

Load Balancing, Auto Scaling & Availability Zones

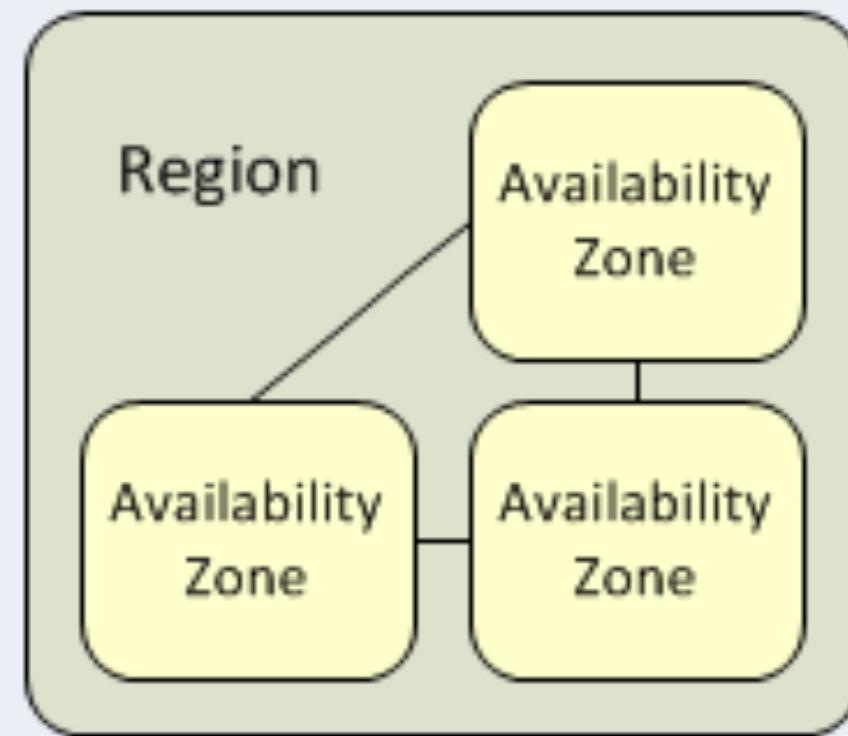
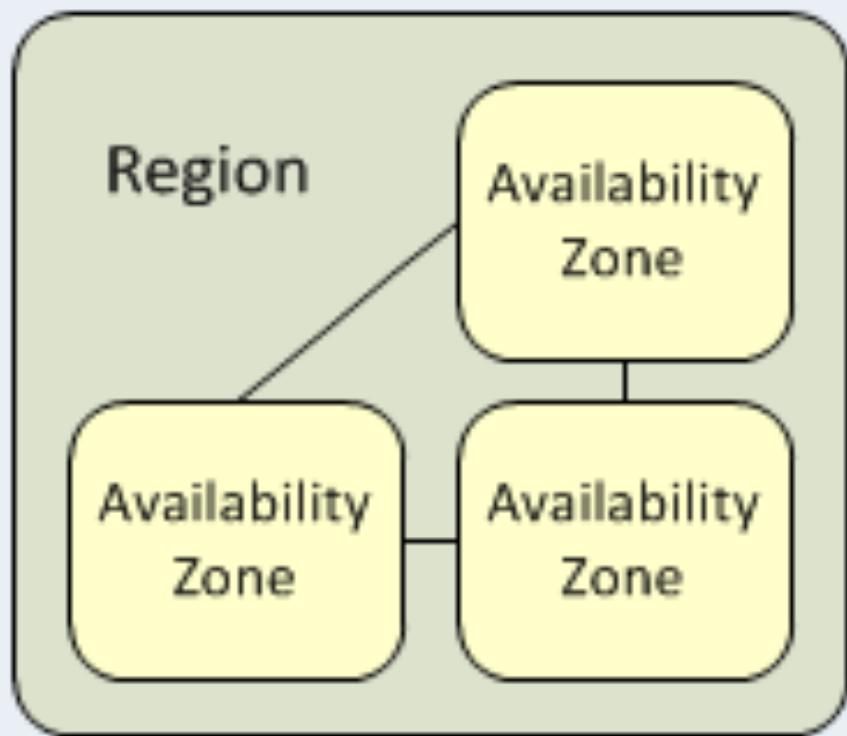
Tejas Parikh (t.parikh@northeastern.edu)

