

# Week 1. Getting started with AWS

- 6个 cloud computing basic:

- Pay as you go. 不需要先建好 data center, hardware is in.

需要你有能力随时启动或停止.

- Benefit from massive economies of scale. 因为云 computing

能 aggregate 在 cloud, 能高 economy of scale, 因此云 cost, pay-as-you-go price 也便宜.

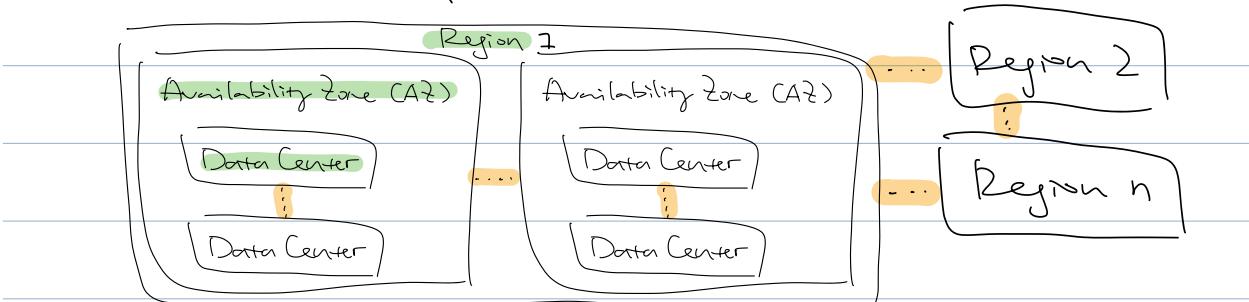
- Stop guessing capacity. 做好 capacity decision 之后, 避免 w/ setting on expensive idle resources & dealing with limited capacity 问题. 因 cloud 可以 scale up or down as required.

- Increase speed and agility.

- No money running / maintaining data center.

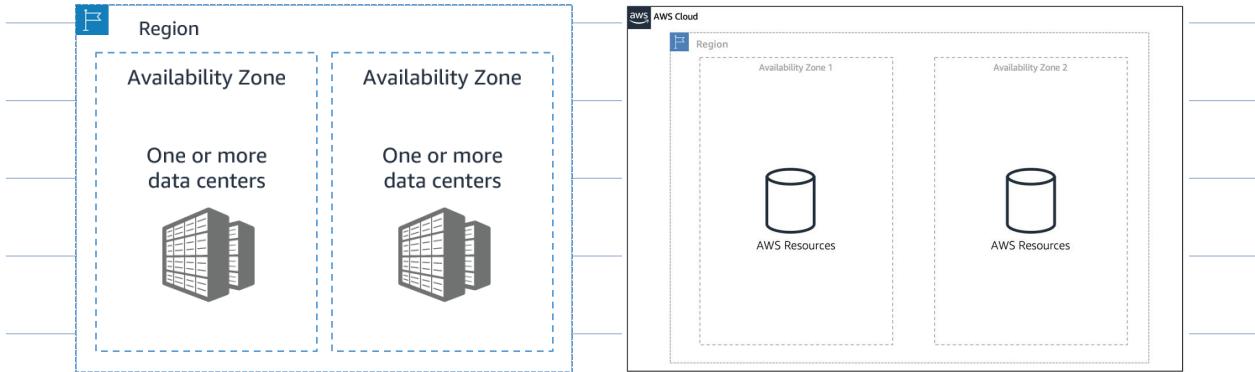
- Go global in minutes.

- AWS Global Infrastructure.



选择 Region 的标准: ① compliance ② latency ③ pricing  
④ service availability

技术上 latency 很大问题。用 edge cache 来解决。



## Interacting with AWS

- 不要自己物理管理，而是通过 API (Application Programming Interface)。

- 3 main ways to interact with AWS API.

