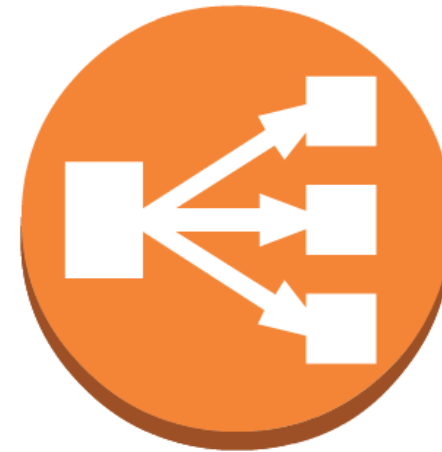




AWS ELB & Auto Scaling

Agenda

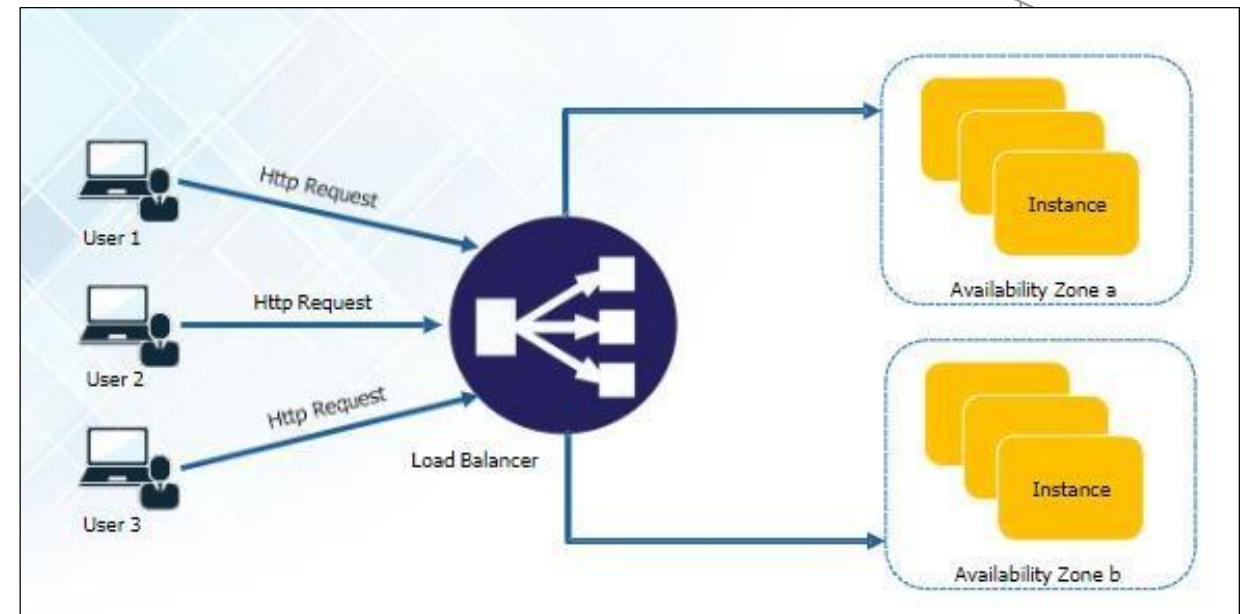
- ❖ What is AWS ELB
- ❖ Classic Load Balancer
 - Features
 - Health Check Configuration
 - Cross-Zone
 - Connection Draining
 - Sticky Sessions
 - Access Logs
 - Limitation
- ❖ Application Load Balancer
 - What is Application ELB
 - Features
 - Application Flow
 - Limitation
- ❖ Network Load Balancer
 - What is Network ELB
- ❖ AWS Autoscaling
- ❖ Hands-On Lab



