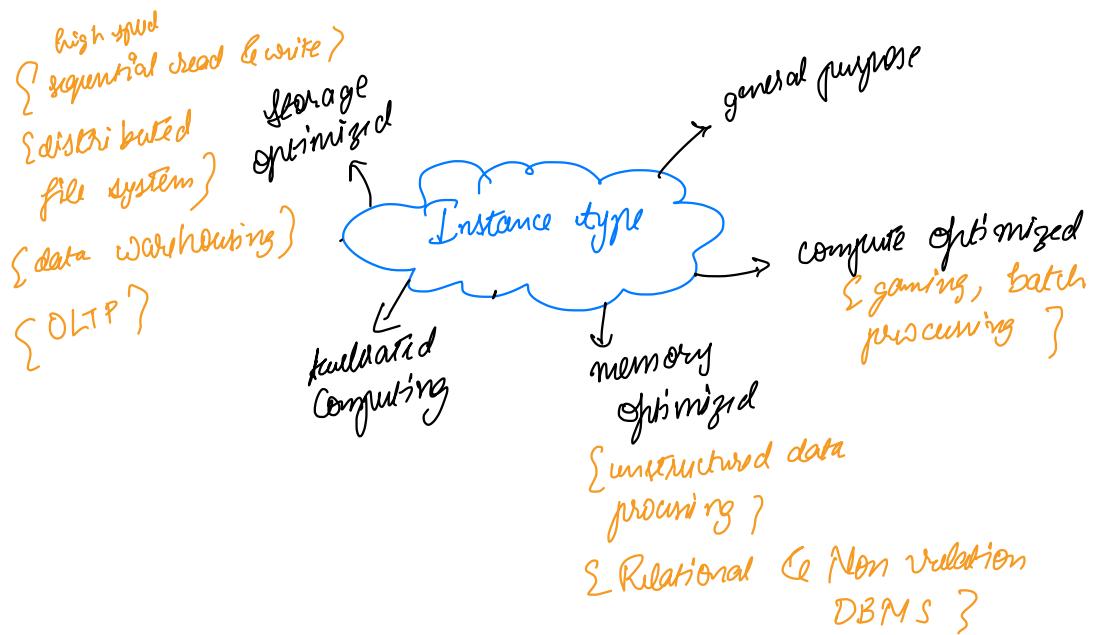
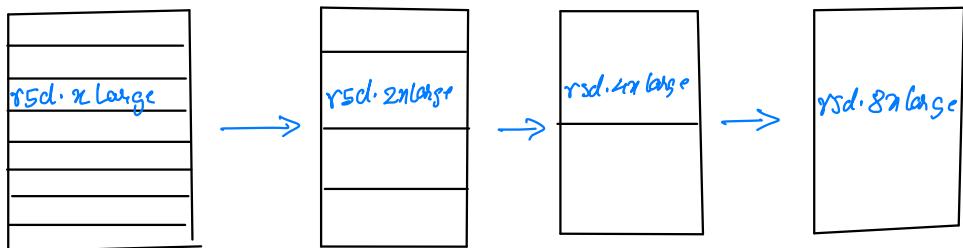
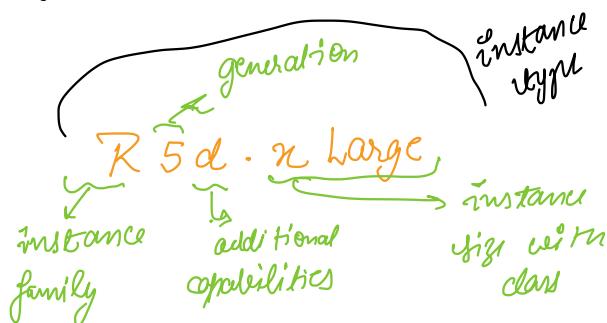
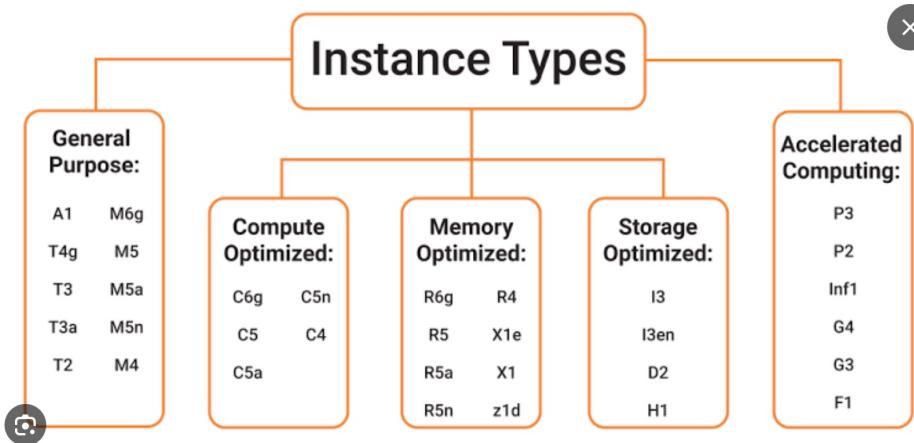


