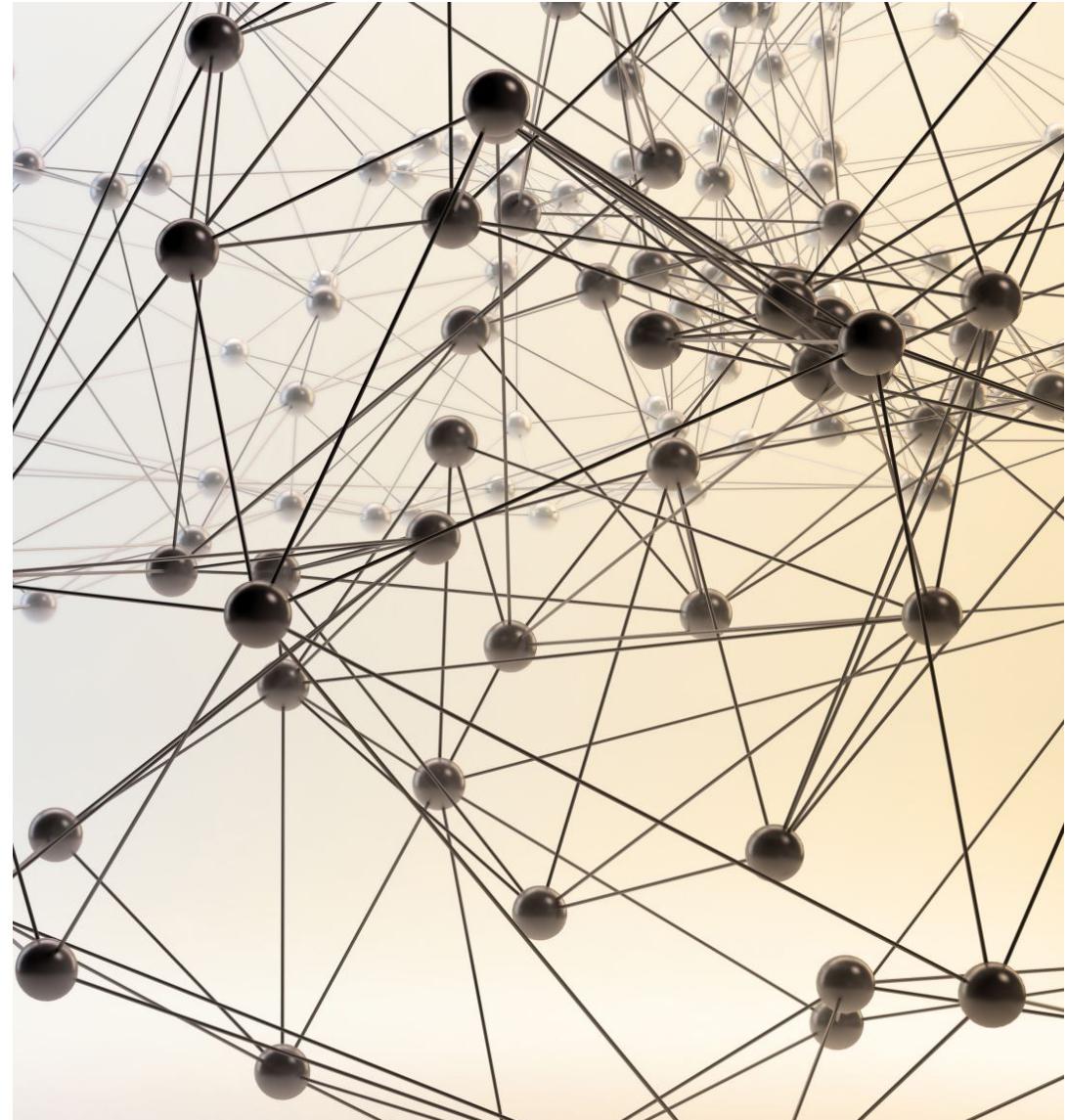




Real-time Streaming Bus Data with Kafka and Spark

by Andrea Zhang
2022 Aug



Project Motivation

Objective: The primary focus of this project is to build a real-time application that collects GPS data steamed by IoT devices located on Toronto TTC buses.

For this project, multiple big data applications were used, including Apache NiFi, CDC(Debezium), Kafka (MSK), Spark Structural Streaming, Docker, MySQL and the analytics layer using Amazon Athena.

