

DataEng: Data Maintenance In-class Assignment

This week you will gain hands-on experience with Data Maintenance by constructing an automated data archiver that compresses, encrypts and stores pipelined data into a low-cost, high-capacity GCP Storage Bucket.

Submit: Make a copy of this document and use it to record your results. Store a PDF copy of the document in your git repository along with any needed code before submitting for this week.

Your job is to develop a new, separate python Kafka consumer similar to the consumers that you have created multiple times for this class. This new consumer should be called `archive.py` and should receive all data from a test Kafka topic, compress the data, encrypt the data (optional) and store the compressed data in a [GCP Storage Bucket](#).

Note that each member of your team should build their own archive mechanism. As always, it is OK to help one another, but each person should develop their own python program for archiving.

Discussion Question for Your Entire Group (do this first)

When archiving data for a data pipeline we could (a) compress, (b) encrypt and/or (c) reduce the data. Here, “reducing the data” refers to the process of transforming detailed data, such as 5 second breadcrumbs for all buses on all trips, into coarser data that contains, for example, only contains a subset of the original data such as only some buses, some trips or possibly fewer breadcrumbs per trip.

Under what circumstances might each of these transformations (compress, encrypt, reduce) be desirable for a data archival feature?

For some log files, a properly compressed archive can be done. This compression feature, data and index compression, and backup compression, reduces the amount of disk space required for database environments.

Part 1: Data Archival to Cold Storage

A. Create test topic

Create new Kafka producer and consumer programs as you did with the Data Transport in-class activity ([link to Transport activity](#)). Create a new Kafka topic that is separate from the topic(s) used for your project. Call it “archivetest” or something similar. As with the Data Transport activity you should initially have your new producer produce test data and have a single consumer that consumes any/all data sent to the Kafka topic.

B. Create separate consumer groups

Similar to part H in the Transport activity, create two separate Consumer Groups for your new Kafka topic. Run a separate consumer for each group and verify that each consumer consumes all of the data sent by the producer.

Your first consumer should simply print all data that it receives. This consumer simulates the main branch of your data pipeline. Typically, the main branch of your pipeline would validate, transform, integrate, load, etc. the data, but for this in-class assignment it only needs to print the data.

C. Archive the data in a GCP Storage Bucket

Your second consumer (call it archive.py) should store all received data into a [GCP Storage Bucket](#). You will need to create and configure a Storage Bucket for this purpose. You are free to choose any of the available storage classes. We recommend using the Nearline Storage class.

D. Compress

Modify your archive.py program to compress the data before it stores the data to the storage bucket. Use [zlib compression](#) which is provided by default by python. How large is the archived data compared to the original?

E. Encrypt (Optional)

Next, modify your archive.py to encrypt the data prior to writing it to the Storage Bucket. Your archive.py program should encrypt after compressing the data. Use RSA encryption as described here: [link](#) There is no need to manage your private encryption keys securely for this assignment, and you may keep your private key in a file or within your python code.

Be sure to test your archiver by decrypting and decompressing the data stored in the Storage Bucket. We suggest that you create a separate python program for this purpose.

How large is the archived data?

F. Add Archiving to your class project (Optional)

Add your archive.py program (or something like it) to your class project's pipeline(s). Mention this in your final project presentation video. Because your project is shared among your project team members you will need to coordinate the adding of new Kafka consumer groups so that each team member may safely add their own archiving service. Again, it is not necessary to securely manage your RSA private encryption key, and it is OK to keep it in a file or in your python source code.

Part 2: VM Monitoring and Resiliency

G. Virtual Machine Resource Usage Monitoring

This section covers how to monitor your virtual machine's health. Walk through this activity together with your project group since you have the same project VM. In particular we'll be looking into how to check the resource usage on your host. As this is a virtual machine, you have the choice of using traditional Linux command line tools like top, or using GCP UI tools.

For the questions about the pipeline, you can either use the pipeline from the activity above, or modify your course project pipeline. Be sure to set `autocommit = False` when setting up the database connection so it doesn't actually commit the data to your tables and mess up your data. Use the new endpoint: <http://www.psudataeng.com:8000/getBreadCrumbDataV2> This endpoint just serves the data from the start of the course again. It'll be shut down at the end of day Friday.

