# ECE566 Enterprise Storage Architecture Lab #1: Drives and RAID

Now we'll play with those hard drives you have and learn ways to establish redundancy for them.

NOTE: This assignment assumes familiarity with basic UNIX concepts such as IO redirection and piping; if you need to brush up on this content, you can use <u>this Linux command line course</u>.

Directions:

- This assignment will be completed in your <u>groups</u>. However, **every member of the group must be fully aware of every part of the assignment**. Further, while you can discuss concepts with other groups, actual steps and answers should not be shared between groups.
- The assignment will ask for short answers or screenshots; this material should be collected in a PDF and submitted via GradeScope. Anything you need to include in this document is highlighted in cyan.
- At times, you will be asked to speculate as to *why* something is happening. **If you aren't confident in your answers, see the instructor to discuss it!** Talking through things you're unsure of is part of the class.

## 1 Check hard drives for existing data

In Lab 0, you created dummy virtual devices for each physical device in the hardware RAID controller, allowing the OS to see the SSD and each HDD.

The SSDs you have are either new or lightly used last year, but the HDDs are used stock purchased on ebay (and many were used last year, too). We're going to check what kind of stuff is on these disks, then we'll overwrite them (and as a side-effect, this will also verify they're good).

Note: In lab 0, some groups had to manually zero their drives during install. Proceed with these steps regardless, as they'll still be a good lesson for you and a good read test for your drives.

## **1.1 Examine block devices**

#### Using the lsblk command, take a look at your drives. Screenshot.

In Linux, drives are represented by /dev/sda, /dev/sdb, etc. You should see your SSD (distinguished by its size of around 120GB or 240GB) plus your HDDs (each around 73 or 146GB depending on your model). Partitions that have been made are indicated with numbers, so /dev/sda1 is the first partition of the first drive. Based on where you installed your SSD (in the first drive slot) and how you partitioned it, I'd expect /dev/sda1 to be the partition that contains your root directory (though it's okay if it's not).

#### Let's see what's mounted: simply run 'mount' without arguments. Screenshot.

You should see something like "/dev/sda1 on / type ext4 (rw)", which indicates that your root filesystem is on /dev/sda1. Make note of which drive your OS is on (/dev/sda in this example), as we won't want to overwrite that. Note your OS drive in the write-up.

All other /dev/sdX drives will be overwritten.

### 1.2 Check the content of the non-OS drives

Let's glance at one of these drives. We'll get into data forensics later, but you're probably curious if there's anything here. Our main question should be "were these drives *zeroed* before being given to us?". To "zero" means to overwrite all bytes with zeroes, and is the standard method of cleaning a drive. We can quickly check by running "hd" (hex dump) on the drive. If non-zero bytes come up, there's probably something there to find. If it just prints a row of zeroes and seems to wait, it's indicating that all subsequent bytes so far are also zeroes, and we likely have a blank drive.

#### Quickly check each drive with a command like: sudo hd /dev/sda

Example of a zeroed drive:

(It sits there reading, so I get bored and press Ctrl+C after a few seconds)

#### Example of a drive with stuff on it:

Note which drives contain data and which appear to be zeroed in your write-up.

### 1.3 Zero the non-OS drives

Let's zero out every drive other than the OS drive.

This turns out to be quite easy. As we've seen in class, the kernel provides a special virtual file that contains infinite zeroes: /dev/zero. We just need to write that content onto our drives. The easiest way is with dd (introduced in class, manpage here, tutorial here).

Using dd, with a block size of 1MB, copy /dev/zero to one of your non-OS drives. In your write-up, show the command used.

# WARNING: Be sure you know which drives are your hard disk drives and which is your SSD with the operating system. If you zero out the OS drive, your OS will stop working and you'll have to reinstall it.

Once you get this working, you can go ahead and zero *all* the drives in parallel (even ones that appeared to be zeroed already). This process will take a while, so if getting disconnected from SSH is an issue, consider learning to use <u>screen</u> or <u>tmux</u> to allow your process to run in the background.

