# Best Practices for Building Serverless Big Data Applications

## Ben Snively, Solutions Architect
## Data and Analytics

aws

Serverless technologies let you build and scale applications and services rapidly without the need to provision or manage servers. In this session, we show you how to incorporate serverless concepts into your big data architectures. We explore the concepts behind and benefits of serverless architectures for big data, looking at design patterns to ingest, store, process, and visualize your data. Along the way, we explain when and how you can use serverless technologies to streamline data processing, minimize infrastructure management, and improve agility and robustness and share a reference architecture using a combination of cloud and open source technologies to solve your big data problems. Topics include: use cases and best practices for serverless big data applications; leveraging AWS technologies such as Amazon DynamoDB, Amazon S3, Amazon Kinesis, AWS Lambda, Amazon Athena, and Amazon EMR; and serverless ETL, event processing, ad hoc analysis, and real-time analytics.

## Agenda

Serverless – what and why?

Serverless – Which Service When?

Common Big Data Applications

Fitting Serverless into Big Data Applications.

Next Steps…

# Serverless Analytics Evolution...



Virtualized ➡ Managed ➡ Serverless

Provision servers · Configure Clusters · Run Analytics

# Serverless characteristics



No servers to provision or manage

Scales with usage

Never pay for idle

Availability and fault tolerance built in

## Serverless nicely fits into big data platforms

- Mix and Match Serverless, Managed, and Virtualized
- Leverage Services to easily
  - Rapidly ingest, categorize, and discover your data
  - Allow easy query and analysis of your data
  - Transform and Load data
  - Provide custom event based handlers
- Serverless allows you to focuses more analytics and not on infrastructure or servers

## Serverless Compute

AWS Lambda

- **Run your code in the cloud - fully managed and highly-available**
- **Triggered through API or state changes in your setup**
- **Scales automatically to match the incoming event rate**
- **Node.js (JavaScript), Python, Java, and C#**
- **Charged per 100ms execution time**

# Serverless Interactive Query Service



Amazon
Athena

- Query directly from Amazon S3
- Use ANSI SQL
- Serverless
- Multiple Data Formats
- Pay per query

# Serverless Catalog and ETL/ELT Service



AWS Glue

**Data Catalog**

Crawl, Discover and Organize Data
Integration with Managed and Serverless Analytics

**Job Authoring**

Serverless ETL – Pay for what you consume

**Job Execution**

# Serverless Streaming Made Easy
## Services make it easy to capture, deliver and process streams on AWS

| Amazon Kinesis Streams | Amazon Kinesis Firehose | Amazon Kinesis Analytics |
|---|---|---|
| • For Technical Developers<br>• Build your own custom applications that process or analyze streaming data | • For all developers, data scientists<br>• Easily load massive volumes of streaming data into S3, Amazon Redshift and Amazon Elasticsearch | • For all developers, data scientists<br>• Easily analyze data streams using standard SQL queries |

---

# Applying Serverless to Big Data Applications?

Data processing | Data warehousing | Reporting | Real-time processing | Predictive analytics | Artificial Intelligence

---

# Characteristics of a Big Data Applications

Collect Anything | Dive in Anywhere | Flexible Access | Future Proof

**Components of Big Data Applications**

Catalog & Search

Ingest & Store — Prepare & Transform — Analyze & Reason — Access & User Interface

Protect & Secure



**Components of Big Data Applications**

Glue Data Catalog

Amazon Kinesis, S3 — Ingest & Store

AWS Lambda, AWS Glue — Prepare & Transform

Kinesis Analytics, Amazon Athena — Analyze & Reason

Amazon QuickSight — Access & User Interface

Protect & Secure

**AWS Glue** has the capability to be a metadata catalog that helps you bring in and maintain new dataset. **Lambda** allows you to write transformation logic for your data as they come in, **Glue** also allows you to write transformations of your data also. **Kinesis Analytics** provides SQL based access on top of data streams while Athena provides SQL based access on top of your data stored on **S3**.

**Components of Big Data Applications**

This is a very common real time analytical flow, it is very easy to setup to start capturing, transforming, and analyzing data



**Serverless Real-time Analytics**



**Demonstration**

This is a Kinesis Firehose, we have specified an S3 bucket with a prefix *raw/*.  We want to capture app messages and analyzed the data