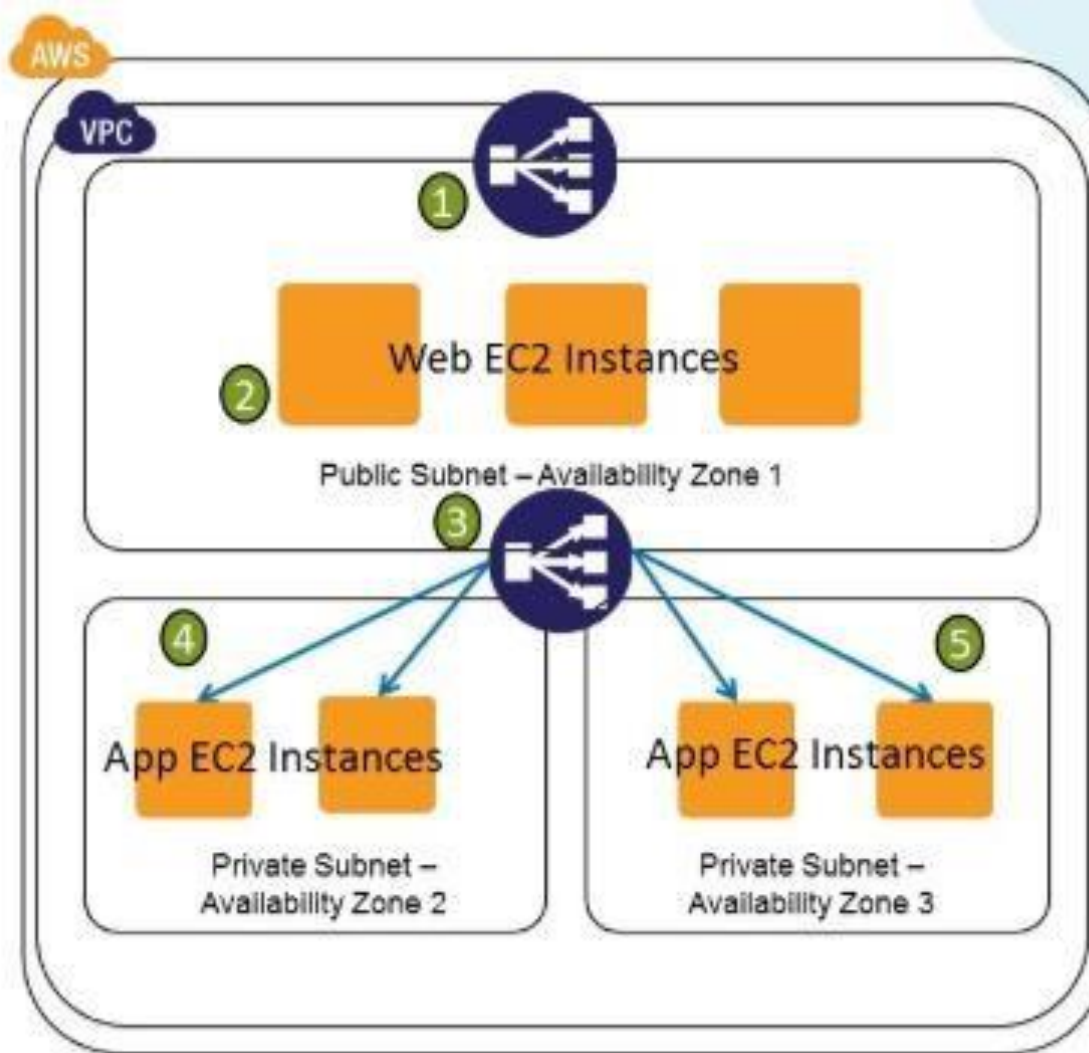
AWS Elastic Load Balancer

What is AWS ELB

- ❑ ELB increases the fault tolerance of your applications.
- ❑ The load balancer serves as a single point of contact for clients.
- ❑ Enable health checks.
- ❑ Types of load balancers:
 - Application Load Balancers
 - Network Load Balancers
 - Classic Load Balancers



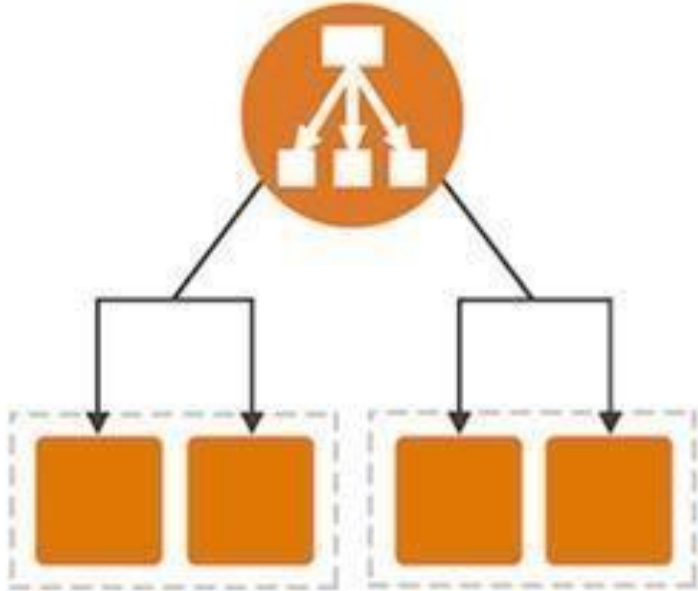
External and Internal Load Balancer



- 1 Front end Elastic Load Balancer exposed to internet accepting web requests
- 2 Scalable web EC2 Instances behind the ELB
- 3 Internal Elastic Load Balancer load balancing only the backend app Instances. Not exposed to internet
- 4 App tier EC2 Instances in AZ1
- 5 App tier EC2 Instances in AZ2

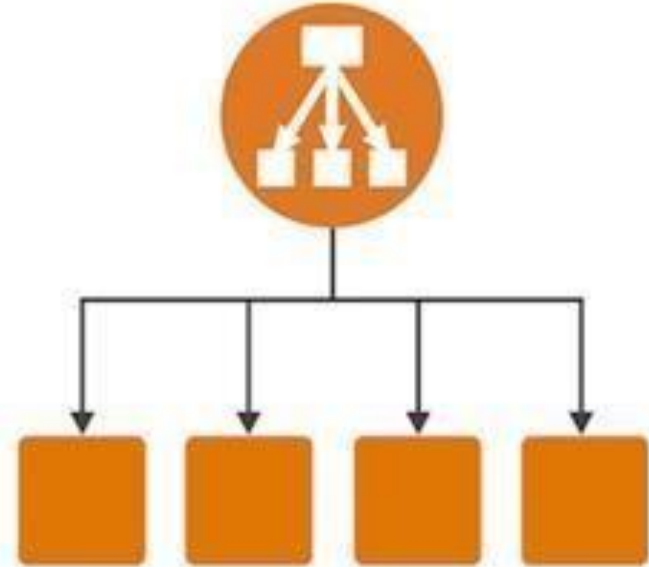
AWS Load BalancerTypes

Application load balancer



An Application load balancer makes routing decisions at the application layer (HTTP/HTTPS), supports path-based routing, and can route requests to one or more ports on each EC2 instance or container instance in your VPC.

Classic load balancer



A Classic load balancer makes routing decisions at either the transport layer (TCP/SSL) or the application layer (HTTP/HTTPS), and supports either EC2-Classic or a VPC.

