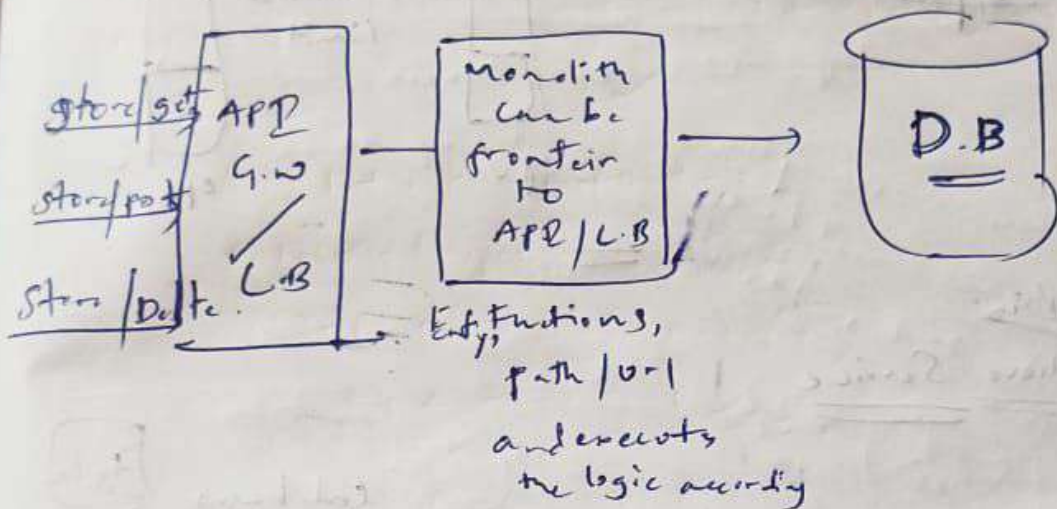


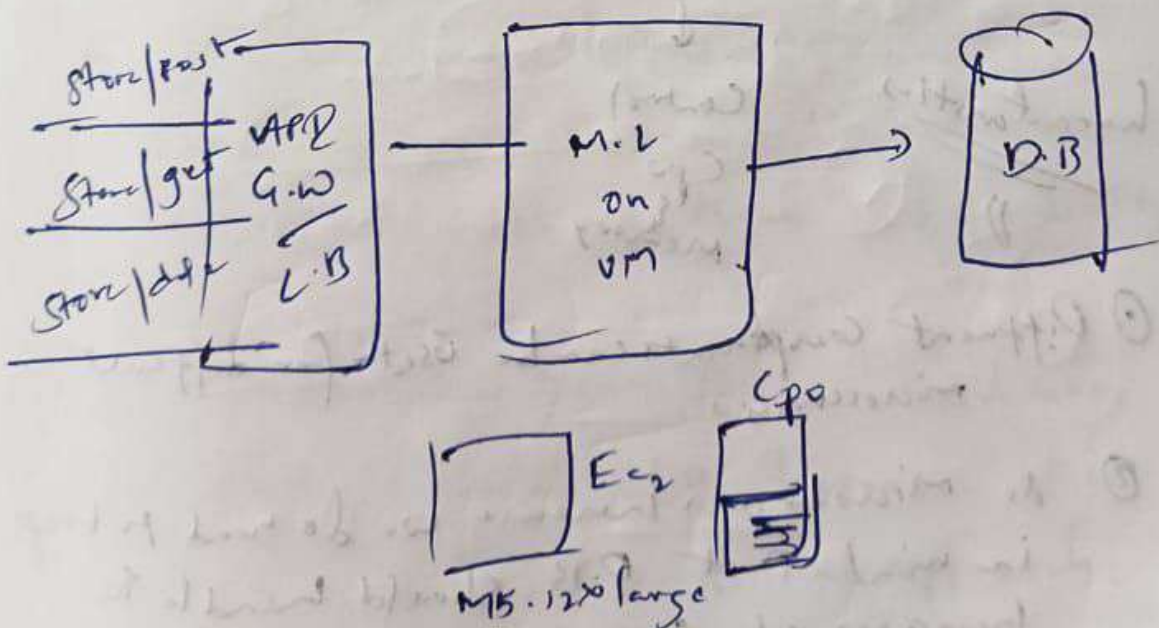
Can you use API with Monoliths

Absolutely yes

API's have nothing to deal with
monoliths or Microservices



Scaling of Monolith : Issues



Generally there are lots of techniques to optimize the reading from the DB for replication catching etc.

