# O'REILLY®

**AWS Core Architecture Concepts**

# What we will cover:

- Fundamentals of AWS: architecture, terminology and concepts

- Virtual Private Cloud (VPC): Networking services

- Elastic Compute Cloud (EC2): Instance deployment and configuration

- Storage solutions: Elastic Block Storage (EBS) and snapshot management

- Simple Storage Service (S3): Object storage
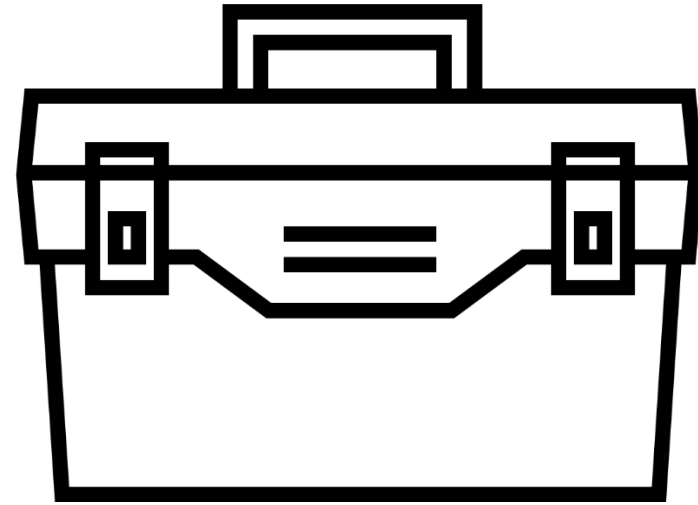
- S3 Glacier: Archive storage

# AWS Core Cloud Services

- AWS Administration – Management portal

- Compute Services – Elastic compute cloud

- Networking Services – Virtual private cloud

- Auto Scaling – Scale EC2 compute automatically

- Elastic Load Balancing – Distribute traffic across EC2 instances or containers

- Elastic Block Storage – Virtual hard drives

- S3 – Durable and scalable object storage

- S3 Glacier – Long-term data archiving

# AWS Services are Managed Services

- Managed services: AWS does most of the setup

- Less managed services: You can do whatever you want

- You to do most of the setup, management, and monitoring (VPC, EC2)

- The reality – there are no completely unmanaged services at AWS
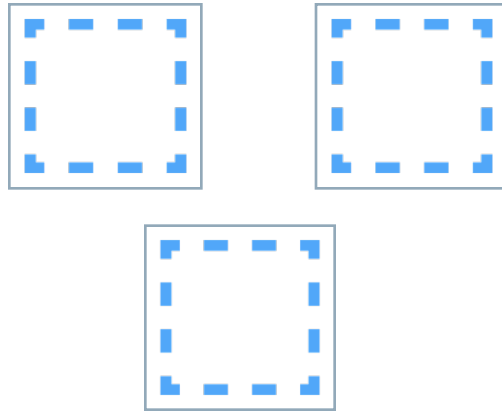
# Demo: Management Services

# AWS Regions

# AWS Regions

Regions start off as independent

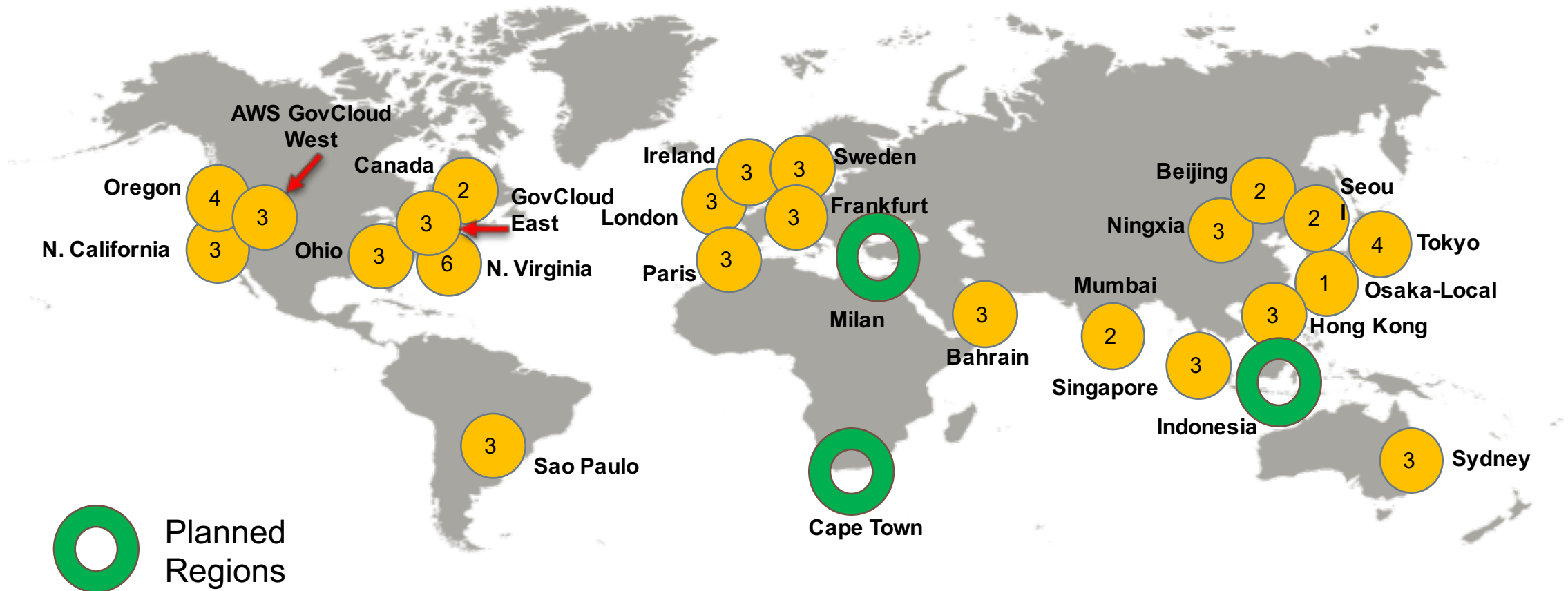Regions have (multiple) Availability Zones

Data transfer charges apply across regions

Resources are not automatically replicated between regions by default

# Choosing an AWS Region

**Latency – to on-premise customer location**

**Costs are different per Region, and AZ**

**Features are different per region**
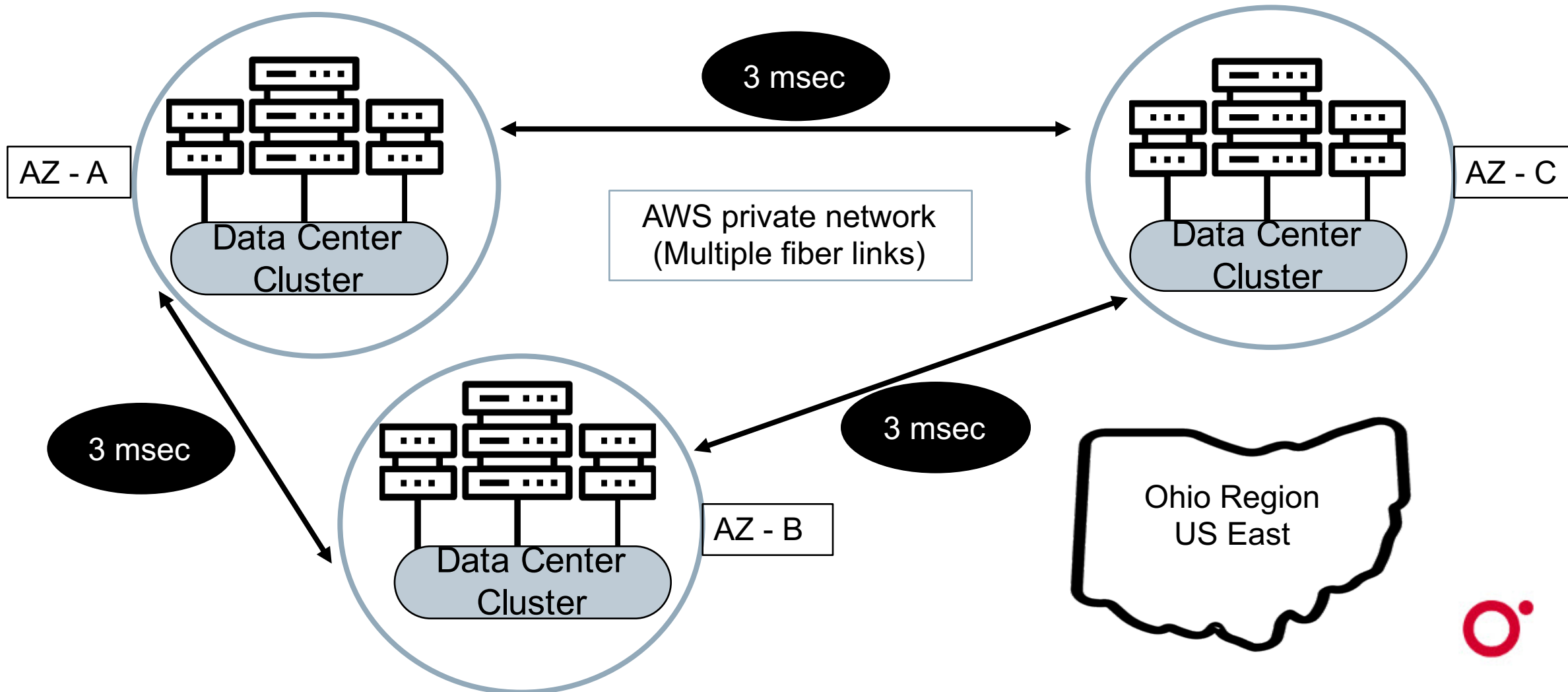
**Compliance: Industry, country, and business**
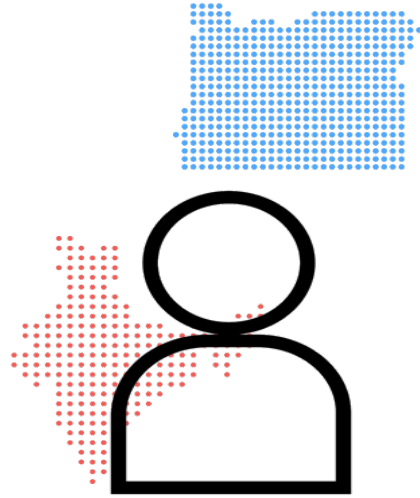
# Demo: Regions

