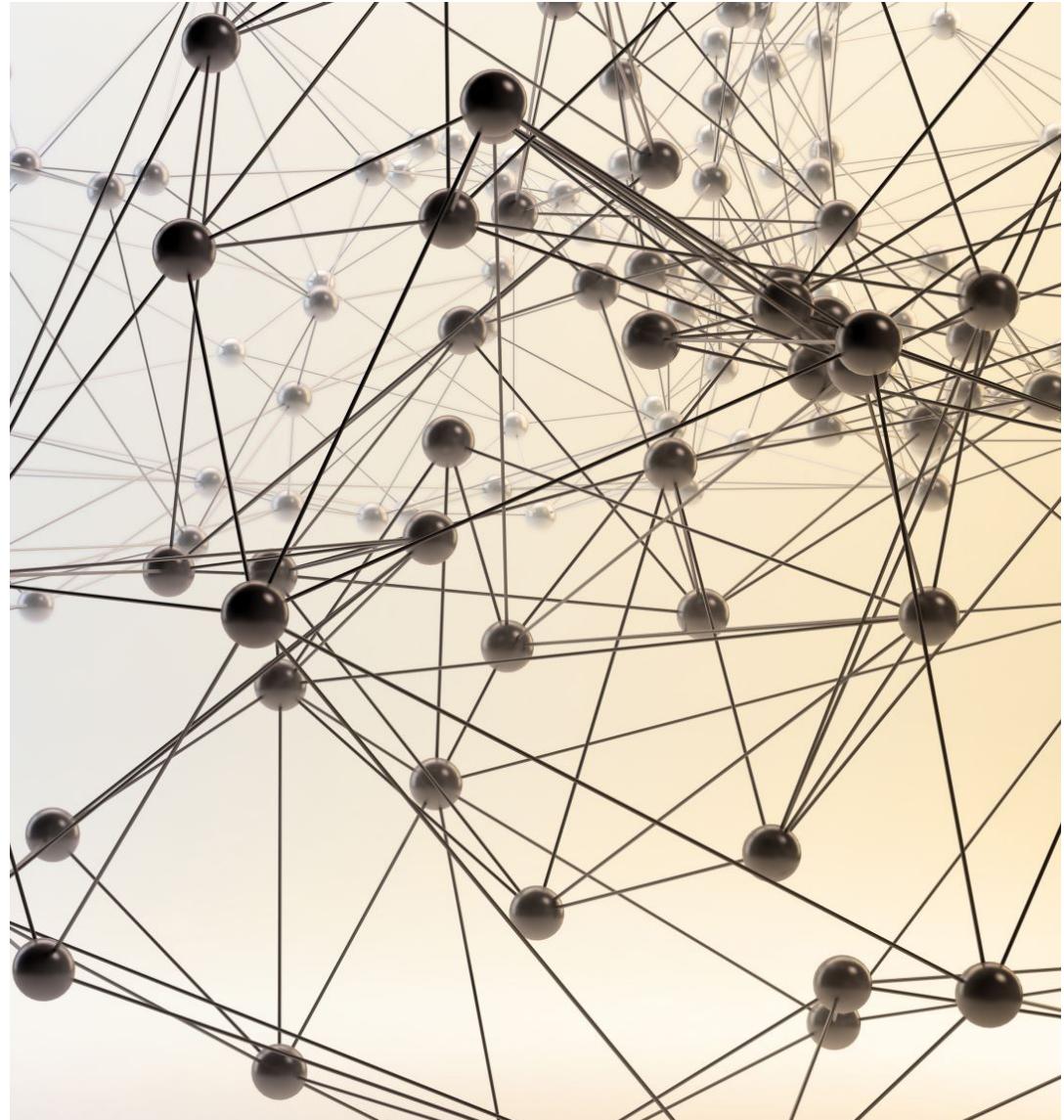




# Real-time Streaming Bus Data with Kafka and Spark

---

by Andrea Zhang