{"sentiment_value": 0.0, "name": "oceane.batz", "sentiment": "neutral", "timestamp": 1511711811578, "longitude": 20.168331146240234, "country": "Albania", "latitude": 41.1533317565918, "message": "The group scavenged a surgical clinic", "channel": "default"}
{"sentiment_value": 0.6124, "name": "mohammed.paucek", "sentiment": "positive", "timestamp": 1511711811679, "longitude": 20.168331146240234, "country": "Albania", "latitude": 41.1533317565918, "message": "The group stumbled upon the nearby urgent care", "channel": "default"}
{"sentiment_value": 0.0, "name": "rodrick.russel", "sentiment": "neutral", "timestamp": 1511711811779, "longitude": 20.168331146240234, "country": "Albania", "latitude": 41.1533317565918, "message": "I discovered a mellow Sugar", "channel": "default"}
{"sentiment_value": 0.6124, "name": "scottie.hirthe", "sentiment": "positive", "timestamp": 1511711811879, "longitude": 20.168331146240234, "country": "Albania", "latitude": 41.1533317565918, "message": "I hunted down the local urgent care", "channel": "default"}
{"sentiment_value": 0.0, "name": "friedrich.larson", "sentiment": "neutral", "timestamp": 1511711811980, "longitude": 20.168331146240234, "country": "Albania", "latitude": 41.1533317565918, "message": "This morning the group came upon a naval doctor's office", "channel": "default"}
{"sentiment_value": 0.0, "name": "evert.satterfield", "sentiment": "neutral", "timestamp": 1511711812085, "longitude": 20.168331146240234, "country": "Albania", "latitude": 41.1533317565918, "message": "The group signed a truce with a golden swing, reported a municipal clam", "channel": "default"}
{"sentiment_value": 0.0, "name": "willard.bogan", "sentiment": "neutral", "timestamp": 1511711812185, "longitude": 20.168331146240234, "country": "Albania", "latitude": 41.1533317565918, "message": "This morning i spoke to the crib, reported a famous elephant", "channel": "default"}
{"sentiment_value": 0.4588, "name": "sheldon.cummerata", "sentiment": "positive", "timestamp": 1511711812285, "longitude": 20.168331146240234, "country": "Albania", "latitude": 41.1533317565918, "message": "I ate the sweet Oysters", "channel": "default"}
{"sentiment_value": 0.0, "name": "felicita.turcotte", "sentiment": "neutral", "timestamp": 1511711812474, "longitude": 20.168331146240234, "country": "Albania", "latitude": 41.1533317565918, "message": "I was converted into a red stove", "channel": "default"}
{"sentiment_value": -0.34, "name": "madie.schowalter", "sentiment": "negative", "timestamp": 1511711812574, "longitude": 20.168331146240234, "country": "Albania", "latitude": 41.1533317565918, "message": "I plans to cease fire against a swing", "channel": "default"}
{"sentiment_value": 0.0, "name": "josie.schamberger", "sentiment": "neutral", "timestamp": 1511711812674, "longitude": 20.168331146240234, "country": "Albania", "latitude": 41.1533317565918, "message": "I scavenged a makeshift doctor's office", "channel": "default"}
{"sentiment_value": 0.0, "name": "lorine.stanton", "sentiment": "neutral", "timestamp": 1511711812775, "longitude": 20.168331146240234, "country": "Albania", "latitude": 41.1533317565918, "message": "I drank the juicy Loquats", "channel": "default"}
{"sentiment_value": 0.0, "name": "coby.greenholt", "sentiment": "neutral", "timestamp": 1511711812875, "longitude":

aws    Services ▾    Resource Groups ▾    ★            △   Ben Snively ▾   N. Virginia ▾   Support ▾

**Amazon Kinesis** ‹

Dashboard

Data Streams

**Data Firehose**

Data Analytics

Firehose delivery streams  >  ZombieDataLake-IngestionFirehoseStream-1GPTKD0DHCJBD

▸ Test with demo data

Use the tabs below to view, edit and monitor your delivery stream.

