

Christopher M. Judd













WWW.BITSTRIPS.COM

INTRODUCTION



What Is Apache Hadoop?

The Apache[™] Hadoop® project develops open-source software for reliable, scalable, distributed computing.

The Apache Hadoop software library is a framework that allows for the distributed processing of large data sets across clusters of computers using simple programming models. It is designed to scale up from single servers to thousands of machines, each offering local computation and storage. Rather than rely on hardware to deliver high-availability, the library itself is designed to detect and handle failures at the application layer, so delivering a highly-available service on top of a cluster of computers, each of which may be prone to failures.

The project includes these modules:

- Hadoop Common: The common utilities that support the other Hadoop modules.
- Hadoop Distributed File System (HDFS™): A distributed file system that provides high-throughput access to application data.
- Hadoop YARN: A framework for job scheduling and cluster resource management.
- Hadoop MapReduce: A YARN-based system for parallel processing of large data sets.

http://hadoop.apache.org/



















Hadoop Approach

- scale-out
- share nothing
- expect failure
- smart software, dumb hardware
- move processing, not data
- build applications, not infrastructure



Hadoop

What is Hadoop good for?

Don't use Hadoop - your data isn't that big

<100 mb - Excel</p>
100 mb > 10 gb - Add memory and use Pandas
100 gb > 1 TB - Buy big hard drive and use Postgres
> 5 TB - life sucks consider Hadoop

http://www.chrisstucchio.com/blog/2013/hadoop_hatred.html

if your data fits in RAM

it is not Big Data

Adoop is an evolving project

Adoop is an evolving project

old api

new api

org.apache.hadoop.mapred

org.apache.hadoop.mapreduce

Adoop is an evolving project

MapReduce 1

Classic MapReduce

MapReduce 2

YARN

SETUP





Hadoop Tutorial user/fun4all hduser/hduser /opt/data

Configure SSH

		Н	adoop Tu	itorial Pr	actice – Ne	etwork			
			\bigcirc						
General	System	Display	Storage	Audio	Network	Ports	Shared Folders	5	
		Adapter 1	Adap	ter 2	Adapter 3	Ada	apter 4		
🗹 En	able Netv	work Adap	ter						
	At	tached to:	NAT		÷				
Name:							A T		
	▽.	Advanced							
	pter Type	Intel	Intel PRO/1000 MT Desktop (82540EM)						
	Promiscu	ous Mode	Deny					A T	
	MAC Address:				080027BA08C1				
			🗹 Cab	le Conn	ected				
Port Forwarding									
?							Cancel	ОК	

attach to NAT

Configure SSH

800	user@user-VirtualBox: ~
user@use eth0	er-VirtualBox:~\$ ifconfig Link encap:Ethernet HWaddr 08:00:27:ba:08:c1 inet addr 10.0.2.15 Bcast:10.0.2.255 Mask:255.255.255.0 inet6 addr: fe80::a00:27ff:feba:8c1/64 Scope:Link UP BROADCAST RUNNING MULTICAST MTU:1500 Metric:1 RX packets:186 errors:0 dropped:0 overruns:0 frame:0 TX packets:210 errors:0 dropped:0 overruns:0 carrier:0 collisions:0 txqueuelen:1000 RX bytes:50163 (50.1 KB) TX bytes:26221 (26.2 KB)
lo	Link encap:Local Loopback inet addr:127.0.0.1 Mask:255.0.0.0 inet6 addr: ::1/128 Scope:Host UP LOOPBACK RUNNING MTU:65536 Metric:1 RX packets:102 errors:0 dropped:0 overruns:0 frame:0 TX packets:102 errors:0 dropped:0 overruns:0 carrier:0 collisions:0 txqueuelen:0 RX bytes:16666 (16.6 KB) TX bytes:16666 (16.6 KB)
user@use	er-VirtualBox:~\$

Configure SSH

Name Protoco		Host IP	Host Port	Guest IP	Guest Port		
ssh	ТСР	127.0.0.1	3022	10.0.2.15	22		
		Cable Conn	ected				
		Port Forwarding					

add port forwarding rule

SSH'ing

\$ ssh -p 3022 user@localhost

user@127.0.0.1's password: Welcome to Ubuntu 12.04.3 LTS (GNU/Linux 3.8.0-29-generic x86_64)

* Documentation: https://help.ubuntu.com/

Last login: Wed Jun 25 22:53:10 2014 from 10.0.2.2
user@user-VirtualBox:~\$

HADDDP



http://www.mapr.com/



Elastic MapReduce





Add Hadoop User and Group



Install Hadoop

\$ sudo mkdir -p /opt/hadoop

\$ sudo tar vxzf /opt/data/hadoop-2.2.0.tar.gz -C /opt/hadoop

\$ sudo chown -R hduser:hadoop /opt/hadoop/hadoop-2.2.0

\$ vim .bashrc

other stuff
java variables
export JAVA_HOME=/usr/lib/jvm/java-7-openjdk-amd64

hadoop variables export HADOOP_HOME=/opt/hadoop/hadoop-2.2.0 export PATH=\$PATH:\$HADOOP_HOME/bin export PATH=\$PATH:\$HADOOP_HOME/sbin export HADOOP_MAPRED_HOME=\$HADOOP_HOME export HADOOP_COMMON_HOME=\$HADOOP_HOME export HADOOP_HDFS_HOME=\$HADOOP_HOME export YARN HOME=\$HADOOP_HOME

\$ source .bashrc
\$ hadoop version

Run Hadoop Job

\$ hadoop jar \$HADOOP_HOME/share/hadoop/mapreduce/hadoop-mapreduce-examples-2.2.0.jar pi 4 1000

Quasi-Monte Carlo method

From Wikipedia, the free encyclopedia

In numerical analysis, **quasi-Monte Carlo method** is a method for numerical integration and solving some other problems using lowdiscrepancy sequences (also called quasi-random sequences or subrandom sequences). This is in contrast to the regular Monte Carlo method or Monte Carlo integration, which are based on sequences of pseudorandom numbers.

Monte Carlo and quasi-Monte Carlo methods are stated in a similar way. The problem is to approximate the integral of a function f as the average of the function evaluated at a set of points $x_1, ..., x_N$:

$$\int_{[0,1]^s} f(u) \, \mathrm{d}u \approx \frac{1}{N} \sum_{i=1}^N f(x_i).$$

Since we are integrating over the *s*-dimensional unit cube, each x_i is a vector of *s* elements. The difference between quasi-Monte Carlo and Monte Carlo is the way the x_i are chosen. Quasi-Monte Carlo uses a low-discrepancy sequence such as the Halton sequence, the Sobol sequence, or the Faure sequence, whereas Monte Carlo uses a pseudorandom



[Pseudorandom sequence]

[Low-discrepancy sequence (Sobol sequence)]

256 points from a pseudorandom number source, Halton sequence, and Sobol sequence (red=1,...,10, blue=11,...,100, green=101,...,256). Points from Sobol sequence are more evenly distributed.

sequence. The advantage of using low-discrepancy sequences is a faster rate of convergence. Quasi-Monte Carlo has a rate of convergence close to O(1/N), whereas the rate for the Monte Carlo method is $O(N^{-0.5})$.^[1]

The Quasi-Monte Carlo method recently became popular in the area of mathematical finance or computational finance.^[1] In these areas, highdimensional numerical integrals, where the integral should be evaluated within a threshold ε , occur frequently. Hence, the Monte Carlo method and the quasi-Monte Carlo method are beneficial in these situations.

