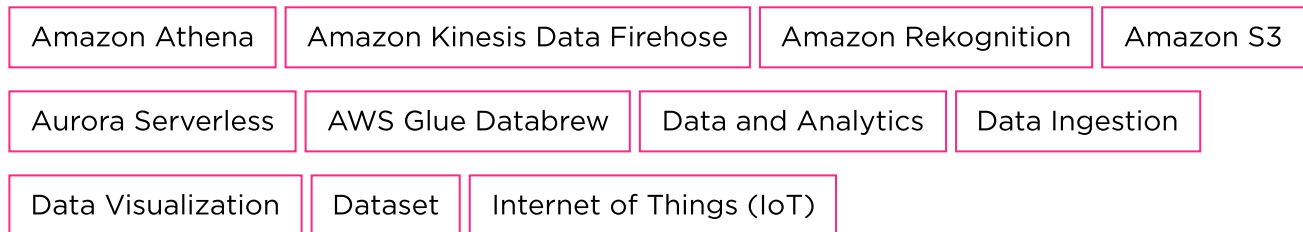
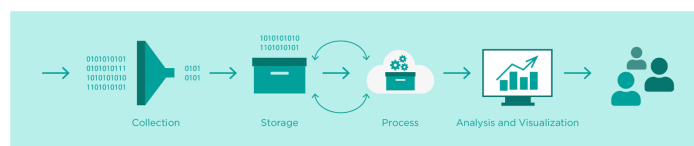


USING AWS TO INGEST AND ANALYZE DATA FROM AN IOT DEVICE: A SIMPLE EXAMPLE WITH AURORA AND ATHENA



beSharp | 11 December 2020

With the Internet of Things quickly becoming a thing of the present (rather of the future...) the number of devices sending collected on the field is increasing exponentially and so does the amount of data, thus data ingestion and analysis has become of the hottest topics of the current IT landscape. AWS offers a wide range of services which allow us to ingest, collect, store, analyze and visualize huge amounts of data quickly and efficiently.



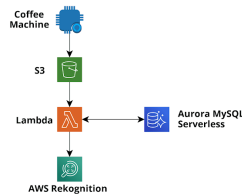
In this brief article, we would like to present a very simple real-world application we developed as a proof of concept demonstration to show the data ingestion and analysis pipeline in AWS and IoT events and conferences.

We customized an existing Nespresso coffee machine to take photos of people making coffees using custom electronics, a Raspberry Pi Zero and a micro camera. The image is immediately uploaded to S3 and an AWS Lambda triggered by the upload analyzes the image using Amazon Rekognition. After the analysis of the image is complete, if the image contains the face of a person, a record is written by the Lambda function in an Aurora MySQL serverless together with metadata output from the Rekognition ML algorithm: does the person have eyeglass? beard? mustaches? is she/he smiling? Finally, a very simple web application was developed and connected to the database to show statistics.

Furthermore, an AWS Athena query cleans the data and moves them to a new S3 bucket as parquet files.

Obviously for our trivial application many of these steps are redundant but they aim to demonstrate the power of AWS building blocks in creating very complex data pipelines.

A scheme of the proposed infrastructure is shown below.



Hereafter we describe all the steps of a common data ingestion and transformation and how we are doing them in our trivial application. Let's dive deep!

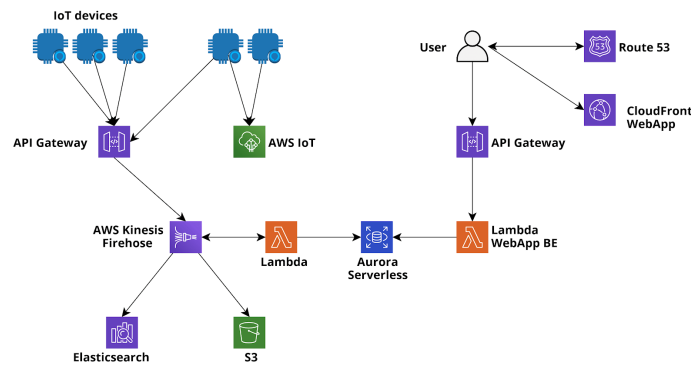
The Ingestion/Storage Step

A very common AWS data ingestion flow is to use AWS IoT Core (Secure MQTT) or Api Gateway (REST APIs or Websocket) as the data entry point, directly connect the data entry point to Kinesis Firehose (using IoT rules or Api gateway Service integrations) and finally leverage the powerful Firehose features for data buffering, buffer transformation (AWS Lambda functions), stream encryption (AWS KMS), data compression (GZIP) and data delivery of compressed and automatically encrypted message batches to both long term object storage (S3) and/or to a data warehouse (AWS Redshift) for complex analytical queries on the huge amount of data collected.