① AWS Management Console [https://aws.amazon.com/console/](#)

② AWS Command Line Interface (CLI)

③ AWS SDK Developer Kits (SDKs)

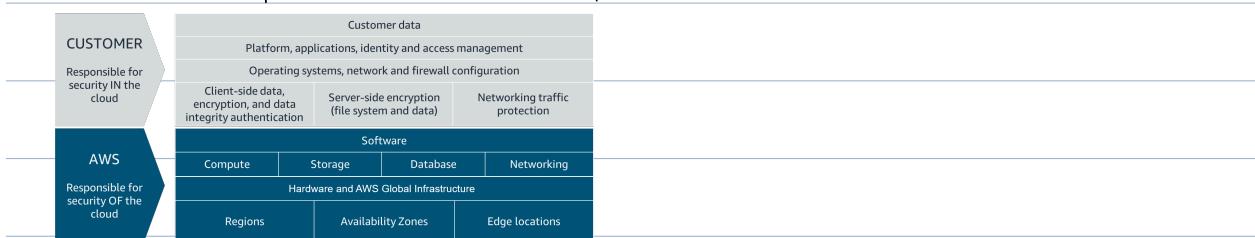
## Security in AWS Clouds

用户对 AWS 云负责其安全责任。

'Shared Responsibility Model'

AWS : security of the cloud

user : .. in ..



Category	Examples of AWS Services in the Category	AWS Responsibility	Category	AWS Responsibility	Customer Responsibility
Infrastructure services	Compute services, such as Amazon Elastic Compute Cloud (Amazon EC2)	AWS manages the underlying infrastructure and foundation services.	Infrastructure services	AWS manages the infrastructure and foundation services.	You control the operating system and application platform, as well as encrypting, protecting, and managing customer data.
Container services	Services that require less management from the customer, such as Amazon Relational Database Service (Amazon RDS)	AWS manages the underlying infrastructure and foundation services, operating system, and application platform.	Container services	AWS manages the infrastructure and foundation services, operating system, and application platform.	You are responsible for customer data, encrypting that data, and protecting it through network firewalls and backups.
Abstracted services	Services that require very little management from the customer, such as Amazon Simple Storage Service (Amazon S3)	AWS operates the infrastructure layer, operating system, and platforms, as well as server-side encryption and data protection.	Abstracted services	AWS operates the infrastructure layer, operating system, and platforms, as well as server-side encryption and data protection.	You are responsible for managing customer data and protecting it through client-side encryption.

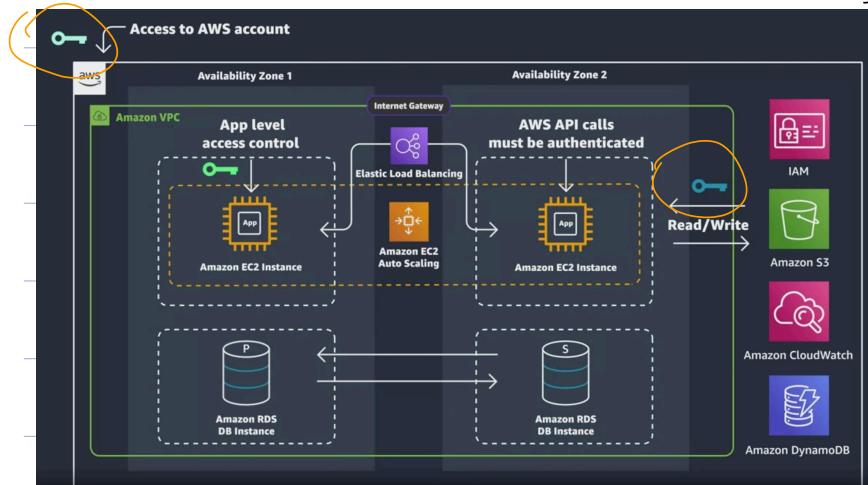
## ⑥ Protecting AWS Rust Users

- When sign up - You own account / You are not user.  
→ log in

→ great power = great responsibility

- ② MFA (multi-factor authentication) 要有 2 个

## ⑦ AWS Identity and Access Management (IAM)



• 這就是 IAM 啟

其實 AWS 在 action, 都是 API call.

