

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



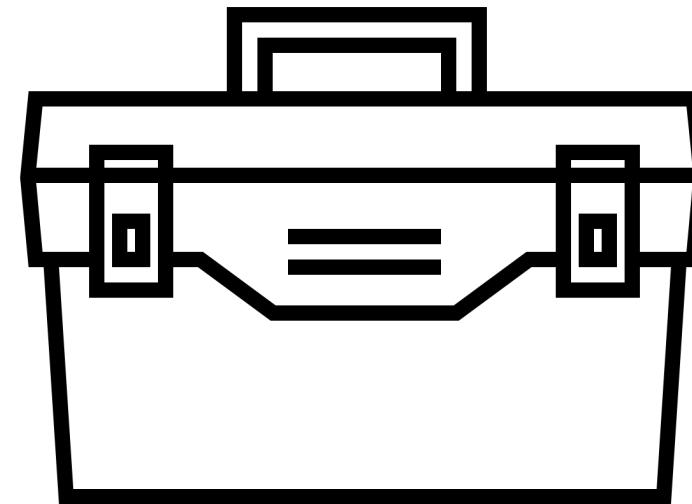
AWS Core Cloud Services

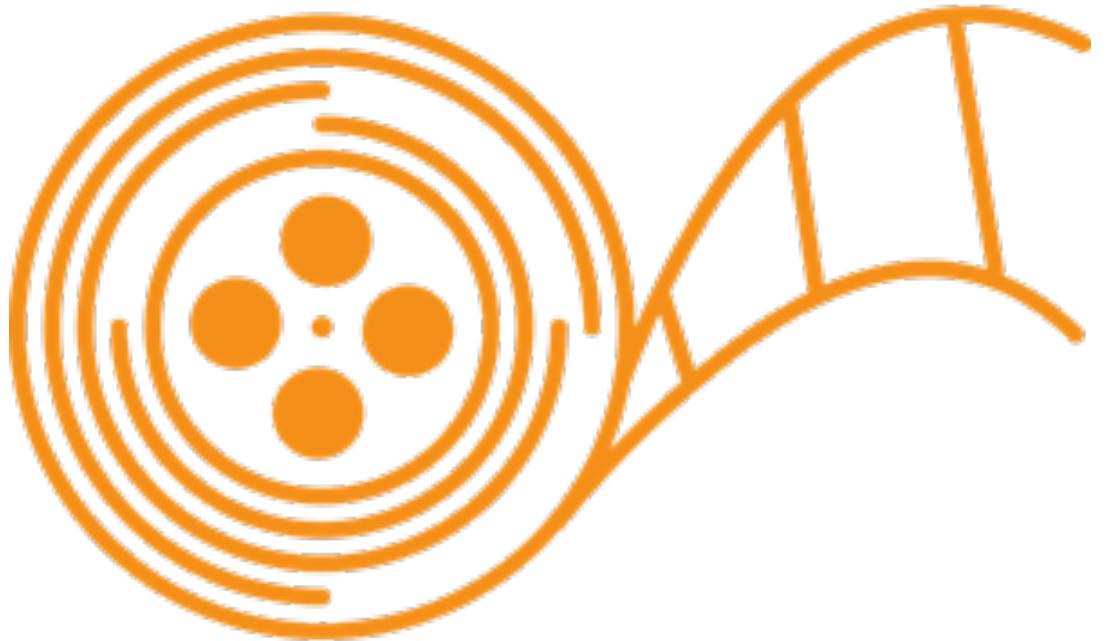
- AWS Administration – Management portal
- Compute Services – Elastic compute cloud
- Networking Services – Virtual private cloud
- Auto Scaling – Scale EC2 compute automatically
- Elastic Load Balancing – Distribute traffic across EC2 instances or containers
- Elastic Block Storage – Virtual hard drives
- 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 do more of the setup
- You do most of the setup, management, and monitoring (VPC, EC2)
- There are no completely unmanaged services at AWS





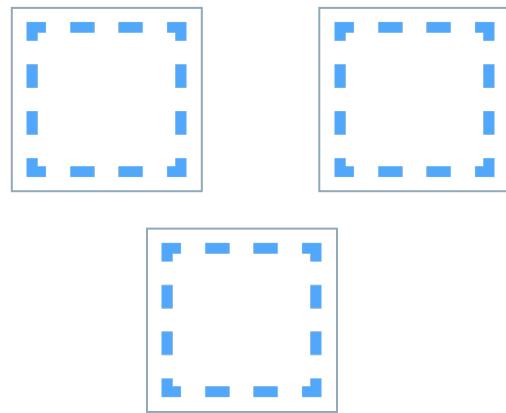
Demo:
Management
Services

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

Regions and Availability Zones



AWS Regions

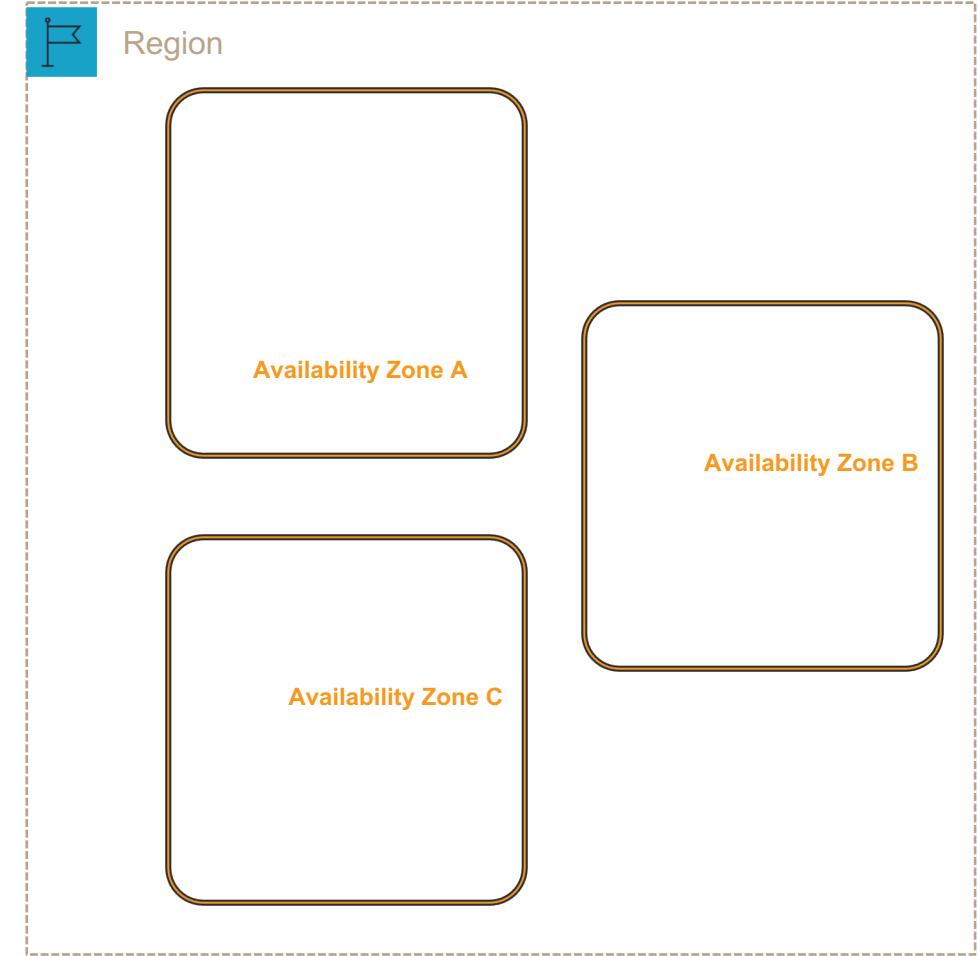
- Areas of the world where Amazon offers AWS cloud services
- Each region is a geographical location
 - Where do you operate?
 - Where your customers?
 - Where are you allowed to operate?
- Each region is completely independent and isolated
- Pricing differences depending on geographical location
- AWS services are not typically replicated across regions unless customers choose to do the setup and replication
- Traffic sent across AWS regions faces additional charges for ingress and egress traffic



Availability Zones

Availability Zones

- Each availability zone contain at least one data center
- Most availability zones contain multiple data centers
- Each availability zone has inexpensive low latency network connectivity to the other availability zones in the same region
- Designing with two AZ's is best practice



Choosing an AWS Region



Latency to on-premise location



Costs are different per Region



Features are different per Region



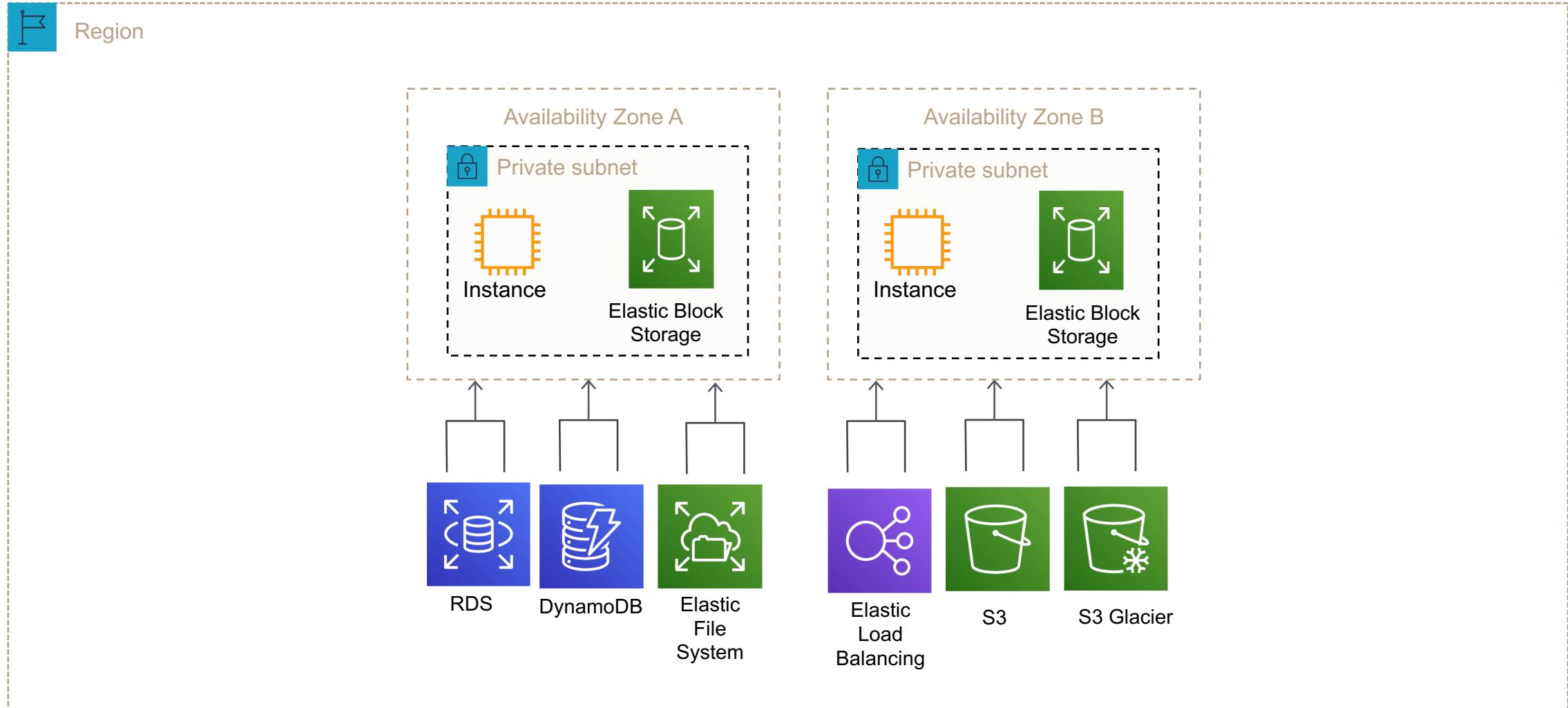
Compliance rules and regulations

AWS Regions in Operation



- Most AWS services are designed to operate within a single AWS region
- Some resources (Snapshots, S3 buckets, and Amazon Machine Images) can be copied from one region to another
- RDS Aurora can also operate multi-AZ, and multi-region
- DynamoDB operates across multiple availability zones, and optionally using global tables that span multiple regions

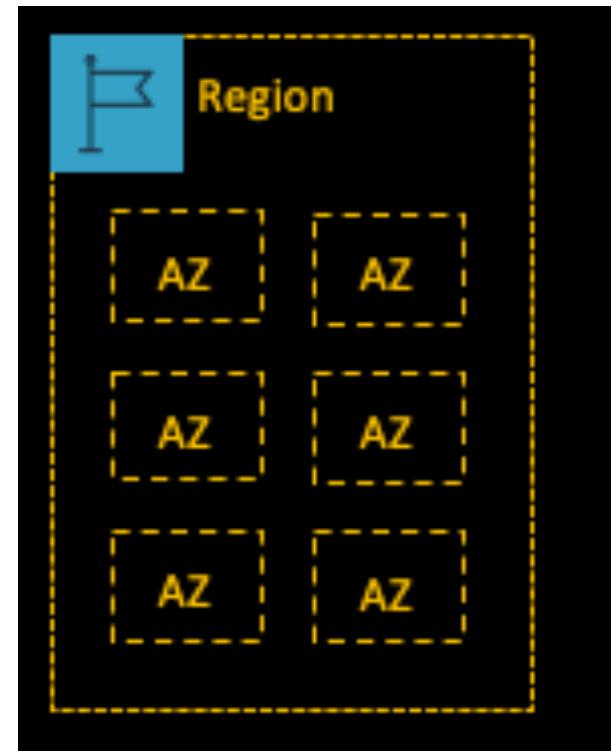
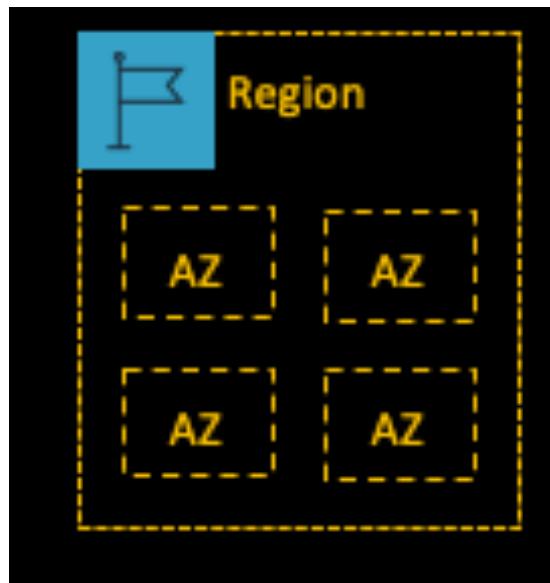
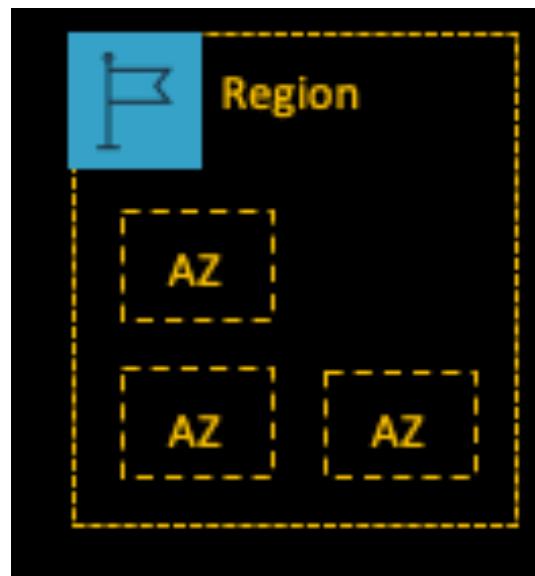
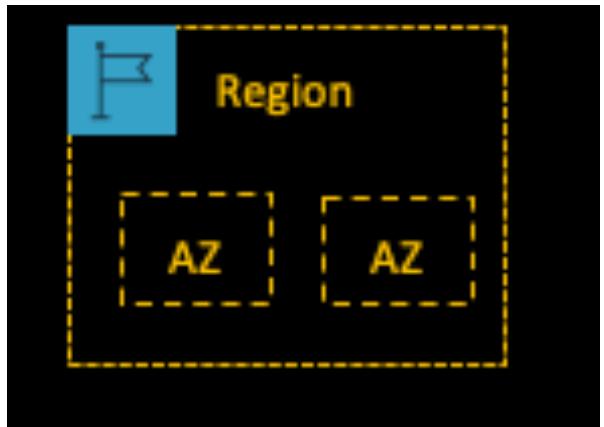
Regional Services



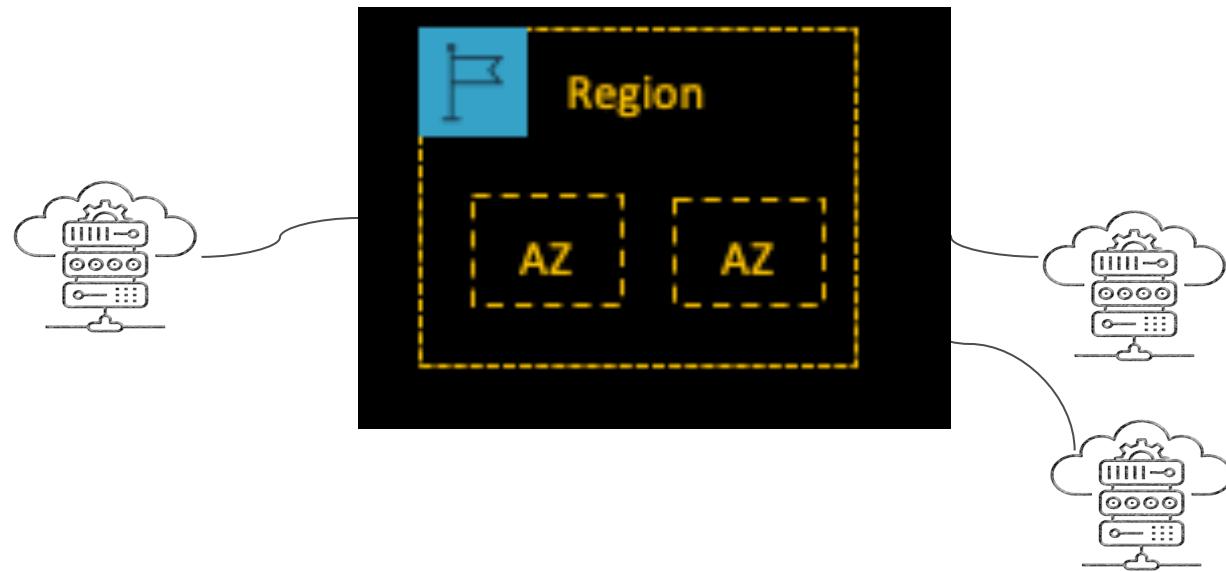


Demo:
Regions

Availability Zones



Availability Zones

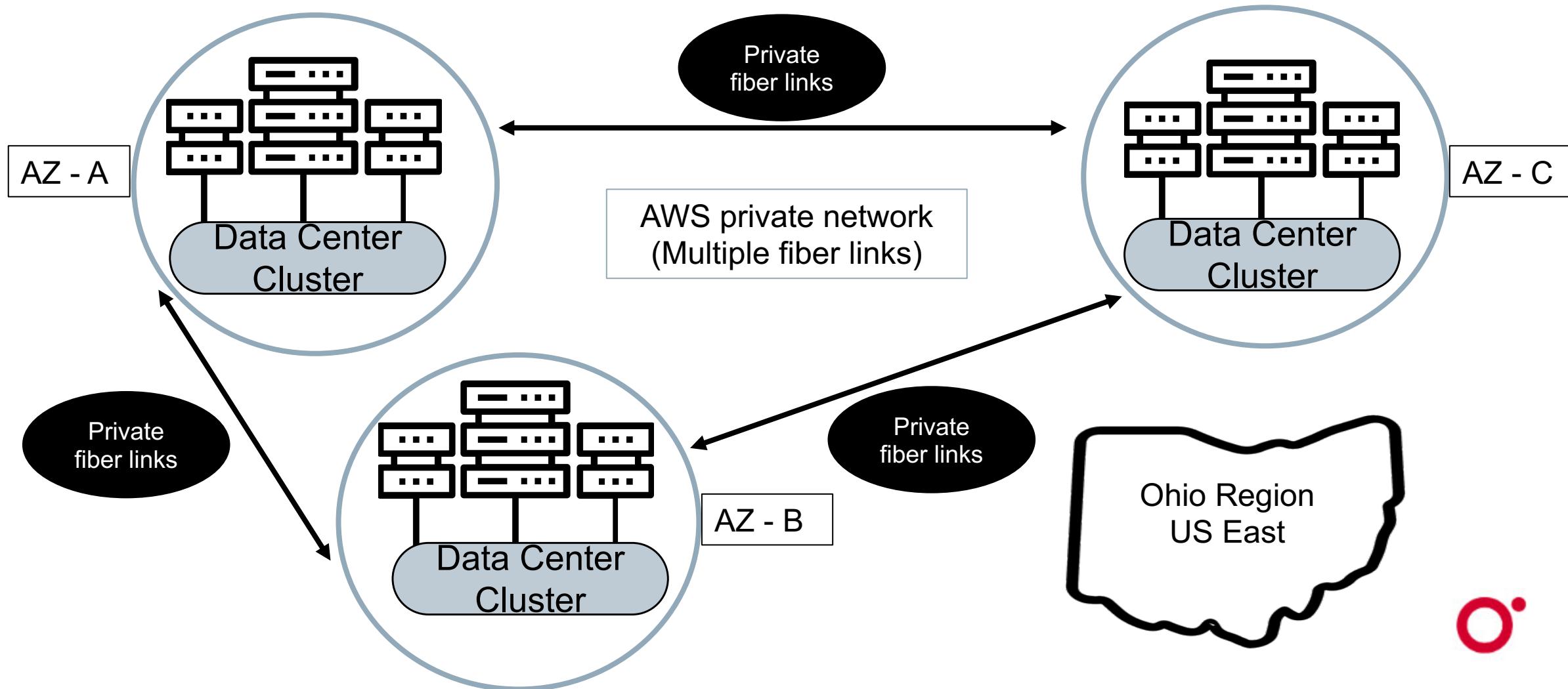


Availability Zones in Operation



- EC2 instances can launch across multiple subnets hosted in multiple availability zones
- ELB can target EC2 instances across multiple availability zones
- EC2 Auto Scaling can scale instances across multiple availability zones
- RDS solutions are replicated across multiple availability zones
- Amazon Aurora can also be multi-AZ
- DynamoDB is replicated across multiple availability zones

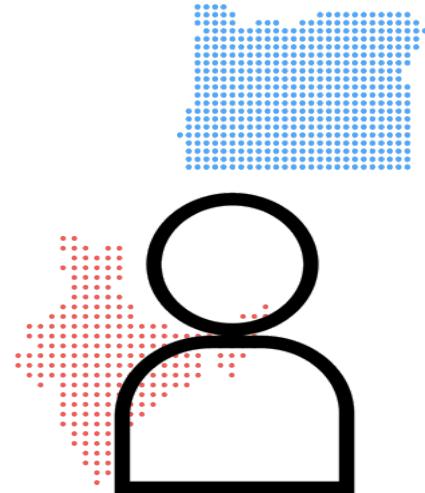
Availability Zones (AZ)



Availability Zones (AZ)



