# Quintrix
A MINDLANCE COMPANY

# Intro to AWS

## MODULE 2 – AWS ESSENTIAL TRAINING

**1095 Morris Ave, SUITE 101, Union NJ 07083**
**646-762-0305 | INFO@QUINTRIXSOLUTIONS.COM**

# Cloud Concepts

**Client-server model**

To review, a **client** can be a web browser or desktop application that a person interacts with to make requests to computer servers. A **server** can be services such as Amazon Elastic Compute Cloud (Amazon EC2), a type of virtual server.

Cloud-Based Deployment

- Run all parts of the application in the cloud.
- Migrate existing applications to the cloud.
- Design and build new applications in the cloud.

In a **cloud-based deployment** model, you can migrate existing applications to the cloud, or you can design and build new applications in the cloud. You can build those applications on low-level infrastructure that requires your IT staff to manage them. Alternatively, you can build them using higher-level services that reduce the management, architecting, and scaling requirements of the core infrastructure.

**Quintrix**
A MINDLANCE COMPANY

# Amazon Web Services (AWS)

Amazon Web Services (AWS) is the world's most comprehensive and broadly adopted cloud platform, offering over 200 fully featured services from data centers globally.

**AWS Developer Tools**
Easily develop applications on AWS in the programming language of your choice with familiar tools.

| Developer Tool | Description |
|---|---|
| Web Console | Simple web interface for Amazon Web Services |
| Command Line Tool | Control your AWS services from the command line and automate service management with scripts |
| Integrated Development Environment (IDE) | Write, run, debug, and deploy applications on AWS using familiar Integrated Development Environments (IDE) |
| Software Development Kit (SDK) | Simplify coding with language-specific abstracted APIs for AWS services |
| Infrastructure as Code | Define cloud infrastructure using familiar programming languages |

**Quintrix**
A MINDLANCE COMPANY

# Amazon Elastic Compute Cloud (Amazon EC2)

Amazon Elastic Compute Cloud (Amazon EC2) provides secure, resizable compute capacity in the cloud as Amazon EC2 instances.

If you are responsible for the architecture of your company's resources and need to support new websites. With traditional on-premises resources, you have to do the following:

- Spend money upfront to purchase hardware.
- Wait for the servers to be delivered to you.
- Install the servers in your physical data center.
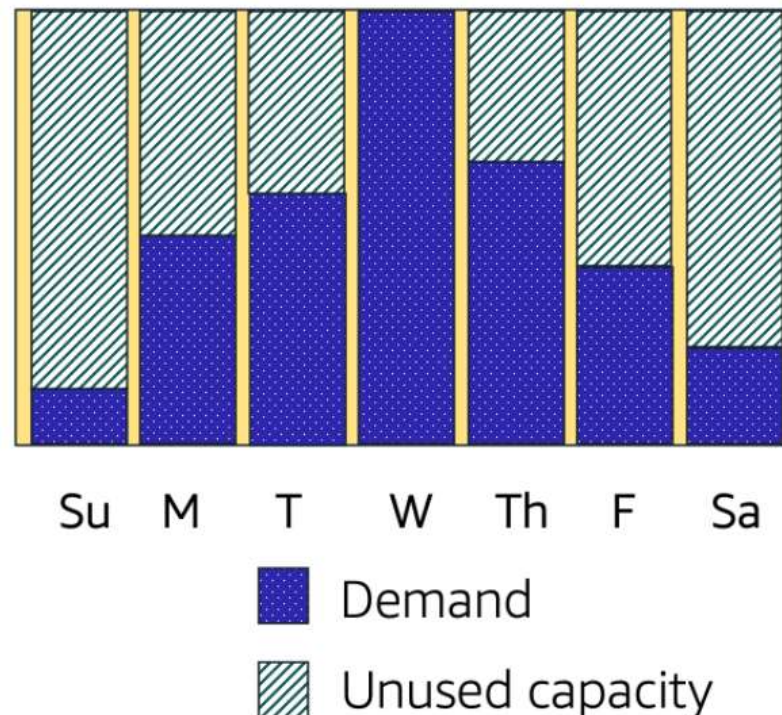- Make all the necessary configurations.

By comparison, with an Amazon EC2 instance you can use a virtual server to run applications in the AWS Cloud.

- You can provision and launch an Amazon EC2 instance within minutes.
- You can stop using it when you have finished running a workload.
- You pay only for the compute time you use when an instance is running, not when it is stopped or terminated.
- You can save costs by paying only for server capacity that you need or want.

**Quintrix**
A MINDLANCE COMPANY

# EC2 Scalability

**Scalability** involves beginning with only the resources you need and designing your architecture to automatically respond to changing demand by scaling out or in. As a result, you pay for only the resources you use. You don't have to worry about a lack of computing capacity to meet your customers' needs.

If you wanted the scaling process to happen automatically, which AWS service would you use? The AWS service that provides this functionality for Amazon EC2 instances is **Amazon EC2 Auto Scaling**.



Su   M   T   W   Th   F   Sa

■ Demand
▨ Unused capacity

**Quintrix**
A MINDLANCE COMPANY

# EC2 Scalability

**Amazon EC2 Auto Scaling**

If you've tried to access a website that wouldn't load and frequently timed out, the website might have received more requests than it was able to handle. This situation is similar to waiting in a long line at a coffee shop, when there is only one barista present to take orders from customers.

Amazon EC2 Auto Scaling enables you to automatically add or remove Amazon EC2 instances in response to changing application demand. By automatically scaling your instances in and out as needed, you are able to maintain a greater sense of application availability.

Within Amazon EC2 Auto Scaling, you can use two approaches: dynamic scaling and predictive scaling.

- *Dynamic scaling* responds to changing demand.
- *Predictive scaling* automatically schedules the right number of Amazon EC2 instances based on predicted demand.

**Quintrix**
A MINDLANCE COMPANY

# Amazon Simple Notification Service (Amazon SNS)

**Amazon Simple Notification Service (Amazon SNS)** is a publish/subscribe service. Using Amazon SNS topics, a publisher publishes messages to subscribers. This is similar to the coffee shop; the cashier provides coffee orders to the barista who makes the drinks.

In Amazon SNS, subscribers can be web servers, email addresses, AWS Lambda functions, or several other options.

