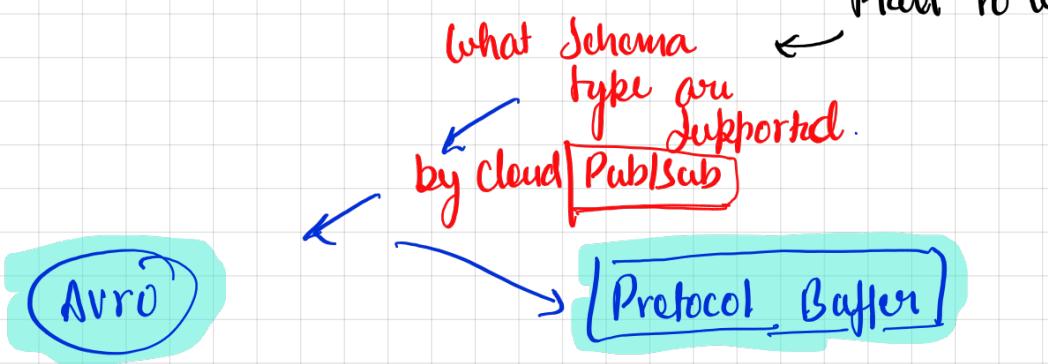


↳ You developing → Distributed System → want to decouple two Services

↳ Also messages uses a standard format ↴



Parquet → is open source file format

↳ used in Hadoop.

↳ Cloud Function → Need to regularly backup → Backup ↴

what  
CMD we can  
use

Backup files need to be  
stored in Cloud Storage

run in the  
background

→ Backend data store **Export** gs://fileStore-backup -aoyne

↓  
is the Cmd for → to save data to Cloud storage Bucket

↳ Gsutil → Only use to manage Cloud Storage

↳ not Cloud datastore

↳ Materialized View → has →

**Bigquery ka**

And got to know  
higher charges ↵

as per aspect

- ↳ frequency of Material View refresh
  - ↳ Data Store
- These 2 thing ↑ Cost

Number of users → reading Data → from → Materialized View  
will not effect cost

But total amount of data Scanned

↳ Migrating a data → from data warehouse → On-Premises to google Cloud.

use of data warehouse

↳ BIGQUERY BI → BiEngine → is an ↑ fast → in memory analysis source

↳ Cloud Memory Stor → is a cache and better suited

for storing  
key-value data

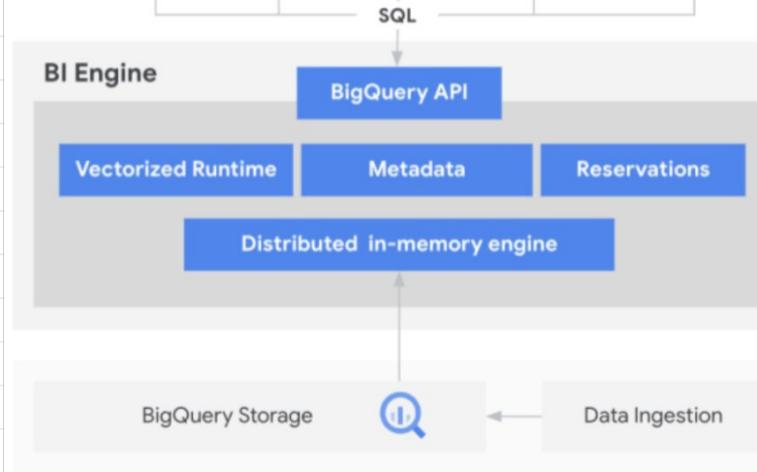
for application

which needs low latency access to data

There is no bigtable Bi Engine Service

↳ BiEngine can accelerate SQL queries from any source  
including thesis

Written by data visualization tools  
 Can manage Cached tables  
 ↓  
 for On going optimization



- Bi Engine might Not fit for
- Rely heavily on wildcards of queries
  - BigQuery's feature with BI Engine doesn't support.
  - Then run unsupported features
    - ↳ Including External tables
    - ↳ Non SQL user-defined functions

↳ New workload → deployed → to → Cloud DataProc  
 ↳ configured with an autoscaling policy.