# Availability Zones (AZ)
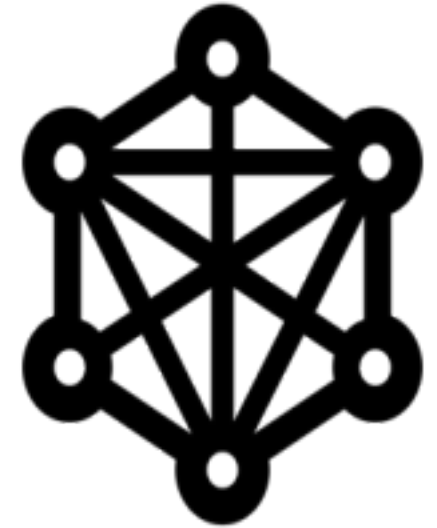
# Availability Zones  (AZ)

Designed as an independent failure zone: Separate power sources

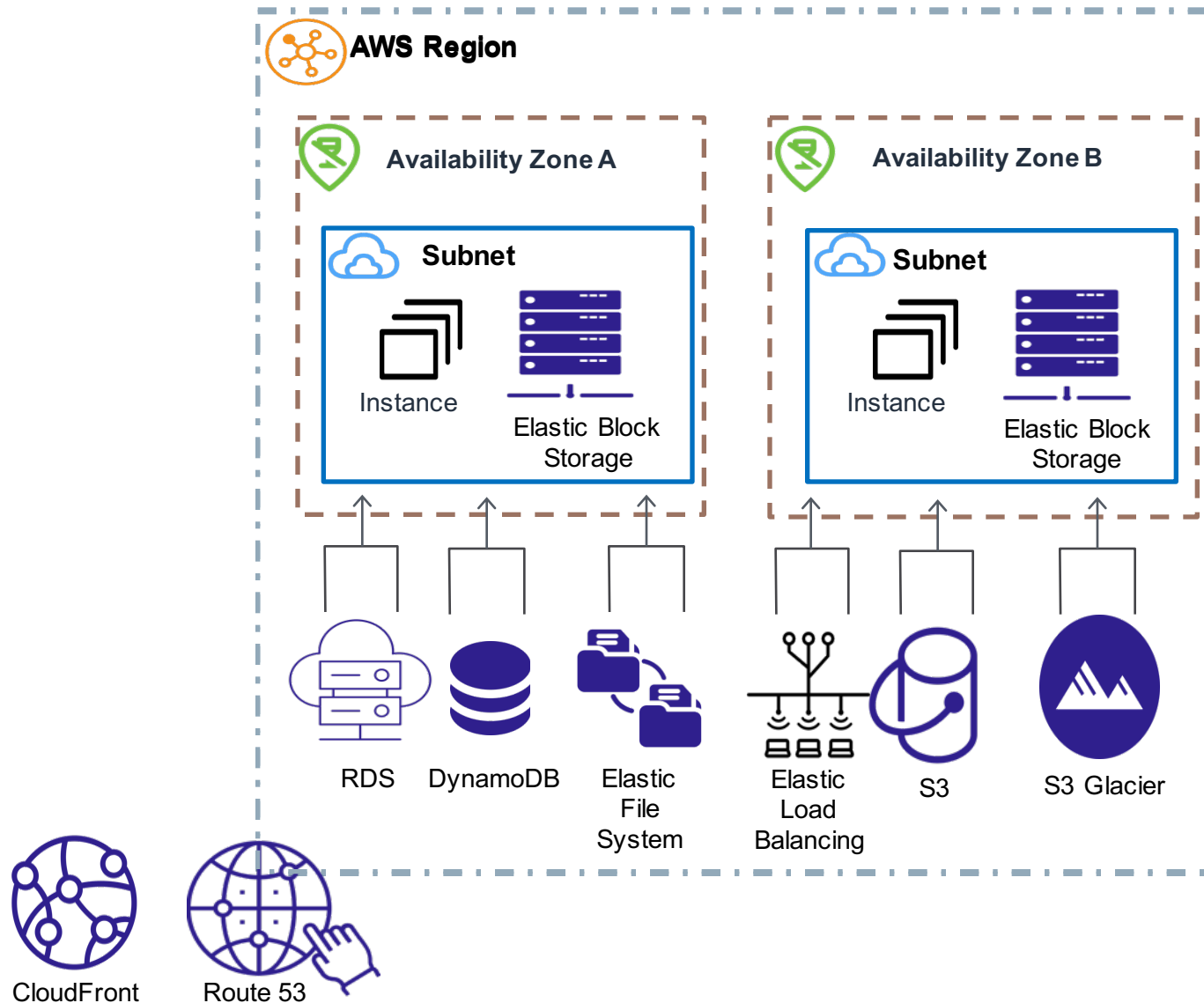AWS account has access to all regions and associated AZ's

Data transfer charges for outbound traffic

Redundant Tier-1 transit private fiber connections

# Regional Services

# Single or Multi-AZ Design

| Single - AZ | Multi - AZ |
|---|---|
| ▪No recovery or failover when disaster happens in a single datacenter<br><br>▪No high availability for instances<br><br>▪No failover in single datacenter<br><br>▪All AWS regions have at least 2 availability zones<br><br>▪Each AZ has at least one datacenter | ▪Better high availability design options<br><br>▪Designing applications hosted across AZ's provides HA options<br><br>▪Load balancing (ELB) supports targeting instances in multiple availability zones<br><br>▪EC2 auto scaling supports multiple AZ's<br><br>▪Use Route 53 (DNS) to balance across multiple AWS regions |

# Demo: Availability Zones

# Edge Locations

# Edge Locations

180 edge locations

11 Regional Edge caches in 73 cities

# Edge Locations @ AWS

# Services at the Edge

| | | |
|---|---|---|
| Caches your content request | Delivers your request to closest edge location | Filters incoming public traffic to the edge |
| CloudFront | Route 53 | WAF |

User (in Singapore)

Website

Origin Server (U.S.A.)

CloudFront

Edge Location

# Demo: Edge Services

# Virtual Private Cloud

# What's a VPC?

- Networking layer at AWS

- A logical and isolated data-center (virtual private cloud)

- Launch EC2 instances and various AWS resources into your virtual network

- Logically isolated from all other virtual networks hosted in the AWS cloud

- Instances run in a virtual private cloud that is logically isolated to your AWS account
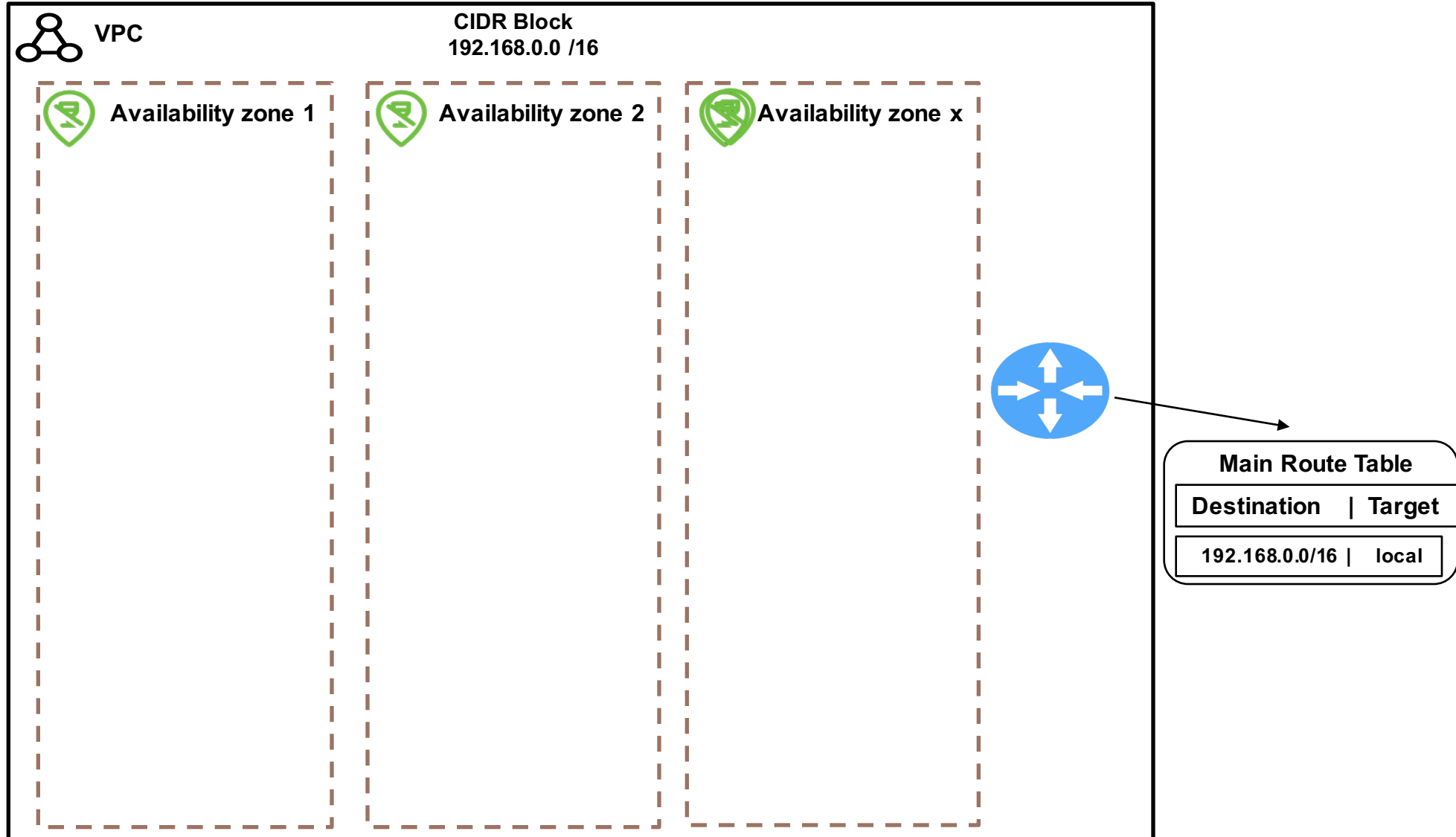
# Creating a New VPC

