

Interview Techm

dbt

classmate

Date _____

Page _____

Data Flow, data proc, federate in bq,
External tables in bq, partition,
Clustering, challenge in project, pax
scenarios, queries, primary key in bq
pub/sub, DAT tool, bq slot, data mod
elling, difference b/w External and Inte
nal table. Ways to export from Bq.

GCP Services for Data Analytics

Bigquery → It is a Serverless and cost-effective
enterprise data warehouse that works
across the cloud and scales with your
data.

Looker → It is a Business Intelligence tool and
data analysis platform for data exploration,
modelling, transformation and visualization.

DataProc → It is a fully managed and highly
scalable service for running Apache Hadoop
and Apache Spark.

Dataflow → It is a fully managed streaming
analytics service that minimizes laten

processing time, and cost through auto-scaling. Dataform → It helps analytic teams transform data in Big using SQL. With dataform data engineers, data analyst can develop Big-table definitions.

Normalization

It is a process to reduce redundancy and dependency. The primary goal of normalization is to eliminate data anomalies and ensure data integrity.

which all anomalies are solved in Normalization is realized in following forms:

- ① Insertion Anomaly, No incomplete or contradictory info in insertion.
- ② Deletion Anomaly, Prevention of loss of valuable info. when deleting data.
- ③ Update Anomaly, update are made in away to maintain data integrity.

Q) Jenkins is a CI/CD Tool that works

A)

Checkout → Build → Test → Deployment

or

Plan → Code → Push → Build → Test → Deploy

Works on groovy syntax

Minute, hour, dayofmonth, Month, dayofyear

0 2 * * *

Means it will trigger daily at 2AM

Q) Global Variable in python is defined outside all the function that is accessed by all the function.

If you want to modify the value of global variable within a function you need to use global keyword:

global variable = 0

Round function syntax

ROUND([number], [ndigs]).

Google Cloud Logging

It is a service by GCP that allow you to store, search, analyse and monitor log generated by your applications and infrastructure.

Window function

These function in SQL are used to perform calculation over a specified range of rows related to the current row.

Ranking → Select employee_id, salary, RANK()
OVER (ORDER BY salary DESC) AS salary_ran
from employee;

Running total → Select o_date, o_units, sum
JUM(COALESCE(o_qty, 0)) OVER (ORDER BY o_date)
Running total from orders;

Three types of Window functions are:-

- ① Aggregate window function;
- ② Ranking window function
- ③ Value window function

Aggregate window function :-
divide the table into n buckets.

Value window function

4 LAG() → To access previous value & value before the current row value.

Syntax : LAG(ColName, 1, 0) OVER [R] OVER(PARTITION BY org.) ORDER BY [year]

LEAD() → It is used to access the subsequent row along with data of current row.

e.g., lead(Order) OVER(ORDER BY [Order])

Benefits of Compute Engine:

- 4 Storage Efficiency.
- 4 Stability.
- 4 Easy Integration.
- 4 Security.
- 4 Compute Globally as per requirement.

- 4 Default bucket location is US region.
- 4 Compute Engine lets you create VM on Google's Infrastructure.
- 4 The data in VM depends on type of disk used.
 - If persistent disk the data is retained even if instance stops. In case of local SSD the data is not retained.
- 4 IAM (Identity and access Management).

4 Storage and Data Services.

Object Cloud Storage (GCS)	Relational Cloud SQL	Non-Relational Cloud Spanner	Wave Cloud Firestore	Cloud Bigtable
----------------------------	----------------------	------------------------------	----------------------	----------------

4. Three Nos of database in OCP.

- ① Database
- ② Cloud Storage
- ③ Cloud Bigtable

Difference b/w these -

Count(*), Count(1)

Select Count(*) from Emp;

Select Count(1) from Emp;

Select Count(DISTINCT 2) from Emp;

Select 1 from Emp;

Deloitte Interview

classmate

Date _____

Page _____

Object storage	Relational	Non-relational	Warehouse
GCS	Cloud SQL	Cloud Firestore	BigQuery
	Cloud Spanner	Cloud Bigtable	

Cloud SQL → It is fully managed Relational database service that supports MySQL, PostgreSQL and SQL Server.

Features → Auto backups, replication of data for high availability, automated patch management.

Cloud Bigtable → Cloud Bigtable is a fully managed highly scalable NoSQL database service that is designed to handle large amounts of data with high read/write speed.

Dataflow, it is a managed service for executing a wide variety of data processing patterns.

Dataproc → It is a fully managed service used to run Apache Hadoop and Spark clusters.

Databricks and Dataproc do the same large-scale data processing.

DataLab → It is an interactive tool for analysing, visualizing data on GCP. It is designed to work with big data technologies such as Big Data, Apache Spark. It provides Jupyter Notebook environment that integrates with various GCP services.

These Notebooks are stored in Google Source Repository, anyone can collaborate, visualize on the code. We can install the DataLab using "gcloud component install datalab".

DataPrep → It is a cloud-based data preparation platform that helps organization clean, enrich, and transform raw, messy and unstructured data into a format suitable for analysis and reporting.

Helps in streamlining the data preparation process for technical and non-technical users to work with data.

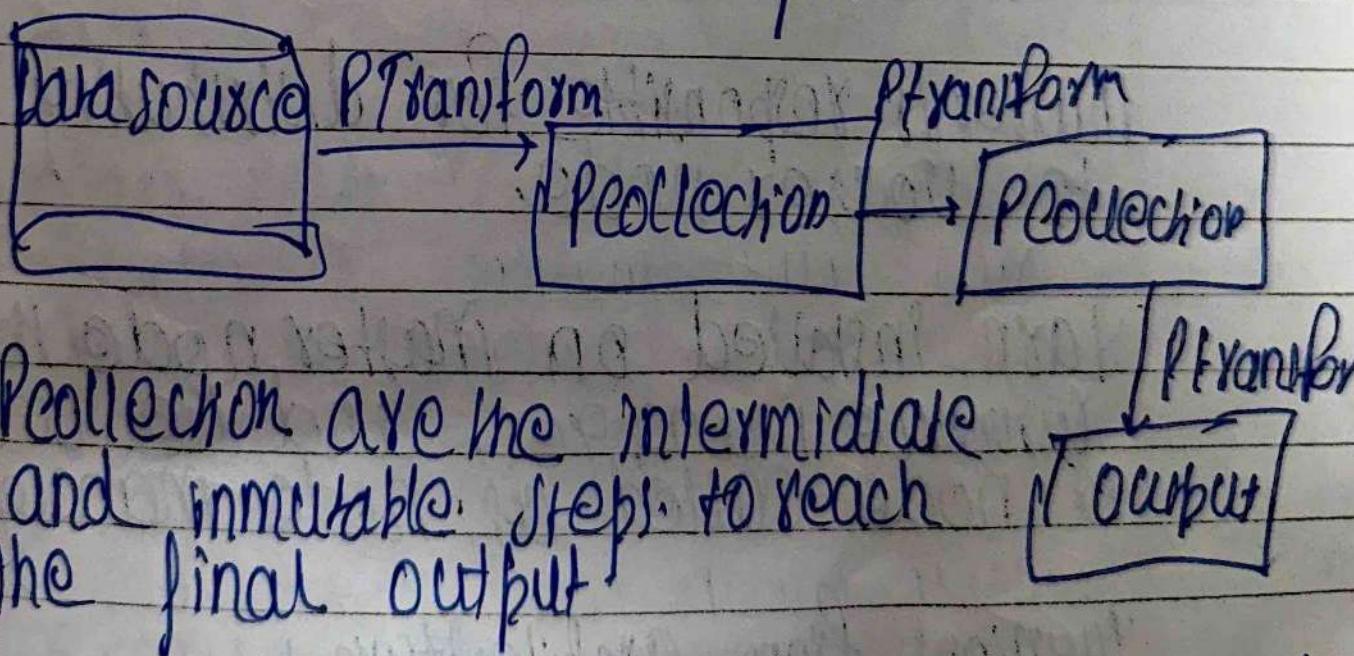
It is managed by Trifacta, and suggests

Intelligent data transformation methods help collaborate with other users.

Pub/Sub → It is a fully managed cloud, real-time messaging service that allows you to send and receive messages b/w independent applications.

Dataflow

It is a managed service to run data processing in a wide variety of patterns with the available default templates.



PCollection are the intermediate and immutable steps to reach the final output.

e.g. To read the word count on AC and post the result on the AC.

BigQuery, It's fully managed serverless data warehouse and analytical platform provided by google cloud platform. It's used for reporting and analytics purpose.

SPARK Architecture

4. Multiple machine connected together called as Cluster.

Works on master-slave architecture, the slave also called as worker node.

Master is responsible for distributing job to worker node.

YARN installed on master node that is resource manager and the worker node called as node manager.

Questions from architecture -

Drivers → The main program that creates spark context, coordinating the distributed

Operations :

1. Spark Context → The entry point to any spark functionality.

Cluster manager, Manages resources across spark application on cluster such as YARN, standalone.