- ② Another advantage is as microservices have their own backends and are independent with each other these can be controlled by different teams

python

nodejs

Go

which is

so

polyglot

Advantages / Characteristics of Microservices

1) The core characteristic of microservices they can be, microservices is independent to each other and they can scale independently to each other.

- ② You can have different Governance, Security guidelines and they can be deployed each of them independently.

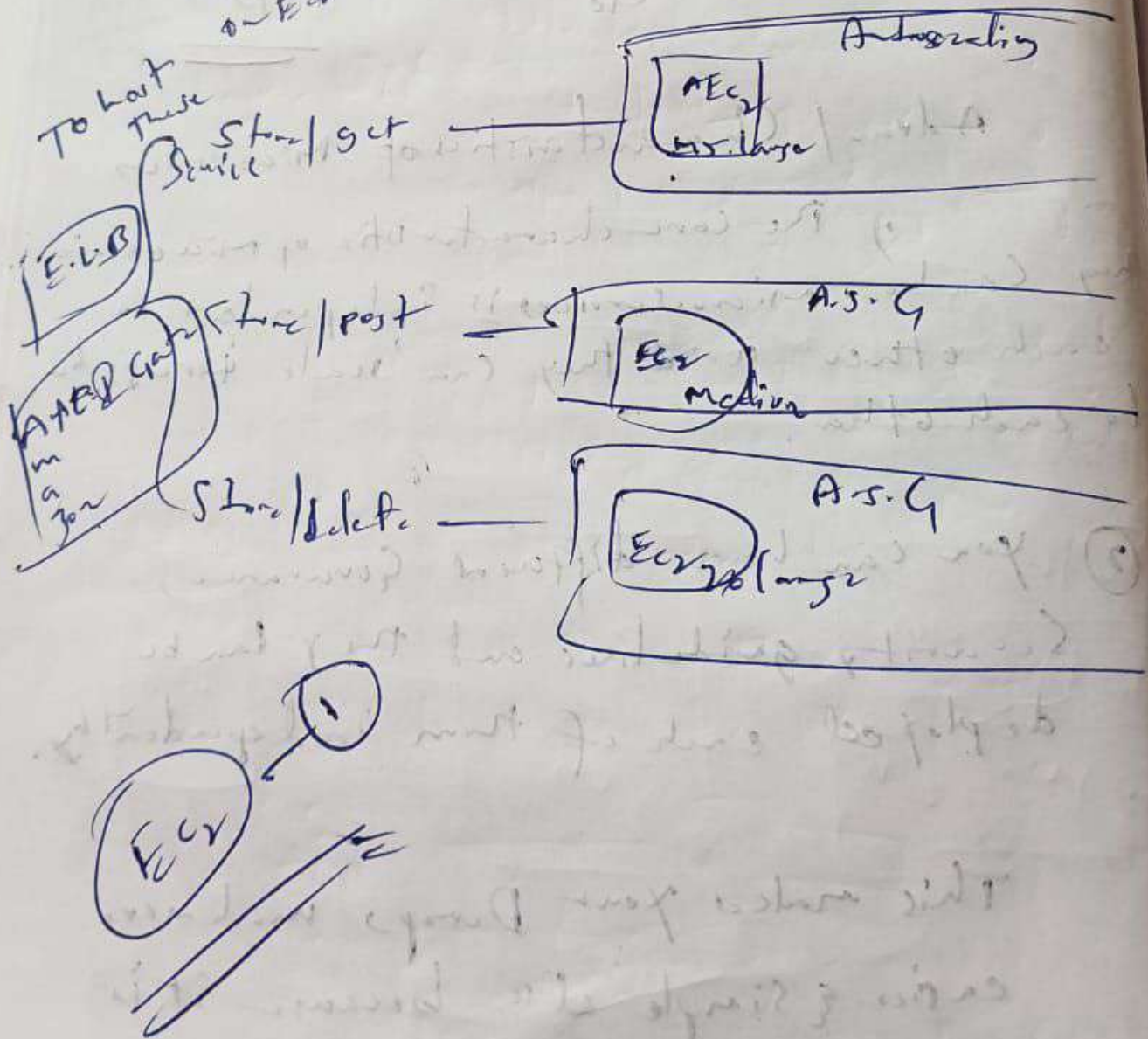
This makes your Deploys much more easier & simple also because it is no longer all or nothing approach.