- When a VPC is created, it spans all the availability zones within the region

- Subnets can be created in each availability zone
  - Each subnet is defined by a CIDR block which is a subset of the VPC CIDR block

- Each subnet is assigned with a default route table that enables local routing throughout the VPC

# VPC's have Multiple AZ's

VPC

**CIDR Block**
**192.168.0.0 /16**

Availability zone 1

Availability zone 2

Availability zone x

**Main Route Table**

| Destination | Target |
|---|---|
| 192.168.0.0/16 | local |

# VPC Core Components

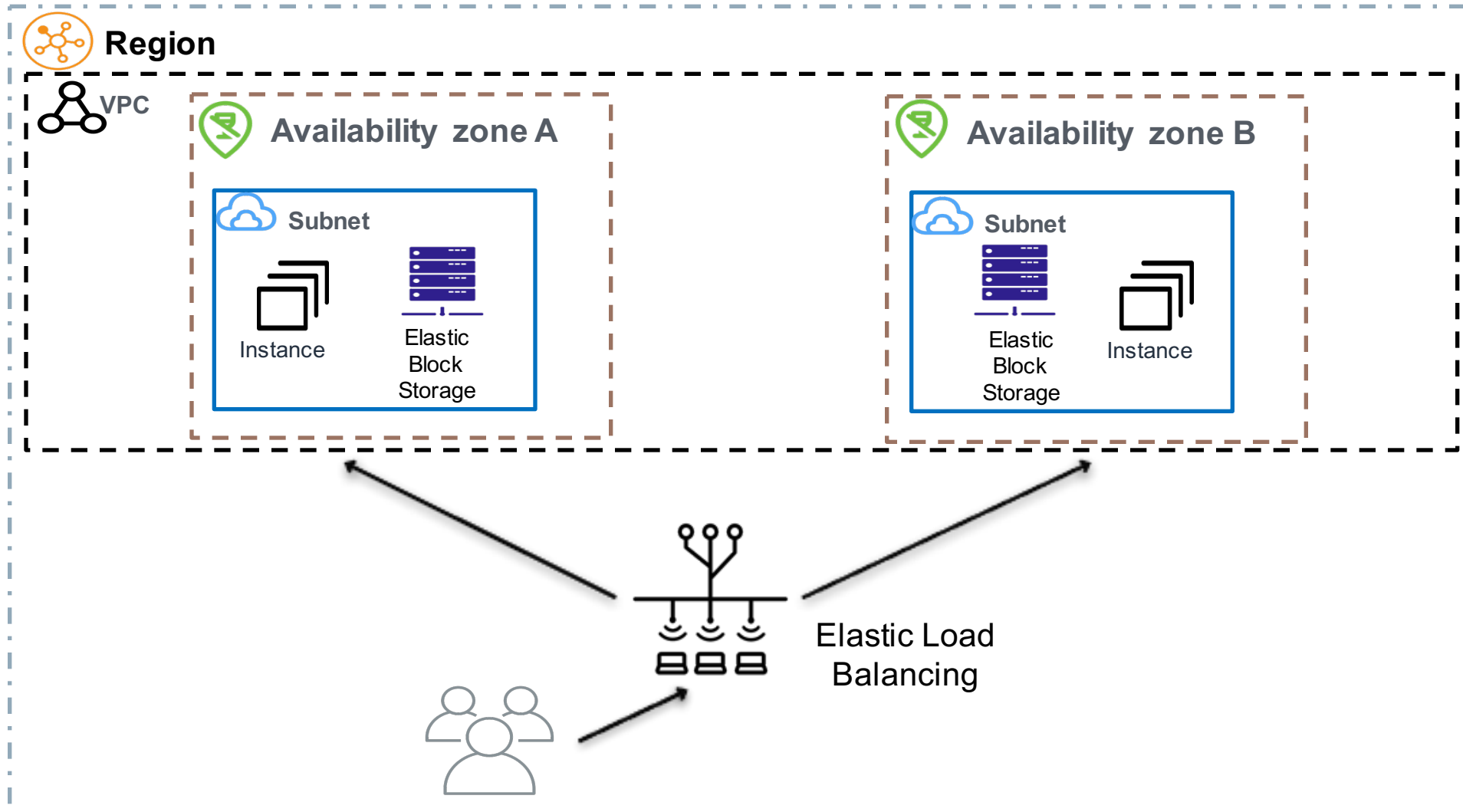| | |
|---|---|
| ▪ Subnets<br><br>▪ Route tables<br><br>▪ Security groups (SG)<br><br>▪ Network access control list (NACLs)<br><br>▪ Internet gateway (IGW) | ▪Virtual private gateway (VGW)<br><br>▪Private endpoints<br><br>▪Peering connections<br><br>▪NAT gateway services<br><br>▪Transit gateway |

# Failover Possibilities with AZ's

# Demo: Create a VPC

# The Default VPC

- /20 CIDR Block is assigned by default

- An internet gateway is connected to the default VPC

- Default route table sends internet traffic to the internet gateway

- Default security group

- Default network access control list

- Default subnets

- Instances are assigned both a private and public IPv4 address

# Internet Gateway (IGW)

- Allows communication between instances or services hosted on public subnets and the Internet

- The internet gateway is a managed AWS service providing Internet access from public subnets

- To enable access to the Internet you must:
  - Order an IGW
  - Attach the IGW to your VPC
  - Add route table entry pointing to the IGW for the public subnet
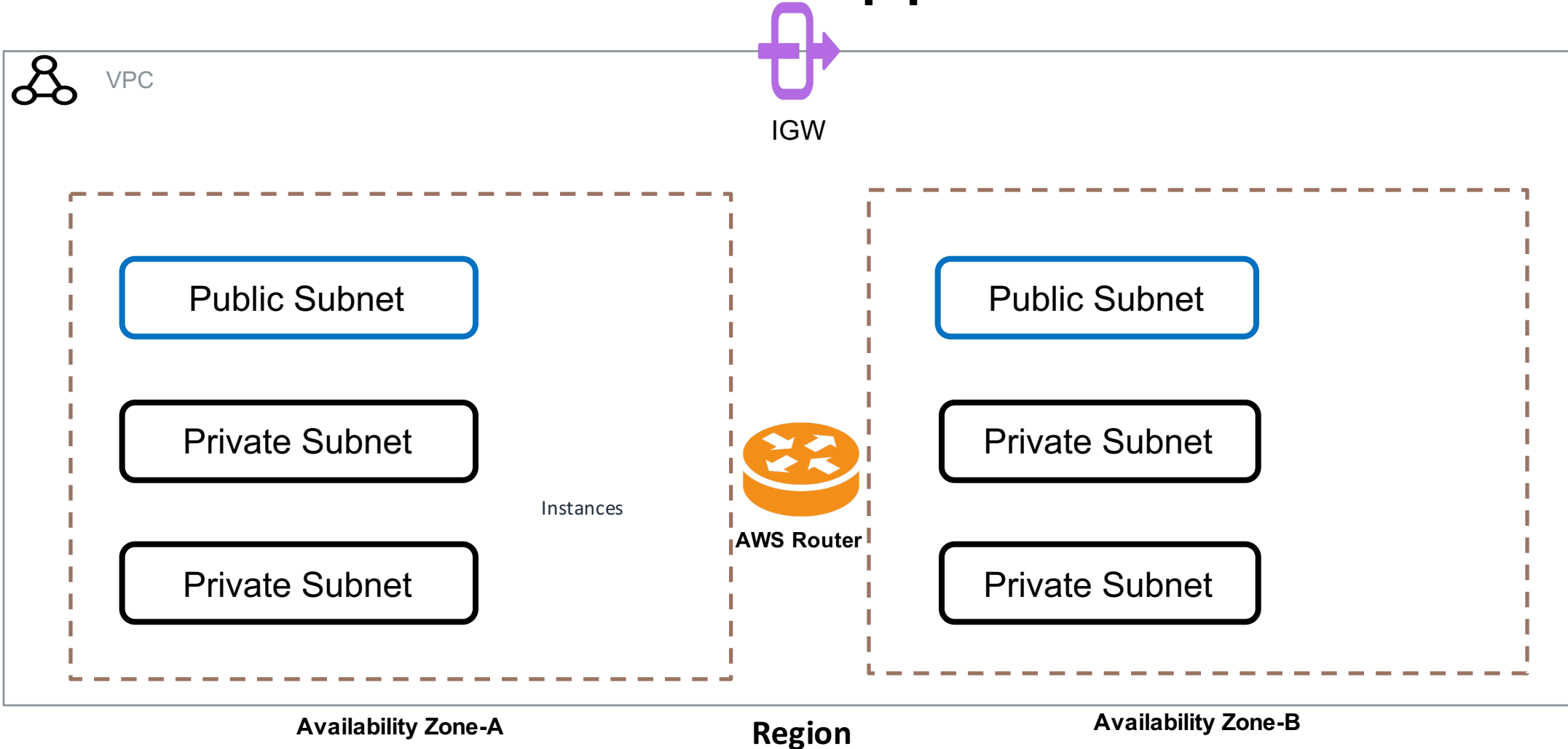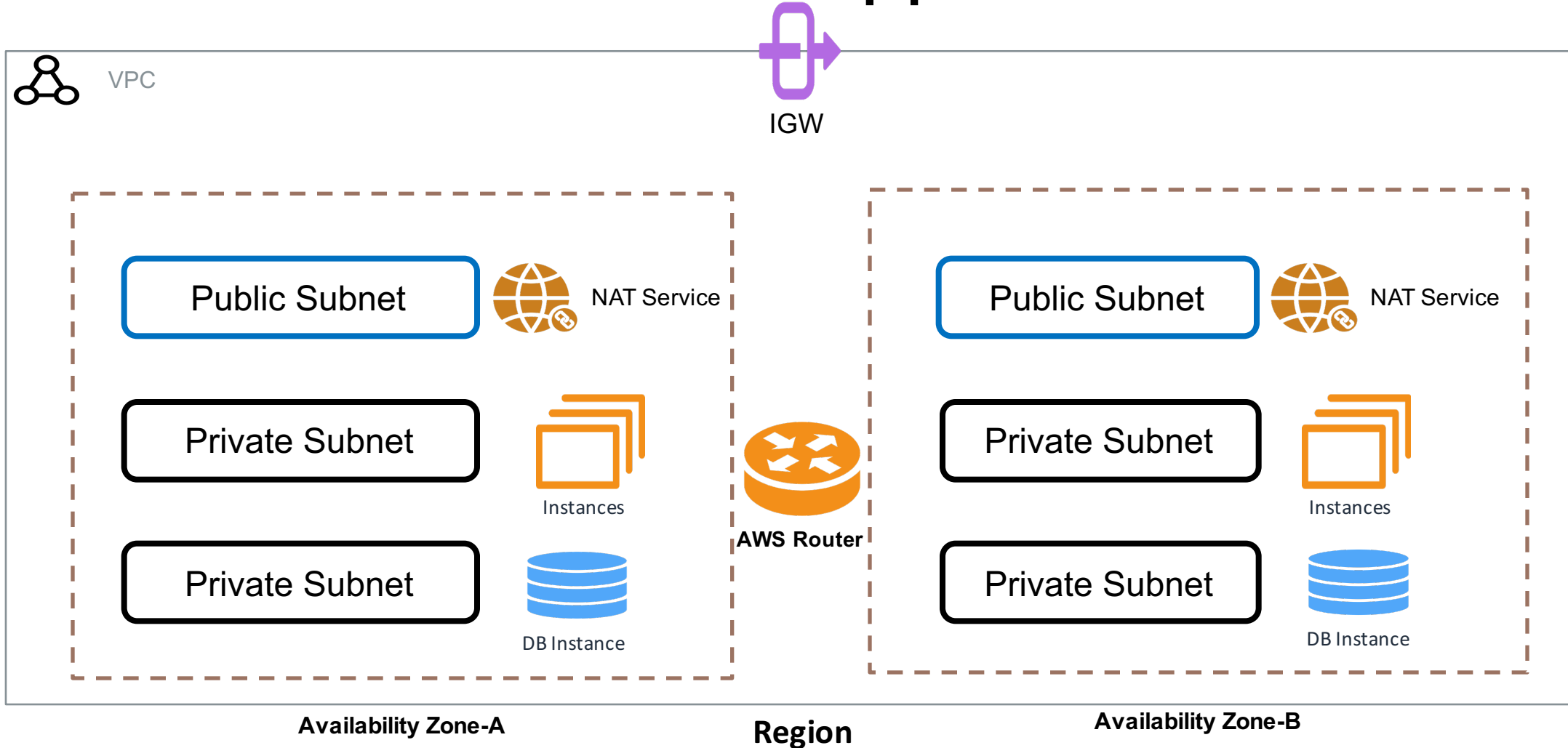  - Instances must have a Public IP address

