# HYBRID CLOUD SERVICES

Mentors:

-Xu (Simon) Chen

Team Members:

- Akash Singh
- Jaison Babu
- Surekha Jadhwani
- Vignesh Shanmuganathan

-

# 1. Introduction:

## 1.1 Problem Statement:

Imagine a case where data is regarding flight journeys, flight launch etc. Some of the data also belong to military airbases. The data of military airbases is something that must not be compromised, however, other data which could be treated public can be worked upon on a public cloud. One can even take the case of credit card no. of credit cards to not be compromised or there are so many cases where some portion of data based on any text like company name etc. or some regex is sensitive.

Our idea is to come up with a text based mechanism where once user specifies different text string to be treated as sensitive, the data containing those text tags won't be compromised. The sensitive information will remain on private cloud whereas all public data can be sent to the public cloud, thus making it possible for private firms to be able to scale to public cloud without compromising any of their sensitive information.

For our case, we are running the text based segregation on live twitter feed. Our reasons for doing this is to demonstrate that the idea can work on a very large scale (we worked with about 20 GB of data) and something which is applicable to live streaming given that there are several applications which have live streaming of data as their backbone. Given these two aspects of twitter feed, we built our framework around it and our framework could easily be adopted to any data which could be either from a warehouse or from some live feed applicable to enumerable scenarios enabling private firms and users to be able to exploit the public cloud. The overall goal is to be able to attract more users to the public cloud along with giving them an assurance that none of their sensitive information will be compromised.
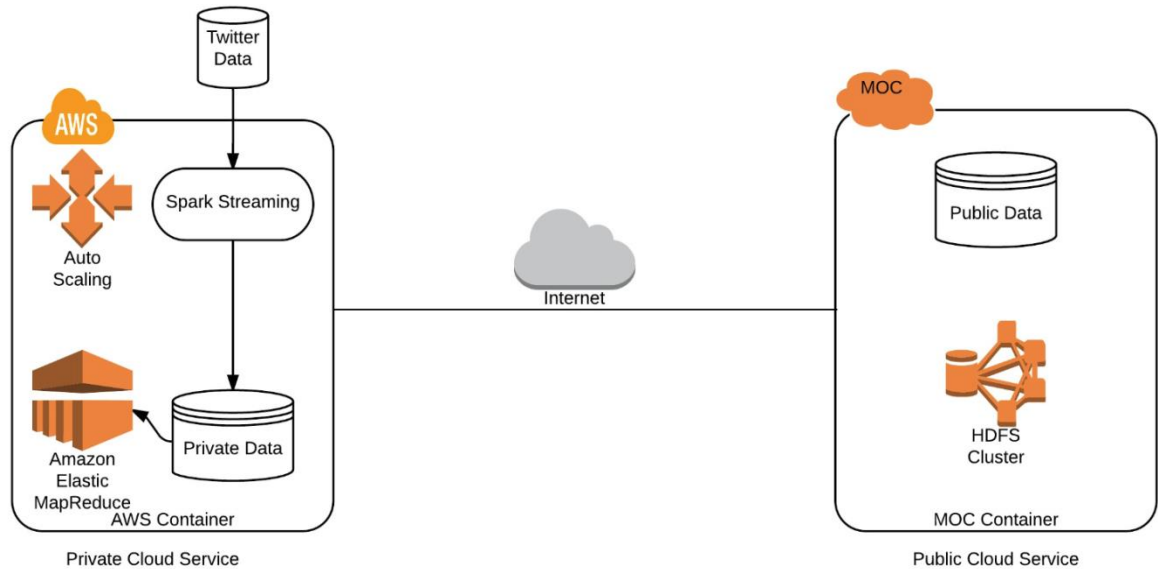
### 1.1.1 Post Segregation:

Data segregation is not enough, if analytics job cannot be run on this data. Our scope is to even run different analytics jobs like Geo-Location wise data distribution using time-zone tag, distribution of tweets based upon language, total counts of tweet in public and private cloud etc on this hybrid setup.

We even added an auto-scaling mechanism which makes assumption that public cloud capacity can be increased as and when needed easily and that by setting space limitations on the usage of private cloud, data can be shifted to public cloud automatically and securely, once space usage is exceeded. This should ideally happen in a VPN and our project makes an assumption that VPN could be setup between the private and public cloud.

## 2. Architecture of our implementation:

**Note:** We have considered MOC as public cloud and AWS as private cloud given the cost constraints.



*Figure 1: Hybrid cloud architecture*

### 2.1 Scope:

a. Creating a public cloud in MOC that contains non-sensitive data.
b. Creating a private cloud in AWS that contains sensitive data.
c. Set-up of Hadoop clusters in AWS and MOC.
d. Establishing communication between AWS and MOC for data transfers.
e. A spark streaming service running on AWS will gather real time data from twitter and partition data according to user configuration for classifying data as sensitive/non-sensitive. The proposed data classifiers include hash tag mechanism and linguistics.
f. Based on the hash tags, the created classifier will distribute data across AWS and MOC.
g. Running jobs in the following setups:
    i. Running data analysis job on public cloud spanning the complete dataset (sensitive and non-sensitive data).
    ii. Running data analysis job on private cloud spanning the complete dataset.
    iii. Running data analysis job on non-sensitive data on public cloud and on sensitive data on private cloud, in parallel.
h. Performing benchmark analysis on the above proposed configurations.

