Justin Schopick
Kushagra Singh
Chris Sultzbaugh
CS 172
May 4, 2018

# Twitter Crawler Report

## Collaboration:

We took a collaborative approach to this project. We all helped with each part to a certain extent. We all helped research the implementation of the crawler and storage. Kushagra did the actual coding of the crawler, Justin did the coding of the storage element, and Chris did the report write up. We all participated in the research and then split up the typing.

## System Details:

Our crawler uses the Twitter Streaming API in order to collect geolocated tweets. The Twitter Stream functionality uses a "listener" to detect when tweets come in. We set up the system to only gather data of geo-tagged tweets.The data is then stored in files of 10 MB each.

The Twitter API is necessary to gain access to the tweets. Our crawler gathers raw Json tweet data and returns it in correct formatting for storage-use. We used Jsoup to find the title of links within the tweets we parse through. We then use Gson to edit the title of the link onto the end of tweet being parsed to add that title to the data collected.

This system is set up within one package called "crawler." Inside this crawler package, we import the needed API and external libraries. We use one class called "streamTest" that has one main function and a couple of helper functions within it.

Our storage system takes the Json data and appends it to a new or current Json file as long as it is less than 10 MB. If no files under 10 MB exists, the system automatically initializes a new file with the proper formatting. As the Twitter Stream continues, our system continuously adds this data into the now-existing files.

The system uses a global variable to count which document we are on; this avoids any chance of overwriting data or causing storage malfunctions. Whenever a document is greater than or equal to 10MB, the system will append the end formatting to the file. After this, it increments the global variable by 1 to move on to the next document. The way our system is set up, it will add to existing files as well, rather than creating new files every time it begins.

**System Limitations:**

Using the Twitter Streaming API, limits us to only collect live tweet info. We currently do not have the capability to search older tweet info. Twitter streaming also limits the speed of our data collection as we cannot crawl for geolocated tweets. We have to wait for geo-tagged tweets to come in from the streamer, which is much slower because not everyone puts a location tag on their tweets.

We also have not yet implemented a stopping mechanism. That is once we start the collection process, unless someone or some outside factor stops the collection it continues. We plan to implement a max data collected parameter to correct this oversight.

We also had a problem we fixed that had us stuck for a bit. Our storage system was not working properly when it came to making a new file when data limit was reached for individual files. This caused us to create the global variable mentioned previously to help keep track of the file being created. This helped avoid creating duplicate file names and overwriting files.

**Deployment:**

Software needed to run our crawler:

- Java Development Kit must be installed in order to compile and run the program

- http://www.oracle.com/technetwork/java/javase/downloads/jdk8-downloads-2133151.html

External Packages Needed(these are all inside of the submission folder):

- jsoup-1.11.3.jar

- gson-2.6.2.jar

- twitter4j-async-4.0.4.jar

- twitter4j-core-4.0.4.jar

- twitter4j-examples-4.0.4.jar

- twitter4j-media-support-4.0.4.jar

- twitter4j-stream-4.0.4.jar

To Run(must be on unix/linux):

- Download the cs172part1_twitter zip file and unzip it to your desired location

- Open your terminal and find the location you stored the cs172part1_twitter directory in

- Open the cs172part1_twitter

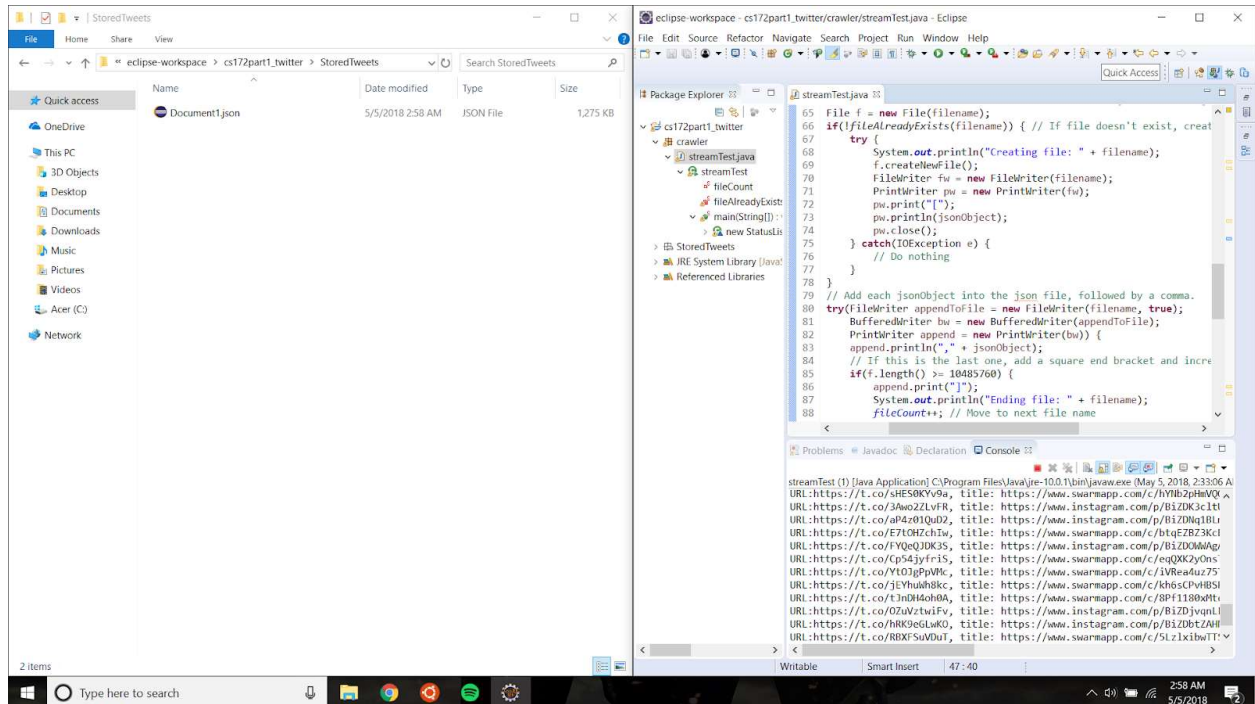- While in the cs172part1_twitter directory, type "./run.sh" into the terminal without quotation marks

