# SENTIMENT ANALYSIS IN MOVIE COMMENTS
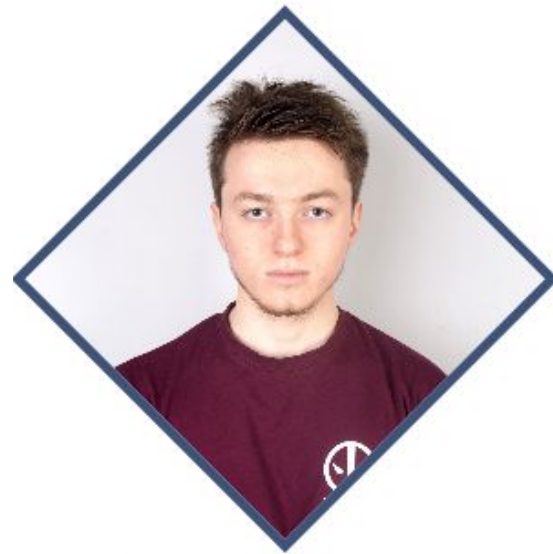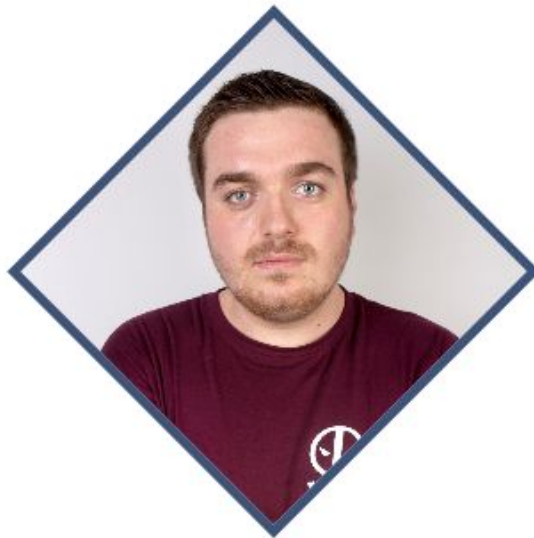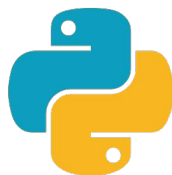
# HELLO!

## Denis Castéran
## Hugo Rybinski

SCIA 2018

# GOAL

» Fetch movie datas (especially reviews)
» Analyze the datas to determine if it is a positive or a negative review
» Extract informations from the datas and display it