AWSELB:Features

**Availability
Zone**

Cross-Zone

**Request
Routing**

**Connection
Draining**

**Internet-facing Load
Balancer**

**Internal
Load Balancer**

**Pay-Only
WhatYou Use**

AWSELB:Health Check Configuration

Ping Protocol

**Ping
Port**

**Ping
Path**

**Response
Timeout**

**HealthCheck
Interval**

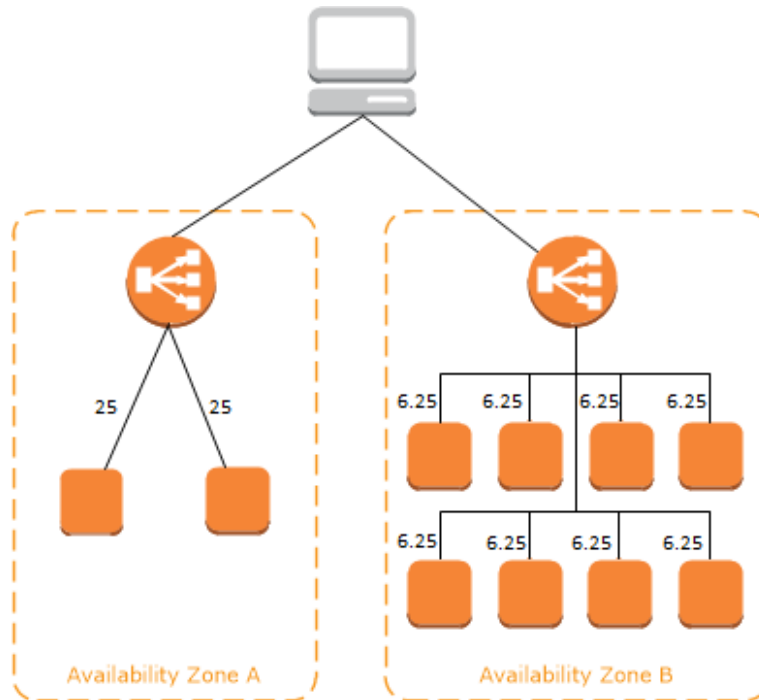
**Unhealthy
Threshold**

**Healthy
Threshold**

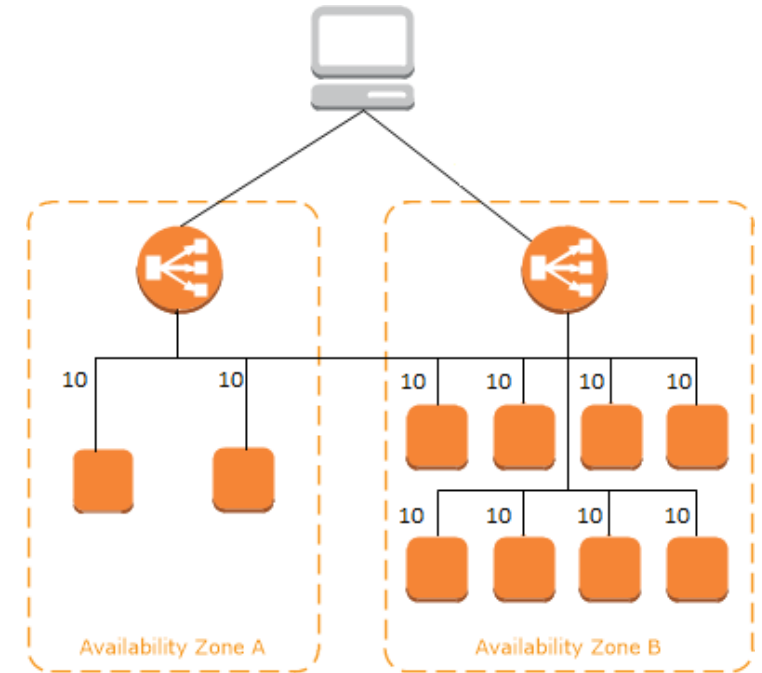
© 2018 CHAITANYA GAJULA ALL RIGHTS RESERVED

AWSELB:Cross-Zone

Cross-Zone Disabled



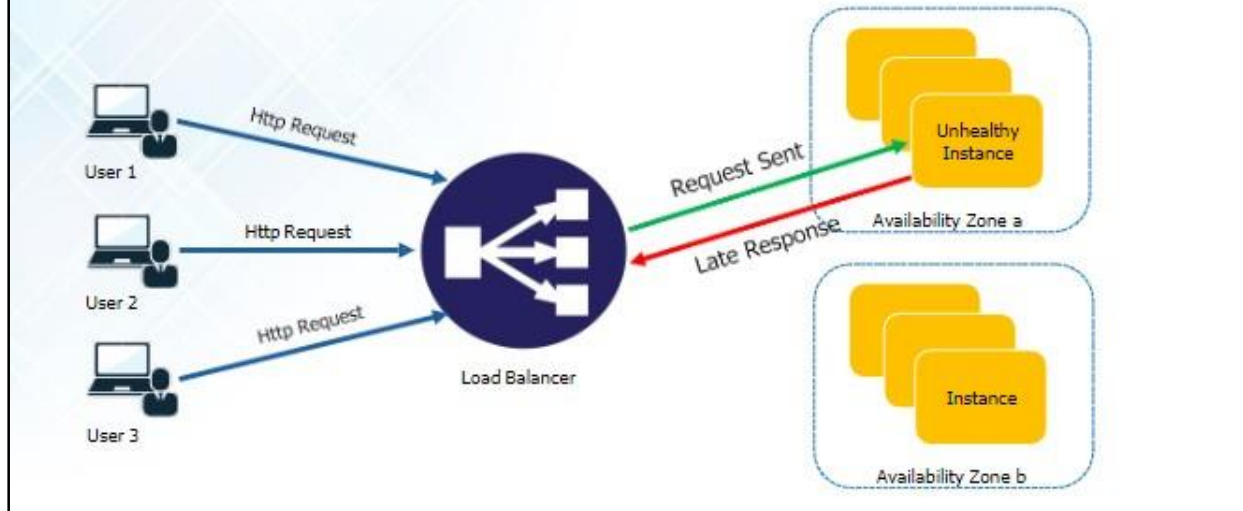
Cross-Zone Enabled



- ❑ Cross-zone load balancing distribute incoming requests evenly across the Availability Zones enabled for your load balancer.
 - Example, if you have 2 instances in Availability Zone us-west-2a and 10 instances in us-west-2b, the requests are distributed evenly across all 12 instances if cross-zone load balancing is enabled.
 - Otherwise, the 2 instances in us-west-2b serve the same number of requests as the 10 instances in us-west-2a.

AWSELB:ConnectionDraining

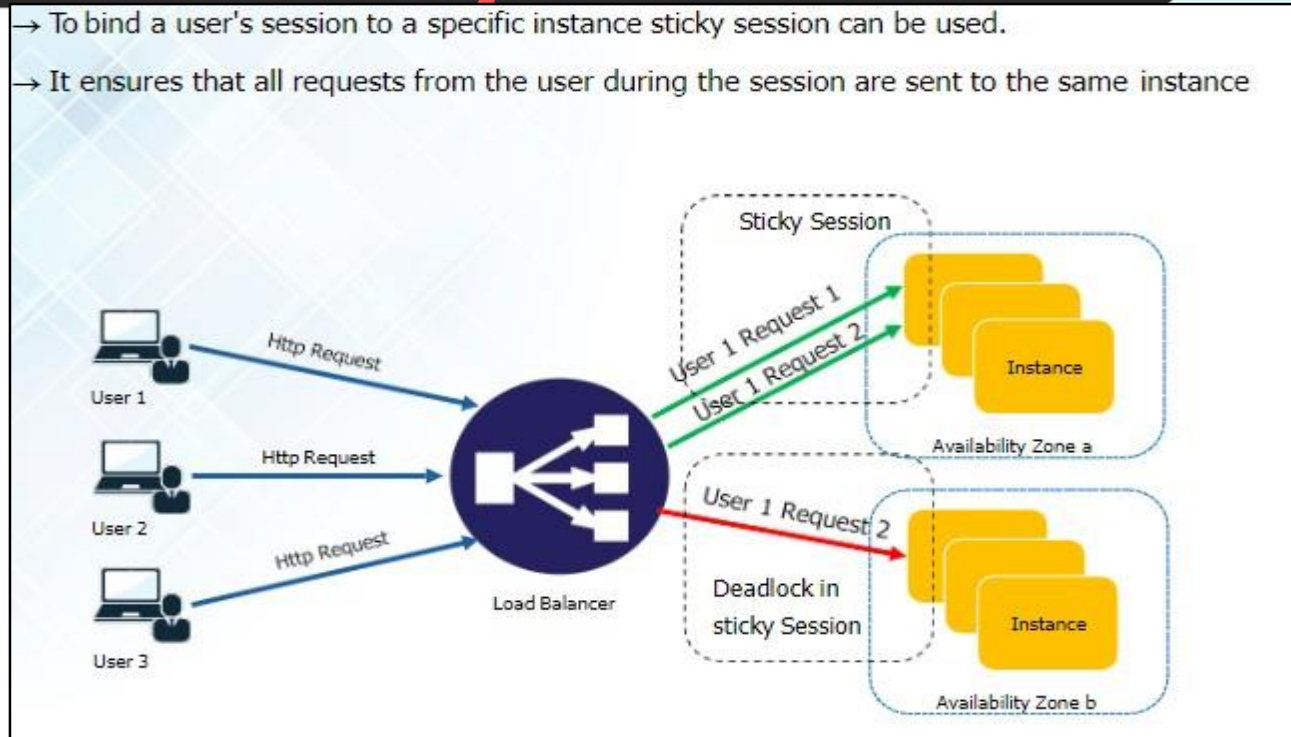
- Connection Draining is used to ensure that a Load Balancer stops sending requests to instances that are de-registering or unhealthy, while keeping the existing connections open
- When connection draining is enabled, a maximum time can be specified for the load balancer to keep connections alive, before reporting the instance as unhealthy