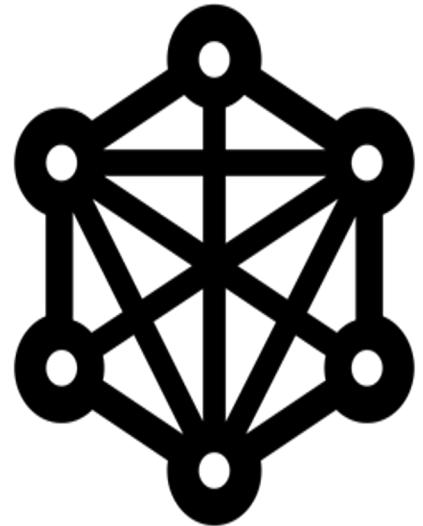
Designed as an independent failure zone: Separate power sources



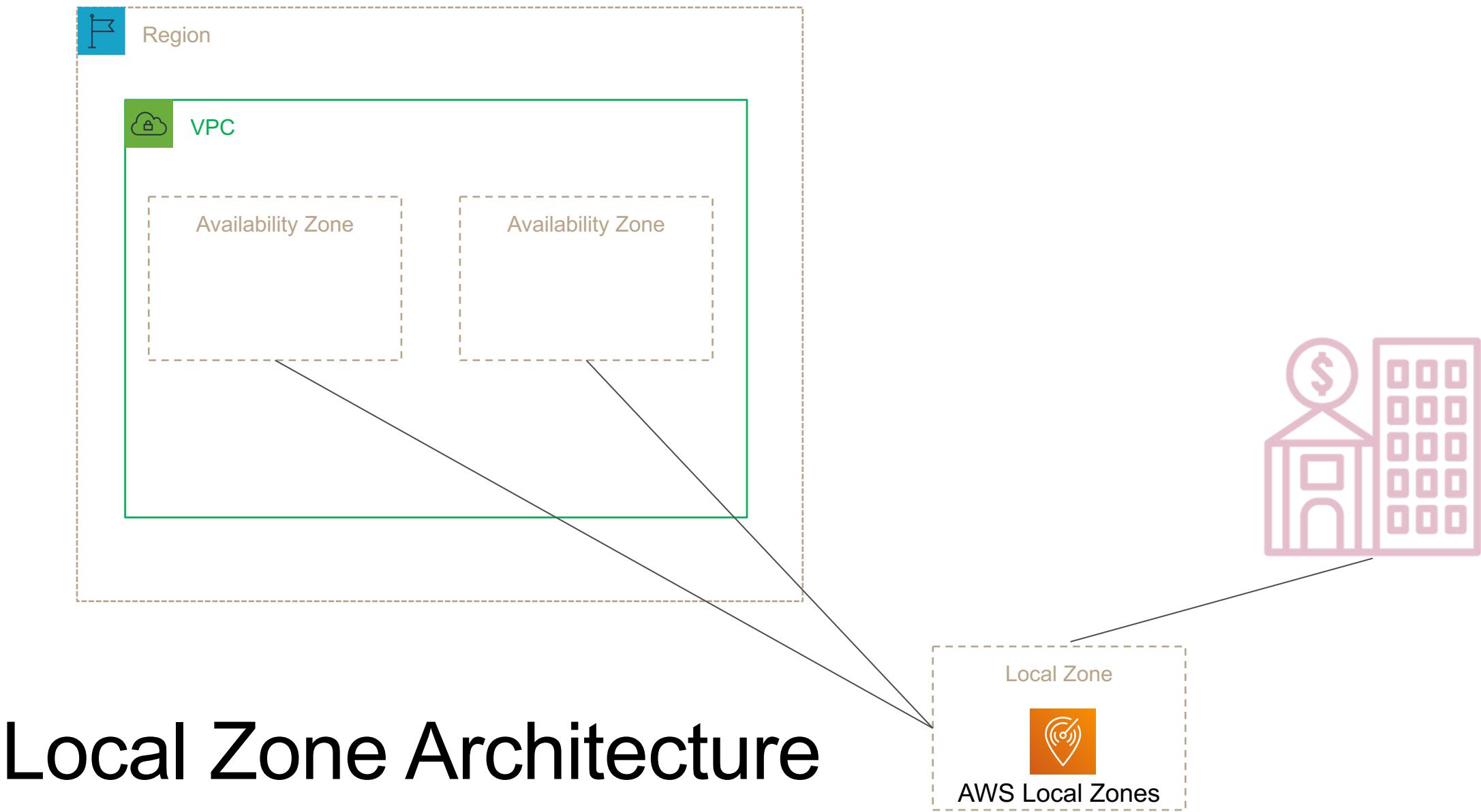
AWS account has access to all regions and associated AZ's



Data transfer charges for outbound traffic



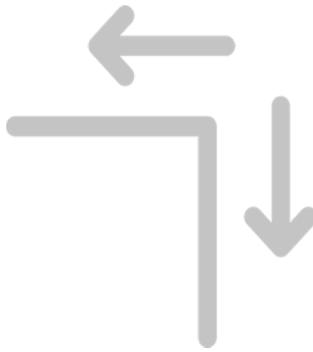
Redundant Tier-1 transit private fiber connections



Local Zones



Enable AWS Local Zones



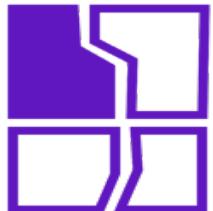
Extend VPC to AWS Local Zone



Build EC2 instances in Local Zone

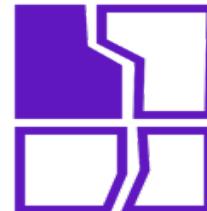
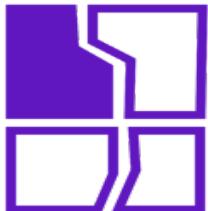
Single Availability Zone

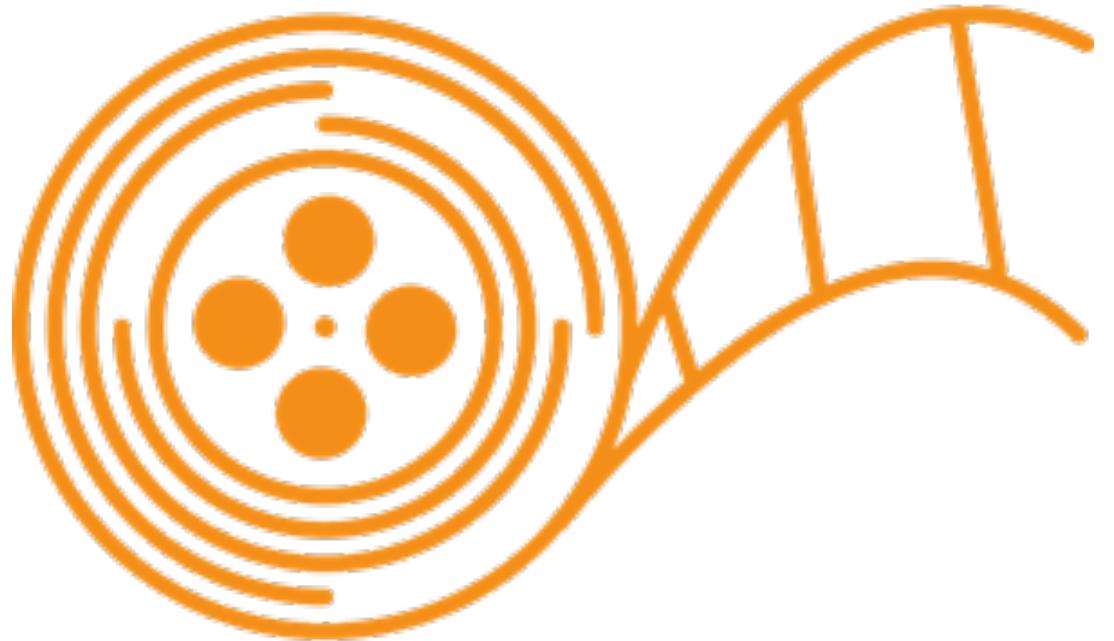
- No recovery or failover when disaster happens in a single datacenter
- No high availability for instances
- No failover in single datacenter
- All AWS regions have at least 2 availability zones
- Each AZ has at least one datacenter



Multiple Availability Zones

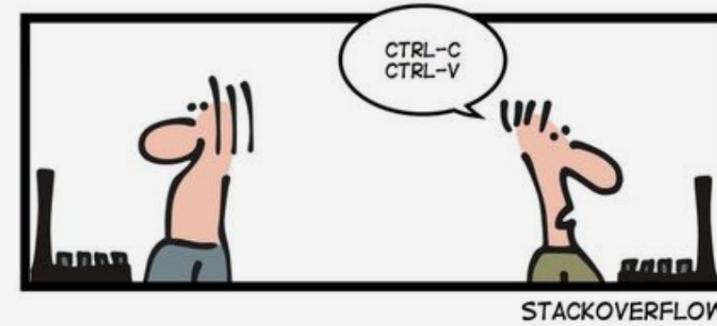
- Better high availability design options
- Designing applications hosted across AZ's provides HA options
- Load balancing (ELB) supports EC2 instances in multiple availability zones)
- EC2 auto scaling supports multiple AZ's
- Route 53 balances resource stacks across multiple AWS regions





Demo:
Availability
Zones

Edge Locations



Edge Locations



Over 200 edge locations



11 Regional Edge caches
in 73 cities

Edge Locations @ AWS



Services at the Edge

Caches your
static and
dynamic content

CloudFront



Delivers your
request from
closest edge
location

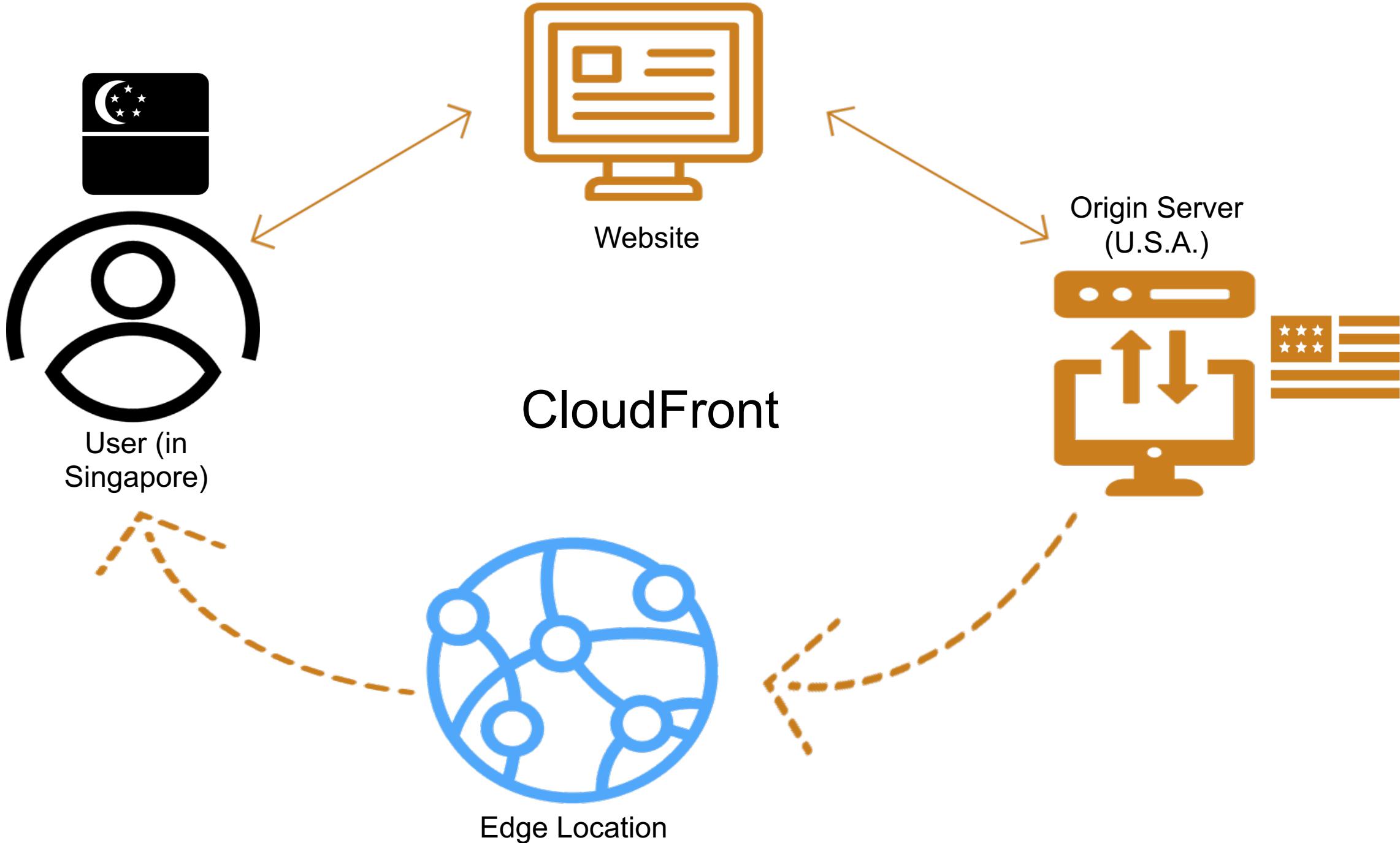
Route 53

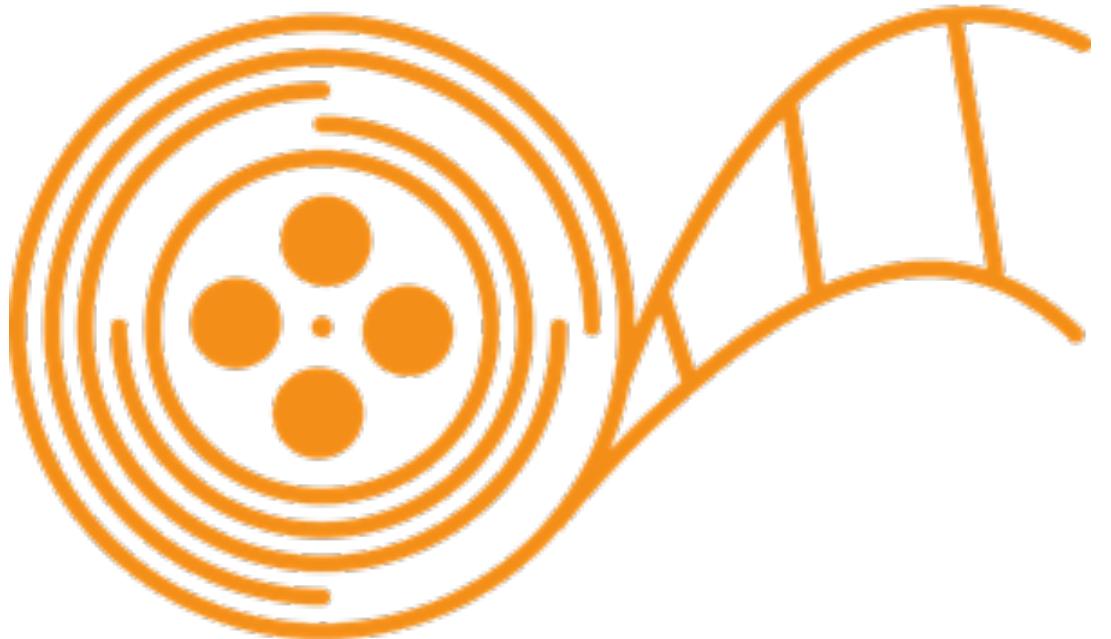


Filters incoming
public traffic at
the edge

WAF







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
- EC2 instances run in a virtual private cloud that is logically isolated to your AWS account

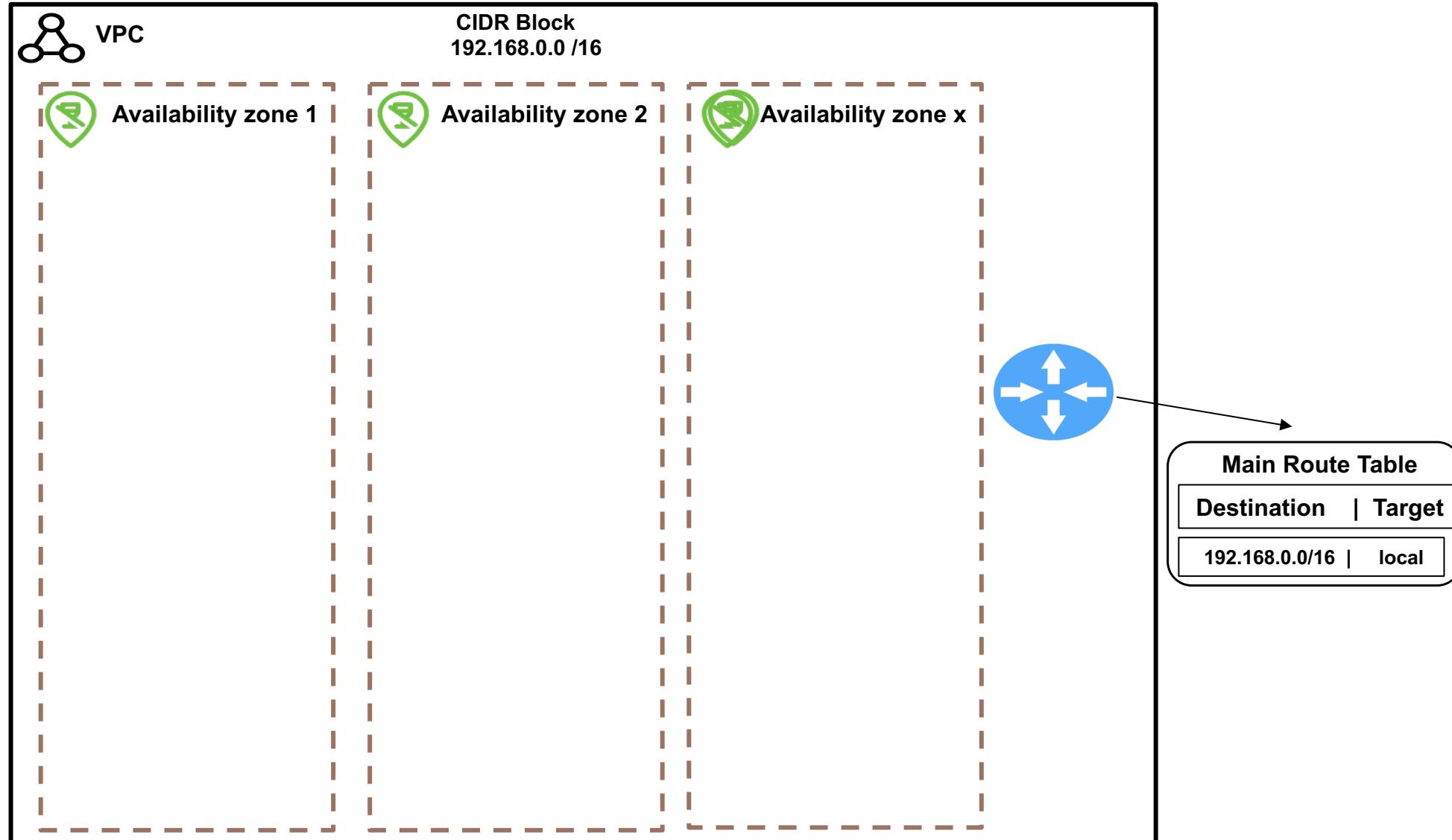


Creating a New VPC

- 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 a default route table that enables local routing throughout the VPC



VPC's have Multiple AZ's

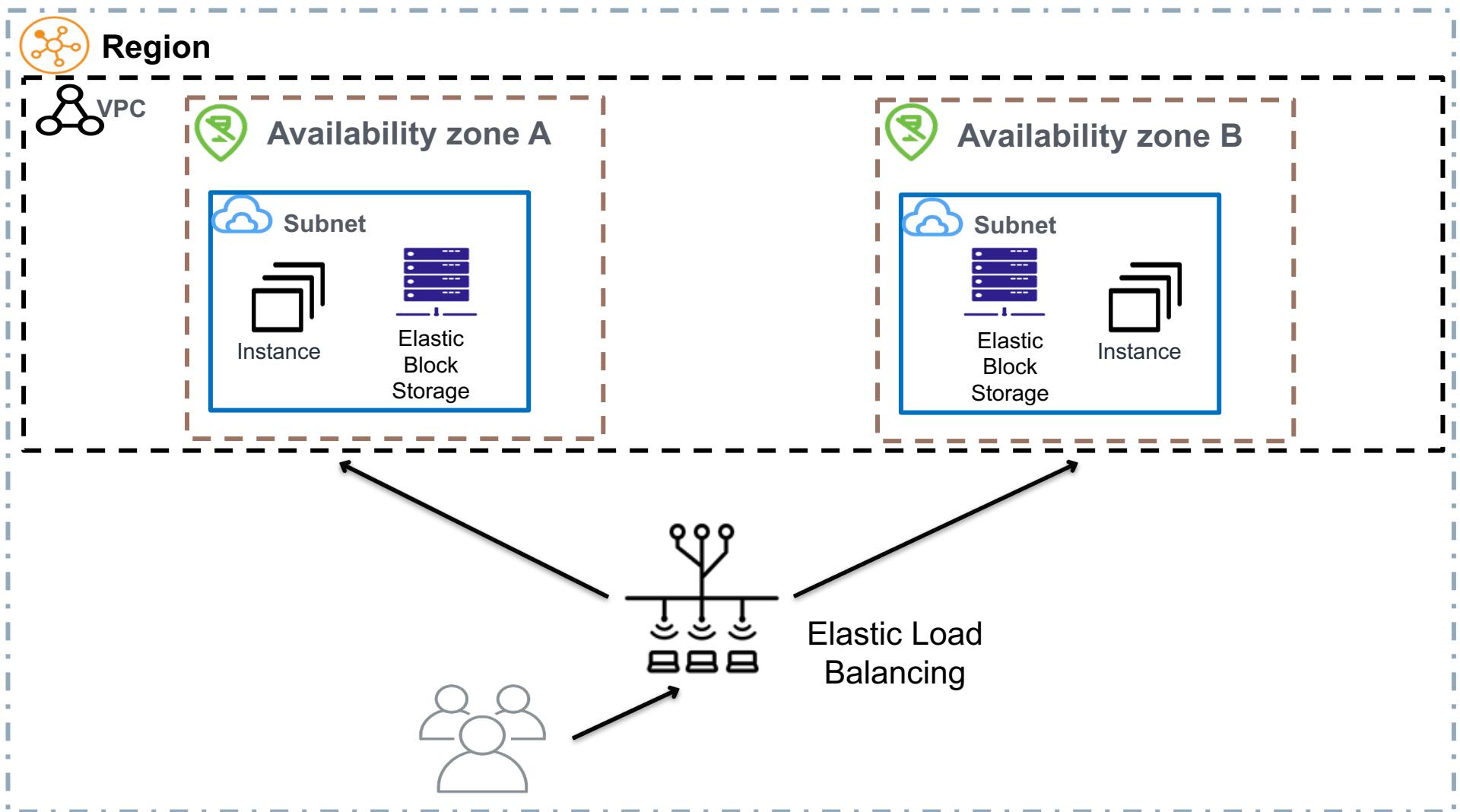


VPC Core Components

- Subnets
- Route tables
- Security groups (SG)
- Network access control list (NACLs)
- Internet gateway (IGW)
- VPN gateway (VGW)
- Private endpoints
- Peering connections
- NAT gateway services
- Transit gateway



Failover Possibilities with AZ's



CIDR Notation for a VPC

| Address (Host or Network) | Netmask (i.e. 24) | Netmask for sub/supernet (optional) |
|--|-------------------------------------|---|
| <input type="text" value="192.168.0.0"/> | / <input type="text" value="16"/> | move to: <input type="text" value="255.255.255.0"/> |
| <input type="button" value="Calculate"/> | <input type="button" value="Help"/> | |