Fetch failed Exception  
 ↳ got

It happen when → Shuffle Data is lost  
↳ Nodes are decommissioned.

To avoid hot-spotting → in your Bigtable Cluster

we are using the prefix of UUID prefix

Issue will be not working as we expected

There is hot-  
spotting  
writing  
Data  
to Bigtable

Cause Can be

UUID → that are sequentially generated

It should use Version 4  
and generate random number generator

Column families don't accept hotSpotting

Cost of Running → some data pipeline → for large Batch Jobs

gets Resource Scheduling

Data flow like RS

which Reduces  
Batch Processing  
Cost → will Reduce Cost

By using → Scheduling Technique

Premittable VMs along with regular VMs

Dataflow Streaming Engine: → Is for Stream processing.

Dataflow Shuffles → Provides fast Execution of Batch Jobs

↳ But don't reduce necessarily cost

Apache Beam → runner → require more management.

↳ Migration Several data warehouse → BigQuery

↳ Cloud storage for

Machine learning Data

ML Engineer

Data Analyst

are having difficult  
finding  
data set

↳ Cloud data Catalog → Can automatically Extract metadata

↳ from sources including

2

- Cloud storage
- Big Query
- Cloud Bigtable
- Cloud Pub/Sub
- Google Sheet

Cloud logging → used for recording data about Events

↳ not best way to Collate metadata.

Cloud fusion → is an ETL tool  
not a metadata

## ↳ Extraction tool

Monitoring System → hai → Jo → Apache Beam Runner  
Pe hai



↳ Temperature → received over past hour is analyzed

↳ If any temperature reading ↴

is more than  
standard deviation

↳ Consistent time

Interval

So it is Sliding Window / hopping window

Consistent time interval in  
a stream.

↳ Data warehouse team is concerned that some data ↴

Sources  
may have

Dont want to  
Bring

incorrect Data or invalid Data

Poor quality of Data



Into Data Warehouse

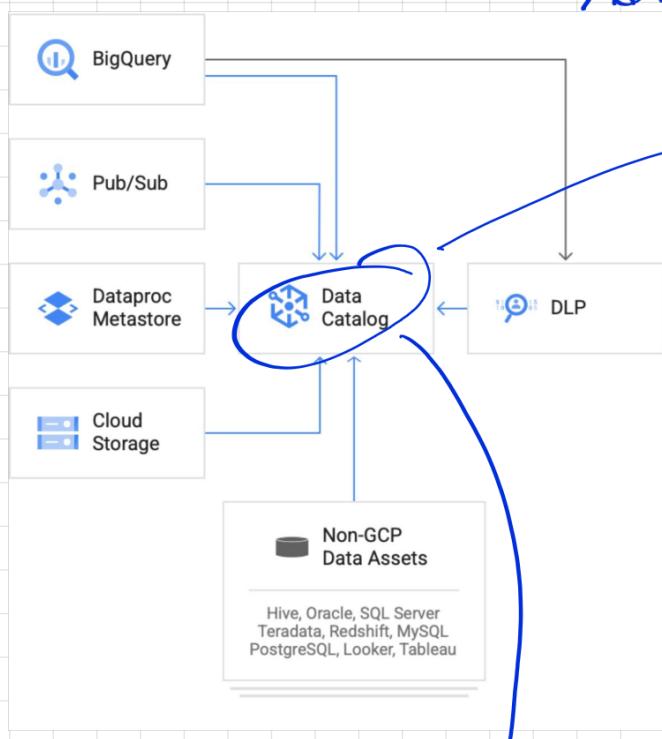
For this we need  
to do ↴

# data quality assessment

## Source System

From the source Data after it Extracted & from

- ↳ Check of ranges
- ↳ distribution Values
- ↳ Attributes Count the number.



