# Architecture

## Architecture diagram



tweets spout     parse bolt     word count bolt

postgres DB
tweetwordcount table

tweetwordcount topology

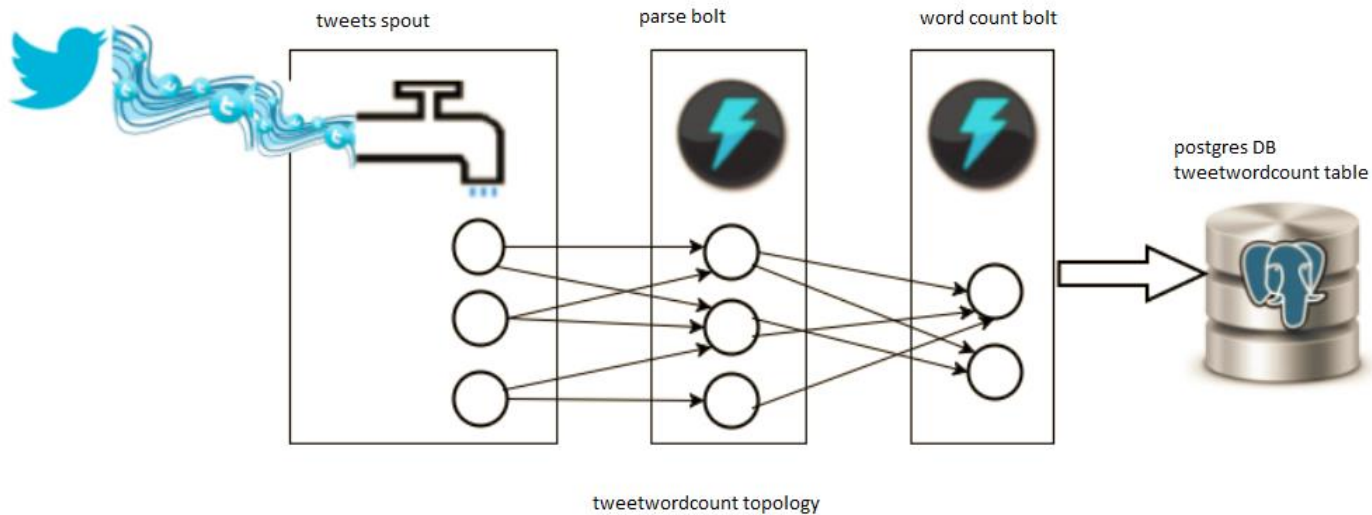## Folder structure

This PC > Documents > GitHub > w205_exercise_1_and_2 > exercise_2 > tweetwordcount

| Name | Date modified | Type | Size |
|------|---------------|------|------|
| src | 12/17/2015 8:51 AM | File folder | |
| topologies | 12/17/2015 8:51 AM | File folder | |
| virtualenvs | 12/17/2015 8:51 AM | File folder | |
| config.json | 12/4/2015 5:08 AM | JSON File | 1 KB |
| fabfile.py | 12/4/2015 5:08 AM | Python source file | 1 KB |
| project.clj | 12/4/2015 5:08 AM | CLJ File | 1 KB |
| README.md | 12/4/2015 5:08 AM | MD File | 0 KB |
| tasks.py | 12/4/2015 5:08 AM | Python source file | 1 KB |

# Submission 1 and 2

I tried to create DB named Tcount, but for some reason I cannot access it from sparse, so I just used default DB named postgres, and also used default folder tweetwordcount.



# Submission 3

When word count bolt counts words from twitter feed, it also reads from postgres DB. If this word already existed, it will read the existing count and increment using the local word count value. If the word does not exist, it will insert a key using current word and count.

