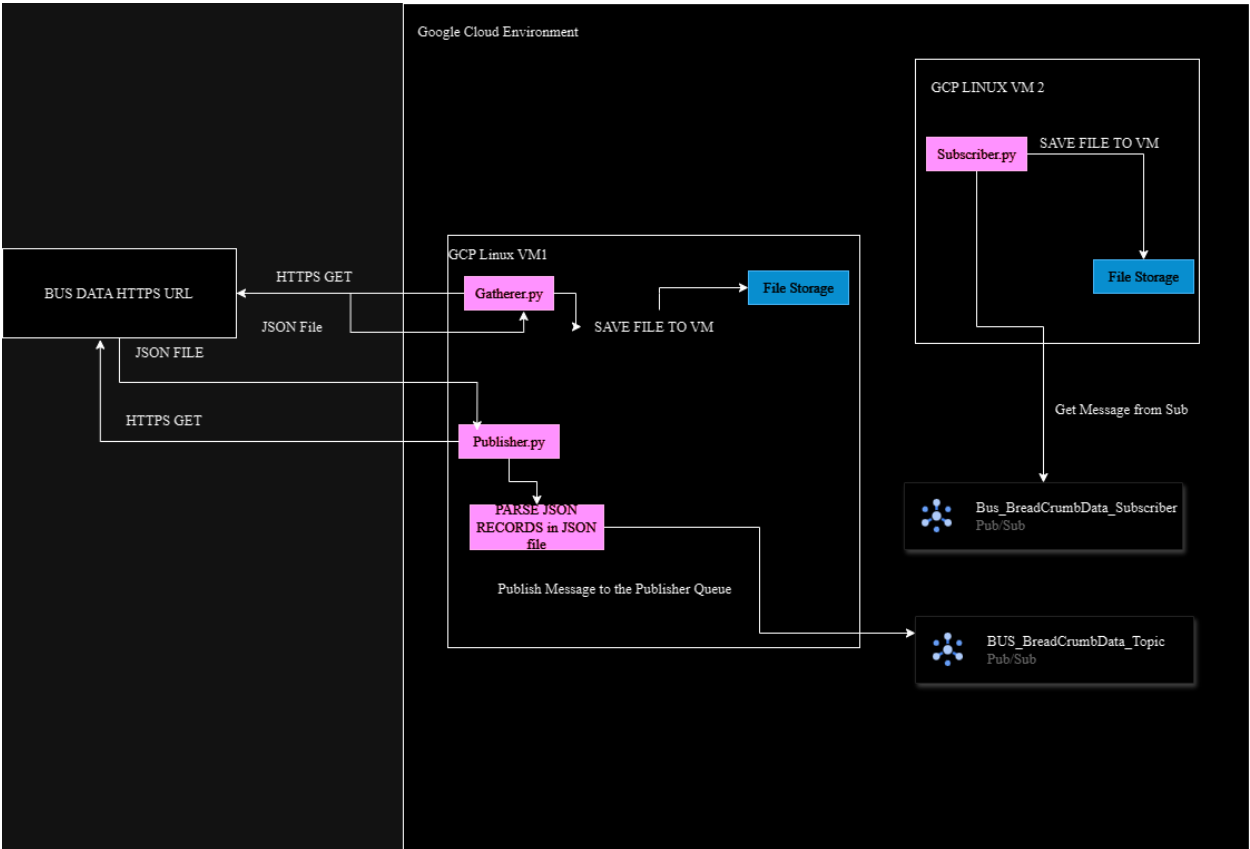# DataEng S24: Project Assignment 1

Gather and Transport

Current Planned Architecture Diagram:



# A. Create, Configure and Run Your Virtual Machine

**VMs created for our project in GCP is as follows :**

| Name ↑ | Zone | Recommendations | In use by | Internal IP | External IP | Connect | | |
|---|---|---|---|---|---|---|---|---|
| instance-20240415-190244 | us-west1-b | | | 10.138.0.5 (nic0) | | SSH | ▾ | ⋮ |
| instance-20240421-224350 | us-west1-b | | | 10.138.0.6 (nic0) | | SSH | ▾ | ⋮ |

# B. Initial Python Data Gatherer

Refer to the file **DataEnggProjectExtractionInitial.py** in our Project Repo

# C. Run the Data Gatherer Daily

```
srirams@instance-20240415-190244:~$ crontab -l
# Edit this file to introduce tasks to be run by cron.
#
# Each task to run has to be defined through a single line
# indicating with different fields when the task will be run
# and what command to run for the task
#
# To define the time you can provide concrete values for
# minute (m), hour (h), day of month (dom), month (mon),
# and day of week (dow) or use '*' in these fields (for 'any').
#
# Notice that tasks will be started based on the cron's system
# daemon's notion of time and timezones.
#
# Output of the crontab jobs (including errors) is sent through
# email to the user the crontab file belongs to (unless redirected).
#
# For example, you can run a backup of all your user accounts
# at 5 a.m every week with:
# 0 5 * * 1 tar -zcf /var/backups/home.tgz /home/
#
# For more information see the manual pages of crontab(5) and cron(8)
#
# m h  dom mon dow   command
30 12 * * * python3 DataEnggProjectExtraction.py
30 12 * * * python3 DataEnggProjectPublisher.py
srirams@instance-20240415-190244:~$
```

Extra Credit:

Updated my Code to make it check for Storage space and for IDs where download fails due to any issue, one Notification email will be sent with the list of vehicle IDs.

