

HỆ THỐNG LƯU TRỮ VÀ XỬ LÝ DỮ LIỆU THÔNG TIN CÁC CHUYẾN BAY

Môn học: Lưu trữ và xử lý dữ liệu lớn

Nhóm sinh viên

1. Vũ Trung Nghĩa - 20173284
2. Lê Vũ Lợi - 20173240
3. Đặng Lâm San - 20170111

Ngày 23 tháng 12 năm 2020

- 1 Phát biểu bài toán
- 2 Kiến trúc hệ thống
- 3 Kết quả chạy chương trình
 - Kết quả chạy nhánh Batch Processing
 - Kết quả chạy nhánh Speed Processing
- 4 Đánh giá hệ thống
 - KAFKA
 - HDFS
 - SPARK
 - Hạn chế và hướng phát triển tiếp theo
- 5 Tài liệu tham khảo

1 Phát biểu bài toán

2 Kiến trúc hệ thống

3 Kết quả chạy chương trình

- Kết quả chạy nhánh Batch Processing
- Kết quả chạy nhánh Speed Processing

4 Đánh giá hệ thống

- KAFKA
- HDFS
- SPARK
- Hạn chế và hướng phát triển tiếp theo

5 Tài liệu tham khảo



Đặt vấn đề

- Dữ liệu các chuyến bay trong thực tế thường đến dưới dạng bản ghi (PNR - Passenger Name Record) theo time series với số lượng rất lớn, vì vậy các hệ thống lưu trữ và xử lý dữ liệu truyền thống không còn hiệu quả.
- Cần xây dựng một hệ thống lưu trữ và xử lý dữ liệu phân tán cho dữ liệu lớn với các ràng buộc:
 - Đáp ứng các yêu cầu về tính toán và lưu trữ
 - Có khả năng chịu lỗi tốt khi các máy trong hệ thống không tin cậy
 - Dễ dàng mở rộng
 - Cho phép người dùng có thể xem dữ liệu real-time và thực hiện các truy vấn không real-time như visualize thống kê hay anomaly detection



Đặc tả một bản ghi

ArrivalTime – local time of arrival

BusinessLeisure – if the trip is for business or leisure

CabinCategory – cabin class

CreationDate – PNR creation date (Julian day)

CurrencyCode – 3-letter currency code of payment

DepartureTime – local time of departure

Destination – IATA code of arrival airport

OfficeIdCountry – country code of office placing the reser

Origin – IATA code of departure airport

TotalAmount – total reservation cost

nPAX – number of passengers

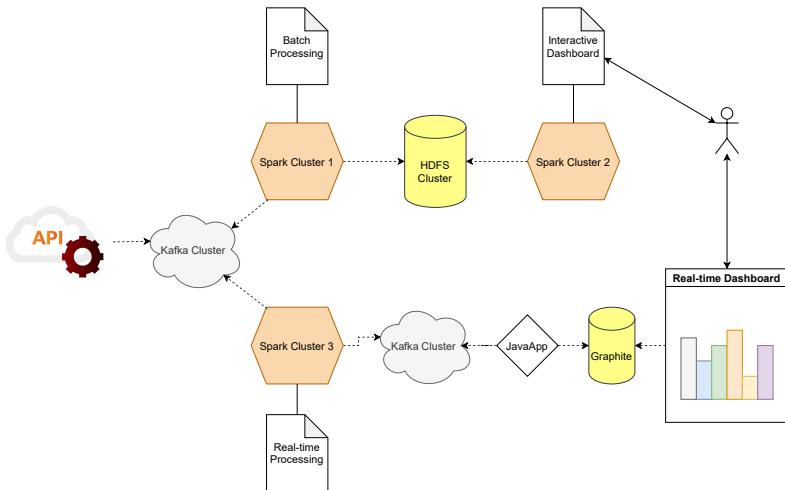
```
{
  "ID": 1188,
  "ArrivalTime": "1453042816",
  "BusinessLeisure": "B",
  "CabinCategory": "40",
  "CreationDate": "2457373",
  "CurrencyCode": "nan",
  "DepartureTime": "1452892672",
  "Destination": "TRD",
  "OfficeIdCountry": "NO",
  "Origin": "ALC",
  "TotalAmount": "nan",
  "nPAX": "1"
}
```



- 1 Phát biểu bài toán
- 2 Kiến trúc hệ thống
- 3 Kết quả chạy chương trình
 - Kết quả chạy nhánh Batch Processing
 - Kết quả chạy nhánh Speed Processing
- 4 Đánh giá hệ thống
 - KAFKA
 - HDFS
 - SPARK
 - Hạn chế và hướng phát triển tiếp theo
- 5 Tài liệu tham khảo



Kiến trúc tổng quan



HDFS Cluster

```
namenode:
  image: bde2020/hadoop-namenode:2.0.0-hadoop2.7.4-java8
  container_name: namenode
  environment:
    - CLUSTER_NAME=test
  env_file:
    - ./hadoop.env
  ports:
    - 8020:8020
    - 50070:50070
datanode-1:
  image: bde2020/hadoop-datanode:2.0.0-hadoop2.7.4-java8
  container_name: datanode-1
  environment:
    SERVICE_PRECONDITION: "namenode:50070"
  env_file:
    - ./hadoop.env
  ports:
    - 50075:50075
datanode-2:
  image: bde2020/hadoop-datanode:2.0.0-hadoop2.7.4-java8
  container_name: datanode-2
  environment:
    SERVICE_PRECONDITION: "namenode:50070"
  env_file:
    - ./hadoop.env
  ports:
    - 50076:50075
```



Kafka Cluster

```

zookeeper:
  image: zookeeper:3.4.10
  container_name: zookeeper
  environment:
    ZOO_MY_ID: 1
    ZOO_SERVERS: server.1=0.0.0.0:2888:3888
    ZOO_TICK_TIME: 15000
  ports:
    - 2181:2181

kafka-broker-1:
  image: wurstmeister/kafka:2.12-2.2.1
  container_name: kafka-broker-1
  depends_on:
    - zookeeper
  ports:
    - "9092:9092"
  environment:
    - KAFKA_ZOOKEEPER_CONNECT=zookeeper:2181
    - ALLOW_PLAINTEXT_LISTENER=yes
    - KAFKA_ADVERTISED_LISTENERS=INSIDE://kafka-broker-1:9093,OUTSIDE://localhost:9092
    - KAFKA_LISTENER_SECURITY_PROTOCOL_MAP=INSIDE:PLAINTEXT,OUTSIDE:PLAINTEXT
    - KAFKA_LISTENERS=INSIDE://kafka-broker-1:9093,OUTSIDE://0.0.0.0:9092
    - KAFKA_INTER_BROKER_LISTENER_NAME=INSIDE

kafka-broker-2:
  image: wurstmeister/kafka:2.12-2.2.1
  container_name: kafka-broker-2
  depends_on:
    - zookeeper
  ports:
    - "9094:9094"
  environment:
    - KAFKA_ZOOKEEPER_CONNECT=zookeeper:2181
    - ALLOW_PLAINTEXT_LISTENER=yes
    - KAFKA_ADVERTISED_LISTENERS=INSIDE://kafka-broker-2:9093,OUTSIDE://localhost:9094
    - KAFKA_LISTENER_SECURITY_PROTOCOL_MAP=INSIDE:PLAINTEXT,OUTSIDE:PLAINTEXT
    - KAFKA_LISTENERS=INSIDE://kafka-broker-2:9093,OUTSIDE://0.0.0.0:9094
    - KAFKA_INTER_BROKER_LISTENER_NAME=INSIDE
  