After zeroing is complete, check for errors encountered:

- See if each dd command gave any error output. Even if you see no error output, issues may have been seen by the kernel but not reported to dd, so check the kernel log. Run "dmesg" to view this log, and look for any entries that appear to deal with hardware-level errors on drives (see <u>here</u> for an example of kernel log messages relating to drive hardware errors). If you see any, document the errors in your write-up, including which drive(s) it was, messages received, etc.
- If errors are found, contact the instructor to swap drives with a spare; document the replacement in your write-up.

## 2 HDD vs. SSD performance

## 2.1 Benchmarking

Let's do a very simple performance test on our drives. We want to read each device and determine the IOs Per Second (IOPS) for random IO and throughput (MB/s) for sequential IO. The IO in this case will be block-oriented: we'll be reading a chunk of data at a time (e.g. 512 or 4096 bytes). This is because IO operations have a fairly high setup cost (including kernel interrupt and scheduling of the actual IO traffic), so we want to spread that cost over a significant amount of data.

The benchmark software we'll be using is <u>the simple "iops" tool written by Benjamin Schweizer</u>. Download this tool onto your server. It is just a Python script, so just enable it for execution (chmod +x) and see the usage message by running it without arguments.

We're going to run it directly on our bare drives (bypassing the filesystem layer), so you'll need to run it as root. It only does read operations, so the test can be done safely even on your OS drive.

Run it in both sequential and random mode for each drive you have (/dev/sd?), and run it with block size 4kB and 1MB. For sequential IO, keep the number of parallel threads at 1. For random IO, use 8

threads. Using parallelism on your random IO test is important, as that's *how* RAID improves random IO performance: it can't make a single seek go faster, but it can allow multiple seeks to occur in parallel – but only if there are multiple IO requests to schedule! Interested readers can use scripting to automate these tests.

In your writeup, include your results as a table, showing the test inputs, the IOPS, and the throughput (MB/s). As always in such things, do multiple repeated runs, show the average, and note the variance (e.g. as standard deviation)!

Based on these numbers, how long would it take to read your first HDD entirely in sequential order? How long if you did it in random order? How does this make you feel?

Do the same computation for your SSD. Wow, SSDs sure are neat, huh?

Observe that sequential IO is still faster than random IO on the SSD, despite it having no mechanical parts. Why is this?

How did increasing the block size change things? Why might that be?

## 2.2 Comparing to the specs

Look up the model of HDDs and SSDs you're using. For the HDDs, look up any 2.5" SAS drive that is the same RPM and capacity as yours. For the SSD, you have a Kingston 120GB A400 SATA-3, model SA400S37/120G, or a Patriot Memory Burst Elite SATA 3 240GB, model PBE240GS25SSDR. If you can't find that, look for specs for a comparable drive (e.g. another 120GB or 240GB SATA-3 SSD).

In your write-up, share links to the specs you find, and summarize the performance stats (random IOPS and sequential throughput for each drive type).

Note: it may be hard to find IOPS numbers from the datasheet, as manufacturers would usually rather talk about seek times and data rates. In this case, you may take the reciprocal of seek time and note this in your writeup. Also, some datasheets may refer to sequential throughput as "sustained data rate".

Compare the rated random IOPS and sequential throughput to your measurements from the previous question. I suspect the rated ones will be higher (as in, you measured less performance than the manufacturer says you can expect). In this case, how might the manufacturer test differ from ours? In the unlikely case that you actually measured higher performance than the rated specs, why might this be?

## **3 Software RAID**

RAID can be done either in software by the operating system or in hardware by a dedicated RAID card. In Lab 0, we configured the hardware RAID card to simply pass each drive through directly to the OS. This will allow us to apply software RAID. In Linux, software RAID is called handled by the "md" subsystem. Ensure the "mdadm" tool is installed.

NOTE: I wrote some of the documentation below without access to a machine with real drives I could trash, so I faked some of the sample output using simulated drives; apologies if my forgeries are imperfect.

## 3.1 Set up software RAID

Noting this <u>RAID setup guide</u>, create a three-disk RAID5 called /dev/md0 using your HDDs, with the fourth one serving as a spare device. Note: you can skip partitioning and instead use the full drives unpartitioned. Document the commands you used to set it up in the write-up.