\$ hadoop jar \$HADOOP_HOME/share/hadoop/mapreduce/hadoop-mapreduce-examples-2.2.0.jar

aggregatewordcount: An Aggregate based map/reduce program that counts the words in the input files.

aggregatewordhist: An Aggregate based map/reduce program that computes the histogram of the words in the input files.

bbp: A map/reduce program that uses Bailey-Borwein-Plouffe to compute exact digits of Pi.

dbcount: An example job that count the pageview counts from a database.

distbbp: A map/reduce program that uses a BBP-type formula to compute exact bits of Pi.

grep: A map/reduce program that counts the matches of a regex in the input.

join: A job that effects a join over sorted, equally partitioned datasets

multifilewc: A job that counts words from several files.

pentomino: A map/reduce tile laying program to find solutions to pentomino problems.

pi: A map/reduce program that estimates Pi using a quasi-Monte Carlo method.

randomtextwriter: A map/reduce program that writes 10GB of random textual data per node.

randomwriter: A map/reduce program that writes I0GB of random data per node.

secondarysort: An example defining a secondary sort to the reduce.

sort: A map/reduce program that sorts the data written by the random writer.

sudoku: A sudoku solver.

teragen: Generate data for the terasort

terasort: Run the terasort

teravalidate: Checking results of terasort

wordcount: A map/reduce program that counts the words in the input files.

wordmean: A map/reduce program that counts the average length of the words in the input files.

wordmedian: A map/reduce program that counts the median length of the words in the input files.

wordstandarddeviation: A map/reduce program that counts the standard deviation of the length of the words in the input files.

\$ hadoop jar \$HADOOP_HOME/share/hadoop/mapreduce/hadoop-mapreduce-examples-2.2.0.jar pi
Usage: org.apache.hadoop.examples.QuasiMonteCarlo <nMaps> <nSamples>

Lab I

- I. Create Hadoop user and group
- 2. Install Hadoop
- 3. Run example Hadoop job such as pi

- Local Standalone mode
- Pseudo-distributed mode
- Fully distributed mode

HADDDP Pseudd-Distributed




O'REILLY" YAROOL HERE THE WINY

Configure YARN

\$ sudo vim etc/hadoop/yarn-site.xml

```
<?xml version="1.0"?>
<configuration>
<property>
        <name>yarn.nodemanager.aux-services</name>
        <value>mapreduce_shuffle</value>
        </property>
        <property>
        <name>yarn.nodemanager.aux-services.mapreduce.shuffle.class</name>
        <value>org.apache.hadoop.mapred.ShuffleHandler</value>
        </property>
</configuration>
```

Configure Map Reduce

\$ sudo mv etc/hadoop/mapred-site.xml.template etc/hadoop/mapred-site.xml
\$ sudo vim etc/hadoop/mapred-site.xml

<configuration> <property> <name>mapreduce.framework.name</name> <value>yarn</value> </property> </configuration>

Configure passwordless login

\$ ssh-keygen -t rsa -P ''
\$ cat ~/.ssh/id_rsa.pub >> ~/.ssh/authorized_keys
\$ ssh localhost
\$ exit

Configure JAVA_HOME

\$ vim etc/hadoop/hadoop-env.sh

other stuff

export JAVA_HOME=/usr/lib/jvm/java-7-openjdk-amd64

more stuff

Start YARN

\$ start-yarn.sh



Run Hadoop Job

\$ hadoop jar \$HADOOP_HOME/share/hadoop/mapreduce/hadoop-mapreduce-examples-2.2.0.jar pi 4 1000

2/noc ×	27
1:8042/node	☆ 🔎 ≡
	Logged in as: dr.who
	NodeManager information
Total Vmem allocated for Containers	16.80 GB
Vmem enforcement enabled	true
Total Pmem allocated for Container	8 GB
Pmem enforcement enabled	true
NodeHealthyStatus	true
LastNodeHealthTime NodeHealthReport	Wed Jan 08 00:23:37 EST 2014
Node Manager Version:	2.2.0 from 1529768 by hortonmu source checksum 6afffa66f656213479c75e45dcfd6e0 on 2013- 10-07T06:34Z
Hadoop Version:	2.2.0 from 1529768 by hortonmu source checksum 79e53ce7994d1628b240f09af91e1af4 on 2013-10-07T06:28Z
	2/noc × 1:8042/node Total Vmem allocated for Containers Vmem enforcement enabled Total Pmem allocated for Container Pmem enforcement enabled NodeHealthyStatus LastNodeHealthTime NodeHealthReport Node Manager Version: Hadoop Version:

About Apache Hadoop

http://localhost:8042/node

Lab 2

- I. Configure YARN
- 2. Configure Map Reduce
- 3. Configure passwordless login
- 4. Configure JAVA_HOME
- 5. Start YARN
- 6. Run pi job

HDF5



Writing Data





ORELLY YAROOL HERE THE WIND

Reading Data





ORELLY YAROOL HIM THE WIN

Configure HDFS

- \$ sudo mkdir -p /opt/hdfs/namenode
- \$ sudo mkdir -p /opt/hdfs/datanode
- \$ sudo chmod -R 777 /opt/hdfs
- \$ sudo chown -R hduser:hadoop /opt/hdfs
- \$ cd /opt/hadoop/hadoop-2.2.0
- \$ sudo vim etc/hadoop/hdfs-site.xml

```
<?xml version="1.0" encoding="UTF-8"?>
<?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
```

```
<configuration>
<property>
<name>dfs.replication</name>
<value>1</value>
</property>
<name>dfs.namenode.name.dir</name>
<value>file:/opt/hdfs/namenode</value>
</property>
<name>dfs.datanode.data.dir</name>
<value>file:/opt/hdfs/datanode</value>
</property>
<configuration>
```

Format HDFS

\$ hdfs namenode -format

Configure Core

\$ sudo vim etc/hadoop/core-site.xml

<property>
<property>
<property>
<property>
<property>.default.name</name>
<property>
</property>
</configuration>

Start HDFS

\$ start-dfs.sh



Use HDFS commands

\$ hdfs dfs -ls /
\$ hdfs dfs -mkdir /books
\$ hdfs dfs -ls /
\$ hdfs dfs -ls /
\$ hdfs dfs -ls /books
\$ hdfs dfs -copyFromLocal /opt/data/moby_dick.txt /books
\$ hdfs dfs -cat /books/moby_dick.txt

- appendToFile
- cat
- chgrp
- chmod
- chown
- copyFromLocal
- copyToLocal
- count
- ср
- du
- get •

- Is
- lsr
- mkdir
- moveFromLocal
- moveToLocal
- mv
- put
- rm
- rmr

tail

• stat

- test
- text
- touchz

http://hadoop.apache.org/docs/r2.2.0/hadoop-project-dist/hadoop-common/FileSystemShell.html

Hadoop NameNode localho ×

C 192.168.56.101:50070/dfshealth.jsp

NameNode 'localhost:9000' (active)

Started:	Tue Jan 07 13:09:59 EST 2014
Version:	2.2.0, 1529768
Compiled:	2013-10-07T06:28Z by hortonmu from branch-2.2.0
Cluster ID:	CID-4a58fa87-92c5-4f73-9a49-3e6fd7d87a69
Block Pool ID:	BP-1094298621-127.0.1.1-1389117580724

Browse the filesystem

NameNode Logs

00

Cluster Summary

Security is OFF

34 files and directories, 16 blocks = 50 total.

Heap Memory used 33.51 MB is 72% of Committed Heap Memory 46.04 MB. Max Heap Memory is 966.69 MB. Non Heap Memory used 28.40 MB is 95% of Committed Non Heap Memory 29.63 MB. Max Non Heap Memory is 214 MB.