# SOFTWARES & TECHNOLOGIES USED

**Languages:**

» Scala
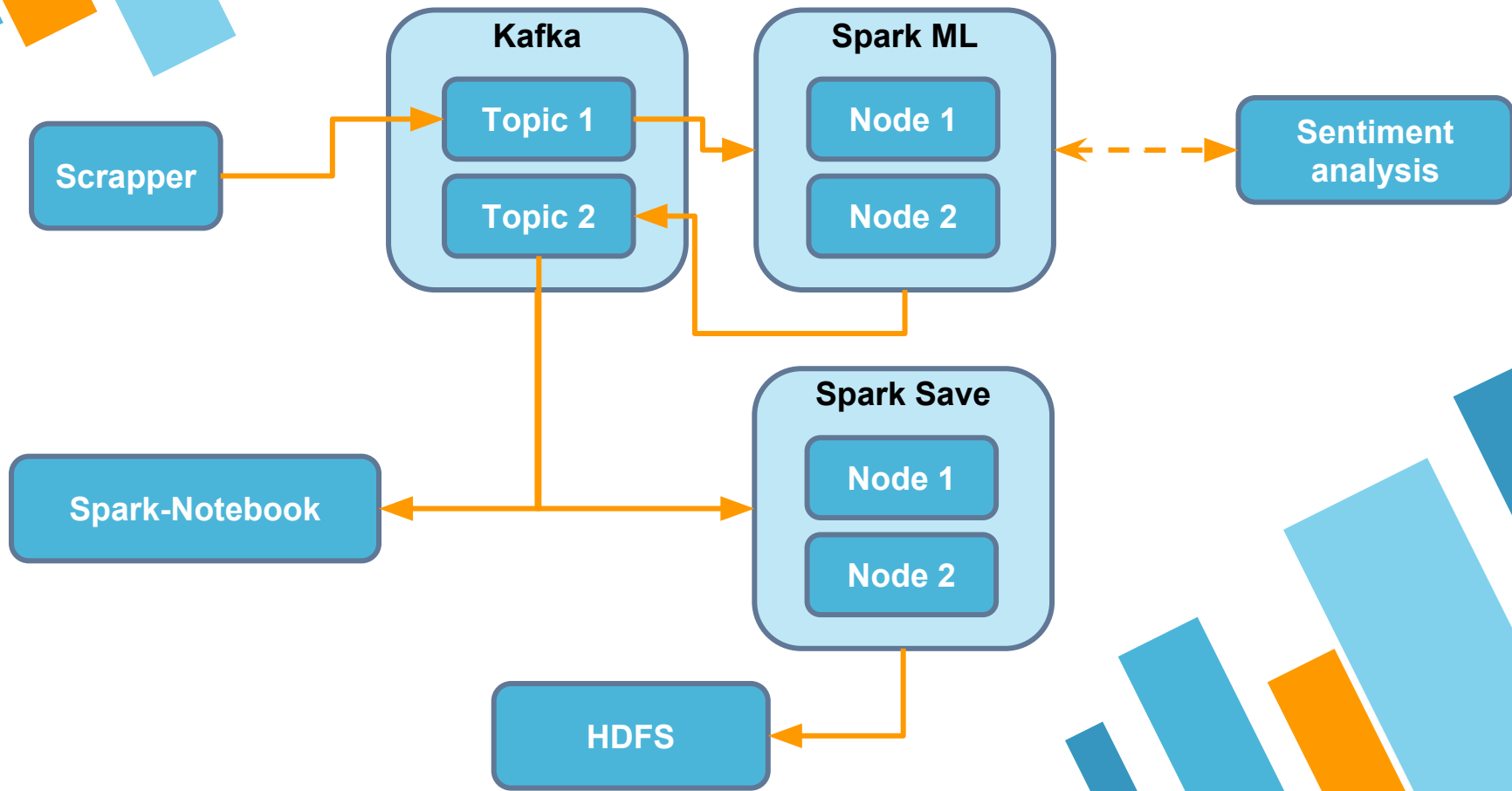
