

# Acquiring Twitter Data

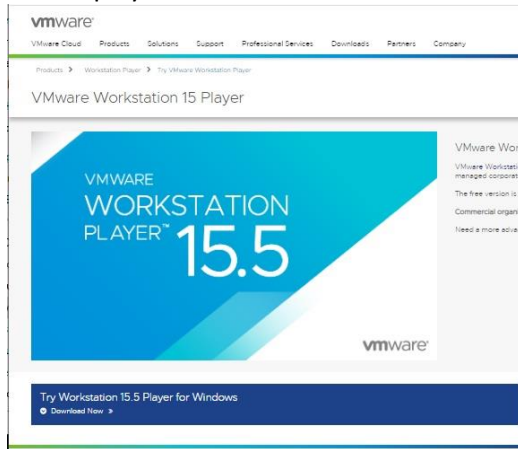
## Using [NiFi](#) & [Kafka](#)

Hortonworks [HDF Sandbox](#)

Introduction: This article is an updated version of the tutorial [here](#). It includes fixes to few errors in the original implementation. This tutorial works on Windows 10 with vMWare workstation player 15. It was not tested on Mac or with Oracle [VirtualBox VM](#).

### Pre-requisitea:

1. vMplayer to install the Hortonworks sandnox, Download and install [vMware workstation player](#)



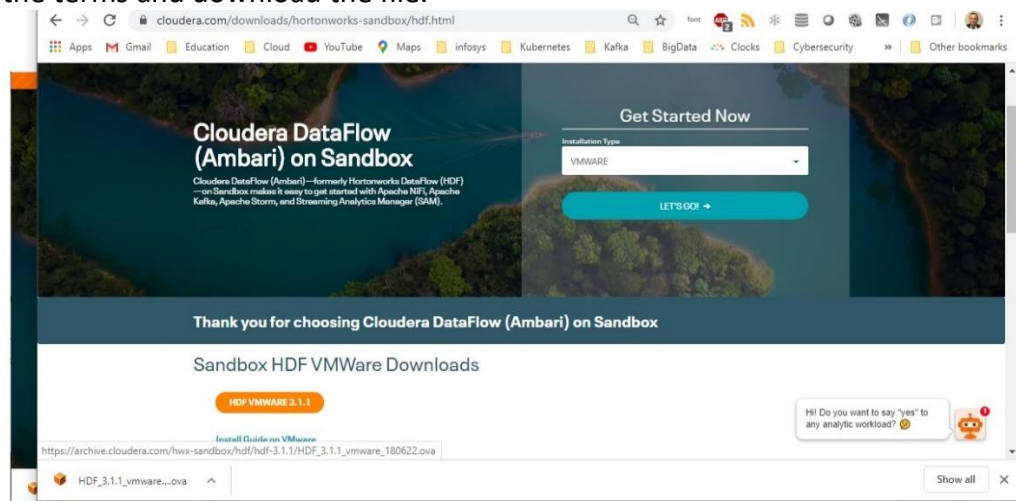
2. Twitter API developer account.

Follow [instructions here](#) in order to create a Twitter developer account.

Save the 4 keys/token that are provided by Twitter: **consumer API Key**, **consumer API Secret key**, **access Token** and **access Token Secret**

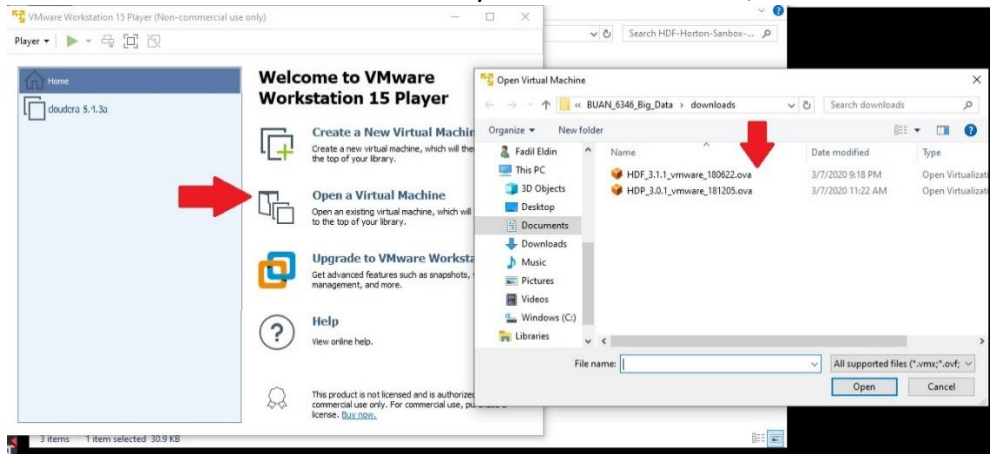
### Install and run the HDF Hortonworks sandbox:

Download the image ova file of the Horton Works HDF from [this page](#), select vMware, accept the terms and download the file.

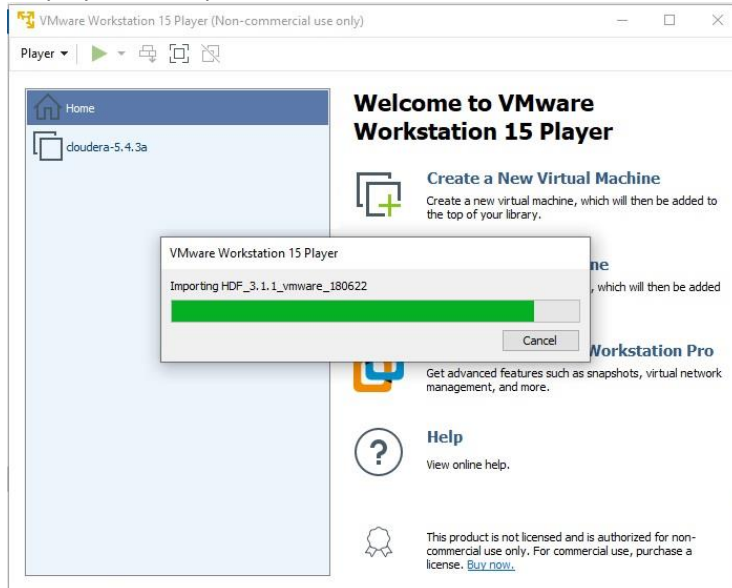


At the time of the creation of this tutorial, the name of the downloaded files was **HDF\_3.1.1\_vmware\_180622.ova**