# EC2 { Elastic Cloud Compute }

"The servers provisioned in cloud"





	Type	Description	Mnemonic
<b>General Purpose</b>	a1	Good for scale-out workloads, supported by Arm	a is for Arm processor – or as light as A1 steak sauce
	t-family: t3, t3a, t2	Burstable, good for changing workloads	t is for tiny or turbo
	m-family: m6g, m5, m5a, m5n, m4	Balanced, good for consistent workloads	m is for main or happy medium
<b>Compute Optimized</b>	c-family: c5, c5n, c4	High ratio of compute to memory	c is for compute
<b>Memory Optimized</b>	r-family: r5, r5a, r5n, r4	Good for in-memory databases	r is for RAM
	x1-family: x1e, x1	Good for full in-memory applications	x is for xtreme
	High memory	Good for large in-memory databases	High memory is for... high memory.
	z1d	Both high compute and high memory	z is for zippy
<b>Accelerated Computing</b>	p-family: p3, p2	Good for graphics processing and other GPU uses	p is for pictures
	Inf1	Support machine learning inference applications	Inf is for inference
	g-family: g4, g3	Accelerate machine learning inference and graphics-intensive workloads	g is for graphics
	f1	Customizable hardware acceleration with field programmable gate arrays (FPGAs)	f is for FPGA or feel as in hardware
<b>Storage Optimized</b>	i-family: i3, i3en	SDD-backed, balance of compute and memory	i is for IOPS
	d2	Highest disk ratio	d is for dense
	h1	HDD-backed, balance of compute and memory	H is for HDD

## # Security groups

↳ fundamental to

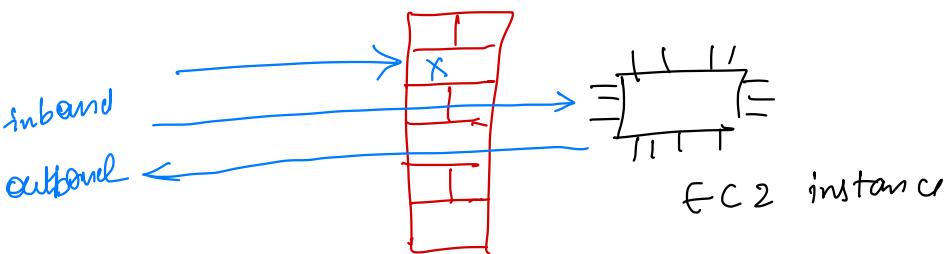
\* like a firewall cell

network security for an EC2 instance

\* **scoped by region!**

\* controls inbound & outbound traffic

{ by default all  
outbound traffic is  
allowed }



\* only contains "Allow rules"

\* lives outside EC2 instances & are like  
policies can be attached to multiple instances

security  
groups are  
stateful!



