

# Amazon Elastic MapReduce

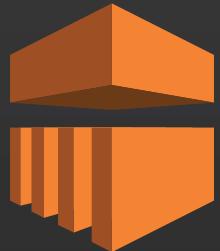


Provides a managed Hadoop framework  
Quickly & cost-effectively process vast amounts of data  
Makes it easy, fast & cost-effective for you to process data  
Run other popular distributed frameworks such as Spark

Easy to Use

Low Cost

Elastic

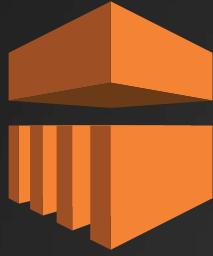


# Amazon EMR

Flexible

Reliable

Secure



# Amazon EMR: Example Use Cases

## Clickstream Analysis

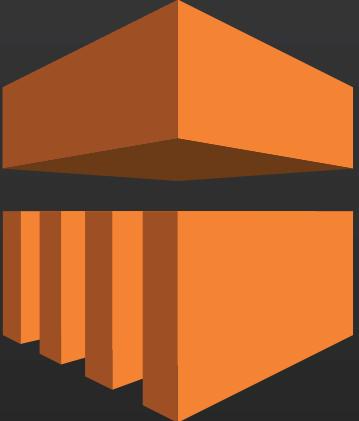
Amazon EMR can be used to analyze click stream data in order to segment users and understand user preferences. Advertisers can also analyze click streams and advertising impression logs to deliver more effective ads.

## Genomics

Amazon EMR can be used to process vast amounts of genomic data and other large scientific data sets quickly and efficiently. Researchers can access genomic data hosted for free on AWS.

## Log Processing

Amazon EMR can be used to process logs generated by web and mobile applications. Amazon EMR helps customers turn petabytes of un-structured or semi-structured data into useful insights about their applications or users.



# Agenda

Hadoop Fundamentals

Core Features of Amazon EMR

How to Get Started with Amazon EMR

Supported Hadoop Tools

Additional EMR Features

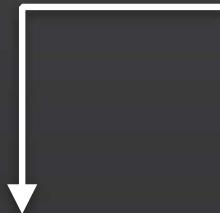
Third Party Tools

Resources where you can learn more

# HADOOP FUNDAMENTALS

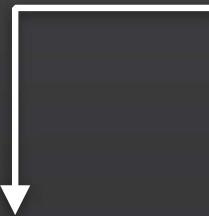
Very large  
clickstream  
logging data  
(e.g TBs)

Lots of actions by  
John Smith



Very large  
clickstream  
logging data  
(e.g TBs)

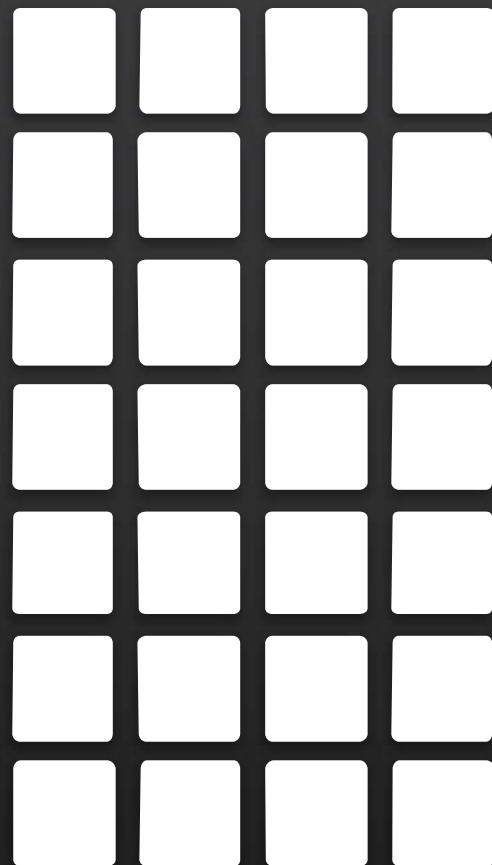
Lots of actions by  
John Smith



Very large  
clickstream  
logging data  
(e.g TBs)



Split the log  
into many  
small pieces

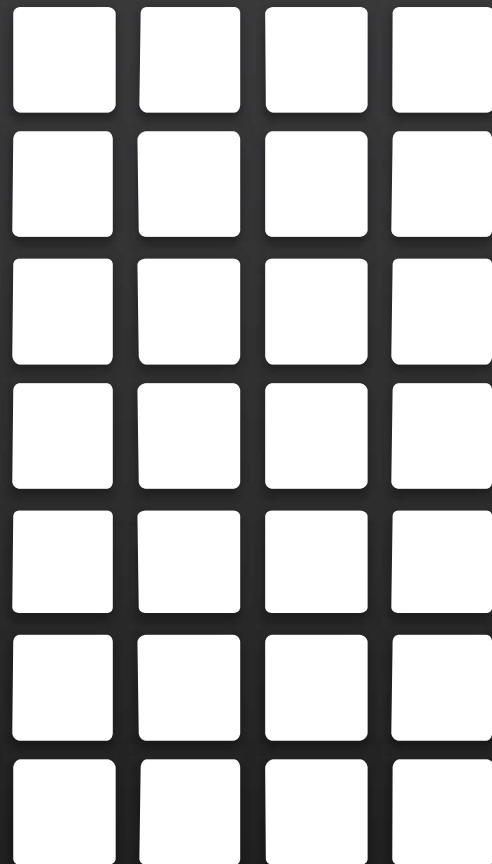


Very large  
clickstream  
logging data  
(e.g TBs)

Lots of actions by  
John Smith



Split the log  
into many  
small pieces

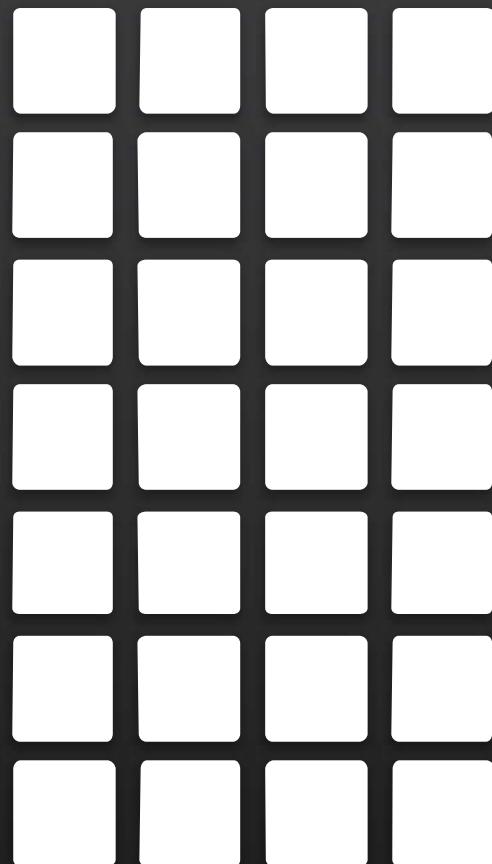


Process in an EMR  
cluster

Very large clickstream logging data (e.g TBs)

