

Module 1

AWS INTRODUCTION AND HISTORY

Jeff Bezos incorporated the company in 1994 and Amazon.com was launched in 1995 as an online bookstore. Amazon.com, Inc. is an American multinational electronic commerce company with its headquarters in Seattle, Washington. It is the world's largest online retailer. Amazon has continued to grow and officially launched Amazon Web Services (AWS) in 2006. More came after, including Amazon Publishing, the Kindle, Amazon Game Studios, and Amazon Art. After over a decade of building and running the highly scalable web application, Amazon.com, the company realized that it had developed a core competency in operating massive scale technology infrastructure and data centers, and embarked on a much broader mission of serving a new customer segment—developers and businesses—with a platform of web services they can use to build sophisticated, scalable applications. Today, AWS is the fastest-growing multi-billion dollar enterprise IT vendor in the world.

Amazon Web Services (AWS)

Enable business and developers to use web services to build scalable, sophisticated applications

Amazon Web Services is 10+ years in the making. Amazon Web Services, also abbreviated to AWS, is a collection of remote computing services called web services. These web services make up a cloud computing platform offered via the Internet. We deliver web-based cloud services for storage, computing, networking, databases, and more.

The AWS mission is to enable businesses and developers to use web services to build scalable, sophisticated applications. Web services is another name for what people now call "the cloud."

AWS has been continually expanding its services to support virtually any cloud workload. It now has more than 50 services that range from compute, storage, networking, database, analytics, application services, deployment, management and mobile. In 2015, AWS launched 722 new features and/or services for a total of 1,950 new features and/or services since its inception in 2006.

As of 1 February 2016, AWS has launched 1,950 new services as well as features and updates to existing services.

AWS Customers

Enterprise Customers:

Enterprise cloud computing with AWS can help IT increase innovation, agility, and resiliency; all while reducing cost. With AWS, you can build enterprise cloud solutions quickly and without a big up-front investment. The free tier allows you to prototype virtually any application for free.

Startup Customers:

Our innovations free you to scale quickly, go to market faster, control costs, and stay lean. AWS Activate is a free program with resources for startups to get the most out of AWS from day one.

Public Sector Customers:

AWS offers scalable, cost-effective cloud services that public sector customers can use to meet mandates, reduce costs, drive efficiencies, and accelerate innovation.

Advantages & Benefits of AWS

Cloud computing provides a simple way to access servers, storage, databases, and a broad set of application services over the Internet. AWS owns and maintains the network-connected hardware required for these application services, while you provision and use what you need via a web application.

Trade capital expense for variable expense:

Instead of having to invest heavily in data centers and servers before you know how you're going to use them, you can pay only when you consume computing resources, and pay only for how much you consume.

Benefit from massive economies of scale:

By using cloud computing, you can achieve a lower variable cost than you can get on your own. Because usage from hundreds of thousands of customers is aggregated in the cloud, providers such as Amazon Web Services can achieve higher economies of scale, which translates into lower pay as you go prices.

Stop guessing capacity:

Eliminate guessing on your infrastructure capacity needs. When you make a capacity decision prior to deploying an application, you often either end up sitting on expensive idle resources or dealing with limited capacity. With cloud computing, these problems go away. You can access as much or as little as you need, and scale up and down as required with only a few minutes' notice.

Increase speed and agility:

In a cloud computing environment, new IT resources are only ever a click away, which means you reduce the time it takes to make those resources available to your developers from weeks to just minutes. This results in a dramatic increase in agility for the organization, since the cost and time it takes to experiment and develop is significantly lower.

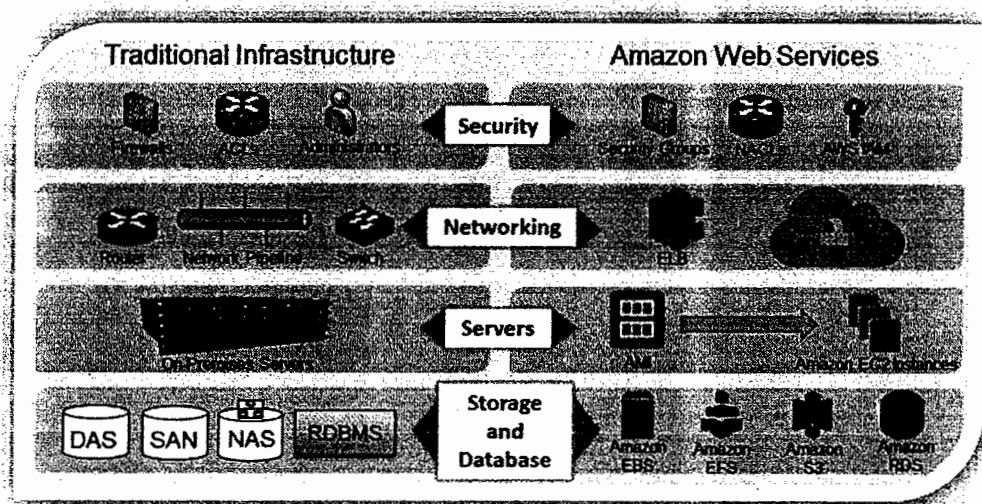
Stop spending money on running and maintaining data centers:

Focus on projects that differentiate your business, not the infrastructure. Cloud computing lets you focus on your own customers, rather than on the heavy lifting of racking, stacking and powering servers.

Go global in minutes:

Easily deploy your application in multiple regions around the world with just a few clicks. This means you can provide a lower latency and better experience for your customers simply and at minimal cost.

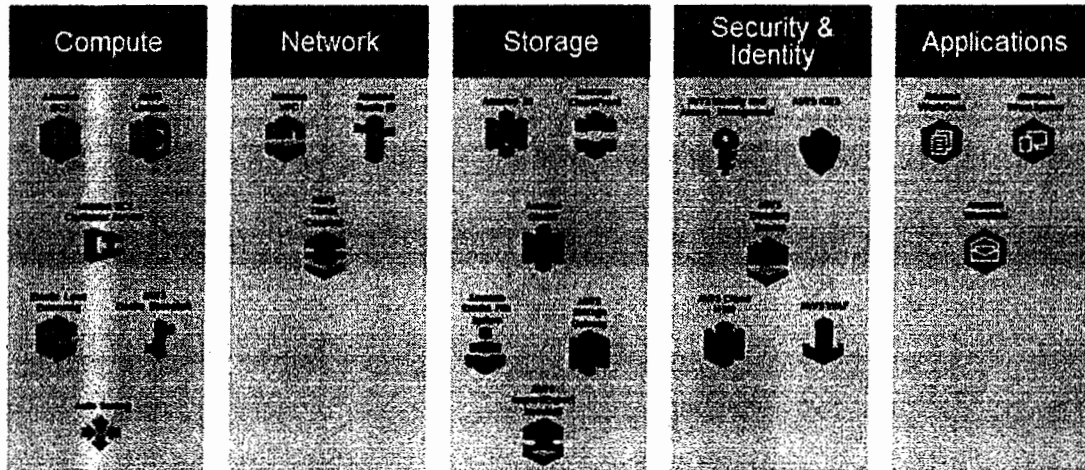
AWS Core Infrastructure and Services



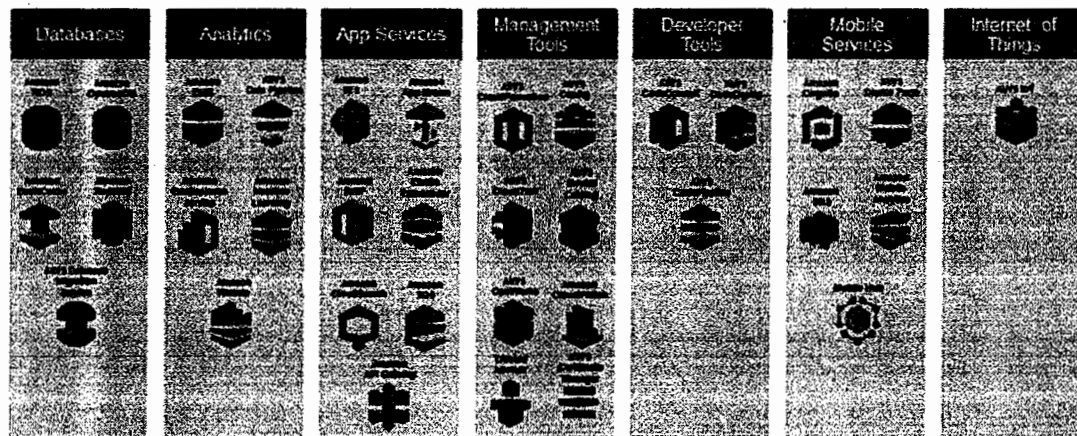
Many of our services have analogs in the traditional IT space and terminology. This side-by-side comparison shows how AWS products and services relate to a traditional infrastructure.

AWS cloud computing provides a simple way to access servers, storage, databases and a broad set of application services over the Internet. AWS owns and maintains the network-connected hardware required for these application services, while you provision and use what you need.

AWS Foundation Services



AWS Platform Services



AWS Global Infrastructure

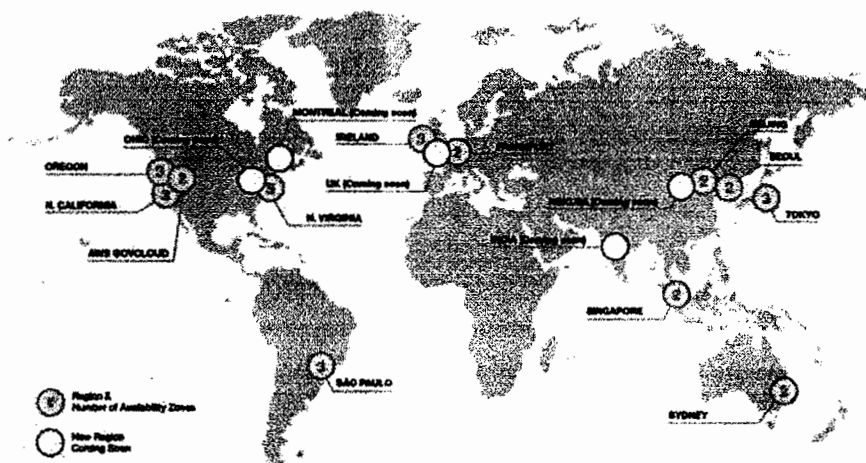
Regions:

- ✓ Geographical Locations
- ✓ Consists of at least two Availability Zones (AZs)

Availability Zones:

- ✓ Clusters of Data centers
- ✓ Isolated from failures in other Availability Zones

AWS Regions are geographic locations that contain multiple Availability Zones (AZs). Availability Zones consist of data centers clustered in a region. Each Availability Zone is engineered to be isolated from failures in other Availability Zones.



AWS is steadily expanding its global infrastructure to help customers achieve lower latency and higher throughput, and to ensure that your data resides only in the region you specify. As you and all customers grow their businesses, AWS will continue to provide infrastructure that meets your global requirements.

As of January 2016, AWS has 12 geographic Regions with 32 Availability Zones. The AWS GovCloud (US) Region is an isolated region designed to allow US government agencies and customers to move sensitive workloads into the cloud by addressing their specific regulatory and compliance requirements. AWS products and services are available by region so you may not see all regions available for a given service.

You can run applications and workloads from a region to reduce latency to end- users while avoiding the up-front expenses, long-term commitments, and scaling challenges associated with maintaining and operating a global infrastructure. In 2016, the AWS Global Infrastructure will expand with at least 10 new Availability Zones in new geographic Regions including Ohio in North America, Ningxia in China, India, Korea, and the United Kingdom.

At least 2 AZs per region.

US East (N. Virginia)

- ✓ us-east-1a
- ✓ us-east-1b
- ✓ us-east-1c
- ✓ us-east-1d
- ✓ us-east-1e

Asia Pacific (Tokyo)

- ✓ ap-northeast-1a
- ✓ ap-northeast-1b
- ✓ Ap-northeast-1c

Each region is a separate geographic area that has multiple locations isolated from each other known as Availability Zones (AZ). Each AZ is isolated, but the AZs in a region are connected through low-latency links. Where natural disasters or fault lines are a consideration, AWS isolates its Availability Zones so that they are not easily affected at the same time. For example, where earthquakes are a problem AWS would not build two AZs on the same fault line. When you launch an instance, you can select an AZ or let AWS choose one for you. If you distribute your instances across multiple AZs and one instance fails, you can design your application so that an instance in another AZ can handle requests.

Achieving High Availability Using Multi-AZ

AWS highly recommends provisioning your compute resources across multiple Availability Zones. If you have multiple instances, you can run them across more than one AZ and get added redundancy. If a single AZ has a problem, all assets in your second AZ will be unaffected.

AWS Global Infrastructure

50+ AWS Edge locations:

Local points-of-presence commonly supporting AWS services like:

- ✓ Amazon Route 53
- ✓ Amazon CloudFront

Edge locations help lower latency and improve performance for end users.

CLOUD COMPUTING CONCEPTS

What is Cloud Computing?

Cloud Computing is on-demand delivery of IT resources and applications via internet with pay-as-you-go pricing.

Cloud computing is a common term for a variety of computing concepts that involve large numbers of computers that are connected through a real-time communication network, like the Internet.

Cloud computing is the use of computing resources (hardware and software) that are delivered as a service over a network (typically the Internet). The name comes from the use of a cloud-shaped symbol as an abstraction in system diagrams for the complex infrastructure it contains. Cloud computing entrusts remote services with a user's data, software, and computation. Cloud computing allows you to access as many resources as you need, almost instantly, and only pay for what you use.

Instead of buying, owning, and maintaining your own data centers and servers, organizations can acquire technology such as compute power, storage, databases, and other services on an as-needed basis.

Essential Characteristics of Cloud Computing?

Cloud computing is characterized by five characteristics:

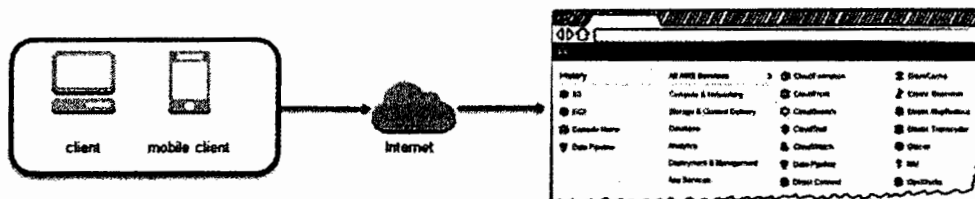
On-demand self-services, broad network access, resource pooling, rapid elasticity, and measured service.

- ✓ On-demand Self Services
- ✓ Broad Network Access
- ✓ Resource Pooling
- ✓ Rapid Elasticity
- ✓ Measured Services

On-Demand Self Services & Broad Network Access

- ✓ User Provisions computing resources as needed.
- ✓ User Interacts with cloud service provider through an online control panel.

✓ Clear solutions are available through a variety of network-connected devices and over varying platforms.

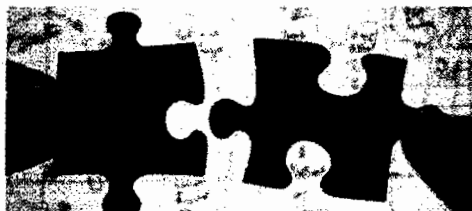


On-demand self-services are an essential characteristic of cloud computing. The user provisions computing resources as needed and interacts with the cloud service provider through an online control panel, such as the AWS Management Console.

Broad network access is another characteristic of cloud computing. Clear solutions are available through a variety of network-connected devices and over varying platforms.

Resource Pooling

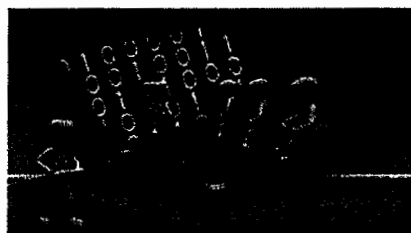
Securely separate resources to service multiple customers.




A third characteristic of cloud computing is resource pooling. Cloud computing solutions securely separate resources to service multiple customers.

Rapid Elasticity

Resources are quickly scalable and flexible based on business needs.



Measured Service



What does My AWS Cloud Look Like?

Module 2

AWS INFRASTRUCUTRE

Amazon Elastic Compute Cloud (EC2)

Understand Amazon Elastic Compute Cloud (EC2) concepts including:

- ✓ Instances vs. servers
- ✓ Types and families
- ✓ Ephemeral vs. persistent storage (root instance volumes)
- ✓ Amazon Machine Images (AMIs)
- ✓ Bootstrapping/user data

Amazon EC2 instances are virtualized servers in Amazon's data centers.

Amazon EC2 is designed to make web-scale computing easier for developers. Amazon EC2's simple web service interface allows you to obtain and configure capacity with minimal friction. It provides you with complete control of your computing resources and allows you to run on Amazon's proven computing environment.

Amazon EC2 reduces the time required to obtain and boot new server instances, allowing you to quickly scale capacity as your computing requirements change. Amazon EC2 changes the economics of computing by allowing you to pay only for capacity that you actually use. Amazon EC2 provides the tools for you to build failure-resilient applications and isolate the applications from common failure scenarios.

Amazon EC2 Facts

- ✓ Scale capacity as your computing requirements change
- ✓ Pay only for capacity that you actually use
- ✓ Choose Linux or Windows
- ✓ Deploy across AWS Regions and Availability Zones for reliability

Amazon EC2 presents a true virtual computing environment, allowing you to use web service interfaces to launch instances with a variety of operating systems, load them with your custom application environment, manage your network's access permissions, and run your image using as many or few systems as you need.

You have the ability to programmatically scale your computing capacity as your requirements change. You pay only for capacity that you actually use and can choose Linux or Windows. You can leverage the AWS global infrastructure to deploy across regions and Availability Zones (AZs) for reliability.

Launching an Amazon EC2 Instance via the web console

- ✓ Determine the AWS Region in which you want to launch the Amazon EC2 instance.
- ✓ Launch an Amazon EC2 instance from a pre-configured Amazon Machine Image (AMI).
- ✓ Choose an instance type based on CPU, memory, storage, and network requirements.
- ✓ Configure network, IP address, security groups, storage volume, tags, and key pair.

Before you create your first Amazon EC2 instance think about which region you want to have that instance in. The AMI comes pre-installed with many AWS API tools as well as CloudInit. AWS API tools enable scripting of important provisioning tasks from within an Amazon EC2 instance. AMIs are like building blocks of EC2 instances. They are templates of a computer's volumes. AMIs can have public or private access. You can also create gold master images of your Amazon EC2 infrastructure, which allow you to decrease your boot times.

AMI Details

An AMI includes the following:

- ✓ A template for the root volume for the instance (for example, an operating system, an application server, and applications).
- ✓ Launch permissions that control which AWS accounts can use the AMI to launch instances.
- ✓ A block device mapping that specifies the volumes to attach to the instance when it's launched.

An AMI is a template that contains a software configuration such as an operating system, application server, and applications. You use an AMI to launch an instance, which is the copy of the AMI running as a virtual server on a host computer in Amazon's data center. You can launch as many instances as you want from an AMI. You can also launch instances from as many AMIs as you need.

You can create your own AMI by customizing the instance that you launch from a public AMI and then saving the configuration as a custom AMI for your own use. You can also buy, share, and sell AMIs.

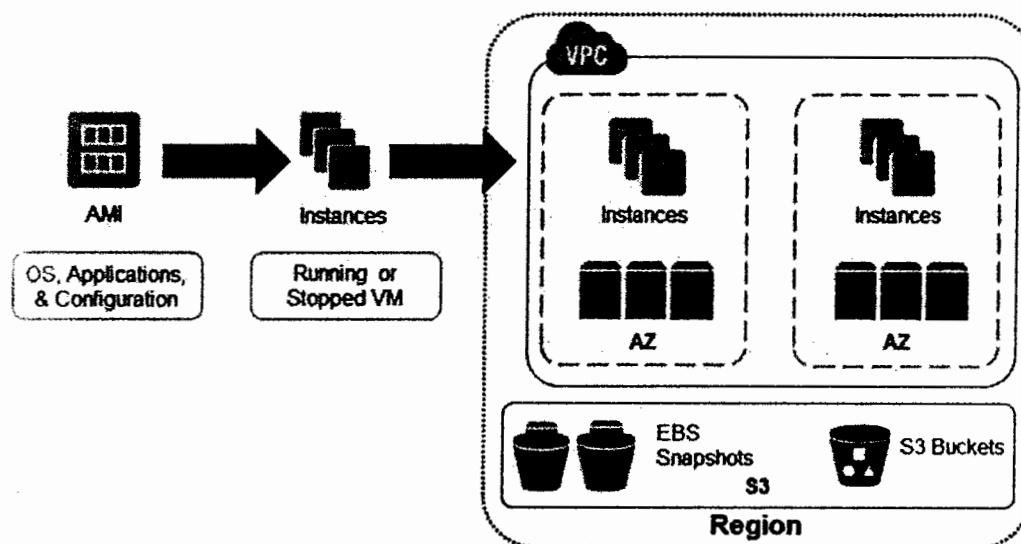
Instances and AMIs

Select an AMIs based on:

- ✓ Region
- ✓ Operating System
- ✓ Architecture (32-bit or 64-bit)
- ✓ Launch permissions
- ✓ Storage for the root device

You select an AMI based on region, operating system, architecture, launch permissions, and storage for the root device. Launch permissions determine availability of an AMI and are either public (the owner grants launch permissions to all AWS accounts), explicit (the owner grants launch permissions to specific AWS accounts), or implicit (the owner has implicit launch permissions for an AMI).

