



# Analyse Tweets using Flume 1.4, Hadoop 2.7 and Hive

**May 2015**

Dr.Thanachart Numnonda  
Certified Java Programmer  
thanachart@imcinstitute.com

Danairat T.  
Certified Java Programmer, TOGAF – Silver  
danairat@gmail.com



---

# Lecture: Understanding Flume

---

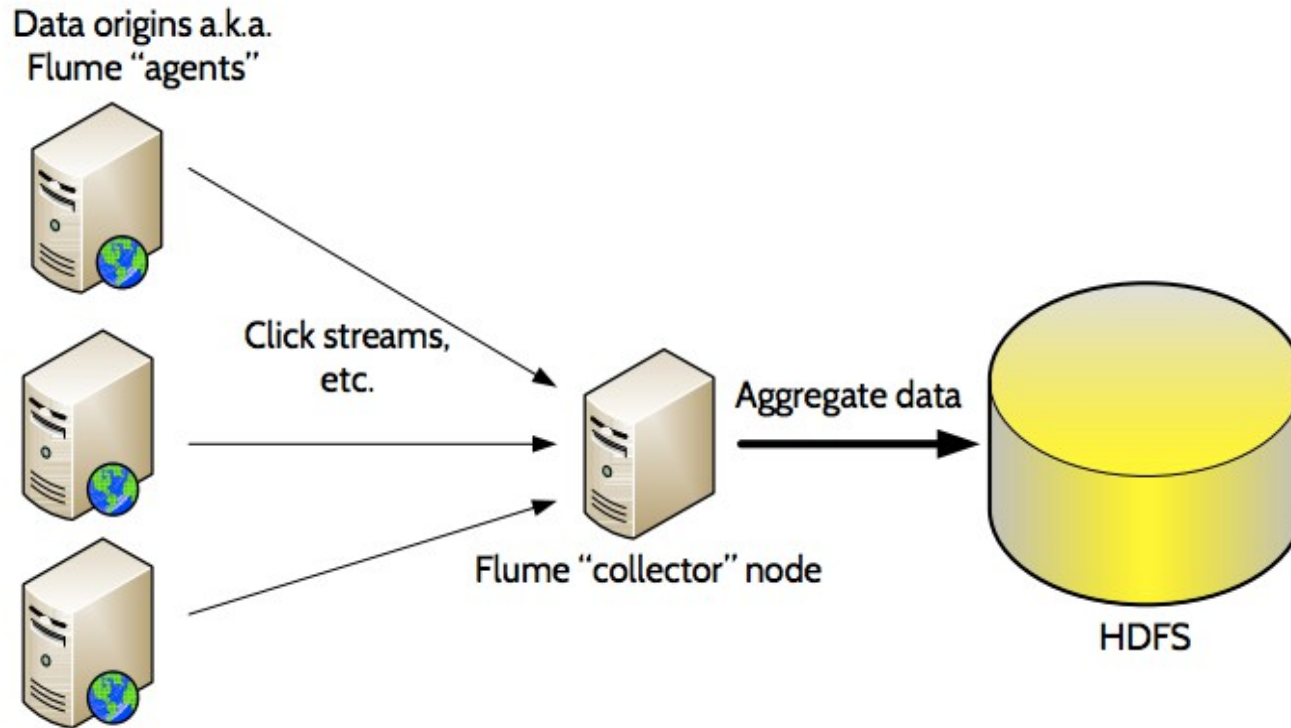
# Introduction



**Apache Flume is:**

- **A distributed data transport and aggregation system for event- or log-structured data**
- **Principally designed for continuous data ingestion into Hadoop... But more flexible than that**