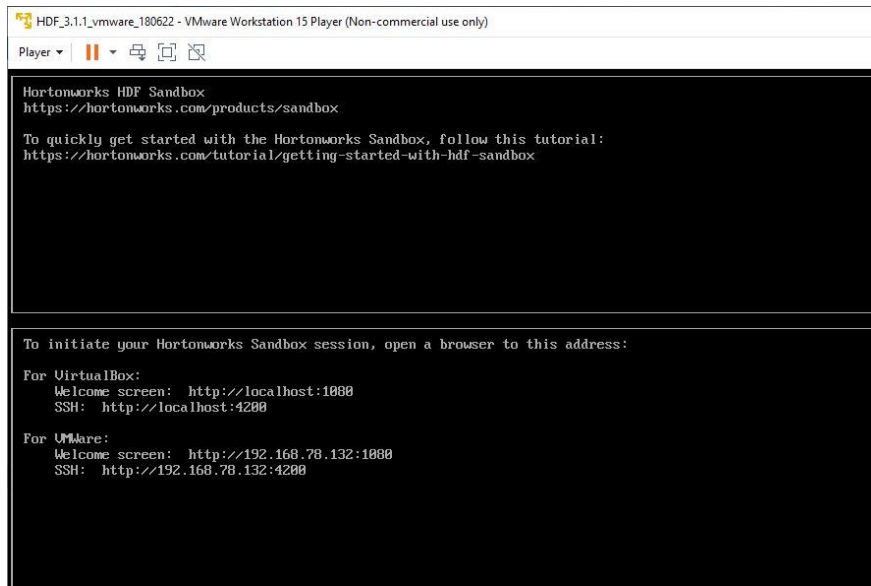
Open (vMware workstation 15 Player) on the right pane, click **“Open A Virtual Machine”** take the Windows file browser to where you save the .ova file, select the file and click **open**



vMplayer will import and install the sandbox



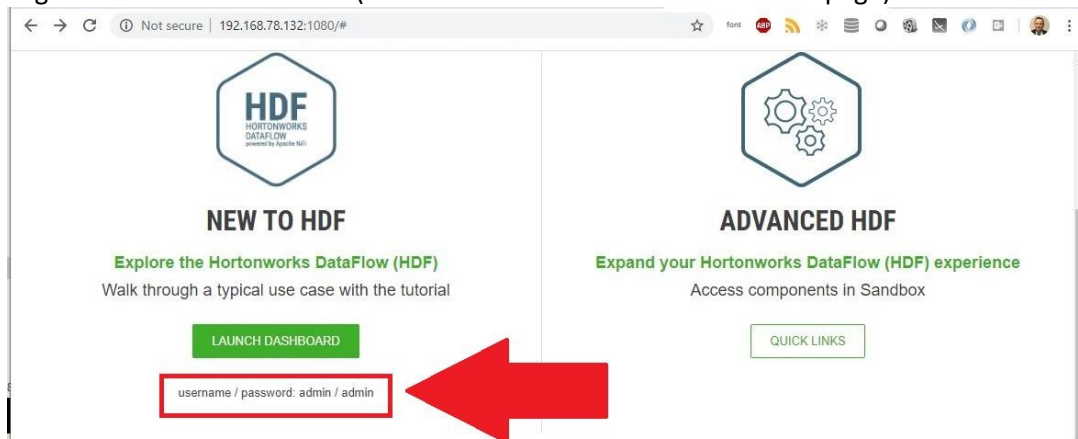
At the end of a successful installation, you will see a page like this:



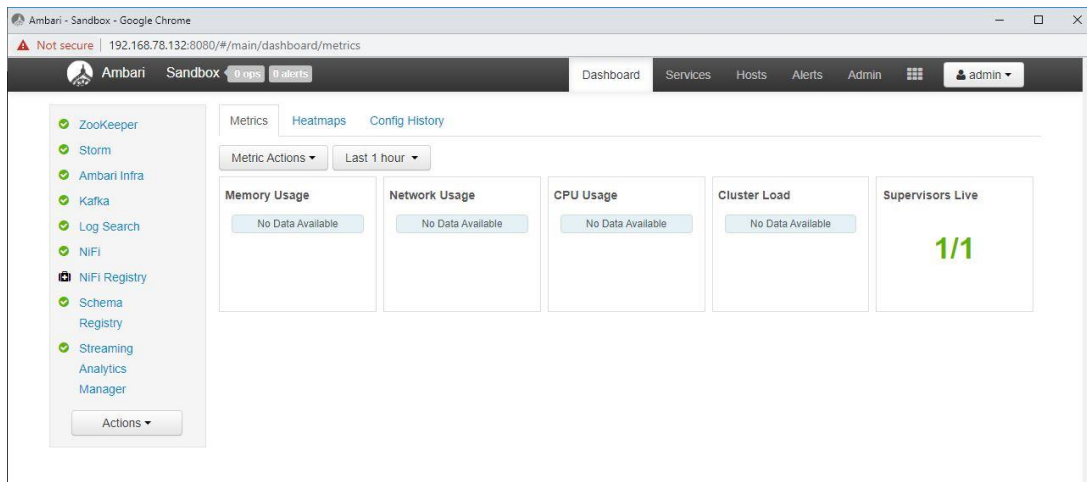
The default username/Password for the welcome screen (See vMware, port 1080) is **admin/admin**  
And the default for the SSH (port 4200) is **root/Hadoop**