All inbound traffic is blocked  
All outbound traffic is allowed

- any connection timeout is because of SG
- "connection refused" is an application error

- \* SSH into a EC2 server requires public IP address of EC2 & not private
- \* Amazon Linux OS is a Cent-OS based & Red Hat based operating system, so to install packages we need to do "yum install" & not "apt-get" or "apt"

→ when trying to connect an EC2 instance via SSH we may encounter access denied due to many permission issues. go to properties > security & remove all users & add only yourself as the only user who is allowed to read & write the .pem file.

```
SSH -i .\\EC2.pem ec2-user@3.110.223.5
```

filename      user      public ip

## # different ways to connect to EC2 instance

Instance connect

{ cloud ⚡ &  
instance 22  
port } { }



SSH using  
port 22

{ requires access keys  
& secret access keys }



SSM session  
manager

{ gives access to EC2  
w/o SSH { port 22 }  
or bastion host  
or SSH access keys }

∴ port 22 is not  
required hence more  
secure

also logs the session  
data to S3 or  
CloudWatch logs.

## \* Older version of windows need PUTTY

w/o SSH  
but windows 10 onwards it's easy.

## # EC2 purchasing options

① On-demand pay as you go

② Reserved instances 72% discount

cheapest.

1 time payment 1y / 3y with options of upfront,  
partial upfront, monthly { regional or zonal scope }  
{ Convertible Reserv. Instan. → can change { 86% off } }  
family, type, OS, scope & tenancy

③ Savings plan 72% discount

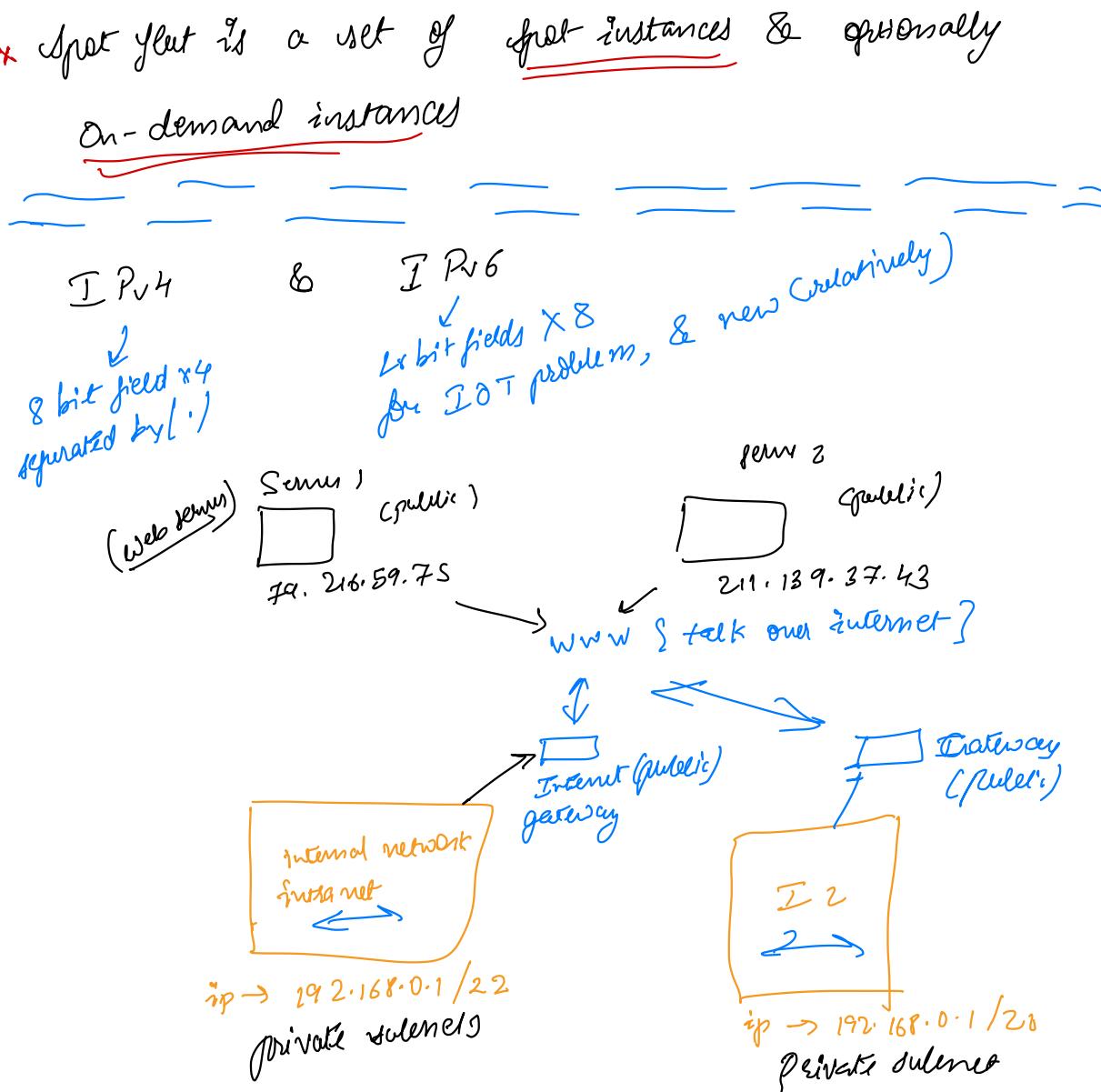
Commit to an amount of usage { 1y or 2y }  
At specific + family specific.  
If usage exceeds then considered on demand. { Shaded to a family }  
{ flexible in family }