Amazon EC2 Instances



You can launch multiple instances of different types from a single AMI when launching an EC2 instance. An instance type essentially determines the hardware of the host computer used for your instance. Each instance type offers different compute and memory capabilities. Select an instance type based on the amount of memory and computing power that you need for the application or software that you plan to run on the instance. Your instance keeps running until you stop or terminate it, or until it fails. Instances are deployed in the Amazon EC2 public cloud or the Amazon Virtual Private Cloud in an Availability Zone (AZ) within a Region. You can configure security and network access on your Amazon EC2 instance.

Customers can deploy to multiple AZs within a Region. You choose which instance types you want, and then start, terminate, and monitor as many instances of your AMI as needed, using the web service APIs or the variety of management tools provided.

Amazon EC2 instances can leverage Amazon Elastic Block Store volumes in each Availability Zone. Determine whether you want to run in multiple locations, utilize static IP endpoints, or attach persistent block storage to your instances. Amazon EBS volumes can be saved using "snapshots." Additionally, Amazon S3 buckets can be used to store data objects needed by

Amazon EC2 instances. Pay only for the resources that you actually consume, like instance-hours or data transfer.

Amazon EBS vs. Amazon EC2 Instance Store

Amazon EBS	Amazon EC2 Instance Store
Data stored on an Amazon EBS volume can persist independently of the life of the instance.	Data stored on a local instance store persists only as long as the instance is alive.
Storage is persistent.	Storage is ephemeral.

Use the local instance store only for temporary data. For data requiring a higher level of durability, use Amazon EBS volumes or back up the data to Amazon S3, topics that are both discussed later in this module. If you are using an Amazon EBS volume as a root partition, set the Delete on termination flag to "No" if you want your Amazon EBS volume to persist outside the life of the instance.

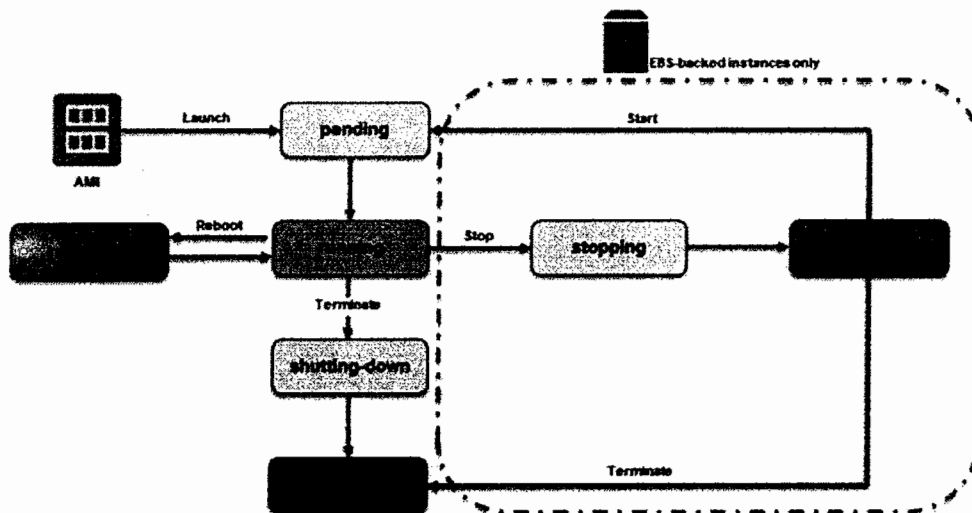
AMI Types – Storage for the Root Device

Characteristic	Amazon EBS-Backed	Amazon Instance Store-Backed
Boot time	Usually < 1 minute	Usually < 5 minutes
Size limit	16 TiB	10 GiB
Data persistence	The root volume is deleted when the instance terminates. Data on any other Amazon EBS volumes persists after instance termination.	Data on any instance store volumes persists only during the life of the instance.
Charges	Instance usage, Amazon EBS volume usage, and storing your AMI as an Amazon EBS snapshot	Instance usage and storing your AMI in Amazon S3
Stopped state	Can be stopped	Cannot be stopped

AMIs are either Amazon Elastic Block Storage (EBS)-backed or backed by instance store. When an AMI is EBS-backed, this means that the root device for an instance is an EBS volume created from an EBS snapshot. When an AMI is instance-store backed, this means that the root device for the instance was created from a template stored in Amazon S3. Key differences between both categories of AMIs are shown in the slide.

Storage will be discussed more extensively later in the modules.

Instance Lifecycle



The lifecycle of an instance launched from an AMI. Note that you can only stop and start instances that are Amazon Elastic Block Store (EBS)-backed.

An EC2 instance can be in one of the following states:

Pending: When you launch an instance, it enters the pending state and the instance moves to a new host computer. The instance type specified at launch determines the hardware of the host computer for your instance.

Running: AWS uses the AMI specified at launch to boot the instance. Once the instance is ready for you, it enters the running state. You can connect to your running instance and use it as you would a computer sitting in front of you. As soon as your instance is in the running state, you're billed for each hour or partial hour that you keep the instance running. You are billed for all running instances, even if they are idle and not being connected to.

Rebooting: You can reboot your instance through the Amazon EC2 console, Amazon EC2 CLI, and the Amazon EC2 API. It is recommended that you reboot your EC2 instance rather than running the operating system reboot from the instance. When an instance is rebooted, it remains on the same host computer and maintains its public DNS name, private IP address, and any data on its instance store volumes. Rebooting an instance doesn't start a new instance billing hour.

Shutting Down: When you've decided you no longer need an instance, you can terminate it. The instance will enter the shutting-down state. You will stop incurring charges as soon as the instance enters shutting-down or terminated states.

Terminated: A terminated instance remains visible in the console for a while before it is deleted. You cannot connect to or recover a terminated instance.

Stopping: Amazon EBS-backed instances can be stopped. When you stop an instance, it enters the stopping state.

Stopped: Amazon EBS-backed instances in the stopped state are no longer eligible for hourly usage or data transfer fees. AWS does charge for the storage of EBS volumes on stopped instances. You can modify certain attributes of stopped instances, including the instance type. Starting a stopped instance puts it back into the pending state, which moves the instance to a new host machine. When you stop and start an instance, you lose any data on the instance store volumes on the previous host computer.

AWS Marketplace – IT Software Optimized for the Cloud

AWS and Oracle have worked together to offer customers convenient options for deploying enterprise applications in the cloud. Customers can not only build enterprise-grade solutions hosted by Amazon Web Services using database and middleware software by Oracle, but they can also launch entire enterprise software stacks from Oracle on Amazon EC2.

New and existing SAP customers can deploy their SAP solutions on SAP-certified Amazon EC2 instances in production environments knowing that SAP and AWS have tested the performance of the underlying AWS resources, verified their performance, and certified them against the same standards that apply to servers and virtual platforms.

AWS also provides infrastructure services that allow customers to easily run Microsoft Windows Server applications in the cloud, without the cost and complexity of having to purchase or manage servers or data centers. AMIs are available; allowing customers to start running fully supported Windows Server virtual machine instances in minutes.

Customers may also rely on the global infrastructure of AWS to power everything from custom .NET applications to enterprise deployments of Microsoft Exchange Server, SQL Server, or SharePoint Server.

Software launched from AWS customers automatically deploys onto Amazon Elastic Compute Cloud (EC2), which is the AWS compute service. AWS customers use 143 million hours a month of Amazon EC2 for AWS Marketplace software products.

AWS Marketplace Benefits

- ✓ Easy product discovery
- ✓ Streamlined buying experience
- ✓ Simplified billing- software charges on AWS bill
- ✓ Expedited deployment cycles
- ✓ Optimized software capacity
- ✓ Matched spend to actual usage-Opex
- ✓ Trust vetted and scanned products

AWS Test Drive provides a private IT sandbox environment containing preconfigured server based solutions. In under an hour, and using a step-by-step lab manual and video, launch, login and learn about these popular 3rd party IT solutions, powered by AWS and AWS CloudFormation.

Choosing the Right Amazon EC2 Instance

Choosing the right EC2 instance type matters. Selecting an appropriate instance type for your workload can save time and money. AWS has a wide variety of EC2 compute instance types to choose from. Each instance type or family (T2, M3, C4, C3, G2, R3, and so on) is optimized for different workloads or use cases. Within an EC2 family, you can choose from different sizes: for example, micro, small, medium, large, xlarge, and 2xlarge. AWS uses Intel® Xeon® processors for the EC2 instances to provide customers with high performance and value for their computing needs.

When you choose your instance type you should consider the several different attributes of each family, such as number of cores, amount of memory, amount and type of storage, network performance, and processor technologies.

Another important consideration is total cost of ownership (TCO). A lowest-price- per-hour instance is not necessarily a money saver; a larger compute instance can sometimes save both money and time. It is important to evaluate all the options to see what is best for your workload.

General Purpose Instances

T2 instances are a low-cost, burstable performance instance type that provide a baseline level of CPU performance with the ability to burst above the baseline. They offer a balance of compute, memory, and network resources for workloads that occasionally need to burst, such as web servers, build servers, and development environments.

M3 and M4 instances provide a balance of compute, memory, and network resources. These instances are ideal for applications that require high CPU and memory performance, such as encoding applications, high traffic content management systems, and memcached applications.

Compute-Optimized Instances

C3 and C4 instances are optimized for compute-intensive workloads. These instances have proportionally more CPU than memory (RAM). They are well suited to applications such as high performance web servers, batch processing, and high- performance scientific and engineering applications.

Memory-Optimized Instances

R3 instances are optimized for memory-intensive workloads. These instances offer large memory sizes for high throughput applications such as high performance databases, distributed memory caches, in-memory analytics, and large enterprise deployments of software such as SAP.

GPU Instances

G2 instances are optimized for graphics and graphic processing unit (GPU) compute applications, such as machine learning, video encoding, and interactive streaming applications.

Storage-Optimized Instances

I2 instances are optimized for storage and high random I/O performance, such as NoSQL databases, scale-out transactional databases, data warehousing, Hadoop, and cluster file systems.

D2 instances are optimized for storage and delivering high disk throughput. D2 instances are suitable for Massively Parallel Processing (MPP) data warehousing, MapReduce and Hadoop distributed computing, distributed file systems, and data processing applications.

Get the Intel® Advantage

AWS recently launched C4 compute-optimized instances which utilize Intel's latest 22nm Haswell microarchitecture. C4 instances use custom Intel® Xeon® v3 processors designed and built especially for AWS.

Through its relationship with Intel®, AWS provides its customers with the latest and greatest Intel® Xeon® processors that help in delivering the highest level of processor performance in Amazon EC2.

Intel® Processor Technologies

Intel® Xeon® processors have several other important technology features that can be leveraged by EC2 Instances.

- ✓ Intel® AVX is perfect for highly parallel HPC workloads such as life sciences or financial analysis.
- ✓ Intel® AES-NI accelerates encryption/decryption of data and therefore reduces the performance penalty that usually comes with encryption.
- ✓ Intel® Turbo Boost Technology automatically gives you more computing power when your workloads are not fully utilizing all CPU cores. Think of it as automatic overclocking when you have thermal headroom.

EC2 Instances with Intel® Technologies

	Burstable	Balanced	Compute	Memory	GPU	I/O	Storage
AWS Instance Type	T2	M4	C4	R3	G2	I2	D2
Intel® processor	Intel® Xeon® family	Intel® Xeon® E5-2676 v3	Intel® Xeon® E5-2668 v3	Intel® Xeon® E5-2670 v2	Intel® Xeon® E5-2670	Intel® Xeon® E5-2670 v2	Intel® Xeon® E5-2676 v3
Intel® process technology	22nm	22nm Haswell	22nm Haswell	22nm Ivy Bridge	32nm Sandy Bridge	22nm Ivy Bridge	22nm Haswell
Intel® AVX	●	●	●	●	●	●	●
Intel® AVX2		●	●				●
Intel® Turbo Boost	●	●	●	●	●	●	●
Storage	EBS only	EBS only	EBS only	SSD	SSD	SSD	HDD

The matrix on the slide highlights the individual Intel® technologies that were discussed previously and the EC2 instance family that can leverage each of these technologies.

Current Generation Instances

Instance Family	Some Use Cases
General purpose (t2, m4, m3)	<ul style="list-style-type: none">• Low-traffic websites and web applications• Small databases and mid-size databases
Compute optimized (c4, c3)	<ul style="list-style-type: none">• High performance front-end fleets• Video-encoding
Memory optimized (r3)	<ul style="list-style-type: none">• High performance databases• Distributed memory caches
Storage optimized (i2, d2)	<ul style="list-style-type: none">• Data warehousing• Log or data-processing applications
GPU instances (g2)	<ul style="list-style-type: none">• 3D application streaming• Machine learning

Each vCPU is a hyperthread of an Intel Xeon core for M4, M3, C4, C3, R3, HS1, G2, I2, and D2. Amazon EC2 lets you choose from a number of different instance types to meet your computing needs. Each instance provides a predictable amount of dedicated compute capacity and is charged per instance-hour consumed. First-generation (M1) general purpose instances provide a balanced set of resources and a low-cost platform that is well suited for a wide variety of applications. Second-generation (M3) general purpose instances provide a balanced set of resources and a higher level of processing performance compared to first-generation general purpose instances. Instances in this family are ideal for applications that require higher absolute CPU and memory performance. Applications that can benefit from the performance of second-generation general purpose instances include encoding applications, high traffic content management systems, and Memcached applications. High-memory instances offer large memory sizes for high throughput applications, including database and memory caching applications. High-CPU instances have proportionally more CPU resources than memory (RAM) and are well suited for compute-intensive applications. There are also various high-storage and cluster-computer instance types available.

Instance Metadata & User Data

Instance Metadata:

- ✓ Is data about your instance.
- ✓ Can be used to configure or manage a running instance.

Instance User Data:

- ✓ Can be passed to the instance at launch.

- ✓ Can be used to perform common automated configuration tasks.
- ✓ Runs scripts after the instance starts.

Although you can only access instance metadata and user data from within the instance itself, the data is not protected by cryptographic methods. Anyone who can access the instance can view its metadata. Therefore, you should take suitable precautions to protect sensitive data (such as long-lived encryption keys). You should not store sensitive data, such as passwords, as user data.

Retrieving Instance Metadata

- 🔗 To view all categories of instance metadata from within a running instance, use the following URI:
<http://169.254.169.254/latest/meta-data/>

- 🔗 On a Linux instance, you can use:

```
> $ curl http://169.254.169.254/latest/meta-data/
> $ GET http://169.254.169.254/latest/meta-data/
```

- 🔗 All metadata is returned as text (content type text/plain).



```
ami-id
ami-launch-index
ami-manifest-path
block-device-mapping/
hostname
instance-action
instance-id
instance-type
local-hostname
local-ipv4
mac
metrics/
network/
placement/
profile
public-hostname
public-ipv4
public-keys/
reservation-id
security-groups
services/
```

Because your instance metadata is available from your running instance, you do not need to use the Amazon EC2 console or the AWS CLI. This can be helpful when you're writing scripts to run from your instance. Note that you are not billed for HTTP requests used to retrieve instance metadata and user data.

Adding User Data

- ✓ You can specify user data when launching an instance.
- ✓ User data can be:
 - ✓ Linux Script-executed by **Cloud-init**
 - ✓ Windows batch or PowerShell Scripts-executed by **EC2Config** Service
- ✓ User data scripts run once per instance-id by default.

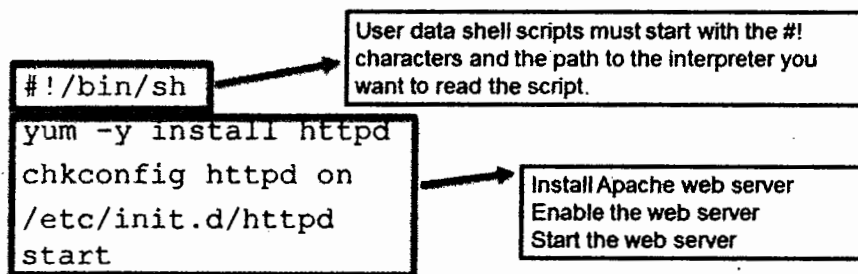
You can specify user data to configure an instance during launch, or to run a configuration script. To attach a file, select the As file option and browse for the file to attach. The cloud-init

package is an open source application built by Canonical that is used to bootstrap Linux images in a cloud computing environment, like Amazon EC2.

User data information:

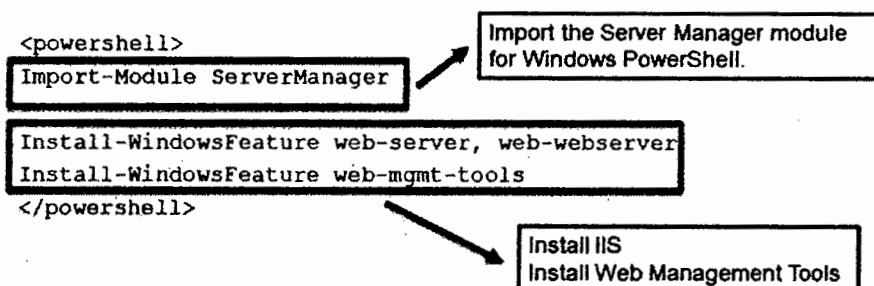
- ✓ User data is treated as opaque data: what you give is what you get back. It is up to the instance to be able to interpret it.
- ✓ User data is limited to 16 KB. This limit applies to the data in raw form, not base64-encoded form.
- ✓ User data must be base64-encoded before being submitted to the API. The Amazon EC2 command line tools perform the base64 encoding for you. The data is decoded before being presented to the instance.
- ✓ User data is executed only at launch. If you stop an instance, modify the user data, and start the instance, the new user data is not executed automatically by default.

User Data Example Linux



The slide shows an example of user data on Linux. You can also provide user data to an instance on Linux by using the `#cloud-config` directive – a format defined by cloud-init.

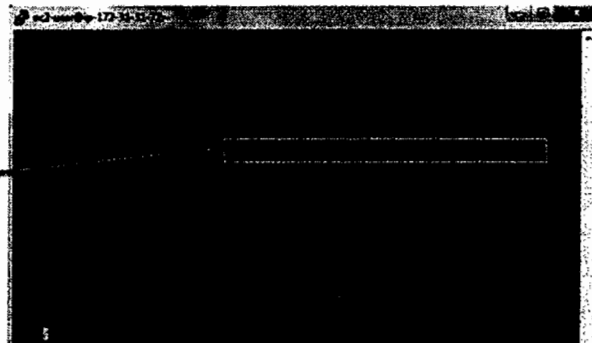
User Data Example Windows



You can send user data to a Windows instance with a PowerShell script (shown in slide) or with a set of Windows batch commands.

