# Getting Started with AWS

## Analyzing Big Data

# Getting Started with AWS: Analyzing Big Data

# Table of Contents

# Analyzing Big Data with Amazon Web Services

The following tutorials show you ways to use Amazon Web Services to process big data:

## Key AWS Services for Big Data

With AWS, you pay only for the resources you use. Instead of maintaining a cluster of physical servers and storage devices that are standing by for possible use, you can create resources when you need them. AWS also supports popular tools like Hadoop, Hive, and Pig, and makes it easy to provision, configure, and monitor clusters for running those tools.

The following table shows how AWS can help address common big-data challenges.

| Challenge | Solution |
|---|---|
| Data sets can be very large. Storage can become expensive, and data corruption and loss can have far-reaching implications. | Amazon S3 can store large amounts of data, and its capacity can grow to meet your needs. It is highly redundant and secure, protecting against data loss and unauthorized use. Amazon S3 also has an intentionally small feature set to keep its costs low. |
| Maintaining a cluster of physical servers to process data is expensive and time-consuming. | When you run an application on a virtual Amazon EC2 server, you pay for the server only while the application is running, and you can increase the number of servers — within minutes, not hours or days — to meet the processing needs of your application. |