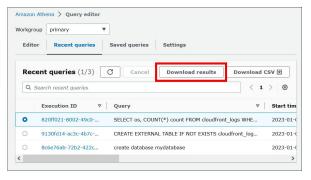
Amazon Athena

William Hay November 25th, 2025

Amazon Athena

SELECT os, COUNT(*) count
FROM cloudfront_logs
WHERE date BETWEEN date '2014-07-05' AND date '2014-08-05'
GROUP BY os



What is it?

- Serverless no infrastructure to manage
- SQL-based interactive query service for Amazon S3
- Uses Facebook Presto for the SQL query engine

What can you do?

- Using S3 and standard SQL, Athena helps analyze unstructured, semi-structured, and structured data that is stored in an S3 bucket
 - Unstructured data still needs to be defined through a schema
- Integration with Amazon QuickSight, Glue Data Catalog, Redshift, EMR
- Download recent queries for further analysis
 - Recent queries are retained for 45 days

History - Who

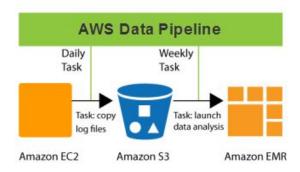
General availability announced in 2016 at the AWS re:Invent event

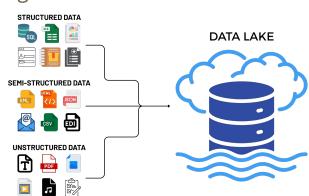
Technology	Туре	Storage	Pricing	Performance	Best For
Amazon Athena	Serverless SQL engine	S3	Per TB scanned	Varies depending on file format	Ad-hoc queries over S3 data lakes
Amazon Redshift	Warehouse	Internal + External	On-demand cluster	High, steady	High-volume analytics
Google BigQuery	Serverless warehouse	Internal + External	Per TB scanned	Very high, consistent	Large-scale data analytics
Snowflake	Warehouse	Internal + External	Compute + storage	Very high	Recurring analytics
Azure Synapse Analytics	Serverless "SQL pool"	Azure Data Lake Storage	Per TB scanned	High	Azure-native workloads

3

History - Why

- Massive growth of data stored in S3
- Need to analyze data without ETL or loading into databases
 - Before Athena, organizations analyzed large datasets by loading data into a database or build and maintain Hadoop or Spark clusters, requiring large amounts of maintenance
- Pay-as-you-go pricing, no idle computing costs
- Makes querying data lakes very easy





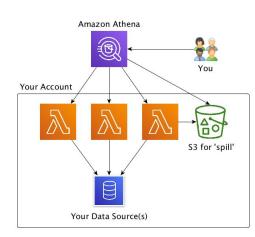
-

History - How

Since its release in 2016...

- Added support for columnar formats (e.g. Parquet)
- Integrated with AWS Glue Data Catalog
- Partition projection
 - Athena will understand S3 partitions logically
- Current engine is Athena engine version 3
- Support added for federated queries across multiple AWS data sources

These changes were driven by a shift towards serverless analytics and increasing size of data lakes in Amazon S3.

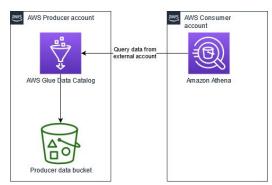


5

Features - 1



- Serverless SQL query engine
- Queries data directly stored in Amazon S3
- Based on Facebook Presto distributed compute engine
- Supports structured and semi-structured formats
 - o CSV, JSON, Parquet, and more...
- Integrates with multiple Amazon services (AWS Glue Data Catalog)
- Can query petabyte-scale datasets without provisioning any clusters
- Supports ANSI SQL
 - o Window functions, complex joins, database and table creation, etc.
- Results can be exported to either S3 or into an ETL pipeline



- Integrates with S3 lifecycle tiers for cost-optimized storage
- Ability to save queries
- Automatically executes queries in parallel
- Supports federated queries
- Invoke SageMaker machine learning models to use an Athena SQL query to run inference
 - o Can use SQL experience to train ML models
- Ability to query from Azure Synapse Analytics data and visualize with Amazon Quicksight