```



Spark Cluster

```

post-batch-processing-spark-master:
  image: vutrungnghia99/spark-master:spark2.4.1-python3.7-hadoop2.7
  container_name: post-batch-processing-spark-master
  ports:
    - "8083:8080"
    - "7078:7077"
  environment:
    - INIT_DAEMON_STEP=setup_spark
post-batch-processing-spark-worker-1:
  image: vutrungnghia99/spark-worker:spark2.4.1-python3.7-hadoop2.7
  container_name: post-batch-processing-spark-worker-1
  depends_on:
    - post-batch-processing-spark-master
  environment:
    - "SPARK_MASTER=spark://post-batch-processing-spark-master:7077"
    - "SPARK_WORKER_CORES=1"
    - "SPARK_WORKER_MEMORY=1G"
    - "SPARK_DRIVER_MEMORY=128m"
    - "SPARK_EXECUTOR_MEMORY=256m"
post-batch-processing-spark-worker-2:
  image: vutrungnghia99/spark-worker:spark2.4.1-python3.7-hadoop2.7
  container_name: post-batch-processing-spark-worker-2
  depends_on:
    - post-batch-processing-spark-master
  environment:
    - "SPARK_MASTER=spark://post-batch-processing-spark-master:7077"
    - "SPARK_WORKER_CORES=1"
    - "SPARK_WORKER_MEMORY=1G"
    - "SPARK_DRIVER_MEMORY=128m"
    - "SPARK_EXECUTOR_MEMORY=256m"

```



System manager - Graphite - Grafana

```
##### manager #####
```

```
system-manager:
```

```
  image: vutrungnghia99/system-manager:spark2.4.1-python3.7-hadoop2.7-kafka2.7.0
```

```
  container_name: system-manager
```

```
  ports:
```

```
    - "8888:8888"
```

```
  volumes:
```

```
    - $PWD/src:/home/jovyan/work
```

```
  environment:
```

```
    - JUPYTER_TOKEN=admin
```

```
##### Graphite and Grafana #####
```

```
graphite:
```

```
  image: vutrungnghia99/graphite:1.1.7-6
```

```
  container_name: graphite
```

```
  ports:
```

```
    - "80:80"
```

```
    - "2003:2003"
```

```
    - "2004:2004"
```

```
grafana:
```

```
  image: grafana/grafana:latest
```

```
  container_name: grafana
```

```
  ports:
```

```
    - "3000:3000"
```



- 1 Phát biểu bài toán
- 2 Kiến trúc hệ thống
- 3 **Kết quả chạy chương trình**
 - Kết quả chạy nhánh Batch Processing
 - Kết quả chạy nhánh Speed Processing
- 4 Đánh giá hệ thống
 - KAFKA
 - HDFS
 - SPARK
 - Hạn chế và hướng phát triển tiếp theo
- 5 Tài liệu tham khảo



- 1 Phát biểu bài toán
- 2 Kiến trúc hệ thống
- 3 Kết quả chạy chương trình
 - Kết quả chạy nhánh Batch Processing
 - Kết quả chạy nhánh Speed Processing
- 4 Đánh giá hệ thống
 - KAFKA
 - HDFS
 - SPARK
 - Hạn chế và hướng phát triển tiếp theo
- 5 Tài liệu tham khảo



Danh sách containers [16]

IMAGE

```
vutrungnghia99/spark-worker:spark2.4.1-python3.7-hadoop2.7
vutrungnghia99/spark-worker:spark2.4.1-python3.7-hadoop2.7
vutrungnghia99/spark-worker:spark2.4.1-python3.7-hadoop2.7
wurstmeister/kafka:2.12-2.2.1
wurstmeister/kafka:2.12-2.2.1
vutrungnghia99/spark-worker:spark2.4.1-python3.7-hadoop2.7
vutrungnghia99/spark-master:spark2.4.1-python3.7-hadoop2.7
zookeeper:3.4.10
vutrungnghia99/system-manager:spark2.4.1-python3.7-hadoop2.7-kafka2.7.0
vutrungnghia99/spark-master:spark2.4.1-python3.7-hadoop2.7
bde2020/hadoop-datanode:2.0.0-hadoop2.7.4-java8
grafana/grafana:latest
vutrungnghia99/spark-master:spark2.4.1-python3.7-hadoop2.7
bde2020/hadoop-namenode:2.0.0-hadoop2.7.4-java8
vutrungnghia99/graphite:1.1.7-6
bde2020/hadoop-datanode:2.0.0-hadoop2.7.4-java8
```

NAMES

```
post-batch-processing-spark-worker-2
post-batch-processing-spark-worker-1
speed-processing-spark-worker-1
kafka-broker-1
kafka-broker-2
pre-batch-processing-spark-worker-1
speed-processing-spark-master
zookeeper
system-manager
post-batch-processing-spark-master
datanode-2
grafana
pre-batch-processing-spark-master
namenode
graphite
datanode-1
```

Spark cluster 1



Spark Master at spark://7a6c33b85408:7077

URL: spark://7a6c33b85408:7077

Alive Workers: 1

Cores in use: 1 Total, 0 Used

Memory in use: 1024.0 MB Total, 0.0 B Used

Applications: 0 Running, 0 Completed

Drivers: 0 Running, 0 Completed

Status: ALIVE

Workers (1)

Worker Id	Address
worker-20201223081126-172.18.0.12-34753	172.18.0.12:34753

Running Applications (0)

Application ID	Name	Cores	Memory per Executor
----------------	------	-------	---------------------

Completed Applications (0)

Application ID	Name	Cores	Memory per Executor
----------------	------	-------	---------------------



Chương trình BatchProcessing

```
def get_categorical(x, m):
    if str(x) == 'nan':
        return 0.0
    else:
        v = m['mapping'][str(x)]
        return (v - m['statistic']['mean']) / m['statistic']['std']

def json_to_processed_data(s):
    t = json.loads(s)
    return [
        t['ID'],
        get_continuous(t['ArrivalTime'], mapping_and_statistic['ArrivalTime']),
        get_categorical(t['BusinessLeisure'], mapping_and_statistic['BusinessLeisure']),
        get_categorical(t['CabinCategory'], mapping_and_statistic['CabinCategory']),
        get_continuous(t['CreationDate'], mapping_and_statistic['CreationDate']),
        get_categorical(t['CurrencyCode'], mapping_and_statistic['CurrencyCode']),
        get_continuous(t['DepartureTime'], mapping_and_statistic['DepartureTime']),
        get_categorical(t['Destination'], mapping_and_statistic['Destination']),
        get_categorical(t['OfficeIdCountry'], mapping_and_statistic['OfficeIdCountry']),
        get_categorical(t['Origin'], mapping_and_statistic['Origin']),
        get_continuous(t['TotalAmount'], mapping_and_statistic['TotalAmount']),
        get_continuous(t['nPAX'], mapping_and_statistic['nPAX']),
        s
    ]
```

```
In [*]: ks = KafkaUtils.createDirectStream(
        ssc, [t['trips'], {'metadata.broker.list': 'kafka-broker-1:9093,kafka-broker-2:9093'}])
        lines = ks.map(lambda x: x[1])

        transform1 = lines.map(lambda tripInfo: json_to_list(tripInfo))
        transform1.foreachRDD(handle_rdd1)

        transform2 = lines.map(lambda tripInfo: json_to_processed_data(tripInfo))
        transform2.foreachRDD(handle_rdd2)

        ssc.start()
        ssc.awaitTermination()
```