# 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.

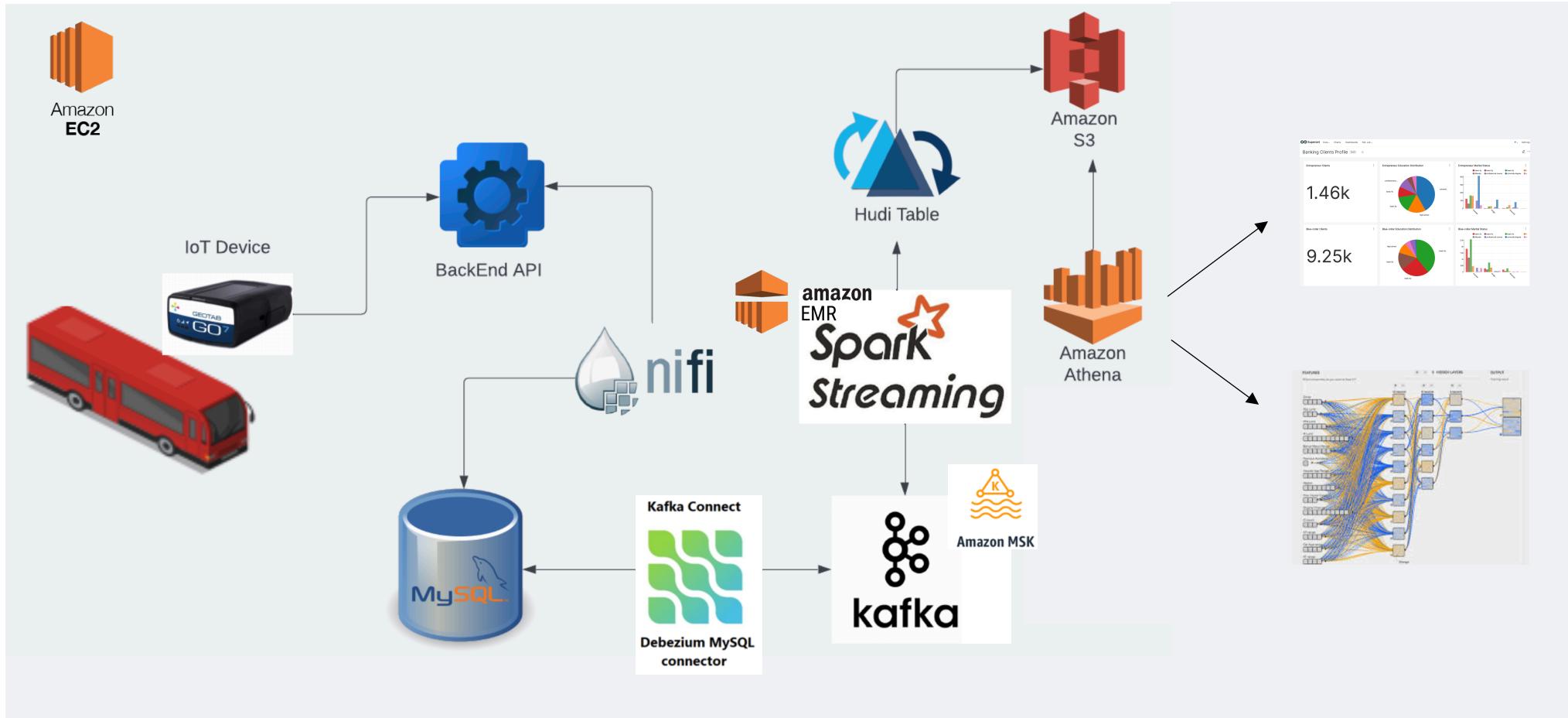
# Conceptual Data Streaming Architecture