RDD → The Resilient Distributed Dataset is the fundamental data structure in spark, representing a distributed collection of objects that can be processed in parallel.

4. Spark SQL → This module of spark for structured data processing. It provides a SQL wrapper and supports various sources like Parquet, Avro, JSON etc. Bridges gap b/w RDD and relational db.

4. Executor → These are worker nodes responsible for executing the tasks on spark cluster.

Narrow and wide Transformations.

- 4 Narrow Transformations involve one-to-one mapping and do not require shuffling of data.
- 4 Wide Transformation involve shuffling and redistribution of data across partitions.
- 4 Spark driver program is responsible for creating Spark Context (Entry point for spark job), Coordinating tasks, Managing the overall execution of spark application.
- 4 It runs main function and communicates with cluster manager to allocate resources and execute tasks.
- 4 Transformations are operations applied to RDDs to create new RDD.

Narrow - filter, union, flatMap, map, flatmap
Wide - groupby, reduceby, sort

*Date _____
Page _____*

Google Transfer Appliance → It provides an easy and secure way to physically ship large amounts of data to the cloud.

It is a briefcase sized appliance shipped from Google to the organization to upload data and return back to Google so it can be uploaded to the cloud storage. It comes in 2 size capacity 40TB and 300TB.

Cloud Transfer Service → It is a managed service of cloud to transfer data from GCP or any other cloud platform to GCS. We can migrate the data from the data transfer service job without writing any code.

We can schedule it as per requirement that one time or recurring basis.

BigQuery Data Transfer Service → It enables you to move data into BigQuery from variety of sources. Transfers can be both one-time as well as ongoing or scheduled.

Docker is a platform that allows you to package and run applications in containers. A container is like a lightweight, standalone and executable package that contains everything needed to run a piece of software, including the code, runtime, software, libraries. It ensures the code or application runs consistently across different environments. And solves the problem of 'This runs on my machine'.

Kubernetes

Kubernetes is like conductor for your software. It helps in managing all the parts of complex application, it helps to coordinate and orchestrate your software component ensuring they play together seamlessly.

Kubernetes helps in:

- ① Deploy applications
- ② Scale Applications
- ③ handle failures
- ④ Update Applications

Lazy Execution in Spark

Lazy Execution refers to strategy where transformations on RDD (operations) or Dataframe are not immediately executed. Instead, they are recorded and the actual computation is deferred until action is invoked.

It is a key optimization feature in spark that helps in improving the efficiency of distributed data processing.

Window Functions

functions that allow to perform calculations on or across range of rows related to the current row.

↳ Row_Number()

↳ Rank()

↳ Dense_RANK()

↳ NTILE()

↳ SUM()

↳ AVG()

↳ MIN()

↳ MAX()

~~False~~

Sales

P_id	sale_date	revenue
1	=	=

classmate

Date

Page

- ① calculate running total of revenue for each product over time.

Select P_id, sale_date, revenue,
 $\text{sum}(\text{revenue}) \text{ over}(\text{partition by } p_id \text{ order by}$
 $\text{sale_date})$ AS total from sales.

- ② Assign a rank to each product base on their total revenue.

Select p_id, revenue, RANK() OVER(ORDER BY
 $\text{sum}(\text{revenue}) \text{ desc})$ AS revenue_rank.
from sales group_by p_id.

- ③ Row_Number assign unique number to each employee based on their salary in descending order.

Select e_id, e_name, salary,
ROW_NUMBER() OVER(ORDER BY salary desc)
AS salary_rank from employee.

- ④ give me the rank of employee based on salary in a emp table, using CTE

h with Ranked Employee As (

Select e_id, ename, salary,
RANK() OVER(ORDER BY salary DESC) AS salary-
rank
from employee

4 Select e_id, c_name, salary, salary_rank
from ranked_employee,

⑤ Hire and highest salary using window fun
with Ranked Employee As (

Select e_id, ename, salary,
DENSE_RANK() OVER(ORDER BY salary DESC)
rank from employee)

Select ename, salary from Ranked Employee
where salary_rank = 2.

Optimizing SQL Queries

We need to address the performance bottleneck to make the query more efficient.

- ① USE indexing. (Columns in where, join, order by) are indexed.
- ② Analyze execution plan (explain before the query).
- ③ limit the result set.
- ④ Partitioning large table.
- ⑤ Avoid select * as it fetches unnecessary columns.
- ⑥ check for blocking or locking queries (the query might be blocked by another transaction).
- ⑦ Review and Optimize JOIN conditions.

