

Data Warehouse

Index

- OLAP vs OLTP
- What is data warehouse
- BigQuery
 - Cost
 - Partitions and Clustering
 - Best practices
 - Internals
 - ML in BQ

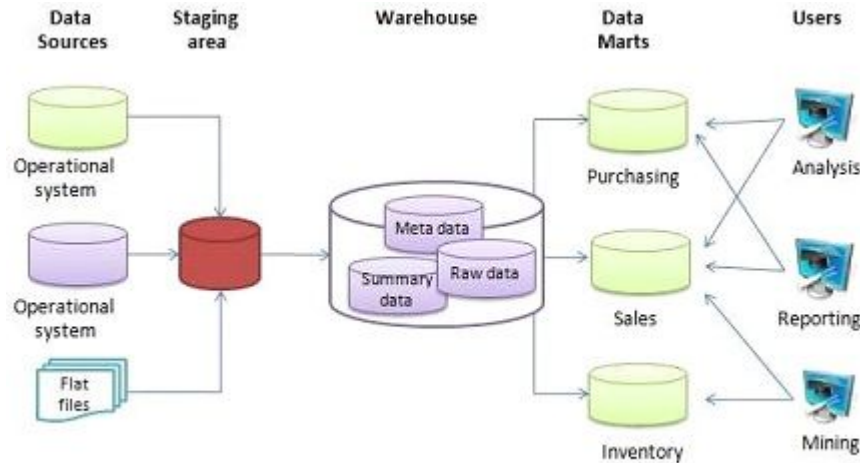
OLAP vs OLTP

	OLTP	OLAP
Purpose	Control and run essential business operations in real time	Plan, solve problems, support decisions, discover hidden insights
Data updates	Short, fast updates initiated by user	Data periodically refreshed with scheduled, long-running batch jobs
Database design	Normalized databases for efficiency	Denormalized databases for analysis
Space requirements	Generally small if historical data is archived	Generally large due to aggregating large datasets

	OLTP	OLAP
Backup and recovery	Regular backups required to ensure business continuity and meet legal and governance requirements	Lost data can be reloaded from OLTP database as needed in lieu of regular backups
Productivity	Increases productivity of end users	Increases productivity of business managers, data analysts, and executives
Data view	Lists day-to-day business transactions	Multi-dimensional view of enterprise data
User examples	Customer-facing personnel, clerks, online shoppers	Knowledge workers such as data analysts, business analysts, and executives

