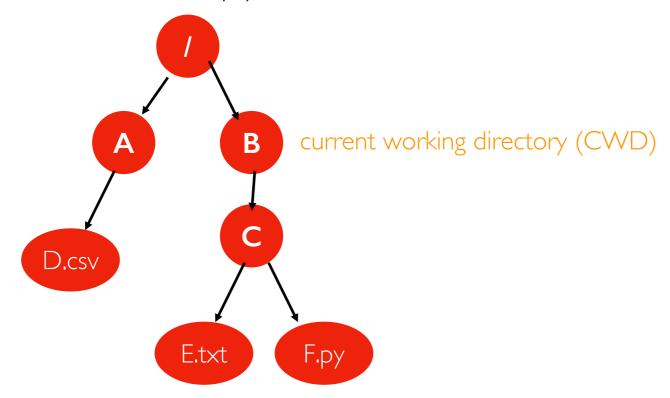
Directories

- special files containing name => inode mappings
- the same inode could be in multiple directories
- each file system has a "root" directory from which you can reach everything else recursively
- formatting creates initial structures (like the root directory)

File System Trees

Nesting of directories and files logically create "trees"

- technically DAGs (directed acyclic graphs) because the same inode number can have multiple names in different directories
- leaves: files and empty directories



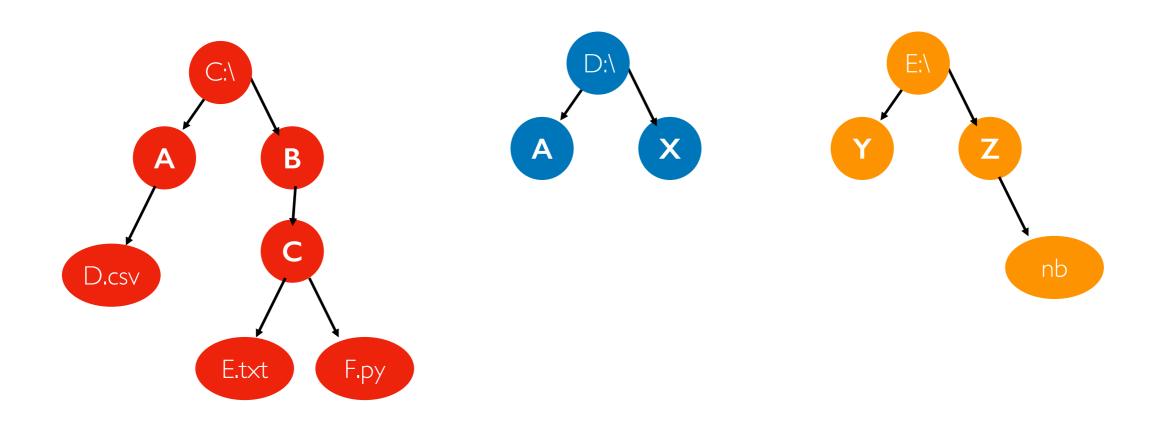
relative path to E.txt: C/E.txt

absolute path to E.txt: /B/C/E.txt relative path to D.csv: ../A/D.csv

absolute path to D.csv: TopHat

Multiple File Systems: Windows Approach

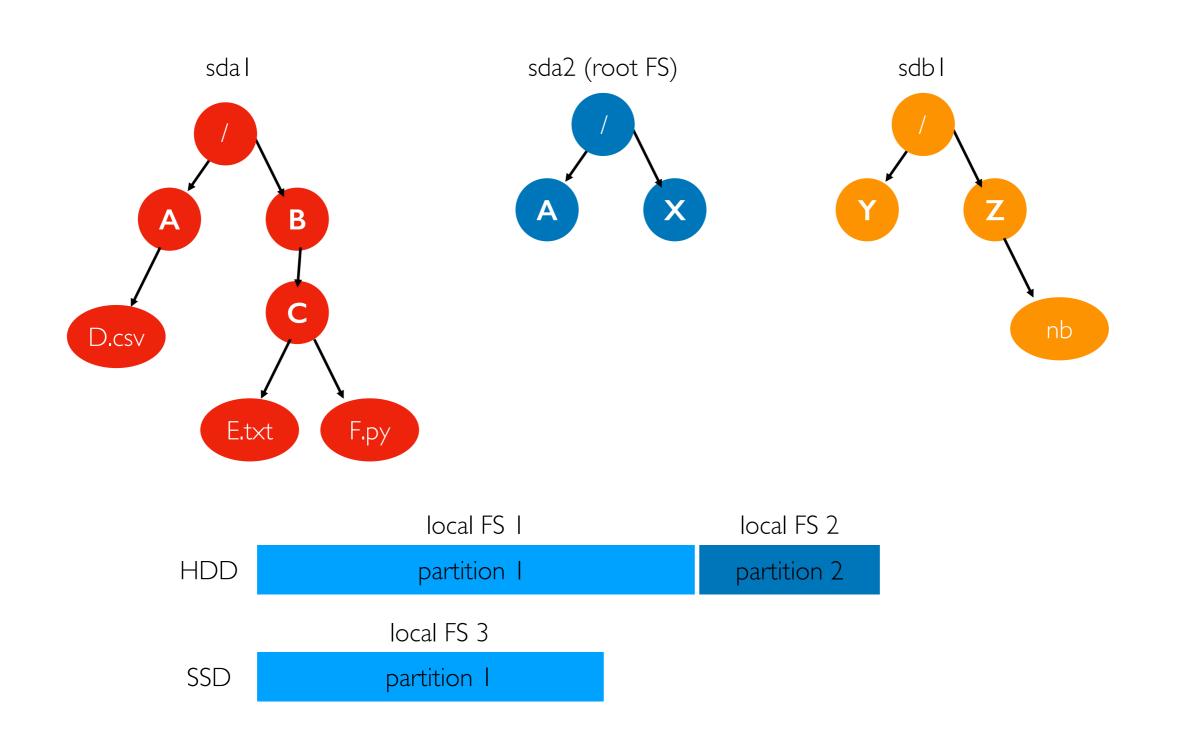
have multiple trees (each is a "drive")

