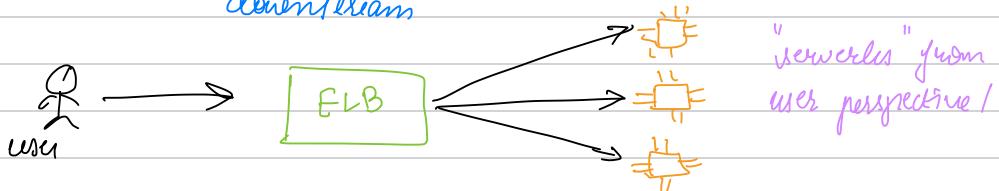


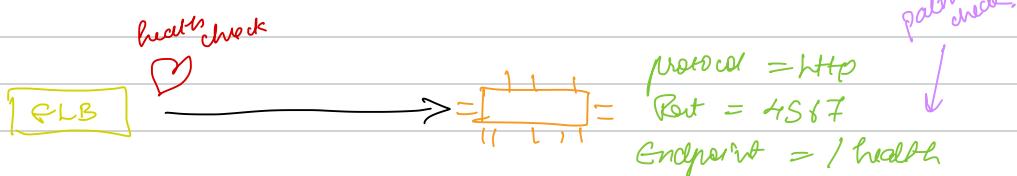
Elastic Load Balances

→ forward traffic to multiple instances downstream



- static DNS for application
- expose a single point of access to application
- do health checks
- provide SSL termination
- high availability
- separate public from private traffic

* ELB is AWS managed! So has ton of integration with other AWS services.



If the response message is not OK (status code ≠ 200) user marked as unhealthy & traffic will be load balanced to another instance.

* 4 types of load balancer can be used for private as well as public

(1) CLB (Classic) { everything { # gone } deprecated }

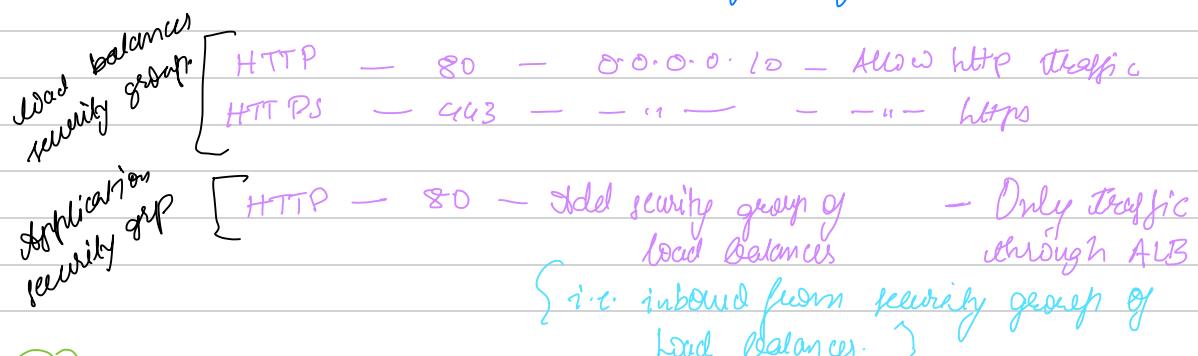
(2) ALB (Application) { HTTP, HTTPS, Web Socket }

(3) NLB (Network) { TCP, TLS ; UDP }

(4) GLB (Gateway) { packet level , IP protocol } (AFN/BYF)
(GWLB)



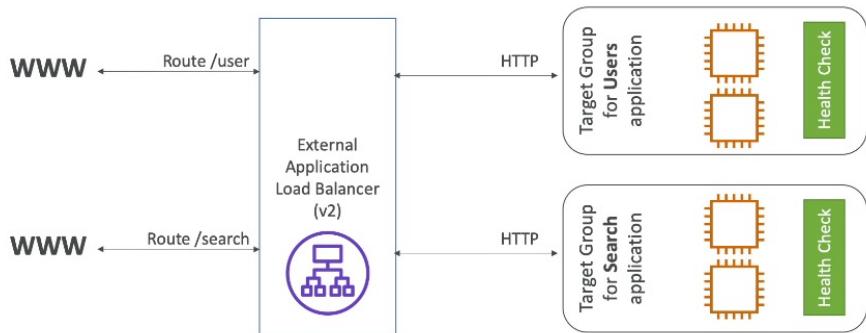
→ int'd we can attack security group of EC2 instance to ELB and restrict EC2 instance to allow request from the ELB & nowhere else. hence tightening network sec.