3. Kết quả chạy chương trình

3.1. Kết quả chạy nhánh Batch Processing

Spark cluster 2



Spark Master at spark://3d40737aa3f2:7077

URL: spark://3d40737aa3f2:7077

Alive Workers: 2

Cores in use: 2 Total, 0 Used

Memory in use: 2.0 GB Total, 0.0 B Used

Applications: 0 Running, 0 Completed

Drivers: 0 Running, 0 Completed

Status: ALIVE

Workers (2)

Worker Id	Address
worker-20201223081126-172.18.0.16-38487	172.18.0.16:38487
worker-20201223081126-172.18.0.17-37513	172.18.0.17:37513

Running Applications (0)

Application ID	Name	Cores	Memory per Executor
----------------	------	-------	---------------------

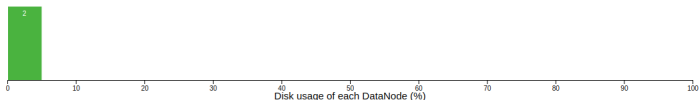
Completed Applications (0)



HDFS

Datanode Information

Datanode usage histogram



In operation

Node	Last contact	Admin State	Capacity	Used	Non DFS Used	Remaining	Blocks	Block pool used	Failed Volumes	Version
c65069a49aa9:50010 (172.18.0.3:50010)	1	In Service	355.6 GB	24 KB	83.51 GB	253.96 GB	0	24 KB (0%)	0	2.7.4
5abe514e7d51:50010 (172.18.0.8:50010)	0	In Service	355.6 GB	24 KB	83.51 GB	253.96 GB	0	24 KB (0%)	0	2.7.4

Decommissioning

Node	Last contact	Under replicated blocks	Blocks with no live replicas	Under Replicated Blocks In files under construction
------	--------------	-------------------------	------------------------------	--

HDFS

[Hadoop](#) [Overview](#) [Datanodes](#) [Snapshot](#) [Startup Progress](#) [Utilities](#) -

Browse Directory

Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
Hadoop, 2017.							



3.1. Kết quả chạy nhánh Batch Processing

```
Topic: trips      PartitionCount: 2      ReplicationFactor: 2      Configs: segment.bytes=1073741824
      Topic: trips      Partition: 0      Leader: 1002      Replicas: 1002,1001      Isr: 1002,1001
      Topic: trips      Partition: 1      Leader: 1001      Replicas: 1001,1002      Isr: 1001,1002
```

late event
late event
late event
late event
late event
late event

100% |

Sent 10000 records in 23.760411262512207 seconds

Sending rate: 420.86813605694687 records/s



Partition 0 trên hai brokers

```
ation": "LHR", "OfficeIdCountry": "GB", "Origin": "LHR", "TotalAmount": "nan",  
    "nPAx": "1"}}, {"ID": 9981, "ArrivalTime": "144346  
9312", "BusinessLeisure": "nan", "CabinCategory": "40", "CreationDate": "24571  
98", "CurrencyCode": "nan", "DepartureTime": "1443745408", "Destination": "KUL  
", "OfficeIdCountry": "DE", "Origin": "MUC", "TotalAmount": "nan", "nPAx": "1"  
} }, {"ID": 9982, "ArrivalTime": "1448134912", "Bu  
sinessLeisure": "nan", "CabinCategory": "40", "CreationDate": "2457292", "Cur  
rencyCode": "nan", "DepartureTime": "1445481088", "Destination": "BCN", "Office  
IdCountry": "SE", "Origin": "ARN", "TotalAmount": "nan", "nPAx": "1"}}, {"ID":  
9983, "ArrivalTime": "1437361280", "BusinessLeisure": "nan", "CabinCategory": "  
40", "CreationDate": "2457162", "CurrencyCode": "nan", "DepartureTime": "1436  
33312", "Destination": "MLA", "OfficeIdCountry": "FR", "Origin": "CDG", "Tota  
lAmount": "nan", "nPAx": "1"}}, {"ID": 9985, "Arr  
ivalTime": "1444902656", "BusinessLeisure": "nan", "CabinCategory": "40", "Cre  
ationDate": "2457162", "CurrencyCode": "nan", "DepartureTime": "1445120256", "  
Destination": "JNB", "OfficeIdCountry": "PL", "Origin": "PRG", "TotalAmount":  
"nan", "nPAx": "1"}}, {"ID": 9986, "ArrivalTime":  
"143966384", "BusinessLeisure": "nan", "CabinCategory": "40", "CreationDate":  
"2457148", "CurrencyCode": "SAR", "DepartureTime": "1435167232", "Destination  
": "KHI", "OfficeIdCountry": "SA", "Origin": "RUH", "TotalAmount": "0.0", "nPA  
X": "1"} } bash-4.4#
```

```
ation": "LHR", "OfficeIdCountry": "GB", "Origin": "LHR", "TotalAmount": "nan",  
    "nPAx": "1"}H0vooooooooooooooooooooo{"ID": 9981, "ArrivalTime": "144346  
9312", "BusinessLeisure": "nan", "CabinCategory": "40", "CreationDate": "24571  
98", "CurrencyCode": "nan", "DepartureTime": "1443745408", "Destination": "KL"  
    , "OfficeIdCountry": "DE", "Origin": "MUC", "TotalAmount": "nan", "nPAx": "1"},  
    {"ID": 9982, "ArrivalTime": "1448134912", "Bu  
ssinessLeisure": "nan", "CabinCategory": "40", "CreationDate": "2457292", "Cur  
urrencyCode": "nan", "DepartureTime": "1445481088", "Destination": "BCN", "Office  
IdCountry": "SE", "Origin": "ARN", "TotalAmount": "nan", "nPAx": "1"}{{"ID":  
9983, "ArrivalTime": "1437361280", "BusinessLeisure": "nan", "CabinCategory":  
40", "CreationDate": "2457162", "CurrencyCode": "nan", "DepartureTime": "1436  
33312", "Destination": "MLA", "OfficeIdCountry": "FR", "Origin": "CDG", "Tota  
lAmount": "nan", "nPAx": "1"}H0vooooooooooooooooooooo{"ID": 9985, "Arr  
ivalTime": "1444902656", "BusinessLeisure": "nan", "CabinCategory": "40", "Cre  
ationDate": "2457162", "CurrencyCode": "nan", "DepartureTime": "1445120256", "  
Destination": "JNB", "OfficeIdCountry": "PL", "Origin": "PRG", "TotalAmount":  
"nan", "nPAx": "1"}H0vooooooooooooooooooooo{"ID": 9986, "ArrivalTime":  
"1439663804", "BusinessLeisure": "nan", "CabinCategory": "40", "CreationDate":  
"2457148", "CurrencyCode": "SAR", "DepartureTime": "1435167232", "Destinat  
ion": "KHI", "OfficeIdCountry": "SA", "Origin": "RUH", "TotalAmount": "0.0", "nPAX":  
"1"}bash-4.4#
```