CSYE 6225, Fall 2019

Northeastern University

Regions & Availability Zone

Amazon Web Services



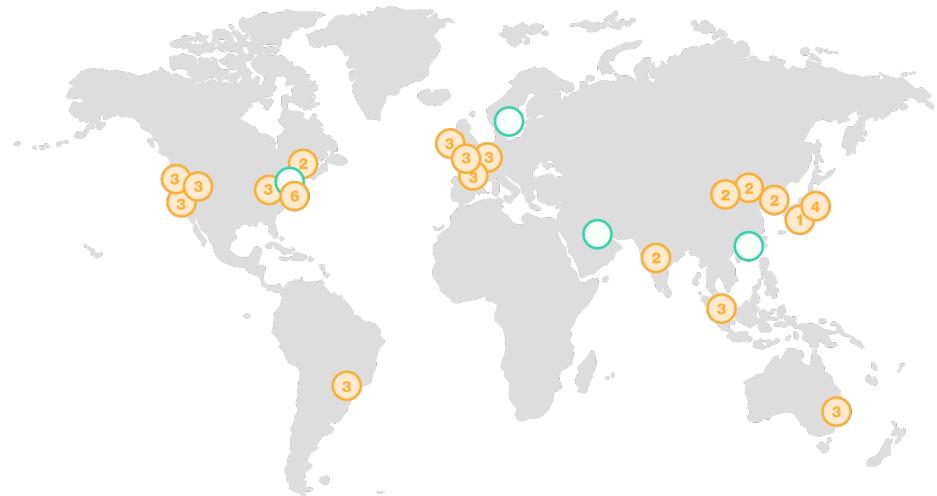
Regions

- Each Amazon EC2 region is designed to be completely isolated from the other Amazon EC2 regions. This achieves the greatest possible fault tolerance and stability.
- When you view your resources, you'll only see the resources tied to the region you've specified. This is because regions are isolated from each other, and AWS does not automatically replicate resources across regions.

Availability Zones

- When you launch an instance, you can select an Availability Zone or let AWS choose one for you. If you distribute your instances across multiple Availability Zones and one instance fails, you can design your application so that an instance in another Availability Zone can handle requests.
- You can also use Elastic IP addresses to mask the failure of an instance in one Availability Zone by rapidly remapping the address to an instance in another Availability Zone.
- An Availability Zone is represented by a region code followed by a letter identifier; for example, us-east-1a. To ensure that resources are distributed across the Availability Zones for a region, AWS independently map Availability Zones to identifiers for each account. For example, your Availability Zone us-east-1a might not be the same location as us-east-1a for another account. There's no way for you to coordinate Availability Zones between accounts.

AWS Regions



Region & Number of Availability Zones

US East	China
N. Virginia (6),	Beijing (2),
Ohio (3)	Ningxia (2)
US West	Europe
N. California (3),	Frankfurt (3),
Oregon (3)	Ireland (3),
Asia Pacific	London (3),
Mumbai (2),	Paris (3)
Seoul (2),	
Singapore (3),	South America
Sydney (3),	São Paulo (3)
Tokyo (4),	
Osaka-Local (1) ¹	AWS GovCloud (US-West) (3)
Canada	
Central (2)	



New Region (coming soon)

Bahrain
Hong Kong
SAR, China
Sweden
AWS GovCloud (US-East)

Code	Name
us-east-1	US East (N. Virginia)
us-east-2	US East (Ohio)
us-west-1	US West (N. California)
us-west-2	US West (Oregon)
ca-central-1	Canada (Central)
eu-central-1	EU (Frankfurt)
eu-west-1	EU (Ireland)
eu-west-2	EU (London)
eu-west-3	EU (Paris)
ap-northeast-1	Asia Pacific (Tokyo)
ap-northeast-2	Asia Pacific (Seoul)
ap-northeast-3	Asia Pacific (Osaka-Local)
ap-southeast-1	Asia Pacific (Singapore)
ap-southeast-2	Asia Pacific (Sydney)
ap-south-1	Asia Pacific (Mumbai)
sa-east-1	South America (São Paulo)

Load Balancing



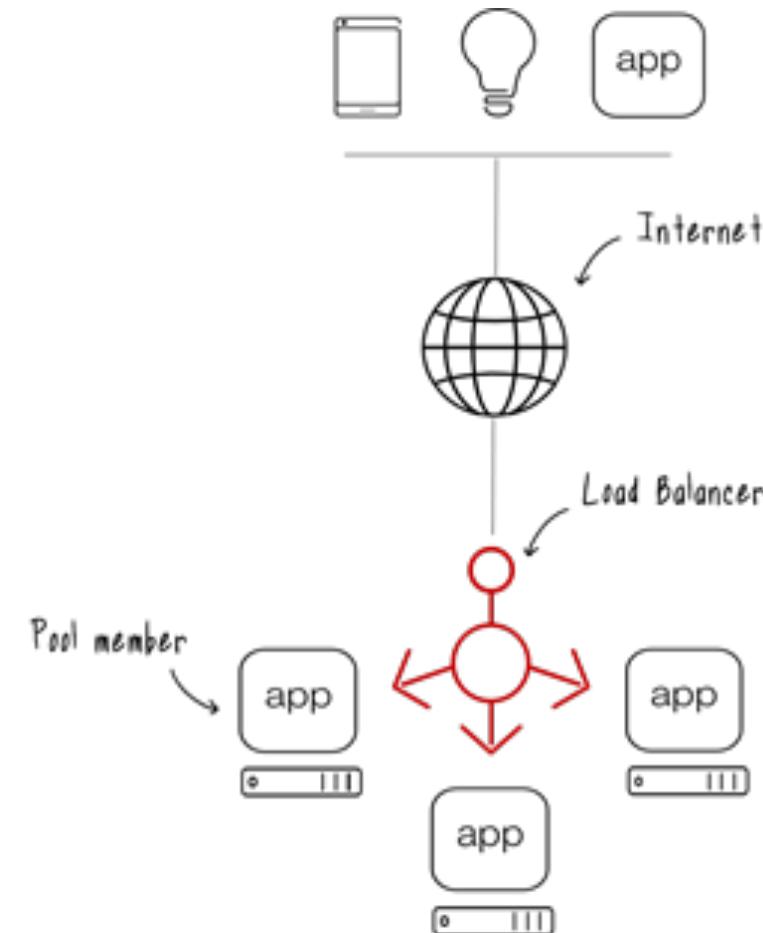
Can you check the load balancer?

