# Lambda Architecture for Twitter real-time sentiment analysis

Lorenzo Agnolucci

Università degli Studi di Firenze
Dipartimento di Ingegneria dell'Informazione
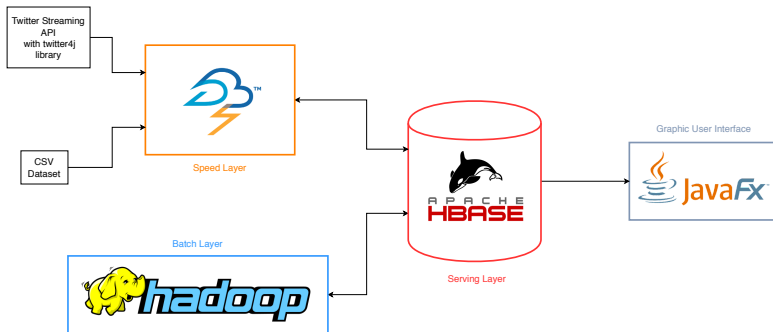
Firenze, 20 Aprile 2019

# Outline

UNIVERSITÀ
DEGLI STUDI
FIRENZE

**DINFO**
Dipartimento di
Ingegneria dell'Informazione

- Big Data requires to find ways to analyze a large amount of data

- Lambda Architecture is a particular approach composed by:

  - *batch layer* : applies batch-oriented technologies (like MapReduce) on a master database. It is effective but it has a high latency

  - *serving layer* : specialized distributed database that supports batch updates and random reads

  - *speed layer* : only looks at recent data and uses low-latency techniques to update real-time views. It compensate for the high latency of the batch layer

- Sentiment analysis is a type of data mining applied to Big Data with some useful applications

- The main goal was not a perfect sentiment classification but the implementation of the architecture
- Speed layer is started first with the keywords as arguments. It creates the speed layer tables at the start of the execution
- Tweets of a dataset are added as if they belonged to the real-time stream to increase the number of tweets
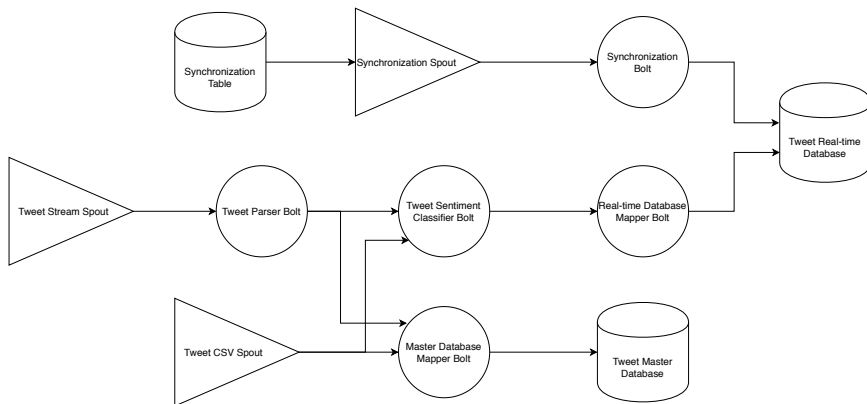
UNIVERSITA
DEGLI STUDI
FIRENZE

**DINFO**
Dipartimento di
Ingegneria dell'Informazione

- Developed with *LingPipe* library

- Trained on 1.6 millions tweets [1]

- Classifies English text with 2 categories: positive and negative

- Decent 0.71 accuracy

---

[1] A. Go, R. Bhayani, and L. Huang. Twitter sentiment classification using distant supervision. CS224N project report, Stanford, 1(12):2009, 2009. https://www.kaggle.com/kazanova/sentiment140

UNIVERSITA
DEGLI STUDI
FIRENZE

**DINFO**
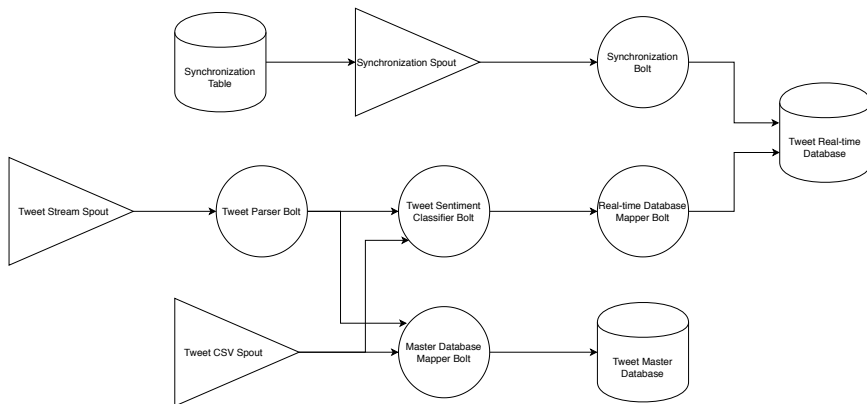Dipartimento di
Ingegneria dell'Informazione

Based on Apache HBase and composed by 4 tables:

- *tweet master database*: master database of the Lambda Architecture

- *tweet real-time database*: stores the tweets on which the real-time view is based

- *batch view*: result of the batch processing

- *synchronization table*: contains the start and the end timestamps of the batch processing