↳ Collect all metadata

Provide Power Based  
Predicate Based Solution

## Search Solution

↳ Por technical and business  
metadata solution

Data catalog → not index a Data  
within a Data Entry

↳ Data Catalog Only  
Indexes the

→ help for  
↳ protect  
↳ discover  
↳ classify → most sensitive  
Data

↳ Data Loss Prevention is to

→ help for  
↳ protect  
↳ discover  
↳ classify → most sensitive  
Data

we have Compute Engine → attached GPU

↳ But GPUs not used when you train  
tensorflow model

→ how to ensure the GPU can be used by

training model

Install GPU Driver

Deploying VM image

which already have  
installed GPU

⇒ MongoDB replacement → [in GCP]

[Cloud firestore]

[Bigtable]

wide column No SQL databases → not good  
replacement

For MongoDB

Load Data from AVRO → Bigtable

Dataflow → Cloud Storage Bigtable template

↳ Cloud SQL Auth Proxy

Recommended way to connect to Cloud SQL



For Deploying a Cloud SQL database to  
Prediction

Strong Encryption is used to Protect

the Confidential

and Integrated of

Not to Perform  
Authentication

Data

## Type of [ Data Model ]

Document Model → Support → Semi Structured Schemas

Star Schema

Snowflake Schema

are denormalized relational model

used in Data Warehousing

That frequently change

Network Model → is used as → graph-like structure

such as transportation network

\* Cloud pub/sub → [ Delivery at least Once ]

Cloud Audit logs

this contain logs till 30 days

to keep them longer than that

→ we need to preserve the logs for

)  
Cloud Storage

Cloud logging → doesn't keep logs beyond 30 days

↳ doesn't support retention policy

Cloud monitoring

→ collects and displays metrics

↳ it doesn't store logs

Bigtable → is not a good option for logs

↳ low latency writes at high volumes  
and provides for key lookup.  
queries.

That requires range scanning

## Overfitting and Underfitting

→ with respect to Data

Poor performance → ML

↳ for example → Supervised learning → ML is best understood as approximately

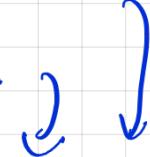
target function ( $f$ )

$$y = f(x)$$

Characterization  
describe

→ Range of Classification

↳ Prediction Problem



Machine algorithm

That can be used to address them.

# Cloud Spanner

→ Database having Performance issue)

Query Insight

Where this is  
Pop up

or wo issue is  
with form of  
due to ↴  
Problematic  
query

↳ Set of Prebuilt Dashboard

going to be faster  
for resolving issue

which makes it Very Easy to find  
spikes in ↑↑ in performances



Query is use more  
then this spikes will  
Pop up

↳ CPU usage

↳ Top query using the CPU in that time window

Filter	Enter property name or value	Query or request tag	CPU utilisation	CPU (%)	Execution count	Avg latency (ms)	Avg. rows scanned
<input type="checkbox"/>	FPRINT	Query or request tag	<div style="width: 100%;"> </div>	88.92	578	431.89	38.28
<input type="checkbox"/>	29457682432492560	SELECT playerUUID FROM (SELECT play...	<div style="width: 100%;"> </div>	8.26	183	119.63	77,487.25
<input type="checkbox"/>	7531924216159805264	SELECT gameUUID FROM (SELECT gam...	<div style="width: 100%;"> </div>	0.93	183	10.67	121.16
<input type="checkbox"/>	-438027238268179866	SELECT PlayerUUID, Stats, Current_game...	<div style="width: 100%;"> </div>	0.61	1,844	11.61	0
<input type="checkbox"/>	-9057436179270446999	INSERT players (playerUUID, player_nam...	<div style="width: 100%;"> </div>	0.00	0	0.00	0