## Preparing the SandBox:

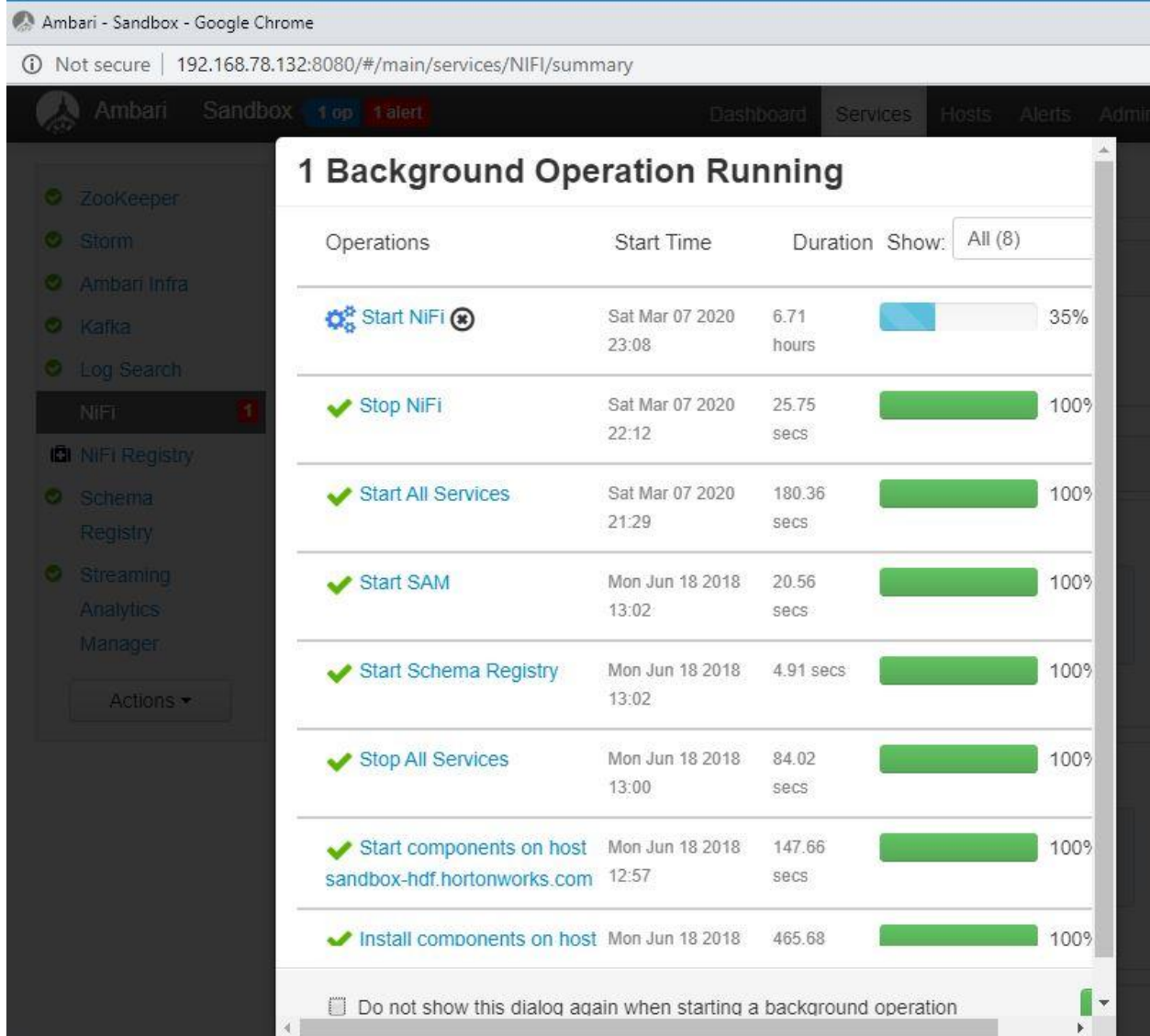
Login to the welcome screen (click Launch dashboard from the 1080 page)



and make sure all services are running

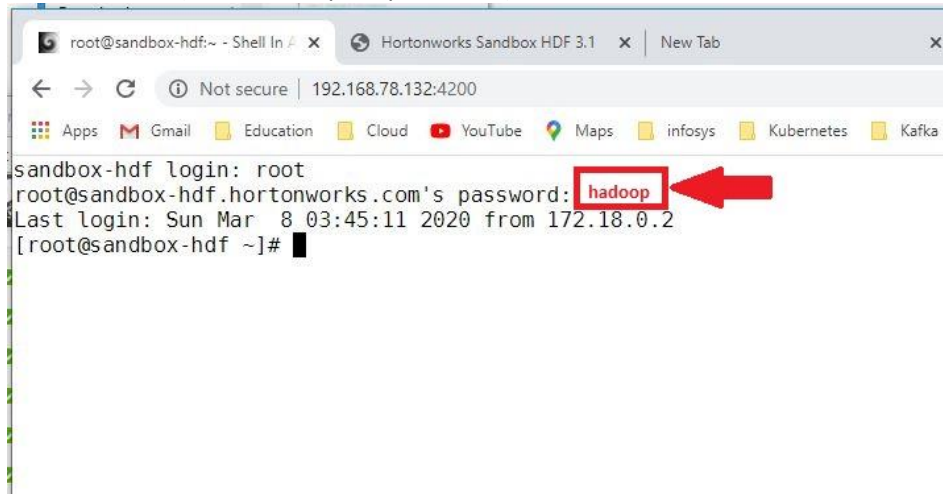


We need **NiFi** and **Kafka**, if they are not running, click and start



Sync the VM clock, this is needed to run Twitter, if the data on your VM is in the past then Twitter will reject the API call

ssh to the VM (root/Hadoop on port 4200)