Vehicle Data Download Failures  External  Inbox ×

metrometricsmavericks@gmail.com                                    5:41PM (0 minutes ago)
to me ▾

The following vehicle IDs failed to download data: 4027, 3507, 2904, 3018, 2935, 4526, 4520, 4502, 3258, 2924, 3719, 4505, 4210, 3325, 3054, 3802, 3035, 3530, 4007, 3630, 3503, 4236, 4518, 3805, 4207, 3524, 3023

Refer to file **DataEnggProjectExtraction.py** in our project Repo.

# D. Configure Google Cloud Pub/Sub

We have Configured Pub/Sub for our project as follows:

| | Topic ID ↑ | Encryption key | Topic name | Retention | Ingestion source |
|---|---|---|---|---|---|
| ☐ | busBreadCrumbData | Google-managed | projects/focus-surfer-420318/topics/busBreadCrumbData | — | — |

| | State | Subscription ID ↑ | Delivery type | Topic name | Ack deadline | Retention | Message ordering | Exactly once delivery | Expiration |
|---|---|---|---|---|---|---|---|---|---|
| ☐ | ✅ | busBreadCrumbData-sub | Pull | projects/focus-surfer-420318/topics/busBreadCrumbData | 10 seconds | 7 days | Disabled | Disabled | 31 days |

# E. Parse JSON into Individual Breadcrumb Records

# F. Send Breadcrumb Records to Pub/Sub Topic

Refer to the file **DataEnggProjectPublisher.py** in our Project Repo

# G. Pub/Sub Receiver

Refer to the file **DataEnggProjectSubscriber.py** in our Project Repo

# H. Configure Linux to Run Pub/Sub Receiver Continually

**Service is enabled for Receiver and we have configured the service to restart on failures and this is the status screenshot.**

```
Apr 21 20:20:27 instance-20240421-224350.c.focus-surfer-420318.internal systemd[1]: Started reciever.service - Pub/Sub Subscriber Service.
root@instance-20240421-224350:/home/srirams# sudo systemctl status reciever
● reciever.service - Pub/Sub Subscriber Service
     Loaded: loaded (/etc/systemd/system/reciever.service; enabled; preset: enabled)
     Active: active (running) since Sun 2024-04-21 20:20:27 PDT; 34s ago
   Main PID: 7015 (python3)
      Tasks: 16 (limit: 4686)
     Memory: 38.8M
        CPU: 593ms
     CGroup: /system.slice/reciever.service
             └─7015 /usr/bin/python3 /home/srirams/DataEnggProjectSubscriber.py

Apr 21 20:20:27 instance-20240421-224350.c.focus-surfer-420318.internal systemd[1]: Started reciever.service - Pub/Sub Subscriber Service.
root@instance-20240421-224350:/home/srirams#
```

# I. Schedule your VM to start and stop automatically

**We have setup the Instance for our project as follows**

Instances schedules

Filter Enter property name or value

| | Name ↑ | Region | Start schedule | Stop schedule | Time zone | Initiation date | Expiration date |
|---|---|---|---|---|---|---|---|
| ☐ | myinstancefordataengg | us-west1 | 12:08PM, every day | 12:45PM, every day | America/Los_Angeles | | |

**DataEng Project Assignment 1 Submission Document**

Construct a table showing each day for which your pipeline successfully, automatically processed one complete day's worth of sensor readings.

| Date | Day of Week | Approximate Time of day for your data access | # Sensor Readings | Total Data Saved (KBs) | # Pub/Sub messages published and received |
|---|---|---|---|---|---|
| 04/14/2024 | Sunday | 12:30 PM PDT | { <br><br> "EVENT_NO_TRIP":220385292, <br><br> "EVENT_NO_STOP":220385302, <br><br> "OPD_DATE":"15DEC2022:00:00:00", <br>   "VEHICLE_ID":2902, <br>   "METERS":19021, <br>   "ACT_TIME":23580, <br>   "GPS_LONGITUDE":-122.69081, <br>   "GPS_LATITUDE":45.535613, <br>   "GPS_SATELLITES":8.0, <br>   "GPS_HDOP":1.1 <br>} | 7898 5328 | 311034 |
| 04/15/2024 | Monday | 12:30 PM PDT | { | 7847 2540 | 325960 |