Quintrix
A MINDLANCE COMPANY

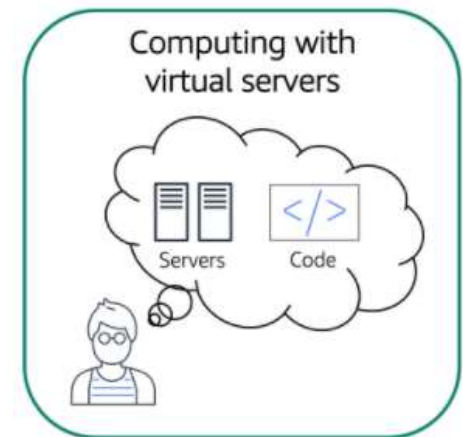# Amazon Simple Queue Service (Amazon SQS)

**Amazon Simple Queue Service (Amazon SQS)** is a message queuing service.

Using Amazon SQS, you can send, store, and receive messages between software components, without losing messages or requiring other services to be available. In Amazon SQS, an application sends messages into a queue. A user or service retrieves a message from the queue, processes it, and then deletes it from the queue.
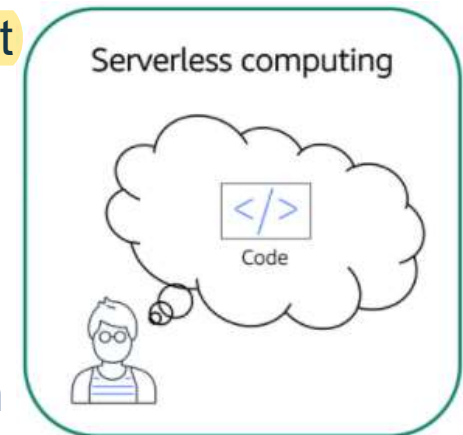
Quintrix
A MINDLANCE COMPANY

# Serverless computing

Amazon EC2, a service that lets you run virtual servers in the cloud. If you have applications that you want to run in Amazon EC2, you must do the following:

1. Provision instances (virtual servers).
2. Upload your code.
3. Continue to manage the instances while your application is running.



Computing with virtual servers

The term "serverless" means that your code runs on servers, but you do not need to provision or manage these servers. With serverless computing, you can focus more on innovating new products and features instead of maintaining servers.

Another benefit of serverless computing is the flexibility to scale serverless applications automatically. Serverless computing can adjust the applications' capacity by modifying the units of consumptions, such as throughput and memory.
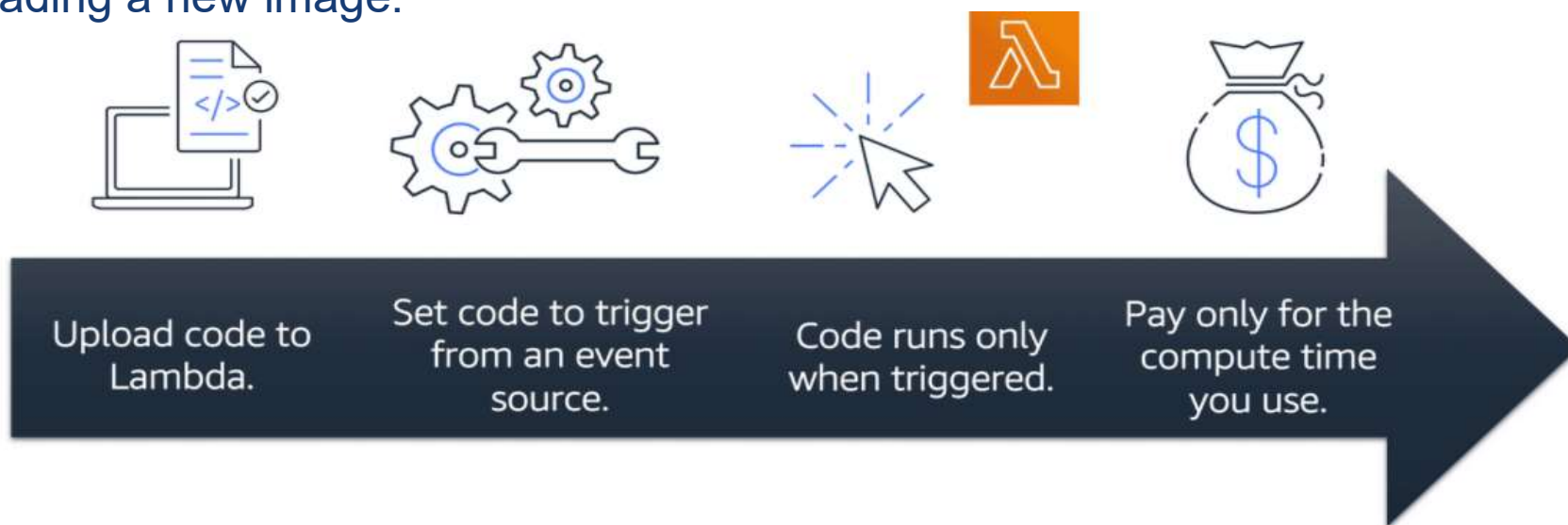


Serverless computing

An AWS service for serverless computing is **AWS Lambda**.

# AWS Lambda

**AWS Lambda** is a service that lets you run code without needing to provision or manage servers.

While using AWS Lambda, you pay only for the compute time that you consume. Charges apply only when your code is running. You can also run code for virtually any type of application or backend service, all with zero administration.

For example, a simple Lambda function might involve automatically resizing uploaded images to the AWS Cloud. In this case, the function triggers when uploading a new image.
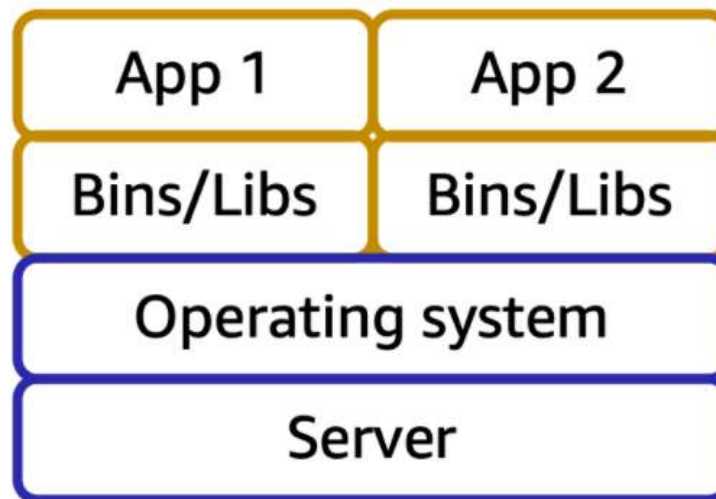


Upload code to Lambda.

Set code to trigger from an event source.

Code runs only when triggered.

Pay only for the compute time you use.

**Quintrix**
A MINDLANCE COMPANY

# Containers

**Containers** provide you with a standard way to package your application's code and dependencies into a single object. You can also use containers for processes and workflows in which there are essential requirements for security, reliability, and scalability.

Suppose that a company's application developer has an environment on their computer that is different from the environment on the computers used by the IT operations staff. The developer wants to ensure that the application's environment remains consistent regardless of deployment, so they use a containerized approach. This helps to reduce time spent debugging applications and diagnosing differences in computing environments.

One host with multiple containers

| App 1 | App 2 |
|---|---|
| Bins/Libs | Bins/Libs |
| Operating system ||
| Server ||

**Quintrix**
A MINDLANCE COMPANY

# Amazon Elastic Container Service (Amazon ECS)

**Amazon Elastic Container Service (Amazon ECS)** is a highly scalable, high-performance container management system that enables you to run and scale containerized applications on AWS.

