

MIDS W205

Exercise #	2	Title	Introduction to the Elements of a Streaming application.
Related Module(s)	8,9	Goal	Implement an end to end streaming app.
Last Updated	3/12/16	Expected duration	15-25 hours.

Introduction

Streaming applications may seem complex but understanding how they operate will be critical for a data scientist. In this exercise, we will explore a streaming application that analyzes the Twitter data. In order to allow you to explore a more complex implementation in a short period of time, you will be developing codes by using an existing code base. You will use Streamparse as seen in Lab 6 with a given topology. The application reads the stream of tweets from the Twitter streaming API, parses them, counts the number of each word in the stream of tweets, and writes the final results back to a Postgres database.

Scope and Goal

In this exercise, you will capture and process the live Twitter data stream and learn the following tasks:

- How to capturing the live data
- How to setup a stream processing pipeline
- How to process and get insights
- How to store the final processing results to a relational database management system

Technologies used in this exercise:

Apache Storm, Amazon EC2, python, Twitter API, Streamparse, Postgres, and Psycopg

Instructions, Resources, and Prerequisites

In the following table you will find all the references related to the tools/libraries and programs that are used in this exercise.

Resource	What
http://storm.apache.org/documentation.html	Apache Storm Documentation
https://streamparse.readthedocs.org/en/latest/api.html	Streamparse Documentation
https://dev.twitter.com/streaming/overview	Twitter Stream API
http://docs.tweepy.org/en/v3.4.0/	Tweepy Documentation
http://docs.tweepy.org/en/v3.4.0/streaming_how_to.html	
http://initd.org/psycopg/	psycopg

Note: You are going to use streamparse and install psycopg. We therefore suggest that you do not use your saved instance image with the attached volume from the previous labs. Instead, we suggest you either use the AMI that you created for Lab 6 or create a new image from UCB AMI specifically for this exercise to avoid broken links or unnecessary package upgrades.

Use Case:

Capturing and analyzing live twitter data around your business interest area can give you a deeper understanding of the current social trends and demands. Historical data can give you information on the mainstream trends over a certain period of time, but live data can provide immediate and real-time insight. For example, a person who manages the ads for the TV programs can capture the live Twitter trends at the time of the live show to engage more TV viewers during a popular TV program (i.e contextual advertising).

So, in this exercise, you will capture live tweets that show people's live interests, process them in real-time to get insights and aggregate the results in a database. Figure 1 shows the overall architecture of the application. Figure 1 also shows the storm topology that you need to develop as part of the application. Using Tweepy library, the application reads the live stream of tweets from twitter in the **Tweet-spout** component. The **Parse-tweet-bolt** parses the tweets, extracts the words from each parsed tweet and emits the words to the next bolt component (i.e **Count-bolt**) in the topology. **Count-bolt** counts the number of each word in the received tuples and updates the counts associated with each words in the **Tweetwordcount** table inside the **Tcount** database. **Tcount** is a postgres database.

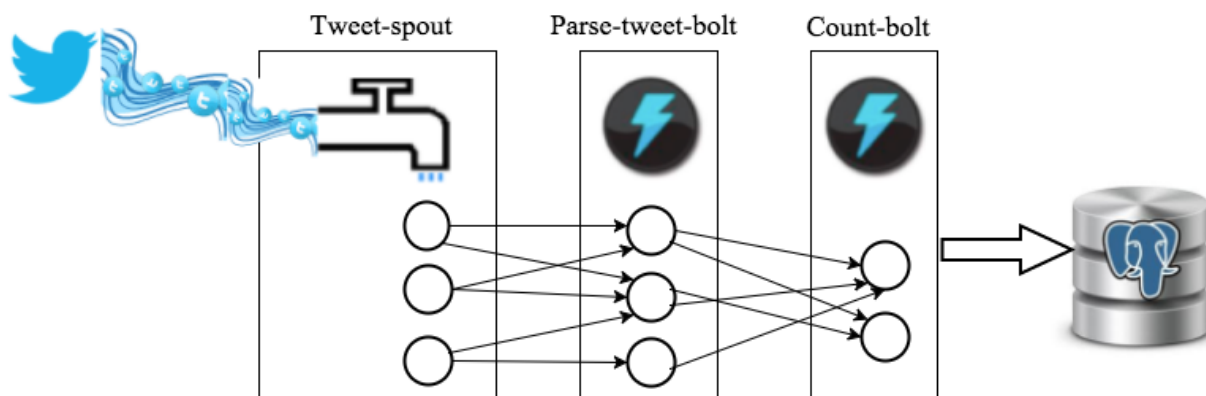


Figure 1: Application Topology

Overall Guidelines for all steps

In this section we provide the overall guidelines for implementing the real time system as shown in Figure 1. You need to follow the below instructions for each of the steps of your implementation. You would use **UCB W205 Spring 2016** AMI for creating your own EC2 server for this exercise. You use your github account to store your programs, data, etc.