Address: 192.168.0.0 11000000.10101000 .00000000.00000000
Netmask: 255.255.0.0 = 16 11111111.11111111 .00000000.00000000
Wildcard: 0.0.255.255 00000000.00000000 .11111111.11111111
=>
Network: 192.168.0.0/16 11000000.10101000 .00000000.00000000 (Class C)
Broadcast: 192.168.255.255 11000000.10101000 .11111111.11111111
HostMin: 192.168.0.1 11000000.10101000 .00000000.00000001
HostMax: 192.168.255.254 11000000.10101000 .11111111.11111110
Hosts/Net: 65534 (Private Internet)



CIDR Notation for Subnets

Subnets

| | | |
|------------|--------------------|--|
| Netmask: | 255.255.255.0 = 24 | 11111111.11111111.11111111 .00000000 |
| Wildcard: | 0.0.0.255 | 00000000.00000000.00000000 .11111111 |
| Network: | 192.168.0.0/24 | 11000000.10101000.00000000 .00000000 (Class C) |
| Broadcast: | 192.168.0.255 | 11000000.10101000.00000000 .11111111 |
| HostMin: | 192.168.0.1 | 11000000.10101000.00000000 .00000001 |
| HostMax: | 192.168.0.254 | 11000000.10101000.00000000 .11111110 |
| Hosts/Net: | 254 | (Private Internet) |
| Network: | 192.168.1.0/24 | 11000000.10101000.00000001 .00000000 (Class C) |
| Broadcast: | 192.168.1.255 | 11000000.10101000.00000001 .11111111 |
| HostMin: | 192.168.1.1 | 11000000.10101000.00000001 .00000001 |
| HostMax: | 192.168.1.254 | 11000000.10101000.00000001 .11111110 |
| Hosts/Net: | 254 | (Private Internet) |
| Network: | 192.168.2.0/24 | 11000000.10101000.00000010 .00000000 (Class C) |
| Broadcast: | 192.168.2.255 | 11000000.10101000.00000010 .11111111 |
| HostMin: | 192.168.2.1 | 11000000.10101000.00000010 .00000001 |
| HostMax: | 192.168.2.254 | 11000000.10101000.00000010 .11111110 |
| Hosts/Net: | 254 | (Private Internet) |



Supported CIDR ranges

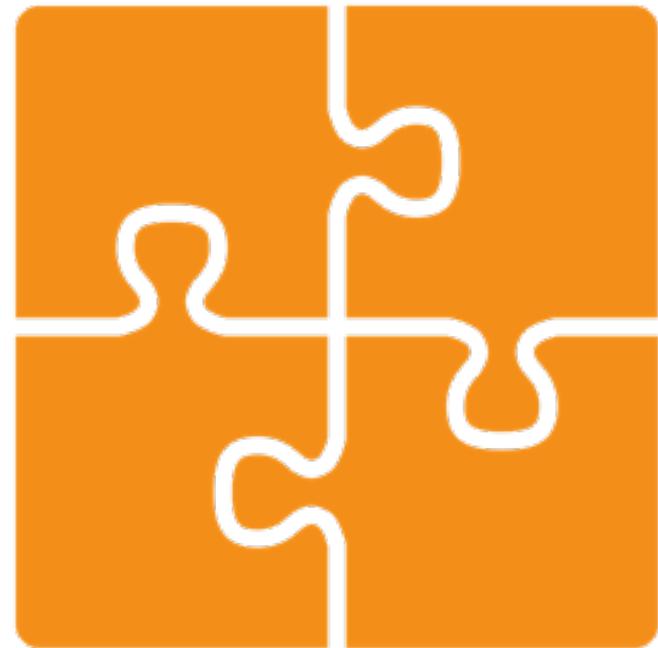
| | Addresses | Hosts | Netmask | Amount of a Class C |
|-----|------------------|--------------|-----------------|----------------------------|
| /30 | 4 | 2 | 255.255.255.252 | 1/64 |
| /29 | 8 | 6 | 255.255.255.248 | 1/32 |
| /28 | 16 | 14 | 255.255.255.240 | 1/16 |
| /27 | 32 | 30 | 255.255.255.224 | 1/8 |
| /26 | 64 | 62 | 255.255.255.192 | 1/4 |
| /25 | 128 | 126 | 255.255.255.128 | 1/2 |
| /24 | 256 | 254 | 255.255.255.0 | 1 |
| /23 | 512 | 510 | 255.255.254.0 | 2 |
| /22 | 1024 | 1022 | 255.255.252.0 | 4 |
| /21 | 2048 | 2046 | 255.255.248.0 | 8 |
| /20 | 4096 | 4094 | 255.255.240.0 | 16 |
| /19 | 8192 | 8190 | 255.255.224.0 | 32 |
| /18 | 16384 | 16382 | 255.255.192.0 | 64 |
| /17 | 32768 | 32766 | 255.255.128.0 | 128 |
| /16 | 65536 | 65534 | 255.255.0.0 | 256 |



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



Subnets

Creating Subnets

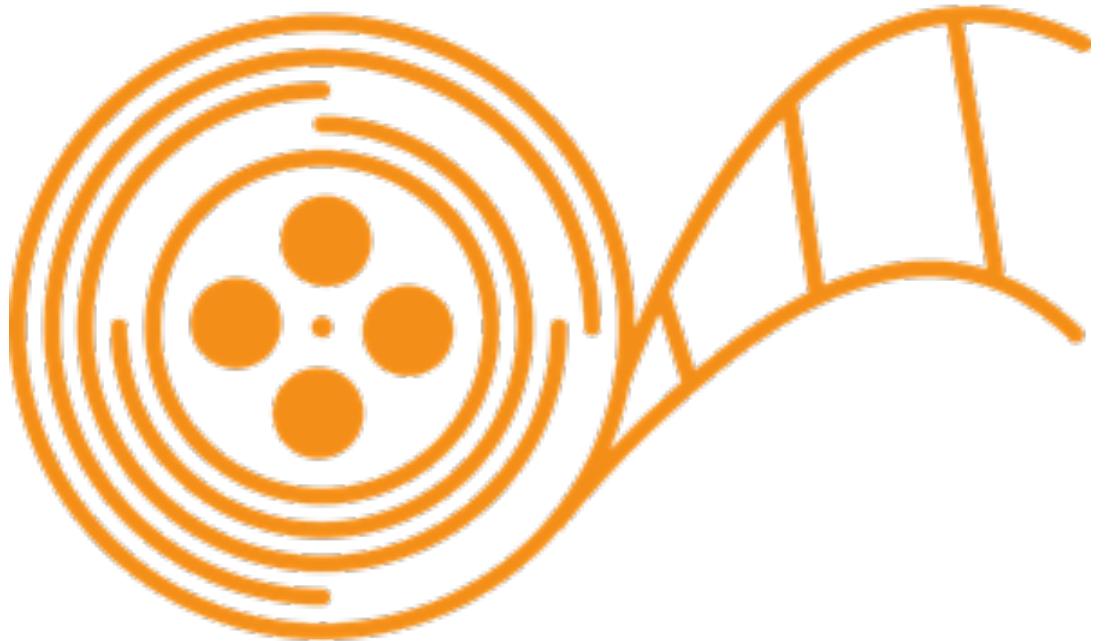
- Public or private subnets can be created in each availability zone
- Subnets cannot span across multiple availability zones
- A subnet that doesn't route to an internet gateway is a private subnet
- If a subnet has traffic routed to an internet gateway it is defined as a public subnet
- EC2 instances in a public subnet must have a public IP address to be able to communicate with the Internet gateway



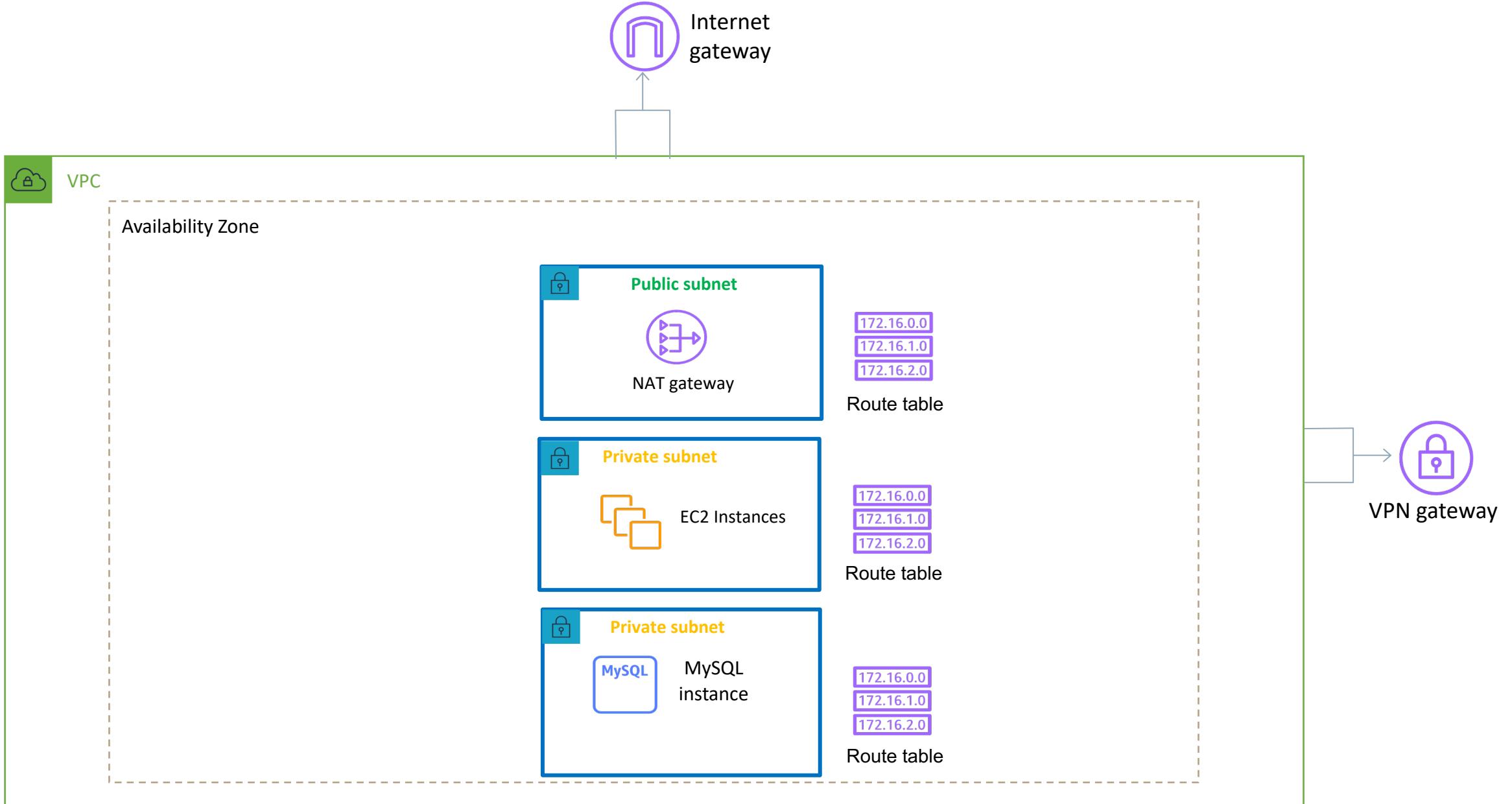
Using Subnets

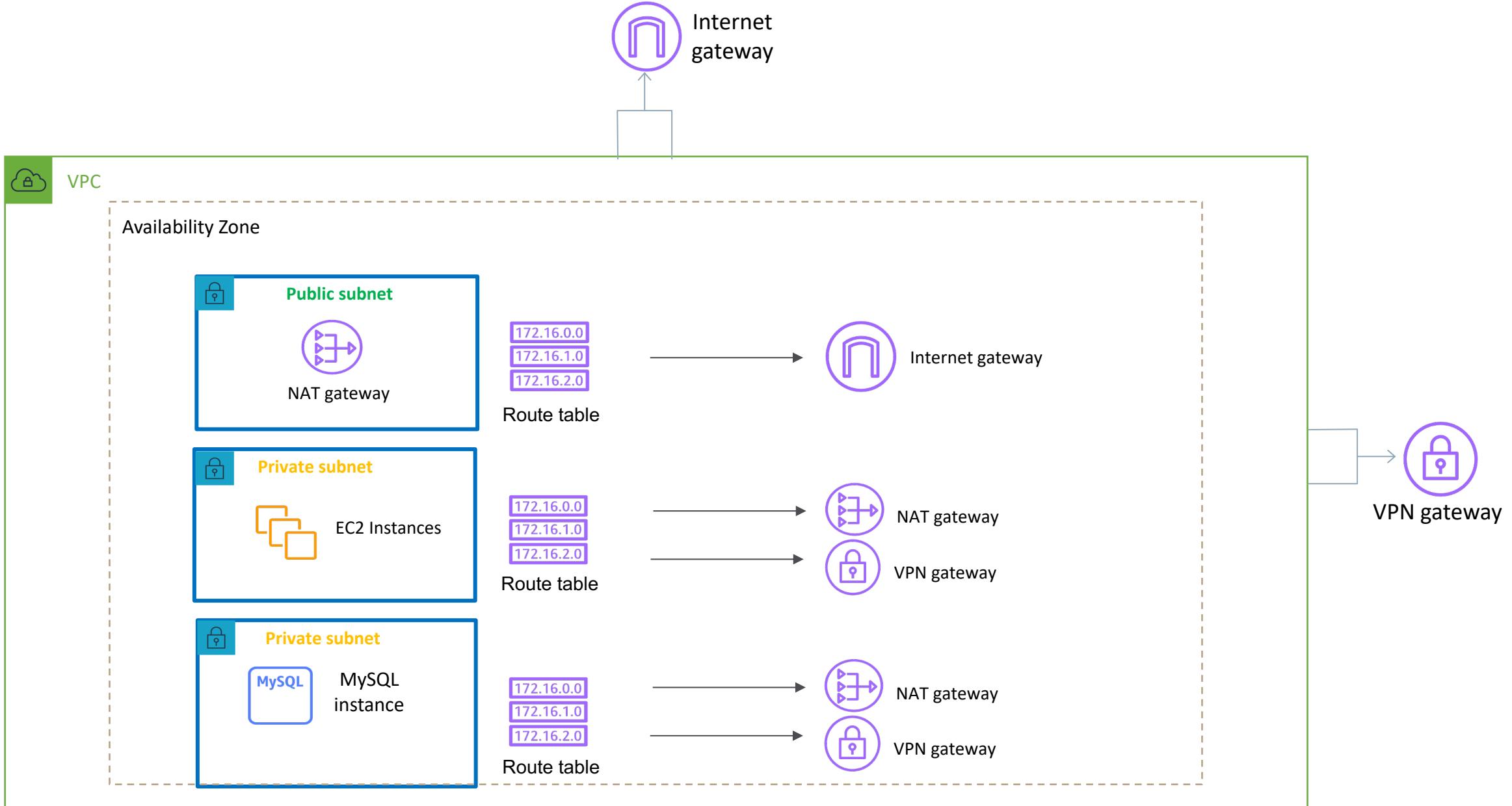
- Instances and AWS services are launched into subnets
- Public subnets can be used for resources that need Internet access (IGW, ELB)
- Private subnets host resources that don't directly connect to the Internet (EC2 Instances, RDS instances)
- Protect subnet access using optional network access control lists (NACLs)
- Network ACL operates at the subnet level and supports allow and deny rules



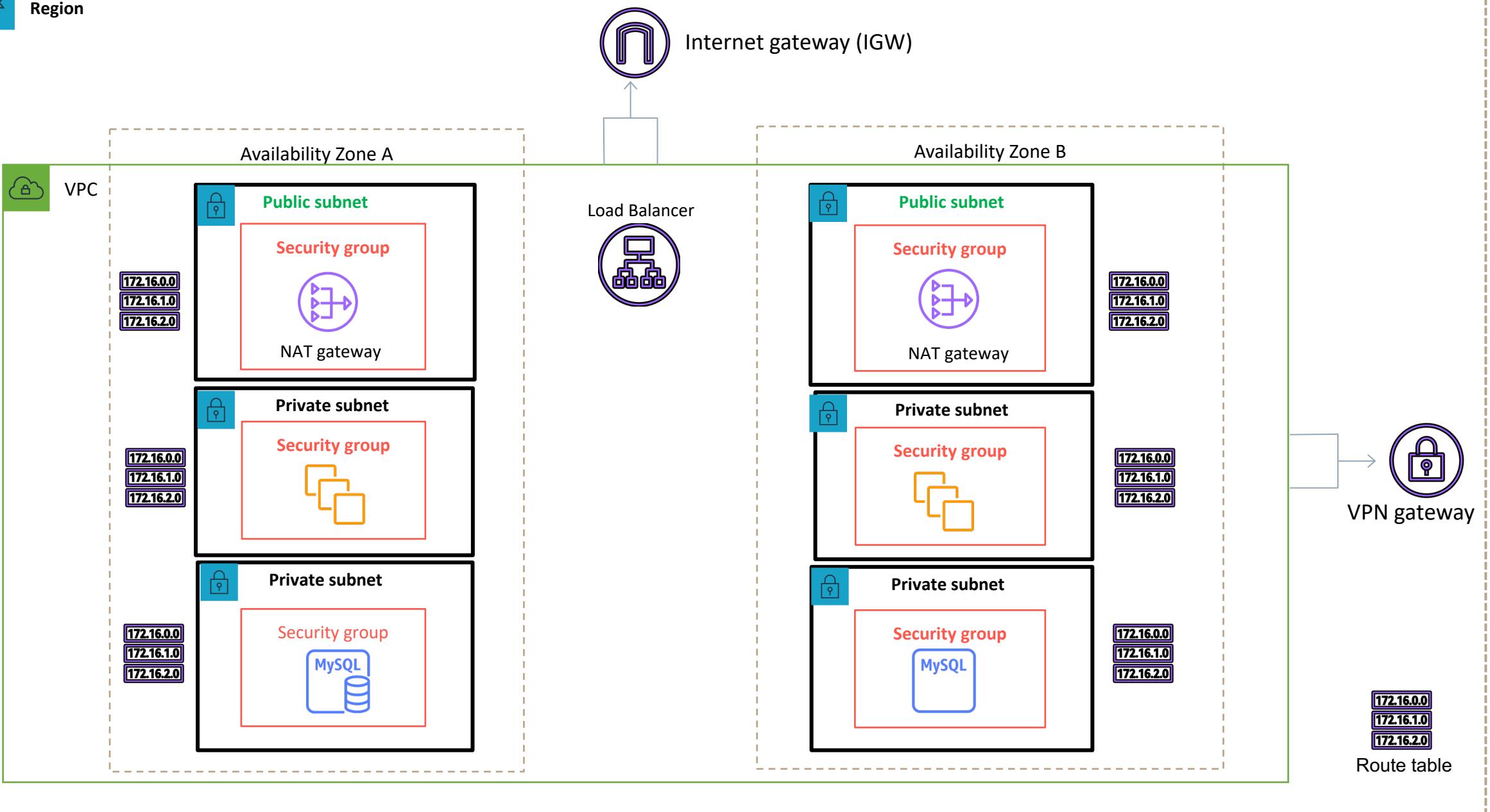


Demo:
Subnets





Region



IP Addresses

IP Addresses

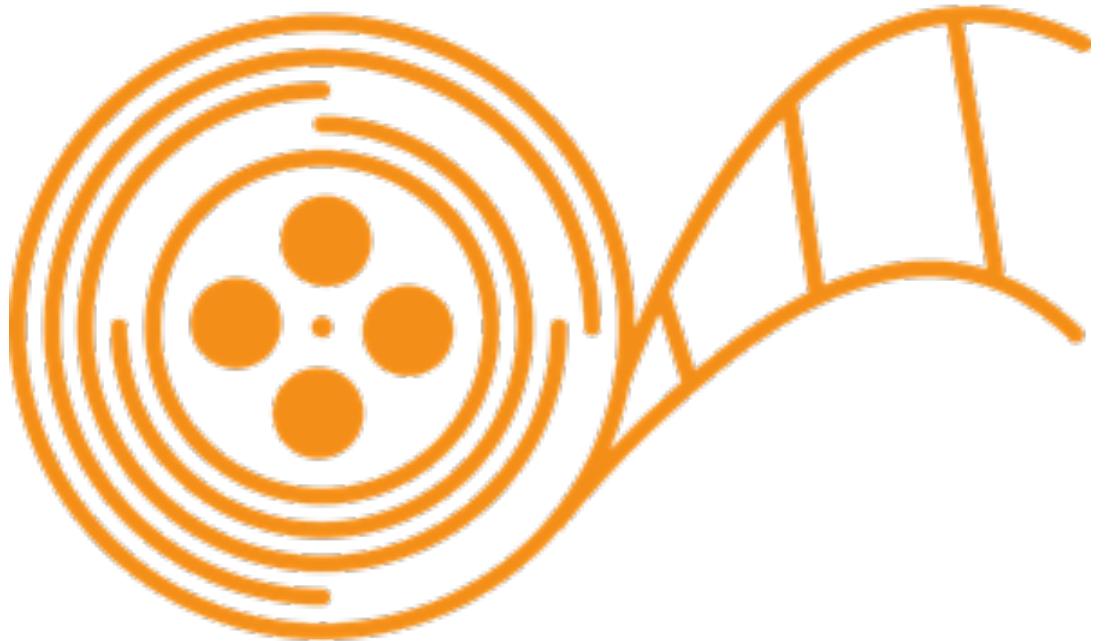
- Each EC2 instance is assigned a private DNS host name associated with its private IP address
- Both IP version 4 and IP version 6 addressing is supported
- IPv4 is default and required; IPv6 is optional
- Address types for EC2 instances:
 - Private IP version 4
 - Public IP version 4
 - Elastic IP address (Static public IP)
 - IP version 6 address (Public)



Private / 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 to the EC2 instance
- Public IP addresses are assigned to your instance from Amazon's pool of public IPv4 addresses
 - When an EC2 instance is stopped the AWS public IP address is unassigned (removed)
 - When an EC2 instance is started a new IP address is assigned
- Assigning an Elastic IP address to an EC2 instance assigns a static public address that is linked to your AWS account