基本上都是 JSON language.

## @ IAM

### - Defi :

- IAM is a web service for managing access to AWS account and resources.
- It provides centralized control with authentication and authorization.
  - access identification
  - permission
- It allows to share access without sharing credentials  
并統括の  
and offers granular access control to resources.

### - Features :

- global, not specific to any region.
- integrated with many AWS service.  
(e.g. 用 fb 登 x).
- supports MFA, identity federation.

### - IAM user credentials.

- access to AWS Management Console.
- Programmatic access to AWS Command Line Interface (AWS CLI) and AWS Application Programming Interface (AWS API)

### - IAM group . 組の権限管理 .

### - IAM Policy Example

```
{  
  "Version": "2012-10-17",  
  "Statement": [  
    {  
      "Effect": "Allow",  
      "Action": "*",  
      "Resource": "*"  
    }]  
}
```

Vision : policy language 格式

Effect : 這個 statement 會 allow / deny

Action : + file 開戶. 要不然就寫清楚

Resource : 這個 寶 告 over in obj

Ty

```
{"Version": "2012-10-17",
 "Statement": [
     {
         "Effect": "Allow",
         "Action": [
             "iam: ChangePassword",
             "iam: GetUser"
         ],
         "Resource": "arn:aws:iam::123456789012:user/${aws:username}"
     }
 ]}
```

只许用这个 resource 为  
这个文件

• Role based access in AWS

权限像给个临时许可证

## Week 2: AWS compute & networking

### • AWS Compute

- Host application  $\xrightarrow{\text{HTTP}}$  server.
- server  $\xrightarrow{\text{HTTP}}$  handle HTTP request & send response to user  
HTTP service  $\xrightarrow{\text{Ports}}$ :
  - Windows option (E.g. Internet Information Service)  $\xrightarrow{\text{IIS}}$
  - Linux option (E.g. Apache HTTP web server, Nginx..)



To use AWS to run HTTP,

From the AWS Management Console  $\xrightarrow{\text{AWS}}$

With computer power as service ( $\xrightarrow{\text{AWS}}$ )

### • EC2: Amazon Elastic Compute Cloud

- EC2, a web service provides secure, resizable compute capacity in cloud. Allows to provision virtual service called EC2 instance.

- To create EC2 instance, need to define:

storage ...)

① hardware specifications (E.g. CPU, memory, network)

② logical configuration (E.g. Networking location, firewall rules, authentication, OS).

When launching EC2 instance  $\xrightarrow{\text{AWS}}$ , it's like creating an Amazon Machine Image (AMI) with the OS.

Image (AMI) 里选用的 OS.

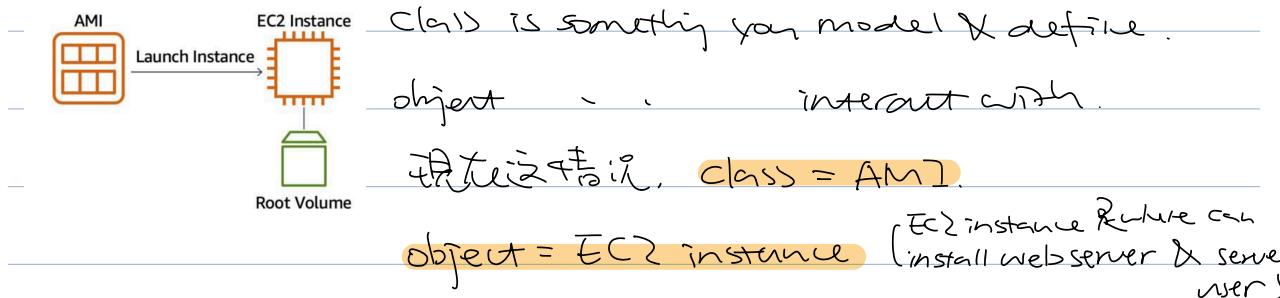
## • AMI

- The AWS cloud 中 OS installation 不是 user 的任务，而是 OS 已经被建在 AMI 中。为什么安装 OS，the traditional infrastructure 中，为了 spinning up server.

