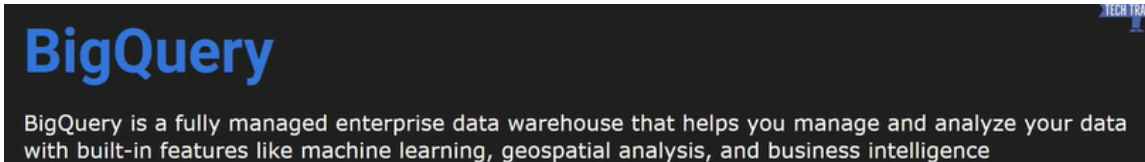Introduction

1. WHat is bigquery?
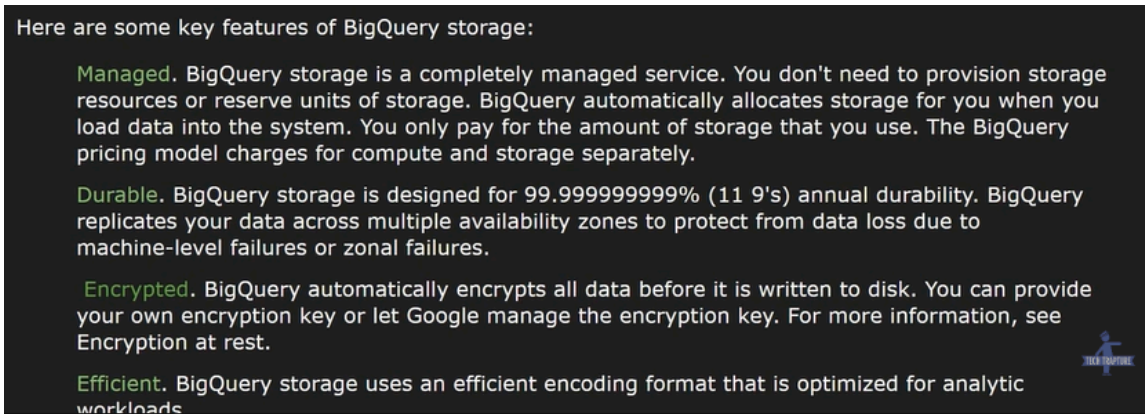


**BigQuery**

BigQuery is a fully managed enterprise data warehouse that helps you manage and analyze your data with built-in features like machine learning, geospatial analysis, and business intelligence

So here we dont need to manage any compute instance…we create the dataset and table…and load the data…and rest of infrastructure needs will be taken care by google
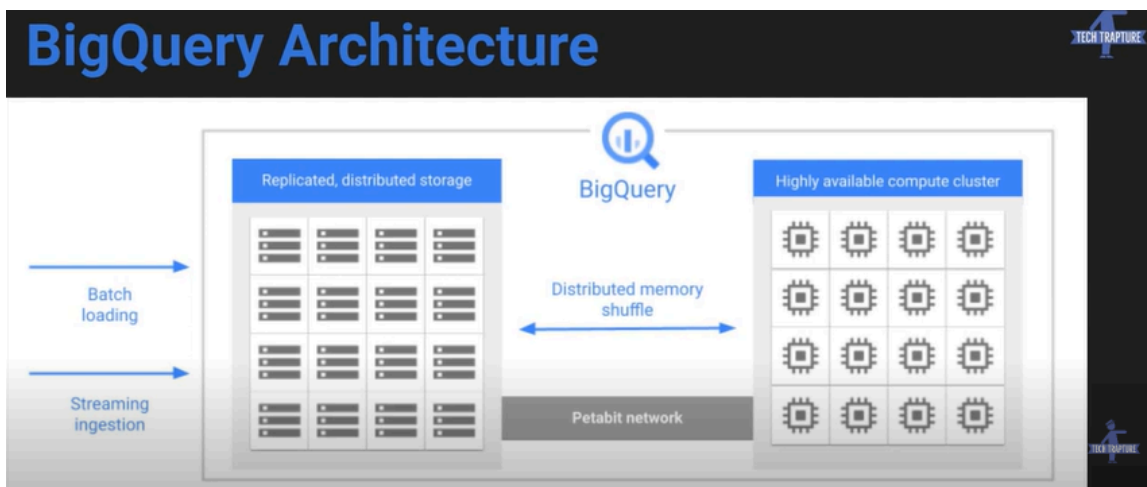
2. Key features



Here are some key features of BigQuery storage:

Managed. BigQuery storage is a completely managed service. You don't need to provision storage resources or reserve units of storage. BigQuery automatically allocates storage for you when you load data into the system. You only pay for the amount of storage that you use. The BigQuery pricing model charges for compute and storage separately.