wordcount.py ☒

```python
from __future__ import absolute_import, print_function, unicode_literals

from collections import Counter
from streamparse.bolt import Bolt

import psycopg2

class WordCounter(Bolt):

    def initialize(self, conf, ctx):
        self.counts = Counter()
        self.conn = psycopg2.connect(database="postgres", user="postgres", password="pass", host="localhost", port="5432")


    def process(self, tup):
        word = tup.values[0]

        # Write codes to increment the word count in Postgres
        # Use psycopg to interact with Postgres
        # Database name: postgres
        # Table name: Tweetwordcount
        # you need to create both the database and the table in advance.
        cur = self.conn.cursor()

        uWord = word
        uCount = self.counts[word] + 1

        cur.execute("SELECT word, count from Tweetwordcount WHERE word=%s",[uWord]);
        records = cur.fetchall()
        #update or insert
        if len(records) > 0:
            self.log('found %s records, before update:' % (len(records)))
            for rec in records:
                self.log('word = %s' % (rec[0]))
                self.log('count = %s' % (rec[1]))
            cur.execute("UPDATE Tweetwordcount SET count=%s WHERE word=%s", (uCount, uWord));
```

```
37
38    else:
39        self.log('record does not exist, try insert')
40        cur.execute("INSERT INTO Tweetwordcount (word,count) VALUES (%s, %s)", (uWord, uCount));
41
42
43    #Select
44    cur.execute("SELECT word, count from Tweetwordcount WHERE word=%s",[uWord]);
45    records = cur.fetchall()
46    self.log('found %s records, after update:' % (len(records)))
47    for rec in records:
48        self.log('word = %s' % (rec[0]))
49        self.log('count = %s' % (rec[1]))
50
51
52    # Increment the local count
53    self.counts[word] += 1
54    self.emit([word, self.counts[word]])
55
56    # Log the count - just to see the topology running
57    self.log('%s: %d' % (word, self.counts[word]))
58
59
```

After running for a few minutes, the table tweetwordcount got updated

## Submission 4

Sample query result for word hello.

Sample query without supplying arg, it will return all word counts, sorted by word in alphabetical ascending order.

Query histogram, using lower limit and upper limit, inclusive on both sides. Parameter format is: lower,upper, such as 3,8

Total number of occurences of "Transponder": 3.
Total number of occurences of "chance": 6.
Total number of occurences of "concert": 3.
Total number of occurences of "ya": 3.
Total number of occurences of "calling": 6.
Total number of occurences of "you'd": 3.
Total number of occurences of "ily": 3.
Total number of occurences of "Go": 6.
Total number of occurences of "Niall": 3.
Total number of occurences of "Snail!": 3.
Total number of occurences of "Johor": 3.
Total number of occurences of "defending": 3.
Total number of occurences of "won": 3.
Total number of occurences of "Boy": 3.
Total number of occurences of "award!": 3.
Total number of occurences of "trey": 3.
Total number of occurences of "wonderful": 3.
Total number of occurences of "battle": 3.
Total number of occurences of "bunch": 3.
Total number of occurences of "thing": 6.
(27env)[root@ip-172-31-50-135 ~]# python exercise_2/histogram.py 100,110
lowerLim is 100
upperLim is 110
Total number of occurences of "a": 108.
Total number of occurences of "go": 106.
Total number of occurences of "will": 107.
(27env)[root@ip-172-31-50-135 ~]#

to                    | 242
Teen                  | 210
just                  | 193
the                   | 192
as                    | 192
when                  | 191
we                    | 173
postgres=#
Using username "root".
Authenticating with public key "imported-openssh-key"
Last login: Sat Dec 19 09:52:11 2015 from 36.101.25.50

Welcome to a virtual machine image brought to you by RightScale!

(27env)[root@ip-172-31-50-135 ~]#

root@ip-172-31-50-135:~

C:\Users\linya\Documents\GitHub\w205_exercise_1_and_2\exercise_2

File | Home | Share | View

Documents > GitHub > w205_exercise_1_and_2 > exercise_2

Search exercise_2

| Name | Date modified | Type | Size |
|---|---|---|---|
| Exercise-2-Subject-205-Real Time Dat... | 12/16/2015 11:19 ... | Adobe Acrobat D... | 752 KB |
| finalresults.py | 12/19/2015 6:13 PM | Python source file | 1 KB |
| hello-stream-twitter.py | 12/4/2015 5:08 AM | Python source file | 2 KB |
| project2.docx | 12/19/2015 6:15 PM | Microsoft Word D... | 856 KB |
| psycopg-sample.py | 12/19/2015 5:57 PM | Python source file | 3 KB |
| README.md | 12/4/2015 5:08 AM | MD File | 1 KB |
| Twittercredentials.py | 12/15/2015 9:18 AM | Python source file | 1 KB |
| histogram.py | 12/19/2015 6:25 PM | Python source file | 1 KB |

12 items    1 item selected 663 bytes

exercise_2 - M3/root@54.152.233.179 - WinSCP

File Commands Mark Session View Help

Address /root/exercise_2

M3/root@54.152.233.179 | New Session

EX2Tweetw... | tweetword... | wordcount | Exercise-2-... Time Data ... | finalresults... | hello-strea... | histogram.py | psycopg-...

README.md | Twittercre... | Twittercre...

0 B of 758 KB in 0 of 11    1 hidden    SFTP-3    0:34:46