



Cloud Platform Comparison

AWS Athena, Google BigQuery & Snowflake

Outline

- Our Progress
- Dataset Overview
- Benchmark Query Results
- Intro to Google BigQuery
- Intro to Snowflake
- Query Comparison
- Intro to AWS Glue
- Pricing

Our Progress

- Loaded ~ 2.5G TB dataset onto AWS, Google and Snowflake platforms
- Ran benchmark queries on three platforms
- Researched and compared additional features
- Experimented with AWS Glue for Parquet conversion

Dataset Overview

The GDELT Project

- Global Data on Events, Location and Tone
- Data on the world's broadcast, print, and web news
- Database is updated every 15 mins

Dataset Overview

The GDELT Project

Tables

- Events (74.1GB) : over 300 categories of physical activities around the world
- Mentions (93GB): every mention of an event over time
- Global Knowledge Graph (GKG) (2.38TB): every person, organization, company, location and several million themes and thousands of emotions from every news report

Queries

We ran a total of 8 queries that involved:

- Table Joins
- Recasting data types and aggregating
- String Matching
- Common Table Expressions ('with' statements)

Benchmark Query Results

Query Speed

Snowflake > BQ Native > Athena Column > BQ External > Athena Row

Data Scanned

Snowflake < Athena Column < BQ Native < BQ External = Athena Row

Cost

Snowflake < Athena Column < BQ Native < BQ External = Athena Row

Remarks

For Row-Based Data

- BQ is generally faster than Athena, but Athena sometimes outperforms
- Athena seems to run faster on all strings dataset, even if query requires recasting data types

For Columnar Data

- BQ and Athena are comparable, but Athena's data is smaller and hence queries cheaper

Remarks

Snowflake

- Partitions, compresses and columnarly stores its data. We could not use it to query data in “raw” format.
- Prices queries based on time with a 1-minute per query minimum. This is competitive for longer-running queries but could be costlier for quicker ones.

Exhausted Resources Error

- Query that produced “exhausted resources error” on Athena (both formats) ran successfully on BigQuery and Snowflake in minutes

Introduction to Google BigQuery

- Serverless: pay for storage + data scanned per query
- Supports data load from cloud storage or local drive
- Parquet files not supported in current version
- Can query external tables (raw data format) or native tables (columnarly stored in BQ)
- Stricter error handling
- User-friendly, detailed error messaging

BigQuery Demo

Google BigQuery

Table creation interface

COMPOSE QUERY

Query History

Job History

Filter by ID or label

gsb-circlerss

GDELT

actortype

actortype_native

eventcodes

eventcodes_native

events_native

gkg_native

mentions_native

v2events

v2gkg

v2mentions

bigquery-public-data

nyc-tlc:yellow

Create Table

Source Data

Create from source

Create empty table

Repeat job

Select Previous Job

Location

Google Cloud Storage

gs://rss_gdelt_data/events/20150218*

File format

CSV

View Files

Destination Table

Table name

GDELT . events_20150218

Table type

Native table

Schema

Automatically detect

Name	Type	Mode
	STRING	NULLABLE

Add Field

Edit as Text

Options

Field delimiter

Comma

Tab

Pipe

Other

Header rows to skip

0

Number of errors allowed

0

Allow quoted newlines

Allow jagged rows

Ignore unknown values

Write preference

Write if empty

Partitioning

None

Create Table

Displays existing dataset and tables

BigQuery Demo

Create Table

Source Data ☒ Create from source ☐ Create empty table

Repeat job

Select Previous Job



Location points to files in google cloud storage, can use wildcard.

Location

Google Cloud Storage

gs://rss_gdelt_data/events/20150218*



File format

CSV

[View Files](#)

Destination Table

Table name

GDEI . events_20150218



Table type

Native table



Native table

External table

With the option to load as an external table or native table

Schema

☐ Automatically detect



Schema auto-detect feature available

Name

Type

Mode

STRING

NULLABLE



Add Field

[Edit as Text](#)

BigQuery Demo

Options

Field delimiter

☐ Comma ☒ Tab ☐ Pipe ☐ Other ?

Header rows to skip

?

Number of errors allowed

? →

The maximum number of bad records that can be ignored before aborting the job.

Allow quoted newlines

☒ ?

Allow jagged rows

☒ ?

Ignore unknown values

☒ ?


Write preference

Write if empty ?