Durable. BigQuery storage is designed for 99.999999999% (11 9's) annual durability. BigQuery replicates your data across multiple availability zones to protect from data loss due to machine-level failures or zonal failures.

Encrypted. BigQuery automatically encrypts all data before it is written to disk. You can provide your own encryption key or let Google manage the encryption key. For more information, see Encryption at rest.

Efficient. BigQuery storage uses an efficient encoding format that is optimized for analytic workloads

3. Architecture



**BigQuery Architecture**

4. Terminology

5. BigQuery Columnar storage



Dataset and table in BigQuery

1. Bigquery is a google manages cloud datawarehouse

2. Click on your project and then we can create dataset



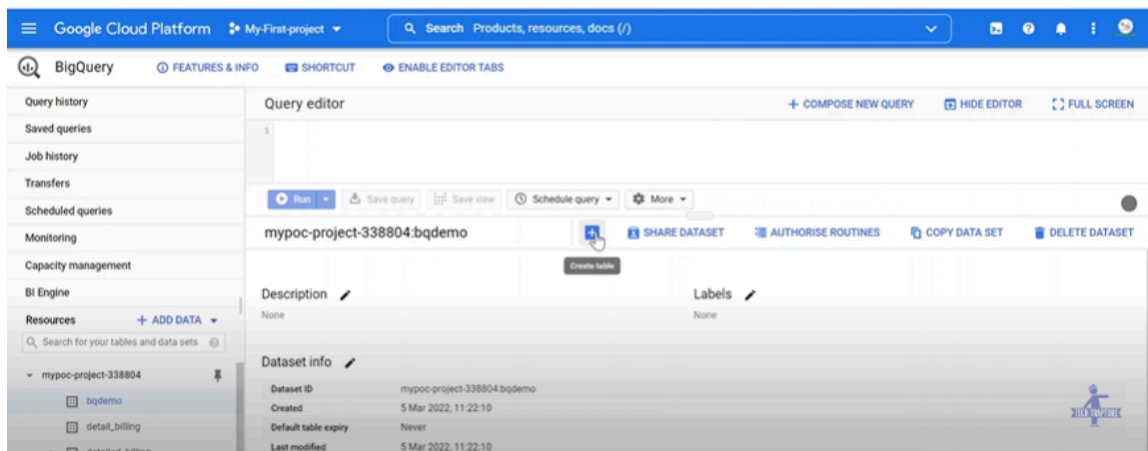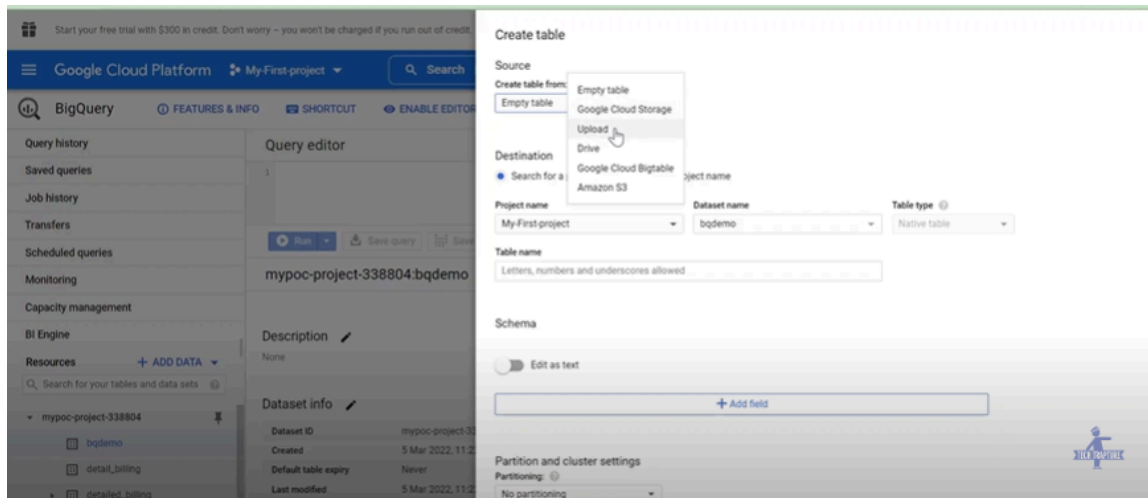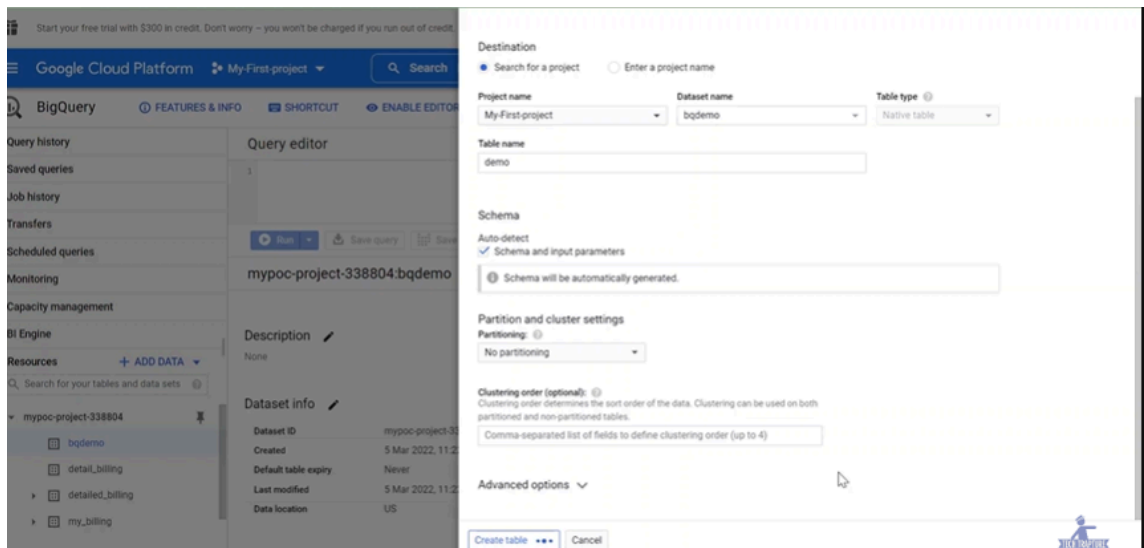3. Click on create dataset..and give a name with default configs
4. After that we create a table in dataset