## • Relationship btw AMI & EC2 instance.

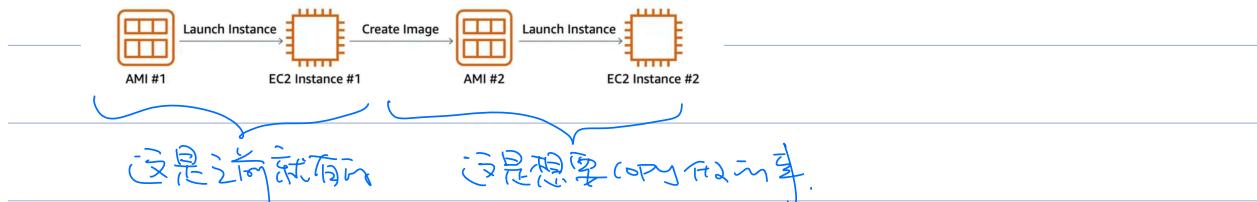
EC2 instance are live instantiations of what is defined

in an AMI, like the relationship btw class and object.

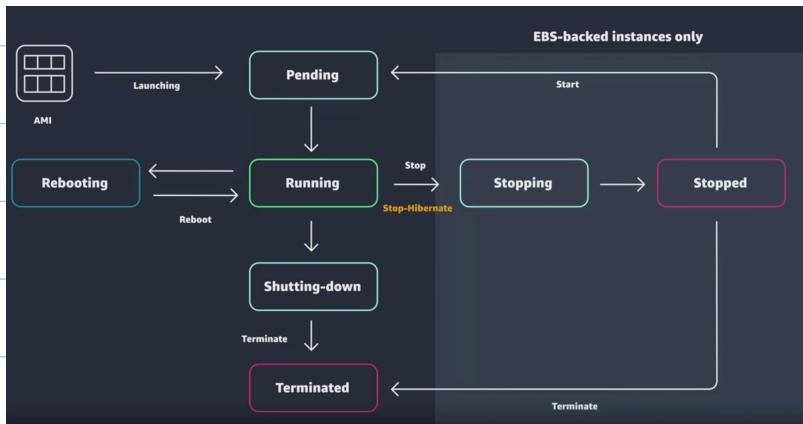


- When you launch a new instance, AWS allocates a virtual machine that runs on a hypervisor. Info to the AMI you selected is copied to the **root device volume**, which contains the image used to boot the volume. From there you get a server you can connect to and install packages and any additional SW.

## - AMI is reusable.



## Amazon EC2 instance lifecycle



- terminate 不一定是坏事, 可能是 SW update, 先 copy-4 例有  
在 EC2 instance, 在这个上进行 update, 然后才正式 terminate  
结束那个老的。

- Change to running state in EC2 instance.
- EC2 instance components:

virtual processors (vCPU) + memory + network

(+ instance storage + graphic processing units (GPU))

- instance type = prefix identifying + size

Instance Family	Description	Use Cases
General purpose	Provides a balance of compute, memory, and networking resources, and can be used for a variety of workloads.	Scale-out workloads such as web servers, containerized microservices, caching fleets, distributed data stores, and development environments.
Compute optimized	Ideal for compute-bound applications that benefit from high-performance processors.	High-performance web servers, scientific modeling, batch processing, distributed analytics, high-performance computing (HPC), machine/deep learning, ad serving, highly scalable multiplayer gaming.
Memory optimized	Designed to deliver fast performance for workloads that process large data sets in memory.	Memory-intensive applications such as high-performance databases, distributed web-scale in-memory caches, mid-size in-memory databases, real-time big-data analytics, and other enterprise applications.
Accelerated computing	Use hardware accelerators or co-processors to perform functions such as floating-point number calculations, graphics processing, or data pattern matching more efficiently than is possible with conventional CPUs.	3D visualizations, graphics-intensive remote workstations, 3D rendering, application streaming, video encoding, and other server-side graphics workloads.
Storage optimized	Designed for workloads that require high sequential read and write access to large data sets on local storage. They are optimized to deliver tens of thousands of low-latency random I/O operations per second (IOPS) to applications that replicate their data across different instances.	NoSQL databases, such as Cassandra, MongoDB, and Redis, in-memory databases, scale-out transactional databases, data warehousing, Elasticsearch, and analytics.