# Demo: Internet Gateway

# Two Tier Application

# Two Tier Application

# NAT Services

- NAT service enable instances in a private subnet to connect to the Internet to get updates

- Traffic requests from the instance are forwarded to the NAT service hosted in the public subnet

- Internet response is sent back to the private instance that made the request

- NAT Options:
  - NAT gateway service – hosted NAT services provided by AWS
  - NAT instance – compute instance created with a NAT AMI

# Updates using NAT Service

# Updates using NAT Instances

Demo: Nat Gateway

# Subnets

- Public or private subnets can be created in each availability zone

- Subnets cannot span across multiple availability zones

- A subnet that doesn't route to an internet gateway is a private subnet

- If a subnet has traffic routed to an internet gateway it is defined as a public subnet

- EC2 instances in a public subnet must have a public IP address, or an Elastic IP address to be able to communicate with an internet gateway

# Subnets

- Instances and AWS services are launched into subnets

- Public subnets can be used for resources that need Internet access (IGW, ELB)

- Private subnets host resources that don't directly connect to the Internet (Instances, RDS)

- Protect subnet access using optional network access control lists (NACLs)

- Network ACL operates at the subnet level and supports allow and deny rules

# Demo: Subnets

# IP Addresses

- Each EC2 instance is assigned a private DNS host name associated with it's private IP address

- Both IP version 4 and IP version 6 addressing is supported

- IPv4 is default and required; IPv6 is optional

- Address types for EC2 instances:
  - Private IP version 4
  - Public IP version 4
  - Elastic IP address (Static public IP)
  - IP version 6 address (Public)

# IP Address Assignment

- A private IPv4 address is not directly reachable from the Internet

- When an EC2 instance is stopped and restarted, the private IP address remains assigned

- A public IP address is reachable from the Internet

- When an EC2 instance is stopped and restarted, the AWS public IP address is unassigned, and a new IP address is assigned

- Public IP addresses are assigned to your instance from Amazon's pool of public IPv4 addresses, as an EIP, or BYOIP

# Elastic IP Addresses (EIPs)

- An elastic IP address is a static IP address (public or private)

- Elastic IP addresses are assigned to your AWS account from AWS or BYOIP ranges

- A public EIP is first allocated for use within a VPC; then assigned to a specific EC2 instance or network interface within your AWS Account

- EIPs can be remapped from one EC2 instance to another EC2 instance within the same or another VPC

# Bring Your Own IP

- Bring your public IPv4 address range to your AWS account

- AWS advertises your public IP address range on the Internet

- Your public IP address range once moved, will appear in your AWS account as a public address pool

- /24 address ranges supported

- Create elastic IP addresses from your public address pool

- Used for EC2 instances, NAT gateways, and Elastic Load Balancers

Demo: IP Addresses

# Route Tables

- Each subnet must be associated with a route table

- External traffic patterns are defined by adding additional routes

- Each route specifies a destination and a target

- Route table rules allow VPC traffic to connect to an Internet gateway (IGW), a Virtual private gateway (VGW), or to a NAT gateway service

- When you create a VPC, it automatically is associated with a default main route table

# Route Tables

- The main route table controls the routing for all subnets that are not explicitly associated with any other route table

- Each route table contains a default entry displayed as "local route" that enables local communication within each VPC

- IP traffic that is routed within the VPC uses the local route

# Route Tables

▪Each new subnet is automatically associated with the default route table that was created when the VPC was first created

▪Subnets can be associated explicitly with a custom route table, or explicitly / implicitly with the main route table

▪Multiple subnets can be associated with the same route table

# Route Tables

Internet gateway

**Availability Zone A**

EC2 Instances

**Subnet 1   10.0.0.1/24**

**Custom Route Table**

| Destination | Target |
|---|---|
| 10.0.0.0/16 | local |
| 0.0.0.0/0 | IGW |

**Availability Zone B**

EC2 Instances

**Subnet 2   10.0.0.2/24**

**Main Route Table**

| Destination | Target |
|---|---|
| 10.0.0.0/16 | local |

**Availability Zone C**

**Subnet 3   10.0.0.3/24**

VGW

VPN Connection

**Custom Route Table**

| Destination | Target |
|---|---|
| 10.0.0.0/16 | local |
| 0.0.0.0.0/0 | VGW |

VPC 10.0.0.0/16

# Demo: Route Tables

# Security Groups

- Security groups are defined as "virtual firewall' protecting EC2 instance's inbound and outbound traffic

- Security groups contain rules that control the inbound and outbound traffic to an instance

- Each instance launched into a VPC can have up to 5 security groups

- Each SG can have 50 inbound / outbound rules

- Each VPC can have up to 500 Security Groups

- When security groups are created, they are linked to a VPC
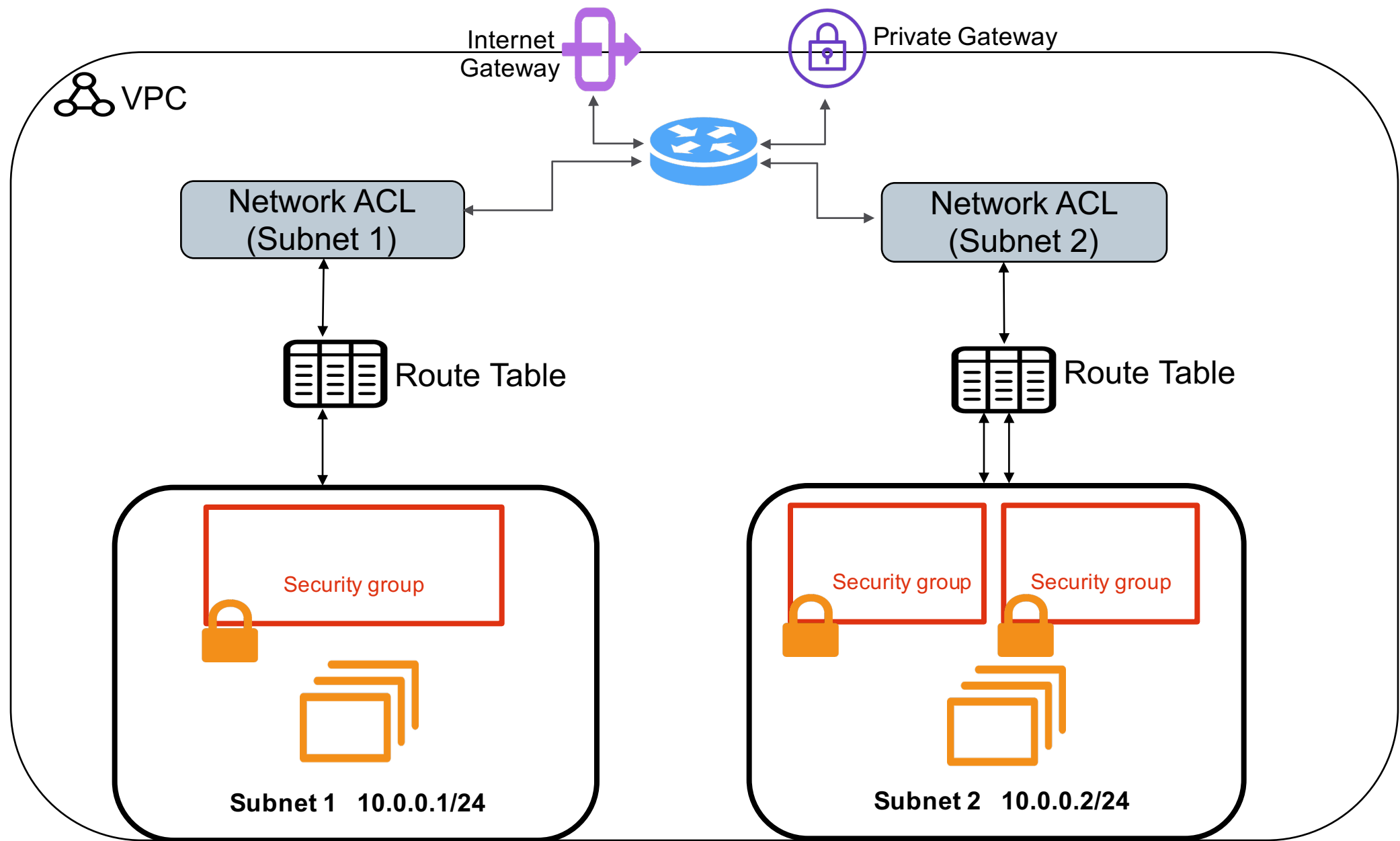
# Security Group Operation