5. Here we can upload the tables from our local computer or from google cloud storage etc
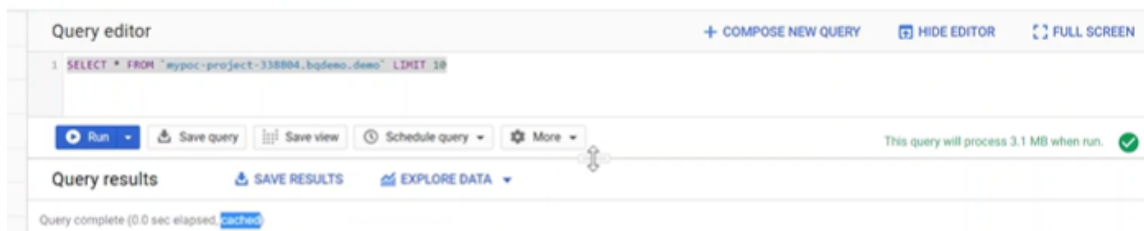
6. Next we upload the file and give a new table name



and create it
7. After that we can query our data
8. Here if we run the same query again..then it retrieve the output from the cache



Partitioning in Bigquery

1. Here we see …how we can optimie the performance using partitioning and clustering
2. Inside our project here..we have multiple datasets and tables