| | | | | | |
|---|---|---|---|---|---|
| | | | {<br><br>"EVENT_NO_TRIP":221431889,<br><br>"EVENT_NO_STOP":221431906,<br><br>"OPD_DATE":"16DEC2022:00:00:00",<br>    "VEHICLE_ID":2902,<br>    "METERS":8487,<br>    "ACT_TIME":47927,<br>    "GPS_LONGITUDE":-122.802103,<br>    "GPS_LATITUDE":45.463523,<br>    "GPS_SATELLITES":12.0,<br>    "GPS_HDOP":0.8<br>  } | | |
| 04/16/2024 | Tuesday | 12:30 PM PDT | {<br><br>"EVENT_NO_TRIP":221510274,<br><br>"EVENT_NO_STOP":221510352,<br><br>"OPD_DATE":"17DEC2022:00:00:00",<br>    "VEHICLE_ID":3023,<br>    "METERS":27683,<br>    "ACT_TIME":66544,<br>    "GPS_LONGITUDE":-122.58922,<br>    "GPS_LATITUDE":45.523058,<br>    "GPS_SATELLITES":12.0,<br>    "GPS_HDOP":0.7<br>  } | 7614 3201 | 290505 |
| 04/17/2024 | Wednesd ay | 12:30 PM PDT | {<br><br>"EVENT_NO_TRIP":222064319,<br><br>"EVENT_NO_STOP":222064321,<br><br>"OPD_DATE":"18DEC2022:00:00:00",<br>    "VEHICLE_ID":3059,<br>    "METERS":104866,<br>    "ACT_TIME":47402,<br>    "GPS_LONGITUDE":-122.572777,<br>    "GPS_LATITUDE":45.43827,<br>    "GPS_SATELLITES":12.0,<br>    "GPS_HDOP":0.7 | 7808 5328 | 310035 |

| | | | | | |
|---|---|---|---|---|---|
| | | | } | | |
| 04/18/2024 | Thursday | 12:30 PM PDT | {<br><br>"EVENT_NO_TRIP":222781137,<br><br>"EVENT_NO_STOP":222781139,<br><br>"OPD_DATE":"19DEC2022:00:00:00",<br>    "VEHICLE_ID":2902,<br>    "METERS":887,<br>    "ACT_TIME":21816,<br>    "GPS_LONGITUDE":-122.84017,<br>    "GPS_LATITUDE":45.509218,<br>    "GPS_SATELLITES":12.0,<br>    "GPS_HDOP":0.7<br>  } | 8148 5328 | 310250 |
| 04/19/2024 | Friday | 12:30 PM PDT | {<br><br>"EVENT_NO_TRIP":223461587,<br><br>"EVENT_NO_STOP":223461601,<br><br>"OPD_DATE":"20DEC2022:00:00:00",<br>    "VEHICLE_ID":2935,<br>    "METERS":12758,<br>    "ACT_TIME":21665,<br>    "GPS_LONGITUDE":-122.988017,<br>    "GPS_LATITUDE":45.521155,<br>    "GPS_SATELLITES":12.0,<br>    "GPS_HDOP":0.8<br>  } | 8268 8306 | 330365 |
| 04/20/2024 | Saturday | 12:30 PM PDT | {<br><br>"EVENT_NO_TRIP":224100636,<br><br>"EVENT_NO_STOP":224100637,<br><br>"OPD_DATE":"21DEC2022:00:00:00",<br>    "VEHICLE_ID":2902,<br>    "METERS":73,<br>    "ACT_TIME":25211,<br>    "GPS_LONGITUDE":- | 8237 2950 | 329324 |

| | | | | | |
|---|---|---|---|---|---|
| | | | 122.844165,<br>  "GPS_LATITUDE":45.503545,<br>  "GPS_SATELLITES":12.0,<br>  "GPS_HDOP":0.9<br> } | | |
| 04/21/2024 | Sunday | 12:30 PM PDT | {<br><br>"EVENT_NO_TRIP":224411073,<br><br>"EVENT_NO_STOP":224411153,<br><br>"OPD_DATE":"22DEC2022:00:00:00",<br>  "VEHICLE_ID":2902,<br>  "METERS":56812,<br>  "ACT_TIME":30645,<br>  "GPS_LONGITUDE":-122.78018,<br>  "GPS_LATITUDE":45.427715,<br>  "GPS_SATELLITES":10.0,<br>  "GPS_HDOP":1.0<br> } | 8337 3964 | 339476 |
| | | | | | |

Additionally, include screenshots for the parts C, H and I

1. Output of crontab -l: Your scheduled cron jobs.
2. systemctl status: This will show the status of your receiver program.
3. VM instance schedule: This will display the schedule settings for your GCP VM instance.