You can always check the current status of all md RAID devices by reading /proc/mdstat. When I built mine, I checked this file immediately after to see it "repairing" the new array (calculating and storing parity). This may take some time, but the array is usable while it occurs, albeit at a lower performance level. Therefore, **wait for it to finish before proceeding with benchmarking**. Below, you can see mdstat showing the rebuild underway, and then later showing the array as ready:

Note the "(S)" flag on sdb: this indicates that this device is a hot spare.

Include similar output of /proc/mdstat either during rebuild or afterward in your write-up.

## 3.2 Benchmark

Based on the theory behind RAID, compute the logical size your RAID array should be, show your work and note the result in your write-up. Then, using Isblk, confirm the size of the RAID device /dev/md0 on the system, showing the output in your write-up. Using the exact same procedure as in section 2.1, benchmark the RAID drive /dev/md0.

In your write-up, include your results as a table, showing the test settings, the IOPS, and the throughput (MB/s). Again, be sure to do multiple repetitions and show the average and variance.

In your write-up, compare all the performance numbers to the average performance of a single one of the drives. Does the performance improvement match the theory?

## 3.3 Set up a filesystem

Our RAID array is represented by /dev/md0. Now we need a file system on top of it. There are many file systems available, and even a basic install of Linux supports many of them. We'll learn more about the design of file systems later in the course, but for now, we just need to create one so that a directory structure can be accessed on our RAID array for testing. To create the initial filesystem, we use the <u>mkfs</u> ("make file system") utility, which has variants for each supported filesystem (e.g. ext2, ext3, ntfs, etc.). Let's just use the common Linux choice of ext4. **Use "mkfs.ext4" to prepare a filesystem on /dev/md0;** provide a screenshot of this.

Now we need a *mount point*, a directory into which we can attach the filesystem. Using <u>mkdir</u>, create a directory called "/x".

Now use the <u>mount</u> utility to mount our /dev/md0 filesystem to /x. Note that while the mount utility does support a lot of options, none are needed for this basic operation – the tool will have no trouble auto-detecting the type of filesystem present and using reasonable default settings when mounting it. When done, you can verify by typing "mount" by itself to list mounted filesystems; you can also type "df" to check free space. Example:

```
root@xub1404dt:~ # mount
/dev/sda1 on / type ext4 (rw,errors=remount-ro)
proc on /proc type proc (rw,noexec,nosuid,nodev)
sysfs on /sys type sysfs (rw,noexec,nosuid,nodev)
...
none on /sys/fs/pstore type pstore (rw)
systemd on /sys/fs/cgroup/systemd type cgroup
(rw,noexec,nosuid,nodev,none,name=systemd)
/dev/md0 on /x type ext4 (rw)
root@xub1404dt:~ # df
```

IOOCEAUDITOTUC. * # CI					
Filesystem	1K-blocks	Used	Available	Use%	Mounted on
/dev/sda1	19478204	4191128	14274596	23%	/
none	4	0	4	0 응	/sys/fs/cgroup
udev	489424	4	489420	18	/dev
tmpfs	100020	1456	98564	28	/run
none	5120	0	5120	0 응	/run/lock
none	500096	152	499944	1%	/run/shm
none	102400	40	102360	1%	/run/user
/dev/md0	203772	2504	201268	2%	/ x

Provide screenshots or a terminal log similar to the above in your write-up.

### 3.4 Use the file system

By whatever means you wish, populate the filesystem at least a few kilobytes of content. At least one file should be a text file you can verify the contents of. When done, you should be able to cd to your filesystem and see your stuff, e.g.:

```
root@xub1404dt:/x # ls -l
total 781
-rwxrwxrwx 1 root root 797772 Sep 30 17:14 kern.log
-rwxrwxrwx 1 root root 21 Sep 30 17:15 test.txt
root@xub1404dt:/x # cat test.txt
This is my text file
```

Provide screenshots or a terminal log of these steps in your write-up. Be sure to include a file listing and "cat" of your small text file, as shown above.

## 3.5 Damage the RAID

We'll now simulate failure of one of the disks. As soon as you do, the array will begin repairing itself using the hot spare we provided.