» Python

**From the Hadoop stack:**

» Kafka

» Spark
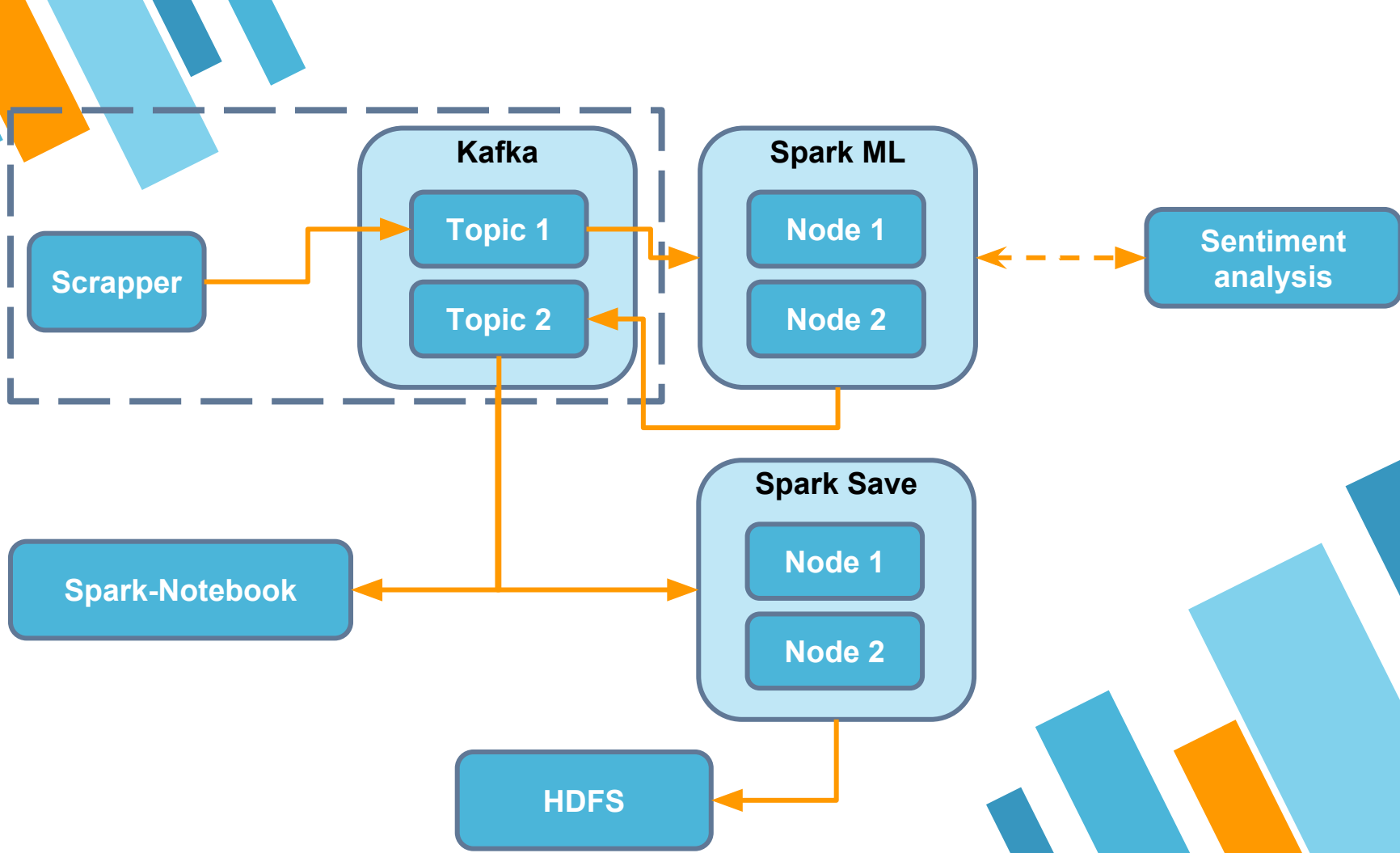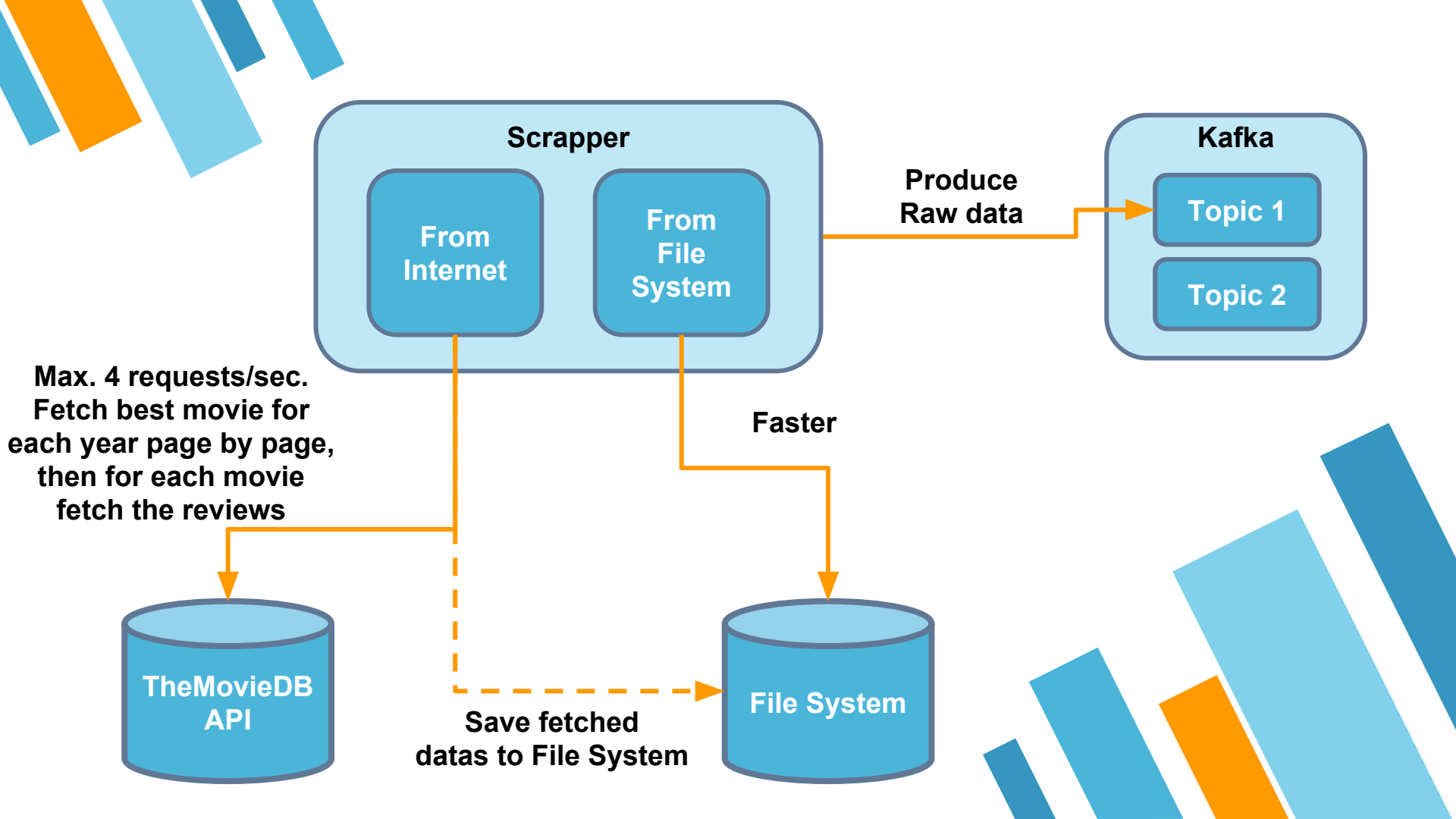
**Misc:**