Always having all the ingested data saved in AWS S3 is an essential step, not only as a lifesaver in case of problems with other hotter storages but also to create a shared data lake which can be later analysed with Athena EMR, Glue Jobs, Glue databrew and also external tools.

Furthermore you can use Firehose to directly deliver data to AWS Elasticsearch for real time analysis and if needed it is also very easy to deliver the batches of ingested data to a relational database (e.g. Aurora Serverless Postgres/MySQL) using either AWS Data Migration tasks or event based Lambda functions. Migrating inserted data (or an aggregation of inserted data) to an existing relational database is often quite useful if you need them to enrich an existing legacy application already using the database.

If you decide to use Lambda functions to move the ingested data to Aurora, which is usually faster and more scalable, you can either use the Firehose transform functions directly or a different function triggered each time Firehose writes an object to S3.



The beauty of Firehose is that you can also add it as a subsequent step! In our simple coffee application we are not using it and images and analyses are saved directly in S3 and AWS Aurora Serverless MySQL by Lambda Functions. Anyway, if the app grows bigger we can integrate it flawlessly!

Analysis Step

Once your data is in a storage, it is time to analyse them. Methodologies can differ greatly and common examples range from simple queries run in relational databases to long and complex analytical jobs run in Redshift data warehouses and to near real time processing using Kinesis connected EMR or ElasticSearch.

In our case we can just run simple queries using our web application backend and display the results in the browser.

However in the future we may be interested in running much more advanced queries on our data and maybe do some data quality inspection and machine learning training. So we need to have these data out of Aurora and into S3 in order to analyze them with Glue jobs and Databrew and if needed to load them easily with Apache Spark either from Glue or AWS EMR. To do this, we can follow several paths: for example we could use AWS DataMigration service to move the data to S3 as Parquet files or maybe we could create a Glue Job, load the data using Glue Connection to RDS and Spark and then write them into S3.

After this would need to run a Glue crawler in order to create DataCatalog that will be used by Athena and Glue for queries and jobs.

Here however we will show you a different and sometimes **much flexible path** to export, clean and catalogue our data from a relational database: **Athena custom data source**.

By default, Athena comes with S3 - Glue data Catalog integrations but AWS recently added the possibility to add customized data sources such as JDBC connected databases, AWS CloudWatch or to query S3 but using a custom Apache Hive metastore. In our case we are interested in connecting to MySQL Aurora Serverless so we need to go to Athena Home, configure a workgroup named AmazonAthenaPreviewFunctionality and then add an S3 query output path:

Athena Query editor Saved queries History Data sources Workgroup: AmazonAthena...			
Workgroups Use workgroups to separate users, teams, applications, or workloads, and to set limits on amount of data each query or the entire workgroup can process. You can also view query-related metrics in AWS CloudWatch. Learn more			
Create workgroup View details Switch workgroup			
Name	Description	Creation time	Workgroup status
awsathenaengineconnector	This functionality is a part of the Athena Preview features and should be run in the workgroup named awsathenaengineconnector.	2020/12/11 11:47:02 UTC+1	Enabled
primary		2020/03/09 10:54:02 UTC+1	Enabled

After this, we can go back to athena home and select Connect Data Source:

AthenaQuery editorSaved queriesHistory

Data sourceConnect data source

AwsDataCatalog

We are presented with a web page where we need to select the type of data source: we go for Query a data source (beta) ad MySQL:

AthenaQuery editorSaved queriesHistoryData sourcesWorkgroup: AmazonAthena...

Connect data source

Step 1: Choose a data source

Step 2: Connection details

Choose where your data is located

Athena queries data where it is. Data is not loaded or moved. [Learn more](#)

Query data in Amazon S3

Choose an external catalog.

Query a data source (beta)

Configure a connector for common data sources.

Choose a data source (beta)

Choose the data source to query with Athena. After you choose a data source, you will configure a Lambda function to handle the connection. [Learn more](#)

Amazon CloudWatch Logs

Amazon DocumentDB

Amazon Redshift

MySQL

Amazon CloudWatch Metrics

Amazon DynamoDB

Apache HBase

PostgreSQL

After that you are requested to enter the name and description of the new catalog and to select or create a Lambda Function to manage the connection. Choose the name you like the most and click Configure new AWS Lambda Function.