(VPC)

- By default, EC2 instance is in **Amazon Virtual Private Cloud**.  
有錢可以各種 transmission cost, 當熟悉 AWS 後, 可以降低成本

**customer VPC**, 這樣能 restrict access with additional routing and connectivity mechanism.

- **Architect for High Availability**:

The network, instance state, and availability zone.

The AZ 中的 AWS Service 必須 be architect of high availability.

EG: Multiple EC2 instance component 中有 size.

- An EC2 instance transmission between different states 具有彈性。

- Reverse Capacity with **Reserved Instance (RI)**

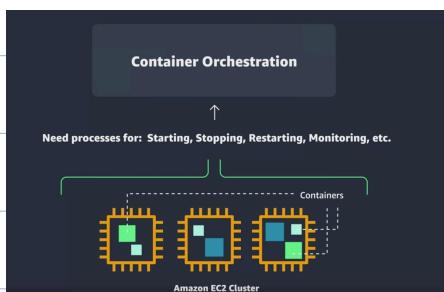
RI & on-Demand instance 重複。

- Container Services on AWS

- container orchestration services:

- Amazon Elastic Container Service (Amazon ECS)

- - - - - Kubernetes (Amazon EKS)



當 container 數量太多時 宜

很難 manage → EKS

## ☞ Container

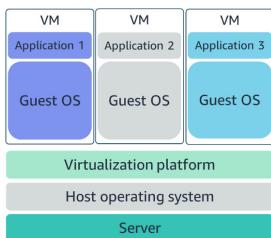
• defi: containers are used as a solution to problems of traditional compute, including the issue of getting SW to run reliably when it moves from one compute environment to another. It's a standardized unit that packages up your code and all of its dependents.

## ☞ Docker

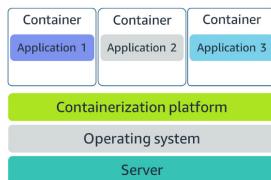
Docker is a popular container runtime that simplifies the management of the entire OS stack needed for container isolation, including networking and storage. Docker makes it easy to create, package, deploy and run.

### ⑥ WHAT IS THE DIFFERENCE BETWEEN CONTAINERS AND VMs?

Virtual machines (VMs)



Containers



## ☞ Orchestrate Containers

In AWS, containers run on EC2 instances. 常常向 instance 中 run 多个 containers. 每个 instance 可以跑多个容器, 但 lack of availability and scalability. - 公司希望多个 container 能够在不同 instance 上运行, across AZ.

↳ manage compute at a large scale ↗, need to know:

① How to place containers on instance.

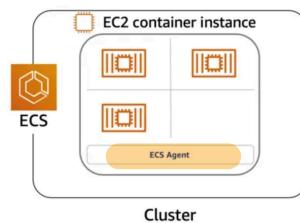
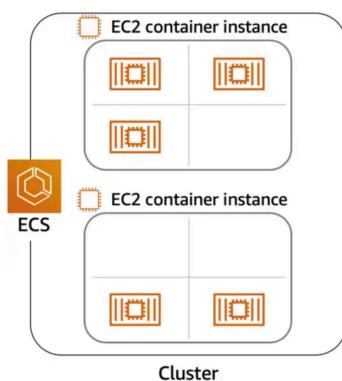
② What happens if container fails.

③ ' ' ' instance - .

④ How to monitor deployments of containers.