- ④ spot instances { 90% discount }  
Bid on EC2 capacity, highly volatile, availability not guaranteed,  
AWS can reclaim an instance with a 2min warning.  
For jobs that can be interrupted. { Eg: Batch processing, analysis }  
\* highly volatile

- ⑤ Dedicated Host  
A physical server for an AWS account. Can launch multiple EC2 instances }  
Allow to address compliance requirements & use existing server bound software licenses (per-socket, per-core, per-VM)  
useful for softwares that have BYOL { Bring your own license } licensing model.  
visibility to low level hardware }

- ⑥ Dedicated Instances  
no other customer will share the hardware  
may share in same account.

- ⑦ Capacity reservations  
Reserve capacity in a specific AZ for any duration  
{ unlike RI, this pay as you go. for the time we keep the reservation active }  
for a period book the instances & pay full price even if you don't use it ?



\* private network ~~don't~~ machines talk to internet  
 standing behind a NAT + internet gateway  
 which acts as a "proxy."  
 network address Translation

# Elastic IPs  
whenever we stop & start an instance via <sup>public</sup> IP changes,  
to fix the public ip of our EC2 instance we need  
elastic IP. Can be assigned to instances & network  
interfaces.  
\* only have 5 Elastic IPs in our account. Can be true?

→ if any instance goes down we can map the traffic to  
another instance.

→ Consider a poor architectural design!

↳ instead use a random public IP & register  
a DNS name to it.  
or

use Load Balancer !

Q Problem with changing public IP is that it changes  
every time an instance is stopped & started {not  
reboot} -

→ Make an elastic IP & associate it with the  
EC2 or dissociate to detach,  
but make sure to remove the Elastic IP  
because it costs!

# Placement groups

how should our EC2 inst. be placed. in AWS

\* 3 diff. strategies

① Cluster → a low-latency group in a single AZ  
{risky?}

② spread → spread instances across different hardware { max 7 inst. per group per AZ }  
high availability

③ partition → spread instances across different logical partitions { only on diff. set of racks }  
upto 7 partitions per AZ too's of GC2 in a group! { for distributed workloads, Kafka, Cassandra, Hadoop }  
connectivity +