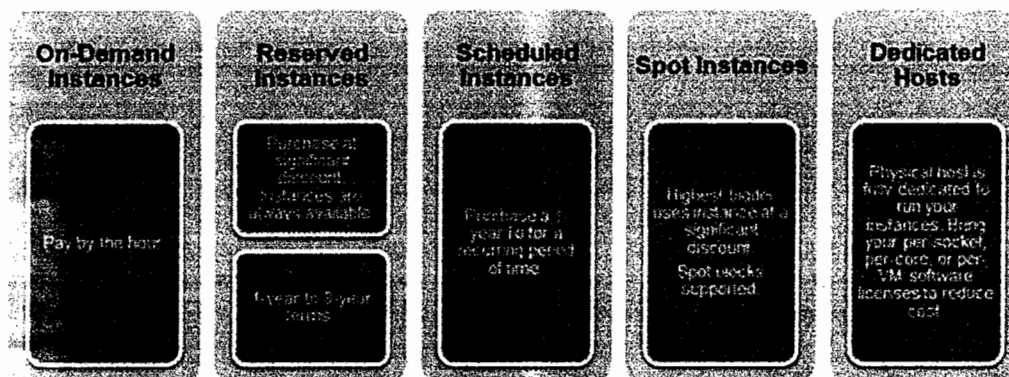
Retrieving User Data

- To retrieve user data, use the following URI:
`http://169.254.169.254/latest/user-data`
- On a Linux instance, you can use:
 - > `$ curl http://169.254.169.254/latest/user-data/`
 - > `$ GET http://169.254.169.254/latest/user-data/`



To retrieve instance user data, use the following URI: `http://169.254.169.254/latest/user-data`

Amazon EC2 Purchasing Options



On-Demand instances are free-tier eligible. It has the lowest up-front cost and the most flexibility. You pay for an hour at a time with no up-front commitments or long-term contracts. This is great for applications with short-term, spiky, or unpredictable workloads.

Amazon EC2 Reserved instance pricing allows you to reserve computing capacity for 1-year to 3-year terms at a significantly discounted hourly rate. Reserved instances are a billing discount and capacity reservation that is applied to instances to lower hourly running costs. A Reserved instance is not a physical instance. The discounted usage price is fixed for as long as you own the Reserved instance, allowing you to predict compute costs over the term of the reservation. If you are expecting consistent, heavy, use, Reserved instances can provide substantial savings over owning your own hardware or running only On-Demand instances.

Scheduled Reserved instances enable you to purchase capacity reservations that recur on a daily, weekly, or monthly basis, with a specified duration, for a 1-year term. You reserve the capacity in advance, so that you know it is available when you need it. You pay for the time the instances are scheduled, even if you do not use them. Scheduled instances are a good choice for workloads that do not run continuously, but do run on a regular schedule and take a finite time to complete.

Spot instances enable you to bid on unused EC2 instances, which can lower your Amazon EC2 costs significantly. The hourly price for a Spot instance (of each instance type in each Availability Zone) is set by Amazon EC2, and fluctuates depending on the supply of and demand for Spot instances. Your Spot instance runs whenever your bid exceeds the current market price.

Spot instances are a cost-effective choice if you can be flexible about when your applications run and if your applications can be interrupted. Amazon EC2 does not terminate Spot instances with a specified duration (also known as Spot blocks) when the Spot price changes. This makes them ideal for jobs that take a finite time to complete, such as batch processing, encoding and rendering, modeling and analysis, and continuous integration.

An Amazon EC2 Dedicated Host is a physical server with EC2 instance capacity fully dedicated to your use. Dedicated Hosts allow you to use your existing per-socket, per-core, or per-VM software licenses, including Microsoft Windows Server, Microsoft SQL Server, SUSE, Linux Enterprise Server, and so on. Dedicated Hosts and Dedicated instances can both be used to launch Amazon EC2 instances onto physical servers that are dedicated for your use. There are no performance, security, or physical differences between Dedicated instances and instances on Dedicated Hosts. However, Dedicated Hosts give you additional visibility and control over how instances are placed on a physical server.

STORAGE SERVICES

Amazon S3 and Amazon EBS

Storage Services

Understand storage options including:

- ✓ Amazon S3 (Requests, Buckets, Objects, Access, Protecting Data, Notifications, Replication, Request Routing, Optimization, Lifecycle Management with Glacier)
- ✓ Amazon EBS (Volumes, Snapshots, Optimization, Encryption, Performance)

Amazon Simple Storage Services

- ✓ Storage for the Internet
- ✓ Natively online, HTTP access
- ✓ Store and retrieve any amount of data, any time, from anywhere on the web
- ✓ Highly scalable, reliable, fast and durable

Amazon S3 is designed to make web-scale computing easier for developers. Amazon S3 provides a simple web services interface that can be used to store and retrieve any amount of data, at any time, from anywhere on the web. It gives any developer access to the same highly scalable, reliable, secure, fast, inexpensive infrastructure that Amazon uses to run its own global network of websites.

Amazon S3 Facts

Here are some facts about Amazon S3. There is a 100-bucket limit per account. You can store unlimited number of objects in a bucket. The size of an object can be up to 5 TB and there is no limit to the size of a bucket. Amazon S3 is designed for 99.999999999% durability and 99.99% availability of objects over a given year. You can use HTTP or HTTPs endpoints to store and retrieve any amount of data, at any time, from anywhere on the web. Most importantly, Amazon S3 is highly scalable, reliable, fast, and inexpensive.

Common use scenarios

- ✓ Storage and Backup
- ✓ Application file Hosting
- ✓ Media Hosting
- ✓ Software Delivery
- ✓ Store AMIs and Snapshots

Advanced use scenarios:

Using Amazon DevPay with Amazon S3: Amazon DevPay enables you to charge customers for using your Amazon S3 product through Amazon's authentication and billing infrastructure. You can charge any amount for your product including usage charges (storage, transactions, and bandwidth), monthly fixed charges, and a one- time charge.

- ✓ **Publishing Content Using Amazon S3 and BitTorrent:** You can direct your clients to your BitTorrent accessible objects by giving them the .torrent file directly or by publishing a link to the BitTorrent URL of your object.
- ✓ **Hosting a Static Website on Amazon S3:** You can host a static website on Amazon S3 by configuring a bucket for website hosting and then uploading your website content to the bucket.

Amazon S3 Pricing

Amazon S3 pricing is based on capacity and bandwidth actually used. Since Amazon S3 is an Internet-scale service that runs natively across an entire region; it can handle significant request throughput and bandwidth output. All bandwidth into Amazon S3 is free, but AWS charges a rate on bandwidth out. Most importantly, since Amazon S3 can handle any amount of data, it is important to note that you only pay for the amount of space you use. Prices are based on a prorated GB per month.

There is also a pricing calculator online as a reference. Note that pricing listed is in the US East (N. Virginia) Region at the time this training was developed.

Amazon S3 Concepts

To get the most out of Amazon S3, you need to understand a few simple concepts. First, Amazon S3 stores data as objects within buckets.

An object is composed of a file, and any metadata that describes that file. To store an object in Amazon S3, you upload the file you want to store into a bucket. When you upload a file, you can set permission on the object as well as any metadata.

Buckets are logical containers for objects. You can have one or more buckets in your account. For each bucket, you can control access, in other words, who can create, delete and list objects in the bucket. You can also view access logs for the bucket, and its objects, and choose the geographical region where Amazon S3 will store the bucket and its contents.

Amazon S3 Buckets

- ✓ Organize the Amazon S3 namespace at the highest level.
- ✓ Identify the account responsible for storage and data transfer charges.
- ✓ Play a role in access control.
- ✓ Serve as the unit of aggregation for usage reporting.
- ✓ Have globally unique bucket names, regardless of the AWS region in which they were created.

A bucket is a logical container for objects stored in Amazon S3. Every object is contained in a bucket. Buckets serve several purposes: They organize the Amazon S3 namespace at the highest level, they identify the account responsible for storage and data transfer charges, they play a role in access control, and they serve as the unit of aggregation for usage reporting. Amazon S3 bucket names are globally unique, regardless of the AWS Region in which you create the bucket. You specify the name at the time you create the bucket.

Object Keys

An object key is the identifier for an object in a bucket.

`http://doc.s3.amazonaws.com/2006-03-01/AmazonS3.html`

Bucket Object/Key

Because the combination of a bucket, key, and version ID uniquely identify each object, Amazon S3 can be thought of as a basic data map between "bucket + key + version" and the object itself. Every object in Amazon S3 can be uniquely addressed through the combination of the web service endpoint, bucket name, key, and optionally, a version.

Amazon S3 Security

- ✓ You can control access to buckets and objects with:
 - ✓ Access Control Lists (ACLs)
 - ✓ Bucket policies
 - ✓ Identity and Access Management (IAM) policies
- ✓ You can upload or download data to Amazon S3 via SSL encrypted endpoints.
- ✓ You can encrypt data using AWS SDKs.

Data Access:

- ✓ **IAM policies:** With IAM policies, you can only grant users within your own AWS account permission to access your Amazon S3 resources.
- ✓ **ACLs:** With ACLs, you can only grant other AWS accounts (not specific users) access to your Amazon S3 resources.
- ✓ **Bucket Policies:** Bucket policies in Amazon S3 can be used to add or deny permissions across some or all of the objects within a single bucket. Policies can be attached to users, groups, or Amazon S3 buckets, enabling centralized management of permissions. With bucket policies, you can grant users within your AWS account or another AWS account access to your Amazon S3 resources.

Data Transfer:

For maximum security, you can securely upload/download data to Amazon S3 via the SSL encrypted endpoints. The encrypted endpoints are accessible both from the Internet and from within Amazon EC2 so that data is transferred securely both within AWS and to and from sources outside of AWS.

Data Storage:

Amazon S3 provides multiple options for protecting data at rest. Customers who prefer to manage their own encryption keys can use a client encryption library like the Amazon S3 Encryption Client to encrypt data before uploading to Amazon S3.

Alternatively, you can use Amazon S3 Server Side Encryption (SSE) if you prefer to have Amazon S3 manage encryption keys for you. With Amazon S3 SSE, you can encrypt data on upload simply by adding an additional request header when writing the object. Decryption happens automatically when data is retrieved.

Amazon S3 SSE uses one of the strongest block ciphers available: 256-bit Advanced Encryption Standard (AES-256). With Amazon S3 SSE, every protected object is encrypted with a unique encryption key. This object key itself is then encrypted with a regularly rotated master key. Amazon S3 SSE provides additional security by storing the encrypted data and encryption keys in different hosts. Amazon S3 SSE also makes it possible for you to enforce encryption requirements. For example, you can create and apply bucket policies that require that only encrypted data can be uploaded to your buckets.

Instead of using Amazon S3 SSE, you also have the option of encrypting your data before sending it to Amazon S3. You can build your own library that encrypts your object data on the client side before uploading it to Amazon S3. Optionally, you can use an AWS SDK to automatically encrypt your data before uploading it to Amazon S3.

Amazon S3 Versioning

- ✓ Protects from accidental overwrites and deletes with no performance penalty.
- ✓ Generates a new version with every upload.
- ✓ Allows easily retrieval of deleted objects or roll back to previous versions.
- ✓ Three states of an Amazon S3 bucket
 - ✓ Un-versioned (default)
 - ✓ Versioning-enabled
 - ✓ Versioning-suspended

Versioning is a means of keeping multiple variants of an object in the same bucket. You can use versioning to preserve, retrieve, and restore every version of every object stored in your Amazon S3 bucket. With versioning, you can easily recover from both unintended user actions and application failures.

In one bucket, for example, you can have two objects with the same key, but different version IDs, such as photo.gif (version 111111) and photo.gif (version 121212).

Once you version-enable a bucket, it can never return to an unversioned state. You can, however, suspend versioning on that bucket.

Amazon S3 Storage Classes

Storage Class	Durability	Availability	Other Considerations
Amazon S3 Standard	99.999999999%	99.99%	None
Amazon S3 Standard - Infrequent Access (IA)	99.999999999%	99.99%	<ul style="list-style-type: none"> • Retrieval fee associated with objects • Most suitable for infrequently accessed data
Glacier	99.999999999%	99.99% (after you restore objects)	<ul style="list-style-type: none"> • Not available for real-time access • Must restore objects before you can access them

Each object in S3 has a storage class associated with it.

S3 Standard is ideal for performance-sensitive use cases and frequently used data. Standard is the default storage class in S3.

S3 Infrequent Access (IA) is optimized for long-lived and less frequently accessed data such as backups and older data that are accessed less but still require high performance.

Glacier is suitable for archiving data where access is infrequent and a retrieval time of several hours is acceptable. Archived objects are not available for real-time access – they must be restored before they can be accessed. The Glacier storage class is very low-cost.

S3 Reduced Redundancy Storage (RSS) is designed for noncritical, reproducible data stored at lower levels of redundancy standards than the Standard or IA classes, reducing cost.

Amazon S3 Object Lifecycle

Lifecycle management defines how Amazon S3 manages objects during their lifetime. Some objects that you store in an Amazon S3 bucket might have a well-defined lifecycle:

- ✓ Log files
- ✓ Archive documents
- ✓ Digital media archives
- ✓ Financial and healthcare records
- ✓ Raw genomics sequence data
- ✓ Long-term database backups
- ✓ Data that must be retained for regulatory compliance

Lifecycle management defines how Amazon S3 manages objects during their lifetime. Some objects that you store in an Amazon S3 bucket might have a well-defined lifecycle: if you are uploading periodic logs to your bucket, your application might need these logs for a week or a month after creation, and after that you might want to delete them. Some documents are frequently accessed for a limited period of time. After that, you might not need real-time access to these objects, but your organization might require you to archive them for a longer period and then optionally delete them.

Digital media archives, financial and healthcare records, raw genomics sequence data, long-term database backups, and data that must be retained for regulatory compliance are some of the kinds of objects that you might upload to Amazon S3 primarily for archival purposes.

When you configure a lifecycle rule, you specify the storage class you want to transition the object to and the number of days after object creation to transition it. You can transition objects to the Standard – Infrequent Access (IA) storage class, archive them to Amazon Glacier, or have them permanently deleted. Standard - IA is useful for data such as backups and other older, infrequently accessed data where high performance continues to be a requirement. It is suitable for objects greater than 128 kilobytes that you want to keep for at least 30 days. There is a retrieval fee associated with Standard - IA objects.

Amazon Glacier

- ✓ Long term low-cost archiving service
- ✓ Optimal for infrequently accessed data
- ✓ Designed for 99.999999999% durability
- ✓ 3-5 hours retrieval time
- ✓ Less than \$0.01 per GB / month (depending on region)

Sound Cloud Case Study

SoundCloud:

- ✓ Operates worldwide.
- ✓ Enables users to upload 12 hours of audio material to its platform every minute.
 - ✓ Each audio file must be transcoded and stored in multiple formats
 - ✓ Logs and analyzes billions of events.

The AWS Solution:

- ✓ SoundCloud uses a storage solution comprised of:
 - ✓ Amazon S3
 - ✓ Amazon Glacier
- ✓ The audio files are:
 - ✓ Placed in Amazon S3.
 - ✓ Distributed from Amazon S3 via the SoundCloud website.
 - ✓ Copied to Amazon Glacier.

The company currently stores 2.5 PB of data on Amazon Glacier.

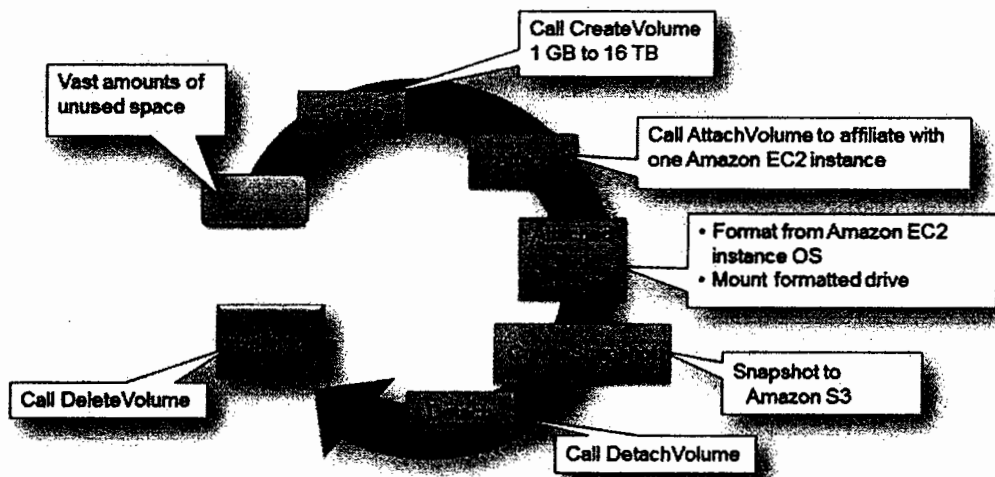
The SoundCloud case study demonstrates one of many ways in which Amazon S3 can be leveraged to securely store data.

Amazon Elastic Block Store (EBS)

- ✓ **Persistent block level storage** volumes offering consistent and low-latency performance
- ✓ Automatically replicated within its Availability Zone
- ✓ Snapshots stored durably in Amazon S3

Amazon Elastic Block Store, also known as Amazon EBS, provides persistent block-level storage volumes for use with Amazon EC2 instance offering consistent and low-latency performance. Amazon EBS is particularly suited for applications that require a database, file system, or access to raw block-level storage. Amazon EBS snapshots are durable and automatically replicated within their Availability Zone. Snapshots can be stored in Amazon S3.

Amazon EBS Lifecycle



Amazon EBS provides block-level storage volumes for use with Amazon EC2 instances. Amazon EBS volumes are highly available and reliable storage volumes that can be attached to any running instance in the same Availability Zone. The Amazon EBS volumes attached to an Amazon EC2 instance are exposed as storage volumes that persist independently from the life of the instance. When the volumes are not attached to an EC2 instance, you pay only for the cost of storage.

Amazon EBS Facts

- ✓ You can create:
 - ✓ **EBS Magnetic** volumes from 1 GiB to 1 TiB in size.
 - ✓ **EBS General Purpose** (SSD) and **Provisioned IOPS** (SSD) volumes up to 16 TiB in size.
- ✓ You can use encrypted EBS volumes to meet a wide range of data at-rest encryption requirements for regulated/audited data and applications.
- ✓ You can create point-in-time snapshots of EBS volumes, which are persisted to Amazon S3.

Amazon Elastic Block Store (Amazon EBS) provides block level storage volumes for use with EC2 instances. EBS volumes are highly available and reliable storage volumes that can be attached to any running instance that is in the same Availability Zone. EBS volumes that are attached to an EC2 instance are exposed as storage volumes that persist independently from the life of the instance. With Amazon EBS, you pay only for what you use.

You can mount multiple volumes on the same instance, but each volume can be attached to only one instance at a time. EBS volumes behave like raw, unformatted block devices. You can create a file system on top of these volumes, or use them in any other way you would use a block device (like a hard drive).

Amazon EBS is recommended when data changes frequently and requires long-term persistence. EBS volumes are particularly well-suited for use as the primary storage for file systems, databases, or for any applications that require fine granular updates and access to raw, unformatted, block-level storage. Amazon EBS is particularly helpful for database-style applications that frequently encounter many random reads and writes across the data set.

Amazon EBS Use Cases

- ✓ OS — Use for boot/root volume, secondary volumes
- ✓ Databases — Scales with your performance needs
- ✓ Enterprise applications— Provides reliable block storage to run mission-critical applications
- ✓ Business continuity — Minimize data loss and recovery time by regularly backing up using EBS Snapshots
- ✓ Applications — Install and persist any application

The Amazon EBS service is simply a virtual hard drive. So, a great use case for Amazon EBS is when you want the hard drive to persist past the life of the Amazon EC2 instance. Before Amazon EBS existed as a service, AWS only used physical local attached hard drives called ephemeral storage. The problem with that was that you couldn't stop an Amazon EC2 instance without losing all your data, because of the temporary nature of local storage.

That's why we created Amazon EBS to decouple the lifecycle of data persistence from the lifecycle of an EC2 instance. Amazon EBS volumes are ideal for root volumes you need to store and have block-level access to your operating system, database storage, and datasets that are smaller than 1 TB. Given its simple snapshot mechanism Amazon EBS is a great use case for simplifying distributed backups as well.

Amazon EBS Pricing

Pay for what you provision:

- ✓ Pricing based on region
- ✓ AWS GovCloud (US) Pricing page
- ✓ Review Pricing Calculator online
- ✓ Pricing is available as:
 - ✓ Storage
 - ✓ IOPS

Amazon EBS pricing is based on allocated storage, whether you use it or not. This is unlike Amazon S3, whose pricing is based on space actually in use. Prices may vary based on region or for IOPS.

Amazon EBS Scope

- ✓ Amazon EBS Volumes are in a Single Availability Zone
- ✓ Volume data is replicated across multiple servers in an Availability Zone.

Amazon EBS volumes are designed to be highly available and reliable. Amazon EBS volume data is replicated across multiple servers in an Availability Zone to prevent the loss of data from the failure of any single component.

The durability of your volume depends on both the size of your volume and the percentage of the data that has changed since your last snapshot.

Amazon EBS volumes are designed for an annual failure rate (AFR) of between 0.1% - 0.2%, where failure refers to a complete or partial loss of the volume, depending on the size and performance of the volume. This is compared with commodity hard disks that will typically fail with an AFR of around 4 percent, making EBS volumes 10 times more reliable than typical commodity.