Lots of actions by John Smith

Split the log into many small pieces



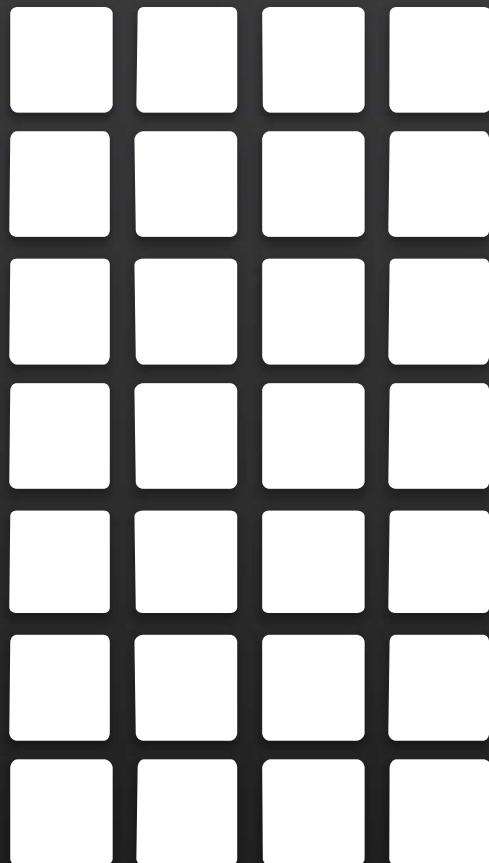
Process in an EMR cluster

Aggregate the results from all the nodes

Very large clickstream logging data (e.g TBs)

Lots of actions by John Smith

Split the log into many small pieces



Process in an EMR cluster

Aggregate the results from all the nodes

What John Smith did

Very large  
clickstream  
logging data  
(e.g TBs)

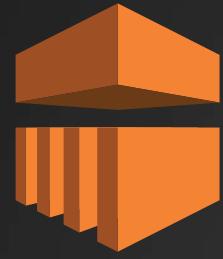


Insight in a fraction of the time

What John  
Smith did

# CORE FEATURES OF AMAZON EMR

# ELASTIC



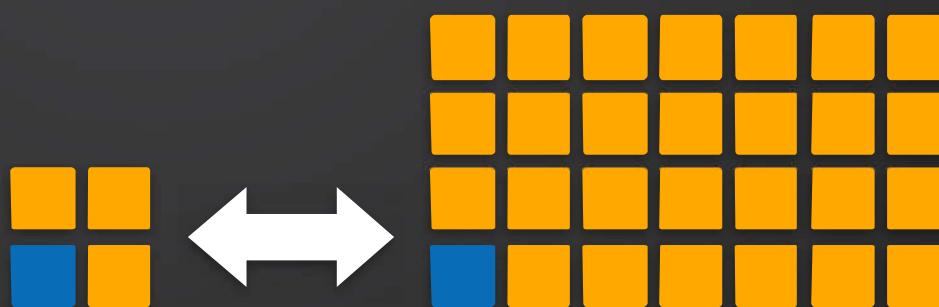
# Elastic

Provision as much capacity as you need  
Add or remove capacity at any time

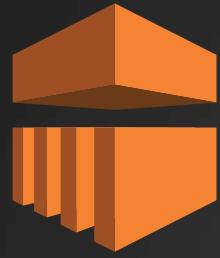
Deploy Multiple Clusters



Resize a Running Cluster



**LOW COST**



# Low Cost

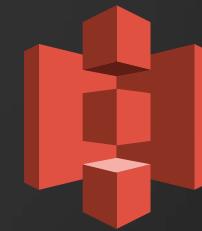
Low Hourly Pricing

Amazon EC2 Spot Integration

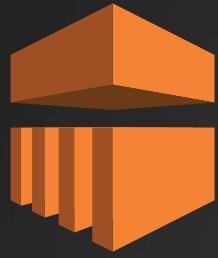


Amazon EC2 Reserved Instance Integration

Elasticity



Amazon S3 Integration



# Low Cost

Accenture Hadoop Study:

Amazon EMR ‘offers better price-performance’

The advertisement features a white cloud icon with a black cable extending from its bottom. A blue arrow points from the text 'High performance. Delivered.' towards the cable. The text 'Accenture Technology Labs' is at the top, followed by 'Where to deploy your Hadoop clusters?' and 'Executive Summary'. The Accenture logo is at the bottom right.

High performance. Delivered.

Accenture Technology Labs  
Where to deploy your Hadoop clusters?  
Executive Summary

accenture

consulting | technology | outsourcing

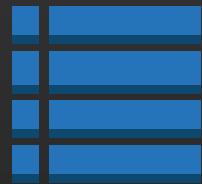
# FLEXIBLE DATA STORES



Amazon  
S3



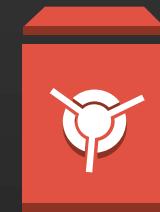
Hadoop Distributed  
File System



Amazon  
DynamoDB



Amazon  
Redshift



Amazon  
Glacier

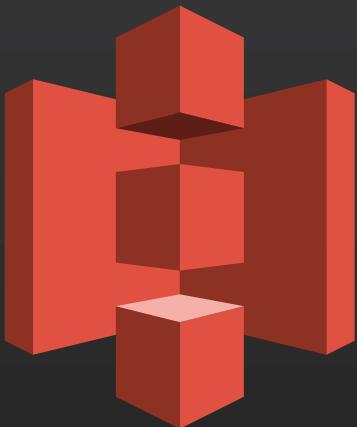


Amazon Relational  
Database Service



Amazon  
EMR

# Amazon S3 + Amazon EMR



- Allows you to decouple storage and computing resources
- Use Amazon S3 features such as server-side encryption
- When you launch your cluster, EMR streams data from S3
- Multiple clusters can process the same data concurrently

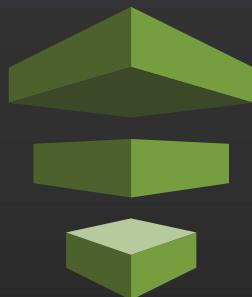
Hadoop Distributed  
File System (HDFS)



