

Coronavirus Playback



Replay the development of Coronavirus with real time simulation

Professor: Mrudula Mukadam

Team:

Le Chi Nhan, Ngo

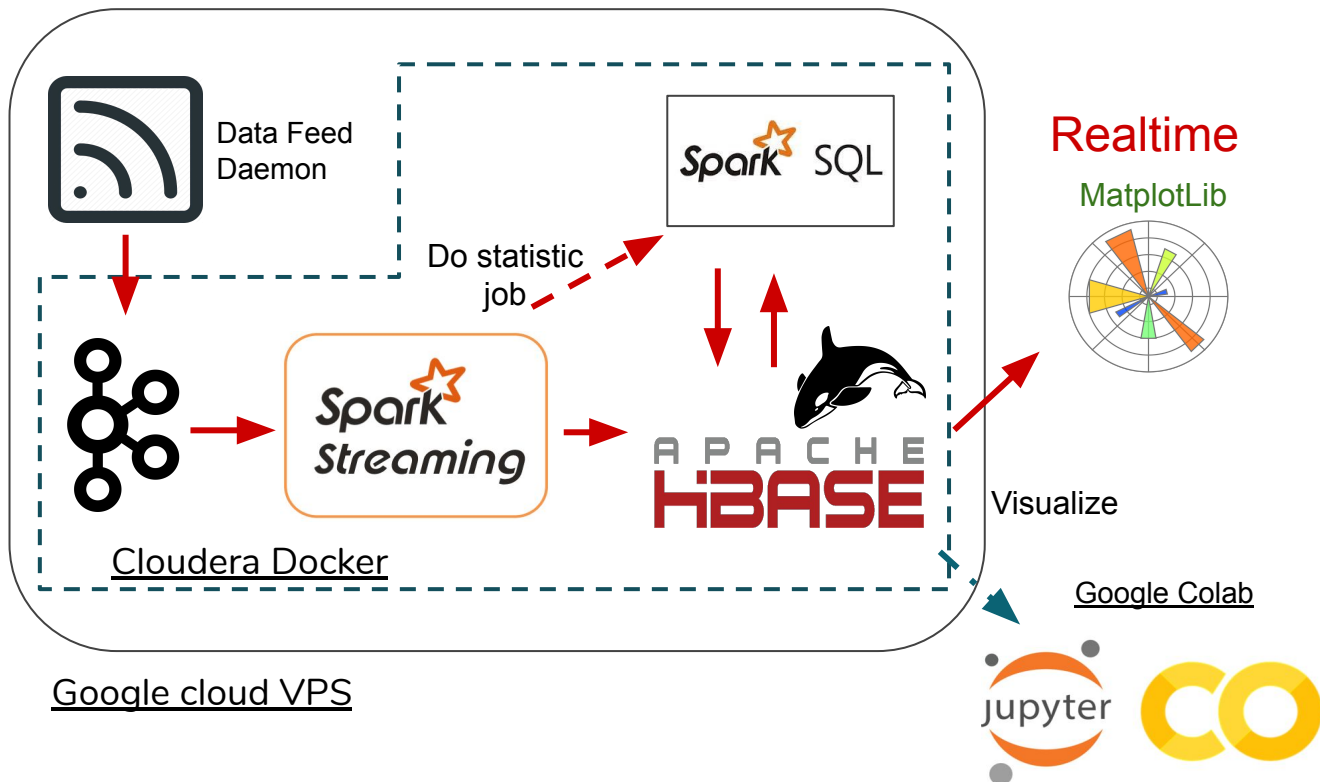
Tien Khanh, Nguyen



Agenda

- Overview
- Data Feed use `python` and `kafka` lib
- Kafka
- Spark Streaming
- HBase
- Spark SQL
- Visualize -
 - Colab - Google
 - Realtime chart with `Matplotlib` and `python`

Overview



Dataset

1. <https://github.com/CSSEGISandData/COVID-19>

	Province/State	Country/Region	Last Update	Confirmed	Deaths	Recovered	Latitude	Longitude
1								
2	Hubei	China	2020-03-15T18:20:18	67794	3085	54288	30.9756	112.2707
3		Italy	2020-03-14T20:13:16	24747	1809	2335	41.8719	12.5674
4		Iran	2020-03-15T18:20:18	13938	724	4590	32.4279	53.688
5		Korea, South	2020-03-15T18:20:18	8162	75	510	35.9078	127.7669
6		Spain	2020-03-15T18:20:18	7798	289	517	40.4637	-3.7492
7		Germany	2020-03-15T18:20:18	5795	11	46	51.1657	10.4515
8	France	France	2020-03-15T18:20:18	4499	91	12	46.2276	2.2137
9		Switzerland	2020-03-15T18:20:18	2200	14	4	46.8182	8.2275