J

It will give more insight such as

How many Sights → Rows  
Column

↳

Retrieved

↳ Enabled by default

↳ It is worth → Only queries  
 ↳ DML ↳ Supported by  
 ↳ Query insight

Spanner mutation → Not available

↳ Cloud Spanner → update data → Data manipulation language

DML

↳ Spanner Mutation

```
Document instance = getDocument(dbClient, customerId, accountIds, lastEvent);
List<Mutation> mutations =
    Arrays.asList(
        Mutation.newInsertBuilder( table: "Ledger" ) Mutation.WriteBuilder
            .set("CustomerId") ValueBinder<Mutation.WriteBuilder>
            .to(CustomerId) Mutation.WriteBuilder
            .set("AccountId") ValueBinder<Mutation.WriteBuilder>
            .to(AccountId) Mutation.WriteBuilder
            .set("TransactionId") ValueBinder<Mutation.WriteBuilder>
            .to(TransactionId) Mutation.WriteBuilder
            .set("Date") ValueBinder<Mutation.WriteBuilder>
            .to(java.time.LocalDate.now().toString()) Mutation.WriteBuilder
            .set("Amount") ValueBinder<Mutation.WriteBuilder>
            .to(Amount) Mutation.WriteBuilder
            .set("Details") ValueBinder<Mutation.WriteBuilder>
            .to(Details) Mutation.WriteBuilder
            .build(),
        Mutation.newUpdateBuilder( table: "Account" ) Mutation.WriteBuilder
            .set("CustomerId") ValueBinder<Mutation.WriteBuilder>
            .to(CustomerId) Mutation.WriteBuilder
            .set("AccountId") ValueBinder<Mutation.WriteBuilder>
            .to(AccountId) Mutation.WriteBuilder
            .set("Balance") ValueBinder<Mutation.WriteBuilder>
            .to(Balance) Mutation.WriteBuilder
            .set("LastTransactionTime") ValueBinder<Mutation.WriteBuilder>
```

```
public static void createAccount(DatabaseClient dbClient, String customerId){
    dbClient
        .readWriteTransaction()
        .run(transaction -> {
            String accountId = String.valueOf(UUID.randomUUID());
            String createdOn = java.time.LocalDate.now().toString();
            String balance = "0";
            String sql =
                String.format("INSERT INTO Account (CustomerId, AccountId, CreatedOn, Balance) "
                    + " VALUES ('%s', '%s', '%s', %s)"
                    , customerId, accountId, createdOn, balance);
            long rowCount = transaction.executeUpdate(Statement.of(sql));
            System.out.printf("%d record inserted with AccountId %.\\n", rowCount, accountId);
            return null;
        });
}
```