7

Use Cases

- Ad-hoc querying for analysis
- Build serverless data lake analysis
- Query information to display visually (e.g. using Amazon QuickSight)

Query parratives data

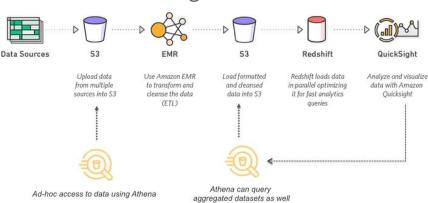
Calculate distance sparse matrix
Run clustering algorithm to detect clusters

· Compare between experiments

· Enrich detected botnets with other data

· Store data as compressed CSV objects in Amazon S3

Machine learning botnet detection





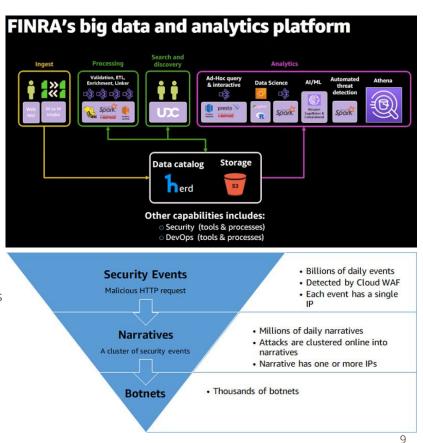
/

18/

8

Industry Deployments

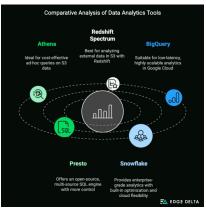
- FINRA
 - Processes 425+ billion market events each day (~10-20 TB of data)
 - Regulates 3,700 security firms and 630k+ brokers
 - Uses Athena for real time analysis
 - Athena eliminated problems with configuring and maintaining physical servers
- Imperva
 - Protects hundreds of thousands of websites
 - Blocks billions of security events each day
 - Uses Athena, SageMaker, and QuickSight to develop a ML clustering algorithm
 - Detects botnets, which contribute to attacks such as DDoS



Advantages

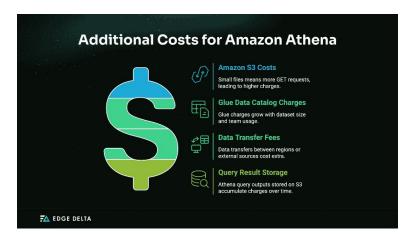
- Serverless
- Build interactive and advanced analytics applications
- Build with preferred tools and frameworks
- Simple and predictive pricing
- Pay only for data scanned
- Highly scalable
- Ideal for intermittent workloads
- Reduces operational costs and ETL costs





Disadvantages

- Performance varies based on data layout within S3
- Can be costly if file compression is not used
- Not ideal for transactional workloads
- Concurrency limits
- Latency has been observed to be unpredictable
- Limited in terms of tuning compared to other options







11

Usability

- If you have SQL knowledge, the hardest part is the setup
 - Always double check to make sure Athena settings are configured correctly for whatever data formats you will be using
- To get Athena working at a minimum, an individual can set it up in about 5 minutes without much knowledge about Athena
- Lots of documentation, tutorials, and settings to mess around with

- ▶ What is Amazon Athena?
- ▼ Use Athena SQL

Tables, databases, and data catalogs

- Get started
- ► Connect to data sources
- Connect to Amazon Athena with ODBC and JDBC drivers
- Create databases and tables
- Create a table from query results (CTAS)
- Use SerDes
- Run queries
- Use ACID transactions
- Security
- ▶ Workload management
- ▶ Athena engine versioning
- SQL reference for Athena
- ▶ Troubleshoot issues
- Code samples
- Use Apache Spark
 Release notes