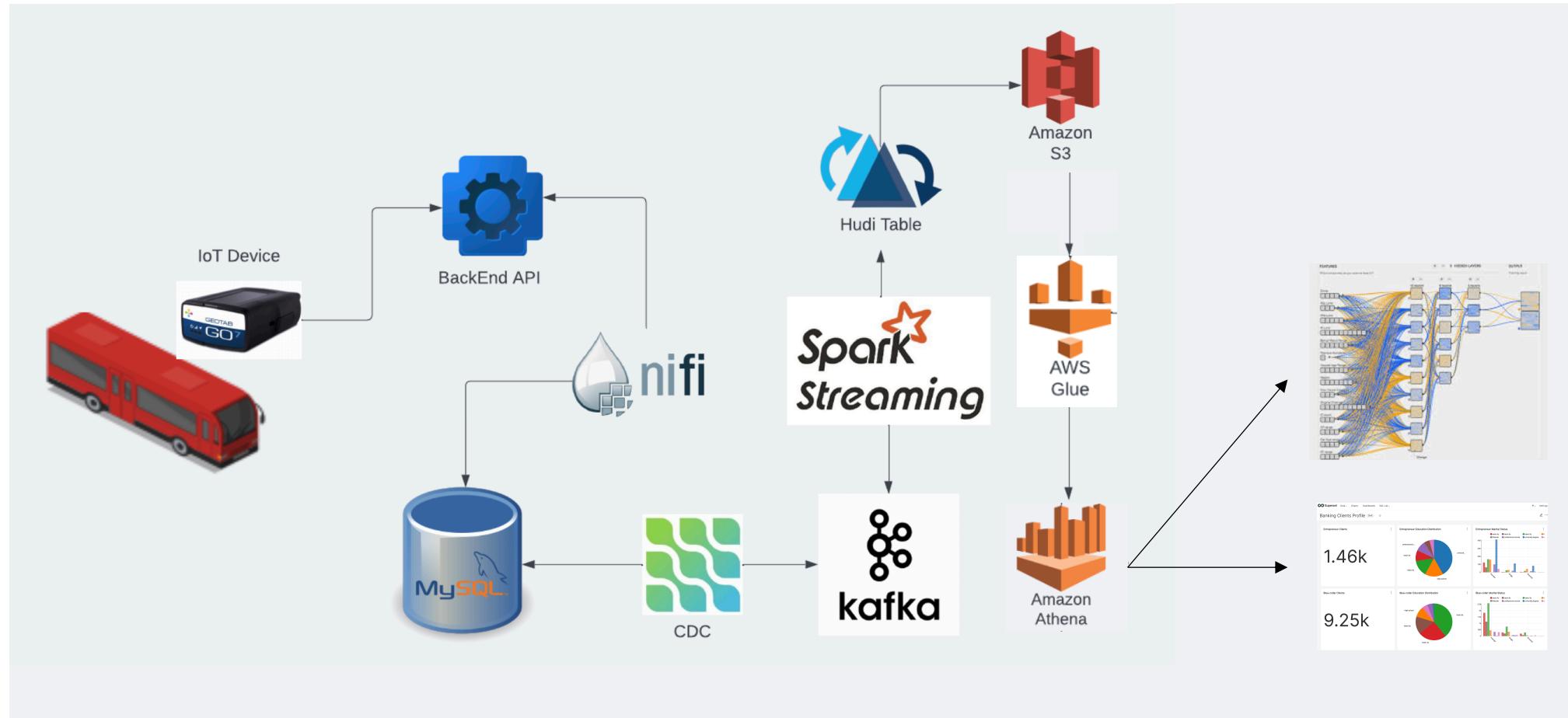
Conceptual Data Streaming Architecture

Data Collection

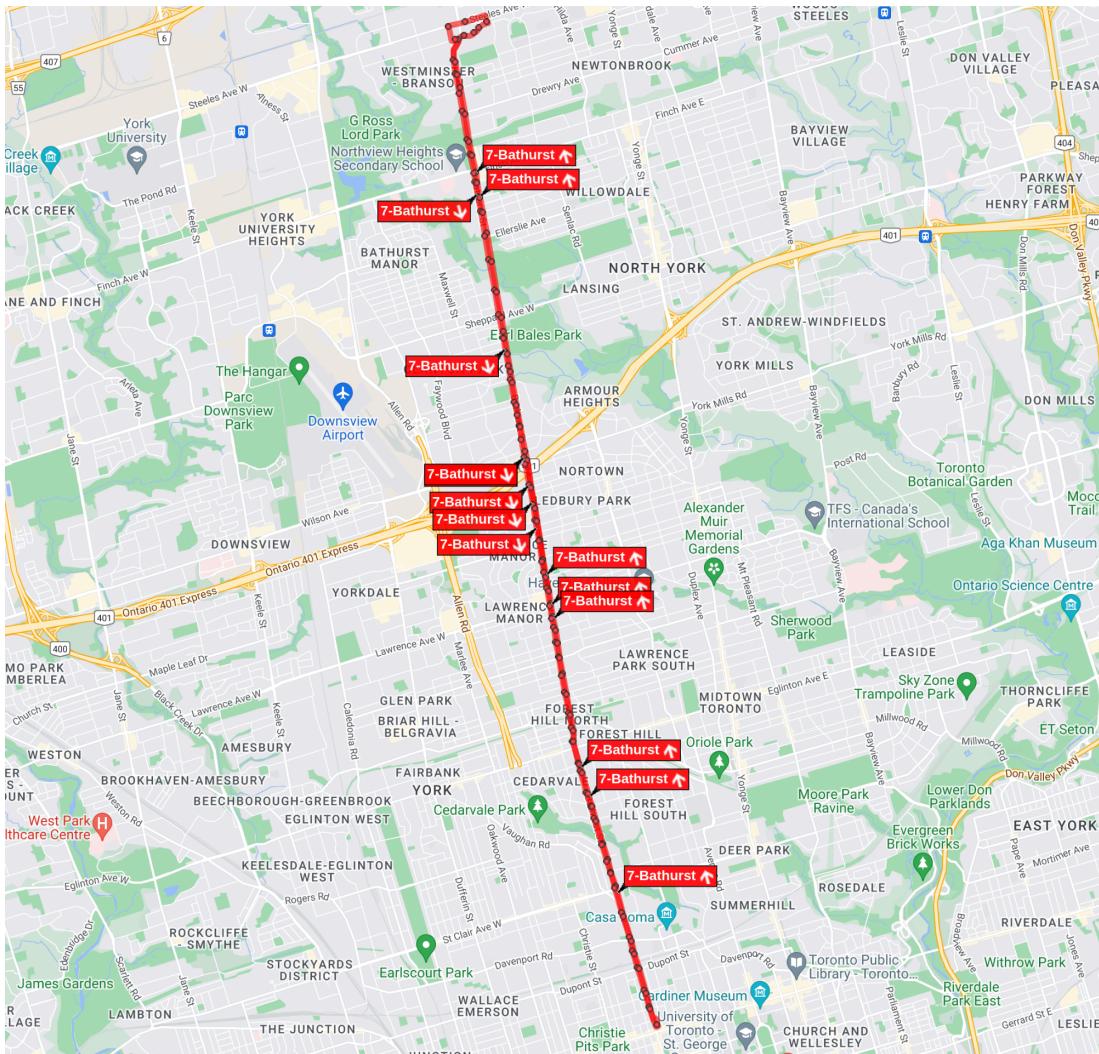
Transformation

Storage

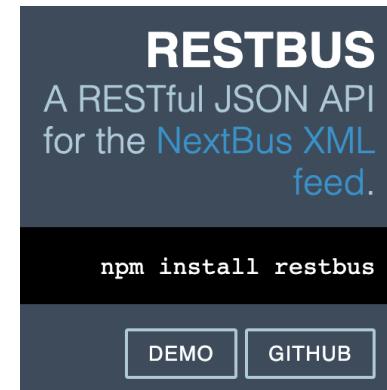
Analytics/ML



Data Collection - RestBus API



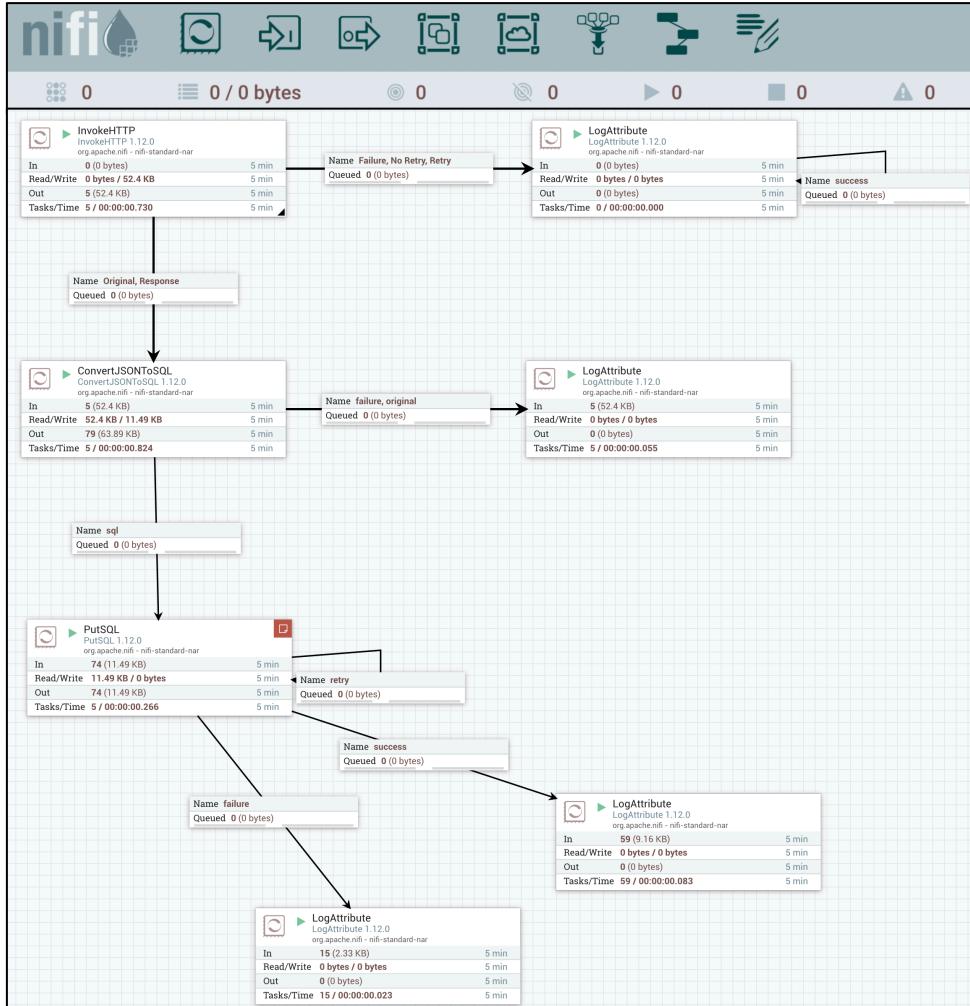
The data source used in the project is TTC bus data from an open API called [Rest Bus](#).



Sample Data

```
[{"id":7,"title":"7-Bathurst","_links": {"self":{"href":"http://restbus.info/api/agencies/ttc/routes/7","type":"application/json","rel":"http://restbus.info/_links/rel/full","rt":"route","title":"Full configuration for ttc route 7-Bathurst."