Partitioning

None

BigQuery Demo

 **Load** gs://rss_gdelt_data/errorfile/20170918204500.gkg.csv to gsb-circlerss:GDELT.gkg_error

Repeat Load Job

Errors:

gs://rss_gdelt_data/errorfile/20170918204500.gkg.csv: CSV table encountered too many errors, giving up. Rows: 517; errors: 1. (error code: [invalid](#))

gs://rss_gdelt_data/errorfile/20170918204500.gkg.csv: Too many values in row starting at position: 6094997. (error code: [invalid](#))

Job ID gsb-circlerss:bquijob_d825234_15fc1ebf095

Creation Time Nov 15, 2017, 3:03:06 PM

Start Time Nov 15, 2017, 3:03:10 PM

End Time Nov 15, 2017, 3:03:17 PM

Destination Table gsb-circlerss:GDELT.gkg_error

Write Preference Write if empty

Source Format CSV

Source URI gs://rss_gdelt_data/errorfile/20170918204500.gkg.csv ([Open in GCS](#))

Autodetect Schema true

Repeat Load Job

Cancel Job

Gives details about the type of error and in which file it is encountered.

BigQuery Demo

New Query ?

Query Editor

UDF Editor



SQL

```
1 SELECT
2   *
3 FROM
4   `gsb-circlerss.GDELT.events_native`
5 WHERE
6   MonthYear = '201710';
```

Valid: This query will process 74.1 GB when run.

Destination Table

Select Table

gsb-circlerss:GDELT.events_1710 X

Write Preference

☒ Write if empty ☐ Append to table ☐ Overwrite table

Results Size

☒ Allow Large Results ?

Results Schema

☒ Flatten Results ?

Query Caching

☐ Use Cached Results ?

Query Priority

☒ Interactive ☐ Batch ?

UDF Source URIs

Edit ?

Maximum Bytes Billed

Project Default ?

SQL Dialect

☐ Use Legacy SQL ?

RUN QUERY

Save Query

Save View

Format Query

Select a destination table to save the query results

Default is Legacy SQL.
Need to uncheck the box to use standard SQL



BigQuery Demo

✓ SELECT * FROM `gsb-circlerss.GDELT.events_native` WHERE MonthYear = '201710';

```
1 SELECT
2 *
3 FROM
4 `gsb-circlerss.GDELT.events_native`
5 WHERE
6   MonthYear = '201710';
```

Job ID	gsb-circlerss:bquijob_2cc3952a_15ffeed95cf
Creation Time	Nov 27, 2017, 11:21:45 AM
Start Time	Nov 27, 2017, 11:21:45 AM
End Time	Nov 27, 2017, 11:22:16 AM
Bytes Processed	74.1 GB
Bytes Billed	74.1 GB
Destination Table	gsb-circlerss:GDELT.events_1710
Allow Large Results	true
Use Legacy SQL	false

Open Query

If destination table is specified, the table will appear in the existing table list after query completion.

gsb-circlerss

▼ GDELT

- actortype
- actortype_native
- eventcodes
- eventcodes_native
- events_1710**
- events_native
- gkg_native
- mentions_native
- v2events
- v2gkg
- v2mentions

BigQuery Demo

Query results saved in a native table can be exported as files.

Table Details: events_1710

Refresh

Query Table

Copy Table

Export Table

Delete Table

Schema

Details

Preview

Row	GLOBALEVENTID	Day	MonthYear	Year	FractionDate	Actor1Code	Actor1Name	Actor1CountryC
1	699036722	20171018	201710	2017	2017.7890	MED	GAZETTE	null
2	698807030	20171018	201710	2017	2017.7890	CAN	ONTARIO	CAN
3	698782867	20171018	201710	2017	2017.7890	USACOP	SEATTLE	USA
4	698215539	20171016	201710	2017	2017.7836	EUR	EUROPEAN	EUR
5	698177739	20171016	201710	2017	2017.7836	LAB	WORKER	null
6	698111132	20171015	201710	2017	2017.7808	UAE	GLIMMER	null

Export to Google Cloud Storage

Export format

CSV

Compression

☒ None ☐ GZIP ?

Google Cloud Storage URI

gs://

[View Files](#)

OK

Cancel

Google Cloud Storage URIs begin with "gs://" and specify the bucket and object you wish to import.

Example:

gs://mybucket/path/to/mydata.csv

Maximum table export size is 1GB. For exports greater than 1GB, specify file pattern with a *

Example:

*gs://mybucket/path/to/mylargedata**

BigQuery Demo

New Query ?

Query Editor

UDF Editor



```
1 SELECT
2   AVG(CAST(AvgTone AS FLOAT64))
3 FROM
4   `gsb-circlerss.GDELT.events_1710`
5 GROUP BY
6   Day
```

SQL

Valid: This query will process 140 MB when run.