Partition 0,1 trên broker 0

```

ation": "LHR", "OfficeIdCountry": "GB", "Origin": "LHR", "TotalAmount": "nan",
  "nPAX": "1"}oHo, oVooooVooooooooooooooooooooo{"ID": 9981, "ArrivalTime": "144346
9312", "BusinessLeisure": "nan", "CabinCategory": "40", "CreationDate": "24571
98", "CurrencyCode": "nan", "DepartureTime": "1443745408", "Destination": "KUL
", "OfficeIdCountry": "DE", "Origin": "MUC", "TotalAmount": "nan", "nPAX": "1"
}o o 6oVoooo!Vooo{"ID": 9982, "ArrivalTime": "1448134912", "Bu
sinessLeisure": "nan", "CabinCategory": "40", "CreationDate": "2457292", "Curr
encyCode": "nan", "DepartureTime": "1445481088", "Destination": "BCN", "Office
IdCountry": "SE", "Origin": "ARN", "TotalAmount": "nan", "nPAX": "1"}oo{"ID":
9983, "ArrivalTime": "1437361280", "BusinessLeisure": "nan", "CabinCategory":
"40", "CreationDate": "2457162", "CurrencyCode": "nan", "DepartureTime": "1436
333312", "Destination": "MLA", "OfficeIdCountry": "FR", "Origin": "CDG", "Tota
lAmount": "nan", "nPAX": "1"}oHsoVooooVooooooooooooooooooooo{"ID": 9985, "Arr
ivalTime": "1444902656", "BusinessLeisure": "nan", "CabinCategory": "40", "Cre
ationDate": "2457162", "CurrencyCode": "nan", "DepartureTime": "1445120256", "
Destination": "JNB", "OfficeIdCountry": "PL", "Origin": "PRG", "TotalAmount":
"nan", "nPAX": "1"}oH-oVooooVooooooooooooooooooooo{"ID": 9986, "ArrivalTime":
"1439666304", "BusinessLeisure": "nan", "CabinCategory": "40", "CreationDate":
"2457148", "CurrencyCode": "SAR", "DepartureTime": "1435167232", "Destination
": "KHI", "OfficeIdCountry": "SA", "Origin": "RUH", "TotalAmount": "0.0", "nP
AX": "1"}bash-4.4# █

oooo{"ID": 9994, "ArrivalTime": "1453793152", "BusinessLeisure": "nan", "Cabi
nCategory": "40", "CreationDate": "2457333", "CurrencyCode": "nan", "Departure
Time": "1452837888", "Destination": "UIO", "OfficeIdCountry": "FR", "Origin":
"ORY", "TotalAmount": "nan", "nPAX": "1"}oHos(oVooooVooooooooooooooooooooo{"ID"
: 9995, "ArrivalTime": "1425617408", "BusinessLeisure": "nan", "CabinCategory"
: "40", "CreationDate": "2457058", "CurrencyCode": "nan", "DepartureTime": "14
26283392", "Destination": "MAD", "OfficeIdCountry": "BE", "Origin": "BIO", "To
talAmount": "nan", "nPAX": "1"}pHRo%VoooIvoooIooooooooooooooooooooo{"ID": 9997, "A
rrivalTime": "1452245376", "BusinessLeisure": "nan", "CabinCategory": "40", "C
reationDate": "2457297", "CurrencyCode": "nan", "DepartureTime": "1449478144",
"Destination": "CCS", "OfficeIdCountry": "DE", "Origin": "FRA", "TotalAmount"
: "nan", "nPAX": "1"}qF-%oVooooVooooooooooooooooooooo{"ID": 9998, "ArrivalTime
": "1417529216", "BusinessLeisure": "L", "CabinCategory": "40", "CreationDate"
: "2456918", "CurrencyCode": "nan", "DepartureTime": "1415165568", "Destinatio
n": "MLE", "OfficeIdCountry": "DK", "Origin": "CPH", "TotalAmount": "nan", "nP
AX": "2"}rH

CoVooooVooooooooooooooooooooo{"ID": 9999, "ArrivalTime": "1438053120
", "BusinessLeisure": "nan", "CabinCategory": "40", "CreationDate": "2457169",
"CurrencyCode": "nan", "DepartureTime": "1438486912", "Destination": "TNR", "
OfficeIdCountry": "KE", "Origin": "LPA", "TotalAmount": "nan", "nPAX": "1"}bas
h-4.4# █

```



Kết quả xử lý và nhận trên Batch Processing Notebook

```
In [*]: ks = KafkaUtils.createDirectStream(
          ssc, ['trips'], {'metadata.broker.list': 'kafka-broker-1:9093,kafka-broker-2:9093'})
          lines = ks.map(lambda x: x[1])

          transform1 = lines.map(lambda tripInfo: json_to_list(tripInfo))
          transform1.foreachRDD(handle_rdd1)

          transform2 = lines.map(lambda tripInfo: json_to_processed_data(tripInfo))
          transform2.foreachRDD(handle_rdd2)

          ssc.start()
          ssc.awaitTermination()
```

```
Recieved 606 records - transform 1
Recieved 606 records - transform 2
Recieved 5832 records - transform 1
Recieved 5832 records - transform 2
Recieved 3562 records - transform 1
Recieved 3562 records - transform 2
```

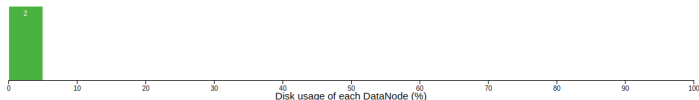


Dữ liệu lưu trong HDFS



Datanode Information

Datanode usage histogram



In operation

Node	Last contact	Admin State	Capacity	Used	Non DFS Used	Remaining	Blocks	Block pool used	Failed Volumes	Version
c65069a49aa9:50010 (172.18.0.3:50010)	2	In Service	355.6 GB	1.91 MB	84.55 GB	252.92 GB	24	1.91 MB (0%)	0	2.7.4
5abe514e7d51:50010 (172.18.0.8:50010)	1	In Service	355.6 GB	1.91 MB	84.55 GB	252.92 GB	24	1.91 MB (0%)	0	2.7.4



Dữ liệu lưu trong HDFS

[Hadoop](#) [Overview](#) [Datanodes](#) [Snapshot](#) [Startup Progress](#) [Utilities](#) [...](#)

Browse Directory

Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
drwxr-xr-x	jovyan	supergroup	0 B	12/23/2020, 3:29:46 PM	0	0 B	processed_trips.parquet
drwxr-xr-x	jovyan	supergroup	0 B	12/23/2020, 3:29:45 PM	0	0 B	trips.parquet

Hadoop, 2017.