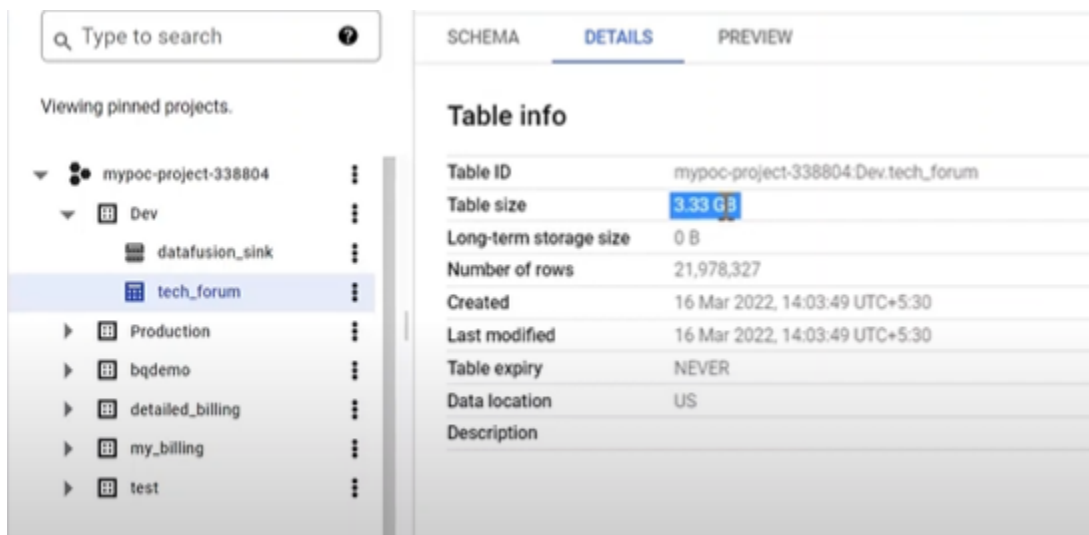3. Here we'll work with techforum table




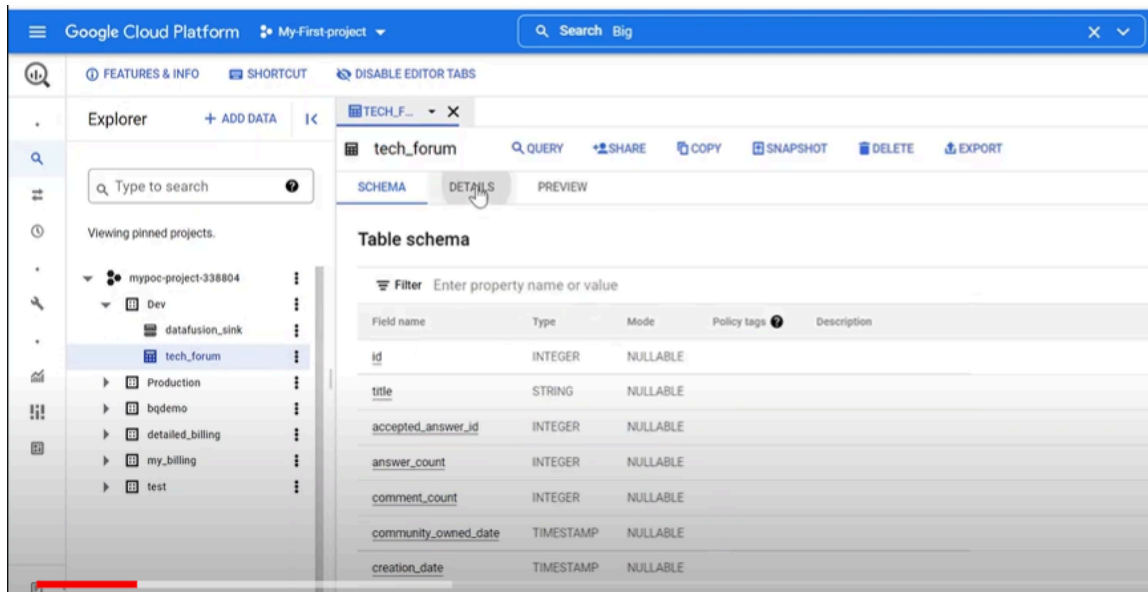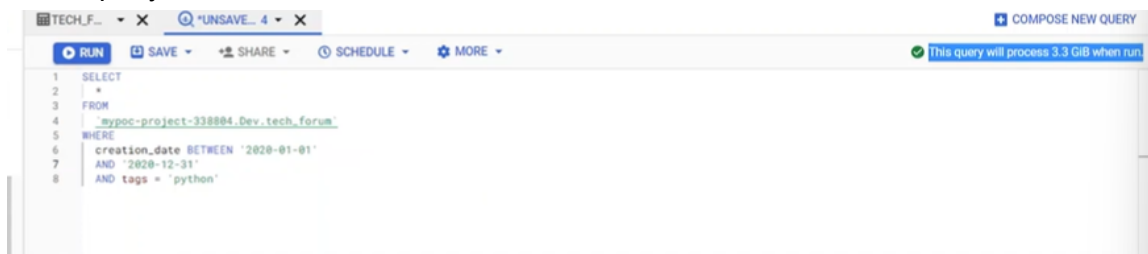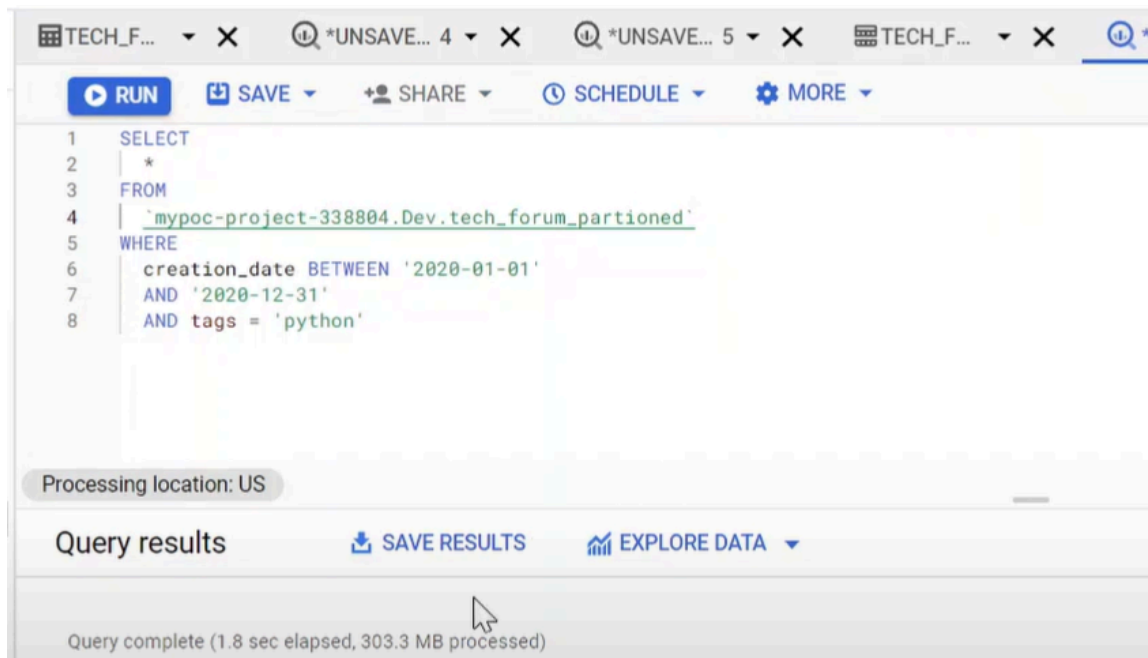
Table size is 3.3GB
4. Now lets query this table
5. As the data is stored in columnar format in bigquery…if we use select * from table limit 10….it read the entire file…even though we need only 10rows
6. Practical example
7. In this query



