# Social Media for Emergency Management

Vimala Nunavath

vimala.nunavath@uia.no

# Agenda

- Recap
- Social media
- Social media use in disasters
- What is streaming?
- Why Spark streaming?
- Spark streaming components
- Social network graphs
- Example: Real-time twitter data stream processing with Apache Spark
- Presentations by you!!

# Recap

- Central Themes:
  - Introduction to Emergency management (Session 1)
  - Introduction to Big data (Session 2)
  - Big Data sources for EM (Session 3)
  - Sensors/IOT for EM (Session 4)
  - **Social media for EM (Session 5)**
  - **Natural language processing (NLP) for EM and visualization/Dashboards (Session 5)**
  - **Machine learning for EM (Session 5)**
- Essay and student programming project

# Social media

# Social media in general

- The term "social media" refers to Internet-based applications that enable people to communicate and share resources and information.

- e.g.,



- Huge volumes of data are generated every

minute, a phenomenon commonly referred to by

researchers as **big data**, **information overload** or

**data deluge**.

- Evolving phenomenon
- New technologies have enabled people to interact and

share information through media.

# Social media in disasters

- Social media (SM) plays a vital role in disaster response and recovery by providing response information before, during and after disasters.

- Social media are changing the way people communicate not only in their day-to-day lives, but also during disasters that threaten public health.

- Engaging with and using emerging social media may well place the emergency-management community, including medical and public health professionals, in a better position to respond to disasters.

- The effectiveness of public emergency system relies on routine attention to preparedness, agility in responding to daily stresses and catastrophes, and the resilience that promotes rapid recovery. Social media can enhance each of these component efforts.

# Social media in disasters

- The use of social media for emergencies and disasters on an organizational level may be conceived of as two broad categories:

  - To disseminate information and receive user feedback via incoming messages, wall posts, and polls.

  - An emergency management tool. Systematic usage might include:

    1) using the medium to conduct emergency communications and issue warnings;

    2) using social media to receive victim requests for assistance;

    3) monitoring user activities and postings to establish situational awareness; and

    4) using uploaded images to create damage estimates, among others.

# Social media in disasters

- For instance:
  - 2018 Indonesia Earthquake
  - 2012 Hurricane Sandy
  - 2012 Utøya bombing

# SM for Situational Awareness

- Social media could be used to alert emergency managers and officials to certain situations by monitoring the flow of information from different sources during an incident.

- Monitoring information flows could help establish situational awareness.

- Situational awareness: the ability to identify, process, and comprehend critical elements of an incident or situation.

- Obtaining real-time information as an incident unfolds can help officials determine where people are located, assess victim needs, and alert citizens and first responders to changing conditions and new threats.

# Challenges with Social media data

- Providing inaccurate and false information
- Malicious use of social media during disasters
- Technological limitations
- Privacy issues

# Accessing Social Media data

# Streaming Twitter data with Spark

# What is Streaming?

- **Data streaming** is a technique for transforming data so that it can be processed as a **steady** and **continuous** stream.

- Streaming technologies are becoming increasingly important with the growth of the internet.

# Why Spark Streaming



**Spark Streaming** is used to stream real-time data from various sources like twitter, Facebook, and geographical systems and perform powerful analytics to help during disasters.

# Spark streaming features

# What is Spark Streaming?

- It is an extension of the *core Spark API* that enables

   - Scalable, high-throughput, fault-tolerant *stream processing of live data streams.*



   - Data can be ingested from many sources

   **Figure: Data from a variety of sources to various storage systems**

   e.g., Kafka, Flume, Kinesis, or TCP sockets

   - It can be processed using complex algorithms expressed with high-level functions *like map, reduce, join and window.*

   **-** processed data can be pushed out *to filesystems, databases, and live dashboards.*

# Spark streaming overview



Figure: Overview Of Spark Streaming

# Spark Streaming

- Spark Streaming receives live input data streams and divides the data into batches, which are then processed by the Spark engine to generate the final stream of results in batches.