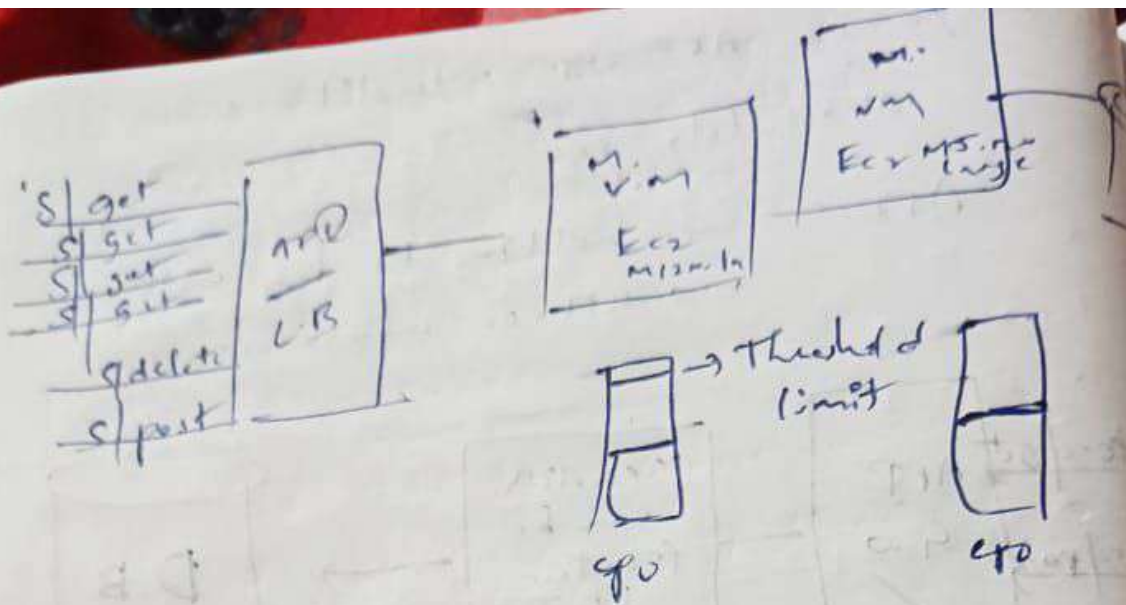
you can test this mis independent to each other and they should have different functionalities.

Deploying Microservices in AWS

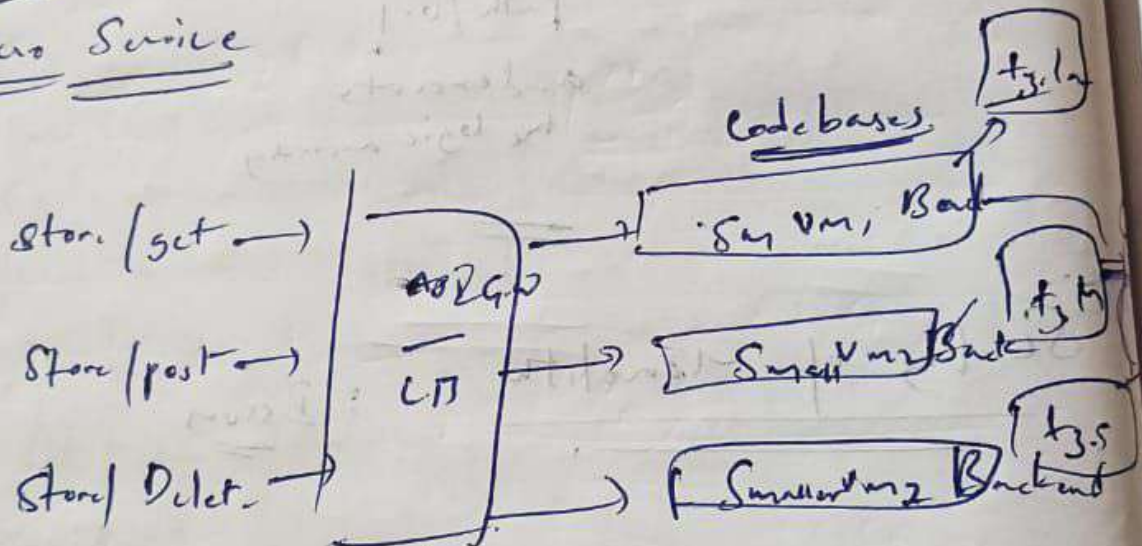
We have one service answer to every thing ECR

there is misconception
mis cannot be run
on ECR. It is wrong





APP's in Micro Service



Characteristics

- ① Different components can be used for different microservices.
- ② As Microservices increase we do need to keep in mind that D.B. should handle the increases of connections.