Data specs:

- 01/22 to current, about 85 daily CS
- Total ~75k records in ~80 csv files ,

```
huestarc@huestarc:~/data/csse_covid_19_daily_reports$ ls
01-22-2020.csv 02-03-2020.csv 02-15-2020.csv 02-27-2020.csv 03-10-2020.csv 03-22-2020.csv 04-03-2020.csv
01-23-2020.csv 02-04-2020.csv 02-16-2020.csv 02-28-2020.csv 03-11-2020.csv 03-23-2020.csv 04-04-2020.csv
01-24-2020.csv 02-05-2020.csv 02-17-2020.csv 02-29-2020.csv 03-12-2020.csv 03-24-2020.csv 04-05-2020.csv
01-25-2020.csv 02-06-2020.csv 02-18-2020.csv 03-01-2020.csv 03-13-2020.csv 03-25-2020.csv 04-06-2020.csv
01-26-2020.csv 02-07-2020.csv 02-19-2020.csv 03-02-2020.csv 03-14-2020.csv 03-26-2020.csv 04-07-2020.csv
01-27-2020.csv 02-08-2020.csv 02-20-2020.csv 03-03-2020.csv 03-15-2020.csv 03-27-2020.csv 04-08-2020.csv
01-28-2020.csv 02-09-2020.csv 02-21-2020.csv 03-04-2020.csv 03-16-2020.csv 03-28-2020.csv 04-09-2020.csv
01-29-2020.csv 02-10-2020.csv 02-22-2020.csv 03-05-2020.csv 03-17-2020.csv 03-29-2020.csv 04-10-2020.csv
01-30-2020.csv 02-11-2020.csv 02-23-2020.csv 03-06-2020.csv 03-18-2020.csv 03-30-2020.csv 04-11-2020.csv
01-31-2020.csv 02-12-2020.csv 02-24-2020.csv 03-07-2020.csv 03-19-2020.csv 03-31-2020.csv 04-12-2020.csv
02-01-2020.csv 02-13-2020.csv 02-25-2020.csv 03-08-2020.csv 03-20-2020.csv 04-01-2020.csv README.md
02-02-2020.csv 02-14-2020.csv 02-26-2020.csv 03-09-2020.csv 03-21-2020.csv 04-02-2020.csv
```

Challengings:

- Uncertain columns set / pattern -> Optional column picking
- Various time format -> Optional parsing
- Mistaken date -> Apply date of data file



Data Feed cron job

To simulate the realtime processing we get data file from oldest to newest.

A **bash** + **python** when run read **csv** and feed to **kafka**

Interval: Every 2 seconds

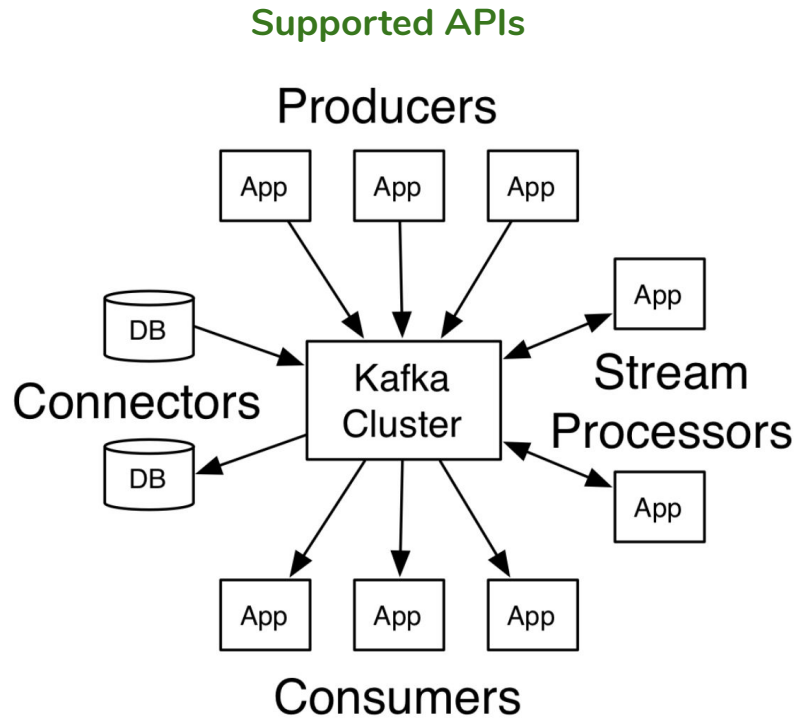
Input: dataset CSV

Periodically publish data into **Kafka**

Kafka - Why Kafka?



- A distributed streaming platform
- Publish and subscribe to streams of records
- Store streams of records in a fault-tolerant durable way.
- Process streams of records as they occur.



Kafka - Consumer implementation (java)

- Kafka Streaming
- Spark Streaming
- HBase

```
@Log4j
public class KConsumer {

    public static void startConsumer() throws InterruptedException {
        JavaSparkContext sc = SparkConfig.getSparkContext();
        Map<String, String> kafkaParams = KafkaConfig.generateKafkaParams();
        Set<String> topicName = Collections.singleton(KafkaConfig.TOPIC_NAME);
        Configuration hadoopConf = sc.hadoopConfiguration();
        HBaseRepository repo = HBaseRepository.getInstance();

        try (JavaStreamingContext streamingContext = new JavaStreamingContext(sc, new Duration(5000))) {
            JavaPairInputDStream<String, String> kafkaSparkPairInputDStream = KafkaUtils.createDirectStream(
                streamingContext, String.class,
                String.class, StringDecoder.class, StringDecoder.class, kafkaParams, topicName);

            JavaDStream<CoronaRecord> recoredRDDs = kafkaSparkPairInputDStream
                .map(RecordParser::parse).filter(r -> r != null);

            recoredRDDs.foreachRDD(rdd -> {
                if (!rdd.isEmpty()) {
                    repo.save(hadoopConf, rdd);

                    // Refresh the visualization
                    Log.info("===== Refreshing the visualization =====");
                    CoronaAnalysisApp.init();
                    CoronaAnalysisApp.generateTotalCasesPilot();
                }
            });

            streamingContext.start();
            streamingContext.awaitTermination();
        }
    }
}
```



HBase tables - corona_cases

ROW	COLUMN+CELL
Australia 20200123	column=cc:confirmedCases, timestamp=1587003705186, value=\x00\x00\x00\x00
Australia 20200123	column=cc:country, timestamp=1587003705186, value=Australia
Australia 20200123	column=cc:county, timestamp=1587003705186, value=
Australia 20200123	column=cc:date, timestamp=1587003705186, value=2020/01/23
Australia 20200123	column=cc:deathCases, timestamp=1587003705186, value=\x00\x00\x00\x00
Australia 20200123	column=cc:recoveredCases, timestamp=1587003705186, value=\x00\x00\x00\x00
Australia 20200123	column=cc:state, timestamp=1587003705186, value=

- Storage Corona cases data
- **Key:** country|state|date
- **Column Family:** cc
- **Column:** country, state, date, deathCases, recoveredCases, confirmedCases



HBase tables - pilot

- Storage Visualization data
- **Key:** yyyyMMdd
- **Column Family:** country
- **Column:** country value
- **Value:** confirmedCases

ROW

20200122
20200122
20200122
20200122
20200123
20200123
20200123
20200123
20200124
20200124
20200124
20200124
20200124

COLUMN+CELL

column=country:CN, timestamp=1587009922682, value=547
column=country:JP, timestamp=1587009922682, value=2
column=country:KR, timestamp=1587009922682, value=1
column=country:US, timestamp=1587009922682, value=1
column=country:CN, timestamp=1587009922682, value=639
column=country:JP, timestamp=1587009922682, value=1
column=country:KR, timestamp=1587009922682, value=1
column=country:US, timestamp=1587009922682, value=1
column=country:CN, timestamp=1587009922682, value=916
column=country:FR, timestamp=1587009922682, value=2
column=country:JP, timestamp=1587009922682, value=2
column=country:KR, timestamp=1587009922682, value=2
column=country:US, timestamp=1587009922682, value=2



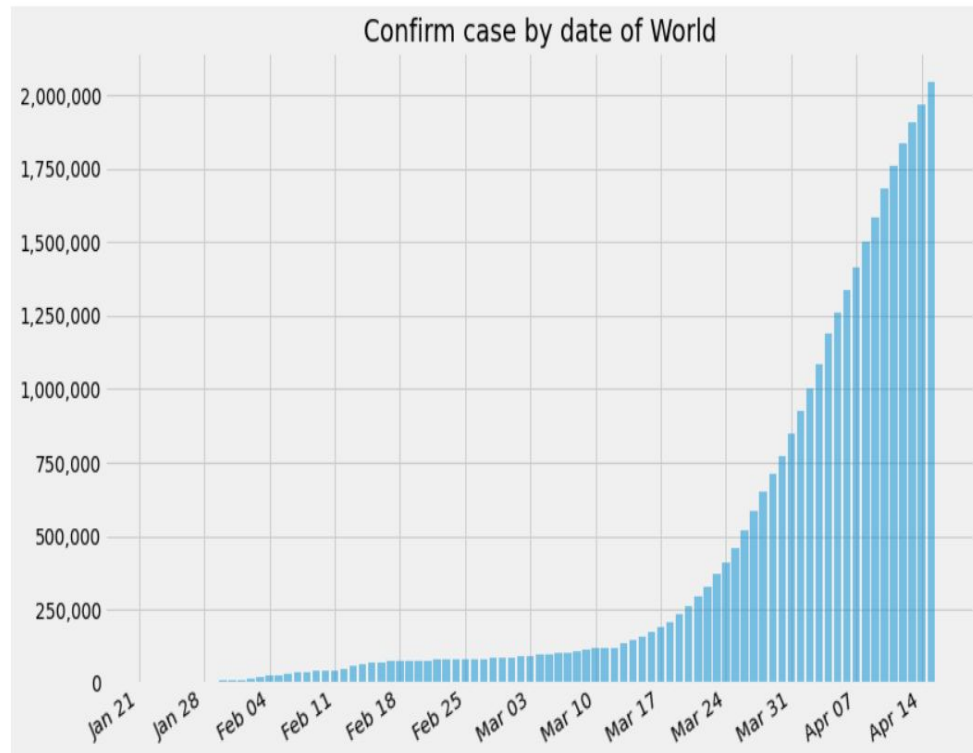
SparkSQL - Performs analysis

1. Scan records in **HBase**
2. Create **SparkSession**
3. Create DataFrame from the **SparkSession** & **HBase** records
4. Query records in DataFrame with **SparkSQL**
5. Display in console/ save data to **HBase** table



Visualize - Colab

1. Bar chart of confirm case of the World



Visualize Realtime - Matplotlib

- Repeat get data from **Hbase** using **happybase**
- Column extract with **panda** and **numpy**
- Render using **FuncAnimation**

```
df = pd.DataFrame(table.scan(), columns=['date', 'countries'])
df2 = pd.concat(
    [
        df.drop(['countries'], axis=1),
        df.countries.apply(pd.Series)
    ], axis=1)

dates = pd.to_datetime(df2['date'].str.decode("utf-8"))
debug(dates)
# print('Last', dates[dates.index[-1]])
if isZeroLast(df2)==True:
    return

plt.cla()
setAxisTick(ax)

for (index, row) in metas.iterrows():
    colName = row['col']
    # print(row['col'], row['id'])

    if colName in df2:
        # print(colName)
        cases = df2[colName].str.decode("utf-8").fillna("0").apply(lambda x: int(x))
        debug(cases)
        ax.plot(dates, cases, label=row['title'])
```



Use happybase to fetch data

```
def connect(tname):
    CONNECTION = happybase.Connection('104.154.17.226', 9090)
    CONNECTION.open()
    # print(CONNECTION.tables())

    _table = CONNECTION.table(tname)
    print(_table)
    return _table

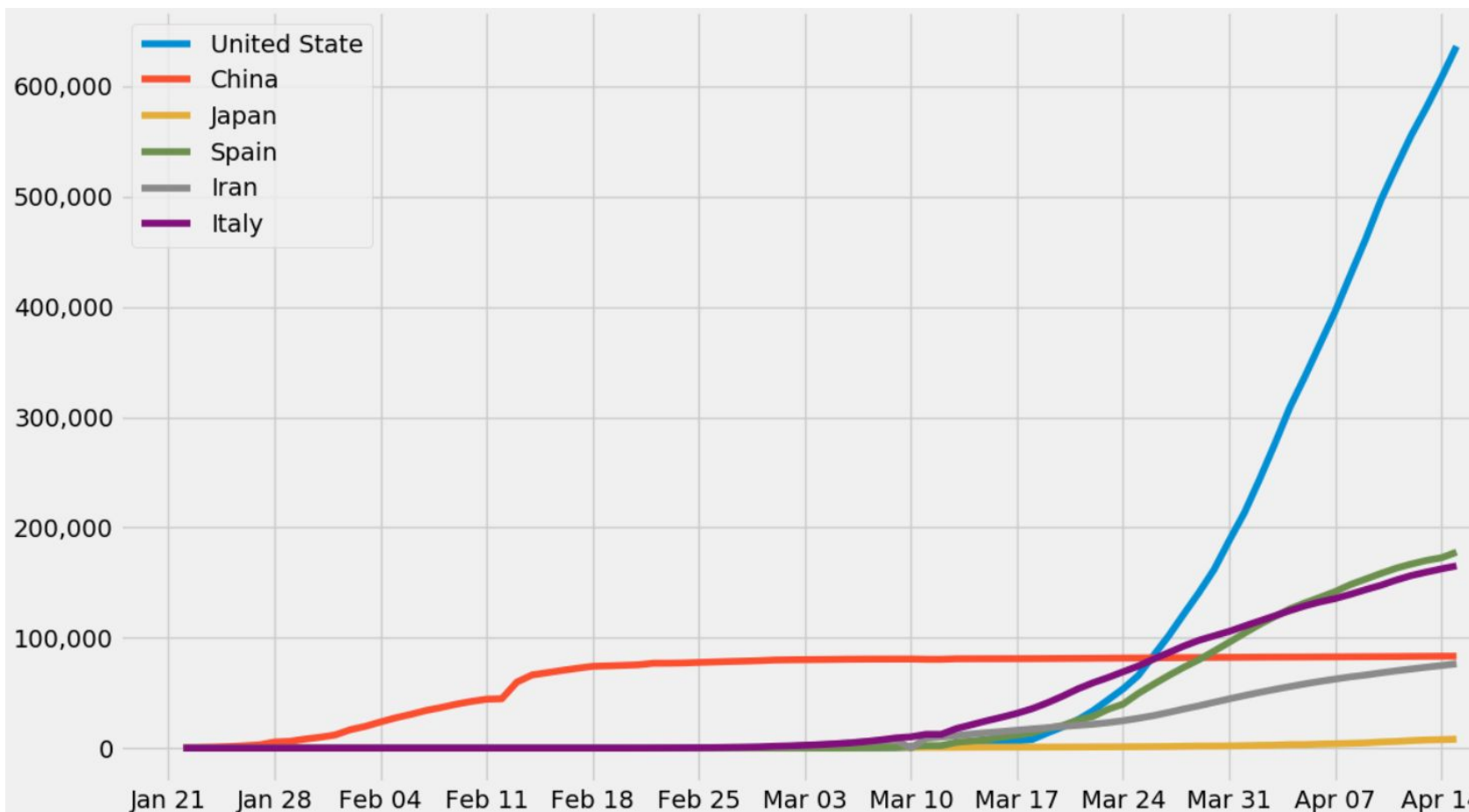
table = connect('pilot')

metas = pd.DataFrame({'col':[
    b'country:US',
    b'country:CN',
    b'country:JP',
    b'country:ES',
    b'country:IR',
    b'country:IT'
],
    'title':['United State', 'China', 'Japan', 'Spain', 'Iran', 'Italy'] })
```



Demo checklist - Realtime

- **Kafka** server ON ?
- Topic created ?
- **Streaming** app started ?
- **Hbase** truncated ?
- **Plot** started ?
- `./sendallfile`
- Go go go!





Sources list

- /proj
 - /kafka
 - **sendkafka.py** -> receive input from stdin and send to kafka
 - **sendallfile.sh, sendfile.sh:** -> scan folder and send csv as input std
 - **./start-kafka.sh, ./recreate-topic.sh,..**
 - /plot
 - **plot.py** -> realtime plot
 - **plothisto.py** -> plot bar chart
 - /cs532
 - /config: **HBaseConfig.java | KafkaConfig.java | SparkConfig.java**
 - /kafka
 - **KConsumer.java**
 - /hbase
 - **HBaseRepository.java**
 - /model
 - **CoronaRecord.java**