Document history

Cost

- Athena charges \$5 per TB of data scanned
- Save up to 90% per query and get better performance by compressing, partitioning, and converting data into columnar formats
- When using SQL queries on data sources, you are charged for the number of bytes per scanned query, rounded to the nearest megabyte
 - o Minimum 10 MB per query unless Provisioned Capacity is used
 - o Data definition language (DDL) statements have no charges
 - CREATE, ALTER, or DROP TABLE
- With 100k queries a day, with a minimum of 10MB: \$145.04 monthly
- With 1m queries a day, with a minimum of 100MB: \$14,503.80 monthly

13

Cost Example

FINRA:

Athena pricing Pay only for the queries you run vs. cost of maintaining long running cluster Athena charges \$5 per TB of data scanned. For \$5 we can execute: 100,000 queries with an average of 10MB data scanning 10,000 queries with an average of 100MB data scanning 1,000 queries with an average of 1,000MB (1GB) data scanning select rec_value from cATFOLA where trade_dt= date '2020-11-11' and cat_lifecycle_id = 1234567890 10.89 MB Datascan

Athena Benchmarks from Upsolver:

Athena vs BigQuery Benchmarks

	Google BigQuery	Amazon Athena O upsolver	Amazon Athena 5-minute Parquet	Amazon Athena 1-minute Parquet
Query Time	57.72	60.98	113.57	292.58
(Seconds)		(+5.6%)	(+96.7%)	(+506%)
Total Cost	1.97	0.42	0.47	0.49
(USD)		(-78.7%)	(-76.2%)	(-75.13%)

* Aggregated results for 9 SQL queries against NY Taxi Rides dataset

Amazon Athena Conclusion

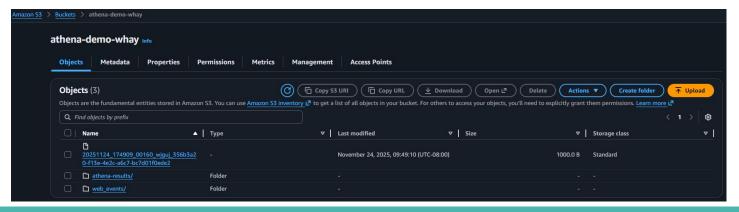


- Serverless SQL query service for S3 analytics or ML inference
- Integrates with multiple AWS services
 - o AWS Glue, CloudFormation, CloudTrail, EMR, S3, IAM
- Ideal for ad-hoc queries and scales to analyze petabyte-scale data
- Best cost-effective method is to optimize and compress data
- Automatically scales and completes queries in parallel
- Built to support standard SQL and federated queries
- Save queries and download results

Demo

To start, we create an S3 bucket to point Amazon Athena to our input/output source

Make sure that the user you are executing queries with has proper permissions (s3:GetObject & s3:ListBucket)



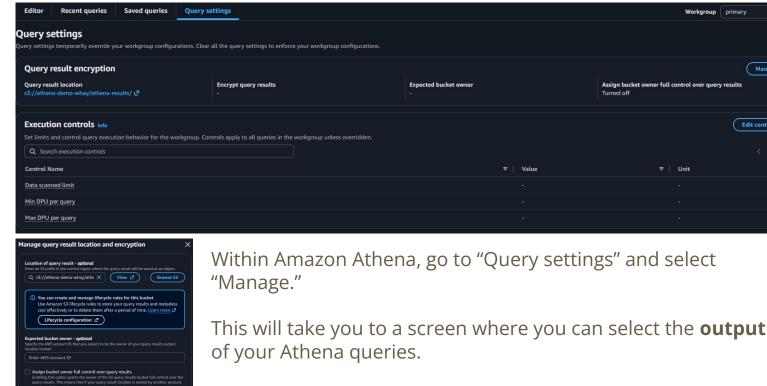
15



18

Workgroup primary

Manage



reate table from S3 bucket data Table details Query in Amazon SageMaker Unified Studio Use your IAM role to analyze and build with your existing AWS ro Saved queries Editor **Recent queries Query settings** Database configuration Info Query 1 : X Query Data C < CREATE TABLE dem Data source WITH (AwsDataCatalog format = 'PARQ external_locat Catalog X View L² Browse S3 None FROM web_events Database WHERE session du demo_db **Tables and views** Create ▲) 🕸 Create a table from data source Q Filter tables and view 53 bucket data Tables (2) Remove AWS Glue Crawler ∠7 Views (0) Create with SQL CREATE TABLE CREATE TABLE AS SELECT Select an S3 bucket as the **input** for your