even if we filtered the data…it is reading complete table(3.3GB)

8. SO to avoid this..we can use partitions
9. Now here we partitioned the table by day
10. Now if run the same query with partition..then time complexity would br

```
    RUN      SAVE ▾    + SHARE ▾    SCHEDULE ▾    MORE ▾
1   SELECT
2    *
3   FROM
4     `mypoc-project-338804.Dev.tech_forum_partioned`
5   WHERE
6     creation_date BETWEEN '2020-01-01'
7     AND '2020-12-31'
8     AND tags = 'python'
```

Processing location: US

Query results        SAVE RESULTS        EXPLORE DATA ▾

Query complete (1.8 sec elapsed, 303.3 MB processed)

BigQuery Federated Queries

# BigQuery External tables

An external table is a table that acts like a standard BigQuery table. The table metadata, including the table schema, is stored in BigQuery storage, but the data itself resides in the external source.

You can use external tables with the following data sources:

Bigtable

Cloud Storage

Drive

Amazon S3

Azure Blob Storage

1.
2. Next we'll open a bigquery…and create a external tables

3. Here we'll create a external table from GCS



and we create table

4. So here the bigquery will only contain the metadata…but the actual data will be in the external source itself



BigQuery external tables let you query data stored in other Google Cloud services directly, without having to copy the data into BigQuery itself. This saves storage space and avoids the need for additional data movement. Here's a breakdown and an example to illustrate:

**Benefits of BigQuery External Tables:**

- **Reduced Storage Costs:** You only pay for storing the data in its original location, not in BigQuery's internal storage.
- **Faster Analytics:** Querying data in its native location can sometimes be faster than copying it to BigQuery first.
- **Simplified Workflow:** Manage and analyze data without needing to load it into BigQuery, streamlining your process.

5. Federated tables

**BigQuery Federated Queries**

A federated query is a way to send a query statement to an external database and get the result back as a temporary table.

Federated queries use the BigQuery Connection API to establish a connection with the external database.

In your standard SQL query, you use the EXTERNAL_QUERY function to send a query statement to the external database, using that database's SQL dialect. The results are converted to BigQuery standard SQL data types.

You can use federated queries with the following external databases:

    Cloud Spanner

    Cloud SQL

practical

BigQuery federated queries empower you to directly query data stored in other Google Cloud databases without physically moving the data into BigQuery. This enables you to analyze data residing in separate services seamlessly, offering a unified view for insightful exploration. Here's a breakdown of federated queries along with an example to illustrate its functionality:

6.