# Architecture Overview

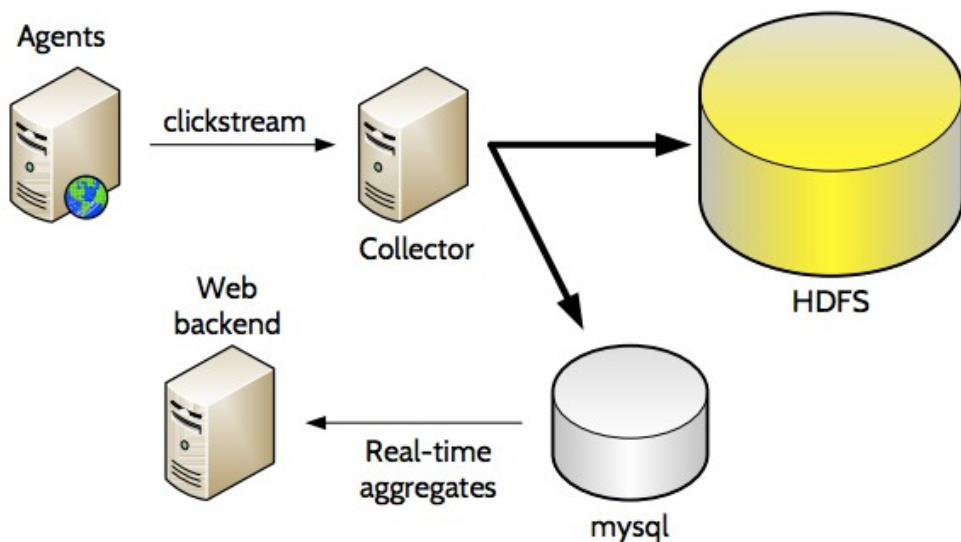


# Flume terminology

- Every machine in Flume is a **node**
- Each node has a **source** and a **sink**
- Some sinks send data to **collector** nodes, which aggregate data from many agents before writing to HDFS
- All Flume nodes heartbeat to/receive config from master
- Events enter Flume within seconds of generation

# Flume isn't an analytic system

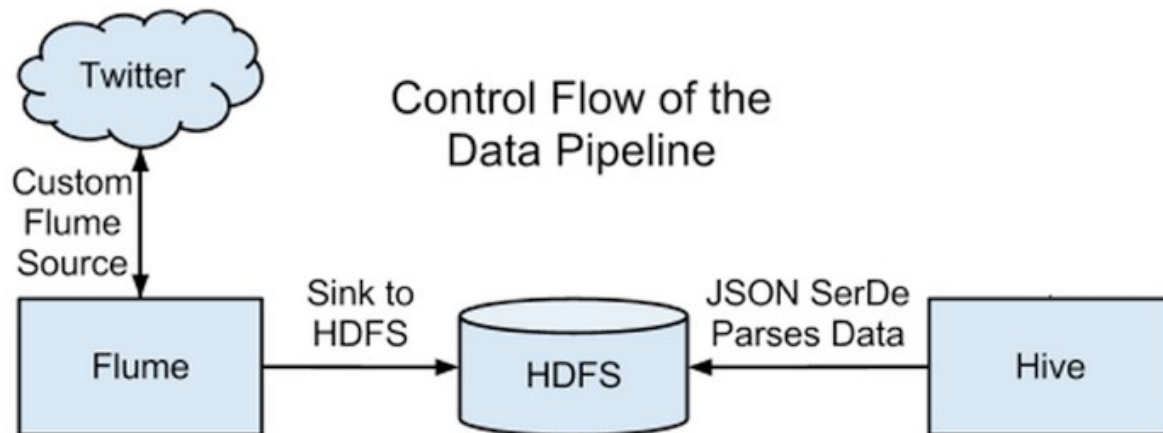
- No ability to inspect message bodies
- No notion of aggregates, rolling counters, etc



# Hands-On: Loading Twitter Data to Hadoop HDFS

---

# Exercise Overview





# 1. Installing Flume

## Install Flume binary file

```
$ wget  
http://apache.mirrors.hoobly.com/flume/1.4.0/apache  
-flume-1.4.0-bin.tar.gz  
  
$ tar -xvzf apache-flume-1.4.0-bin.tar.gz  
  
$ sudo mv apache-flume-1.4.0-bin flume  
  
$ sudo mv flume /usr/local  
  
$ rm apache-flume-1.4.0-bin.tar.gz
```

# 1. Installing Flume (cont.)

Edit \$HOME ./bashrc

```
$ sudo vi $HOME/./bashrc
```

```
export FLUME_HOME=/usr/local/flume
export PATH=$PATH:$HADOOP_HOME/bin:$HADOOP_HOME/sbin:$JAVA_HOME/bin:$HIVE_HOME/
bin:$PIG_HOME/bin:$SQOOP_HOME/bin:$FLUME_HOME/bin
```

```
$ exec bash
```