Amazon ECS supports Docker containers. Docker is a software platform that enables you to build, test, and deploy applications quickly. AWS supports the use of open-source Docker Community Edition and subscription-based Docker Enterprise Edition. With Amazon ECS, you can use API calls to launch and stop Docker-enabled applications.

# Amazon Elastic Kubernetes Service (Amazon EKS)

**Amazon Elastic Kubernetes Service (Amazon EKS)** is a fully managed service that you can use to run Kubernetes on AWS.

Kubernetes is open-source software that enables you to deploy and manage containerized applications at scale. A large community of volunteers maintains Kubernetes, and AWS actively works together with the Kubernetes community. As new features and functionalities release for Kubernetes applications, you can easily apply these updates to your applications managed by Amazon EKS.

**Quintrix**
A MINDLANCE COMPANY

# AWS Fargate

**AWS Fargate** is a serverless compute engine for containers. It works with both Amazon ECS and Amazon EKS.

When using AWS Fargate, you do not need to provision or manage servers. AWS Fargate manages your server infrastructure for you. You can focus more on innovating and developing your applications, and you pay only for the resources that are required to run your containers.

Quintrix
A MINDLANCE COMPANY

# Creating an AWS account

**AWS account:**

**https://aws.amazon.com/resources/create-account/**

Click on **Create an AWS Account** in the upper right-hand corner

This gives you Administrator access
One account per student

**Quintrix**
A MINDLANCE COMPANY

# Selecting a Region

**Compliance with data governance and legal requirements**

Depending on your company and location, you might need to run your data out of specific areas. For example, if your company requires all of its data to reside within the boundaries of the UK, you would choose the London Region.

Not all companies have location-specific data regulations, so you might need to focus more on the other three factors.

**Proximity to your customers**

Selecting a Region that is close to your customers will help you to get content to them faster. For example, your company is based in Washington, DC, and many of your customers live in Singapore. You might consider running your infrastructure in the Northern Virginia Region to be close to company headquarters, and run your applications from the Singapore Region.

**Quintrix**
A MINDLANCE COMPANY

# Selecting a Region

**Available services within a Region**

Sometimes, the closest Region might not have all the features that you want to offer to customers. AWS is frequently innovating by creating new services and expanding on features within existing services. However, making new services available around the world sometimes requires AWS to build out physical hardware one Region at a time.
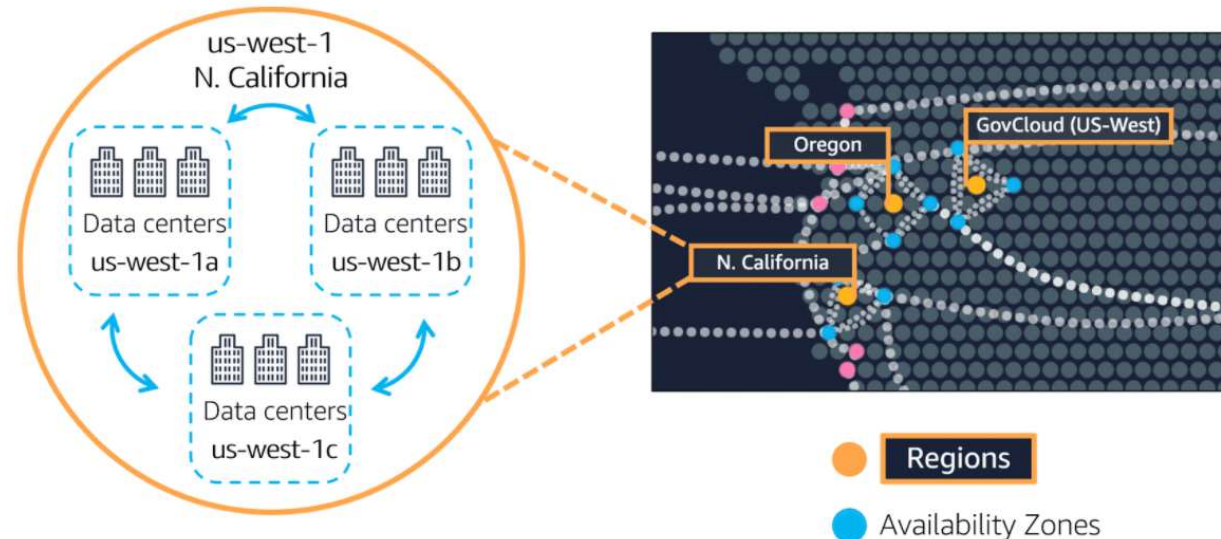
Suppose that your developers want to build an application that uses Amazon Braket (AWS quantum computing platform). At this time Amazon Braket is not yet available in every AWS Region around the world, so your developers would have to run it in one of the Regions that already offers it.

**Pricing**

Suppose that you are considering running applications in both the United States and Brazil. The way Brazil's tax structure is set up, it might cost 50% more to run the same workload out of the São Paulo Region compared to the Oregon Region. You will learn in more detail that several factors determine pricing, but for now know that the cost of services can vary from Region to Region.
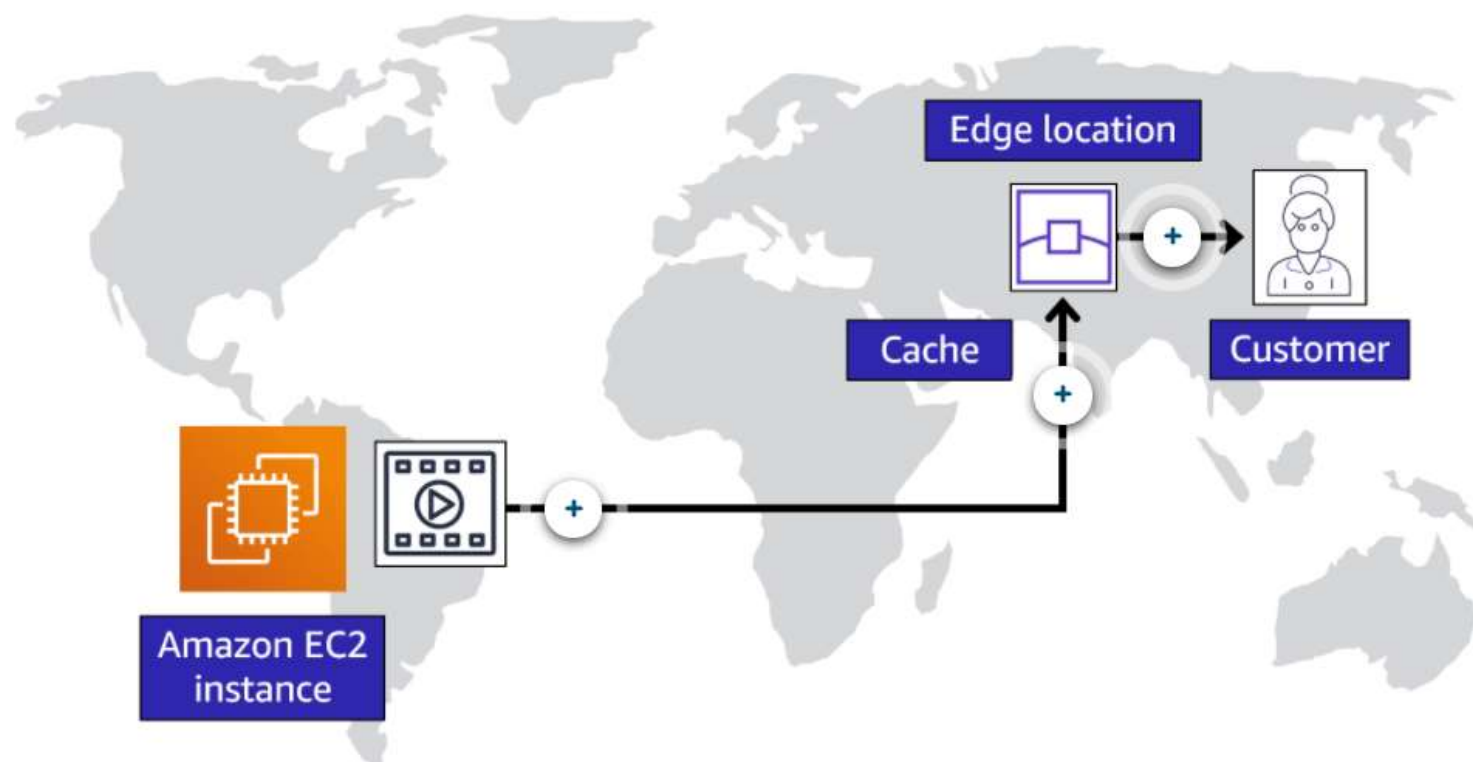
17