| Details | Monitoring | S3 Logs |

Delete Delivery Stream

Edit

Delivery stream name*  ZombieDataLake-IngestionFirehoseStream-1GPTKD0DHCJBD

Source  Direct PUT

S3 bucket  ☑ zombiedatalake-zombiebucket-md5turcapouy

S3 prefix  raw/
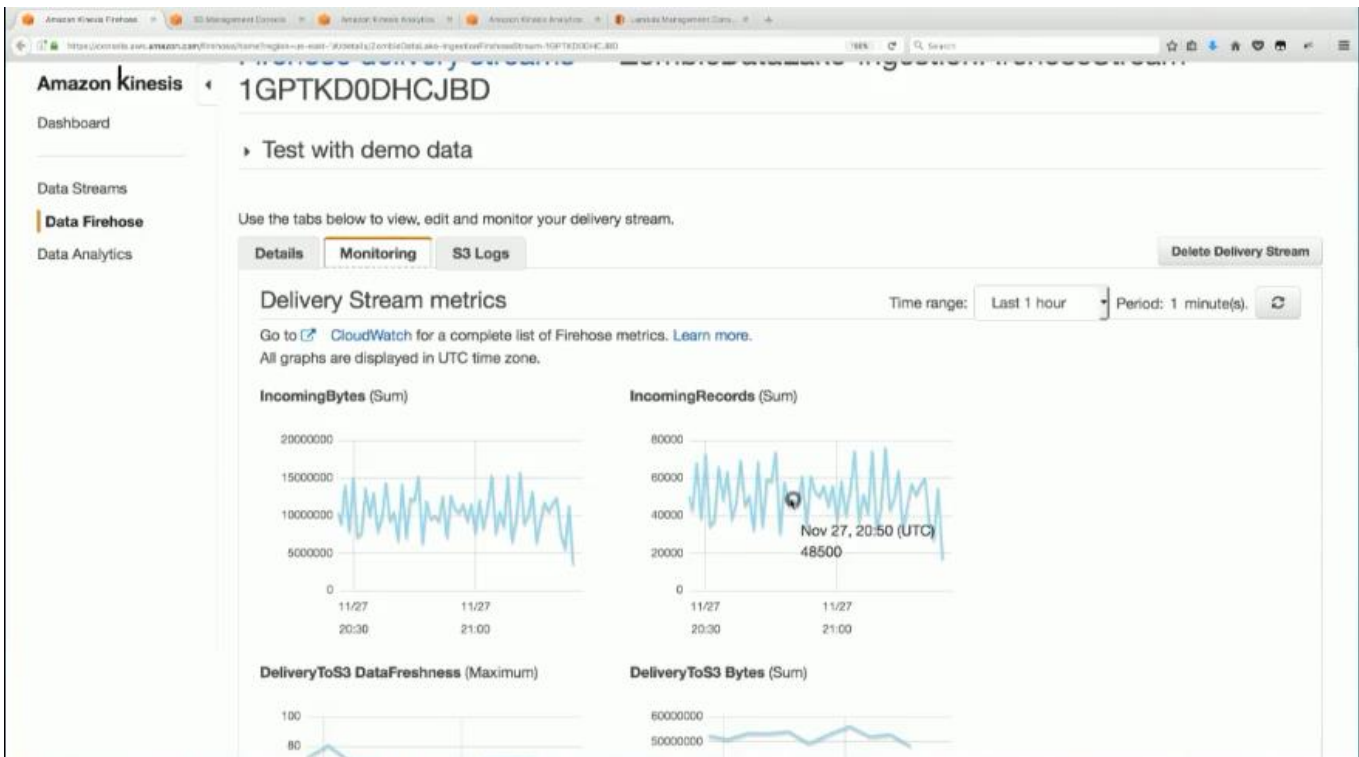
IAM role*  ZombieDataLake-IngestionFirehoseRole-13LC2RV1TNLD9

Data transformation*  Enabled

We can easily query this data using Athena or Kinesis Analytics as shown below



We write a SQL query, we have a SQL statement that does our real time analytics and start getting results

This defines the output using a sliding window



We are now generating the data into the stream, this is doing a 30 second window of data

Add and run SQL queries to continuously analyze source data in real-time. Then, optionally, connect the in-application stream to a destination to deliver results.

**Add SQL from templates**                              **Download SQL**

```
 9      MESSAGE_COUNT  INTEGER);
10
11   CREATE OR REPLACE PUMP "STREAM_PUMP" AS INSERT INTO "DESTINATION_SQL_STREAM"
12   SELECT STREAM STEP("SOURCE_SQL_STREAM_001".ROWTIME BY INTERVAL '30' SECOND) AS "ingest_time",
13       "country", "latitude", "longitude",
14       COUNT(CASE WHEN "sentiment" = 'negative' THEN 1 ELSE NULL END) as NEGATIVE_COUNT,
15       COUNT(CASE WHEN "sentiment" <> 'negative' THEN 1 ELSE NULL END) as POSITIVE_COUNT,
16       COUNT(*) AS MESSAGE_COUNT
17   FROM "SOURCE_SQL_STREAM_001"
18   GROUP BY "country", "latitude", "longitude",
19       STEP("SOURCE_SQL_STREAM_001".ROWTIME BY INTERVAL '30' SECOND);
20
```

**Exit (done editing)**     **Save and run SQL**

| Source data | Real-time analytics | Destination |     **Application status:** RUNNING

ZombieDataLake-
IngestionFirehoseStream-                     **Refresh stream sample**    **Download CSV**
1GPTKD0DHCJBD:

SOURCE_SQL_STREAM_001          **▼** Filter by column name          **Edit schema**

| | | | | | |
|---|---|---|---|---|---|
| 2017-11-27 21:18:55.077 | 0.6124 | alysha.oconnell | positive | 1511719148984 | 69. |
| 2017-11-27 21:18:55.077 | 0.0 | jimmie.dicki | neutral | 1511717289084 | 69. |
| 2017-11-27 21:18:55.077 | 0.0 | nathen.towne | neutral | 1511716561008 | 69. |
| 2017-11-27 21:18:55.077 | 0.0 | jerrod.harber | neutral | 1511718789284 | 69. |
| 2017-11-27 21:18:55.077 | 0.0 | jamar.herman | neutral | 1511720289384 | 69. |

