

S. No.	Date	Title	Page No.	Teacher's Sign / Remarks
		AWS		
1)		Cloud Introduction		
2)		EC2		
3		AMI		
4		Load balancer		
5		cloudwatch		
6		SNS		
7		Auto scaling Groups		
8		EBS		
9		S3		
10		Cloudfront		
11		IAM		
12		SES		
13		RDS		
		aws cli		
14		cloudformation (awscli)		
15		Route53		
16.		(λ) Lambda function		
17		Elastic Beanstalk		

* what is cloud?

→ Without buying the hardware, if a customer can use H/W resources of vendor, on demand of customer, can access globally through internet is called as cloud.

e.g.: - google drive

→ The cloud is all about "pay for what you use"

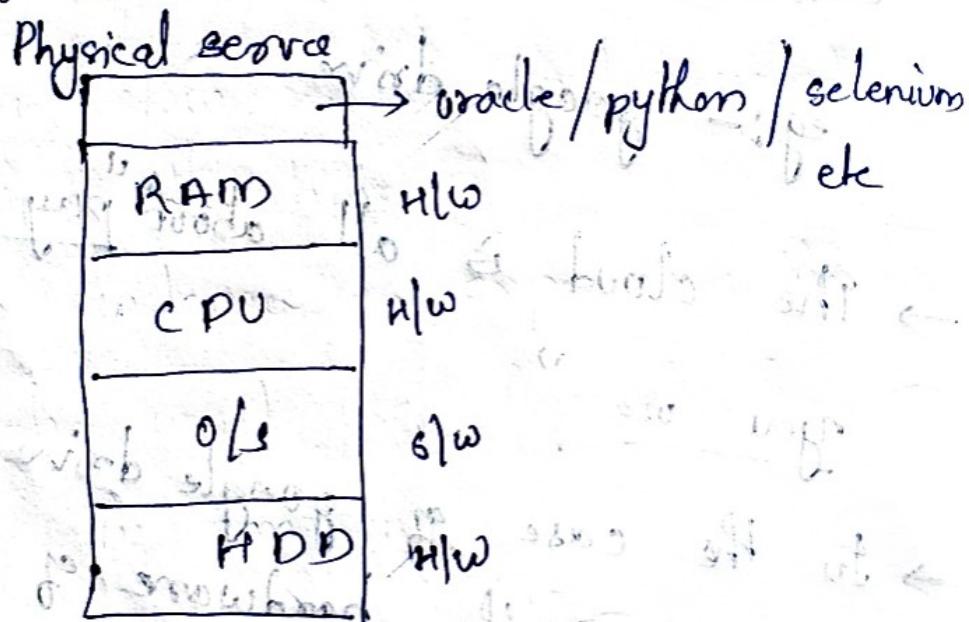
→ In the case of google drive we are not buying the hardware of google drive by just utilizing the H/W of we are just creating account, upto 15GB is free, beyond that you have to pay;

→ Google drive is a storage service of vendor (google)

Server:-

→ It is serving something (an application) to you. It is called server. It will take a request from you and process it and give response to you.

H/w peripheral



- If it is oracle, it is called database server.
- If it is python, it is called development server.
- If it is selenium, it is called testing server.

→ The purpose of the server you use
on it will define the server of
that purpose.

Physical Server in real time

- It needs atleast 64 GB of RAM,
it is not free, let's say 50K rupees
- CPU → atleast 2-3 Lakh
- O/S → software licensing → 2-3 Lakh
- HDD → 50K
- Oracle → 3 Lakh for enterprise addition
of database
- We are spending Lakhs of rupees to
setup one server for our business.
- This total process is called H/W
procurement (buying/getting)
- It will be taking more than 2-3 weeks of time to setup a realtime
server.

→ To install what we need different beans in our organisation.

Let's say,

We setup a server for flipcart which can process 100 user per second.

→ On certain day of flipcart sale, 175 user are hitting the server, the service becomes slow or it won't respond at all. It is called downtime.

→ If the customer is experiencing the downtime, he won't use flipcart, this is loss to the company.

→ He has to maintain the no. of servers depending the workload.

→ Let's say on the flipcart sale day he needs 100 servers and remaining days he needs only 10 servers. In this case 90 servers are sitting idle.

→ These all factors makes the loss to the company because to buy and setup 100 servers, is time taking and lots of money he has to spend to run a business.

→ To avoid this problem, the service provider is providing the same infrastructure on rent i.e. pay for what you use.

→ In a physical server, we can create no. of services for our use and delete them after completing your work, this creation of services in a physical server is called virtual Server, and it is given by EC2 service in AWS.

→ Advantages :-

→ Lots of money saving

→ mainly time saving

→ We can easily setup a virtual server with in < 1 min of time. and we can also delete.

Advantages:-

1) cost benefits

→ No upfront cost

→ pay for what you use

→ purchasing options

2) Automation options

→ Aws cli

→ Aws cloud formation

3) Global availability

4) Low total cost of IT

Cloud Providers:-

→ amazon webservices

→ Microsoft Azure

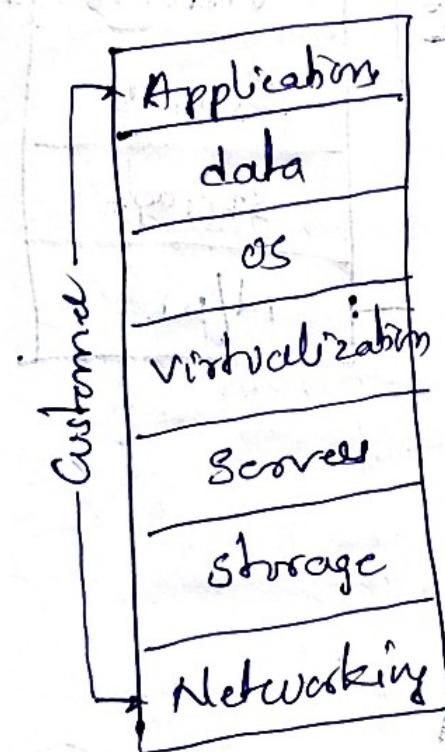
→ Google cloud Platform

→ IBM Bluemix

Service Models:

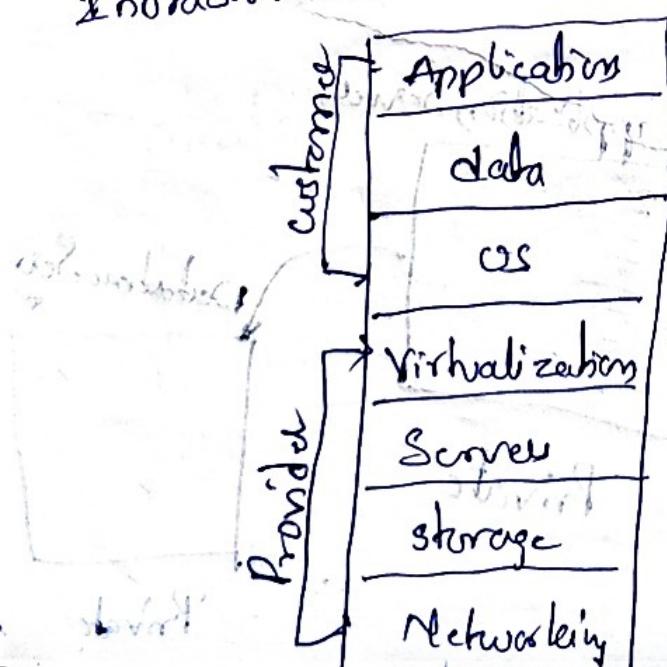
- IaaS (Infrastructure as a Service)
- PaaS (Platform as a Service)
- SaaS (Software as a Service)

Private Cloud



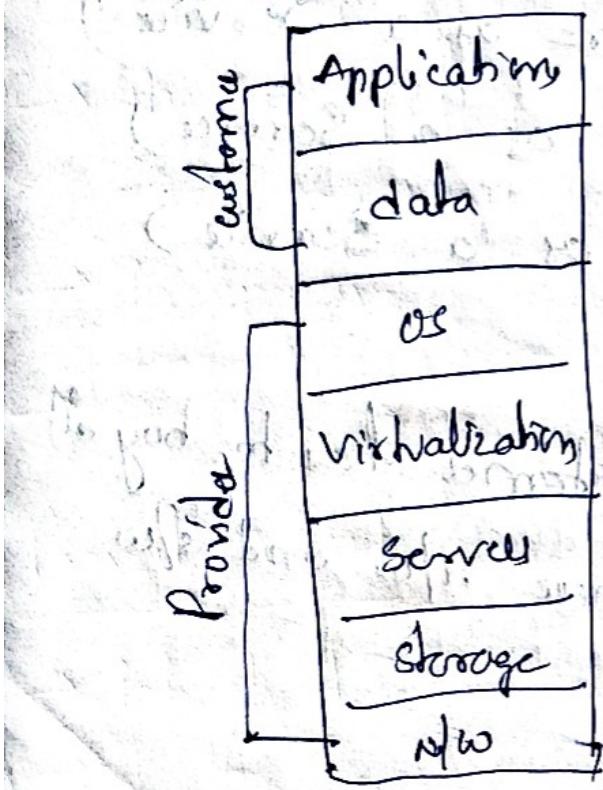
→ customer has to buy all these H/w and SW

Infrastructure as cloud

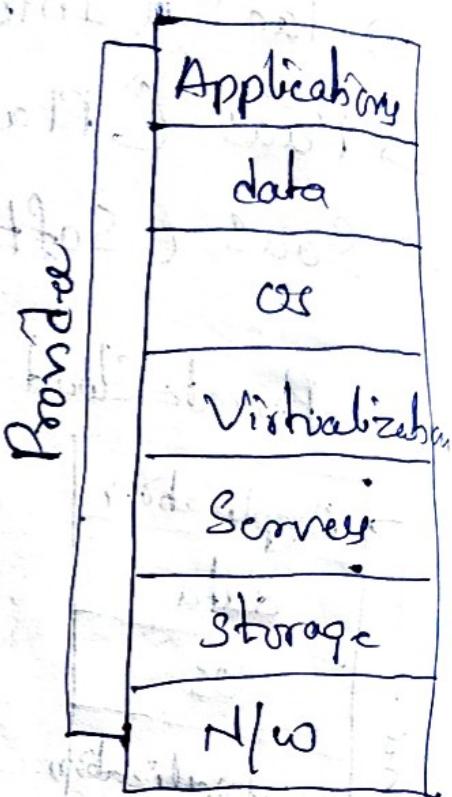


→ In this case, from n/w to virtualization, this H/w is provided by service provider, we are cutting our cost for this.

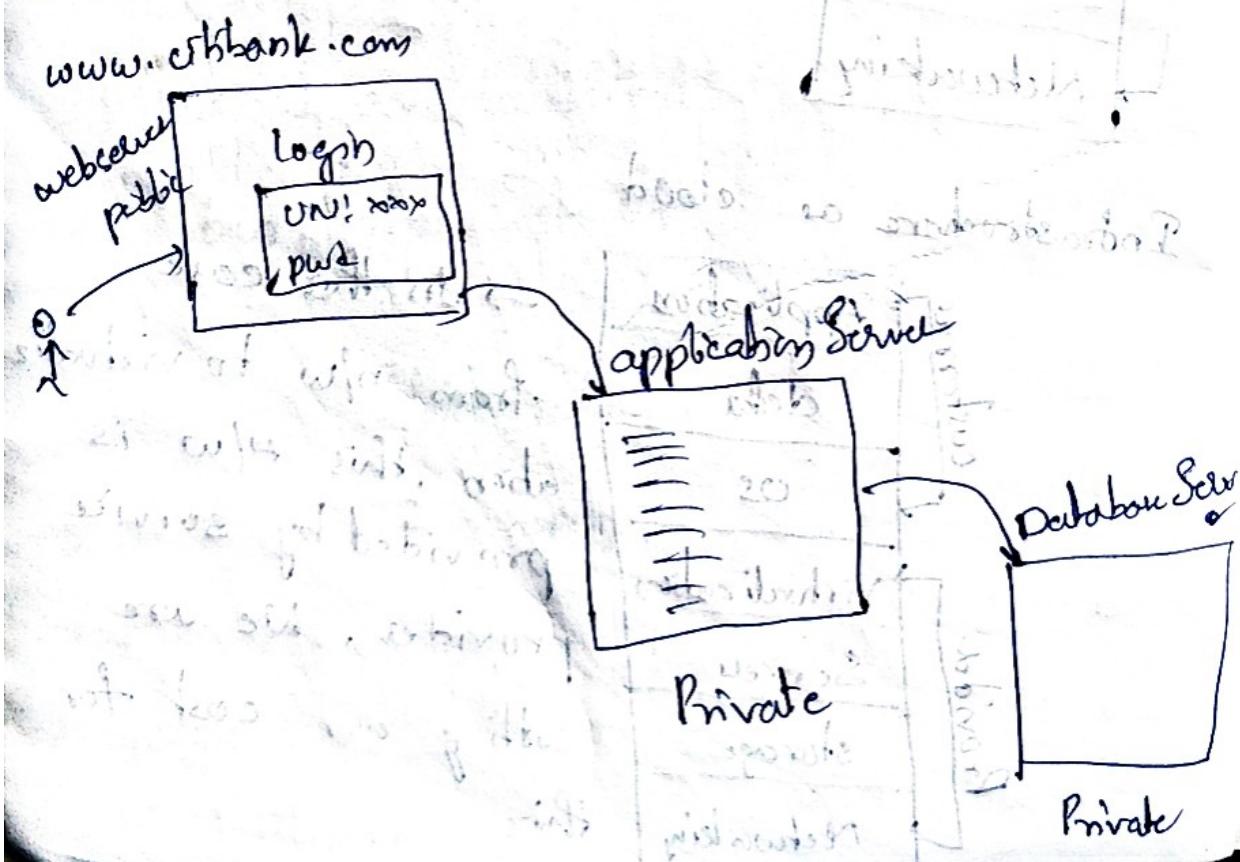
Platform as a Service



Software as a Service



3 tier presentation



- We can't expose application service (it has code for executing the application)
- Database servers are also not exposing to the customers; these two are private servers
- ~~Web~~ web servers are more in number for any application, whereas application servers and database servers are less in number
- If you are combining both, public and private servers is called hybrid cloud platform.
- Databases and application servers are created on premises, are also on public cloud platforms.

Public Cloud:-

- Infrastructure managed by cloud vendor
- The customer has no visibility and

- control over where the computing infrastructure is hosted
- The computing infrastructure is shared across organisations
- Inexpensive
- Secured

Private Cloud

- Infrastructure is dedicated to a particular organisation
- Most expensive and more secured
- customer has to manage the infrastructure and s/w.
- Build your own cloud using open source or commercial softwares like
 - (i) apache cloud stack
 - (ii) Eucalyptus
 - (iii) openstack

Hybrid cloud

- Organisations may host critical applications in private clouds and applications with relatively less security concerns on the public cloud
- The usage of both, private and public clouds together is called hybrid cloud
- Use databases and application servers on private and web servers on public cloud, this will for normal usage use private cloud and high/peak load use public clouds.

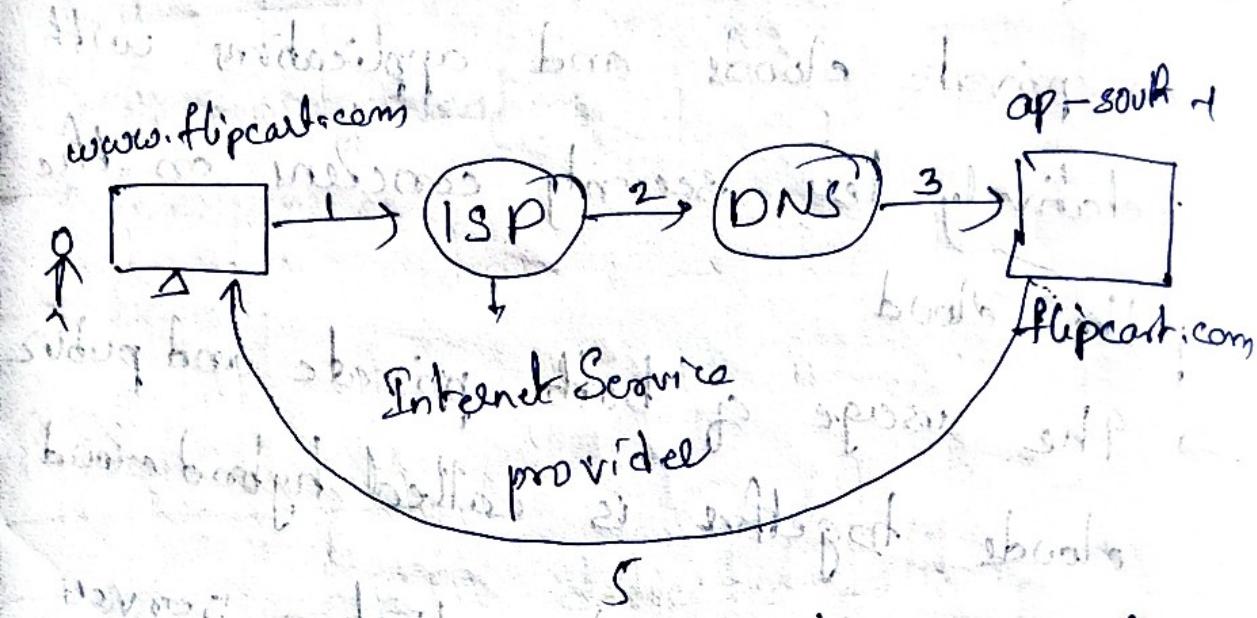
Aws Global Infrastructure

Region:-

It is a geographical location where Amazon is maintaining their data centers

- Now it has 23 regions

latency :- The time taken to process a request and give response is called latency.