- It provides a high-level abstraction called **discretized stream** or **Dstream**.
- We can write Spark Streaming programs in Scala, Java or Python



*Figure: Incoming streams of data divided into batches*

# Streaming fundamentals

# Streaming Context



Figure: Spark Streaming Context

- Consumes a stream of data in Spark.
- Registers an **InputDstream** to produce a **Reciever** object.
- It is the main entry point for Spark functionality.
- Spark provides a number of default implementations of sources like Twitter, Akka Actor, and ZeroMQ that are accessible from the context.



Figure: Default Implementation Sources

# Streaming Context - Initialization

- A StreamingContext object can be created from a SparkContext object.

- A SparkContext represents the connection to a Sprak cluster and can be used to create RDDs, accumulators and broadcast variables on that cluster.

```
import org.apache.spark._
import org.apache.spark.streaming._
var ssc = new StreamingContext(sc,Seconds(1))
```

# Dstream

- Discretized stream (Dstream) *is the basic* abstraction provided by Spark Streaming.

- It represents a continuous stream of data.

- It is received from source or from a processed data stream generated by transforming the input stream.

- Internally, a DStream is represented by a continuous series of Resilient Distributed Datasets (RDDs). Each RDD contains data from a certain interval.



**Figure: Input data streams divided into discrete chunks of data**

# Discretized Stream Processing

- Dstreams: Run a streaming computation as a series of very small, deterministic batch jobs

  - Chop up the live stream into batches of X seconds

  - Spark treats each batch of data as RDDs and processes them using RDD operations

  - Finally, the processed results of the RDD operations are returned in batches



live data stream

Spark Streaming

batches of X seconds

Spark

processed results

# Discretized Stream Processing



- Dstreams: Run a streaming computation as a series of very small, deterministic batch jobs

  - Batch sizes as low as ½ second, latency of about 1 second

  - Potential for combining batch processing and streaming processing in the same system

# Example – Get hashtags from Twitter

```
val tweets = ssc.twitterStream()
```

**DStream:** a sequence of RDDs representing a stream of data

Twitter Streaming API    batch @ t    batch @ t+1    batch @ t+2 →

tweets DStream

stored in memory as an RDD
(immutable, distributed)

# Dstream operation

- Any operation applied on a DStream translates to operations on the underlying RDDs.

- For example, in the example of converting a stream of lines to words, ***the flatMap operation*** is applied on each RDD in the lines DStream to generate the RDDs of the words Dstream.



**Figure: Extracting words from an Inputstream**

# Example – Get hashtags from Twitter

```
val tweets = ssc.twitterStream()

val hashTags = tweets.flatMap (status => getTags(status))
```

**new DStream**

**transformation:** modify data in one DStream to create another DStream

batch @ t     batch @ t+1     batch @ t+2

tweets DStream

flatMap    flatMap    flatMap

hashTags Dstream
[#cat, #dog, ... ]

new RDDs created for every batch

# Input DStreams

- Input DStreams are DStreams representing the stream of input data received from streaming sources.



Figure: Input datastream divided into discrete chunks of data

- In the quick example, ***lines*** was an input DStream as it represented the stream of data received from the netcat server.

# Input DStreams



- If we want to receive multiple streams of data in parallel in our streaming application, then we can create multiple input Dstreams.
- This will create multiple receivers which will simultaneously receive multiple data streams.

# Receiver

- Every input DStream is associated with a Receiver object which receives the data from a source and stores it in Spark's memory for processing.



**Figure:** *The Receiver sends data onto the DStream where each Batch contains RDDs*

# Transformations on Dstreams

- Transformations allow the data from the inputDstream to be modified similar to RDDs. Dstreams support many of the transfromations available on normal Spark RDDs.

Most Popular Spark Streaming Transformations

- map
- flatMap
- filter
- reduce
- groupBy



DStream 1    Transform    DStream 2

Drop split point

**Figure:** *DStream Transformations*

# Transformations on Dstreams

- Map(func):

  - It returns a new Dstream by passing each element of the source Dstream through a function func.



Figure: Input DStream being converted through map(func)

Figure: Map Function

# Transformations on Dstreams

- flatMap(func):

  - It is similar to map(func), but each input item can be mapped to 0 or more output items and returns a new Dstream by passing each source element through a function func.



Figure: Input DStream being converted through flatMap(func)



Figure: flatMap Function

# Transformations on Dstreams

- <u>Filter(func):</u>

  - It returns a new <span style="color:#00B0F0">Dstream</span> by <span style="color:#00B0F0">selecting</span> only the records of the source Dstream on which <span style="color:#00B0F0">func</span> returns true.