① Application Load Balancer (V2)

- Layer 7 only
- multiple HTTP application across machines (or on the same machine)
- Supports redirect (from http to https) machine in containers
परन्तु यह http:// का पोर्ट नहीं बदल सकता।
- Supports route routing. { building tables? }
- Based on { URL path { /error }
hostnames
Query string & HTTP headers }
- great for microservices & container based applic'n
- has port mapping feature to redirect to a dynamic port in ECS
- support protocols like HTTP, HTTPS & WebSockets

* Same load balancer for different web applications!



Target groups

EC2 (HTTP)

ECS (HTTP)

A functions (HTTP)

ip address (private)

* ALB can route to

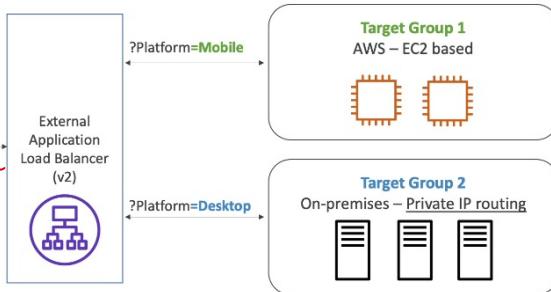
multiple tg.

* health checks @ tg level

translated to JSON events

Application Load Balancer (v2) Query Strings/Parameters Routing

Done in the
listeners of the ALB
add rules
for the listener
to direct incoming
to our audience

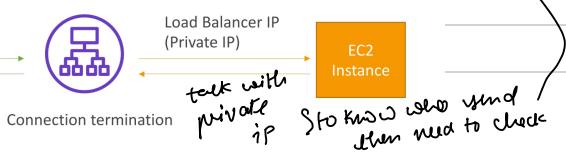


Application Load Balancer (v2) Good to Know

About the load balancer!

- Fixed hostname (XXX.region.elb.amazonaws.com)
- The application servers don't see the IP of the client directly
 - The true IP of the client is inserted in the header X-Forwarded-For
 - We can also get Port (X-Forwarded-Port) and proto (X-Forwarded-Proto)

Client IP
12.34.56.78



Step 2
Define rule conditions

Step 3
Define rule actions

Step 4
Set rule priority

Step 5
Review and create

Listener details: HTTP:80

Actions

Forward to target groups Redirect to URL Return fixed response

Return fixed response Info
Use fixed-response actions to drop client requests and return a custom HTTP response. When a fixed-response action is taken, the action and the URL of the redirect target are recorded in the access logs.

Response code
The type of message you want to send.
404
2xx, 4xx, 5xx

Content type - optional
The format of your message.
text/plain

Response body - optional
Enter your response message.
Not found, custom error!