3. Kết quả chạy chương trình

3.1. Kết quả chạy nhánh Batch Processing

Dữ liệu lưu trong HDFS

Hadoop Overview Datanodes Snapshot Startup Progress Utilities -								
Browse Directory								
/trips/trips.parquet								Go!
Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name	
-rw-r--r--	joyvan	supergroup	0 B	12/23/2020, 3:29:45 PM	3	128 MB	_SUCCESS	
-rw-r--r--	joyvan	supergroup	88.21 KB	12/23/2020, 3:29:41 PM	3	128 MB	part-00000-19bb1bf0-592d-492d-9595-8487c8c2bcc7-c000.snappy.parquet	
-rw-r--r--	joyvan	supergroup	30.03 KB	12/23/2020, 3:29:45 PM	3	128 MB	part-00000-29e54c4e-d45c-4c86-8dbf-5d6dc4cc17d-c000.snappy.parquet	
-rw-r--r--	joyvan	supergroup	87.09 KB	12/23/2020, 3:29:33 PM	3	128 MB	part-00000-4d646859-ab52-4db8-ad0a-aaec01a80740-c000.snappy.parquet	
-rw-r--r--	joyvan	supergroup	91.86 KB	12/23/2020, 3:29:38 PM	3	128 MB	part-00000-5ad763d6-a1ff-4f85-993c-4c36eb4580f6-c000.snappy.parquet	
-rw-r--r--	joyvan	supergroup	90.76 KB	12/23/2020, 3:29:36 PM	3	128 MB	part-00000-74231608-9c1f-4f52-9d4f-e1ea64568ebb-c000.snappy.parquet	
-rw-r--r--	joyvan	supergroup	42.25 KB	12/23/2020, 3:29:30 PM	3	128 MB	part-00000-c90b3fa1-d391-4476-93b1-ed7d511ed249-c000.snappy.parquet	
-rw-r--r--	joyvan	supergroup	81.64 KB	12/23/2020, 3:29:41 PM	3	128 MB	part-00001-19bb1bf0-592d-492d-9595-8487c8c2bcc7-c000.snappy.parquet	
-rw-r--r--	joyvan	supergroup	29.43 KB	12/23/2020, 3:29:45 PM	3	128 MB	part-00001-29e54c4e-d45c-4c86-8dbf-5d6dc4cc17d-c000.snappy.parquet	
-rw-r--r--	joyvan	supergroup	91.71 KB	12/23/2020, 3:29:33 PM	3	128 MB	part-00001-4d646859-ab52-4db8-ad0a-aaec01a80740-c000.snappy.parquet	
-rw-r--r--	joyvan	supergroup	88.01 KB	12/23/2020, 3:29:39 PM	3	128 MB	part-00001-5ad763d6-a1ff-4f85-993c-4c36eb4580f6-c000.snappy.parquet	



Đọc dữ liệu từ HDFS và visualize

Read data from parquet file "trips.parquet" in hdfs

```
df = ss.read.parquet("hdfs://namenode:9000/trips/trips.parquet")
```

```
print(f"Number of records: {df.count()}")
df = df.sort('ArrivalTime')
```

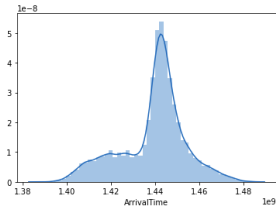
Number of records: 10000

Data Mining

Distribution of trips over time

```
arrivalTime = df.select('ArrivalTime').toPandas()['ArrivalTime'].astype('int64')
sns.distplot(arrivalTime)
```

<matplotlib.axes._subplots.AxesSubplot at 0x7f27d875bf90>



- 1 Phát biểu bài toán
- 2 Kiến trúc hệ thống
- 3 Kết quả chạy chương trình
 - Kết quả chạy nhánh Batch Processing
 - Kết quả chạy nhánh Speed Processing
- 4 Đánh giá hệ thống
 - KAFKA
 - HDFS
 - SPARK
 - Hạn chế và hướng phát triển tiếp theo
- 5 Tài liệu tham khảo



Đưa dữ liệu real time từ kafka => kafka

- Chương trình spark đọc dữ liệu từ kafka, xử lý và đẩy vào kafka

```

ss = SparkSession.Builder() \
    .appName("SparkBatchStreamingKafka") \
    .master("spark://speed-processing-spark-master:7077") \
    .config("spark.jars", "./spark-streaming-kafka-0-10-assembly_2.11-2.4.1.jar,./kafka-clients-0.10.1.0.jar,./spark-sql-kafka-0-10_2.11-2.4.1.jar") \
    .config("spark.sql.warehouse.dir", "hdfs://namenode:9000/") \
    .getOrCreate()

df = ss \
    .readStream \
    .format("kafka") \
    .option("kafka.bootstrap.servers", "kafka-broker-1:9093,kafka-broker-2:9093,kafka-broker-3:9093") \
    .option("partition.assignment.strategy", "none") \
    .option("subscribe", "trips") \
    .load()

import random

def transform_window(s):
    """
    s = Row(start=datetime.datetime(2020, 12, 21, 17, 9, 30), end=datetime.datetime(2020, 12, 21, 17, 9, 40))
    """
    return str(int(s.end.timestamp()))

def transform_count(s):
    """
    s = 941
    """
    return str(s)

udf_transform_window = udf(transform_window)
udf_transform_count = udf(transform_count)

query = df.withWatermark("timestamp", "15 seconds") \
    .groupBy(window("timestamp", "5 seconds", "5 seconds")) \
    .count() \
    .withColumn("count", udf_transform_count("count")) \
    .withColumn("window", udf_transform_window("window")) \
    .withColumn("value", sf.concat(sf.col('window'),sf.lit('_'), sf.col('count')))) \
    .writeStream \
    .format("kafka") \
    .option("kafka.bootstrap.servers", "kafka-broker-1:9093,kafka-broker-2:9093,kafka-broker-3:9093") \
    .option("topic", "real-time-statistic") \
    .option("checkpointLocation", "/tmp/checkpoint") \
    .outputMode("append") \
    .option("truncate", False) \
    .start()

querv.awaitTermination()

```



Đưa dữ liệu real time từ kafka vào graphite

- Chương trình chạy java để đẩy dữ liệu từ kafka vào graphite

```

① localhost:terminal$
jupyter

group.id = bigdata
heartbeat.interval.ms = 3000
interceptor.classes = null
key.deserializer = class org.apache.kafka.common.serialization.StringDeserializer
max.partition.fetch.bytes = 2840576
max.poll.interval.ms = 300000
max.poll.records = 500
metadata.max.age.ms = 300000
metric.reporters = []
metrics.num.samples = 2
metrics.recording.level = INFO
metrics.sample.window.ms = 30000
partition.assignment.strategy = [class org.apache.kafka.clients.consumer.RangeAssignor]
receive.buffer.bytes = 65536
reconnect.backoff.ms = 60
request.timeout.ms = 305000
retry.backoff.ms = 100
sasl.jaas.config = null
sasl.kerberos.kinit.cmd = /usr/bin/kinit
sasl.kerberos.min.time.before.relogin = 60000
sasl.kerberos.service.name = null
sasl.kerberos.ticket.renew.jitter = 0.05
sasl.kerberos.ticket.renew.window.factor = 0.8
sasl.mechanism = GSSAPI
security.protocol = PLAINTEXT
send.buffer.bytes = 512072
session.timeout.ms = 180000
ssl.cipher.suites = null
ssl.enabled.protocols = [TLSv1.2, TLSv1.1, TLSv1]
ssl.endpoint.identification.algorithm = null
ssl.key.password = null
ssl.keymanager.algorithm = SunX509
ssl.keystore.location = null
ssl.keystore.password = null
ssl.keystore.type = JKS
ssl.protocol = TLS
ssl.provider = null
ssl.secure.random.implementation = null
ssl.trustmanager.algorithm = PKIX
ssl.truststore.location = null
ssl.truststore.password = null
ssl.truststore.type = JKS
value.deserializer = class org.apache.kafka.common.serialization.StringDeserializer

08:26:17.711 [main] INFO o.a.kafka.common.utils.AppInfoParser - Kafka version: 0.10.2.0
08:26:17.711 [main] INFO o.a.kafka.common.utils.AppInfoParser - Kafka commitId: 974093abdc0cf421
08:26:17.761 [main] INFO o.a.k.c.c.i.AbstractCoordinator - Discovered coordinator kafka-broker-1:9093 (id: 2147482646 rack
null) for group bigdata
08:26:17.764 [main] INFO o.a.k.c.c.i.ConsumerCoordinator - Revoking previously assigned partitions [] for group bigdata
08:26:17.760 [main] INFO o.a.k.c.c.i.AbstractCoordinator - (Re-)joining group bigdata
08:26:17.817 [main] INFO o.a.k.c.c.i.AbstractCoordinator - Successfully joined group bigdata with generation 1
08:26:17.838 [main] INFO o.a.k.c.c.i.ConsumerCoordinator - Setting newly assigned partitions [real-time-statistic-1, real-
time-statistic-0] for group bigdata

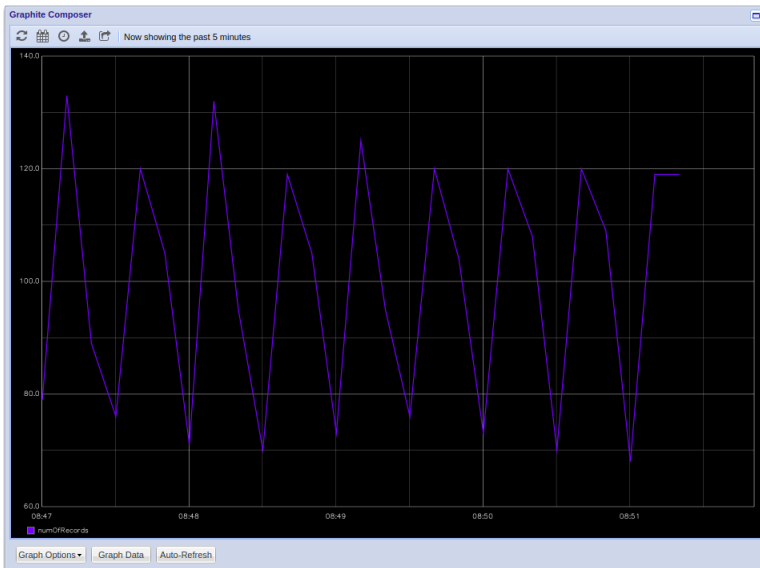
```