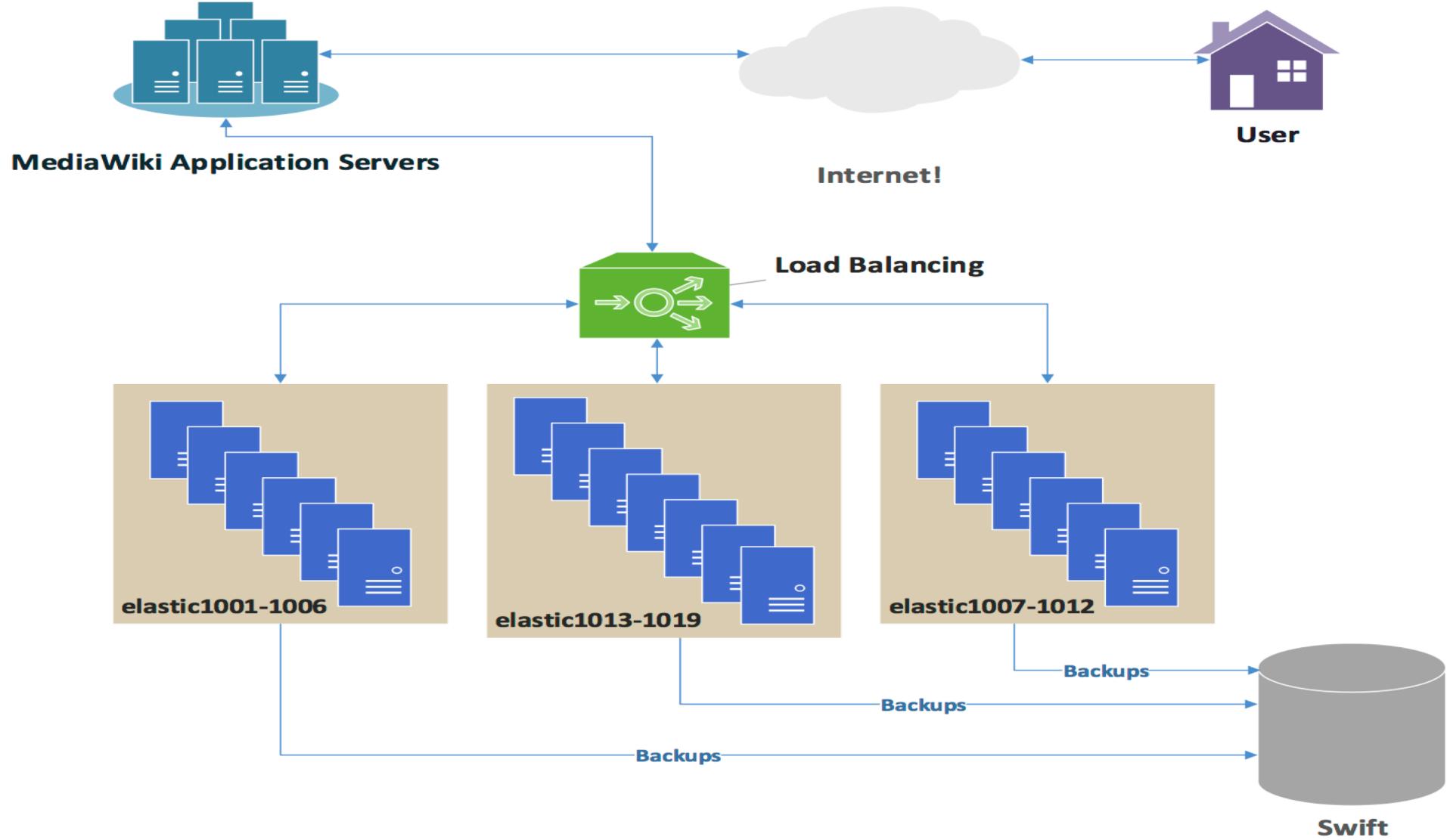
What is Load Balancer?

A load balancer is a device that acts as a reverse proxy and distributes network or application traffic across a number of servers.

Why Use Load Balancers?

- Load balancers are used to increase capacity (concurrent users) and reliability of applications.
- They improve the overall performance of applications by decreasing the burden on servers associated with managing and maintaining application and network sessions, as well as by performing application-specific tasks.
- Load balancers ensure reliability and availability by monitoring the "health" of applications and only sending requests to servers and applications that can respond in a timely manner.





Types of Load Balancers

- **Layer 4 Load Balancers** - Layer 4 load balancers act upon data found in network and transport layer protocols (IP, TCP, FTP, UDP).
- **Layer 7 Load Balancers** - Layer 7 load balancers distribute requests based upon data found in application layer protocols such as HTTP.