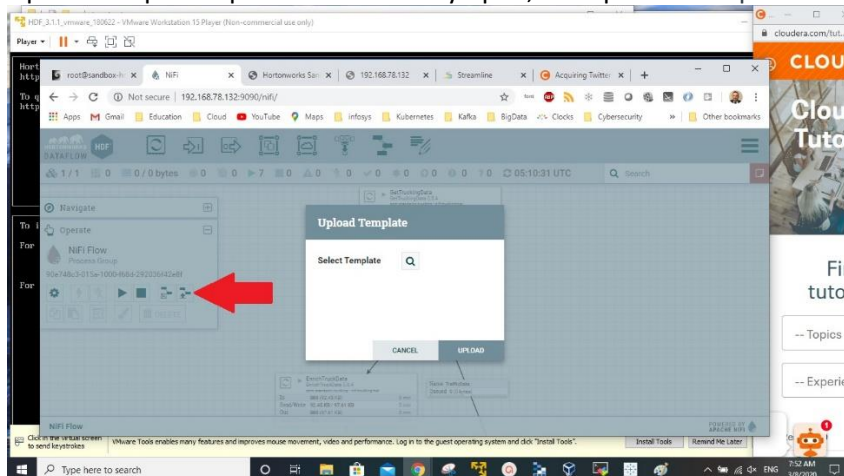
Run the following commands

```
sudo yum install -y ntp
sudo service ntpd stop
sudo ntpdate pool.ntp.org
sudo service ntpd start
```

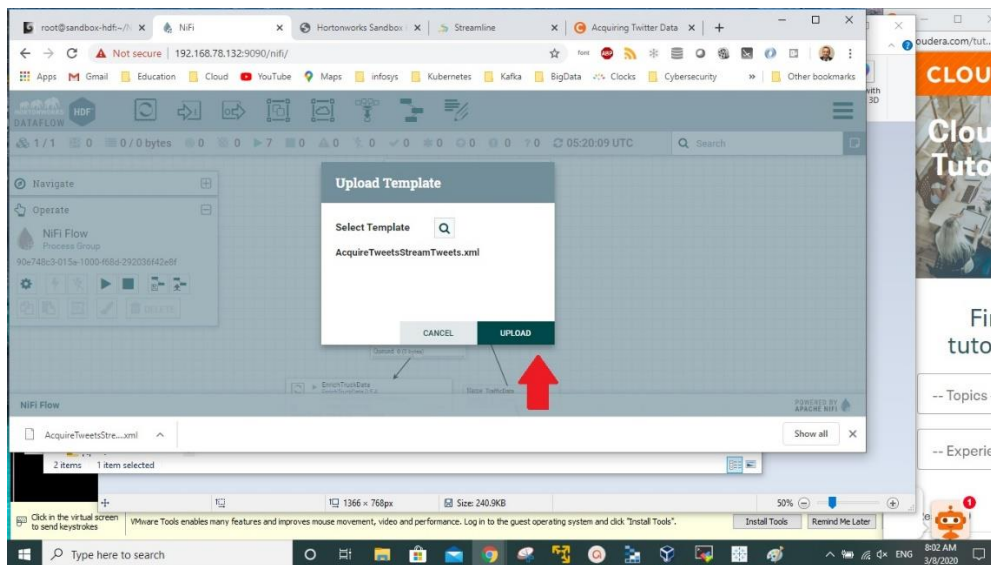
## Import NiFi Flow For Twitter Acquisition and storage

Download the NiFi xml template file from [AcquireTweetsStreamTweets.xml](#) to your local computer. open HDF NiFi UI at <http://YOUR-IP:9090/nifi>.

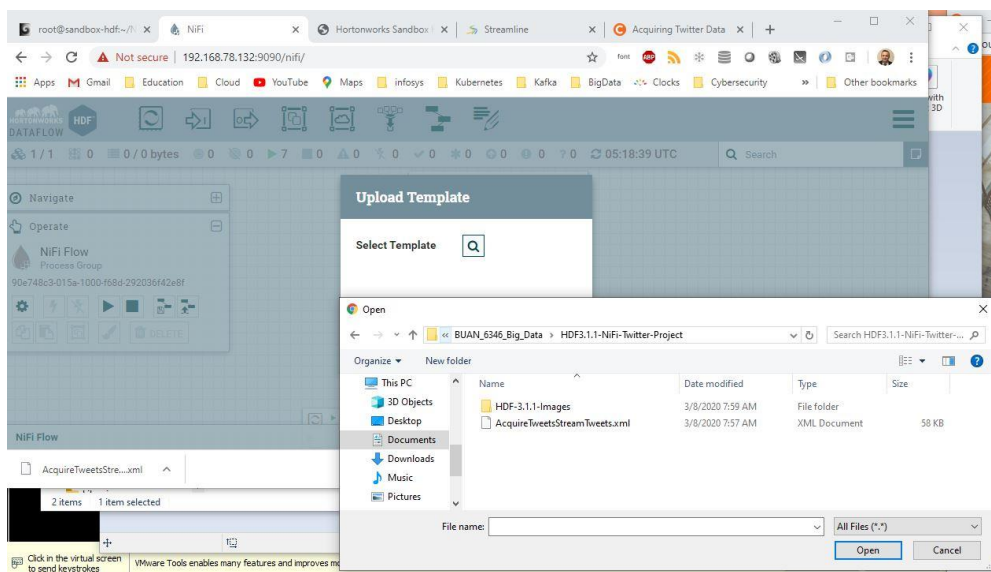
Open the Operate panel if not already open, then press the Upload Template icon .



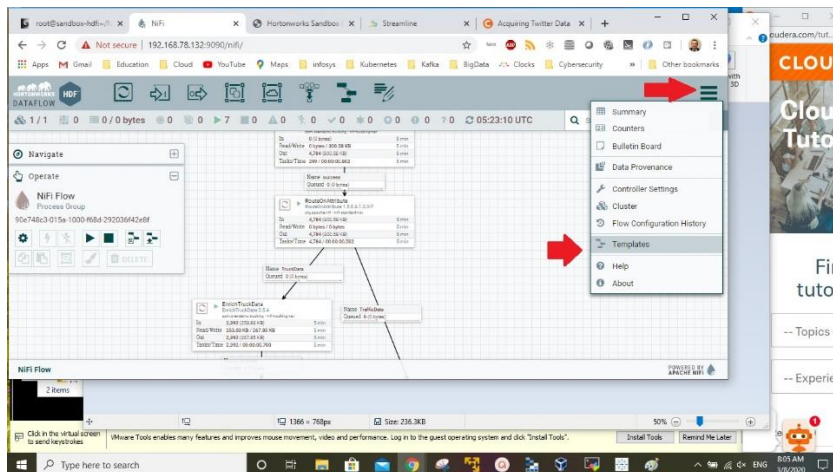
Press on Select Template icon .