→ If the server is near to you, the latency is less.

→ By creating more servers to reduce the latency we can just enable cache servers (in this servers frequently used

things are stored). The cache is calling things are stored as edge locations

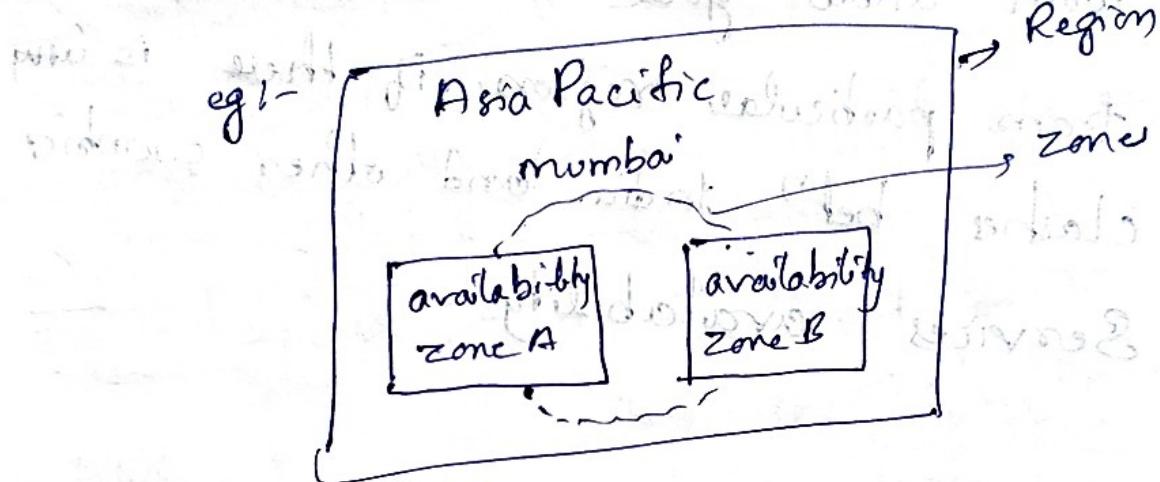
This cache server to enable we

have a service called cloud front,

they are calling this cache servers as edge locations.

Availability zones:-

The region in which servers are kept in a location is called zone.



→ If a client is in India, the server region to be selected near by him

Advantages:-

- Global reach to customers in minutes
- High availability across regions
- Disaster recovery sites

The region to use / on what basis region
is choosed.

- ① Proximity :— (Closest location to customer)
- ② Legal Compliances → this legal team won't allow you to use the servers from particular regions if there is any clashes betn India and other countries
- ③ Services availability

(EC2) Elastic Compute Cloud

- It is service provided by AWS, to create virtual machines on AWS cloud platform.
- Based on the client requirement create the instances
- EC2 is giving us pre installed o/s images
- Marketplace :- It is like playstore of AWS, we can select the o/s we want and install it on the instance

EC2 Instances = Virtual Server

- ① Resizable compute capacity in 200+ instance types.
- ② Reduces the time required to obtain and book new server by minutes or seconds.

③ scale capacity as your computing requirement

- to change

④ pay only for capacity that you actually use

⑤ choose Linux or Windows

⑥ Deploy across Regions and Availability zones for reliability

⑦ Flexible networking (NAT/elastic, VPC)

(Elastic IPs)

⑧ Support for virtual n/w interfaces that can be attached to EC2 instances in

your VPC

→ Based on the clients requirement,

the AWS has created instance family,

there are 200+ instance families, are

they

→ Similar kind of works that a instance can do, we group them as a Family

eg:- general purpose family

cpu optimised family

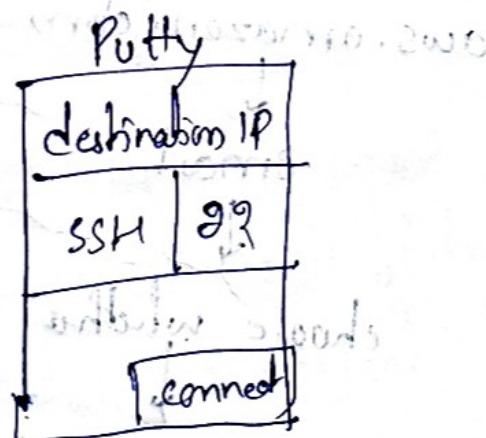
graphic card "

memory based " " etc

→ The default port no. of SSH is 22

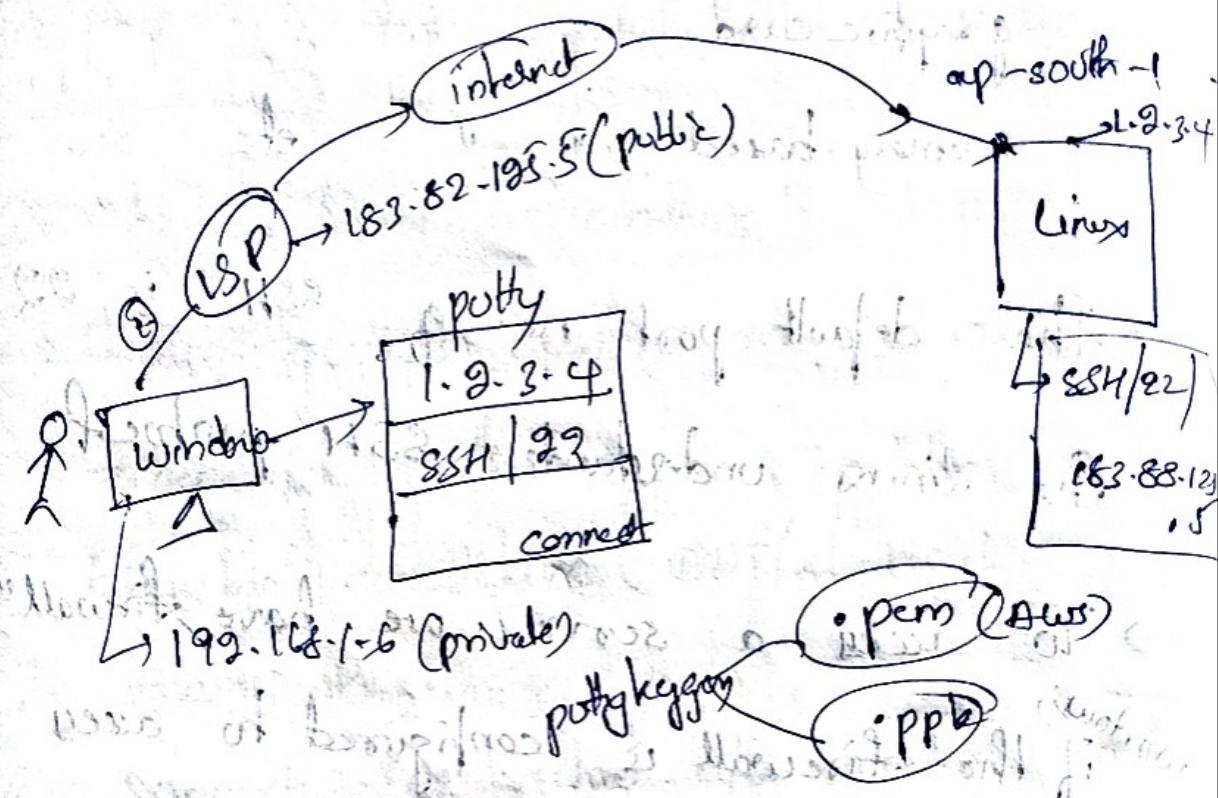
The Linux understands SSH protocol

→ To access a service we have firewall,
Windows,
if the firewall is configured to access
then only we can access where as in
Linux it is called as Security group



→ AWS will give the public key, the
extension is .pem (public encryption
module)

→ To convert .pem key into .ppk file
we need to have puttykeygen



→ To create AWS account, there are
5 steps

aws.amazon.com

↓
email

choose whether ① Personal ② Business

↓
payment

↓
personal information

↓
select plan

↓
complete sign up

→ To know the ip address of your laptop [ip config] in putty only

To connect to Linux server from windows laptop

→ We need to have putty and puttygen

to be installed
→ with the help of puttygen convert the public key pair .pem file into .ppk

→ In putty configuration

give public ip of your server

→ ssh

→ 22, port no

→ then select expand ssh which is present in left hand side

→ in that just select Auth

→ then provide private key file

→ open/connect.

Public ip :- It is used to access the server

from outside the AWS network.

Private ip :- It is used to access the server from within the network.

→ When you start a service, the AWS community will assign a public IP address to it. At the moment you stop the instance, it releases the IP address, because assigning a public IP to a stopped instance is no use.

→ Whenever you stop and start the instance we are getting new IP for the instances.

→ To avoid that problem we have

a service called Elastic ip address,

→ It won't change whenever you stop or start.

→ These elastic ip addresses are chargeable

Key pair :-

- the scope of a keypair is within the region.
- You can launch in no. of instances with one keypair in the same region only.
- The same keypair will not work in another region.

To create a elastic IP and add it to the server

- Network and Security
- select Elastic IP's
- click on Allocate Elastic IP address (this option is present on right corner)
- click on Amazon's pool of IPv4 addresses
- Then click on Allocate button on side
- Then an elastic ip is created.

Now attach that ip address to the server.

→ select the elastic ip address

→ Now go to Actions

→ select Associate Elastic IP address

→ click on instance and select the instance

→ Associate

For deleting elastic ip address

→ select the elastic ip address

→ go to Actions

→ select Release elastic ip address

→ In organisation, to change anything we have to raise a request like CR (change request)

Interview

Sir on the daily basis we face some issues

happening we will be raising a CR to AWS team, then they will solve it.

- yum is a installer tool in Linux
- It will download, install, verify the package. (automate)

Note:-

- In production we use only "rpm" to download and install any package.
- yum installs whatever dependencies are in rpm. the package also where as in rpm. download whatever packages we select.

rpm is also installed tool (manual)

- rpm is also installed tool (manual)
- How to host/configure/install a website

- on Linux server:
- install httpd (to configure a webserver)
- by default httpd will be in inactive state
- To check the status of httpd
- service httpd status/start/stop/restart

service : service name , status

→ Now start the httpd

service httpd start

→ All the webpage are created in a

location called /var/www/html

→ Create a index.html file for webpage

→ After changing anything we have to

restart the service

service httpd restart

→ chkconfig httpd on → whenever anybody reboot the httpd will go in dead mode

→ To see the website, go to security tab

group and add a rule like

http [80] anywhere

→ From browser give public-ip : 80

How to automate websaver configuration

→ To automate we use shell script in linux.

→ shell script:- sequence of linux commands with some logic is called shell script

Rules or conditions to write shell script

• File ends .sh

① file must end with .sh

e.g. - websaver.sh

② The first line of a script must and

should start with

#!/bin/bash → shabang line

We are telling our o/s, that I want to

write many linux commands, the commands

library is present in /bin/bash only. Then

o/s will picks them and use it

vim webscrevd.sh

#!/bin/bash → no space in the firstline

yum install -y httpd git

service httpd start

chkconfig httpd on

cd /var/www/html

git clone pastur •

service httpd restart

→ write above script in notepad and save with extension as .sh

→ while launching the server attach that file to it, that's all, then take public ip and port no. and access it from browser

→ while launching server, you find advanced details section

- go to user data section → select As file and launches the server.
- At the time of booting the server we are installing the script.
- it is called bootstrapping
bootstrapping = attaching

Scenarios

- ① We deployed a website on prod server and it is running successfully.
- ② Now client came and asking I need to include few more modules in the same website, what will you do?
- Never add those two modules directly on prod server because there may be compatibility issues, there may cause the on going business to go down.

→ Then go for testing in low production servers, as QA/deployment servers.

→ Now development teams is asking prod manager, that create a copy of prod server and give it to me, we will ^{add} two more modules and test it, if it is going well, we will let you know.

→ Then prod team will create a Image out of that server ie AM2.

→ Now the development team can ^{create} server from that AM2.

AM2 → s/w library + softwares installed
direct as image, can be copied but direct as server can't be copied

→ If we have AM2, we can create a server.

AMIs creation

- How to copy an instance? → AMI
- How to backup of instance? → AMI
- How can you do backup or restore prod server? → AMI.
- How to migrate an instance? → AMI

Scenarios!

- ① If prod servers are in same aws account and same region, → there no need of copying → he can create AMI out of it and use it
- ② Both are in same aws account and both are in different regions?
 - We need to copy the AMI
- ③ Imagine that prod servers are maintaining separate AWS account, non-prod servers are maintaining separate AWS account?

lets say TCS is maintaining AWS account
and Bajaj insurance has gave project to
implement and give it me and he is
also maintaining AWS account.
→ then TCS will create a AMI and it
will be copied into Bajaj companies
account.

Creation of AMI

→ let say you have prod server.
→ select that server and go to Actions
→ select image and template → create image
→ give some name
~~No reboot~~
 enable
when you enable it, then it wont reboot
production servers, reboot will be done
by authorization or planned downtime

→ create.

You can find the AMI in Images section

Copy a AMI into other region

→ select the AMI

→ go to Actions → select copy AMI.

→ select the destination region and specify your regions

→ It will take 30 to 40 mins depending on the configurations etc of that AMI.

copy AMI from one AWS ac to another

→ select the AMI

→ go to Actions → select Edit AMI permissions

→ select private

→ In Shared accounts → select Add account

and the AWS account ID of others

persons

Elastic Load Balancer:-

Scenario:-

- In production environment, we have two servers, if the load is increasing on them either of the server may crush or it will go down.
- So that case we have to create one more server and do all necessary configurations, this may take time and business will go down.
- To overcome such problems, elastic load balancer is used, it will distribute load to the servers equally.
- If we have more than one server, we must have the loadbalancer.
- Now create two servers and one load balancer

parameters to be configured in Load Balancer

see
Ping

① protocol : http

② port : 80

③ path : /index.html

④ Interval : 30 sec

→ Load Balancer will ping the servers in every 30 sec

⑤ Response time out : 5 sec

→ It will ~~receive~~ receive the response

for every 5 sec

⑥ Unhealthy threshold : 2

→ If the first server or second server

is not responding for first 5 sec

it will show other server as it

service and next time after 5 sec

it again not responding , then

Load balancer shows that ~~the~~ server

is in out of service

→ Then the requests of user are not sent

to unhealthy server and all are forwarded to next server.

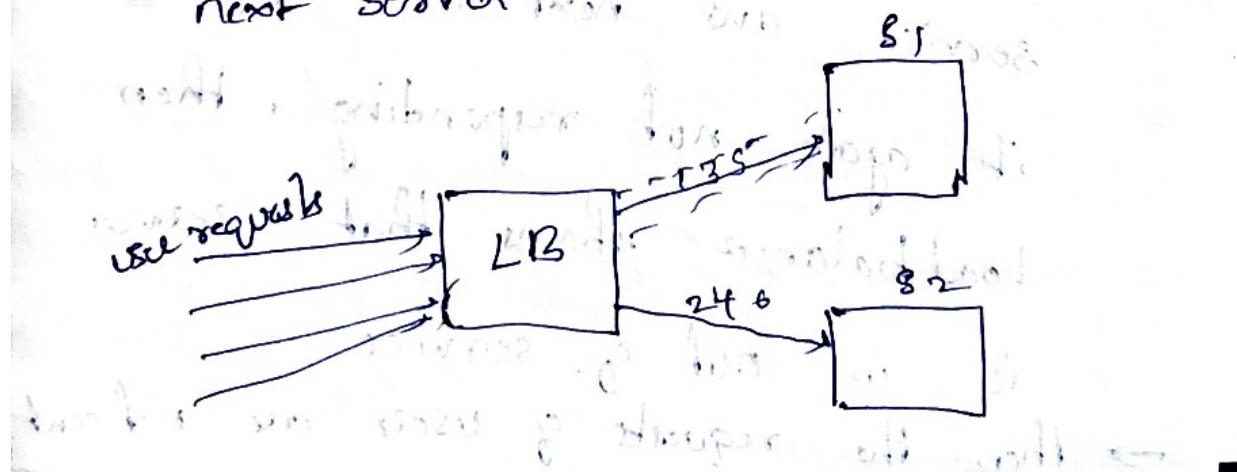
⑦ healthy threshold: 2

→ Load balancer keeps on pinging it when the server comes into working state, it will show as in-service after 2 thresholds.

⑧ connection draining

enable it:

→ It will send less no. of requests to unhealthy server in just 5 sec, then it shows as out of service and all the requests are redirected to next server.



① Cross zone load balancing

- one service may be in 1a zone
- Second service may be in 1b zone
- But both must be in one region only
and load balancer too.
- ~~We can add~~ ~~number~~ services to a load balancer
- If both services are in different zones
~~Protocol~~ enable it.
- ② On what basis this distribution of load happens?