, ,			-		
Configured Capacity	:	5.78 GB			
DFS Used	:	2.80 MB			
Non DFS Used	:	3.36 GB			
DFS Remaining	:	2.42 GB			
DFS Used%	:	0.05%			
DFS Remaining%	:	41.79%			
Block Pool Used	:	2.80 MB			
Block Pool Used%	:	0.05%			
DataNodes usages	:	Min %	Median %	Max %	stdev %
		0.05%	0.05%	0.05%	0.00%
Live Nodes	:	1 (Decommissioned: 0)			
Dead Nodes	:	0 (Decommissioned: 0)			
Decommissioning Nodes	:	0			
Number of Under-Replicated Blocks	:	0			
	Configured Capacity DFS Used Non DFS Used DFS Remaining DFS Used% DFS Remaining% Block Pool Used Block Pool Used Block Pool Used% DataNodes usages Live Nodes Dead Nodes Decommissioning Nodes Number of Under-Replicated Blocks	Configured Capacity : DFS Used : Non DFS Used : DFS Remaining : DFS Remaining% : DFS Remaining% : Block Pool Used : Block Pool Used% : DataNodes usages : Live Nodes : Dead Nodes : Number of Under-Replicated Blocks :	Configured Capacity:5.78 GBDFS Used:2.80 MBNon DFS Used:3.36 GBDFS Remaining:2.42 GBDFS Used%:0.05%DFS Remaining%:41.79%Block Pool Used:2.80 MBBlock Pool Used%:0.05%DataNodes usages:Min %Live Nodes:1 (Decommissioned: 0)Dead Nodes:0 (Decommissioned: 0)Decommissioning Nodes:0	Configured Capacity:5.78 GBDFS Used:2.80 MBNon DFS Used:3.36 GBDFS Remaining:2.42 GBDFS Used%:0.05%DFS Remaining%:41.79%Block Pool Used:2.80 MBBlock Pool Used%:0.05%DataNodes usages:Min %Live Nodes:1 (Decommissioned: 0)Dead Nodes:0 (Decommissioned: 0)Decommissioning Nodes:0	Configured Capacity:5.78 GBDFS Used:2.80 MBNon DFS Used:3.36 GBDFS Remaining:2.42 GBDFS Remaining%:0.05%DFS Remaining%:41.79%Block Pool Used:0.05%DataNodes usages:Min %Median %Max %Min %0.05%Live Nodes:1 (Decommissioned: 0)Dead Nodes:0Number of Under-Replicated Blocks:0

NameNode Journal Status:

Current transaction ID: 384

Journal Manager	State	
FileJournalManager(root=/opt/hdfs/namenode)	EditLogFileOutputStream(/opt/hdfs/namenode/current/edits_inprogress_000000000000000384)	

http://localhost:50070/dfshealth.jsp

H.

☆ 🎤 🔳

Lab 3

- I. Configure HDFS
- 2. Format HDFS
- 3. Configure Core
- 4. Start HDFS
- 5. Experiment HDFS commands (ls, mkdir, copyFromLocal, cat)

COMBINE HADOOP & HDFS

Run Hadoop Job

\$ hadoop jar share/hadoop/mapreduce/hadoop-mapreduce-examples-2.2.0.jar wordcount /books out
\$ hdfs dfs -ls out

\$ hdfs dfs -cat out/_SUCCESS

\$ hdfs dfs -cat out/part-r-00000

young-armed	1
young; 2	
younger 2	
youngest 1	
youngish 1	
your251	
your@login	1
yours 5	
yours? 1	
yourselbs	1
yourself 14	
yourself,	5
yourself,"	1
yourself.'	1
yourself;	4
yourself?	1
yourselves	1
yourselves!	3
yourselves,	1
yourselves,"	1
yourselves;	1
youth 5	
youth, 2	
youth. 1	
youth; 1	
youthful 1	

Lab 4

- I. Run wordcount job
- 2. Review output
- 3. Run wordcount job again with same parameters

WRITING MAP REDUCE JOBS

{**N**1,**V**1}



The MapReduce Pipeline

A mapper receives (Key, Value) & outputs (Key, Value) A reducer receives (Key, Iterable[Value]) and outputs (Key, Value) Partitioning / Sorting / Grouping provides the Iterable[Value] & Scaling

$\{KI,VI\} \longrightarrow \{K2, List < V2 >\} \rightarrow \{K3,V3\}$

MOBY DICK; OR THE WHALE

By Herman Melville

CHAPTER 1. Loomings.

Call me Ishmael. Some years ago--never mind how long precisely--having little or no money in my purse, and nothing particular to interest me on shore, I thought I would sail about a little and see the watery part of the world. It is a way I have of driving off the spleen and regulating the circulation. Whenever I find myself growing grim about the mouth; whenever it is a damp, drizzly November in my soul; whenever I find myself involuntarily pausing before coffin warehouses, and bringing up the rear of every funeral I meet; and especially whenever my hypos get such an upper hand of me, that it requires a strong moral principle to prevent me from deliberately stepping into the street, and methodically knocking people's hats off--then, I account it high time to get to sea as soon as I can. This is my substitute for pistol and ball. With a philosophical flourish Cato throws himself upon his sword; I quietly take to the ship. There is nothing surprising in this. If they but knew it, almost all men in their degree, some time or other, cherish very nearly the same feelings towards the ocean with me.

There now is your insular city of the Manhattoes, belted round by wharves as Indian isles by coral reefs--commerce surrounds it with her surf. Right and left, the streets take you waterward. Its extreme downtown is the battery, where that noble mole is washed by waves, and cooled by breezes, which a few hours previous were out of sight of land. Look at the crowds of water-gazers there.



ΚV

1 Call me Ishmael. Some years ago--never mind how long precisely--having

2 little or no money in my purse, and nothing particular to interest me on

3 shore, I thought I would sail about a little and see the watery part of

4 the world. It is a way I have of driving off the spleen and regulating

5 the circulation.



Mapper

```
package com.manifestcorp.hadoop.wc;
import java.io.IOException;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;
                                                ΚI
                                                    VI K2
public class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
   private static final String SPACE = " ";
   private static final IntWritable ONE = new IntWritable(1);
   private Text word = new Text();
                                KI VI
   public void map(LongWritable key, Text value, Context context)
                                            throws IOException, InterruptedException {
      String[] words = value.toString().split(SPACE);
      for (String str: words) {
         word.set(str); K2 V2
         context.write(word, ONE);
      }
   }
```



ΚV

1 Call me Ishmael. Some years ago--never mind how long precisely--having

2 $% \left({{\left({{\left({{\left({1 \right)}} \right)}} \right)}} \right)$ little or no money in my purse, and nothing particular to interest me on

3 shore, I thought I would sail about a little and see the watery part of

4 the world. It is a way I have of driving off the spleen and regulating

5 the circulation.

	Κ	V	К	V		Κ	V
	Call	1	agonever	1		agonever	1
	me	1	Call	1		Call	1
mah 🦛	Ishmael.	1	how	1		how	1
map	Some	1	Ishmael.	1		Ishmael.	1
	years	1 cort	me	1	aroub	me	1
	agonever	1 SUIL	little	1	group	little	1
	mind	1	long	1	0 1	long	1
	how	1	mind	1		mind	1
	of	1	of	1		of	1,1,1
	long	1	of	1		or	1
	little	1	of	1		Some	1
	of	1	or	1		years	1
	or	1	Some	1			
	of	1	years	1			



Reducer

```
package com.manifestcorp.hadoop.wc;
import java.io.IOException;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;
                                               K2
                                                                  K3
                                                      V2
                                                                         V3
public class WordCountReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
                                                       V2
                           К2
   public void reduce(Text key, Iterable<IntWritable> values, Context context)
                                        throws IOException, InterruptedException {
      int total = 0;
      for (IntWritable value : values) {
          total++;
       }
                     K3
      context.write(key, new IntWritable(total));
   }
}
```



ΚV

1 Call me Ishmael. Some years ago--never mind how long precisely--having

2 $% \left({{\left({{\left({{\left({1 \right)}} \right)}} \right)}} \right)$ in my purse, and nothing particular to interest me on

3 shore, I thought I would sail about a little and see the watery part of

4 the world. It is a way I have of driving off the spleen and regulating

5 the circulation.