## 2. Installing a jar file

### Copy a jar file and edit conf file

```
$ wget http://files.cloudera.com/samples/flume-sources-1.0-SNAPSHOT.jar  
  
$ sudo mv flume-sources-1.0-SNAPSHOT.jar  
/usr/local/flume/lib/  
  
$ cd /usr/local/flume/conf/  
  
$ sudo cp flume-env.sh.template flume-env.sh  
  
$ sudo vi flume-env.sh
```

```
JAVA_HOME=/usr/lib/jvm/java-7-openjdk-amd64
```

```
# Give Flume more memory and pre-allocate, enable remote monitoring via JMX  
#JAVA_OPTS="-Xms100m -Xmx200m -Dcom.sun.management.jmxremote"
```

```
# Note that the Flume conf directory is always included in the classpath.  
FLUME_CLASSPATH="/usr/local/flume/lib/flume-sources-1.0-SNAPSHOT.jar"
```

# 3. Create a new Twitter App

Login to your Twitter @ twitter.com

Home

Notifications

Messages

Discover

Search Twitter

imcinstitute

@imcinstitute

TWEETS

88

FOLLOWING

9

FOLLOWERS

23

Get more from Twitter

Sign up

Follow 5 accounts

Complete your profile

What's happening?

สำนักข่าวเนชั่น

NNA

NATION NEWS AGENCY

สำนักข่าวเนชั่น @nnanews · 1h  
นายกฯ สั่งห้ามมือปล่อยน้ำเสีย ทำปลาในแม่น้ำปลาสดตายยกกระชัง พร้อมให้ทหาร เร่งนำปลาที่ตายขึ้นจากน้ำ #nna

3

3

hp

HP OpenNFV @hpnfv

Where would we be without the carrier networks? Follow @hpnfv to learn more about what's next for telecom.

HP OpenNFV

Promoted

Follow

Pongsuk Hiranprueck @nuishow · 2h

Facebook เริ่มทดสอบการเชื่อมต่อระหว่าง WhatsApp กับ Facebook บน Android แล้ว buff.ly/1xULvS9 #beartai

6

3

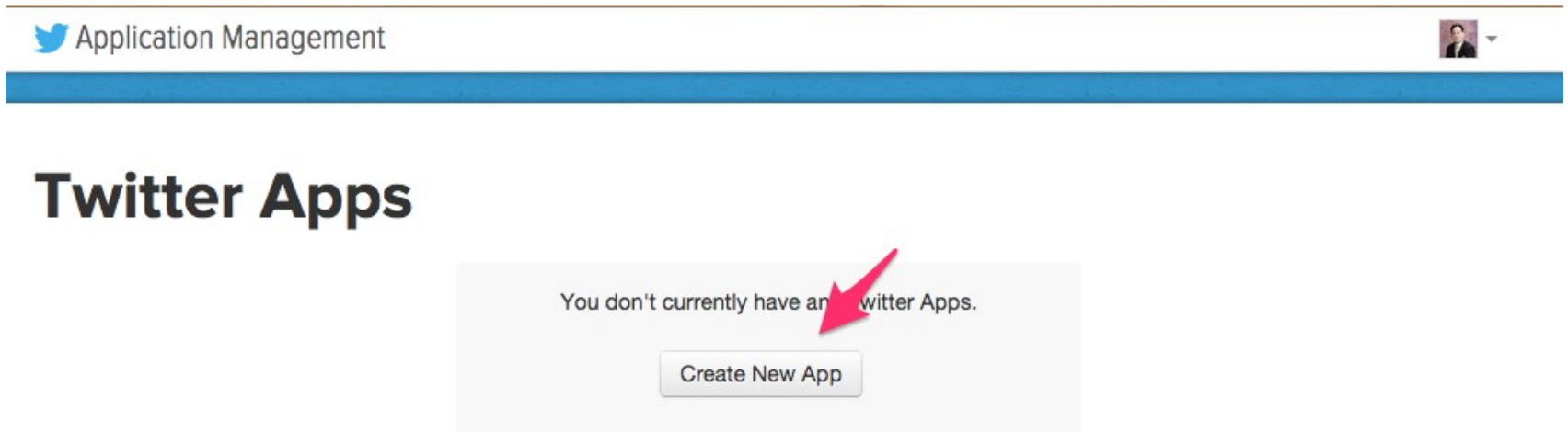
View summary

Big Data using Hadoop workshop

Danairat T., danairat@gmail.com: Thanachart Numnonda, thanachart@imcinstitute.com Apr 2015