The file browser on your local computer will appear, find AcquireTweetsStreamTweets.xml template you just downloaded, then press Open, then press UPLOAD.

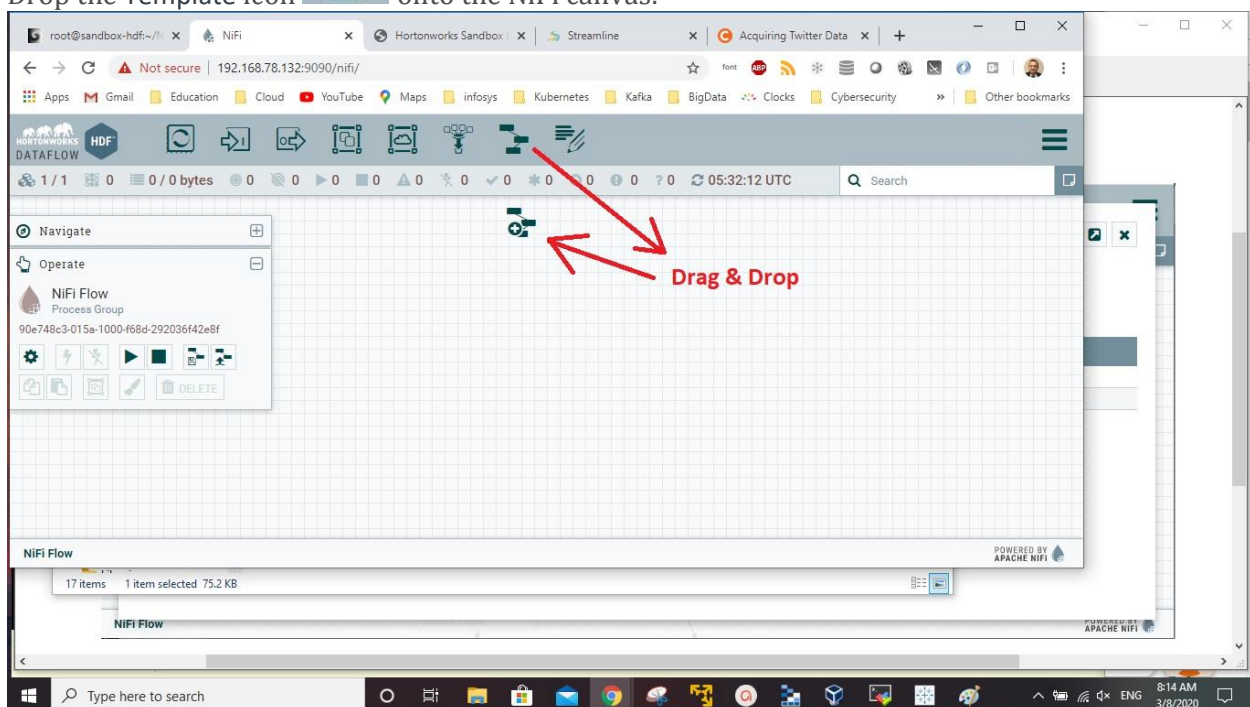




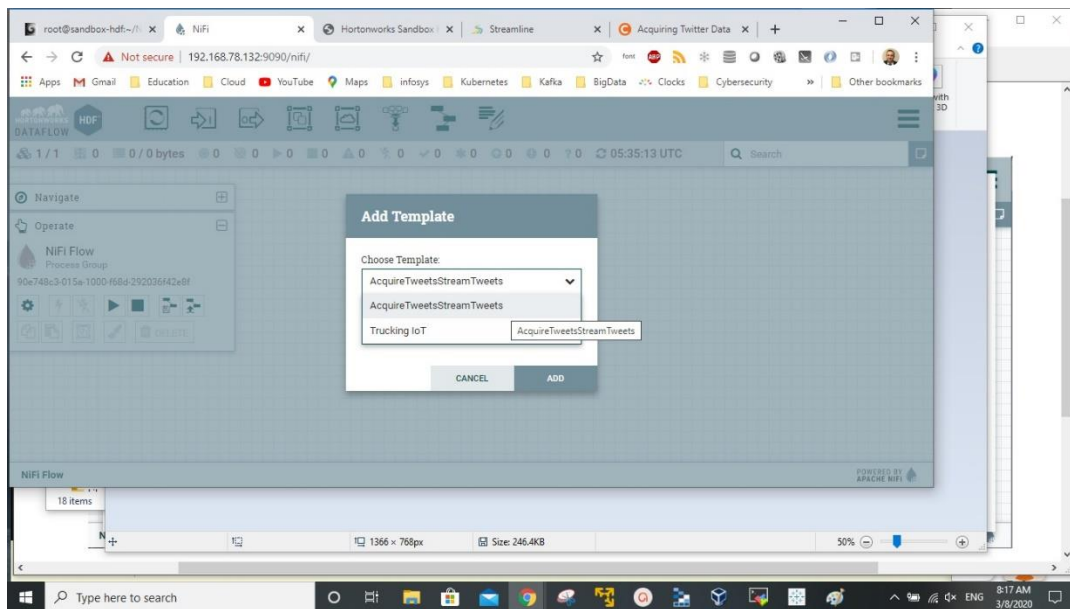


You should receive a notification that the Template successfully imported. Press OK to acknowledge.

Drop the Template icon  onto the NiFi canvas.