Figure: Input DStream being converted through filter(func)

Figure: Filter Function For Even Numbers

# Transformations on Dstreams

- Reduce(func):

    - It returns a new Dstream of single-element RDDs by aggregating the elements in each RDD of the source Dstream using a function func.



Figure: Input DStream being converted through reduce(func)



Figure: Reduce Function To Get Cumulative Sum

# Transformations on Dstreams

- groupBy(func):
    - It returns the new RDD which basically is made up with a key and corresponding list of items of that group.



Figure: Input DStream being converted through groupBy(func)



Figure: Grouping By First Letters

# Dstream Window Operations

- Spark Streaming also provides windowed computations, which allow you to apply transformations over a sliding window of data.



**Figure:** *DStream Window Transformation*

# Window-based Transformations

```
val tweets = ssc.twitterStream()
val hashTags = tweets.flatMap (status => getTags(status))
val tagCounts = hashTags.window(Minutes(1), Seconds(5)).countByValue()
```

- sliding window operation
- window length
- sliding interval

window length

DStream of data

sliding interval

# Output Operations on DStreams

- It allows DStream's data to be pushed out to external systems like databases or file systems.

- Output operations trigger the actual execution of all the DStream transformations.



**Figure:** *Output Operations on DStreams*

# Output Operations on DStreams

- Currently, the following output operations are defined:

| | | Output Operations | | |
|---|---|---|---|---|
| print() | saveAsTextFiles (prefix, [suffix]) | saveAsObjectFiles (prefix, [suffix]) | saveAsHadoopFiles (prefix, [suffix]) | foreachRDD(func) |

# Example – Get hashtags from Twitter

```
val tweets = ssc.twitterStream()
val hashTags = tweets.flatMap (status => getTags(status))
hashTags.saveAsHadoopFiles("hdfs://...")
```

**output operation:** to push data to external storage



tweets DStream

batch @ t    batch @ t+1    batch @ t+2

flatMap    flatMap    flatMap

hashTags DStream

save    save    save

every batch saved to HDFS

INFO319, autumn 2018, session 5         source: http://ampcamp.berkeley.edu/

# Design Patterns for using foreachRDD

- dstream.foreachRDD is a powerful primitive that allows data to be sent out to external systems.

- The lazy evaluation achieves the most efficient transfer of data.

```
dstream.foreachRDD { rdd =>
 rdd.foreachPartition { partitionOfRecords =>

// ConnectionPool is a static, lazily initialized pool of connections
val connection = ConnectionPool.getConnection()
 partitionOfRecords.foreach(record => connection.send(record))

// Return to the pool for future reuse
ConnectionPool.returnConnection(connection)
 }
}
```

# Caching and persistent

- Dstreams allow developers to cache/persist the stream's data in memory. This is useful if the data in the Dstream will be computed multiple times.

- This can be done using the persist() method on a Dstream.

- For input streams that receive data over the network (such as Kafka, Flume, sockets, etc), the default persistence level is set to replicate the data to two nodes for fault-tolerance.



**Figure:** *Caching Into 2 Nodes*

# Accumulators/Broadcast/Variables/Checkpoints

- Accumulators are variables that are only added through an associative and commutative operation.

- They are used to implement counters or sums.



**Figure:** *Accumulators In Spark Streaming*

- Tracking accumulators in the UI can be useful for understanding the progress of running stages.

- Spark natively supports numeric accumulators. We can create named or unnamed accumulators.

# Accumulators/Broadcast/Variables/Checkpoints

- Broadcast variables allow the programmer to keep a read-only variable cached on each machine rather than shipping a copy of it with tasks.

- They can be used to give every node a copy of a large input dataset in an efficient manner.

- Spark also attempts to distribute broadcast variables using efficient broadcast algorithms to reduce communication post.

Figure: Broadcasting A Value To Executors

Figure: SparkContext and Broadcasting

# Accumulators/Broadcast/Variables/Checkpoints

- Checkpoints are similar to checkpoints in gaming. They make it run 24/7 and make it resilient to failures unrelated to the application logic.



