# Elastic Map Reduce

Shadi Khalifa
Database Systems Laboratory (DSL)
khalifa@cs.queensu.ca

# The Amazon Web Services Universe

# Infrastructure Services



Infrastructure Services

EC2
http://aws.amazon.com/ec2/

VPC
http://aws.amazon.com/vpc/

S3
http://aws.amazon.com/s3/

EBS
http://aws.amazon.com/ebs/

**S3**
Scalable Storage in the Cloud

# Amazon Simple Storage Service (S3)

- Amazon S3 provides a simple web services interface that can be used to store and retrieve any amount of data, at any time, from anywhere on the web.

- Write, read, and delete objects containing from 1 byte to 5 terabytes of data each. The number of objects you can store is unlimited.

- Each object is stored in a bucket and retrieved via a unique, developer-assigned key.
  - A bucket can be stored in one of several Regions.
  - You can choose a Region to optimize for latency, minimize costs, or address regulatory requirements.
  - Objects stored in a Region never leave the Region unless you transfer them out.

- Authentication mechanisms are provided to ensure that data is kept secure from unauthorized access.
  - Objects can be made private or public, and rights can be granted to specific users.

- S3 charges based on per GB-month AND per I/O requests AND per data modification requests.

**All Buckets**

| Name |
| --- |
| 🔍  aws-test-grep |
| 🔍  aws-test-HE |
| 🔍  cloud-final |
| 🔍  cloud-final-results |
| 🔍  cloudfinal-out1 |
| 🔍  dasbucket1 |
| 🔍  data432 |
| 🔍  eh-test |
| 🔍  EH_test |
| 🔍  final-report1 |
| 🔍  final-reprot1 |
| 🔍  finalreport |
| 🔍  hanaEsi_final |
| 🔍  mmyawsbucket |
| 🔍  mmyawsbucket1 |
| 🔍  ms499-mapredcapri-test |
| 🔍  ms499-wordcount-test |

✕

💬 **S3 Console**

Let us know what you think, with the Feedback button at the bottom of the page.

Feedback

# Platform Services

EMR

http://aws.amazon.com/elasticmapreduce/

RDS

http://aws.amazon.com/rds/

Platform Services

Not Only SQL

DynamoDB

http://aws.amazon.com/dynamodb/

Beanstalk

http://aws.amazon.com/elasticbeanstalk/

**Elastic MapReduce**
Managed Hadoop Framework

# Amazon Elastic MapReduce (EMR)

- Amazon EMR is a web service that makes it easy to quickly and cost-effectively process vast amounts of data using Hadoop.

- Amazon EMR distribute the data and processing across a resizable cluster of Amazon EC2 instances.

- With Amazon EMR you can launch a persistent cluster that stays up indefinitely or a temporary cluster that terminates after the analysis is complete.

- Amazon EMR supports a variety of Amazon EC2 instance types and Amazon EC2 pricing options (On-Demand, Reserved, and Spot).

- When launching an Amazon EMR cluster (also called a "job flow"), you choose how many and what type of Amazon EC2 Instances to provision.

- The Amazon EMR price is in addition to the Amazon EC2 price.

- Amazon EMR is used in a variety of applications, including log analysis, web indexing, data warehousing, machine learning, financial analysis, scientific simulation, and bioinformatics.

# What is Hadoop?

# Example: Word Count

- **Objective**: Count the number of occurrences of each word in the provided input files.

- How it works
  - In the Map phase the text is tokenized into words then we form a key value pair with these words where the key being the word itself and value is set to '1'.
  - In the reduce phase the keys are grouped together and the values for the same key are added.

**Input Files**

Apple Orange Mango
Orange Grapes Plum

Apple Plum Mango
Apple Apple Plum

# Example 2: Joins with MapReduce

**UserDetails.txt**

123 456, Jim

456 123, Tom

789 123, Harry

789 456, Richa

**DeliveryDetails.txt**

123 456, Delivered

456 123, Pending

789 123, Failed

789 456, Resend

**Expected Output**

Jim, Delivered

Tom, Pending

Harry, Failed

Richa, Resend

- We have 2 input files as follows:
  - **UserDetails.txt** : Every record is of the format 'mobile number , consumer name'.
  - **DeliveryDetails.txt**: Every record is of the format 'mobile number, delivery status'.

- **objective**: Associate the customer name with the delivery status.

# Formulating as a MapReduce

Since we have two input files, we need 2 Map functions (one for each input file)

**UserDetails.txt**
<123 456, CD~Jim>
<456 123, CD~Tom>
<789 123, CD~Harry>
<789 456, CD~Richa>

**DeliveryDetails.txt**
<123 456, DR~Delivered>
<456 123, DR~Pending>
<789 123, DR~Failed>
<789 456, DR~Resend>

**Reduce Input**
<123 456, CD~Jim>
<123 456, DR~Delivered>

<456 123, CD~Tom>
<456 123, DR~Pending>

- value: 'DR' for the

**Reduce Output**
<Jim, Delivered>

<Tom, Pending>

- on the reducer, every            s (for simplicity) one with prefix 'CD' and other 'DR'.