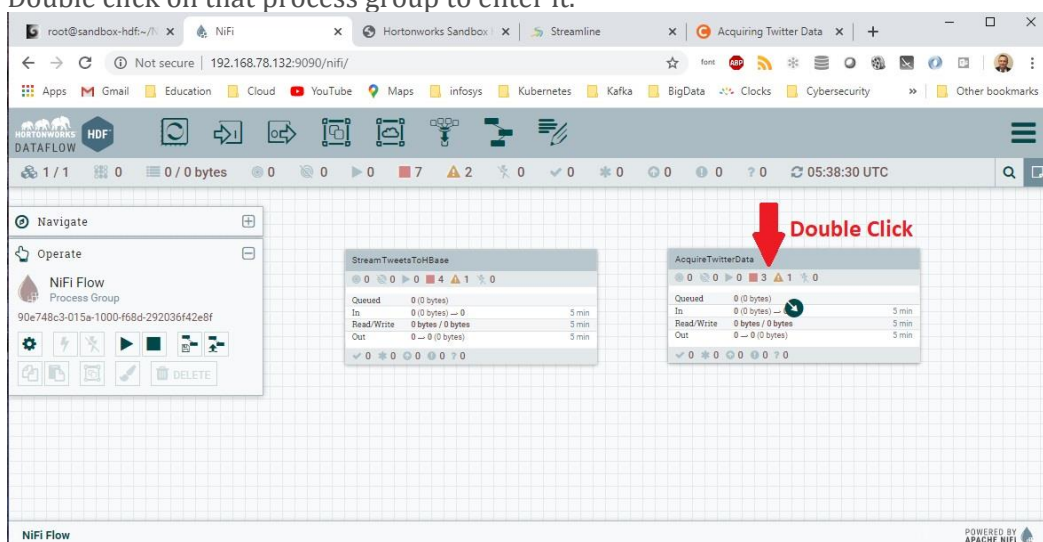


Add Template called AcquireTweetsStreamTweets.



You will notice on the process group called AcquireTwitterData, there is one yellow warning.

Double click on that process group to enter it.



Zoom in if needed. GrabGardenHose processor has the warning. The reason is that we need to update the Consumer API Key and Consumer API Secret Key and the Access Token and Access Token Secret in the processor's properties table for the warning to go away. Modify what terms you want to filter On.



SETTINGS

SCHEDULING

PROPERTIES

COMMENTS

Required field +


Property	Value
Twitter Endpoint	Sample Endpoint
Consumer Key	
Consumer Secret	Sensitive value set
Access Token	
Access Token Secret	Sensitive value set
Languages	en
Terms to Filter On	coronavirus, COVID-2019, COVID2019, COVID19, COVID-...
IDs to Follow	No value set
Locations to Filter On	No value set

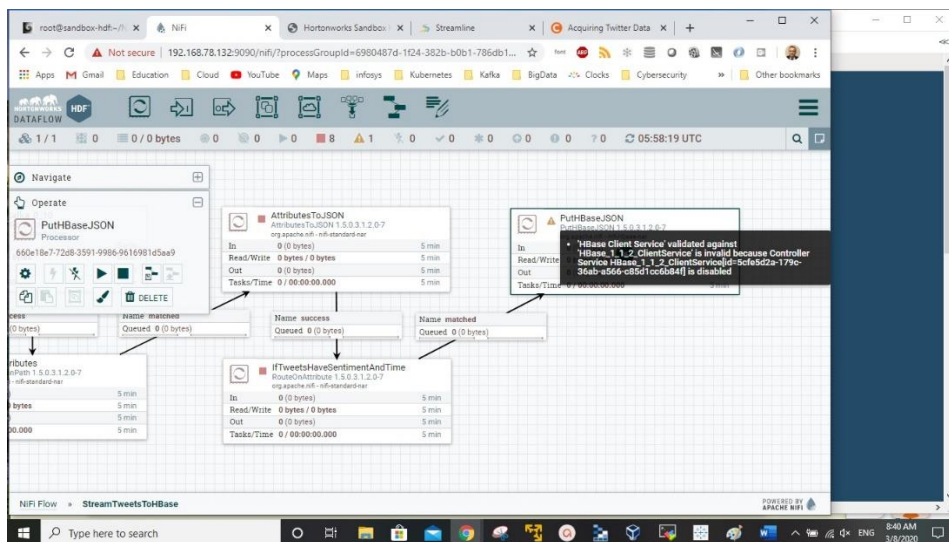
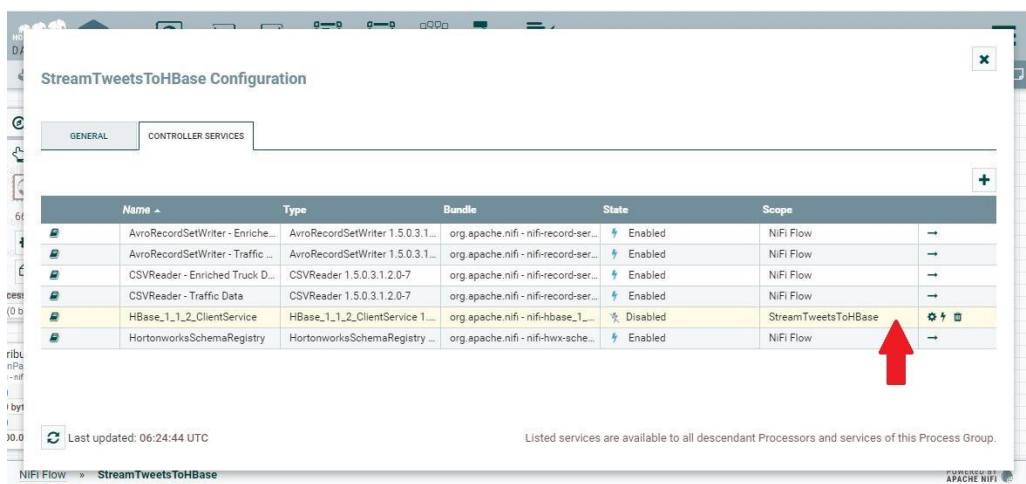
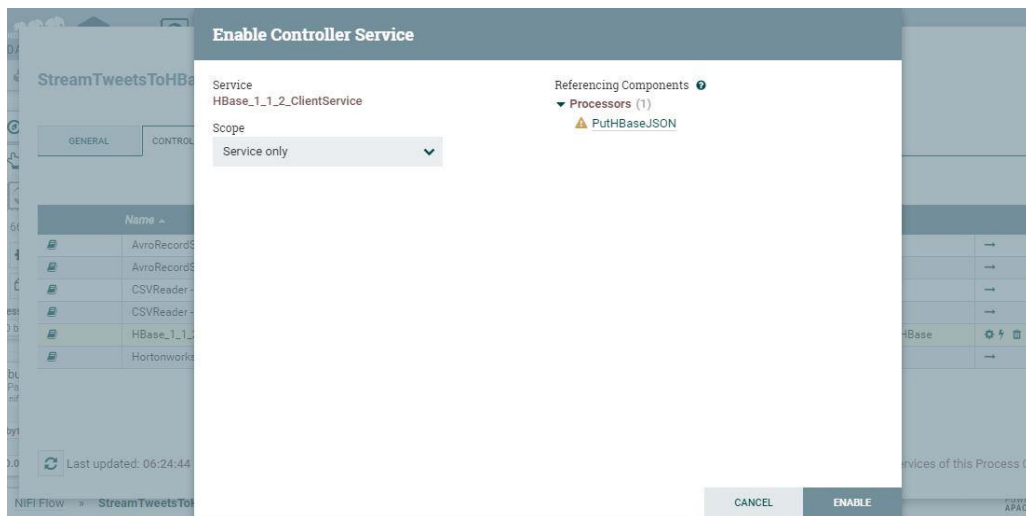
CANCEL