Since Amazon EBS servers are replicated within a single Availability Zone, mirroring data across multiple Amazon EBS volumes in the same Availability Zone will not significantly improve volume durability.

For those interested in even more durability, with Amazon EBS you can create point-in-time consistent snapshots of your volumes that are then stored in Amazon S3, and automatically replicated across multiple Availability Zones.

Taking frequent snapshots of your volume is a convenient and cost-effective way to increase the long-term durability of your data. In the unlikely event that your Amazon EBS volume does fail, all snapshots of that volume will remain intact, and will allow you to recreate your volume from the last snapshot point.

Amazon EBS and Amazon S3

	Amazon EBS	Amazon S3
Paradigm	Block storage with file system	Object store
Performance	Very fast	Fast
Redundancy	Across multiple servers in an Availability Zone	Across multiple facilities in a Region
Security	EBS Encryption – Data volumes and Snapshots	Encryption
Access from the Internet?	No (1)	Yes (2)
Typical use case	It is a disk drive	Online storage

(1) Accessible from the Internet if mounted to server and set up as FTP, etc.

(2) Only with proper credentials, unless ACLs are world-readable

This table demonstrates significant differences between Amazon S3 and Amazon EBS. Amazon EBS volumes are network-attached hard drives that can be written to or read from at a block level. Amazon S3 is an object-level storage medium.

This means that you must write whole objects at a time. If you change one small part of a file, you must still rewrite the entire file in order to commit the change to Amazon S3. This can be very time-consuming if you have frequent writes to the same object.

Amazon S3 is optimized for write once/read many use cases. The other major difference is cost. With Amazon S3 you pay for what you use, and with Amazon EBS you pay for what you provision.

Amazon EC2 Instance Storage

- ✓ Local, complimentary direct attached block storage resource.
- ✓ Availability, number of disks, and size is based on EC2 instance type.
- ✓ Storage optimized instances for up to 365,000 Read IOPS and 315,000 First Write IOPS.
- ✓ SSD or magnetic.
- ✓ No persistence.
- ✓ All data is automatically deleted when an EC2 instance stops, fails or is terminated.

An instance store provides temporary block-level storage for your instance. This storage is located on disks that are physically attached to the host computer. Instance store is ideal for temporary storage of information that changes frequently, such as buffers, caches, scratch data, and other temporary content, or for data that is replicated across a fleet of instances, such as a load-balanced pool of web servers.

Reboot vs. Stop vs. Terminate

Characteristic	Reboot	Stop/Start (if EBS-backed instances only)	Terminate
Host computer	The instance stays on the same host computer.	The instance runs on a new host computer.	N/A
Private and public IP addresses	Stay the same.	Instance keeps its private IP address and gets a new public IP address.	N/A
Elastic IP addresses (EIP)	EIP remains associated with the instance.	EIP remains associated with the instance.	The EIP is disassociated from the instance.
Instance store volumes	The data is preserved.	The data is erased.	The data is erased.
EBS volume	The volume is preserved.	The volume is preserved.	The volume is deleted by default.
Billing	Instance billing hour doesn't change.	You stop incurring charges as soon as state is changed to stopping.	You stop incurring charges as soon as state is changed to shutting-down.

NETWORKING

Amazon VPC

Amazon Virtual Private Cloud (VPC)

- ✓ Provision a private, isolated virtual network on the AWS cloud.
- ✓ Have complete control over your virtual networking environment.

With Amazon Virtual Private Cloud (VPC), you can define a virtual network topology that closely resembles a traditional network that you might operate in your own data center. You have complete control over your virtual networking environment, and you can easily customize the network configuration for your Amazon VPC such as selection of IP address range, creation of subnets, configuration of route tables, and network gateways.

VPCs and Subnets

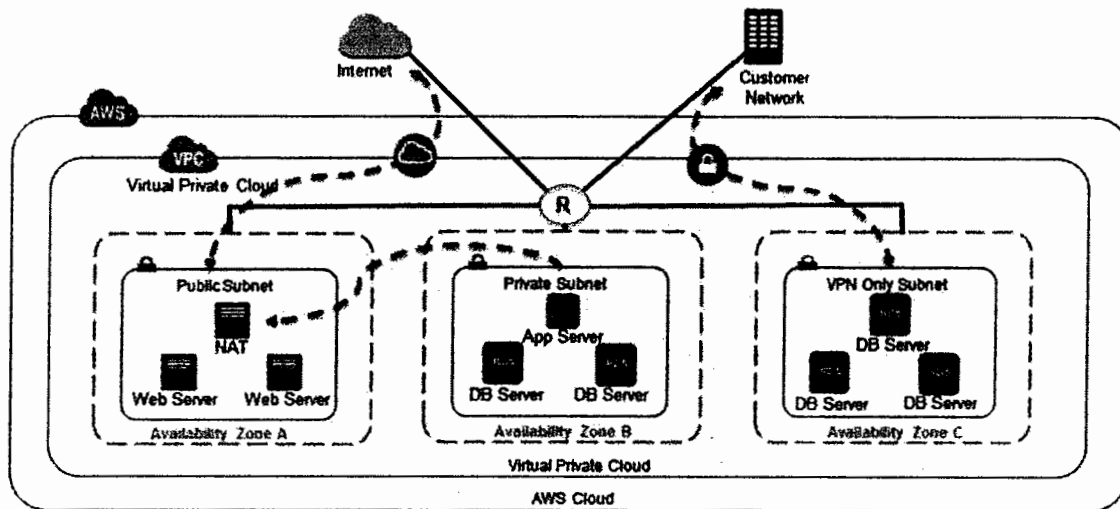
- ✓ A subnet defines a range of IP addresses in your VPC.
- ✓ You can launch AWS resources into a subnet that you select.
- ✓ A private subnet should be used for resources that won't be accessible over the Internet.
- ✓ A public subnet should be used for resources that will be accessed over the Internet.
- ✓ Each subnet must reside entirely within one Availability Zone and cannot span zones.

AWS assigns a unique ID to each subnet. Regardless of the type of subnet (public or private), the internal IP address range of the subnet is always private.

A public subnet has a route to an internet gateway (i.e., for a web server accessible from the Internet).

A private subnet has no route to an internet gateway (i.e., for a database server only accessed within the VPC).

Amazon VPC Example



Amazon Virtual Private Cloud also known as Amazon VPC, allows you provision a logically isolated section of the AWS cloud where you can launch AWS resources in a virtual network that you define. You have complete control over your virtual networking environment, including selection of your own IP address range, creation of subnets, configuration of route tables, network access control lists, and network gateways. You can easily customize the network and configuration for your Amazon VPC instance. For example, you can create a public-facing subnet for your web servers that require access to the Internet, and place your back end systems such as databases or application servers in a private-facing subnet with no Internet access. You can leverage multiple layers of security, including security groups and network access control lists, to help control access to Amazon EC2 instances in each subnet. Additionally, you can create a hardware virtual private network (VPN) connection between your corporate data center and your VPC, allowing you to leverage the AWS cloud as an extension of your corporate data center.

Security in your VPC

- ✓ Security groups
- ✓ Network access control lists (ACLs)

Amazon VPC provides three features that you can use to increase and monitor the security for your VPC:

- ✓ Security groups act as a firewall for associated Amazon EC2 instances, controlling both inbound and outbound traffic at the instance level.
- ✓ Network access controls lists (ACLs) act as a firewall for associated subnets, controlling both inbound and outbound traffic at the subnet level.

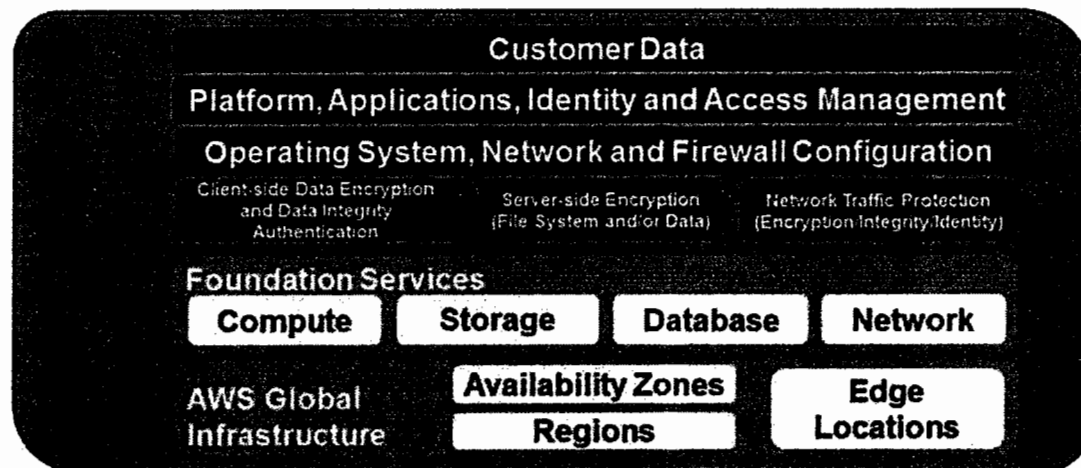
VPN Connections

VPN Connectivity option	Description
AWS Hardware VPN	You can create an IPsec, hardware VPN connection between your VPC and your remote network.
AWS Direct Connect	AWS Direct Connect provides a dedicated private connection from a remote network to your VPC.
AWS VPN CloudHub	You can create multiple AWS hardware VPN connections via your VPC to enable communications between various remote networks.
Software VPN	You can create a VPN connection to your remote network by using an Amazon EC2 instance in your VPC that's running a software VPN appliance.

Module 3

SECURITY, IDENTITY AND ACCESS MANAGEMENT

Shared Responsibility – AWS



When we talk about cloud security, we like to start with a discussion of the Shared Security Responsibility Model. While AWS takes care of provisioning and maintaining the underlying cloud infrastructure, you will still need to perform several security configuration tasks to ensure that you stay safe in the cloud. AWS's responsibility goes from the ground up to the hypervisor.

AWS secures the hardware, software, facilities, and networks that run all products and services. Customers are responsible for securely configuring the services they sign up for as well as anything they put on those services.

AWS also performs the following responsibilities:

- ✓ Obtaining industry certifications and independent third party attestations
- ✓ Publishing information about AWS security and control practices in whitepapers and web site content
- ✓ Providing certificates, reports, and other documentation directly to AWS customers under NDA (as required).

The amount of security configuration work you have to do varies, depending on how sensitive your data is and which services you select. For example, AWS services such as Amazon EC2 and Amazon S3 are completely under your control and require you to perform all of the necessary security configuration and management tasks. In the case of Amazon EC2, you are responsible for management of the guest OS (including updates and security patches), any application software or utilities you install on the instances, as well as the configuration of the AWS-provided firewall (called a security group) on each instance.

When you use any of AWS's managed services like Amazon RDS, Amazon RedShift, or Amazon WorkDocs, you don't have to worry about launching and maintaining instances or patching the guest OS or applications—AWS handles that for you. For these managed services, basic security configuration tasks like data backups, database replication, and firewall configuration happen automatically.

However, there are certain security features—such as IAM user accounts and credentials, SSL for data transmissions, and user activity logging—that you should configure no matter which AWS service you use.

AWS Support provides a highly personalized level of service for customers seeking technical help.

Physical Security

- ✓ 24/7 trained security staff
- ✓ AWS data centers in nondescript and undisclosed facilities
- ✓ Two-factor authentication for authorized staff
- ✓ Authorization for data center access

One of the main security responsibilities of AWS is the physical security of the data centers that house the AWS cloud infrastructure. Amazon has many years of experience designing, constructing, and operating large-scale data centers.

The physical security measures that protect these data centers are some of the most comprehensive in the industry and include: 24/7 trained security guards; locations in nondescript, undisclosed facilities; two-factor authentication for ingress; authorization for data center access only for an approved, specific need; and continuous monitoring, logging, and auditing of physical access controls.

Hardware, Software, and Network

- ✓ Automated change-control process
- ✓ Bastion servers that record all access attempts
- ✓ Firewall and other boundary devices
- ✓ AWS monitoring tools

The hardware and software that supports AWS cloud services has been architected to be not only highly available and redundant, but also extremely secure. All changes to AWS hardware

and software are managed through a centralized, automated change control process, and all access to hardware or software must be authorized.

Privileged access to software and systems requires SSH logon and is allowed only through bastion servers that record all access attempts. AWS network devices, including firewall and other boundary devices, monitor and control communications at the external boundary of the network and at key internal boundaries.

AWS monitoring tools are designed to detect unusual or unauthorized activities and conditions at ingress and egress communication points. These tools monitor server and network usage, port scanning activities, application usage, and unauthorized intrusion attempts. AWS security monitoring tools help identify several types of denial of service (DoS) attacks, including distributed, flooding, and software/logic attacks.

Certifications and Accreditations



AWS has successfully completed multiple audits, attestations, and certifications. AWS publishes a Service Organization Controls SOC 1 report, published under both the SSAE 16 and the ISAE 3402 professional standards, as SOC 2-Security and SOC 3 Report.

In addition, AWS has achieved ISO 9001, ISO 27001, ISO 27017 and ISO 27018 certifications, has been successfully validated as a Level 1 service provider under the Payment Card Industry (PCI) Data Security Standard (DSS), and currently offers HIPAA Business Associate Agreements to covered entities and their business associates subject to HIPAA.

In the realm of public sector certifications, AWS has achieved FedRAMP compliance, has received authorization from the U.S. General Services Administration to operate at the FISMA Moderate level, and is also the platform for applications with Authorities to Operate (ATOs) under the Defense Information Assurance Certification and Accreditation Program (DIACAP).

NIST, FIPS 140-2, CJIS, DoD SRG Levels 2 and 4 are some of the other certifications AWS has received.

SSL Endpoints

SSL Endpoints	Security Groups	VPC
Secure Transmission Establish secure communication sessions (HTTPS) using SSL/TLS.	Instance Firewalls Configure firewall rules for instances using Security Groups.	Network Control In your Virtual Private Cloud, create low-level networking constraints for resource access. Public and private subnets, NAT and VPN support.

AWS provides customer access points, also called API endpoints, that allow HTTPS access so that you can establish secure communication sessions with your AWS services including SSL and TLS. SSL encrypts the transmission, protecting each request or the response from being viewed in transit.

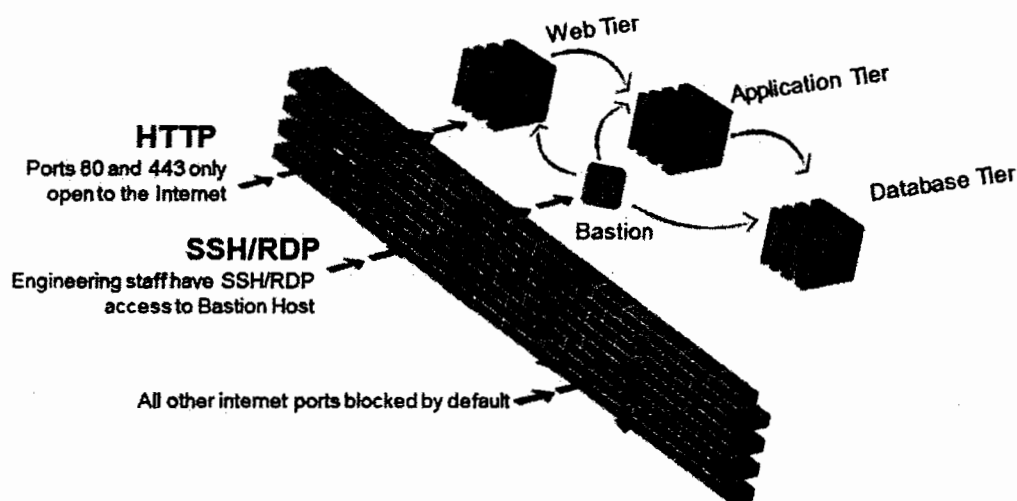
Security Groups

SSL Endpoints	Security Groups	VPC
Secure Transmission Establish secure communication sessions (HTTPS) using SSL/TLS.	Instance Firewalls Configure firewall rules for instances using Security Groups.	Network Control In your Virtual Private Cloud, create low-level networking constraints for resource access. Public and private subnets, NAT and VPN support.

AWS also provides security groups, which act like built-in firewalls for your virtual servers. You can control how accessible your instances are by configuring security group rules--from totally public to completely private, or somewhere in between. And when your instances reside within a Virtual Private Cloud (VPC) subnet, you can control egress as well as ingress traffic.

Security Groups can also be used by AWS services such as Amazon RDS, Amazon Redshift, Amazon EMR and Amazon ElastiCache.

AWS Multi-Tier Security Groups



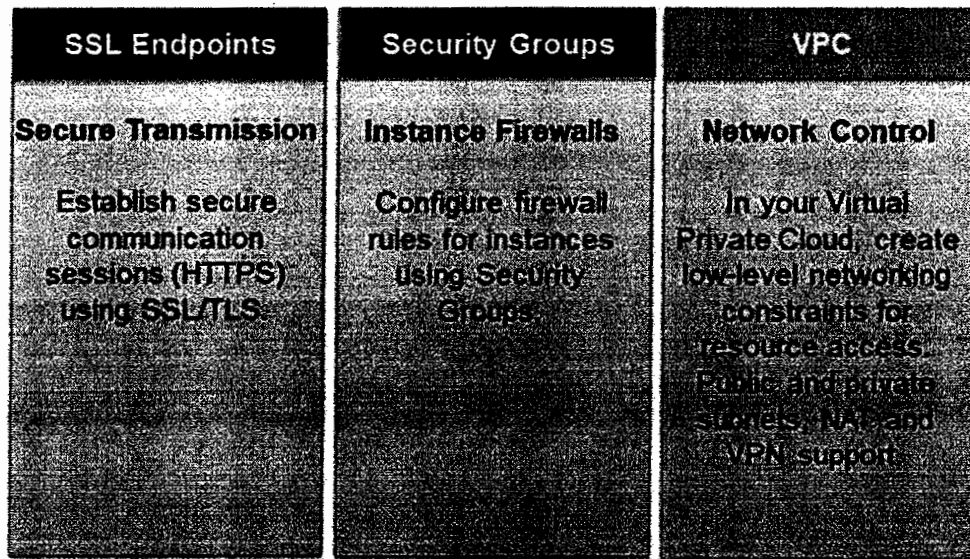
You can set up security group rules for your EC2 instances to create a traditional multi-tiered web architecture:

The web tier security group can accept traffic on port 80/443 from anywhere on the Internet if you select source 0.0.0.0/0. Alternatively, it might make more sense to only accept traffic from a load balancer so that individual clients cannot overload a single server and the load balancer can perform its job.

Similarly, the app tier can only accept traffic from the web tier, and the DB tier can only accept traffic from the app tier.

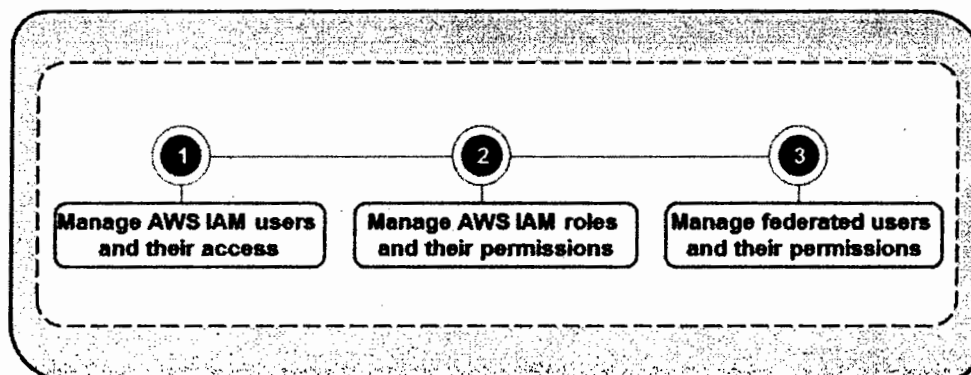
Lastly, we have also added a set of rules to allow remote administration over SSH port 22. We have restricted remote access by funneling all traffic through the app tier and allowing access only from a specific IP. After you use SSH to access an app tier server, you can then connect to machines on the web and DB security groups.

