

Lab: Real-Time Clickstream Anomaly Detection Kinesis Analytics

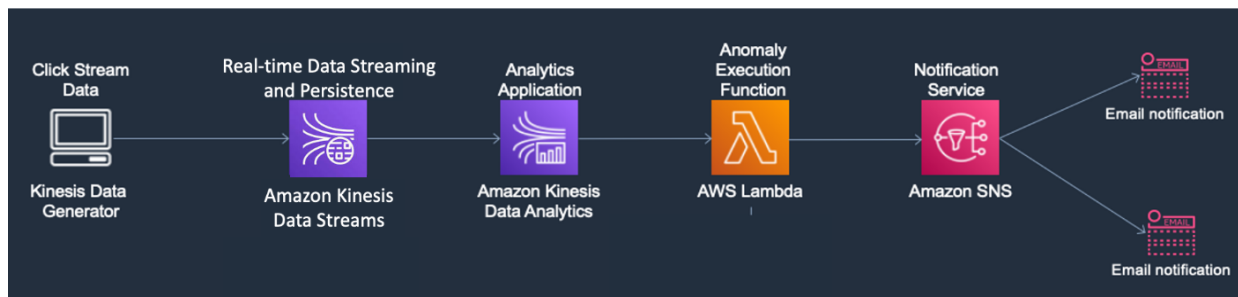
Introduction

This guide helps you complete Real-Time Clickstream Anomaly Detection using Amazon Kinesis Data Analytics. Analyzing web log traffic to gain insights that drive business decisions has historically been performed using batch processing. Although effective, this approach results in delayed responses to emerging trends and user activities. There are solutions that process data in real-time using streaming and micro-batching technologies, but they can be complex to set up and maintain. Amazon Kinesis Data Analytics is a managed service that makes it easy to identify and respond to changes in data behavior in real-time.

Steps:

- [Set up an Amazon Analytics Studio Application through CloudFormation stack deployment](#)
- [Generate real time website traffic using Amazon Kinesis Data Generator \(KDG\)](#)
- [Perform real-time Data Analytics](#)
- [Environment Cleanup](#)
- [Appendix: Anomaly Detection Scripts](#)

In the [Kinesis prelab setup](#), you fulfilled the prerequisites for this lab. In this lab, you will create the following Kinesis Data Analytics pipeline.



Set up an Amazon Analytics Studio Application through CloudFormation stack deployment

1. Click the [Deploy to AWS](#) link here to stand up the pre-lab workshop infrastructure
2. The button above will open a “Quick create stack” form, please accept the default parameters, select the checkbox to acknowledge new IAM role creation and select Create Stack to run the Amazon CloudFormation Template

Lab: Real-Time Clickstream Anomaly Detection Kinesis Analytics

Quick create stack

Template

Template URL:
https://kda-flink-pre-lab-cfn-template.s3.amazonaws.com/1-Lab-Kinesis-Clickstream_CFN.yaml

Stack description:
Supporting elements for the Kinesis Analytics click stream lab with Flink

Stack name

Stack name:
kda-flink-pre-lab
Stack names can include letters (a-z and A-Z), numbers (0-9), and dashes (-).

Parameters

Parameters are defined in your template and allow you to input custom values when you create or update a stack.

Kinesis Flink Pre Lab set up


Firehose version:
1.15.4

GitOrganizationName:
aws-samples

Release:
master

CloudWatchLogGroupName:
prelab_data_db

Capabilities


 The following resource(s) require capabilities: [AWS-IAM:elastic]. This template contains Identity and Access Management (IAM) resources that might provide entities access to make changes to your AWS account. Check that you want to create each of these resources and that they have the minimum required permissions. [Learn more](#)

☒ I acknowledge that AWS CloudFormation might create IAM resources.

[Cancel](#) [Create change set](#) [Create stack](#)

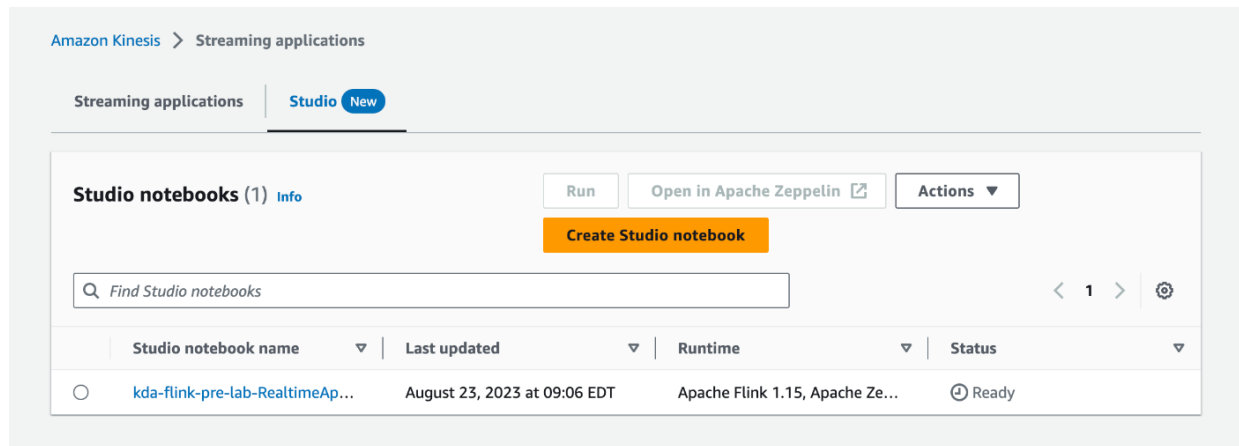
3. The stack will create six Amazon Kinesis Data Streams in the [Amazon Kinesis Console](#) -

