# Experiment: Data Engineering Pipeline
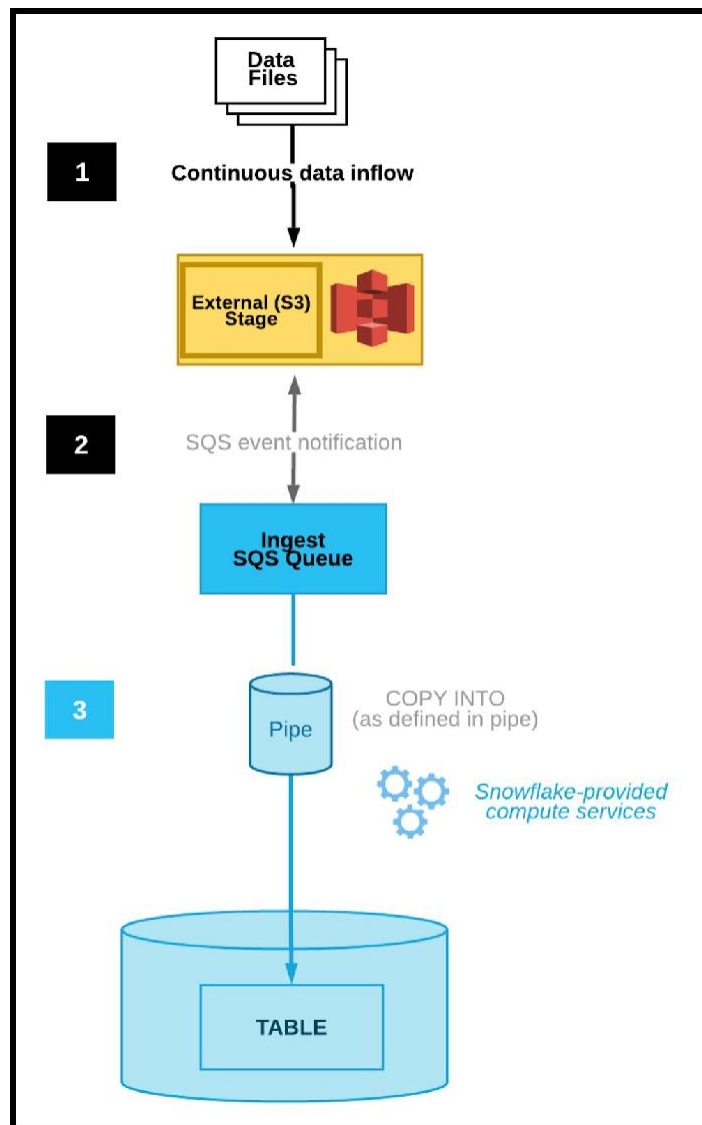## Overview

This experiment demonstrates Snowflake features specifically aligned to "Data Engineering" workload to build modern data pipelines. Data pipelines automate many of the manual steps involved in transforming and optimizing continuous data loads. Frequently, the "raw" data is firstloaded temporarily into a staging table used for interim storage and then transformed using a series of SQL statements before it is inserted into the destination reporting tables. The most efficient workflow for this process involves transforming only data that is new or modified.

## Experiment 22:  Snowpipe

### 22.1 Leverage Snowpipe to Populate Table using S3 Event Notifications

This section describes the most common option for triggering Snowpipe data loads automatically using Amazon SQS (Simple Queue Service) notifications for an S3 bucket. The following diagram shows the Snowpipe auto-ingest process flow:

1. Data files are loaded in a stage.
2. An S3 event notification informs Snowpipe via an SQS queue that files are ready to load. Snowpipe copies the files into a queue.
3. A Snowflake-provided virtual warehouse loads data from the queued files into the target table based on parameters defined in the specified pipe.

**22.1.1** Create a table that Snowpipe will use to write the incoming data. The data being loaded is JSON formatted so create a table using the VARIANT data type.

```
create or replace table trips_raw (v variant);
```

**22.1.2** Create a PIPE definition to tell Snowpipe to load the JSON records in the STAGE into the table you just created.

```
create or replace pipe trips_pipe auto_ingest=true as copy into
trips_raw from@streaming_data/;
```

**22.1.3** Confirm PIPE definition to tell Snowpipe to load the JSON records in the STAGE into the table you just created.

```
show pipes;
```

Make a note of the ARN of the SQS queue for the stage in the *notification_channel* column. Copy the ARN to a convenient location.



**22.1.4** Configuring Event Notifications. Complete the below configuration in your S3 bucket to automatically invoke Snowpipe when files are written there. This is the mechanism that is used to automatically load the streaming data files. Snowpipe relies on notifications from your S3 bucket to achieve this. Instructions for completing this configuration are found in the Snowflake documentation link below. ***Perform only Step 4- Configure Event Notifications and come back to this document***.

https://docs.snowflake.com/en/user-guide/data-load-snowpipe-auto-s3.html#step-4-configure-event-notifications

From the Snowflake console when we execute "show pipes" it will provide the SQS queue that is the notification channel.