Data Collection

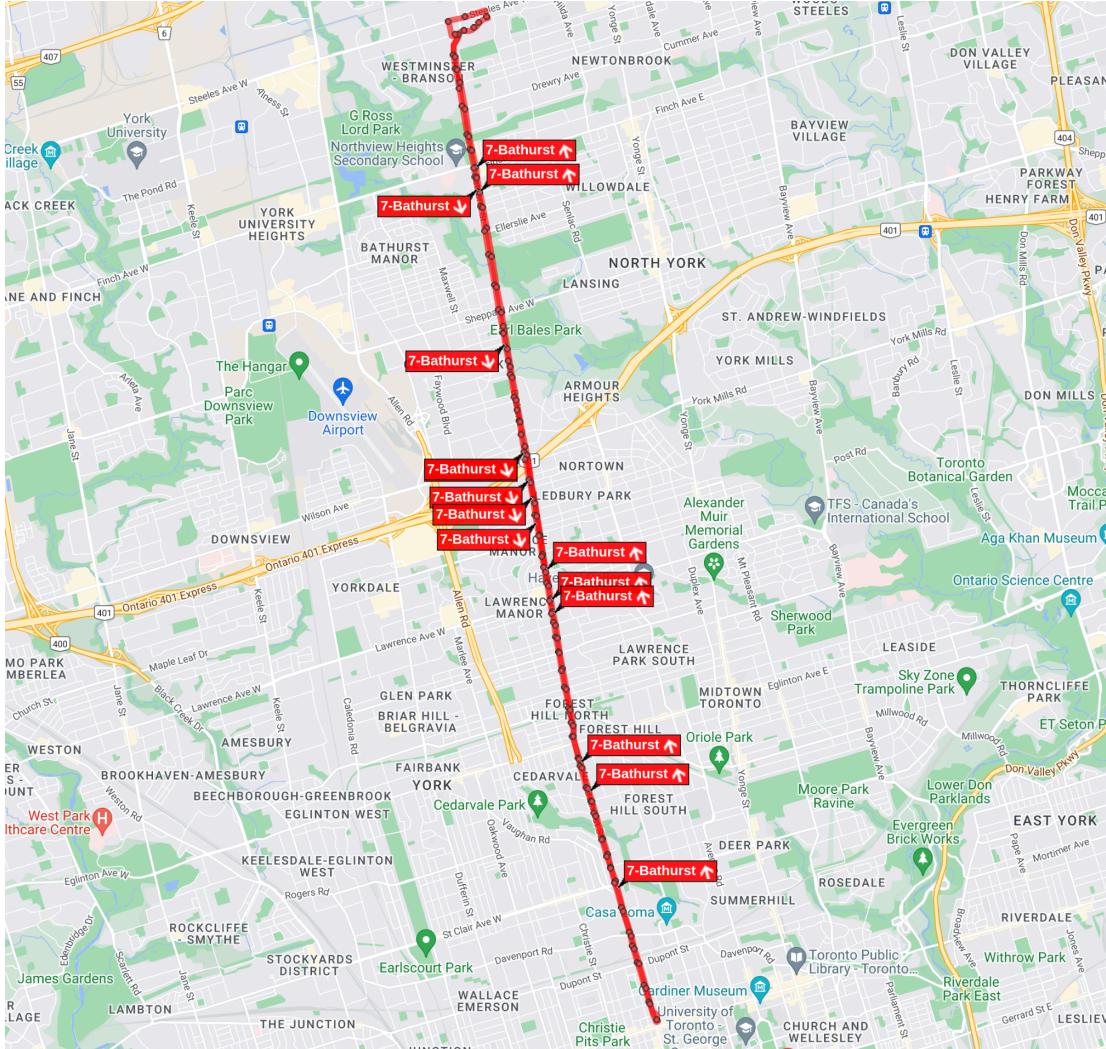
Transformation

Storage

Analytics/ML



# Data Collection - RestBus API



In order to collect the public bus data in real time, an open API called [Rest Bus](#) is used to import the XML Feed data using URLs with parameters specified in the query string.

Data output is in a JSON format. This data need to be transferred to an analytical database for further consumption.