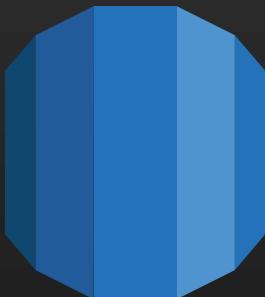
Amazon  
DynamoDB



AWS  
Data Pipeline

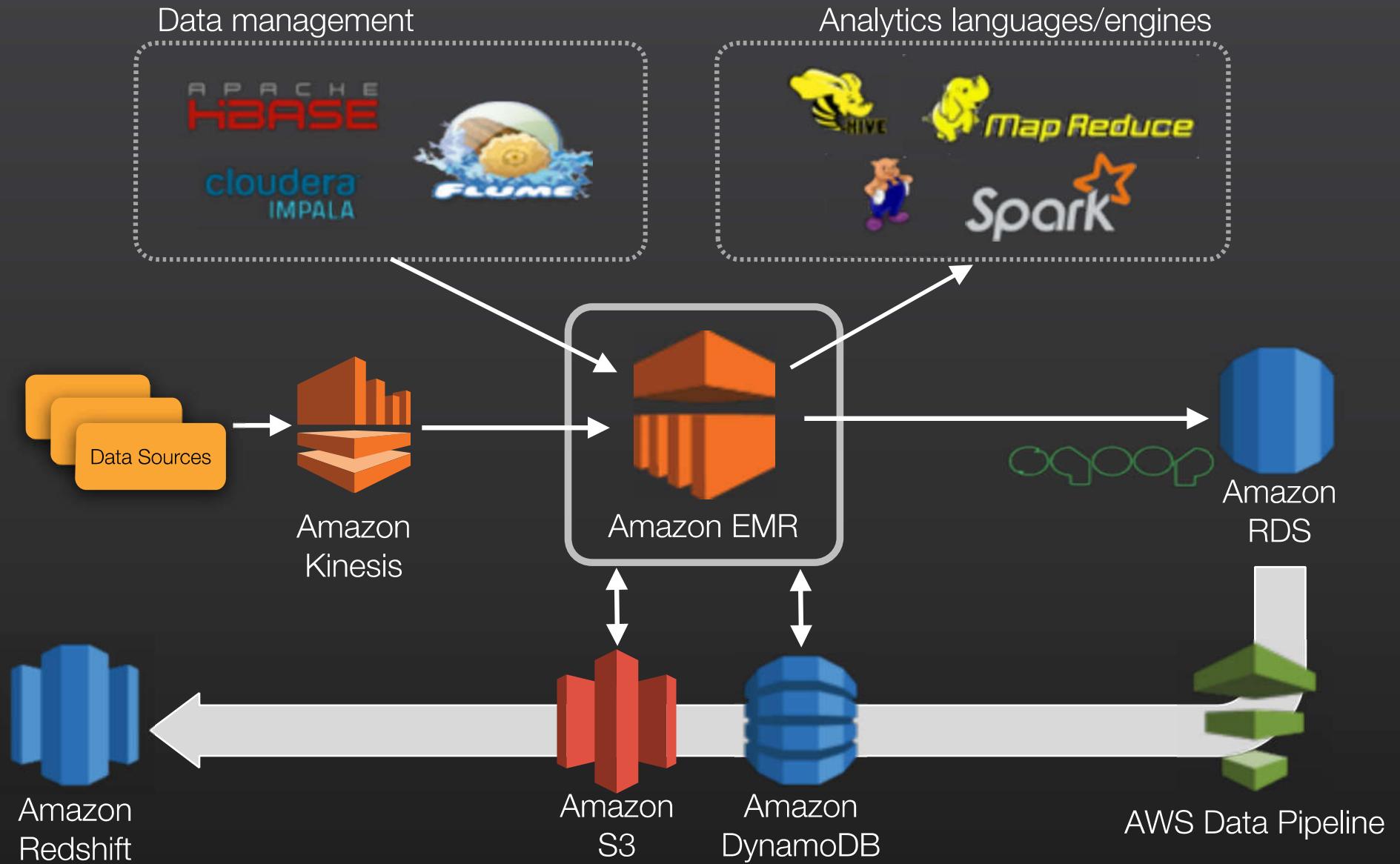


Amazon  
RDS



Amazon  
Redshift





# GETTING STARTED WITH AMAZON ELASTIC MAPREDUCE

# Develop your data processing application

The screenshot shows a list of articles under the 'Articles & Tutorials' section. The articles include:

- Run Search and Sort Job on Amazon Elastic MapReduce**: An article about creating and running search and sort jobs on Amazon EMR.
- Common Administration using Apache Hadoop and Amazon EMR**: An article about common administration tasks on Amazon EMR.
- Logistics for Amazon CloudFront**: An article about using Amazon CloudFront with Amazon CloudWatch Metrics.
- Using Hadoop with Amazon Elastic MapReduce**: An article about using Hadoop with Amazon EMR.
- Process and Analyze Big Data Using Hive on Amazon EMR and Microsoft Big Data**: An article about using Microsoft Big Data services with Amazon EMR.
- Analyze Log Data with Apache Hadoop, Windows PowerShell, and Amazon EMR**: An article about using Windows PowerShell with Amazon EMR to analyze log data.
- Using Cascading MapReduce with EMR**: An article about using Cascading with Amazon EMR.
- Amazon EMR with the MapReduce API for Tez**: An article about using the MapReduce API for Tez with Amazon EMR.
- Processing Streaming Data using Google Beam to Process Data and Apache Hadoop on Amazon Redshift**: An article about using Google Beam with Amazon Redshift.
- Building Stream with the AWS SDK for Java API**: An article about using the AWS SDK for Java API to build a stream.
- Operating on Data Warehouses with EMR, Amazon Redshift, MapReduce, and Amazon DMS**: An article about operating on data warehouses with Amazon Redshift, MapReduce, and Amazon DMS.
- Finding Similar Items with Amazon Elastic MapReduce, Python, and Amazon Swami**: An article about finding similar items using Amazon Swami.
- Word Count Example**: An example showing how to use Amazon Swami to count the number of lines that words occur within a text file.
- CloudFront**: An article about using CloudFront with Amazon Swami.

<http://aws.amazon.com/articles/Elastic-MapReduce>

Develop your data processing application

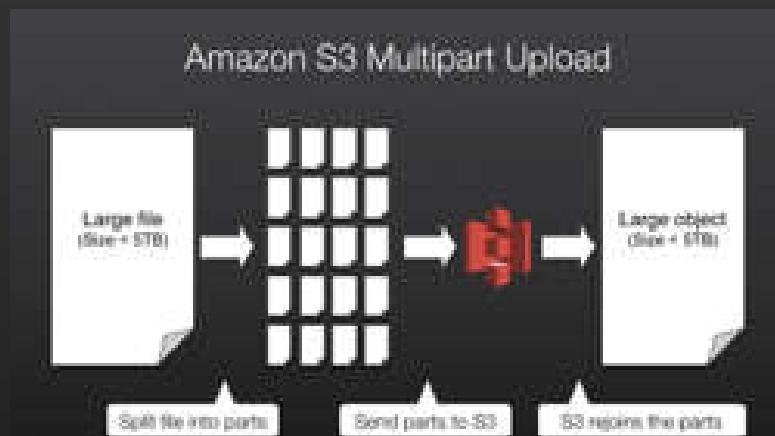


Upload your application and data to Amazon S3

Develop your data processing application



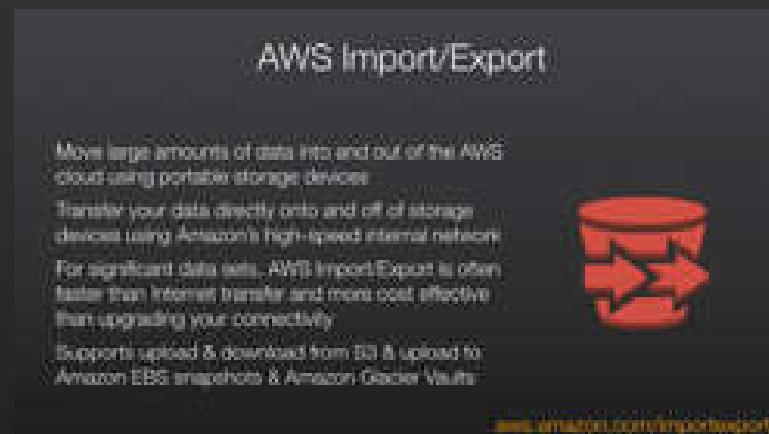
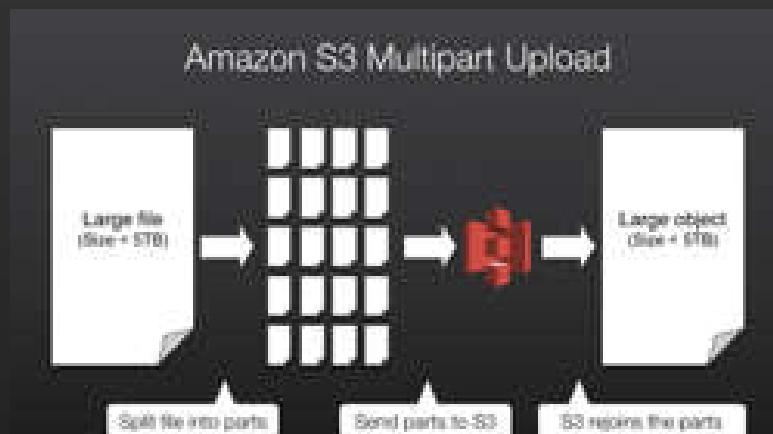
Upload your application and data to Amazon S3