What is a data warehouse

- OLAP solution
- Used for reporting and data analysis



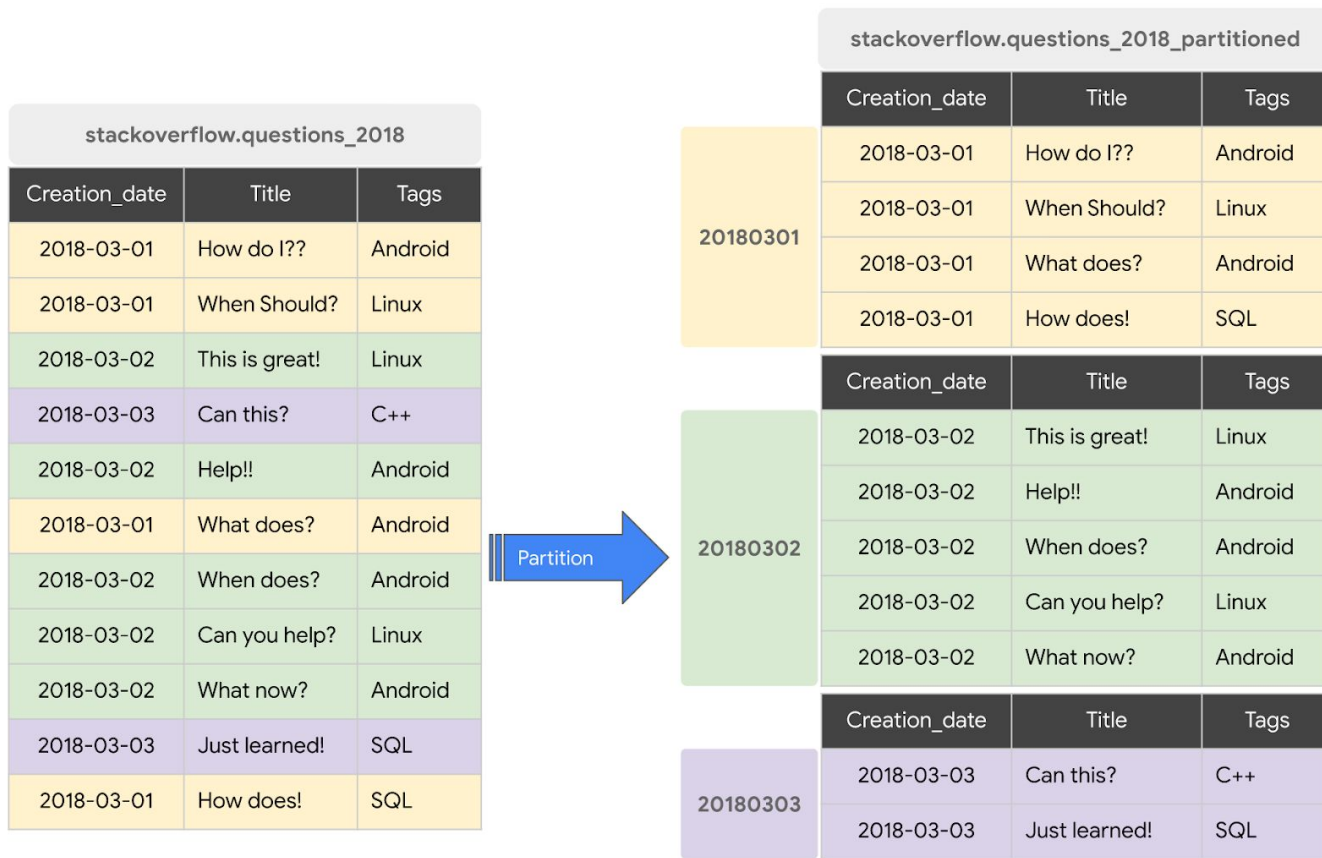
BigQuery

- Serverless data warehouse
 - There are no servers to manage or database software to install
- Software as well as infrastructure including
 - **scalability** and **high-availability**
- Built-in features like
 - machine learning
 - geospatial analysis
 - business intelligence
- BigQuery maximizes flexibility by separating the compute engine that analyzes your data from your storage

BigQuery Cost

- On demand pricing
 - 1 TB of data processed is \$5
- Flat rate pricing
 - Based on number of pre requested slots
 - 100 slots → \$2,000/month = 400 TB data processed on demand pricing

Partition in BQ



Clustering in BigQuery

BigQuery Partitioning & Clustering

Stack_Questions

Date	Title	Tags
2019-03-01	How do I??	Android	
2019-03-01	When Should?	Linux	
2019-03-02	This is great!	Linux	
2019-03-03	Can this?	C++	
2019-03-02	Help!!	Android	
2019-03-01	What does?	Android	
2019-03-02	When does?	Android	
2019-03-02	Can you help?	Linux	
2019-03-02	What now?	Android	
2019-03-03	Just learned!	SQL	
2019-03-01	How does?	SQL	

Stack_Questions_2019_03_01

Date	Tags	Title	...
2019-03-01	Android	How do I??	
2019-03-01	Android	What does?	
2019-03-01	Linux	When Should?	
2019-03-01	SQL	How Does?	

Stack_Questions_2019_03_02

Date	Tags	Title	...
2019-03-02	Android	Help!!	
2019-03-02	Android	When does?	
2019-03-02	Android	What now?	
2019-03-02	Linux	This is great!	
2019-03-02	Linux	Can you help?	

Stack_Questions_2019_03_03

Date	Tags	Title	...
2019-03-03	SQL	Just learned!	
2019-03-03	C++	Can this?	

BigQuery partition

- Time-unit column
- Ingestion time (`_PARTITIONTIME`)
- Integer range partitioning
- When using Time unit or ingestion time
 - Daily (Default)
 - Hourly
 - Monthly or yearly
- Number of partitions limit is 4000

Resource: <https://cloud.google.com/bigquery/docs/partitioned-tables>

BigQuery Clustering

- Columns you specify are used to colocate related data
- Order of the column is important
- The order of the specified columns determines the sort order of the data.
- Clustering improves
 - Filter queries
 - Aggregate queries
- Table with data size < 1 GB, don't show significant improvement with partitioning and clustering
- You can specify up to four clustering columns

BigQuery Clustering

Clustering columns must be top-level, non-repeated columns

- DATE
- BOOL
- GEOGRAPHY
- INT64
- NUMERIC
- BIGNUMERIC
- STRING
- TIMESTAMP
- DATETIME

Partitioning vs Clustering

Clustering	Partitoning
Cost benefit unknown	Cost known upfront
You need more granularity than partitioning alone allows	You need partition-level management.
Your queries commonly use filters or aggregation against multiple particular columns	Filter or aggregate on single column
The cardinality of the number of values in a column or group of columns is large	

