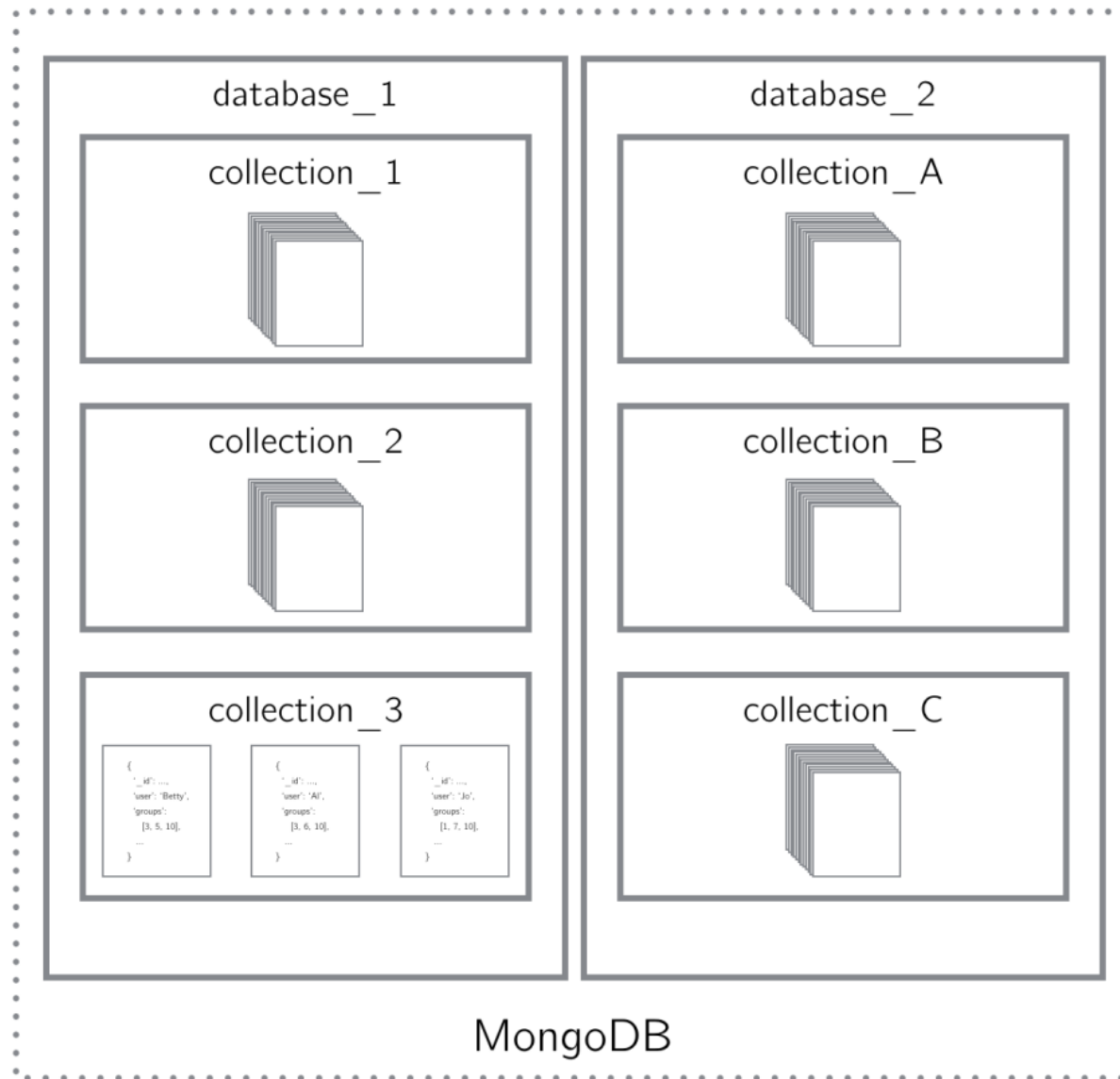


# Twitter, S3, and MongoDB

# MongoDB

- A Database Server
- Contains Databases that contain Collections that contain Documents
- Documents are stored using a JSON-like format; play very well with `dict` objects

# MongoDB



# NoSQL

- MongoDB is a NoSQL database
- No schema required
- Just add a document (`dict`) to a collection