3. Kết quả chạy chương trình

3.2. Kết quả chạy nhánh Speed Processing

Dữ liệu sau khi được đưa vào graphite



3. Kết quả chạy chương trình

3.2. Kết quả chạy nhánh Speed Processing

Visualize dữ liệu của graphite với grafana



- 1 Phát biểu bài toán
- 2 Kiến trúc hệ thống
- 3 Kết quả chạy chương trình
 - Kết quả chạy nhánh Batch Processing
 - Kết quả chạy nhánh Speed Processing
- 4 **Đánh giá hệ thống**
 - KAFKA
 - HDFS
 - SPARK
 - Hạn chế và hướng phát triển tiếp theo
- 5 Tài liệu tham khảo



- 1 Phát biểu bài toán
- 2 Kiến trúc hệ thống
- 3 Kết quả chạy chương trình
 - Kết quả chạy nhánh Batch Processing
 - Kết quả chạy nhánh Speed Processing
- 4 Đánh giá hệ thống
 - KAFKA
 - HDFS
 - SPARK
 - Hạn chế và hướng phát triển tiếp theo
- 5 Tài liệu tham khảo



Đánh giá khả năng chịu tải của kafka

	2 brokers			1 brokers		
Delay (s)	Tốc độ gửi (records/s)	Số gói tin nhận	Số gói tin bị mất	Tốc độ gửi (records/s)	Số gói tin nhận	Số gói tin bị mất
0	4501.5	10000	0	4990.3	10000	0
	5045.5	10000	0	4945.3	10000	0
	5120.8	10000	0	5003.4	9973	27
0.0000001	1525.6	10000	0	1572.5	10000	0
	1621.1	10000	0	1962.5	10000	0
	1551.4	10000	0	1558.5	10000	0
0.000001	1335.4	10000	0	1126.4	10000	0
	1596.8	10000	0	1866.5	10000	0
	1597.2	10000	0	1581.1	10000	0

Đánh giá khả năng chịu lỗi của kafka

```
(base) vutrungnghia@Lusheeta:~/kafka_2.13-2.6.0$ bin/kafka-topics.sh --describe --topic trips --bootstrap
p-server localhost:9092,localhost:9094
Topic: trips      PartitionCount: 2      ReplicationFactor: 2      Configs: segment.bytes=1073741824
      Topic: trips Partition: 0      Leader: 1002      Replicas: 1002,1001      Isr: 1002,1001
      Topic: trips Partition: 1      Leader: 1002      Replicas: 1001,1002      Isr: 1002,1001
(base) vutrungnghia@Lusheeta:~/kafka_2.13-2.6.0$ bin/kafka-topics.sh --describe --topic trips --bootstrap
p-server localhost:9092,localhost:9094
Topic: trips      PartitionCount: 2      ReplicationFactor: 2      Configs: segment.bytes=1073741824
      Topic: trips Partition: 0      Leader: 1002      Replicas: 1002,1001      Isr: 1002
      Topic: trips Partition: 1      Leader: 1002      Replicas: 1001,1002      Isr: 1002
(base) vutrungnghia@Lusheeta:~/kafka_2.13-2.6.0$ █
```

Read data from parquet file "trips.parquet" in hdfs

```
df = ss.read.parquet("hdfs://namenode:9000/trips/trips.parquet")
```

```
print(f"Number of records: {df.count()}")
df = df.sort('ArrivalTime')
```

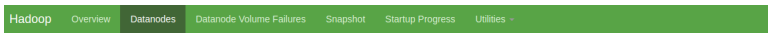
Number of records: 10000



- 1 Phát biểu bài toán
- 2 Kiến trúc hệ thống
- 3 Kết quả chạy chương trình
 - Kết quả chạy nhánh Batch Processing
 - Kết quả chạy nhánh Speed Processing
- 4 Đánh giá hệ thống
 - KAFKA
 - **HDFS**
 - SPARK
 - Hạn chế và hướng phát triển tiếp theo
- 5 Tài liệu tham khảo

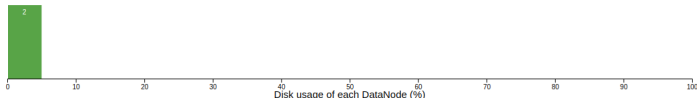


Đánh giá khả năng chịu lỗi của HDFS



Datanode Information

Datanode usage histogram



In operation

Node	Last contact	Admin State	Capacity	Used	Non DFS Used	Remaining	Blocks	Block pool used	Failed Volumes	Version
f4a27057f907:50010 (172.18.0.3:50010)	1	In Service	355.6 GB	1.97 MB	88.41 GB	249.06 GB	60	1.97 MB (0%)	0	2.7.4
3e4473579333:50010 (172.18.0.8:50010)	1	In Service	355.6 GB	1.97 MB	88.41 GB	249.06 GB	60	1.97 MB (0%)	0	2.7.4