- Allow rules can be specified

- Explicit deny rules can't be specified

- Inbound rules define the source of the traffic, and the destination port, port range, or security group

- Any TCP protocol that is defined with a standard port number is supported

- Outbound rules can define the destination for the traffic and the destination port or security group

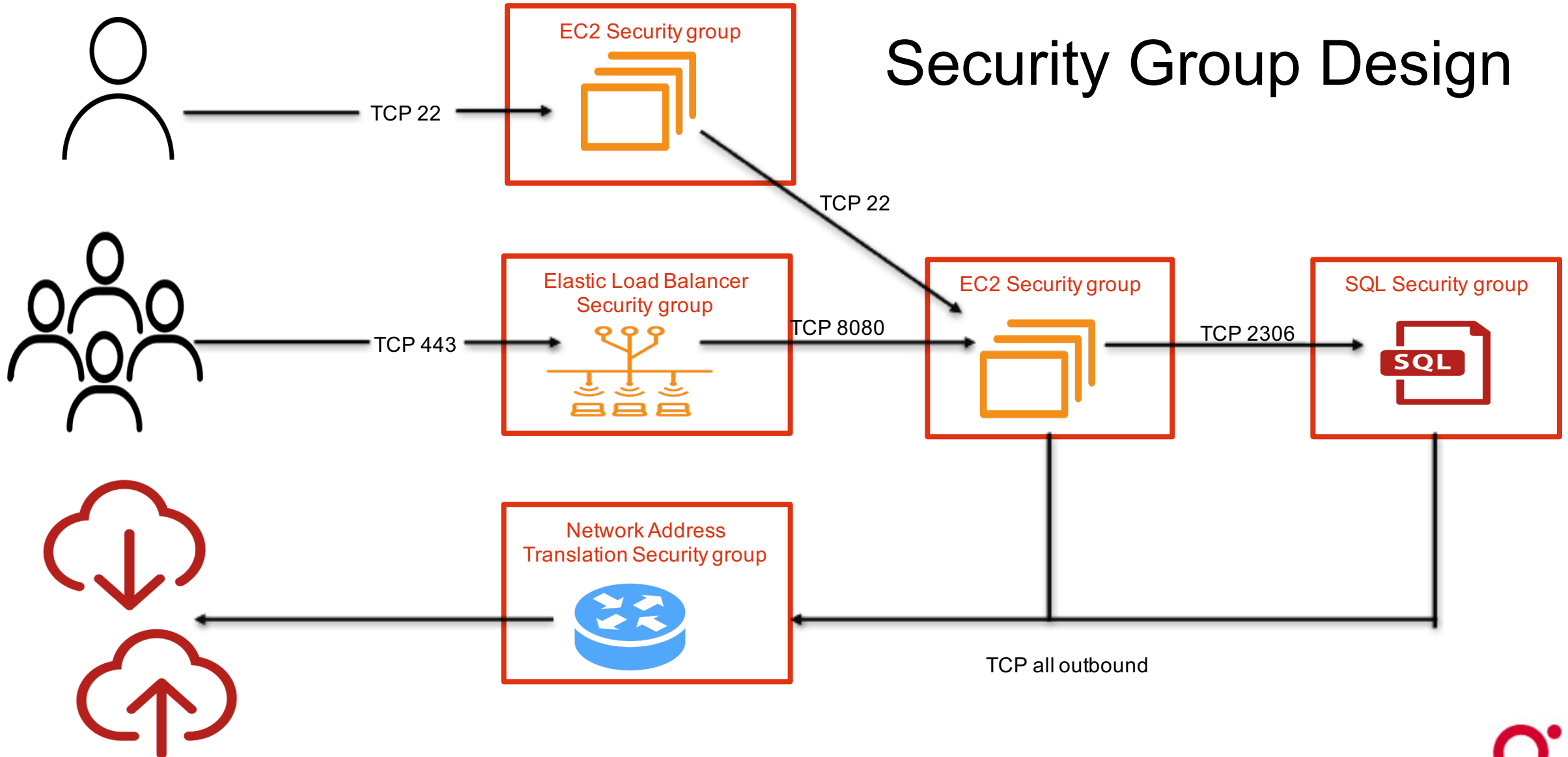- Separate allow rules can be defined for both inbound and outbound traffic

# Security Group Operation

- Security groups are defined as stateful – if a request is made inbound to an EC2 instance, the response traffic for the incoming request is allowed outbound

- Traffic is restricted to:

- IP protocol

- Service port

- Source / destination IP address, or address block

- Security group

# Demo: Security Group

# Network ACL's

- NACLs are an optional security control for subnets

- NACLs act as an "subnet firewall" for controlling traffic in and out of each subnet

- The default network ACL for a VPC allows all inbound and outbound IPv4 traffic

- Each subnet is associated with a single NACL

- A single network ACL can be associated with multiple subnets within the same VPC

# Network ACL Rules

- Inbound Rule

- Allow or deny for the specified traffic pattern

- Outbound Rule

- Allow or deny for the specified traffic pattern

# Network ACL Operation

- NACL rules are defined as stateless

- Rules are evaluated in order until a match is found

- Evaluation starts with the lowest numbered rule to determine if traffic is allowed in or out of the subnet associated with the network ACL

- Create rules in multiples of 10, so adding new rules doesn't cause problems in the future

# Security Groups vs NACLs

| Security Groups | NACLs |
|---|---|
| ▪Operates at the EC2 instance level<br><br>▪Allow rules only supported<br><br>▪Stateful: Return traffic is automatically allowed<br><br>▪All rules are processed before traffic decisions are made<br><br>▪Applied to the selected EC2 instance elastic network adapter | ▪Operates at the subnet level<br><br>▪Allow and deny rules supported<br><br>▪Stateless: Return traffic must be explicitly allowed by a rule<br><br>▪Rules are processed in numerical order before traffic decisions are made<br><br>▪Applied to a subnet |

# VPC Flow Logs

- Flow logs can be created for a VPC, a subnet, or a network interface

- Logs IP traffic to and from network interfaces in a VPC (accepted / rejected)

- Each NIC has a unique log stream

- Flow log data is published to a log group stored as a CloudWatch log group, or S3 Bucket

- Does not capture DNS, license, metadata, or default VPC router traffic

# Demo: Flow Logs

# VPC Private Endpoints

- Privately connect your VPC to AWS services

- Interface endpoint – network interface with private IP address to hosted AWS service

- Gateway endpoint – S3 bucket or Dynamo DB table

# Peering VPC's

- Networking connection between two VPC's

- Peer your VPC's or between other account holders VPC's using a private IP address

- Peering is a one-to-one relationship

- Peering connections are not transitive

- CIDR blocks can't overlap in a peering relationship

- Peering connections can be created between VPCs in the same region

- Peering connections can be created between VPCs in different regions

# Demo: Private Connections

# EC2 Instances

# Instance Hosting since 2017

Performance of storage, networking, management improved with custom chipsets

Hardware replaces software emulation speeding up EC2 communications

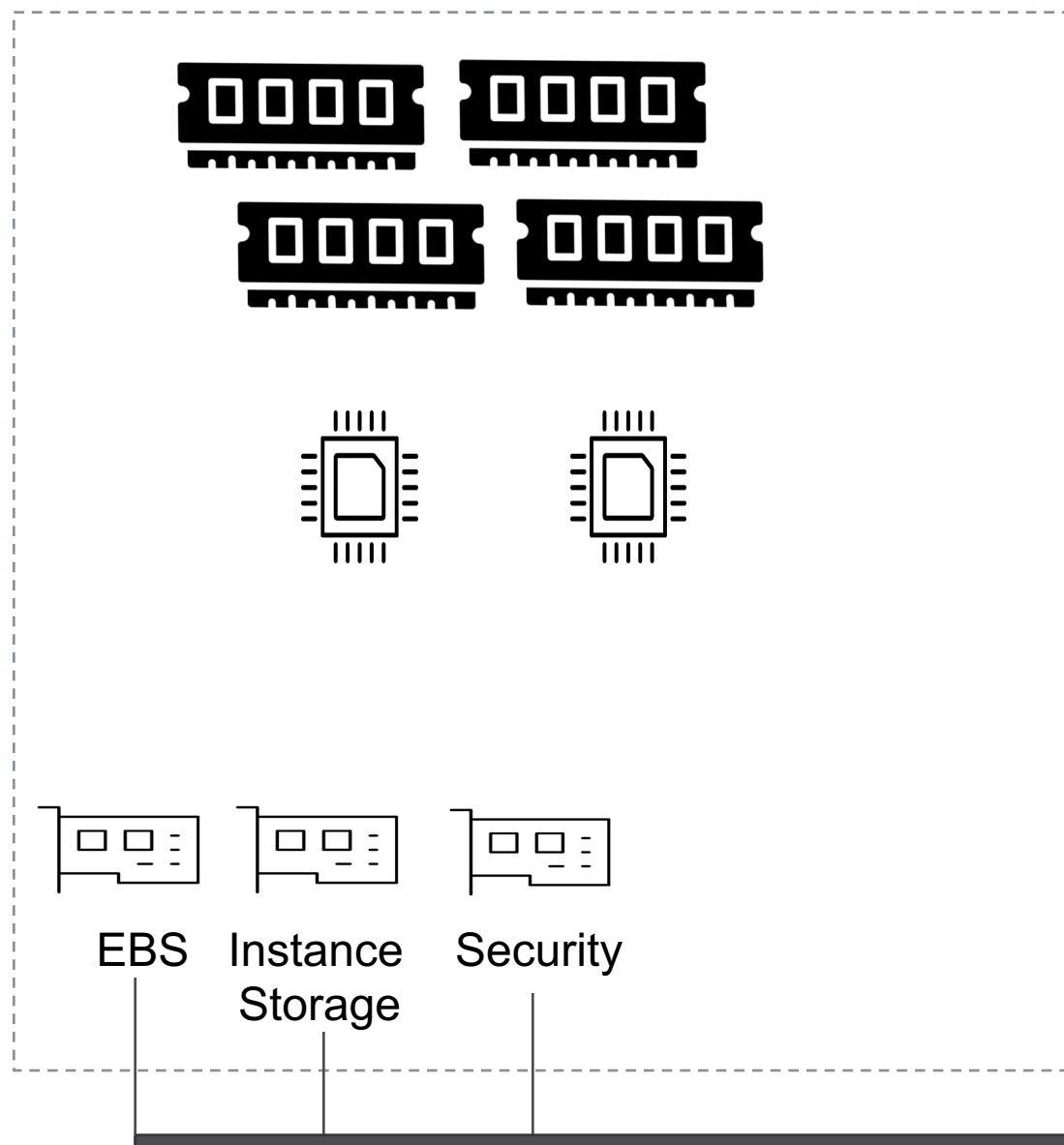Replacing Xen hypervisor with lightweight hypervisor called Nitro