Optimizing Spark SQL Jobs

- ① Choose appropriate file formats, as Parquet is often more efficient than CSV or JSON.
- ② Leverage Partitioning and bucketing to organize data on the disk as it reduces amounts of data to be scanned.
- ③ Using Broadcast join for smaller datasets or dataframes.
- ④ Using a broadcast variable as a local copy if stored to each worker node that reduces data transport.
- ⑤ Tuning shuffle operation we can adjust the partitions based on size of data - default is 200. Set `spark.sql.shuffle.partitions`.
- ⑥ Dynamic partition pruning. → It is a technique of spark to push the unnecessary data and push down the required data for better data retrieval.

Set `spark.sql.adaptive.enabled = True`

Performance Tuning

classmate

Date _____

Page _____

- ⑦ Sufficient hardware resources allocated to cluster.
- ⑧ Spark UI and Spark history server are used to identify the bottlenecks and optimize your queries accordingly.
- ⑨ There is a complex query with multiple subqueries and join condition. How can you optimize it?
 - ① Examine the execution plan.
 - ② Ensure Join conditions and where clause have indexed column.
 - ③ Review the subqueries and evaluate if they can be optimized or replaced with Join or EXISTS clause.
 - ④ Limit the result set.
 - ⑤ Consider materialized views if the query is having heavy aggregation or computation.
 - ⑥ USE CTE (Common Table Expression).

⑦ Partitioning table

Materialized View → It is a type of view which stores the data to avoid recomputing the results of making the heavy queries easy to execute.

Create materialized view to Abs as
select bid, sale from sales;

In DWH, we create materialized views as a summary table by which we can improve the performance of analytical queries.

Why we use Views, Virtual table

- ① Views encapsulate complex queries into single name.
- ② Views can be used to control access to underlying table.
- ③ Views promote code readability. Instead of writing complex queries again.

Cloud Composer, It is fully managed Apache Airflow service by which we can create, schedule, monitor, and manage our work-flows.

In data Engineering, workflow represents a series of tasks for ingesting, transforming, analyzing, or utilizing data. Airflow creates workflows using DAG's.

Despite of several connectors available to migrate data, why we don't use in live environment?

- ① Downtime Concern.
- ② Data Volume and Velocity.
- ③ Network reliability.
- ④ Custom logic or m^o. transformations
- ⑤ Fear of data loss.

Teladate architecture Vs BO Architecture

Bigquery architecture is decoupled architecture that is storage and processing are separated making it less expensive.

Colossus → It is a successor of AFS (Google file system) that is used for storage.

Dremel → It is Compute Engine to process the query to give output.

It has Root Server, Mixers, Leaf Nodes. where Leaf Node is used to read data from Colossus. Mixers does the heavy aggregate part and Root Server co-ordinates everything.

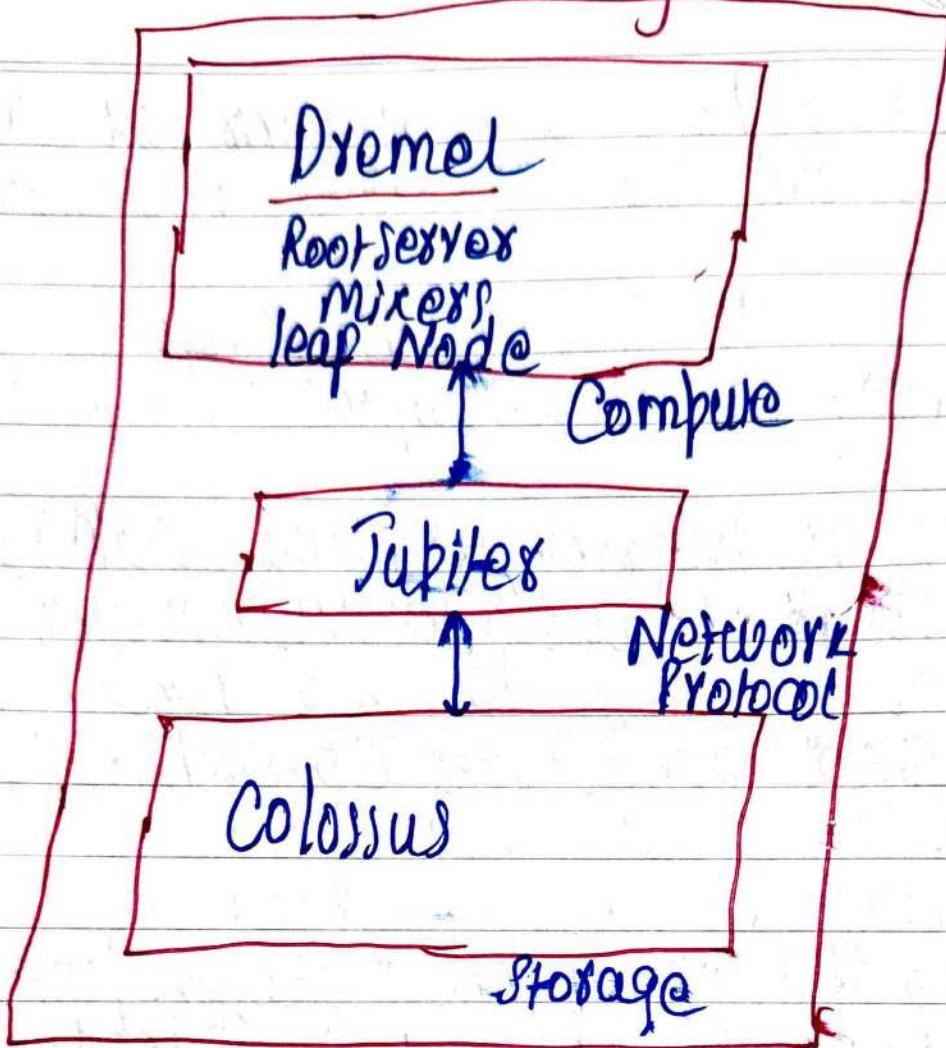
Jupiter → It is high speed network protocol to connect Dremel and Colossus. Called as perabit Network protocol.