- ❑ Complete in-flight requests made to instances that are de-registering or unhealthy.
- ❑ Specify a maximum time for the load balancer to keep connections alive
- ❑ State:
 - InService: Instance deregistration currently inprogress
 - OutOfService: Instance isnot currently registered with theLoadBalancer

AWSELB: Sticky Sessions

- To bind a user's session to a specific instance sticky session can be used.
- It ensures that all requests from the user during the session are sent to the same instance



- You can use the *sticky session* feature (also known as *session affinity*), which enables the load balancer to bind a user's session to a specific instance.
- The stickiness policy configuration defines a cookie expiration, which establishes the duration of validity for each cookie.
- After a cookie expires, the session is no longer sticky.
- If an instance fails or becomes unhealthy, the load balancer stops routing requests to that instance.

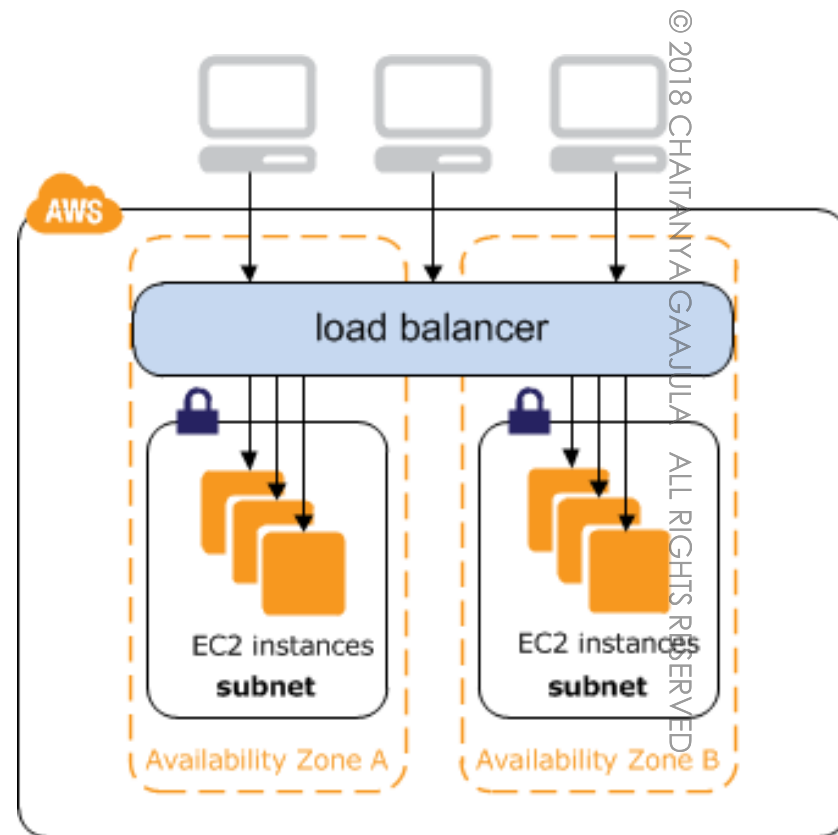
AWSELB: Access Logs

- ❑ ELB provides access logs that capture detailed information about requests sent to ELB.
- ❑ Each log contains information such as time the request, the client's IP address, etc.
- ❑ ELB captures the logs and stores them in the Amazon S3 bucket.
- ❑ You can use these access logs to troubleshoot issues.
- ❑ Access logging is an optional feature of Elastic Load Balancing that is disabled by default.
- ❑ Syntax
 - Each log entry contains the details of a single request made to the load balancer.
 - timestamp elb client:port backend:port request_processing_time backend_processing_time response_processing_time elb_status_code backend_status_code received_bytes sent_bytes "request" "user_agent" ssl_cipher ssl_protocol

AWS Classic ELB: Limitation

Resource	Default Limit
Load balancers per region	20
Listeners per load balancer	100
Security groups per load balancer	5
Subnets per Availability Zone per load balancer	1

Application Load Balancer



What is AWS Application ELB

- ❑ An Application Load Balancer functions at Interconnection (OSI) model.
- ❑ It evaluates the listener rules to determine which rule to apply, and then selects a target from the target group for the rule action.
- ❑ Configure listener rules to route requests to different target groups based on the content of the application traffic.
- ❑ Configure health checks, which are used to monitor the health of the registered targets.
- ❑ Listeners support the HTTP/HTTPS protocols.