(DML) → Insert / update / Delete SQL

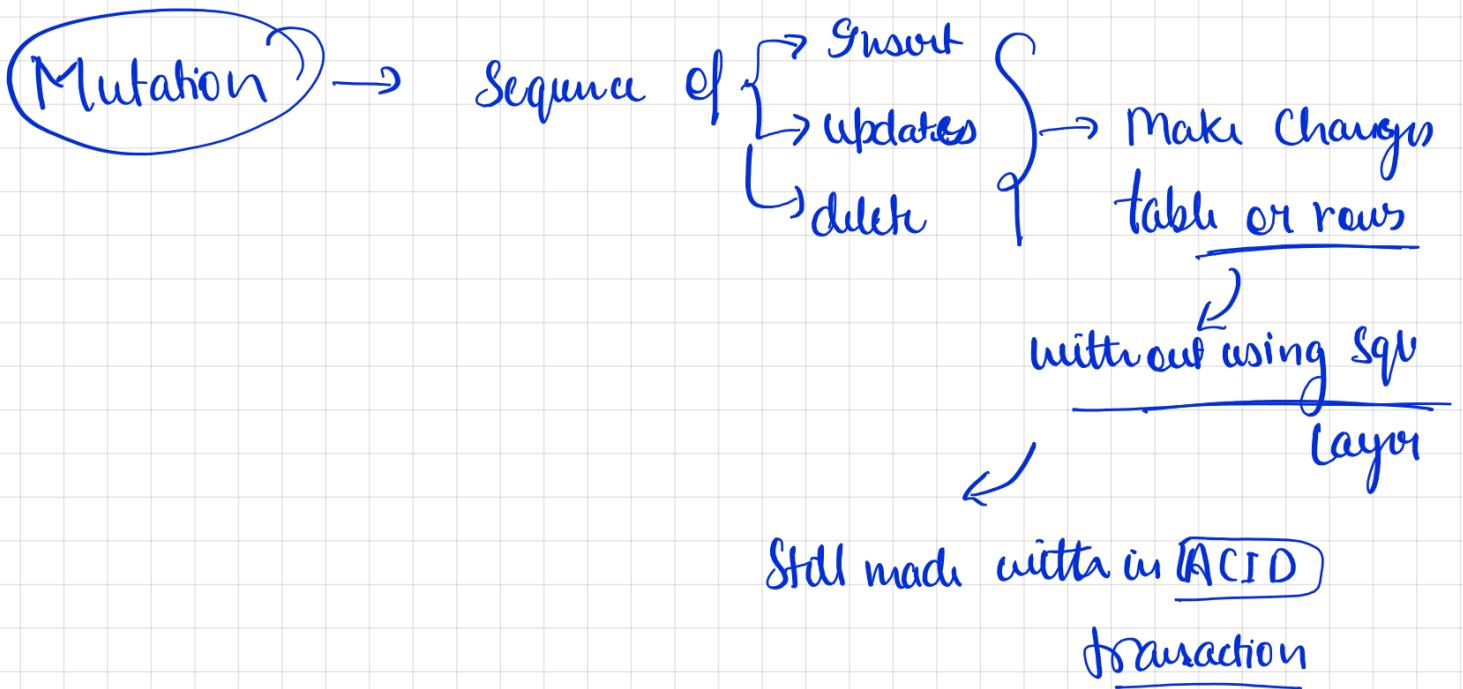
Standard DML

↓  
great for single statement

↳ Special implementation  
for bulk data changes

(DML) → If you change a Data → That Change will remain to later  
Statement

Within the same transaction



Changes  $\rightarrow$  arr  $\rightarrow$  done  $\rightarrow$  Tise order mai specified hai

↳ Buffered in Client locally and then sent,

To spanner

When transaction is finished

Due to that

Mutation don't support read after write

Within same transaction

Constrain is → check at commit time

Rather after each mutation

DML → great for OLTP users

Large Change Operation use → Mutation  
 ↓ Partition DML

↳ If you want to check → Read | Change ) within same transaction

Or constraint should be check after each statement

use DML

DML or Mutations?

Operations	DML	Mutations
Insert Data	Supported	Supported
Delete Data	Supported	Supported
Update Data	Supported	Supported
Read Your Writes (RYW)	Supported	Unsupported
Upsert (insert or update)	Unsupported	Supported
SQL Syntax	Supported	Unsupported
Constraint Checking	After every statement	At commit time

For Example

financial Bank hai

Jab account bata kar raha hai

Single query hoga

Then DML

```
public static void createAccount(DatabaseClient dbClient, String customerId){
    dbClient
        .readWriteTransaction()
        .run(transaction -> {

            String accountId = String.valueOf(UUID.randomUUID());
            String createdOn = java.time.LocalDate.now().toString();
            String balance = "0";
            String sql =
                String.format("INSERT INTO Account (CustomerId, AccountId, CreatedOn, Balance) "
                    + " VALUES ('%s', '%s', '%s', %s)"
                    , customerId, accountId, createdOn, balance);
            long rowCount = transaction.executeUpdate(Statement.of(sql));
            System.out.printf("%d record inserted with AccountId %s.\n", rowCount, accountId);
            return null;
        });
}
```