Lightweight: hypervisor tasks performed by the new custom hardware chipsets

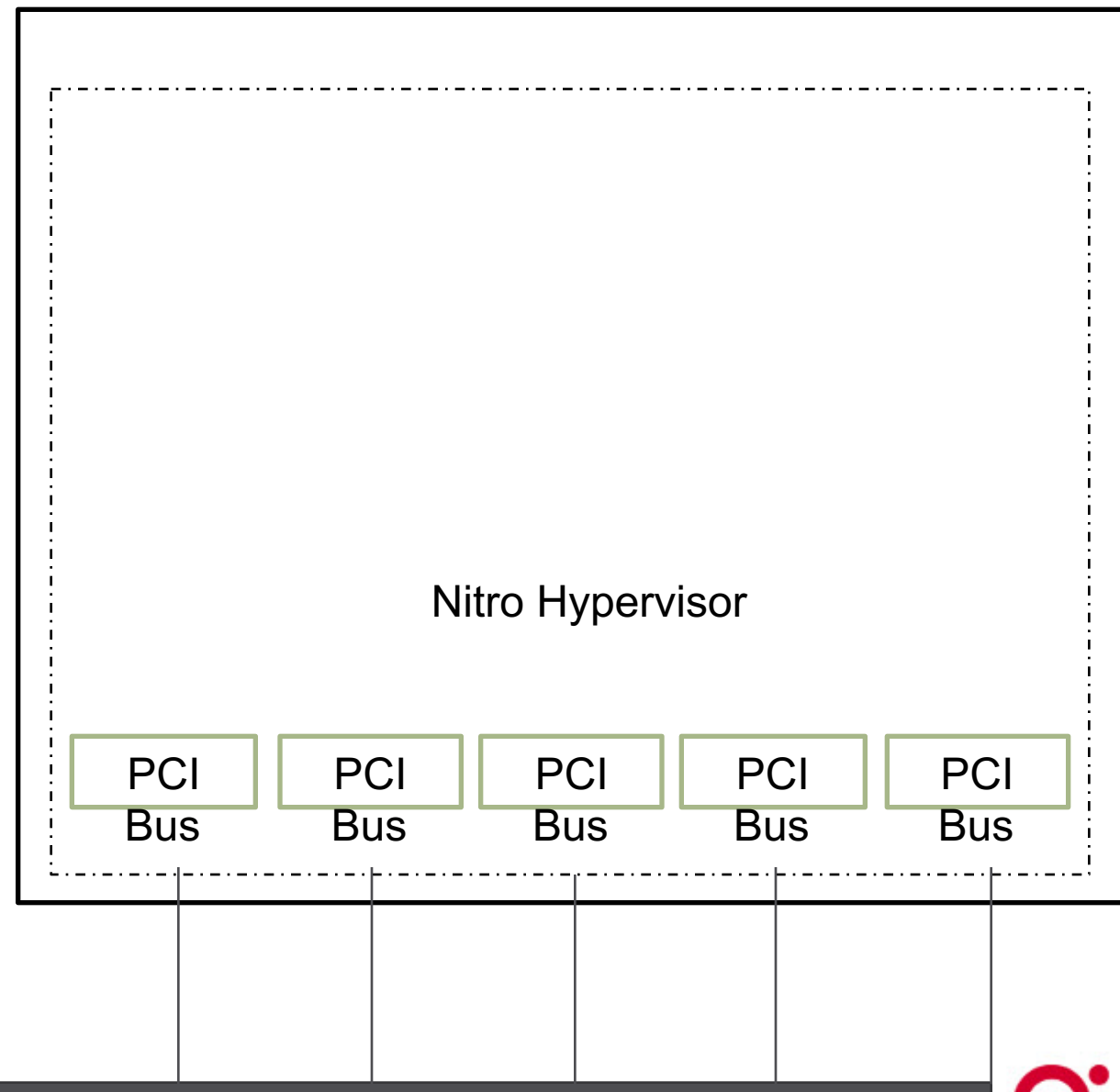The C5 instance was the first instance to use the Nitro hypervisor

Networking, storage, and encryption duties offloaded to custom hardware chipsets

Custom Hardware

Nitro Hypervisor

| PCI Bus | PCI Bus | PCI Bus | PCI Bus | PCI Bus |

EBS  Instance Storage  Security

C5 Architecture

# Current Generation vs All Generations

"Current generation" means the latest and greatest virtualized choices available

Changing view to "All generations" reveals that the older virtualization choices still exist

Choosing para-virtualization, or older instance types means your applications could run slowly at AWS

Long-term, do the work of upgrading your on-premise operating system versions to the latest versions

Then you can start with the current generation instance types and HVM AMI's

AWS Management Console

All instance types
Micro instances
General purpose
Compute optimized
FPGA instances
GPU graphics
GPU instances
GPU compute
Memory optimized
Storage optimized

Current generation
All generations

Filter by:     All instance types  ▼     Current generation  ▼     **Show/Hide Columns**

**Currently selected:** t2.micro (Variable ECUs, 1 vCPUs, 2.5 GHz, Intel Xeon Family, 1 GiB memory, EBS only)

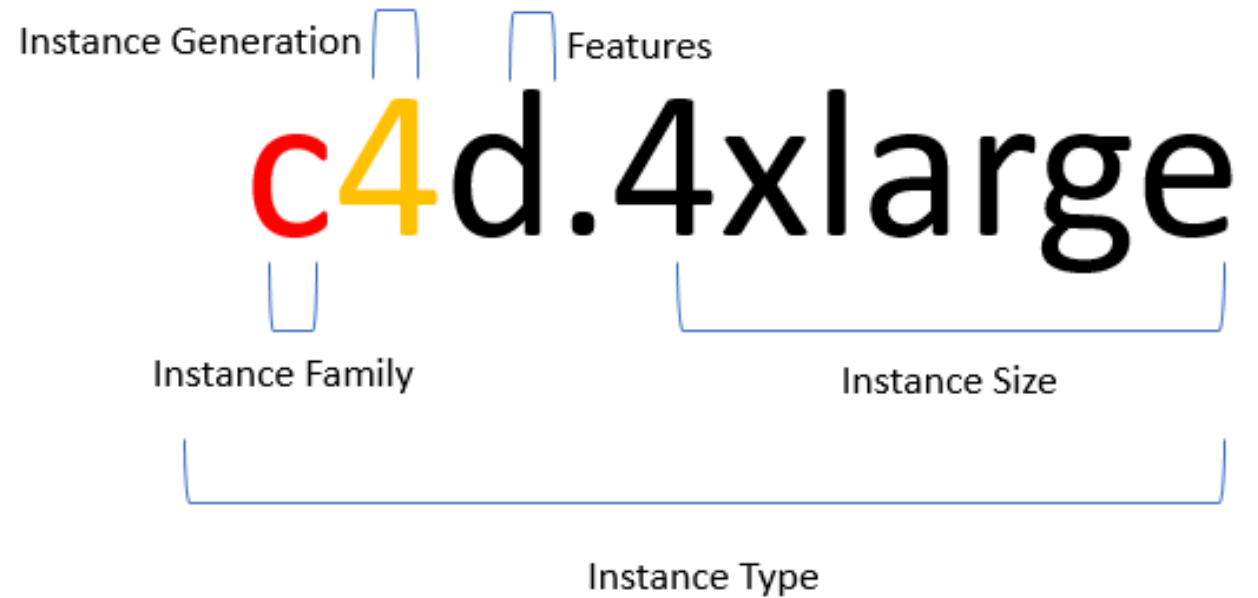| | Family ▼ | Type ▼ | vCPUs ⓘ ▼ | Physical Processor ▼ | Clock Speed ▼ | Memory (GiB) ▼ | Instance Storage (GB) ⓘ ▼ | EBS-Optimized Available ⓘ ▼ | Network Performance ⓘ ▲ |
|---|---|---|---|---|---|---|---|---|---|
| ☐ | Compute optimized | c5.18xlarge | 72 | Intel Xeon Platinum 8124M | 3 GHz | 144 | EBS only | Yes | 25 Gigabit |
| ☐ | Compute optimized | c5d.18xlarge | 72 | Intel Xeon Platinum 8124M | 3 GHz | 144 | 2 x 900 (SSD) | Yes | 25 Gigabit |
| ☐ | Compute optimized | c5n.18xlarge | 72 | Intel Xeon Platinum 8124M | 3 GHz | 192 | EBS only | Yes | 100 Gigabit |
| ☐ | Compute optimized | c5n.2xlarge | 8 | Intel Xeon Platinum 8124M | 3 GHz | 21 | EBS only | Yes | Up to 25 Gigabit |