Amazon Virtual Private Cloud (VPC)



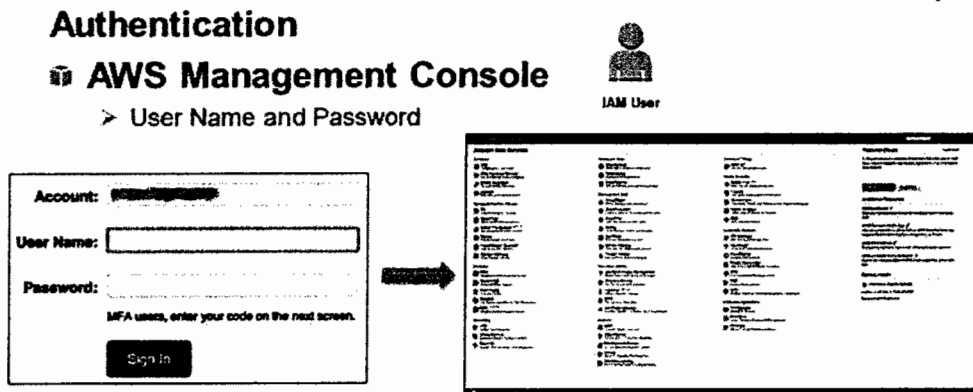
The Amazon Virtual Private Cloud (VPC) service allows you to add another layer of network security to your instances by creating private subnets and even adding an IPsec VPN tunnel between your network and your VPC. Amazon VPC allows you to define your own network topology, including definitions for subnets, network access control lists, Internet gateways, routing tables, and virtual private gateways. The subnets that you create can be defined as either private or public.

AWS Identity and Access Management (IAM)



Using IAM you can create and manage AWS users and groups and use permissions to allow and deny their access to AWS resources. You can use existing corporate identities to grant secure access to AWS resources, such as Amazon S3 buckets, without creating any new AWS identities.

AWS IAM Authentication

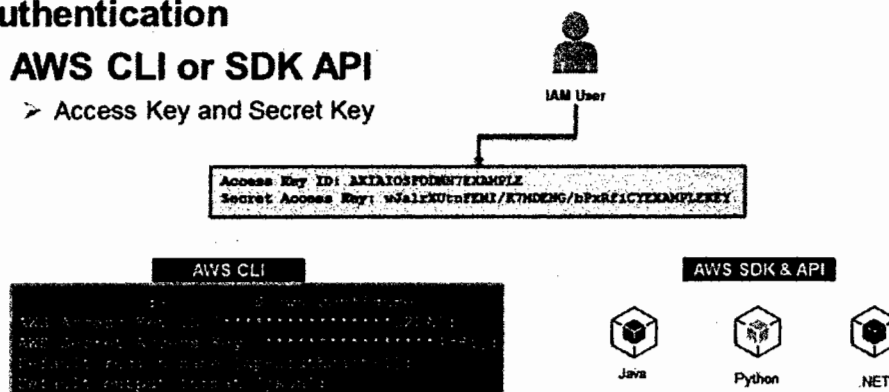


AWS services and resources can be accessed using the AWS Management Console, AWS CLI or through SDKs and APIs from a wide range of supported platforms. Users and systems have to be authenticated before they can access AWS services and resources.

The AWS Management Console provides a web-based way to administer AWS services. If you're the account owner, you can sign in to the console directly using the Root Account. It is, however, advisable to create individual IAM users for each user and login using individual credentials.

IAM is a complimentary service.

Authentication
AWS CLI or SDK API
➤ Access Key and Secret Key

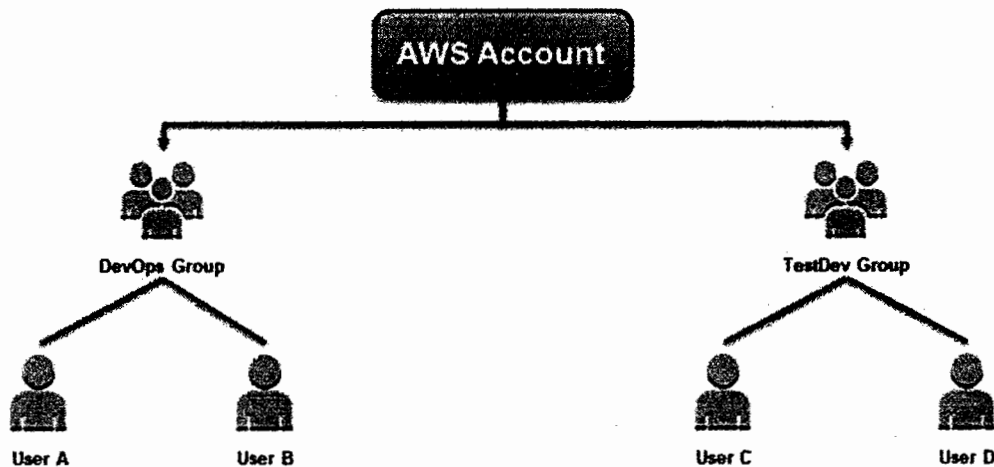


The AWS Command Line Interface is a unified tool to manage your AWS services. With AWS CLI, you can control multiple AWS services from the command line and automate them through scripts.

AWS CLI is supported on Windows, Linux, OS X, and Unix platforms.

AWS offers support for a wide variety of programming platforms including .NET, Java, Python etc.

AWS IAM User Management - Groups



As the number of users managing your AWS environment increases, it is helpful to manage permissions for multiple IAM users using IAM groups.

AWS IAM Authentication

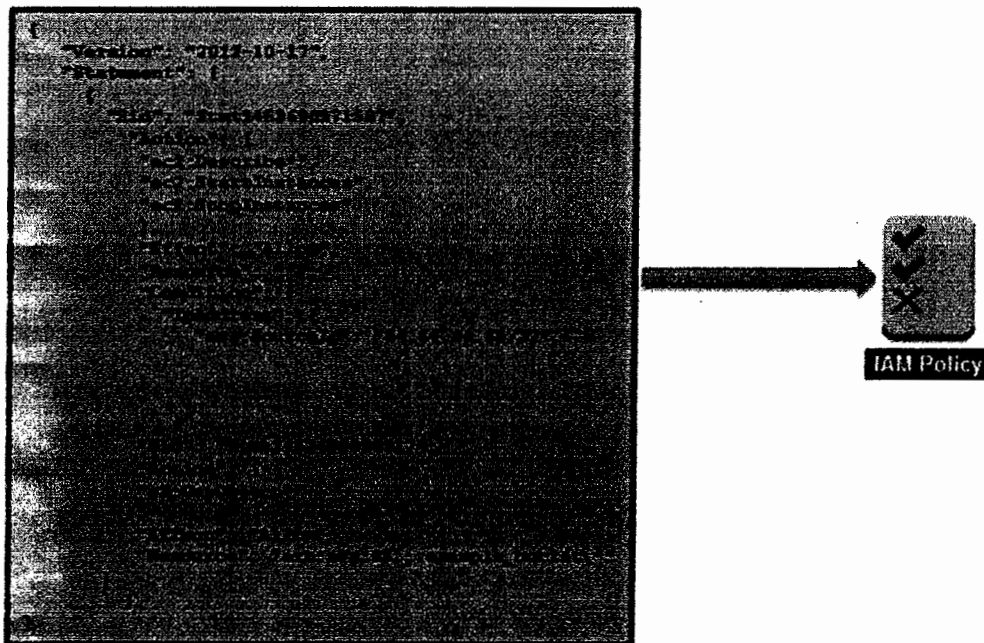
Policies:

- ✓ Are JSON documents to describe permissions.
- ✓ Are assigned to Users, Groups or Roles.

After a user or system has been authenticated, they have to be authorized to access AWS services. To assign permissions to a user, group, role, or resource, you create a policy, which is a document that explicitly lists permissions.

An IAM role is similar to a user, in that it is an AWS identity with permission policies that determine what the identity can and cannot do in AWS. However, instead of being uniquely associated with one person, a role is intended to be assumable by anyone who needs it. Also, a role does not have any credentials (password or access keys) associated with it. Instead, if a user is assigned to a role, access keys are created dynamically and provided to the user.

AWS IAM Policy Elements



Policies are documents that are created using JavaScript Object Notation (JSON). A policy consists of one or more statements, each of which describes one set of permissions.

An IAM policy may consist of:

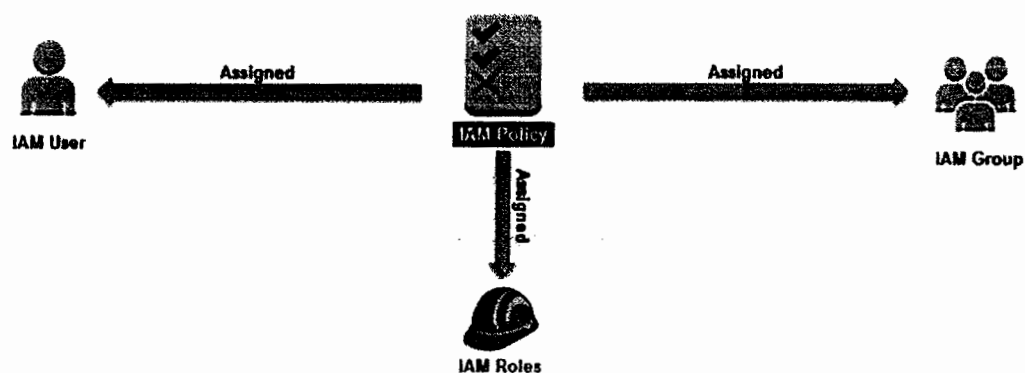
- ✓ **Version**
- ✓ **Id**
- ✓ **Statement**
- ✓ **Sid**
- ✓ **Effect:** Defines what the effect will be when the user requests access—either allow or deny. Because the default is that resources are denied to users, you typically specify that you will allow users access to resource.
- ✓ **Principal**
- ✓ **NotPrincipal**
- ✓ **Actions:** Defines what actions you want to allow. Each AWS service has its own set of actions. Any actions that you do not explicitly allow are denied.
- ✓ **NotAction**
- ✓ **Resources:** Defines which resources you allow the action on. Users cannot access any resources that you have not explicitly granted permissions to.
- ✓ **NotResource**
- ✓ **Condition**
- ✓ **Supported Data Types**

- ✓ **AWS Policy Generator:** You can use the AWS Policy Generator to generate policies at ease.
- ✓ **AWS Policy Validator:** Policy Validator automatically examines your existing IAM access control policies to ensure that they comply with the IAM policy grammar.
- ✓ **AWS Policy Simulator:** The simulator evaluates the policies that you choose and determines the effective permissions for each of the actions that you specify. The simulator uses the same policy evaluation engine that is used during real requests to AWS services.
- ✓ **Managed policies** – Are standalone policies that you can attach to multiple users, groups, and roles in your AWS account. Managed policies apply only to identities (users, groups, and roles) - not resources. You can use two types of managed policies:
 - ✓ **AWS managed policies** – Managed policies that are created and managed by AWS. If you are new to using policies, it is recommended that you start by using AWS managed policies.
 - ✓ **Customer managed policies** – Managed policies that you create and manage in your AWS account. Using customer managed policies, you have more precise control over your policies than when using AWS managed policies.
- ✓ **Inline policies** – Policies that you create and manage, and that are embedded directly into a single user, group, or role.

AWS IAM Policy Assignment



IAM Policies are assigned to IAM users and Groups. These users are bound by the permissions defined in the IAM Policy.



IAM Policies may also be assigned to an IAM Role.

An IAM role is similar to a user, in that it is an AWS identity with permission policies that determine what the identity can and cannot do in AWS. However, instead of being uniquely associated with one person, a role is intended to be assumable by anyone who needs it. Also, a role does not have any credentials (password or access keys) associated with it. Instead, if a user is assigned to a role, access keys are created dynamically and provided to the user.

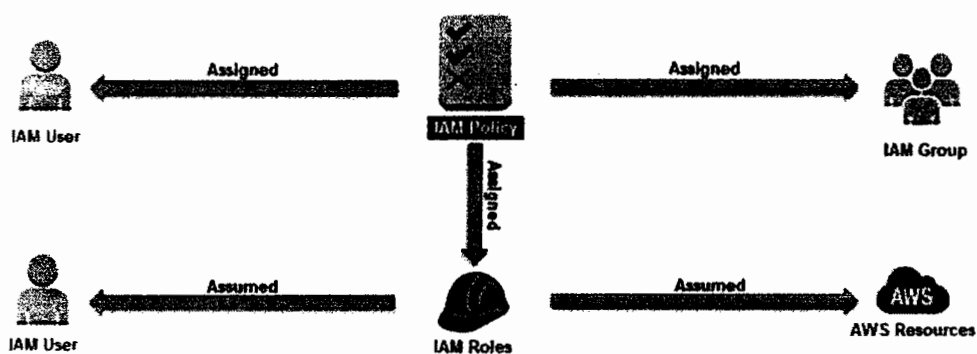
AWS IAM Roles

- ✓ An IAM role uses a policy.
- ✓ An IAM role has no associated credentials.
- ✓ IAM users, applications, and services may assume IAM roles.

IAM Policies may also be assigned to an IAM role.

An IAM role is similar to a user, in that it is an AWS identity with permission policies that determine what the identity can and cannot do in AWS. However, instead of being uniquely associated with one person, a role is intended to be assumable by anyone who needs it. Also, a role does not have any credentials (password or access keys) associated with it. Instead, if a user is assigned to a role, access keys are created dynamically and provided to the user.

AWS IAM Policy Assignment



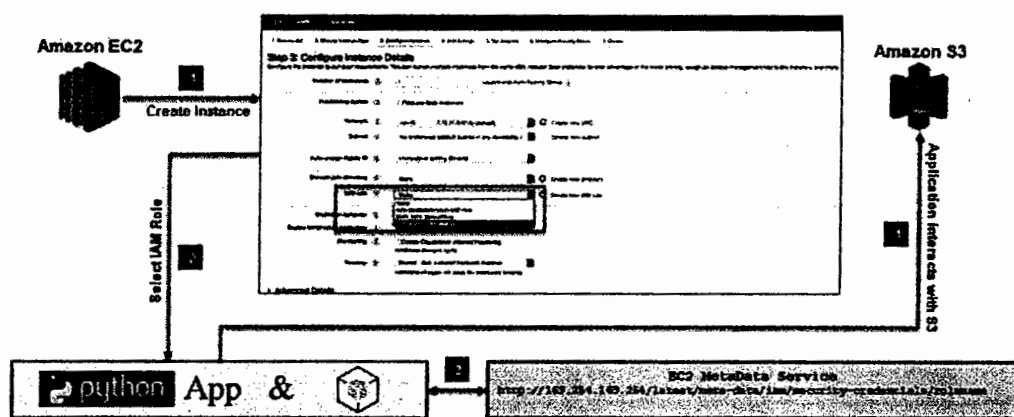
You can use roles to delegate access to users, applications, or services that don't normally have access to your AWS resources.

Application Access to AWS Resources

- ✓ Python application hosted on an Amazon EC2 Instance needs to interact with Amazon S3.
- ✓ AWS credentials are required:
 - ✓ Option-1 Store AWS Credentials on the Amazon EG2 instance.
 - ✓ Option 2: Securely distribute AWS credentials to AWS Services and Applications.

In the example above a custom application written in Python and hosted on an Amazon EC2 instance needs to interact with objects stored in an Amazon S3 bucket. Applications may access AWS resources in multiple ways. One way is to embed your AWS access key ID and secret access key in the application code or in a Config file supported by the application. However, doing so may compromise the user's credentials. Changing or rotating the user's credentials would require an update in the code each time. This approach is not secure and feasible in many cases. The alternate and secure option is to use an IAM role to pass temporary security credentials as part of an instance profile.

AWS IAM Roles – Instance Profiles



An instance profile is a container for an IAM role that you can use to pass role information to an EC2 instance when the instance starts.

In the example, an IAM role named `PythonInEC2AccessS3` is created by an IAM user. The role grants access to an Amazon S3 bucket.

Step 1: An application developer selects the PythonInEC2AccessS3 role while creating the Amazon EC2 instance. The instance would host a Python application which would need access to an Amazon S3 bucket.

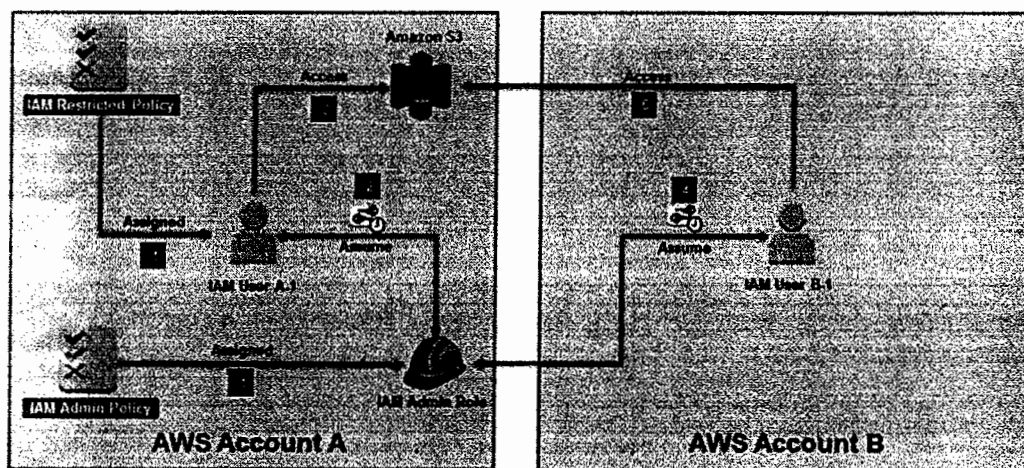
Note: An IAM role may be associated with an EC2 instance only during creation. The policy associated with the role may be modified at any time. A user launching an EC2 instance also needs appropriate permissions to associate an IAM role to the EC2 instance.

Step 2: Post instance creation the Python application is installed on the EC2 instance. AWS SDK for Python (Boto3) is also installed on the instance. The application tries to access an Amazon S3 bucket. However AWS credentials are not available on the instance.

Step 3: The Python application then uses the EC2 metadata service to gain access to Temporary Security Credentials. Temporary Security Credentials will be discussed later.

Step 4: The application interacts with the Amazon S3 bucket specified in the PythonInEC2AccessS3 role.

AWS IAM Roles – Assume Role



IAM roles may also be associated with users.

In the above example, there are two AWS accounts A and B. IAM User A-1 is part of Account A and IAM User B-1 is part of Account B.

Step 1: An IAM policy named IAM Admin Policy with access to an Amazon S3 bucket is associated to an IAM role named IAM Admin Role. User A-1 has an IAM policy with restricted access. This is done as User A-1 does not normally need administrative privileges. However, User A-1 may sometimes have to perform tasks that require administrative privileges.

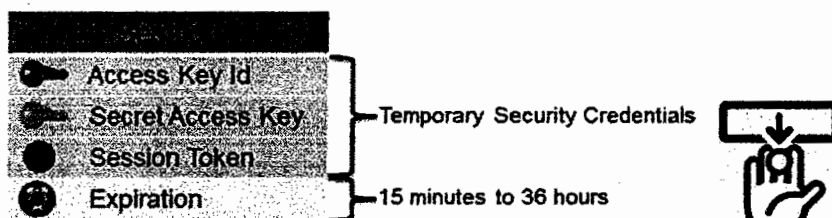
Step 2: When required User A-1 assumes the IAM Admin Role. Doing so gives User A-1 access to the S3 bucket. A user who assumes a role temporarily gives up his or her own permissions and instead takes on the permissions of the role. When the user exits, or stops using the role, the original user permissions are restored. It is therefore helpful to use IAM roles instead of changing the user's policies each time a change is required.

Note: User A-1's policy must contain permissions to assume the role. Step 3: User A-1 gains access to the Amazon S3 bucket.

Step 4: With IAM roles, you can establish trust relationships between your trusting account and other AWS trusted accounts. The trusting account owns the resource to be accessed and the trusted account contains the users who need access to the resource. User B-1 from Account B assumes the IAM Admin Role from Account A.

Step 5: User B-1 gains access to the Amazon S3 bucket owned by Account A.

Temporary Security Credentials (AWS STS)



Use Cases:

- ✓ Cross account access
- ✓ Federation
- ✓ Mobile Users
- ✓ Key rotation for Amazon

AWS Security Token Service (AWS STS) provides trusted users with temporary security credentials that can control access to your AWS resources. These credentials are short-term and work almost identically to the long-term access key credentials. These credentials are generated dynamically and provided to the user when requested.

A session established with AWS STS consists of an access key ID, secret access key, a session token, and an expiration time. The expiration time could last between 15 minutes to 36 hours. The keys are used to sign API requests and pass in the token as an additional parameter, which AWS uses to verify that the temporary access keys are valid.