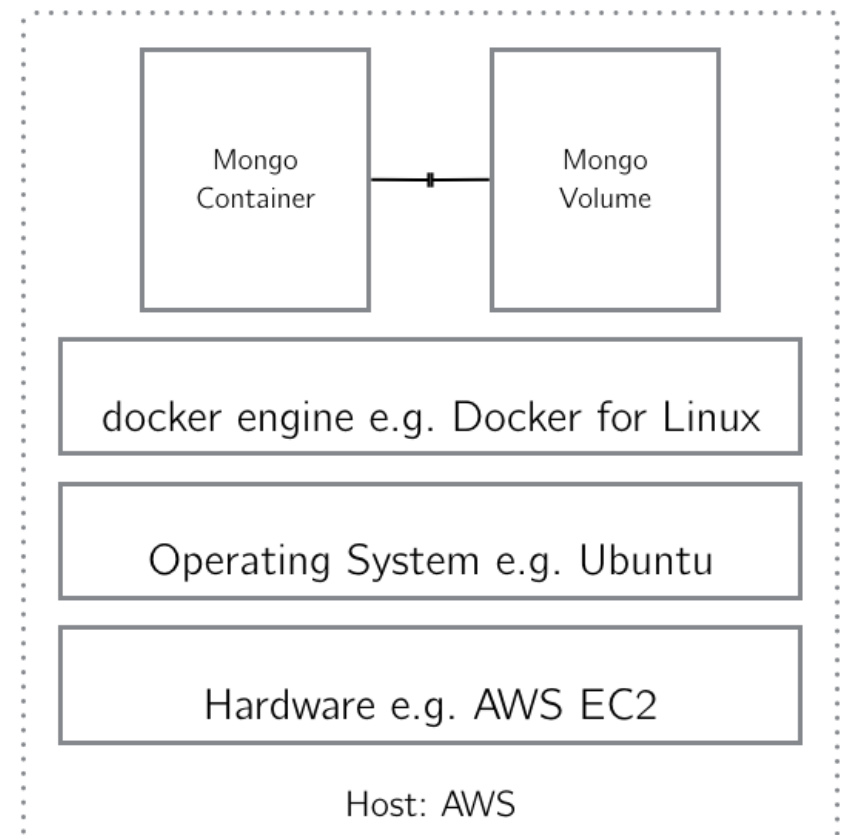
# NoSQL



You == Han Solo  
Your data == Chewbacca

# Configuring Mongo on AWS

- Considerations:
  - Networking
    - Solve via AWS Security Groups
  - Data Persistence
    - Solve via Docker Volumes



# Configuring Mongo on AWS

- Note that this configuration will use Mongo on a separate instance from the instance on which you are running Jupyter.
- If you are running Jupyter on AWS, you will need a second `t2.micro`.
- Accessing a database managed by Docker from a different AWS instance is actually easier than accessing a database managed by Docker on the same system.

# Set up a new t2.micro

- From the AWS EC2 Dashboard, select “Launch Instance.”
- On the Choose AMI tab, choose Ubuntu Server 16.04.
- On the Choose Instance Type tab, choose t2.micro.
- On the Add Storage tab, use the default setting of 8GB.
- On the Configure Security Group tab, choose “Create a new security group.”
- a. Confirm that inbound SSH traffic can be accepted over port 22 from anywhere.
- b. Add a rule that accepts inbound traffic over port 2376 from anywhere. This port will allow you to pull images from Docker Hub.
- c. Add a rule that accepts inbound traffic over port 27016 from anywhere. This is the default port for accessing MongoDB.
- Review and launch an instance, taking care to confirm that you have access to the SSH keys stored with your AWS account.



# Configure the new t2.micro

- Take note of the IP address of the newly configured AWS instance.
- SSH into the instance using that IP address.
- Install Docker via a shell script.
- Add the ubuntu user to the docker group.
- Log out and back in.

```
(local) $ ssh ubuntu@255.255.255.255  
(remote) $ curl -sSL https://get.docker.com/ | sh  
(remote) $ sudo usermod -aG docker ubuntu
```

# Run Mongo via Docker

- Pull the mongo image

```
$ docker pull mongo
```

- Create a New Data Volume

```
$ docker volume create mongo-dbstore
```

- Launch MongoDB as a Persistent Service

```
$ docker run -d --name this_mongo \  
    -v mongo-dbstore:/data/db \  
    -p 27017:27017 mongo
```

# Verify MongoDB Installation

- You can verify that you are running the mongo service by connecting to the running MongoDB via the MongoDB client, `mongo`, issued via `docker exec`.
- To do this, connect and then insert a trivial document to a mongo collection. You are inserting the JSON object `{"foo":1}` into the collection `test`. You then search for the document you inserted using the `.find()` command.

```
$ docker exec -it this_mongo mongo
> db.test.insert({"foo":1})
> db.test.find()
```

# Using MongoDB with Jupyter

- You will need to install the necessary Python library, pymongo

```
!pip install pymongo
```

- This should be run from a Jupyter server that is not on the same AWS instance as your Mongo server.