}, "to": [{"href":"http://restbus.info/api/agencies/ttc/routes/7","type":"application/json","rel":"http://restbus.info/_links/rel/full","rt":"route","title":"Full configuration for ttc route 7-Bathurst."}], "from": [{"href":"http://restbus.info/api/agencies/ttc","type":"application/json","rel":"via","rt":"agency","title":"Transit agency ttc details."}]}
```

Data Transformation - Apache NiFi



Apache NiFi extracts data from the Bus API and save it to a MySQL database to make data available for the further consumption.

- Apache NiFi was set up using a Docker Container on an AWS EC2 instance
- ✓ **InvokeHTTP** processor collects data in JSON blobs. It is scheduled to get information from the Bus API every 30 seconds
- ✓ **ConvertJSONToSQL** processor converts JSON blobs into SQL inserts to update MySQL table
- ✓ **PutSQL** processor executes SQL statements that are generated by the ConvertJSONToSQL processor
- ✓ **LogAttribute** processors handle edge cases

Data Storage - MySQL

```
CREATE TABLE bus_status (
    record_id INT NOT NULL AUTO_INCREMENT,
    id INT NOT NULL,
    routeId INT NOT NULL,
    directionId VARCHAR(40),
    predictable BOOLEAN,
    secsSinceReport INT NOT NULL,
    kph INT NOT NULL,
    heading INT,
    lat REAL NOT NULL,
    lon REAL NOT NULL,
    leadingVehicleId INT,
    event_time DATETIME DEFAULT NOW(),
    PRIMARY KEY (record_id)
);
```

Bus data was constantly loaded into a MySQL table.