Application Authentication



AWS IAM is not appropriate for OS and application authentication.

AWS IAM Authentication and Authorization

IAM is a powerful service to authenticate and authorize users and AWS resources.

AWS IAM Best Practices

- ✓ Delete AWS account (root) access keys.
- ✓ Create individual IAM users.
- ✓ Use groups to assign permissions to IAM users.
- ✓ Grant least privilege.
- ✓ Configure a strong password policy.
- ✓ Enable MFA for privileged users.
- ✓ Use roles for applications that run on Amazon EC2 instances.
- ✓ Delegate by using roles instead of by sharing credentials.
- ✓ Rotate credentials regularly.
- ✓ Remove unnecessary users and credentials.
- ✓ Use policy conditions for extra security.
- ✓ Monitor activity in your AWS account.

AWS Resources – Based Policies

- ✓ Are an alternative to IAM and supported by some services.
- ✓ Grant cross-account access to your resources.
- ✓ Use a principal to uniquely identify account in the policy.
- ✓ Supported AWS services include :
 - ✓ Amazon S3 Bucket Policy
 - ✓ Amazon SNS Topic Policy
 - ✓ Amazon SQS Queue Policy
 - ✓ Amazon Glacier Vault Policy
 - ✓ AWS OpsWorks Stack Policy
 - ✓ AWS Lambda Function Policy

For some AWS services, you can grant cross-account access to your resources. To do this, you attach a policy directly to the resource that you want to share, instead of using a role as a proxy. The resource that you want to share must support resource-based policies. Unlike a user-based policy, a resource-based policy specifies who (in the form of a list of AWS account ID numbers) can access that resource. Cross-account access with a resource-based policy has an advantage over a role. With a resource that is accessed through a resource-based policy, the user still works in the trusted account and does not have to give up his or her user permissions in place of the role permissions. In other words, the user continues to have access to resources in the trusted account at the same time as he or she has access to the resource in the trusting account. This is useful for tasks such as copying information to or from the shared resource in the other account.

Principal: This element defines an account in a policy. In a resource-based policy, the principal may refer to the same account or another account.

Module 4

DATABASES

Databases

Understand the concepts of fundamental AWS database services including:

- ✓ **Amazon Relational Database Service (RDS)**

- ✓ DB Instances
- ✓ Security Groups
- ✓ DB Parameter Groups
- ✓ DB Option Groups
- ✓ RDS Interfaces

- ✓ **Amazon DynamoDB**

- ✓ DynamoDB Data Model
- ✓ Supported Operations
- ✓ Provisioned Throughput
- ✓ Accessing DynamoDB

SQL and NoSQL Databases

	SQL	NoSQL
Data Storage	Rows and Columns	Key-Value
Schemas	Fixed	Dynamic
Querying	Using SQL	Focused on collection of documents
Scalability	Vertical	Horizontal

SQL

ISBN	Title	Author	Format
9182932465265	Cloud Computing Concepts	Wilson, Joe	Paperback
3142536475869	The Database Guru	Gomez, Maria	eBook

NoSQL

```
{
  ISBN: 9182932465265,
  Title: "Cloud Computing Concepts",
  Author: "Wilson, Joe",
  Format: "Paperback"
}
```

A SQL database stores data in rows and columns. Rows contain all the information about one entry and columns are the attributes that separate the data points. A SQL database schema is fixed – columns must be locked before data entry. Schemas can be amended if the database is altered entirely and taken offline. Data in SQL databases is queried using SQL (Structure Query Language), which can allow for complex queries. SQL databases scale vertically, by increasing hardware power.

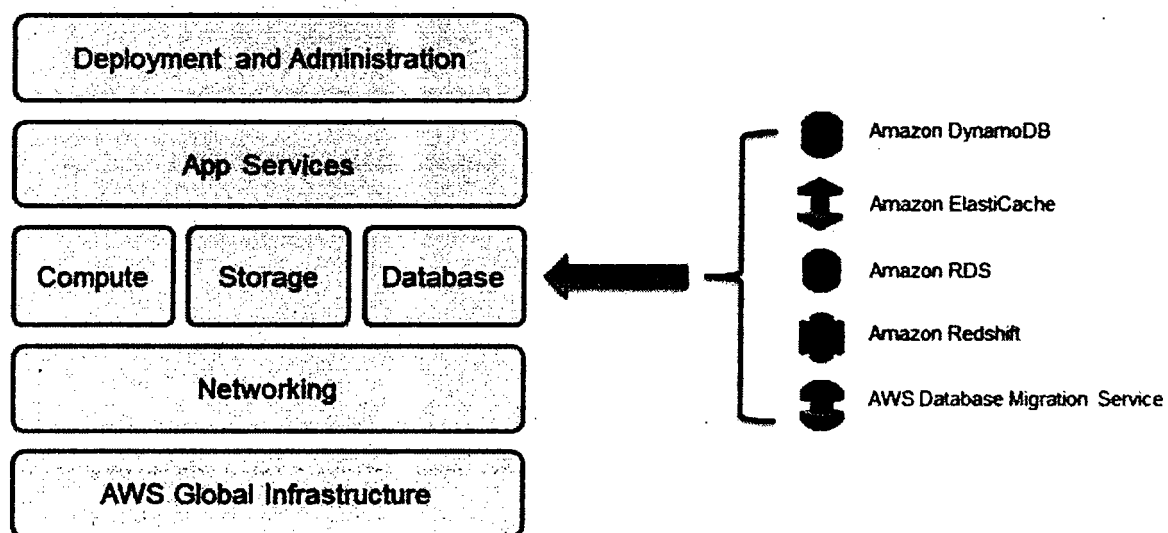
NoSQL databases store data using one of many storage models including key-value pairs, documents, and graphs. NoSQL schemas are dynamic and information can be added on the fly. Each 'row' doesn't have to contain data for each 'column'. Data in NoSQL databases is queried by focusing on collections of documents. NoSQL databases scale horizontally, by increasing servers.

Data storage Considerations

- ✓ No one size fits all.
- ✓ Analyze your data requirements by considering:
 - ✓ Data formats
 - ✓ Data size Query frequency
 - ✓ Data access speed
 - ✓ Data retention period

No one size fits all when considering database types. You must take into consideration your data requirements such as data formats, data size, and the frequency of your queries, how quickly you need your data, and for how long you need to keep it.

AWS Managed Database Services



Amazon Relational Database Service (RDS)

- ✓ Cost-efficient and **resizable capacity**
- ✓ Manages time-consuming **database administration** tasks
- ✓ Access to the full capabilities of **Amazon Aurora, MySQL, MariaDB, Microsoft SQL Server, Oracle, and PostgreSQL** databases

With Amazon RDS, you can access the full capabilities of a familiar MySQL, MariaDB, Microsoft SQL Server, Oracle, or PostgreSQL database. In addition, Amazon RDS for MySQL provides two distinct, but complementary, replication features: Multi-AZ deployments and read replicas that can be used in conjunction with each other to gain enhanced database availability, protect your latest database updates against unplanned outages, and scale beyond the capacity constraints of a single DB instance for read-heavy database workloads.

Amazon Aurora is a MySQL-compatible relational database engine that is part of Amazon RDS.

Amazon RDS Use Case

- ✓ Flipboard is an online magazine with millions of users and billions of "flips" per month.
- ✓ Flipboard is one of the world's first social media magazines.
- ✓ Flipboard uses Amazon RDS and its Multi-AZ capabilities to store **mission critical user data**.

Amazon RDS

- ✓ Simple and fast to deploy
- ✓ Manages common database administrative tasks
- ✓ Compatible with your applications
- ✓ Fast, predictable performance
- ✓ Simple and fast to scale
- ✓ Secure
- ✓ Cost-effective

Amazon RDS is a web service that makes it easy to set up, operate, and scale a relational database in the cloud. It provides cost-efficient and resizable capacity while managing time-consuming database administration tasks, freeing you up to focus on your applications and business. Amazon RDS gives you access to the full capabilities of a MySQL, Oracle, SQL Server, or Amazon Aurora database engine. This means that the code, applications, and tools you already use today with your existing databases can be used with Amazon RDS. Amazon RDS automatically patches the database software and backs up your database, storing the backups

for a user-defined retention period and enabling point-in-time recovery. You benefit from the flexibility of being able to scale the compute resources or storage capacity associated with your relational database instance via a single API call.

DB Instances

- ✓ DB Instances are the basic building blocks of Amazon RDS.
- ✓ They are an isolated database environment in the cloud.
- ✓ They can contain multiple user-created databases.

The basic building block of Amazon RDS is the DB instance. A DB instance is an isolated database environment in the cloud. A DB instance can contain multiple user-created databases, and you can access it by using the same tools and applications that you use with a stand-alone database instance. You can create and modify a DB instance by using the AWS Management Console, Amazon AWS command line interface, or the Amazon RDS API.

How Amazon RDS Backups Work

Automatic Backups:

- ✓ Restore your database to a point in time.
- ✓ Are enabled by default.
- ✓ Let you choose a retention period up to 35 days.

Manual Snapshots:

- ✓ Let you build a new database instance from a snapshot.
- ✓ Are initiated by the user.
- ✓ Persist until the user deletes them.
- ✓ Are stored in Amazon S3.

When automated backups are turned on for your DB Instance, Amazon RDS automatically performs a full daily snapshot of your data (during your preferred backup window) and captures transaction logs (as updates to your DB Instance are made). When you initiate a point-in-time recovery, transaction logs are applied to the most appropriate daily backup in order to restore your DB Instance to the specific time you requested. Amazon RDS retains backups of a DB Instance for a limited, user-specified period of time called the retention period, which by default is one day but can be set to up to thirty five days.

Manual database snapshots are user-initiated and enable you to back up your DB Instance in a known state as frequently as you want, and then restore to that specific state at any time. DB

Snapshots can be created with the AWS Management Console or CreateDBSnapshot API and are kept until you explicitly delete them with the Console or DeleteDBSnapshot API.

Manual database snapshots are kept in Amazon Simple Storage Service (Amazon S3). There is no additional charge for backup storage up to 100% of your consumed database storage for an active DB Instance.

Cross-Region Snapshots

- ✓ Are a copy of a database snapshot stored in a different AWS Region.
- ✓ Provide a backup for disaster recovery.
- ✓ Can be used as a base for migration to a different region.

Cross-region snapshot copy is available for all Amazon RDS engines. You can copy snapshots of any size. Copies can be moved between any of the public AWS Regions, and you can copy the same snapshot to multiple regions simultaneously by initiating more than one transfer. There is no charge for the copy operation itself; you pay only for the data transfer out of the source region and for the data storage in the destination region.

Amazon RDS Security

- ✓ Run your DB instance in an **Amazon VPC**.
- ✓ Use **IAM policies** to grant access to Amazon RDS resources.
- ✓ Use security groups.
- ✓ Use Secure Socket Layer (**SSL**) connections with DB instances (Amazon Aurora, Oracle, MySQL, MariaDB, PostgreSQL, Microsoft SQL Server).
- ✓ Use Amazon RDS **encryption** to secure your RDS instances and snapshots at rest.
- ✓ Use network encryption and transparent data encryption (**TDE**) with Oracle DB and Microsoft SQL Server instances.
- ✓ Use the security features of your DB engine to control access to your DB instance.

You can manage access to your Amazon Relational Database Service (Amazon RDS) resources and your databases on a DB instance. The method you use to manage access depends on what type of task the user needs to perform with Amazon RDS.

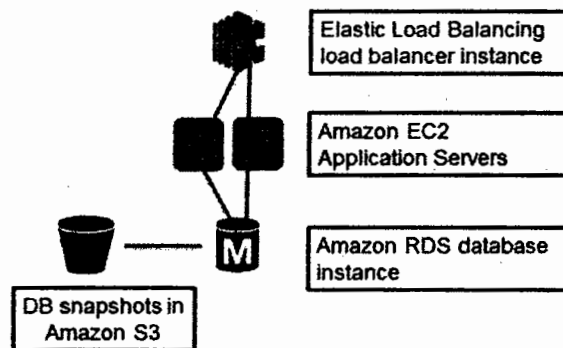
Run your DB instance in an Amazon Virtual Private Cloud (VPC) for the greatest possible network access control.

- ✓ Use AWS Identity and Access Management (IAM) policies to assign permissions that determine who is allowed to manage RDS resources. For example, you can use IAM to

determine who is allowed to create, describe, modify, and delete DB instances, tag resources, or modify DB security groups.

- ✓ Use security groups to control which IP addresses or EC2 instances can connect to your databases on a DB instance. When you first create a DB instance, its firewall prevents any database access except through rules specified by an associated security group.
- ✓ Use Secure Socket Layer (SSL) connections with DB instances running the MySQL, MariaDB, PostgreSQL, or Microsoft SQL Server database engines.
- ✓ Use Amazon RDS encryption to secure your RDS DB instances and snapshots at rest. Amazon RDS encryption uses the industry standard AES-256 encryption algorithm to encrypt your data on the server that hosts your RDS DB instance.
- ✓ Use network encryption and transparent data encryption with Oracle DB instances.
- ✓ Use the security features of your DB engine to control who can log in to the databases on a DB instance, just as you would if the database was on your local network.

A Simple Application Architecture



A simple application stack with an application running in an Amazon EC2 instance supported by a master database running in an Amazon RDS database instance. Presenting the application behind an Elastic Load Balancer allows for compute resiliency and scaling features such as Auto Scaling and ELB groups to be adopted in the future.

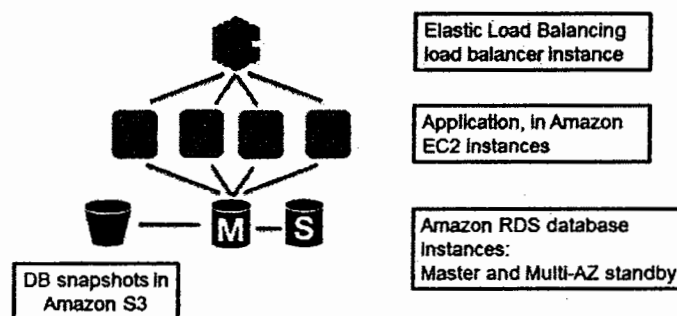
Multi-AZ RDS Deployment

- ✓ With Multi-AZ operation, your database is synchronously replicated to another AZ in the same AWS Region.
- ✓ Failover automatically occurs to the standby in case of master database failure.
- ✓ Planned maintenance is applied first to standby databases.

Amazon RDS Multi-AZ deployments provide enhanced availability and durability for Database (DB) instances, making them a natural fit for production database workloads. When you

provision a Multi-AZ DB instance, Amazon RDS automatically creates a primary DB instance and synchronously replicates the data to a standby instance in a different Availability Zone (AZ). Each AZ runs on its own physically distinct, independent infrastructure, and is engineered to be highly reliable. In case of an infrastructure failure (for example, instance hardware failure, storage failure, or network disruption), Amazon RDS performs an automatic failover to the standby, so that you can resume database operations as soon as the failover is complete. Since the endpoint for your DB instance remains the same after a failover, your application can resume database operation without the need for manual administrative intervention.

A Resilient, Durable Application Architecture



An application stack that uses AWS reliability and durability features. An ELB group of Amazon EC2 instances supports the application logic. The instances use a Multi-AZ Amazon RDS deployment. In the event of infrastructure failure, the database fails over to a standby instance. The application logic retries its database connections, to the same endpoint as before, and the service resumes using the new master. Meanwhile, a new standby is instantiated.

In addition to Amazon RDS's automatic backups, the database snapshot feature is used to ensure that backups are durably retained. You can create a new database instance from a database snapshot whenever you want.

Amazon RDS Best Practices

- ✓ **Monitor** your memory, CPU, and storage usage.
- ✓ Use **Multi-AZ** deployments to automatically provision and maintain a synchronous standby in a different Availability Zone.
- ✓ Enable automatic backups.
- ✓ Set the **backup window** to occur during the daily low in WriteIOPS.
- ✓ To increase the I/O capacity of a DB instance:
 - ✓ Migrate to a DB instance class with high I/O capacity.

- ✓ Convert from standard storage to provisioned IOPS storage and use a DB instance class optimized for provisioned IOPS.
- ✓ Provision additional throughput capacity (if using provisioned IOPS storage).
- ✓ If your client application is caching the DNS data of your DB instances, set a TTL of less than 30 seconds.
- ✓ **Test** failover for your DB instance.
- ✓ Monitor your memory, CPU, and storage usage. Amazon CloudWatch can be set up to notify you when usage patterns change or when you approach the capacity of your deployment, so that you can maintain system performance and availability.
- ✓ Use Multi-AZ deployments to automatically provision and maintain a synchronous standby replica in a different Availability Zone.
- ✓ Enable automatic backups and set the backup window to occur during the daily low in WriteIOPS.
- ✓ On a MySQL DB instance:
 - ✓ Do not create more than 10,000 tables using provisioned IOPS (IOPS are input/output operations per second) or 1000 tables using standard storage. Large numbers of tables will significantly increase database recovery time after a failover or database crash. If you need to create more tables than recommended, set the `innodb_file_per_table` parameter to 0.
 - ✓ Avoid tables in your database growing too large. Underlying file system constraints restrict the maximum size of a MySQL table file to 2 TB. Instead, partition your large tables so that file sizes are well under the 2 TB limit. This approach can also improve performance and recovery time.
- ✓ If your database workload requires more I/O than you have provisioned, recovery after a failover or database failure will be slow. To increase the I/O capacity of a DB instance, do any or all of the following:
 - ✓ Migrate to a DB instance class with high I/O capacity.
 - ✓ Convert from standard storage to provisioned IOPS storage, and use a DB instance class that is optimized for provisioned IOPS.
- ✓ If you are already using provisioned IOPS storage, provision additional throughput capacity.
- ✓ If your client application is caching the DNS data of your DB instances, set a time to live (TTL) of less than 30 seconds. Because the underlying IP address of a DB instance can change after a failover, caching the DNS data for an extended time can lead to connection failures if your application tries to connect to an IP address that no longer is in service.
- ✓ Test failover for your DB instance to understand how long the process takes for your use case and to ensure that the application that accesses your DB instance can automatically connect to the new DB instance after failover.

Amazon DynamoDB

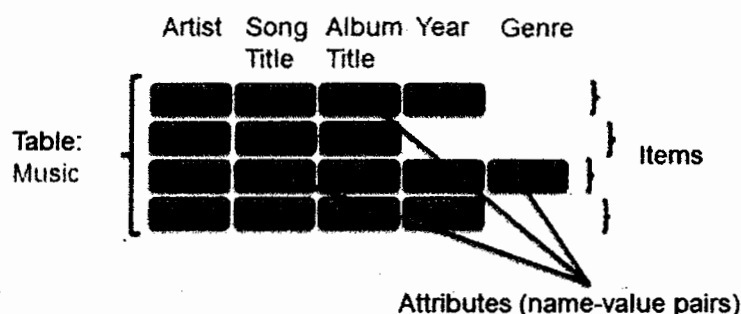
- ✓ Store any amount of data with **no limits**
- ✓ Fast, predictable performance using **SSDs**
- ✓ Easily provision and change the **request capacity** needed for each table
- ✓ **Fully managed, NoSQL** database service

Amazon DynamoDB is a fully-managed NoSQL database service that offers high performance, predictable throughput and low cost. It is easy to set up, operate, and scale. With Amazon DynamoDB, you can start small, specify the throughput and storage you need, and easily scale your capacity requirements in seconds, as needed. It automatically partitions data over a number of servers to meet your requested capacity. In addition, Amazon DynamoDB automatically replicates your data synchronously across multiple Availability Zones within an AWS Region to ensure high availability and data durability.

DynamoDB Use Case

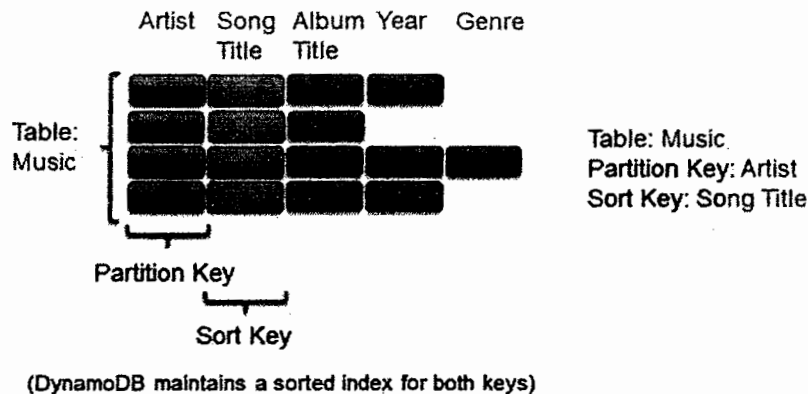
- ✓ Adroll Uses AWS to grow by more than 15,000% in a year
- ✓ Needed high-performance, flexible platform to swiftly sync data for worldwide audience
- ✓ Processes 50 TB of data a day
- ✓ Serves 50 billion impressions a day
- ✓ Stores 1.5 PB of data
- ✓ Worldwide deployment minimizes latency

DynamoDB Data Model



In Amazon DynamoDB, a table is a collection of items and each item is a collection of attributes. Each attribute in an item is a name-value pair. An attribute can be a scalar (single-valued), a JSON document, or a set.