1024 character maximum

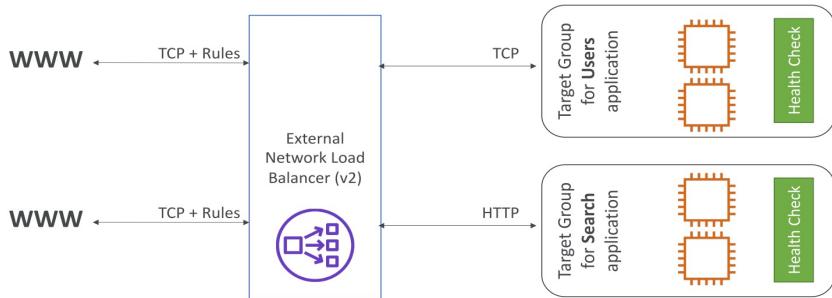
Cancel Previous Next

2) Network Load balancer (v2) *(Not you?)*

↳ Layer 4

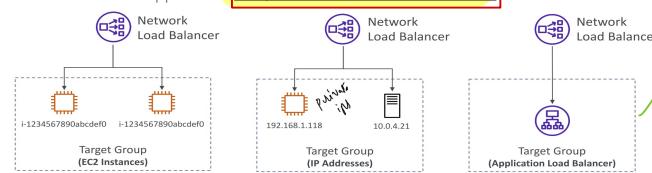
- forward TCP & UDP traffic, extreme performance
- low latency (~100ms) [~400ms for ALB]
- has only 1 static IP per AZ & supports assigning Elastic IP *ALB1 can't be attached with elastic IPs*

Network Load Balancer (v2) TCP (Layer 4) Based Traffic



Network Load Balancer – Target Groups

- EC2 instances
- IP Addresses – must be private IPs
- Application Load Balancer
- Health Checks support the **TCP, HTTP and HTTPS Protocols**

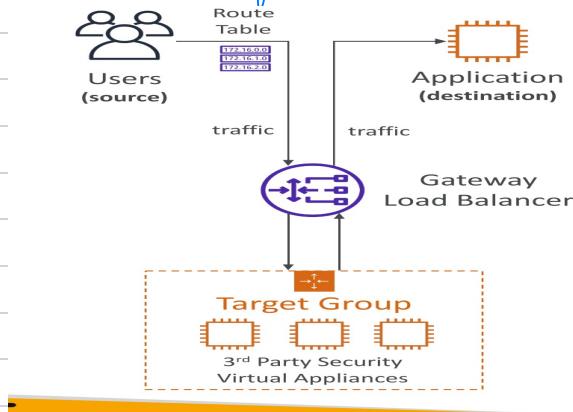


get our own rule for traffic handling with self

③ Gateway Load Balancer

→ Advanced, deploy, scale & manage of fleet of 3rd party network virtual appliances in AWS

Ex: Firewalls, IDPS, Deep pkt inspection systems, payload manipulation.



* It analyzes the network layer traffic and then re-routes the traffic

* Layer 3 (Network layer) - IP pkts

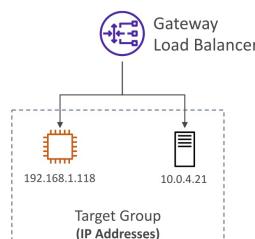
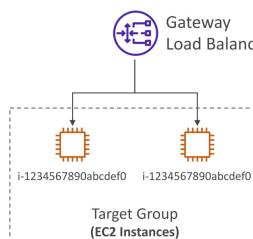
* Uses GENEVE protocol on port 6081

2 func
Transparent Network gateway
(single entry/exit for all traffic)

Load balancer
(distributes target to virtual appliances)

Gateway Load Balancer – Target Groups

- EC2 instances
- IP Addresses – must be private IPs

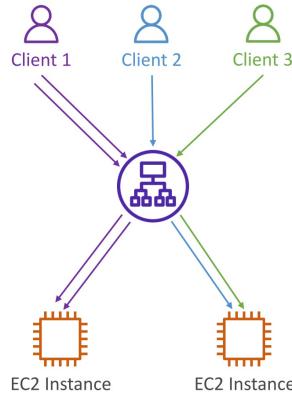


Stick Session (Session affinity)

- It is possible to implement stickiness so that the same client is always redirected to the same instance behind a load balancer
- This works for Classic Load Balancer, Application Load Balancer, and Network Load Balancer
- For both CLB & ALB, the "cookie" used for stickiness has an expiration date you control
- Use case: make sure the user doesn't lose his session data

* Enabling stickiness may bring imbalance to the load over the backend EC2 instances

* NLB doesn't use cookie!



Cookie → Application based
↓
duration based

- Application-based Cookies
 - Custom cookie
 - Generated by the target
 - Can include any custom attributes required by the application
 - Cookie name must be specified individually for each target group
 - Don't use AWSALB, AWSALBAPP, or AWSALBTG (reserved for use by the ELB)
 - Application cookie
 - Generated by the load balancer
 - Cookie name is AWSALBAPP
- Duration-based Cookies
 - Cookie generated by the load balancer
 - Cookie name is AWSALB for ALB, AWSELB for CLB

* Can introduce stickiness during configuration.