**Quintrix**
A MINDLANCE COMPANY

# Availability Zones

An **Availability Zone** is a single data center or a group of data centers within a Region. Availability Zones are located tens of miles apart from each other. This is close enough to have low latency (the time between when content requested and received) between Availability Zones. However, if a disaster occurs in one part of the Region, they are distant enough to reduce the chance that multiple Availability Zones are affected.

Quintrix
A MINDLANCE COMPANY

# Edge Locations

An **edge location** is a site that Amazon CloudFront uses to store cached copies of your content closer to your customers for faster delivery.
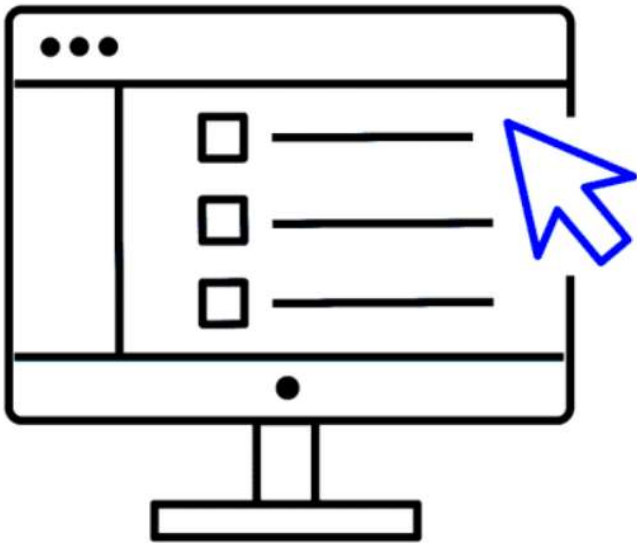
# Ways to interact with AWS

**AWS Management Console**
Web-based interface for accessing and managing AWS services. You can quickly access recently used services and search for other services by name, keyword, or acronym. The console includes wizards and automated workflows that can simplify the process of completing tasks.

You can also use the AWS Console mobile application to perform tasks such as monitoring resources, viewing alarms, and accessing billing information. Multiple identities can stay logged into the AWS Console mobile app at the same time.

**Quintrix**
A MINDLANCE COMPANY

# Ways to interact with AWS

**AWS Command Line Interface**
To save time when making API requests, you can use the **AWS Command Line Interface (AWS CLI)**. AWS CLI enables you to control multiple AWS services directly from the command line within one tool. AWS CLI is available for users on Windows, macOS, and Linux.

By using AWS CLI, you can automate the actions that your services and applications perform through scripts. For example, you can use commands to launch an Amazon EC2 instance, connect an Amazon EC2 instance to a specific Auto Scaling group, and more.
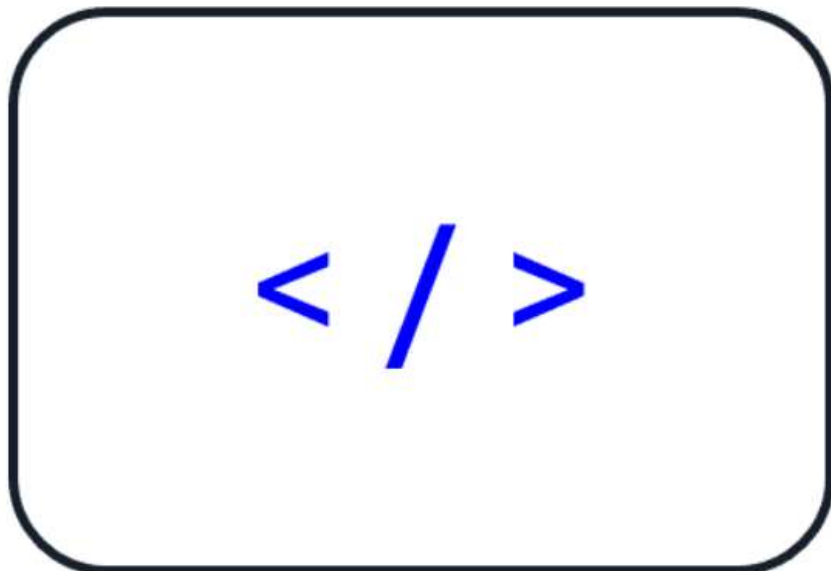
aws> _

Quintrix
A MINDLANCE COMPANY

# Ways to interact with AWS

**Software Development Kits**
Another option for accessing and managing AWS services is the **software development kits (SDKs)**. SDKs make it easier for you to use AWS services through an API designed for your programming language or platform. SDKs enable you to use AWS services with your existing applications or create entirely new applications that will run on AWS.

To help you get started with using SDKs, AWS provides documentation and sample code for each supported programming language. Supported programming languages include C++, Java, .NET, and more.

< / >

**Quintrix**
A MINDLANCE COMPANY

# AWS Elastic Beanstalk

With **AWS Elastic Beanstalk**, you provide code and configuration settings, and Elastic Beanstalk deploys the resources necessary to perform the following tasks:

- Adjust capacity
- Load balancing
- Automatic scaling
- Application health monitoring

# AWS CloudFormation

With **AWS CloudFormation**, you can treat your infrastructure as code. This means that you can build an environment by writing lines of code instead of using the AWS Management Console to individually provision resources.

AWS CloudFormation provisions your resources in a safe, repeatable manner, enabling you to frequently build your infrastructure and applications without having to perform manual actions or write custom scripts. It determines the right operations to perform when managing your stack and rolls back changes automatically if it detects errors.