Primary Keys

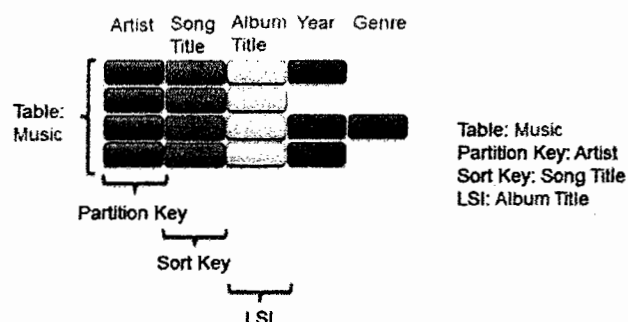


When you create a table, in addition to the table name, you must specify the primary key of the table. As in other databases, a primary key in DynamoDB uniquely identifies each item in the table, so that no two items can have the same key. When you add, update, or delete an item in the table, you must specify the primary key attribute values for that item.

DynamoDB supports two different kinds of primary keys:

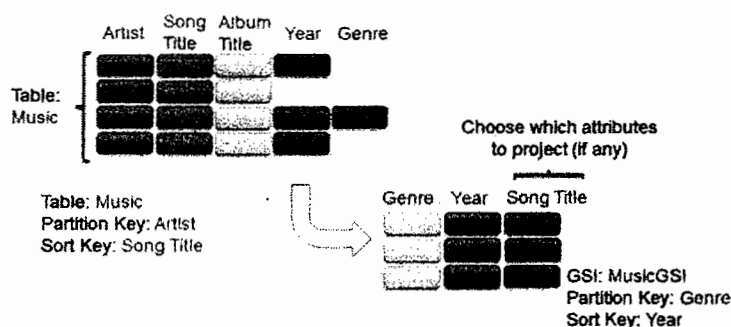
- ✓ **Partition Key** – A simple primary key, composed of one attribute known as the partition key. DynamoDB uses the partition key's value as input to an internal hash function; the output from the hash function determines the partition where the item is stored. No two items in a table can have the same partition key value.
- ✓ **Partition Key and Sort Key** – A composite primary key, composed of two attributes. The first attribute is the partition key, and the second attribute is the sort key. DynamoDB uses the partition key value as input to an internal hash function; the output from the hash function determines the partition where the item is stored. All items with the same partition key are stored together, in sorted order by sort key value. It is possible for two items to have the same partition key value, but those two items must have different sort key values.

Local Secondary Index



If you want to read the data using non-key attributes, you can use a secondary index to do this. A local secondary index is an index that has the same partition key as the table, but a different sort key.

Global Secondary Index



A Global secondary index is an index with a partition key and sort key that can be different from those on the table. They can be thought of as "pivot charts" for your table.

Provisioned Throughput

- ✓ You specify how much provisioned throughput capacity you need for reads and writes.
- ✓ Amazon DynamoDB allocates the necessary machine resources to meet your needs.
- ✓ Read capacity unit:
 - ✓ One strongly consistent read per second for items as large as 4 KB
 - ✓ Two eventually consistent reads per second for items as large as 4 KB.
- ✓ Write capacity unit:
 - ✓ One write per second for items as large as 1 KB.

When you create or update a table, you specify how much provisioned throughput capacity you need for reads and writes. Amazon DynamoDB will allocate the necessary machine resources to meet your throughput needs while ensuring consistent, low-latency performance.

A unit of read capacity represents one strongly consistent read per second (or two eventually consistent reads per second) for items as large as 4 KB. A unit of write capacity represents one write per second for items as large as 1 KB.

Supported Operations

- ✓ **Query:**
 - ✓ Query a table using the partition key and an optional sort key filter.
 - ✓ If the table has a secondary index, query using its key.
 - ✓ It is the most efficient way to retrieve items from a table or secondary index.
- ✓ **Scan:**
 - ✓ You can scan a table or secondary index.
 - ✓ Scan reads every item - slower than querying.
- ✓ You can use conditional expressions in both Query and Scan operations.

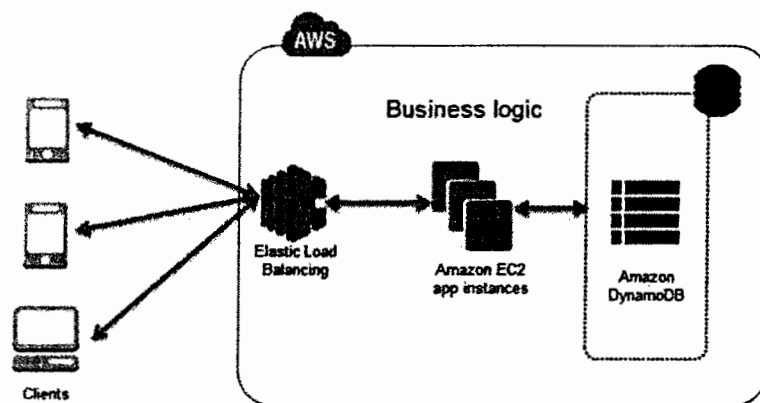
The Query operation enables you to query a table using the partition key and an optional sort key filter. If the table has a secondary index, you can also Query the index using its key. You can query only tables that have a composite primary key (partition key and sort key). You can also query any secondary index on such tables. Query is the most efficient way to retrieve items from a table or a secondary index.

Amazon DynamoDB also supports a Scan operation, which you can use on a table or a secondary index. The Scan operation reads every item in the table or secondary index. For large tables and secondary indexes, a Scan can consume a large amount of resources; for this reason, we recommend that you design your applications so that you can use the Query operation mostly, and use Scan only where appropriate.

You can use conditional expressions in both the Query and Scan operations to control which items are returned.

Simple Application Architecture

A simple application architecture using Amazon DynamoDB to store data processed by applications on Amazon EC2 instances.



Amazon RDS and Amazon DynamoDB

Factors	Relational (Amazon RDS)	NoSQL (Amazon DynamoDB)
Application Type	<ul style="list-style-type: none"> Existing database apps Business process-centric apps 	<ul style="list-style-type: none"> New web-scale applications Large number of small writes and reads
Application Characteristics	<ul style="list-style-type: none"> Relational data models, transactions Complex queries, joins, and updates 	<ul style="list-style-type: none"> Simple data models, transactions Range queries, simple updates
Scaling	Application or DBA-architected (clustering, partitions, sharding)	Seamless, on-demand scaling based on application requirements
QoS	<ul style="list-style-type: none"> Performance—depends on data model, indexing, query, and storage optimization Reliability and availability Durability 	<ul style="list-style-type: none"> Performance—Automatically optimized by the system Reliability and availability Durability

Database Considerations

If You Need	Consider Using
A relational database service with minimal administration	Amazon RDS <ul style="list-style-type: none"> Choice of Amazon Aurora, MySQL, MariaDB, Microsoft SQL Server, Oracle, or PostgreSQL database engines Scale compute and storage Multi-AZ availability
A fast, highly scalable NoSQL database service	Amazon DynamoDB <ul style="list-style-type: none"> Extremely fast performance Seamless scalability and reliability Low cost
A database you can manage on your own	Your choice of AMIs on Amazon EC2 and Amazon EBS that provide scale compute and storage, complete control over instances, and more.

AWS provides a number of database alternatives for developers. You can run fully managed relational and NoSQL services, or you can operate your own database in the cloud on Amazon EC2 and Amazon EBS. If you need a relational database service with minimal administration, consider using Amazon RDS. If you need a fast, highly scalable NoSQL database service, consider using Amazon DynamoDB. If you need a relational database you can manage on your own, consider using your choice of relational AMIs.

Module 5

AWS ELASTICITY AND MANAGEMENT TOOLS

Auto Scaling

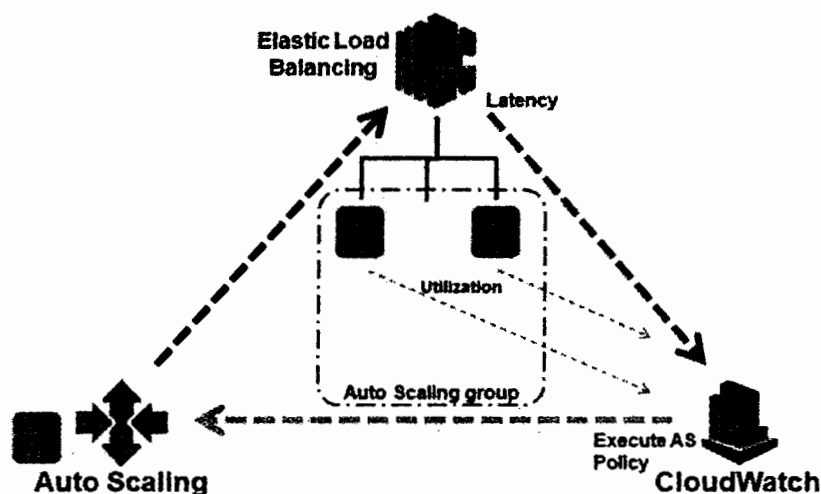
- ✓ Scale your Amazon EC2 capacity automatically
- ✓ Well-suited for applications that experience variability in usage
- ✓ Available at no additional charge

Understand Auto Scaling concepts including:

- ✓ Launch Configurations
- ✓ Auto Scaling Groups
- ✓ Scaling Plans
- ✓ Auto Scaling Lifecycle
- ✓ Auto Scaling Limits

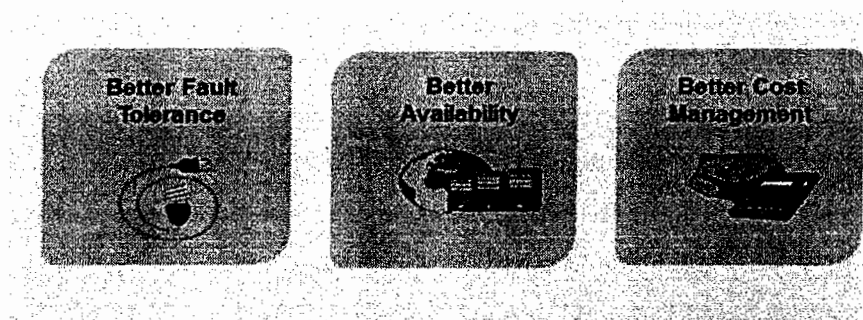
Auto Scaling helps you ensure that you have the correct number of EC2 instances available to handle the load for your application. Auto Scaling is particularly well suited for applications that experience hourly, daily, or weekly variability in usage.

Trio of Services



Auto Scaling works as a triad of services working in sync. Elastic Load Balancing and EC2 instances feed metrics to Amazon CloudWatch. Auto Scaling defines a group with launch configurations and Auto Scaling policies. Amazon CloudWatch alarms execute Auto Scaling policies to affect the size of your fleet. All of these services work well individually, but together they become more powerful and increase the control and flexibility our customers demand.

Auto Scaling Benefits



Adding Auto Scaling to your application architecture is one way to maximize the benefits of the AWS cloud. When you use Auto Scaling, your applications gain the following benefits:

Better fault tolerance: Auto Scaling can detect when an instance is unhealthy, terminate it, and launch an instance to replace it. You can also configure Auto Scaling to use multiple Availability Zones. If one Availability Zone becomes unavailable, Auto Scaling can launch instances in another one to compensate.

Better availability: Auto Scaling can help you ensure that your application always has the right amount of capacity to handle the current traffic demands.

Better cost management: Auto Scaling can dynamically increase and decrease capacity as needed. Because you pay for the EC2 instances you use, you save money by launching instances when they are actually needed and terminating them when they aren't needed.

Launch Configurations

- ✓ A launch configuration is a template that an Auto Scaling group uses to launch EC2 instances.
- ✓ When you create a launch configuration, you can specify:
 - ✓ AMI ID
 - ✓ Instance type
 - ✓ Key pair
 - ✓ Security groups
 - ✓ Block device mapping
 - ✓ User data

When you create an Auto Scaling group, you must specify a launch configuration. You can specify your launch configuration with multiple Auto Scaling groups. However, you can only specify one launch configuration for an Auto Scaling group at a time, and you can't modify a

launch configuration after you've created it. If you want to change the launch configuration for your Auto Scaling group, you must create a new launch configuration and then update your Auto Scaling group with the new launch configuration. When you change the launch configuration for your Auto Scaling group, any new instances are launched using the new configuration parameters, but existing instances are not affected.

Auto Scaling Groups

- ✓ Contain a collection of EC2 instances that share similar characteristics.
- ✓ Instances in an Auto Scaling group are treated as a logical grouping for the purpose of instance scaling and management.

You can create collections of EC2 instances, called Auto Scaling groups. You can specify the minimum number of instances in each Auto Scaling group, and Auto Scaling ensures that your group never goes below this size. You can specify the maximum number of instances in each Auto Scaling group, and Auto Scaling ensures that your group never goes above this size. If you specify the desired capacity, either when you create the group or at any time thereafter, Auto Scaling ensures that your group has this many instances. If you specify scaling policies, then Auto Scaling can launch or terminate instances as demand on your application increases or decreases.

Dynamic Scaling

You can create a scaling policy that uses CloudWatch alarms to determine:

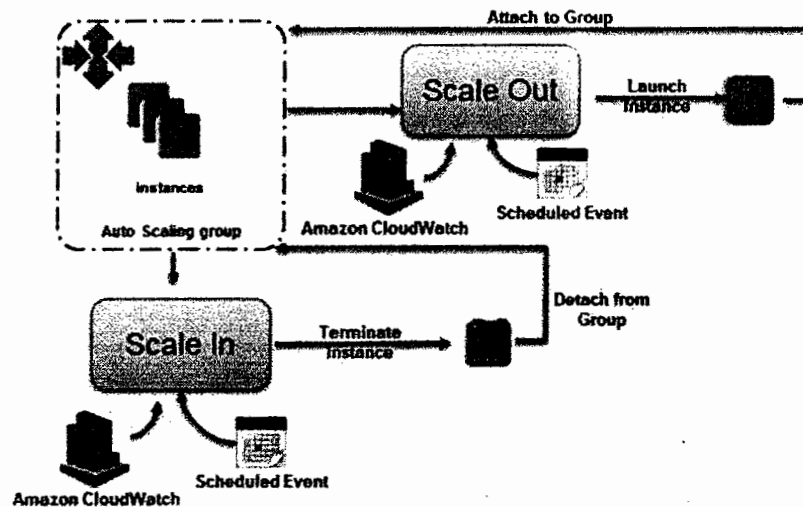
- ✓ When your Auto Scaling group should scale out.
- ✓ When your Auto Scaling group should scale in.

You can use alarms to monitor:

- ✓ Any of the metrics that AWS services send to Amazon CloudWatch.
- ✓ Your own custom metrics.

Each CloudWatch alarm watches a single metric and sends messages to Auto Scaling when the metric breaches a threshold that you specify in your policy.

Auto Scaling Basic Lifecycle



The basic lifecycle of instances within an Auto Scaling Group.

- ✓ The Scaling Group has a desired capacity of three instances.
- ✓ A Cloud Watch alarm trigger scaling events and policies scale the group at specific dates and times.
- ✓ The scaling policy launches an instance and attaches it to the Auto Scaling Group.
- ✓ A health check fails and triggers an alarm similar to scaling out.
- ✓ The instance is terminated.
- ✓ The instance is detached from the Auto Scaling Group.

Elastic Load Balancing

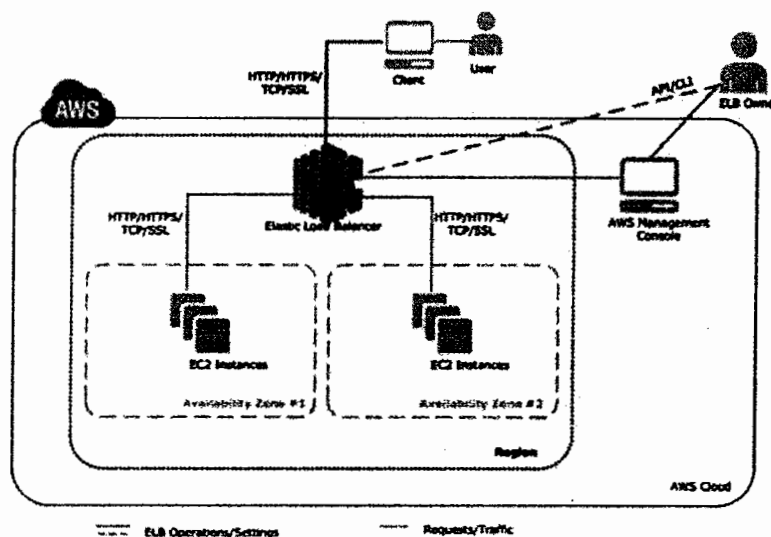
- ✓ Distributes traffic across multiple instances
- ✓ Supports health checks to detect unhealthy Amazon EC2 instances
- ✓ Supports the routing and load balancing of HTTP, HTTPS, and TCP traffic to Amazon EC2 instances

Understand Elastic Load Balancing (ELB) concepts including:

- ✓ Request Routing
- ✓ Internet-facing vs. internal vs. http load balancers
- ✓ Back-end instances
- ✓ Listeners

Elastic Load Balancing automatically distributes incoming application traffic across multiple Amazon EC2 instances. It enables you to achieve greater fault tolerance in your applications, seamlessly providing the amount of load balancing capacity needed in response to incoming application traffic. Elastic Load Balancing detects unhealthy instances within a pool and automatically reroutes traffic to healthy instances until the unhealthy instances have been restored. You can enable Elastic Load Balancing within a single Availability Zone or across multiple zones for even more consistent application performance.

Elastic Load Balancing Example



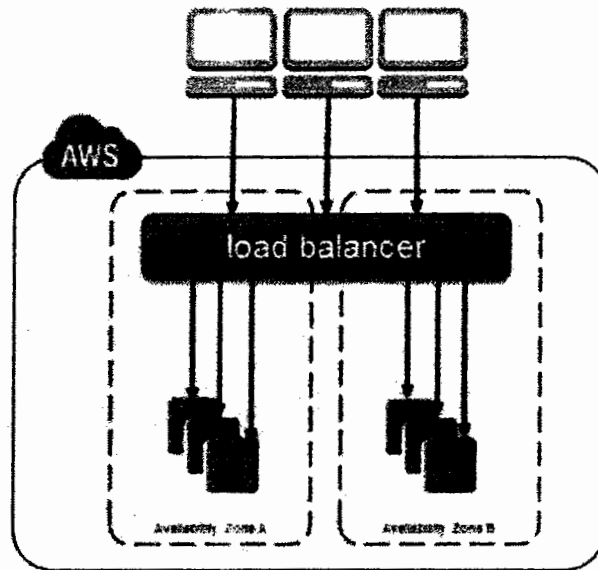
Elastic Load Balancing automatically scales its request handling capacity in response to incoming traffic. The diagram shows how the various components of Elastic Load Balancing work together. You can access and work with your load balancer using one of the following interfaces:

- ✓ AWS Management Console - A simple web browser interface that you can use to create and manage your load balancers without using additional software or tools.
- ✓ Command Line Interfaces - A Java-based command line client that wraps the SOAP API.

