# CSCE 5300

# Spark Streaming

# What is Streaming?

- Data Streaming is a technique for transferring data so that it can be processed as a steady and continuous stream

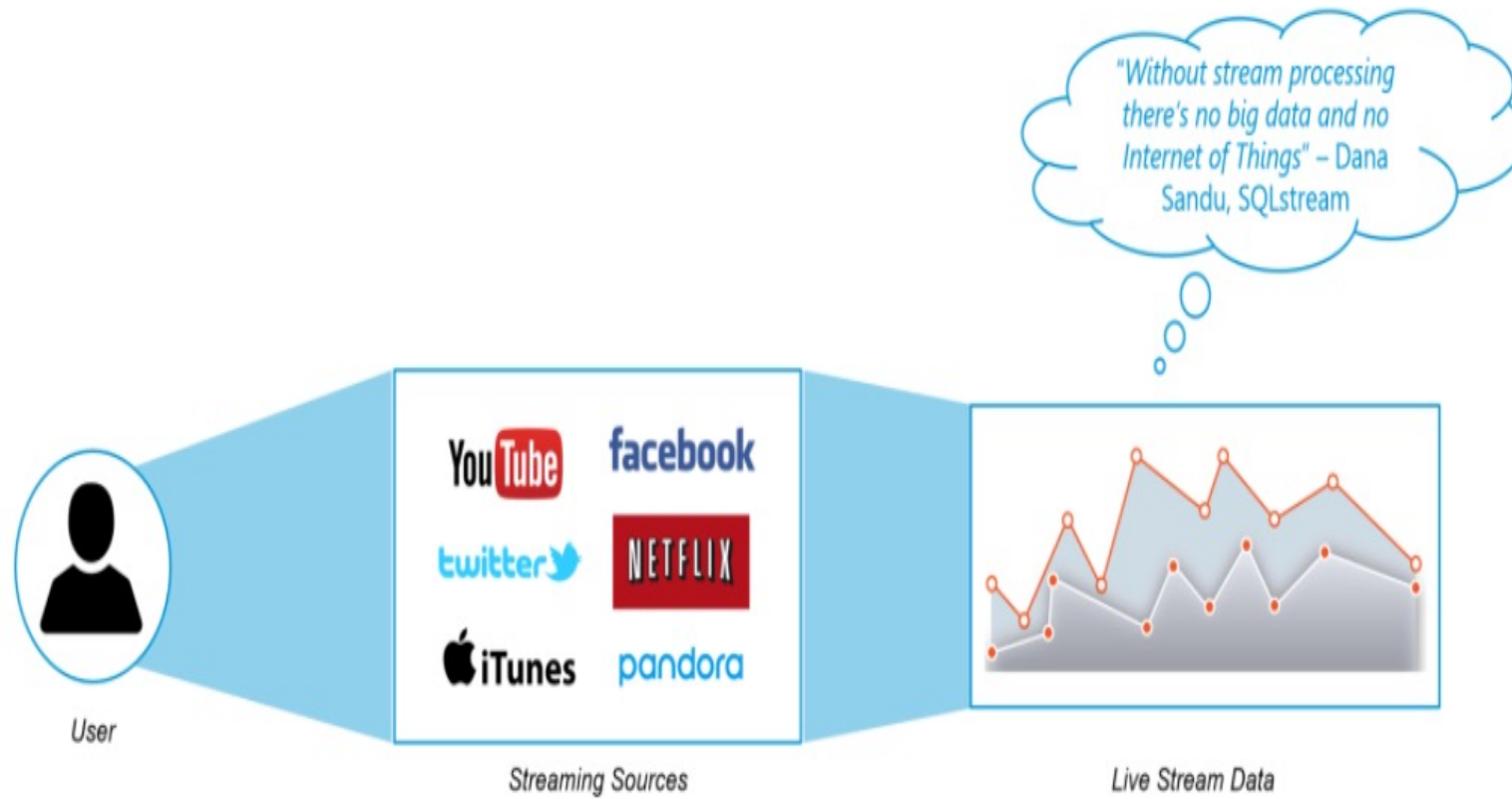- Streaming technologies are becoming increasingly important with the growth of the Internet

**Figure:** *What is Streaming?*

# Why Spark Streaming?



Spark Streaming is used to stream real-time data from various sources like Twitter, Stock Market and Geographical Systems and perform powerful analytics to help businesses.
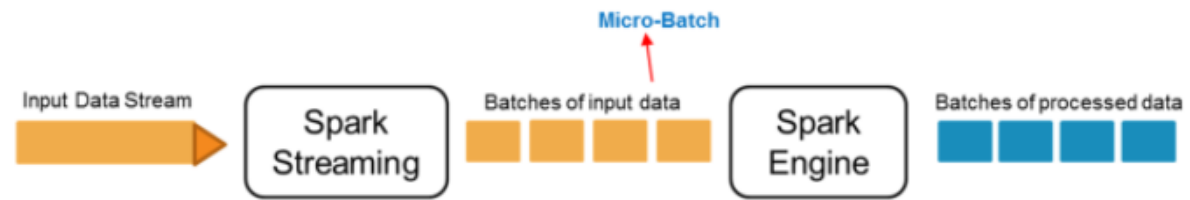
**Figure:** *Why Spark Streaming?*