Ans - Round Robin

- even distribution happens

Load balance configuration

- Create two services and one load balancer
- Now for our understanding, create two

different websites and deploy one in
webserver1 and another in webserver2

```
#!/bin/bash
yum install httpd -y
service httpd start
chkconfig httpd on
cd /var/www/html
echo "welcome to webserver1" > index.html
service httpd restart
→ and do bootstrapping for one server
```

```
#!/bin/bash
yum install httpd -y
service httpd start
chkconfig httpd on
cd /var/www/html
echo "welcome to webserver2" > index.html
service httpd restart
```

→ do bootstrapping in server 2.

→ In lefthand side search for load balance in AWS. and create load balance. and configure as required.

Scenario:

→ one user is accessing a sbi bank's server and doing transactions, he came to 1 page and waiting for otp and enter, meanwhile the service has gone down,

→ In that case the request will be redirected to next service, but the user doesn't know where the stopped user wants to is seeing for otp, then starts from starting of the server page, we have to fill all the details and go upto otp page, and so on.

→ To avoid this AWS has given

"stickiness", enable it then next service

starts from the same page ~~to~~ where
the user has stopped to do
stickiness. —

- enable it in loadbalance;
- when you enable, the AWS is maintaining buffer area, in that all the data is stored,
- when one server is going down the next server will look into buffer area and starts from the same page where the ~~other~~ other servers has stopped.
- go to instances and scroll down, you will find the stickiness → enable it depending on the application.
- You can ^{have} merge no. of servers in the same load balancer.

④ We have configured on load balance and running successfully, at any time a service may down or go in out of service, then how cloud engineer know that the service has gone down?

→ He can know when login into this account, and in the Loadbalance section he can see instances of status, all he can see instances of status.

⑤ Is this possible to know the status 24/7 manually?

→ It is not possible,

→ Then we depend on monitoring tools

called (i) Nagios

(ii) bmcpcatof

(iii) splunk

(iv) prometheus

→ These tools are not free,

→ If there is any error, it generates Alarms

Cloud Watch

- It is a service given by AWS to monitor the various resources in AWS.
- It is a managed service.

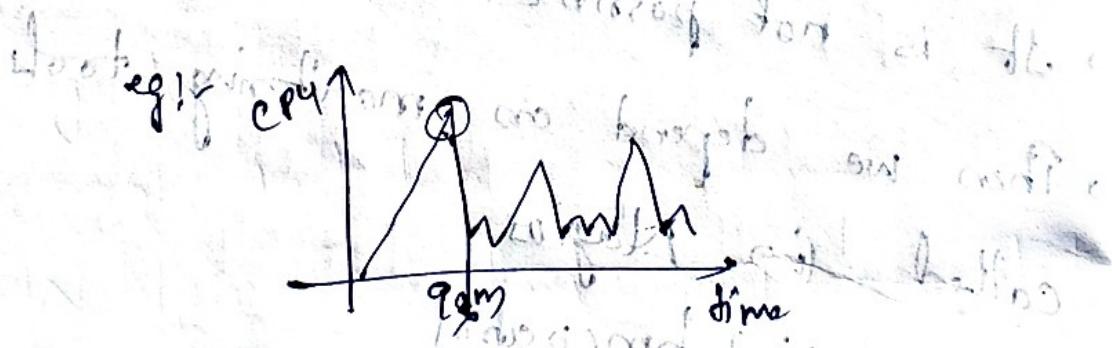
① How to monitor a Linux server?

- based on a metric

② Collecting information of a server

- all the information is stored in log files

③ Use visualization tool



take 9am time before 5wins and
after 5win to 9am, and analyse the
problem and troubleshoot it.

④ Root Cause Action

- Understand the problem and take action

Cloud watch configuration

→ Write condition for this, as

whenever $\boxed{}$ of CPU utilization $\geq 70\%$ for
min
max
avg
sum

- $\boxed{1}$ interval of $\boxed{15 \text{ min}}$

→ It is collecting info and storing in datapoint



→ min means it takes 40 and compare with threshold 70%, if the condition satisfied it shows the status as

$\boxed{\text{OK}}$

→ if not satisfied, the status is

$\boxed{\text{ON-ALARM}}$

→ If the service it self crashes, it shows

$\boxed{\text{INSUFFICIENT}}$

Note:-

→ In real time we use avg or min

→ better set avg

* If alarm is generated, how do you know?

→ For this I must get a notification as email or sms

Simple Notification Service (SNS)

Scenario:-

In e-commerce portal, let's say electronic goods,

→ event: Instock

→ **Notify me** → whenever product is in stock

email id

mobile no.

confirm
subscription

→ All these details will be given,

→ whenever the concerned product is in stock

093314-1107

it will send a email or alarm to the customer.

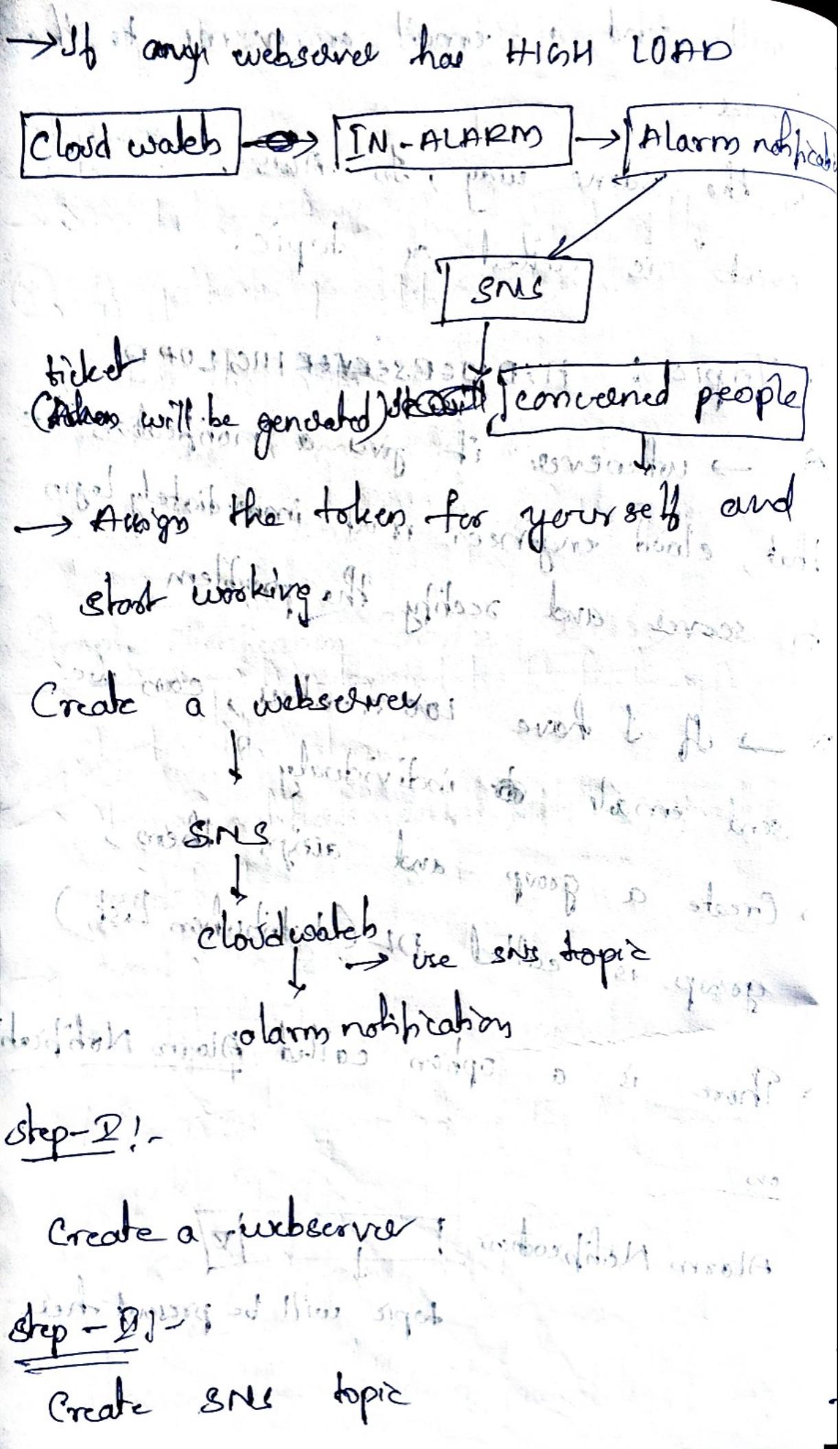
→ In the same way, in AWS, "event" is called as topic.

① Topic : HYD WEB SERVER HIGH LOAD
→ whenever it gives a notification like that, cloud engineer must immediately login to sevice and rectify the problem.

② → If I have 100 members, can we send emails individually?
→ Create a group and assign them, group is called DL (Distribution List)

→ There is a option called Alarms Notification

on Alarms Notification []
topic will be present here



- go to SNS
- click on create topic
- give topic name
- display name (with what name the message is to be displayed)
- In real time enable encryption
- In Access policy → choose basic ✓
If you are programme choose Advanced ✓
- Define who can publish messages to the topic
or everyone
- Define who can subscribe to the topic
or everyone
- Delivery retry policy:-
How many times you can send message whenever you send the message their won't be network in that area.
- enter max no. of times you want to send
- Now create subscription
- Protocol → select email option

→ endpoint → give email address

Step-IV:-

→ Create a alarm in cloud watch

→ go to cloud watch

→ select the metric

(or)

→ select the EC2 on which you are setting alarm

→ Actions → Monitors and troubleshooting → Manage Cloud Watch Alarms

→ Alarms notification → select the SNS topic you created

Alarms action:-

Be careful for enabling this.

Note:- Don't enable it without knowing the problems to do.

→ Without analysis don't enable it.

Select the parameters

select 5 minutes → Interview

→ 5 minutes is free upto 10 alarms, after that they are chargeable

To check the CPU utilization increase → go to EC2 server and install stress package.

To install stress package, we need to have Amazon Linux - extra install epel, sudo yum install stress epel → extra packages for enterprise Linux.

* top command will use task manager in Linux

Now execute the command

stress -c 100 --timeout 1800 &

no. of times

to run in

Now see the top command background and see CPU utilization

Note:- Don't execute stress command in real time

→ Never execute apt-get update command on production servers, because update takes more time, until the update is executed server won't be available for the business

→ mr-rf command to stop the process, we use kill command, if the 100 processes are running each process has id, we can't give 100 id's to do so

killall stress

killall process-name

Scenario:-

- whenever a application is running on a group of servers.
- Let's say 3 servers, then if one server needs to be replaced, then to create a server and replace it, is a time consuming process in the mean while my business will get down.
- The creation of server must be automatically happen.

Requirements:-

- whenever the load increase, the server must be created automatically.
- whenever the load decreases, the servers must be decreased automatically.
- whenever the service goes down loadbalance will sense it.
- For this we need to have
 - (i) Create AMI

(ii) Launch configurations

→ Server (AMI), config, EBS, tag, Security group, key

(iii) Load balanced:

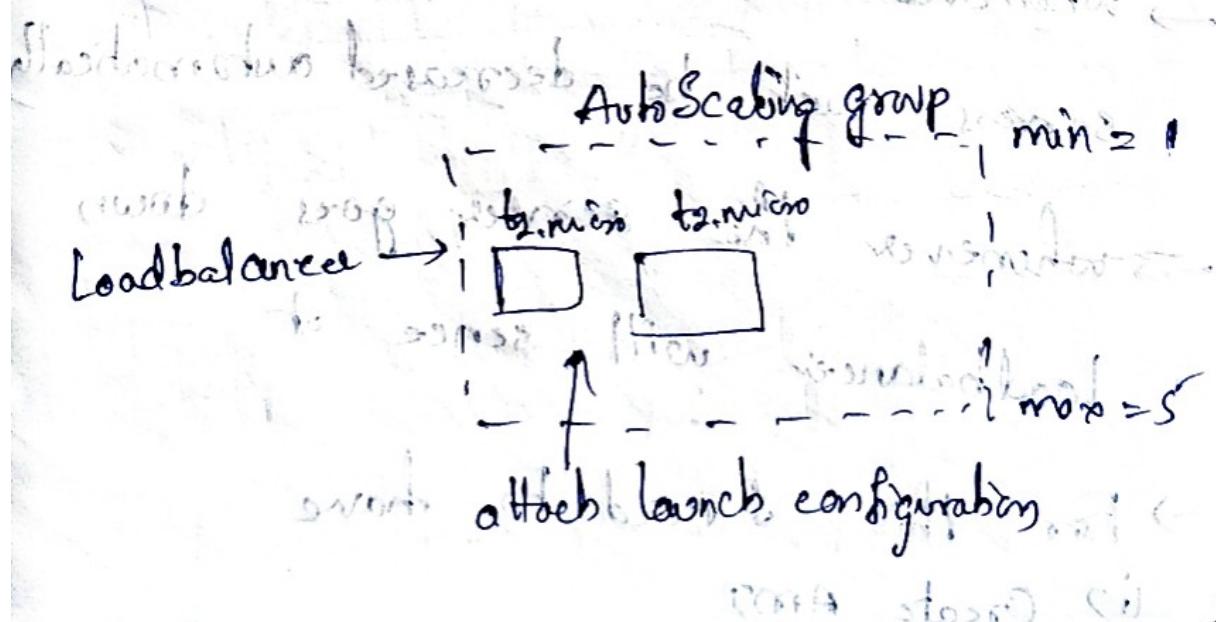
→ It must be empty at starting point.

→ The addition of servers in load balancer must be automatic

(iv) We need auto scaling group
→ whenever load increases or decreases
the servers will come up automatically
is called Auto scaling.

(v) scaling policies

Create auto scaling group



→ write scaling policies, i.e. write condition for increase and decrease the servers

② Policies - (i) if demand is high
(ii) increase group size i.e. adding more no. of servers

(iii) CPU utilization $\geq 75\%$ for 1 hour

- Take Action after 5 min
- (i) Add instances
- (ii) Add alarms, SNS

(iv) decrease group size, write condition whenever CPU utilization increases in my case it adds servers

→ when the addition of servers reaches the max limit you specified,

→ The autoscaling group is trying to add

more servers but addition of servers reached max limit; then it shows

ERROR i.e. failed to launch

↓ Interview

→ decrease group size whenever large g. CPU ≤ 30 to 1 g size.

Takes Action remove 9 Instances
or your choice

→ whenever it reaches min limit while removing, it throws an ERROR
i.e. failed to terminate

~~Procedure~~ How to upgrade auto scaling group servers?

→ launch configuration

→ Now change values min=0 max=0

then old servers will be removed, downtime: only

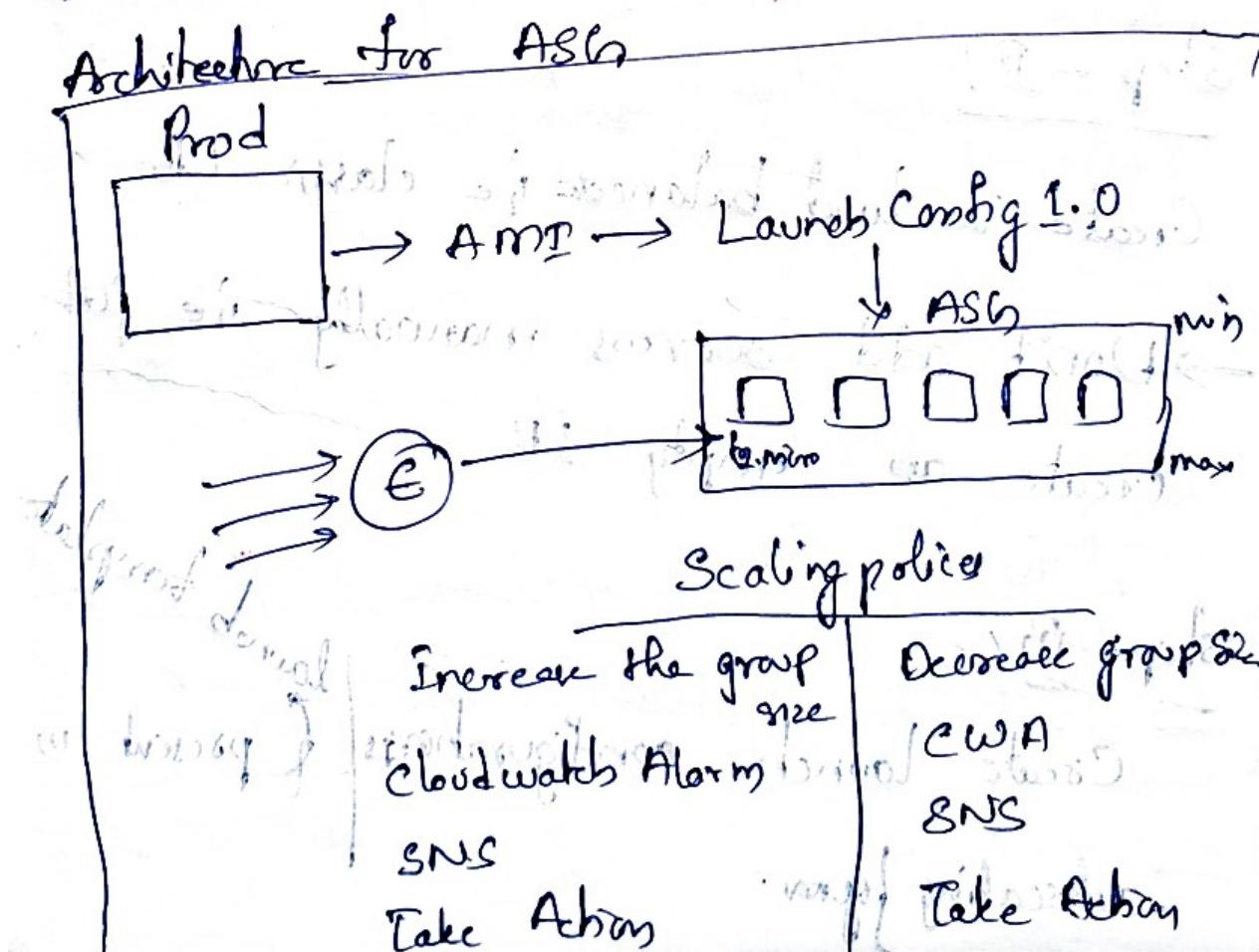
do this in planned do this in planned

attach new launch configuration.