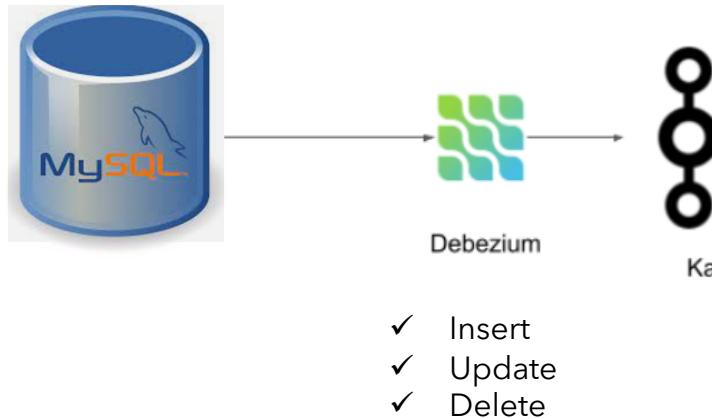
- MySQL was set up using a Docker Container on the AWS EC2 instance
- A MySQL table was created with a pre-defined schema to receive data
- ✓ The Bus data is imported from NiFi and stored in the MySQL table.
- ✓ The data in SQL database constantly gets updated and stored for further downstream processes.

mysql> select * from bus_status;												
record_id	id	routeId	directionId	predictable	secsSinceReport	kph	heading	lat	lon	leadingVehicleId	event_time	
1	8156	7	7_1_7	1	28	18	348	43.73719140	-79.43381820	NULL	2022-07-24 18:13:55	
2	8178	7	7_0_7	1	28	0	168	43.70559560	-79.42656040	NULL	2022-07-24 18:13:55	
3	8132	7	7_1_7	1	28	33	73	43.79323970	-79.44185970	NULL	2022-07-24 18:13:55	
4	8138	7	7_1_7	1	28	28	349	43.77753710	-79.44385090	NULL	2022-07-24 18:13:55	
5	8358	7	7_1_7	1	28	24	344	43.68229960	-79.41795270	NULL	2022-07-24 18:13:55	
6	8389	7	7_0_7	1	28	0	181	43.70171730	-79.42577360	NULL	2022-07-24 18:13:55	
7	8130	7	7_0_7	1	28	0	181	43.70195440	-79.42571090	NULL	2022-07-24 18:13:55	
8	8384	7	7_1_7	1	28	0	230	43.79187920	-79.44193110	NULL	2022-07-24 18:13:55	
9	8150	7	7_1_7	1	28	0	342	43.69939040	-79.42483870	NULL	2022-07-24 18:13:55	
10	8360	7	7_0_7	1	28	0	263	43.79083100	-79.44400010	NULL	2022-07-24 18:13:55	
11	8195	7	7_0_7	1	28	6	167	43.70161430	-79.42572780	NULL	2022-07-24 18:13:55	
12	8382	7	7_0_7	1	28	5	75	43.73743050	-79.43408200	NULL	2022-07-24 18:13:55	
13	8119	7	7_1_7	1	29	41	350	43.77991930	-79.44451140	NULL	2022-07-24 18:13:55	
14	8156	7	7_1_7	1	28	3	251	43.73764020	-79.43405470	NULL	2022-07-24 18:14:25	

Data Distribution - CDC (Debezium)

Run Debezium using a Kafka tool called Connect to enable MySQL and Kafka Connection. Debezium helps monitoring changes in both MySQL target table and Kafka topics.

- ✓ Debezium reads data from MySQL databases and turn it into event streams, so applications can quickly react to each row-level change in the databases, such as inserts, updates, and deletes.
 - ✓ Debezium pushes the data into a Kafka topic and monitors data changes in Kafka topics.



Kafka Topics

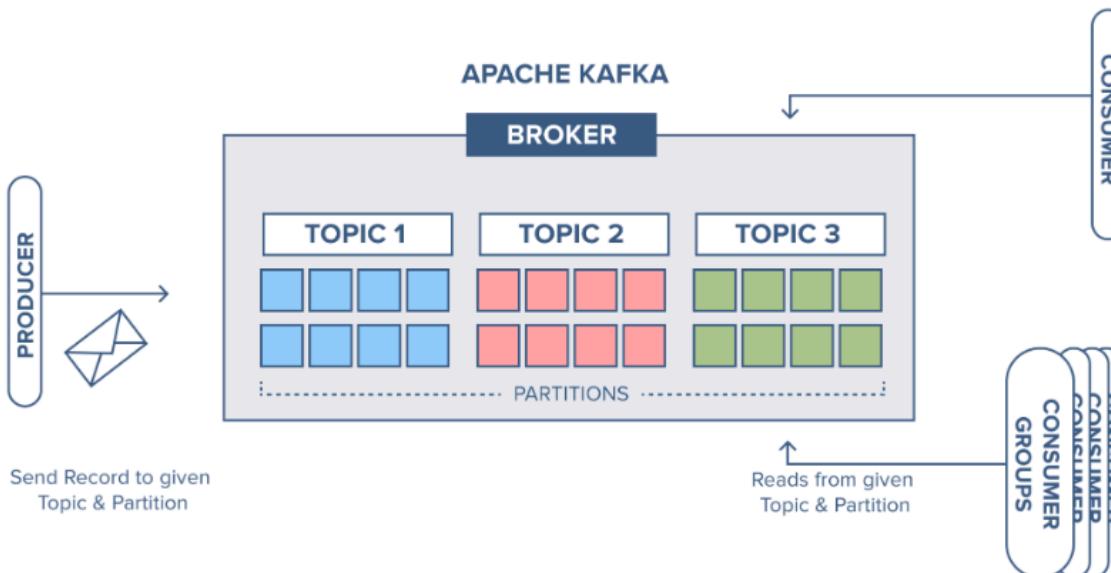
```
__amazon_msk_canary
__consumer_offsets
dbhistory.demo
dbserver1
dbserver1.demo.bus_status
my_connect_config
my_connect_offsets
my_connect_statuses
```