Hello World from ip-172-31-7-176.eu-central-1.compute.internal

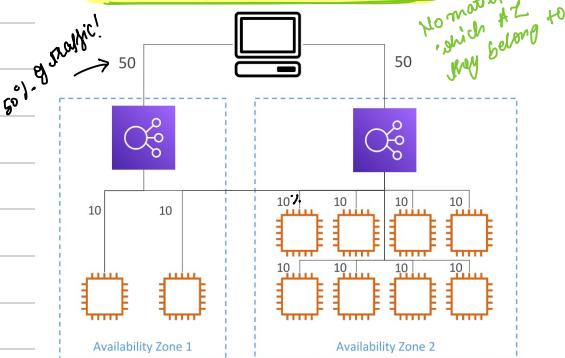
duration based cookie!

Name	Value	Expires
AWSELB-COOKIE	value="XUUTLjR3Q0hpmcYtIhdH4QfTfH0L2gjCn88TTfDAnhJd8m8AA0vYzJ8URPkg/zaedCT4eK7P8L027357lpmuXe6TrBuHe6C3abfHQ2NThyvsl"	2021-07-08T08:22:47.000Z

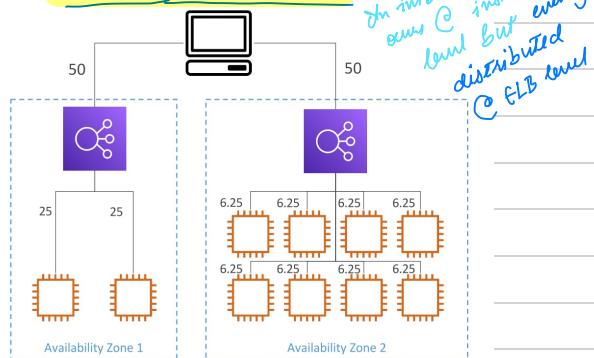
Cross-Zone Load Balancing

Cross-Zone Load Balancing

With Cross Zone Load Balancing:
each load balancer instance distributes evenly across all registered instances in all AZ



Without Cross Zone Load Balancing:
Requests are distributed in the instances of the node of the Elastic Load Balancer



Cross-Zone Load Balancing

- Application Load Balancer
 - Enabled by default (can be disabled at the Target Group level)
 - No charges for inter AZ data
- Network Load Balancer & Gateway Load Balancer
 - Disabled by default
 - You pay charges (\$) for inter AZ data if enabled
- Classic Load Balancer
 - Disabled by default
 - No charges for inter AZ data if enabled

↑ Generally AWS are charged to be paid for inter AZ data transfers!

* Cross Zone LB is only default enabled for ALBs

LB > Attributes > C Z LB Enabled? > edit > enable

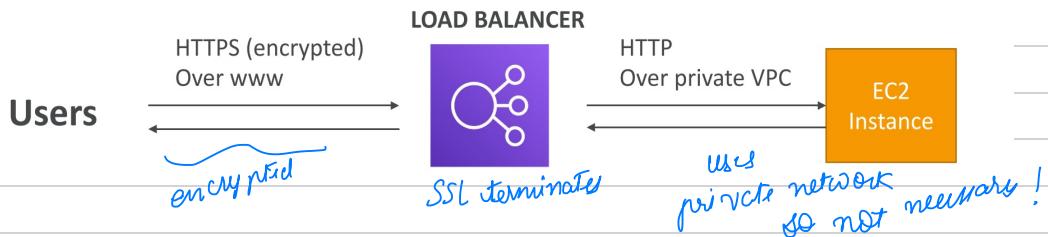
SSL / TLS

by adding listeners to the ELB

SSL allows in flight encryption b/w client & LB
Secure socket layer

TLS → newer version } Transport Layer Security }

public SSL certificates issued by Certificate Authorities (CA)
they need to be renewed from time to time because
they have an expiration date



- The load balancer uses an X.509 certificate (SSL/TLS server certificate)
- You can manage certificates using ACM (AWS Certificate Manager)
- You can create/upload your own certificates alternatively
- HTTPS listener:
 - You must specify a default certificate
 - You can add an optional list of certs to support multiple domains
 - Clients can use SNI (Server Name Indication) to specify the hostname they reach
 - Ability to specify a security policy to support older versions of SSL / TLS (legacy clients)

SNI { Server name indication }

→ allows to expose multiple HTTP(s) application each with its own SSL certificate on the same port .