## # Elastic Network Interface { ENI }

logical component in a VPC that represents a virtual network card

instances → internet access, to FW or network card creation to dev't

→ has <sup>1</sup> primary IPv4 & one or more secondary IPv4

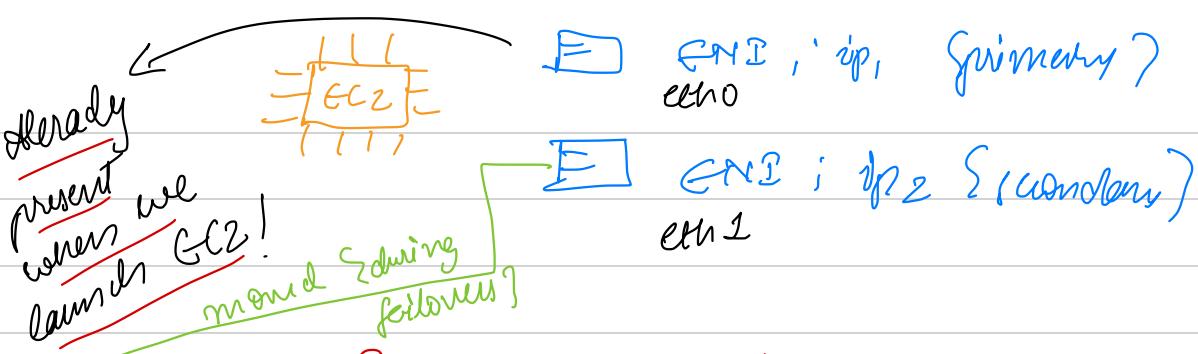
→ gives the private IPv4 to an GC2 instance

→ 1 Elastic IPv4 per private IPv4

→ 1 Public IPv4 auto assigned only to primary ENI (eth0)

→ 1 or more security groups

→ 1 MAC address



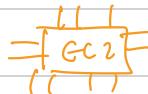
\* ENIs can be created independently & attached on fly.

\* AZ specific { Etho will also give itv? public IPv4 }

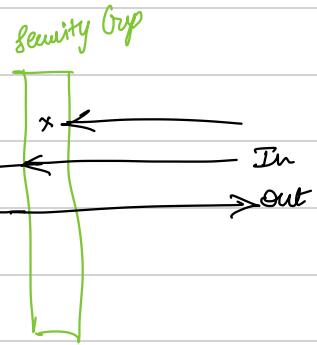


Eth0 - primary

So basically,



ENI  
Etho  
{ public ip  
private ip }



\* ENIs give a more control over the networking

## # EC2 hibernate

what actually happens?

loop → The instance is just not running  
 ALL network volumes (EBS) are intact  
 terminate → all attached network volumes are destroyed [if flags set?]

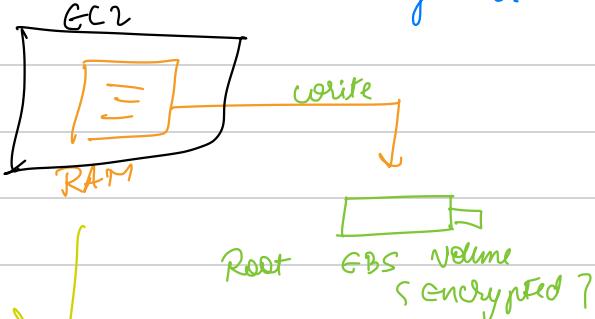
Starting

First start :- OS boots + EC2 user data script

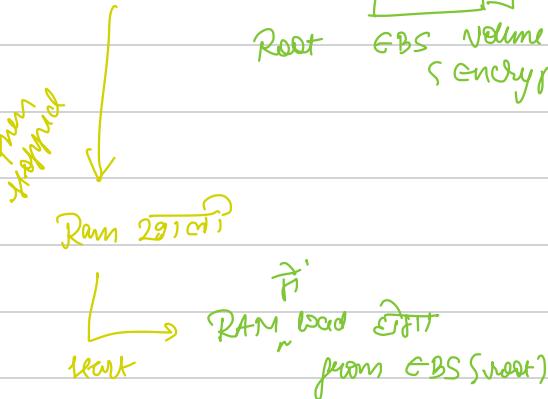
Following starts :- OS boots up

Starting up of application, caches local, etc.