---

Data Firehose
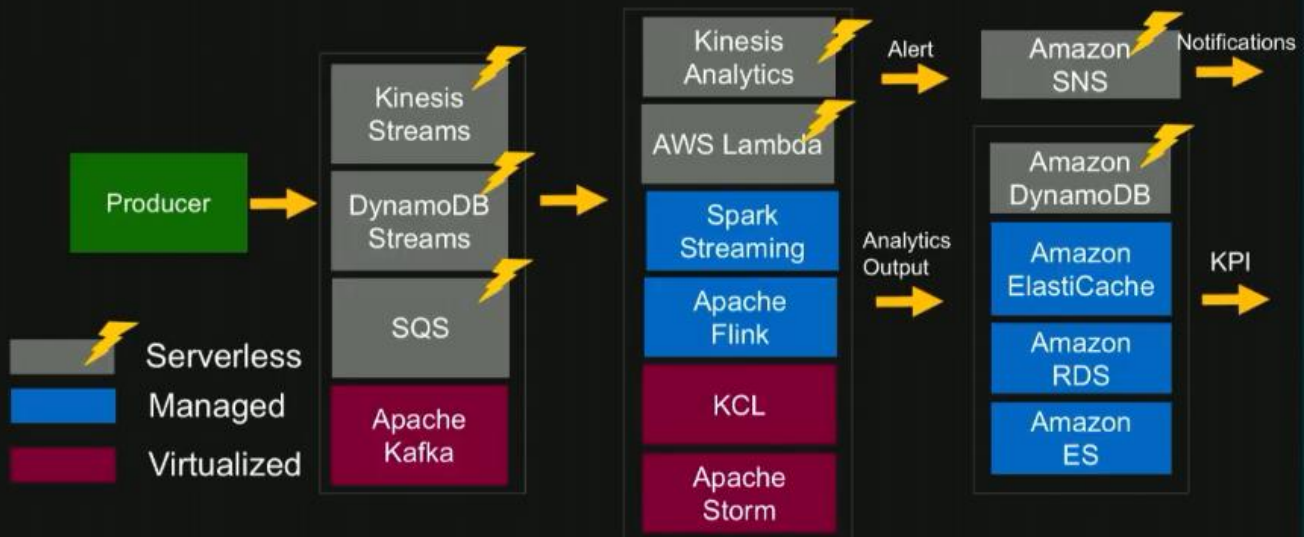
**Data Analytics**

```
 9      MESSAGE_COUNT  INTEGER);
10
11   CREATE OR REPLACE PUMP "STREAM_PUMP" AS INSERT INTO "DESTINATION_SQL_STREAM"
12   SELECT STREAM STEP("SOURCE_SQL_STREAM_001".ROWTIME BY INTERVAL '30' SECOND) AS "ingest_time",
13       "country", "latitude", "longitude",
14       COUNT(CASE WHEN "sentiment" = 'negative' THEN 1 ELSE NULL END) as NEGATIVE_COUNT,
15       COUNT(CASE WHEN "sentiment" <> 'negative' THEN 1 ELSE NULL END) as POSITIVE_COUNT,
16       COUNT(*) AS MESSAGE_COUNT
17   FROM "SOURCE_SQL_STREAM_001"
18   GROUP BY "country", "latitude", "longitude",
19       STEP("SOURCE_SQL_STREAM_001".ROWTIME BY INTERVAL '30' SECOND);
20
```

**Exit (done editing)**     **Save and run SQL**

| Source data | Real-time analytics | Destination |     **Application status:** RUNNING

In-application streams:        **Pause results**  ⟳ New results are added every 2-10 seconds. The results below are sampled. ⓘ

DESTINATION_SQL_STREAM          ☐  Scroll to bottom when new results arrive.

error_stream                   **▼** Filter by column name

| ROWTIME | INGEST_TIME | COUNTRY | LATITUDE | LONGITUDE |
|---|---|---|---|---|
| 2017-11-27 21:19:00.0 | 2017-11-27 21:18:30.0 | Pakistan | 30.375320434570316 | 69.345115661 |
| 2017-11-27 21:19:30.0 | 2017-11-27 21:19:00.0 | Australia | -25.274398803710938 | 133.77513122 |
| 2017-11-27 21:19:30.0 | 2017-11-27 21:19:00.0 | Uzbekistan | 41.37749099731445 | 64.585258483 |
| 2017-11-27 21:19:30.0 | 2017-11-27 21:19:00.0 | Saint Lucia | 13.909443855285645 | -60.97889328 |
| 2017-11-27 21:19:30.0 | 2017-11-27 21:19:00.0 | Niger | 17.6077880859375 | 8.0816659927 |

We can specify a destination like S3 for the analyzed result

## Source

Connect to an existing Kinesis stream or Firehose delivery stream, or easily create and connect to a new demo Kinesis stream. The limit is one streaming source for each application. Learn more.

| | Source | In-application stream name | ID ⓘ | Record pre-processing ⓘ |
|---|---|---|---|---|
| ✏ | Firehose delivery stream ZombieDataLake-IngestionFirehoseStream-1GPTKD0DHCJBD ☐ | SOURCE_SQL_STREAM_001 | 2.1 | ZombieDataLake-DataTransformationLambda-1DM4YR |

## Real time analytics

Continuously analyzing your source data with SQL. Learn more

**Go to SQL results**

## Destination

(Optional) Connect an in-application stream to a Kinesis stream, or to a Firehose delivery stream, to continuously deliver SQL results to AWS destinations. The limit is three destinations for each application.