Check Data in Kafka Topic

Data Distributed Streaming - Kafka (MSK)

Kafka was used to transfer the Bus data from MySQL to Spark Streaming in a fast, fault-tolerance and high scalable way.

- ✓ The Kafka topic receives incoming messages in JSON blobs.
- ✓ Built-in partitioning
- ✓ Replication



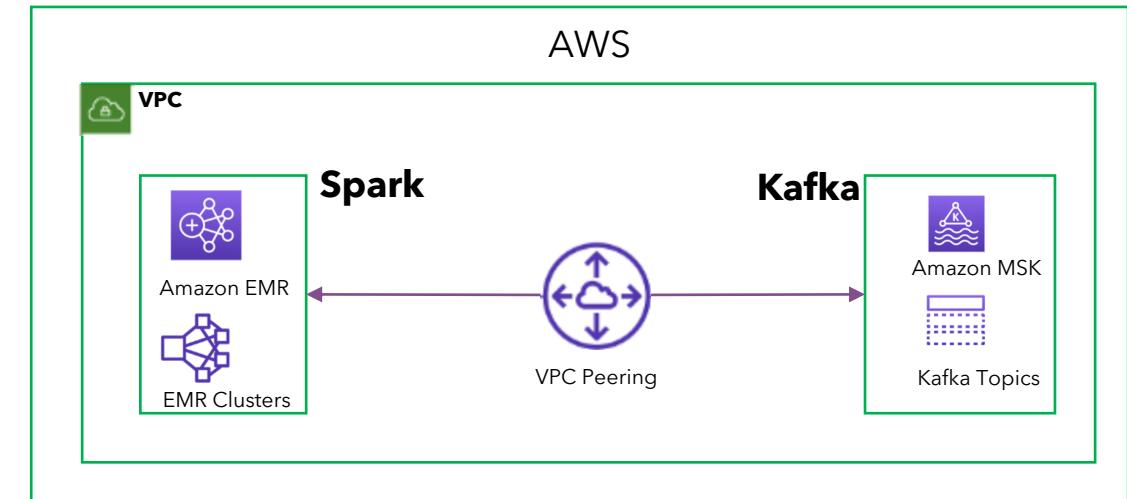
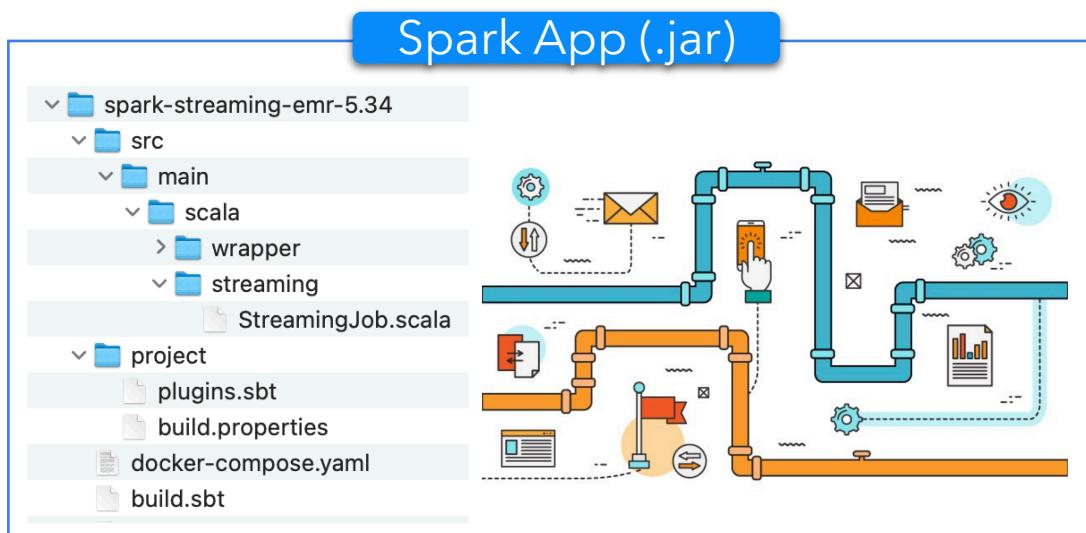
Kafka Topic Description

```
$ bin/kafka-topics.sh --describe my_connect_config --bootstrap-servers mp3s.1qo03c.c23.kafka.us-east-1.amazonaws.com:9092
topic:my_connect_config PartitionCount: 1 ReplicationFactor: 2
topic:my_connect_config Partition: 0 Leader: 2 Replicas: 2
topic:my_connect_config PartitionCount: 5 ReplicationFactor: 2
topic:my_connect_config Partition: 0 Leader: 2 Replicas: 2
topic:my_connect_config Partition: 1 Leader: 1 Replicas: 2
topic:my_connect_config Partition: 2 Leader: 3 Replicas: 2
topic:my_connect_config Partition: 3 Leader: 2 Replicas: 2
topic:my_connect_config Partition: 4 Leader: 1 Replicas: 2
topic:my_connect_config Partition: 5 Leader: 0 Replicas: 2
```