Demo: IP
Addresses

Route Tables

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

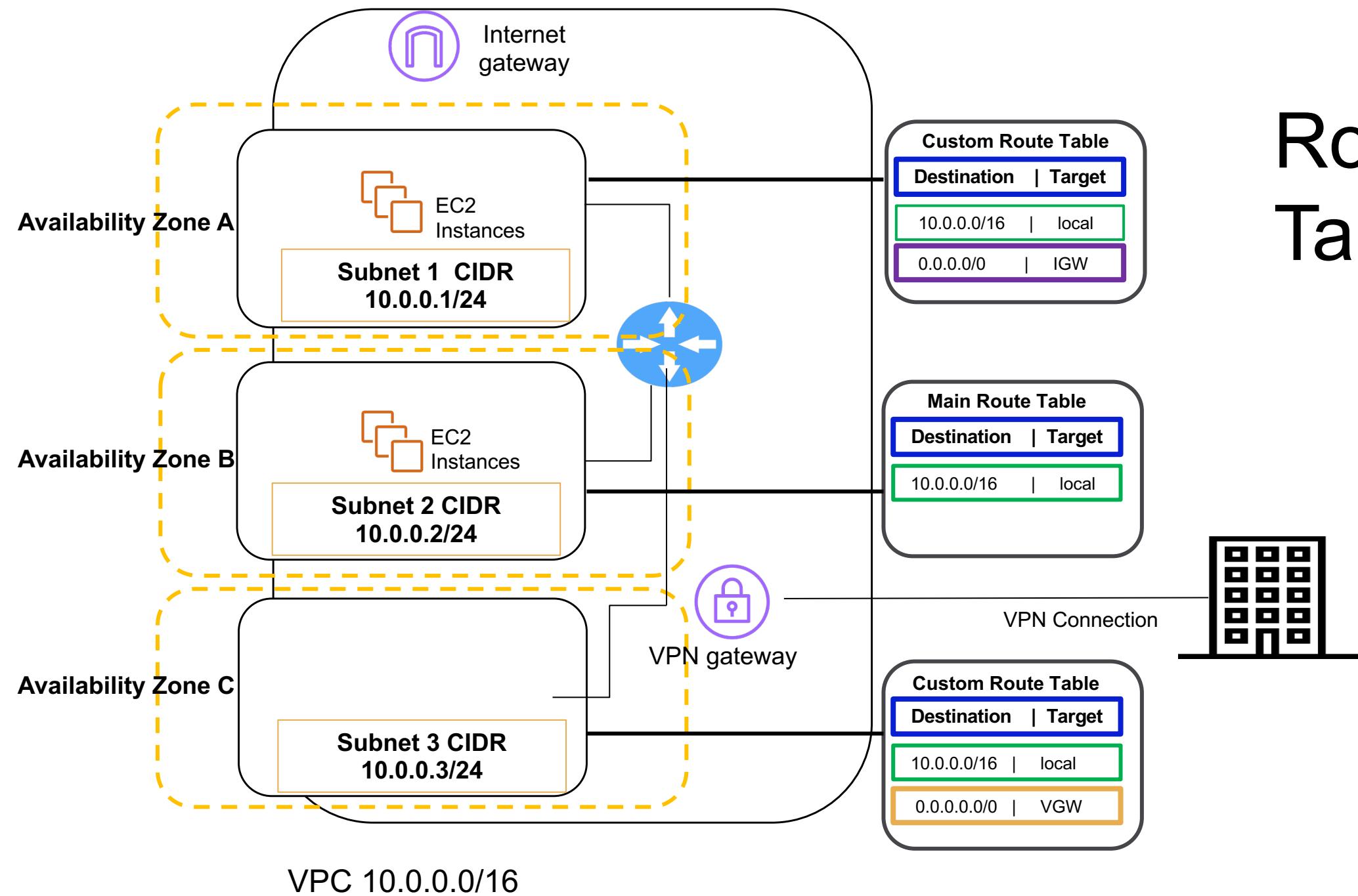


Route Tables

- Each new subnet is automatically associated with the default route table that was created when the VPC was first created
- The main route table controls the routing for all subnets that are not explicitly associated with any other route table
- 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

Internet Gateway (IGW)

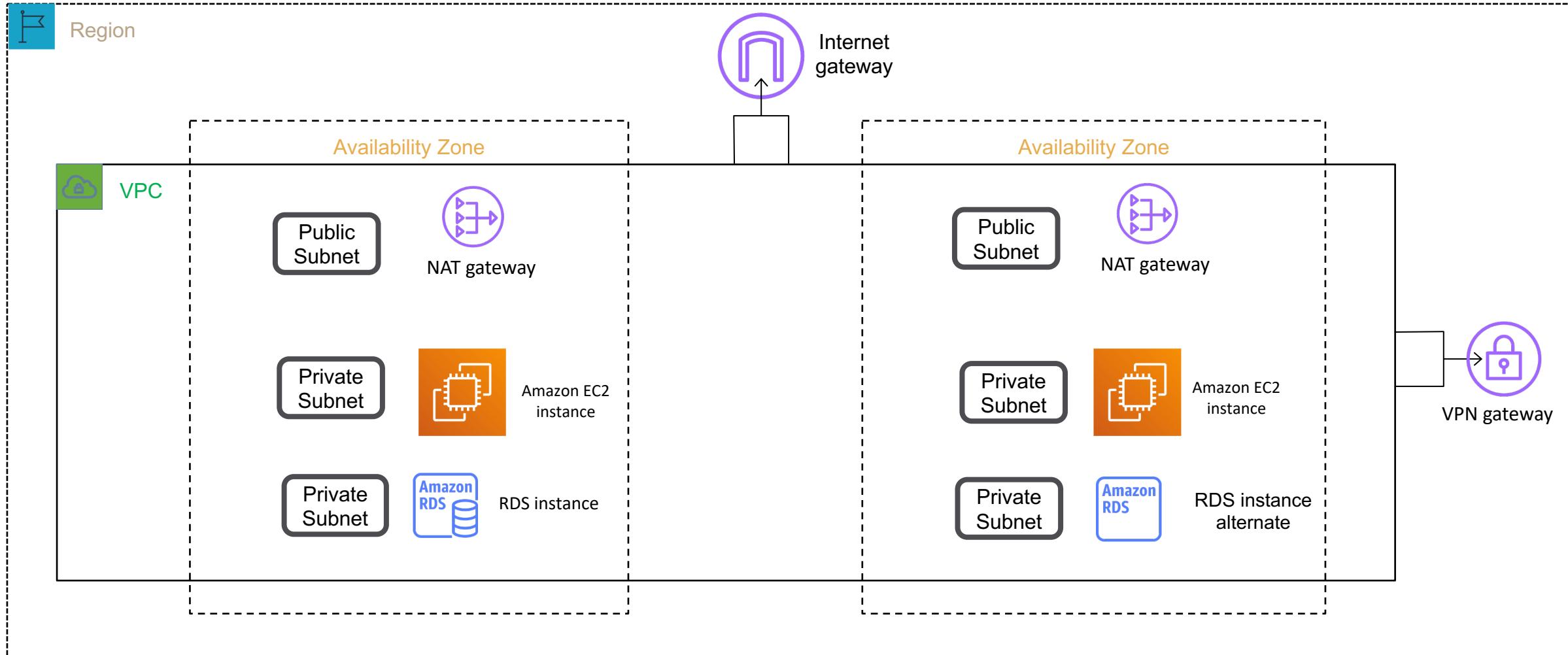
- Allows communication between instances or services hosted on public subnets and the Internet
- The internet gateway is a managed AWS service providing Internet access from public subnets
- To enable access to the Internet you must:
 1. Order an IGW
 2. Attach the IGW to your VPC
 3. Add route table entry pointing to the IGW for the public subnet
 4. Instances must have a public IP address





Demo:
Internet
Gateway

Two-tier Application



Internet gateway



Public Door

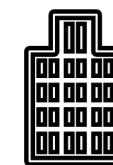
Internet

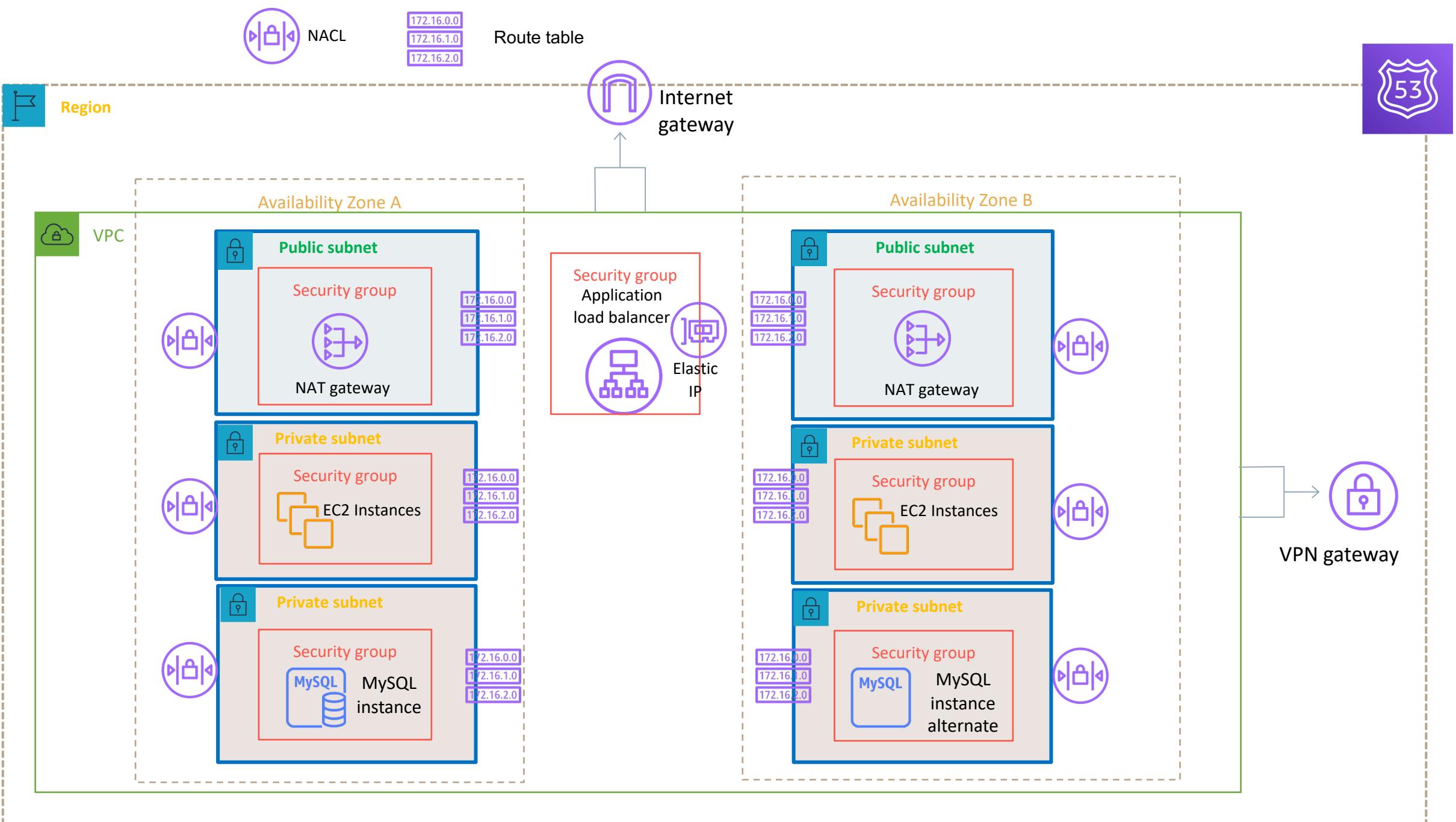


VPN gateway



Private
Door





NAT Gateway Service

NAT Services

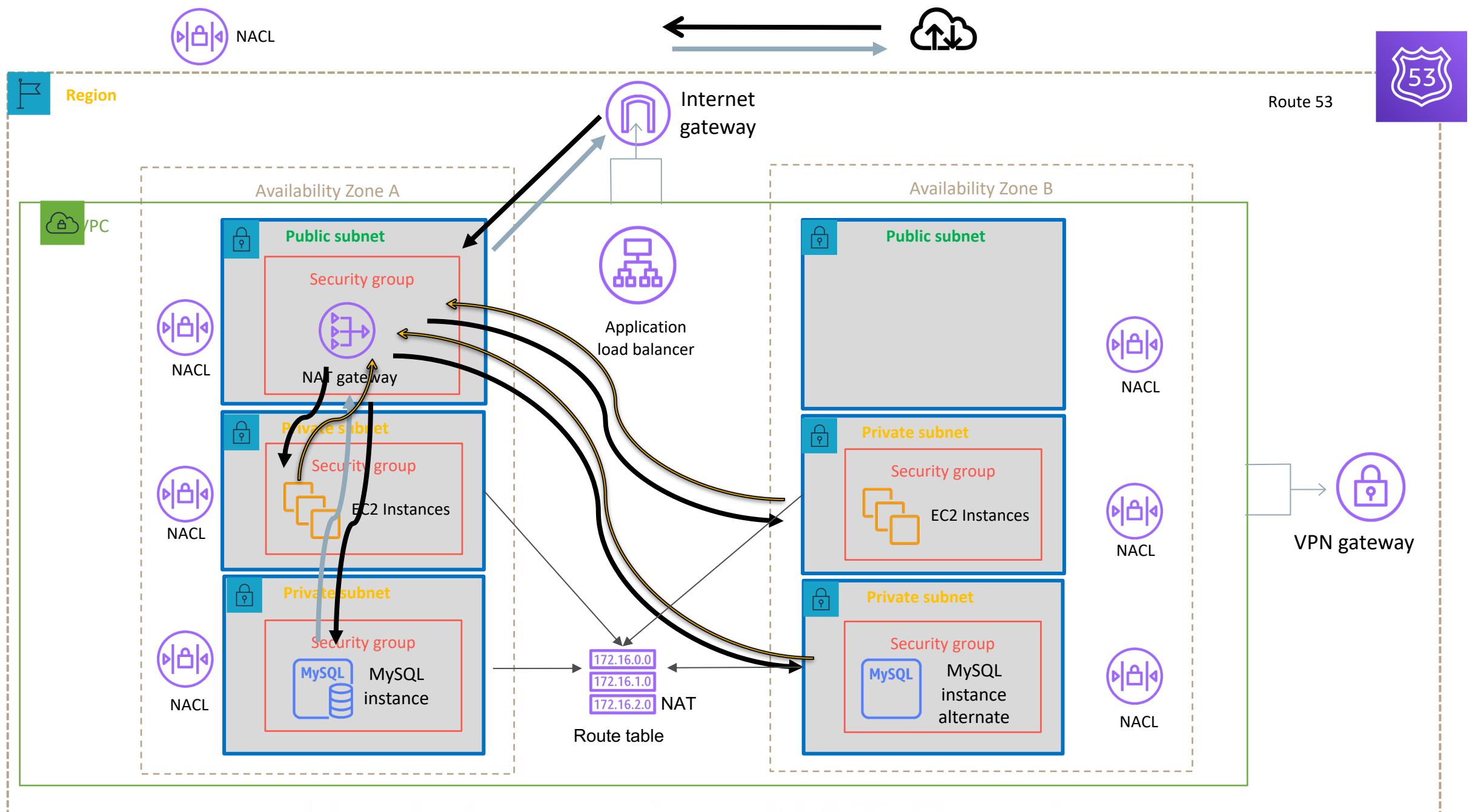
- NAT services enable instances in a private subnet to connect to the Internet to get updates
- Traffic requests from the instance are forwarded to the NAT service hosted in the public subnet
- Internet response is sent back to the private instance that made the request
- NAT Options:
 - NAT gateway service – hosted NAT services provided by AWS
 - NAT instance – compute instance created with a NAT AMI

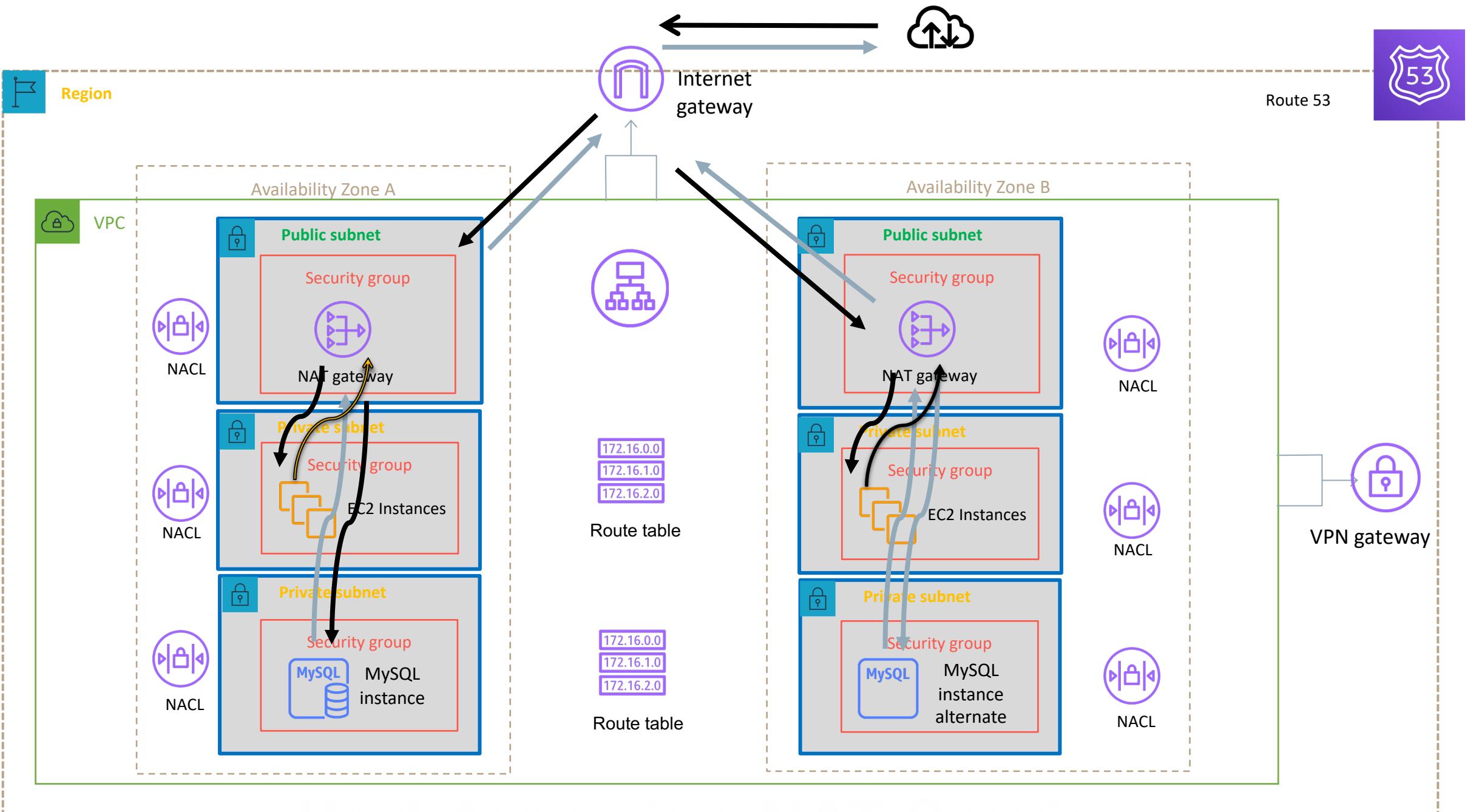


NAT Gateway Services



- 5 Gbps baseline : scales automatically to 45 Gbps
- One Elastic IP address can be assigned
- Security Groups are not supported
- Network ACL's are supported
- Supports up to 55,000 simultaneous connections to unique destinations