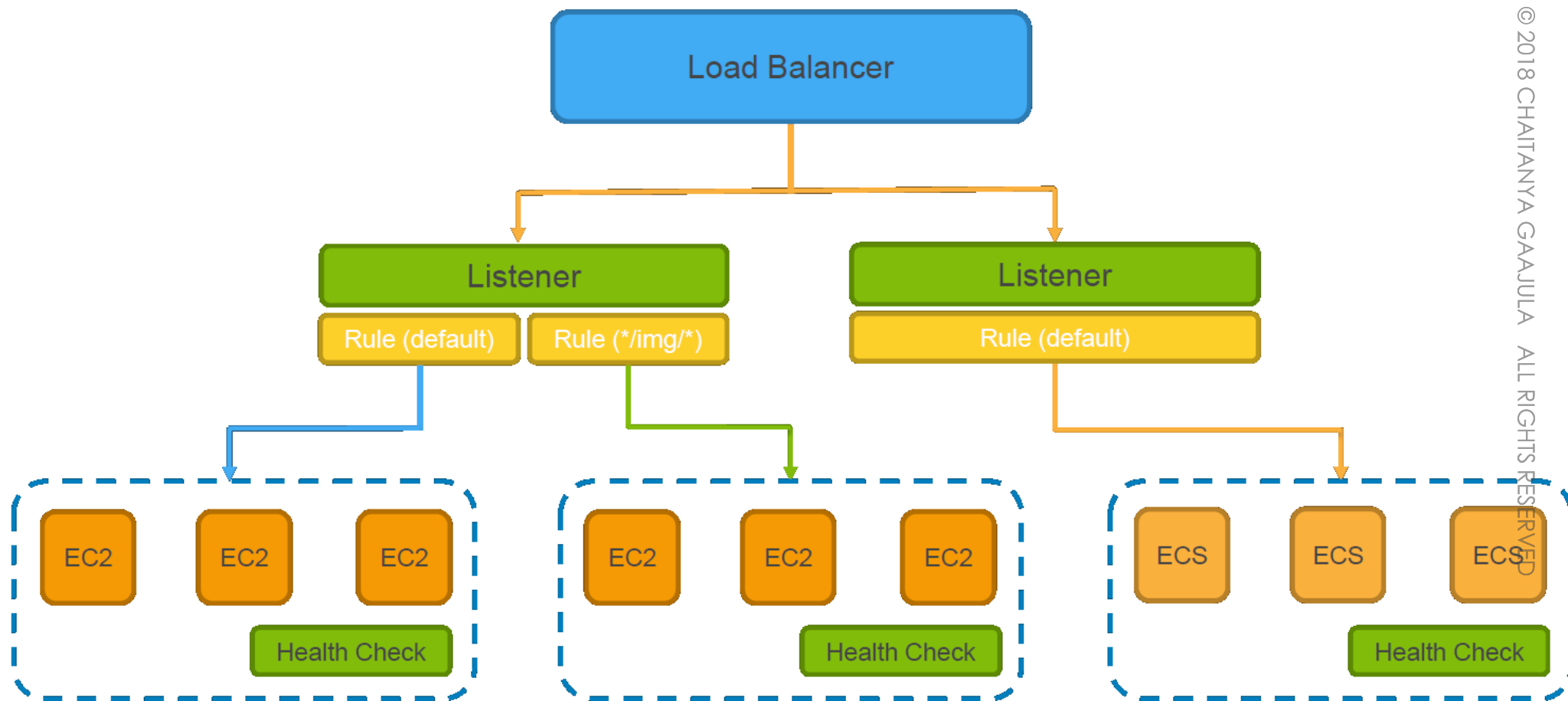
Application Load Balancer Components

Load balancer : Distributes incoming application traffic across multiple targets, such as EC2 instances, in multiple Availability Zones. This increases the availability of your application. One or more listeners can be added to the load balancer.

Listener : Checks for connection requests from clients, using the protocol and port that is configured, and forwards requests to one or more target groups (Content of the request), based on the rules that are defined. Each rule specifies a target group, condition, and priority. When the condition is met, the traffic is forwarded to the target group.

Target group : Routes requests to one or more registered targets, such as EC2 instances, using the protocol and port number that is specified. Multiple target groups can be registered and health checks can be configured on a per target group basis. Health checks are performed on all targets registered to a target group that is specified in a listener rule for your load balancer.

AWSELB: Application Flow



AWSELB:Features

**Path-based
Routing**

**Host-based
Routing**

**Path-based
Routing**

**Multiple applications
on a single EC2**

**Registering targets
by
IP address**

**Pay-Only
What You Use**

AWSELB: Limitation

Resource	Default Limit
Load balancers per region	20
Target groups per region	3000
Load balancers per target group	1
Targets per load balancer	1000
Targets per target group	1000
Listeners per load balancer	50
Rules per load balancer	100
Number of times a target can be registered per load balancer	100
Security groups per load balancer	5
Subnets per Availability Zone per load balancer	1
Certificates per listener	1
Conditions per rule (one host condition, one path condition)	2
Actions per rule	1
Target groups per action	1



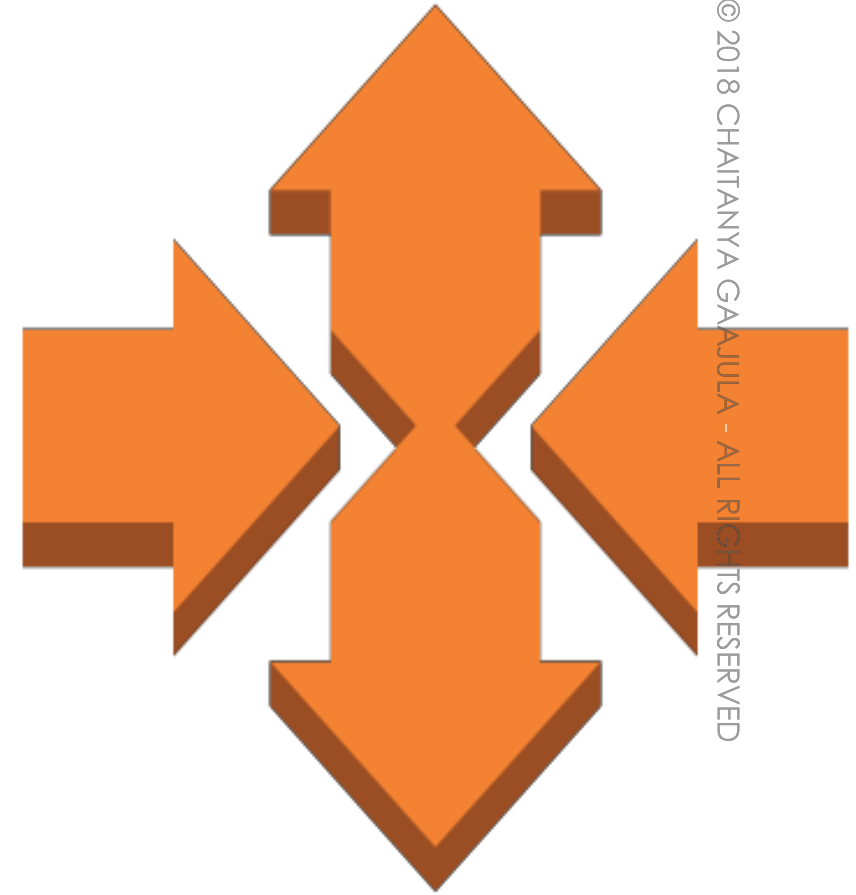
Network Load Balancer

- ❑ An Network Load Balancer functions at the fourth layer of the Open Systems Interconnection (OSI) model.
- ❑ Ability to handle volatile workloads and scale to millions of requests per second.
- ❑ Capable of handling millions of requests per second.
- ❑ Static IP Address can be assigned for Load balancer
- ❑ Listeners support the TCP protocols.