### 3. Create a new Twitter App (cont.)

Create a new Twitter App @ apps.twitter.com



### 3. Create a new Twitter App (cont.)

Enter all the details in the application:

 Application Management



## Create an application

### Application Details

**Name \***



Your application name. This is used to attribute the source of a tweet and in user-facing authorization screens. 32 characters max.

**Description \***

Your application description, which will be shown in user-facing authorization screens. Between 10 and 200 characters max.









**Website \***


Your application's publicly accessible home page, where users can go to download, make use of, or find out more information about your application. This fully-qualified URL is used in the source attribution for tweets created by your application and will be shown in user-facing authorization screens.

(If you don't have a URL yet, just put a placeholder here but remember to change it later.)

### 3. Create a new Twitter App (cont.)

Your application will be created:

 <https://apps.twitter.com/app/8158163>       


Application Management 

Your application has been created. Please take a moment to review and adjust your application's settings.

## IMC\_Institute\_App

Test OAuth

Details Settings Keys and Access Tokens Permissions

 IMC Institute Demo App  
<http://www.imcinstitute.com>

### Organization

Information about the organization or company associated with your application. This information is optional.

|                      |      |
|----------------------|------|
| Organization         | None |
| Organization website | None |

### Application Settings

### 3. Create a new Twitter App (cont.)

Click on Keys and Access Tokens:

 Application Management



## IMC\_Institute\_App

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)       | MjpswndxVj27yInpOoSBrnfLX                          |
| Consumer Secret (API Secret) | QYmuBO1smD5Yc3zE0ZF9ByCgeEQxnxUmhRVCIsAvPFudYVJC4a |
| Access Level                 | Read and write (modify app permissions)            |
| Owner                        | imcinstitute                                       |
| Owner ID                     | 921172807  |



### 3. Create a new Twitter App (cont.)

**Click on Keys and Access Tokens:**

#### Application Actions

Regenerate Consumer Key and Secret

Change App Permissions

#### Your Access Token

*You haven't authorized this application for your own account yet.*

*By creating your access token here, you will have everything you need to make API calls right away. The access token generated will be assigned your application's current permission level.*

#### Token Actions

Create my access token



### 3. Create a new Twitter App (cont.)

**Your Access token got created:**

#### 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        | 921172807-EfMXJj6as2dFECdH1vDe5goyTHcxPrF1RIJozqgx |
| Access Token Secret | HbpZEVip3D5j80GP21a37HxA4y10dH9BHcgEFXUNcA9xy      |
| Access Level        | Read and write                                     |
| Owner               | imcinstitute                                       |
| Owner ID            | 921172807  |

#### Token Actions

Regenerate My Access Token and Token Secret

Revoke Token Access

## 4. Configuring the Flume Agent

Copy the flume.conf file from the following url:

<https://github.com/cloudera/cdh-twitter-example/blob/master/flume-sources/flume.conf>

```
$ vi /usr/local/flume/conf/flume.conf
```