Batch layer

UNIVERSITA
DEGLI STUDI
FIRENZE

DINFO
Dipartimento di
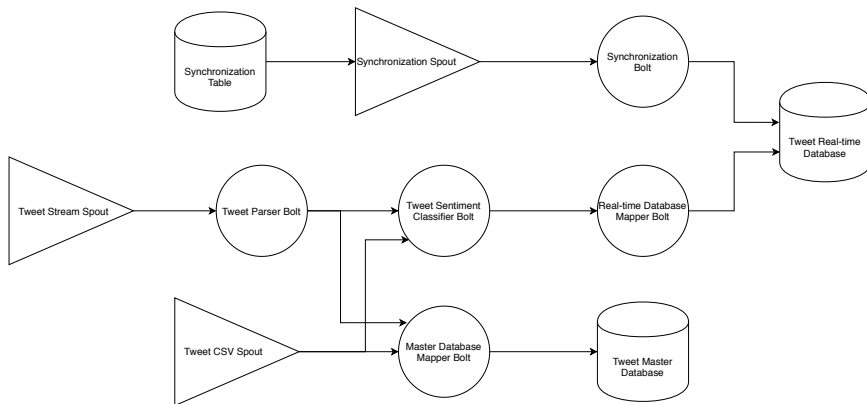Ingegneria dell'Informazione

- Represented by Apache Hadoop

- Computes a MapReduce job on *tweet master database* in a infinite loop

- Writes its results from scratch in *batch view*

- Writes the start and the end timestamps of the computation in *synchronization table*

- Mapper takes a tweet in input and outputs a $< Keyword, Sentiment >$ tuple

- Reducer takes a tuple in input and increment the corresponding cell in *batch view*
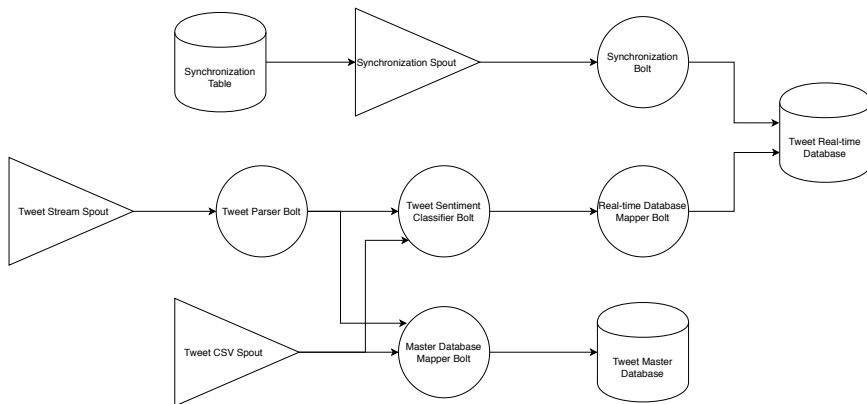
- *tweet stream spout*: gets a real-time stream of tweets with *Twitter4j* library and filters them

- *tweet parser bolt*: parse a tweet object to a tuple

UNIVERSITA
DEGLI STUDI
FIRENZE

**DINFO**
Dipartimento di
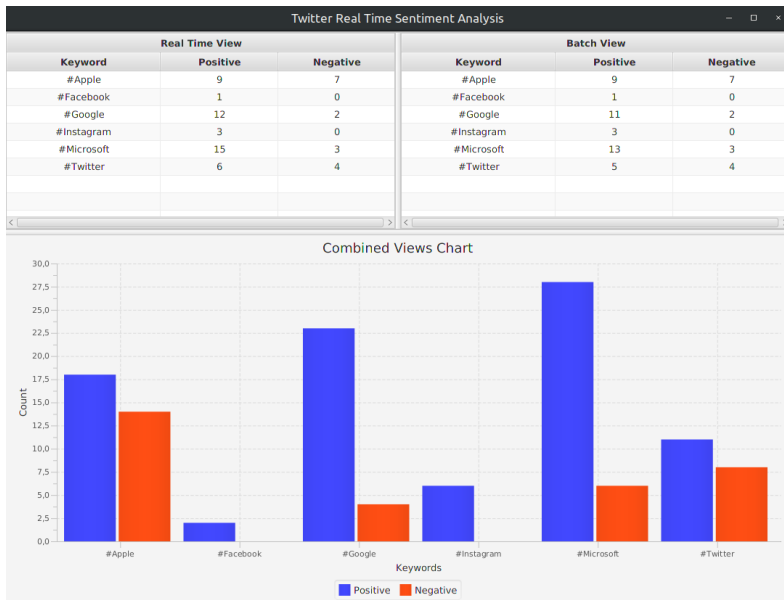Ingegneria dell'Informazione



- *tweet CSV spout*: outputs a tuple for each tweet of the dataset
- *master database mapper bolt*: inserts tweets in *tweet master database*

- *tweet sentiment classifier bolt*: classifies the sentiment of the tweet
- *real-time database mapper bolt*: inserts tuples in *tweet real-time database*

- *synchronization spout*: checks when batch processing ends
- *synchronization bolt*: deletes already processed tweets from *tweet real-time database*

# GUI

# Conclusions

- It has been shown an implementation of a Lambda Architecture capable of getting sentiment analysis statistics of real-time tweets

- The GUI that was developed lets to visualize how the different parts of the architecture work together

- As a future development a neutral category could be added to the sentiment classifier