

# Cloudera Enterprise Reference Architecture for AWS deployments



## Table of Content

Cloudera on AWS	3
Amazon Web Services Overview	4
Elastic Compute Cloud (EC2)	4
Simple Storage Service (S3)	4
Relational Database Service (RDS)	4
Elastic Block Store (EBS)	5
Direct Connect	5
Virtual Private Cloud	5
Deployment Architecture	5
Deployment Topologies	5
Workloads, Roles and Instance types	8
Regions and Availability Zones	10
Networking, connectivity and security	10
Supported AMIs	12
Storage options and configuration	12
Capacity planning	13
Relational Databases	13
Installation and Software Configuration	14
Provisioning instances	14
Preparation	15
Deploying Cloudera Enterprise	15
Cloudera Enterprise configuration considerations	15
Summary	16
References	16
Cloudera Enterprise	16
Amazon Web Services	16

## Abstract

Organizations' requirement for a big data solution is simple: The ability to acquire and combine any amount or type of data in its original fidelity, in one place, for as long as necessary, and deliver insights to all kinds of users, as fast as possible.

Cloudera, an enterprise data management company, introduced the concept of enterprise data hub, a single central system to store and work with all data. The enterprise data hub (EDH) has the flexibility to run a variety of enterprise workloads (i.e. batch processing, interactive SQL, enterprise search, and advanced analytics) while meeting enterprise requirements such as integrations to existing systems, robust security, governance, data protection, and management. The EDH is the emerging and necessary center of enterprise data management. EDH builds on [Cloudera Enterprise](#), which consists of the open source CDH, a suite of management software and enterprise class support.

In addition to needing an enterprise data hub, enterprises are also looking to move or add this powerful data management infrastructure to the Cloud to gain benefits such as operation efficiency, cost reduction, and compute/capacity flexibility, and speeds and agility.

As organizations are looking to embrace Hadoop-powered big data deployments in cloud environments, they also want features such as enterprise-grade security, management tools, and technical support which are a part Cloudera Enterprise.

Customers of Cloudera and Amazon Web Services (AWS) now have the ability to run the enterprise data hub in the AWS public cloud, leveraging the power of the Cloudera Enterprise platform and the flexibility of the AWS cloud together.

## Cloudera on AWS

Cloudera delivers on that objective with Cloudera Enterprise and now makes it possible for organizations to deploy the Cloudera solution as an enterprise data hub in the Amazon Web Services (AWS) cloud. This joint solution combines Cloudera's expertise in large-scale data management and analytics, along with AWS's expertise in cloud computing.

This joint solution offers benefits, including:

[Flexible Deployment, Faster Time to Insight](#) - Running Cloudera Enterprise on AWS provides customers the greatest flexibility in how they deploy Hadoop, and can now bypass prolonged infrastructure selection and procurement processes, to rapidly put Cloudera's Platform for Big Data to work to start realizing tangible business value from their data immediately. Hadoop excels at large scale data management and the AWS cloud focuses on providing infrastructure services on demand. Combining these allows customers to be able to leverage the power of Hadoop much faster and on-demand.

**Scalable Data Management** - At many large organizations, it can take weeks or even months to add new nodes into a traditional data cluster. By deploying Cloudera Enterprise in AWS, enterprises can effectively shorten rest-to-growth cycles to scale their data hubs as their business grows.

**On-demand Processing Power** - While Hadoop focus on collocating compute to disk, there are many processes that benefit from increased compute power. Deploying Hadoop on Amazon allows a fast ramp-up / ramp-down based on the needs of specific workloads, a flexibility that does not come easy with on-premise deployment.

**Improved Efficiency, Increased Cost Savings** - Deploying in AWS eliminates the need for organizations to dedicate resources toward maintaining a traditional data center, enabling them to focus instead on core competencies. As annual data growth for the average enterprise continues to skyrocket, even relatively new data management systems may experience strain under the demands of modern high performance workloads. By moving their data management platform to the cloud, enterprises can now offset or avoid the need to make costly annual investments in their on-premises data infrastructure to support new enterprise data growth, applications and workloads.