7. To establish a connection..we click on add data…and create an external

**External data source**

Connection type
Cloud Spanner ▾

Connection ID *
fed-demo

Data location
us (multiple regions in United States) ▾ ❓

Friendly name
demo

Description

Database name *
projects/{project}/instances/{instance}/databases/{database} ❓

☐ Read data in parallel ❓

CREATE CONNECTION    CANCEL

connection…and give details

Database name *
projects/{project}/instances/{instance}/databases/{database} ❓

☐ Read data in parallel ❓

8. Database name should be in this format
   click on create connection

9. After that we can query the table using external query

 it gives all the tables in external DB

10. Using federated query we have executed queries of external DB
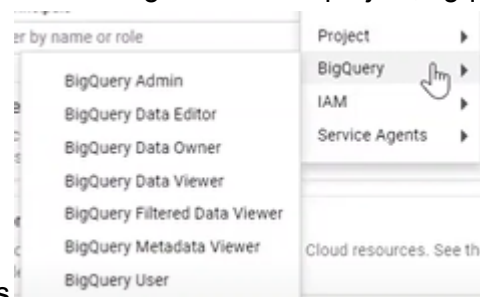


BigQuery Access control

1. Here we'll learn access controls
2. Lets suppose we have 2 accounts in GCP…first account has is owner and have access to all the files including production

3. 2nd user is employee and has only access to dev
4. Now to grant the access to 2nd user…we have to add dataset permissions and add principals(2nd user email) and select a role
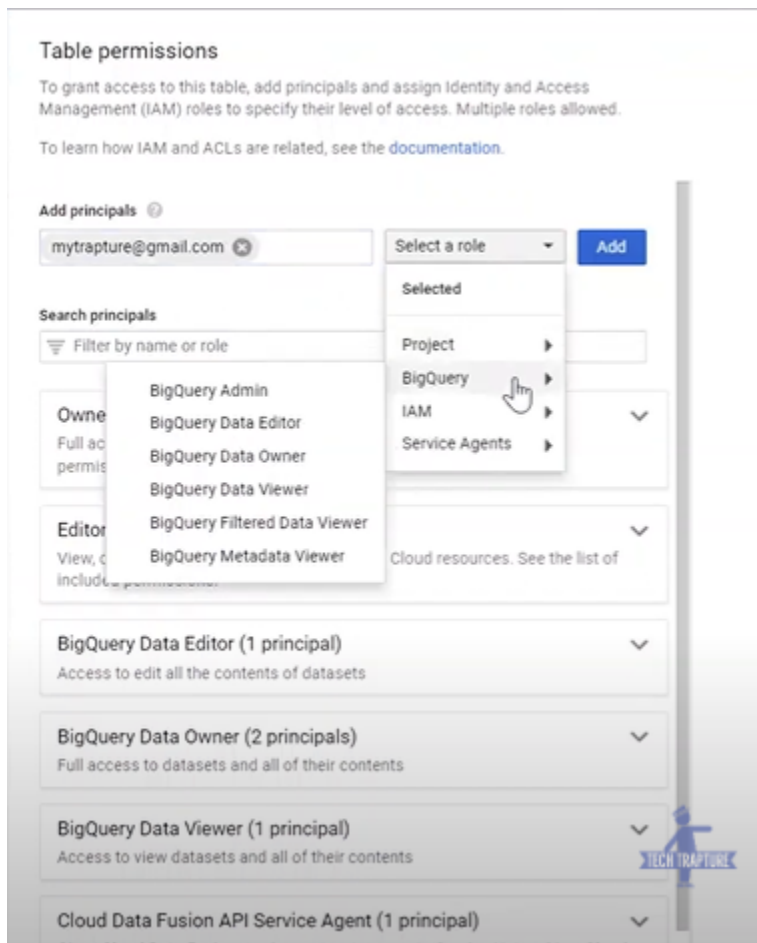