Data Streaming - Spark Streaming

Ingesting Data From Kafka with Spark Streaming

- Create an AWS EMR for Spark and enable EMR and MSK communication through CPV Perring.
- Compiling Spark Streaming Jobs Built in Scala. Upload .jar file to an S3 bucket.
- Spark-submit using the .jar file to the Spark (EMR cluster) master node.
- The Spark application ingests data from Kafka, transform the data according to a predefined schema (using an S3 file), and export data to an S3 bucket.



Data Streaming - Spark Streaming

Micro batches of data coming from the Kafka topic.

Micro-batch processing: divides incoming streams into small batches for processing.

Batch: 189										
ldirectionId	event_time	lheadingId	lkphi	lat	lleadingVehicleId	lon	lpredictable	lrecord_id	lrouteId	lsecsSinceReport
17_0_7	1656546597000 281	1900410	143.790825	null	-79.444496	true	14706	17	126	
17_1_7	1656546597000 349	1900211	143.76817	null	-79.441475	true	14707	17	126	
17_1_7	1656546597000 172	1900610	143.66606	null	-79.411064	true	14708	17	127	
17_0_7	1656546597000 169	1831310	143.737576	null	-79.434235	true	14709	17	126	
17_1_7	1656546597000 349	18395131	143.7115	null	-79.42788	true	14710	17	126	
17_0_7	1656546597000 173	19000124	143.719265	null	-79.4298	true	14711	17	126	
lnull	1656546597000 142	18309114	143.741898	null	-79.45461	false	14712	17	127	
17_0_7	1656546597000 172	1810210	143.666378	null	-79.411125	true	14713	17	126	
17_1_7	1656546597000 349	1901210	143.73826	null	-79.43412	true	14714	17	126	
17_0_7	1656546597000 216	1816510	143.790855	null	-79.44464	true	14715	17	126	
17_1_7	1656546597000 337	19018116	143.68691	null	-79.419716	true	14716	17	125	
17_0_7	1656546597000 163	19016122	143.768547	null	-79.44178	true	14717	17	126	
17_0_7	1656546597000 169	18302125	143.70999	null	-79.42762	true	14718	17	126	
17_1_7	1656546597000 349	1370139	143.77536	null	-79.443245	true	14719	17	126	
17_0_7	1656546627000 281	1900410	143.790825	null	-79.444496	true	14720	17	122	
17_1_7	1656546627000 349	19002134	143.770443	null	-79.44208	true	14721	17	122	
17_1_7	1656546627000 172	1900610	143.66606	null	-79.411064	true	14722	17	123	
17_0_7	1656546627000 167	18313146	143.736576	null	-79.43387	true	14723	17	122	
17_1_7	1656546627000 350	1839510	143.71225	null	-79.42798	true	14724	17	122	
17_0_7	1656546627000 169	1900014	143.71797	null	-79.42948	true	14725	17	123	