**Connect new destination**     Disconnect destination

| | | Destination | In-application stream name | ID ⓘ |
|---|---|---|---|---|
| ☐ | ✏ | Firehose delivery stream ZombieAnnihilationSentimentAggregationStream ☐ | DESTINATION_SQL_STREAM | 4.1 |

---



aws    Services ▾    Resource Groups ▾    ✦          🔔 Ben Snively ▾    N. Virginia ▾    Support ▾

## AWS Lambda    ✕

Dashboard
**Functions**

Lambda > Functions > ZombieDataLake-DataTransformationLambda-1DM4YRG5RMA27 > $LATEST

ARN - arn:aws:lambda:us-east-1:783526147575:function:ZombieDataLake-DataTransformationLambda-1DM4YRG5RMA27:$LATEST

# ZombieDataLake-
# DataTransformationLambda-1DM4YRG5RMA27:$LATEST

| Version: $LATEST ▾ | Actions ▾ | Select a test event.. ▾ | **Test** |

ⓘ You are now viewing version $LATEST and associated code, config, and event sources.    ✕

**Configuration**    Triggers    Monitoring

▼ **Function code**

⚠ This function contains external libraries. Uploading a new file will override these libraries.    ✕

| Code entry type | Runtime | Handler  Info |
|---|---|---|
| Edit code inline ▾ | Python 2.7 ▾ | index.lambda_handler |

We can run some custom lambda code to transform or enrich the data coming into the stream before storing it

```
2017-11-27 16:25:00.000,Fiji,-16.578193664550/8,179.41441345214844,1,249,250
2017-11-27 16:25:00.000,India,20.59368324279785,78.96288299560547,5,245,250
2017-11-27 16:25:00.000,Japan,36.2048225402832,138.2529296875,5,245,250
2017-11-27 16:25:00.000,Saint Lucia,13.909443855285645,-60.978893280029304,5,245,250
2017-11-27 16:25:00.000,Turkey,38.9637451171875,35.24332046508789,2,248,250
2017-11-27 16:25:00.000,Burkina Faso,12.238332748413088,-1.5615930557250977,2,248,250
2017-11-27 16:25:00.000,Cambodia,12.565678596496584,104.990966796875,1,249,250
2017-11-27 16:25:00.000,Spain,40.46366882324219,-3.74921989440918,2,248,250
2017-11-27 16:25:00.000,Australia,-25.274398803710938,133.77513122558594,0,250,250
2017-11-27 16:25:00.000,Niger,17.6077880859375,8.081665992736816,0,250,250
2017-11-27 16:25:00.000,Uruguay,-32.52277755737305,-55.76583480834961,0,250,250
2017-11-27 16:25:00.000,Cayman Islands,19.513469696044922,-80.56695556640625,1,249,250
2017-11-27 16:25:00.000,El Salvador,13.79418468475342,-88.89653015136719,10,240,250
2017-11-27 16:25:00.000,Ethiopia,9.145000457763672,40.48967361450195,1,249,250
2017-11-27 16:25:00.000,Costa Rica,9.748916625976562,-83.75342559814453,2,248,250
2017-11-27 16:25:00.000,Bhutan,27.514162063598633,90.43360137939452,0,250,250
2017-11-27 16:25:00.000,Martinique,14.641528129577637,-61.024173736572266,7,243,250
2017-11-27 16:25:00.000,Egypt,26.820552825927738,30.80249786376953,1,249,250
2017-11-27 16:25:00.000,Nepal,28.39485740661621,84.12400817871094,4,246,250
2017-11-27 16:25:00.000,Togo,8.619543075561523,0.8247820138931273,3,247,250
2017-11-27 16:25:00.000,Portugal,39.399871826171875,-8.224453926086426,1,249,250
2017-11-27 16:25:00.000,Chile,-35.675148010253906,-71.54296875,0,250,250
2017-11-27 16:25:00.000,Bahrain,25.9304141998291,50.63777160644531,4,246,250
2017-11-27 16:25:00.000,Madagascar,-18.76694679260254,46.86910629272461,0,250,250
2017-11-27 16:25:00.000,Greece,39.0742073059082,21.82431221008301,1,249,250
2017-11-27 16:25:00.000,Mauritius,-20.348403930664066,57.55215072631836,10,240,250
2017-11-27 16:25:00.000,Montserrat,16.74249839782715,-62.1873664855957,6,244,250
2017-11-27 16:25:00.000,Mozambique,-18.66569519042969,35.529563903808594,0,250,250
2017-11-27 16:25:00.000,Norway,60.472023010253906,8.46894645690918,0,250,250
2017-11-27 16:25:00.000,Venezuela,6.423749923706056,-66.58972930908203,0,250,250
2017-11-27 16:25:00.000,Saint Kitts and ,17.35782241821289,-62.78299713134766,6,244,250
2017 11 27 16:25:00 000 Qatar 25.254825072510746 51.1828026660210 0 250 250
```