→ Then attach new launch configuration to the auto scaling group and now

change values min=1 max=25

- ~~Batchup~~
- * You have set it for CPU utilization but the service goes down due to another reason i.e. a server crashes what is the autoscaling group do?
- The autoscaling group has ~~self healing~~ mechanism i.e. the server will come up, if a server is going unhealthy the autoscaling group will sense it and replace it with healthy server
- It is called ~~horizontal Scaling~~



Target Tracking :- It is a service given by AWS to write policies automatically → just add metric and threshold.
→ no. of servers to add and remove, adding alarm etc there all be taken care by target tracking.

Step - 2 :- Select PROD servers and create AMI out of it. I don't think nothing is left.

Step - 3 :- Create a load balancer i.e classic LB
→ Don't add servers manually i.e just create an empty LB.

Step - 4 :- Create launch configurations | launch template { present in autoscaling group.

- Click on "Create launch config".
- Select AMI for which you have to create.
- choose instance type.

Monitoring:-

- Enable EC2 instance detailed monitoring within cloud watch.
- EBS-optimized instance
- Launch as EBS-optimized instance

→ These all are chargeable.

~~Additional~~ If you have webserver AMI,
① Already you have webserver AMI.
How can we add additional s/w's

on top of same AMI?

→ While creating launch config, in advanced settings details → write a script or attach the s/w file in a script or attach the s/w file in user data (will give it to each instance).

- Like this you can install more s/w on top of AMI.
- select all the requirements you want and create launch config

Step - IV:

Create ASG

- select Launch configured → Actions → create Auto Scaling group

→ Name

* You are creating min 10 servers in ASG, should call the servers in one zone.

→ No, we should not create 1st zone,

because to have high availability.

* Select atleast two or more zones

→ Attach to a existing load balancer

→ select classic LB

Health checks

ELB

Additional settings

Group size

Desired capacity :- It is the value between min and max.

→ We can't go beyond the desired capacity,

min

max

desired capacity

e.g. - desired capacity

2

/ or give value between them as you wish

min

1

max

5

Scaling policies:-

Target tracking scaling policy

None

- If you select none, you have to write policies
- If you select target, just select metric and threshold, other policies are written by AWS

~~Warms up instance~~ ^{Interval} seconds

- The time taken by the service to get ready to work
- give the value depending upon how big is your application

Instance scale in protection

whatever instances added, you don't want to remove them, then enable it

scale in :- decreasing/removing

scale out :- adding/increasing

Add notifications :-

select the events

Add notifications one for each event

Tags :-

→ It is important, give some tag name

key	value
Name	ASG-Survey

Dynamic target policy :-

→ write your own policies

→ Name

→ select step scaling

→ Take the Action

Add survey

→ Depending on client use target or

dynamic policy

→ No. of instances - survey

→ survey out of baseline info

→ out of baseline info

→ baseline under different env

Snapshots - Storage:-

Q How many types of hard drives?

- Tape drives, cold ~~hot~~ HDD, throughput optimised
- SSD's → 2000 i/p, o/p/s/sec
- In SSD's → different types of SSD's.

Scenarios:-

Two teams in the organisation require 5GB each to store their data,
→ To allocate this type of volumes to the teams is done by EBS.

EBS (Elastic Block Store)

→ It gives types of storage devices, select the type and use it.

→ In /dev/xvda → it has all the s/w's required to the server, therefore don't store your application data in it.

xvda → extendable volume virtual device a

→ It is the first volume attached to any server

Note:-

The instance and the volume should be in the same region and same availability zone.

→ It is a disadvantage

→ first volume name is /dev/xvda

→ second volume name is /dev/xvdf

Q) How can we do partition of 10GB and allocate to each team

→ fdisk is used to do partition

fdisk path name of volume/hdd

e.g:- fdisk /dev/xvdf

: n → new partitions

: l → it will ask for partition number, if

starts from 1

: t → size you enter

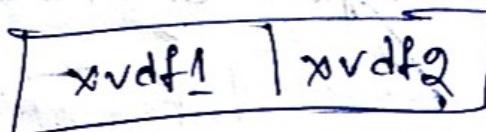
→ For next 5GB block add file 3.

: 0 ↘

: 2 ↘

blocks after which standard off
will be used to come out

: w ↘ → to write 10GB capacity same as off



② format :-

To make any file system to know the OS
as ~~FS~~ → extensible file system

mkfs -t xfs /dev/xvdf1

mkfs -t xfs /dev/xvdf2

→ Now OS knows what kind of file to
store in this

③ mount points:-

→ There are nothing but directories

mkdir /shiftA /shiftB

application name/give any name

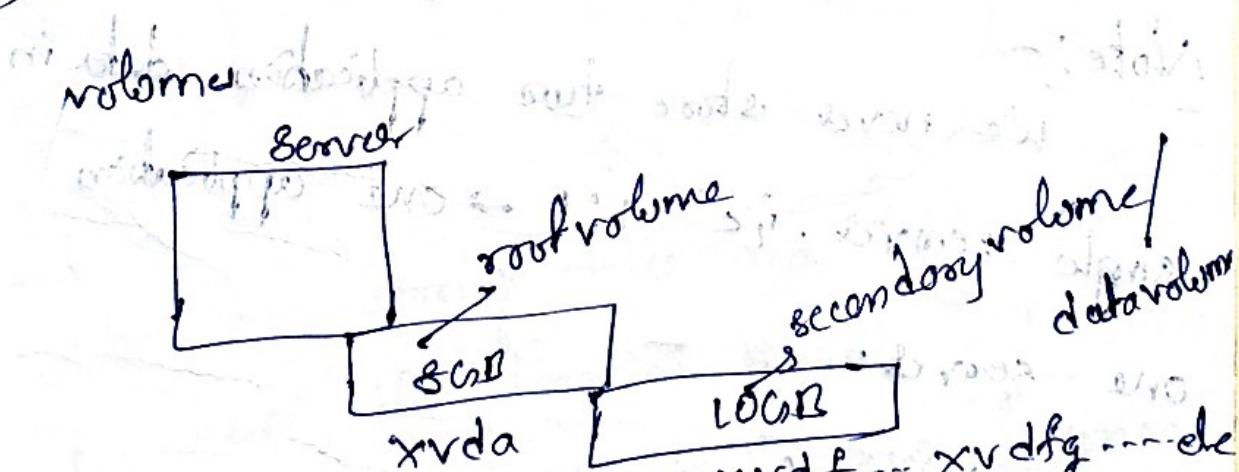
→ Now, give shiftA directory to one team
and shiftB directory to another team

④ mounting :- attaching these folders to
teams

```
mount /dev/xvdf1 /shiftA  
mount /dev/xvdf2 /shiftB
```

→ We have service called ebs, with
which we are creating volumes and
attaching the volume to the server

→ while launching the server we create



→ We can attach 12 volumes maximum

for a server

- The max size of root volume is 2TB
- The max size of secondary volume is 16TB

Step 00 - Root Volume

- To see the disk information

[fdisk -l]

- Select the ~~server~~ and go to ~~CBS~~ service
- first create a volume and attach to the server

Note:-
We never store two application data in single server, i.e. \rightarrow one application one server.

- Create volume.
- After that select the volume \rightarrow Actions
- attach volume \rightarrow select the server

→ Now go inside the screen and partition

① Partition

stick pathname of HDD

[fdisk /dev/xvdf] ↪

: n ↪

partition type p : ↪

first sector : ↪

: +5G ↪



→ We are storing files and folders in the HDD, the location of file and folder path/address is stored in reserved space.

For 16 GB pendrive, we get 14.83 gb

free and remaining gb is reserved.

: w ← for saving all three partitions

② Formatting:

mkfs -t ext3 /dev/xvdf1

mkfs -t xfs /dev/xvdf2

label
DF-NT

③ mounting points:

mkdir /shifta /shiftb

④ mounting:

mount /dev/xvdf1 /shifta

mount /dev/xvdf2 /shiftb

→ Now go inside the mount point and

add your data

cd shifta

→ touch file{1..100}

→ It will create 100 files

→ go to another mount point shiftb and create files

lslblk → To see the files added to volume.

→ umount is used to unmount the files.

→ To keep the mount point permanently → go to /etc/fstab and write as follows:

Snapshot → /dev/xvda1 /filename xfs defaults

Taking back up of a volume is

called snapshot.

→ We can read only the data but we can't write.

→ To read write on the data, create a volume and attach it.

→ Create a volume out of snapshot and attach it to the server.

Note:- After attaching the volume from

snapshot, don't create partitions, if you

do so data will be lost.

→ partitions X → do only mount

→ formating X

→ when you mount only you can see the data.

Intent: - transfer of data and migration

① How to copy data → snapshot

② How to backup and restore data → snapshot

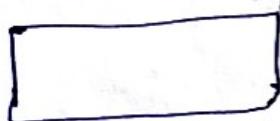
③ How to migrate the data → snapshot

Scenarios:

① Same AWS ac, same region

② same AWS ac, different region

③ Different accounts



① create a snapshot

PROD-DAT-VOL

(snap1)

create volume for QA server

create QA server ?

attach Q/A volume

mount on QA db

X adding
X removing

- Now you can copy into another region
- To copy snapshot into other AWS account

→ click Actions → modify permissions →
give **AC DP**

- We are using EBS volumes

dynamic data

- To store huge static data, we have

a service called S3

S3: (Simple Storage Service)

- Huge static data can be stored in it.

e.g.: google drive

- It is called as "bucket"

- It can store unlimited data without cost

- file is called as "object"

→ Each file can be of 0 - 5TB, like this we can upload unlimited files.
eg:- netflix, gaana.com etc

* How many buckets can be created in AWS ac by default?

→ 100

* How can we create more than 100 buckets?

→ Yes

→ By raising a support ticket.. feature called increase limit

→ 1:1 bucket is required i.e 1 application 1 bucket.

Scenario:

→ whenever you upload a movie in netflix; many users are accessing it,

→ There shouldn't be latency.

- It must be accessed to everybody without latency.
 - For that we have "standard storage class".
 - storage class:-
 - standard is a default storage class.
 - They are using SSD's only (Solid state drive).
 - They will maintain 9 copies of a file by default.
 - It is managed by AWS.
- Scenario:-
- For the first 30 days everybody is accessing the data/movie.
 - After 30 days very less no. of users are accessing it, then why to maintain 9 copies and pay for them.
 - So I want to reduce copies and cost cutting also.
 - For that we have a feature called Reduced Redundancy Storage (RRS).

Reduced Redundancy storage - (180 days)

Redundancy means duplicate

- RRS means less no. of duplicate copies
- By default it is maintaining 2 copies
- Therefore cost is reduced
- Thus, write one policy i.e. if my data age is more than 30 days charge the class

Scenario!

- If my end users are accessing it very rarely, why to maintain 2 copies and pay for it.
- For that we have a class called glacier

Glacier - (slowest movement / slowest storage)

- It maintains only 1 copy
- It is very slow to access the data.

- Cost will be reduced, archive your data
- Tapes are used in this case.
- It takes (3-4) hrs to access the data
- Raise a request, After 365 days delete it or not upto you

DLMP → Data Life cycle Management Policy

- 90 days (standard) — 180 days (RRS) — 365 days (Glossary) → delete it
- Writer DLMP → based on the wish of client

To overcome the DLMP, AWS given standard only sub clause

standard Infrequent Access
6 copies

standard

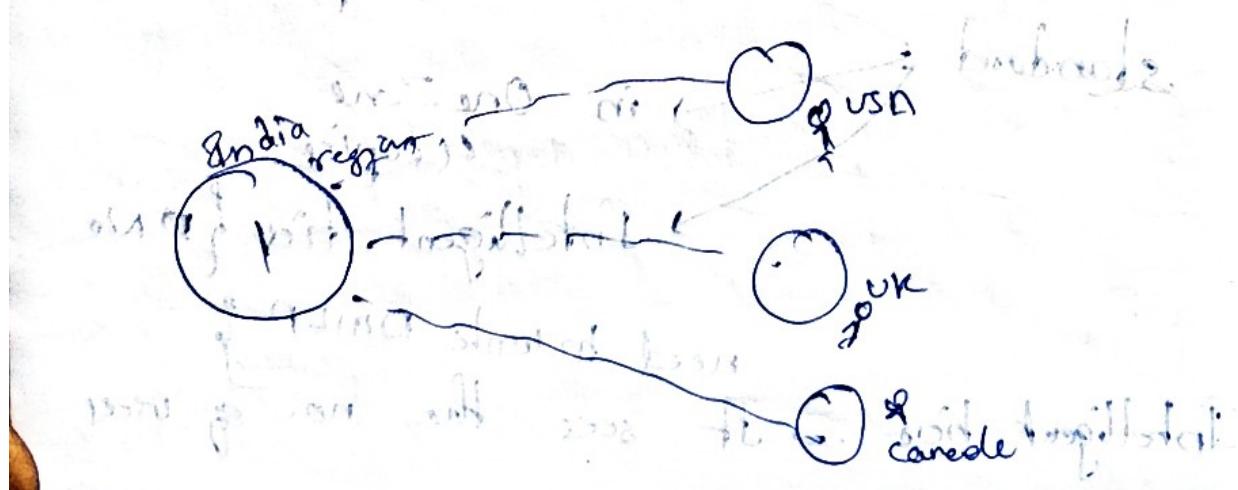
in One Zone
4 copies

Intelligent tiering { No

need to write DLMP

Intelligent tiering → It sees the no. of users

- occurring it, if more no. of users are accessing it, it changes to standard IP, if no body occurring it changes to glasse
- Here we no need of writing DMIP, depending on the user's load, it changes automatically.
- store 180 days in intelligent file
- 1 day store in standard, and next 180 days store in Intelligent file, then IT will take care, billing is reduced
- Intelligent file, use it in organisation to reduce the cost, therefore no need to reduced again and again.



→ The objects are called static object store.

Security:-

→ In S3 bucket, there is encryption for

→ Two types of encryption:

① S3 service is giving encryption algorithm

algorithm

→ S3 is additionally charging you.

② Customer can manage encryption algorithm

ie he can buy encryption algorithms

from third party and store, with that

only he can encrypt his data and store

→ Encryption certificate can be uploaded

→ KMS / customer managed keys

→ The organisations which are using second

option, every week you have to update

the algorithm SSL (Secure Sockets Layer)

③ Property:-

→ We can also host a static website on S3, without going with EC2.

- enable the static website ○ by default it will be disabled
- It will give you URL, not the IP address

① index page

index.html

(or)

homepage.html

② Error page

error.html

→ Save

③ All the data of your website upload in the S3

→ With this we are saving lot of money

S3 property:-

→ I want to get notification, if anybody doing something in S3 bucket

→ For that we have Event notification

②. Property! - Versioning

- The older and later versions of your data are maintained by s3
- enable versioning while creating s3 bucket
- Once you enable, there is no option to disable, if you disable it, the older objects will be lost
- enabling is chargeable
- We can suspend the versioning,
- whenever you click on suspend button, the older versions will retain and next object onwards there is no versioning.
- S3 service is global; i.e., we can access it from anywhere.
- S3 bucket creation: how to do?
 - click on create bucket
 - Bucket name

Interview

→ while creating new bucket, can you copy existing bucket settings?

→ Yes, click on copy existing bucket settings and click on the bucket you want.

→ SACL (Access Control List)

→ Block all public access.

→ Bucket Versioning

enable

→ Default encryption

Disable

enable

→ In real time enable encryption

Interview

→ If anybody doing deleting the objects in

s3, how can you block it?

→ Yes, by object lock.

→ create bucket.

Permissions -

- Others to access the bucket give the permissions
- unblock public access permissions and write a policy in json format
- In real-time scenarios don't enable public access
- give the ARN of your bucket and don't change anything

```
{  
  "Version": "2012-10-17",  
  "Statement": [{  
    "Sid": "Public ReadGetObject",  
    "Effect": "Allow",  
    "Principal": "*",  
    "Action": ["s3:GetObject"],  
    "Resource": ["arn:aws:s3:::nithu  
      ckethyd/tut/  
      /folder1/*"]  
  }]}  
}
```

- On which folder, file you want to give access, then specify the folder, filename in the code
- Now you can open the file through url given in the bucket.

→ To suspend the bucket versioning →

Go to properties of the bucket



Create event Notification

→ select whatever option you want, this will send a notification

→ create SNS topic

Transfer acceleration

→ To transfer data faster or, enable it



Static website hosting

→ By default, it is disabled

- enable it and host static website
- enter two files
 - (i) index.html
 - (ii) error.html
- It gives an ~~401~~, it throws an error which means your project folders are not added, 404 Not found error.
- Then upload your project files in the S3 bucket.
- Then access it with url, and you can see the static website.