hibernate → in-memory (RAM) is preserved  
 ⇒ faster boot



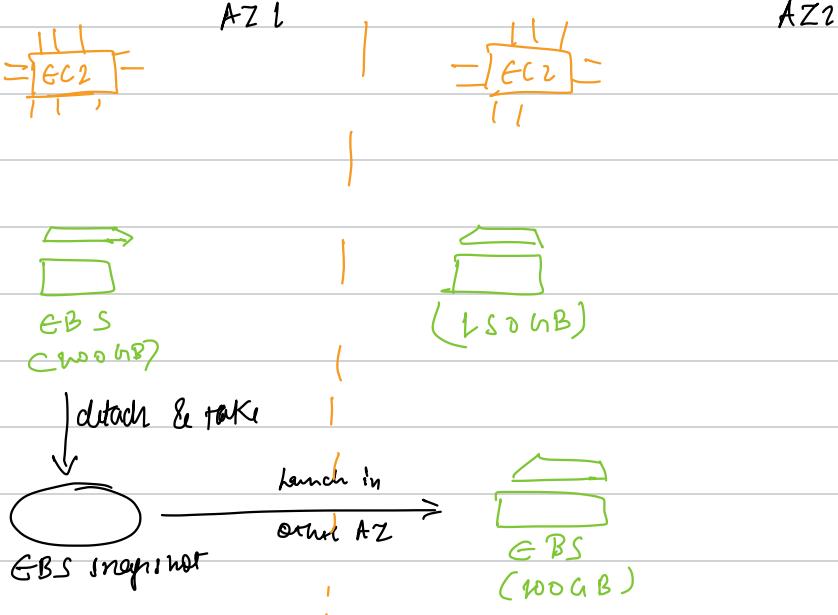
\* provided root EBS volume has enough spaces to store.



- \* instance RAM size < 150 GB
- \* Cannot hibernate for more than 60 days
- \* bare metal family does not support

# # EBS volumes {Elastic Block Storage}

- a "network USB stick" {a network drive}
- allows data to persist even after termination.
- can be attached to multiple instances
- AZ specific if need in another AZ then take EBS snapshot & then make an EBS volume out of that snapshot in another AZ



- has a "delete on termination" attribute
  - default "True" for root volume
  - "False" for other —

EBS Snapshot → it point in time snapshot of an EBS.

pointing of snapshots  
based on the size of  
the snapshot

Fast detach  
EBS be your  
snapshot!  
taken & reattach!

- ① EBS snapshot archive {1-3 days}
  - {FSI v. chapter? {24-72 hrs vs retain}}
- ② Recycle Bin
  - {Specify retention period?}
- ③ Fast Snapshot Restore (FSR)
  - {full ini of snapshot to home no latency}
  - {Good for big snapshot?}

## # AMI {Amazon Machine Image}

Snapshot  
of entire EC2  
instance! EBS, and  
with FSI, configuration  
all the

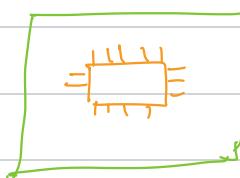
↓  
customizations  
of EC2 instance  
on our we care

{Can be found @}  
Market place

→ created for a specific AZ  
but can be copied across  
regions.

- \* public → AWS provided
- \* your own
- \* AWS market place

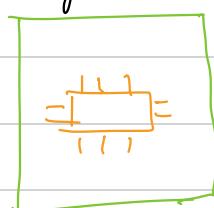
AZ1 (us-east-1a)



custom AMI



AZ2  
(us-east-1b)

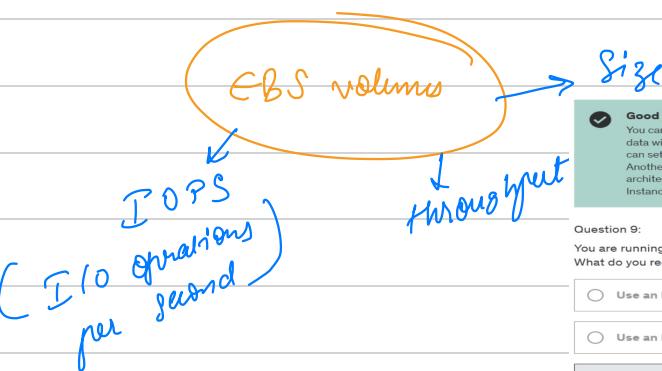


## # EC2 instance store

- hardware disk attached to EC2 instance
- attached to physical server
- extremely fast I/O
- not durable long term, will get deleted on termination, or may fail if EC2 fails.
- Backup & replication is user responsibility.
- good for cache / buffer etc & low time storage?