Develop your data processing application



Upload your application and data to Amazon S3



Develop your data processing application



Upload your application and data to Amazon S3



Develop your data processing application



Upload your application and data to Amazon S3



Configure and launch your cluster

Configure and launch your cluster

Amazon EMR Cluster

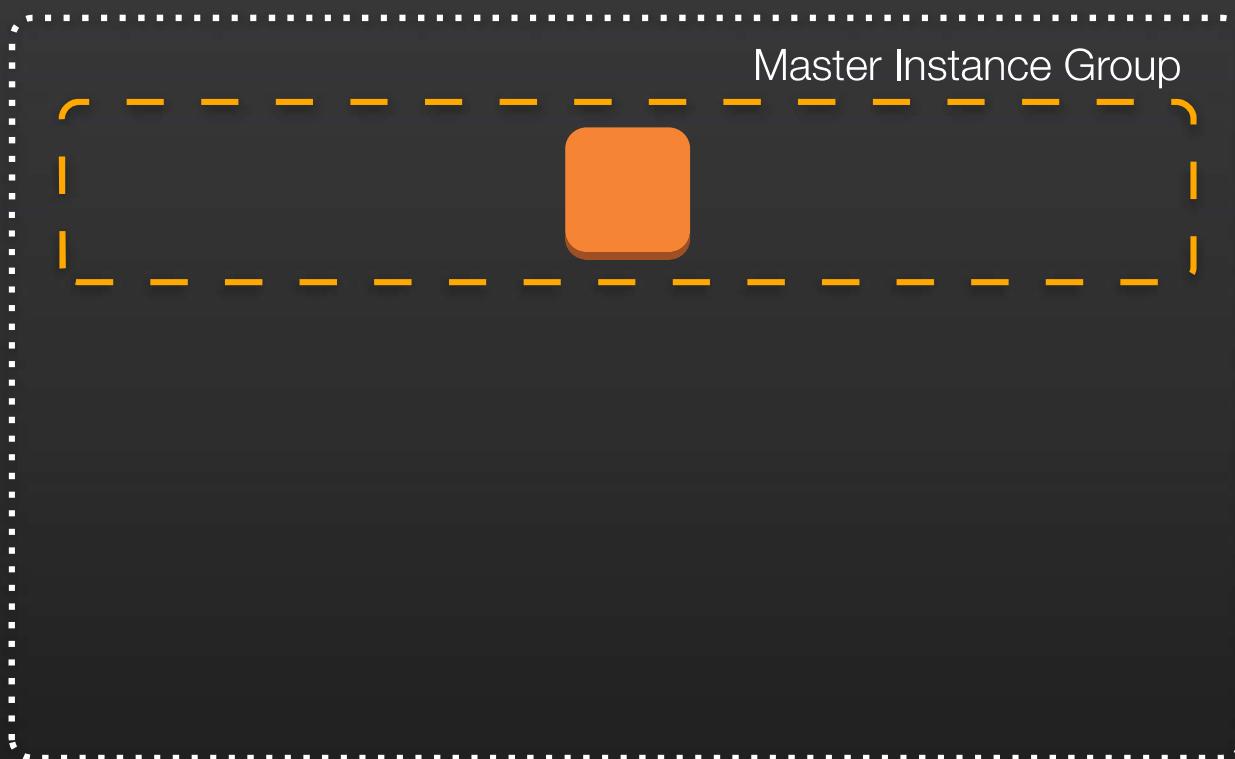
Start an EMR cluster  
using console, CLI tools  
or an AWS SDK