	K	V		К	V		К	V	Κ	V
	Call	1		agonever	1		agonever	1	agonever	1
_	me	1		Call	1		Call	1	Call	1
mah 🧖	Ishmael.	1		how	1		how	1	how	1
map	Some	1		Ishmael.	1		Ishmael.	1	Ishmael.	1
	years		v+ -	me	1	groub	me	1 roduco	me	1
	agonever	1 SO	יונ	little	1	group	little		little	1
	mind	1		long	1	0 1	long	1	long	1
	how	1		mind	1		mind	1	mind	1
	of	1		of	1		of	1,1,1	of	3
	long	1		of	1		or	1	or	1
	little	1		of	1		Some	1	Some	1
	of	1		or	1		years	1	years	1
	or	1		Some	1					
	of	1		years	1					



Driver

```
package com.manifestcorp.hadoop.wc;
```

}

```
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
```

```
public class MyWordCount {
   public static void main(String[] args) throws Exception {
       Job job = new Job();
       job.setJobName("my word count");
       job.setJarByClass(MyWordCount.class);
      FileInputFormat.addInputPath(job, new Path(args[0]));
      FileOutputFormat.setOutputPath(job, new Path(args[1]));
       job.setMapperClass(WordCountMapper.class);
       job.setReducerClass(WordCountReducer.class);
       job.setOutputKeyClass(Text.class);
       job.setOutputValueClass(IntWritable.class);
       System.exit(job.waitForCompletion(true) ? 0 : 1);
   }
```

pom.xml

<properties> <hadoop.version>2.2.0</hadoop.version> </properties>

<build>

```
<plugins>
    <plugin>
      <artifactId>maven-compiler-plugin</artifactId>
      <version>2.3.2</version>
      <configuration>
        <source>1.7</source>
        <target>1.7</target>
      </configuration>
    </plugin>
  </plugins>
</build>
<dependencies>
  <dependency>
    <groupId>org.apache.hadoop</groupId>
    <artifactId>hadoop-client</artifactId>
    <version>${hadoop.version}</version>
    <scope>provided</scope>
```

```
</dependency> </dependencies>
```

```
</project>
```

Run Hadoop Job

\$ hadoop jar target/hadoop-mywordcount-0.0.1-SNAPSHOT.jar com.manifestcorp.hadoop.wc.MyWordCount /books out

Lab 5

- I. Unzip /opt/data/hadoop-mywordcount-start.zip
- 2. Write Mapper class
- 3. Write Reducer class
- 4. Write Driver class
- 5. Build (mvn clean package)
- 6. Run mywordcount job
- 7. Review output

UNIT TESTING


https://mrunit.apache.org/

```
<dependency>
```

<groupId>junit</groupId>
 <artifactId>junit</artifactId>
 <version>4.11</version>
 <scope>test</scope>
</dependency>

```
<dependency>
```

```
<groupId>org.apache.mrunit</groupId>
<artifactId>mrunit</artifactId>
<version>1.1.0</version>
<scope>test</scope>
<classifier>hadoop2</classifier>
</dependency>
```

```
package com.manifestcorp.hadoop.wc;
```

```
import java.io.IOException;
import java.util.ArrayList;
import java.util.List;
```

```
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mrunit.mapreduce.MapDriver;
import org.apache.hadoop.mrunit.mapreduce.MapReduceDriver;
import org.apache.hadoop.mrunit.mapreduce.ReduceDriver;
import org.junit.Before;
import org.junit.Test;
```

public class WordCountTest {

```
MapReduceDriver<LongWritable, Text, Text, IntWritable, Text, IntWritable> mapReduceDriver;
MapDriver<LongWritable, Text, Text, IntWritable> mapDriver;
ReduceDriver<Text, IntWritable, Text, IntWritable> reduceDriver;
```

```
@Before
```

```
public void setUp() {
```

```
WordCountMapper mapper = new WordCountMapper();
WordCountReducer reducer = new WordCountReducer();
```

```
mapDriver = new MapDriver<>(mapper);
reduceDriver = new ReduceDriver<>(reducer);
```

```
mapReduceDriver = new MapReduceDriver<>(mapper, reducer);
```

```
}
```

```
@Test
```

```
public void testMapper() throws IOException {
    mapDriver.withInput(new LongWritable(1), new Text("java hadoop java"));
    mapDriver.withOutput(new Text("java"), new IntWritable(1));
    mapDriver.withOutput(new Text("hadoop"), new IntWritable(1));
    mapDriver.withOutput(new Text("java"), new IntWritable(1));
    mapDriver.runTest();
}
// ...
```

```
public class WordCountTest {
  MapReduceDriver<LongWritable, Text, Text, IntWritable, Text, IntWritable> mapReduceDriver;
  MapDriver<LongWritable, Text, Text, IntWritable> mapDriver;
  ReduceDriver<Text, IntWritable, Text, IntWritable> reduceDriver;
  @Before
  public void setUp() {
    WordCountMapper mapper = new WordCountMapper();
    WordCountReducer reducer = new WordCountReducer();
   mapDriver = new MapDriver<>(mapper);
    reduceDriver = new ReduceDriver<>(reducer);
   mapReduceDriver = new MapReduceDriver<>(mapper, reducer);
  @Test
  public void testMapper() throws IOException {
   mapDriver.withInput(new LongWritable(1), new Text("java hadoop java"));
   mapDriver.withOutput(new Text("java"), new IntWritable(1));
   mapDriver.withOutput(new Text("hadoop"), new IntWritable(1));
    mapDriver.withOutput(new Text("java"), new IntWritable(1));
   mapDriver.runTest();
  @Test
  public void testReducer() throws IOException {
   List<IntWritable> values = new ArrayList<>();
    values.add(new IntWritable(1));
    values.add(new IntWritable(1));
    reduceDriver.withInput(new Text("java"), values);
    reduceDriver.withOutput(new Text("java"), new IntWritable(2));
    reduceDriver.runTest();
  }
  @Test
  public void testMapReduce() throws IOException {
   mapReduceDriver.withInput(new LongWritable(1), new Text("java hadoop java"));
   mapReduceDriver.withInput(new LongWritable(2), new Text("java spring java"));
    mapReduceDriver.addOutput(new Text("hadoop"), new IntWritable(1));
    mapReduceDriver.addOutput(new Text("java"), new IntWritable(4));
   mapReduceDriver.addOutput(new Text("spring"), new IntWritable(1));
   mapReduceDriver.runTest();
```

HADDOP IN THE CLOUD



Elastic MapReduce





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

T AWS - Services - Edit -		cmj @ 5637-0073-6850	✓ Global ✓ Support ✓
Upload Create Folder Actions ~ All Buckets / cmj-emr		None Properties	Transfers C ⁴
Name	Storage Class	Size	Last Modified
hadoop-mywordcount-0.0.1-SNAPSHOT.jar	Standard	5.4 KB	Tue Jan 06 22:01:56 GMT-500 2015

AWS - Services - Edit -

Welcome to Amazon Elastic MapReduce

Amazor Elastic MapReduce (Amazon EMR) is a web service that enables businesses, researchers, data analysts, and developers to easily and cost-effectively process vast amounts of data.

You on the prear to have any clusters. Create one now:

Create cluster

Υ

How Elastic MapReduce Works

Upload

Create



Upload your data and processing application to S3.