Some handy links are listed below. We recommend skimming or searching them as a reference during the rest of this activity. You are not required to read them. Remember to Google anything that isn't immediately clear!

- Linux command line monitoring tools: <https://www.makeuseof.com/best-cli-tools-to-monitor-linux-performance-terminal/>
 - My personal favorites: dstat, top, free -h, df -h, du -h -d1, ps aux
- GCP Monitoring overview: <https://cloud.google.com/monitoring/docs/monitoring-overview>
- GCP VM monitoring dashboards: https://console.cloud.google.com/monitoring/dashboards/resourceList/gce_instance?tab=overview
- GCP Monitoring Processes: <https://cloud.google.com/monitoring/agent/process-metrics>

- GCP Alerts doc: <https://cloud.google.com/monitoring/monitor-compute-engine-virtual-machine>
- GCP Alerts UI: https://console.cloud.google.com/monitoring/alerting?walkthrough_id=monitoring_quicks_tart_compute_uptime

Do you have any other monitoring tools you prefer?

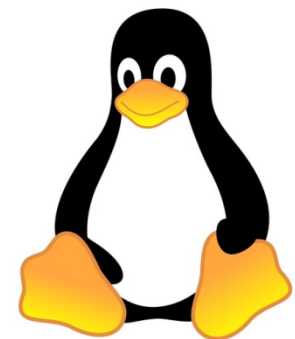
No. I prefer Linux command line monitoring tools because I am more familiar with Linux.

H. Getting to know your VM

Which [Linux distribution](#) are you using? (Distro and year: Ubuntu 20)

X86_64 && Linux

Linux version 5.13.0-1024-gcp (buildd@lcy02-amd64-110) (gcc (Ubuntu 9.4.0-1ubuntu1~20.04.1) 9.4.0, GNU ld (GNU Binutils for Ubuntu) 2.34)



What are your system's resource allocations? How much of each of the following resources do you have on the VM?

- CPUs
- CPU cores
- Memory
- Swap space
- Disk

Are there other resources you're used to monitoring? When are those resources important to keep an eye on?