Configure and launch your cluster

Amazon EMR Cluster

Master Instance Group

Master instance group  
created that controls the  
cluster



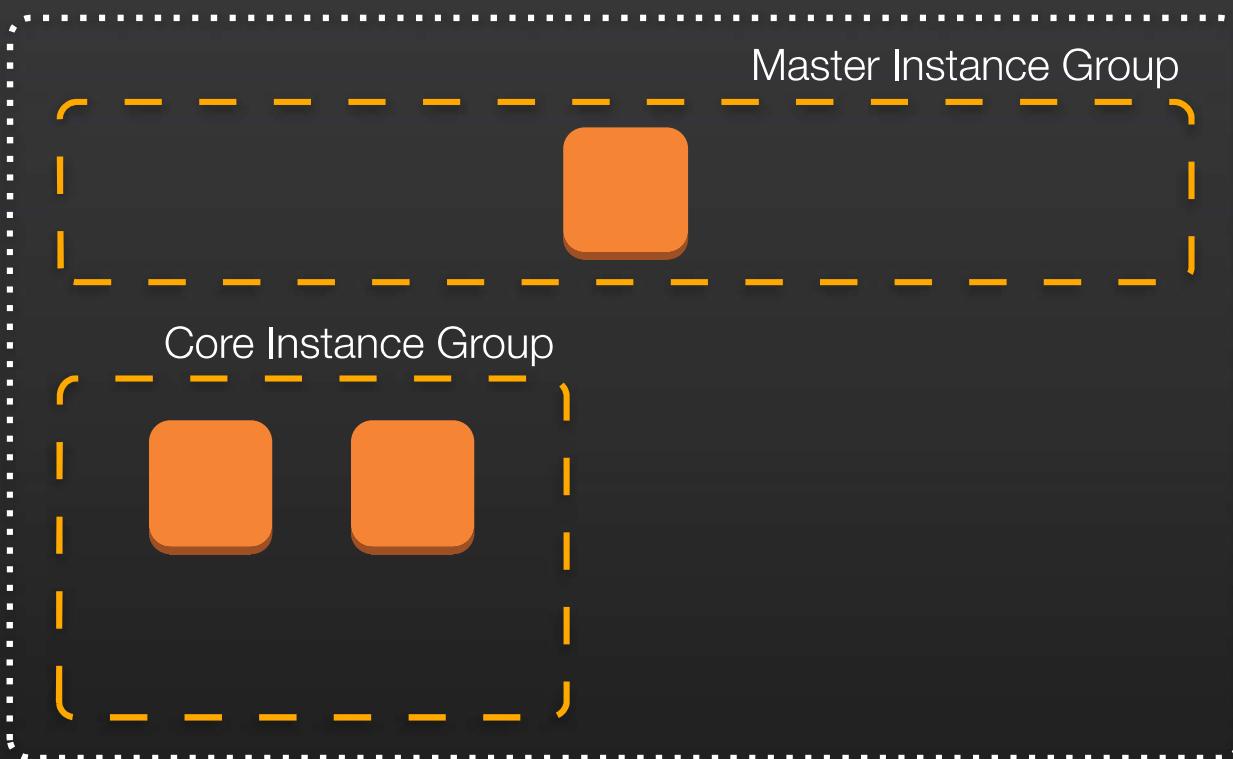
Configure and launch your cluster

Amazon EMR Cluster

Master Instance Group

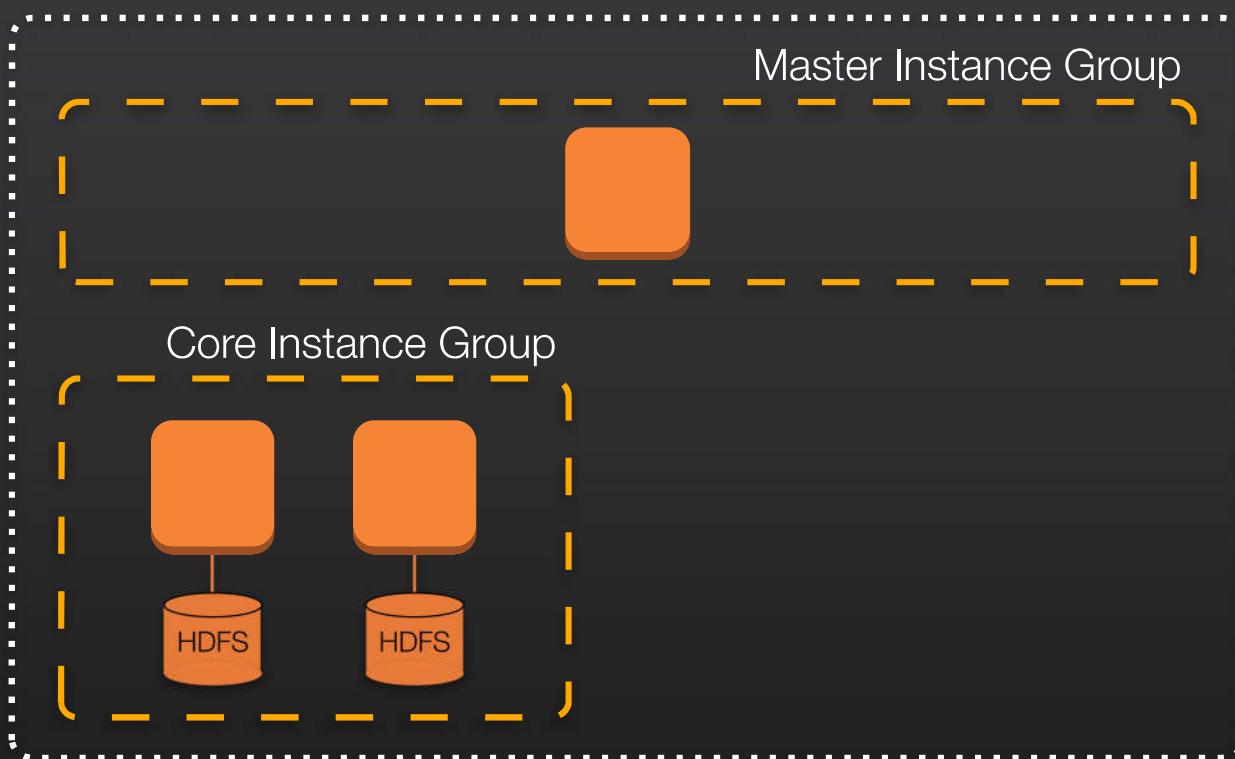
Core Instance Group

Core instance group  
created for life of cluster



Configure and launch your cluster

Amazon EMR Cluster



Core instance group  
created for life of cluster

Core instances run  
DataNode and  
TaskTracker daemons

Configure and launch your cluster

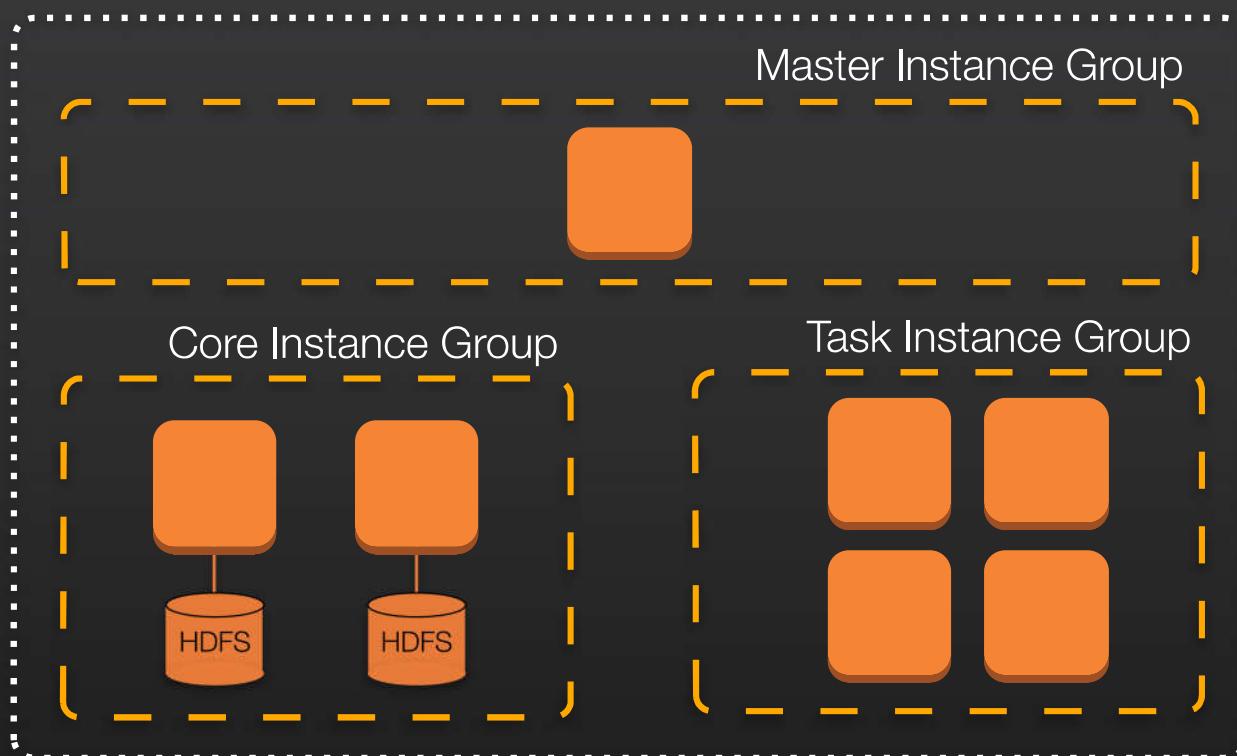
Amazon EMR Cluster

Master Instance Group

Core Instance Group

Task Instance Group

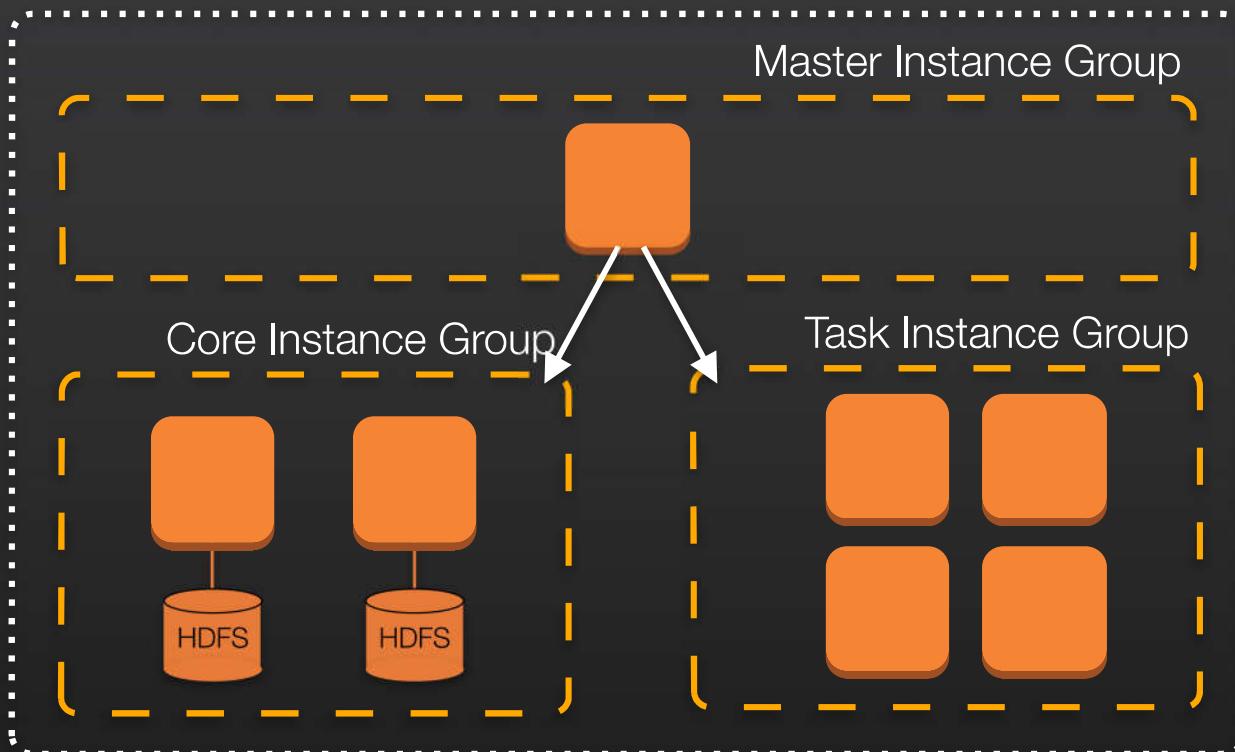
Optional task instances  
can be added or  
subtracted