Clustering over partitioning

- Partitioning results in a small amount of data per partition (approximately less than 1 GB)
- Partitioning results in a large number of partitions beyond the limits on partitioned tables
- Partitioning results in your mutation operations modifying the majority of partitions in the table frequently (for example, every few minutes)

Automatic reclustering

As data is added to a clustered table

- the newly inserted data can be written to blocks that contain key ranges that overlap with the key ranges in previously written blocks
- These overlapping keys weaken the sort property of the table

To maintain the performance characteristics of a clustered table

- BigQuery performs automatic re-clustering in the background to restore the sort property of the table
- For partitioned tables, clustering is maintained for data within the scope of each partition.

BigQuery-Best Practice

- Cost reduction
 - Avoid SELECT *
 - Price your queries before running them
 - Use clustered or partitioned tables
 - Use streaming inserts with caution
 - Materialize query results in stages

BigQuery-Best Practice

- Query performance
 - Filter on partitioned columns
 - Denormalizing data
 - Use nested or repeated columns
 - Use external data sources appropriately
 - Don't use it, in case u want a high query performance
 - Reduce data before using a JOIN
 - Do not treat WITH clauses as prepared statements
 - Avoid oversharding tables

BigQuery-Best Practice

- Query performance
 - Avoid JavaScript user-defined functions
 - Use approximate aggregation functions (HyperLogLog++)
 - Order Last, for query operations to maximize performance
 - Optimize your join patterns
 - As a best practice, place the table with the largest number of rows first, followed by the table with the fewest rows, and then place the remaining tables by decreasing size.

Internals

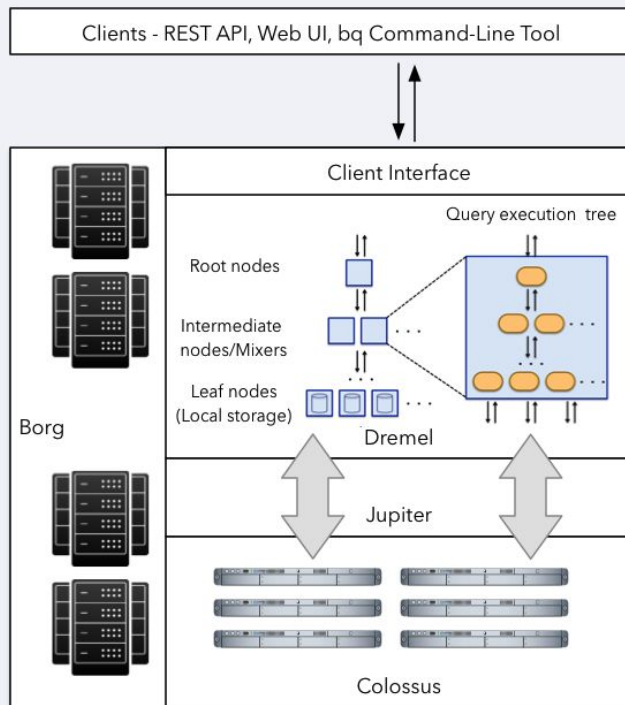
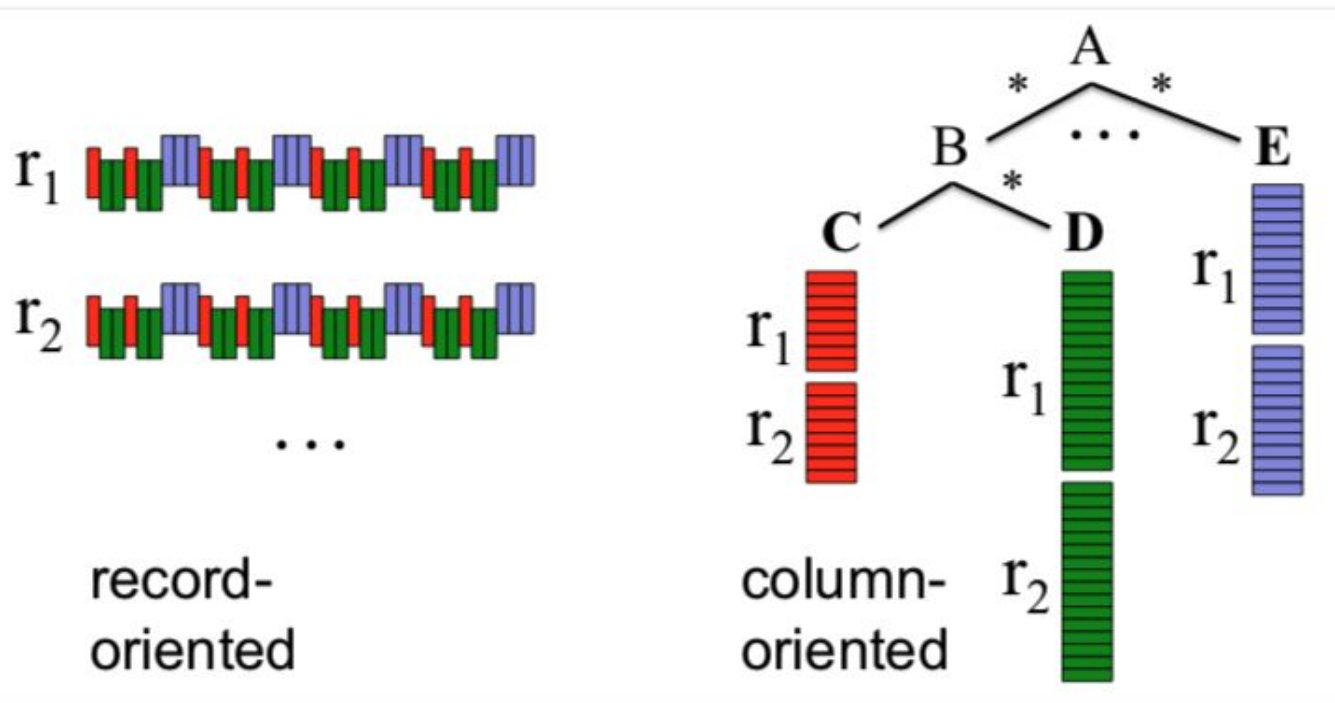