And whole architecture is orchestrated by Borg. that co-ordinates and schedules the overall architecture.

BQ is column oriented Database or storage.

Borg.

Date _____
Page _____



Read a CSV, infer schema then filter them
and write to a table

from pyspark.sql import SparkSession
spark = SparkSession.builder.appName('m')
getOrCreate()

my_schema = StructType([

StructField('id', IntegerType(), True),

`df = spark.read.format('csv').options('mode':
 'option', 'schema': my-schema)
 .load('path')`

`transformed_df = df.filter(col('id') > 100)`

`transformed_df.write.format('parquet').
 saveAsTable('db.tablename').`

HSBC

- 4 SQL Transformers in GCP
- 4 life cycle rules in buckets
- 4 Decoupled BD architecture
- 4 Storage classes
- 4 Python code to find hobbies from with T
in a file list of dictionaries
- 4 Create a cluster table
- 4 hadoop coupled system
- 4 BD table is clustered on year, month,
day and hour, all are int. Find
the data in date range (use partition
concept).
- 4 find percentage like in table of emp
- 4 how BD stores and compute data

PwC

- ↳ Introduction
- ↳ past Experience
- ↳ Questions related to project (Data loading problem, how to validate data)
- ↳ Volume and size of objects
- ↳ SQL query to get the latest and older record from SCD table.
- ↳ What is SCD and how we apply SCD on table. What is its advantages.
- ↳ Python code to find frequency of each word in string.
- ↳ SQL to get the 2nd highest of the salary.
- ↳ Apache airflow basics, dags, RDD.
- ↳ BQ architecture
- ↳ Python code to read a csv and give the top 5 apps from one dataset of playstore.

Infosys (2ounds).

1st round

- ↳ Qn180 and role in project
- ↳ questions on projects
- ↳ BQ architecture
- ↳ nested and repeated fields in BQ.
- ↳ ways to export data from table in BQ
- ↳ Python basic questions, list, dictionary
- ↳ what is set and tuple
- ↳ a scenario where we can use tuple
- ↳ spark architecture
- ↳ migration experience questions

and bound

- ↳ Qn180 and project responsibility
- ↳ Python proficiency
- ↳ SQL question on rank and dense_rank
- ↳ GenAI basic knowledge
- ↳ Any example of GenAI

Public Sapient

- ↳ Q. Have an input-df with two column
 ↳ I want to add new column in the df with True or False.
- ↳ Write full syntax to read a csv file in spark
- ↳ Write syntax to declare the schema in spark.
- ↳ What is inferSchema, header
- ↳ What if inferSchema is false.
- ↳ Input → name, math, english
 Sam 70 30
- Out → name subject math
 Sam Math 70
 Sam Eng 30

↳ Table join Count Questions

↳ bank, grade score

scid	a	1
hdfe	a	3
rnb	b	2

out →	bank	grade	Score	Competitor score
	pcici	a	1	2
	Hdfc	a	3	2
	Pnb	b	2	2

4 Com_avg → ii) Competitor average score.

- 4 Questions related to project loading data.
- 4 how loading was done
- 4 many counter questions on loading volume
- 4 why used manual approach to load data
- 4 Transfers services.