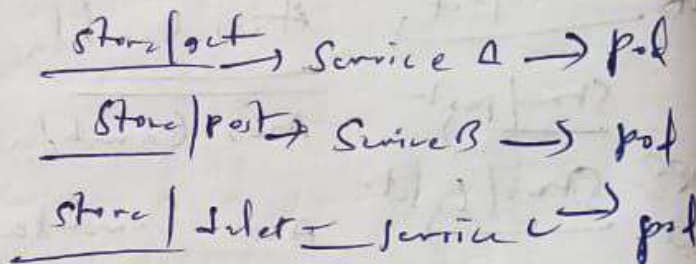
kubernetes part 6

pod's smallest deployable unit

- 1) Elastic LB
- 2) Amazon API gateway
- 3) Application Load Balancer
- 4) ingress ALB.

=

Ingress ALB

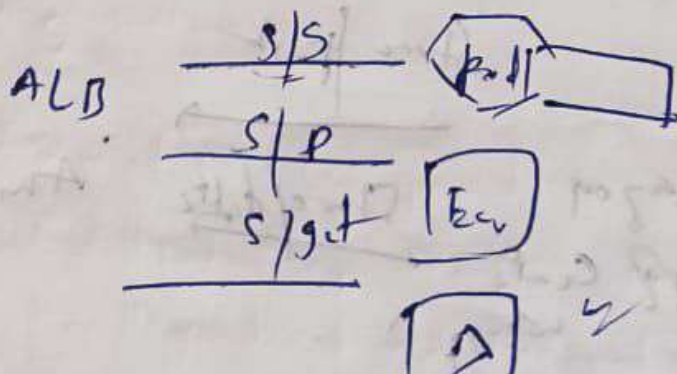


A. Flat
Kuber
ser

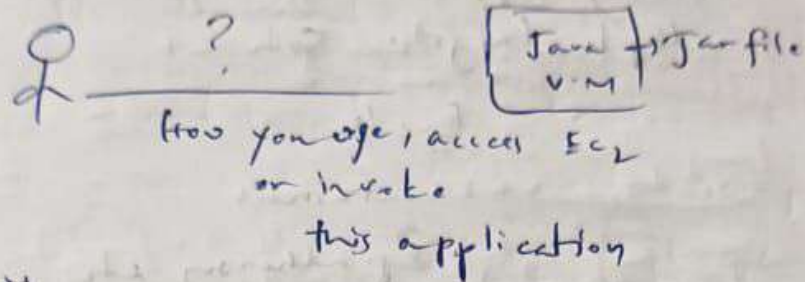
Since it independent with each other
& written different
programming language

Conte
stor
code
=

we can mix & match all these services



Load Balancer



ways

- 1) Each Amazon EC2 has IP address so you can directly invoke this IP address to invoke the application

② But there are some caveats to it?

① What if this EC2 goes down and another EC2 comes up with different IP address

② How will you discover the new IP address

③ What if Application Scales and now instead of one EC2 you have 3 Amazon EC2 what if even more EC2's keep getting added

④ How will you discover these EC2's as well as invoke the application and ensure that traffic is balanced b/w these EC2's

That's where Load Balancer comes to the picture

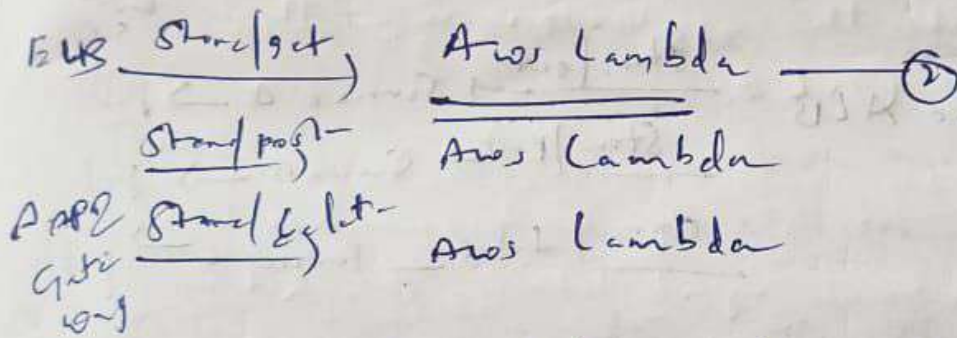
In modern ^{API} services we have In modern application development we can use different AWS Lambda for each microservices backend