» Spark-notebook

» Docker

# 1.

## DATABASE SCRAPING & KAFKA

Get the datas

# KAFKA SETUP

- Simple setup with two Docker images:
    - Zookeper
    - Kafka
- Possibility to add broker on other machines:
    - Just start the Kafka image somewhere else, but change the address and port of the Zookeeper server
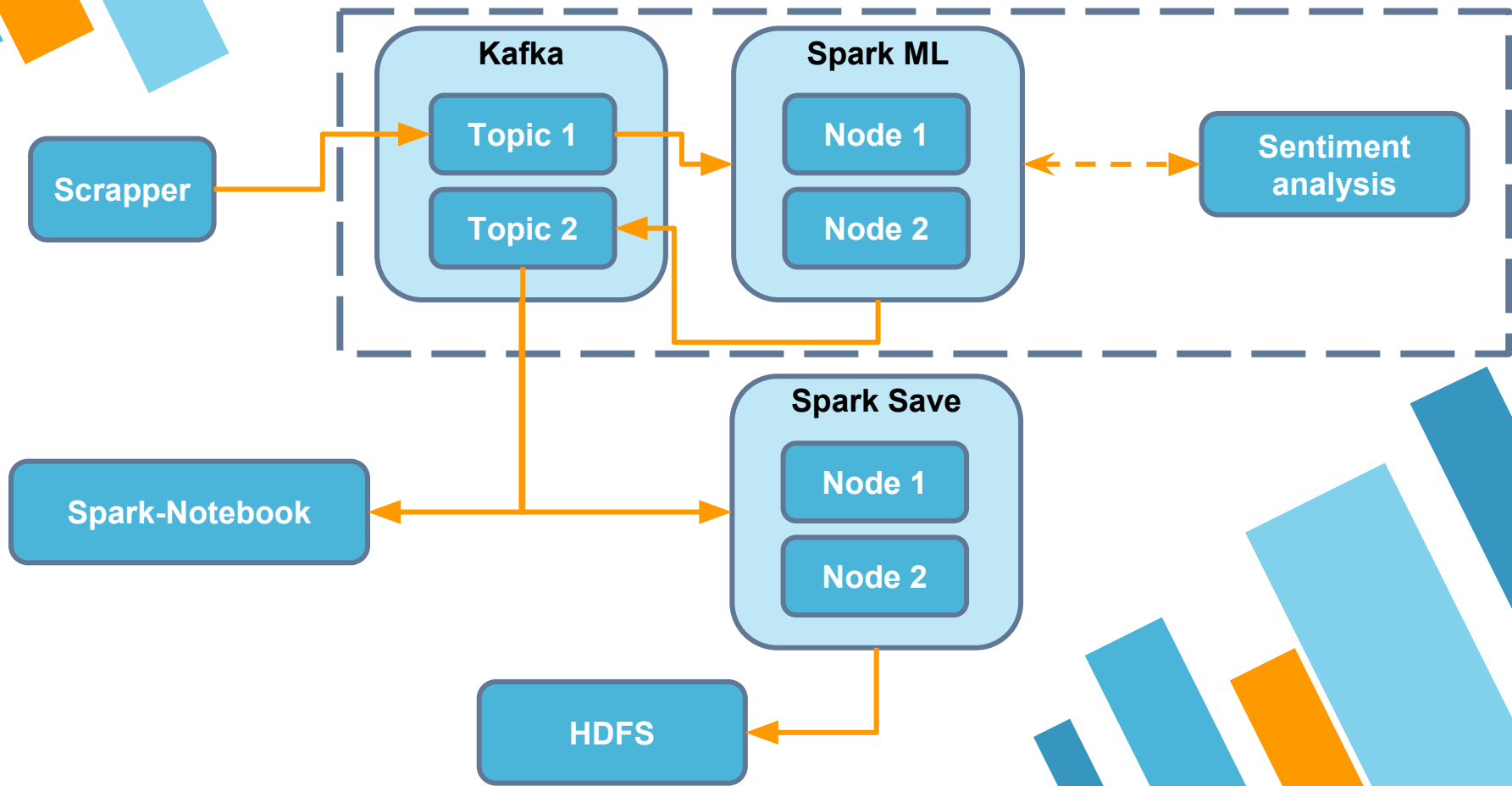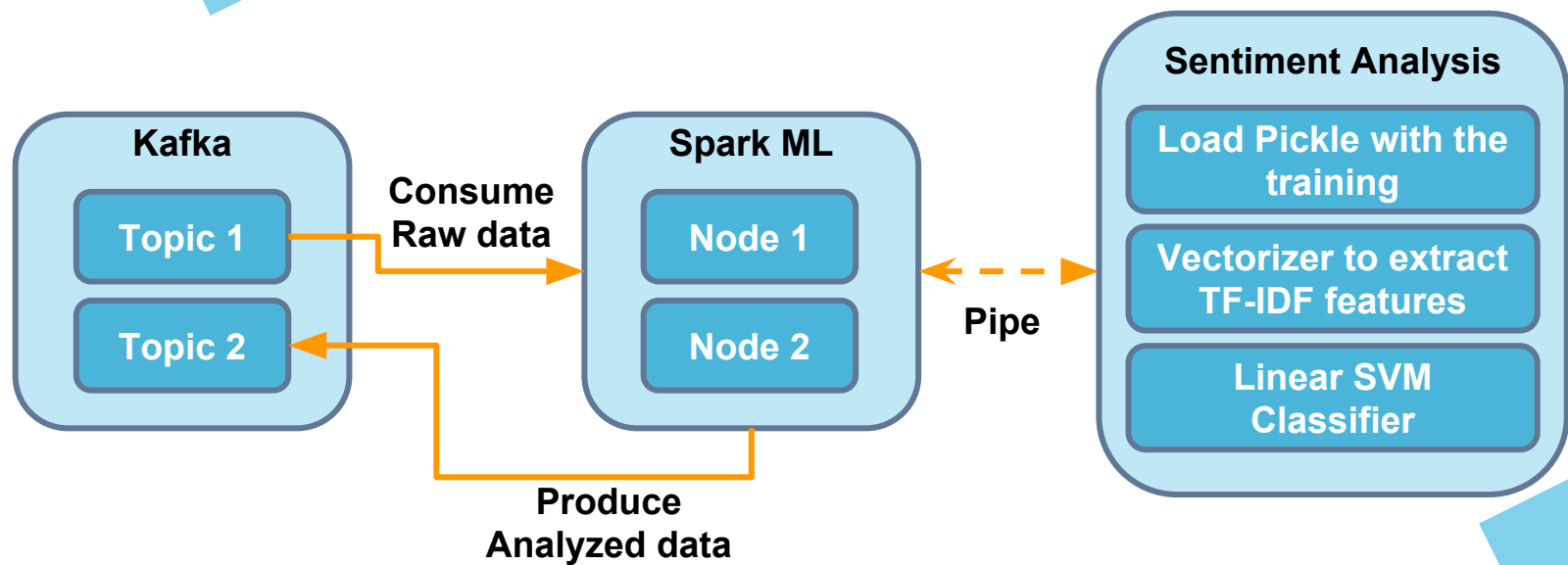- Easy to manage with kafka-manager