# pymongo

- pymongo is a Python module containing the MongoDB tools recommended for working with the database.
- You begin by instantiating a connection to MongoDB using `pymongo.MongoClient`.
- Here, you use the IP address of your AWS instance on which MongoDB is running.

```
from pymongo import MongoClient  
client = MongoClient('255.255.255.255', 27016)
```

# pymongo

- pymongo has a very useful “get or create” mechanism for both databases and collections.

```
client.database_names()
```

# pymongo

- Databases and collections are accessed using either attribute-style (`client.database_name`) or dictionary-style (`client['test-database']`).
- If the database exists, this method will return a reference to the existing database or collection (“get”). If the database does not exist, this method will create the database or collection and then return a reference to it (“create”).

# pymongo

- The creation happens at the time of insertion of a document.

```
db_ref = client.my_database  
client.database_names()
```

```
coll_ref = db_ref.my_collection  
client.database_names(), db_ref.collection_names()
```

```
sample_doc = {"name": "Joshua", "message": "Hi!",  
'my_array' : [1,2,3,4,5,6,7,9]}  
coll_ref.insert_one(sample_doc)
```

```
client.database_names(), db_ref.collection_names()
```



# Mongo and Twitter

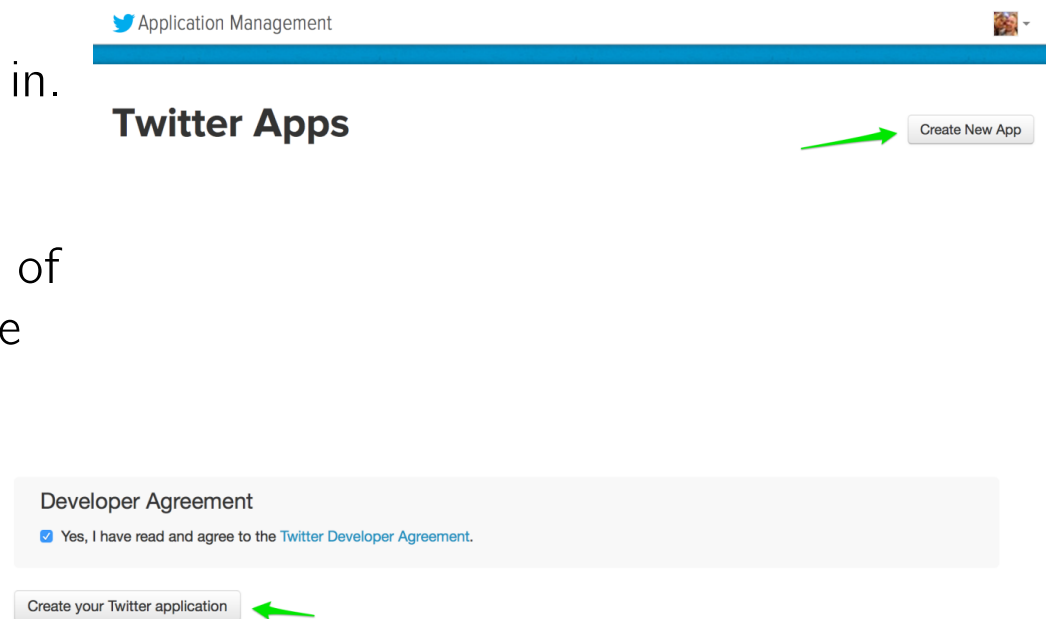
- To demonstrate a simple usage for MongoDB with Jupyter, you will implement a basic Twitter streamer that inserts captured tweets into a MongoDB collection.
- Twitter data represents an ideal use case for the NoSQL MongoDB.
- Each tweet obtained via the Twitter API is received as an unstructured nested JSON object.

# Mongo and Twitter

- Adding such an object to a SQL database would be a non-trivial task by any measure involving numerous foreign keys and JoinTables as the user seeks to manage each of the one-to-one, one-to-many, and many-to-one relationships built into the tweet.
- Adding such an object to Mongo, on the other hand, is a trivial task.
- MongoDB's native Binary JSON (BSON) format was designed precisely to accept such an object.