Note: If and only if static website choose S3 to host it.

Management:- (Lifecycle rules)

- create lifecycle rule

- name

Lifecycle rule actions

→ select all actions

Transition current versions of objects between storage class

choose storage class transition days after object creation

Intelligent-Tiering 20

Transition non current versions

These are older versions of objects

Glacier shadow 180 5

Expire current versions of objects

→ want to delete the current versions after 8 years

days
730

Permanently delete non current versions

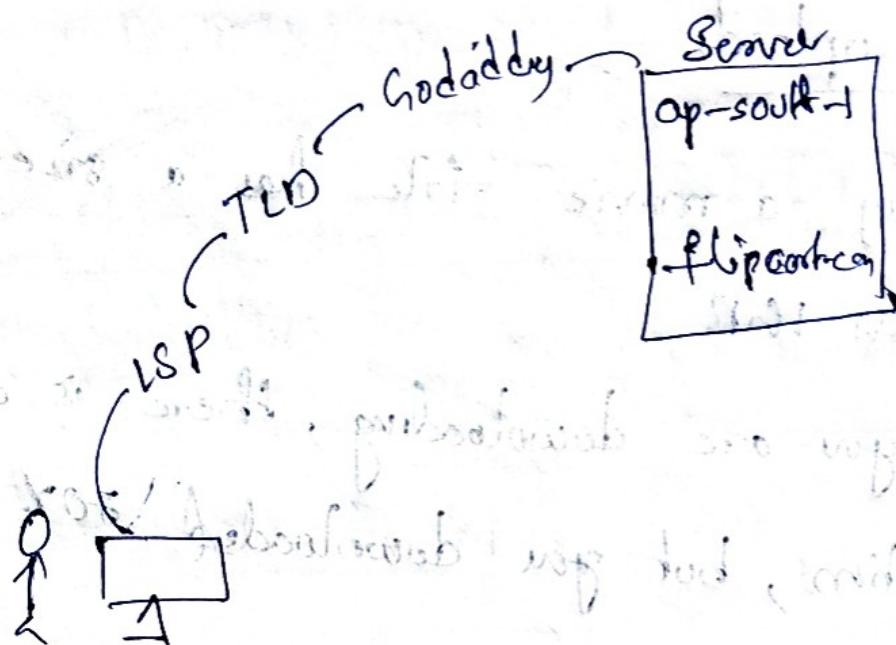
day
1365

5 → no. of copies to retain

Delete expired object delete marker or incomplete multipart uploads

- let's say a movie file has a size of ~~random~~ 1GB
 - while you are downloading, there is an interruption, but you downloaded 70% i.e. 700mb,
 - with that 700mb you can't play the movie you need full file to play that movie.
 - with 70% of file no use, and you have to pay for it, but it is not in use of this process
 - Next time you can resume the same activity
 - But these types of incomplete parts are retained upto 7 days
- Delete incomplete multipart uploads
days
7

Global Infrastructure



- If the same person and server is in the same region, latency will be less.
- If the user is in USA, accessing the server in India, then that request has to travel long distance. Latency will be more.
- To overcome this problem if we create more instances in all the regions, cost will be more.
- To overcome that we have to enable cache servers (frequently accessed ^{static} data is stored in it).

→ For this type of problem, In AWS we have a service called CloudFront

CloudFront:-

It is a service, to create cache servers in AWS.

- Just enable edge locations.
- This can be done through distribution i.e. we are distributing our static data/content over the network.

CloudFront → edge location → distribution

① origin location

② origin location { loadbalanced or S3 bucket }

select it

③ price class

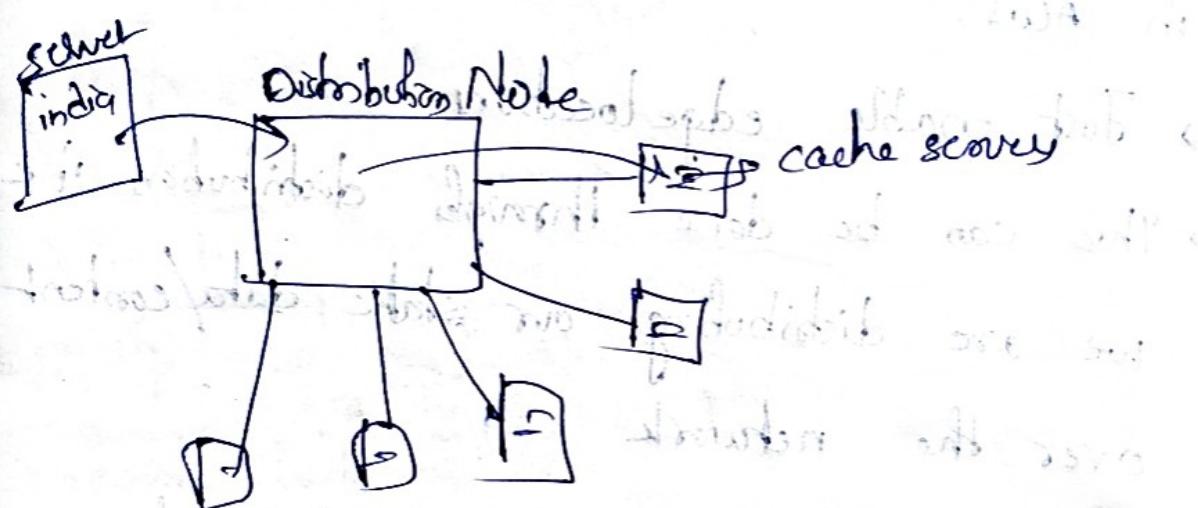
(i) USA, Europe, Canada

(ii) USA, Europe, Asia, Africa, Middle East

(iii) All edge locations

→ More money will be spent, depends on the option selected.

- Here we are saving considerable amount of money for client
- AWS is maintaining no. of edge location in all the regions.
- They are called distribution Nodes



- All the edge locations/cache servers are connected to distribution node,
- The distribution node will forward all the request to nearest cache server of the user.
- If the cache server gone down, the distribution node will forward the same request to the next cache server.
- DNode will pull the information from the edge locations.

server and push that information into the subsequent cache servers, it will maintain same data.

→ To update the server we have a parameter

TTL → Time To Live

→ We are setting the time to live static.

data inside cache servers

$$\boxed{\text{TTL} = 7 \times 86400} \text{ — for 1 week}$$

↓
7 days

→ Because of this, the existing content in the nodes will be expired after 7 days and new content will be pulled again.

→ Every company will be having their own update policy. (for 30 days, 7 days, 3 months)

→ To pull the content from the source to the node, latency will be there,

because it has to go through internet,

to overcome that AWS is using its own network

→ For the distribution it has secured n/w, and latency will be lower.

Ques

* How do you know that secured n/w is reducing the latency

→ least congested algorithm

Confusion

Advantages

i) latency solved

ii) cost reduced

iii) security

iv) least congested algorithm

CloudFront - To load the websites faster

without creating virtual machines across the globe, we use cloudfront, We solved

the problems like

(i) latency

(ii) cost cutting

(iii) security

- Ques:
- (*) what is CloudFront?
 - (*) what is cache servers?
 - (*) what are edge locations?
 - (*) what is Content Delivery Network?

Ans: - Cloudfront

Scenario:

- Create a S3 bucket and host a website in it.
- Then, create a CloudFront and create a distribution and use the S3 as origin in the CloudFront.
- Add new identity in H2O table.
 - (i) No → customer has to write policy
 - (ii) Yes → CloudFront has to write policy
- Now create CloudFront via CDN Content Delivery Network
- We are enabling the cache servers to reduce the latency of my website.

Create OAI (Origin Access Identity)

- ⑧ Yes use OAI
- Create OAI
- Bucket policy
- ⑧ Yes

*settings:-

Price class

- ⑧ use all edge locations (best performance)

Ques:-

- ⑧ My business is in India, Can I select USA as edge location?

- No,
- We have only 3 options, from that

only you have to select

Default root object optional

index.html

→ Cloud front gives SSL certificate, give secure access

Point 1:
* All of a sudden, if one of the caches goes down, who will be taking care, who will be knowing that?

→ DR (Distribution Node)

→ The popular cache service is "AKAMAI".

→ In AWS we have cloudfront

→ YouTube data is handled by "AKAMAI"

→ We can use AKAMAI's services, management

→ -ent has to enable it, get the AKAMAI

people come and install it.

* Instead of S3 bucket as origin, use

load balancer as origin.

Facebook uses load balancer as origin.

Facebook uses load balancer as origin.

Facebook uses load balancer as origin.

IAM

Identity Access Management

- It is nothing but creating credentials to the employee.
- Creating users and giving permissions on AWS account.
- e.g. Create a dev environment, and create group of developers and give them access to in the dev environment.
- We can manage them by giving concerned group permissions to them.

In windows

- We have Active Directory to create and manage them.
- AD → user management

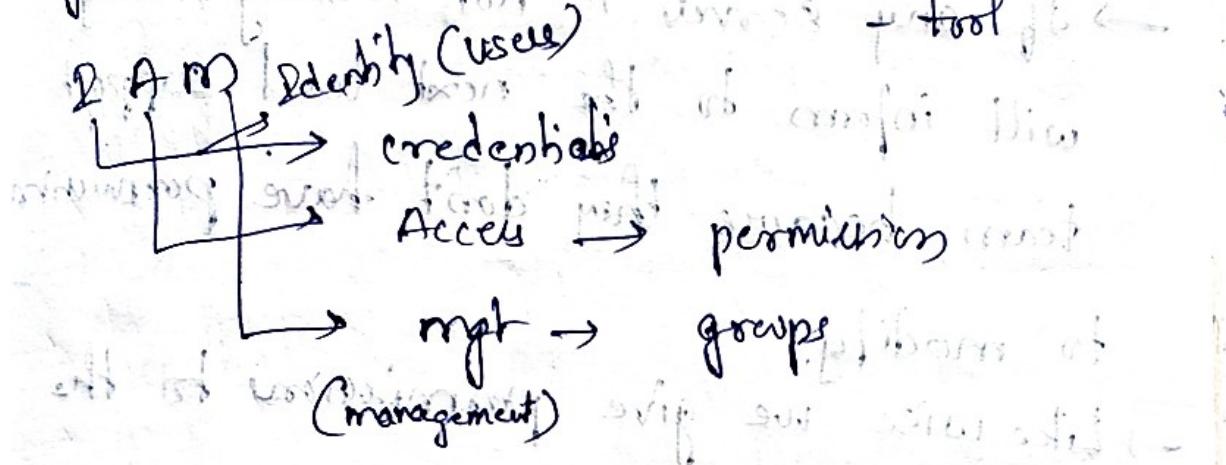
In Linux

- Lightweight Directory Access Protocol (LDAP)

→ Employee ID is called as LDAP only

In AWS:-

We have "IAM", is a service given by AWS, it is User Management



(i) Creating a user in IAM

User Access Types :-

i) AWS mgt console Access (GUI)

Username: user1

prod!

ii) AWS programmatic Access (CLI)

iii) Access :-

→ Provide permissions to the User.

→ Permissions are given on Service level only

eg! - LI support team.

→ They will just check the health of servers.

→ They have only read only permissions.

→ If any server is not running, they will inform to the next level support team, because they don't have permission to modify.

→ Likewise we give permissions to the concerned team.

eg! On EC2 service, we give read only permissions.

★ ROT (Rule of thumb) → Always provide least privileged Access

→ give only the permissions, the user is working on.

→ Good set of understanding about the user's needs, no writing and reading of all the data.

Permission / policy :-

(i) AWS managed policy :-

- It is default policy, created by AWS.
- No modifications
- These are listed in the policies.
- attach to any user
- On every service we have two access they are (i) read only access
(ii) Full access

(ii) Customer managed Policy

- modification is possible
- Listed
- attach to any user

(iii) In Line policy:

- modification is possible
- not listed
- only for a particular user only

Management :-

- If I am having lots user in a organisation, they are working on same operations.
- eg:- LI support teams has 100 users.
 - Create a policy and attach to the individual user is very time taking.
 - Instead of that,
 - (i) create a group
 - (ii) Create a policy and attach to that group
 - (iii) Now add users to the group.
 - Then the permissions of the group will be applied to all the users in the group.
 - One user can be a part of multiple groups.

- Scenario:
- Let's say a website has login page, signup button will be there, We have to fill all our details and upload photo etc.
- The client says that store all the details in a database (RDS) and photos in S3 bucket.
- For this there must be communication between RDS and S3.
- Website from EC2 service is uploading photos in S3 and details in RDS, i.e. service to service communication is happening.
- We can't attach policy to the service
- * policies can be attached to the User/group
- Roles are attached to the service
- Role: - It is nothing but a policy giving this service credentials to other service to login into it.

polices are for Users

Roles are for Services

No Ifs

e.g. -

→ If all persons is wearing a kaki uniform what is the role name she has to control

law and order

→ The uniform defines his role

→ In the same way I am creating a role for ec2 instance, if the ec2

instance has s3 full access and RDS full access

→ If the ec2 instance has the role on it, then he can access s3 and RDS

MFA: - Multi Factor Authentication

→ This is an additional layer of password on top of credentials.

Scanned with ACE Scanner

- Users are global
- Access management Service is force

Creating User

→ click on Add User

g. Username : Sudarshan@ibm.com

depends upon your company naming convention/policy

Select AWS access type :-

① Programmable access

② AWS management Console Access

Console password :-

→ click on custom password

Require password reset :-

User must rotate password

→ Password :-

→ EC2 readonly Access

→ Every service has readonly access & Full access

- Create user
- send the details to the user, login credentials to him via email.
- Create alias to Account ID, because the ID should not be known to user.

Customer managed Policy

→ Give the user just start, stop, restart of instance on top of EC2 randomly and then write the policy in JSON.

- To create a policy
- go to policies in the left hand side of AWS ac page
 - click on policies
 - There will be visual editor and json
 - click on visual editor and follow the steps

Inline Policy :-

- To create it, go inside the user and click on Add inline policy

MFA to User :-

- go inside user
- click on Security credentials
- select Virtual MFA device
- select one of the list of compatible applications
- scan QR code, then assign MFA

Role :-

- giving permissions to services
- eg:- Role for EC2
- select AWS service
- Permission policies
 - (i) S3 full access
 - (ii) RDS full access
- give Role name → create a role

→ Then attach it any particular EC2 instance.

User groups :-

- click on user groups
- Give group name
eg:- L1 support
- create a new policy to that group or add existing policy to that group
- create group

HWS :-

- ① Create 3 buckets $S3_1$ $S3_2$ $S3_3$
- ② Create a user
- ③ Create a policy that, the User has access only on $S3_1$ and $S3_3$, but not on $S3_2$

Scenario! How can we send bulk emails/marketing of IT products?

- We need to have EC2 instance
- OS and its software
- SMTP (Simple mail transfer protocol)
- Cost increases
- Email marketing is used in IT
- For this we have SES in AWS

Simple Email Service

It is a service given by AWS, to send bulk emails for marketing.

Advantages!

- We are not creating EC2 instance
- not installing SMTP
- Not using third party subscription plan
- We are not using third party s/w

- Cost is reduced for all these
- But SES is not free, but it can reduce the considerable cost
- It is disabled by default
- When you enable it, it can give you free tier i.e. only 200 emails/24h

File download/see

whenever you enable, it shows sandbox environment, it means you can do whatever you want in this, If it is getting success and running good, it is like

POC (Proof of Concept)

→ When you give POC to the client, then it satisfies to the client, that will be deployed in ~~testing~~

→ Then that will be deployed in ~~testing~~ prod environment

→ We can increase the limit

→ upgrade to production & Acy

- increase the sending quota i.e. for eg:- 50,000 emails/day
 - It is cheaper than any other third party service
 - If you are upgrading to production account they are giving 62,000 emails/month is free.
- Suppose you are sending 70,000 emails/month
How they will charge you?
- For 62,000 it is free
 - next 12,000 emails, they will charge
 - fee for every 1000 emails \rightarrow 0.10 \$ only

Sandbox:-

- It has some conditions, to use this
- sender email and receive email has to be verified. (must and should)
- It is not possible, so organisations are upgrading to production account

Creation of SES

- Search for SES in AWS account
- Amazon Simple Email Service
- click on ~~Amazon SES~~ → ~~Get Started~~
- Create Identity
- For organization select Domains
- For single email select Email address
- Now we are going with email address
- provide your email id to the person.
- who is sending emails
- click on Create Identity
- Now send email to the employer / customer
- select the sender email and choose a Send email from me
- In that select Formatted text
- In Scenario → select Custom

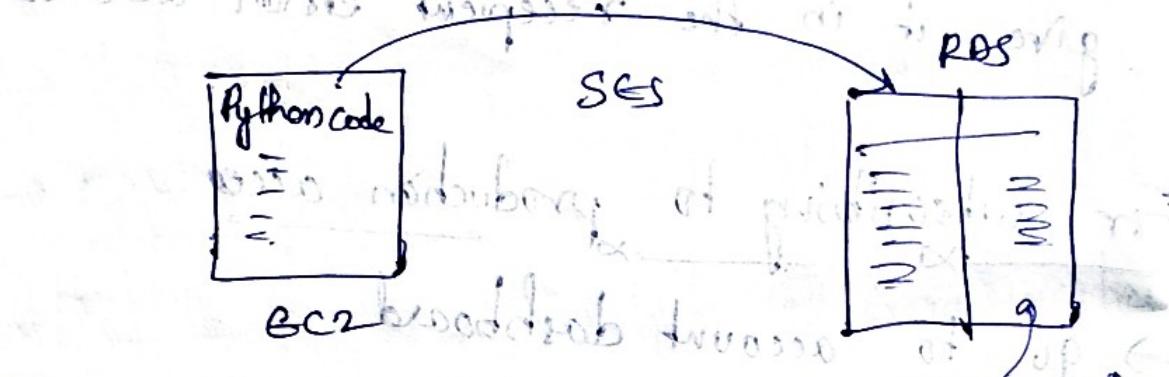
- The provide recipient email address
- send test email.
- To send a email to group of user
create a DL(Distribution list), and
give it in the recipient email address

For subscribing to production access

- go to account dashboard
 - Request production access
 - Then they will talk to you, and they will upgrade your request.
- Scenario:
- For the end users email information is not available, as distribution list.
 - whenever users are registering in the website, they give their details, all the details are registered in a separate database in the form of table.