Auto Scaling

Agenda

- ❖ What is AWS Auto Scaling
- ❖ Auto Scaling Components
- ❖ Auto Scaling Group
- ❖ Auto Scaling Launch Configuration
- ❖ Auto Scaling Benefits
- ❖ Auto Scaling Lifecycle
- ❖ Auto Scaling Plans
- ❖ Manual Scaling
- ❖ Schedule Scaling
- ❖ Dynamic Scaling
- ❖ Auto Scaling StepAdjustment
- ❖ Auto Scaling TerminationPolicy
- ❖ Default TerminationPolicy
- ❖ Health Check
- ❖ Quiz
- ❖ Hands-On Lab



What is AWS Auto Scaling

Scalability is the capability of a system, network, or process to handle a growing amount of work, or its potential to be enlarged to accommodate that growth.

- ❑ Types of scaling:
 - **Horizontal Scaling** [scaling out and scaling in]
 - **Vertical Scaling** [scaling up and scaling down]

AWS AutoScaling: Components



Groups

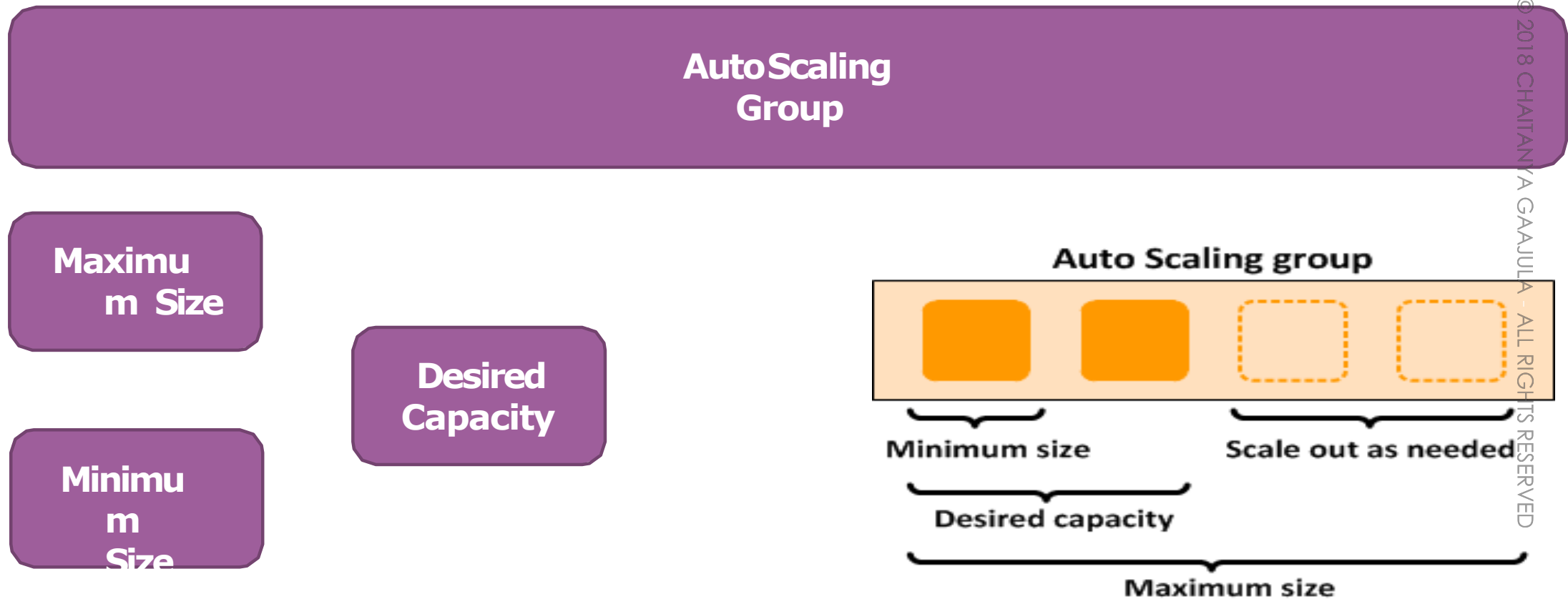


**Launch
Configuration**



**Scaling
Plans**

AWS AutoScaling: Group



AWS AutoScaling: LaunchConfiguration

Creating Launch Configuration from Scratch

Amazon Machine Image

Key Pair

Security Group

Instance Type

Block Device Mapping

Creating Launch Configuration from a running EC2 Instance