Demo:
Nat
Gateway

Network ACLs

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





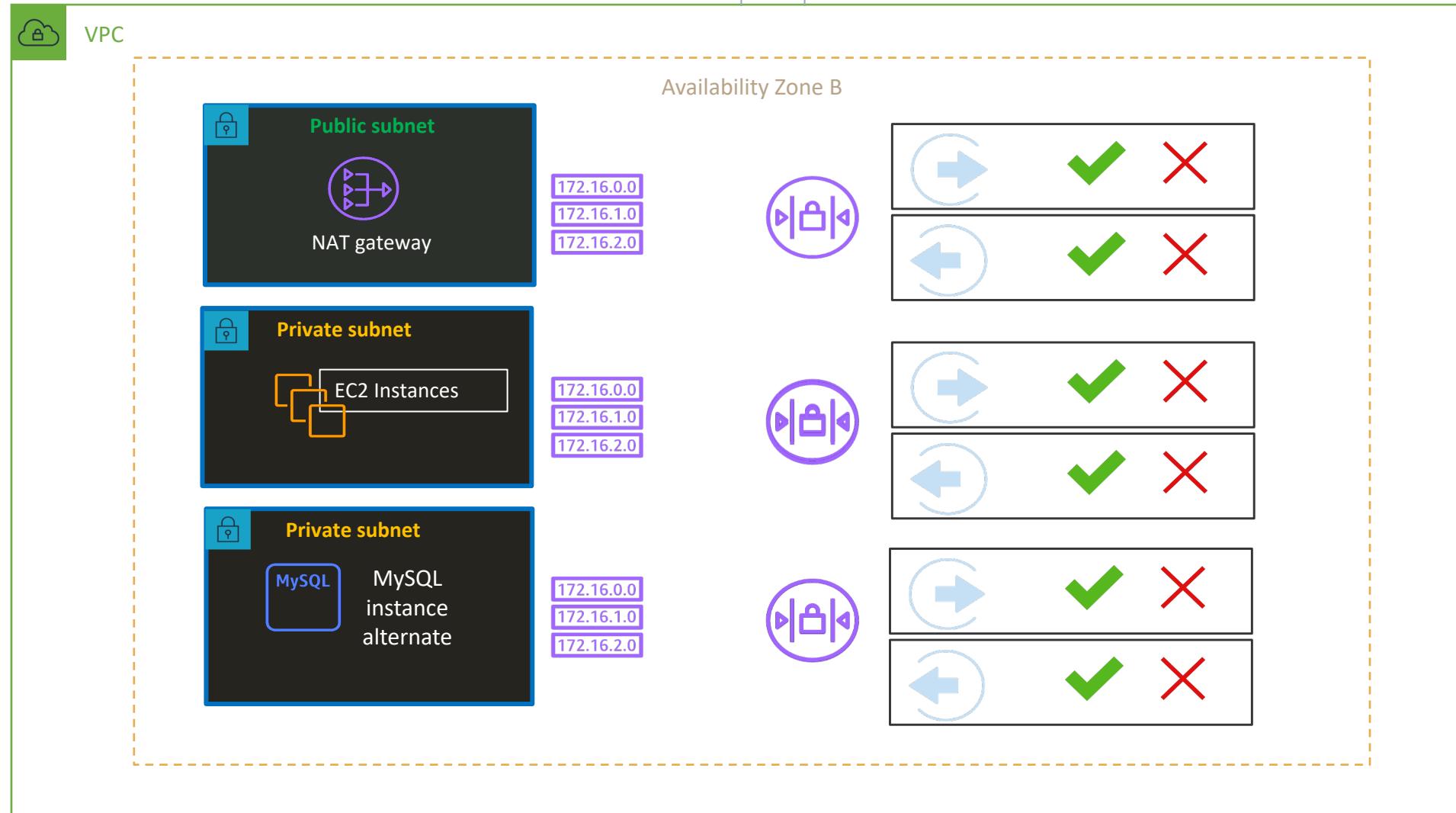
NACL

| |
|------------|
| 172.16.0.0 |
| 172.16.1.0 |
| 172.16.2.0 |

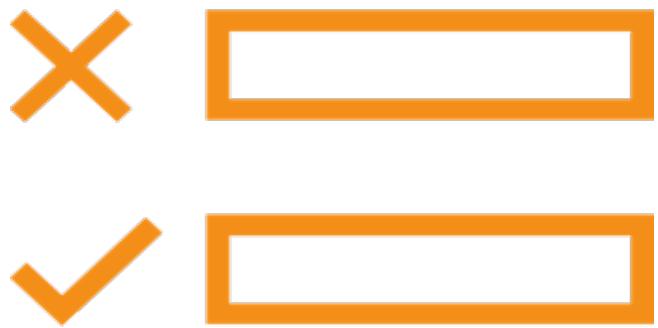
Route table



Internet gateway



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

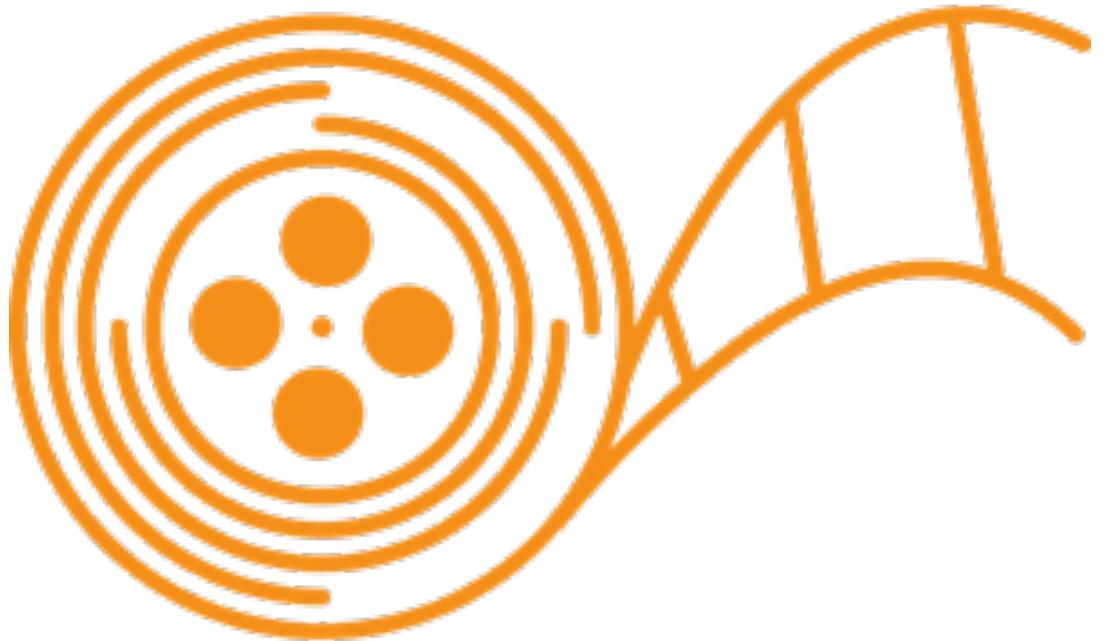


VPC Flow Logs

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





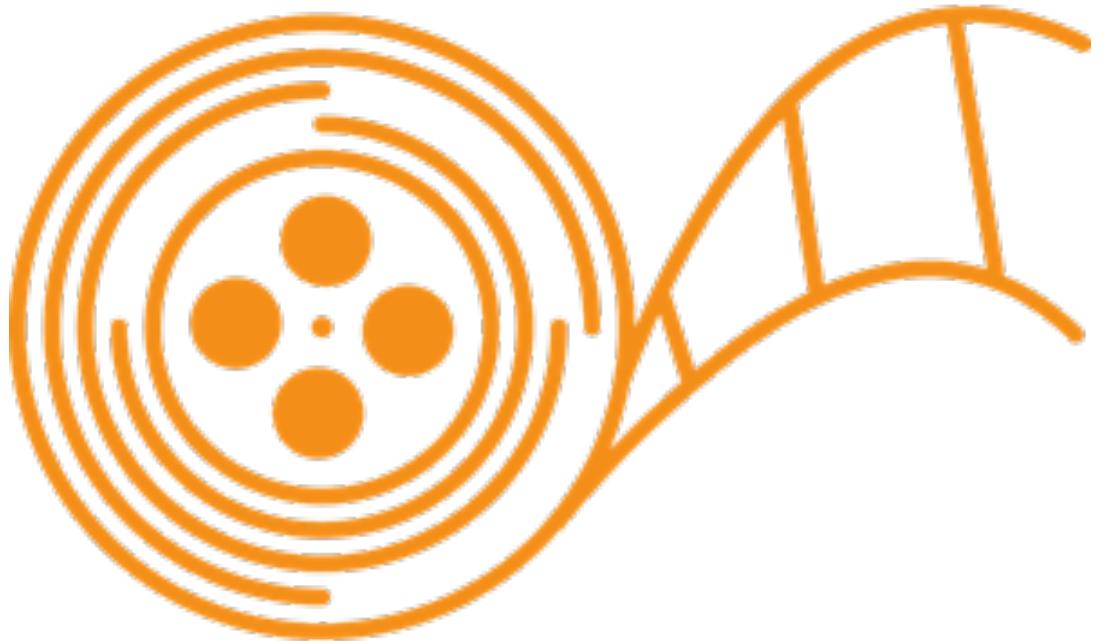
Demo:
Flow
Logs

Private Connections

VPC Private Endpoints

- Privately connect your VPC to AWS services
- Interface endpoint – network interface with private IP address to hosted AWS service
- Gateway endpoint – Dynamo DB table

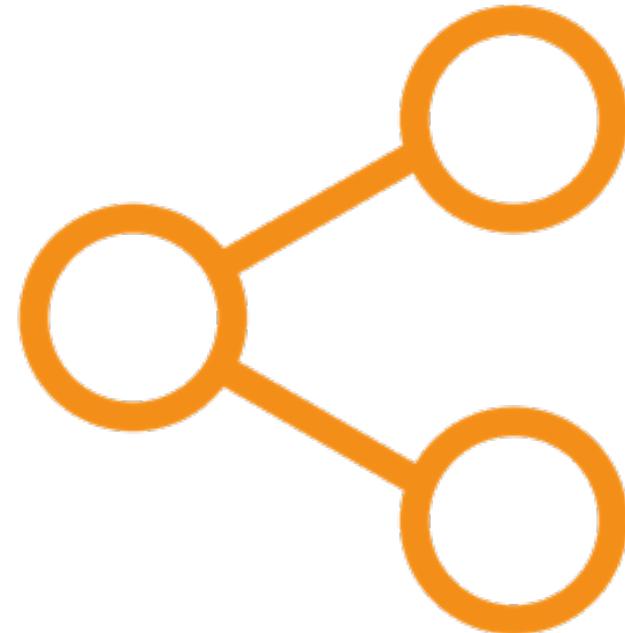




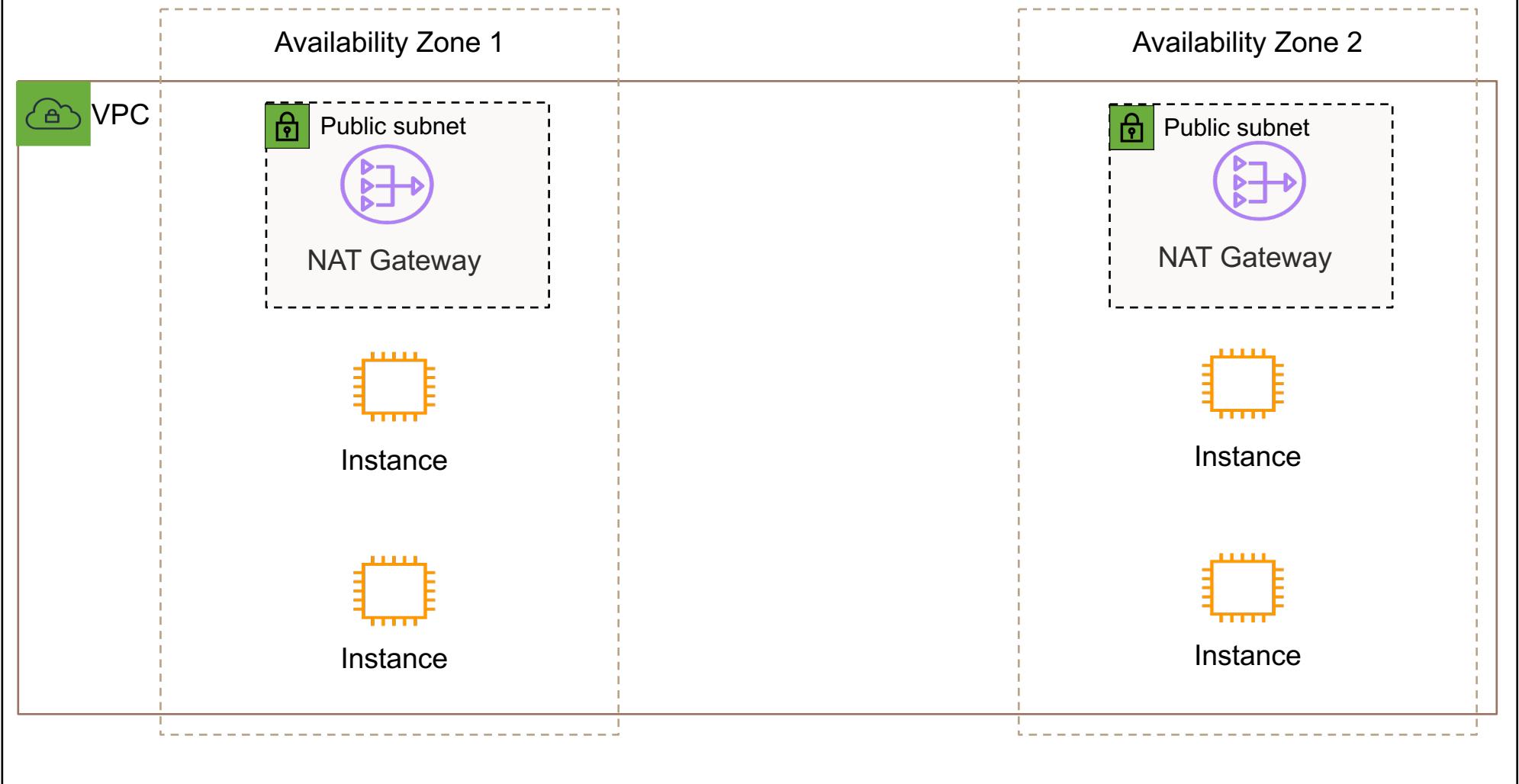
Demo:
Private
Connections

Peering VPC's

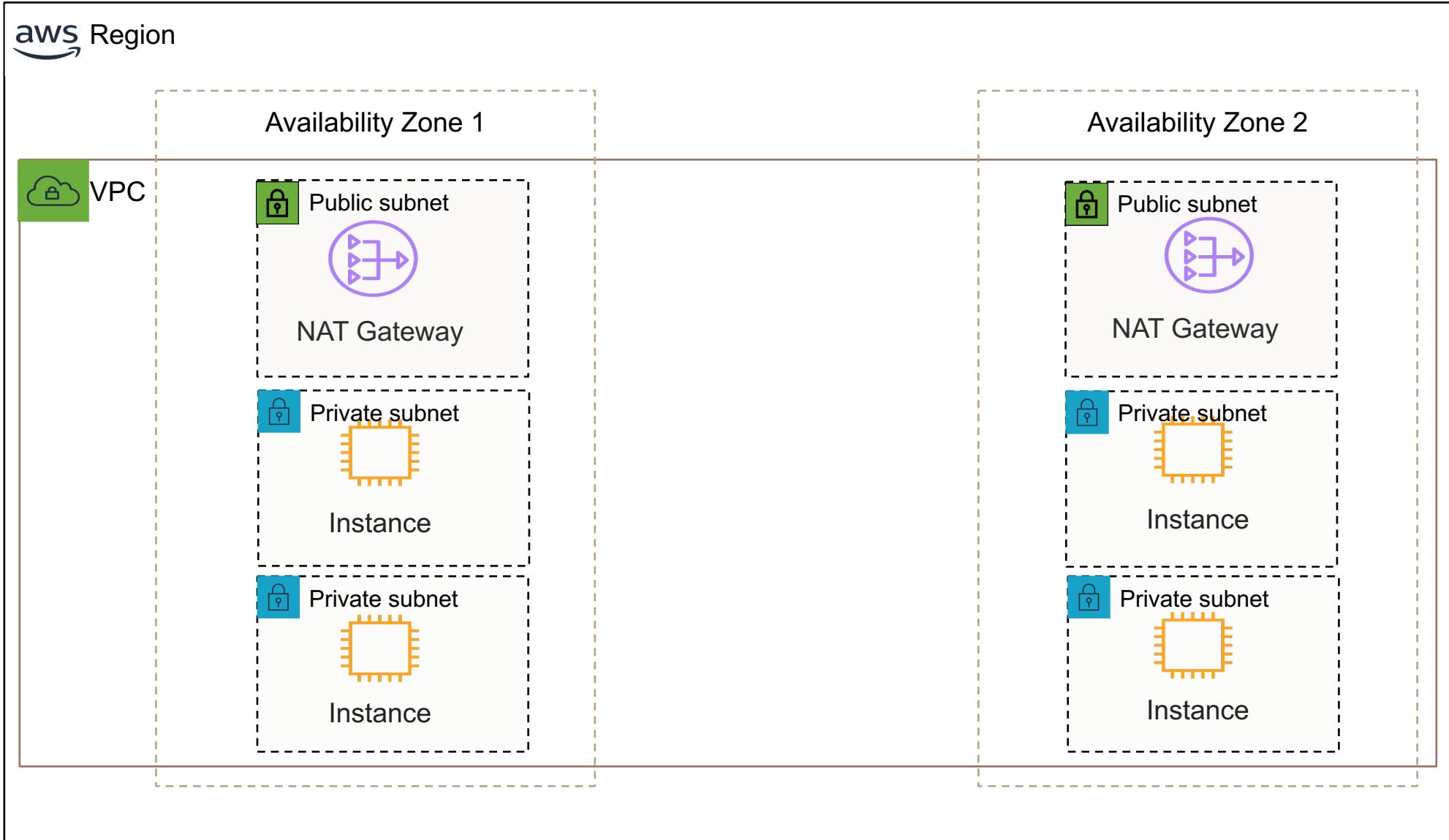
- Networking connection between two VPC's
- Peer your VPC's or between other account holders VPC's using a private IP address
- Peering is a one-to-one relationship
- Peering connections are not transitive
- CIDR blocks can't overlap in a peering relationship
- Peering connections can be created between VPCs in the same region
- Peering connections can be created between VPCs in different regions

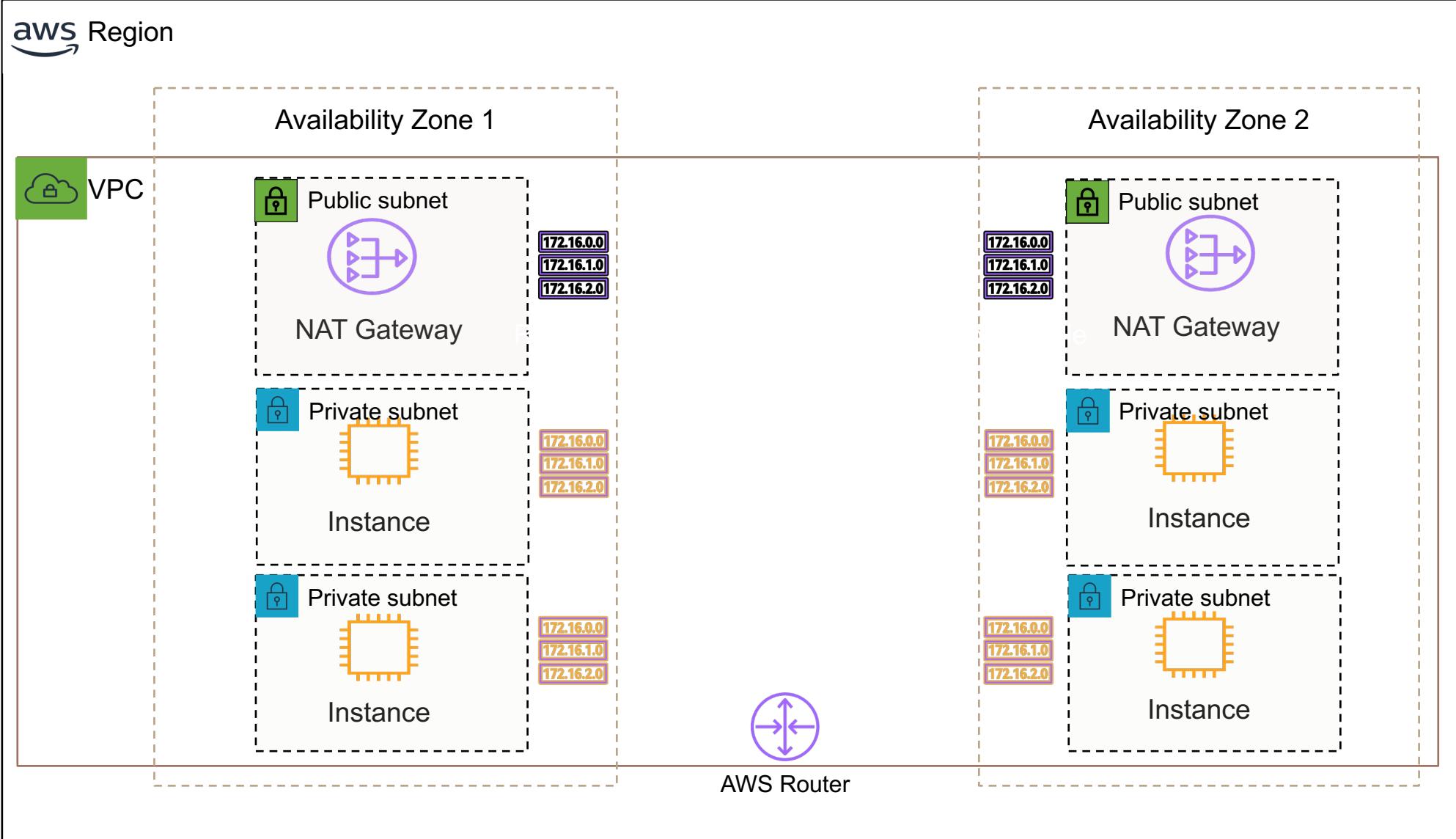


aws Region



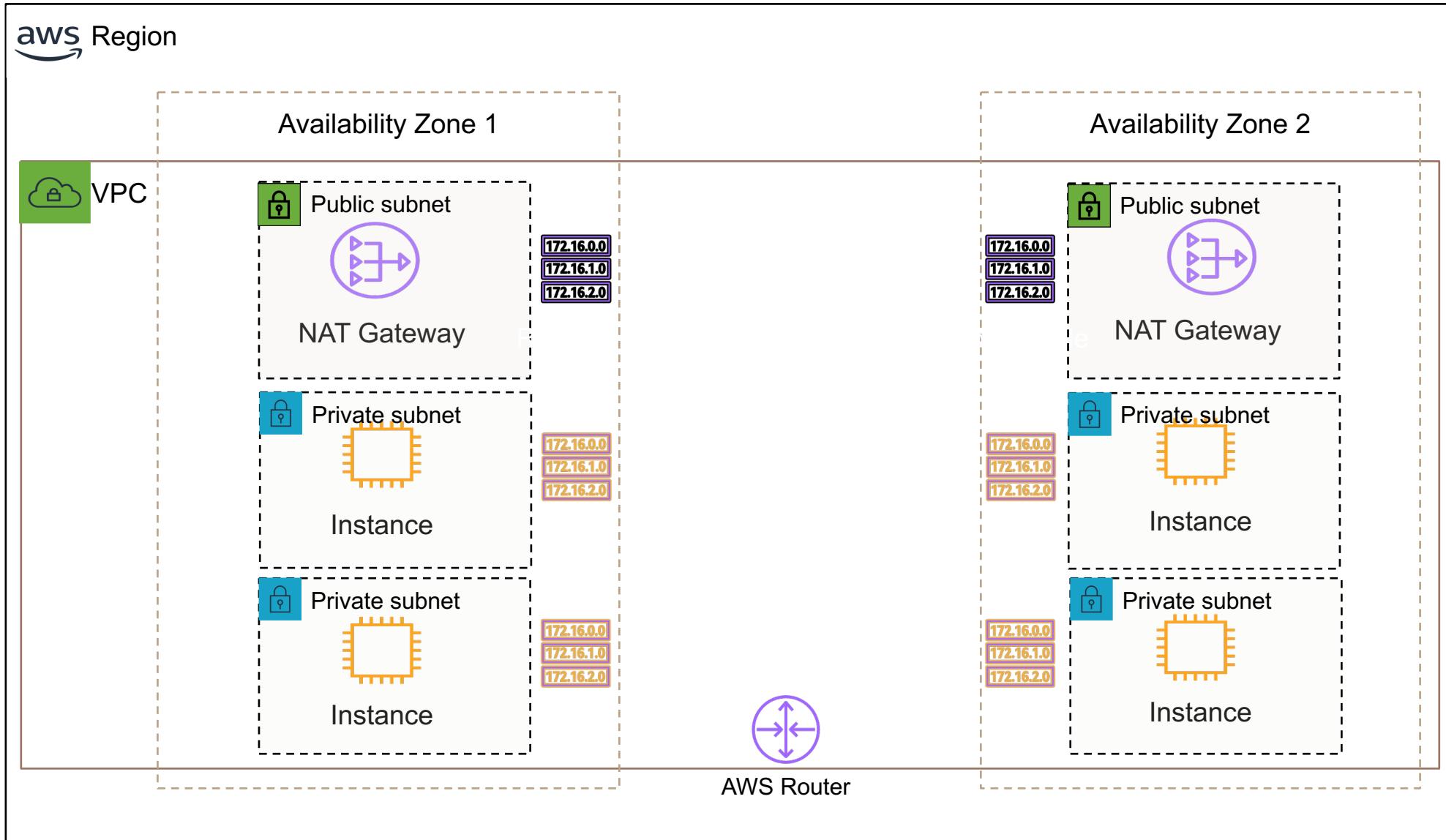
aws Region





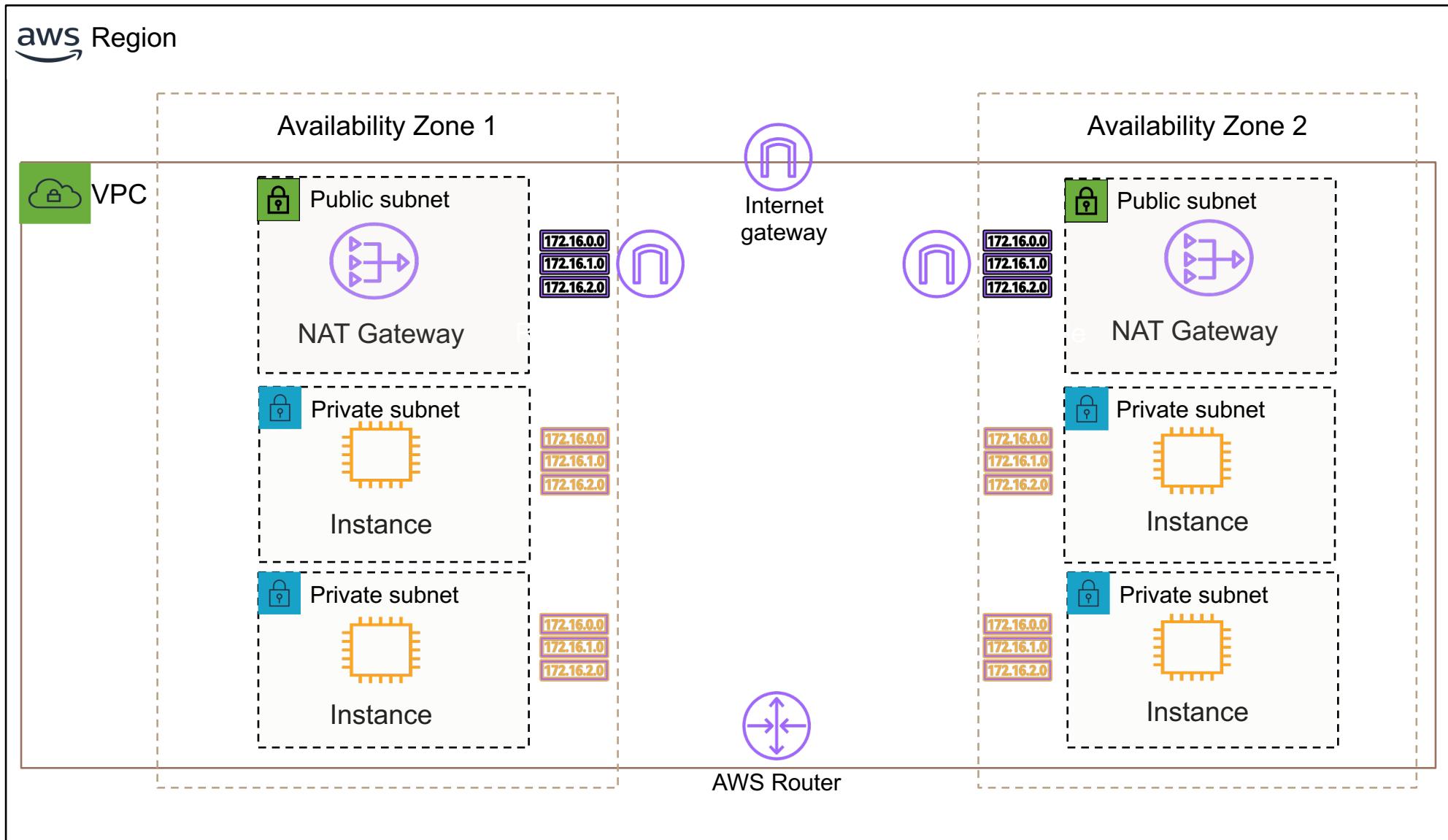
Custom
172.16.0.0
172.16.1.0
172.16.2.0
Table