Đánh giá khả năng chịu lỗi của HDFS

```

20/12/23 17:29:34 INFO FileScanRDD: Reading File path: hdfs://namenode:9000/trips/trips.parquet/part-00001-2t
range: 0-166672, partition values: [empty row]
20/12/23 17:29:52 WARN BlockReaderFactory: I/O error constructing remote block reader.
java.net.NoRouteToHostException: No route to host
    at sun.nio.ch.SocketChannelImpl.checkConnect(Native Method)
    at sun.nio.ch.SocketChannelImpl.finishConnect(SocketChannelImpl.java:716)
    at org.apache.hadoop.net.SocketIOWithTimeout.connect(SocketIOWithTimeout.java:206)
    at org.apache.hadoop.net.NetUtils.connect(NetUtils.java:531)
    at org.apache.hadoop.hdfs.DFSClient.newConnectedPeer(DFSClient.java:3436)
    at org.apache.hadoop.hdfs.BlockReaderFactory.nextTcpPeer(BlockReaderFactory.java:777)
    at org.apache.hadoop.hdfs.BlockReaderFactory.getRemoteBlockReaderFromTcp(BlockReaderFactory.java:694)
    at org.apache.hadoop.hdfs.BlockReaderFactory.build(BlockReaderFactory.java:355)
    at org.apache.hadoop.hdfs.DFSInputStream.blockSeekTo(DFSInputStream.java:673)
    at org.apache.hadoop.hdfs.DFSInputStream.readWithStrategy(DFSInputStream.java:882)
    at org.apache.hadoop.hdfs.DFSInputStream.read(DFSInputStream.java:934)
    at org.apache.hadoop.hdfs.DFSInputStream.read(DFSInputStream.java:735)
    at java.io.FilterInputStream.read(FilterInputStream.java:83)
    at org.apache.parquet.io.DelegatingSeekableInputStream.read(DelegatingSeekableInputStream.java:61)
    at org.apache.parquet.bytes.BytesUtils.readIntLittleEndian(BytesUtils.java:80)
    at org.apache.parquet.hadoop.ParquetFileReader.readFooter(ParquetFileReader.java:520)
    at org.apache.parquet.hadoop.ParquetFileReader.readFooter(ParquetFileReader.java:505)
    at org.apache.parquet.hadoop.ParquetFileReader.readFooter(ParquetFileReader.java:499)
    at org.apache.parquet.hadoop.ParquetFileReader.readFooter(ParquetFileReader.java:448)
    at

```



Đánh giá khả năng chịu lỗi của HDFS

```
20/12/23 17:31:00 WARN DFSCClient: Failed to connect to /172.18.0.3:50010 for block, add to deadNodes
java.net.NoRouteToHostException: No route to host
    at sun.nio.ch.SocketChannelImpl.checkConnect(Native Method)
    at sun.nio.ch.SocketChannelImpl.finishConnect(SocketChannelImpl.java:716)
    at org.apache.hadoop.net.SocketIOWithTimeout.connect(SocketIOWithTimeout.java:206)
    at org.apache.hadoop.net.NetUtils.connect(NetUtils.java:531)
    at org.apache.hadoop.hdfs.DFSCClient.newConnectedPeer(DFSCClient.java:3436)
    at org.apache.hadoop.hdfs.BlockReaderFactory.nextTcpPeer(BlockReaderFactory.java:777)
    at org.apache.hadoop.hdfs.BlockReaderFactory.getRemoteBlockReaderFromTcp(BlockReaderFactory.java:777)
    at org.apache.hadoop.hdfs.BlockReaderFactory.build(BlockReaderFactory.java:355)
    at org.apache.hadoop.hdfs.DFSInputStream.blockSeekTo(DFSInputStream.java:673)
    at org.apache.hadoop.hdfs.DFSInputStream.readWithStrategy(DFSInputStream.java:882)
    at org.apache.hadoop.hdfs.DFSInputStream.read(DFSInputStream.java:934)
    at org.apache.hadoop.hdfs.DFSInputStream.read(DFSInputStream.java:735)
    at java.io.FilterInputStream.read(FilterInputStream.java:83)
    at org.apache.parquet.io.DelegatingSeekableInputStream.read(DelegatingSeekableInputStream.java:80)
    at org.apache.parquet.bytes.BytesUtils.readIntLittleEndian(BytesUtils.java:80)
    at org.apache.parquet.hadoop.ParquetFileReader.readFooter(ParquetFileReader.java:520)
    at org.apache.parquet.hadoop.ParquetFileReader.readFooter(ParquetFileReader.java:505)
    at org.apache.parquet.hadoop.ParquetFileReader.readFooter(ParquetFileReader.java:499)
    at org.apache.parquet.hadoop.ParquetFileReader.readFooter(ParquetFileReader.java:448)
    at
```



Đánh giá khả năng chịu lỗi của HDFS

```

    at org.apache.spark.sql.execution.BufferedRowIterator.hasNext(BufferedRowIterator.java:43)
    at org.apache.spark.sql.execution.WholeStageCodegenExec$$anonfun$13$$anon$1.hasNext(WholeStageCode
    at scala.collection.Iterator$$anon$11.hasNext(Iterator.scala:409)
    at org.apache.spark.shuffle.sort.BypassMergeSortShuffleWriter.write(BypassMergeSortShuffleWriter.j
    at org.apache.spark.scheduler.ShuffleMapTask.runTask(ShuffleMapTask.scala:99)
    at org.apache.spark.scheduler.ShuffleMapTask.runTask(ShuffleMapTask.scala:55)
    at org.apache.spark.scheduler.Task.run(Task.scala:121)
    at org.apache.spark.executor.Executor$TaskRunner$$anonfun$10.apply(Executor.scala:403)
    at org.apache.spark.util.Utils$.tryWithSafeFinally(Utils.scala:1360)
    at org.apache.spark.executor.Executor$TaskRunner.run(Executor.scala:409)
    at java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPoolExecutor.java:1149)
    at java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoolExecutor.java:624)
    at java.lang.Thread.run(Thread.java:748)
20/12/23 17:31:00 INFO DFSClnt: Successfully connected to /172.18.0.8:50010 for BP-1621182885-172.18.0.7
20/12/23 17:31:00 INFO Executor: Finished task 0.0 in stage 4.0 (TID 5). 1716 bytes result sent to driver
20/12/23 17:31:18 INFO CoarseGrainedExecutorBackend: Got assigned task 7
20/12/23 17:31:18 INFO Executor: Running task 0.0 in stage 5.0 (TID 7)
20/12/23 17:31:18 INFO MapOutputTrackerWorker: Updating epoch to 2 and clearing cache
20/12/23 17:31:18 INFO TorrentBroadcast: Started reading broadcast variable 7
20/12/23 17:31:18 INFO MemoryStore: Block broadcast_7_piece0 stored as bytes in memory (estimated size 3.8
20/12/23 17:31:18 INFO TorrentBroadcast: Reading broadcast variable 7 took 13 ms
20/12/23 17:31:18 INFO MemoryStore: Block broadcast_7 stored as values in memory (estimated size 7.1 KB, f
20/12/23 17:31:18 INFO MapOutputTrackerWorker: Don't have map outputs for shuffle 1, fetching them
20/12/23 17:31:18 INFO MapOutputTrackerWorker: Doing the fetch; tracker endpoint = NettyRpcEndpointRef(spa
20/12/23 17:31:18 INFO MapOutputTrackerWorker: Got the output locations
20/12/23 17:31:18 INFO ShuffleBlockFetcherIterator: Getting 2 non-empty blocks including 1 local blocks an
20/12/23 17:31:18 INFO ShuffleBlockFetcherIterator: Started 1 remote fetches in 1 ms
20/12/23 17:31:18 INFO Executor: Finished task 0.0 in stage 5.0 (TID 7). 1782 bytes result sent to driver

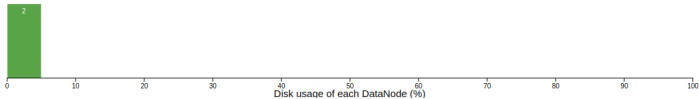
```