Then

→ update the account

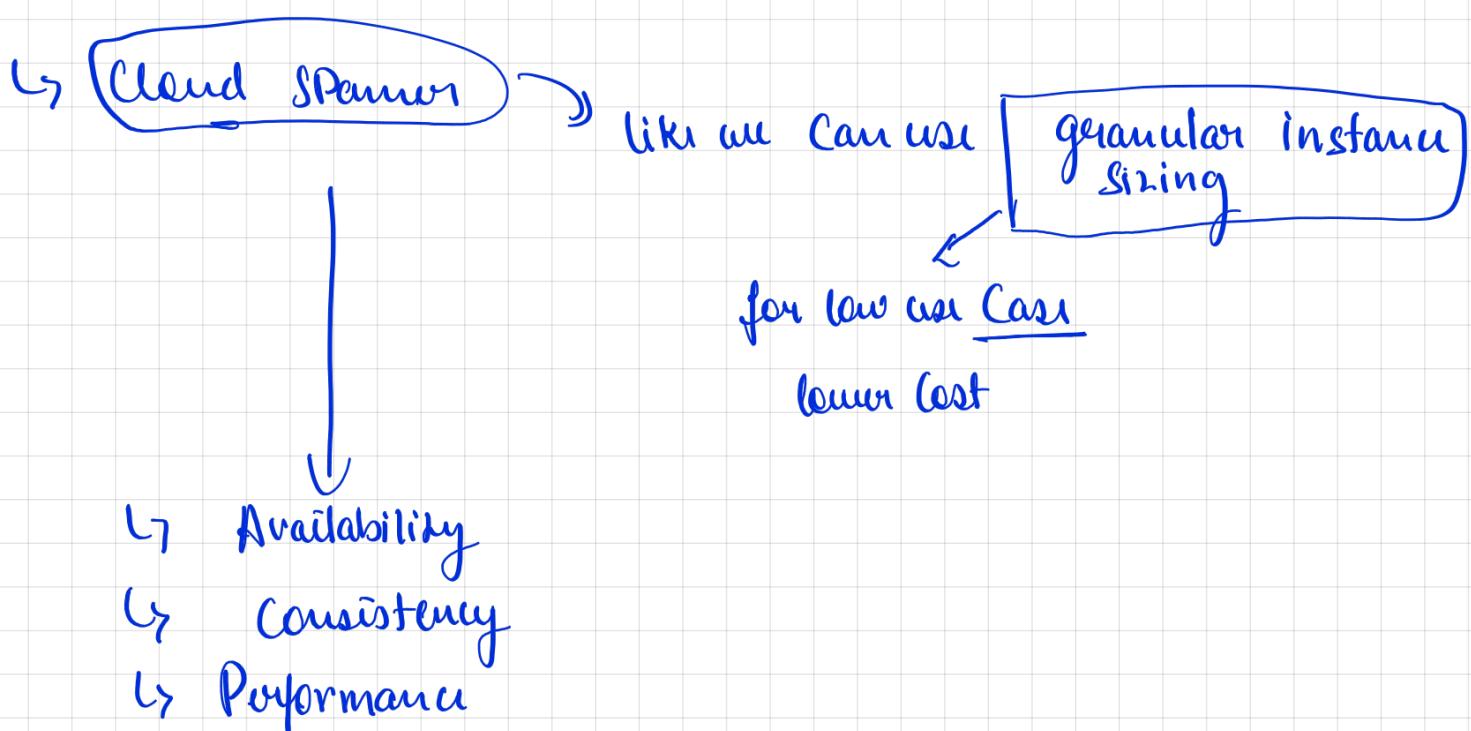
↳ Enter the Details into ledger

Mutation

As we don't need to

Read the Data after the Changes.

(DML) Statement And Mutation Statement  
shouldn't be mixed  
within the Same transaction



Spanner → Easily → Manage load → without require complex sharding.

↳ support Scalability by → Separating → support  
↳ storage

↳ **BigQuery** → support materialized Views  
↳ Performance ←  
↳ efficiency ↗  
↳ Precomputed Views that  
Periodically Cache the result of a query of increased.

Queries that use → Materialized Views → generally faster  
that retrieve the same data only from base tables.  
↓  
Consume fewer resources than queries.

Materialized Views → Key characteristic → zero maintenance  
↓  
push Data  
Smart tuning

Process get Benefits

↓ from Views → OLAP Online Analytics processing  
↓ ETL  
↓ BI

↓ Improve query performance → Preaggregate Data :> Aggregation and Streaming Data

↓ Pre-filter Data  
Run query Only in Particular subset  
Prjoined Data

Query Joins , Espically between Large and small table