Custom
172.16.0.0
172.16.1.0
172.16.2.0
Table



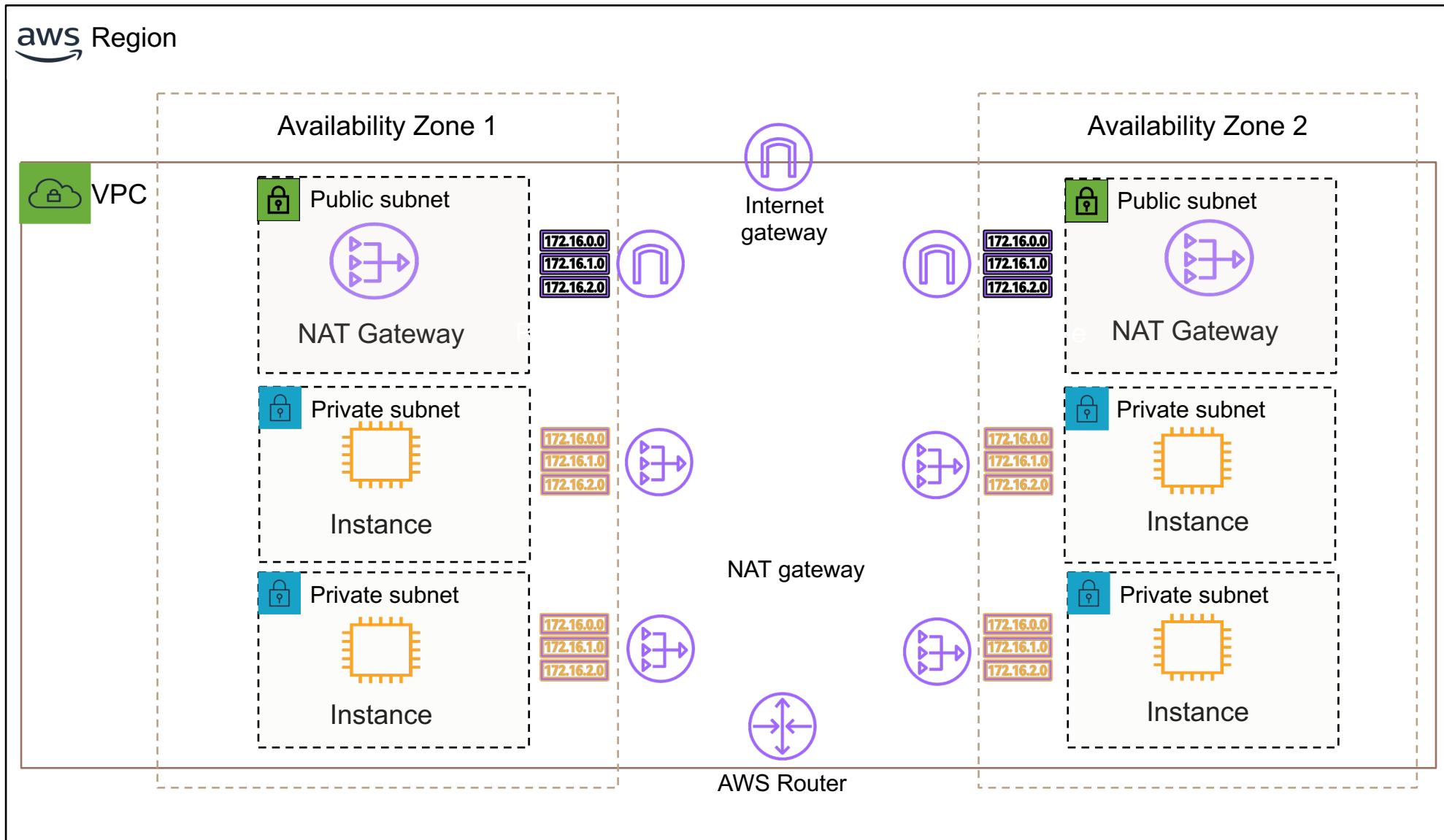
172.16.0.0 Custom
172.16.1.0 Public Route
172.16.2.0 Table

172.16.0.0 Custom
172.16.1.0 Private Route
172.16.2.0 Table



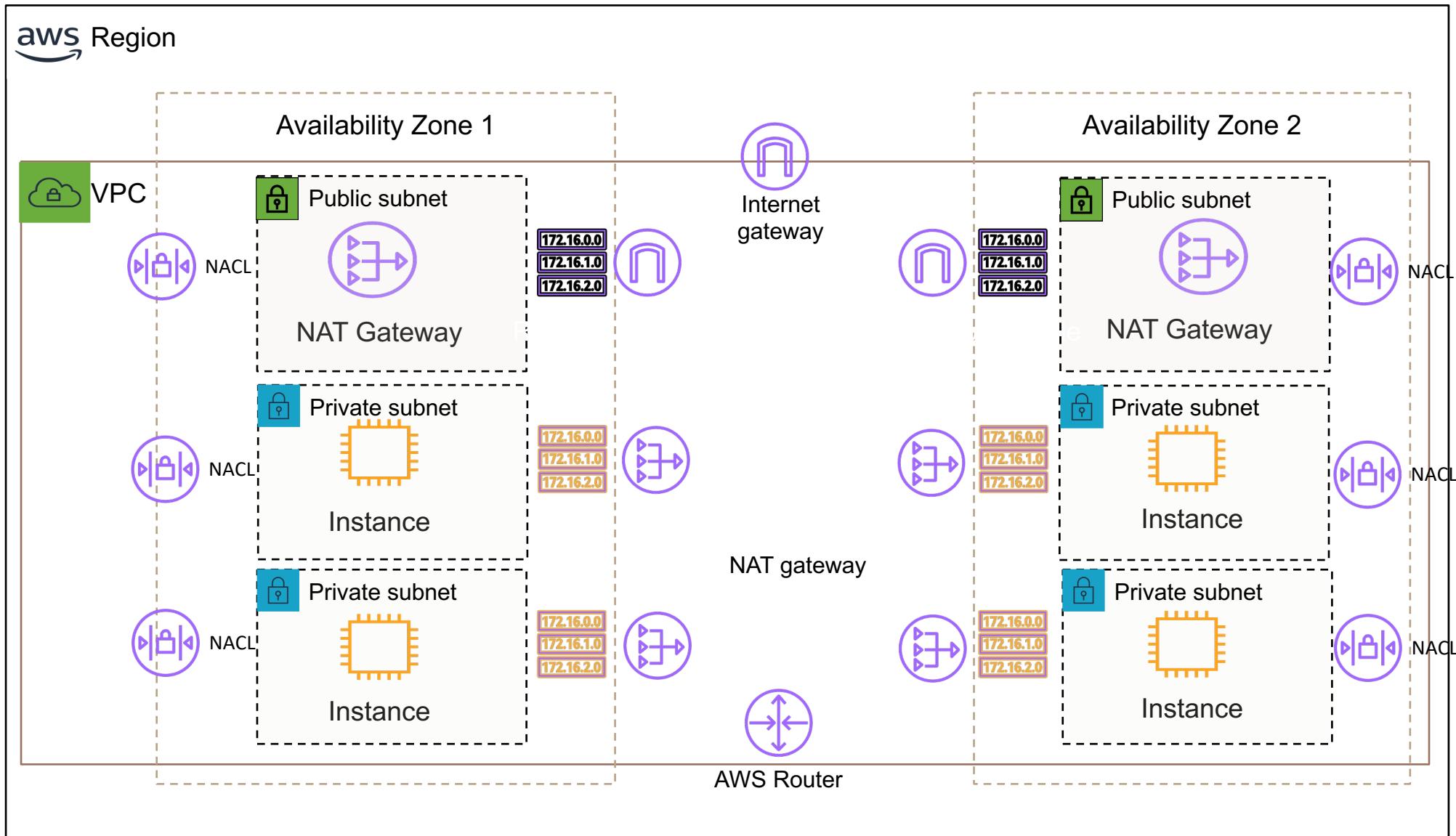
172.16.0.0 Custom
172.16.1.0 Public Route
172.16.2.0 Table

172.16.0.0 Custom
172.16.1.0 Private Route
172.16.2.0 Table



Custom
172.16.0.0
172.16.1.0
172.16.2.0
Table

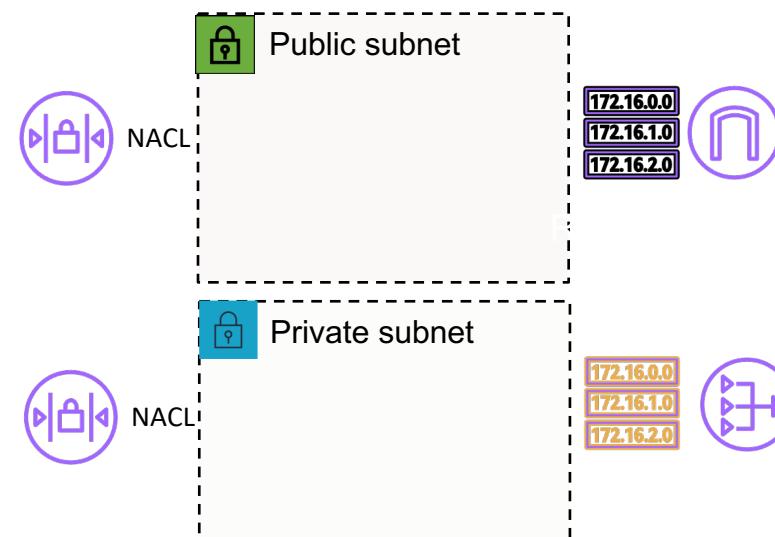
Custom
172.16.0.0
172.16.1.0
172.16.2.0
Table



Subnet Security Controls

172.16.0.0
172.16.1.0
172.16.2.0
Custom
Public Route
Table

172.16.0.0
172.16.1.0
172.16.2.0
Custom
Private Route
Table



EC2 Instances

Instance Hosting since 2017

Performance of storage, networking, management improved with custom chipsets

Hardware replaces software emulation speeding up EC2 communications

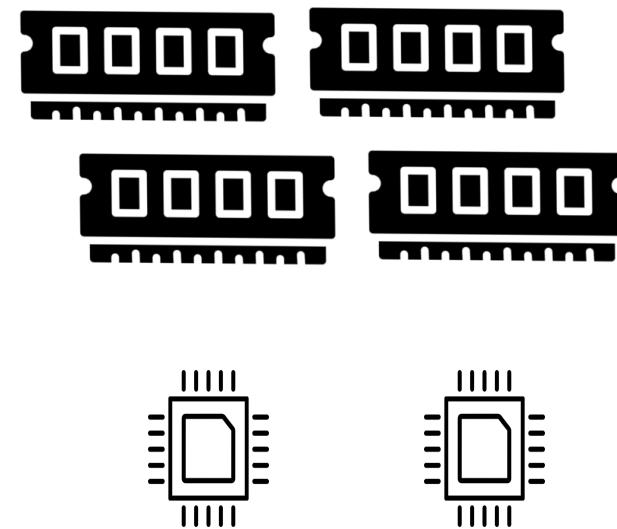
Replacing Xen hypervisor with lightweight hypervisor called Nitro

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



EBS Instance Storage Security

Nitro Hypervisor



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



The screenshot shows the AWS Management Console interface for selecting instance types. On the left, a sidebar lists various instance types categories: All instance types, Micro instances, General purpose, Compute optimized, FPGA instances, GPU graphics, GPU instances, GPU compute, Memory optimized, and Storage optimized. Below this is a 'Filter by:' dropdown set to 'All instance types'. In the center, there are two buttons: 'Current generation' (which is highlighted) and 'All generations'. At the bottom, there are 'Show/Hide Columns' and another 'Current generation' dropdown.

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 |
|--------------------------|-------------------|--------------|-------|---------------------------|-------------|--------------|-----------------------|-------------------------|---------------------|
| <input type="checkbox"/> | Compute optimized | c5.18xlarge | 72 | Intel Xeon Platinum 8124M | 3 GHz | 144 | EBS only | Yes | 25 Gigabit |
| <input type="checkbox"/> | Compute optimized | c5d.18xlarge | 72 | Intel Xeon Platinum 8124M | 3 GHz | 144 | 2 x 900 (SSD) | Yes | 25 Gigabit |
| <input type="checkbox"/> | Compute optimized | c5n.18xlarge | 72 | Intel Xeon Platinum 8124M | 3 GHz | 192 | EBS only | Yes | 100 Gigabit |
| <input type="checkbox"/> | Compute optimized | c5n.2xlarge | 8 | Intel Xeon Platinum 8124M | 3 GHz | 21 | EBS only | Yes | Up to 25 Gigabit |

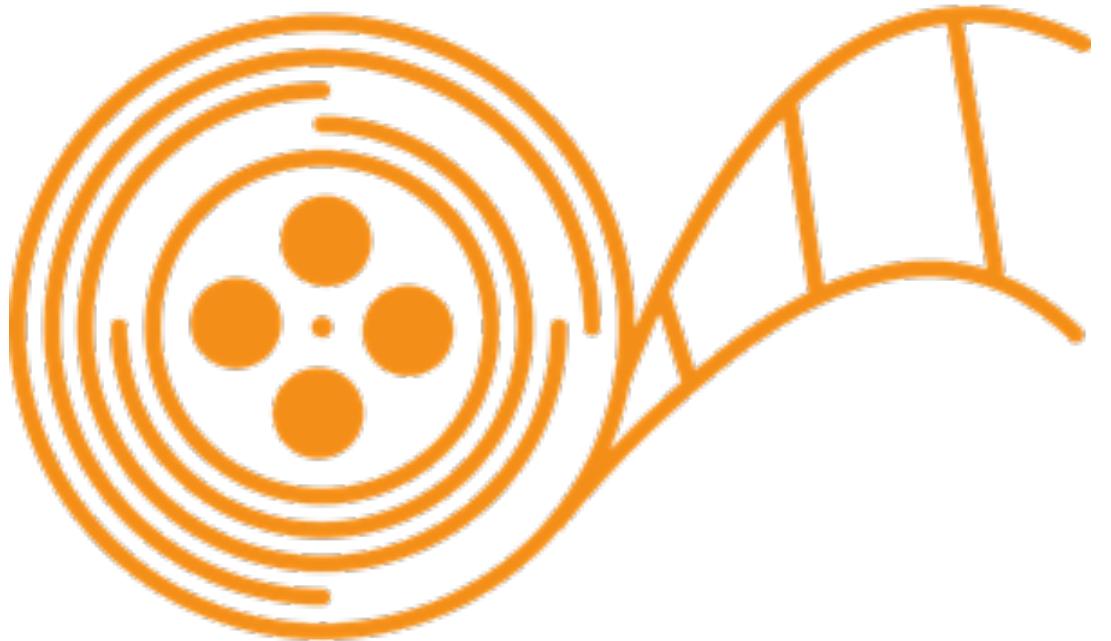
EC2 Instance FYI

- 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



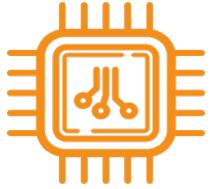
Instance Families at AWS

Instance Generation Features
Instance Family Instance Size
c4d.4xlarge
Instance Type



Demo:
Instance
Families

Instance Types



micro



Micro Instances – there's only one instance type in this class; the t1.micro with an unidentified processor

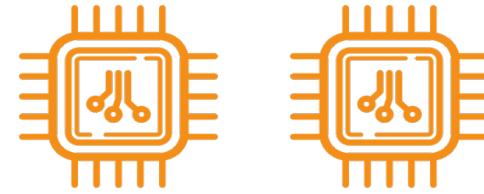


The clock speed is not identified, but you have .613 GiB of memory with very low networking performance



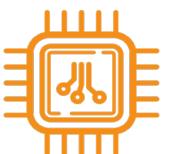
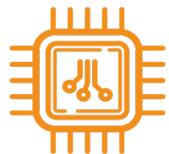
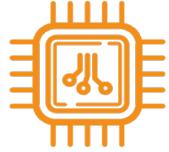
It supports 32 or 64-bit workloads and only shows up in the management console if you search for the All Generation types of instances

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 |



| 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 |

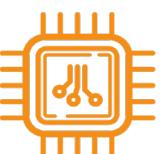
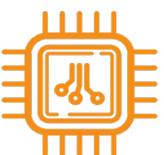
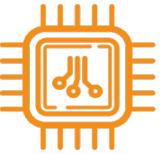
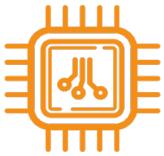
storage optimized

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

memory optimized

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

| 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 | |



Primary Network Interface

Every instance in a VPC has a default network interface, called the primary network interface (eth0)

The primary network interface can not be detached from an instance

When a network interface is created, it inherits the public IPv4 addressing attribute from the subnet it is assigned

When EC2 instances first launch, the primary network interface is assigned a reserved private IP address from the default DHCP pool assigned to the VPC

This private IP address stays assigned to the network interface until the EC2 instance is deleted



Elastic Network Interfaces

- An *elastic network interface* represents a virtual network card
- A primary private IPv4 address from the IPv4 address range of your VPC
- One or more secondary private IPv4 addresses from the IPv4 address range of your VPC
- One Elastic IP address (IPv4) per private IPv4 address
- One public IPv4 address
- One or more IPv6 addresses
- One or more security groups
- A MAC address
- A source/destination check flag
- A description

Elastic Network Interfaces (ENI)

Attach a network interface to an instance when it's running (hot attach), when it's stopped (warm attach), or when the instance is being launched (cold attach)

You can't detach the primary network interface from an EC2 instance

You can move an ENI network interface from one EC2 instance to another

The instances must be in the same Availability Zone and VPC but located on different subnets

Enhanced Networking (ENA)

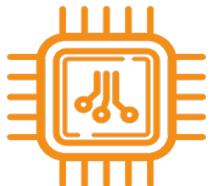
- Enhanced networking uses single root I/O virtualization (SR-IOV) to provide high-performance networking
- ENA supports network speeds of up to 200 Gbp
- The Intel 82599 Virtual Function interface supports network speeds of up to 10 Gbps
- The number and the type of network interfaces depends on the instance type

EC2 Auto Recovery

- Create an Amazon CloudWatch alarm that monitors an Amazon EC2 instance and automatically recovers the instance if it becomes impaired due to an underlying hardware failure or a problem
- A recovered instance is identical to the original instance, including the instance ID, private and public / Elastic IP addresses, and all instance metadata.
- StatusCheckFailed
- Network connectivity, loss of system power, software or hardware issues on the physical host

Burst Credits

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 it 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



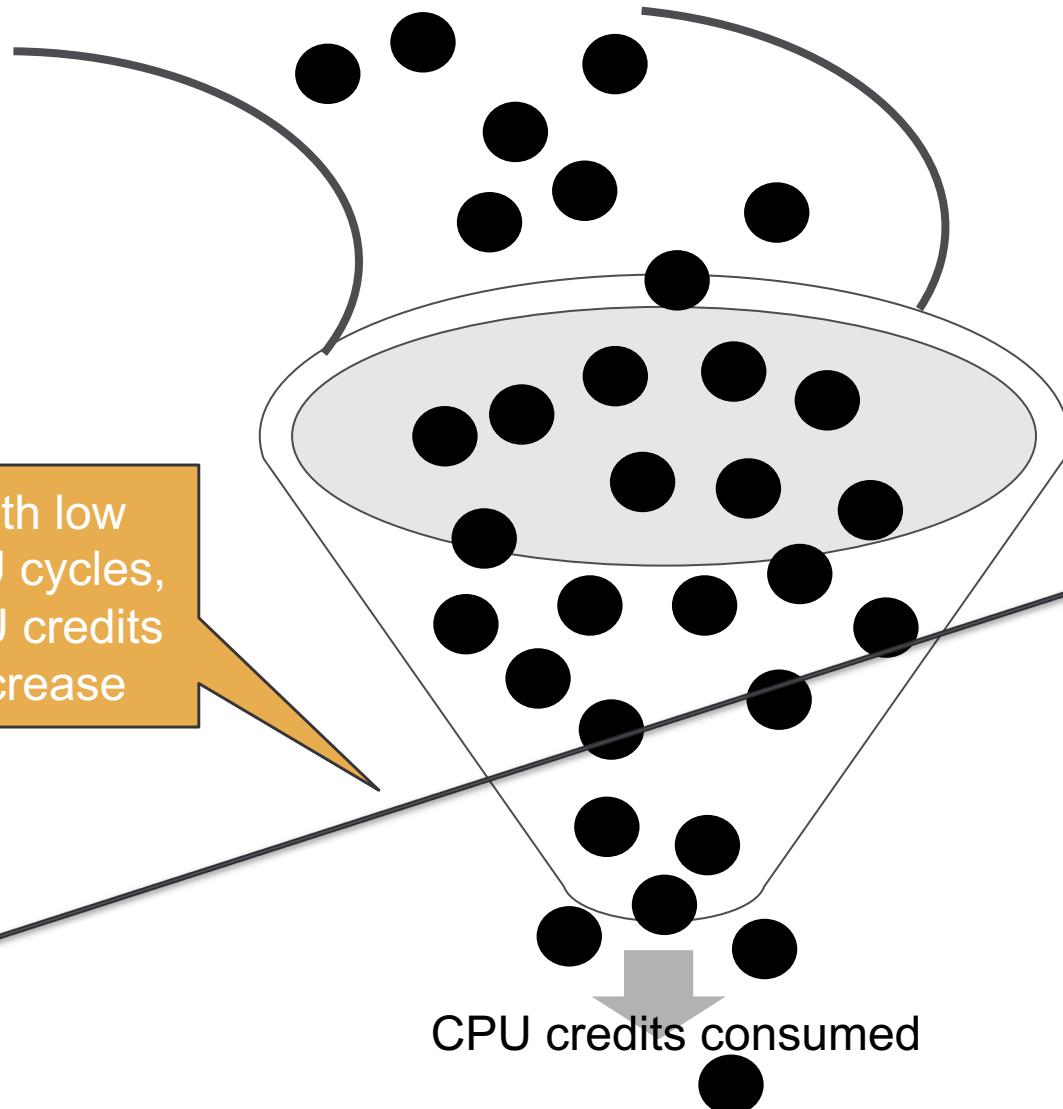
After a T2 instance is powered on it earns CPU credits at a defined steady rate up to a defined maximum value.