This is what the analyzed result looks like in S3

# Lambda Pre-processing



Transform, Enrich, Filter

# Amazon Kinesis Analytics - SQL

In this example:
19 Lines of SQL = Serverless Realtime Analytics

```sql
11  CREATE OR REPLACE PUMP "STREAM_PUMP" AS INSERT INTO "DESTINATION_SQL_STREAM"
12  SELECT STREAM STEP("SOURCE_SQL_STREAM_001".ROWTIME BY INTERVAL '30' SECOND) AS "ingest_time",
13      "country", "latitude", "longitude",
14      COUNT(CASE WHEN "sentiment" = 'negative' THEN 1 ELSE NULL END) as NEGATIVE_COUNT,
15      COUNT(CASE WHEN "sentiment" <> 'negative' THEN 1 ELSE NULL END) as POSITIVE_COUNT,
16      COUNT(*) AS MESSAGE_COUNT
17  FROM "SOURCE_SQL_STREAM_001"
18  GROUP BY "country", "latitude", "longitude",
19      STEP("SOURCE_SQL_STREAM_001".ROWTIME BY INTERVAL '30' SECOND);
```

# Larger Picture of what we showed:



# Fitting into existing Real-time Analytics

Components of Big Data Applications



Serverless Interactive Analytics



Demonstration

Let us now do some interactive analytics on some of our data at rest in S3.

We are starting with the Glue catalog to help describe our data on S3



There are actually several data type available in the S3 location, we can filter the data by type