## # EBS volume types <sup>new gen</sup>

- ① gp2 / gp3 (SSD) general purpose <sup>{cont system}</sup>
  - ② io1 / io2 Block Express (SSD) high perform.
  - ③ ST1 (HDD) frequently access & throughput intensive
  - ④ SC1 (HDD) less frequent
- used as boot volumes*



Good job!

You can run a database on an EC2 instance that uses an Instance Store, but you'll have a problem that the data will be lost if the EC2 instance is stopped (it can be restarted without problems). One solution is that you can set up a replication mechanism on another EC2 instance with an Instance Store to have a standby copy. Another solution is to set up backup mechanisms for your data. It's all up to you how you want to set up your architecture to validate your requirements. In this use case, it's around IOPS, so we have to choose an EC2 Instance Store.

Question 9:

You are running a high-performance database that requires an IOPS of 310,000 for its underlying storage. What do you recommend?

Use an EBS gp2 drive

Use an EBS io1 drive

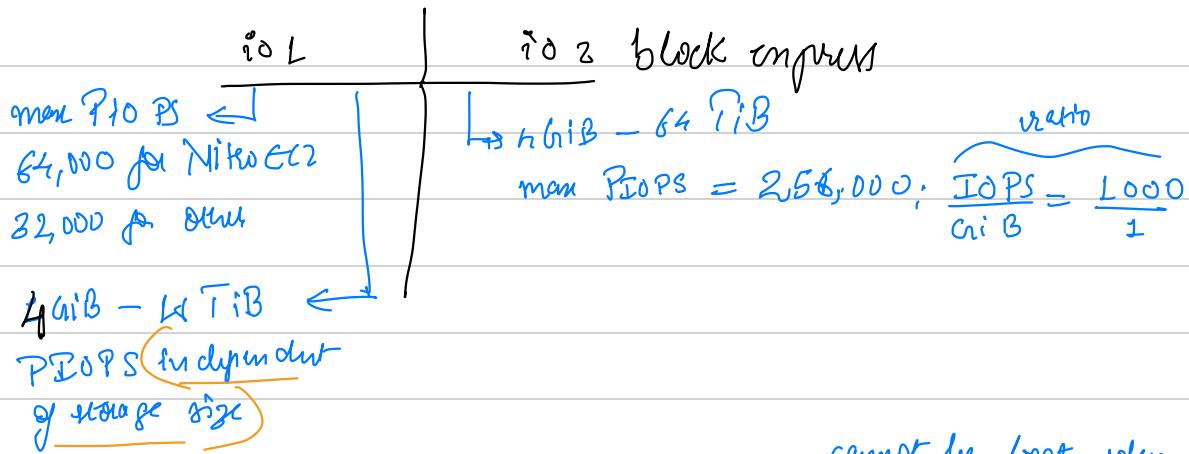
Use an EC2 Instance Store

Use an EBS io2 Block Express drive

General purpose SSDs  $\rightarrow$  cost effective, low latency  
 $3,000 < \text{IOPS} < 16,000$

gp 2 / gp 3 storage (newer) independent IOPS  
 { to indicate if storage increases  
 if storage increases  
 { 3 IOPS per GiB } (PIOPs)  
 Provisioned IOPS  $\rightarrow$  supports multi attach  
 business critical  
 $\sum_{i=1}^3 2500 = 7500$   
 $\frac{7500}{3} = 2500$  GiB needed for max IOPS

database workloads



flared disk drives (HDD)  $\leftarrow$  125 GiB - 16 TiB cannot be boot volume

(throughput optimized) STL  $\leftarrow$  SCL (cold HDD)

max throughput = 250 MiBs  $\downarrow$   
 250 IOPS

\* data warehousing

\* Log processing

\* Big data

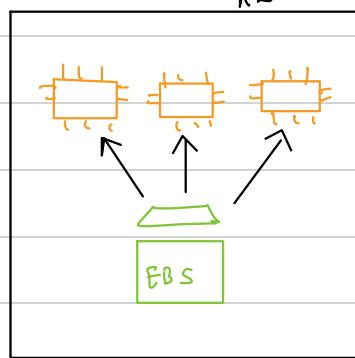
\* infrequently accessed data  
 \* lower cost