That SQS ARN is available by clicking the Copy button on the results pane and pasting that into a text file or other editor.
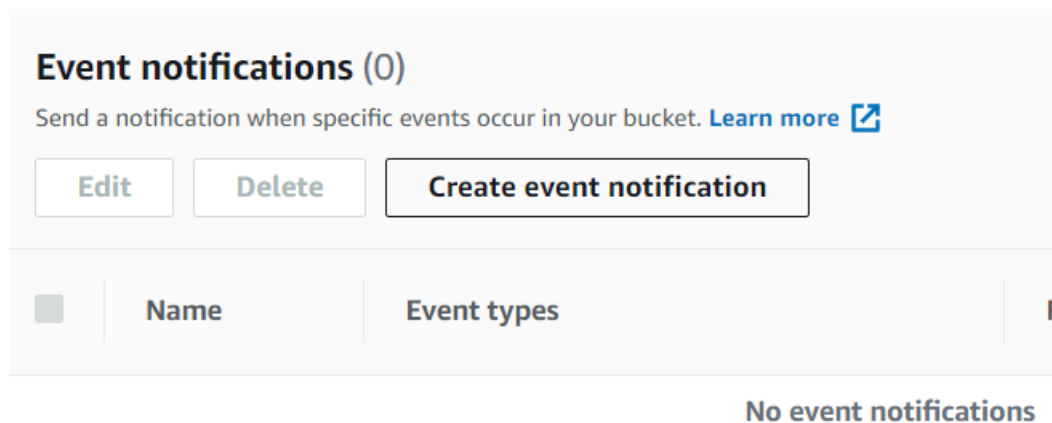
**arn:aws:sqs:us-east-1:941093095223:sf-snowpipe-AIDA5WHMFA43YOTSXNQ7J--ZFTKNPCWky2PyrpnYtzsg**

Sample screenshots of the S3 bucket enabling & configuring of event notifications below

View the S3 Console and choose the bucket that we created in our previous work

Amazon S3 > snowflake-data-pipeline-session1-george

# snowflake-data-pipeline-session1-george Info

| Objects | Properties | Permissions | Metrics | Management | Access Points |

Choose the **Properties** tab and Scroll down to the **Event Notifications** section

## Event notifications (0)

Send a notification when specific events occur in your bucket. Learn more [↗]

| Edit | Delete | **Create event notification** |

| | Name | Event types | F |

No event notifications

Choose **Create event notification**

# Create event notification Info

To enable notifications, you must first add a notification configuration that identifies the eve
publish and the destinations where you want Amazon S3 to send the notifications.

## General configuration

Event name

> Auto-ingest Snowflake

Event name can contain up to 255 characters.

Prefix - *optional*
Limit the notifications to objects with key starting with specified characters.

> *images/*

Suffix - *optional*
Limit the notifications to objects with key ending with specified characters.

> *.jpg*

Enter **Auto-ingest Snowflake** for the Event name

## Event types

Specify at least one event for which you want to receive notifications. For each group, you can choose an event t
can choose one or more individual events.

## Object creation

☑ All object create events
s3:ObjectCreated:*

☐ Put
s3:ObjectCreated:Put

☐ Post
s3:ObjectCreated:Post

☐ Copy
s3:ObjectCreated:Copy

☐ Multipart upload completed
s3:ObjectCreated:CompleteMultipartUpload

22.2.1 Under Event Types and Object creation check the box for **All object create events**

22.2.2 Under Destination choose radio for SQS Queue and SQS Queue ARN and past the arn for the SQS
queue that we retrieved with the "show pipes" result.

### 22.2.3 Choose **Save changes**



### 22.2.4 Under Event Types and Object creation check the box for **All object create events**

**22.1.5** Confirm the setup is working by running the below steps.

Load 1 day of data to test the ingestion by calling the STREAM_DATA stored procedure.

```
call stream_data('2018-01-01', '2018-01-01');
```

**22.1.6** Verify if there are any files already in the bucket/folder.

```
list @streaming_data;
```



**22.1.7** Review the contents of the external stage

```
select $1 from @streaming_data limit 100;
```

## 22.1.8 Check the status of the Snow Pipe



## 22.1.9 Further check the status

```
select system$pipe_status('trips_pipe');
```

Snowpipe copies the data into our raw table. Review the table data.

```
select count(*) from
trips_raw; select * from
trips_raw limit 100;
```

Clean up the stage by calling the purge_files function.



```
call purge_files('trips_raw', '@streaming_data/');
```

Truncate table.
```
truncate table trips_raw;
```

## 22.2 Troubleshooting Snowpipe

If data isn't written to the table you can use the following commands to see where the processis broken.

Verify data was in fact written to your STAGE by the stored procedure.

```
list @streaming_data;
```

Next check the status of the PIPE using the system$pipe_status call. Make sure theresult shows status 'RUNNING' and there are no errors present.

### Details

```
1  {"executionState":"RUNNING","pendingFileCount":0,"notificationChannelN
   ame":"arn:aws:sqs:us-west-2:235488214722:sf-snowpipe-
   AIDATNVBQ73BDZ5EBJD3V-
   VZ7EhFqXekKdtfRTZjsKkA","numOutstandingMessagesOnChannel":0,"lastRecei
   vedMessageTimestamp":"2020-11-
   08T23:10:22.319Z","lastForwardedMessageTimestamp":"2020-11-
   08T23:10:22.389Z"}
```

Done

```
select system$pipe_status('trips_pipe');
```

Lastly, check the information_schema.copy_history to see if there were any problems encountered while writing records to the trips_raw table. Parsing errors or permissions errors could accidentally affect your load jobs. You can find the details for each of the files you are loading. Play with the dateadd function to adjust the size of the window when looking at processed files.

```
select *
from
table(information_schema.copy_history(
table_name=>'citibike_pipelines.public.tri
ps_raw',start_time=>dateadd(hour, -1,
current_timestamp)));
```

**Note:** *the above commands aren't strictly used for troubleshooting purposes. You can use these anytime to gain insight into the status or jobs and track the velocity of arriving data.*