Remember Lambda scales

automatically, so

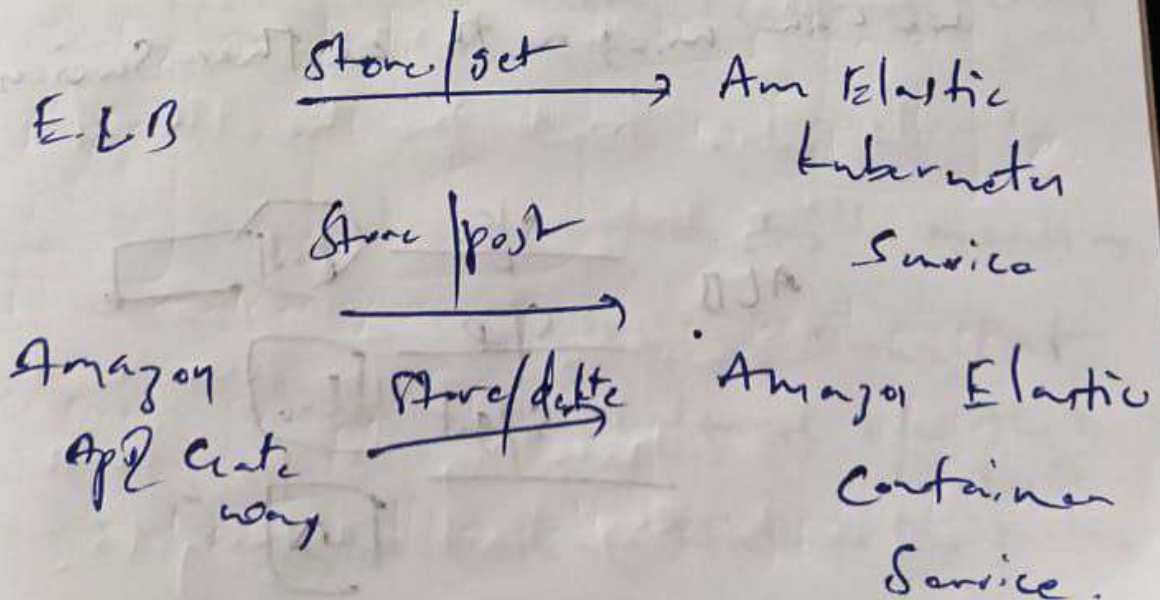
or automatically,

And Lambda can be fronted by ELB
or Amazon API Gateway



Lambda scales automatically

Popular choice
Another choice is Containers — ②



LB with Kubernetes pods

Similarly Kubernetes pods Load Balancer can distribute traffic across multiple pods and in this case

LB acts as ingress / service.

LB with Lambda

Similarly even if your application is running in lambda load balancer can distribute traffic to AWS lambda (those as well)

What application is exposed to website that you already have such as
low no. storage.com:

How you solve that?

E.L.B gives you DNS. Will you ^{can} use

ELB with a DNS services such as Amazon Route53.

Whenever user invokes www.stores.com
The traffic first goes to Amazon Route 53
and from there ↓

User types URL

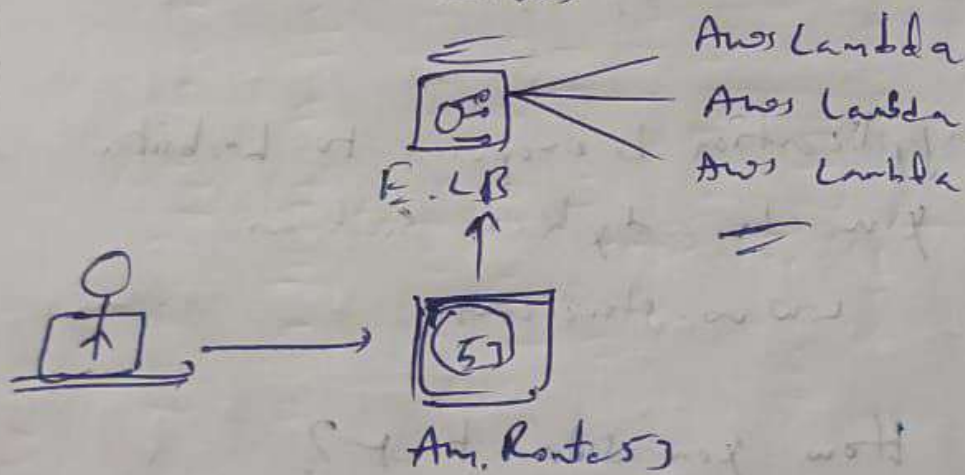
↓
DNS Resolution starts

↓
Resolution queries authoritative DNS

↓
DNS returns ELB IP

↓
Browser send HTTP Request to
ELB

ELB routes traffic to backend
Servers



Load Balancer automatically distributes incoming application traffic across multiple targets. Such as

EC2's

Containers

IP addresses etc

Each Elastic L.O has Unique DNS name so that your user can invoke this application using these DNS URLs.