Configure and launch your cluster

Amazon EMR Cluster

Master node  
coordinates distribution  
of work and manages  
cluster state



Develop your data processing application



Upload your application and data to Amazon S3



Configure and launch your cluster



Optionally, monitor the cluster

Develop your data processing application



Upload your application and data to Amazon S3



Configure and launch your cluster



Optionally, monitor the cluster



Retrieve the output

Retrieve the output

Amazon EMR Cluster

Master Instance Group

Core Instance Group

Task Instance Group

HDFS

HDFS

S3 can be used as  
underlying ‘file system’  
for input/output data



**DEMO:**

**GETTING STARTED WITH  
AMAZON EMR USING A SAMPLE  
HADOOP STREAMING APPLICATION**

# Hadoop Streaming

Utility that comes with the Hadoop distribution

Allows you to create and run Map/Reduce jobs with any executable or script as the mapper and/or the reducer

Reads the input from standard input and the reducer outputs data through standard output

By default, each line of input/output represents a record with tab separated key/value

# Job Flow for Sample Application

Steps			
<p><span>Info</span> A step is a unit of work you submit to the cluster. A step might contain one or more Hadoop jobs, or contain instructions to install or configure an application. You can submit up to 256 steps to a cluster. <a href="#">Learn more</a></p>			
Name	Action on failure	JAR location	Arguments
Word count	Terminate cluster	/home/hadoop/contrib/streaming/hadoop-streaming.jar	<code>-files s3://eu-west-1.elasticmapreduce/samples/wordcount/wordSplitter.py -mapper wordSplitter.py -reducer aggregate -input s3://eu-west-1.elasticmapreduce/samples/wordcount/input -output s3://lambda-aws-emr/intermediate/</code>
Streaming program	Terminate cluster	/home/hadoop/contrib/streaming/hadoop-streaming.jar	<code>-mapper /bin/cat -reducer org.apache.hadoop.mapred.lib.IdentityReducer -input s3://lambda-aws-emr/intermediate/ -output s3://lambda-aws-emr/output -jobconf mapred.reduce.tasks=1</code>

# Job Flow: Step 1

JAR location: /home/hadoop/contrib/streaming/hadoop-streaming.jar

Arguments:

```
-files s3://eu-west-1.elasticmapreduce/samples/wordcount/wordSplitter.py  
-mapper wordSplitter.py  
-reducer aggregate  
-input s3://eu-west-1.elasticmapreduce/samples/wordcount/input  
-output s3://ianmas-aws-emr/intermediate/
```

# Step 1: mapper: wordSplitter.py

```
#!/usr/bin/python
import sys
import re

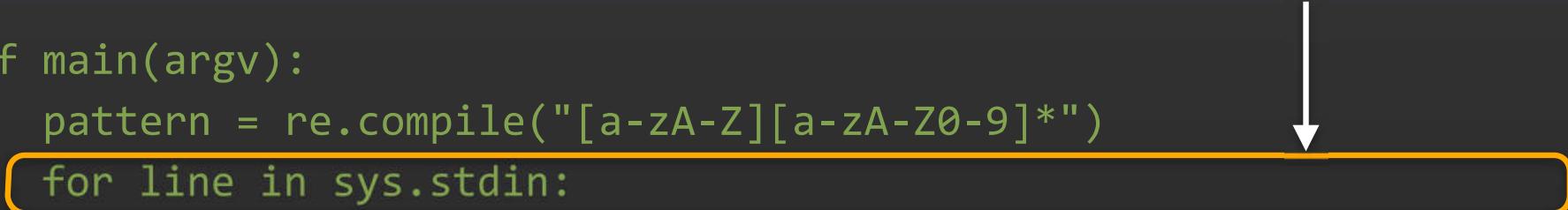
def main(argv):
    pattern = re.compile("[a-zA-Z][a-zA-Z0-9]*")
    for line in sys.stdin:
        for word in pattern.findall(line):
            print "LongValueSum:" + word.lower() + "\t" + "1"

if __name__ == "__main__":
    main(sys.argv)
```

# Step 1: mapper: wordSplitter.py

```
#!/usr/bin/python
import sys
import re
def main(argv):
    pattern = re.compile("[a-zA-Z][a-zA-Z0-9]*")
    for line in sys.stdin:
        for word in pattern.findall(line):
            print "LongValueSum:" + word.lower() + "\t" + "1"
if __name__ == "__main__":
    main(sys.argv)
```

Read words from StdIn line by line



# Step 1: mapper: wordSplitter.py

```
#!/usr/bin/python
import sys
import re

def main(argv):
    pattern = re.compile("[a-zA-Z][a-zA-Z0-9]*")
    for line in sys.stdin:
        for word in pattern.findall(line):
            print "LongValueSum:" + word.lower() + "\t" + "1"

if __name__ == "__main__":
    main(sys.argv)
```

Output to StdOut tab delimited records  
in the format “LongValueSum:abacus 1”



# Step 1: reducer: aggregate

Sorts inputs and adds up totals:

“Abacus 1”

“Abacus 1”

“Abacus 1”

becomes

“Abacus 3”

# Step 1: input/output

The input is all the objects in the S3 bucket/prefix:

`s3://eu-west-1.elasticmapreduce/samples/wordcount/input`

Output is written to the following S3 bucket/prefix to be used as input for the next step in the job flow:

`s3://ianmas-aws-emr/intermediate/`

One output object is created for each reducer (generally one per core)

# Job Flow: Step 2

JAR location: /home/hadoop/contrib/streaming/hadoop-streaming.jar

Arguments:

Accept anything and return as text

```
-mapper /bin/cat  
-reducer org.apache.hadoop.mapred.lib.IdentityReducer  
-input s3://ianmas-aws-emr/intermediate/  
-output s3://ianmas-aws-emr/output  
-jobconf mapred.reduce.tasks=1
```



# Job Flow: Step 2

JAR location: /home/hadoop/contrib/streaming/hadoop-streaming.jar

Arguments:

```
-mapper /bin/cat  
-reducer org.apache.hadoop.mapred.lib.IdentityReducer  
-input s3://ianmas-aws-emr/intermediate/  
-output s3://ianmas-aws-emr/output  
-jobconf mapred.reduce.tasks=1
```

Sort  
↓

# Job Flow: Step 2

JAR location: /home/hadoop/contrib/streaming/hadoop-streaming.jar

Arguments:

- mapper /bin/cat
- reducer org.apache.hadoop.mapred.lib.IdentityReducer
- input s3://ianmas-aws-emr/intermediate/**
- output s3://ianmas-aws-emr/output
- jobconf mapred.reduce.tasks=1

Take previous output as input



# Job Flow: Step 2

JAR location: /home/hadoop/contrib/streaming/hadoop-streaming.jar

Arguments:

- mapper /bin/cat
- reducer org.apache.hadoop.mapred.lib.IdentityReducer
- input s3://ianmas-aws-emr/intermediate/
- output s3://ianmas-aws-emr/output**
- jobconf mapred.reduce.tasks=1

Output location



# Job Flow: Step 2

JAR location: /home/hadoop/contrib/streaming/hadoop-streaming.jar

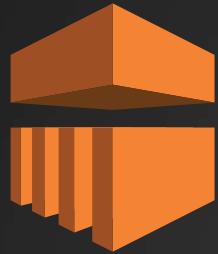
Arguments:

- mapper /bin/cat
- reducer org.apache.hadoop.mapred.lib.IdentityReducer
- input s3://ianmas-aws-emr/intermediate/
- output s3://ianmas-aws-emr/output
- jobconf mapred.reduce.tasks=1**

Use a single reduce task  
to get a single output object



# SUPPORTED HADOOP TOOLS



# Supported Hadoop Tools

Hive



An open source data warehouse & analytics package that runs on top of Hadoop. Operated by Hive QL, a SQL-based language which allows users to structure, summarize, and query data.

