

# Analytics and Visualization of Big Data

*Fadel M. Megahed*

## ***Lecture 04: Map-Reduce and an Introduction to AWS***



**AUBURN UNIVERSITY**

SAMUEL GINN  
COLLEGE OF ENGINEERING

***Department of Industrial and Systems Engineering***

*Spring 13*

## Chapter 02: Map-Reduce (Cont.)



# What is a Distributed File System?

- Motivation for a Distributed File System
- Physical Organization of Compute Nodes
- Large-Scale File-System Organization

# Map-Reduce

- **The Map Tasks**
- **Grouping and Aggregation**
- **The Reduce Tasks**
- **Combiners**
- **Coordination**
- **Coping with Node Failures**



## Motivation Example: Reading 5M Books with a Program 😊 4

- Counting the number of times each distinct word or phrase appears in a large collection of documents.

Source: [http://www.ted.com/playlists/56/making\\_sense\\_of\\_too\\_much\\_data.html](http://www.ted.com/playlists/56/making_sense_of_too_much_data.html)

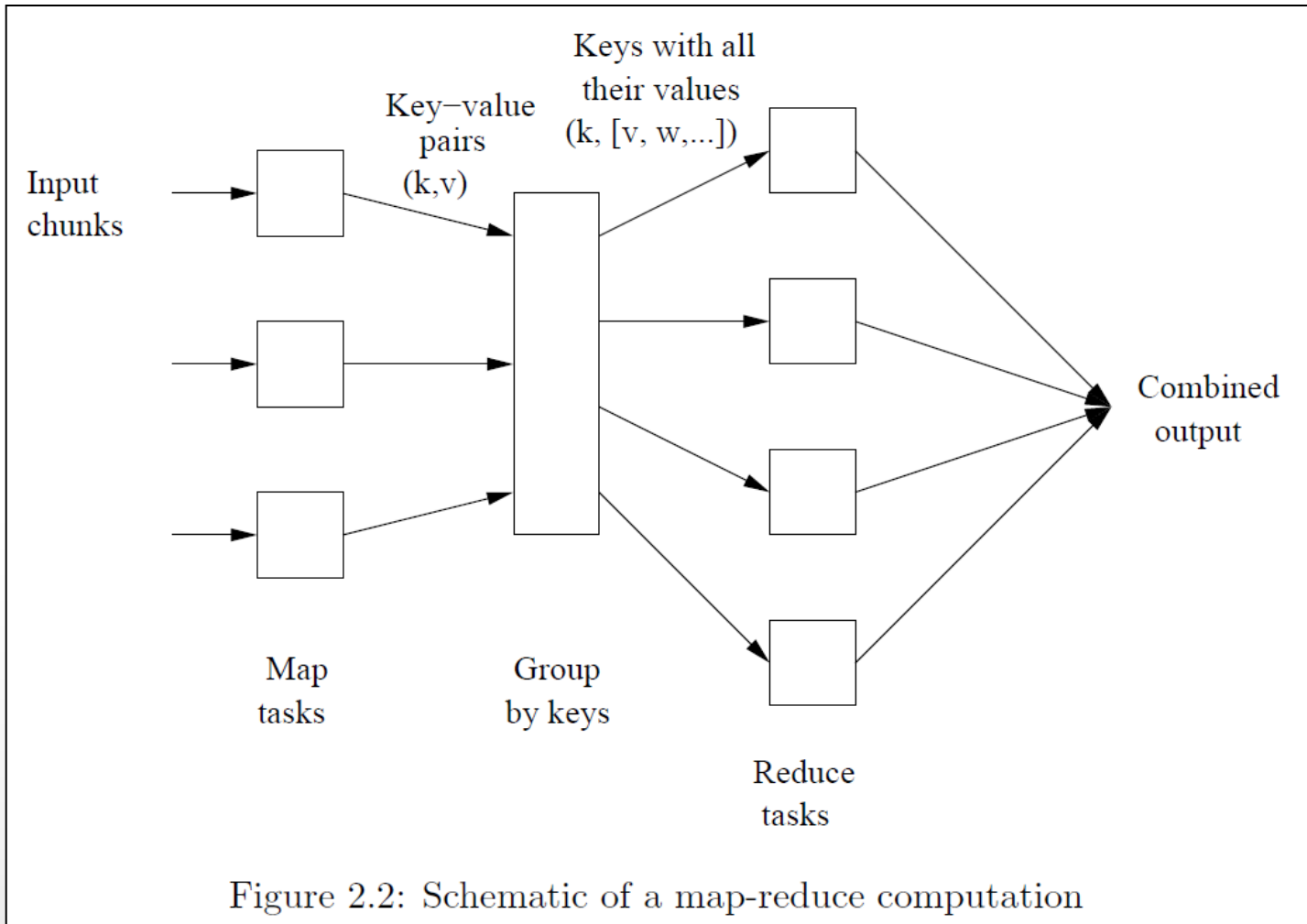


- It is a style of computing that has been implemented in several systems:
  - Google's own implementation → map-reduce
  - Several popular open-source implementations → Hadoop
- Map-reduce can manage large-scale computations in a way that is tolerant to hardware failures.
- All you need to write are two functions **Map** and **Reduce**, while the system:
  - Manages the parallel execution
  - Coordinates the tasks that execute Map and Reduce
  - Deals with the possibility that one of the tasks will fail