It is the saving of the information defining the streaming computation

Metadata Checkpoints

Checkpoints

Data Checkpoints

It is saving of the generated RDDs to reliable storage

# Social network graphs

# Social network graphs

- Social Networks as graphs

    - Social networks are naturally modeled as graphs, which we sometimes refer to as a social graph.

- What is a social network??

    - When we think of a social network, we think of Facebook, Twitter, Google+, or another website that is called a "social network," and indeed this kind of network is representative of the broader class of networks called "social."

# Social network graphs

- In social networks, the entities are the nodes, and an edge connects two nodes if the nodes are related by the relationship that characterizes the network.

- If there is a degree associated with the relationship, this degree is represented by labeling the edges.

- Often, social graphs are undirected, as for the Facebook friends' graph. But they can be directed graphs, as for example the graphs of followers on Twitter or Google+.

# An example social network



Figure 10.1: Example of a small social network

# SNA example 1:

# SNA Example 2:



INFO319, autumn 2018, session 5

# Example 3



Figure 2. Information communication and tracking Network

# Social network analysis

- Social network analysis (SNA) is "the process of investigating social structures through the use of networks and graph theory".

- Under the SNA model, social networks are represented by "nodes" and "edges," or the elements (e.g. individuals, organizations, ideas) of a given social network and the various relationships that connect them.

- In the case of social networks on Twitter, nodes may represent different users while edges may represent any of the interactions that connect these individual accounts (e.g. likes, replies, retweets, mentions etc.).

# Twitter data on SNG



Who is retweeted by whom?

#allertameteoTOS – is.retweeted by

#allertameteoPIE – is.retweeted by

#allertameteoLIG – is.retweeted by

# Use Case- Twitter analysis

# Apache Spark

```scala
import org.apache.spark.SparkConf
import org.apache.spark.streaming.StreamingContext
import org.apache.spark.streaming.Seconds
import twitter4j.conf.ConfigurationBuilder
import twitter4j.auth.OAuthAuthorization
import twitter4j.Status
import org.apache.spark.streaming.twitter.TwitterUtils
object TwitterData {
  def main(args: Array[String]) {
    if (args.length < 4) {
      System.err.println("Usage: TwitterData <ConsumerKey><ConsumerSecret><accessToken><access
        "[<filters>]")
      System.exit(1)
    }
    val appName = "TwitterData"
    val conf = new SparkConf()
    conf.setAppName(appName).setMaster("local[3]")
    val ssc = new StreamingContext(conf, Seconds(5))
    val Array(consumerKey, consumerSecret, accessToken, accessTokenSecret) = args.take(4)
    val filters = args.takeRight(args.length - 4)
    val cb = new ConfigurationBuilder
    cb.setDebugEnabled(true).setOAuthConsumerKey(consumerKey)
      .setOAuthConsumerSecret(consumerSecret)
      .setOAuthAccessToken(accessToken)
      .setOAuthAccessTokenSecret(accessTokenSecret)
    val auth = new OAuthAuthorization(cb.build)
    val tweets = TwitterUtils.createStream(ssc, Some(auth))
    val englishTweets = tweets.filter(_.getLang() == "en")
    englishTweets .saveAsTextFiles("tweets", "json")
    ssc.start()
    ssc.awaitTermination()
  }
}
```

# Setting up the Spark streaming context

- We need to set the Spark streaming context as follows:

```
val ssc = new StreamingContext(conf, Seconds(5))
```

# Authenticate twitter user

- val cb = new ConfigurationBuilder
-     cb.setDebugEnabled(true).setOAuthConsumerKey(consumerKey)
-         .setOAuthConsumerSecret(consumerSecret)
-         .setOAuthAccessToken(accessToken)
-         .setOAuthAccessTokenSecret(accessTokenSecret)

<br>

- Authentication:
        val auth = new OauthAuthorization(cb.build)

# Starting the spark streaming

- val tweets = TwitterUtils.createStream(ssc, Some(auth))

- val englishTweets = tweets.filter(_.getLang() == "en")

- Run the class file and then console window is appeared.

# Spark streaming

# Conclusions

- Social media

- Social media use in disasters

- What streaming is?

- Why Spark streaming is important?

- Different Spark streaming components

- Example: Real-time twitter data stream processing with Apache Spark