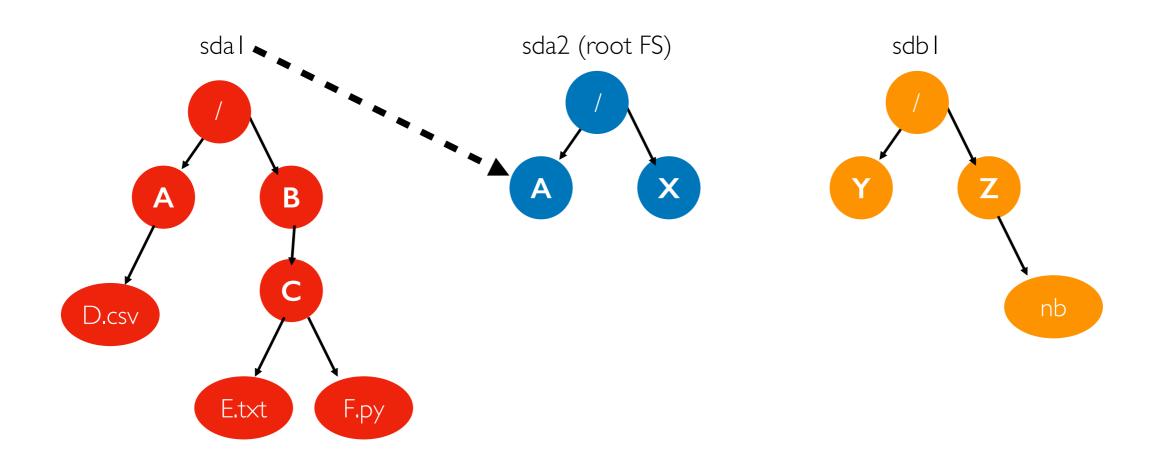


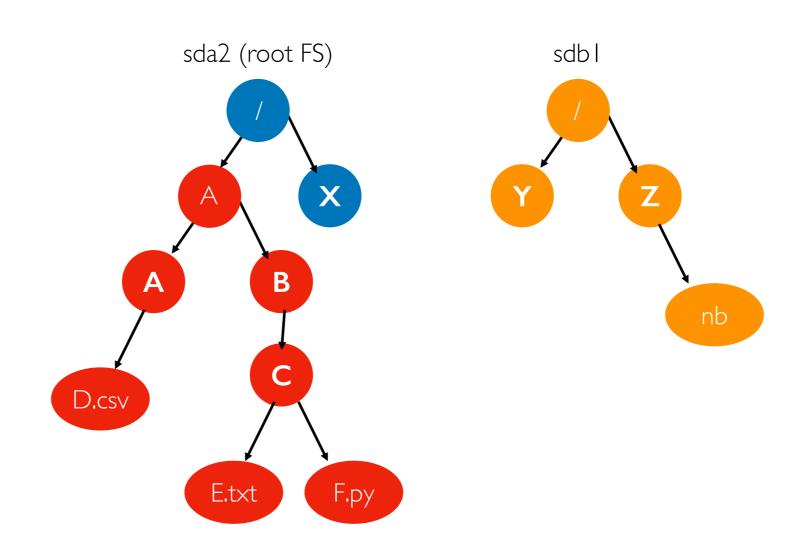




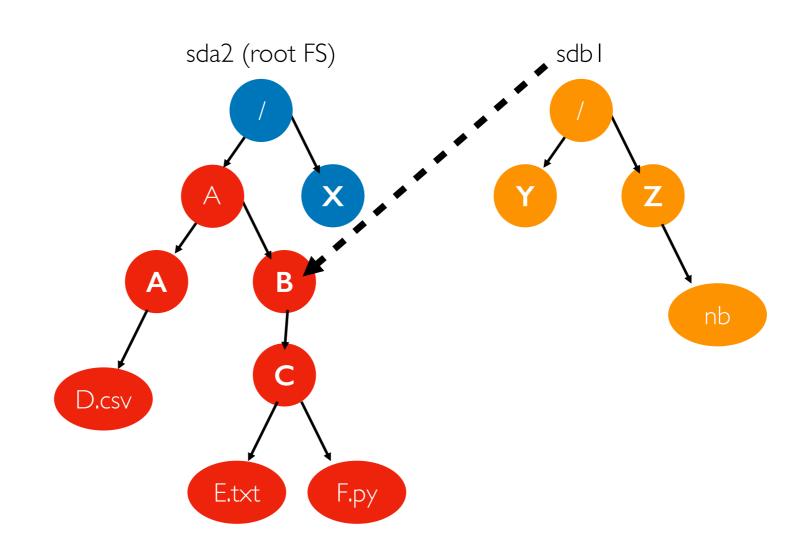
https://www.brit.co/fruit-salad-tree/





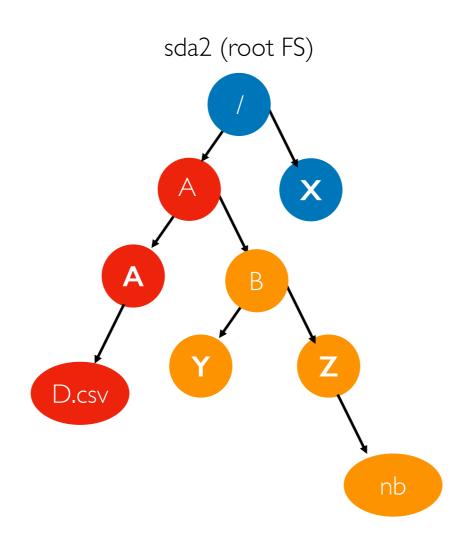


mount file systems over directories of other file systems to make one big tree



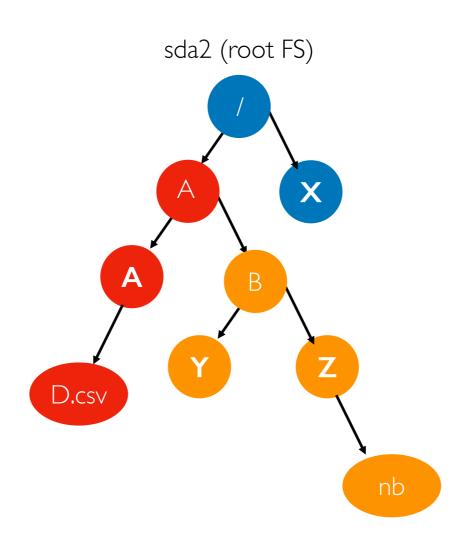
mount /dev/sda1 /A
mount /dev/sdb1 /A/B

mount file systems over directories of other file systems to make one big tree