To Terminate the Program:

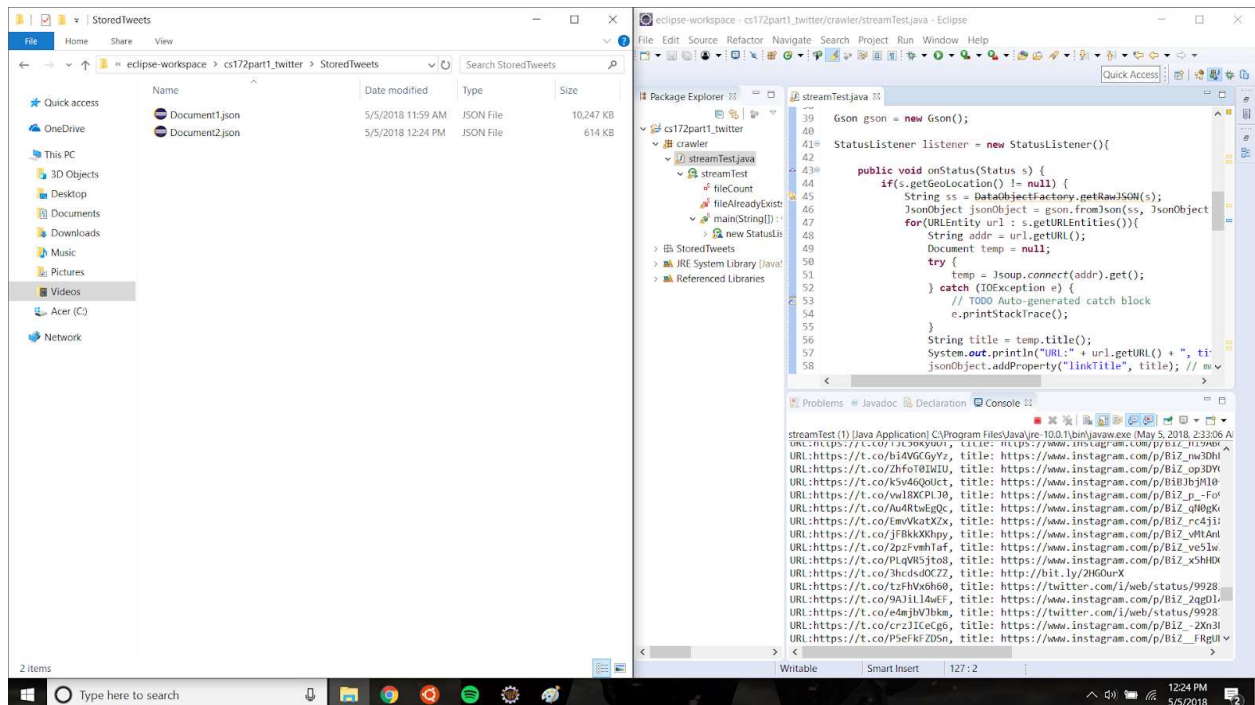- On Windows: Control+C

- On Mac: Command+C

**Screenshots:**

<u>Using the Eclipse IDE</u>

After Running for ~2 Hours:



After Running for ~12 Hours:

# Using the Terminal

## After Running for ~18 Hours: