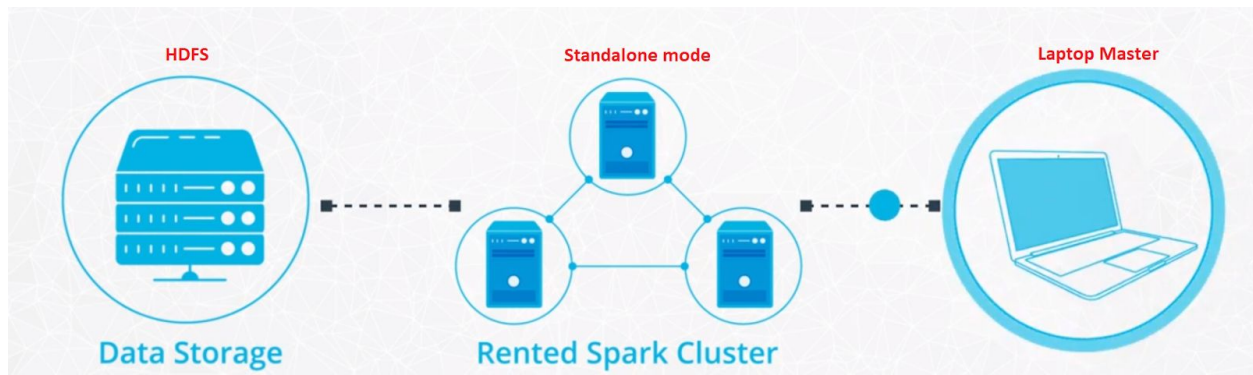# BÁO CÁO LAB 4
## CÀI ĐẶT CỤM SPARK VÀ CHẠY CHƯƠNG TRÌNH WORDCOUNT
## NHÓM: GIỮA CHÚNG TA

Tiếp nối bài cài đặt HDFS + YARN, tính WordCount bằng Hadoop MapReduce, trong bài này chúng ta sẽ cài đặt cụm spark và chạy trương trình WordCount.

**Nguyên lý hoạt động:** Cài đặt Spark Cluster ở Standalone mode. Master gửi job, load dữ liệu từ HDFS, phân tích dữ liệu bằng Spark. Dữ liệu khi phân tích xong được lưu lại vào HDFS, hiển thị kết quả.



hadoopuser@hadoop-master: Master + Worker
hadoopuser@hadoop-slave1: Worker
hadoopuser@hadoop-slave2: Worker

1. **Tải spark về**

   wget http://apache.claz.org/spark/spark-2.4.0/spark-2.4.0-bin-hadoop2.7.tgz

   tar -xzf spark-2.4.0-bin-hadoop2.7.tgz

   mv spark-2.4.0-bin-hadoop2.7 spark

2. **Tạo các biến môi trường**

## 3. Tạo file slave ghi các địa chỉ slave



## 4. Start spark

```
$ cd $SPARK_HOME/sbin
$ ./start-all.sh
```

**5. Check trên từng máy xem đã thành công chưa**

- Với hadoopuser@hadoop-master: Master + Worker



- Với hadoopuser@hadoop-slave1: Worker



- Với hadoopuser@hadoop-slave2: Worker

## 6. Check trên từng máy thành công, kiểm tra tiếp trên webUI



## 7. Chạy chương trình WordCount

- Trước tiên kiểm tra file từ cụm HDFS



- Thực thi chương trình WordCount và up lên cụm HDFS

- Kết quả thu được trên cụm HDFS