Learn more

Learn more

Configure and create your cluster by

specifying data inputs, outputs,

cluster size, security settings, etc.

Monitor



Monitor the health and progress of your cluster. Retrieve the output in \$3.

Learn more

Additional Information

More about Elastic MapReduce

EMR overview FAQs Pricing

More Help Using Elastic MapReduce

Forum

Documentation Developer Guide Quick Reference Card API Reference EMR on GitHub Help portal

RapReduce < Create Cluster			anj	
Cluster Configuration Configure sample app Cluster name Word count Termination protection • Yes • No Prevents accidential termination of the cluster, some must turn of termination of the cluster some must turn of termination of the cluster, some must turn of termination of the cluster, some must turn of termination of the cluster, some must turn of termination of the cluster some must turn o	Reduce - Create Cluste	er		
Cluster configuration Cluster name Vest could Prevets accident termination of the cluster to shup or the former to cluster is on the	luster Configuration			Configure sample application
Termination protection Yes No Prevents accidental termination protection Logging Enabled Copy the cluster's log files automatically to 53. Lear more Copy the cluster's log files automatically to 53. Lear more Copy the cluster's log files automatically to 53. Lear more Copy the cluster's log files automatically to 53. Lear more Copy the cluster's log files automatically to 53. Lear more Copy the cluster's log files automatically to 53. Lear more Tags Obtional: Add up to 10 tags to your EMR cluster. A tag consists of a case-sensitive key-value pair. Tags on EMR clusters are propage to the underlying EO2 instances. Learn more about tagging your Amazon EMR clusters. Key Value (optional) Add a key to create a tag Detromises the base configuration Hadoop distribution Amazon MapR Use MapR's Hadoop distribution. Learn more MapR Use MapR's Hadoop distribution. Learn more Additional applications Sect an application File System Configuration Configure and add File System Configuration Configure and add File System Configuration Server-side encryption and consistent view for EMR's Learn more Server-side encryption Consigure and add <td>Cluster name</td> <td>Word count</td> <td></td> <td></td>	Cluster name	Word count		
Consistent view No Logging Enabled Log folder S3 location Software Software Copy the cluster's log files automatically to S3. Last more Tags Index logs to enable corrected debugging functionality (requires logging). Learn more Tags Optional: Add up to 10 tags to your EMR cluster. A tag consists of a case-sensitive key-value pair. Tags on EMR clusters are propage to the underlying EC2 instances. Learn more about tagging your Amazon EMR clusters. Key Value (optional) Add a key to oraste a tag Optional: Add up to 10 tags to your EMR cluster. A tag consists of a case-sensitive key-value pair. Tags on EMR clusters are propage to the underlying EC2 instances. Learn more about tagging your Amazon EMR clusters. Key Value (optional) Add a key to oraste a tag Optional: Add up to 10 tags to your EMR cluster. A tag consists of a case-sensitive key-value pair. Tags on EMR clusters are propage to the underlying EC2 instances. Learn more about tagging your Amazon EMR cluster, including the Hadoop distribution. Learn more Mil version Optional: Add up to restrict a tag Software Configuration Use Amazon's Hadoop distribution. Learn more MapR Use MapPris Hadoop distribution. Learn more Applications to be installed Version Additional applications Select an application The EMR Fis	Termination protection	Ves	Prevents accidental ter	mination of the cluster: to shut
Logging Enabled Copy the cluster's log files automatically to S3. Leamore Log folder S3 location s3://-cbucket-names/-dolders/ Debugging Enabled Index logs to enable console debugging functionality frequines logging). Leam more Tags Optional: Add up to 10 tags to your EMR cluster. A tag consists of a case-sensitive key-value pair. Tags on EMR clusters are propage to the underlying EO2 instances. Learn more about tagging your Amazon EMR clusters. Key Value (optional) /.dds key to create a tag	renning of processor	No	down the cluster, you r protection. Learn mor	nust turn off termination
Log folder S3 location st/strice-envisor st/str	Logging	Enabled	Copy the cluster's log	iles automatically to S3. Learn
		Log folder S3 location	india.	
Software Configuration Additional applications Sever-side encryption Enabled Index logs to enable console debugging functionality (requires logging). Learn more Index logs to enable console debugging functionality (requires logging). Learn more Index logs to enable console debugging functionality (requires logging). Learn more Index logs to the underlying EC2 Instances. Learn more about tagging your Amazon EMR clusters. Key Value (optional) Add a key to create a tag Software Configuration Hadoop distribution Amazon Use Amazon's Hadoop distribution. Learn more Additional applications File System Configuration The EMR File System (EMRFS) and the Hadoop Distributed File System (HDFS) are both installed on your EMR cluster. HDFS stores on EMR Cluster, while EMRFS allows EMR clusters to store data on S3. You can enable server-side encryption and consistent view fe EMRFs Learn more Consistent view Enabled Marine Enable Marine Enable Marine Enable		s3://cmi-emr/logs/	5	
Debugging Enabled Index logs to enable console debugging functionality (requires logging). Learn more Tags Optional: Add up to 10 tags to your EMR cluster. A tag consists of a case-sensitive key-value pair. Tags on EMR clusters are propage to the underlying EC2 instances. Learn more about tagging your Amazon EMR clusters. Key Value (optional) Add a key to create a tag Software Configuration Macop distribution Amazon Use Amazon's Hadoop distribution. Learn more Add a key to create a tag Determines the base configuration of the instances in your cluster, including the Hadoop version. Learn more MapR Use MapPi's Hadoop distribution. Learn more Applications to be installed Version Configure and add File System Configuration The EMR File System (EMRFS) and the Hadoop Distributed File System (HDFS) are both installed on your EMR cluster. HDFS stores on EMR cluster, while EMRFS allows EMR clusters to store data on S3. You can enable server-side encryption and consistent view for EMRFs. Learn more Server-side encryption Enabled Uses S3 server-wile encryption on files written to S3 EMRFs. Learn more Consistent view Enabled Monitorn its tan read-after-write (for new puts) Monitorn its an read-after-write (for new puts) Monitorn		s3:// <bucket-name>/<folder>/</folder></bucket-name>	_	
Tags • Optional: Add up to 10 tags to your EMR cluster. A tag consists of a case-sensitive key-value pair. Tags on EMR clusters are propage to the underlying EC2 Instances. Learn more about tagging your Amazon EMR clusters. • May to create a lag Software Configuration Hadoop distribution AMI version 3.1 your cluster, including the Hadoop version. Learn more AMI version AMI version Aution MapR Use MapR's Hadoop distribution. Learn more Additional applications Flie System Configuration The EMR File System (EMRFS) and the Hadoop Distributed File System (HDFS) are both installed on your EMR cluster. HDFS stores on an EMR cluster, while EMRFS allows EMR clusters to store data on S3. You can enable server-side encryption and consistent view for EMRFS below, or use a bootstrap action to configure additional settings for EMRFS. Server-side encryption on files written to S3 EMRFS below, or use a bootstrap action to configure additional settings for EMRFS. Server-side encryption end listers to store data on S3. You can enable server-side encryption and consistent vie	Debugging	Enabled	Index logs to enable co (requires logging). Let	nsole debugging functionality m more
Tags • Optional: Add up to 10 tags to your EMR cluster. A tag consists of a case-sensitive key-value pair. Tags on EMR clusters are propage to the underlying EG2 instances. Learn more about tagging your Amazon EMR clusters. Key Value (optional) Add a key to create a tag Software Configuration Macop distribution Hadoop distribution Amazon Maloop distribution Amazon MapR Determines the base configuration of the instances in your cluster, including the Hadoop viersion Applications to be installed Version Additional applications Select an application Configure and add File System Configuration The EMR File System (EMRFS) and the Hadoop Distributed File System (HDFS) are both installed on your EMR cluster, HDFS stores on an EMR cluster, while EMRFS allows EMR clusters to store data on S3. You can enable server-side encryption and consistent view for EMRFS. Learn more EMRFS below, or use a bootstrap action to configure additional settings for EMRFS. Server-side encryption on files written to S3 EMRFS. Learn more				