5. Here we can assign access on project,bigquery etc…we gave big query(data viewer)



access

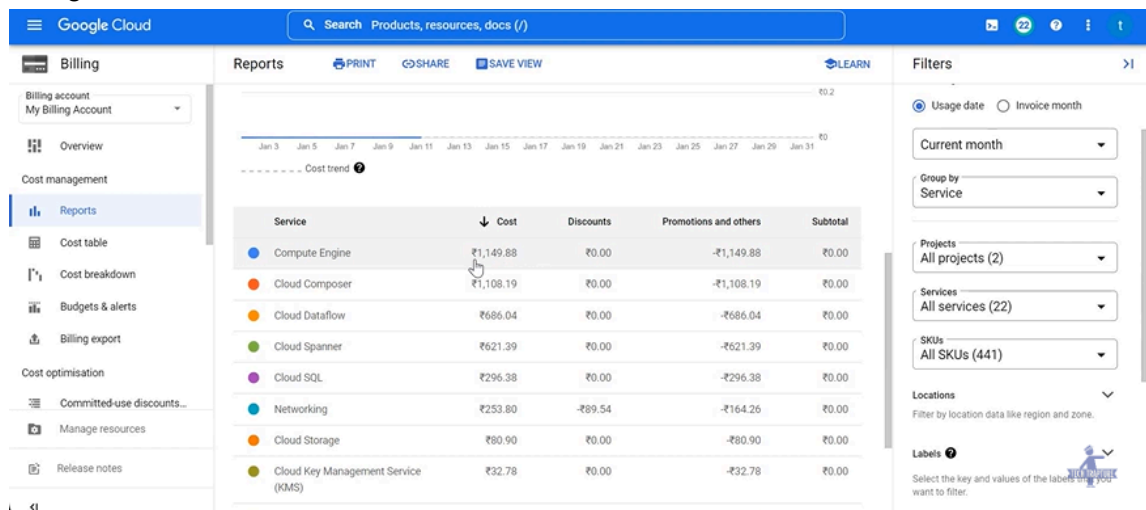6. We can also give access levels on tables too

7. For that we go to the tables in our database and select the table and click on share table..then add principal and select role
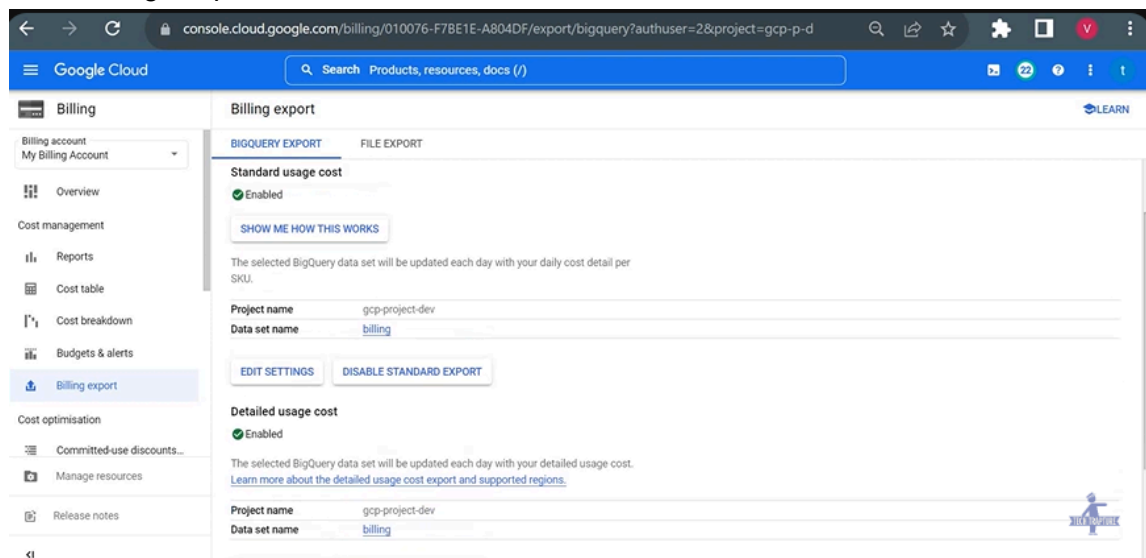


Export Billing data to Bigquery

1. We'll learn how can we export billings data from big query and create some dashboards
2. The billings permissions will only have who recreate the account
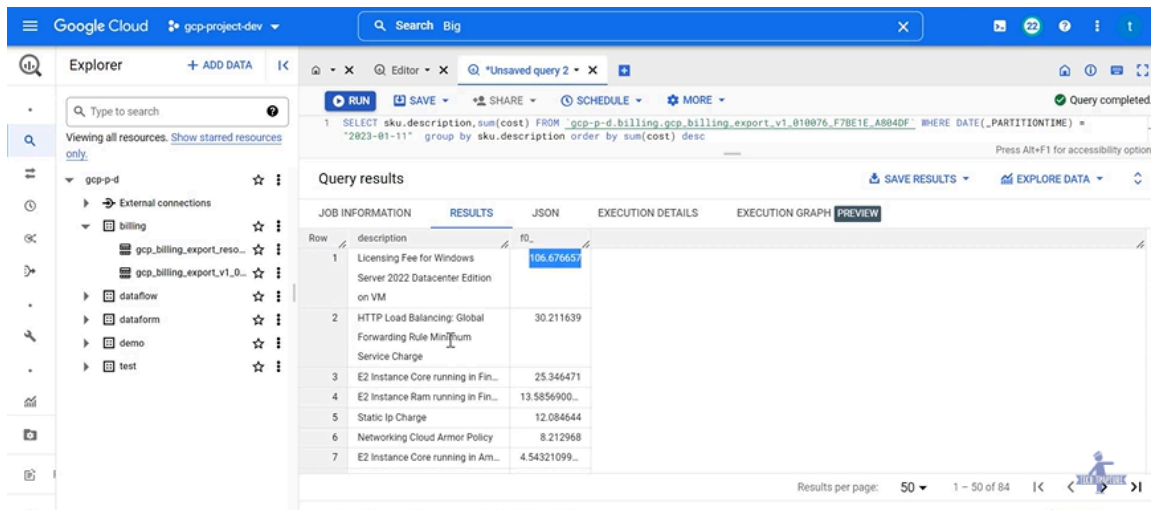
3. Billings of our services