- SNI solves the problem of loading **multiple SSL certificates onto one web server** (to serve multiple websites)
- It's a "newer" protocol, and requires the client to **indicate the hostname of the target server** in the initial SSL handshake
- The server will then find the correct certificate, or return the default one

Note:

- Only works for ALB & NLB (newer generation), CloudFront
- Does not work for CLB (older gen)

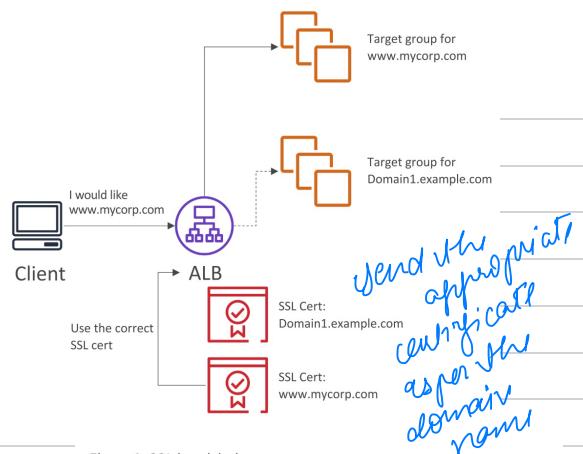
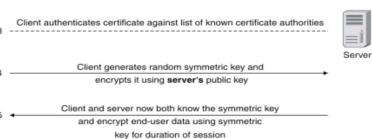
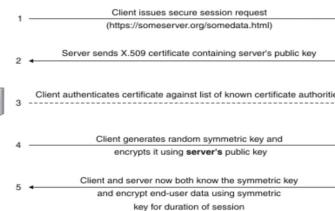


Figure 1. SSL handshake



Elastic Load Balancers – SSL Certificates

- Classic Load Balancer (v1)
 - **Support only one SSL certificate**
 - Must use multiple CLB for multiple hostname with multiple SSL certificates
- Application Load Balancer (v2)
 - Supports **multiple listeners with multiple SSL certificates**
 - **Uses Server Name Indication (SNI)** to make it work
- Network Load Balancer (v2)
 - **Supports multiple listeners with multiple SSL certificates**
 - **Uses Server Name Indication (SNI)** to make it work

ELB > Listener > Add listener > { HTTPS { 443 } > Forward to tg > ELB security policy > ACM { default } or import it }

Connection Draining

Retaining the connection while the instance is de-registering itself or going unhealthy.

existing req.

के साथ रहें

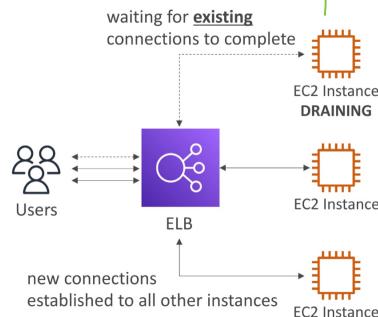
तो यह उपी !

जब तक यह उपी !

Connection Draining

- Feature naming
 - Connection Draining – for CLB
 - Deregistration Delay – for ALB & NLB
- Time to complete "in-flight requests" while the instance is de-registering or unhealthy
- Stops sending new requests to the EC2 instance which is de-registering
- Between 1 to 3600 seconds (default: 300 seconds)
- Can be disabled (set value to 0)
- Set to a low value if your requests are short

* * in draining state!



Auto Scaling Groups

get rid of servers very quickly!

scale out { increase the no. of instances } or
scale in { decrease the no. of instances }

* ASG is free { but have to pay for the instance created }

minimum capacity
desired - - -
maximum - - -



→ When ALB frontend's a ASG & suppose that one of the EC2 instance fails the health check, then the ASG will automatically terminate that instance & relaunch a new EC2 instance.