- tickerstream – the raw stream to send the initial traffic
- clickstream – captures the number of clicks
- impressionstream – captures the number of impressions
- ctrstream – captures the calculated click through rate
- destinationstream – captures the anomaly scores
- anomalydetectionstream – captures the records with anomaly score greater than 1

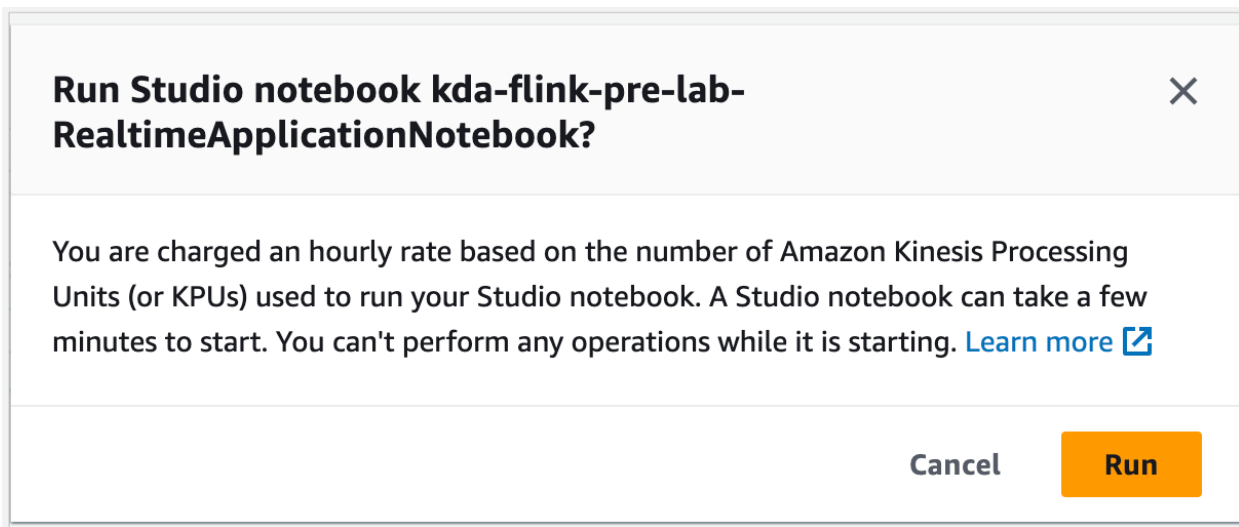
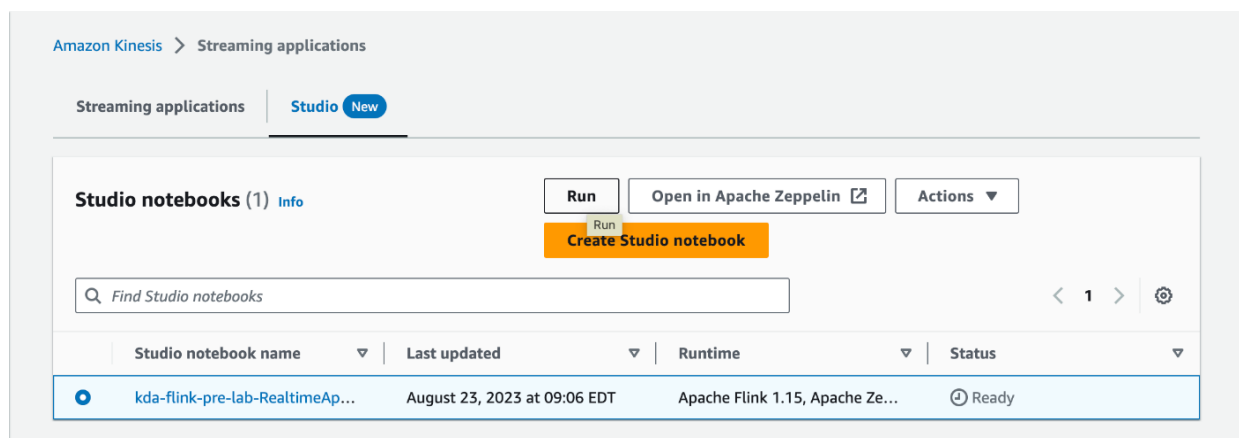
Data streams (8) Info							
Process data in real time		Create a Firehose delivery stream		Actions ▼		Create data stream	
<input type="text" value="Q :"/>		8 streams		< 1 >			
<input type="checkbox"/>	Name ▲	Status ▼	Capacity mode ▼	Provisioned shards ▼	Data retention period ▼	Encryption ▼	Consumers with enhanced fan-out
<input type="checkbox"/>	anomalydetectionstream	✔ Active	Provisioned	1	1 day	Disabled	0
<input type="checkbox"/>	clickstream	✔ Active	Provisioned	1	1 day	Disabled	0
<input type="checkbox"/>	ctrstream	✔ Active	Provisioned	1	1 day	Disabled	0
<input type="checkbox"/>	destinationstream	✔ Active	Provisioned	1	1 day	Disabled	0
<input type="checkbox"/>	impressionstream	✔ Active	Provisioned	1	1 day	Disabled	0
<input type="checkbox"/>	tickerstream	✔ Active	Provisioned	1	1 day	Disabled	0

4. The template would also create a Amazon Kinesis Data Analytics Studio application called kda-flink-prelab-RealtimeApplicationNotebook in the [Amazon Kinesis Console](#) → Studio tab. We will write interactive Studio Notebook in Apache Zeppelin for real time data analysis.