**Figure:** *Streams in Spark Streaming*

# Spark Streaming Features

- **Scaling:** Spark Streaming can easily scale to hundreds of nodes.
- **Speed:** It achieves low latency.
- **Fault Tolerance:** Spark has the ability to efficiently recover from failures.
- **Integration:** Spark integrates with batch and real-time processing.
- **Business Analysis:** Spark Streaming is used to track the behavior of customers which can be used in business analysis
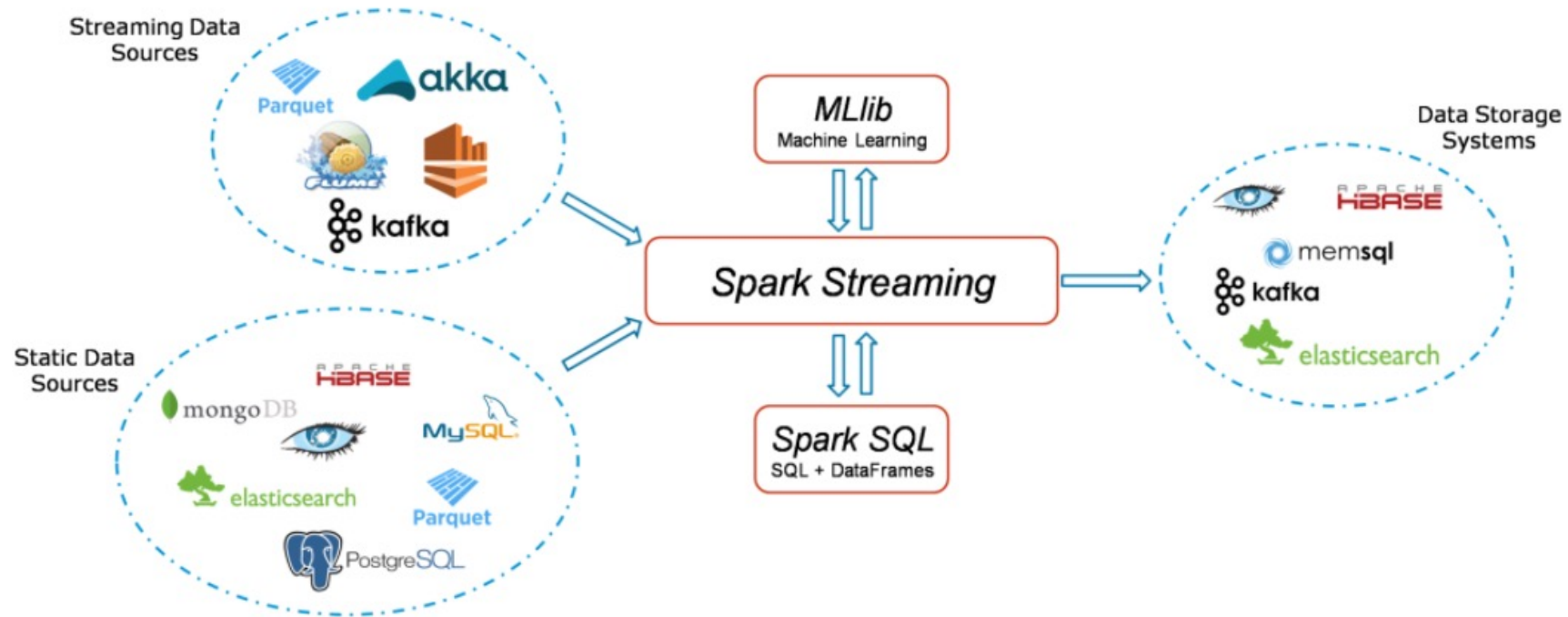
# Spark Streaming Workflow



**Figure:** *Overview Of Spark Streaming*

# Spark Streaming Fundamentals

- Streaming Context

- DStream

- Caching

- Accumulators, Broadcast Variables and Checkpoints

# Streaming Context

- *Streaming Context* consumes a stream of data in Spark.

- It registers an *Input DStream* to produce a *Receiver* 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:** *Spark Streaming Context*



**Figure:** *Default Implementation Sources*

# DStream

- *Discretized Stream* (DStream) is the basic abstraction provided by Spark Streaming.

- It is a continuous stream of data.

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

**Figure:** *Extracting words from an Input DStream*

# Transformations on DStreams

| | |
|---|---|
| map(*func*) | map(*func*) returns a new DStream by passing each element of the source DStream through a function *func*. |
| flatMap(*func*) | flatMap(*func*) 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*. |
| filter(*func*) | filter(*func*) returns a new DStream by selecting only the records of the source DStream on which *func* returns true. |
| reduce(*func*) | reduce(*func*) returns a new DStream of single-element RDDs by aggregating the elements in each RDD of the source DStream using a function *func*. |
| groupBy(*func*) | groupBy(*func*) returns the new RDD which basically is made up with a key and corresponding list of items of that group. |

# Accumulators, Broadcast Variables and Checkpoints

- **Accumulators:**

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

➢ They are used to implement counters or sums. 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.

- **Broadcast Variables:**

➤ *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 cost.

- **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

Data Checkpoints

It is saving of the generated RDDs to reliable storage

Checkpoints

**Figure:** *Features of Checkpoints*

# Spark Streaming

# File Generation

# Download Netcat

- Go https://joncraton.org/files/nc111nt.zip

- Unzip the files and use nc as password

- Go to environment variables

- Select Path and Add new path as you location to netcat folder.

- Open command line

- Type nc to check either netcat is working or not

# Running Results

# References

- [https://spark.apache.org/docs/2.2.0/streaming-programming-guide.html](https://spark.apache.org/docs/2.2.0/streaming-programming-guide.html)
- [https://www.edureka.co/blog/spark-streaming/](https://www.edureka.co/blog/spark-streaming/)