Auto Scaling Group Attributes

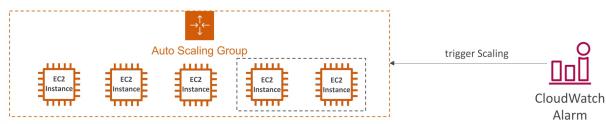
- A Launch Template (older "Launch Configurations" are deprecated)
 - AMI + Instance Type
 - EC2 User Data
 - EBS Volumes
 - Security Groups
 - SSH Key Pair
 - IAM Roles for your EC2 Instances
 - Network + Subnets Information
 - Load Balancer Information
- Min Size / Max Size / Initial Capacity
- Scaling Policies



information about the instances
to be launched!

Auto Scaling - CloudWatch Alarms & Scaling

- It is possible to scale an ASG based on CloudWatch alarms
- An alarm monitors a metric (such as Average CPU, or a custom metric)
- Metrics such as Average CPU are computed for the overall ASG instances
- Based on the alarm:
 - We can create scale-out policies (increase the number of instances)
 - We can create scale-in policies (decrease the number of instances)



Scaling policies for ASG |

i) Dynamic scaling

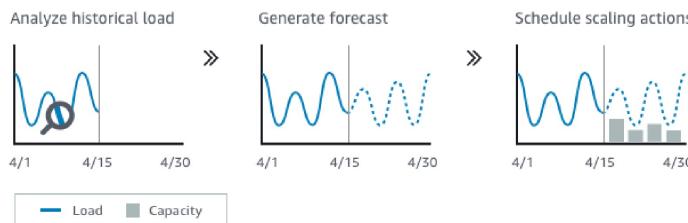
target tracking scaling (Simpler)
Scaling CPU utilization to 70% (example?)

simple / step scaling
if CPU > 70%. then add 2 units
if CPU < 30%. then remove 1 unit

ii) Scheduled scaling

Schedule a scaling based on known usage patterns
check time for this it's late to AM to 6 PM (1 min capacity)

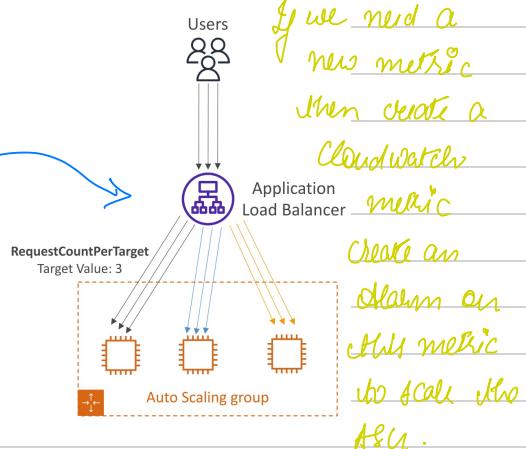
3) Predictive Scaling { continuously forecast load & schedule scaling ahead}



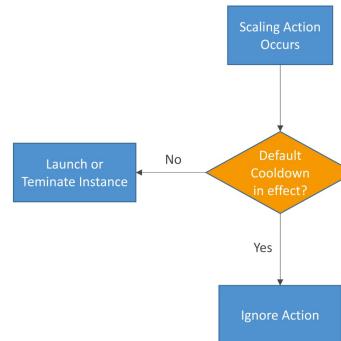
Good metrics to scale on

Always present
just choose
useful metrics.

- CPUUtilization: Average CPU utilization across your instances
- RequestCountPerTarget: to make sure the number of requests per EC2 instances is stable
- Average Network In / Out (if you're application is network bound)
- Any custom metric (that you push using CloudWatch)



- After a scaling activity happens, you are in the **cooldown period** (default 300 seconds)
- During the cooldown period, the ASG **will not launch or terminate** additional instances (to allow for metrics to stabilize)
- Advice: Use a ready-to-use AMI to reduce configuration time in order to be serving request faster and reduce the cooldown period



ASG > attributes { automatic scaling } > add scaling policies