**Quintrix**
A MINDLANCE COMPANY

# AWS Networking

**Amazon Virtual Private Cloud (Amazon VPC)**
Imagine the millions of customers who use AWS services. Also, imagine the millions of resources that these customers have created, such as Amazon EC2 instances. Without boundaries around all of these resources, network traffic would be able to flow between them unrestricted.

A networking service that you can use to establish boundaries around your AWS resources is **Amazon Virtual Private Cloud (Amazon VPC)**.
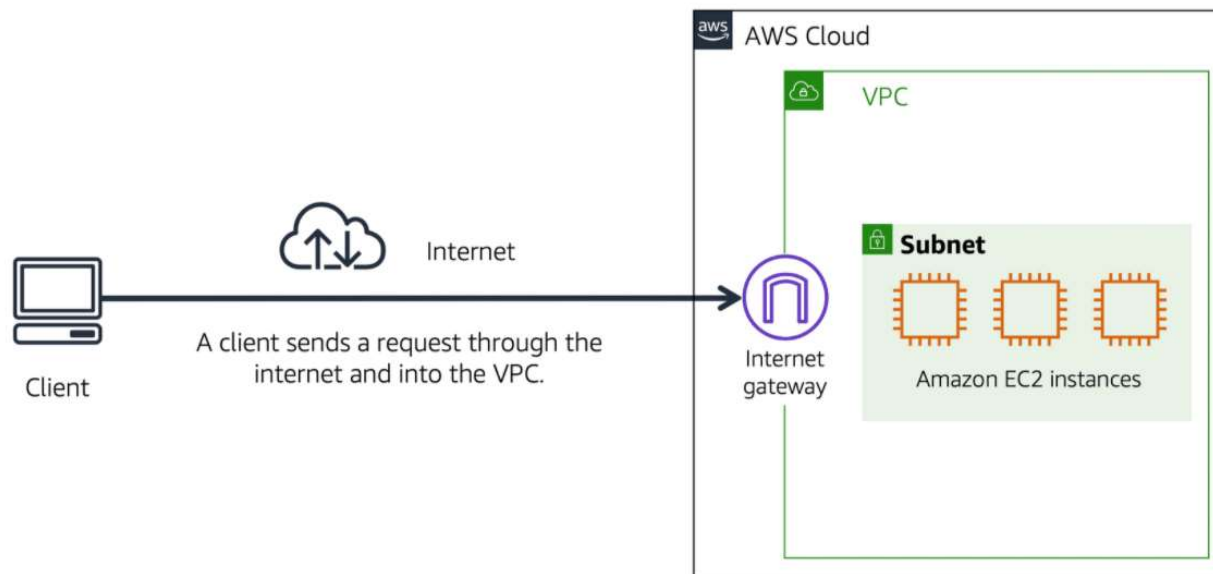Amazon VPC enables you to provision an isolated section of the AWS Cloud. In this isolated section, you can launch resources in a virtual network that you define. Within a virtual private cloud (VPC), you can organize your resources into subnets.

A **subnet** is a section of a VPC that can contain resources such as Amazon EC2 instances.

**Quintrix**
A MINDLANCE COMPANY

# AWS Networking

**Internet gateway**

To allow public traffic from the internet to access your VPC, you attach an **internet gateway** to the VPC.
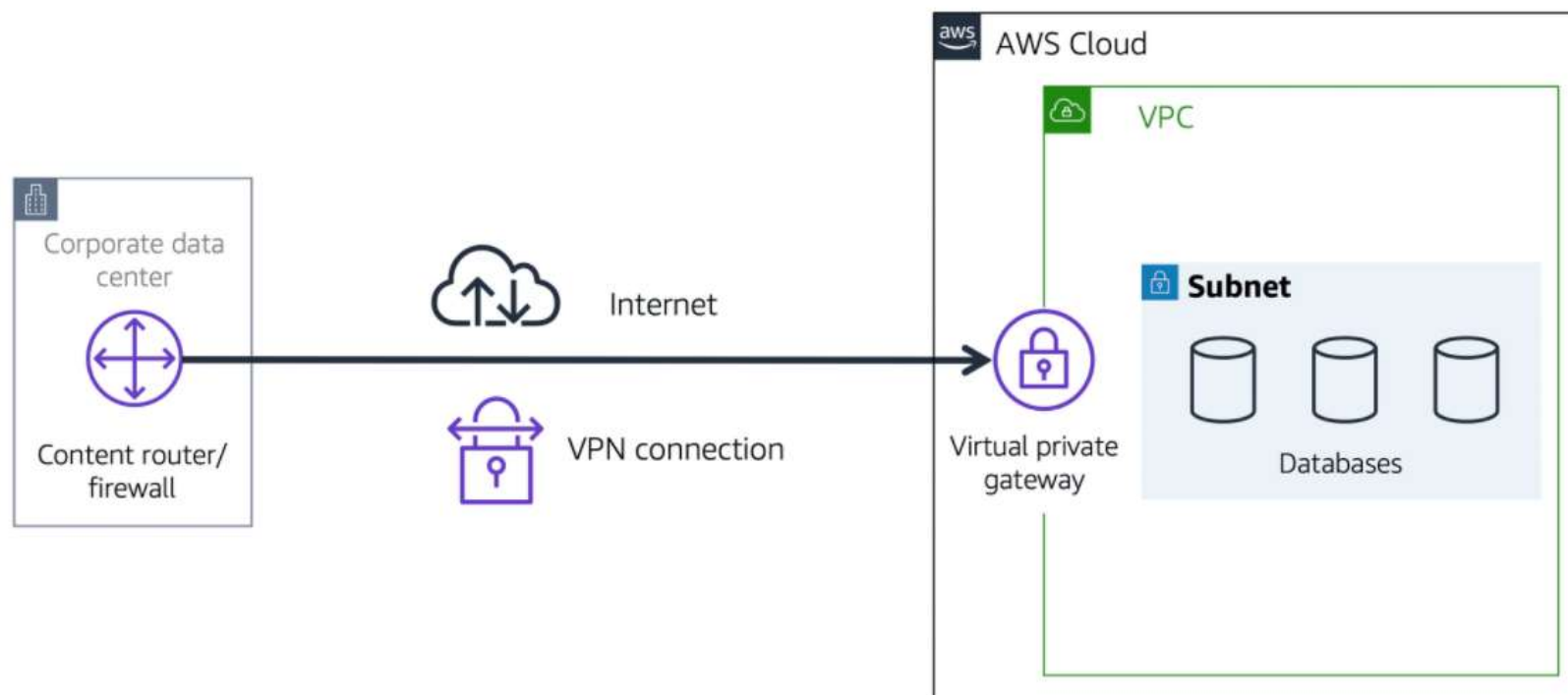


An internet gateway is a connection between a VPC and the internet. You can think of an internet gateway as being similar to a doorway that customers use to enter the coffee shop. Without an internet gateway, no one can access the resources within your VPC.