OSI Model				
Layer		Protocol data unit (PDU)	Function ^[3]	
Host layers	7	Application	Data	
	6	Presentation	High-level APIs, including resource sharing, remote file access	
	5	Session	Translation of data between a networking service and an application; including character encoding, data compression and encryption/decryption	
	4	Transport	Managing communication sessions, i.e. continuous exchange of information in the form of multiple back-and-forth transmissions between two nodes	
Media layers	3	Network	Segment, Datagram	Reliable transmission of data segments between points on a network, including segmentation, acknowledgement and multiplexing
	2	Data link	Packet	Structuring and managing a multi-node network, including addressing, routing and traffic control
	1	Physical	Frame	Reliable transmission of data frames between two nodes connected by a physical layer
			Symbol	Transmission and reception of raw bit streams over a physical medium

Request Distribution

- Requests are received by both types of load balancers and they are distributed to a particular server based on a configured algorithm. Some industry standard algorithms are:
 - Round robin
 - Weighted round robin
 - Least connections
 - Least response time
- Layer 7 load balancers can further distribute requests based on application specific data such as HTTP headers, cookies, or data within the application message itself, such as the value of a specific parameter.

Classic Load Balancer

AWS Classic Load Balancer

The Classic Load Balancer routes traffic based on application or network level information and is ideal for simple load balancing of traffic across multiple EC2 instances where high availability, automatic scaling, and robust security are required.

AWS Classic Load Balancer Features

- High Availability
- Health Checks
- SSL offloading
- Sticky Sessions
- Layer 4 or 7 Load Balancing
- Access Logs
- Operational Monitoring

Application Load Balancer

AWS Application Load Balancer

An Application Load Balancer is a load balancing option for the Elastic Load Balancing service that operates at the application layer and allows you to define routing rules based on content across multiple services or containers running on one or more Amazon Elastic Compute Cloud (Amazon EC2) instances.

AWS Application Load Balancer Features

- Host-based routing
- Path-based routing
- HTTP/2 Support
- WebSockets Support
- Health Checks
- Layer 7 Load Balancing
- SSL Offloading
- Request Tracking
- Web Application Firewall
- Containerized Application Support

Load Balancer Features

- Asymmetric Load
- Distributed Denial of Service (DDoS) attack protection - load balancers can provide features such as SYN cookies and delayed-binding (the back-end servers don't see the client until it finishes its TCP handshake) to mitigate SYN flood attacks and generally offload work from the servers to a more efficient platform.
- HTTP Compression
- See
[https://en.wikipedia.org/wiki/Load_balancing_\(computing\)#Load_balancer_features](https://en.wikipedia.org/wiki/Load_balancing_(computing)#Load_balancer_features) for more features.

Network Load Balancer

Network Load Balancer

- Network Load Balancer operates at the connection level (Layer 4), routing connections to targets - Amazon EC2 instances, containers and IP addresses based on IP protocol data.
- Ideal for load balancing of TCP traffic, Network Load Balancer is capable of handling millions of requests per second while maintaining ultra-low latencies.
- Network Load Balancer is optimized to handle sudden and volatile traffic patterns while using a single static IP address per Availability Zone.
- It is integrated with other popular AWS services such as Auto Scaling, Amazon EC2 Container Service (ECS), and Amazon CloudFormation.

Network Load Balancer Key Features

- Connection-based Load Balancing - You can load balance TCP traffic, routing connections to targets - Amazon EC2 instances, microservices and containers, and IP addresses.
- Preserve source IP address - Network Load Balancer preserves the client side source IP allowing the back-end to see the IP address of the client.
- Static IP support - Network Load Balancer automatically provides a static IP per Availability Zone (subnet) that can be used by applications as the front-end IP of the load balancer.
- Long-lived TCP Connections - Network Load Balancer supports long-lived TCP connections that are ideal for WebSocket type of applications.
- Central API Support - Network Load Balancer uses the same API as Application Load Balancer. This will enable you to work with target groups, health checks, and load balance across multiple ports on the same Amazon EC2 instance to support containerized applications.

Auto Scaling

What is Autoscaling

- Autoscaling is a method used in cloud computing, whereby the amount of computational resources in a server farm, typically measured in terms of the number of active servers, scales automatically based on the load on the farm.
- It is closely related to, and builds upon, the idea of load balancing.

How does auto-scaling work?

Auto Scaling detects impaired compute instances and unhealthy applications, and replace the instances without human intervention ensuring that your application is getting the compute capacity it needs.

AWS Autoscaling

Auto Scaling will perform three main functions to automate fleet management for EC2 instances:

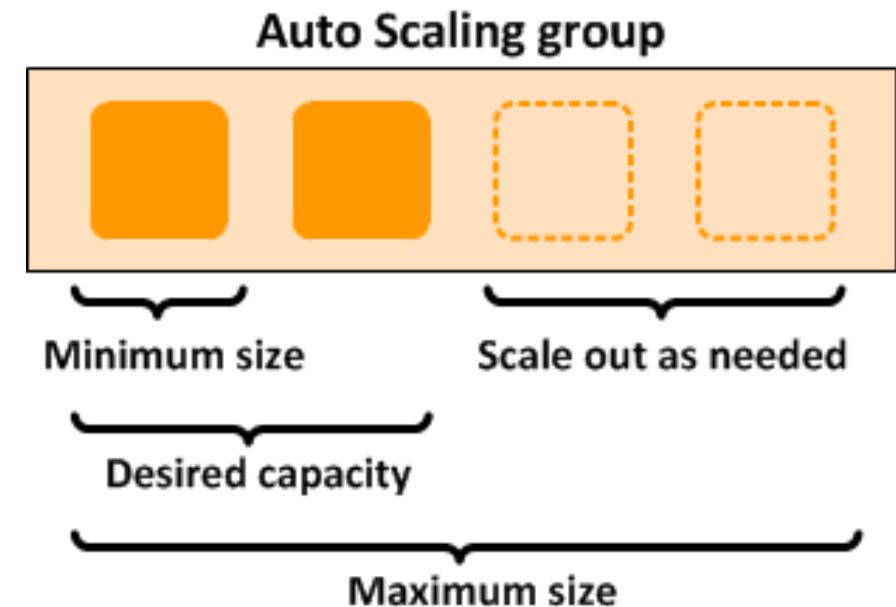
- 1. Monitor the health of running instances** - Amazon EC2 Auto Scaling ensures that your application is able to receive traffic and that EC2 instances are working properly. Amazon EC2 Auto Scaling periodically performs health checks to identify any instances that are unhealthy.
- 2. Replace impaired instances automatically** - When an impaired instance fails a health check, Amazon EC2 Auto Scaling automatically terminates it and replaces it with a new one. That means that you don't need to respond manually when an instance needs replacing.
- 3. Balance capacity across Availability Zones** - Amazon EC2 Auto Scaling can automatically balance instances across zones, and always launches new instances so that they are balanced between zones as evenly as possible across your entire fleet.

Dynamic Scaling

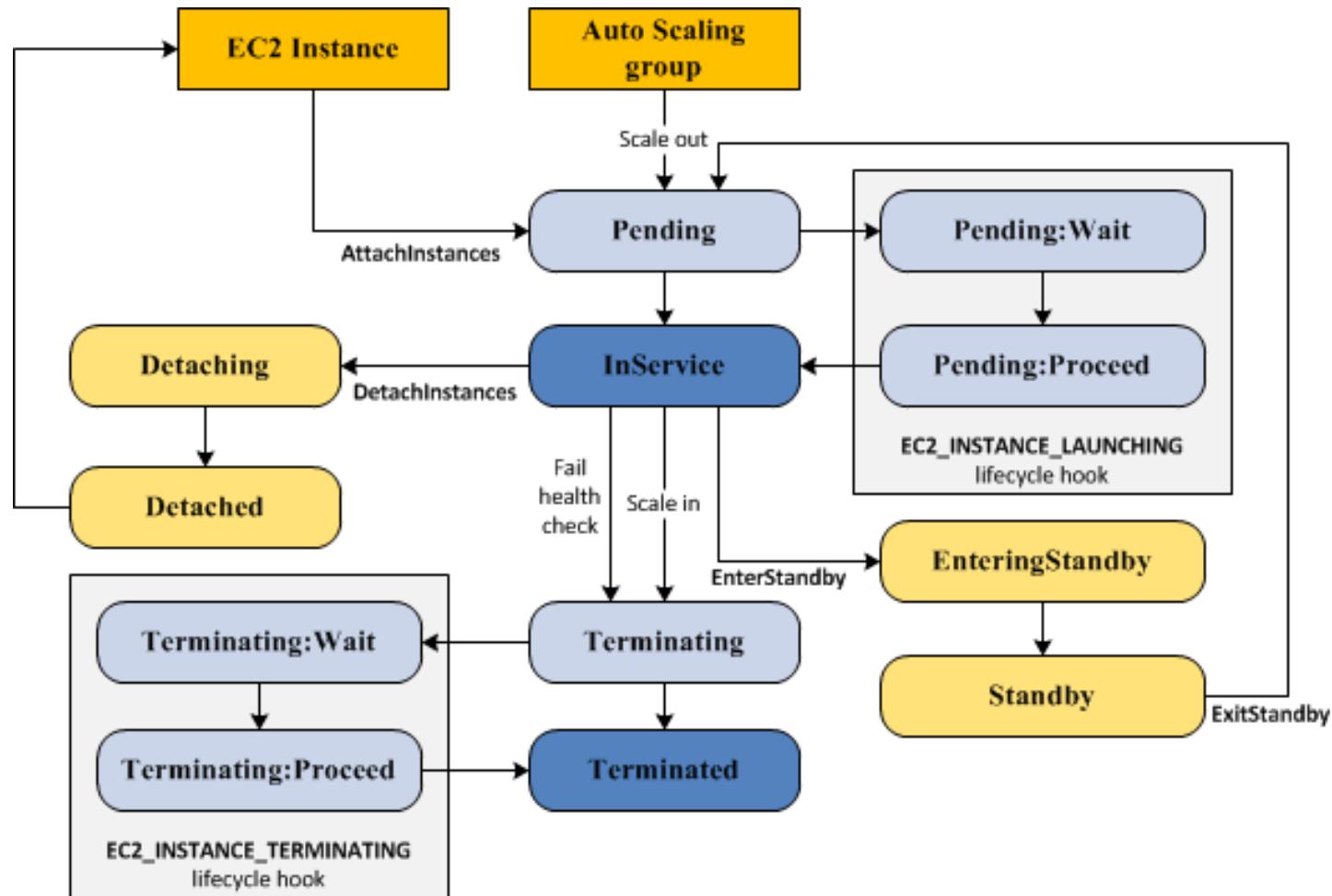
- Amazon EC2 Auto Scaling enables you to follow the demand curve for your applications closely, reducing the need to manually provision Amazon EC2 capacity in advance. For example, you can use target tracking scaling policies to select a load metric for your application, such as CPU utilization. Or, you could set a target value using the new “Request Count Per Target” metric from Application Load Balancer, a load balancing option for the Elastic Load Balancing service. Amazon EC2 Auto Scaling will then automatically adjust the number of EC2 instances as needed to maintain your target.
- You can also use simple scaling policies to set a condition to add new Amazon EC2 instances in increments when the average utilization of your Amazon EC2 fleet is high, and similarly, you can set a condition to remove instances in the same increments when CPU utilization is low. If you have predictable load changes, you can also set a schedule through Amazon EC2 Auto Scaling to plan your scaling activities.
- Amazon EC2 Auto Scaling can also be used with Amazon CloudWatch, which can send alarms to trigger scaling activities, and Elastic Load Balancing to help distribute traffic to your instances within EC2 Auto Scaling groups.
- You can also use Amazon EC2 Auto Scaling in combination with AWS Auto Scaling to scale multiple services.