# EC2 Instance FYI

- Instances are members of compute families

- For each instance's name, the first letter is the instance family that it belongs to

- The letter describes the resources allocated to the instance

- The workloads that the instance has been designed for

- The letter stands for something; for example, the letter C stands for compute, R for RAM and I for IOPS

- The resources (vCPUs, memory, and network bandwidth) are assigned to your account and are never shared with any other AWS customer
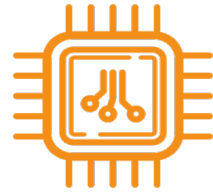
# Demo: Instance Families

# Instance Types

# micro

**Micro Instances** – there's only one instance type in this class; the t1.micro with an unidentified processor
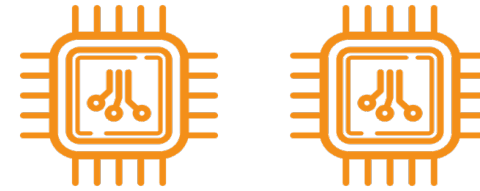
The clock speed is not identified, but you have .613 GiB of memory with very low networking performance

It supports 32 or 64-bit workloads and only shows up in the management console if you search for the All Generation types of instances
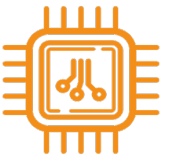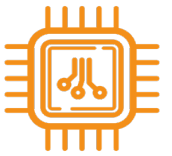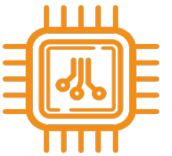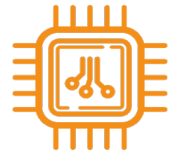
# compute optimized

- Compute optimized instances are designed for batch processing workloads, media transcoding and high-performance application or Web servers

- The C5 architecture takes advantage of the Nitro system components for enhanced networking

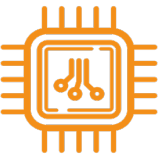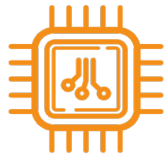| Maximum Size | Storage | AMI |
|---|---|---|
| **Maximum 72 vCPUs, 144 GiB of memory, enhanced networking up to 25 Gbps** | EBS optimized storage with dedicated bandwidth up to 4,000 Mbps | 64-bit HVM AMIs that include drivers for enhanced networking and NVMe storage |

# memory optimized

| Instance Type | Maximum Size | Storage |
|---|---|---|
| **r5** | 96 vCPUs, 769 GiG of memory, enhanced networking speeds up to 25 Gbps | EBS optimized storage with dedicated EBS bandwidth up to 4,000 Mbps |
| **r4** | 16 vCPUs, 488 GiG of memory, enhanced networking speeds up to 25 Gbps | Local instance storage using NVMe SSD |
| **x1** | 128 vCPUs, 1952 GiG of memory, enhanced networking speeds up to 25 Gbps | 14,000 Mbps of EBS optimized storage bandwidth |
| **x1e** | 128 vCPUs, 3904 GiG of memory, enhanced networking speeds up to 10 Gbps | |

Workloads that need to process vast data sets hosted in memory such as MYSQL or NoSQL databases

# storage optimized

| Instance Type | Maximum Size |
|---|---|
| **h1** | 64 vCPUs, 256 GiB memory, 4 x 200 GiB of instance storage, up to 25 Gbps enhanced networking |
| **d2** | 36 vCPUs, 244 GiB memory, 24 x 2048 GiB of instance storage,10 Gbps enhanced networking |
| **i3** | 36 vCPUs, 244 GiB memory, 8 X 1900 GiB of instance storage, up to 25 Gbps enhanced networking |

Storage optimized instance are designed for workloads that require local storage for large data sets

# t instances

When you launch a T2 or T3 instance, depending on the size, you will get a baseline of CPU performance

The use case for these instances could include applications where the CPU processing time is infrequent

t instances are designed with the ability to burst above an initial CPU baseline of performance

# Burst credits

The design of a t instance provides you with CPU credits for the time that your CPU is idle

Banking your CPU credits allows you to use them, when your application needs to burst above the baseline that has been assigned to your instance

The typical server doesn't run flat out at 100%; instead its has peaks and valleys in its performance needs

When performance is needed; banked CPU credits are first used

# Burst credits in operation

At launch, there are enough CPU credits allocated to carry out the initial tasks of booting the operating system and running the application

A single CPU credit has the performance of one full CPU core running at 100 % for one minute

After a T2 instance is powered on it earns CPU credits at a defined steady rate; the larger the instance the more CPU credits are earned up to a defined maximum value.

Earned credits expire after 24 hours; if the CPU credits were not used, then, they weren't needed by the running application

Demo: t instances

# Amazon Machine Images

- The precise definition of an AMI is a template that contains the desired software configuration for an instance:
  - Operating system
  - Optionally an application
  - Additional supporting software
  - Root device boot volume

- After selecting an AMI, you then choose the instance type where your selected AMI will be installed

Each AMI contains the necessary technical information required to launch an instance

You must use an AMI to launch an AWS instance

There are two AMI storage options to consider:

[1] An EBS backed AMI providing a persistent block storage boot volume

[2] A local instance store backed AMI, which provides temporary local block storage

# AMI FYI

# AMI Components

**Boot Volume** – describes what will be used as the root boot volume for the instance – either an EBS boot volume, or a local instance storage volume

**Launch permissions** – define the AWS accounts that are permitted to use the AMI to launch the instances

**Volumes to attach** – the volumes to attach to the instance at launch are contained in a block device mapping document

**Default location** – AMI's are region specific; can be manually copied

**Operating system** – Linux or Windows

Demo:  AMI Creation

# EC2 Instances: Pricing Options

**On-Demand**

Pay by per hour/ second

Short-term, unpredictable workloads

**Reserved Instances**
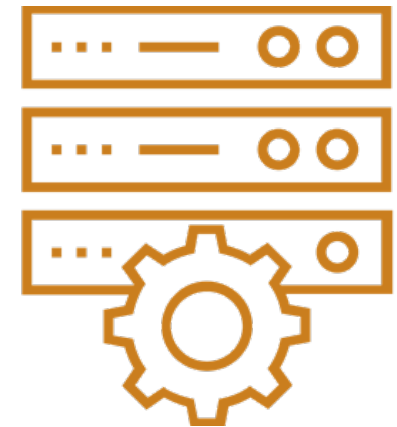
Discount for 1 - 3-year commitment

Applications with consistent usage

**Spot Requests**

Spare AWS capacity > 90% discount

Applications with flexible start and end times

**Dedicated Hosts**

Physical server dedicated to customer

Compliance requirements for applications

# on-demand instances

There are over 160 instance types to choose

Pricing is per second for Linux instances

Pricing is per minutes for Windows instances

The number of on-demand instances each AWS account can launch in each region are bound by a soft limit

Soft limits can be increased upon request

# Reserved instance pricing (ri)

A reserved instance is a billing discount applied to the on-demand instances currently being used, or that will be used

RI pricing is applied when you launch a new instance with the same specifications as your reserved instance pricing

The number of regional reserved instances that you can purchase depends on your current soft limit

Reserved instances have limits based on the region, and the number of availability zones within the region itself

# Reserved instance: Standard reservation

Provides the biggest discount and can be purchased as repeatable one-year terms, or a three-year term.

After purchasing a standard reserved reservation you can make some changes within the reservation:

- Availability Zone
- Instance size
- Networking type

# Reserved instance: Convertible reservation

Change instance types, operating systems, or switch from multi-tenancy to single tenancy compute operation

The convertible reserved discount could be over 50%; and the term can be a one, or three-year term

You can also request reserved EC2 capacity if you need to guarantee that on-demand instances are always available for use in a specific availability zone

After you've created a capacity reservation, you will be charged for the capacity reservation whether you use the instances or not

Without a capacity reservation, there is no guarantee that instances will be available when you need them

# On-demand Capacity Reservations

# Regional Pricing Compared

# Demo: Simple Monthly Calculator

# Spot Requests

A spot instance is spare compute capacity that AWS is not currently using

Save up to 90% of the purchase price; when AWS wants your instance back, a two-minute warning is provided  (CloudWatch alert)

Spot instance pricing is based on supply and demand; the instance will run until another AWS client offers a higher spot price for the same type of spot instance

To counteract this possibility, you can define a maximum spot price that you're willing to pay

# Spot instances

Spot instances can also be hibernated or stopped when it is interrupted

When spot instances are hibernated, the RAM contents is stored on the root EBS drive and your private IP address is held

Choose a spot instance price based on a guaranteed term of 1 to 6 hours

Spot Pool - a current number of unused EC2 instances of the same instance type

Spot Fleet – number of spot instances that are launched based on type and price

# EBS

# Elastic Block Storage

**SSD**

- General purpose:
  Boot volumes, low latency
  applications
- Provisioned IOPS: Databases
  with sustained IOPS

**HDD**

- Throughput optimized:
  High-throughput sequential
  workloads
- Cold: Logging and minimal
  needs

# EC2 Instance Stores

- Local disks attached to the bare metal server that hosts your instance(s)

- Called "Ephemeral storage"

- Temporary storage – buffers , cache, etc.