If the traffic goes up and your backend needs to scale up and let's say it adds another EC2.

L.B automatically discover this new ~~EC2~~ instance and start distributing traffic to it.

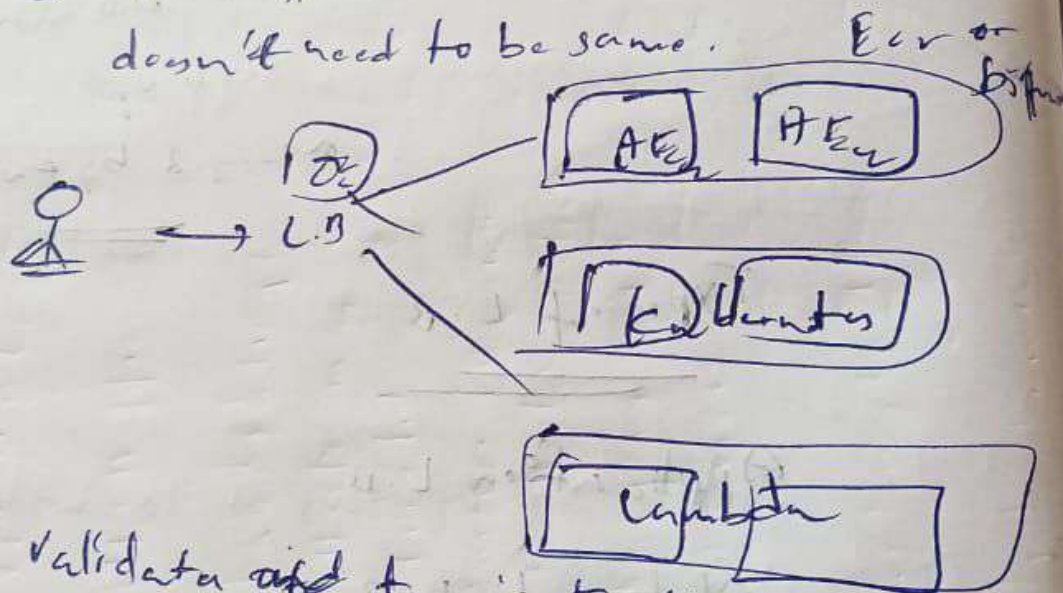
It also tracks the health of all ~~EC2~~ instances.

Even if one EC2 goes down, and it will start distributing traffic to the rest of the healthy back end targets.

we can distribute the traffic to backends with paths to the traffic groups.

Properties of LB

- ① Operates on OSI Layer 7
- ② Routes traffic based on url path
- ③ For different url, the backend doesn't need to be same.



- ④ Validates and terminates SSL



- Operates on OS layer 4
- Route traffic based on protocol and port of incoming traffic

NLB cannot send traffic to backend based on different URL because URL path is Application layer property

- SSL passthrough by default

ALB or NLB ? What to use

① NLB handles spiky traffic better

② NLB exposes static IP address

③ ALB need Global acceleration
Influenced by choices we have

• APP Gateway REST API

integration with backend

like EC2, ELB, etc

Integrate with private link

Integrates with NLB Support EC2 instance

and IP address in Backend

connect target group. It does not

support NLB support as B.K

Support NLB

Compass Through &

Terminate as well

Properties of Load Balancer

- Automatically distributes incoming traffic across multiple targets
- Monitor health of targets
- Integrates with SSL
- Elastic - traffic goes up LB also will also scale up.

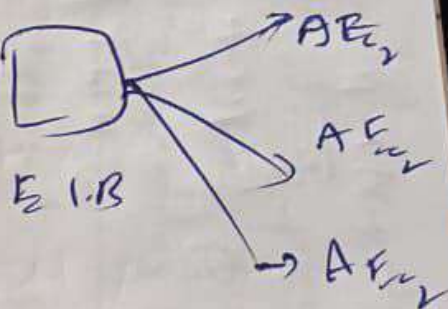
Managed by AWS

Types of L.B

Application L.B

Network L.B

ALB



Let's say our instance website `www.shom.com`.
I want different instances and backend for different operation.