# Obtain Twitter Credentials

- In order to follow along, you must obtain API credentials for accessing the Twitter API. This is done by creating a Twitter application.
- In order to do this, follow these steps:
  1. Visit <https://apps.twitter.com> and sign in.
  2. Choose “Create New App”.
  3. Give the new app a name, description, and website. For your purposes, the values of these responses are irrelevant, although the website will need to have a valid URL structure.
  4. Agree to the Developer Agreement and click “Create your Twitter Application”.



The screenshot displays the Twitter Application Management page. At the top, there is a header with the Twitter logo and the text 'Application Management'. Below this, the main heading is 'Twitter Apps'. To the right of this heading is a button labeled 'Create New App', which is highlighted by a green arrow. Below the heading, there is a section titled 'Developer Agreement' with a checkbox and the text 'Yes, I have read and agree to the Twitter Developer Agreement.' At the bottom of this section is a button labeled 'Create your Twitter application', which is also highlighted by a green arrow.

# Obtain Twitter Credentials

- Next, you will need to access your credentials on the “Keys and Access Tokens” tab.

- You will need a total of four values:
  1. A consumer key (API Key)
  2. A consumer secret (API Secret)
  3. An access token
  4. An access token secret




## Test\_CJP

Test OAuth

[Details](#) [Settings](#) [Keys and Access Tokens](#) [Permissions](#)




### Application Settings

Keep the “Consumer Secret” a secret. This key should never be human-readable in your application.

Consumer Key (API Key)		←
Consumer Secret (API Secret)		←
Access Level	Read and write ( <a href="#">modify app permissions</a> )	
Owner	joshuacook	
Owner ID		

### Your Access Token

This access token can be used to make API requests on your own account's behalf. Do not share your access token secret with anyone.

Access Token		←
Access Token Secret		←
Access Level	Read and write	
Owner	joshuacook	
Owner ID		

# Load Twitter Credentials

- Load Twitter Credentials as Strings
- Replace this with your credentials:

```
CONSUMER_KEY = None  
CONSUMER_SECRET = None  
ACCESS_TOKEN = None  
ACCESS_SECRET = None
```

# Install the `twitter` library

- I prefer the `twitter` library over `tweepy`. I've found it to be better for streaming. Others have found `tweeps` better for historical data.

```
!pip install twitter
```

# Authentication

- You next instantiate a `twitter.OAuth` object using the Python twitter module and the credentials you have just loaded.
- You will use this object to facilitate your connection to Twitter's API.

```
from twitter import OAuth
oauth = OAuth(ACCESS_TOKEN, ACCESS_SECRET,
              CONSUMER_KEY, CONSUMER_SECRET)
```