Lab: Real-Time Clickstream Anomaly Detection Kinesis Analytics



5. Run the Studio Application by selecting the kda-flink-prelab-RealtimeApplicationNotebook under Studio tab. Select “Run” again on the next screen.



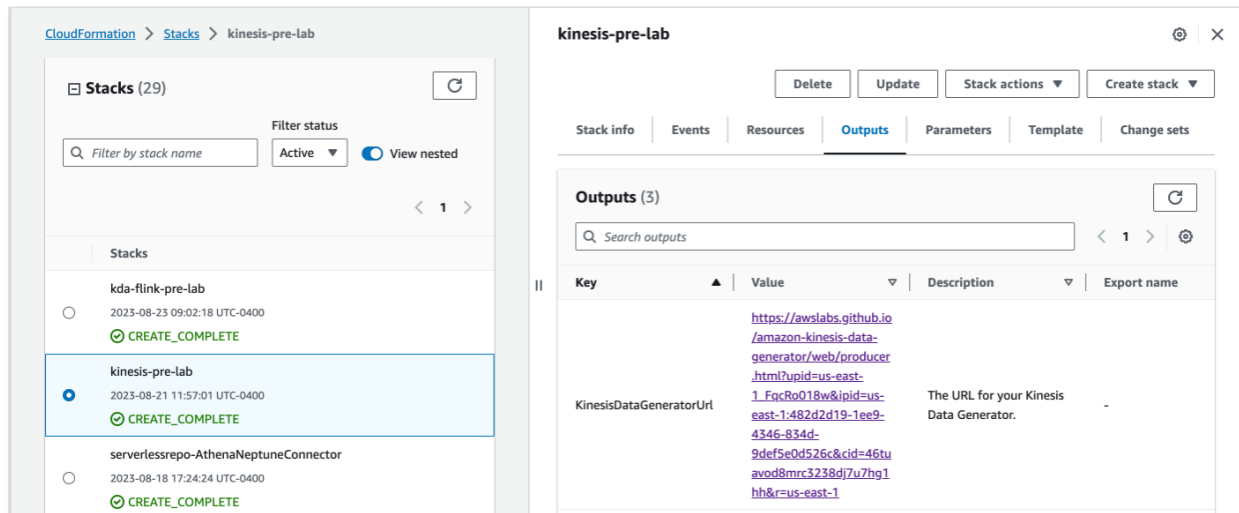
Generate real time website traffic using Amazon Kinesis Data Generator (KDG)

1. Navigate to the [Amazon CloudFormation console](#) in your AWS account, click on the **Kinesis-pre-lab** stack created

Lab: Real-Time Clickstream Anomaly Detection Kinesis Analytics

during the [Streaming Data Analytics Prelab setup](#).

2. Go to the Outputs tab of the stack to get the Kinesis Data Generator link as shown below:



The screenshot shows the AWS CloudFormation console. On the left, the 'Stacks' list shows three stacks: 'kda-flink-pre-lab', 'kinesis-pre-lab' (selected), and 'serverlessrepo-AthenaNeptuneConnector'. The 'kinesis-pre-lab' stack is in the 'CREATE_COMPLETE' state. On the right, the 'Outputs' tab for the 'kinesis-pre-lab' stack is displayed. It shows a single output named 'KinesisDataGeneratorUrl' with a long URL value. The URL is: https://awslabs.github.io/amazon-kinesis-data-generator/web/producer.html?upid=us-east-1_FqcRo018w&ipid=us-east-1:482d2d19-1ee9-4346-834d-9def5e0d526c&cid=46tuavod8mrc3238dj7u7hg1hh&r=us-east-1. The description for this output is 'The URL for your Kinesis Data Generator.'

3. Open two concurrent sessions of the KDG UI in your browser. Sign in using the username and password you entered in the CloudFormation template while creating the [Streaming Data Analytics Prelab setup](#).
 - a. We want to generate more Click messages than Impressions. So, in the first session, send impression messages at rate of one message per second for 30 seconds to the tickerstream, the message body is

```
{"browseraction": "Impression", "site": "https://www.mysite.com"}
```

Lab: Real-Time Clickstream Anomaly Detection Kinesis Analytics

Amazon Kinesis Data Generator

[Configure](#) [Help](#) [Log Out](#)

Region

us-east-1

Stream/delivery stream

tickerstream

Records per second

Constant

Periodic

5

Compress Records [?](#)

☐

Record template [?](#)

Impression Payload

Template 2

Template 3

Template 4

Template 5

Impression Payload

```
{"browseraction": "Impression", "site": "https://www.mysite.com"}
```