# An Overview of Map-Reduce Programming

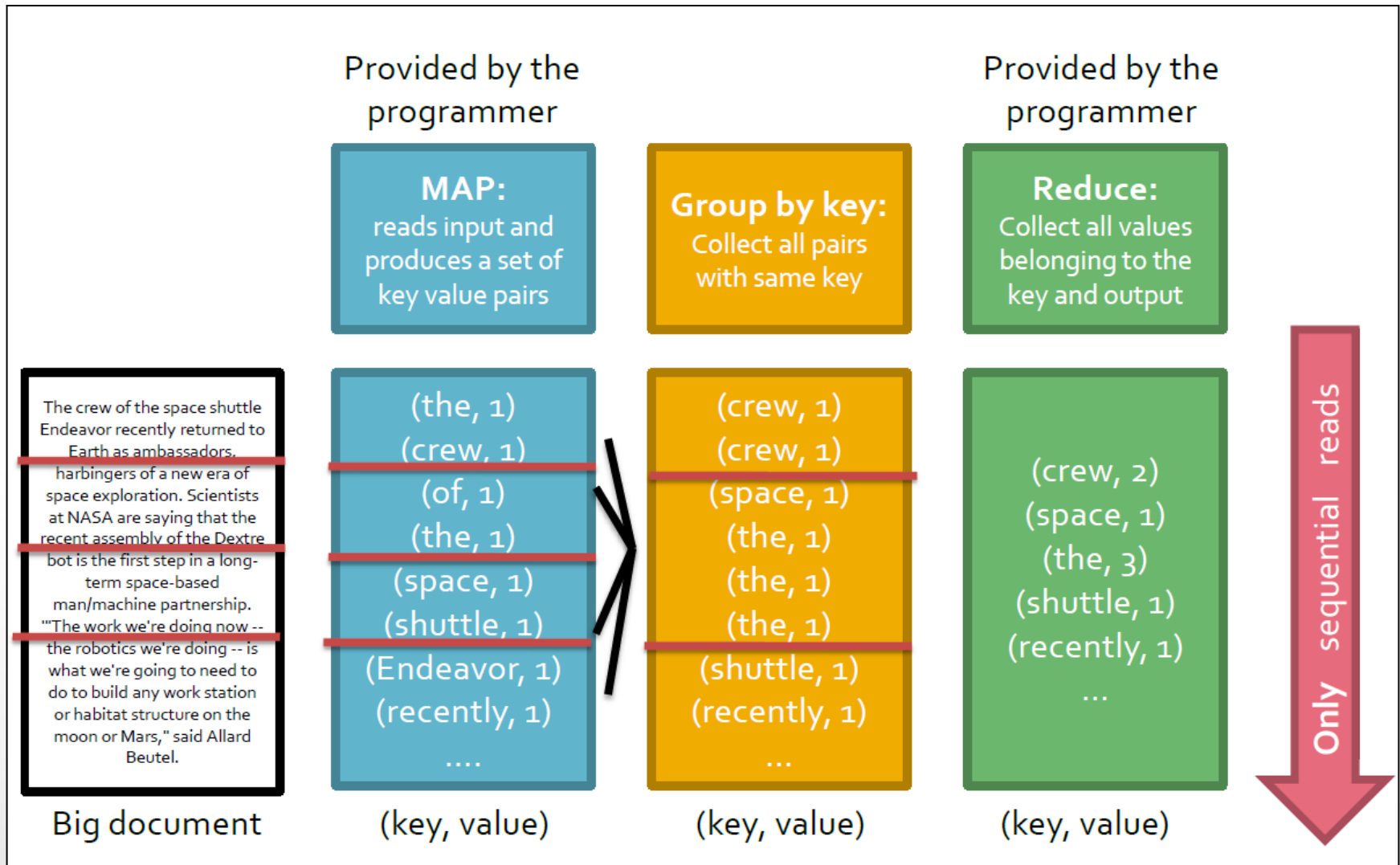
6



Source: A. Rajaraman, J. Leskovec, J.D. Ullman. (2012). "Mining of Massive Datasets". <http://i.stanford.edu/~ullman/mmds.html>

# The Use of Map-Reduce for Word Counting

7



Source: Figure from Jure Leskovic, Stanford CS246, Lecture Notes, see <http://cs246.stanford.edu>