- From CD get the customer name corresponding to the cell number (input key) and from DR get the status.

- The output Key values from the reducer would be as follows
  - Key : Customer Name
  - Value : Status Message

# Tools you will need

- Eclipse IDE for Java EE Developers
  - http://www.eclipse.org/downloads/packages/eclipse-ide-java-ee-developers/keplersr1

- hadoop-core-1.2.1.jar
  - http://mvnrepository.com/artifact/org.apache.hadoop/hadoop-core/1.2.1

- Amazon web Services Account
  - https://console.aws.amazon.com/elasticmapreduce/home?region=us-east-1

# Create a new Java Project

# Add Hadoop jar to the project

- Create a 'lib' folder in the project.

- Copy and paste the Hadoop jar into the 'lib' folder.

- Add the jar to the project build path.

# Map classes: 1) UserFileMapper

```java
public class UserFileMapper extends MapReduceBase implements Mapper<LongWritable, Text, Text, Text>
{
    //varia                          er Details
    private                      tomerName,fileTag="CD~";

    /* map                          nsu            me
pairs
        Key(
        Valu                     r to indicate the source of              e
customer det                    Name
    */
    public void map(LongWritable key, Text value, OutputCollector<Text, Text> output,
Reporter reporter) throws IOException
    {
        //taking one line/record at a time and parsing them into key value pairs
        String line = value.toString();
        String splitarray[] = line.split(",");
        cellNumber = splitarray[0].trim();
        customerName = splitarray[1].trim();

      //sending the key value pair out of mapper
        output.collect(new Text(cellNumber), new Text(fileTag+customerName));
    }
}
```

**UserDetails.txt**

123 456, Jim
456 123, Tom
789 123, Harry
789 456, Richa

**UserDetails.txt**

<123 456, CD~Jim>
<456 123, CD~Tom>
<789 123, CD~Harry>
<789 456, CD~Richa>

# Map classes: 2) DeliveryFileMapper

```java
public class DeliveryFileMapper extends MapReduceBase implements Mapper<LongWritable,
Text, Text. Text>
{
    //                        ry report
    pr                        iveryCode,fileTag="DR~";

    /*                        iv                    e
    Ke
    Va                            o indicate the source of
deliv                         de*/

    public void map(LongWritable key, Text value, OutputCollector<Text, Text> output,
Reporter reporter) throws IOException
    {
        //taking one line/record at a time and parsing them into key value pairs
        String line = value.toString();
        String splitarray[] = line.split(",");
        cellNumber = splitarray[0].trim();
        deliveryCode = splitarray[1].trim();

        //sending the key value pair out of mapper
        output.collect(new Text(cellNumber), new Text(fileTag+deliveryCode));
    }
}
```

DeliveryDetails.txt
123 456, Delivered
456 123, Pending
789 123, Failed
789 456, Resend

DeliveryDetails.txt
<123 456, DR~Delivered>
<456 123, DR~Pending>
<789 123, DR~Failed>
<789 456, DR~Resend>

# Reduce Class

```java
public                             apReduceBase implements Reducer<Text, Text, Text, Text> {

                                    n process
                                  e.deliveryReport;


                            y, Iterator<Text> values,
output                              IOException
    {
        while (values.hasNext())
        {
            String currValue = values.next().toString();
            String valueSplitted[] = currValue.split("~");
            /*identifying the record source that corresponds to a cell number
            and parses the values accordingly*/
            if(valueSplitted[0].equals("CD"))
            {
              customerName=valueSplitted[1].trim();
            }
            else if(valueSplitted[0].equals("DR"))
            {
              deliveryReport = valueSplitted[1].trim();
            }
        }
        output.collect(new Text(customerName), new Text(deliveryReport));
    }
}
```