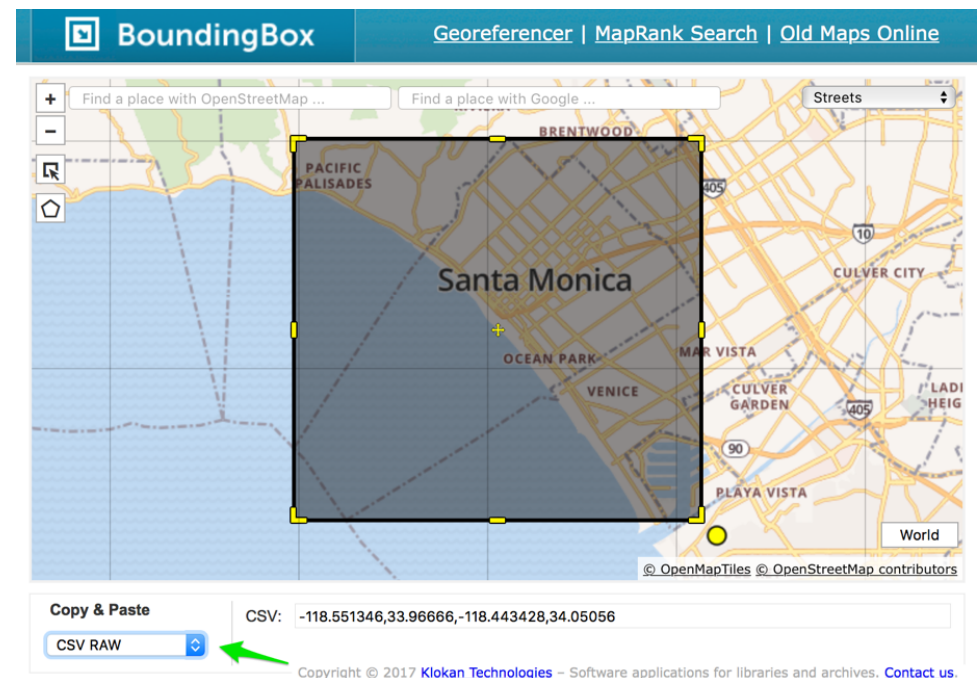
# Collect Tweets by Geolocation

- For this example, you will be using Twitter's Public Stream.
- Applications that are able to connect to a streaming endpoint will receive a sample of public data flowing through Twitter and will be able to do so without polling or concern of API rate limits.
- In other words, the Public Stream is a safe and sanctioned way to collect a sample of live public tweets.
- That said, even this sample will return a great deal of unordered data.



# Collect Tweets by Geolocation

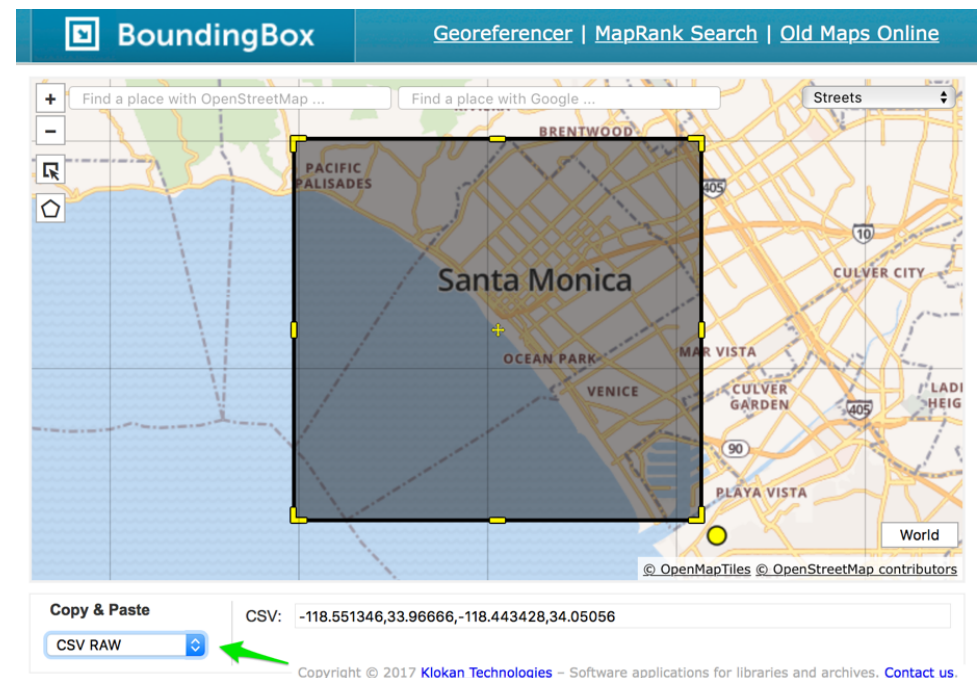
- In order to provide a modicum of order to your Twitter stream, you will restrict incoming tweets using a geolocation bounding box, or bbox. You can easily obtain a bbox for a location of interest using the Klokantech BoundingBox Tool.
- Let's obtain a bbox for Santa Monica, California in the United States, making sure to select CSV Raw as the copy and paste format.



```
los_angeles_bbox = "-118.551346,33.96666,-118.443428,34.05056"
```