# 2.

# SPARK TREATMENT

Parallelize machine learning

# SPARK

- Spark streaming :
    - Allow parallelized computation
    - Quick computation of sentiment analysis
    - Easy to communicate with kafka

# SENTIMENT ANALYSIS

- Built using Scikit-Learn
- Train on local machine once, trained classifier is saved in a Pickle
  - Training dataset : 10 000 labelled reviews from IMDB
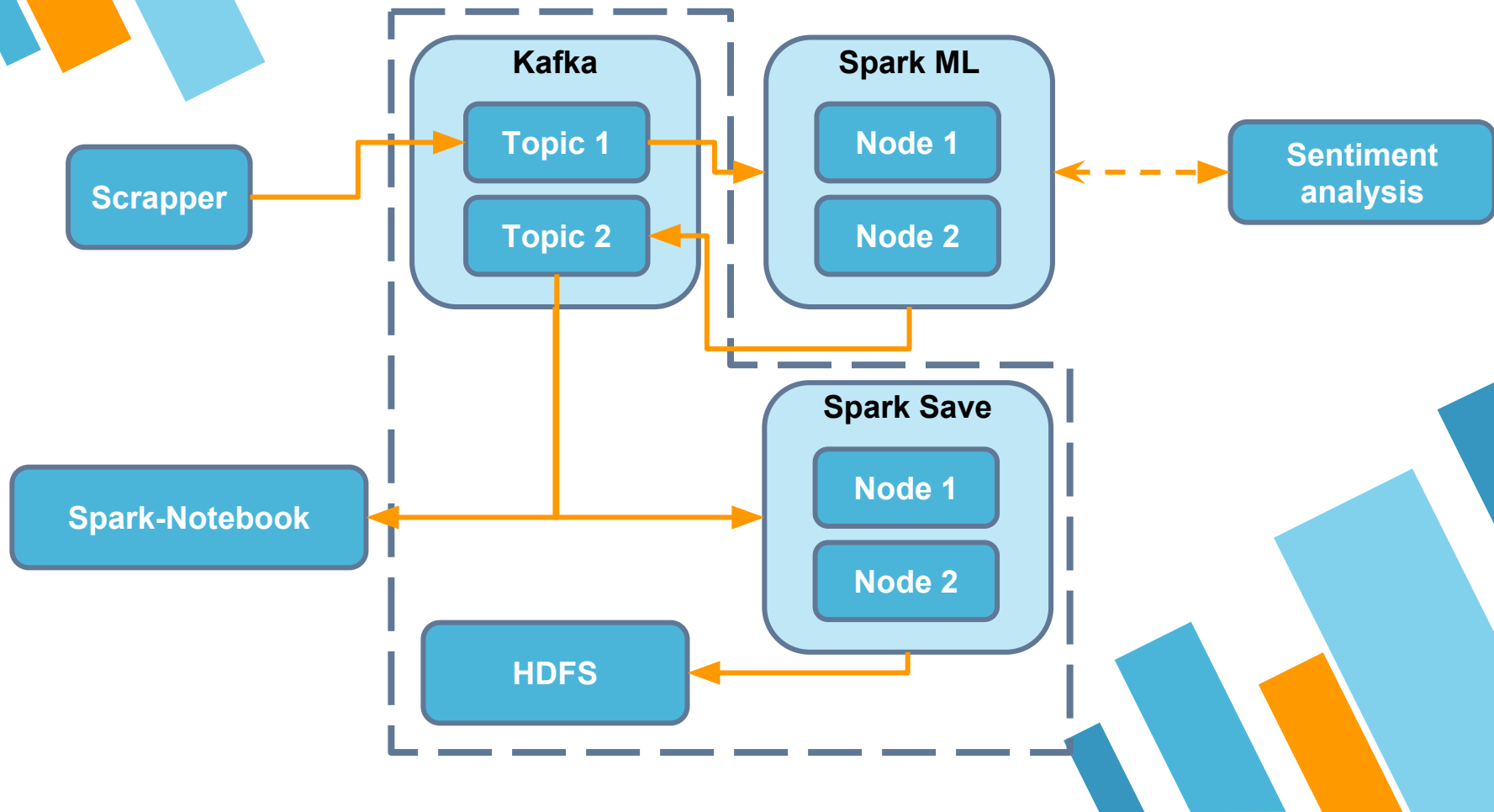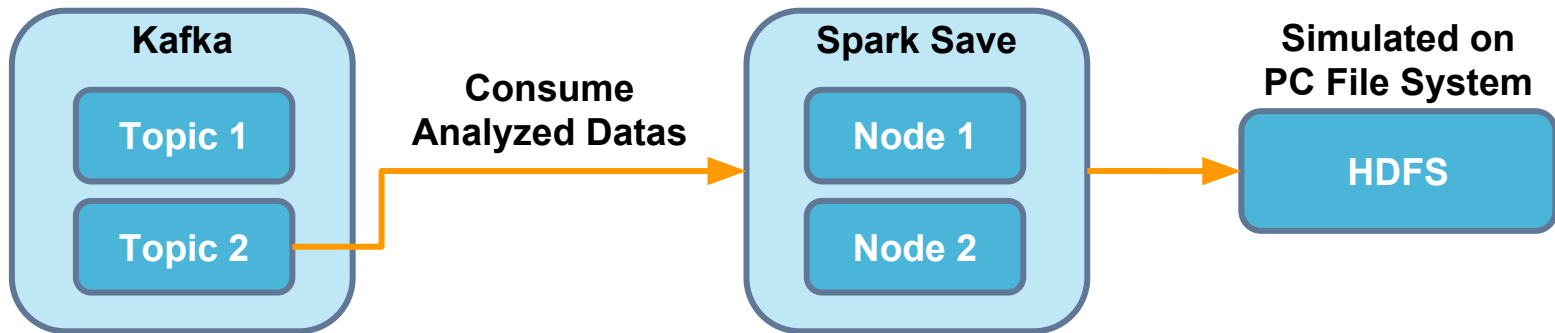  - Testing dataset : 10 000 labelled reviews from IMDB
  - ~80% accuracy