Step-1. Environment and Tool Setup

1. Clone the Github Repository for this exercise from [git@github.com:UC-Berkeley-I-School/w205-spring-16-labs-exercises.git](https://github.com/UC-Berkeley-I-School/w205-spring-16-labs-exercises.git) (If you want to just clone the exercise_2 directory, see the appendix).
2. Create an EC2 instance using **UCB W205 Spring 2016** and install Streamparse on it. You can either use the instructions from Lab 6 to install the Streamparse or use the configured AMI with a Streamparse installation from Lab 6 if you saved it. You may want to save this AMI once you configured it for the first time to save time in future re-launches. (Note: make sure python 2.7 is the default python version)
3. Install **psycopg** by running:

```
$ pip install psycopg2
```
4. Create a project called **EX2Tweetwordcount** in Streamparse.
5. Copy the files from the **tweetwordcount** directory in your cloned repository and paste them in the corresponding folders in the new **EX2Tweetwordcount** project. The description of the files in the code base is provided in Table 1.
6. Modify the **EX2tweetwordcount.clj** to implement the topology in Figure 1.

Note1: You may not want to keep your EC2 server live when you are not working, as you will run out of credit that way. So, you could save your work in github as you progress and when you launch your sever, you can re -pull the code and use.

Name of the program	Location	Description
Tweets.py	exercise_2/tweetwordcount/src/spouts/	Tweet-spout
Parse.py	exercise_2/tweetwordcount/src/bolts/	Parse-tweet-bolt
Wordcount.py	exercise_2/tweetwordcount/src/bolts/	Count-bolt
Twittercredentials.py	exercise_2/	Twitter App Keys
hello-stream-twitter.py	exercise_2/	Sample twitter Stream program
tweetwordcount.clj	/exercise_2/tweetwordcount/topologies/	Topology for the program
psycopg-sample.py	exercise_2/	Sample codes on how to use psycopg

Table 1: Description of the files in the code base

Step-2. Twitter application setup

Once you have the topology configured with the codes in the cloned repository, you need to modify the spout program (i.e Tweets.py) to pull the tweets from twitter streaming API. In order to get the tweets, you need to set up a twitter application and get access keys for pulling the tweets out of the twitter streaming API.

The following instructions will step you through the process of acquiring data from Twitter. The first thing that you need to do is to install Tweepy, which is a python library for accessing the Twitter API.

2.1 Install Tweepy

The easiest way to install Tweepy is by using pip:

```
$pip install tweepy
```

You may also use Git to clone the repository directly from Github and install it manually:

```
$git clone https://github.com/tweepy/tweepy.git
$cd tweepy
$python setup.py install
```

The next step is to get access to twitter data. Twitter data can be accessed over the Web by creating a Twitter application and then using the access keys Twitter provide for the application in your program.

2.2 Create an Application

Note: You will need to have a Twitter account to create an application.

To create an application:

1. Login to Twitter (<https://www.twitter.com/>).
2. Visit <https://apps.twitter.com> and click on "Create New App".

 Application Management



3. Fill in the application name, description, and Website. The name will be listed in your application list when you return this Website.

Create an application

Application Details

Name *

Your application name. This is used to attribute the source of a tweet and in user-facing authorization screens. 32 characters max.

Description *

Your application description, which will be shown in user-facing authorization screens. Between 10 and 200 characters max.

Website *

Your application's publicly accessible home page, where users can go to download, make use of, or find out more information about your application. This fully-qualified URL is used in the source attribution for tweets created by your application and will be shown in user-facing authorization screens.
(If you don't have a URL yet, just put a placeholder here but remember to change it later.)

Callback URL

Where should we return after successfully authenticating? [OAuth 1.0a](#) applications should explicitly specify their `oauth_callback` URL on the request token step, regardless of the value given here. To restrict your application from using callbacks, leave this field blank.

- **Name:** Your app name. It needs to be a unique name across all twitter applications (i.e. no one has used it before)
- **Description:** A short description for your app.
- **Website:** Enter the website address where the app will be hosted.
- **Callback URL:** Ignore this field.