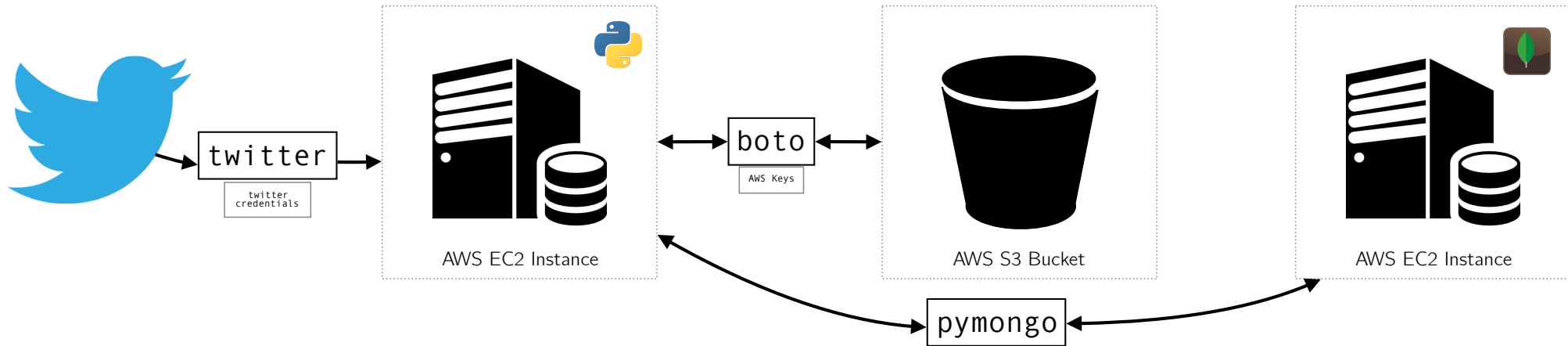
# Collect Tweets by Geolocation

- In order to provide a modicum of order to your Twitter stream, you will restrict incoming tweets using a geolocation bounding box, or bbox. You can easily obtain a bbox for a location of interest using the Klokantech BoundingBox Tool.
- Let's obtain a bbox for Santa Monica, California in the United States, making sure to select CSV Raw as the copy and paste format.



```
los_angeles_bbox = "-118.551346,33.96666,-118.443428,34.05056"
```

# Twitter Analysis Pipeline



- Input of Tweets
  - Use `twitter` to stream live tweets
  - Use `boto` to write json of tweets to S3 bucket
- Process Tweets into Mongo
  - Use `boto` to read json from S3 bucket
  - Use `pymongo` to write tweets to MongoDB
- Analyze Tweets
  - Use `pymongo` to load tweets into Jupyter for analysis

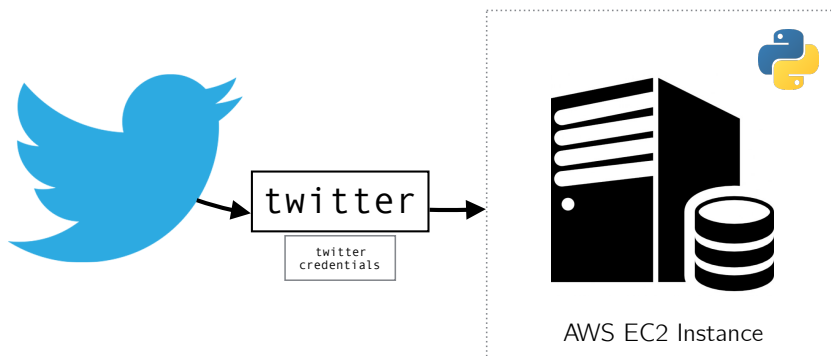
# Write Tweets to S3

## Create a Tweet Iterator

```
from lib import create_tweet_iterator

token          = "YOUR_SECRET_VALUE_HERE"
token_secret   = "YOUR_SECRET_VALUE_HERE"
consumer_key    = "YOUR_SECRET_VALUE_HERE"
consumer_secret = "YOUR_SECRET_VALUE_HERE"
bounding_box = "-118.5137323688,34.0001996344,-118.4702449172,34.0331651696"

tweet_iterator = create_tweet_iterator(token,
                                       token_secret,
                                       consumer_key,
                                       consumer_secret,
                                       bounding_box)
```