## 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."}]}
```

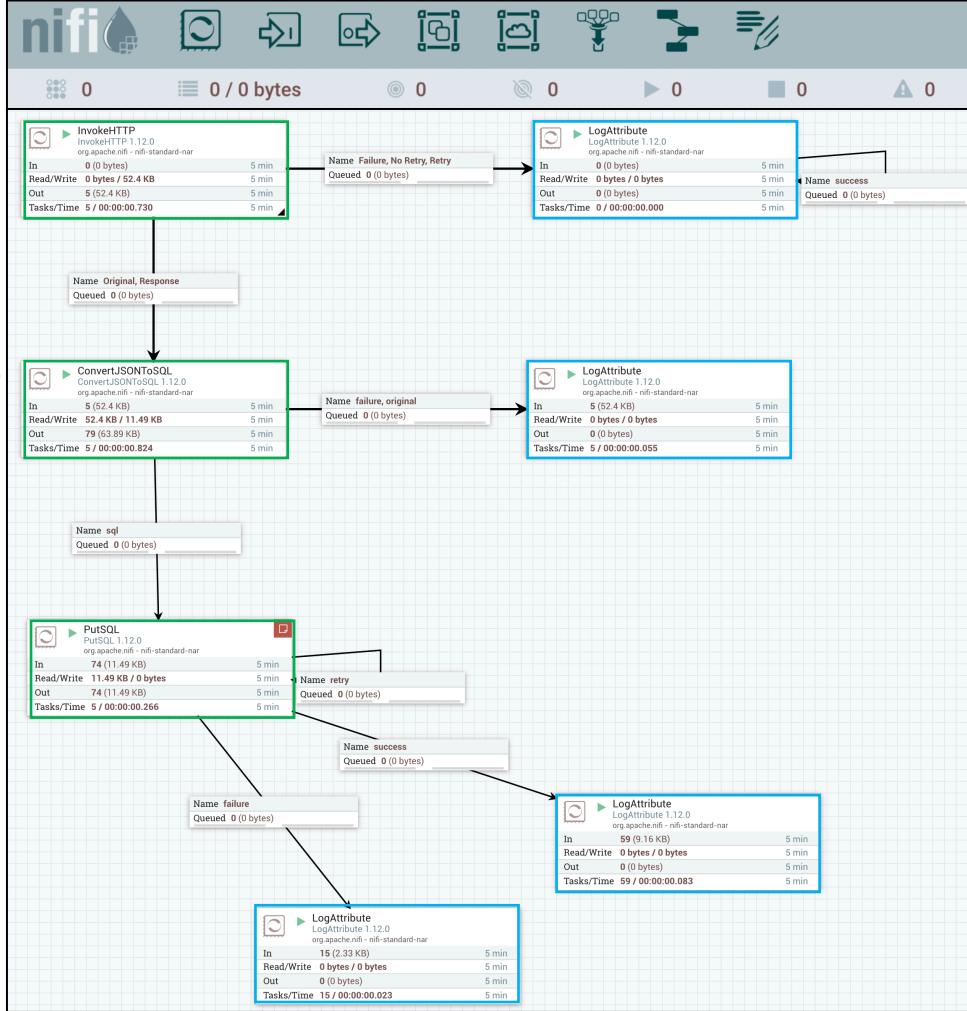
# Data Transformation - Apache NiFi

**InvokeHTTP**  
collects data in JSON blobs with a scheduler of every 30s.

**ConvertJSONToSQL**  
converts data from JSON blobs into SQL insert statements.

**PutSQL**  
executes SQL insert statements.

**LogAttribute**  
Retry upon failure.



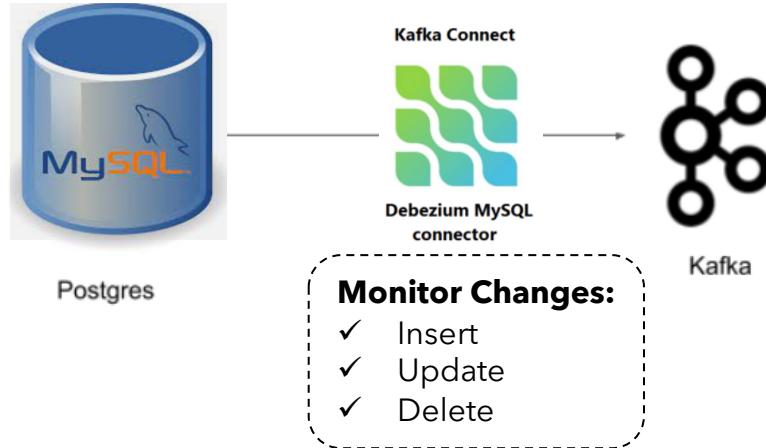
**Apache NiFi is the selected solution for the data transformation problem because it provides a drag-and-drop tool that creates a data flow framework to extract the large amount of data from the bus API and send it to an analytical database with lightning speed effortlessly.**

**Similar to Apache Kafka, but different.**

Pros	Cons
<ul style="list-style-type: none"><li>Easy to use</li><li>Scalability</li><li>Security - pull out the queued data at any instance</li><li>Processor group - better organization of different pipelines</li><li>Template pipeline - quick creation of multiple pipelines processing different bus routes using a template</li></ul>	<ul style="list-style-type: none"><li>Once the