In this white paper, we provide an overview of general best practice for running Cloudera on AWS, leveraging different AWS services such as EC2, S3, and RDS

## Amazon Web Services Overview

AWS (Amazon Web Services) is the leading public cloud infrastructure provider. Their offerings consists of several different kinds of services, ranging from storage to compute to services higher up in the stack for things like automated scaling, messaging and queuing etc. For the purpose of Cloudera Enterprise deployments, the following service offerings are relevant:

### Elastic Compute Cloud (EC2)

**Elastic Compute Cloud (EC2)** is a service where end users can rent virtual machines of different configurations on-demand and pay for the amount of time they use them. For this deployment, EC2 instances are the equivalent of servers that run Hadoop. There are several different **types** of instances that EC2 offers, with different **pricing** options. For Cloudera Enterprise deployments, each individual node in the cluster conceptually maps to an individual server. A list of supported instance types and the roles that they play in a Cloudera Enterprise deployment are highlighted later in the document.

### Simple Storage Service (S3)

**Simple Storage Service (S3)** is a storage service which allows users to store and retrieve arbitrary sized data objects using simple API calls. S3 is designed for 99.999999999% durability and 99.99% availability. When using S3, users only get the ability to store data. There is no compute element to it. The compute service is provided by the EC2 service, which is independent of S3.

### Relational Database Service (RDS)

**Relational Database Service (RDS)** is a service which allows users to provision a managed relational database instance. Users can provision different flavors of relational database instances, including Oracle and MySQL. RDS handles database management tasks, such as backups for a user-defined retention period and enabling point-in-time recovery, patch management, and replication, allowing the user to pursue higher value application development or database refinements.

## Elastic Block Store (EBS)

[Elastic Block Store \(EBS\)](#) provides users block level storage volumes that can be used as network attached disks with EC2 instances. Users can provision volumes of different capacities and IOPS guarantees. Unlike S3, these volumes can be mounted as network attached storage to EC2 instances and have an independent persistence lifecycle, i.e. they can be made to persist even after the EC2 instance has been shut down. At a later point, the same EBS volume can be attached to a different EC2 instance. EBS volumes can also be snapshotted to S3 for higher durability guarantees. EBS is primarily optimized for random access patterns.

## Direct Connect

[Direct Connect](#) is the way to establish direct connectivity between your data center to AWS region. You can configure direct connect links with different bandwidths based on your requirement. This service allows you to logically consider AWS infrastructure as an extension to your data center.

## Virtual Private Cloud

[Virtual Private Cloud \(VPC\)](#) gives you the ability to logically cordon off a section of the AWS cloud and provision services inside of that cordoned off network that you define. VPC is the recommended way to provision services inside AWS and is enabled by default for all new accounts. There are different configuration options for VPC. The difference between various options are in method of accessibility to Internet and other AWS services. You can create public facing subnets in VPC, where the instances have the option of having direct access to the public Internet gateway and other AWS services. Instances can be provisioned in private subnets too, where their access to the Internet and other AWS services can be restricted entirely or done via NAT. RDS instances can be accessed from within a VPC.

# Deployment Architecture

## Deployment Topologies

There are two kinds of Cloudera Enterprise deployments supported in AWS, both of which are within VPC but with different accessibility.

1. Cluster inside Public Subnet in VPC
2. Cluster inside Private Subnet in VPC

The choice between the public subnet and private subnet deployments depends predominantly on the accessibility of the cluster, both inbound and outbound and the bandwidth required for outbound access.

## Public Subnet deployments

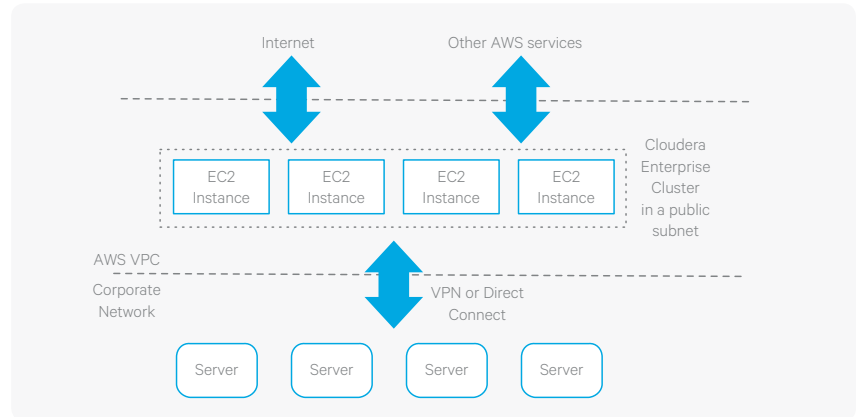
A public subnet in this context is defined as subnet with a route to the Internet gateway. Instances provisioned in public subnets inside VPC can have direct access to the Internet as well as to other AWS services such as RDS and S3. If your requirement is to have the cluster access S3 for data transfers, or ingest from sources on the Internet, your cluster should be deployed in a public subnet. This gives each instance full bandwidth access to the Internet and other AWS services. Unless it's a requirement, we don't recommend opening full access to your cluster from the Internet. The cluster can be configured to have access to other AWS services but not to the Internet. This can be done via security groups (discussed later).

## Private Subnet deployments

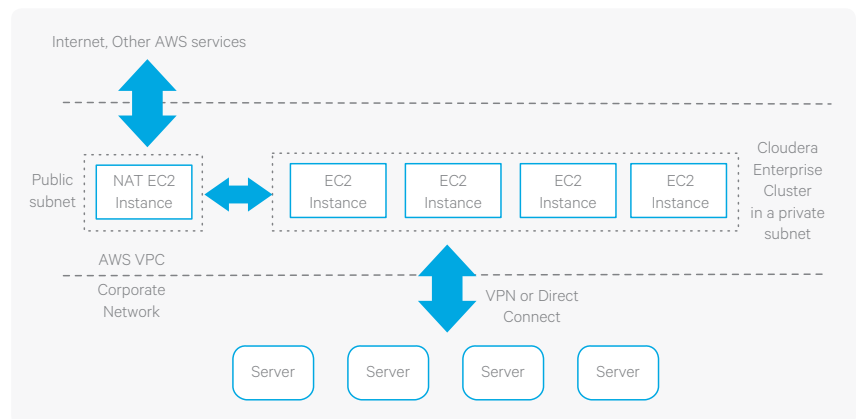
Instances provisioned in private subnets inside VPC don't have direct access to the Internet or to other AWS services. In order to access the Internet, they have to go through a NAT instance in the public subnet. If your cluster does not require full bandwidth access to the Internet or to other AWS services, you should deploy in a private subnet.

In both cases, you can have VPN or Direct Connect setup between your corporate network and AWS. This will make AWS look like an extension to your network and the Cloudera Enterprise deployment will be accessible as if it was on servers in your own data center.

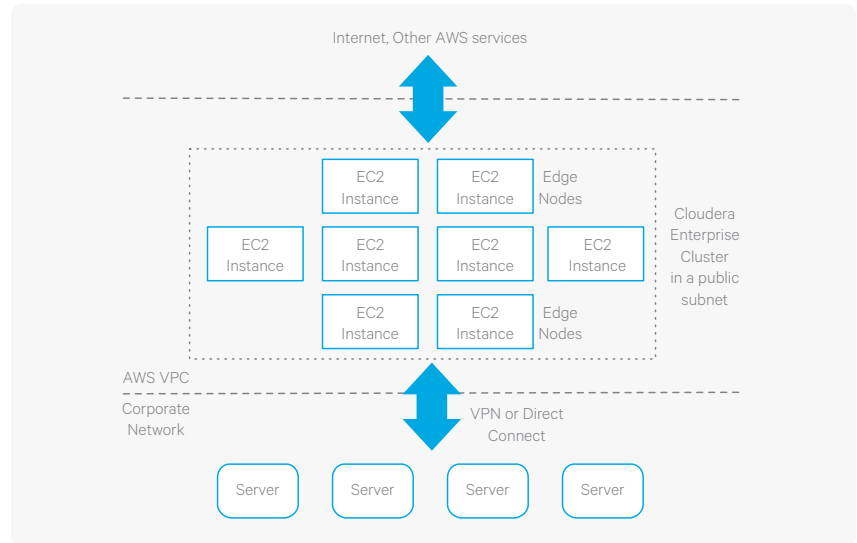
Deployment in the public subnet looks like:



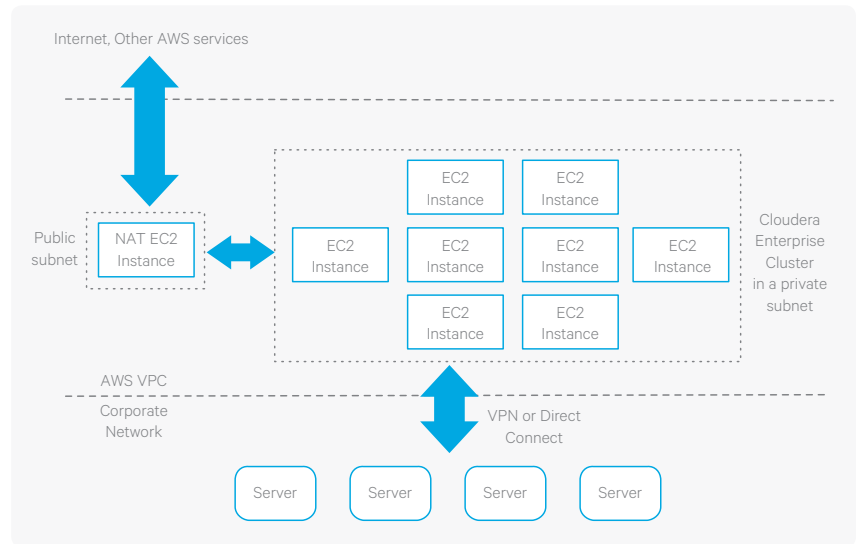
Deployment in the private subnet looks like:



The accessibility of your Cloudera Enterprise cluster is defined by the VPC configuration and depends on the security requirements and the workload. Typically there are edge/client nodes that have direct access to the cluster. Users go through these edge nodes via client applications to interact with the cluster and the data residing there. These edge nodes could be running a web application for real time serving workloads, BI tools, or simply the Hadoop command line client that can be used to submit or interact with HDFS. The public subnet deployment with edge nodes looks like:



Deployment in private subnet with edge nodes looks like:



The edge nodes in case of a private subnet deployment could be in the public subnet, depending on how they have to be accessed. The figure above shows them in the private subnet as one deployment option.

The edge nodes can be EC2 instances in your VPC or servers in your own data center. It's recommended to allow access to the Cloudera Enterprise cluster via edge nodes only. This can be configured in the security groups for the instances that you provision. In the rest of the document, the various options are described in detail.

## Workloads, Roles and Instance types

In this reference architecture, we take into account different kinds of workloads that are run on top of an enterprise data hub and make recommendations on the different kinds of EC2 instances that are suitable for each of these workload types. The recommendations span across new as well as old generation instance types, with storage options including magnetic disks and SSDs. Customers can choose instance types based on the workload they want to run on the cluster. We leave the exercise to do a cost-performance analysis for the customer.

We currently support RHEL 6.4 AMIs, on CDH 4.5+ and CDH5.x.

Matrix of workload categories and services that typically combined for the workload type is as follows:

Workload Type	Typical Services	Comments
Low	MapReduce, YARN, Spark, Hive, Pig, Crunch	Suitable for workloads that are predominantly batch oriented in nature and involved MapReduce or Spark.
Medium	HBase, Solr, Impala	Suitable for higher resource consuming services and production workloads but limited to only one of these running at any time.
High / Full EDH workloads	All CDH services	Full scale production workloads with multiple services running in parallel on a multi-tenant cluster.

## Management Nodes

Management nodes for a Cloudera Enterprise deployment are the ones that run the management services. Management services include:

- > Cloudera Manager
- > JobTracker
- > Standby JobTracker
- > NameNode
- > Standby NameNode
- > JournalNodes
- > HBase Master
- > Zookeeper
- > Oozie

## Worker Nodes

Worker Nodes for a Cloudera Enterprise deployment are the ones that run worker services. These include:

- > DataNode
- > TaskTracker
- > HBase RegionServer
- > Impala Daemons
- > Solr Servers



## Edge Nodes

Edge nodes are where your Hadoop client services are run. They are also known as Gateway Services. These include:

- > Third party tools
- > Hadoop command line client
- > Beeline
- > Impala shell
- > Flume agents
- > Hue Server

Following is a matrix showing the different workload categories, instance types and roles they are suited for in a cluster:

Workload Type	Typical Services	Instances for Management Nodes	Instances for Worker Nodes
Low	MapReduce, YARN, Spark, Hive, Pig, Crunch	> m2.4xlarge	> c3.8xlarge
		> c3.8xlarge	> r3.8xlarge
		> r3.8xlarge	> i2.2xlarge
		> i2.2xlarge	> i2.4xlarge
		> i2.4xlarge	> i2.8xlarge
		> i2.8xlarge	> hs1.8xlarge
		> hs1.8xlarge	> m1.large
		> m1.xlarge	> m1.xlarge
		> m1.large	> c1.xlarge
		> c1.xlarge	> cc2.8xlarge
		> cc2.8xlarge	> m2.4xlarge
		> m2.2xlarge	> hi1.4xlarge
		> hi1.4xlarge	
Medium	HBase, Solr, Impala	> c3.8xlarge	> i2.4xlarge
		> r3.8xlarge	> i2.8xlarge
		> i2.4xlarge	> hs1.8xlarge
		> i2.8xlarge	> cc2.8xlarge
		> hs1.8xlarge	> hi1.4xlarge
		> m1.xlarge	
		> cc2.8xlarge	
		> m2.4xlarge	
> hi1.4xlarge			
High / Full EDH workloads	All CDH services	> i2.2xlarge	> cc2.8xlarge
		> i2.4xlarge	> hs1.8xlarge
		> cc2.8xlarge	> i2.8xlarge
		> hs1.8xlarge	

## Regions and Availability Zones

[Regions](#) are self-contained geographical locations where AWS services are deployed. Regions have their own deployment of each service. Each service within a region has its own [endpoint](#) that you can interact with to use the service.

Within regions there are [availability zones](#). These are isolated locations within a general geographical location. Some regions have more availability zones than others. While provisioning, you can choose specific availability zones or let AWS pick for you.

Cloudera EDH deployments are to be restricted to single availability zones. Clusters spanning availability zones and regions are not supported.

## Networking, connectivity and security

### VPC

There are several different configuration options for VPC. See the [VPC documentation](#) for detailed explanation of the different options and choose based on your networking requirements. You can deploy Cloudera Enterprise clusters either in public subnets or in private subnets, as highlighted above. In both cases, the instances forming the cluster are advised to not be assigned a publicly addressable IP unless the requirement is for them to be accessible from the internet or other AWS services. If you assign public IP addresses to the instances and want to block incoming traffic, you can do so by configuring it in the security groups.

### Connectivity to Internet and other AWS services

Deploying the instances in a public subnet allows them to have access to the internet for outgoing traffic as well as to other AWS services, such as S3, RDS etc. Clusters that need data transfer between other AWS services (especially S3) and HDFS should be deployed in a public subnet and with public IP addresses assigned so that they can directly transfer data to those services. You should configure the security group for the cluster nodes to block incoming connections to the cluster instances.

Clusters that don't need heavy data transfer between other AWS services or the internet and HDFS should be launched in the private subnet. These clusters still might need access to services like RDS or software repositories for updates etc. This can be accomplished by provisioning a NAT instance in the public subnet, allowing access outside the private subnet into the public domain. The NAT instance is not recommended to be used for any large-scale data movement.

If you choose to completely disconnect the cluster from the internet, you will block access for software updates as well as to other AWS services, which makes maintenance activities hard. If the requirement is to completely lock down any external access because of which you don't want to keep the NAT instance running all the time, Cloudera recommends spinning up a NAT instance as and when external access is required and spinning it down once the activities are complete.

### Private Data Center Connectivity

You can establish connectivity between your data center and the VPC hosting your Cloudera Enterprise cluster by using a VPN or using Direct Connect. We recommend using Direct Connect so that there is a dedicated link between the two networks with lower latency, higher bandwidth, security and encryption via IPSec as compared to the public Internet. If you don't need high bandwidth and low latency connectivity between your data center and AWS, connecting to EC2 through the Internet is sufficient and need Direct Connect may not be required.

## Security Groups

**Security Groups** are analogous to firewalls. You can define rules for EC2 instances and define what traffic will be allowed from what IP address and port ranges. Instances can belong to multiple security groups. For Cloudera Enterprise deployments, you need the following security groups:

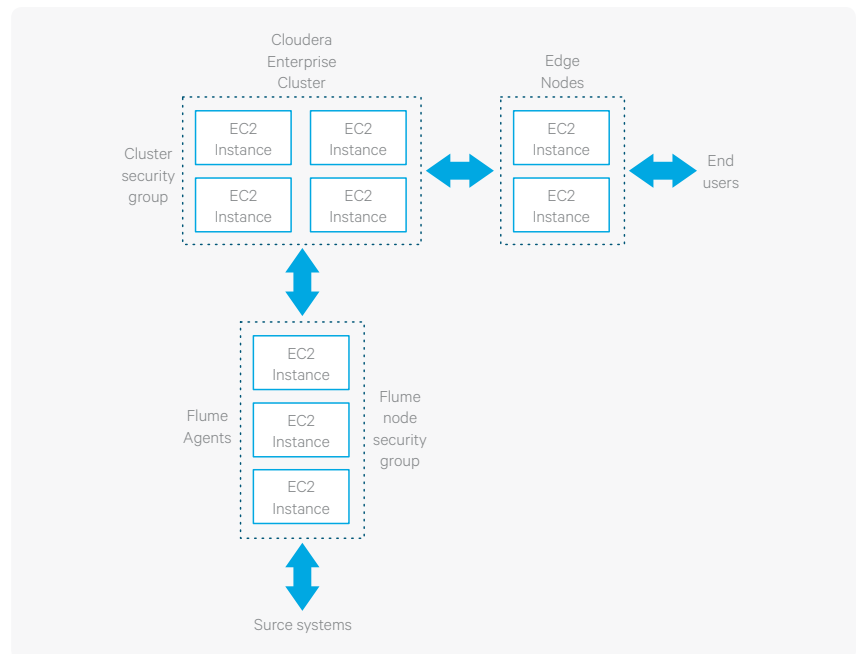
**Cluster** - This security group blocks all inbound traffic except that coming from the security group containing the flume nodes and edge nodes. You can allow outbound traffic for Internet access during installation and upgrade time and disable it thereafter if you choose to do so. You can also allow outbound traffic if you intend to access large volumes of Internet based data sources like S3.