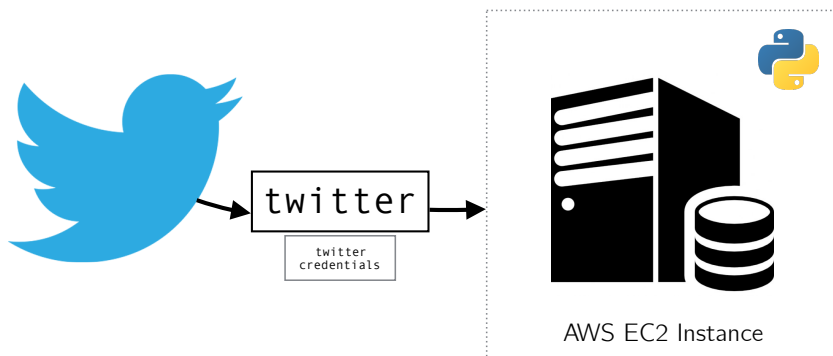


# Write Tweets to S3

## Build an Iterative Process for Collecting Tweets

#TODO

1. Use your `tweet_iterator`
2. Collect as many tweets as your are comfortable with and add them to a list of tweets



# Write Tweets to S3

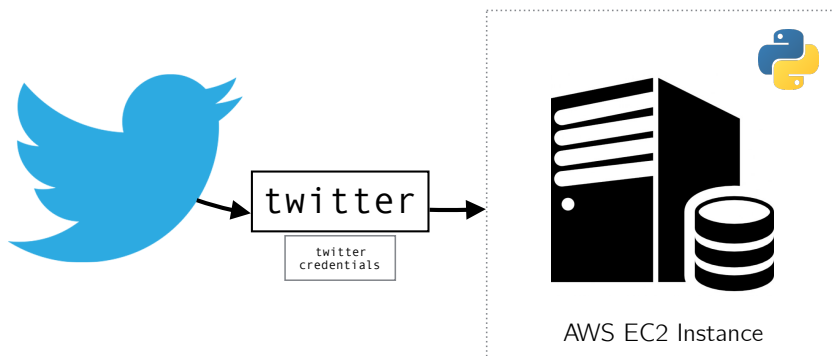
## Write list of tweets to JSON file on disk

`tweets` should be the list of tweets you previously created

```
import json
from lib import create_timestamped_filename

username = 'joshua'
filename = create_timestamped_filename(username)

with open(filename, 'w') as outfile:
    json.dump(tweets, outfile)
```

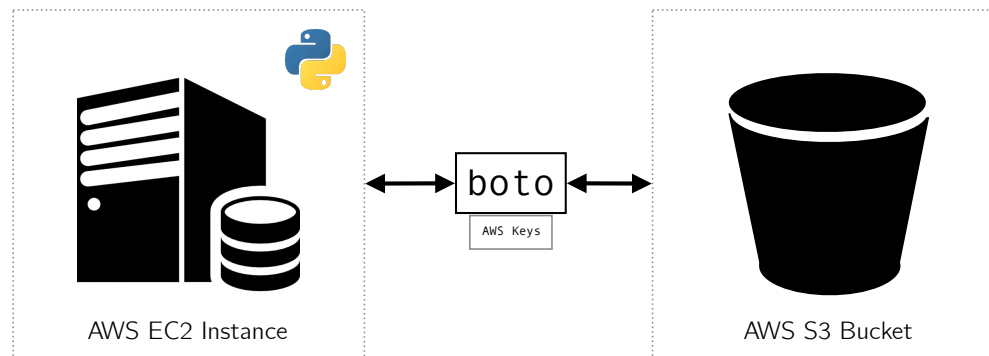


# Write Tweets to S3

## Write Tweet Files to S3, Step 1 - Create a Boto Client to S3

<https://boto3.readthedocs.io>

```
from lib import create_boto_client  
s3_client = create_boto_client()
```



# Write Tweets to S3

## Write Tweet Files to S3, Step 2 - Write a file to S3

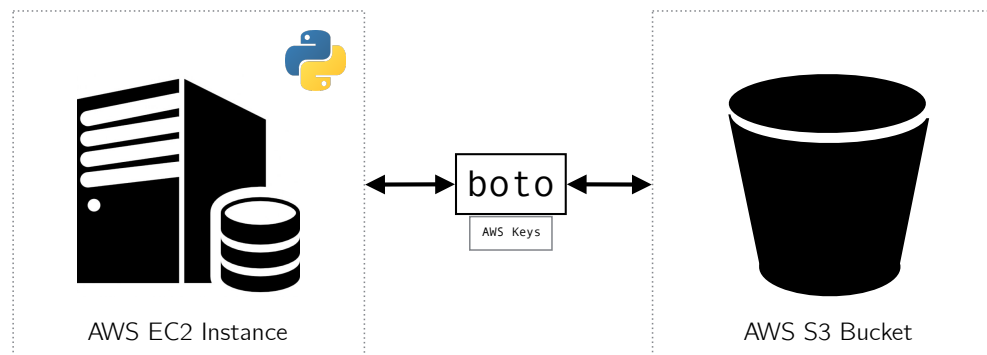
<https://boto3.readthedocs.io>

`filename` should be the name of a tweet file you saved earlier.  
`client` should be the name of the boto client you created in the last step.

```
from lib import write_file_to_S3

filename = 'tweets-joshua-2018-07-21_14-09-23-363540.json'
S3_BUCKET = 'uclax-data-science'

write_file_to_S3(s3_client, filename, S3_BUCKET)
```