node is detached from the cluster, it's difficult to fetch the data from that node. Although you can manually copy the data from that node, it will occupy your time and energy. However, if you are on the receiver end, you cannot find the data from the detached node until the admin follows the above procedure.</li><li>Another limitation is the state persistence issue of the Apache NiFi. The issue makes it difficult for the source to fetch the data. If you experience an issue in the first node, you can expect the impact in other nodes, n as well.</li></ul>

# Data Distribution - MySQL with Debezium

The data from Bus API will be stored in MySQL tables, constantly updated, and sent over the changes to the Kafka topic via Debezium's Kafka Connect.



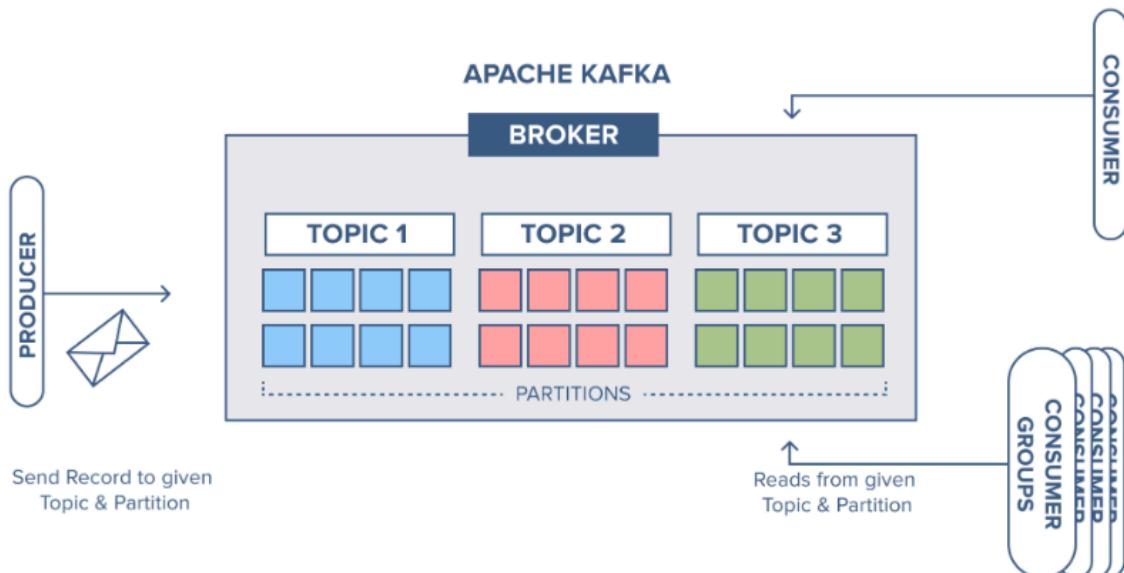
Set up MySQL with Debezium as an extra layer to help to capture each row-level changes in the database and send it over to Kafka.

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 - Apache Kafka (Amazon 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 Topics

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

Topic Description