AWS AutoScaling: **Benefits**

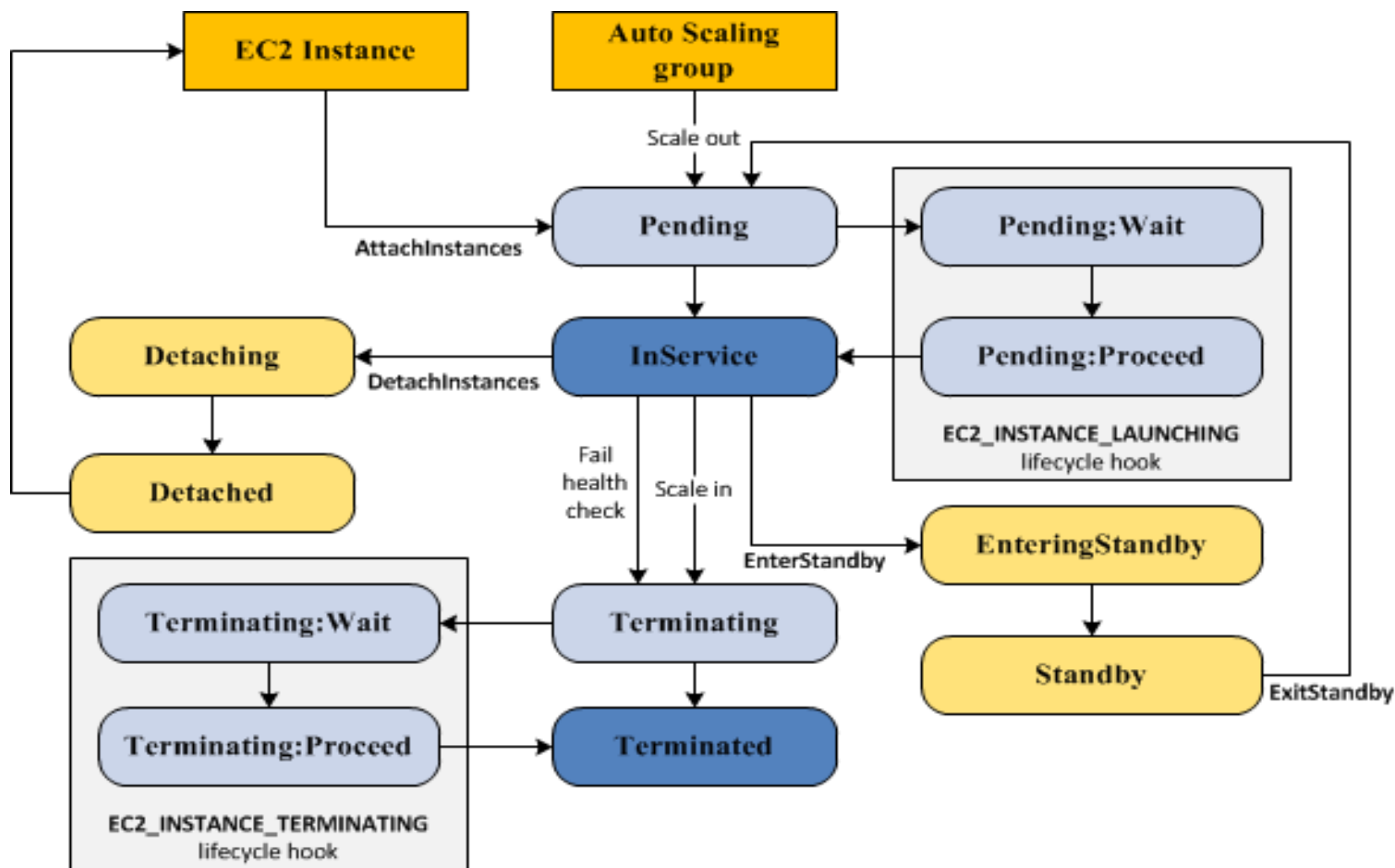
**Fault
Tolerance**

Multiple AZ(s)

Availability

**Cost
Management**

AWS AutoScaling: Lifecycle



AWS AutoScaling: **Plans**

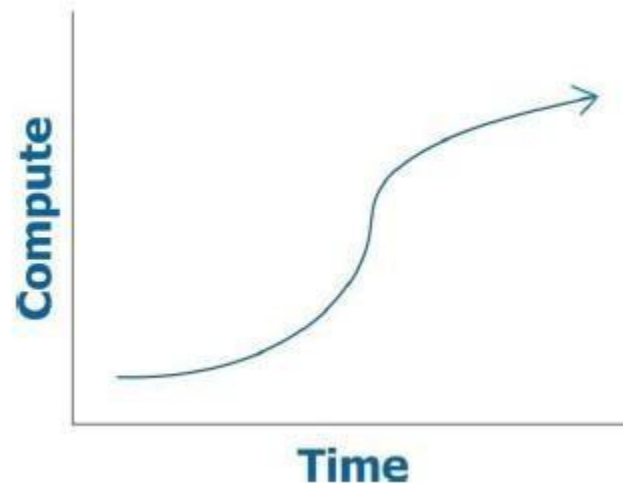
Manual Scaling

Scale based on a Schedule

Scale based on demand

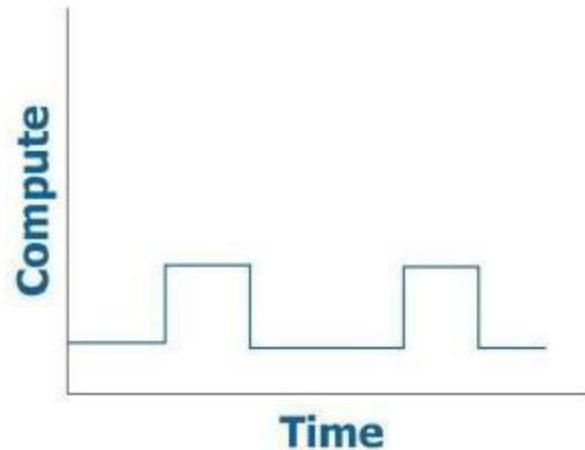
AWS AutoScaling: Manual Scaling

- ❑ At any time, you can change the size of an existing Auto Scaling group.
- ❑ Update the desired capacity of the Auto Scaling group, or update the instances that are attached to the Auto Scaling group.
- ❑ After changes, verify that your Auto Scaling group has launched/ terminated additional instance.



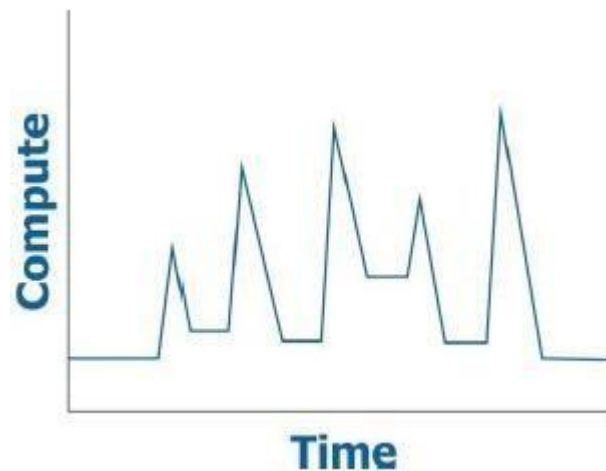
AWS AutoScaling: Schedule Scaling

- ❑ Scaling based on a schedule allows you to scale your application in response to predictable load changes.
- ❑ For example, every week the traffic to your web application starts to increase on Wednesday, remains high on Thursday, and starts to decrease on Friday.
- ❑ You can plan your scaling activities based on the predictable traffic patterns.