Auto Scaling Groups

- You create collections of EC2 instances, called Auto Scaling groups.
- You can specify the minimum number of instances in each Auto Scaling group, and Auto Scaling ensures that your group never goes below this size.
- You can specify the maximum number of instances in each Auto Scaling group, and Auto Scaling ensures that your group never goes above this size.



Auto Scaling Lifecycle



Scaling Policies

- **Manual scaling** - Manual scaling is the most basic way to scale your resources. Specify only the change in the maximum, minimum, or desired capacity of your Auto Scaling group. Amazon EC2 Auto Scaling manages the process of creating or terminating instances to maintain the updated capacity.
- **Scheduled Scaling** - Sometimes you know exactly when you will need to increase or decrease the number of instances in your group, simply because that need arises on a predictable schedule. Scaling by schedule means that scaling actions are performed automatically as a function of time and date.
- **Dynamic Scaling** - A more advanced way to scale your resources, scaling by policy, lets you define parameters that control the scaling process. For example, you can create a policy that calls for enlarging your fleet of EC2 instances whenever the average CPU utilization rate stays above ninety percent for fifteen minutes. This is useful when you can define how you want to scale in response to changing conditions, but you don't know when those conditions will change.

Scale Out

- The following scale out events direct the Auto Scaling group to launch EC2 instances and attach them to the group:
 - You manually increase the size of the group.
 - You create a scaling policy to automatically increase the size of the group based on a specified increase in demand.
 - You set up scaling by schedule to increase the size of the group at a specific time.
- When a scale out event occurs, the Auto Scaling group launches the required number of EC2 instances, using its assigned launch configuration.
- When each instance is fully configured and passes the Amazon EC2 health checks, it is attached to the Auto Scaling group and it enters the InService state. The instance is counted against the desired capacity of the Auto Scaling group.

Scale In

- The following scale in events direct the Auto Scaling group to detach EC2 instances from the group and terminate them:
 - You manually decrease the size of the group.
 - You create a scaling policy to automatically decrease the size of the group based on a specified decrease in demand.
 - You set up scaling by schedule to decrease the size of the group at a specific time.
- When a scale in event occurs, the Auto Scaling group detaches one or more instances. The Auto Scaling group uses its termination policy to determine which instances to terminate. Instances that are in the process of detaching from the Auto Scaling group and shutting down enter the Terminating state, and can't be put back into service. If you add a lifecycle hook to your Auto Scaling group, you can perform a custom action here. Finally, the instances are completely terminated and enter the Terminated state.

Scaling Cooldowns

- The cooldown period is a configurable setting for your Auto Scaling group that helps to ensure that it doesn't launch or terminate additional instances before the previous scaling activity takes effect.
- After the Auto Scaling group dynamically scales using a simple scaling policy, it waits for the cooldown period to complete before resuming scaling activities.
- Cooldown periods are not supported for scheduled scaling.

Content Delivery Network (CDN)

What is CDN?

- A content delivery network or content distribution network (CDN) is a geographically distributed network of proxy servers and their data centers.
- CDN is an umbrella term spanning different types of content delivery services: video streaming, software downloads, web and mobile content acceleration, licensed/managed CDN, transparent caching, and services to measure CDN performance, load balancing, multi-CDN switching and analytics and cloud intelligence.
- A CDN allows for the quick transfer of assets needed for loading Internet content including HTML pages, JavaScript files, stylesheets, images, and videos.
- The popularity of CDN services continues to grow, and today the majority of web traffic is served through CDNs, including traffic from major sites like Facebook, Netflix, and Amazon.

What are the benefits of using a CDN?