```
]$ bin/kafka-topics.sh --describe my_connect_config --bootstrap-server 3.1q003c.c23.kafka.us-east-1.amazonaws.com:9092
Topic: my_connect_config PartitionCount: 1 ReplicationFactor: 2
Leader: 2 Replicas: 2
Partition: 0 Leader: 2 Replicas: 2
PartitionCount: 5 ReplicationFactor: 2
Leader: 1 Replicas: 2
Partition: 0 Leader: 2 Replicas: 2
Partition: 1 Leader: 1 Replicas: 2
Partition: 2 Leader: 3 Replicas: 2
Partition: 3 Leader: 2 Replicas: 2
Partition: 4 Leader: 1 Replicas: 2
```

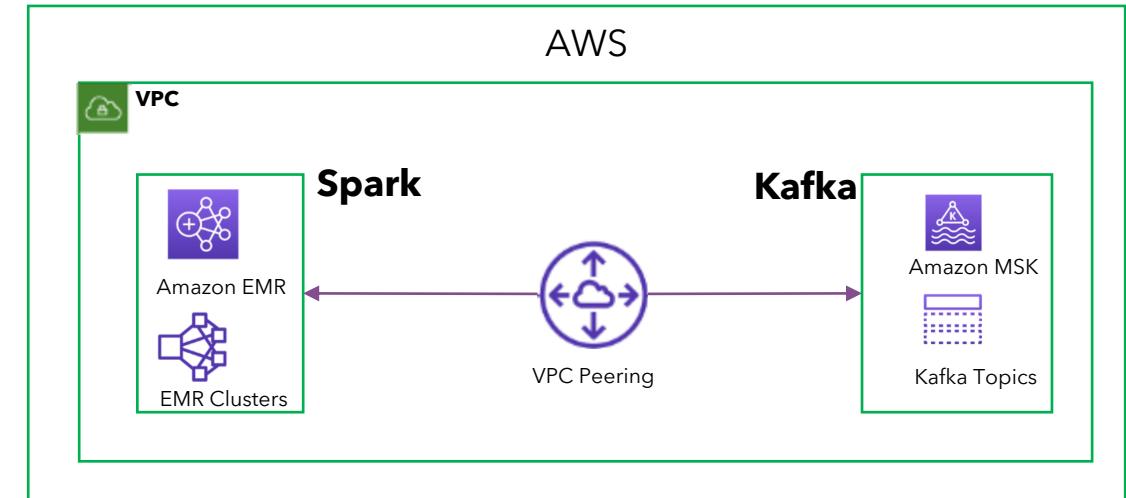
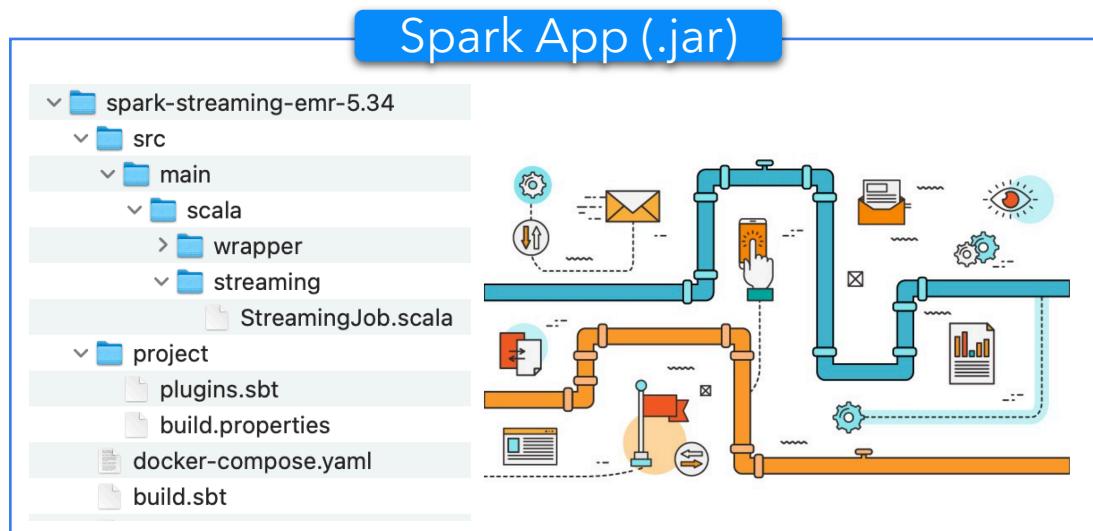
Data in Topic