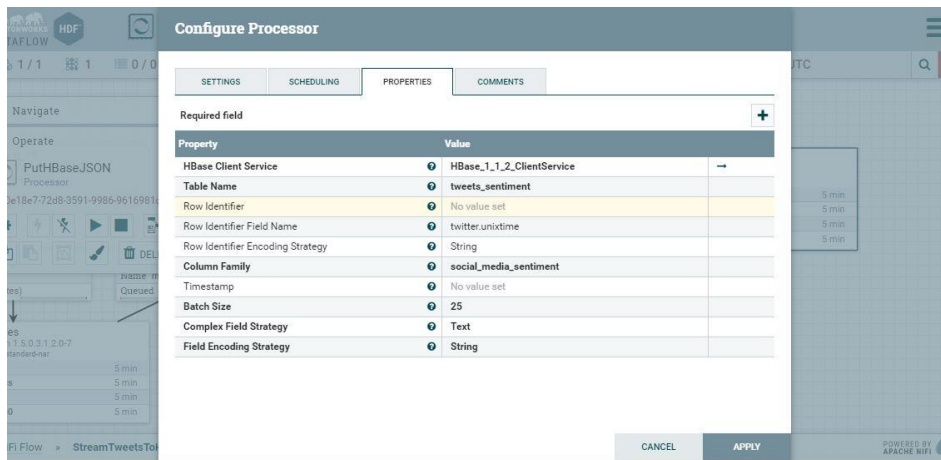
APPLY

Start the NiFi flow. Hold control + mouse click on each process group, then click the start option.

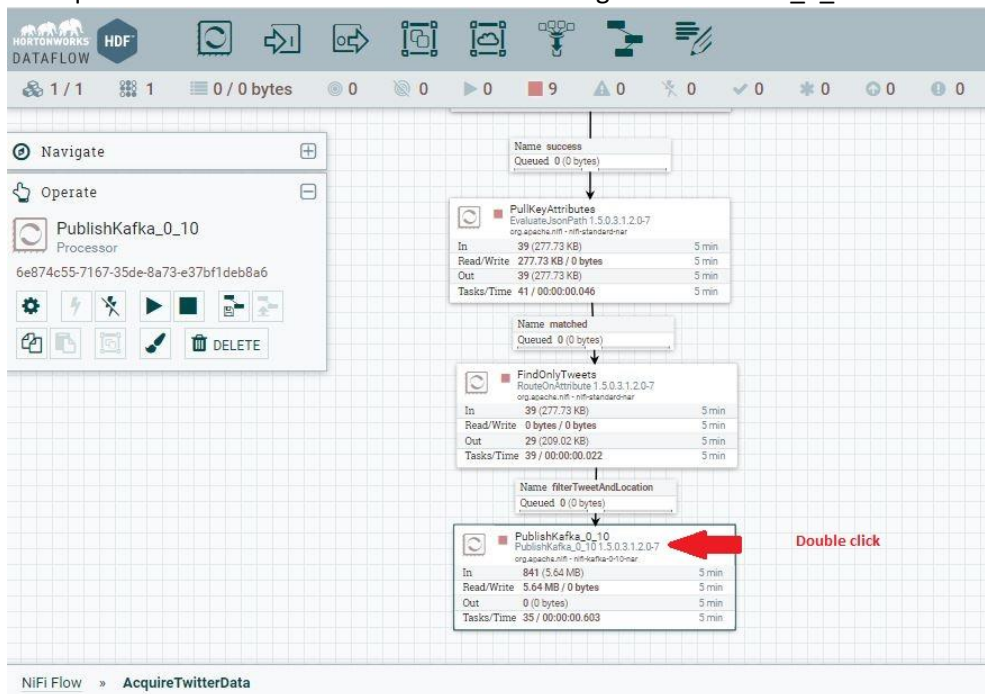
StreamTweetsToHBase	AcquireTwitterData
<div> <div>0</div> <div>0</div> <div>0</div> <div>5</div> <div>0</div> <div>0</div> </div>	<div> <div>0</div> <div>0</div> <div>0</div> <div>4</div> <div>0</div> <div>0</div> </div>
Queued 0 (0 bytes)	Queued 0 (0 bytes)
In 0 (0 bytes) → 0 5 min	In 0 (0 bytes) → 0 5 min
Read/Write 0 bytes / 0 bytes 5 min	Read/Write 0 bytes / 0 bytes 5 min
Out 0 → 0 (0 bytes) 5 min	Out 0 → 0 (0 bytes) 5 min
<div> <div>✓ 0</div> <div>* 0</div> <div>0</div> <div>0</div> <div>?</div> <div>0</div> </div>	<div> <div>✓ 0</div> <div>* 0</div> <div>0</div> <div>0</div> <div>?</div> <div>0</div> </div>