**Map-Reduce environment takes care of:**

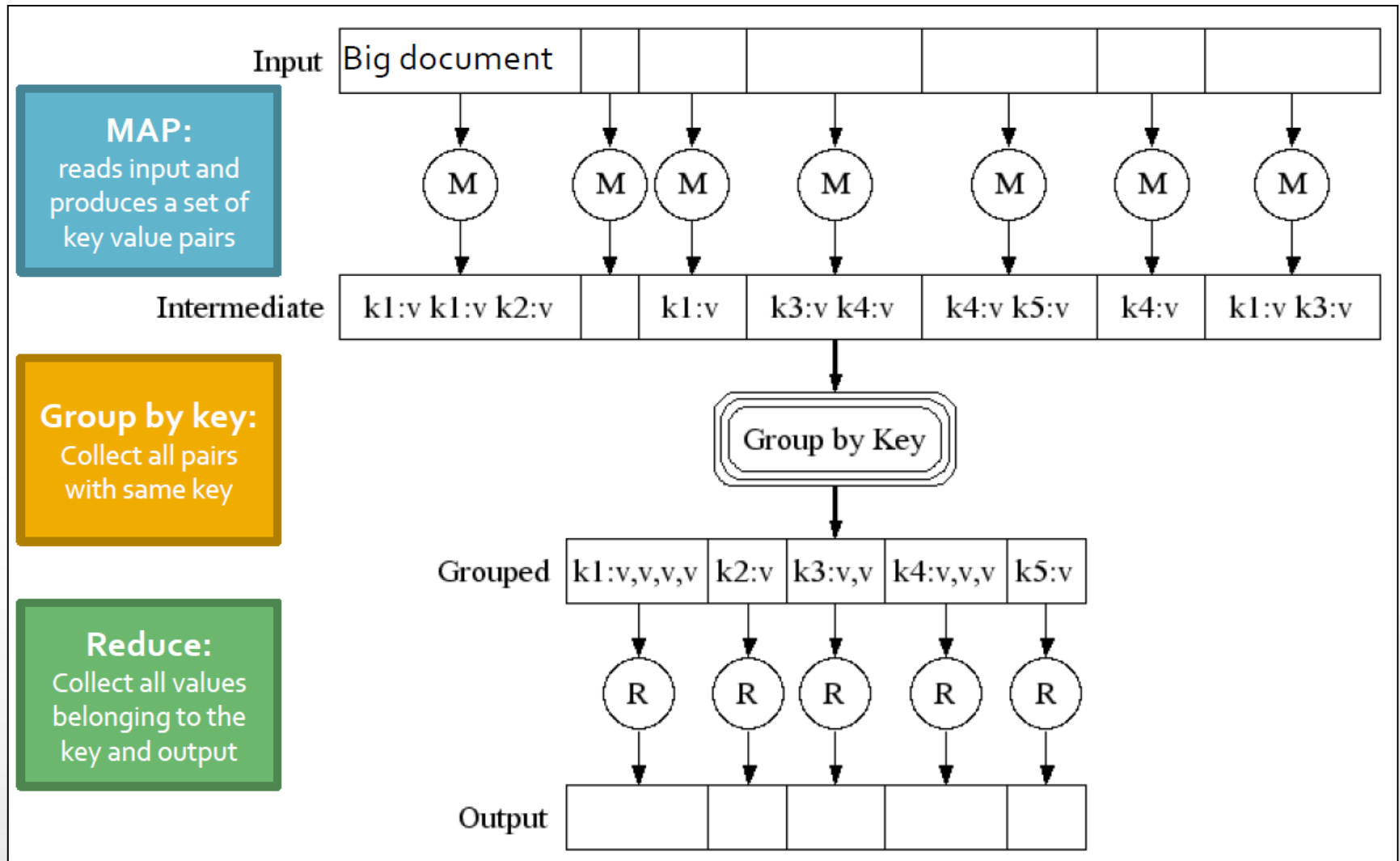
- **Partitioning** the input data
- **Scheduling** the program's execution across a set of machines
- Handling machine **failures**
- Managing required inter-machine **communication**