**Flume nodes** - This security group is for instances running Flume agents. Outbound traffic to the Cluster security group has to be allowed and inbound traffic from sources from which flume is receiving data from needs to be allowed.

**Edge nodes** - This security group is for instances running client applications. Outbound traffic to the Cluster security group has to be allowed and incoming traffic from IP addresses that will interact with the client applications as well the cluster itself needs to be allowed.

Each of these security groups can be put in public subnets or private subnets depending on the access requirements highlighted above.

A full deployment would look like the following:



Source systems are where the data is being ingested from using Flume. You'll have flume sources deployed on those machines.

End users are the end clients that will interact with the applications running on the edge nodes that can interact with the Cloudera Enterprise cluster.

## Placement Groups

As described in the AWS documentation for Placement Groups available [here](#), Placement Groups are a logical grouping of EC2 instances within an availability zone, where instances are provisioned such that the network between them has higher throughput and lower latency. AWS accomplishes this by trying to provision instances as close to each other as possible. This does limit the pool of instances available for provisioning but this is the only way to guarantee uniform network performance. We recommend provisioning the worker nodes of the cluster within a placement group. Master and edge nodes can be outside the placement group unless you need high throughput and low latency between those and the cluster. That would be the case if you are moving large amounts of data or expect low latency responses between the edge nodes and the cluster.

## Supported AMIs

Amazon Machine Images are the virtual machine images that run on EC2 instances. These consist of the operating system and any other software that the AMI creator wanted to bundle into them. For Cloudera Enterprise deployments in AWS, we support [Red Hat AMIs](#). There are HVM and PV AMIs available. For certain instances like cc2.8xlarge, you have to use [Hardware Assisted Virtualization \(HVM\)](#) whereas for instances like m1.xlarge, you have to use [Paravirtualization \(PV\)](#). For instances like hs1.8xlarge where there's a choice, we recommend you use HVM. You can find a list of the Red Hat AMIs for each region [here](#). The only version currently supported is Red Hat 6.4.

## Storage options and configuration

AWS offers different kinds of storage options that you can use. These vary in their performance, durability and cost characteristics.

### Instance storage

EC2 instances have storage attached at the instance level, similar to disks on a physical server. The storage is virtualized and referred to as ephemeral storage. It's called so because the lifetime of the storage is the same as the lifetime of your EC2 instance. If you stop or terminate the EC2 instance, the storage is lost. It's not lost on restarts however. Different EC2 instances come with different amounts of instance storage, as highlighted above. For long running Cloudera Enterprise clusters, the HDFS data directories should use the instance storage. This allows you to benefit from all the benefits of shipping compute close to the storage and not having to read remotely over the network.

### Simple Storage Service

We strongly recommend using S3 to keep a copy of the data you have in HDFS for disaster recovery. The durability and availability guarantees make it a very good place to have a cold backup that you can restore from in case the primary HDFS cluster goes down. If you want a hot backup, you'll need a second HDFS cluster holding a copy of your data.

### Root Device

We recommend using EBS volumes as root devices for the EC2 instances. At the time of instantiating the instances, you can define the root device size. The root device size for Cloudera Enterprise clusters should be at least 500GB, allowing for parcels and logs to be stored. You don't need to use any of the instance storage for the root device. The instance storage can be 100% utilized for HDFS data dirs.