Send data

Test template

b. Then in the second session, send click messages at a rate of five messages per second for 30 seconds to the tickerstream, the message body is

```
{"browseraction": "Click", "site": "https://www.mysite.com"}
```

Lab: Real-Time Clickstream Anomaly Detection Kinesis Analytics

Amazon Kinesis Data Generator

ConfigureHelpLog Out

Region

us-east-1

Stream/delivery stream

tickerstream

Records per second

ConstantPeriodic

5

Compress Records ⓘ

☐

Record template ⓘ

Click PayloadTemplate 2Template 3Template 4Template 5

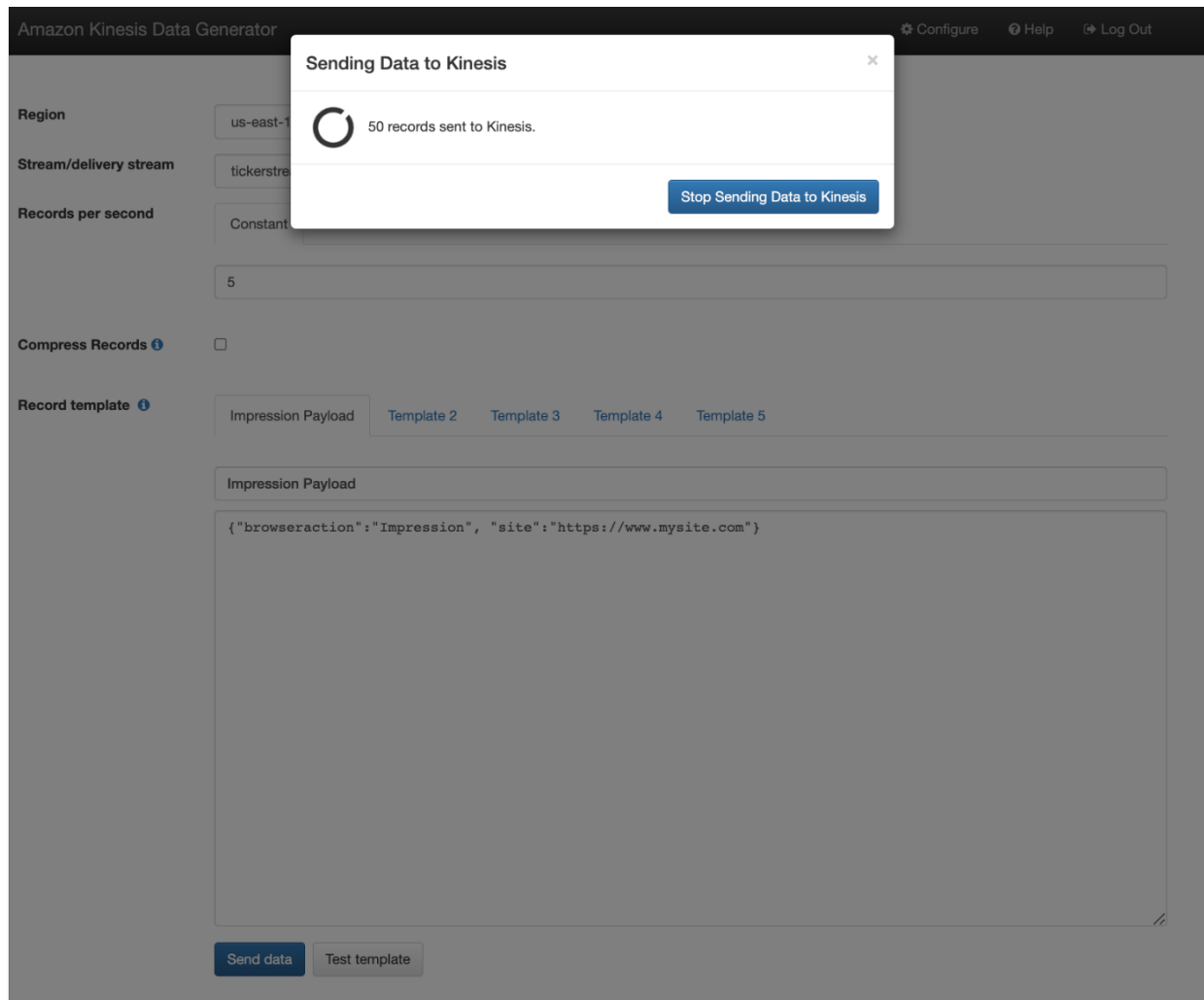
Click Payload

```
{"browseraction": "Click", "site": "https://www.mysite.com"}
```

Send dataTest template

c. You can view the number of messages being sent to the data stream when you start posting the messages

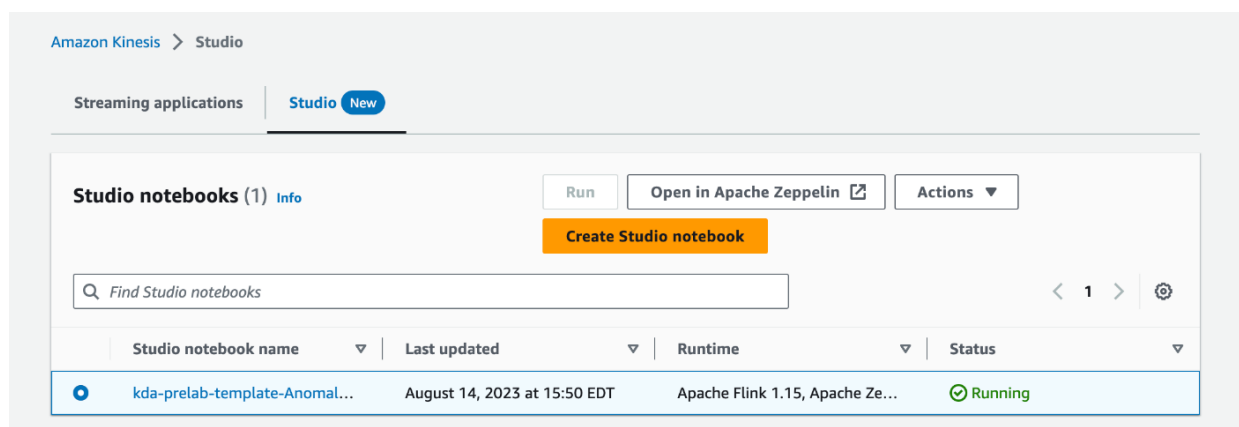
Lab: Real-Time Clickstream Anomaly Detection Kinesis Analytics



d. After 30 seconds, please stop sending both Click and Impression messages

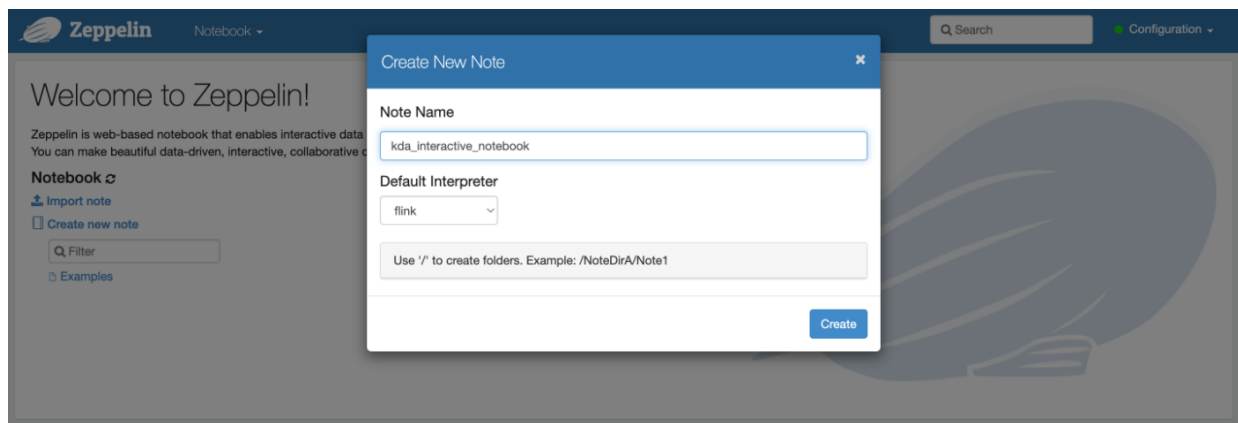
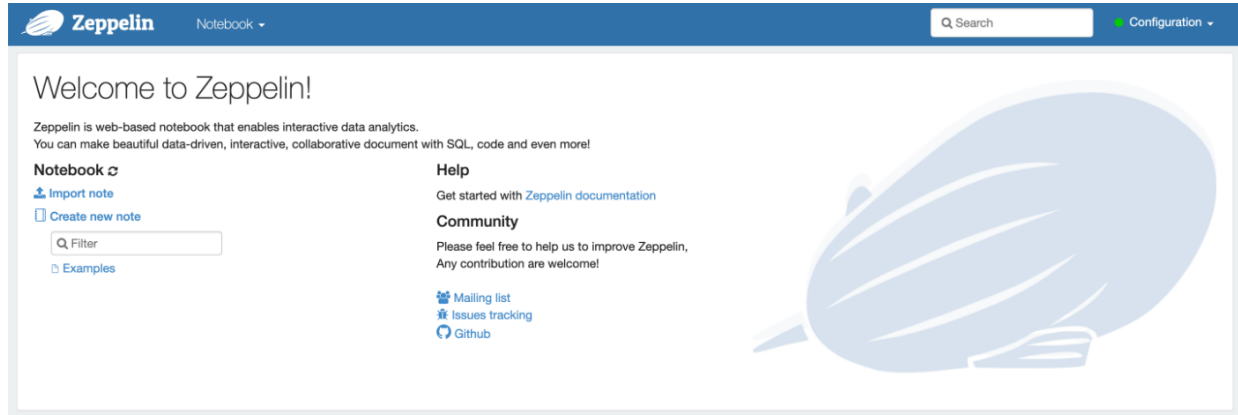
Perform real-time Data Analytics

1. Navigate to the [Amazon Kinesis Console](#) → Analytics Application. Under the Studio tabs, select `kda-prelab-template-RealtimeApplicationNotebook`. Select "Open in Apache Zeppelin".

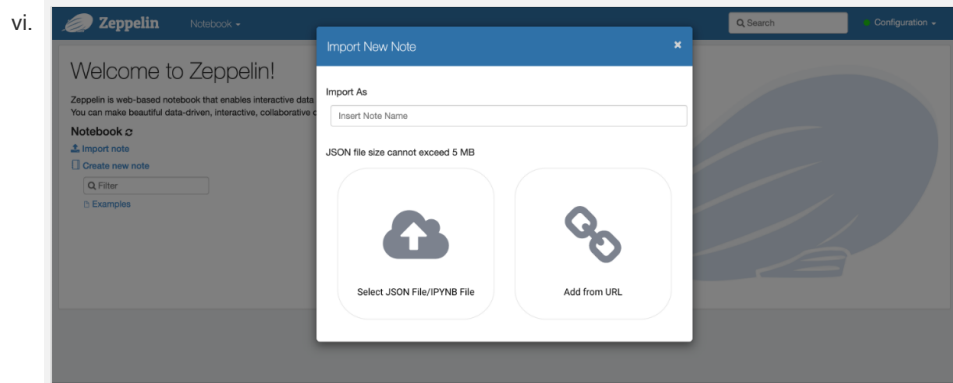


Lab: Real-Time Clickstream Anomaly Detection Kinesis Analytics

2. On the Apache Zeppelin Console, select Create new note. Provide the notebook name as **kda_interactive_notebook**

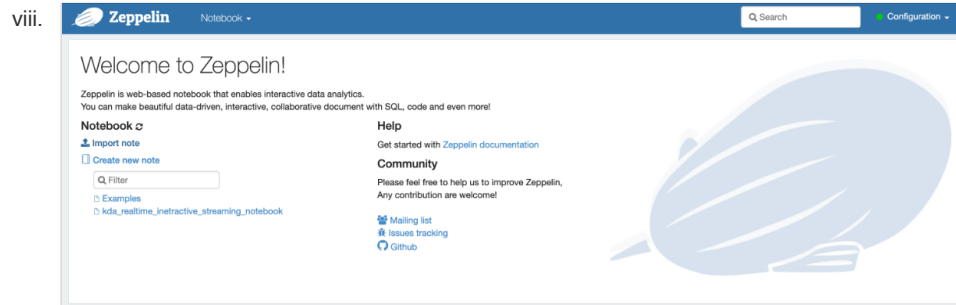


3. Now lets perform real time interactive analytics with Kinesis data streams. We will
 - i. Create Flink tables using Flink SQL Queries
 - ii. Use Flink SQL queries to transform and create new data streams in real-time
 - iii. Perform anomaly detection using Flink User Defined Function and trigger anomaly notification emails in real-time.
 - iv. The scripts are available [here](#).
 - v. A notebook is also available [here](#) which can be downloaded and imported through Apache Zeppelin console.

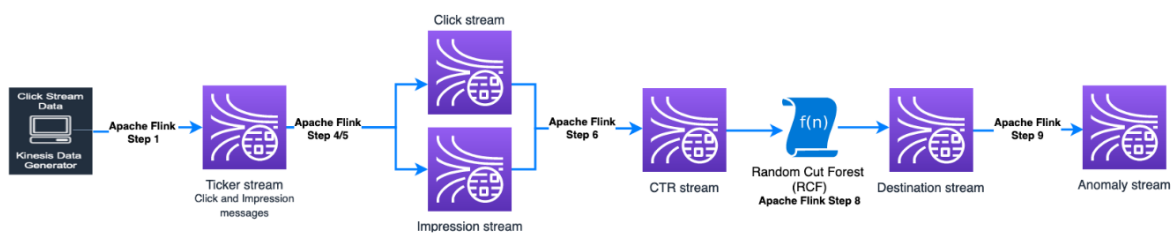


Lab: Real-Time Clickstream Anomaly Detection Kinesis Analytics

vii. You can then open the notebook and run the paragraphs one after the other.



Please refer to the anomaly detection logic in the image below before you start running the remaining steps



a. Run the table creation scripts.

```
#1 Create all the Streaming tables

%Flink.pyFlink
st_env.execute_sql("""DROP TEMPORARY TABLE IF EXISTS tmpclickstream""")
st_env.execute_sql("""DROP TEMPORARY TABLE IF EXISTS tmpimpressionstream""")
st_env.execute_sql("""DROP TEMPORARY TABLE IF EXISTS tmpctrstream""")
st_env.execute_sql("""DROP TEMPORARY TABLE IF EXISTS tmpdestinationstream""")
st_env.execute_sql("""DROP TEMPORARY TABLE IF EXISTS tmpanomalydetectionstream""")
st_env.execute_sql("""CREATE TEMPORARY TABLE tmpclickstream (
  'browseraction' STRING,
  'site' STRING,
  'rowtime' TIMESTAMP(3) METADATA FROM 'timestamp',
  WATERMARK FOR rowtime AS rowtime - INTERVAL '0' SECOND
)
WITH (
  'connector' = 'kinesis',
  'stream' = 'clicker',
  'aws.region' = 'us-east-1',
  'scan.stream.initpos' = 'TRIM_HORIZON',
  'format' = 'json'
)""")
```

b. User Defined Function (UDF) performs Anomaly Detection in real-time using Random Cut Forest algorithm. Run step #2.

```
#2: Create the UDF (User Defined Function)

%Flink(parallelism=1)
import software.amazon.flink.example.RandomCutForestUDF
stenv.registerFunction("RANDOM_CUT_FOREST", new RandomCutForestUDF())

import software.amazon.flink.example.RandomCutForestUDF
```

c. You can view real time data from website visits by running the query in Step #3.