Optional: Add up to 10 tags to your EMR cluster. A tag consists of a case-sensitive key-value pair. Tags on EMR clusters are propage to the underlying EC2 instances. Learn more about tagging your Amazon EMR clusters. Key Value (optional) Add a key to create a tag Software Configuration Hadoop distribution Hadoop distribution AMI version 3.1 Our cluster, including the Hadoop version. Learn more MapR Use MapR's Hadoop distribution. Learn more Additional applications MapR Use MapR's Hadoop distribution. Learn more Additional applications Select an application Configure and add File System (EMRFS) and the Hadoop Distributed File System (HDFS) are both installed on your EMR cluster. HDFS stores on an EMR cluster, while EMRFS allows EMR clusters to store data on S3. You can enable server-side encryption and consistent view for EMRFS. Learn more EMRFS below, or use a bootstrap action to configure additional settings for EMRFS. Server-side encryption Enabled Version Consistent view Enabled	ags			
Key Value (optional) Add a key to create a tag	 Optional: Add up to 10 tags to to the underlying EC2 instances. 	o your EMR cluster. A tag consists of a ca Learn more about tagging your Amazon E	e-sensitive key-value pair. Tags on IR clusters.	EMR clusters are propagated
Add a key to create a tag Software Configuration Hadoop distribution • Amazon Use Amazon's Hadoop distribution. Learn more AMI version 3.3.1 Determines the base configuration of the instances in your cluster, including the Hadoop version. MapR Use MapR's Hadoop distribution. Applications to be installed Version Additional applications Select an application Configure and add File System Configuration The EMR File System (EMRFS) and the Hadoop Distributed File System (HDFS) are both installed on your EMR cluster. HDFS stores on an EMR cluster, while EMRFS allows EMR clusters to store data on S3. You can enable server-side encryption and consistent view figMR's below, or use a bootstrap action to configure additional settings for EMRFS. Server-side encryption Enabled Uses S3 server-side encryption on files written to S3 EMR's. Learn more Consistent view Enabled	Кеу	Value (optional)		
Software Configuration Hadoop distribution Amazon Use Amazon's Hadoop distribution. Learn more AMI version 3.1 Determines the base configuration of the instances in your cluster, including the Hadoop version. Learn more MapR Use MapR's Hadoop distribution. Learn more Applications to be installed Version Additional applications Select an application Configure and add Configure and add	Add a key to create a tag			
AMI version 3.3.1 MapR Use MapR's Hadoop distribution. Learn more Applications to be installed Version Additional applications Select an application Configure and add	oftware Configuration	Amazon	Use Amazon's Hadoop	
3.3.1 Configuration of the instances is your cluster, including the Hadoop version. Learn more MapR Use MapR's Hadoop distribution. Learn more Applications to be installed Version Additional applications Select an application Configure and add Configure and add File System Configuration Configure and add Sever-side encryption Enabled Version Uses S3 server-side encryption and consistent view for EMRFS. Learn more Consistent view Enabled		•		distribution. Learn more
MapR Use MapR's Hadoop distribution. Learn more Applications to be installed Version Additional applications Select an application Configure and add Configure and add File System Configuration Configure and add The EMR File System (EMRFS) and the Hadoop Distributed File System (HDFS) are both installed on your EMR cluster. HDFS stores on an EMR cluster, while EMRFS allows EMR clusters to store data on S3. You can enable server-side encryption and consistent view for EMRFS below, or use a bootstrap action to configure additional settings for EMRFS. Server-side encryption Enabled Uses S3 server-side encryption on files written to S3 EMRFS. Learn more Consistent view Enabled		AMI version		distribution. Learn more
Applications to be installed Version Additional applications Select an application Configure and add Configure and add File System Configuration Configure and add File System Configuration Configure and add Sever-side encryption Enabled Uses S3 server-side encryption on files written to S3 EMRFS. Learn more Consistent view Enabled		AMI version 3.3.1	Determines the base of ware cluster inclusion	onfiguration of the instances in
Additional applications Select an application File System Configuration Image: The EMR File System (EMRFS) and the Hadoop Distributed File System (HDFS) are both installed on your EMR cluster. HDFS stores on an EMR cluster, while EMRFS allows EMR clusters to store data on S3. You can enable server-side encryption and consistent view for EMRFS below, or use a bootstrap action to configure additional settings for EMRFS. Server-side encryption Enabled Uses S3 server-side encryption on files written to S3 EMRFS. Consistent view Enabled Monitors list and read-after-write (for new puts)		AMI version 3.3.1 MapR	Determines the base of your cluster, including Use MapR's Hadoop d	onfiguration of the instances in the Hadoop version. Learn more
Configure and add File System Configuration Image: The EMR File System (EMRFS) and the Hadoop Distributed File System (HDFS) are both installed on your EMR cluster. HDFS stores on an EMR cluster, while EMRFS allows EMR clusters to store data on S3. You can enable server-side encryption and consistent view for EMRFS below, or use a bootstrap action to configure additional settings for EMRFS. Server-side encryption Enabled Uses S3 server-side encryption on files written to S3 EMRFS. Consistent view Enabled Monitors list and read-after-write (for new puts)	Applications to be installed	AMI version 3.3.1 MapR Version	Determines the base of your cluster, including Use MapR's Hadoop d	distribution. Learn more onfiguration of the instances in the Hadoop version. Learn more istribution. Learn more
File System Configuration Image: The EMR File System (EMRFS) and the Hadoop Distributed File System (HDFS) are both installed on your EMR cluster. HDFS stores on an EMR cluster, while EMRFS allows EMR clusters to store data on S3. You can enable server-side encryption and consistent view for EMRFS below, or use a bootstrap action to configure additional settings for EMRFS. Server-side encryption Enabled Uses S3 server-side encryption on files written to S3 EMRFS. Consistent view Enabled Monitors list and read-after-write (for new puts)	Applications to be installed Additional applications	AMI version 3.3.1 MapR Version Select an application	 Determines the base of your cluster, including Use MapR's Hadoop d 	distribution. Learn more onfiguration of the instances in the Hadoop version. Learn more istribution. Learn more
on an EMR cluster, while EMRFS allows EMR clusters to store data on S3. You can enable server-side encryption and consistent view f EMRFS below, or use a bootstrap action to configure additional settings for EMRFS. Server-side encryption Enabled Uses S3 server-side encryption on files written to S3 EMRFS. Learn more Consistent view Enabled Monitors list and read-after-write (for new puts)	Applications to be installed Additional applications	AMI version 3.3.1 MapR Version Select an application Configure and add	 Determines the base of your cluster, including Use MapR's Hadoop d Ise MapR's Hadoop d 	distribution. Learn more onfiguration of the instances in the Hadoop version. Learn more istribution. Learn more
Server-side encryption Enabled Uses S3 server-side encryption on files written to S3 EMRFS. Consistent view Enabled Monitors list and read-after-write (for new puts)	Applications to be installed Additional applications ile System Configuration	AMI version 3.3.1 MapR Version Select an application Configure and add add add add add	 Determines the base of your cluster, including to Use MapR's Hadoop d Use MapR's Hadoop d Use MapR's Hadoop d 	EMR cluster. HDFS stores data
Consistent view Enabled Monitors list and read-after-write (for new puts)	Applications to be installed Additional applications ile System Configuration The EMR File System (EMRFS on an EMR cluster, while EMRFS EMRFS below, or use a bootstrap	AMI version 3.3.1 MapR Version Select an application Configure and add S) and the Hadoop Distributed File System allows EMR clusters to store data on S3. D action to configure additional settings for	Determines the base of your cluster, including Use MapR's Hadoop d C (HDFS) are both installed on your B 'ou can enable server-side encrypt EMRFS.	Additionant Consistent view for
consistency for files in S3. Learn more	Applications to be installed Additional applications ile System Configuration The EMR File System (EMRFS on an EMR cluster, while EMRFS EMRFS below, or use a bootstrap Server-side encryption	AMI version 3.3.1 MapR Version Select an application Configure and add add add add add add add add add ad	Determines the base or your cluster, including Use MapR's Hadoop d Use MapR's Hadoop d (HDFS) are both installed on your B 'ou can enable server-side encrypt EMRFS. Uses S3 server-side en EMRFS. Learn more	EMR cluster. HDFS stores data tion and consistent view for cryption on files written to S3 by