4. Agree to the terms and agreements and click on "Create your Twitter Application"

Once you have successfully created an application, it should take you to the newly created application. Here you must create access keys for subsequent operations by your application. To do so, use the following procedure:

1. Click on the "Keys and Access Tokens" tab.
2. Click on "Create my Access Token" near the bottom of the page.

Your access token

It looks like you haven't authorized this application for your own Twitter account yet. For your convenience, we give you the opportunity to create your OAuth access token here, so you can start signing your requests right away. The access token generated will reflect your application's current permission level.

Create my access token

The response should be relatively immediate.

OAuth settings

Your application's OAuth settings. Keep the "Consumer secret" a secret. This key should never be human-readable in your application.

Access level	Read-only About the application permission model
Consumer key	<code>XX</code>
Consumer secret	<code>XX</code>
Request token URL	<code>https://api.twitter.com/oauth/request_token</code>
Authorize URL	<code>https://api.twitter.com/oauth/authorize</code>
Access token URL	<code>https://api.twitter.com/oauth/access_token</code>
Callback URL	None
Sign in with Twitter	No

Your access token

Use the access token string as your "oauth_token" and the access token secret as your "oauth_token_secret" to sign requests with your own Twitter account. Do not share your oauth_token_secret with anyone.

Access token	<code>XX</code>
Access token secret	<code>XX</code>
Access level	Read-only

[Recreate my access token](#)

Now you have four things (as blurred above):

1. A consumer key that identifies your application.
2. A consumer secret that acts as a "password" for your application.
3. An access token that identifies your authorized access.
4. An access token secret that acts as a "password" for that authorized access.

At any point, you can revoke the access key or regenerated any of these values. To completely disable the application, you must delete the application. This does remove the consumer key, secret, and access tokens from Twitter's system and any program using them will immediately stop working.

2.3 Test your Application

In the code base that you cloned, `hello-stream-twitter.py` is a sample application that pulls tweets from the twitter streaming API. This program uses `Tweepy` to work with the streaming API. Use this program to test your application. Change the code in

`Twittercredentials.py` and insert your consumer key, consumer secret, access token, and access token secret. You should then be able to just run the program and get tweets:

```
$python hello-stream-twitter.py
```

Step-3. Application deployment

Now you have a streamparse project and a twitter application, your task is to write codes to connect necessary pieces together as to create a full stream tweetword count processing application as depicted in Figure 1. Your application has the following pieces:

1. A spout connected to twitter streaming API that pulls the tweets from twitter stream and emits them to the parse bolt. You need to modify the necessary codes to complete this part.
2. A tweet-parse bolt that parses the tweets emitted by the spout and extracts the words out of the received tweets.
3. A tweet word count bolt that counts the number of words emitted by the tweet-parse bolt and updates the total counts for each word in a corresponding table inside a database. In order to be able to update the results in a DB, you need to create a Postgres DB called **Tcount** and a table called **Tweetwordcount** table inside **Tcount** database. You also need to modify the code in the `wordcount.py` in order to update the word count in the **Tweetwordcount** table for each word in the tweet stream. To interact with Postgres, you can use **psycopg** library. You can find sample codes on how to use **psycopg** in `psycopg-sample.py` file.

Step-4. Serving Scripts

In this step, your task is to develop two simple scripts that query the database and return specific results as follows:

1. `finalresults.py`

This script gets a word as an argument and returns the total number of word occurrences in the stream. For example:

```
$ python finalresults.py hello
$ Total number of occurrences of "hello": 10
```

2. Running `finalresults.py` without an argument returns all the words in the stream and their total count of occurrences, sorted alphabetically in an ascending order, one word per line. For example:

```
$ python finalresults.py
$ (<word1>, 2), (<word2>, 8), (<word3>, 6), (<word4>, 1), ...
```

3. histogram.py

This script gets two integers k1,k2 and returns all the words that their total number of occurrences in the stream is more or equal than k1 and less or equal than k2. For example:

```
$ python histogram.py 3,8
$      <word2>: 8
      <word3>: 6
      <word1>: 3
```

Submission Instructions:

General guidelines:

1. All code outlined above **must be** committed and pushed to your GitHub repository.
2. All code must be runnable by your instructor in the **UCB W205 Spring 2016**.
3. Your GitHub repository **must be** shared with your section instructor via pull request.

Your Github repository should include:

1. Complete and fully functional twitter application codes based on the description above.
2. `Architecture.pdf`: a complete documentation (max 4 pages) of your twitter application including directory and file structure, application idea, description of the architecture, file dependencies, any necessary information to run the application, etc.
3. A directory called `screenshots/` that has at least three screenshots of an end-to-end execution of your application of your choice (name screenshots consistently, e.g.: `screenshot-twitterStream.png`, `screenshot-storm-components.png`, `screenshot-extract-results.png`).
4. `Readme.txt`: a file that shows the step-by-step instructions on how to run the application.
5. `Plot.png`: a bar chart showing the top 20 words in your twitter stream.