Let `www.shom.com/get/post/delete`
like this.

Target type

Application level

HTTP status code

Health checks

Application level

(HTTP status code)

ODF - PP

Backend IP address,
private IP

N/Land Top Connections

N/Land Top Connections

Simple Telugu:

ALB = Web master

Browser to URI open chesting

HTTP/HTTPS request ni need
check path or header batri
correct server ki postistadi

NLB: Highway traffic cop NLB ki content
thorani ledu, yavaraini
connect they fast ga
top level lo forward
chestadi

HTTP/HTTPS/path
based routing karate

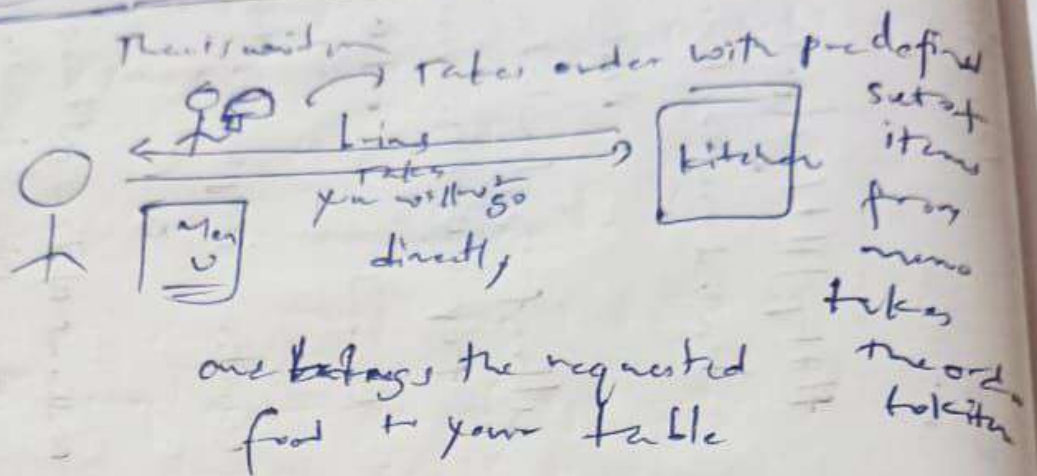
→ ALB

Ultrafast, low latency
top & UDP service
kanali NLB

Content inspection karate
ALB

Millions of low latency
connections, handle
chuppali (gaming, IoT)
real time

APR & why do we need them?



Here waiter acts as an interface to the you / kitchen. So waiter is APR in this case

Backend → Bank / online booking /
finder etc → with APR
waiter → APR

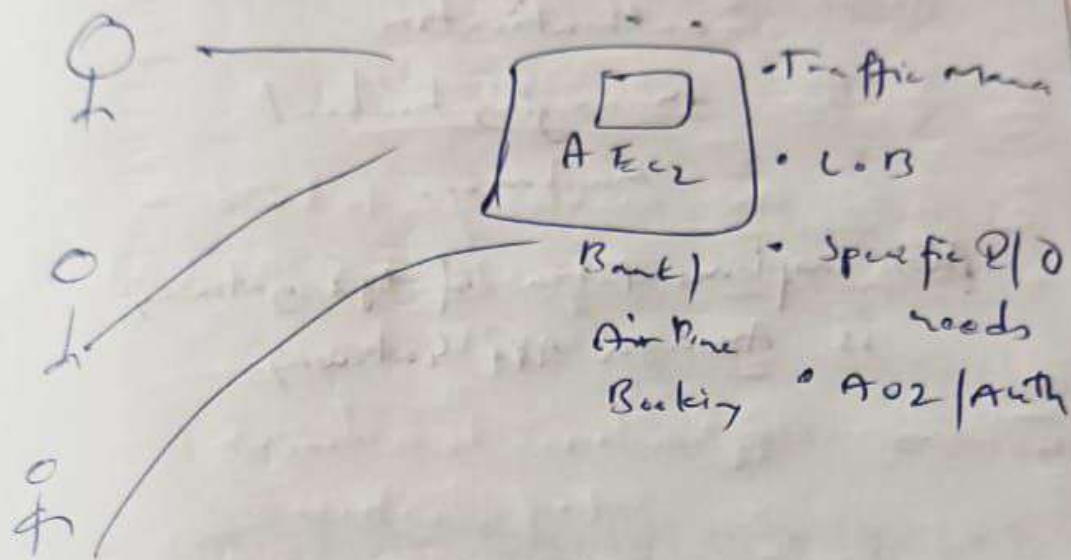
Menu → with predefined input
fields

What is the use of APR? How do you
skip APR directly Backend in realtime?