"htop"

```

1  [ 0.0%] Tasks: 36, 44 srr: 1 running
2  [ 0.0%] Load average: 0.00 0.00 0.00
Mem [|||||] 204M/3.93G Uptime: 30 days, 23:39:34
Cpu [|||||] 0%/0%

#PID#USER#CPU#MEM#RSS#VSZ#ST#TIME#COMMAND#
194200 postgres 20 0 212M 24060 2028 S 0.0 0.6 0:31.01 /usr/lib/postgresql/12/bin/postgres -D /var/lib
1 root 20 0 164M 11012 6456 S 0.0 0.3 2:10.73 /sbin/init
170 root 19 - 76912 21888 2326 S 0.0 0.6 1:04.19 /lib/systemd/systemd-journald
201 root 20 0 8860 4476 3864 S 0.0 0.1 0:06.71 /lib/systemd/systemd-udev
324 root RT 0 273M 18000 8208 S 0.0 0.4 0:27.19 /sbin/multipathd -d -s
325 root RT 0 273M 18000 8208 S 0.0 0.4 0:00.00 /sbin/multipathd -d -s
326 root RT 0 273M 18000 8208 S 0.0 0.4 0:04.29 /sbin/multipathd -d -s
327 root RT 0 273M 18000 8208 S 0.0 0.4 3:18.92 /sbin/multipathd -d -s
328 root RT 0 273M 18000 8208 S 0.0 0.4 0:00.01 /sbin/multipathd -d -s
329 root RT 0 273M 18000 8208 S 0.0 0.4 0:00.00 /sbin/multipathd -d -s
322 root RT 0 273M 18000 8208 S 0.0 0.4 6:36.63 /sbin/multipathd -d -s
446 systemd- 20 0 237M 4812 3872 S 0.0 0.1 0:21.59 /lib/systemd/systemd-networkd
450 systemd- 20 0 238M 9100 5176 S 0.0 0.2 0:20.04 /lib/systemd/systemd-resolved
565 root 20 0 235M 6380 5360 S 0.0 0.2 1:10.57 /usr/lib/accounts-service/accounts-daemon
582 root 20 0 235M 6380 5360 S 0.0 0.2 0:00.02 /usr/lib/accounts-service/accounts-daemon
564 root 20 0 235M 6380 5360 S 0.0 0.2 1:10.70 /usr/lib/accounts-service/accounts-daemon
581 messagebu 20 0 766M 3776 3048 S 0.0 0.1 0:15.66 /usr/bin/dbus-daemon --system --address=system
595 root 20 0 768M 12540 616 S 0.0 0.3 1:02.23 /usr/bin/google_osconfig_agent
596 root 20 0 768M 12540 616 S 0.0 0.3 0:55.35 /usr/bin/google_osconfig_agent
597 root 20 0 768M 12540 616 S 0.0 0.3 0:00.00 /usr/bin/google_osconfig_agent
599 root 20 0 768M 12540 616 S 0.0 0.3 0:00.00 /usr/bin/google_osconfig_agent
605 root 20 0 768M 12540 616 S 0.0 0.3 0:53.56 /usr/bin/google_osconfig_agent
607 root 20 0 768M 12540 616 S 0.0 0.3 0:53.71 /usr/bin/google_osconfig_agent
608 root 20 0 768M 12540 616 S 0.0 0.3 0:51.47 /usr/bin/google_osconfig_agent
610 root 20 0 768M 12540 616 S 0.0 0.3 0:55.04 /usr/bin/google_osconfig_agent
613 root 20 0 768M 12540 616 S 0.0 0.3 0:52.20 /usr/bin/google_osconfig_agent
587 root 20 0 768M 12540 616 S 0.0 0.3 6:24.54 /usr/bin/google_osconfig_agent
612 daemon 20 0 3800 2188 2008 S 0.0 0.1 0:00.19 /usr/sbin/atd -F
762 root 20 0 834M 13008 6664 S 0.0 0.3 0:42.85 /usr/bin/google_guest_agent
763 root 20 0 834M 13008 6664 S 0.0 0.3 0:00.00 /usr/bin/google_guest_agent

```

I. Pipeline Resource usage

When your pipeline is running, what is the typical system-wide usage of these resources:

- CPU
- Load Average
- Memory
- Disk Usage (fullness)
- Disk Utilization (I/O)
- Network (I/O)

When your pipeline is running, what are some of the top processes that are running? Which ones are consuming the most CPU and memory?

Python files are running. Python files are the most CPU and memory consuming.

J. Understanding Linux Monitoring Metrics

What is the difference between disk usage and disk utilization?

Utilization is how often you use and the resources you use. Utilization refers to the space you have at your disposal and use. Disk usage (DU) refers to the portion or percentage of computer storage that is currently in use. It contrasts with disk space or capacity, which is the total amount of space a given disk can store. Disk usage is typically measured in kilobytes (KB), megabytes (MB), gigabytes (GB), and/or terabytes (TB). Disk usage is a key metric for any computing system because it provides users with the information they need not only for storage, but also for software requirements and overall operation. While it usually refers to a computer's hard drive, it may also refer to external storage devices. You can use utilization to estimate the processing time on disk for a transaction that requires a given number of disk transfers.

What is the difference between memory that is free and memory that is available, in the command `free -h`? Is it concerning when there is very little “free” memory, if there is a lot of “available” memory?

```
gen@dataeng:~$ free -h
```

	total	used	free	shared	buff/cache	available
Mem:	968Mi	189Mi	79Mi	12Mi	699Mi	624Mi
Swap:	0B	0B	0B			

Does swap space use memory or disk? Is it a good or bad sign (or both/neither) when it is used? What is a standard amount of swap space to allocate?

Should you be concerned if the load average is 7?

Yes, the previous question is a trick question. The answer to your followup question: the system has 32 cores. Are you concerned now?