## Capacity planning

One of the benefits of AWS is the ability to scale your Cloudera Enterprise cluster up and down easily. If your storage or compute requirements change, you can provision and deprovision instances accordingly and meet your requirements quickly, without waiting to buy physical servers. However, a little bit of planning upfront makes operations easier. The fundamental thing you have to plan for is whether your workloads need high amount of storage capacity or not. The supported EC2 instances have different amounts of memory, storage and compute and deciding which instance type and generation should make up your initial deployment depends on the storage and workload requirement. The operational cost of your cluster will depend on the type and number of instances you choose.

### Low storage density

For use cases with lower storage requirements, using cc2.8xlarge is recommended. It gives you lower amount of storage per instance but has a high amount of compute and memory resources to go with it. The cc.8xlarge instances have 4 x 840GB raw instance storage on them.

### High storage density

For use cases with higher storage requirements, using hs1.8xlarge is recommended. These give you a high amount of storage per instance but the compute is lower than cc2.8xlarge instances. The hs1.8xlarge instances have 24 x 2TB instance storage on them.

### Reserved Instances

AWS offers the ability to [reserve EC2 instances](#) upfront and pay a lower per hour price. This is beneficial for users that are making a decision to use EC2 instances for the foreseeable future and will be keeping them on for majority of the time. Reserving instances can drive down the TCO significantly of long running Cloudera Enterprise clusters. There are different options for reserving instances in terms of the time period of the reservation and the utilization of each instance. Refer to the AWS documentation to plan instance reservation.

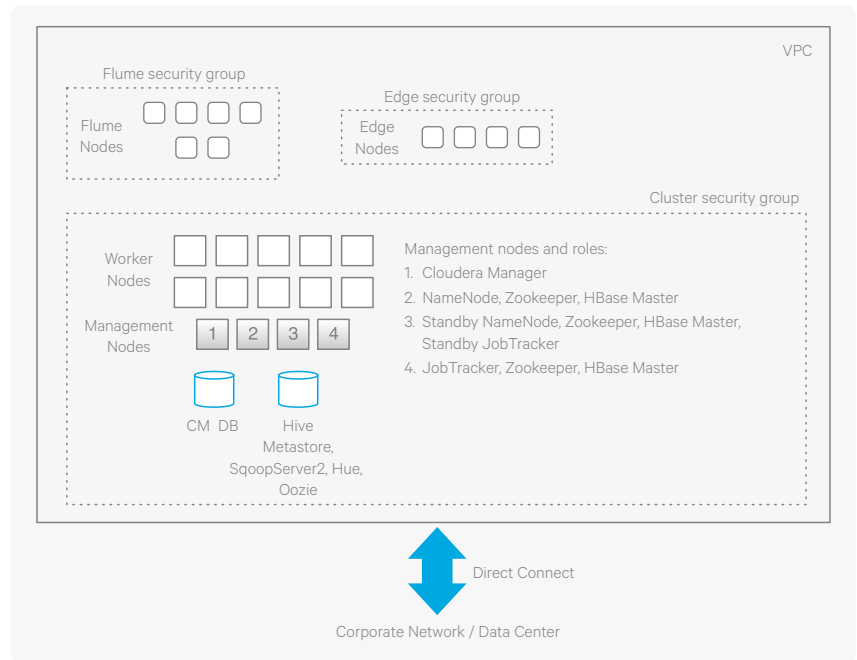
## Relational Databases

Cloudera Enterprise deployments require relational databases for the following components:

- > Cloudera Manager databases
- > Hive and Impala metastore
- > Hue database
- > Oozie database
- > SqoopServer2 database

You have two choices for operating relational databases in AWS. You can either provision EC2 instances and install/manage your own database instances, or you can use RDS. The list of supported database types and versions is available [here](#).

With all the considerations highlighted thus far, a deployment in AWS would look like (for both private and public subnets):



The next section of this whitepaper addresses preparation and installation of the cluster.

## Installation and Software Configuration

### Preparation

#### Provisioning instances

To provision EC2 instances, the first step is to define the VPC configurations based on your requirements on aspects like access to Internet, other AWS services and connectivity to your corporate network. After that, you can either use the [EC2 command line API tool](#) to provision instances or do it through the [AWS management console](#). To provision instances, you must create a keypair with which you will later be able to log into the instances. In Red Hat AMIs, you'll be able to use this keypair to log in as ec2-user, which has sudo privileges for further administrative tasks. While provisioning instances, make sure to specify the following:

- > Red Hat 6.4 AMI
- > Root device size of at least 500GB
- > All ephemeral storage devices to be attached to the instances
- > Tags to indicate the role that the instance will play later on. This makes identifying instances easy later on

Along with provisioning instances, databases must be provisioned (RDS or self managed). If you are provisioning in a public subnet, RDS instances can be accessed directly. If you are deploying in a private subnet, you either need a NAT instance to access RDS instances or you have to setup a database instances on EC2 inside the private subnet. The database credentials will be required during Cloudera Enterprise installation.