We can use Glue to crawl a new S3 location, the crawlers run with an IAM role that we set up for the crawler to allow it crawl our data on S3 to discover all the new tables in there

We now have several tables created for us by the crawler with the data sets discovered



This is the twitter data it discovered, we can now use this in our analytics ad queries

It also discovered the JSON that can represent the twitter data, we can modify this schema as we like

Now we can use Athena to start querying the data using SQL to analyze this dataset using SELECT commands.

We can see that we have over 10 million records available for this month

The results come back very fast

# Demonstration

# Interactive Analytics



| Producer | → | Amazon S3 | → | Amazon Athena / Amazon EMR (Presto, Impala, Spark) / Amazon Redshift | → | QuickSight / kibana / IPython Interactive Computing |

Legend:
- Serverless
- Managed
- Virtualized

# Components of Big Data Applications



Glue Data Catalog

Integrated Pipeline

- Kinesis / S3 — Ingest & Store
- AWS Lambda / AWS Glue — Prepare & Transform
- Kinesis Analytics / Amazon Athena — Analyze & Reason
- QuickSight — Access & User Interface
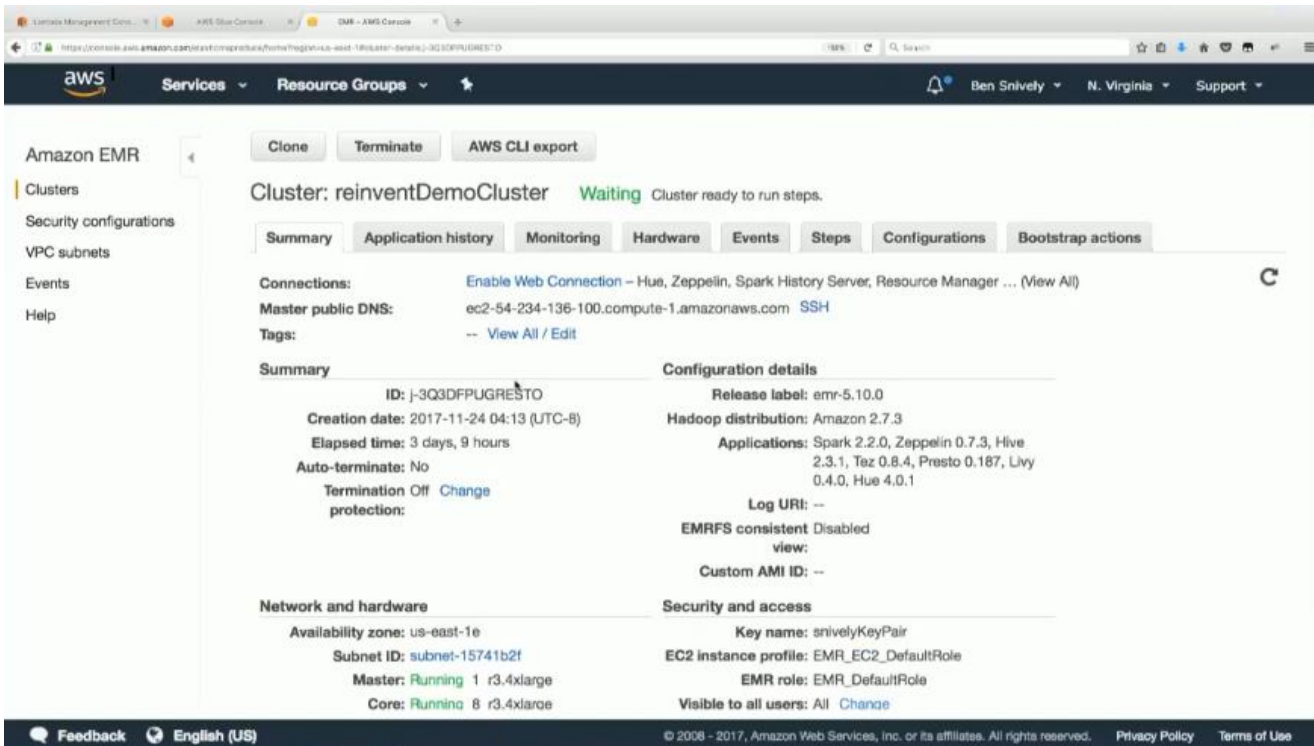
Protect & Secure
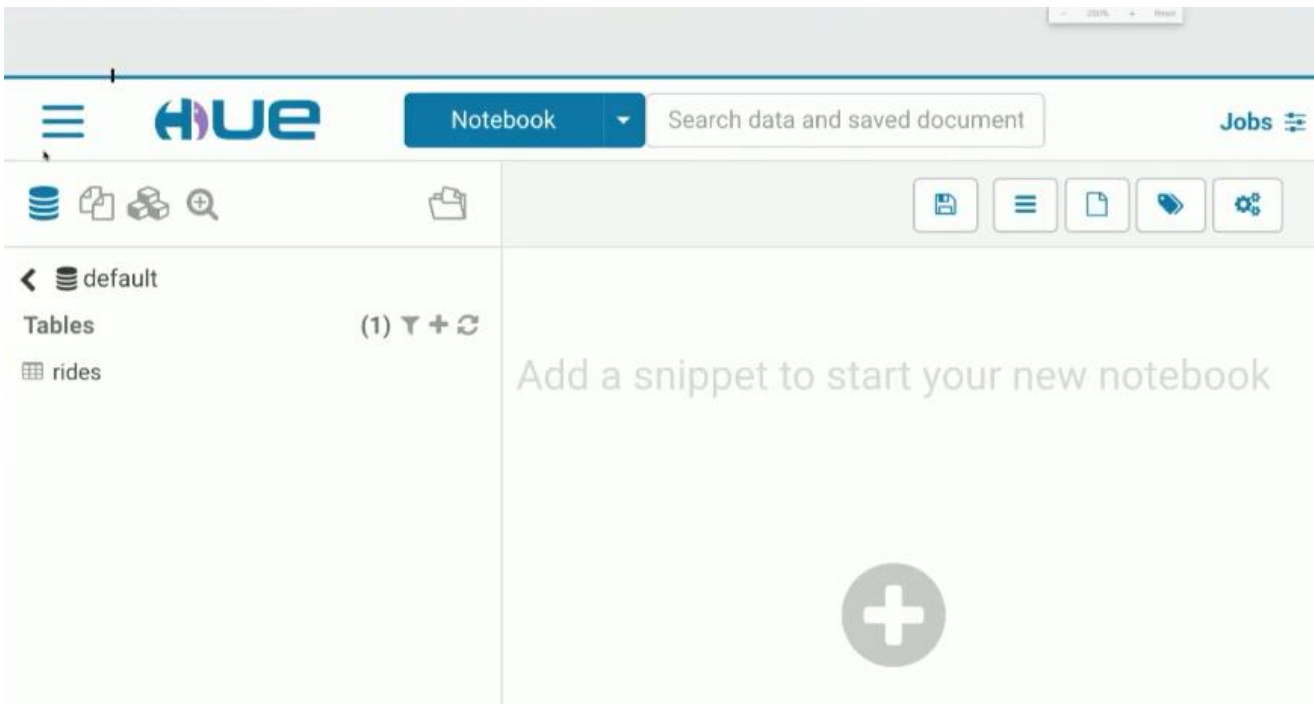
# Data Lake reference architecture