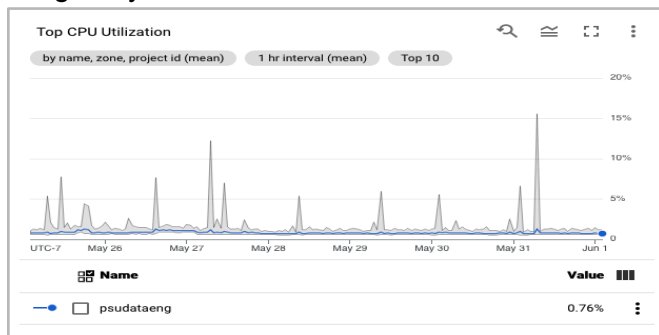
Should you be concerned if the 1 minute load average is really high, but the 5 and 15 minute load averages are normal? Why or why not? What does this mean, is load going up or down?

The average load refers to the average number of processes that the system is in a running state and an uninterruptible state per unit time, that is, the average number of active processes. When judging load averages, you first need to know how many CPUs your system has.

- If the three values of 1 minute, 5 minutes, and 15 minutes are basically the same, or the difference is not large, it means that the system load is very stable.
- However, if the value of 1 minute is much smaller than the value of 15 minutes, it means that the load of the system has decreased in the last minute, but there has been a heavy load in the past 15 minutes.
- Conversely, if the value of 1 minute is much larger than the value of 15 minutes, it means that the load of the last minute has increased. This increase may only be temporary, or it may continue to increase, so continuous observation is required. Once the 1-minute average load is close to or exceeds the number of CPUs, it means that the system is experiencing overload problems. At this time, it is necessary to analyze and investigate what is causing the problem and find ways to optimize it.

K. GCP Monitoring

Where in the GCP cloud console can you see charts like the one below, showing the resource usage of your VM?



Where in the GCP cloud console can you see disk usage (not utilization)?

Click on the instance name in the virtual machine instance and select Observability, then we can choose CPU or DISK, etc.

How can you see GCP billing usage, excluding credits from promotions and course coupons?

Instances

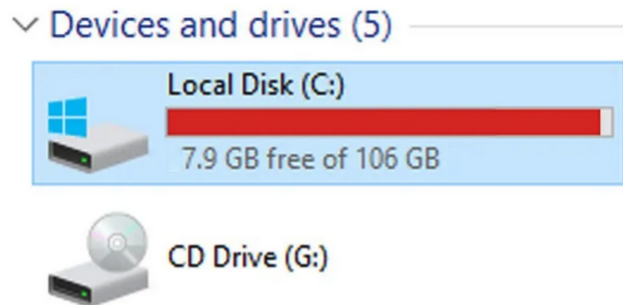
Select billing report in VM instance.

L. Predicting and Preventing Outages

Consider the following questions, and how you can improve the resiliency of your VM and pipeline.

Running out of resources

Judging from your current usage patterns, and if nothing were to change, on what date do you expect to run out of disk space? What about on your file archive GCP storage bucket? When will you run out of GCP credits?



A picture of a full disk on Windows. Linux just doesn't run out of disk space with quite the same flamboyance.

When your pipeline is downloading new data, is the system running out of any resources? Did it run out of memory? Does it start paging/swapping? Was the CPU pegged? What was the normalized CPU usage (total CPU usage, averaged over the count of CPU cores)?

Judging from the last question, which resource would you say is the biggest bottleneck?



Alerting on and resolving problems

Given the typical resource usage patterns we've explored here, if you were to create some resource usage alerts, what types of thresholds would you set? For example, would you get out of bed to fix a disk at 70% full or wait till 90% full? Thinking of typical human behavior, what issues can you think of if alert thresholds are set too high or too low?

If your system were running out of disk and you needed to free up space right away, what files would you delete first? Could you recreate the data from the remaining files? Do you have a safe copy of the data elsewhere?

If your system were running out of memory, and you knew it was caused by your consumer, how can you kill the consumer process? What effect (if any) would this have on the integrity of data in your database? Could you safely start the consumer again once the system is healthy again, or would you need to fix any data first? Can you think of any database features that could make it more resilient to failures in the middle of running the pipeline?

Ctrl-c to close the current process. It does have an impact. Therefore, if something goes wrong I will manually delete the wow city part, or the completed part of the killed process. Using Explore Actifio Go to get it is a good way to do it, but I don't use it in my project and really ignore the easier way to delete or restore data.