# Map-Reduce: Another Representation

9



Source: Figure from Jure Leskovic, Stanford CS246, Lecture Notes, see <http://cs246.stanford.edu>

- **Input and final output** are stored on a **distributed file system (DFS)**:
  - Scheduler tries to schedule map tasks “close” to physical storage location of input data
- **Intermediate results** are stored on **local FS** of map and reduce workers
- Output is often input to another map reduce task

Source: Slide is adapted from Jure Leskovic, Stanford CS246, Lecture Notes, see <http://cs246.stanford.edu>



# Coordination – How compute nodes, tasks and files interact?

- Normally a compute node (worker) handles either map tasks or reduce tasks but not both.
- The master node has many responsibilities:
  - Task status: (idle, in-progress, completed)
  - Idle tasks get scheduled as workers become available
  - When a map task completes, it sends the master the location and sizes of its R intermediate files, one for each reducer
  - Master pushes this info to reducers
  - Master pings workers periodically to check for compute node failures



- **Map worker failure**

- Map tasks completed or in-progress at worker are reset to idle
- Reduce workers are notified when task is rescheduled on another worker

- **Reduce worker failure**

- Only in-progress tasks are reset to idle

- **Master failure**

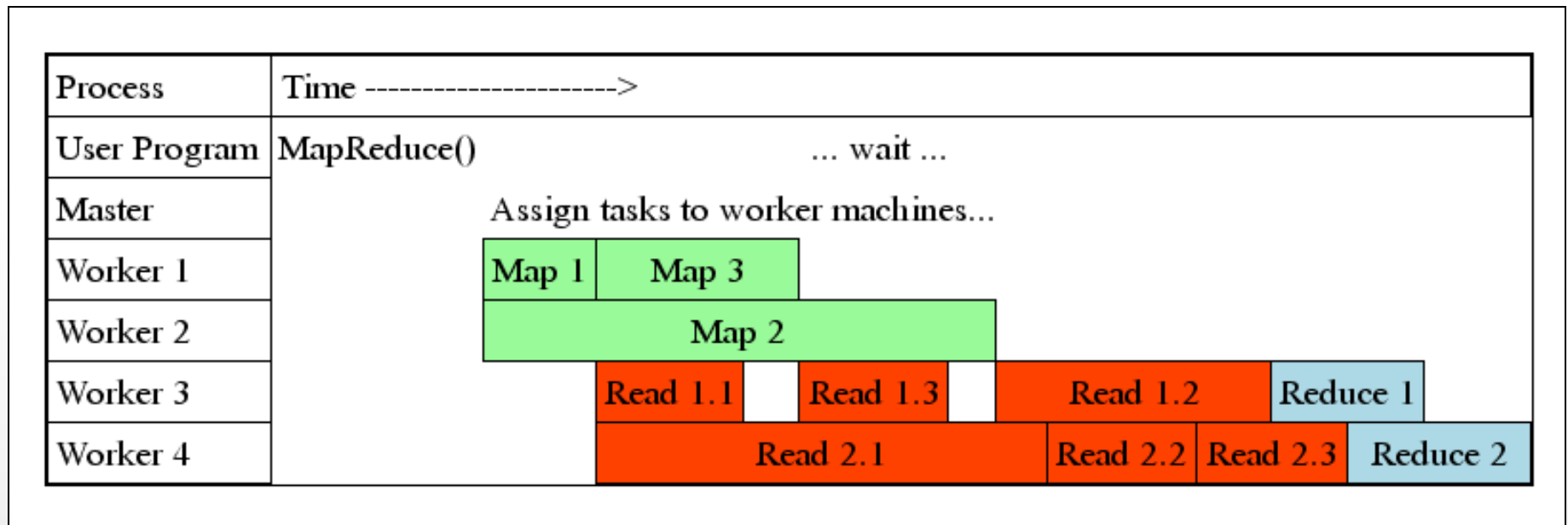
- Map-Reduce task is aborted and client is notified (In case of one master node)



# Similarities of Map-Reduce to MRP???

13

- **Fine granularity tasks:** map tasks  $\gg$  machines
  - Minimizes time for fault recovery
  - Can pipeline shuffling with map execution
  - Better dynamic load balancing



Source: Slide is adapted from Jure Leskovic, Stanford CS246, Lecture Notes, see <http://cs246.stanford.edu>

# **An Introduction to Amazon Web Services**



- AWS) is an '*Infrastructure as a Service*' (IaaS) provider.