Lab: Real-Time Clickstream Anomaly Detection Kinesis Analytics

#3: View the Ticker Stream Data in real-time FLINK JOB ABORT

```
%flink.ssql(type=update)
select * from tmpclickerstream;
```

settings

browseraction	site	rowtime
Click	https://www.mysite.com	2023-08-23 07:30:11.812
Click	https://www.mysite.com	2023-08-23 07:30:13.007
Click	https://www.mysite.com	2023-08-23 07:30:13.009
Click	https://www.mysite.com	2023-08-23 07:30:13.009
Click	https://www.mysite.com	2023-08-23 07:30:13.009
Click	https://www.mysite.com	2023-08-23 07:30:13.009
Click	https://www.mysite.com	2023-08-23 07:30:13.982
Click	https://www.mysite.com	2023-08-23 07:30:13.984
Click	https://www.mysite.com	2023-08-23 07:30:13.984

Note: Use %flink.ssql while querying the table

- d. Create impressionstream by filtering messages from tickerstream.

#4: Generate the Impression Stream Data in real-time based on the browser action value in the Ticker stream FLINK JOB ERROR

```
%flink.ssql(type=update, parallelism=1)
INSERT INTO tmpimpressionstream
SELECT
  TUMBLE_START(rowtime, INTERVAL '10' SECOND) AS rowtime,
  COUNT(*) AS impressioncount
FROM tmpclickerstream
WHERE browseraction='Impression'
GROUP BY
  TUMBLE(rowtime, INTERVAL '10' SECOND);
```

- e. Create clickstream by filtering messages from tickerstream.

#5: Generate the Click Stream Data in real-time based on the browser action value in the Ticker stream FLINK JOB ERROR

```
%flink.ssql(type=update, parallelism=1)
INSERT INTO tmpclickstream
SELECT
  TUMBLE_START(rowtime, INTERVAL '10' SECOND) AS rowtime,
  COUNT(*) AS clickcount
FROM tmpclickerstream
WHERE browseraction='Click'
GROUP BY
  TUMBLE(rowtime, INTERVAL '10' SECOND);
```

- f. Calculate Click Through Rate (CTR) and populate ctrstream.

#10: View and Analyze the Anomalies in real-time FLINK JOB ABORT

```
%flink.ssql
select * from anomalydetectionstream;
```

settings

rowtime	ctrpercent	anomaly_score
2023-08-23 13:46:30.000000	875.0	6.530422766870707
2023-08-23 13:46:40.000000	500.0	6.723916774777987
2023-08-23 13:46:50.000000	500.0	6.530422766870706
2023-08-23 13:47:00.000000	500.0	6.675543272801167
2023-08-23 13:47:10.000000	500.0	6.869037280708447
2023-08-23 13:47:20.000000	500.0	6.869037280708447
2023-08-23 13:47:30.000000	166.7	6.3853022609402466
2023-08-23 13:47:50.000000	500.0	6.482049264893886

- g. You can view the Click Through Rate in real time by executing Step #7.

Lab: Real-Time Clickstream Anomaly Detection Kinesis Analytics

#7: Analyze the Click Through Rate in real-time

```
%flink.sql(type=update)
select * from ctrstream;
```

time	ctr
1692798390	875.0
1692798400	500.0
1692798410	500.0
1692798420	500.0
1692798430	500.0
1692798440	500.0
1692798450	166.7
1692798470	500.0

h. Use the UDF (Random Cut Forest) to generate anomaly scores.

#8: Generate the Anomaly Score in real-time with User Defined Function RANDOM_CUT_FOREST

```
%flink.sql(type=update, parallelism=1)
INSERT INTO tmpdestinationstream
SELECT time, ctr, RANDOM_CUT_FOREST(cast(ctr as float)) as score
FROM ctrstream;
```

i. Populate anomalydetectionstream by executing Step #9.

#9: Create a Stream with the Anomaly data in real-time

```
%flink.sql(type=update, parallelism=1)
INSERT INTO tmpanomalydetectionstream
SELECT to_timestamp(from_unixtime('time')) as rowtime, ctr as ctrpercent, score as anomaly_score
FROM destinationstream
WHERE score>2;
```

Now check the anomaly scores from Random Cut Forest algorithm in real time.

a. #10: View and Analyze the Anomalies in real-time

```
%flink.sql
select * from anomalydetectionstream;
```

rowtime	ctrpercent	anomaly_score
2023-08-23 13:46:30.000000	875.0	6.530422766870707
2023-08-23 13:46:40.000000	500.0	6.723916774777987
2023-08-23 13:46:50.000000	500.0	6.530422766870706
2023-08-23 13:47:00.000000	500.0	6.675543272801167
2023-08-23 13:47:10.000000	500.0	6.869037280708447
2023-08-23 13:47:20.000000	500.0	6.869037280708447
2023-08-23 13:47:30.000000	166.7	6.3853022609402466
2023-08-23 13:47:50.000000	500.0	6.482049264893886

b. You will start receiving notifications in your email when anomalies are detected:

c. **Anomaly Detected**



o ClkStrEv2 <no-reply@sns.amazonaws.com>

Today at 9:53 AM

To: [redacted]

Anomaly detected with a click through rate of 500% and an anomaly score of 6.482049264893886