# Driver Class

```java
public class SmsDriver extends Configured implements Tool
{
        public int run(String[] args) throws Exception {

                //get the configuration parameters and assigns a job name
                JobConf conf = new JobConf(getConf(), SmsDriver.class);
                conf.setJobName("SMS Reports");

                //setting key value types for mapper and reducer outputs
                conf.setOutputKeyClass(Text.class);
                conf.setOutputValueClass(Text.class);

                //specifying the custom reducer class
                conf.setReducerClass(SmsReducer.class);

// If only one Mapper exists
   //    conf.setMapperClass(Mapper.class);
   //    FileInputFormat.addInputPath(conf, newPath(args[0]));
//Specifying the input directories(@ runtime) and Mappers independently for inputs from multiple sources
    MultipleInputs.addInputPath(conf, new Path(args[0]), TextInputFormat.class, UserFileMapper.class);
    MultipleInputs.addInputPath(conf, new Path(args[1]), TextInputFormat.class, DeliveryFileMapper.class);

                //Specifying the output directory @ runtime
                FileOutputFormat.setOutputPath(conf, new Path(args[2]));

                JobClient.runJob(conf);
                return 0;
        }

        public static void main(String[] args) throws Exception {
                int res = ToolRunner.run(new Configuration(), new SmsDriver(), args);
                System.exit(res);
        }
}
```

# Export Project to Jar

# Uploading Project Jar and Data to S3

- https://console.aws.amazon.com/s3/home?region=us-east-1



| Name | Storage Class | Size | Last Modified |
|---|---|---|---|
| DeliveryDetails.txt | Standard | 54 bytes | Sun Sep 29 22:16:50 GMT-400 2013 |
| DeliveryStatusCodes.txt | Standard | 54 bytes | Sun Sep 29 22:16:50 GMT-400 2013 |
| JoinHadoop.jar | Standard | 3.7 MB | Mon Sep 30 00:32:41 GMT-400 2013 |
| UserDetails.txt | Standard | 58 bytes | Sun Sep 29 22:16:50 GMT-400 2013 |
| log | -- | -- | -- |
| out | -- | -- | -- |

# Hadoop Execution Overview

# Using Amazon EMR

https://console.aws.amazon.com/elasticmapreduce/home?region=us-east-1

# Using Amazon EMR

Make sure that the output folder (3rd argument in our example) does NOT exists on S3 (MapReduce will create it)



**Create a New Job Flow**                                                        Cancel ☒

DEFINE JOB FLOW    **SPECIFY PARAMETERS**    CONFIGURE EC2 INSTANCES    ADVANCED OPTIONS    BOOTSTRAP ACTIONS    REVIEW

Specify the location in Amazon S3 of your JAR. Hadoop executes the JAR. You can specify its main class in its manifest. If you don't you must specify a class name as the first argument of the JAR.

**JAR Location\*:** s3n://shadihadoop/JoinHadoop.jar

**JAR Arguments\*:** s3n://shadihadoop/UserDetails.txt
s3n://shadihadoop/DeliveryDetails.txt
s3n://shadihadoop/out

‹ Back                          Continue  ▶                          * Required field

# Using Amazon EMR

Set the number and type of the EC2 instances used to process your application

**Create a New Job Flow**       Cancel ✕

DEFINE JOB FLOW    SPECIFY PARAMETERS    **CONFIGURE EC2 INSTANCES**    ADVANCED OPTIONS    BOOTSTRAP ACTIONS    REVIEW

Specify the master, core and task nodes to run your job flow. For more than 20 instances, complete the limit request form.

**Master Instance Group:** This EC2 instance assigns Hadoop tasks to core and task nodes and monitors their status.

Instance Type: | Small (m1.small) ▾ | ☐ Request Spot Instance

**Core Instance Group:** These EC2 instances run Hadoop tasks and store data using the Hadoop Distributed File System (HDFS). Recommended for capacity needed for the life of your job flow.

Instance Count: | 2

Instance Type: | Small (m1.small) ▾ | ☐ Request Spot Instances

**Task Instance Group (Optional):** These EC2 instances run Hadoop tasks, but do not persist data. Recommended for capacity needed on a temporary basis.

Instance Count: | 0

Instance Type: | Small (m1.small) ▾ | ☐ Request Spot Instances

‹ Back      **Continue** ▶      * Required field

# Using Amazon EMR

Remember to ENABLE the Debugging and provide the log Path

# Performance

Using Amazon cloudwatch (Monitoring Tab) you can  check the performance of the

# Debugging the Job

Debug

For detailed information on the MapReduce progress, click on the syslog link

**Debug a Job Flow**                                                    Close ☒

**Job Flow:** My Job Flow (j-1YPC558KQ9PN6)

View logs for steps, Hadoop jobs, tasks, and task attempts.

**Steps → Jobs → Tasks → Task Attempts**                    ⟳ Refresh List

| Step | Name | State | Start Time | Log Files | Actions |
|------|------|-------|-----------|-----------|---------|
| 1 | Setup Hadoop Debugging | 🔵 COMPLETED | 2013-10-17 17:12 EDT | controller \| stderr \| stdout \| syslog | View Jobs |
| 2 | Streaming Job | 🔵 COMPLETED | 2013-10-17 17:12 EDT | controller \| stderr \| stdout \| syslog | View Jobs |

# Accessing the Results on S3

Elastic Map Reduce

Questions?