4. And we want our team to analyze the data…to have them access this data..we have export this data to bigquery
5. Go to billings export



and enable export to GCP

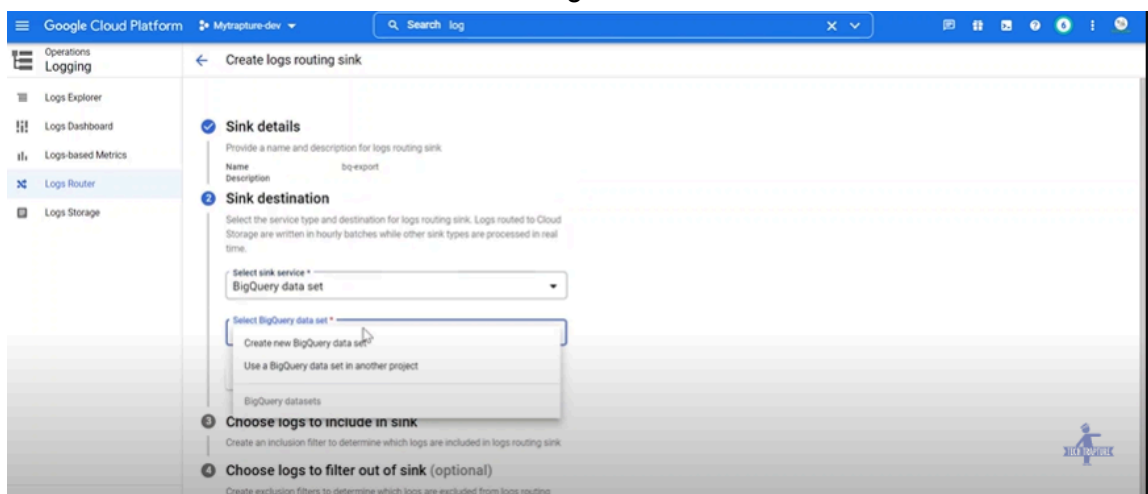6. Lets go to bigquery and we can query the billing data…to see more accurate details



7. Now we can create dashboards in the looker studio…well connect to GCP and share the datasets

## Export GCP logs to BigQuery

1. Search logs and go to logs router
2. Then we have to create sink…but before that we have to create a dataset in the bigquery to store this log tables/data
3. Next we have to create the sink details in log router



select the dataset for logs..and create log sink
4. Now we can logs in our bigquery

Restore deleted data in Bigquery

1. We have time travel features…where we can see our data back in time at particular moment…we can only go backwards upto 7 days

2.

**Time Travel:**

- **Concept:** Allows you to access historical versions of your data within a specific timeframe.
- **Functionality:** Think of it like a rewind button for your data. You can query the state of your data at a specific point in the past, even if it has been updated or deleted since then.
- **Benefits:**
  - **Data Recovery:** Recover from accidental data modifications or deletions.
  - **Auditing:** Analyze how data has changed over time for audit purposes.
  - **Debugging:** Identify the root cause of issues that might be related to historical data changes.

- Use time travel for:
  - Short-term data recovery needs.
  - Analyzing recent data changes.
  - Debugging issues potentially related to historical data modifications.

**Snapshot:**

- **Concept:** Creates a static copy (snapshot) of your data at a specific point in time.
- **Functionality:** Similar to taking a picture of your data at a particular moment. You can then analyze or store this snapshot independently.
- **Benefits:**
  - **Long-term Retention:** Snapshots can be stored indefinitely, allowing access to historical data outside of the time travel window.
  - **Faster Access:** Querying snapshots is typically faster than using time travel as it's just accessing a static copy.

- Use snapshots for:
  - Archiving historical data for long-term analysis.
  - Regulatory or compliance needs that require data retention for extended periods.
  - Situations where faster access to historical data is crucial.

3.