mount /dev/sda1 /A
mount /dev/sdb1 /A/B

mount file systems over directories of other file systems to make one big tree

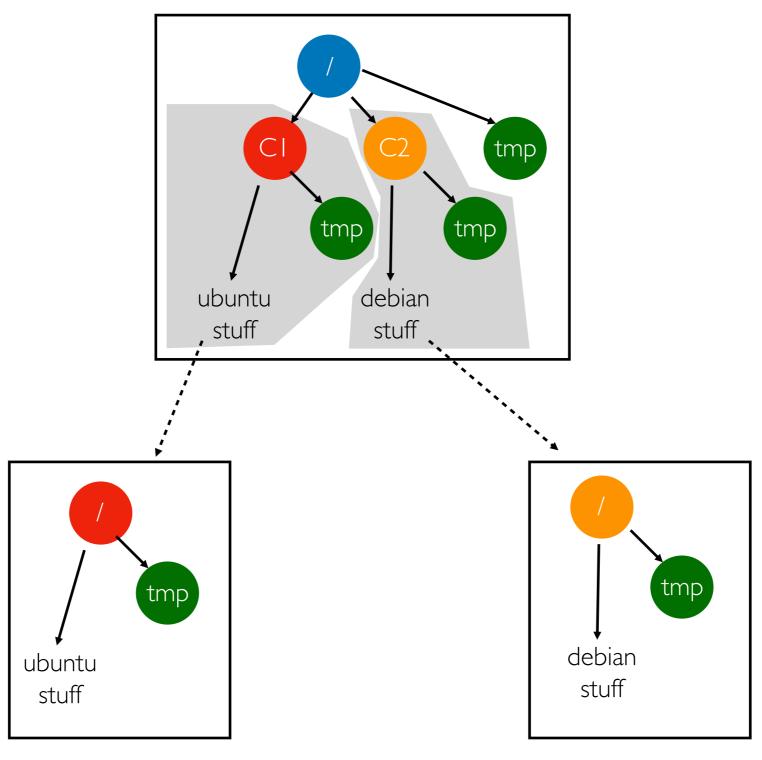


Note: each container has it's own root file system and mount namespace

mount /dev/sda1 /A
mount /dev/sdb1 /A/B

Container File Systems (Simplified)

mount namespace (VM)



Outline

Block Devices (overview, HDD, SSD)

File Systems

Demos