## \* Multi Attach { ior / io2 family }

- LEBS to multiple EC2 instances
- each instance has full read & write permissions. (with T) performance

use  
case

- ① Clustered Linux application  
for higher availability
- ② Concurrent write operations



AZ

\* 16 EC2 instances @ a time

\* must use file system that is cluster aware

## # EBS Encryption

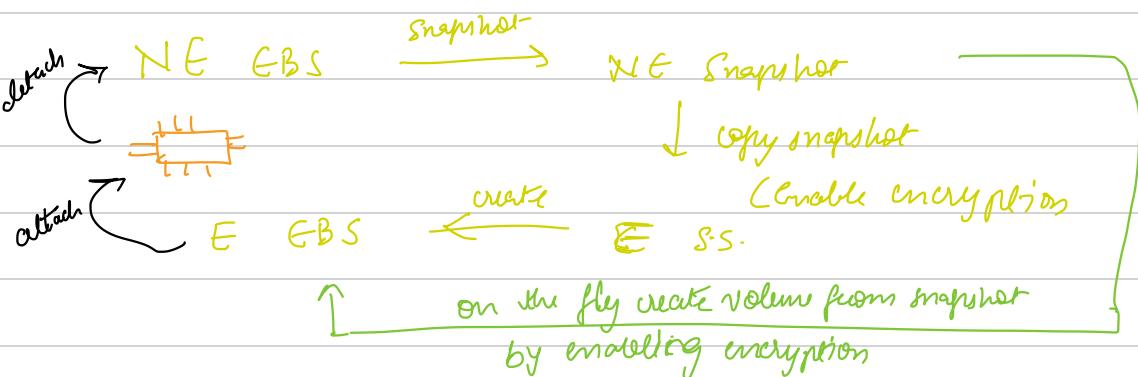
→ when enabled its automatically handled by EC2 & EBS.

when  
volume  
is created  
with EBS  
encryption,

- data in rest encrypted
- data in flight encrypted
- snapshots encrypted
- volumes from snapshots encrypted

- leverages encryption Keys from KMS (AES-256)
- Encrypted volumes & snapshots are encrypted

- Encryption  
of unencrypted  
EBS volumes!
- Create snapshot
  - Encrypt the EBS snapshot (using copy)
  - Create new volume from snapshot
  - attach encrypted EBS



## # Amazon EFS (Elastic File System)

- network file system (managed by AWS)
- mount of many EC2, in multi AZ
- high {
  - ↑ scalability
  - ↑ availability
  - ↑ pay per use

\* ~~Scalability~~

- uses NFS v4.1 protocol
- only compatible with Linux based AMI { Not ~~Windows~~ }
- POSIX { Linux } filesystem
- auto scaling i.e. no capacity planning.

## (1) Performance mode (set @ EFS creation time)

\* default  
general purpose

{ latency sensitive }

More I/O

{ higher latency, throughput & parallel }  
(big data, media proc.)

throughput → amount of bytes or bytes per second by a storage device

IOPS → no. of read/write operations per second.

→ throughput increases with increase in size

## (2) Throughput mode

→ Bursting

{  $LTB = 50 \text{ MiB/s} + \text{burst}$  }  
up to 100 MiB/s

Provisioned

{ set throughput regardless of storage size 1 GiB/s for LTB storage }

{ throughput independent of storage size }

service automatically scale throughput when unpredictable workload !?

## EFS - Storage Classes

storage tiers

1) Standard

2) Infrequent access [GFS-IA]

3) Archive {50% cheaper?}

{ with right option  
90% of cost saving  
can be done. }

→ implement lifecycle policies to move from different tiers {to save money?}

availability

durability

1) Standard {Multi AZ} \* default

2) One Zone {Single AZ} {can also set up Infrequent Access one zone?}

much cheaper