Appendix

In this section we introduce some additional instructions. Some of these are from other labs and you may be able to re-use installations from those.

Installing Streamparse

First make sure you are using python version 2.7.x. If you do not, update to 2.7 by following the instructions in the appendix section.

Install `ez_setup`.

```
$sudo curl -o ez_setup.py https://bootstrap.pypa.io/ez\_setup.py
$sudo python ez_setup.py
```

Use `ez_install` to install `pip`, which is a package manager for Python software.

```
$sudo /usr/bin/easy_install-2.7 pip
```

Then use `pip` to install `virtualenv`. `virtualenv` is a tool to create and manage dependencies for distinct Python environments. Streamparse uses `virtualenv` to manage all dependencies for individual Python Storm projects.

```
$sudo pip install virtualenv
```

Streamparse requires the build tool `lein` to resolve dependencies, so now you will install `lein`.

Note: In the video the installer fails to save to `/usr/bin`, so we have to move the `lein` file there. If the command succeeds you will not need to use the `mv lein /usr/bin` command.

```
$wget --directory-prefix=/usr/bin/
https://raw.githubusercontent.com/technomancy/leiningen/stable/bin/lein
```

If you check the permissions on the `lein` file, you will see it is not executable. This means the shell and operating system will not allow you to run it as a command.

```
$ls -l /usr/bin/lein
-rw-r--r-- 1 root root 12713 Oct 25 17:01 /usr/bin/lein
```

Use the following `chmod` command to turn on the executable permission for all users.

```
$ chmod a+x /usr/bin/lein
```

Check that it looks like what you expected.

```
$ls -l /usr/bin/lein
-rwxr-xr-x 1 root root 12713 Oct 25 17:01 /usr/bin/lein
```

The first time you run lein, it will install itself.

```
$sudo /usr/bin/lein

$lein version
WARNING: You're currently running as root; probably by accident.
Press control-C to abort or Enter to continue as root.
Set LEIN_ROOT to disable this warning.

Leiningen 2.5.3 on Java 1.7.0_79 Java HotSpot(TM) 64-Bit Server
```

Install streamparse:

```
$pip install streamparse
```

Now you have streamparse installed. It will greatly simplify the creation of Python Storm projects and help you get a simple example up and running quickly. The following command creates an installation of the wordcount example.

```
$sparse quickstart wordcount
```

After watching the video and understanding the structure of topology definitions and the actual spout and bolt, run the word-count example use the following commands:

```
$cd wordcount
$sparse run
```

Python version 2.7

If you run into issues running Streamparse, check your python version. You need to use python 2.7. Follow the instructions below, and you will need to rerun the streamparse installation.

Install the required version of Python.

```
$sudo yum install python27-devel -y
```

You can see that the Python in your execution path (\$PATH) is still 2.6.X by checking the version again.

```
$python --version
Python 2.6.6
```

Rename the current version to reflect its correct version.

```
$mv /usr/bin/python /usr/bin/python266
```

Create a symbolic link from the file in the path to the version you want to execute.

```
$ln -s /usr/bin/python2.7 /usr/bin/python
```

Check that the link indeed refers to the intended version of Python.

```
$/usr/bin/python --version
Python 2.7.3
```

Check that the shell picks up the version of Python you intended.

```
$python --version
Python 2.7.3
```

Cloning a subdirectory

Git clones a complete repository. If you only want to clone the `exercise_2` directory, you need to use “sparse checkout”. Follow the below instructions to clone only the `exercise_2` sub directory.

```
$mkdir ex2
$cd ex2/
$git init
$git remote add -f origin git@github.com:UC-Berkeley-I-School/w205-spring-16-labs-exercises.git
$git config core.sparseCheckout true
$echo "exercise_2" >> .git/info/sparse-checkout
$git pull origin master
```

Exercise 2 Grading Guidelines:

Aspect	Description	Grading
Functionality	Complete and fully functional spout and bolts programs based on the description: Max 80 pt	<ul style="list-style-type: none">• Correct and executable spout: 25pt• Correct and executable bolts: 30pt• Correct topology: 5pt• Finalresults.py: 10 pt• Histogram.py: 10 pt
Design/Architecture	Complete documentation of the system: Max 20 pt	<ul style="list-style-type: none">• Correct Project/Folder/File/DB/table names: 3pt• Description/Execution instructions (readme.txt): 3pt• Architecture.pdf : 10pt• Screenshots: 2pt• Plot.png : 2 pt