

# **BigQuery Review Paper**

**Name:** Sharayu Rasal

## **BigQuery Paper Review: Problem Statement Tackled**

Google needed interactive, seconds-level analytics over web-scale, largely read-only datasets used across products such as Search, Gmail, YouTube, and Ads. Traditional databases and batch systems could not answer ad hoc business questions fast enough without heavy indexing or pre-aggregation. Internally, Google built Dremel to run SQL-like queries over trillions of records in seconds.

BigQuery is the externalization of Dremel, exposing similar capabilities via a managed cloud service (Web UI, CLI, REST) without infrastructure management.

## **Key Design Principles**

- Interactive speed at web scale - Big Query provides full-table scans over billions of rows in seconds or tens of seconds, without indexes or pre-aggregates.
- Columnar storage first - It only reads the needed columns, higher compression (roughly 10:1 vs ~3:1 in row stores) dramatically cuts I/O.
- Massively parallel execution- It's a distributed tree architecture that fans queries out to thousands of machines and aggregates results quickly.
- Cloud-managed, zero-ops - The capacity, replication, reliability, and APIs delivered as a service.

## **Architecture & Components**

Columnar Storage - Data is stored column-wise, so queries scan only referenced columns. This yields high compression and throughput but is optimized for read-heavy OLAP/BI (no in-place updates).

Tree Query Execution - Coordinator nodes distribute a query down a multi-level tree and leaf workers scan and partially aggregate the results, which are merged up the tree to deliver low-latency answers.

Externalization as BigQuery - The same underlying ideas are exposed via REST API, Web UI, CLI, access control, schema management, and tight integration with Google Cloud Storage.

## **Query Model & Capabilities**

BigQuery supports SQL-like ad hoc queries over structured schemas, ideal for diagnostics, quick roll-ups, and exploratory trial-and-error analytics. Example given in

the paper: a regex count over ~314M Wikipedia rows (35.7 GB) answered in ~10 seconds with no index.

OLAP - Excellent for aggregations across dimensions (time, region, product, etc.) without building cubes or indexes.

## **Operational Model (Why Cloud Matters)**

As a fully managed service, BigQuery requires no provisioning, replication management, patching, or 24×7 operations. APIs enable dashboards or mobile front-ends. Lower total cost of ownership lets teams focus on analysis rather than infrastructure.

## **Real-World Uses (Inside Google)**

Dremel underpins spam analysis, crash reports, OCR (Books), Android install tracking, Maps debugging, Bigtable tablet migrations, build/test results, data-center resource monitoring, and more illustrating broad utility at "Google speed."

## **Limitations & Scope**

Optimized for read-heavy analytics; no in-place updates to existing records. Best with structured schemas, unstructured or complex iterative ML is often better suited to batch frameworks (then analyzed in BigQuery).

## **Related Work**

BigQuery relates to Google's broader ecosystem: MapReduce (batch processing), Bigtable (storage), and the Dremel research paper on interactive analysis of web-scale datasets. It contrasts with traditional ROLAP/MOLAP by emphasizing scan speed over advance indexing.

## **Conclusion / Summary**

BigQuery is the cloud-powered externalization of Dremel that brings interactive, SQL-based analytics over massive datasets to everyone without managing clusters, indexes, or cubes. It excels at ad hoc OLAP/BI and complements batch frameworks for heavy ETL and complex data mining. By combining columnar storage and tree-based parallel execution on Google's infrastructure, it delivers web-scale performance and strong cost effectiveness enabling teams to analyze data at "Google speed"