In Real world there are
multiple users using apps
Sometimes millions

	<u>ALB</u>	<u>NLB</u>
Layer	Layer 7 Application Layer	Layer 4 Transport Layer
Traffic Type	HTTP/HTTPS traffic	TCP, UDP, TLS level traffic
Routing Logic	Content based Routing (path, Hostname, Headers etc)	Fastest connection based routing, no deep request inspection
Latency	Bit more (Because it checks content)	Ultra low latency (millisecond level)
Use cases	for web apps, API microservices	for high performance apps, real time, etc)
SSL Termination	Ecn, Ecs, Lambda	odd apps
Target type	Application level http status code	Ecn + Ip address, private links N/Level TCP connection

Beyond just thinking the program with some input parameter and getting output, there are a lot of other factors such as



If you want code all things in application program, there are lot of overload because

- ① all this interface and cases have nothing to do with business logic that's running your application
- ② So coding all this into your application makes it difficult to test, develop, modify and deploy your code.
- ③ And inversely that a lot of open interfacing standards that you need to implement and implementing all these with their application code makes it very cumbersome and painful.

App Gateway

1) Can implement rate limiting, bursting for APIs

2) Not possible to get static IP address for endpoint

3) Accepts HTTPS traffic

4) It is able to Request validation Request / response mapping

5) Able to handle Spiky traffic (default Rct. looks at burst rate)

Both Traffic Superspeed

Both integrate with

AWS Web application Firewall (WAF)

ALB's

Rate limiting, bursting Capability

possible to get static IP address for LB endpoint (by using Global accelerator)

Both

Accepts HTTP & HTTPS traffic

Not able to do Request validation Request / response mapping

Able to handle but sometimes But we can protect LB's to avoid delay

Both can handle Traffic Spiky

Both can handle Traffic Spiky

- ① Able to integrate with Lambda ^{from} different regions with different accounts but both were in diff account we would be to connect
- ② We can able to import/export API from openapi, Swagger etc
- ③ Have extensive API/A2 Integration
- API key, OAuth, Cognito user pools etc
- ④ Able to Cache Responses
- ⑤ Timeout limit 30 seconds
- ⑥ Integrates almost all API services
- ⑦ ALB is a regional service
- ⑧ No direct method to support/export for cross platform
- ⑨ Integrates with any OIDC compliant IdP Cognito etc.
- ⑩ Able to cache responses
- ⑪ Timeout limit 4000 seconds
- ⑫ Use ECS Lambda, EC2 address as backend.

ALB (App LB)

Automatically distributes incoming traffic across backend targets
Layer 7 load balancer

Infrastructure managed by AWS
highly available elastic

when traffic goes up, the LB needs to scale up as well. So to account for the increase of rate of traffic

ALB (App LB)

Fully managed and serverless ALB

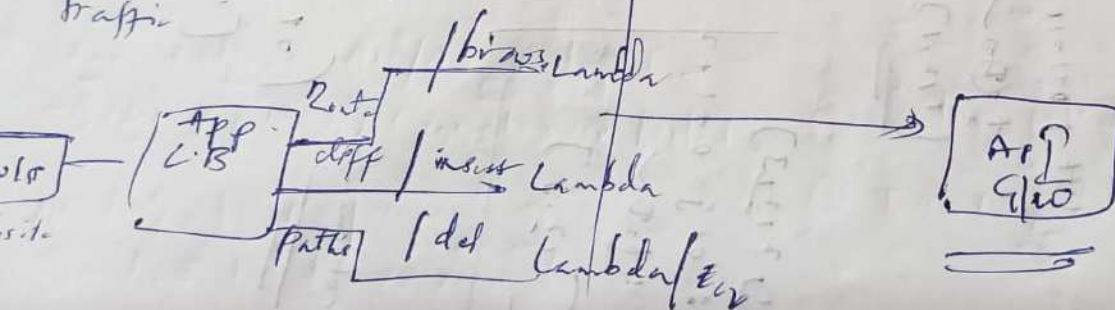
Service from AWS

Automatically scales up and down

Same Option

Application Load Balancer

ALB Gateway



ALB Gateway

Both integrate with

ALB