## Setting up instances

Once the instances are provisioned, there are a few tasks that need to be done in order to get them ready for deploying CE. These are:

- > Disabling IP tables
- > Disabling SELinux
- > Formatting and mounting the instance storage
- > Resizing the root volume if it does not show full capacity

For more information on operating system preparation and configuration, please see the Cloudera Manager installation instructions available [here](#).

## Deploying Cloudera Enterprise

To deploy Cloudera Enterprise, log onto the instance that you have elected to host Cloudera Manager and follow installation instructions available [here](#).

## Cloudera Enterprise configuration considerations

### HDFS

**Durability** for Cloudera Enterprise deployments in AWS, the supported storage option is the ephemeral storage. HDFS on EBS volumes is not a supported configuration. Data stored on ephemeral storage is lost if instances are stopped, terminated or go down for some other reason. Data does persist on restarts however. Data durability in HDFS can be guaranteed by keeping replication at 3. We don't recommend lowering the replication factor.

Secondly, a persistent copy of all data should be maintained in S3 to guard against total cases where you can lose all 3 copies of the data. This can be accomplished by either writing to S3 at ingest time or distcp'ing datasets from HDFS after the fact.

**Availability** of HDFS can be accomplished by deploying the NameNode with High Availability with at least 3 Journal Nodes.

### Zookeeper

We require running at least 3 Zookeeper servers for availability and durability reasons for production clusters.

### Flume

For durability in Flume agents, you can use memory channel or file channel. Flume's memory channel offers increased performance at the cost of no data durability guarantees. File channels offer a higher level of durability guarantee since the data is persisted on disk in the form of files. We support file channel on ephemeral storage as well as EBS. If the EC2 instance goes down, the data on the ephemeral storage will be lost. For guaranteed data delivery, use EBS backed storage for Flume file channel.

## Summary

Cloudera and AWS allow users to deploy and use Cloudera Enterprise on AWS infrastructure, combining the scalability and functionality of the Cloudera Enterprise suite of products with the flexibility and economics of the AWS cloud. This whitepaper provided reference configurations for Cloudera Enterprise deployments in AWS. These configurations leverage different AWS services such as EC2, S3 and RDS.

## References

### Cloudera Enterprise

Cloudera [homepage](#)

Cloudera Enterprise [documentation](#)

Cloudera Enterprise [support](#)

### Amazon Web Services

AWS [homepage](#)

EC2 [homepage](#)

EC2 instance [lifecycle](#)

S3 [homepage](#)

RDS [homepage](#)

VPC [homepage](#)

Direct Connect [homepage](#)

EC2 [networking and security](#)

Red Hat certified [AMIs](#)

AWS [developer tools](#)

AWS [support](#)

## About Cloudera

Cloudera is revolutionizing enterprise data management by offering the first unified Platform for Big Data, an enterprise data hub built on Apache Hadoop. Cloudera offers enterprises one place to store, access, process, secure, and analyze all their data, empowering them to extend the value of existing investments while enabling fundamental new ways to derive value from their data. Cloudera's open source Big Data platform is the most widely adopted in the world, and Cloudera is the most prolific contributor to the open source Hadoop ecosystem. As the leading educator of Hadoop professionals, Cloudera has trained over 22,000 individuals worldwide. Over 1,200 partners and a seasoned professional services team help deliver greater time to value. Finally, only Cloudera provides proactive and predictive support to run an enterprise data hub with confidence. Leading organizations in every industry plus top public sector organizations globally run Cloudera in production. [www.cloudera.com](http://www.cloudera.com).

---

### [cloudera.com](http://cloudera.com)

1-888-789-1488 or 1-650-362-0488

Cloudera, Inc. 1001 Page Mill Road, Palo Alto, CA 94304, USA

© 2014 Cloudera, Inc. All rights reserved. Cloudera and the Cloudera logo are trademarks or registered trademarks of Cloudera Inc. in the USA and other countries. All other trademarks are the property of their respective companies. Information is subject to change without notice.