### flume.conf file

```
TwitterAgent.sources.Twitter.consumerKey =25dhte51roU2s0rBogRh8Zwa5
TwitterAgent.sources.Twitter.consumerSecret =h1PgWYoqoIa9DCryueyStX80tSvrt50J1vWwISN58LthFdm90:
TwitterAgent.sources.Twitter.accessToken =921172807-EzYUG4Tb1AZegzta1NhX27MUsxlzw7VNHGnzPIt5
TwitterAgent.sources.Twitter.accessTokenSecret =mA9d8AQ1Vjapp7MatQESgdrZdyJ55xrrbZv1RjefKs1Uv
```

```
TwitterAgent.sources.Twitter.keywords = hadoop, big data, analytics, bigdata, cloudera, data science, data scientiest, business intelligence, mapreduce, data warehouse, data warehousing, mal out, hbase, nosql, newsql, businessintelligence, cloudcomputing
```

```
TwitterAgent.sinks.HDFS.channel = MemChannel
TwitterAgent.sinks.HDFS.type = hdfs
TwitterAgent.sinks.HDFS.hdfs.path = hdfs://localhost:9000/user/flume/tweets/
TwitterAgent.sinks.HDFS.hdfs.fileType = DataStream
TwitterAgent.sinks.HDFS.hdfs.writeFormat = Text
TwitterAgent.sinks.HDFS.hdfs.batchSize = 1000
TwitterAgent.sinks.HDFS.hdfs.rollSize = 0
TwitterAgent.sinks.HDFS.hdfs.rollCount = 10000
TwitterAgent.channels.MemChannel.type = memory
TwitterAgent.channels.MemChannel.capacity = 10000
TwitterAgent.channels.MemChannel.transactionCapacity = 100
```

# Fixing Bug for running Flume1.4 on Hadoop 2.x

Need to remove file guava-10.0.1.jar and  
protobuf-java-2.4.1.jar

```
$ cd /usr/local/flume/rm
```

```
$ rm guava-10.0.1.jar
```

```
$ rm protobuf-java-2.4.1.jar
```

## 5. Fetching the data from twitter

```
$ flume-ng agent -n TwitterAgent -c conf -f  
/usr/local/flume/conf/flume.conf
```

Wait for 60-90 seconds and let flume stream the data on HDFS, then press Ctrl-c to break the command and stop the streaming. (Ignore the exceptions)

```
15/05/08 04:07:30 INFO hdfs.HDFSDataStream: Serializer = TEXT, UseRawLocalFileSystem = false  
15/05/08 04:07:31 INFO hdfs.BucketWriter: Creating hdfs://localhost:9000/user/flume/tweets//FlumeData.1431058050787.tmp  
15/05/08 04:08:04 INFO hdfs.BucketWriter: Renaming hdfs://localhost:9000/user/flume/tweets/FlumeData.1431058050787.tmp to hdfs://localhost:9000/user/flume/tweets/FlumeData.1431058050787  
15/05/08 04:08:10 INFO hdfs.BucketWriter: Creating hdfs://localhost:9000/user/flume/tweets//FlumeData.1431058050788.tmp  
15/05/08 04:08:40 INFO hdfs.BucketWriter: Renaming hdfs://localhost:9000/user/flume/tweets/FlumeData.1431058050788.tmp to hdfs://localhost:9000/user/flume/tweets/FlumeData.1431058050788
```

## 6. View the straming data

```
$ hdfs dfs -ls /user/flume/tweets
```

```
ubuntu@ip-172-31-2-61:~$ hdfs dfs -ls /user/flume/tweets
Found 7 items
-rw-r--r-- 1 ubuntu supergroup 58224 2015-05-08 04:08 /user/flume/tweets/FlumeData.1431058050787
-rw-r--r-- 1 ubuntu supergroup 187394 2015-05-08 04:08 /user/flume/tweets/FlumeData.1431058050788
-rw-r--r-- 1 ubuntu supergroup 161913 2015-05-08 04:09 /user/flume/tweets/FlumeData.1431058050789
-rw-r--r-- 1 ubuntu supergroup 29949 2015-05-08 04:09 /user/flume/tweets/FlumeData.1431058050790
-rw-r--r-- 1 ubuntu supergroup 70739 2015-05-08 04:10 /user/flume/tweets/FlumeData.1431058050791
-rw-r--r-- 1 ubuntu supergroup 44753 2015-05-08 04:11 /user/flume/tweets/FlumeData.1431058050792
-rw-r--r-- 1 ubuntu supergroup 44426 2015-05-08 04:11 /user/flume/tweets/FlumeData.1431058050793
```

```
$ hdfs dfs -cat /user/flume/tweets/FlumeData.1431058050787
```

```
{"filter_level":"low","retweeted":false,"in_reply_to_screen_name":null,"possibly_sensitive":false,"truncated":false,"lang":"en","in_reply_to_status_id_str":null,"id":596526811430924288,"in_reply_to_user_id_str":null,"timestamp_ms":"1431058049434","in_reply_to_status_id":null,"created_at":"Fri May 08 04:07:29 +0000 2015","favorite_count":0,"place":null,"coordinates":null,"text":"RT @JanJekielek: Thanks! MT @rafat Building a Vertical Media Brand Around Trendlines #business intelligence #skift #inmawc15 http://t.co/qgJ\u0026","contributors":null,"retweeted_status":{"f
```

```
$ hdfs dfs -rm /user/flume/tweets/*.tmp
```