→ Now you can write a python program to automatically access the database from SES. We have to do the following steps:

→ For this we need a role.



- ① We need Role that access the service
 - ② Programmatic access → we need for user
- Programmatic access gives Access key
- Using access key in my program I can access the emails in the database and send bulk emails.

For above scenario we have following steps:

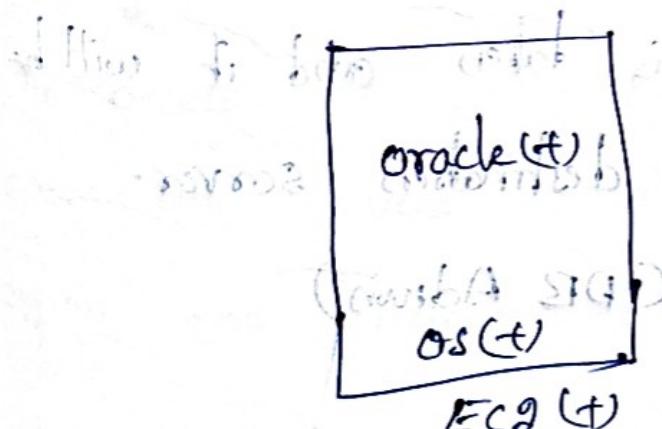
- Create SMTP settings in AWS Lambda function
- Create SMTP credentials with Lambda function
- It is taking IAM role to IAM user
- Create Lambda function
- It will create SMTP credentials

- Copy those credentials and give to developers
- They will use this credentials in their programs, and they will send bulk emails with pictures to billion of people.

DataBase

- It is a s/w, with this raw data is stored in structured (tables, rows, columns) format by indexing.
- We can retrieve the data with the help of indexing very quickly.

challenges of maintaining a DB in our own Virtual machine



- (i) oracle DB admin
- (ii) os admin
 - Cost increase
 - to maintain a DB

The cost are not free, many type

~~Chit~~: ~~Establishing a new project~~
① If the DB of ECG machine is going down we can't access the data

SPOF → Single point of failure

→ To avoid this, we are maintaining another ECG it is called as standby database server

② Backup: -

→ Every organisation maintains a backup policy

with the following frequencies:
hourly
daily
monthly
yearly

③ Restore DB

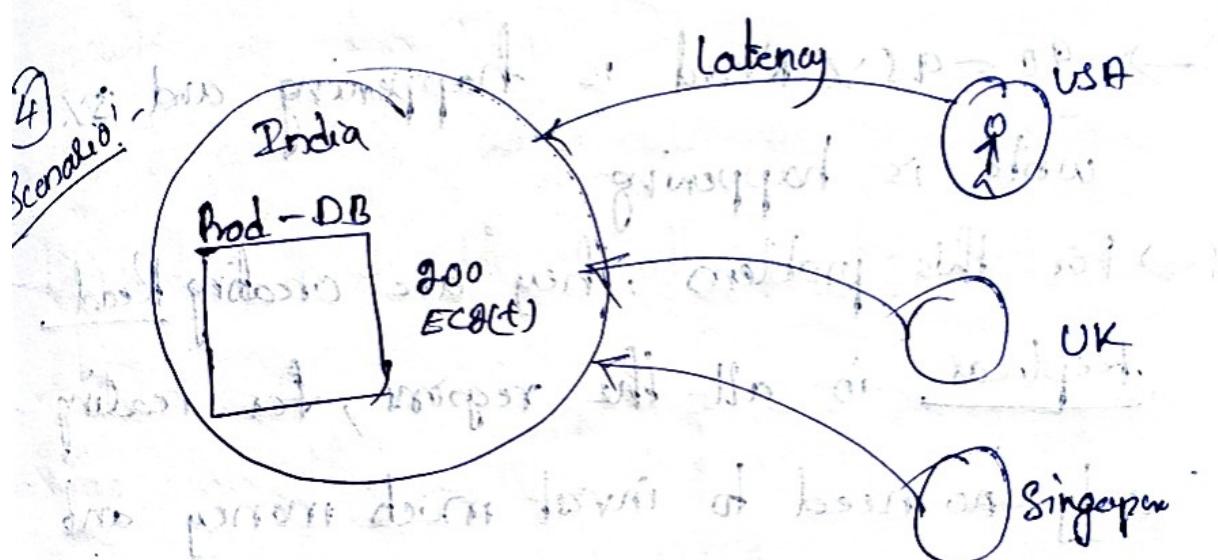
→ The backup is taken and it will be stored in a destination server.

DBA (DB Admin)

a destination of sh

ops team → start doing backup

- These backups will take 5-6-8 hrs, while doing it fails at 99%, again we have to start doing.
- We have to start fresh backup, for this we have to clean backup servers.



- To access the financial DB which is in India from other regions, latency will increase. So performance will be slow, transactions will go down.
- To decrease the latency we have to create same DB servers in all the regions and again cost increase.

~~eg:- ATM~~

- Insert card → It will check with DB and give

the response of it matches \rightarrow Reading DB

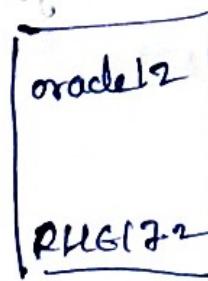
- \rightarrow PIN \rightarrow Read (first) \rightarrow quick statement
- \rightarrow statement \rightarrow Read
- \rightarrow min statement \rightarrow Read
- \rightarrow Balance enquiry \rightarrow Read
- \rightarrow withdraw \rightarrow write

\rightarrow 90% - 95% Read is happening and 5% write is happening

- \rightarrow For this problem, they are creating Read Replicas in all the regions, for reading only no need to invest much money and they will maintain 1 or 2 Primary DB's covering the each region, by this cost is reduced and there will be no latency.

\rightarrow latency $\downarrow \rightarrow$ performance \uparrow

⑥ O/S administration



upgrade \rightarrow o/s patching
 \rightarrow security patching

- ⑥ oracle 12 → upgrade.
 - DB patching
 - DB security patching

⑦ All of a sudden oracle server is not running again it is a problem, Compatibility issues

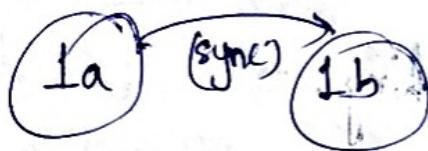
→ As a customer, we are maintaining the DB there will be challenges he has to face

→ For all these problems to maintain we have a service, i.e. RDS

Relational DataBase Service

- It takes care of all the challenges
- select the readily available databases in the RDS,
- It will automatically takes care of compatibility issues

→ check multi Available Zone Deployment



when you click this a DB is created in (1a) zone and another standby DB is created in (1b) zone and data is always synchronous with (1a)^{zone} DB.

→ Create Read Replica → select destination (UK)

→ To take DB backup

automated backup → open windows and specify time for backup.

→ Restore is happens automatically.

→ It is a managed DB service.

* → RDS is a service given by AWS to create database to store raw data in structured format.

Creating RDS

- click on create database
- click on standard create
- because it has more options than easy create
- Engine option
 - select any db engine
 - eg:- oracle , mysql

Templates :-

- ① Production { for production purpose }
- ② Dev/Test { for develop/testing purpose }
- ③ Free tier { for dev free applications and testing purpose }

Settings

DB instance identifier

give name for your db

→ Master username

dbadmin / any name

DB instance class

① Burstable classes (include t class)

include previous generation clause

→ It give t class for lower cost and some are free also, that's why we are choosing this one

select t.micro

storage :-

give some storage based on your requirement

* Storage Autoscaling

while we are creating we created with TCRB, but some transactions are happening we require more gbs than what is there with us.

→ Thus it automatically increases the storage space,

→ It is not a free service.

Availability & durability

→ To create a standby db instance it's created in another zone.

→ It is not free service.

Public access

→ Don't enable it.

mysql } port no is 3306.
mariadb

Amazon Aurora

oracle → 1521

→ Create a new security group and attach to it.

Database authentication

① Password and IAM db authentication

→ for this create a user, and give full

access of RDS.

② Password authentication

→ master password will be for this

Backup retention period → 1 day

→ for every day 1 backup will happen automatically, the older backup will be removed and the latest backup will be kept or retained.

Backup window

at what time you want to take backups
set the time here, depends on your organization.

To give the time in UTC, there will be online converters from ITC to UTC

To give the time in UTC, there will be online converters from ITC to UTC

then specify the time

Maintainance window → The maintenance windows and backup windows

should not overlap

Deletion protection

enable it, always

- To access database from your laptop you have to install mysql workbench
- for different different db's there are different client s/w's
 - for oracle → sql developer, toad
- copy the url of dB and paste in the mysql workbench at the place of host ip, give port no, and password
click on test connect

Scenario:- ~~single environment~~ ~~multiple environments~~ ~~multiple environments~~

→ The initial requirements for the IT business
are extracted from the client requirement

~~client~~

Client → 4 → Dev, QA, Stage and Prod

→ We have 4 environments

→ 30 servers for Dev in 3 different zones
in Mumbai

10 × mumbai × 1a × ubuntu × key × sg
10 × " " 1b × " " × " " × sg
10 × mumbai × 1c × ubuntu × key × sg

(20)
Dev

20 × mumbai × 1a × ubuntu × key × sg

20 × " " 1b × " " × " "

20 × " " 1c × " " × " "

Prod Prod
(150)

50 × USA × 1a × ubuntu × key × sg

50 × " " 1b × " " × " "

50 × " " 1c × " " × " "

Prod
(60)

100 × USA × 1a × ubuntu × key × sg

100 × canada × 1b × " " × " "

100 × africa × 1b × " " × " "

→ Total we require 540 servers

To create a instance, we have to go through
7 steps:

Ans → Bot.type → config → Gbs → tag → sg → key pair

→ To create 540 instance, we have to go
for 12 times all these steps

→ After doing all these things, you have
done a mistake i.e Dev servers in wrong
region, then client won't agree, and client
will say please remove them,

→ The moment you launch Dev servers,
(30 servers) will be in the ownership of
Dev manager, You have to get approval from
him, billing will be there, getting the
approval will take 1 week also,

→ Then manager have to write what is
the impact of error, conduct meeting,
get approval,

→ Then you can remove wrong instances
after the approval only.

→ These all are human errors, manually.

→ To automate this activity AWS gives us awscli

log of over 100 commands stored in awscli

→ The awscli is a programmatic access

→ It uses scripting

→ It has a user interface with editor

→ For user → We go with IAM

→ with IAM → create User → with programme
access → attach a policy (for EC2 full access)

→ In every organization

→ The moment you click on create user
it gives two access key

AK: ***

SK: xyzt1234567890

→ Now open a cliconfig source/boston file
it will be present in every organization

→ After logging into that server, first check whether there is awscli

aws --version

yum install awscli -y

→ After installation, then configure cli user

aws configure

→ AK: paste your access keys here

→ SK:

→ Default Region: ap-south-1

In which region you have to

create servers, specify here

→ output format: text/json/table

→ It asks all the four question, answer it

→ Now to create 30. instances in mumbai

aws ec2 run-instances ami-id ami-0b823425

instance-type t2.micro vpc-id → a-zone ap-south-1a

count 10 --region ap-south-1a bag →

Sq-id ↪ key-name ↪

→ In this way, we can create by a single command.

* Should you give all the parameters in an order step by step?

→ No

Advantages

→ Execution time is less

→ Reusability

Drawbacks!

→ If you are doing errors, to modify it automatically roll-back is not possible

→ We have to delete it manually after getting approved

→ To overcome all these errors, the AWS gives us a service called

Cloud Formation

Cloudformation

It will help us to provision the Infrastructure.

- We need to write 12 steps to create 540 instances.
- It is letting us to write only 1 script to create all 540 instances only one script.

Scenario:
→ suppose, the 1000 servers are running in singapore are prone to disaster, then you have to replace them very soon.

→ In this scenario's we go with these scripts and setup all the 1000 servers just in matter of minutes.

→ To create a same environment, replica environment very quickly, you can create script in cloudformation

- IAC, Infrastructure as a Code
- IAC is a service given by AWS

→ Script languages are allowed in cloudfront; they're of two types

- i) json
- ii) yaml

→ The cloudfront calling the script as template (ie predefined blue print)

Syntax

filename.json

cloudformation version → It is a mandatory line

Resources → mandatory section

Properties → properties

aws::ec2::instance → ami-id

ami-id → instance type

instance type → stars

→ stars → no relationship → same →
stars → private → private is in same

execution methods - Inside down right

① upload the code into S3 bucket

→ then open cloudformation service

→ then open first execution method i.e S3

→ copy object url, copy into cloudform

ation and then click Next, then

code will be executed

② As a local file upload directly into

cloudformation.

* what is the significance of S3 bucket

in case of cloudformation

→ S3 will be working as compiler

→ There will be inbuilt compiler in S3

→ S3 will execute the code,

→ If the code executed upto Dev, QA

and pre in pre-prod the code has

got error and stopped executing further.

→ Then what about the created servers
, 90 servers total after all badge

Dev → 30, QA → 60, → total 90 created

out of 90, 540 failed ways with

These are two options for this problem

(1) Rollback on failure (2) create

→ If a roll's back everything that is
created, removes everything created

Advantages

→ No need of taking approvals from
manager

→ No need of listening the songs of
client
etc

Scenarios

I don't want to delete the 90 servers
which are successfully created

~~(1)~~ preserve successfully created resources &

→ We can create any service, i.e. EC2, S3, RDS etc. for this, we require IAM role is required, and attach it to cloudformation

Timeout value:
→ If there is 200 lines of code, who execute the code if it can't set time, under this time only it has to execute

awscli creation
→ We need a server
→ Now create a IAM user
→ for this we go with programme

→ copy AKS, STK
→ Now enter into a server, then
→ check whether awscli package is present or not

aws --version

(or) yum install awscli

→ Now give command as

aws configure

AKS : paste here

SDK : "

• Default region : give region name

Default o/p format : text

→ lets check, can we access s3 bucket?

aws s3 ls

It shows permission denied because we have only ec2 permissions for this user.

→ We can access EC2 only

→ To see the list of instances

aws ec2 describe-instances --region region-name

→ To see the o/p of above command in different format

aws ec2 describe-instances --region region-name --output [table/json]

→ Create a ec2 instance

```
aws ec2 run-instances --image-id ( )
--count 1 --instance-type t2.micro --key-name
--security-group-ids < >
--subnet-id < > --region < >
```

json script for ec2 formation.

filename.json

```
{
  "AWS::TemplateFormatVersion": "2010-09-09",
  "Description": "Demo ec2 instance creation",
  {
    "Resources": {
      "DEVELOPMENT": {
        "Type": "AWS::Ec2::Instance",
        "Properties": {
          "AvailabilityZone": "",
          "ImageId": "ami-00000000",
          "InstanceType": "t2.micro",
          "KeyName": "key"
        }
      }
    }
  }
}
```

"SecurityGroups": [" "]

() [] } - ~~CloudFormation does not support nested stacks~~
more post - ~~CloudFormation supports nested stacks~~ I found -
< } > ab3-geop-structure

→ Now go to ~~Cloud~~ cloud formation

→ create cloudformation stack

→ Template is ready

→ You have written the template

① Use a sample template.

→ AWS provider &

→ Select any one

② Amazon S3 url ③ upload a template file

→ Next

Permissions!