- **Improving website load times** - By distributing content closer to website visitors by using a nearby CDN server (among other optimizations), visitors experience faster page loading times. As visitors are more inclined to click away from a slow-loading site, a CDN can reduce bounce rates and increase the amount of time that people spend on the site. In other words, a faster website means more visitors will stay and stick around longer.
- **Reducing bandwidth costs** - Bandwidth consumption costs for website hosting is a primary expense for websites. Through caching and other optimizations, CDNs are able to reduce the amount of data an origin server must provide, thus reducing hosting costs for website owners.
- **Increasing content availability and redundancy** - Large amounts of traffic or hardware failures can interrupt normal website function. Thanks to their distributed nature, a CDN can handle more traffic and withstand hardware failure better than many origin servers.
- **Improving website security** - A CDN may improve security by providing DDoS mitigation, improvements to security certificates, and other optimizations.

- 1 vrid-225.core-sw.aus.us.siteprotect.com (216.139.225.1) 0.627 ms
- 2 xe-3-4.brdr-rtr-02.aus.us.siteprotect.com (216.139.253.53) 0.219 ms
- 3 66.113.197.121 0.452 ms
- 4 xe-5-2-0.edge3.Dallas1.Level3.net (4.59.112.37) 4.978 ms
- 5 ae-73-70.ebr3.Dallas1.Level3.net (4.69.145.116) 9.817 ms
- 6 ae-7-7.ebr3.Atlanta2.Level3.net (4.69.134.22) 30.570 ms
- 7 ae-2-2.ebr1.Washington1.Level3.net (4.69.132.86) 38.801 ms
- 8 ae-81-81.csw3.Washington1.Level3.net (4.69.134.138) 41.795 ms
- 9 ae-3-89.edge2.Washington1.Level3.net (4.68.17.145) 39.193 ms
- 10 72.21.222.139 35.767 ms



CDN Usecases

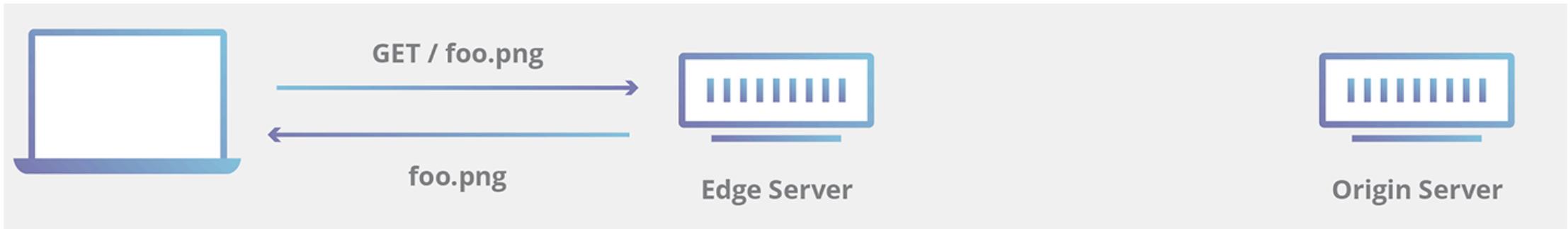
- Static Asset Caching
- Live and On-Demand Video Streaming
- Security and DDoS Protection
- API Acceleration
- Software Distribution

What is an Origin Server?

- The purpose of an origin server is to process and respond to incoming Internet requests from Internet clients.
- The concept of an origin server is typically used in conjunction with the concept of an edge server or caching server.
- At its core, an origin server is a computer running one or more programs that are designed to listen for and process incoming Internet requests.
- An origin server can take on all the responsibility of serving up the content for an Internet property such as a website, provided that the traffic does not extend beyond what the server is capable of processing and latency is not a primary concern.

What is a CDN edge server?

- A CDN edge server is a computer that exists at the logical extreme or “edge” of a network. An edge server often serves as the connection between separate networks.
- A primary purpose of a CDN edge server is to store content as close as possible to a requesting client machine, thereby reducing latency and improving page load times.



What is Anycast?