Pricing for Amazon EMR and Amazon EC2 (On-Demand)

Region: US East (Northern Virginia) +	
	Amazon EC2 Price
Standard On-Demand Instances	
Small (Default)	\$0.06 per hour
Medium	\$0.12 per hour
Large	\$0.24 per hour
Extra Large	\$0.48 per hour
Hi-Memory On-Demand Instances	
Extra Large	\$0.41 per hour
Double Extra Large	\$0.82 per hour
Quadruple Extra Large	\$1.64 per hour

Hi-CPU On-Demand Instances		
Medium	\$0.145 per hour	\$0.03 per hour
Extra Large	\$0.58 per hour	\$0.12 per hour
Cluster Compute On-Demand Instances		
Quadruple Extra Large	\$1.30 per hour	\$0.27 per hour
Eight Extra Large	\$2.40 per hour	\$0.50 per hour
Cluster GPU On-Demand Instances		
Quadruple Extra Large	\$2.10 per hour	\$0.42 per hour
High-I/O On-Demand Instances		
Quadruple Extra Large	\$3.10 per hour	\$0.47 per hour
High-Storage On-Demand Instances		
Eight Extra Large	\$4.60 per hour	\$0.69 per hour

Amazon Elastic

MapReduce Price

\$0.015 per hour

\$0.03 per hour

\$0.06 per hour

\$0.12 per hour

\$0.09 per hour

\$0.21 per hour

\$0.42 per hour

MAKING IT MORE REAL

http://aws.amazon.com/architecture/



AWS Reference Architectures

more details on getting one on

one support for your architecture questions.

> The flexibility of AWS allows you to design your application architectures the way you like. AWS Reference Architecture Datasheets provide you with the architectural guidance you need in order to build an application that takes full advantage of the AWS cloud. Each datasheet includes a visual representation of the architecture and basic description of how each service is used.









Large Scale Processing and Huge Data sets Build high-performance computing systems that involve Big Data (PDF) Ad Serving Build highly-scalable online ad serving solutions (PDF)

Disaster Recovery for Local Applications Build cost-effective Disaster Recovery solutions for on-premises applications (PDF) File Synchronization Build simple file synchronization service (PDF)







Online Games Build powerful online games (PDF)







Financial Services Grid Computing Build highly scalable and elastic grids for the Financial Services Sector (PDF)



E-Commerce Website Part 1: Web Frontend Build elastic Web Frontends for an e-Commerce website (PDF)



E-Commerce Website Part 2: Checkout Pipeline

Build highly scalable checkout pipeline for an e-Commerce website (PDF)



E-Commerce Website Part 3: Marketing and Recommendations Build highly scalable recommendation engine for an e-Commerce website (PDF)

WEB LOG ANALYSIS

Amazon Web Services provides services and infrastructure to build reliable, fault-tolerant, and highly available web applications in the cloud. In production environments, these applications can generate huge amounts of log information.

This data can be an important source of knowledge for any company that is operating web applications. Analyzing logs can reveal information such as traffic patterns, user behavior, marketing profiles, etc.

However, as the web application grows and the number of visitors increases, storing and analyzing web logs becomes increasingly challenging.

This diagram shows how to use Amazon Web Services to build a scalable and reliable large-scale log analytics platform. The core component of this architecture is Amazon Elastic MapReduce, a web service that enables analysts to process large amounts of data easily and cost-effectively using a Hadoop hosted framework.



System Overview

The web front-end servers are running on Amazon Elastic Compute Cloud (Amazon EC2) instances.

Amazon CloudFront is a content delivery network that 2 uses low latency and high data transfer speeds to distribute static files to customers. This service also generates valuable log information.

3 Log files are periodically uploaded to Amazon Simple Storage Service (Amazon S3), a highly available and reliable data store. Data is sent in parallel from multiple web servers or edge locations.

An Amazon Elastic MapReduce cluster processes the data set. Amazon Elastic MapReduce utilizes a hosted Hadoop framework, which processes the data in a parallel job flow.

When Amazon EC2 has unused capacity, it offers EC2 5 instances at a reduced cost, called the Spot Price. This price fluctuates based on availability and demand. If your workload is flexible in terms of time of completion or required capacity, you can dynamically extend the capacity of your cluster using Spot Instances and significantly reduce the cost of running your job flows.

Data processing results are pushed back to a relational 6 database using tools like Apache Hive. The database can be an Amazon Relational Database Service (Amazon RDS) instance. Amazon RDS makes it easy to set up, operate, and scale a relational database in the cloud.

Like many services, Amazon RDS instances are priced on a pay-as-you-go model. After analysis, the database can be backed-up into Amazon S3 as a database snapshot, and then terminated. The database can then be recreated from the snapshot whenever needed.

ADVERTISEMENT SERVING

Internet advertising services need to serve targeted advertising and must do so under limited time. These are just two of multiple technical challenges they face.

Amazon Web Services provides services and infrastructure to build reliable, fault-tolerant, and highly available ad serving platforms in the cloud. In this document, we describe the two main parts of such a system: ad serving infrastructure and click-through collection featuring a data analysis cluster.



System Overview

When visitors load a web page, ad servers return a 1 pointer to the ad resource to be displayed. These servers are running on Amazon Elastic Compute Cloud (Amazon EC2) instances. They query a data set stored in an Amazon DynamoDB table to find relevant ads depending on the user's profile.

Ad files are downloaded from Amazon CloudFront, a content delivery service with low latency, high data-transfer speeds, and no commitments. Log information from displayed ads is stored on Amazon Simple Storage Service (Amazon S3), a highly available data store.

The click-through servers are a group of Amazon EC2 3 Ine click-through servers are a group of through data. This information is contained in the log files of the clickthrough web servers, which are periodically uploaded to Amazon S3.

Ad impression and click-through data are retrieved and processed by an Amazon Elastic MapReduce cluster using a hosted Hadoop framework to process the data in a parallel job flow. The cluster's capacity can be dynamically extended using Spot Instances to reduce the processing time and the cost of running the job flow.

Data processing results are pushed back into Amazon DynamoDB, a fully managed NoSQL database service that provides fast and predictable performance with seamless scalability. Amazon DynamoDB tables can store and retrieve any amount of data, and serve any level of request traffic, both of which are specific requirements for storing and quickly retrieving visitors' profile information.

The high availability and fast performance of Amazon DynamoDB enable ad server front-ends to serve requests with predictable response time, even with high traffic volumes or large profile's data sets.

ONLINE GAMES

Online games back-end infrastructures can be challenging to maintain and operate. Peak usage periods, multiple players, and high volumes of write operations are some of the most common problems that operations teams face.

But the most difficult challenge is ensuring flexibility in the scale of that system. A popular game might suddenly receive millions of users in a matter of hours, yet it must continue to provide a satisfactory player experience. Amazon Web Services provides different tools and services that can be used for building online games that scale under high usage traffic patterns.