TiamRole Name

→ Without creating a role will it execute or

→ "Yes", every cloudformation template will be

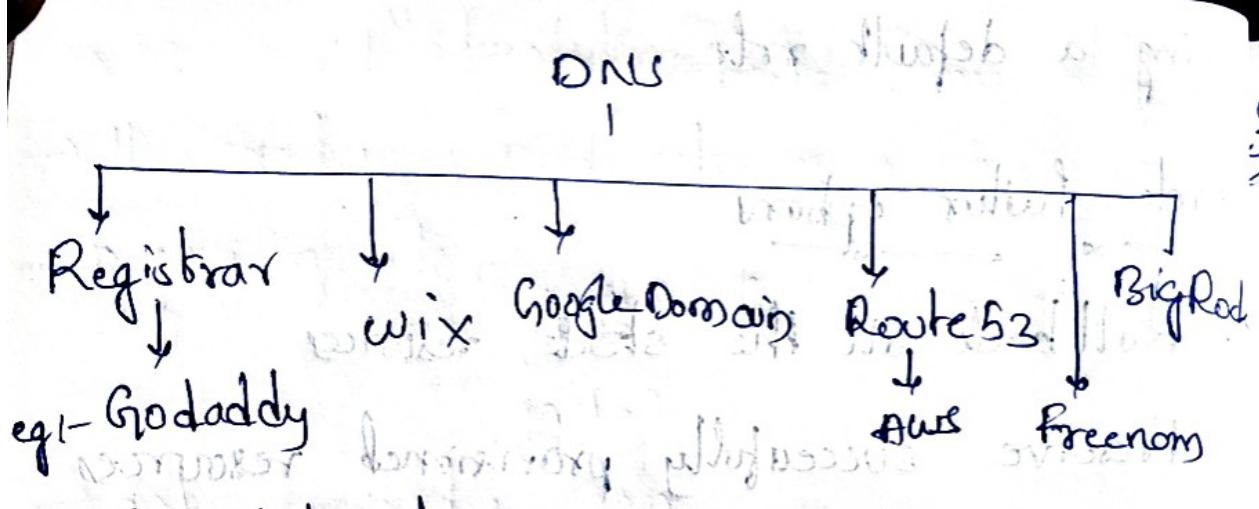
using a default role

stack failure options

- o Roll back all the stack resources
- o Preserve successfully provisioned resources
- Notification of which is available after
- Timeout policy
- Protection termination
- Next

DNS

- It gives a unique name for your ip address
- We can't remember ip's, and load balancer's url, at which we do some changes
- It is Domain Naming System
- We need to provide a meaningful name unique name that can be used across the globe is called DNS
- To maintain all the names, there must be a central body, it is DNS



- The distributor or dealer of domain names inside DNS is called Registrar
- There are many registrars, Very famous one is egi-Godaddy
- In Registrar, You will register your domain name or you will buy your domain name
- We will buy a domain name/ registering a domain name, we are going to attach to our load balanced

eg! - whenever a user types the domain name and hit enter, then

Domain Name → Registrar → Actual Server
 will forward

Scenario-1
e.g.: Infosys has its domain name www.infosys.com.

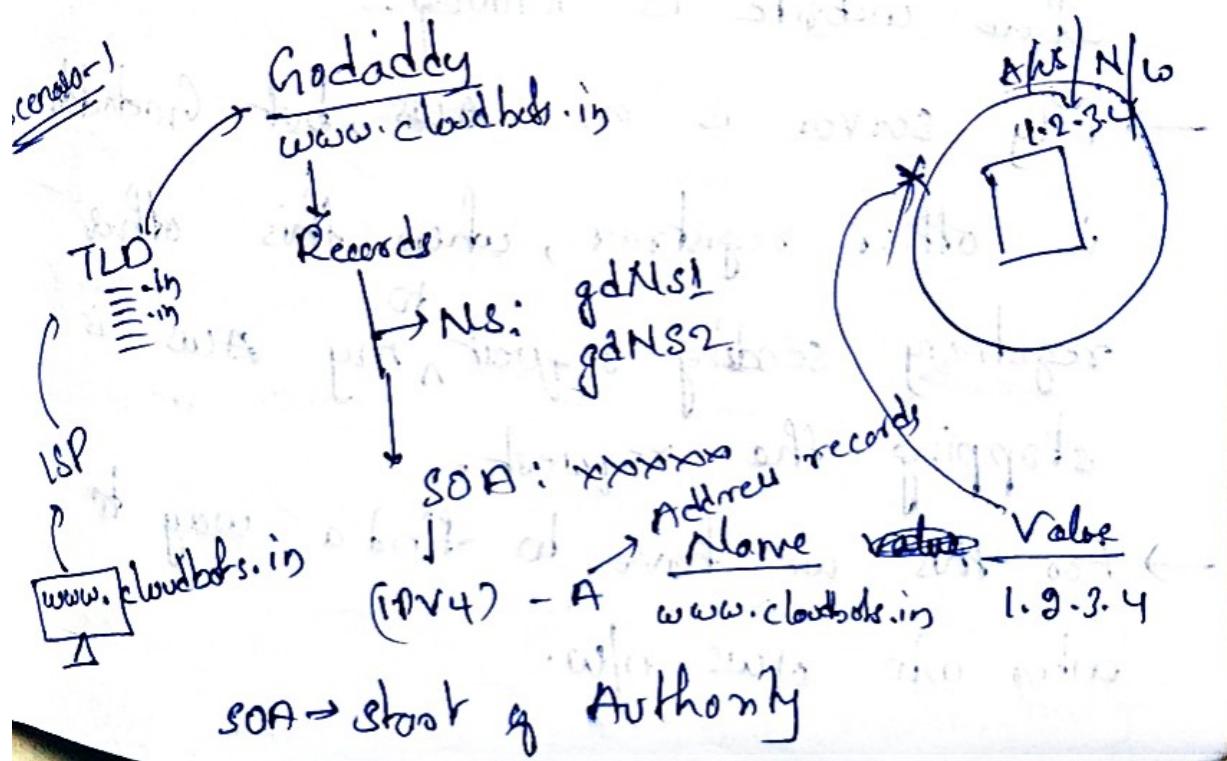
If they want to migrate their business to AWS
How to manage their domain name?

Scenario-2

→ If there is a startup company, without buying a domain name outside, they can buy it from AWS itself.

→ By or from Route 53

→ Every domain name has expiry, & Every year you have to renew, if you are not renewing your domain name, they will release it to others.



- We registered our website name in Godaddy
- It creates record sets of number (two) naming records
eg: gdNS1 (Name server record)
gdNS2
- and it creates SOA: xxxxxxxxx
- whenever a user types www.cloudbees.in
this request reaches to ISP
- Then ISP looks into a TLD in which all the domains ends with .in,
- Then TLD sends the request to Godaddy, then if the name is registered even in it, it will check and send to server where website is running.
- My server is on AWS, but Godaddy is other registrar, when this other registry sending request, my AWS is stopping the request.
- For this we have to find a way to enter into AWS n/w.

Route 53

It is a service given by AWS to configure or to register a domain name for a existing domain name.

→ 53 is the port no. of DNS in Linux

→ It is a basic double check mark OS

→ In this I am going to create a hosted zone with the same exact domain name in Route 53.

→ This Route 53 thinks that it is fresh name and gives four NS records

→ If 3 NS records are going down, we can find answer with 4 NS record, it gives high availability

→ It is faster than others because of 4 NS records

It has SOA: xxxx.x.com also

→ Now still request is going to Godaddy only, how can make that request directly

comes to Route 53

→ For that we replace GoDaddy NS with AWS NS

aws NS1
aws NS2
aws NS3
aws NS4

→ Now create address records in Route 53

create Name value
www.clashbot.in 1.2.3.4

on SOA → where you are buying the domain name, your request should start from them.

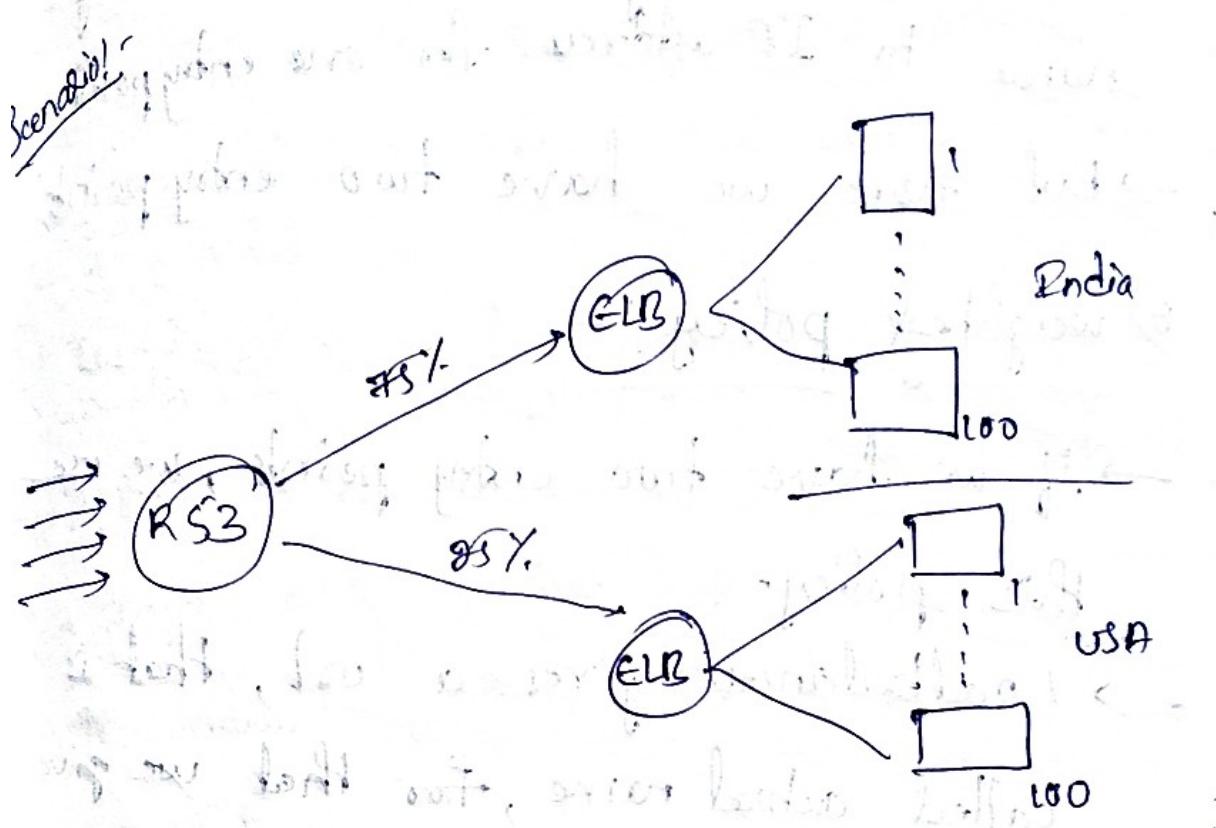
Routing policy!

→ You are having single service or many services for one website, then how can you make entry point for many services at one entry point

① Simple routing!

→ If you have one entry point, you can

use this policy
eg:- single server, loadbalance, s3 bucket



- For my business, 100 servers are running in India with 1 LB and 100 servers in USA with 1 LB.
- Loadbalance has a limit, it can work in one region only, i.e. we can't include USA region servers in India GLB.
- Then RS3 has to forward request to one of the loadbalance;

Here we have two entry points, for this we can't go with simple routing policy.

→ For simple routing, we are resolving name to IP address for one entry point

→ But here we have two entry points

2) weighted policy:-

→ If we have two entry points, we use this policy.

→ Loadbalance gives a val, that is called actual name, for that we give another name i.e alias name.

→ Now I create one alias record for loadbalance in India and one for another alias record for USA. Load balance

→ let's say I have more user in India

but less in USA,

→ For this I create 75% alias record for

India and 85% alias record for USA.
Loadbalance
→ q! - If RS3 is receiving 100 requests
route 53 sends 75 requests to India and
25 requests to USA.

③ Geo location routing policy:-

→ Based on the location, the request should
reach LB of that region only.

q! - Indian request should reach India LB
USA " USA "

→ For every country there is a separate

IP range

→ This information is having with Route 53,

then the moment you type the

url, it reaches to RS3, then it

understands from which region and

sends that request to that region

LB only

④ Latency based policy:-

- If my user is in India and typing my website
- To reach India LB, the latency is 7ms and to reach USA it has latency of 5ms then RS3 sends request to USA LB only even though he is in India, because he choosed this policy.
- The RS3 calculates the Latency based on heartbeat

⑤ Failover routing policy:-

- If India LB is more efficient, then if choose it has Primary LB, then all the requests go to Primary only
- In case Primary LB fails then only the requests are send to secondary LB in USA

- If your organisation is maintaining IPv4, we have to create ~~A~~ Address record of 'A' → IPv4 for IPv6 → "AAAA" record
- Name → IP resolve
IP → Name resolve
PTR → pointer Record
- For an organisation, while buying a domain name, you have to mail service records like eg:- sudarshan@bcs.com
- If you already have mail service we can configure in AWS with 'mx' mail service exchange record values before.
- for mx also → Name Value will be there

Note:- whenever you are creating Record sets and Domain name it will take (4-4t) hrs to reflect into all TLD's and to access it from browser

- Attaching a domain name to IP
- Create instance and it has one IP
 - Launch a website in it
 - Now attach that IP to domain name
 - Now let's say we bought domain name now I have attach that domain name from GoDaddy to R53
 - my domain name is www.cloudbots.in
 - For existing domain name in a registrar, to configure that name into AWS we create hosted zone
 - search for Route 53
 - select create hosted zone
 - Domain Name : cloudbots.in
comment : my org
create
 - Now it will give you two record sets
 - A → ns-1449.awsdns-53.org.
 - NS → ns-1005.awsdns-61.net.
 - NS → ns-206.awsdns-95.com.
 - NS → ns-1592.awsdns-02.co.uk.

SOP: ns-1449.awsdns-53.org., awsdns-hoster
er. amazon

→ Now copy NS names and paste it

~~Godaddy~~ Godaddy NS,

Note:- copy till the name only leave the dot, copy it carefully.

Scenario 2 → don't copy the trailing dot

→ Now create Record Sets in R53

→ click on create Record Set

Name : www.cloudbots.in

Type : A - IP or address

Value : paste IP of your server (pub-ip)

Create.

Note:- copy LB vol, S3 vol in the type

section if and only if you have only one LB and one S3 only.

Scenario 3

→ If my customer is not typing www

only typing cloudbots.in, then we will lose the business

→ For this create an alias Record Set

Name: cloudbots.in

Type: A - IPv4 - address

Alias Yes No

Alias Target: www.cloudbots.in

select it from drop down

Scenario-3

How to create a subdomain name

Name: dev.cloudbots.in

Type: CNAME - Canonical name

value: paste DNS name given by AWS for service

eg: ec2-3-110-169-85.ap-south-1.compute.amazonaws.com

→ Create

Scenario-4

Two entrypoints, one is India and

other in USA

related to me

>Create Record Set for Loadbalancee of India

Name: www.cloudbox.in

Type: A - IPv-4 address

Value: paste ip of LB of India

Routing policy: Weighted

weight: 75

Set ID: IndiaLB

→ Create

Create Record set for Loadbalancee of USA

Name: www.cloudbox.in

Type: A - IPv-4 address

Value: paste ip of LB of USA

Routing policy: Weighted

weight: 95

Set ID: USALB

→ Create (0.0.0.0)

Note:- Create Records based on Geolocation,

latency, failover for the same

website.

(A) LAMBDA functions

Scenario:

- In our organisation, we have 2000 servers and 8000 have moved to cloud, they are telling that we have to check the servers whether it is open to internet ($24 \times 7 \times 365$ day)
- If any server is intermediate to company IP or company IP range?
- We can't check whether the server is open to internet through security groups, inbound rules
- Open the security groups and check in inbound rules, if the IP address is $(0.0.0.0/0)$ → it means open to internet.
If not it must have company IP address eg $10.1.2.4/8$.
- We can't check 8000 security groups manually,

→ let's say we checked 1000 security groups manually, and while we are heading to next 1000 SG's meanwhile anyone may opened the service to internet in above 1000 servers.

→ How to restrict someone opening the servers to internet? manually ~~is it~~ is not possible to do/ check SG always

→ For this we do automation, to do so we need to write program, for this we ~~should~~ have to buy

(i) VPC

(ii) OS, licensing

(iii) Java/python etc

(iv) We need to have programme

with access of AWS account

(v) to automate aws account

using python we need to

import a library called 'boto3'

→ To write a code in a service, it takes
95% of memory, time, Hdd, Cpu,
→ 95-98% of the time these all resources
are sitting idle, for this we have to
pay to AWS

→ To overcome all these problems, they
invented a service called Lambda,
which is serverless programming platform

It is a platform as a service (PaaS)

→ We will assign a user to Lambda
function, he opens it and selects a
program for customizations; they are giving
a readily available IDE

→ Lambda function works to access S3,
DynamoDB, EC2 etc, create role and
S3, attaches to Lambda function

→ whenever it is executing the code only
utilizes the resources we ~~can~~ can cut

cost.
→ I can control executing with timeout.
★ because when the Lambda is executing it
is charged. (charged at runtime).

④ How/when to execute the Lambda function.

- ① scheduled based execution
 - it is like cron jobs, i.e. every
~~hourly~~ hourly scheduling, it
 - At 8:00 am, if am scheduling, it
will run next time at 9:00 am, next
10:00 am, ..., In between can anyone
open my score to internet, then
how can we trace it, there is a
chance of opening the score.
→ Because you should it for hourly,
next 9:00 am only it will execute
and trace if any score is open
and ~~none~~ remediate to secure n/w.

Case-II

→ It will have scheduled it for *.*.*.*.*
ie for every minute, every hour. It is
going to charge because we are
executing for every execution it charges

⑨ Event based execution

e.g:- creating a service it is a event,
starting a service, doing anything in
aws is called event.

→ By default events ^{log files} are not captured
in aws, for that they have a service
which is called Cloud Trail.

Cloud Trail :-

- It is not a free service.
- The moment you enable it, it
creates a S3 bucket and storing
all the 23 regions log files into it
- If anybody is changing the security

group, what is the event for it?