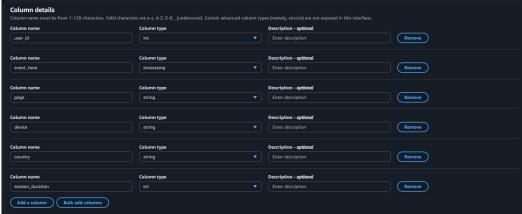
Athena queries.

CREATE TABLE AS SELECT(ICEBERG)

CREATE VIEW

Other than selecting the S3 bucket, no other options are used.

user_id	event_time	page	device	country	session_duration
1	2025-11-20T10:01:00Z	/home	desktop	US	120
2	2025-11-20T10:02:15Z	/products	mobile	CA	45
3	2025-11-20T10:03:10Z	/cart	mobile	US	300
4	2025-11-20T10:05:27Z	/home	desktop	GB	60
5	2025-11-20T10:06:55Z	/checkout	tablet	AU	210
6	2025-11-20T10:08:33Z	/products	mobile	US	95
7	2025-11-20T10:10:42Z	/home	desktop	US	180
8	2025-11-20T10:12:01Z	/checkout	mobile	IN	340
9	2025-11-20T10:14:29Z	/home	tablet	CA	75
10	2025-11-20T10:17:58Z	/cart	mobile	US	420

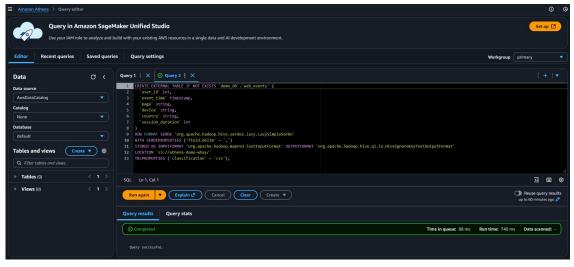


Make sure to insert your test file into your input S3 bucket!

Athena does **not** automatically infer a schema unless you use something like AWS Glue.

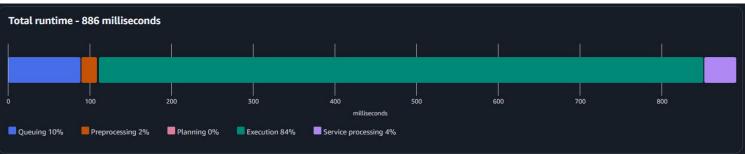
Under "Column details" after selecting our input, we can specify a schema that we want Athena to use.

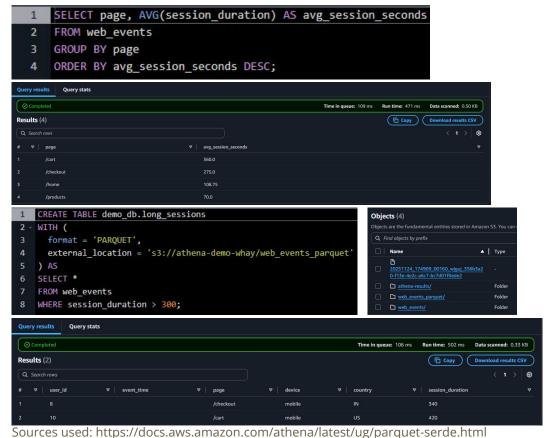
Al was used to generate the test file.



Once you've selected "Create table," you should see an automatically generated query to create the table you specified.

You can also see query stats!





Lets try some queries!

You can create multiple query tabs to save and rerun any queries

Athena can also write files with a specific format, capturing any query and storing it in S3 using that file format

As Athena also supports other formats than CSV, you can also run queries on Parquet files

21

Questions?