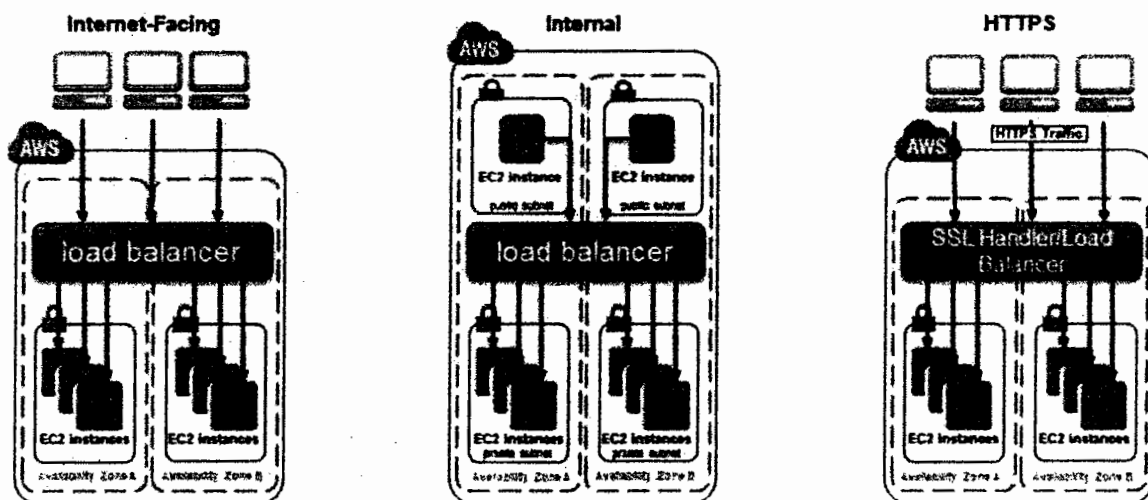
How It Works

A load balancer accepts incoming traffic from clients and routes requests to its registered EC2 instances in one or more Availability Zones. The load balancer also monitors the health of its registered instances and ensures that it routes traffic only to healthy instances. When the load

balancer detects an unhealthy instance, it stops routing traffic to that instance, and then resumes routing traffic to that instance when it detects that the instance is healthy again.



Load Balancer Types



Internet-facing load balancer:

An Internet-facing load balancer takes requests from clients over the Internet and distributes them across the EC2 instances that are registered with the load balancer.

Internal load balancer:

An internal load balancer routes traffic to your EC2 instances in private subnets. The clients must have access to the private subnets.

HTTPS load balancer:

You can create a load balancer that uses the SSL/TLS protocol for encrypted connections (also known as SSL offload). This feature enables traffic encryption between your load balancer and the clients that initiate HTTPS sessions, and for connections between your load balancer and your back-end instances.

Back-end Instances for Load Balancer

- ✓ Health Checks
- ✓ Security Groups
- ✓ Subnets
- ✓ Register
- ✓ De-Register Instances

After you've created your load balancer, you must register your EC2 instances with the load balancer. You can select EC2 instances from a single Availability Zone or multiple Availability Zones within the same region as the load balancer. Elastic Load Balancing routinely performs health checks on registered EC2 instances, and automatically distributes incoming requests to the DNS name of your load balancer across the registered, healthy EC2 instances.

Health checks are periodic pings, attempted connections, or requests sent to EC2 instances by your load balancer to check the availability of your EC2 instances. The load balancer performs health checks on all registered instances, whether the instance is in a healthy state or an unhealthy state. The load balancer routes requests only to the healthy instances. When the load balancer determines that an instance is unhealthy, it stops routing requests to that instance. The load balancer resumes routing requests to the instance when it has been restored to a healthy state.

A security group acts as a firewall that controls the traffic allowed to and from one or more instances. When you launch an EC2 instance, you can associate one or more security groups with the instance. For each security group, you add one or more rules to allow traffic. You can modify the rules for a security group at any time; the new rules are automatically applied to all instances associated with the security group. You must ensure that the security groups for your instances allow the load balancer to communicate with your back-end instances on both the listener port and the health check port. In a VPC, your security groups and network access control lists (ACL) must allow traffic in both directions on these ports.

When you attach a subnet to your load balancer, Elastic Load Balancing creates a load balancer node in the Availability Zone. Load balancer nodes accept traffic from clients and forward requests to the healthy registered instances in one or more Availability Zones. For load

balancers in a VPC, we recommend that you attach one subnet per Availability Zone for at least two Availability Zones. This improves the availability of your load balancer. Note that you can modify the subnets attached to your load balancer at any time.

Registering an EC2 instance adds it to your load balancer. The load balancer continuously monitors the health of its registered instances, and routes requests to the healthy registered instances. If demand on your instances increases, you can register additional instances with the load balancer to handle the demand.

De-registering an EC2 instance removes it from your load balancer. The load balancer stops routing requests to an instance as soon as it is de-registered. If demand decreases, or you need to service your instances, you can de-register instances from the load balancer. A de-registered instance remains running, but no longer receives traffic from the load balancer, and you can register it with the load balancer again when you are ready.

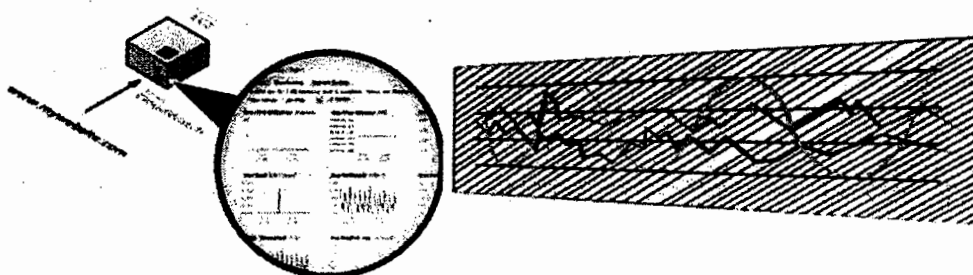
Amazon CloudWatch

- ✓ A monitoring service for AWS cloud resources and the applications you run on AWS
- ✓ Visibility into resource utilization, operational performance, and overall demand patterns
- ✓ Custom application-specific metrics of your own
- ✓ Accessible via AWS Management Console, APIs, SDK, or CLI

CloudWatch lets you view graphs, set alarms to troubleshoot, spot trends, and take automated action based on the state. It is accessible via the AWS Management Console, APIs, SDK or CLI. You can customize with your own metrics or use a sample template found online.

Amazon CloudWatch Facts

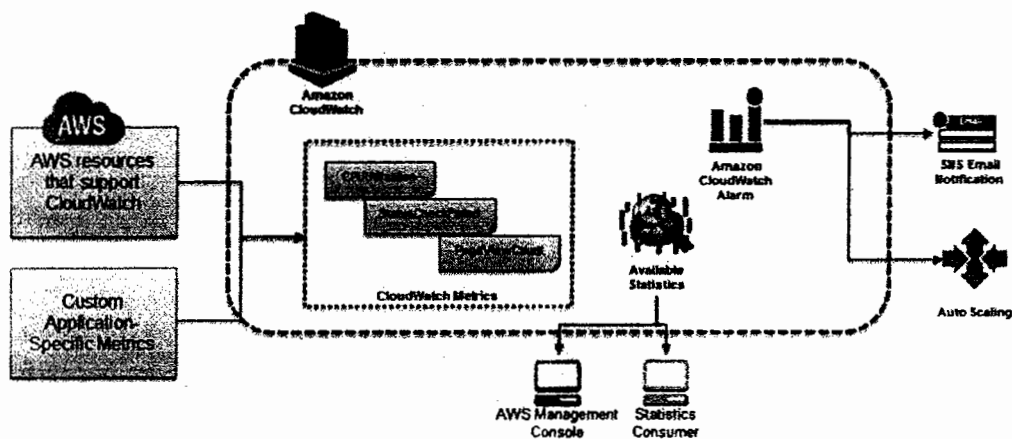
- ✓ Monitor other AWS resources
 - ✓ View graphics and statistics
- ✓ Set Alarms



For Amazon EC2 instances, Amazon CloudWatch basic monitoring collects and reports metrics for CPU utilization, data transfer, and disk usage activity from each Amazon EC2 instance at a five-minute frequency. Amazon CloudWatch detailed monitoring provides these same metrics at one-minute intervals, and also enables data aggregation by Amazon EC2 AMI ID and instance type.

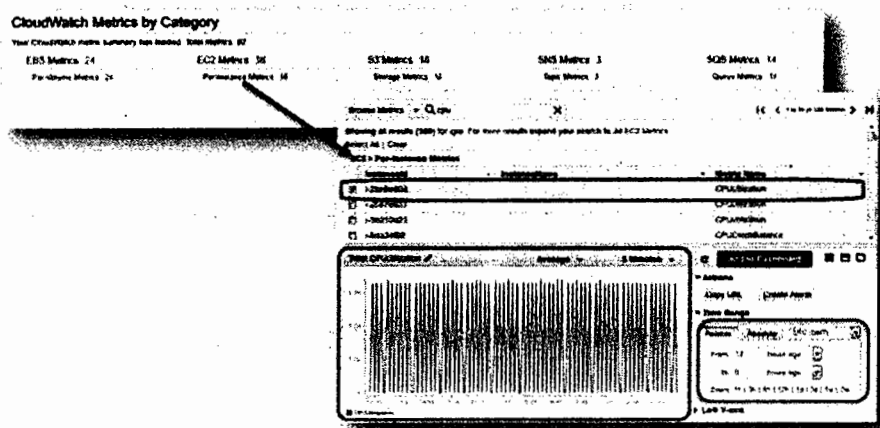
Set alarms on any of your metrics to receive notifications. You can also use Auto Scaling to add or remove Amazon instances.

Amazon CloudWatch Architecture



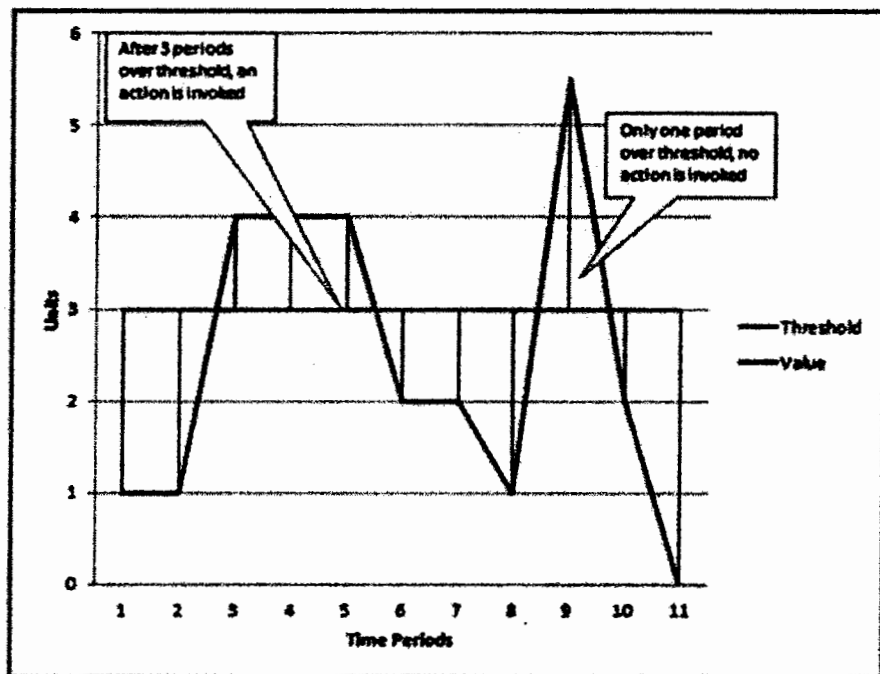
Amazon CloudWatch is a metrics repository. AWS products put metrics into the repository and you retrieve statistics based on the metrics. Statistics can be graphically presented in the CloudWatch console.

CloudWatch Metrics Examples



This screenshots from the Amazon CloudWatch Console. In the example, the user has selected an EC2 per-instance metric of CPU Utilization.

CloudWatch Alarms

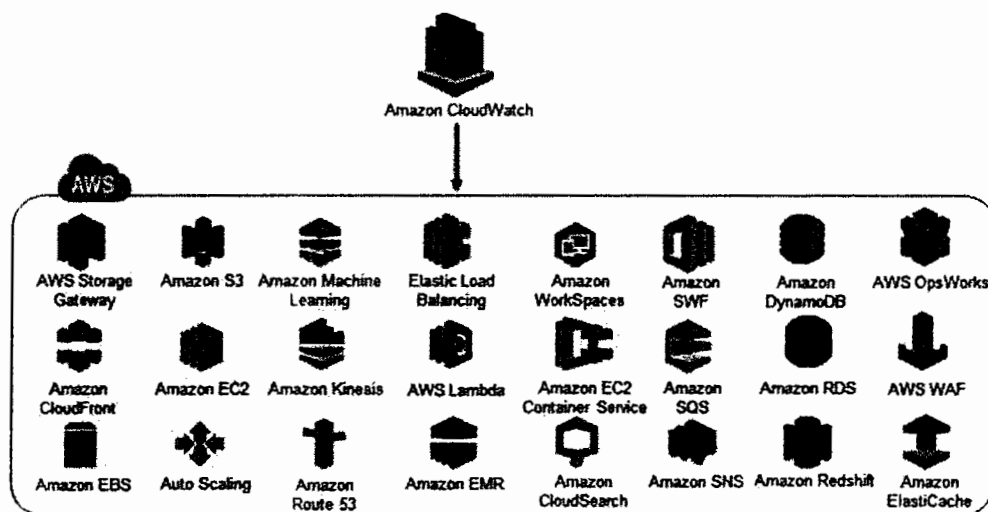


You can create a CloudWatch alarm that sends an Amazon Simple Notification Service (SNS) message when the alarm changes state. An alarm watches a single metric over a time period you specify, and performs one or more actions based on the value of the metric relative to a given threshold over a number of time periods. The action is a notification sent to an Amazon SNS topic or Auto Scaling policy. Alarms invoke actions for sustained state changes only.

CloudWatch alarms will not invoke actions simply because they are in a particular state, the state must have changed and been maintained for a specified number of periods.

In the slide, the alarm threshold is set to 3 and the minimum breach is 3 periods. The alarm invokes its action only when the threshold is breached for 3 consecutive periods. In the figure, this happens with the third through fifth time periods, and the alarm is triggered. At period six, the value dips below the threshold, and the state is set to OK. Later, during the ninth time period, the threshold is breached again, but not for the necessary three consecutive periods. Consequently, the alarm's state remains OK.

Supported AWS Services



AWS Trusted Advisor

- ✓ **Best practice** and recommendation engine.
- ✓ Provides AWS customers with performance and security recommendations in four categories: **cost optimization, security, fault tolerance, and performance improvement.**

The status of checks provided by AWS Trusted Advisor is shown by using color coding on the dashboard page:

Red: action recommended

Yellow: investigation recommended

Green: no problem detected

For each check, you can review a detailed description of the recommended best practice, a set of alert criteria, guidelines for action, and a list of useful resources on the topic.

Cost Optimization

- ✓ Amazon EC2 Reserved Instance Optimization
- ✓ Low Utilization Amazon EC2 Instances
- ✓ Idle Load Balancers
- ✓ Underutilized Amazon EBS Volumes
- ✓ Unassociated Elastic IP Addresses
- ✓ Amazon RDS Idle DB Instances

AWS Trusted Advisor helps you save money on AWS by checking for unused and idle resources and making commitments to reserved capacity.

The following cost optimization checks are available with Trusted Advisor:

- ✓ Amazon EC2 Reserved Instance Optimization: Checks your Amazon Elastic Compute Cloud (Amazon EC2) computing consumption history and calculates an optimal number of Partial Upfront Reserved Instances. Recommendations are based on the previous calendar month's hour-by-hour usage aggregated across all consolidated billing accounts.
- ✓ Low Utilization Amazon EC2 Instances: Checks the Amazon EC2 instances that were running at any time during the last 14 days and alerts you if the daily CPU utilization was 10% or less and network I/O was 5 MB or less on 4 or more days.
- ✓ Idle Load Balancers: Checks your Elastic Load Balancing configuration for load balancers that are not actively used.
- ✓ Underutilized Amazon EBS Volumes: Checks Amazon Elastic Block Store (Amazon EBS) volume configurations and warns when volumes appear to be underused.
- ✓ Unassociated Elastic IP Addresses: Checks for Elastic IP addresses (EIPs) that are not associated with a running Amazon EC2 instance.
- ✓ Amazon RDS Idle DB Instances: Checks the configuration of your Amazon Relational Database Service (Amazon RDS) for any DB instances that appear to be idle. If a DB instance has not had a connection for a prolonged period of time, you can shut down the instance to reduce costs. If persistent storage is needed for data on the instance, you can use lower-cost options such as taking and retaining a DB snapshot.

Security

- ✓ Security Groups
- ✓ AWS IAM Use
- ✓ Amazon S3Bucket Permissions
- ✓ MFA on Root Account
- ✓ AWS IAM Password Policy
- ✓ Amazon RDS Security Group Access Risk

AWS Trusted Advisor helps you improve the security of your application by closing gaps, enabling various AWS security features, and examining your permissions.

The following security checks are available with Trusted Advisor:

- ✓ Security Groups - Specific Ports Unrestricted: Checks security groups for rules that allow unrestricted access (0.0.0.0/0) to specific ports. Unrestricted access increases opportunities for malicious activity (hacking, denial-of-service attacks, loss of data). The

ports with highest risk are flagged red, and those with less risk are flagged yellow. Ports flagged green are typically used by applications that require unrestricted access, such as HTTP and SMTP.

- ✓ Security Groups - Unrestricted Access: Checks security groups for rules that allow unrestricted access to a resource. Unrestricted access increases opportunities for malicious activity (hacking, denial-of-service attacks, loss of data).
- ✓ IAM Use (Free!): Checks for your use of AWS Identity and Access Management (IAM).
- ✓ Amazon S3 Bucket Permissions : Checks buckets in Amazon Simple Storage Service (Amazon S3) that have open access permissions. This check examines explicit bucket permissions, but it does not examine associated bucket policies that might override the bucket permissions.
- ✓ MFA on Root Account: Checks the root account and warns if multi-factor authentication (MFA) is not enabled.
- ✓ IAM Password Policy: Checks the password policy for your account and warns when a password policy is not enabled, or if password content requirements have not been enabled.
- ✓ Amazon RDS Security Group Access Risk: Checks security group configurations for Amazon Relational Database Service (Amazon RDS) and warns when a security group rule might grant overly permissive access to your database.

Fault Tolerance

- ✓ Amazon EBS Snapshots
- ✓ Load Balancer Optimization
- ✓ Auto Scaling Group Resources
- ✓ Amazon RDS Multi-AZ
- ✓ Amazon Route 53 Name Server Delegations
- ✓ ELB Connection Draining

AWS Trusted Advisor helps you increase the availability and redundancy of your AWS application by taking advantage of auto scaling, health checks, multi AZ, and backup capabilities.

The following fault tolerance checks are available with Trusted Advisor:

- ✓ Amazon EBS Snapshots: Checks the age of the snapshots for your Amazon Elastic Block Store (Amazon EBS) volumes (available or in-use).
- ✓ Load Balancer Optimization: Checks your load balancer configuration.
- ✓ Auto Scaling Group Resources: Checks the availability of resources associated with launch configurations and your Auto Scaling groups.
- ✓ Amazon RDS Multi-AZ: Checks for DB instances that are deployed in a single Availability Zone.

- ✓ Amazon Route 53 Name Server Delegations: Checks for Amazon Route 53 hosted zones for which your domain registrar or DNS is not using the correct Route 53 name servers.
- ✓ ELB Connection Draining: Checks for load balancers that do not have connection draining enabled.

Performance Improvement

- ✓ High Utilization Amazon EC2 Instances
- ✓ Service Limits
- ✓ Large Number of Rules in EC2 Security Group
- ✓ Over Utilized Amazon EBS Magnetic Volumes
- ✓ Amazon EC2 to EBS Throughput Optimization
- ✓ Amazon CloudFront Alternate Domain Names

AWS Trusted Advisor helps you improve the performance of your service by checking your service limits, ensuring you take advantage of provisioned throughput, and monitoring for over utilized instances.

The following performance improvement checks are available with Trusted Advisor:

- ✓ High Utilization Amazon EC2 Instances: Checks the Amazon Elastic Compute Cloud (Amazon EC2) instances that were running at any time during the last 14 days and alerts you if the daily CPU utilization was more than 90% on 4 or more days.
- ✓ Service Limits: Checks for usage that is more than 80% of the service limit.
- ✓ Large Number of Rules in EC2 Security Group: Checks each Amazon EC2 security group for an excessive number of rules.
- ✓ Over Utilized Amazon EBS Magnetic Volumes: Checks for Amazon Elastic Block Store (EBS) Magnetic volumes that are potentially over-utilized and might benefit from a more efficient configuration.
- ✓ Amazon EC2 to EBS Throughput Optimization: Checks for Amazon EBS volumes whose performance might be affected by the maximum throughput capability of the Amazon EC2 instance they are attached to.
- ✓ CloudFront Alternate Domain Names: Checks Amazon CloudFront distributions for alternate domain names with incorrectly configured DNS settings.