Figure-1: A high-level architecture for BigQuery service.

Internals



Internals

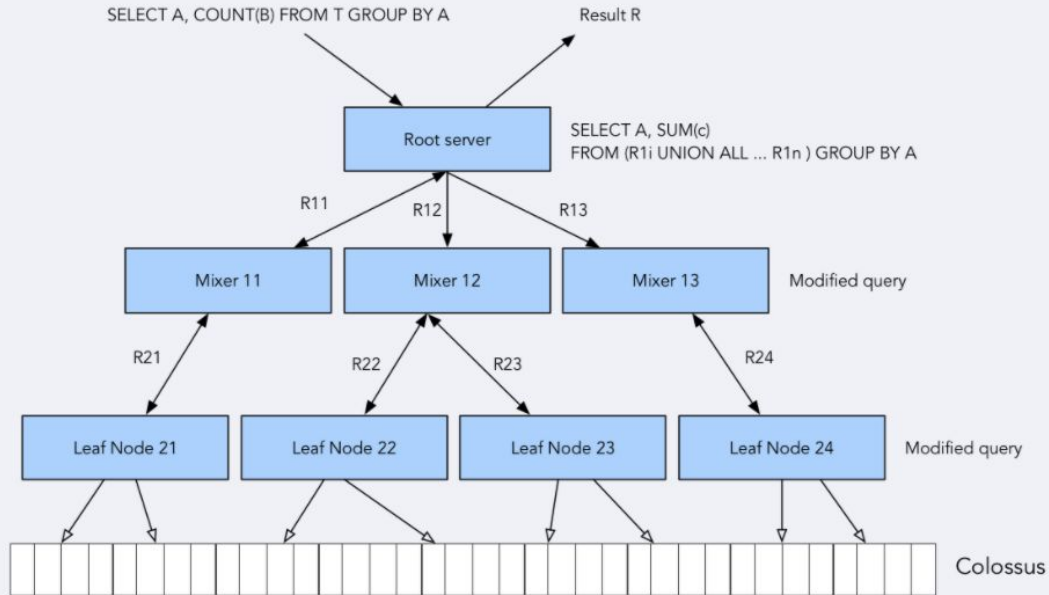


Figure-2: An example of Dremel serving tree.