This document presents a cost-effective online game architecture featuring automatic capacity adjustment, a highly available and high-speed database, and a data processing cluster for player behavior analysis.



System Overview

AWS

Browser games can be represented as client-server applications. The client generally consists of static files, such as images, sounds, flash applications, or Java applets. Those files are hosted on Amazon Simple Storage Service (Amazon S3), a highly available and reliable data store.

2 As the user base grows and becomes more geographically distributed, a high-performance cache Amazon CloudFront can provide substantial improvements in latency, fault tolerance, and cost. By using Amazon S3 as the origin server for the Amazon CloudFront distribution, the game infrastructure benefits from fast network data transfer rates and a simple publishing/caching workflow.

3 Requests from the game application are distributed by Elastic Load Balancing to a group of web servers running on Amazon Elastic Compute Cloud (Amazon EC2) instances. Auto Scaling automatically adjusts the size of this group, depending on rules like network load, CPU usage, and so on.

Player data is persisted on Amazon DynamoDB, a fully managed NoSQL database service. As the player population grows, Amazon DynamoDB provides predictable performance with seamless scalability.

5 Log files generated by each web server are pushed back into Amazon S3 for long-term storage.

6 Managing and analyzing high data volumes produced by online games platforms can be challenging. Amazon Elastic MapReduce (Amazon EMR) is a service that processes vast amounts of data easily. Input data can be retrieved from web server logs stored on Amazon S3 or from player data stored in Amazon DynamoDB tables to run analytics on player behavior, usage patterns, etc. Those results can be stored again on Amazon S3, or inserted in a relational database for further analysis with classic business intelligence tools.

7 Based on the needs of the game, Amazon Simple Email Service (Amazon SES) can be used to send email to players in a cost-effective and scalable way.

FINANCIAL SERVICES **GRID COMPUTING**

Financial services grid computing on the cloud provides dynamic scalability and elasticity for operation when compute jobs are required, and utilizing services for aggregation that simplify the development of grid software.

On demand provisioning of hardware, and template driven deployment, combined with low latency access to existing on-premise data sources make AWS a powerful platform for high performance grid computing systems.



System Overview

AWS

Date sources for market, trade, and counterparties are installed on startup from on premise data sources, or from Amazon Simple Storage Service (Amazon S3).

AWS DirectConnect can be used to establish a low 2 latency and reliable connection between the corporate data center site and AWS, in 1 to 10Gbit increments. For situations with lower bandwidth requirements, a VPN connection to the VPC Gateway can be established.

3 Private subnetworks are specifically created for customer source data, compute grid clients, and the grid controller and engines.

Application and corporate data can be securely stored in the cloud using the Amazon Relational Database Service (Amazon RDS).

Grid controllers and grid engines are running Amazon Elastic Compute Cloud (Amazon EC2) instances started on demand from Amazon Machine Images (AMIs) that contain the operating system and grid software.

6 Static data such as holiday calendars and QA libraries and additional gridlib bootstrapping data can be downloaded on startup by grid engines from Amazon S3.

Grid engine results can be stored in Amazon DynamoDB, a fully managed database providing configurable read and write throughput, allowing scalability on

demand. Results in Amazon DynamoDB are aggregated using a map/reduce job in Amazon Elastic MapReduce (Amazon EMR) and final output is stored in Amazon S3.

The compute grid client collects aggregate results from Amazon S3

Aggregate results can be archived using Amazon Glacier, a low-cost, secure, and durable storage service.

E-COMMERCE WEBSITE

With Amazon Web Services, you can build a recommendation and marketing service to manage targeted marketing campaigns and offer personalized product recommendations to customers who are browsing your e-commerce site.

In order to build such a service, you have to process very large amounts of data from multiple data sources. The resulting user profile information has to be available to deliver real-time product recommendations on your e-commerce website.

The insights that you gain about your customers can also be used to manage personalized marketing campaigns targeted at specific customer segments.

With the tools that AWS provides, you can build highly scalable recommendation services that can be consumed by different channels, such as dynamic product recommendations on the e - commerce website or targeted email campaigns for your customers.



System Overview

Amazon Elastic MapReduce (Amazon EMR) is a hosted Hadoop framework that runs on Amazon Elastic Compute Cloud (Amazon EC2) instances. It aggregates and processes user data from server log files and from the customer's purchase history.

2 An Amazon Relational Database Services (Amazon RDS) Read Replica of customer and order databases is used by Amazon EMR to compute user profiles and by Amazon Simple Email Service (Amazon SES) to send targeted marketing emails to customers.

3 Log files produced by the e-commerce web front end have been stored on Amazon Simple Storage Service (Amazon S3) and are consumed by the Amazon EMR cluster to compute user profiles.

User profile information generated by the Amazon EMR cluster is stored in Amazon DynamoDB, a scalable, high-performance managed NoSQL database that can serve recommendations with low latency.

A recommendation web service used by the web front end is deployed by AWS Elastic Beanstalk. This service uses the profile information stored on Amazon DynamoDB to provide personalized recommendations to be

shown on the e-commerce web front end.

A marketing administration application deployed by AWS Elastic Beanstalk is being used by marketing managers to send targeted email campaigns to customers with specific user profiles. The application reads customer email addresses from an Amazon RDS Read Replica of the customer database.

7 Amazon SES is used to send marketing emails to customers. Amazon SES is based on the scalable technology used by Amazon web sites around the world to send billions of messages a year.

TIME SERIES PROCESSING

When data arrives as a succession of regular measurements, it is known as time series information. Processing of time series information poses systems scaling challenges that the elasticity of AWS services is uniquely positioned to address.

This elasticity is achieved by using Auto Scaling groups for ingest processing, AWS Data Pipeline for scheduled Amazon Elastic MapReduce jobs, AWS Data Pipeline for intersystem data orchestration, and Amazon Redshift for potentially massive-scale analysis. Key architectural throttle points involving Amazon SQS for sensor message buffering and less frequent AWS Data Pipeline scheduling keep the overall solution costs predictable and controlled.



System Overview

AWS

Remote devices such as power meters, mobile clients, ad-network clients, industrial meters, satellites, and environmental meters measure the world around them and send sampled sensor data as messages via HTTP(S) for processing.

2 Send messages to an Amazon Simple Queue Service queue for processing into Amazon DynamoDB using autoscaled Amazon EC2 workers. Or, if the sensor source can do so, post sensor samples directly to Amazon DynamoDB. Try starting with a DynamoDB table that is a week-oriented, time-based table structure.

3 If a Supervisory Control and Data Acquisition (SCADA) system exists, create a flow of samples to or from Amazon DynamoDB to support additional cloud processing or other existing systems, respectively.

Using AWS Data Pipeline, create a pipeline with a regular Amazon Elastic MapReduce job that both calculates expensive sample processing and delivers samples and results.

5 The pipeline places results into Amazon Redshift for additional analysis.

The pipeline exports historical week-oriented sample tables, from Amazon DynamoDB to Amazon Simple Storage Service (Amazon S3)

The pipeline also optionally exports results in a format custom applications can accept.

Amazon Redshift optionally imports historic samples to reside with calculated results.

9 Using in-house or Amazon partner business intelligence solutions, Amazon Redshift supports additional analysis on a potentially massive scale.

Resources















Hadoop Real-World Solutions Cookbook

lealistic, simple code examples to solve problems at scale with ladoop and related technologies

Jonathan R. Owens Jon Lentz (Constant) open source





Resources





Christopher M. Judd



CTO and Partner email: cjudd@juddsolutions.com web: www.juddsolutions.com blog: juddsolutions.blogspot.com twitter: javajudd