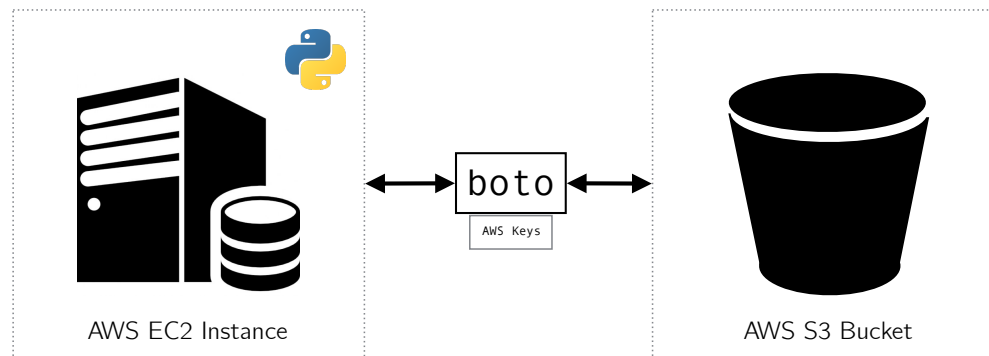




# Write Tweets to S3

## List Files on S3

```
from lib import list_files_in_S3_bucket  
S3_BUCKET = 'uclax-data-science'  
  
list_files_in_S3_bucket(s3_client, S3_BUCKET)
```



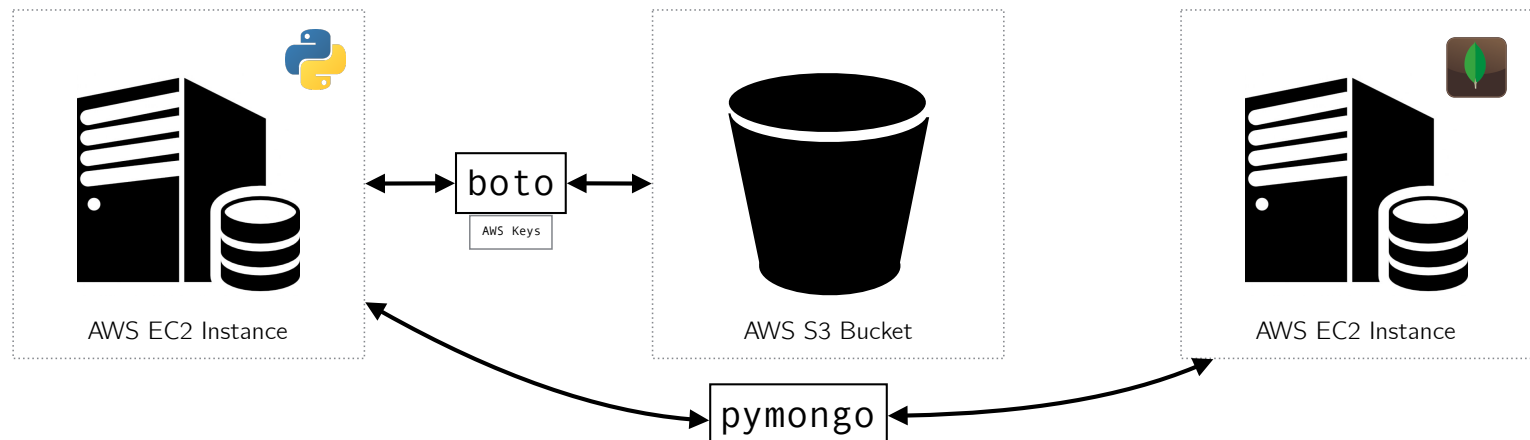
# Write Tweets to Mongo

## Read an object from S3

`key` should be the name of a tweet object on S3.  
`client` should be the name of the boto client you created previously.

```
from lib import read_object_from_S3
S3_BUCKET = 'uclax-data-science'

tweets_from_s3 = read_object_from_S3(s3_client, key, S3_BUCKET)
```



# Write Tweets to Mongo

## Write Tweets to Mongo

You will need to install `pymongo`.

You will need to read each json file from S3.

`tweet` is a single tweet you have read from an S3 object.

```
from lib import create_mongo_client_to_database_collection

collection_client = create_mongo_client_to_database_collection('twitter', 'tweets')
collection_client.insert_one(tweet)
```