Earned credits expire after 24 hours

Baseline CPU credits added

With low
CPU cycles,
CPU credits
increase





Demo: t
instances

Amazon Machine Images

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



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



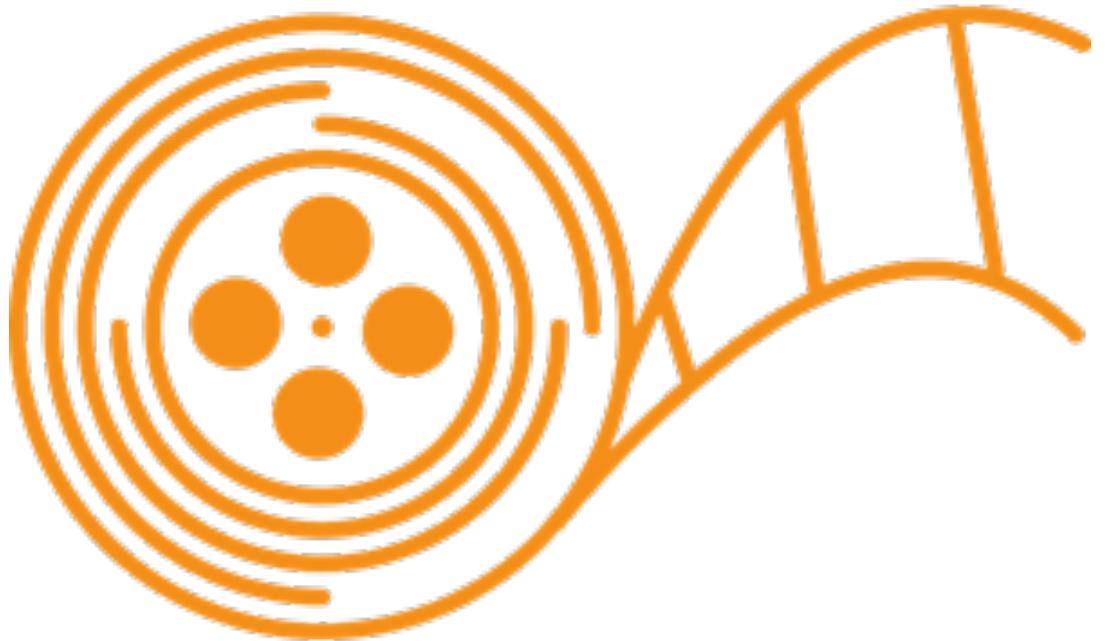
Volumes to attach – the volumes to attach to the instance at launch are contained in a block device mapping document



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



Operating system – Linux or Windows



Demo:
AMI
Creation

Security Groups

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

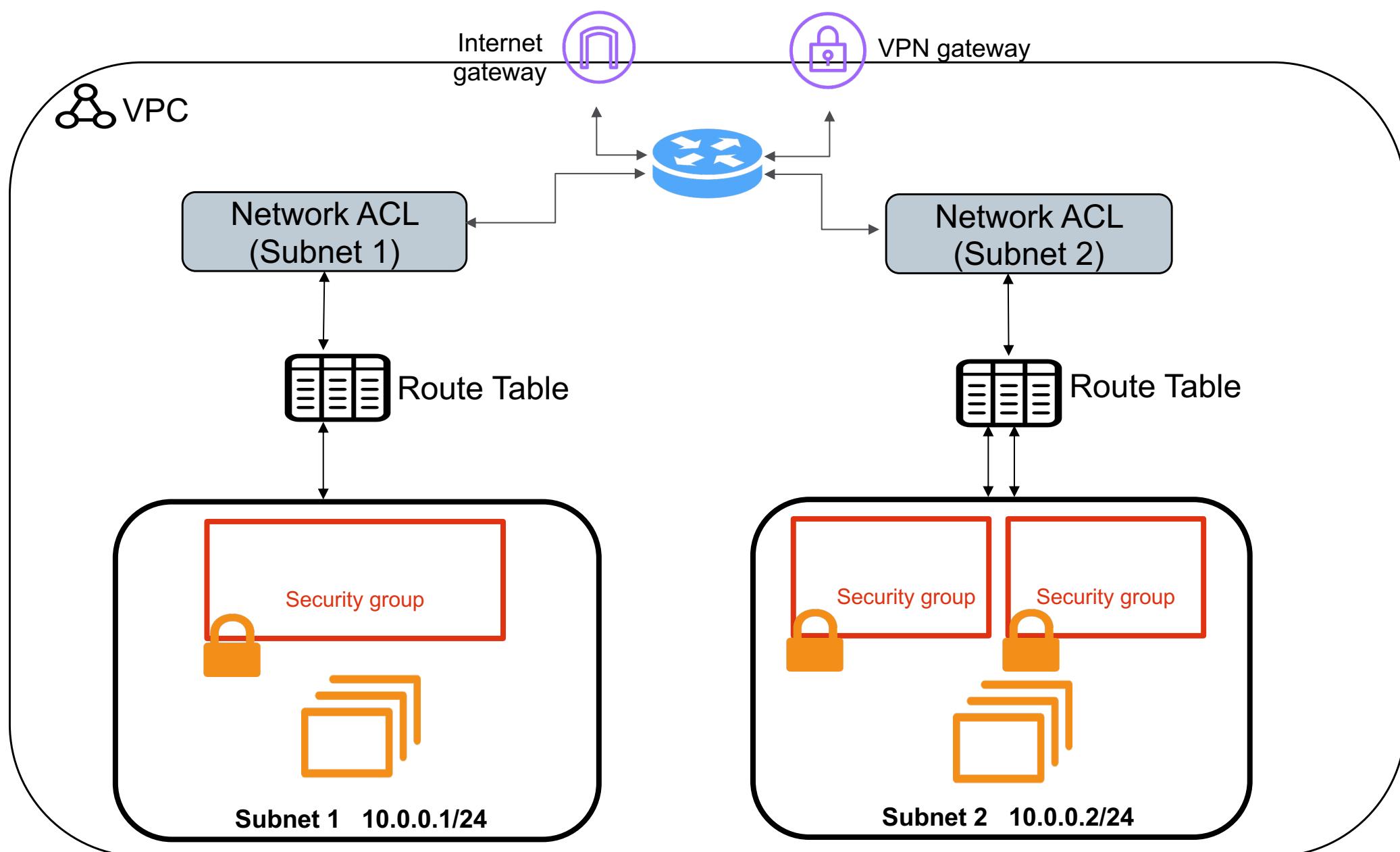
- Allow rules can be specified
- Explicit deny rules can't be specified
- Any TCP protocol that is defined with a standard port number is supported
- Outbound rules can define the destination for the traffic and the destination port or security group
- Separate allow rules can be defined for both inbound 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 can be restricted by:
 - IP protocol
 - Source / destination IP address, or address block
 - Security group





Security Groups vs NACLs

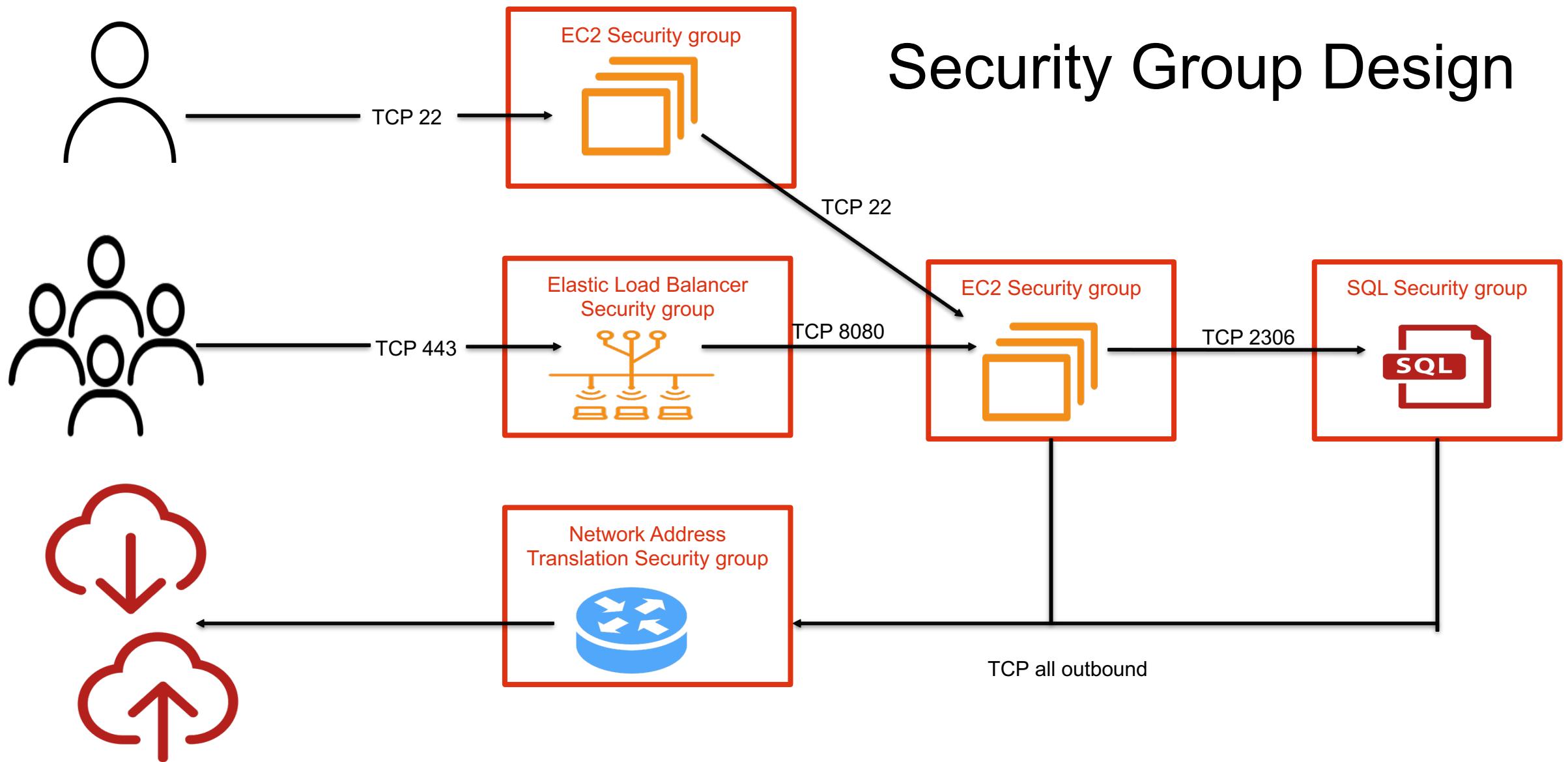
Security Groups

- Operates at the EC2 instance level
- Allow rules only supported
- Stateful: Return traffic is automatically allowed
- All rules are processed before traffic decisions are made
- Applied to the selected EC2 instance elastic network adapter

NACLs

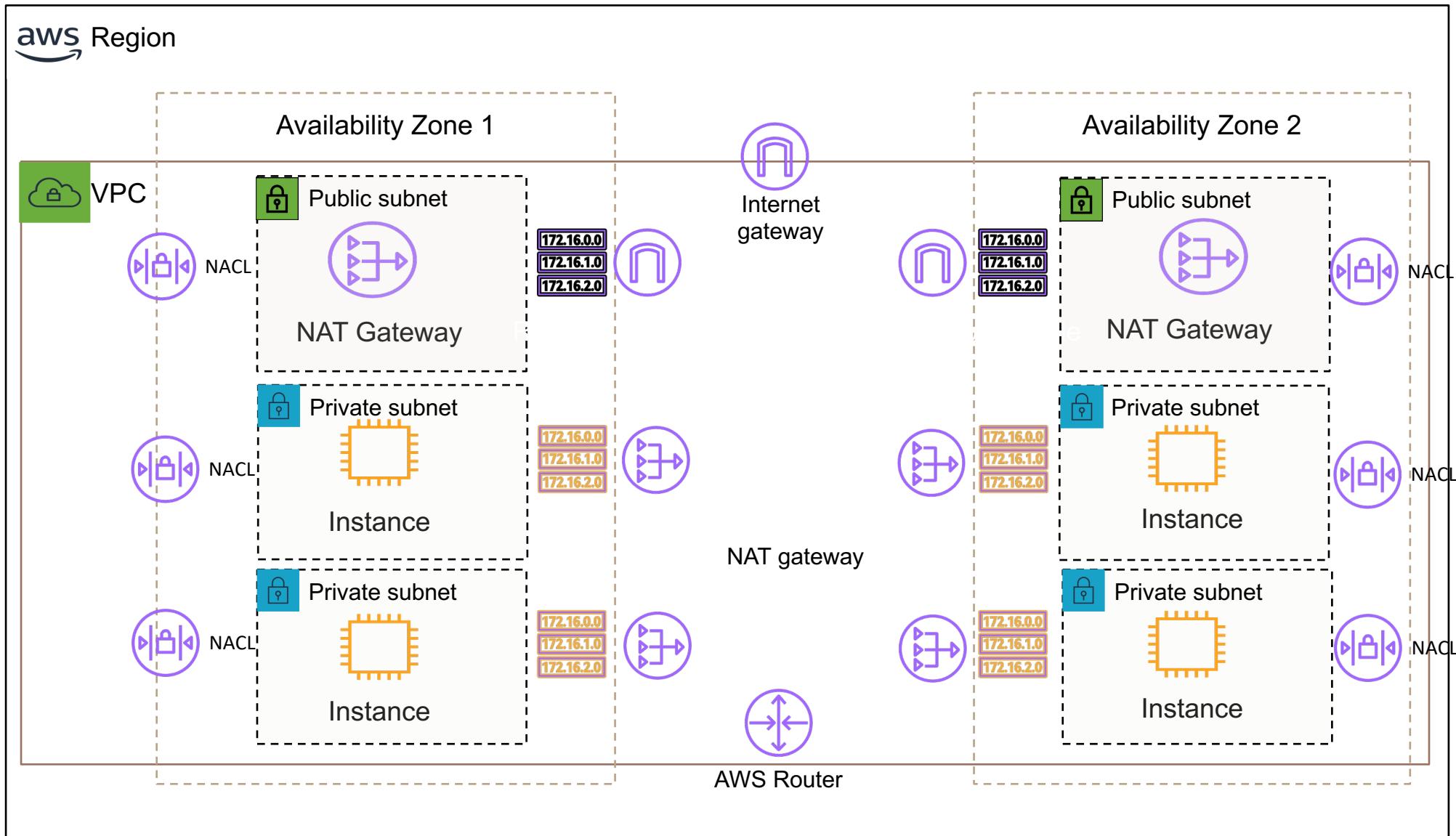
- Operates at the subnet level
- Allow and deny rules supported
- Stateless: Return traffic must be explicitly allowed by a rule
- Rules are processed in numerical order before traffic decisions are made
- Applied to a subnet

Security Group Design



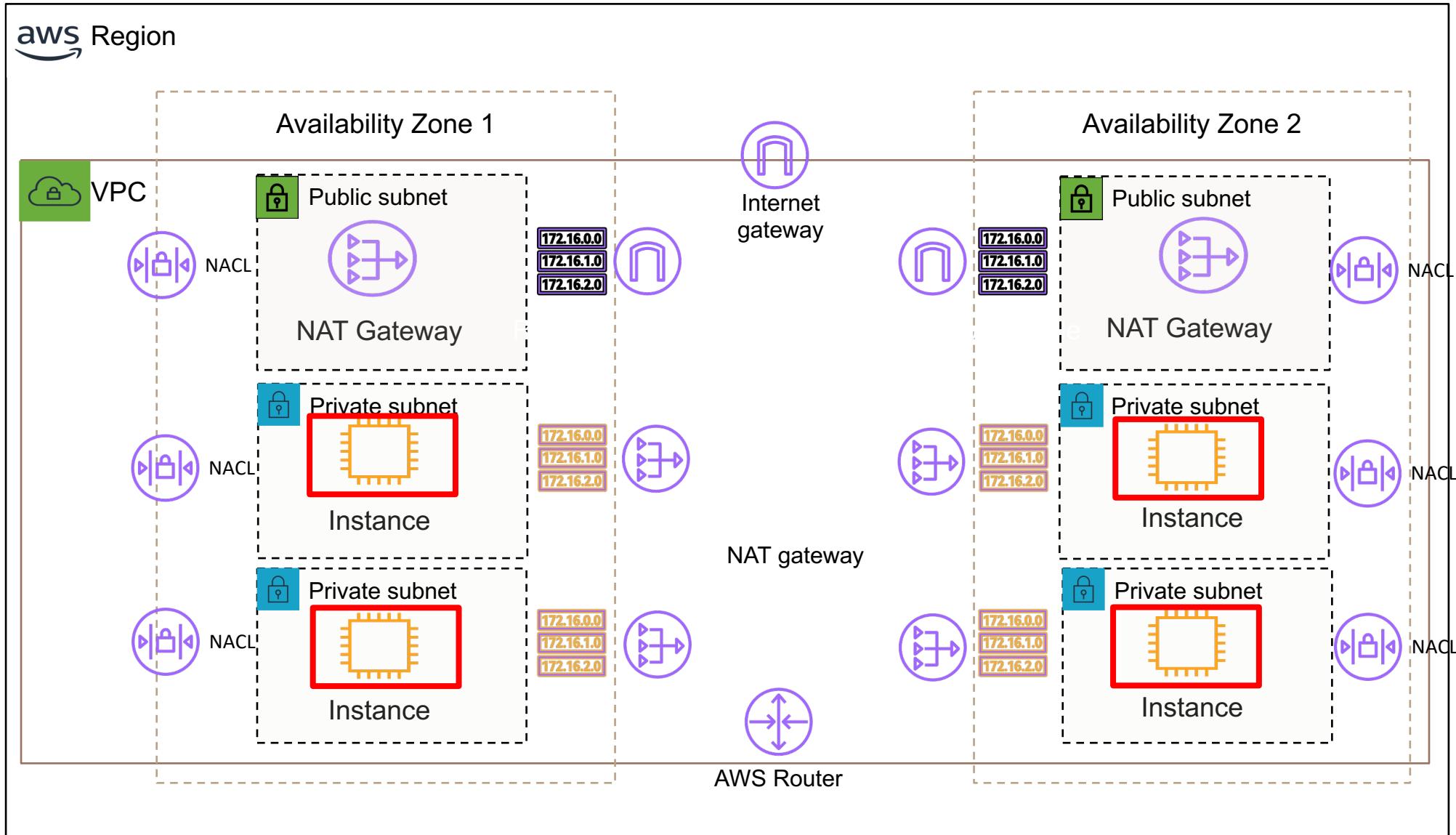
Custom
172.16.0.0
172.16.1.0
172.16.2.0
Table

Custom
172.16.0.0
172.16.1.0
172.16.2.0
Table

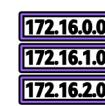


Custom
172.16.0.0
172.16.1.0
172.16.2.0
Table

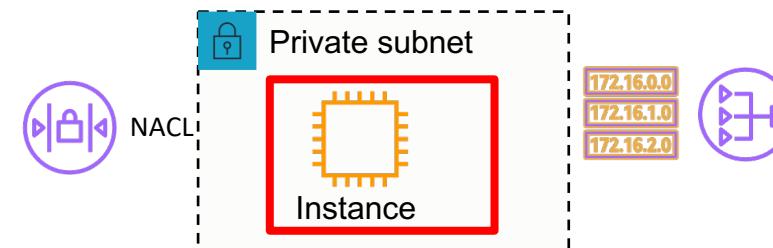
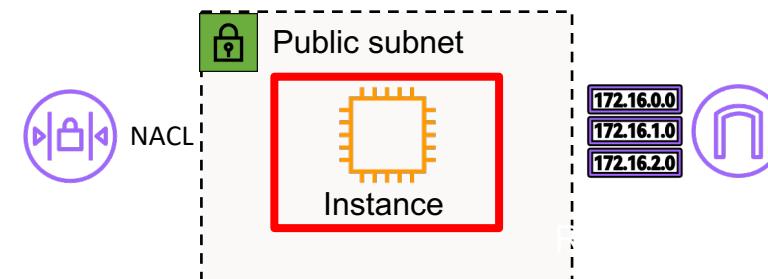
Custom
172.16.0.0
172.16.1.0
172.16.2.0
Table

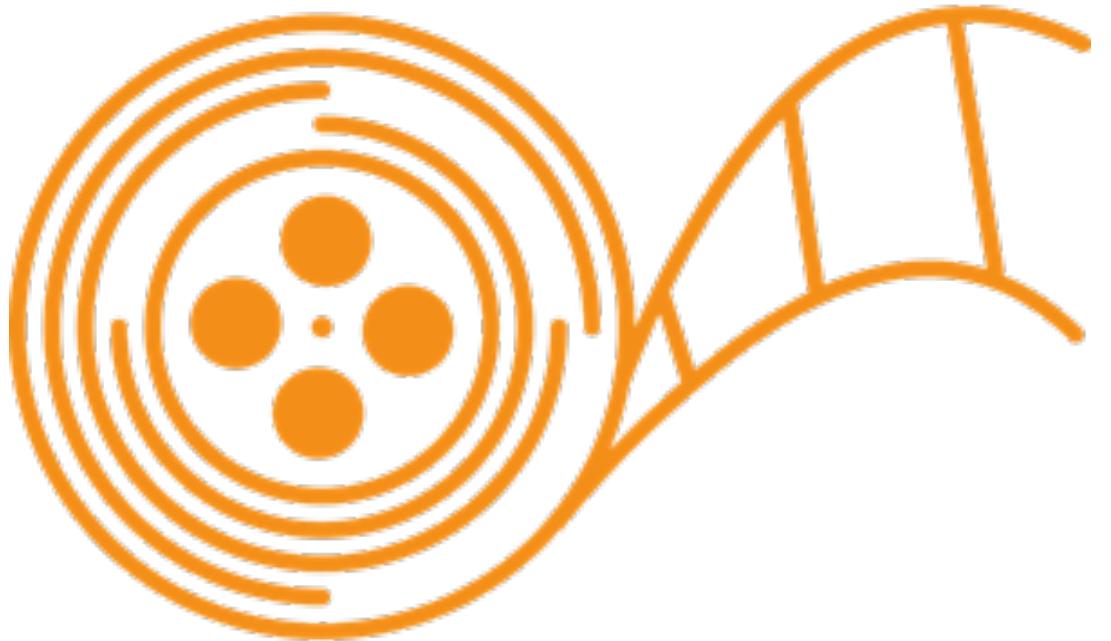


Subnet and Instance Security Controls

 Custom
Public Route
Table
172.16.0.0
172.16.1.0
172.16.2.0

 Custom
Private Route
Table
172.16.0.0
172.16.1.0
172.16.2.0

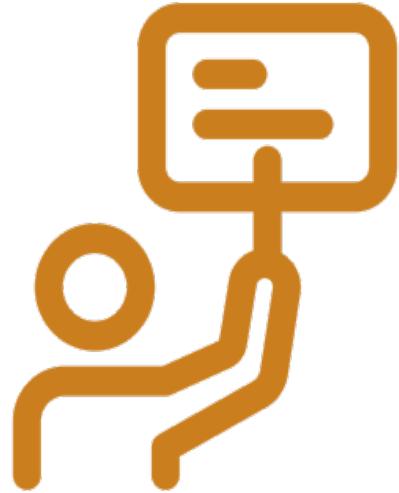




Demo:
Security
Group

Pricing

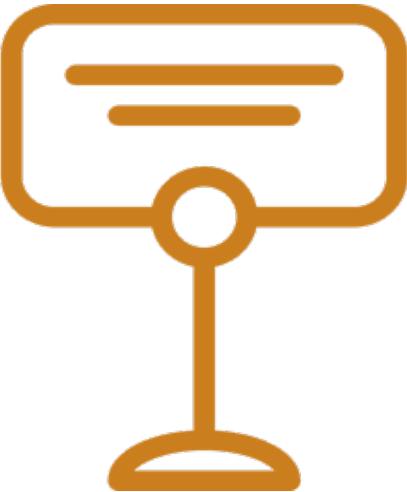
EC2 Instances: Pricing Options



On-Demand

Pay by per hour/
second

Short-term,
unpredictable
workloads



Reserved Instances

Discount for 1 - 3-year
commitment

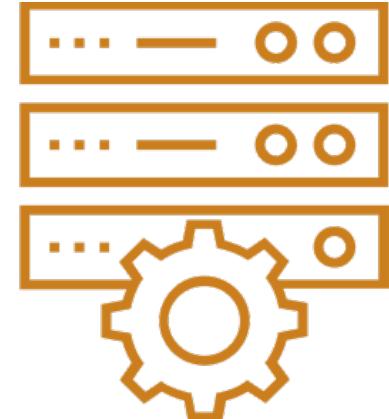
Applications with
consistent usage



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 hundreds of 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 quota