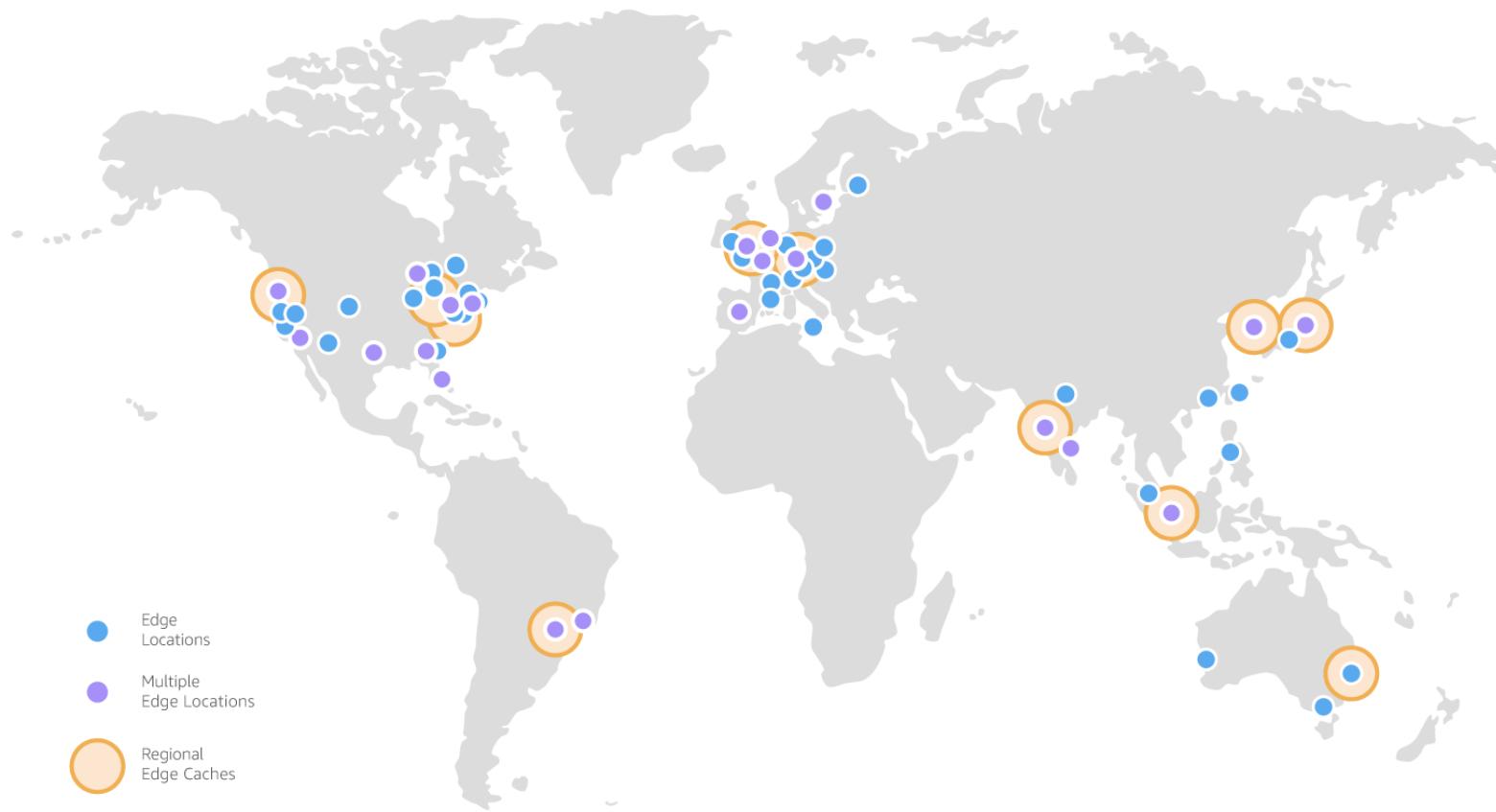
- Anycast is a network addressing and routing method in which incoming requests can be routed to a variety of different locations or “nodes.”
- In the context of a CDN, Anycast typically routes incoming traffic to the nearest data center with the capacity to process the request efficiently.
- Selective routing allows an Anycast network to be resilient in the face of high traffic volume, network congestion, and DDoS attacks.



How CDN Works

- At its core, a CDN is a network of servers linked together with the goal of delivering content as quickly, cheaply, reliably, and securely as possible.
- In order to improve speed and connectivity, a CDN will place servers at the exchange points between different networks.
- These Internet exchange points (IXPs) are the primary locations where different Internet providers connect in order to provide each other access to traffic originating on their different networks.
- By having a connection to these high speed and highly interconnected locations, a CDN provider is able to reduce costs and transit times in high speed data delivery.

The Amazon CloudFront Global Edge Network



CloudFront Regional Edge Cache

A default CloudFront feature that brings more of your content close to your viewers, even when the content is not popular to stay at a particular edge location. This helps improve performance for your viewers while lowering the operational burden and cost of scaling your origin resources.

Type of HTTP requests are supported by Amazon CloudFront

- GET
- HEAD
- POST
- PUT
- PATCH
- DELETE
- OPTIONS

Note: Amazon CloudFront does not cache the responses to POST, PUT, DELETE, and PATCH requests – these requests are proxied back to the origin server.

Latency - How does a CDN improve website load times?

When it comes to websites loading content, users drop off quickly as a site slows down. CDN services can help to reduce load times in the following ways:

- The globally distributed nature of a CDN means reduce distance between users and website resources. Instead of having to connect to wherever a website's origin server may live, a CDN lets users connect to a geographically closer data center. Less travel time means faster service.
- CDNs can reduce the amount of data that's transferred by reducing file sizes using tactics such as minification and file compression. Smaller file sizes mean quicker load times.
- CDNs can also speed up sites which use TLS/SSL certificates by optimizing connection reuse and enabling TLS false start.

Reliability and Redundancy - How does a CDN keep a website always online?

Uptime is a critical component for anyone with an Internet property. Hardware failures and spikes in traffic, as a result of either malicious attacks or just a boost in popularity, have the potential to bring down a web server and prevent users from accessing a site or service. A well-rounded CDN has several features that will minimize downtime:

- Load balancing distributes network traffic evenly across several servers, making it easier to scale rapid boosts in traffic.
- Intelligent failover provides uninterrupted service even if one or more of the CDN servers go offline due to hardware malfunction; the failover can redistribute the traffic to the other operational servers.
- In the event that an entire data center is having technical issues, Anycast routing transfers the traffic to another available data center, ensuring that no users lose access to the website.

Additional Resources

<https://fall2019.csye6225.cloud/>