--

If you wish to stop receiving notifications from this topic, please click or visit the link below to unsubscribe:

[redacted] > [https://sns.us-east-1.amazonaws.com/unsubscribe.html?SubscriptionArn=arn:aws:sns:us-east-1:\[redacted\]:ClickStreamEvent2:5aa1d647-6293-4037-8f03-535f0a7cfcff&Endpoint=\[redacted\]](https://sns.us-east-1.amazonaws.com/unsubscribe.html?SubscriptionArn=arn:aws:sns:us-east-1:[redacted]:ClickStreamEvent2:5aa1d647-6293-4037-8f03-535f0a7cfcff&Endpoint=[redacted])

Please do not reply directly to this email. If you have any questions or comments regarding this email, please contact us at <https://aws.amazon.com/support>

Lab: Real-Time Clickstream Anomaly Detection Kinesis Analytics

- d. If you do not receive the anomaly notification emails on your first attempt
 - i. Open the two concurrent sessions of KDG UI in your browser again.
 - ii. In the first session, send impression messages at rate of one message per second to the tickerstream, the message body is {"browseraction": "Impression", "site": "https://www.mysite.com"}
 - iii. In the second session, send click messages at a rate of five messages per second to the tickerstream, the message body is {"browseraction": "Click", "site": "https://www.mysite.com"}
 - iv. Stop sending messages after 30-40 seconds.
 - v. Now on the Apache Zeppelin Notebook, repeat steps 3 to 10 and you should start receiving email notifications from the second attempt.

Environment Clean Up

1. After completing the lab, click **Actions** → **Stop Application** to stop your application and avoid flood of SMS and e-mails messages.
2. If you would like to delete the entire resource stack, navigate to [Amazon CloudFormation Console](#) → Stacks, select **kda-flink-pre-lab** and click Delete to remove the stack. This step will clean up all the resources created by the stack earlier.

CloudFormation > Stacks

Stacks (29) Refresh Delete Update Stack actions ▼ Create stack ▼

Filter status: Active ▼ View nested < 1 > ⚙

Stack name	Status	Created time	Description
<input checked="" type="radio"/> kda-flink-pre-lab	✔ CREATE_COMPLETE	2023-08-23 11:33:39 UTC-0400	Supporting elements for the Kinesis Analytics click stream lab with Flink
<input type="radio"/> kinesis-pre-lab	✔ CREATE_COMPLETE	2023-08-21 11:57:01 UTC-0400	Supporting elements for the Kinesis Analytics click stream lab

2. Select Delete on the next screen

Delete stack? ×

Delete stack **kda-flink-pre-lab** permanently? This action cannot be undone.

Deleting this stack will delete all stack resources. Resources will be deleted according to their DeletionPolicy. [Learn more](#)

Cancel Delete

Lab: Real-Time Clickstream Anomaly Detection Kinesis Analytics

3. If Delete Stack operation fails due to non-empty S3 buckets, select Delete again. On the pop up, either retain the S3 buckets failing to delete or manually empty the buckets.

The screenshot shows the AWS CloudFormation console's 'Stacks' page. At the top, there are buttons for 'Delete', 'Update', 'Stack actions', and 'Create stack'. Below these is a search bar labeled 'Filter by stack name' and a 'Filter status' dropdown set to 'Active'. A 'View nested' toggle is also present. The main table lists two stacks:

Stack name	Status	Created time	Description
kda-flink-pre-lab	DELETED_FAILED	2023-08-23 11:33:39 UTC-0400	Supporting elements for the Kinesis Analytics click stream lab with Flink
kinesis-pre-lab	CREATE_COMPLETE	2023-08-21 11:57:01 UTC-0400	Supporting elements for the Kinesis Analytics click stream lab

4. This should pop up the list of the non-empty S3 buckets, either retain the S3 buckets failing to delete by hitting the checkbox or manually empty and remove those buckets. Hit "Delete" again once you have performed either of the previously mentioned steps, this time it should remove the stack completely.

The screenshot shows a 'Delete stack?' dialog box. It contains the following text:

Delete stack?

Deleting this stack will delete all stack resources. Resources will be deleted according to their DeletionPolicy. [Learn more](#)

You may retain resources that are failing to delete
This stack previously failed to delete because the following resources failed to delete. If you choose to retain resources, they will be skipped during this delete operation.

Resources to retain - optional
Selected resources will be skipped during the delete stack operation

☒ **ArtifactBucket**
[kda-flink-pre-lab-](#) [redacted] [\[link icon\]](#)

At the bottom, there are 'Cancel' and 'Delete' buttons.

Appendix: Anomaly Detection Scripts

1. Amazon CloudFormation Template: https://aws-bigdata-blog.s3.amazonaws.com/DE-ID-KDA-Lab/1-Lab-Kinesis-Clickstream_CFN.yaml
2. Apache Zeppelin Notebook: https://aws-bigdata-blog.s3.amazonaws.com/DE-ID-KDA-Lab/kda_realtime_inettractive_streaming_notebook_2J9H888TM.zpln
3. Amazon Random Cut Forest User Defined Function: <https://raw.githubusercontent.com/aws-samples/data-engineering-for-aws-immersion-day/master/AnomalyDetection/RandomCutForest/src/main/java/software/amazon/flink/example/RandomCutForestUDF.java>