Reference

<https://cloud.google.com/bigquery/docs/how-to>

<https://research.google/pubs/pub36632/>

<https://panoply.io/data-warehouse-guide/bigquery-architecture/>

http://www.goldsborough.me/distributed-systems/2019/05/18/21-09-00-a_look_at_dremel/

ML in BigQuery

- Target audience Data analysts, managers
- No need for Python or Java knowledge
- No need to export data into a different system

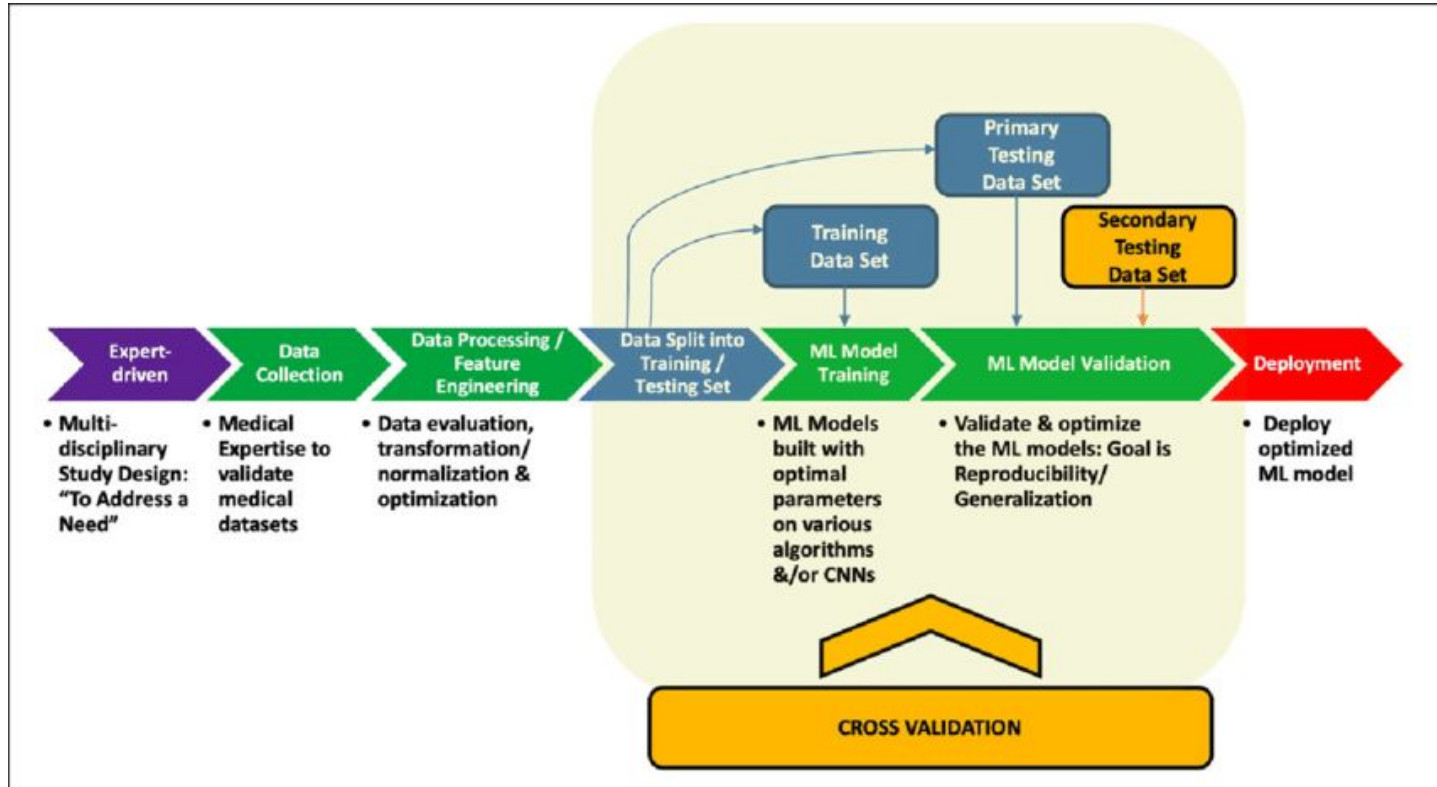
ML in BigQuery pricing

- Free
 - 10 GB per month of data storage
 - 1 TB per month of queries processed
 - ML Create model step: First 10 GB per month is free

ML in BigQuery pricing

Operation	Pricing
Logistic regression model creation ¹	\$250.00 per TB
Linear regression model creation ¹	
K-means clustering model creation ¹	
Time series model creation ^{1, 2}	
AutoML Tables model creation ¹	\$5.00 per TB, plus Vertex AI training cost
DNN model creation ¹	
Boosted tree model creation	

ML in BigQuery



ML in BigQuery