```
{"schema": {"type": "struct", "fields": [{"type": "struct", "fields": [{"type": "int32", "optional": false, "field": "record_offset": true, "field": "directionId"}, {"type": "int16", "optional": true, "field": "predictable"}, {"type": "int32", "optional": true, "field": "heading"}, {"type": "bytes", "optional": false, "name": "org.apache.kafka.connect.data.Decimal", "version": 1, "parameters": {"scale": 8}, "field": "bus_status.Value", "field": "id"}, {"type": "int32", "optional": false, "field": "routeId"}, {"type": "string", "optional": true, "name": "dbserver1.demo.bus_status.Value", "field": "event_time"}, {"type": "int32", "optional": true, "field": "kph"}, {"type": "int32", "optional": true, "field": "heading"}, {"type": "float", "optional": true, "field": "lon"}, {"type": "float", "optional": true, "field": "lat"}, {"type": "bytes", "optional": false, "name": "org.apache.kafka.connect.data.Decimal", "version": 1, "parameters": {"scale": 10}, "field": "lon"}, {"type": "int64", "optional": true, "name": "io.debezium.time.Timestamp", "version": 1, "parameters": {"precision": 10}, "field": "leadingVehicleId"}, {"type": "int64", "optional": true, "name": "io.debezium.time.Timestamp", "version": 1, "parameters": {"precision": 10}, "field": "last_fall_increments"}, {"type": "string", "optional": false, "name": "io.debezium.data.Enum", "version": 1, "parameters": {"allowed": "true, last_fall_increments"}, "field": "sequence"}, {"type": "string", "optional": true, "field": "table"}, {"type": "int64", "optional": false, "field": "server_id"}, {"type": "string", "optional": false, "name": "io.debezium.data.TableOrder", "version": 1, "parameters": {"order": "row"}, "field": "source"}, {"type": "string", "optional": false, "field": "op"}, {"type": "int64", "optional": true, "field": "data_collection_order"}, {"type": "int64", "optional": false, "field": "routeId"}, {"type": "int64", "optional": true, "field": "secsSinceReport"}, {"type": "int32", "optional": false, "field": "row"}, {"type": "int64", "optional": true, "field": "transaction_id"}, {"type": "int64", "optional": true, "field": "ts_ms"}, {"type": "string", "optional": true, "name": "dbserver1", "field": "snapshot"}, {"type": "string", "optional": true, "name": "db", "field": "demo"}, {"type": "string", "optional": true, "name": "sequence", "field": "bus_status"}, {"type": "string", "optional": true, "name": "table", "field": "bus_status"}, {"type": "string", "optional": true, "name": "transaction": null}}}
```

# 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 and produces parquet files in a 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 Analytics - AWS Athena

Spark streaming app point to Athena table 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') and 'Database' (set to 'streaming-project'). Below that is a 'Tables and views' section with a 'Create' button and a 'Filter tables and views' search bar. Under 'Tables and views', there's a 'Tables (1)' section listing a single table named 'output'. On the right, the main area is titled 'Query 3' (with 'Query 2' also visible). It contains a SQL query: 'SELECT \* FROM "streaming-project"."output" limit 10;'. Below the query are buttons for 'Run again', 'Explain', 'Cancel', 'Save', 'Clear', and 'Create'. The 'Query results' tab is selected, showing a table titled 'Results (10)'. The table has columns: #, record\_id, id, routeid, directionid, predictabl, seccssincereport, kp, and headin. The data is as follows:

#	record_id	id	routeid	directionid	predictabl	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