AWS AutoScaling: Dynamic Scaling

- ❑ When you use Auto Scaling to scale dynamically, you must define how you want to scale in response to changing demand.
- ❑ For example, say you have a web application that currently runs on two instances and you do not want the CPU utilization of the Auto Scaling group to exceed 70 percent.
- ❑ Scaling Policy Types:
 - Simple scaling
 - Step scaling
 - Target tracking scaling



AWS AutoScaling: Step Adjustment

Scale Out Policy				
Lower bound	Upper bound	Adjustment	Metric value	Changes
0	10	0	$50 \leq \text{value} < 60$	maintains the desired capacity while the aggregated metric value is less than 60
10	20	10	$60 \leq \text{value} < 70$	increases the desired capacity of the group by 1 instance, to 11 instances (add 10 percent of 10 instances)
20	null	30	$70 \leq \text{value} < +\text{infinity}$	increase the desired capacity by another 3 instances, to 14 instances (add 30 percent of 11 instances, 3.3 instances, rounded down to 3 instances).

AWS AutoScaling: Step Adjustment

Scale In Policy				
Lower bound	Upper bound	Adjustment	Metric value	Changes
-10	0	0	$40 < \text{value} \leq 50$	maintains the desired capacity while the aggregated metric value is greater than 40
-20	-10	-10	$30 < \text{value} \leq 40$	if the metric value gets to 40, decreases the desired capacity of the group by 1 instance, to 13 instances, (remove 10 percent of 14 instances, 1.4 instances, rounded down to 1 instance).
null	-20	-30	$-\infty < \text{value} \leq 30$	if the metric value falls to 30, decreases the desired capacity of the group by another 3 instances, to 10 instances. (remove 30

AWSAutoScaling: TerminationPolicy

**Oldest
Instance**

**Newest
Instance**

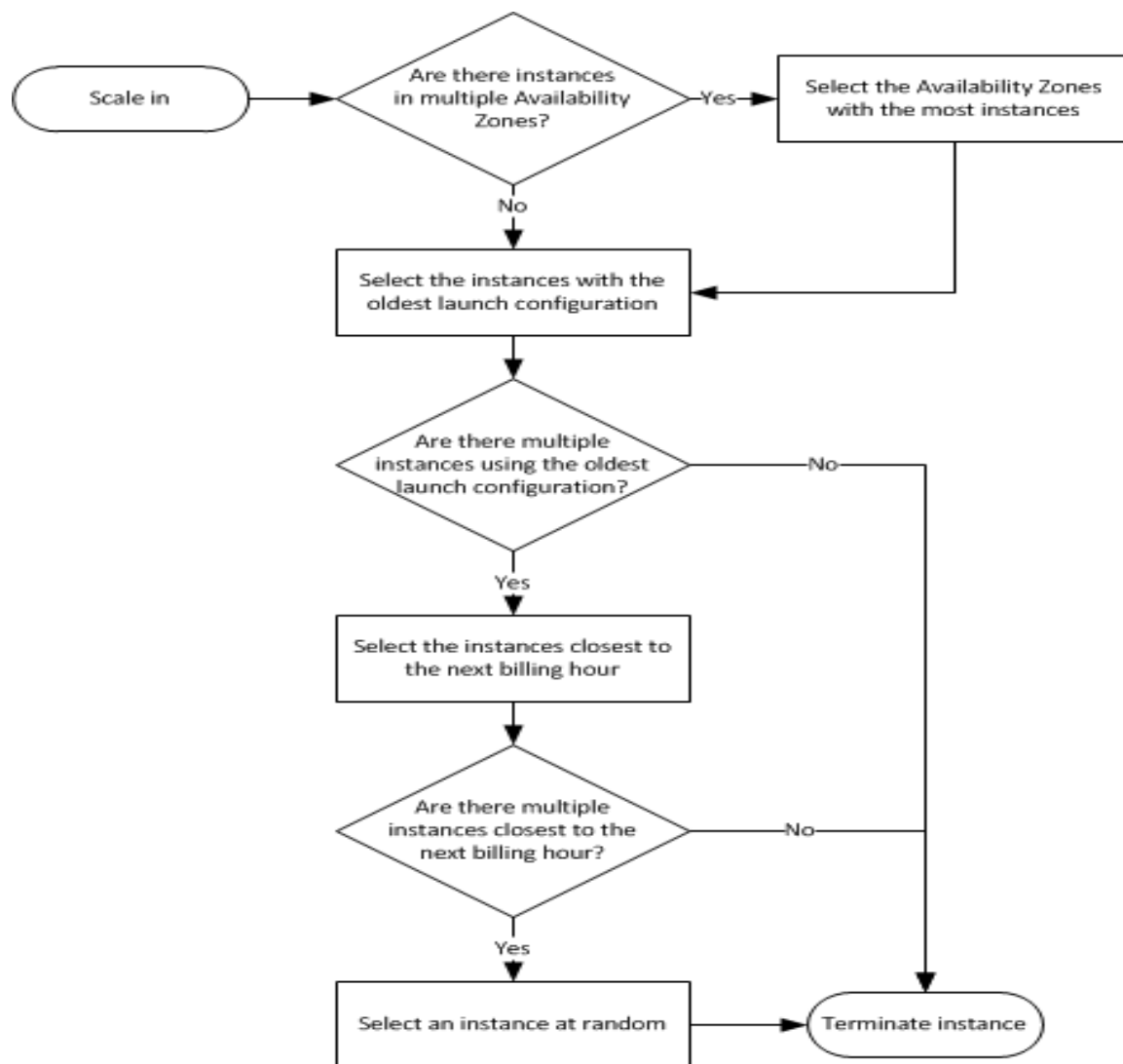
**ClosestToNext
Instance Hour**

**OldestLaunch
Configuration**

Default

**Instance
Protection**

AWS AutoScaling: Default Termination Policy



AWSAutoScaling: HealthCheck

- ❑ Auto Scaling determines the health status of an instance using one or more of the following:
 - Status checks provided by Amazon EC2.
 - Health checks provided by Elastic Load Balancing.
 - Custom health checks.
- ❑ By default, Auto Scaling health checks use the results of the EC2 status checks to determine the health status of an instance.
- ❑ If you attached a load balancer to your Auto Scaling group, you can configure Auto Scaling to mark an instance as unhealthy if Elastic Load Balancing reports the instance as **OutOfService**.

AWS AutoScaling: Limitations

Resource	Default Limit
Launch configurations perregion	100
Auto Scaling groups perregion	20
Scaling policies per Auto Scalinggroup	50
Scheduled actions per Auto Scalinggroup	125
Lifecycle hooks per Auto Scalinggroup	50

© 2018 CHAITANYA GAJULA ALL RIGHTS RESERVED

Hands –On –Lab

- Maintaining High Availability using Auto Scaling