- Up to 24 TB depending on instance type

- Deleted when instance is stopped, or fails

# EBS Volumes

- Instances that use EBS volumes can be stopped and restarted without data loss
  - EBS volumes can be:
  - Root / boot drives
  - Data drives
  - Encrypted

- Replicated with multiple copies within the AZ where the instance is located

# EBS Features

- Persistent data storage
  - Change volume type
  - Change volume size

- Increase or decrease provisioned IOPS

- Designed for 99.999 service availability

# Elastic Block Storage (EBS)

- Single EBS volume attached to one instance
- Multiple EBS volumes can be attached to one instances

- General Purpose SSD – 1 GB to 16 TB
  - ( 3 IOPS per GB) burstable to 10,000 IOPS

- Provisioned IOPS SSD    4GB – to 16 TB
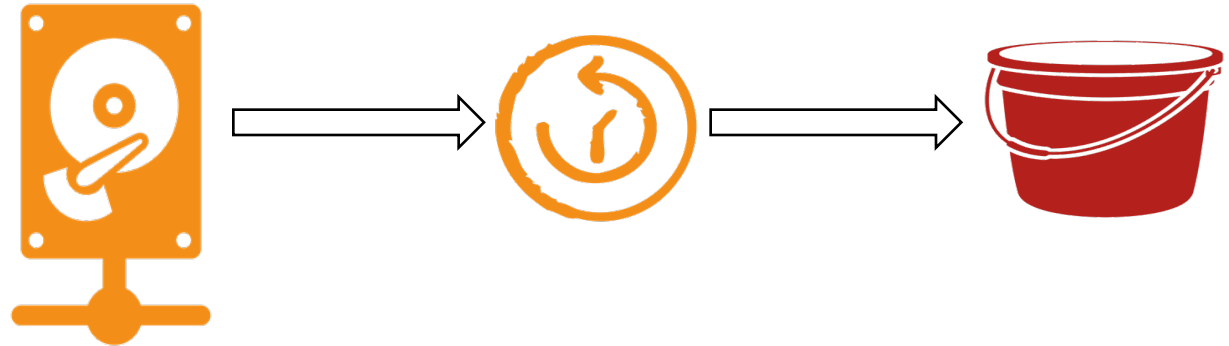  - Minimum  100 IOPS, Max: 64000 IOPS

# Burst Storage Operation

- The baseline for a general purpose SSD is designed with a minimum baseline of 100 to 10,000 IOPS with an average of 3 IOPS per GiB

- Bursting is for a use case where applications have periods of idle time followed by periods of high IOPS

- The smallest gp2 drive can burst to 3000 IOPS, while maintaining a single digit millisecond latency

- As EBS volumes gets larger, your volume also is assigned additional burst credits allowing the drive to burst for a longer time frame

  - 300 GiB volumes can burst up to 40 minutes
  - 500 GiB volumes can burst for almost an hour
  - 900 GiB volumes can burst for almost 10 hours

# Protecting EBS Volumes

- Backup / Recovery snapshots
  - Snapshots are a "Point in time backup"
  - Stored in S3 in AWS "Controlled storage"

- Create a new EBS Volume from an existing snapshot

- EBS volumes can be encrypted – KMS service handles key management

Demo:  EBS Administration

# S3 Storage

# What is S3 Storage ?

- Simple Storage Service
  - Secure, durable and scalable

- Object Storage – Cloud object storage
  - Pay only for the storage you use
  - Each object contains data and metadata

- Accessed over the Internet

- Private endpoint from a subnet hosted in a VPC

- Data is managed as an object using API calls and HTTP verbs (PUT,GET)

- Native interface to S3 using a Restful API (HTTP or HTTPS methods)

# S3 Buckets

- Objects are stored in containers called buckets
  - Buckets are top-level management components

- Bucket names are global, must be unique across all AWS accounts

- Each object is identified, and accessed using a specified unique key

- Each bucket can be divided into folders (delimiters) \
  - Each bucket can hold an unlimited number of objects
  - You can't mount a bucket, install software, host a database

- Highly durable, scalable object store optimized for Reads

# S3 FAQ

- S3 can store any type of data
  - Up to 5 TB max for single object

- Each object has a unique key
  - Key = filename
  - Must be unique within each bucket
  - Multi-part upload for objects greater than 5 GB
  - Bucket contents can be copied to buckets in other regions (additional costs)

- Metadata describes the data
  - System metadata – AWS   date, size, content-type
  - User metadata – tags specified only at the time the object is created

# S3 Storage Classes

- S3 Standard – no minimum storage time

- S3 Intelligent-tiering – monitor and move after 30 days

- S3 Standard-1A – min 30 days

- S3 One Zone-1A – one AZ – min 30 days

- S3 Glacier – min 90 days

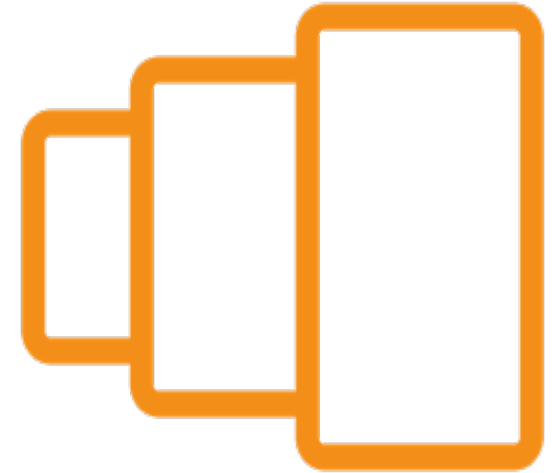- S3 Glacier Deep Archive – min 180 days

Demo: S3 Storage

# S3 Management

# S3 Versioning

- Versioning is enabled at the bucket level

- Versioning allows you to store multiple versions of the same object in one bucket

- Protect yourself from unintended overwrites or deletions

- Once enabled, versioning can't be disabled but can be suspended

# Lifecycle Rules

- Rules defines an action for S3 to apply to a selected group of stored objects

- Rules control the retention of objects
  - Change storage tier, archive, or delete
  - Stored logs: delete after 90 days
  - Documents less frequently accessed: archive to S3 Glacier
  - Delete objects not required after certain date

# Demo: Lifecycle Rules

# S3 Durability

Stored in multiple devices in multiple facilities, within a region

- Designed to sustain concurrent loss of two facilities without loss of data

- Standard
  - 11 9's durability
  - 4 9's availability
  - Over a given year

- Standard 1-A
  - 4 9's durability

# S3 Consistency

- Objects are eventually consistent

- Multiple copies means replicated storage

- PUT's to new objects – read after write consistency

- PUT's to existing object – eventual consistency

# Cross Region Replication

- Asynchronous replication from source bucket in region to bucket in same or another region
- Helps move data closer to end-users
- Compliance / additional durability

# Access Control

- Only owner has access by default
  - Private by default

- Coarse grained – S3 ACL
  - Read / Write / Full Control at object level

- Fine-grained – bucket policies
  - Associated with the bucket / not an IAM security principal
  - Can specify access from where, who can access, and what time of day

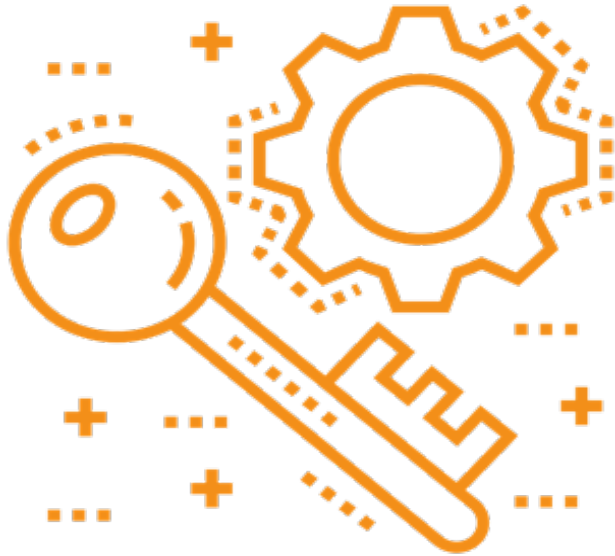- Can be associated with different AWS accounts

# S3 Encryption

- SSE – S3 (AWS Managed Keys)
  - AWS rotates encryption keys
  - Data, encryption, and master keys are stored on separate hosts

- SSE - S3 (AWS KMS Keys) Customer Managed
  - Manage permissions for the master key
  - CloudTrail auditing; view failed attempts

# Key Management Service

**AWS KMS** – Managed service allow you to generate, store, enable / disable and delete symmetric keys

- Customer managed keys – Each CMS is per customer and is used to encrypt and decrypt data
- Data keys – Used to encrypt data objects within data storage

**AWS Cloud HSM** – Secure your cryptographic keys using Hardware Security Modules

- Recommendation is to use two HSM's configured in a highly available configuration

# S3 Notifications

S3 server-access logs track requests to S3 bucket
- Account name and IP address
- Bucket name
- Request time
- Action ( GET PUT LIST)
- Response or error code

Event Notifications
- Response to objects uploaded to S3
- Monitored at the bucket level

Object creation, removal triggers response
- Simple notification service, Simple queue service, transcoding, Lambda

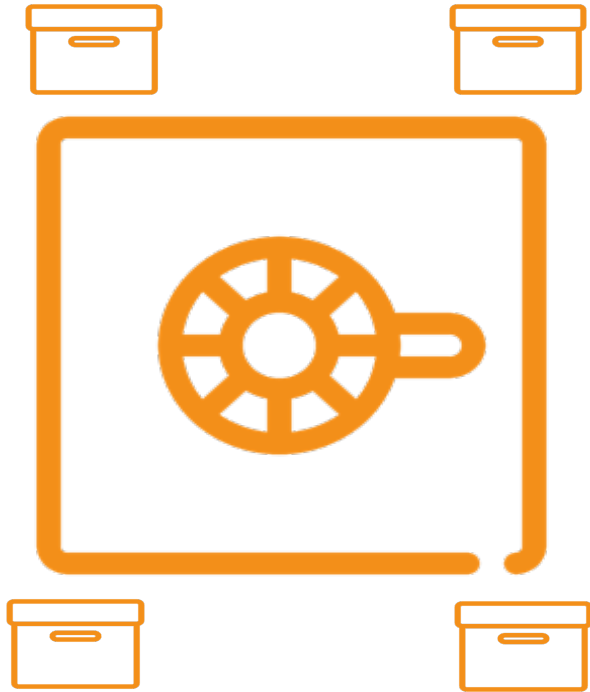# Demo: S3 Notifications

# S3 Glacier

# S3 Glacier Storage

- Low cost archival storage

- Data is stored in archives

- Unlimited # of archives

- 40 TB archive size

- Glacier – Encrypted by default

# S3 Glacier Vaults

- Archives are held in containers called vaults

- Each AWS account can have up to 1,000 vaults

- Compliance controls per vault with a vault lock policy    (WORM)

- Retrieval policy to control data access

# Demo: S3 Glacier

# What we covered:

- Fundamentals of AWS: architecture, terminology and concepts

- Virtual Private Cloud (VPC): Networking services

- Elastic Compute Cloud (EC2): Instance deployment and configuration

- Storage solutions: Elastic Block Storage (EBS) and snapshot management

- Simple Storage Service (S3): Object storage

- S3 Glacier: Archive storage