Quintrix
A MINDLANCE COMPANY

# Virtual Private Gateway

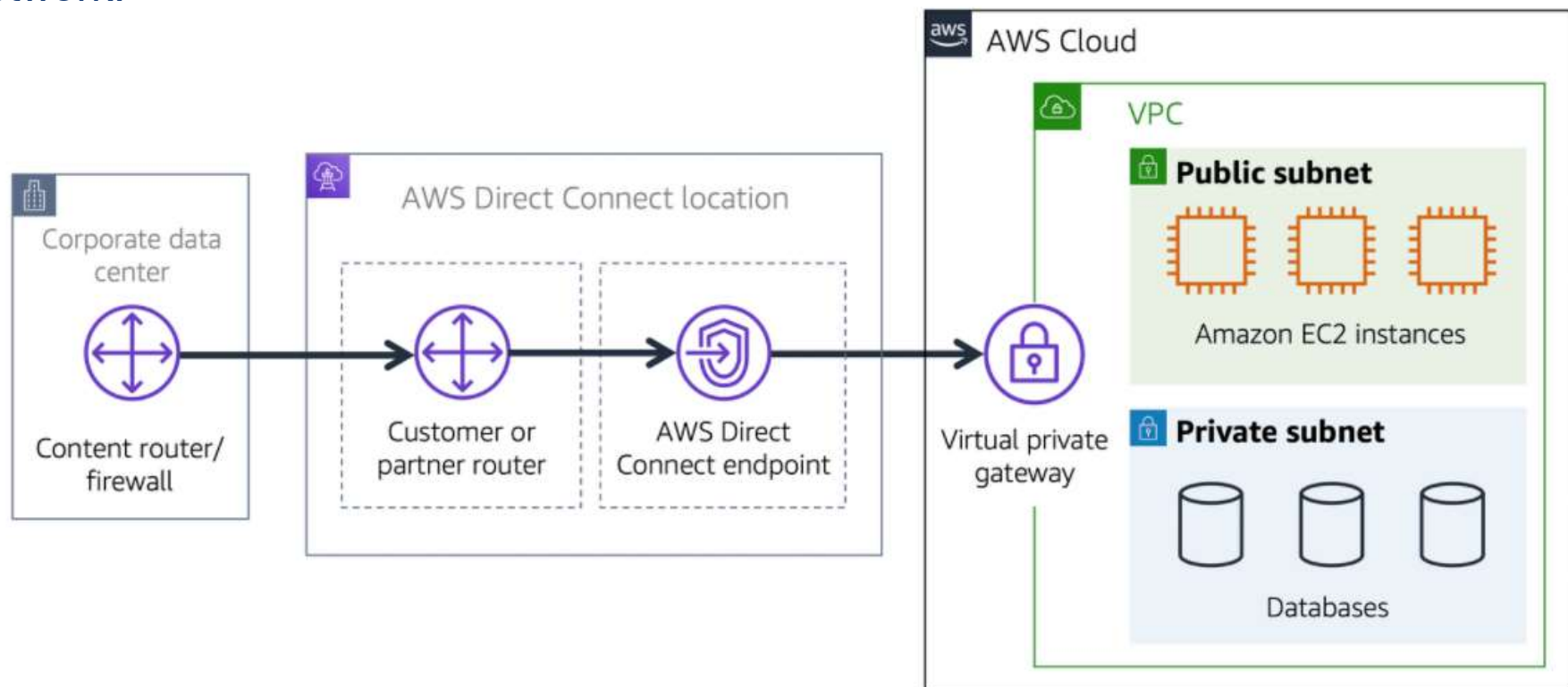To access private resources in a VPC, you can use a **virtual private gateway**.

A virtual private gateway enables you to establish a virtual private network (VPN) connection between your VPC and a private network, such as an on-premises data center or internal corporate network. A virtual private gateway allows traffic into the VPC only if it is coming from an approved network.

# AWS Direct Connect

**AWS Direct Connect** is a service that enables you to establish a dedicated private connection between your data center and a VPC.

The private connection that AWS Direct Connect provides helps you to reduce network costs and increase the amount of bandwidth that can travel through your network.
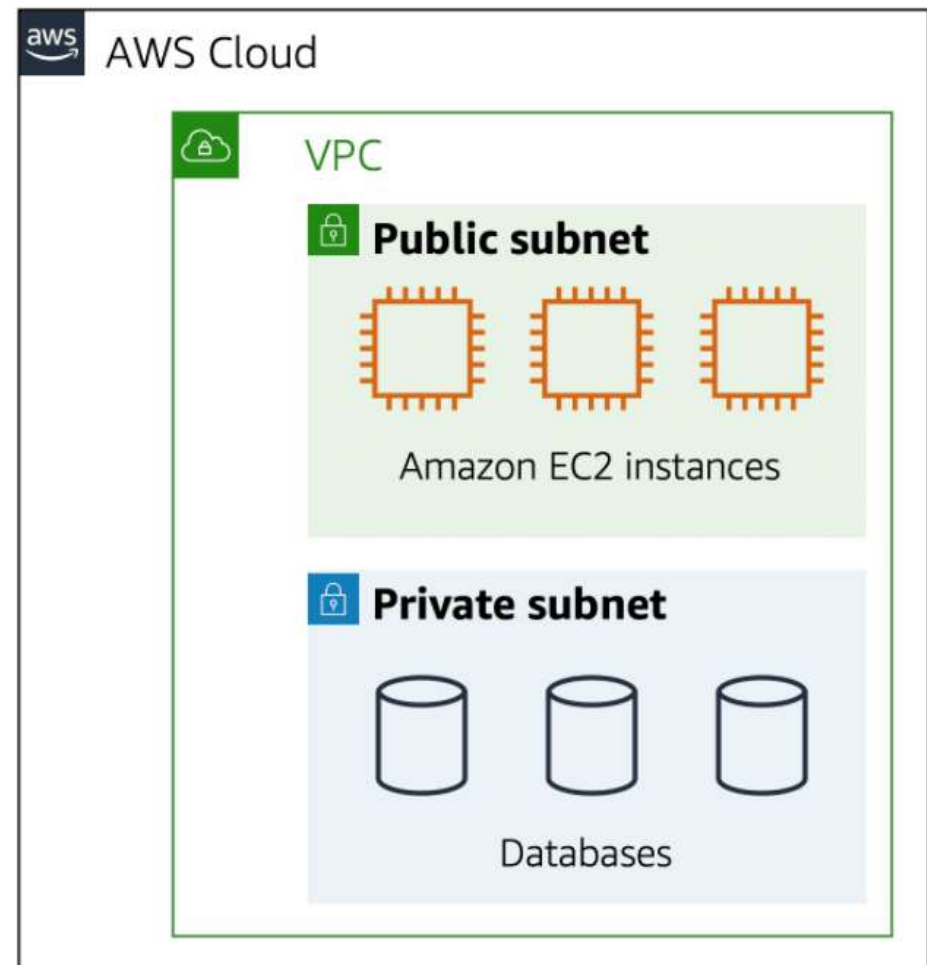
# Subnets

A subnet is a section of a VPC in which you can group resources based on security or operational needs. Subnets can be public or private.