Pig

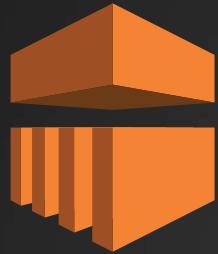


An open source analytics package that runs on top of Hadoop. Pig is operated by Pig Latin, a SQL-like language which allows users to structure, summarize, and query data. Allows processing of complex and unstructured data sources such as text documents and log files.

HBase



Provides you an efficient way of storing large quantities of sparse data using column-based storage. HBase provides fast lookup of data because data is stored in-memory instead of on disk. Optimized for sequential write operations, and it is highly efficient for batch inserts, updates, and deletes.



# Supported Hadoop Tools

Impala



A tool in the Hadoop ecosystem for interactive, ad hoc querying using SQL syntax. It uses a massively parallel processing (MPP) engine similar to that found in a traditional RDBMS.

This lends Impala to interactive, low-latency analytics. You can connect to BI tools through ODBC and JDBC drivers.

Presto



An open source distributed SQL query engine for running interactive analytic queries against data sources of all sizes ranging from gigabytes to petabytes.

Hue



An open source user interface for Hadoop that makes it easier to run and develop Hive queries, manage files in HDFS, run and develop Pig scripts, and manage tables.

# **DEMO: APACHE HUE ON EMR**

[aws.amazon.com/blogs/aws/new-apache-spark-on-amazon-emr/](http://aws.amazon.com/blogs/aws/new-apache-spark-on-amazon-emr/)



AWS Official Blog

## New – Apache Spark on Amazon EMR

by Jeff Barr [on 18 JUN 2013] | in [Amazon EMR](#) | [Permalink](#)

My colleague Jon Fritz wrote the guest post below to introduce a powerful new feature for [Amazon EMR](#).

— Jeff

I'm happy to announce that Amazon EMR now supports [Apache Spark](#). Amazon EMR is a web service that makes it easy for you to process and analyze vast amounts of data using applications in the Hadoop ecosystem, including Hive, Pig, HBase, Presto, Impala, and others. We're delighted to officially add Spark to this list. Although many customers have previously been installing Spark using custom scripts, you can now launch an Amazon EMR cluster with Spark directly from the Amazon EMR Console, CLI, or API.

### [Apache Spark: Beyond Hadoop MapReduce](#)

We have seen great customer successes using Hadoop MapReduce for large scale data processing, batch reporting, ad hoc analysis on unstructured data, and machine learning. Apache Spark, a newer distributed processing framework in the Hadoop ecosystem, is also proving to be an enticing engine by increasing job performance and development velocity for certain workloads.

By using a directed acyclic graph (DAG) execution engine, Spark can create a more efficient query plan for data transformations. Also, Spark uses in-memory, fault-tolerant, resilient distributed datasets (RDDs), keeping intermediates, inputs, and outputs in memory instead of on disk. These two elements of functionality can result in better performance for certain workloads when compared to Hadoop MapReduce, which will force jobs into a sequential map-reduce framework and incur an IO cost from writing intermediates out to disk. Spark's performance enhancements are particularly applicable for iterative workloads, which are common in machine learning and low-latency querying use cases.

Additionally, Spark natively supports Scala, Python, and Java APIs, and it includes libraries for SQL, popular machine learning algorithms, graph processing, and stream processing. With many tightly integrated development options, it can be easier to create and maintain applications for Spark than to work with the various abstractions wrapped around the Hadoop MapReduce API.

### [Introducing Spark on Amazon EMR](#)

Today, we are introducing support for [Apache Spark](#) in Amazon EMR. You can quickly and easily create scalable, managed Spark

clusters on a variety of [Amazon Elastic Compute Cloud \(EC2\)](#) instance types from the Amazon EMR console, [AWS Command Line](#)

and [command-line tools](#), or even directly from the AWS Management Console. To learn more about how to get started with Spark on Amazon EMR, see the [Spark on Amazon EMR Getting Started Guide](#).



# Create a Cluster with Spark

```
$ aws emr create-cluster --name "Spark cluster" \  
  --ami-version 3.8 --applications Name=Spark \  
  --ec2-attributes KeyName=myKey --instance-type m3.xlarge \  
  --instance-count 3 --use-default-roles
```

```
$ ssh -i myKey hadoop@masternode
```

invoke the spark shell with

```
$ spark-shell
```

or

```
$ pyspark
```

pyspark

# Working with the Spark Shell

Counting the occurrences of a string a text file stored in Amazon S3 with spark

```
$ pyspark
>>> sc
<pyspark.context.SparkContext object at 0x7fe7e659fa50>
>>> textfile = sc.textFile("s3://elasticmapreduce/samples/hive-ads/tables/impressions/
dt=2009-04-13-08-05/ec2-0-51-75-39.amazon.com-2009-04-13-08-05.log")
>>> linesWithCartoonNetwork = textfile.filter(lambda line: "cartoonnetwork.com" in
line).count()
15/06/04 17:12:22 INFO lzo.GPLNativeCodeLoader: Loaded native gpl library from the
embedded binaries
<snip>
<Spark program continues>
>>> linesWithCartoonNetwork
9
```

# **ADDITIONAL EMR FEATURES**

# CONTROL NETWORK ACCESS TO YOUR EMR CLUSTER

Using SSH local port forwarding

```
ssh -i EMRKeyPair.pem -N \  
-L 8160:ec2-52-16-143-78.eu-west-1.compute.amazonaws.com:8888 \  
hadoop@ec2-52-16-143-78.eu-west-1.compute.amazonaws.com
```

# MANAGE USERS, PERMISSIONS AND ENCRYPTION

## File System Configuration

**i** The [EMR File System \(EMRFS\)](#) and the Hadoop Distributed File System (HDFS) are both installed on your EMR cluster. HDFS stores data on an EMR cluster, while EMRFS allows EMR clusters to store data on S3. You can enable [S3 server-side encryption](#) or [S3 client-side encryption](#) and [consistent view](#) for EMRFS below, or use a bootstrap action to configure additional settings for EMRFS.