Destination Table

Select Table

No table selected

Write Preference

☒ Write if empty ☐ Append to table ☐ Overwrite

Results Size

☐ Allow Large Results ?

Results Schema

☒ Flatten Results ?

Query Caching

☒ Use Cached Results ?

Query Priority

☒ Interactive ☐ Batch ?

UDF Source URIs

Edit ?

Maximum Bytes Billed

Project Default ?

SQL Dialect

☐ Use Legacy SQL ?

If no destination table is selected, results will be saved to a temporary table and cached for ~ 24hrs

RUN QUERY

Save Query

Save View

Format Query

Hide Options



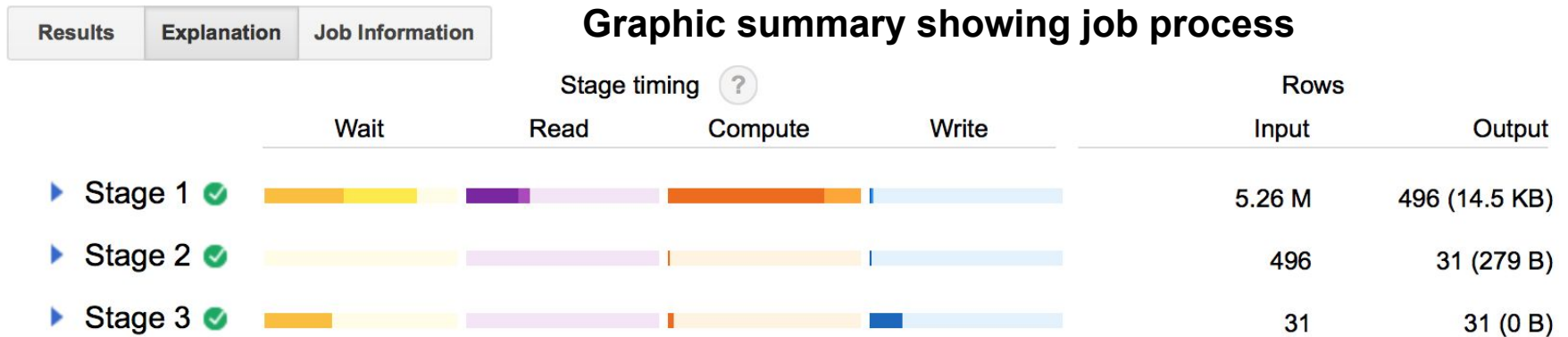
BigQuery Demo

Download option only
available in Chrome

Results	Explanation	Job Information
Download as CSV	Download as JSON	Save as Table
Save to Google Sheets		

Row	f0_
1	-2.026962960405946
2	-2.087468767197474
3	-1.9637318551111413

Graphic summary showing job process



Intro to Snowflake

- Pay for storage + time per query
- Serverless insofar as warehouse can be automatically suspended and resumed just prior to querying
- Local files must be staged on either S3 or Snowflake stage
- Load from S3 is fast since within AWS environment
- Auto-scalability & elasticity: can change warehouse size to improve concurrency and/or processing speed
- Data automatically partitioned and columnar compressed
- Meta-data is extracted to enable efficient query processing

High-Level Comparison

Athena

- Query interface only
- Can read Parquet files
- Manual columnar conversion
- Loose error handling on data load

BigQuery

- Query interface or data warehouse
- Cannot read Parquet files
- Automatic columnar conversion
- Stricter error handling, can be overridden

Snowflake

- Data warehouse
- Can read Parquet files
- Automatic columnar conversion
- Stricter error handling, can be overridden

High-Level Comparison

Athena

- Results automatically saved to S3
- Need to repoint to query output

BigQuery

- Results can be displayed/ cached or written to native table
- Detailed error messaging that indicates file at fault

Snowflake

- Query results can be displayed/ cached or written to table

- See [here](#) for more details

Query Comparison

Select all records where theme includes 'terror'

	Athena Row	Athena Columnar	BQ External	BQ Native	Snowflake
Runtime	82 min	71.98 min	74.6 min	59.9 min	35.22 min
Data Scanned	2.38 TB	2.37 TB	2.38 TB	2.38 TB	924.4 GB
Cost (\$)	11.9	11.85	11.9	11.9	4.696

Query Comparison

Select all events where the confidence score ≥ 100

	Athena Row	Athena Columnar	BQ External	BQ Native	Snowflake
Runtime	82 min	7.88 min	74.6 min	5.1 min	1.32 min
Data Scanned	158 GB	36.06 GB	158 GB	80.6 GB	14.9 GB
Cost (\$)	0.79	0.1803	0.79	0.403	0.352