**Public subnets** contain resources that need to be accessible by the public, such as an online store's website.

**Private subnets** contain resources that should be accessible only through your private network, such as a database that contains customers' personal information and order histories.

In a VPC, subnets can communicate with each other. For example, you might have an application that involves Amazon EC2 instances in a public subnet communicating with databases that are located in a private subnet.



AWS Cloud

VPC

Public subnet

Amazon EC2 instances

Private subnet

Databases

**Quintrix**
A MINDLANCE COMPANY

# Network traffic in a VPC

When a customer requests data from an application hosted in the AWS Cloud, this request is sent as a packet. A **packet** is a unit of data sent over the internet or a network.

It enters into a VPC through an internet gateway. Before a packet can enter into a subnet or exit from a subnet, it checks for permissions. These permissions indicate who sent the packet and how the packet is trying to communicate with the resources in a subnet.
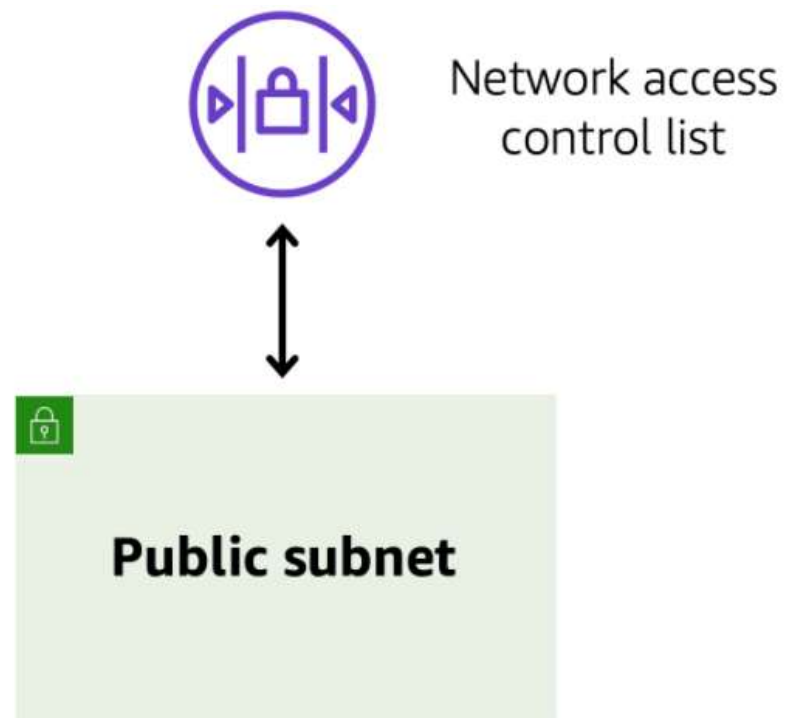
The VPC component that checks packet permissions for subnets is a **network access control list (ACL)**

**Quintrix**
A MINDLANCE COMPANY

# Network Access Control List (ACLs)

A network access control list (ACL) is a virtual firewall that controls inbound and outbound traffic at the subnet level.

Each AWS account includes a default network ACL. When configuring your VPC, you can use your account's default network ACL or create custom network ACLs.
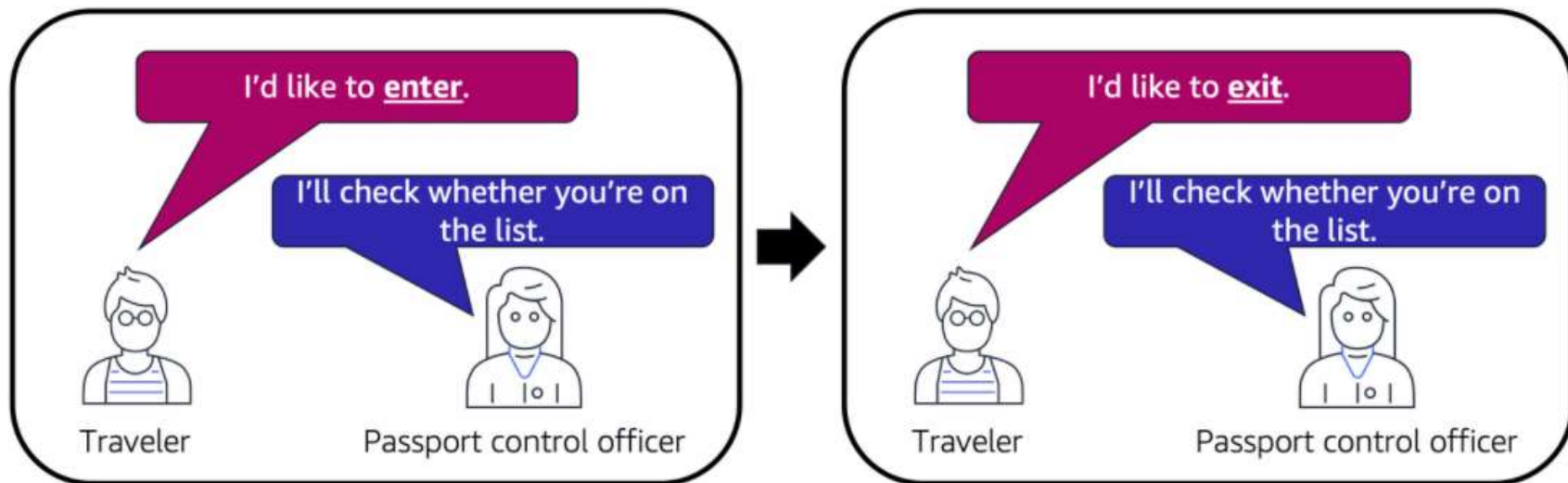
By default, your account's default network ACL allows all inbound and outbound traffic, but you can modify it by adding your own rules. For custom network ACLs, all inbound and outbound traffic is denied until you add rules to specify which traffic to allow. Additionally, all network ACLs have an explicit deny rule. This rule ensures that if a packet doesn't match any of the other rules on the list, the packet is denied.

Network access control list

**Public subnet**

**Quintrix**
A MINDLANCE COMPANY

# Stateless packet filtering

Network ACLs perform **stateless** packet filtering. They remember nothing and check packets that cross the subnet border each way: inbound and outbound.

When a packet response for that request comes back to the subnet, the network ACL does not remember your previous request. The network ACL checks the packet response against its list of rules to determine whether to allow or deny.



After a packet has entered a subnet, it must have its permissions evaluated for resources within the subnet, such as Amazon EC2 instances.
The VPC component that checks packet permissions for an Amazon EC2 instance is a **security group**.
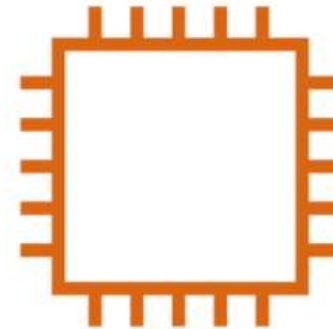
Quintrix
A MINDLANCE COMPANY

# Security groups

A security group is a virtual firewall that controls inbound and outbound traffic for an Amazon EC2 instance.

