

Architecture of W205 Exercise 2

Real Time Data Processing Using Apache Storm

This exercise implements the Apache Storm topology illustrated in Figure 1.

Tweet-spout (tweet.py) uses tweepy library to read sampled live Twitter stream and emits tweets to Parse-tweet-bolt (parse.py). Parse-tweet-bolt cleans up tweets and emits words to Count-bolt (wordcount.py). Count-bolt counts each word and insert or update word count into tweetwordcount table in tcount Postgres database.

Two python scripts (finalresult.py and histogram.py) were created to query tweetwordcount table to show how streamed and processed data can be used by other applications.

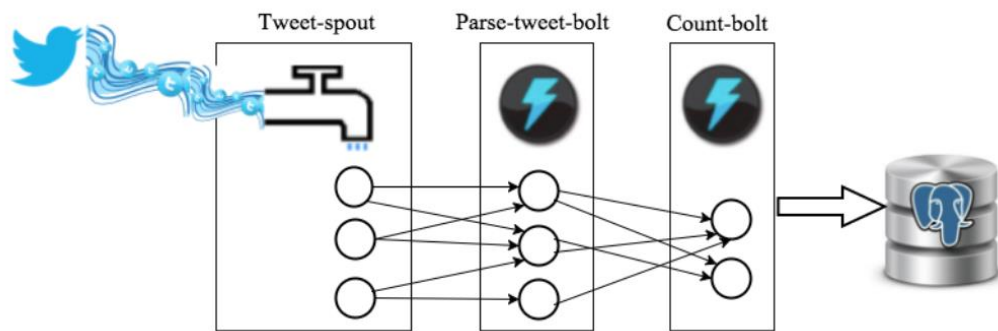


Figure 1: Application Topology

Table 1 shows the directories and files structure of the completed application. Please refer to the Readme.txt file on how to set up the environment and step-by-step instruction to run this application.

<ul style="list-style-type: none"> EXERCISE_2 <ul style="list-style-type: none"> screenshots <ul style="list-style-type: none"> screenshot-finalresults-all.png screenshot-finalresults-word.png screenshot-histogram.png screenshot-storm-components.png screenshot-twitterStream.png serving <ul style="list-style-type: none"> finalresults.py histogram.py tweetwordcount <ul style="list-style-type: none"> src <ul style="list-style-type: none"> bolts <ul style="list-style-type: none"> parse.py wordcount.py spouts <ul style="list-style-type: none"> tweets.py topologies <ul style="list-style-type: none"> tweetwordcount.clj virtualenvs <ul style="list-style-type: none"> tweetwordcount.txt config.json fabfile.py project.clj README.md tasks.py Architecture.pdf db_setup.sh Plot.png Readme.txt 	<p><u>EXERCISE_2</u>: This is the root folder.</p> <p><u>screenshots</u>: Screenshots of serving scripts result and streaming in progress.</p> <p><u>serving</u>: Scripts used to query data stored in database.</p> <p>tweetwordcount: This is the streaming application. Storm spouts and bolts components are under “src” and Storm topology is under “topologies”.</p> <p><u>Architecture.pdf</u>: This document.</p> <p><u>db_setup.sh</u>: Bash script used to set up database and table in Postgres.</p> <p><u>Plot.png</u>: Chart shows the top 20 most frequent words.</p> <p><u>Readme.txt</u>: Step by step instruction to set up and run this application.</p>
---	--

Table 1 Application Directories and File Structure