Connection details: MySQL

Choose a Lambda function that is configured to connect to your data source, or create and configure a Lambda function to handle the connection. [Learn more](#)

Lambda function

Choose or configure a new AWS Lambda function to connect to the data source.

Choose Lambda function

Configure new AWS Lambda function

Catalog name

Create a unique name to specify this data source within a SQL statement.

iotarticolo

Use up to 127 characters and must be unique within your account. It cannot be changed after creation. Valid characters are a-z, A-Z, 0-9, _ (underscores), @ (at) and - (hyphen).

Description

iotarticolo catalog

Use up to 1024 characters.

Cancel

Previous

Connect

You are presented with this page where you need to enter the JDBC connection uri for Aurora and select the subnet and security group for the Lambda function that athena will use to establish the JDBC connection. **Choose them wisely** otherwise the Lambda won't reach the Aurora instance!

Application settings

Application name

The stack name of this application created via AWS CloudFormation

AthenaJdbcConnector

SecretNamePrefix

Used to create resource-based authorization policy for "secretsmanager:GetSecretValue" action. E.g. All Athena JDBC Federation secret names can be prefixed with "athena-jdbc-federation" and authorization policy will allow "arn:aws:secretsmanager:\${AWS::Region}:\${AWS::AccountId}:secret:athena-jdbc-federation". Parameter value in this case should be "athena-jdbc-federation". If you do not have a prefix, you can manually update the IAM policy allow any secret names.

SpillBucket

The name of the bucket where this function can spill data.

JdbcConnectorConfig

DefaultConnectionString

The default connection string is used when catalog is "lambda:\${Lambda:functionName}". Catalog specific Connection Strings can be added later. Format: {DatabaseType}/{NativeJdbcConnectionString}.

DisableSpillEncryption

If set to "false" data spilled to S3 is encrypted with AES GCM

false

LambdaFunctionName

The name you will give to this catalog in Athena. It will also be used as the function name. This name must satisfy the pattern "^[a-z0-9_]{1,64}\$"

LambdaMemory

Lambda memory in MB (min 128 - 3008 max).

3008

LambdaTimeout

Maximum Lambda invocation runtime in seconds. (min 1 - 900 max)

900

SecurityGroups

One or more SecurityGroup IDs corresponding to the SecurityGroup that should be applied to the Lambda function. (e.g. sg-1sg2xg1)

SpillPrefix

The prefix within SpillBucket where this function can spill data.

athena-spill

SubnetIds

One or more Subnet IDs corresponding to the Subnet that the Lambda function can use to access you data source. (e.g. subnet-1subnet2)

☐ I acknowledge that this app creates custom IAM roles.

[Info](#)

Cancel

Previous

Next

Done

Secret Name prefix is used to store the DB creds in AWS Secret Manager. This is essential for a production environment. Leaving it blank means no integration will be created. After you select deploy and the Lambda you just created in the Athena dashboard you'll see a new catalog different from the standard AwsGlueCatalog:

Data source

iotarticolo

Connect data source

Note that at first, Databases and tables won't appear. Fear not: If you go to the lambda functions you'll see failures and in cloudwatch you'll see an error like:

Catalog is not supported in multiplexer. **After** registering the catalog in Athena, must **set** 'iotarticolo_connection_string' environment variable in Lambda. See JDBC connect or README **for** further details.: java.lang.RuntimeException

Go on and set the required Lambda function env variable by using the same JDBC connection string used as DefaultConnection string in the preceding step. After this, the connection will work and you'll be able to query your DB directly from Athena! Sweet!

Data source [Connect data sources](#)

Databases

Tables (3)

- + internal_metadata(Partitioned)
- + coffee(Partitioned)
- + schema_migrations(Partitioned)

New query 1 **New query 2** **New query 3** **New query 4** **New query 5** **New query 6** **New query 7** **New query 8** **New query 9** **New query 10**

```
| SELECT * FROM `internal_metadata`.`last` "coffee" LIMIT 200;
```

Run query **Save as** (Run time 3.28 seconds, Data scanned 0 KB)

Use Ctrl + Enter to save your query. Ctrl + Space to auto-complete.

Results

#	ID	pkts_id *	order *	board *	executable *	glosses *	coffee_bar_n *	partition_n *
1	1	5E_LA_0000	true	false	true	true	2017-04-19 00:00:00	*
2	2	16_A-LA_0000	true	false	true	true	2017-04-19 00:00:00	*
3	3	5E_LA_0000	true	false	true	true	2017-04-19 00:00:00	*
4	4	5E_LA_0000	true	false	true	true	2017-04-19 00:00:00	*
5	5	5E_LA_0000	true	false	true	true	2017-04-19 00:00:00	*
6	6	16_A-LA_0000	true	false	true	true	2017-04-19 00:00:00	*
7	7	https://3.amazonaws.com/artifacts/2018-03-20T17-04-16-03Z.jpg	true	false	false	false	2018-03-20 16:45:45	*
8	8	https://3.amazonaws.com/artifacts/2018-03-20T17-04-16-03Z.jpg	false	true	true	true	2018-03-20 16:45:00	*

However at a closer look we immediately notice that something is afoul with the data: here is a screen of what we can read directly from MySQL:

```
mysql> mysql select * from coffees;
```

id	photo_url	sale	beard	mustache	glasses	coffee_hour	created_at	updated_at
1	url_a_caso	1	1	1	1	2017-06-15 00:00:00	2018-03-26 16:27:08	2018-03-26 16:27:08
2	url_a_caso	1	1	0	1	2017-06-15 08:00:00	2018-03-26 16:27:08	2018-03-26 16:27:08
3	url_a_caso	1	1	0	0	2017-06-15 08:00:00	2018-03-26 16:27:08	2018-03-26 16:27:08
4	url_a_caso	1	1	0	1	2017-06-15 15:00:00	2018-03-26 16:27:08	2018-03-26 16:27:08
5	url_a_caso	1	1	0	1	2017-06-15 15:00:00	2018-03-26 16:27:08	2018-03-26 16:27:08
6	url_a_caso	1	1	0	1	2017-06-15 16:00:00	2018-03-26 16:27:08	2018-03-26 16:27:08
7	url_a_caso	1	1	0	1	2017-06-15 16:00:00	2018-03-26 16:27:08	2018-03-26 16:27:08
8	url_a_caso	1	1	0	1	2018-03-26 16:00:00	2018-03-26 16:53:00	2018-03-26 16:53:00

```
mysql> cp west-1-us-east-1.amazonaws.com:/usr/local/mysql/2018-03-26:16:51:59.9999.jpg
```

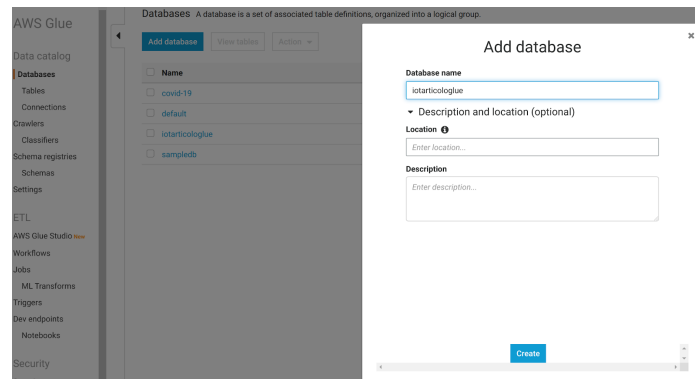
```
mysql> cp west-1-us-east-1.amazonaws.com:/usr/local/mysql/2018-03-26:16:51:59.9999.jpg
```

As you can see, Athena was smart enough to convert tinyint(1) data to bool but could not fetch mysql datetime columns. This is due to a very well known problem with jdbc connector and the easier fix is to just create a new field were the datetime is a string in Java datetime format:

```
UPDATE coffees SET coffees.coffee_hour_str=DATE_FORMAT(coffee_hour, '%Y-%m-%d %H:%i:%s');
ALTER TABLE coffees ADD COLUMN coffee_hour_str VARCHAR(255) AFTER coffee_hour;
```

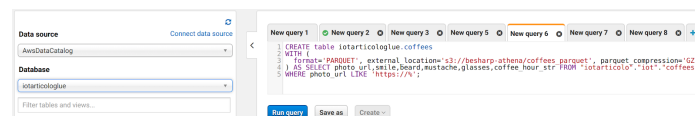
At this point, Athena will be able to read the new field.

And now we are ready for a beautiful trick: let's just go to the Glue dashboard and create a new Database. A database is just a logical container for metadata. You can chose the name you prefer:



At this point we can go back to Athena and run a query like this:

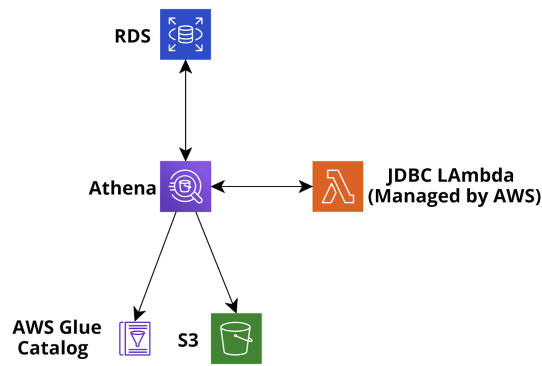
```
CREATE table iotarticologlue.coffees
WITH (
  format='PARQUET', external_location='s3://besharp-athena/coffees_parquet', parquet_compression='GZIP'
) AS SELECT photo_url,smile,beard,mustache,glasses,coffee_hour_str FROM
"iotarticolo"."iot"."coffees"
WHERE photo_url LIKE 'https://%';
```



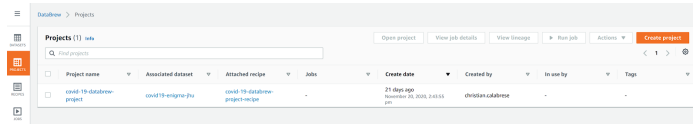
This will create a new Table in the Database we just added to our Glue data catalog and save all the data in S3 as a GZIP Parquet file. Furthermore you could also change the compression (e.g. Snappy or BZIP) if you like.

The query will also filter out the date with a bad S3 url in photo_url!

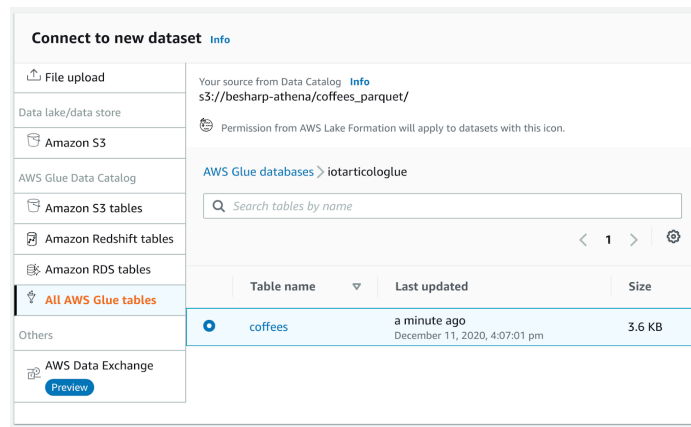
So we now have a super-fast way to export our DB to S3 as parquet while automatically creating the Glue catalog (the query does also that for free under the hood).



And now it is trivial to visualize this new catalog in AWS Glue databrew: got to dashboard and create a new project:

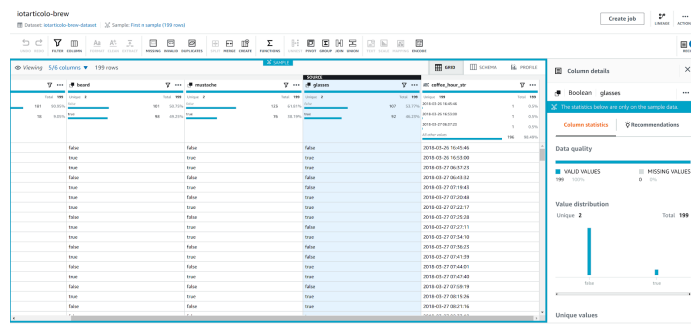


now create a new dataset in the add dataset section:



And create the project! If you encounter an error try to set the object name to parquet in s3 and to crawl again the table with Glue crawlers (Databrew is pretty new too!)

And voilà a beautiful data visualization of our dataset complete with column statistics!



Conclusion

In this article, we described a very simple **IoT application** using **Rekognition** and **Aurora**. We explained how it can be enhanced with firehose and finally we used Athena to transform and clean the collected data and save them very easily to parquet to be analyzed with Glue Databrew, Athena, and other AWS tools such as EMR.

Have you ever tried something similar for your Data Analysis process?

Feel free to write to us about your solutions: we'll be glad to offer you a "connected" coffee ☺

That's all for today.

Keep reading and see you in 14 days on **#Proud2beCloud!**



beSharp

Dal 2011 beSharp guida le aziende italiane sul Cloud. Dalla piccola impresa alla grande multinazionale, dal manifatturiero al terziario avanzato, aiutiamo le realtà più all'avanguardia a realizzare progetti innovativi in campo IT.

Get in touch

beSharp.it
proud2becloud@besharp.it

Copyright © 2011-2021 by beSharp srl - P.IVA IT02415160189