Query Comparison

Count events per year per actor label

	Athena Row	Athena Columnar	BQ External	BQ Native	Snowflake
Runtime	122.81s	11.5 s	126.8s	3.2s	1.7s
Data Scanned	70 GB	40.92 MB	70 GB	1.45 GB	0.16 GB
Cost (\$)	0.35	0.0002	0.35	0.0072	0.132

Query Comparison

Number of rows per combination of 27 column & join of three tables
(Exhausted Resources Error query)

	Athena Row	Athena Columnar	BQ External	BQ Native	Snowflake
Runtime	Failed	Failed	7.05 mins	4 mins	3.88 mins
Data Scanned	-	-	158 GB	20.3 GB	4.2 GB
Cost (\$)	-	-	0.79	0.1015	0.516

- See [here](#) for more details

Cost of Parquet Conversion with EMR

	Data Size	Job Time	Master & Core Nodes	Task Nodes	Price
Events	70 GB	46 mins	1 x m3.xlarge 3 x m3.xlarge	None	\$0.9774
Mentions	93 GB	42 mins	1 x m3.xlarge 3 x m3.xlarge	None	\$0.8925
GKG	2.4 TB	3 hrs 2 mins	1 x m3.xlarge 4 x m3.2xlarge	5 x m3.2xlarge (spot price)	\$13.13

m3.xlarge	On demand	0.266/hr
m3.2xlarge	On demand	0.532/hr
m3.2xlarge	Spot	0.15/hr

Cost Estimate for Data Transfer from AWS

10TBs; ~4000 files (assume 1 file = 1 request)

To BigQuery

- Data Transfer: \$0.02
- Storage Cost: \$200/month

To Snowflake

- Data Transfer: \$0
- Storage Cost:
 - \$400/month on-demand
 - \$230/month pre-pay

Intro to AWS Glue

Data Catalog

- A persistent metadata store (can be used with Athena, Redshift Spectrum and Redshift)
- Uses a crawler to scan a data store and automatically detect schemas or manually input/update schemas
- Can schedule when the crawler is to run and update existing schemas if changes detected

Intro to AWS Glue

ETL (Extract, Transform, Load) - Serverless

- Automatically generates code to extract, transform and load data based on a specified output format (e.g. convert csv to parquet format)
- The code is generated in Python and written for the Apache Spark 2.1 environment
- Issue: failed on converting 70GB dataset
- Reason: cannot overwrite Spark configuration, specifically, 'driver.maxResultSize' parameter.

Pricing - AWS Glue

ETL jobs and Crawler: charged by runtime

- \$0.44 per DPU-Hour, billed per second, with a 10-minute minimum for each ETL job
- Default is 10 DPUs assigned per ETL job; 2 DPUs minimum
- A single DPU: 4 vCPUs compute and 16 GB of memory.

Data Catalog: charged by unit storage and access request

- **Storage:** free for first million objects; \$1/month/100,000 objects stored above 1M
- **Requests:** Free for the first million requests per month; \$1/month/million requests above 1M

Pricing - Athena

Storage

- \$0.023/GB/mo

Compute

- \$5/TB with a 10MB-per-query min

Pricing - BigQuery

Storage

- **External Table**
 - Multi-regional bucket: \$0.026/GB/mo
- **Native Table:**
 - Long-term (table not edited for 90 consecutive days): \$0.01/GB/mo
 - Short-term (eg. growing table): \$0.02/GB/mo

Compute

- \$5/TB with a 10MB-per-query minimum

Pricing - SnowFlake

- Minimum \$25 per month in storage or compute

Storage

- On-Demand: \$40/TB/mo
- Pre-Pay: \$23/TB/mo (negotiable)

Compute

- \$2/credit/hour; queries billed per second with a 1min-per-query minimum
- Number of credits per hour depends on warehouse size:

X-Small	Small	Medium	Large	X-Large	2X-Large	3X-Large	4X-Large
1	2	4	8	16	32	64	128

Pricing - Data Transfers

Google to AWS

- Transferring data from Google to AWS was treated as DOWNLOAD, which also incurs network charges (\$0.02/GB)

AWS to Google

- Transferring data from AWS to Google was treated as COPY request on AWS (\$0.005 per 1,000 requests) - Per Documentation

AWS to Snowflake

- None noticed

Next Steps

Expanding Pipeline

- E.g. incorporating statistical analyses

Additional Features

- Concurrency
- Security
- Backup and Recovery Options

Additional Platforms

- Microsoft Azure
- Amazon Redshift + Redshift Spectrum



change lives. change organizations. change the world.