→ Event name: ~~Security group~~

Security Group Ingress

→ I have to use Lambda function on
Security Group Ingress

⑧ Now how can Lambda service is going
to search in the S3 bucket for SecurityGroup
Ingress

→ For this, after writing the logs (in S3,
Cloud Trail is going to stream (live insert
logs) into a service called Cloudwatch Log.

→ Now after writing logs into Cloudwatch Log,
then I will trigger my Lambda function and

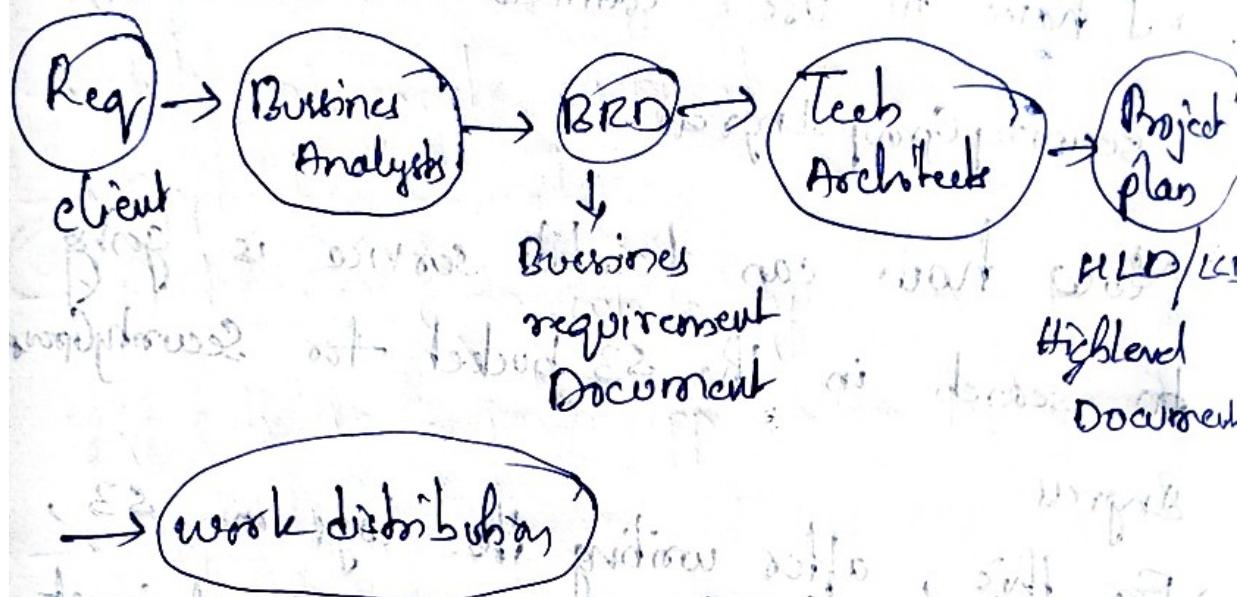
from dropdown Security Group Ingress

and save

→ whenever any event is done, my Lambda
function executes, whenever the Lambda
function executes the error is identified and
solved immediately

Elastic Beanstalk

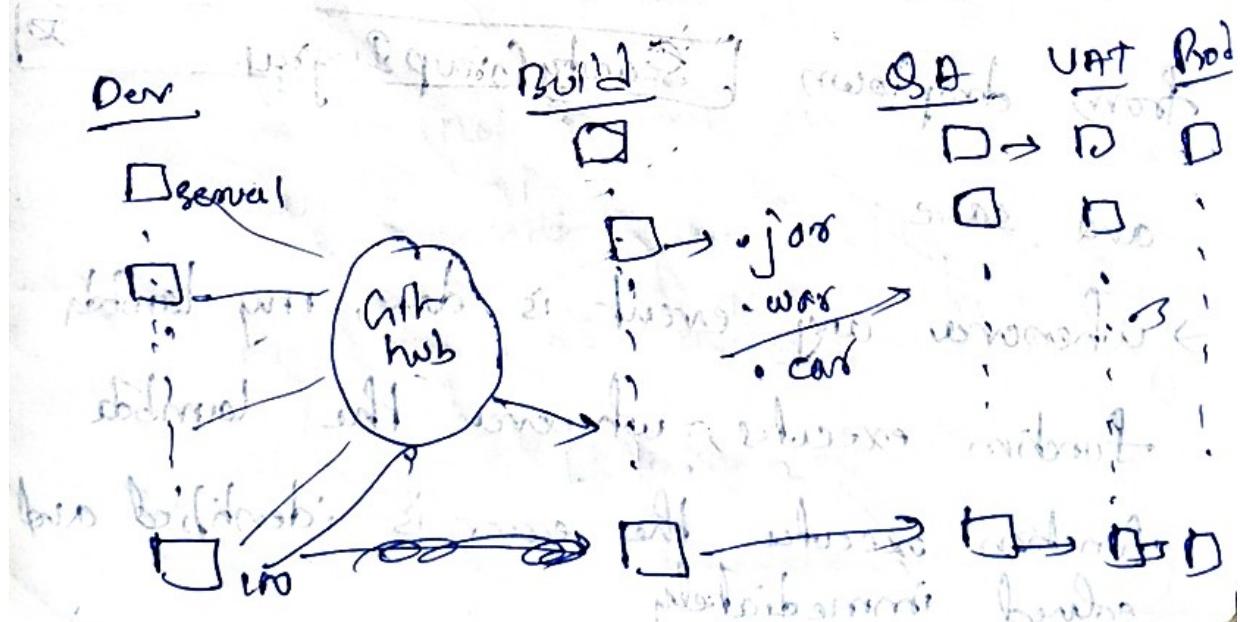
→ For a business input IT, we have a following cycle.



→ work distribution

→ Then cloud engineer has to create Dev, build, QA/deployment, UAT/Prod environments

UAT → User Acceptance Test environment



Roles of cloud engineer

① Create environment

② Deployment

③ Management

→ To solve above challenges, there is a service to quickly setup entire environment is Elastic Beanstalk

★ Note:- This is only for web application.

→ Upload your build file in that service it will deploy in all environments.

→ You can deploy, undeploy, redeploy very quickly.

~~It will not charge for this service but it charges for the resources utilized in the service.~~

→ You simply upload the codes and beanstalk automatically handles the deployment, from

capacity provisioning, load balancing, and automate scaling to web application, health monitoring, with ongoing fully managed patch and security updates

→ click on create application

→ Application name

→ Platform (on which you have to deploy)
e.g. if it is java app → then deploy on tomcat

→ Application code:

① sample application

② upload your own code

→ Source code origin

Version label any name

① localfile

② Public S3 vol

→ If your file is in your machine then

select localfile

→ If it is in S3 bucket, select S3 vol

- Create application
- For redeploying another version of war file just select **upload and deploy**

Then it will remove older version and deploy later version

Configure more options

- You can select **instance types**
- No. of instances
- Deployment policy
 - o Rolling update
 - o All at once {it is dangerous, don't go with it}

Security

→ Add the keypair, it by default won't create any keypair

- Create RDS data base and attach to it
- The cloudformation is advanced compared to elasticbeanstalk

~~Today~~ With the help of cloudformation we can
create no. of resources, no. of environments
and deploy any kind of application

~~But~~ But with elastic beanstalk only 1 environ
ment only for web application.

VPC

Virtual Private Cloud

- It is nothing but a n/w.
- AWS gives default n/w to any account.
- We can create our own n/w.
- The private IP range is ~~given~~ created by the n/w.
- We can create our own IP ranges i.e. private IP through VPC with AWS.
- We are isolating the n/w from AWS given n/w i.e. we are isolating our resources from AWS n/w and keeping them separately in our own n/w.

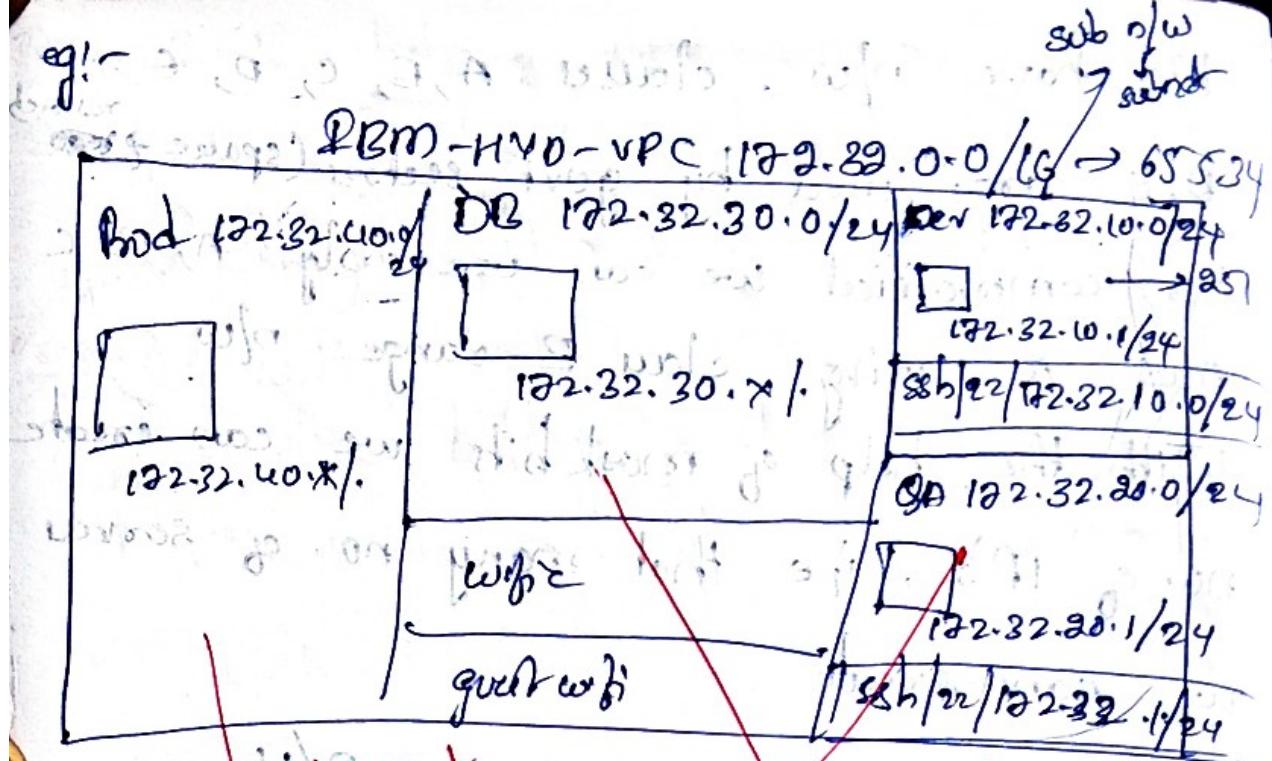
- We have new classes A, B, C, D, E
- D, E are used by govt sectors (space & ^{research})
- For commercial we can use only A, B, C
- AWS is using class B range 172.0.0.0/16
- With the help of mask bits we can create no. of IP's, i.e. that many no. of servers we can create

eg:- AWS IP range 172.31.0.0/16
 ↓
 address bits

- It creates 65534 ip addresses
- depending upon the requirement in your organisation create IP addresses

→ 172.31.0.0/16 → 65534
 ↓
 172.31.0.1 172.31.1.1 172.31.2.1 172.31.3.1
 ↓ ↓ ↓ ↓
 0.255 0.255 0.255 0.255
 ↓ ↓ ↓ ↓
 0.255 0.255 0.255 0.255

eg:- We are creating a VPC for IBM
IBM-HYD-VPC
 ↓
 172.0.0.0/16



→ We have created our own n/w through VPC

→ I am creating a subnetwork in my main n/w, i.e. called as subnet because I don't want all the dev servers accessed by testing team.

→ Like this I can create no. of n/w's in a main n/w and restrict the access by others and outside of this n/w too

→ eg:- 172.32.10.0/24 → 254 IP addresses from 3 IP address are used by internet purpose and only 251 can be utilized

Case 1 To communicate with the Dev servers.

We need the IP range of Dev servers only.

→ To restrict the access from one subnet to another subnet, I am creating a security group.

group

i.e Sg [ssh/22/172.32.10.0/24]

The Dev team members coming through the above Sg only can access the Dev servers.

Case 2

→ In QA environment, I create a subnet of IP range 172.32.20.0/24.

→ All the team members can communicate in QA servers those who have Sg of

[ssh/22/172.32.20.0/24]

Case 3

→ The IP range of 172.32.20.0 can't access

the IP range of 172.32.10.0

→ because the SG is restricting it, it has
no gate to each subnetwork.

Note! - The starting point of your project is
networking, i.e. create n/w and separate
subnetwork for each environment.

Case-4: Prod/live environment.

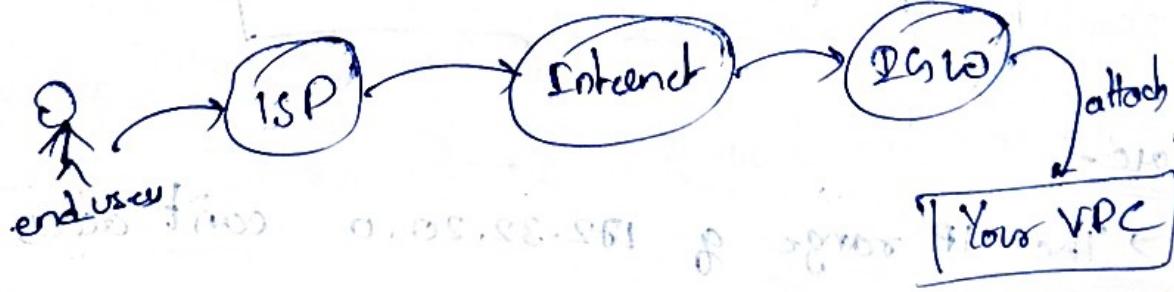
→ These servers are accessed by the end users

→ But my prod server is in my own n/w.

n/w:

→ The end users will be outside the n/w
to access it we need public IP

→ We have to expose prod servers to public
for this we have to enable public IP
for prod servers



Internet Gate Way ! (IGW)

- It acts a gatekeeper to every v/v.
- It has permission that whom to allow / whom to restrict to enter into your VPC.
- Create one IGW and attach it your VPC
- It allows you to connect to your VPC only.

Note! We never expose, Dev, testing, and DB servers to public, for that we never enable public ip for these servers.

- We are not enabling these servers to public ip address, it means there are accessed by our employee only, then we can say it as private subnetwork.

and Prod. environment we can call it as Public subnetwork.

- By default all the sub nets are private, when we are exposing to public, it will become public subnets.

Scenario 1

- * From outside the VPC, we have to login into Dev service and install httpd.
- Now an employee who is doing work from home, login to Dev service and install httpd.
- For this he will enter into Dev service through RDS and loged in.
- Then he is installing httpd for that he again has to connect to internet and download the httpd package and install it.
- For VPC there is only one way to enter into it, and there is no out way for going out and connect to internet. There is a path to enter into it but if there is no path to go out then he can't connect to outside RDS and download httpd.
- For VPC there is only inbound and no outlet.

From inside the n/w, to connect outside n/w

We have routing tables, store best level routes to

Routing table! - (Path)

→ It creates a path from private VPC to internet, it is like a way/outgate to VPC

→ To connect a part of subnet to outside there must be outgate/path

e.g. - Route table

$0.0.0.0/0$ → it means you can

connect anywhere and anyone in the inter-

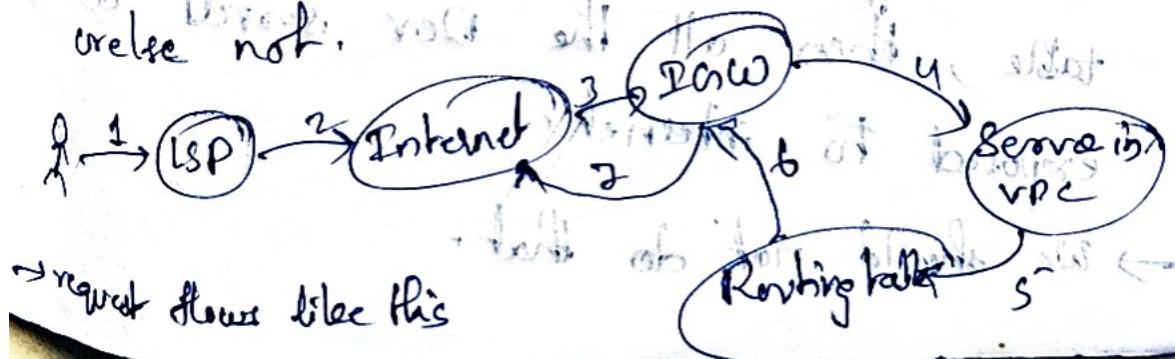
-net, through gw only.

→ All the servers in the subnet, contacting

routing table, if the server has access in

that table they will be connected outside

else not, will go through gw



→ request flows like this

→ Firewall is a security group at instance level and route table is security group at subnet level

Note: → DCSW, allows outsider to connect to the resources in your n/w, and Routable allows inside servers to connect to outside of bridge to a team of a.

Scenario: To connect prod member into Dev team

For that create a Security group in the Dev subnet as [3sh|22|122.32.40.1/24]

Then only prod team servers can communicate with dev team servers

→ From Dev to connect internet create a routing table in the Dev subnet,

→ If we give 1GW per IP in routing

table, then all the Dev servers are exposed to internet

→ We should not do that.