## 7. Analyse data using Hive

### Get a Serde Jar File for parsing JSON file

```
$ wget  
http://files.cloudera.com/samples/hive-serdes-1.0-SNAPSHOT.jar  
  
$ mv hive-serdes-1.0-SNAPSHOT.jar /usr/local/apache-hive-  
1.1.0-bin/lib/  
  
$ hive
```

### Register the Jar file.

```
hive> ADD JAR /usr/local/apache-hive-1.1.0-bin/lib/hive-  
serdes-1.0-SNAPSHOT.jar;
```

## 7. Analyse data using Hive (cont.)

### Running the following hive command

```
1 CREATE EXTERNAL TABLE tweets (  
2     id BIGINT,  
3     created_at STRING,  
4     source STRING,  
5     favorited BOOLEAN,  
6     retweet_count INT,  
7     retweeted_status STRUCT<  
8         text:STRING,  
9         user:STRUCT<screen_name:STRING,name:STRING>>,  
10    entities STRUCT<  
11        urls:ARRAY<STRUCT<expanded_url:STRING>>,  
12        user_mentions:ARRAY<STRUCT<screen_name:STRING,name:STRING>>,  
13        hashtags:ARRAY<STRUCT<text:STRING>>>,  
14    text STRING,  
15    user STRUCT<  
16        screen_name:STRING,  
17        name:STRING,  
18        friends_count:INT,  
19        followers_count:INT,  
20        statuses_count:INT,  
21        verified:BOOLEAN,  
22        utc_offset:INT,  
23        time_zone:STRING>,  
24    in_reply_to_screen_name STRING  
25 )  
26 ROW FORMAT SERDE 'com.cloudera.hive.serde.JSONSerDe'  
27 LOCATION '/user/flume/tweets';
```

<http://www.thecloudavenue.com/2013/03/analyse-tweets-using-flume-hadoop-and.html>



## 7. Analyse data using Hive (cont)

### Finding user who has the most number of followers

```
hive> elect user.screen_name, user.followers_count c from  
tweets order by c desc;
```

```
Starting Job = job_201504051617_0010, Tracking URL = http://localhost:50030/jobdetails.jsp?jobid=job_201504051617_0010  
Kill Command = /usr/local/hadoop/libexec/./bin/hadoop job -kill job_201504051617_0010  
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1  
2015-04-06 15:37:27,782 Stage-1 map = 0%, reduce = 0%  
2015-04-06 15:37:31,837 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 1.27 sec  
2015-04-06 15:37:39,899 Stage-1 map = 100%, reduce = 33%, Cumulative CPU 1.27 sec  
2015-04-06 15:37:40,908 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 2.42 sec  
MapReduce Total cumulative CPU time: 2 seconds 420 msec  
Ended Job = job_201504051617_0010  
MapReduce Jobs Launched:  
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 2.42 sec HDFS Read: 170686 HDFS Write: 687 SUCCESS  
Total MapReduce CPU Time Spent: 2 seconds 420 msec  
OK  
vinnaum 11523  
navchatterji 5485  
HCITExpert 4751  
NWDCScoop 4097  
7wdata 3005  
MotivasiMariaP 2007  
WesleyBackelant 1977  
IFTTMarketing 1307  
jonathangibs 968  
ephraimcohen 914  
feshob 716  
DKajouri 713
```

# Thank you

[www.imcinstitute.com](http://www.imcinstitute.com)

[www.facebook.com/imcinstitute](https://www.facebook.com/imcinstitute)