# 3.

# DATA PERSISTENCE

Save the computed datas

# HDFS

- Simulated on file system :
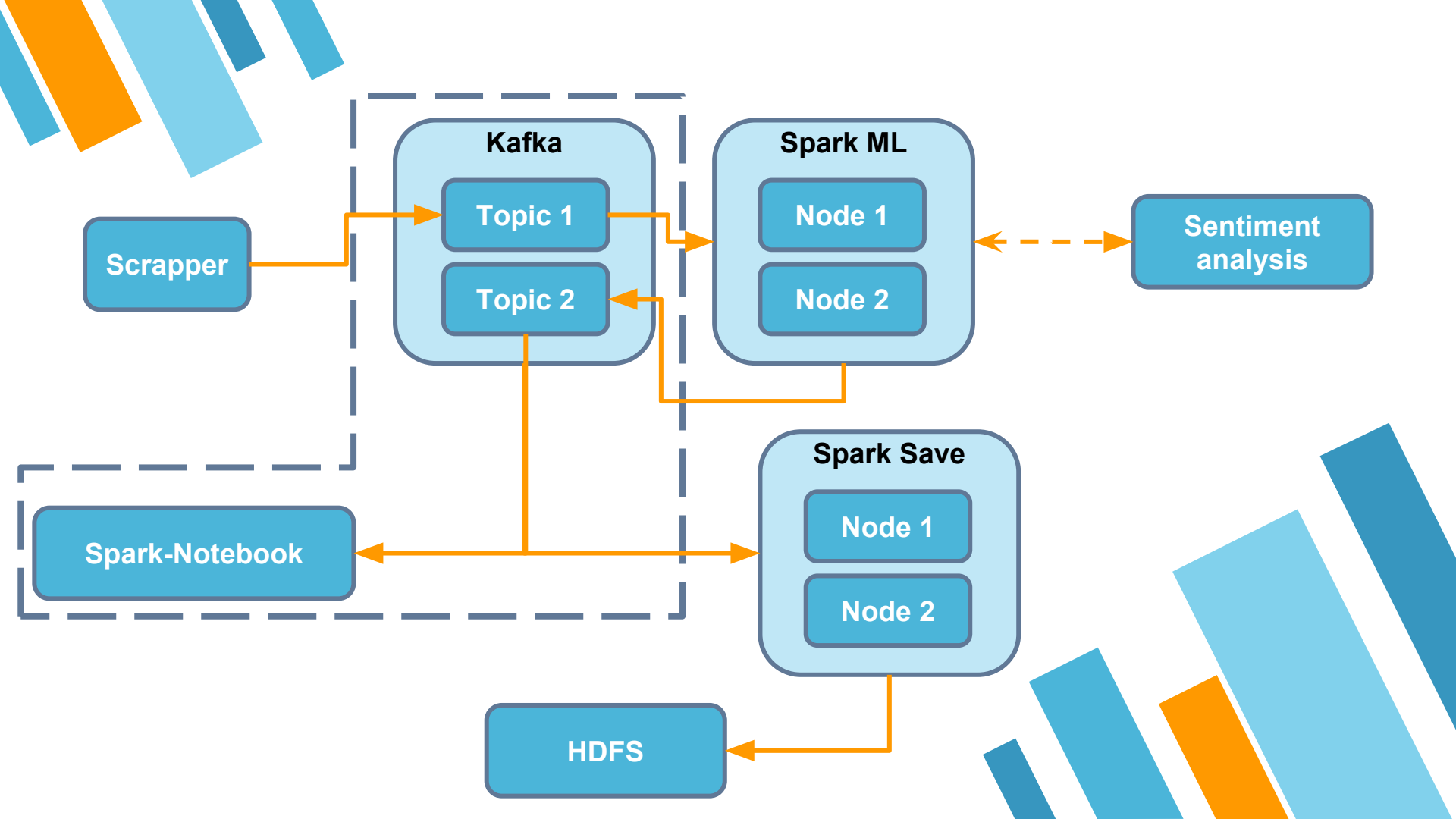    - The dataframe that contains the RDD is saved
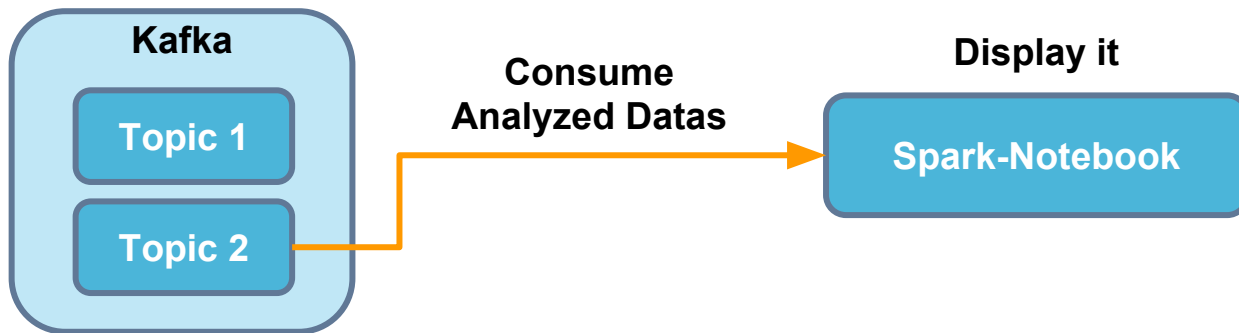
# 4.

## DATA DISPLAY

With spark-notebook

# SPARK-NOTEBOOK

- Easy to the data from kafka with spark streaming
- Prints cool graphics !

# THANKS!

## Any questions?