- AWS' cloud based computing offers dynamically scalable computing, storage, and data access “on demand” over the Internet, with “pay as you go” pricing for the hardware and software that is delivered.
- AWS includes over two dozen cloud-related services, including their popular elastic computing (EC2) and storage (S3) capabilities. Broadly speaking, their cloud services can be grouped into 3 categories, infrastructure services, platform services and supporting services.

Source: <http://aws.amazon.com/>





The screenshot shows the Amazon Web Services homepage. At the top left is the AWS logo. To its right are links for "Sign Up", "My Account / Console", and "English". Below the logo is a navigation bar with "AWS Products & Solutions", a search bar, "Entire Site", "Developers", and "Support". The main banner is titled "Still Managing Your Own Relational Database?". It features logos for MySQL, Microsoft SQL Server, and ORACLE. The text reads: "Amazon RDS makes it easier for you to set up, manage, and scale a relational database in the cloud." Below this is a link: "Discover how Amazon RDS can help your business today »". A diagram shows several server icons connected to a central database icon. At the bottom of the banner is a button "Get Started for Free »" and the text "Pay only for what you use."

amazon web services

Sign Up

My Account / Console

English

AWS Products & Solutions

Entire Site

Developers

Support

## Still Managing Your Own Relational Database?

Amazon RDS makes it easier for you to set up, manage, and scale a relational database in the cloud.

MySQL Microsoft SQL Server ORACLE

Discover how Amazon RDS can help your business today »

Get Started for Free » Pay only for what you use.

Source: <http://aws.amazon.com/>



# Analytics and Visualization of Big Data

*Fadel M. Megahed*

## ***Lecture 04: Map-Reduce and an Introduction to AWS***



**AUBURN UNIVERSITY**

SAMUEL GINN  
COLLEGE OF ENGINEERING

***Department of Industrial and Systems Engineering***

*Spring 13*