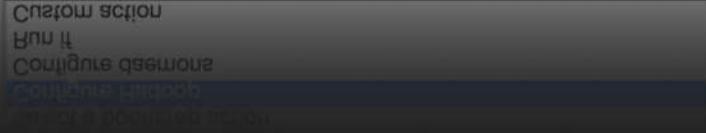
EMRFS S3 Encryption	<input type="radio"/> None <input type="radio"/> None <input checked="" type="radio"/> S3 server-side encryption <input type="radio"/> S3 client-side encryption with AWS Key Management Service (KMS) <input type="radio"/> S3 client-side encryption with custom encryption materials provider	Choose encryption method for objects written to or read from S3 using EMRFS. Please note that this will not encrypt files written to HDFS. <a href="#">Learn more</a>
Consistent view		Monitors list and read-after-write (for new puts) consistency for files in S3. <a href="#">Learn more</a>

# INSTALL ADDITIONAL SOFTWARE WITH BOOTSTRAP ACTIONS

## Bootstrap Actions

Bootstrap actions are scripts that are executed during setup before Hadoop starts on every cluster node. You can use them to install additional software and customize your applications. [Learn more](#)

Bootstrap action type	Name	S3 location	Optional arguments
Add bootstrap action	<input type="text" value="Custom action"/> 	<input type="text" value="Select a bootstrap action"/> 	
		<ul style="list-style-type: none"><li><input checked="" type="checkbox"/> Configure Hadoop</li><li><input type="checkbox"/> Configure daemons</li><li><input type="checkbox"/> Run if</li><li><input type="checkbox"/> Custom action</li></ul>	



# EFFICIENTLY COPY DATA TO EMR FROM AMAZON S3

Run on a cluster master node:

```
$ hadoop jar /home/hadoop/lib/emr-s3distcp-1.0.jar -  
Dmapreduce.job.reduces=30 --src s3://s3bucketname/ --dest hdfs://  
$HADOOP_NAMENODE_HOST:$HADOOP_NAMENODE_PORT/data/ --outputCodec 'none'
```

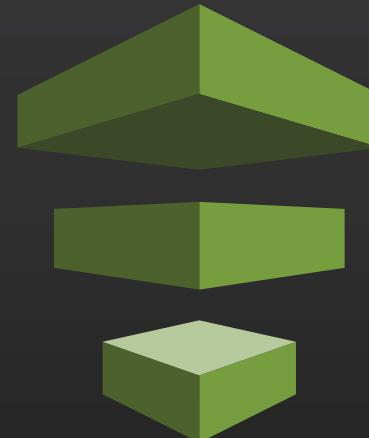
# SCHEDULE RECURRING WORKFLOWS

## AWS Data Pipeline

AWS Data Pipeline is a web service that helps you reliably process and move data between different AWS compute and storage services, as well as on-premise data sources, at specified intervals. With AWS Data Pipeline, you can regularly access your data where it's stored, transform and process it at scale, and efficiently transfer the results to AWS services such as Amazon S3, Amazon RDS, Amazon DynamoDB, and Amazon Elastic MapReduce (EMR).

AWS Data Pipeline helps you easily create complex data processing workloads that are fault tolerant, repeatable, and highly available. You don't have to worry about ensuring resource availability, managing inter-task dependencies, retrying transient failures or timeouts in individual tasks, or creating a failure notification system. AWS Data Pipeline also allows you to move and process data that was previously locked up in on-premise data silos.

soilis stabs esimerkki-ni illä backjoi



[aws.amazon.com/datapipeline](http://aws.amazon.com/datapipeline)

# MONITOR YOUR CLUSTER

[docs.aws.amazon.com/ElasticMapReduce/latest/DeveloperGuide/UsingEMR\\_ViewingMetrics.html](https://docs.aws.amazon.com/ElasticMapReduce/latest/DeveloperGuide/UsingEMR_ViewingMetrics.html)

# DEBUG YOUR APPLICATIONS

Log files generated by EMR Clusters include:

- Step logs
- Hadoop logs
- Bootstrap action logs
- Instance state logs

# USE THE MAPR DISTRIBUTION

## Amazon EMR with the MapR Distribution for Hadoop

Amazon Elastic MapReduce (Amazon EMR) makes it easy to provision and manage Hadoop in the AWS Cloud. Hadoop is available in multiple distributions and Amazon EMR gives you the option of using the Amazon Distribution or the [MapR Distribution](#) for Hadoop.

MapR delivers on the promise of Hadoop with a proven, enterprise-grade platform that supports a broad set of mission-critical and real-time production uses. MapR brings unprecedented dependability, ease-of-use and world-record speed to Hadoop, NoSQL, database and streaming applications in one unified Big Data platform. MapR is used across financial services, retail, media, healthcare, manufacturing, telecommunications and government organizations as well as by leading Fortune 100 and Web 2.0 companies. Investors include Lightspeed Venture Partners, Mayfield Fund, NEA, and Redpoint Ventures. Connect with MapR on [Facebook](#), [LinkedIn](#), and [Twitter](#).



Twitter bus uuluqaklu' lookdesu on Rqsmi nifiv  
Facebook. seilutnuv iulogbaF bus, AEN, bnuU biallyM, israuhpaF  
beedsligl ebuliqin stotsevul: seisndoms 0.5 dev bus 00f enutfo  
gribsel vd ss lew ss snofezzisnglo ihemmvovg bus stotseciuummojelej

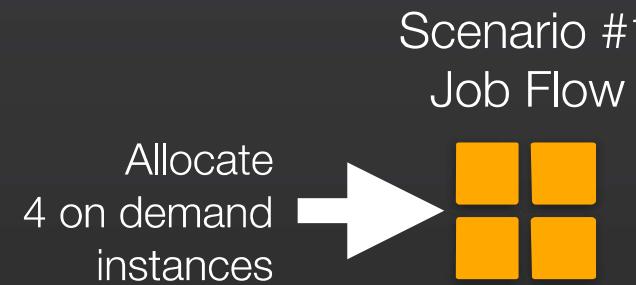
[aws.amazon.com/elasticmapreduce/mapr/](http://aws.amazon.com/elasticmapreduce/mapr/)

# TUNE YOUR CLUSTER FOR COST & PERFORMANCE

Supported EC2 instance types

- General Purpose
- Compute Optimized
- Memory Optimized
- Storage Optimized - D2 instance family   
D2 instances are available in four sizes with 6TB, 12TB, 24TB, and 48TB storage options.
- GPU Instances

# TUNE YOUR CLUSTER FOR COST & PERFORMANCE



Time Savings: 50%  
Cost Savings: ~22%

Duration:

A horizontal progress bar divided into two equal segments. The left segment is blue and labeled "14 hours". The right segment is white. A vertical line is positioned at the center of the bar.

Cost *without* Spot  
4 instances \* 14 hrs \* \$0.50  
Total = \$28



Duration:

A horizontal progress bar divided into two equal segments. The left segment is blue and labeled "7 hours". The right segment is white. A vertical line is positioned at the center of the bar.

Cost *with* Spot  
4 instances \* 7 hrs \* \$0.50 = \$14 +  
5 instances \* 7 hrs \* \$0.25 = \$8.75  
Total = \$22.75

# THIRD PARTY TOOLS

**MicroStrategy**

**MAPR**

**Datameer**

**ATTUNITY**

BI/Visualization

Hadoop Distribution

Graphical IDE

Data Transfer

**MORTAR**

**SAP**  
Business Objects™

**bndy**

**JASPERSOFT**  
THE INTELLIGENCE EDGE

Integration and Analytics

Business Intelligence

Monitoring

BI/Visualization

**talend\***  
open data solutions

**splunk>**

**Compuware**

**+tableau**

Graphical IDE

Data Exploration

Performance Tuning

BI/Visualization

Graphical IDE

Data Exploration

Performance Tuning

BI/Visualization

**RESOURCES YOU CAN USE  
TO LEARN MORE**

[aws.amazon.com/emr](http://aws.amazon.com/emr)

Getting Started with Amazon EMR Tutorial guide:

[docs.aws.amazon.com/ElasticMapReduce/latest/DeveloperGuide/emr-get-started.html](https://docs.aws.amazon.com/ElasticMapReduce/latest/DeveloperGuide/emr-get-started.html)

Customer Case Studies for Big Data Use-Cases

[aws.amazon.com/solutions/case-studies/big-data/](https://aws.amazon.com/solutions/case-studies/big-data/)

Amazon EMR Documentation:

[aws.amazon.com/documentation/emr/](https://aws.amazon.com/documentation/emr/)