Quotas 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

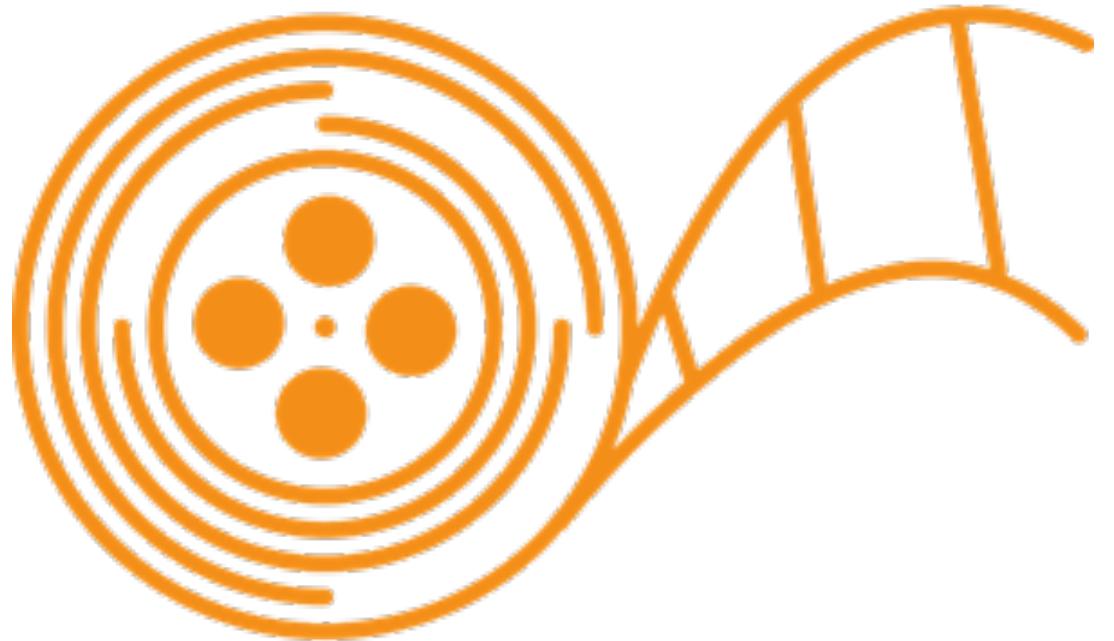
After you've created a capacity reservation, you will be charged for the capacity reservation whether you use the instances or not

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

On-demand Capacity Reservations

Regional Pricing Compared

| Estimate of Your Monthly Bill | | |
|---|--|--------------------------------|
| <input type="checkbox"/> Show First Month's Bill (include all one-time fees, if any) | | |
|  Below you will see an estimate of your monthly bill. Expand each line item to see cost breakout of each service. To save this bill and input values, click on 'Save and Share' button. To remove the service from the estimate, jump back to the service and clear the specific service's form. | | |
| Export to CSV | | Save and Share |
| <input type="checkbox"/> Amazon EC2 Service (US East (N. Virginia)) | | \$ 339.70 |
| Compute: | | \$ 339.70 |
| <input type="checkbox"/> Amazon EC2 Service (South America (Sao Paulo)) | | \$ 3935.30 |
| <input type="checkbox"/> Amazon Elastic Load Balancing (US East (N. Virginia)) | | \$ 16.47 |
| Application LBs: | | \$ 16.47 |
| <input type="checkbox"/> Amazon Elastic Load Balancing (South America (Sao Paulo)) | | \$ 24.89 |



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 EC2 needs your resources back, a two-minute warning is provided (CloudWatch alert)

Spot instance pricing is based on supply and demand; the instance will run until EC2 needs the resources back

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 are stored on the root EBS drive and the assigned 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



Saving Plans

- Savings Plans provides savings of up to 72% on your AWS compute usage
Applies to all Amazon EC2 instances (OS, tenancy or Region) and AWS Fargate and AWS Lambda usage
- You commit to use a specific amount of compute power (measured in \$/hour) for a one, or three-year period

Basic Monitoring

Use CloudWatch to monitor your EC2 instances at no additional charge

CPU load, disk I/O, and network I/O metrics are collected at 5-minute intervals and are stored for two weeks

Detailed Monitoring can also be enabled at 1-minute intervals

Launch Templates

Specify the Amazon machine image (AMI) from which to launch the instances.

→ 1 →

Choose an instance type that is compatible with the AMI you've specified.

→ 2 →

Specify the key pair to use when connecting to instances, for example, using SSH.

3

Add one or more security groups to allow relevant access to the instances from an external network.

→ 4 →

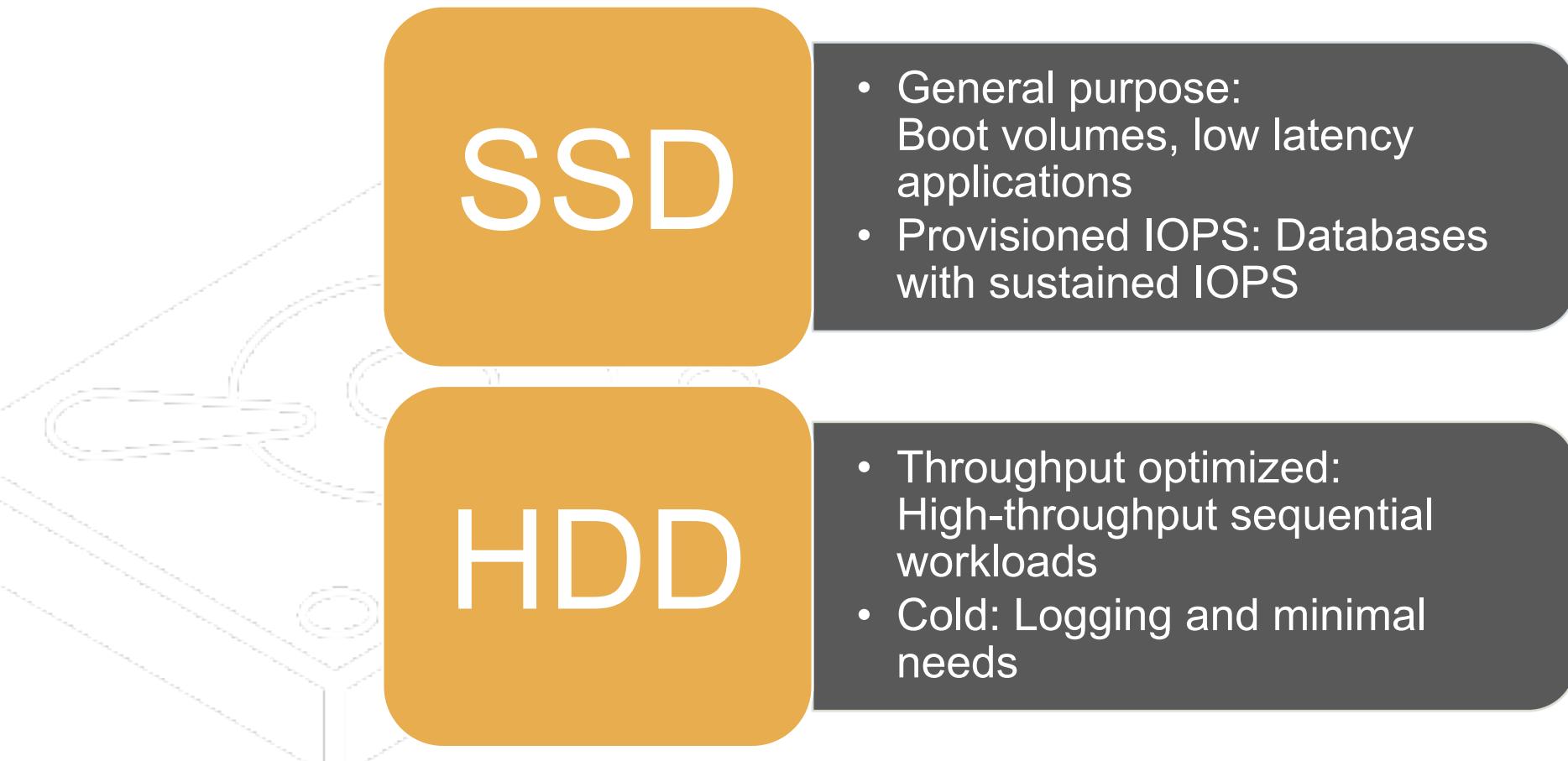
Specify whether to attach additional EBS volumes or instance store volumes to each instance.

→ 5 →

Add custom tags (key-value pairs) to the instances and volumes.

Elastic Block Storage

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

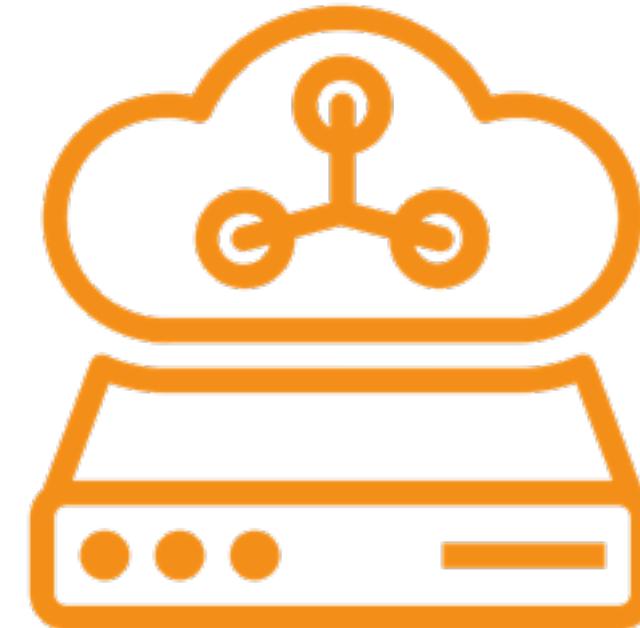
EC2 Instance Stores

- 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
- EBS volumes are replicated with multiple copies within the AZ where the volume 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



* EBS volumes can be shared with up to 16 EC2 instances (Nitro)

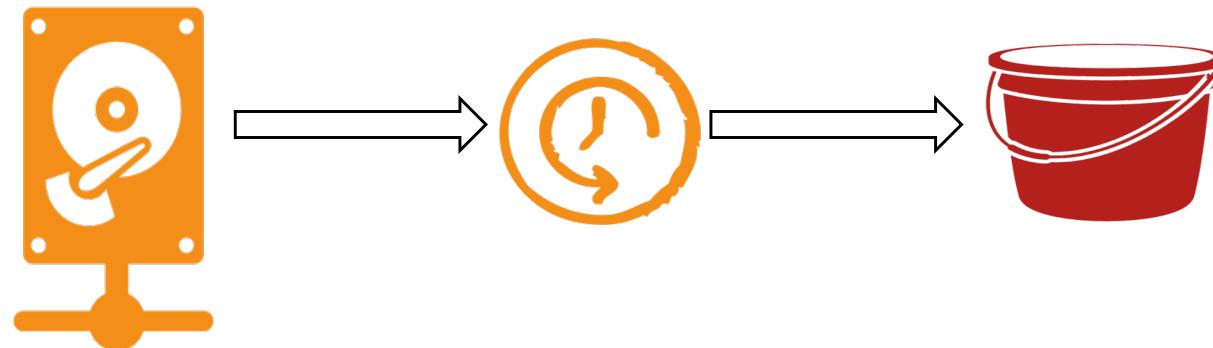
Burst Storage Operation

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



EBS Volume 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 stored objects
 - 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 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
 - 11 9's durability
 - 4 9's durability
 - Minimum storage time of 30 days



S3 Consistency

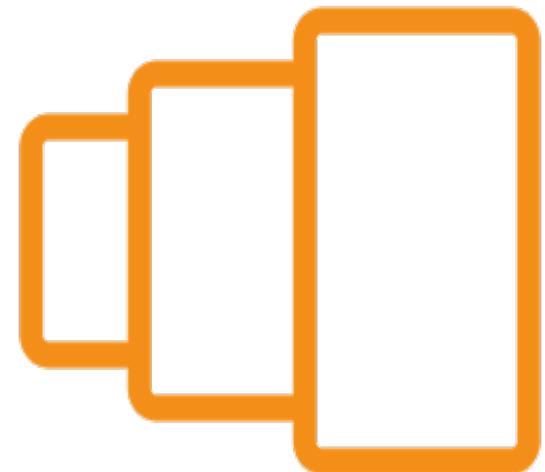
- Objects are eventually consistent
- Multiple copies means replicated storage
- PUT's to new objects – read after write consistency
- PUT's to existing object – eventual consistency



S3 Management

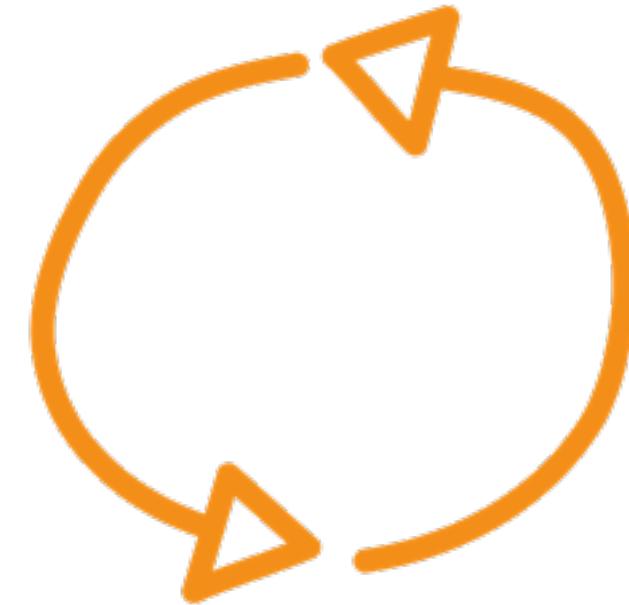
S3 Versioning

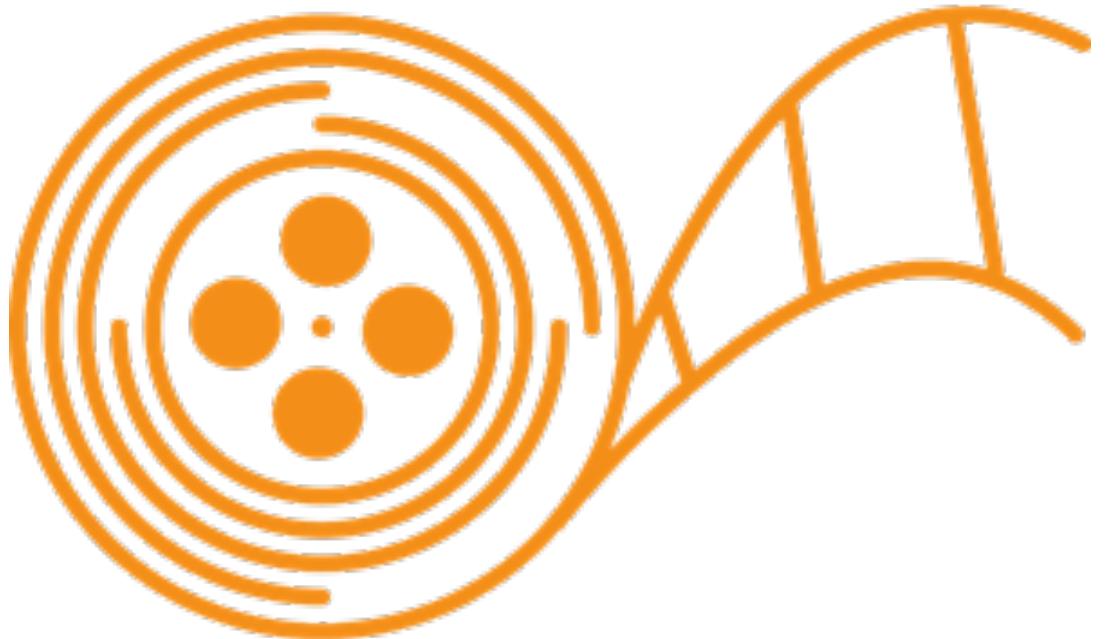
- Versioning is enabled at the bucket level
- Versioning allows you to store multiple versions of the same object in one bucket
- Protect yourself from unintended overwrites or deletions
- Once enabled, versioning can't be disabled but can be suspended



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

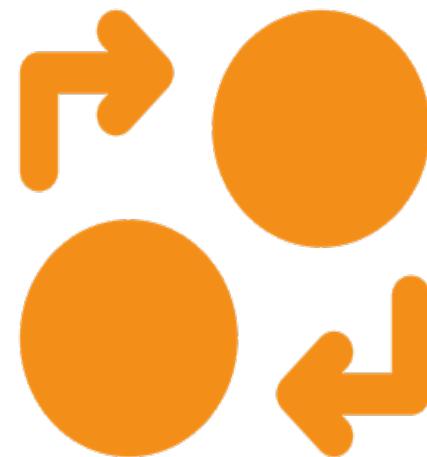




Demo:
Lifecycle
Rules

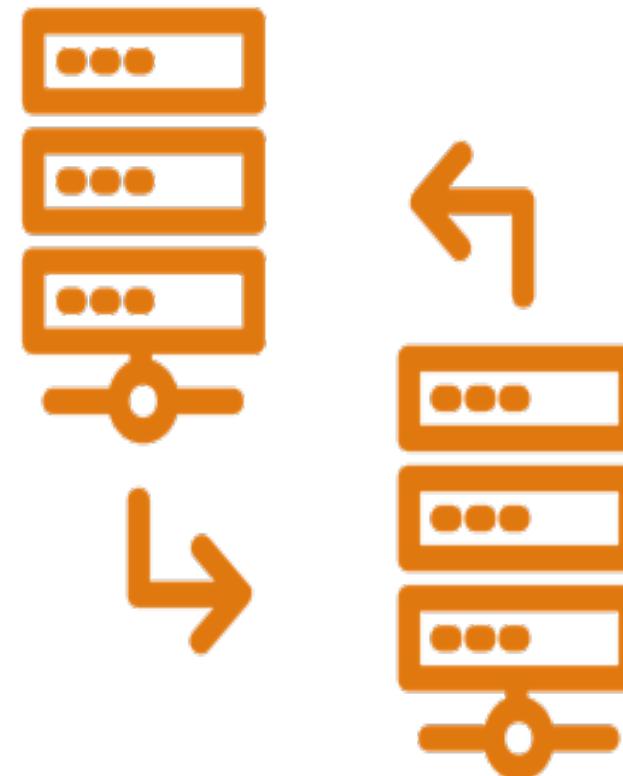
Cross Region Replication (CRR)

- Asynchronous replication from source bucket in AWS region to bucket in another AWS region
- Automatic asynchronous copying of objects
- Versioning must be enabled for both the source and destination buckets
- Separate S3 Lifecycle rules can be configured on both the source and destination buckets
- Helps move data closer to end-users
- Compliance / additional durability



Same Region Replication (SRR)

- Use SRR to replicate objects to a destination bucket within the same AWS region as the source bucket
- Replication is automatic and asynchronous
- Replication can be to any storage class
- Triggers include uploading objects, deleting objects, or changes to the object



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
- Buckets can be associated with different AWS accounts



S3 Encryption Options

- SSE – S3 (AWS Managed Keys)
 - Each object is encrypted with a unique key
 - Encryption key is encrypted with master key.
 - AWS regularly rotate the master key
- SSE-KMS – Server-Side Encryption with AWS KMS keys
 - Can use CMK key, or use your own uploaded key

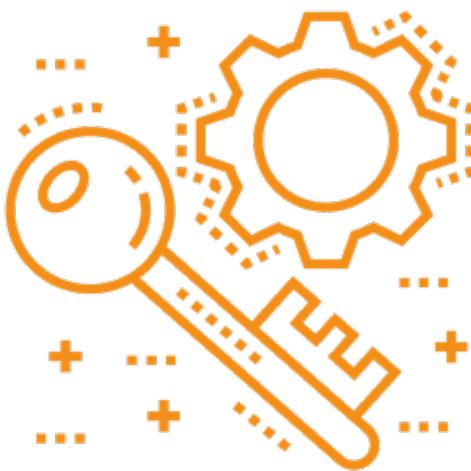


S3 Encryption Options

- SSE-C – Server-Side Encryption with client provided keys
 - Client manages the keys, S3 manages encryption
 - AWS does not store the encryption keys
 - If keys are lost data cannot be decrypted



Key Management Service



AWS KMS – 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

Event Notifications FYI

Notifications can be generated for the following events:

- New object created, object removal or restore events
- Replication events

Event notifications can be sent to the following destinations:

- A Simple Notification Service topic
- Amazon Simple Queue Service queue
- AWS Lambda invoking a Lambda function
- S3 must have permissions to post notifications or to invoke a Lambda function





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 Storage Options

S3 Glacier

- Designed for durability of 99.99999999% of objects in a single Availability Zone in 3 separate facilities
- Designed for 99.99% availability
- Supports SSL for data in transit and encryption of data at rest

S3 Glacier Deep Archive.

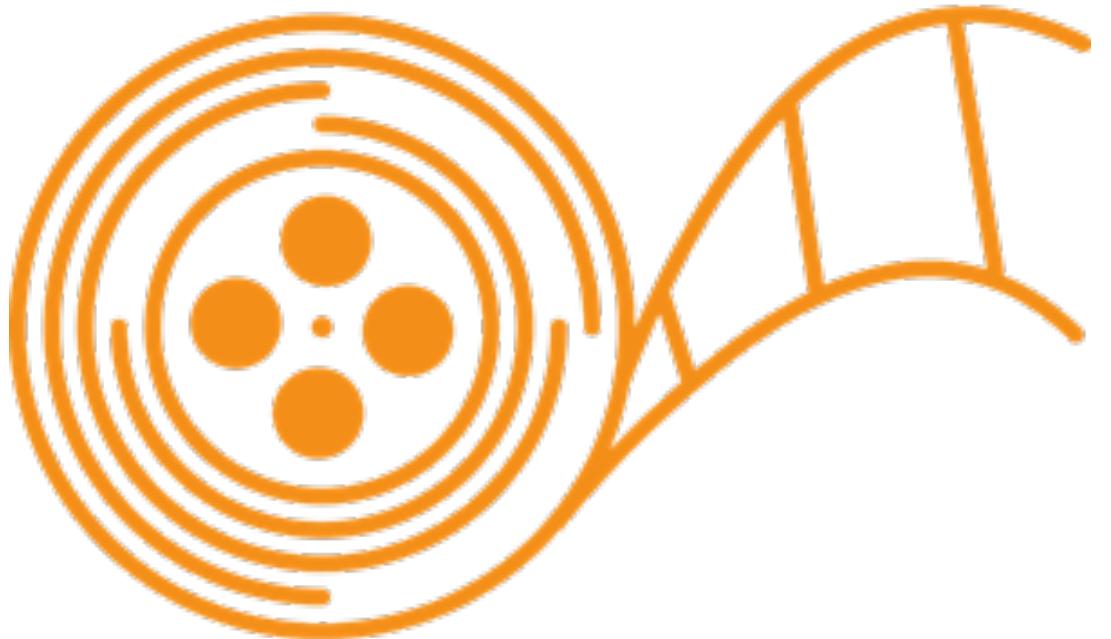
- Designed for durability of 99.99999999% of objects across multiple Availability Zones.
- Data is resilient in the event of one Availability Zone destruction
- Configurable retrieval times, from minutes to hours



S3 Glacier Storage

- Glacier file archives from 100MB up to 40TB can be uploaded to Glacier using the multipart upload API
- 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