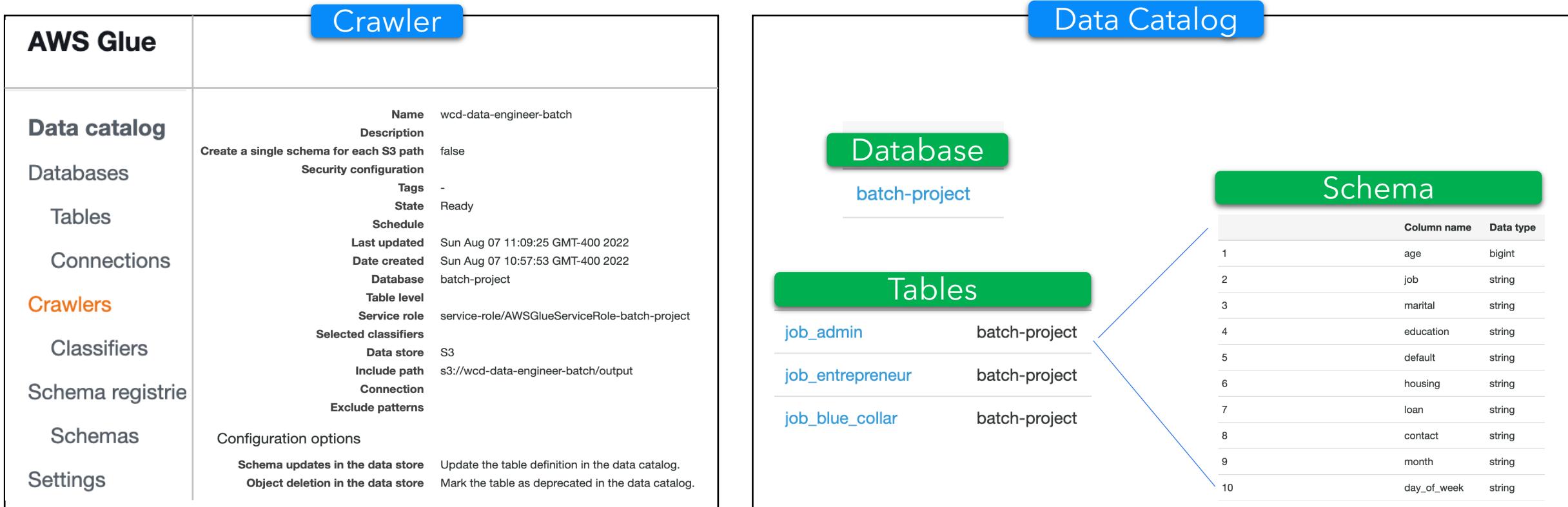
only showing top 20 rows

Data Integration - Glue

Glue provides serverless data integration that prepare and combine data for analytics.

The Glue data Catalog serves as a central metadata repository.

- ✓ Discover and search across multiple data sets using Glue data Catalog.
- ✓ Once the data is cataloged, it is immediately available for search and query using Amazon Athena.



Data Analytics - AWS Athena

Connect Glue Data Catalog and point to the streaming project database and enable querying the bus data instantly.

The screenshot shows the Amazon Athena Query editor interface. On the left, there's a sidebar titled 'Data' with sections for 'Data source' (set to 'AwsDataCatalog'), 'Database' (set to 'streaming-project'), and 'Tables and views'. A table named 'output' is listed under 'Tables (1)'. The main area is the 'Query editor' with two tabs: 'Query 2' (green checkmark) and 'Query 3' (orange checkmark). The SQL query in 'Query 3' is:

```
1 | SELECT * FROM "streaming-project"."output" limit 10;
```

Below the query, the status bar shows 'SQL Ln 1, Col 1'. To the right of the status bar are buttons for 'Run again', 'Explain', 'Cancel', 'Save', 'Clear', and 'Create'. Underneath these buttons are tabs for 'Query results' and 'Query stats'. The 'Query results' tab is selected, showing a table titled 'Results (10)' with 10 rows of data. The columns are: #, record_id, id, routeid, directionid, predictablr, seccssincereport, kp, and headin. The data is as follows:

#	record_id	id	routeid	directionid	predictablr	seccssincereport	kp	headin
1	1	8156	7	7_1_7	1	28	18	348
2	2	8178	7	7_0_7	1	28	0	168
3	3	8132	7	7_1_7	1	28	33	73
4	4	8138	7	7_1_7	1	28	28	349
5	5	8358	7	7_1_7	1	28	24	344
6	6	8389	7	7_0_7	1	28	0	181
7	7	8130	7	7_0_7	1	28	0	181
8	8	8384	7	7_1_7	1	28	0	230
9	9	8150	7	7_1_7	1	28	0	342

Use Cases



Monitoring



Trip Planner



Optimization



ML predictions