Select StreamStweets To HBase, click configuration  and enable StreamTweetsToHBase





In AcquireTwitterData double click the last rectangle PublishKafka\_0\_10



**Configure Processor**

SETTINGS SCHEDULING PROPERTIES COMMENTS

Required field +

Property	Value
Kafka Brokers	sandbox-hdf.hortonworks.com:6667
Security Protocol	PLAINTEXT
Kerberos Service Name	No value set
Kerberos Principal	No value set
Kerberos Keytab	No value set
SSL Context Service	No value set
Topic Name	tweets
Delivery Guarantee	Best Effort
Kafka Key	No value set
Key Attribute Encoding	UTF-8 Encoded
Message Demarcator	No value set
Max Request Size	1 MB
Acknowledgment Wait Time	5 secs
Max Metadata Wait Time	5 sec

CANCEL APPLY

In the properties, change hdp to hdf change the topic name to tweets.

Once NiFi writes tweet data to Kafka on HDP, you can check the provenance events quickly by looking at the PublishKafka\_0\_10 processor inside the AcquireTwitterData process group.

To turn off a process group, you can do so by holding control + mouse click on for instance the AcquireTwitterData process group, then choose stop option.

## Run the flow and inspect the results.

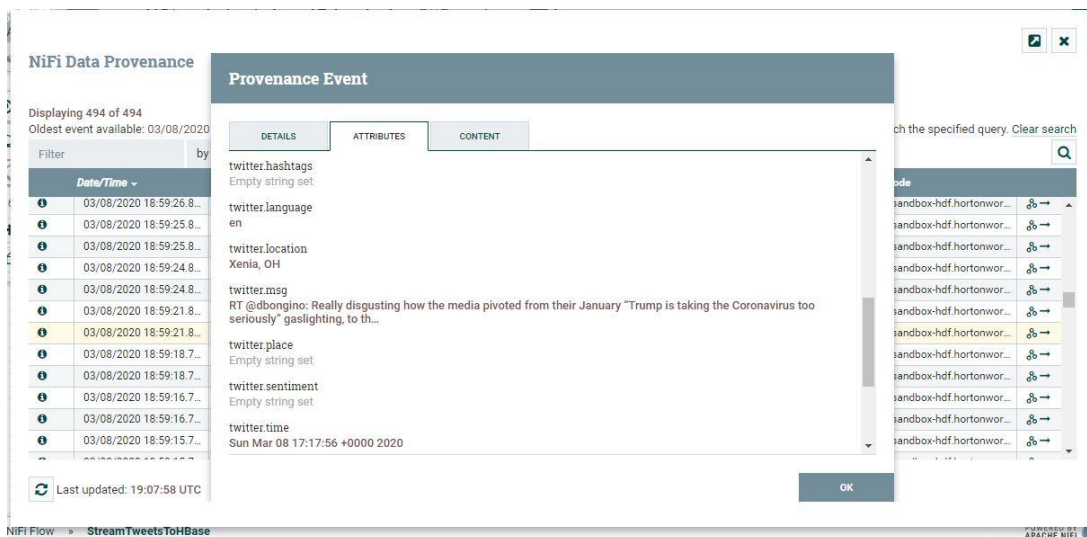
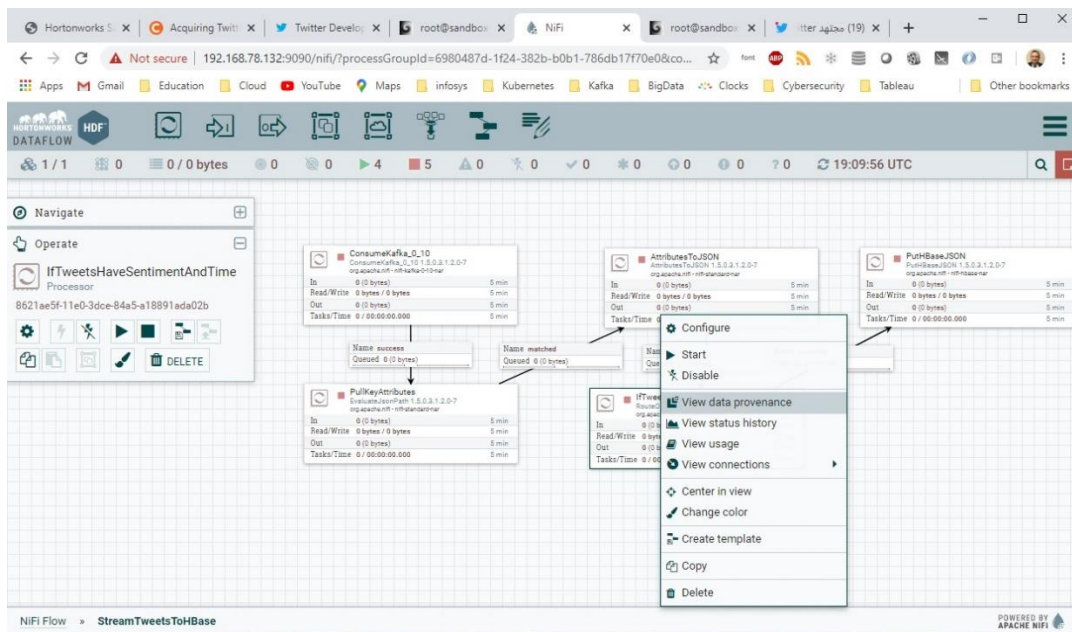
On NiFi flow main window, click CTRL-A and Run the flow.

Monitor as the input and output pipes get filled with tweets,

StreamTweetsToHBase			
0	0	5	0
Queued	2,797 (18.41 MB)		
In	0 (0 bytes) → 0	5 min	
Read/Write	815.36 KB / 18.6 MB	5 min	
Out	0 → 0 (0 bytes)	5 min	
0	0	0	0

They are filtered according to the keywords we put and analyzed.

Here are few sample tweet



This Tutorial on [LinkedIn](#)  
Thanks