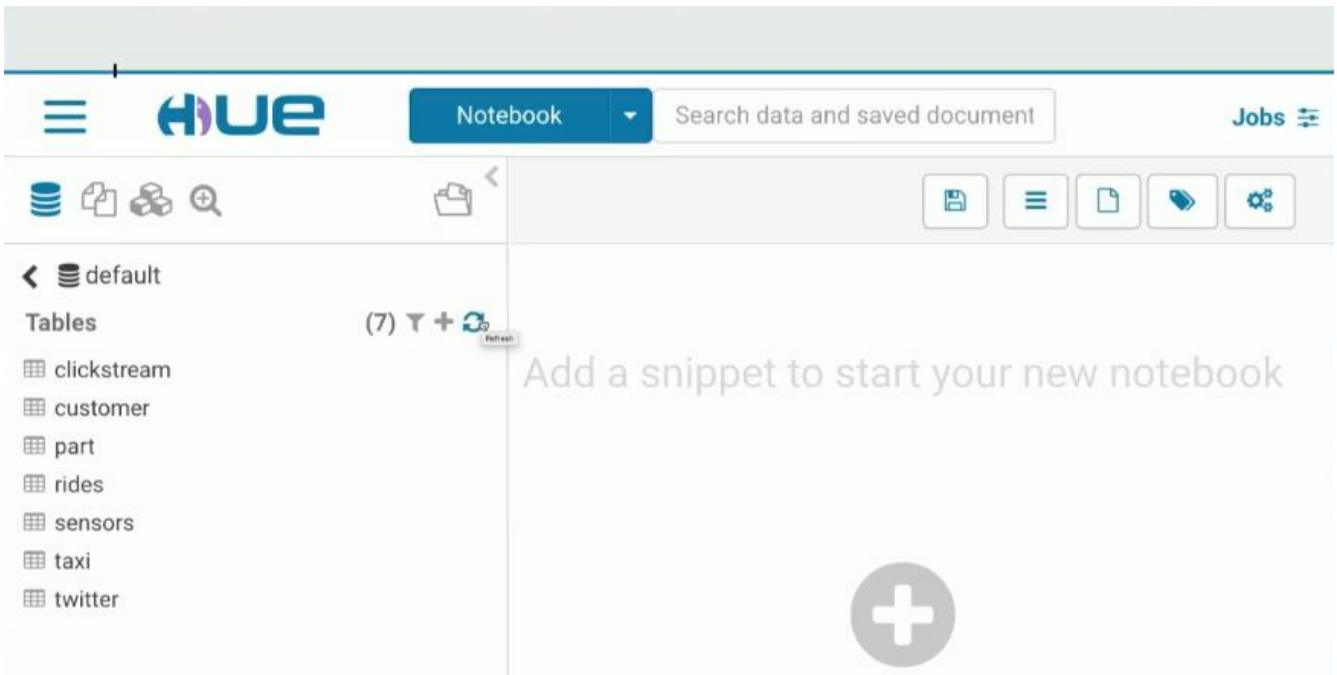# The Right tool for the Right Job



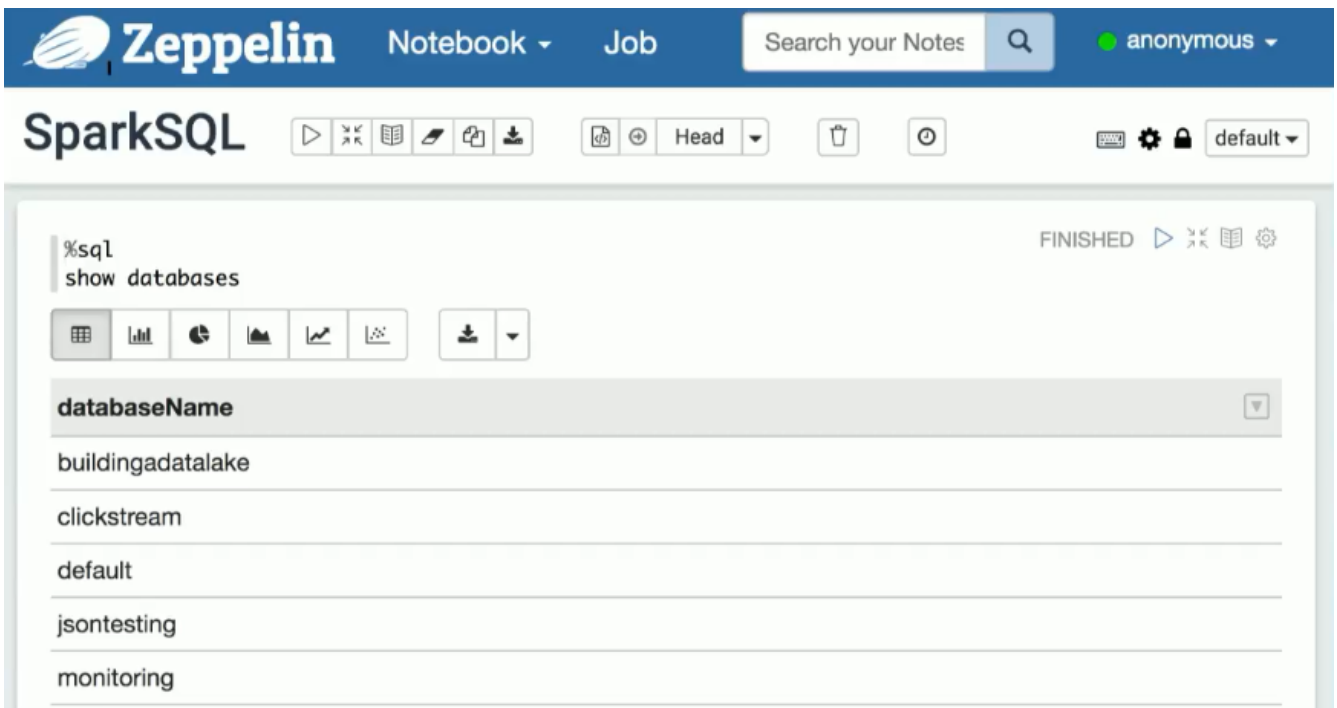# What about existing Hadoop Clusters?

We have our EMR cluster running here, it has several of the ecosystem applications on it, we also specified that it use the Glue data catalog.



Hue is an interactive web interface for interacting with your Hadoop cluster.

These are the data tables we discovered with the Glue crawler



We can also have our data scientist using Zeppelin notebooks running on the EMR cluster look at the data

## What about existing Hadoop Clusters?

## The Right tool for the Right Job

## Serverless nicely fits into big data platforms

- Mix and Match Serverless, Managed, and Virtualized Services
- Leverage Services to easily
    - Rapidly ingest, categorize, and discover your data
    - Allow easy query and analysis of your data
    - Transform and Load data
    - Provide custom event based handlers
- Serverless allows you to focuses more analytics and not on infrastructure or servers

- Pay only for what you use