## 2.2 Features:

a. Easy user configuration by using configuration file.
b. Establishing communication between the cloud clusters.
c. Ability to run data analysis jobs on either/both clouds.
d. Analyzing real world data (twitter feed) on real world load on the hybrid setup.
e. Configuring the setup to sustain load of such load without major failures spanning both clouds.

# 3. Design & Implementation:

The project fundamentally compromises of 3 parts that involves:

a) Private cloud - Amazon Web Services (AWS)
b) Private cloud - Massachusetts Open Cloud (MOC)
c) Data - Live twitter feed streaming using Twitter4j API

## 1.1 Cluster Configurations:

| AWS Configuration | | MC Configurations | |
|---|---|---|---|
| No of Machine: | 2 | No of Machine: | 5 |
| Ram: | 1GB/machine | Ram: | 8GB/machine |
| Storage: | 15GB/machine | Storage: | 80GB/machine |
| Sensitive Data: | 4GB | Non-sensitive data: | 15gb |

## 1.2 Challenges faced:
a) Setting up spark streaming environment in AWS and Hadoop cluster in MOC.
b) Establishing secure communication between clouds.
c) Data segregation using user configuration for twitter data and auto scaling.

## 1.3 New Technology Learnt:
a) Private cloud and public cloud set-up in AWS and MOC.
b) Configuring HDFS in cluster mode.
c) Spark streaming job.
d) Hadoop analytical job.
e) Task tracking tool such as Trello.
f) Twitter 4j API in Java.

## 1.4 Implementation:
### 1.4.1 Data tagging mechanism:
Based on user twitter configuration file, the spark streaming job classifies the inflow of the twitter data into sensitive and non-sensitive data. A tweet that contains a hash tag mentioned in twitter configuration file is considered as sensitive data and is stored in HDFS on AWS (private cloud) and any tweet that does not contain the tag

mentioned in twitter configuration file is considered as non-sensitive data and is stored in MOC (public cloud) HDFS.

### 1.4.2 Streaming job:

A spark streaming job be an in memory continuous operation is opted for running a data gathering task that continuously retrieves twitter data from twitter server and based on data tagging mechanism the data is segregated into private and public data to gather enough sample data set that can be used to run analytical job to perform performance evaluation of different cloud processing jobs.

### 1.4.3 Analytical job:

Different analytical jobs such as classifying tweets based on location, language and time zone is run that combines execution of the job either in standalone mode or clubbed mode to evaluate performance of these jobs to showcase the benefits of the hybrid cloud service over standalone cloud models.

### 1.4.4 Performance benchmarking:

Based on different jobs run in public, private and hybrid model, we evaluate the options of running jobs in the above different models.

### 1.4.5 Auto-scaling:

Based on the user scaling configuration file, that contains 2 parameters (threshold of data in bytes, time in 10 sec intervals) [refer readme for more information properties], the job continuously monitors the private cloud in the mentioned time span if the current sensitive data size has exceeded the threshold mentioned and if true, the entire sensitive data is moved to public cloud to accommodate more data space in private cloud.

# 4. Results:

## 4.1 Segments of result of location wise tweet count : Left portion is result on private cloud and Right portion is result on public cloud.



## 4.2 Portion of result of running language based distribution on data on small dataset.

**4.3 Following is the result of the performance benchmarking we did on analytics job of geo-location wise tweet count.**

## 5. Limitations:

a) Currently the hybrid service has been modelled to work for data from twitter based on provided tweet tags.
b) Type of analytical job depends on the meta-data provided by twitter.
c) Communication security depends on Linux security layer.
d) Makes an assumption that a VPN can be created between MOC (public) and AWS (private) cloud.

## 6. Conclusion:

We were successful in coming up with a framework which can segregate data on the presence of text configuration between private and public cloud to enable use of public cloud to private firm. It overcomes the hurdle of sensitive information compromise, which otherwise is difficult to tackle. We did this on live data feed to demonstrate its abilities to work not just on data in a warehouse, but even on live feed on which several analytics are done these days. Our model can help in giving assurance to firms that they will be able to use public cloud without compromising any sensitive data and attract more people to the public cloud.

## 7. Future Scope:

e) Extend the security layer between private and public cloud by implementing an VPN communication layer
f) Current data segregation method relies on preprocessed text inputs, extend the data segregation for generic data by preprocessing using a spark job.
g) Extend auto-scaling to modify number of private and public cluster nodes based on requirement of analytical job.
h) Provide cluster backup in case of any fatal issue with the cluster.