Use the mdadm command to "fail" one of the active drives in the array (not the spare!), then print the content of /proc/mdstat so you see the array status while it's rebuilding. Check the /proc/mdstat file again after it's done to see the final status. You'll find the failed disk with a "(F)" flag, and the former spare disk is now in the array, so it no longer has the "(S)" flag. Example output of mdstat before, during, and after rebuilding:

```
root@xub1404dt:~ # cat /proc/mdstat
Personalities : [raid6] [raid5] [raid4]
md0 : active raid5 sdd[3](S) sdb[5] sdc[4] sda[0]
      203776 blocks super 1.2 level 5, 512k chunk, algorithm 2 [3/3]
[UUU]
(failure is simulated)
root@xub1404dt:~ # cat /proc/mdstat
Personalities : [raid6] [raid5] [raid4]
md0 : active raid5 sdd[3] sdb[5](F) sdc[4] sda[0]
      203776 blocks super 1.2 level 5, 512k chunk, algorithm 2 [3/2]
[U U]
      [======>....] recovery = 50.0% (51580/101888)
finish=0.0min speed=51580K/sec
unused devices: <none>
root@xub1404dt:~ # cat /proc/mdstat
Personalities : [raid6] [raid5] [raid4]
md0 : active raid5 sdd[3] sdb[5](F) sdc[4] sda[0]
      203776 blocks super 1.2 level 5, 512k chunk, algorithm 2 \left[3/3\right]
[UUU]
unused devices: <none>
```

During reconstruction, verify that the content in the mounted filesystem /x is undamaged, despite a disk failure, e.g.:

```
root@xub1404dt:/x # ls -l
total 781
-rwxrwxrwx 1 root root 797772 Sep 30 17:14 kern.log
-rwxrwxrwx 1 root root 21 Sep 30 17:15 test.txt
root@xub1404dt:/x # cat test.txt
This is my text file
```

Provide screenshots or a terminal log of these steps in your write-up, including evidence of your text file surviving the failure.

## 3.6 Replace the faulty drive

At this point, because the spare was there to allow rapid reconstruction, your array can again withstand a drive failure without data loss, but there would be no more spares to rebuild it immediately in that case. Running with no spares is dangerous, so we'll now simulate replacing the faulty drive. Steps:

- 1. Use the mdadm "remove" option to delete the faulty drive /dev/sdb from the array listing. This is the software equivalent of pulling the drive out of the server.
- 2. At this point you'd replace the drive with a new one, but since our drive isn't actually bad, just imagine doing so in your mind.
- 3. Use the mdadm "add" option to add our 'replacement' drive (the same drive) to our /dev/md0 array.

When you're done, replaced drive should show as a new spare device in /proc/mdstat:

Your array is now capable of immediately rebuilding if another drive failure occurs.

Provide screenshots or a terminal log of these steps in your write-up.

### 3.7 Tear down software RAID

Unmount /c, then use mdadm to stop and destroy your RAID device.

## 4 Hardware RAID

Let's compare the setup and performance of software RAID to hardware RAID.

## 4.1 Set up hardware RAID

Either in-person or via remote access, reboot the server, and press the appropriate key to get into the RAID controller setup interface (like you did in Lab 0). Destroy the virtual devices for the hard disks, *taking extreme care not to harm the virtual device for the SSD*, which is where your OS lives.

Then create a new virtual device that is a three-disk RAID5 with a single hot spare (the same as configuration as our previous software RAID).

Take a screenshot of the firmware interface showing the RAID5 setup for your write-up.

## 4.2 Benchmark

Your new hardware RAID device will show up as a new /dev/sd? device. Figure out which device it is and note the drive size.

Wait for the array rebuild to complete (I believe you can check this in the boot-time menu).

Using the exact same procedure as in section 2.1, benchmark the RAID drive.

In your write-up, include your results as a table, showing the test settings, the IOPS, and the throughput (MB/s). As always, do multiple reps and characterize variance.

In your write-up, compare all the performance numbers to the software RAID. Is one strictly better than the other? Are there tradeoffs? What's the difference?

## 4.3 Tear down hardware RAID and restore single-drive access

Reboot, go back into the RAID controller firmware, destroy the RAID5 virtual device, and in a procedure similar to Lab 0, make pass-through virtual devices for each HDD, basically restoring the system to its previous configuration from before you started this assignment. Again, <u>take extreme care not to harm</u> <u>the virtual device for the SSD</u>, which is where your OS lives.