By default, a security group denies all inbound traffic and allows all outbound traffic. You can add custom rules to configure which traffic to allow or deny.

For this example, suppose that you are in an apartment building with a door attendant who greets guests in the lobby. You can think of the guests as packets and the door attendant as a security group. As guests arrive, the door attendant checks a list to ensure they can enter the building. However, the door attendant does not check the list again when guests are exiting the building

If you have multiple Amazon EC2 instances within a subnet, you can associate them with the same security group or use different security groups for each instance.
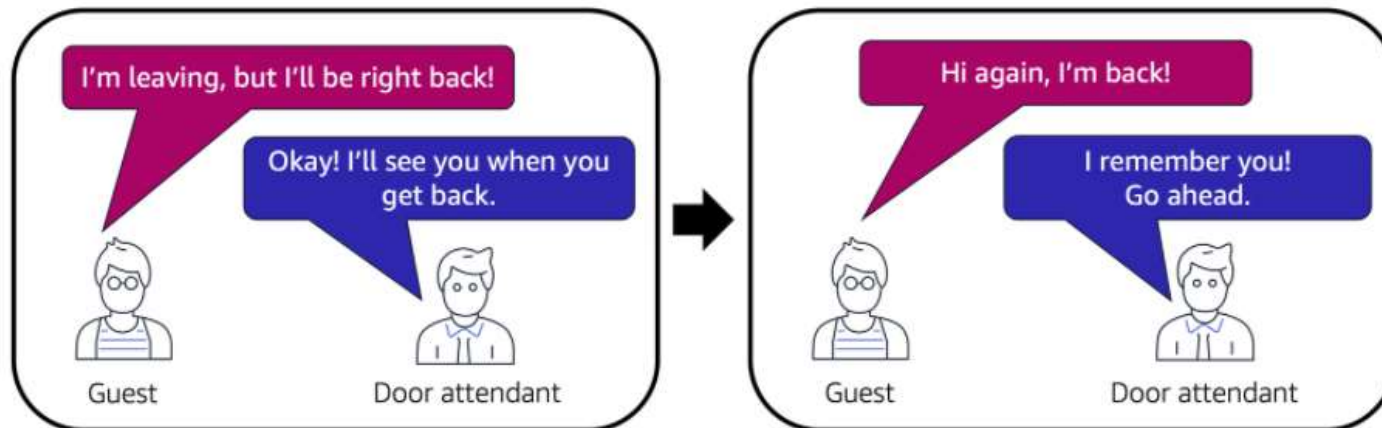
Security group

Amazon EC2 instance

# Stateful packet filtering

Security groups perform **stateful** packet filtering.  Consider the same example of sending a request out from an Amazon EC2 instance to the internet.

When a packet response for that request returns to the instance, the security group remembers your previous request. The security group allows the response to proceed, regardless of inbound security group rules.
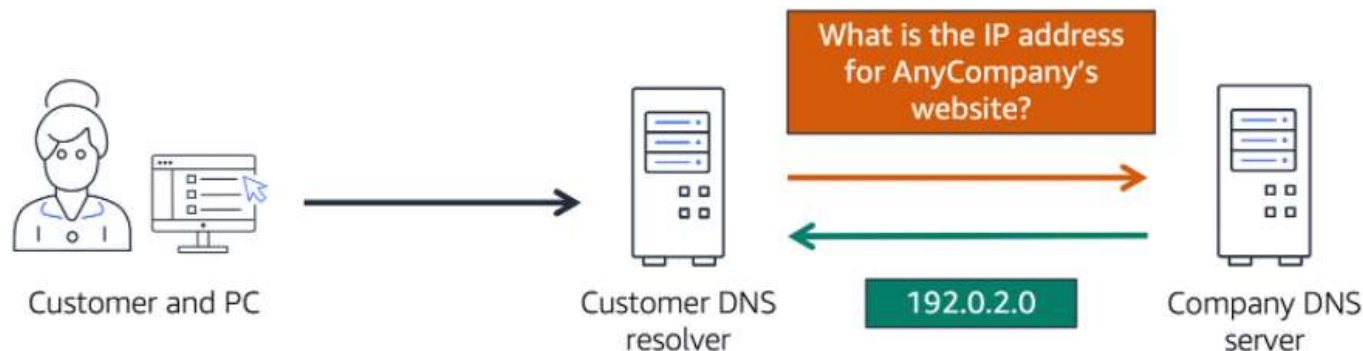


Both network ACLs and security groups enable you to configure custom rules for the traffic in your VPC. As you continue to learn more about AWS security and networking, make sure to understand the differences between network ACLs and security groups.

**Quintrix**
A MINDLANCE COMPANY

# Domain Name System (DNS)

Suppose that AnyCompany has a website hosted in the AWS Cloud. Customers enter the web address into their browser, and they are able to access the website.

This happens because of **Domain Name System (DNS)** resolution. DNS resolution involves a customer DNS resolver communicating with a company DNS server. You can think of DNS as being the phone book of the internet. DNS resolution is the process of translating a domain name to an IP address.

What is the IP address for AnyCompany's website?

Customer and PC

Customer DNS resolver

192.0.2.0

Company DNS server

Quintrix
A MINDLANCE COMPANY

# Domain Name System (DNS)

For example, suppose that you want to visit AnyCompany's website.

1. When you enter the domain name into your browser, this request is sent to a customer DNS resolver.

2. The customer DNS resolver asks the company DNS server for the IP address that corresponds to AnyCompany's website.

3. The company DNS server responds by providing the IP address for AnyCompany's website, 192.0.2.0.



Customer and PC          Customer DNS          What is the IP address          Company DNS
                         resolver              for AnyCompany's                 server
                                               website?
                                               192.0.2.0

Quintrix
A MINDLANCE COMPANY

# Amazon Route 53

**Amazon Route 53** is a DNS web service.

It gives developers and businesses a reliable way to route end users to internet applications hosted in AWS.

Amazon Route 53 connects user requests to infrastructure running in AWS (such as Amazon EC2 instances and load balancers). It can route users to infrastructure outside of AWS.

Another feature of Route 53 is the ability to manage the DNS records for domain names. You can register new domain names directly in Route 53. You can also transfer DNS records for existing domain names managed by other domain registrars. This enables you to manage all of your domain names within a single location.

**Quintrix**
A MINDLANCE COMPANY

# Amazon Route 53

Suppose that AnyCompany's application is running on several Amazon EC2 instances. These instances are in an Auto Scaling group that attaches to an Application Load Balancer.

1. A customer requests data from the application by going to AnyCompany's website.
2. Amazon Route 53 uses DNS resolution to identify AnyCompany.com's corresponding IP address, 192.0.2.0. This information is sent back to the customer.
3. The customer's request is sent to the nearest edge location through Amazon CloudFront.
4. Amazon CloudFront connects to the Application Load Balancer, which sends the incoming packet to an Amazon EC2 instance.

Quintrix
A MINDLANCE COMPANY