Đánh giá khả năng chịu lỗi của HDFS

Datanode Information

Datanode usage histogram



In operation

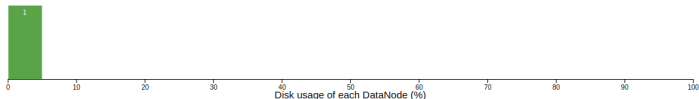
Node	Last contact	Admin State	Capacity	Used	Non DFS Used	Remaining	Blocks	Block pool used	Failed Volumes	Version
f4a27057f907:50010 (172.18.0.3:50010)	470	In Service	355.6 GB	2.38 MB	88.43 GB	249.04 GB	116	2.38 MB (0%)	0	2.7.4
3e4473579333:50010 (172.18.0.8:50010)	1	In Service	355.6 GB	2.98 MB	88.43 GB	249.03 GB	116	2.98 MB (0%)	0	2.7.4



Đánh giá khả năng chịu lỗi của HDFS

Datanode Information

Datanode usage histogram



In operation

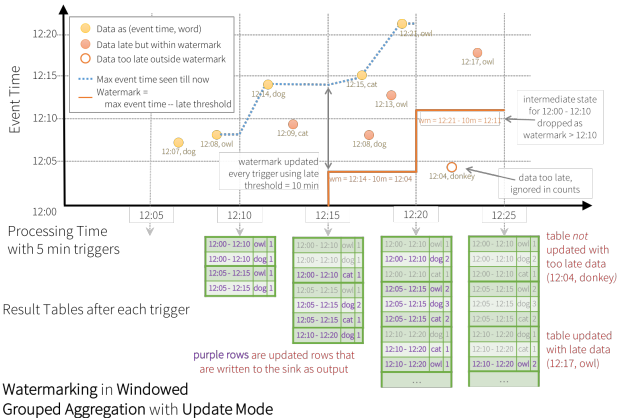
Node	Last contact	Admin State	Capacity	Used	Non DFS Used	Remaining	Blocks	Block pool used	Failed Volumes	Version
3e4473579333:50010 (172.18.0.8:50010)	2	In Service	355.6 GB	2.98 MB	88.44 GB	249.02 GB	116	2.98 MB (0%)	0	2.7.4
14a27057f907:50010 (172.18.0.3:50010)	Thu Dec 24 2020 00:29:21 GMT+0700 (Indochina Time)	Dead	-	-	-	-	-	-	-	-



- 1 Phát biểu bài toán
- 2 Kiến trúc hệ thống
- 3 Kết quả chạy chương trình
 - Kết quả chạy nhánh Batch Processing
 - Kết quả chạy nhánh Speed Processing
- 4 Đánh giá hệ thống
 - KAFKA
 - HDFS
 - **SPARK**
 - Hạn chế và hướng phát triển tiếp theo
- 5 Tài liệu tham khảo



Đánh giá khả năng xử lý gói tin đến sai thứ tự của Spark



(*) source: <https://spark.apache.org>



Đánh giá khả năng xử lý gói tin đến sai thứ tự của Spark

```
query = df.withWatermark("timestamp", "20 seconds") \
    .groupBy(window("timestamp", "10 seconds", "10 seconds")) \
    .count() \
```

- Để đưa các gói tin về đúng khoảng thời gian chính xác mà nó được gửi, Structured Spark sử dụng Window Grouped Aggregation.
- Đối với các dạng event đến quá trễ và trở nên vô giá trị, Watermarking bỏ qua các gói tin đến trễ so với khoảng thời gian cho trước.



Đánh giá khả năng xử lý gói tin đến sai thứ tự của Spark

```

Row(start=datetime.datetime(2020, 12, 22, 4, 22), end=datetime.datetime(2020, 12, 22, 4, 22, 10))_108
Row(start=datetime.datetime(2020, 12, 22, 4, 22, 10), end=datetime.datetime(2020, 12, 22, 4, 22, 20))_138
Row(start=datetime.datetime(2020, 12, 22, 4, 22, 20), end=datetime.datetime(2020, 12, 22, 4, 22, 30))_51
Row(start=datetime.datetime(2020, 12, 22, 4, 22, 30), end=datetime.datetime(2020, 12, 22, 4, 22, 40))_106
Row(start=datetime.datetime(2020, 12, 22, 4, 22, 40), end=datetime.datetime(2020, 12, 22, 4, 22, 50))_135
Row(start=datetime.datetime(2020, 12, 22, 4, 22, 50), end=datetime.datetime(2020, 12, 22, 4, 23))_56
Row(start=datetime.datetime(2020, 12, 22, 4, 23), end=datetime.datetime(2020, 12, 22, 4, 23, 10))_103
Row(start=datetime.datetime(2020, 12, 22, 4, 23, 10), end=datetime.datetime(2020, 12, 22, 4, 23, 20))_136
Row(start=datetime.datetime(2020, 12, 22, 4, 23, 20), end=datetime.datetime(2020, 12, 22, 4, 23, 30))_58
Row(start=datetime.datetime(2020, 12, 22, 4, 23, 30), end=datetime.datetime(2020, 12, 22, 4, 23, 40))_104
Row(start=datetime.datetime(2020, 12, 22, 4, 23, 40), end=datetime.datetime(2020, 12, 22, 4, 23, 50))_134
Row(start=datetime.datetime(2020, 12, 22, 4, 23, 50), end=datetime.datetime(2020, 12, 22, 4, 24))_59

```

Hình: Bài toán đến số lượng bản ghi trong 10s gần nhất

- Sử dụng KafkaProducer để giả lập luồng dữ liệu có các events đến trễ.
- Cứ sau 20s, một nửa các gói tin được gửi đi từ Producer trong 10s tiếp theo được đánh dấu là trễ 10s, nửa còn lại gửi đi theo đúng thời gian thực.



- 1 Phát biểu bài toán
- 2 Kiến trúc hệ thống
- 3 Kết quả chạy chương trình
 - Kết quả chạy nhánh Batch Processing
 - Kết quả chạy nhánh Speed Processing
- 4 Đánh giá hệ thống
 - KAFKA
 - HDFS
 - SPARK
 - Hạn chế và hướng phát triển tiếp theo
- 5 Tài liệu tham khảo



Hạn chế

- Các ứng dụng Graphite, Grafana và JavaApp không cần nhiều tài nguyên nên có thể đóng gói để chạy trên 1 container
- Ngoại trừ Interactive Dashboard Notebook, các notebook còn lại có thể chuyển sang ngôn ngữ Scala hay Java và đóng gói thành 1 ứng dụng, từ đó giúp tăng tốc độ tính toán cũng như phù hợp hơn với môi trường production



Hướng phát triển tiếp theo

- Xử lý các hạn chế
- Xử lý nhiều luồng dữ liệu với các bài toán khác nhau thay vì chỉ 1 luồng dữ liệu cho 1 bài toán
- Xây dựng các mô hình học máy mạnh mẽ hơn cho các bài toán phức tạp
- Triển khai trên cụm máy tính thật



Tài liệu tham khảo

- <https://github.com/haiphucnguyen/BigDataDemo>
- <http://www.diva-portal.org/smash/get/diva2:897808/FULLTEXT01.pdf>
- [https://blog.softwaremill.com/7-mistakes-when-using-ap](https://blog.softwaremill.com/7-mistakes-when-using-apache-spark-2017-01-10/)
- <https://spark.apache.org/docs/1.5.2/sql-programming-guide.html>
- <https://spark.apache.org/docs/latest/structured-streaming-programming-guide.html>
- [https://medium.com/dev-genius/an-in-depth-look-at-zook](https://medium.com/dev-genius/an-in-depth-look-at-zookeeper-2017-01-10/)
- <https://spark.apache.org/docs/latest/streaming-programming-guide.html>