#### MANAGE CONTAINERS WITH AMAZON ELASTIC CONTAINER SERVICE (AMAZON ECS)

Amazon ECS is an end-to-end container orchestration service that allows you to quickly spin up new containers and manage them across a cluster of EC2 instances.



To run & manage,

need a **ECS Agent**.

To prepare your application to run on ECS, need to create task definition in JSON.

Eg. simple task definition runs on Nginx web server.

```
{  
    "family": "webserver",  
    "containerDefinitions": [  
        {  
            "name": "web",  
            "image": "nginx",  
            "memory": "100",  
            "cpu": "99"  
        }],  
        "requiresCompatibilities": [ "FARGATE" ],  
        "networkMode": "awsvpc",  
        "memory": "512",  
        "cpu": "256"  
    ]}
```

→ Kubernetes

⇒ ECS 与 EKS . EKS 与 ECS 是  
conceptually similar, 但它们是:

① EC2 instance with ECS Agent installed  
and configured as container instance.  
With EKS as worker node.

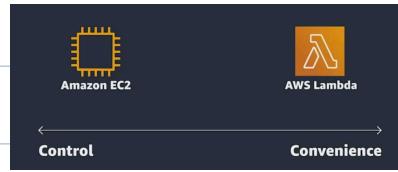
② ECS container as task.

In EKS ecosystem, as pod.

③ ECS runs on AWS native technology.  
EKS runs on top of Kubernetes.

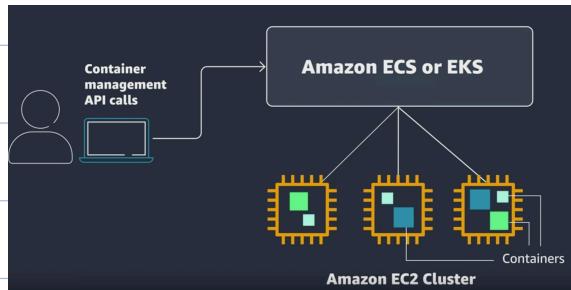
## • Serverless

只需要关注小生意的，剩下的交给 AWS.

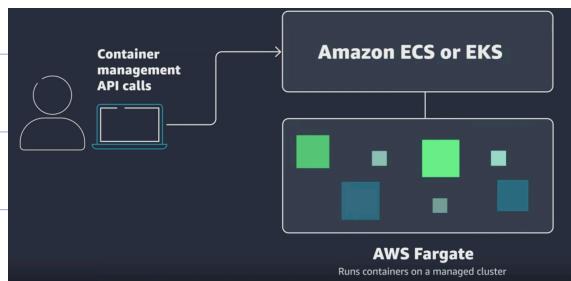


## • Serverless with AWS Fargate

之前：



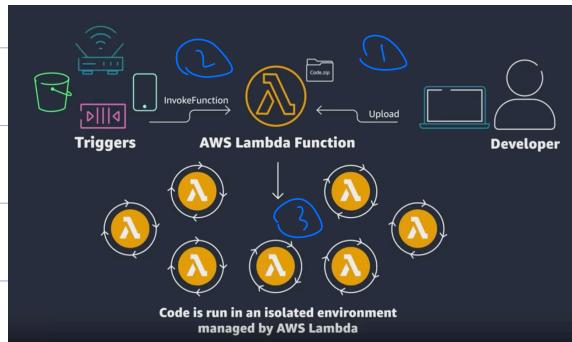
现在：



EC2 肯定不是 serverless

Fargate 是 serverless

## • AWS Lambda



①: 自己写 Lambda function.

②: 当有 trigger 触发后会自动调用

call Lambda function.

③: 可以有多个入 func, 都是 isolated.

Lambda func 不适用于需要长期在线的 app, 而是等到条件满足  
之后再使用。

courseware 里说后使用 sample 来表示。

tiny summary: AWS compute service { EC2, ECS/EKS , Lambda }

## • AWS Networking

### ◦ Networking on AWS

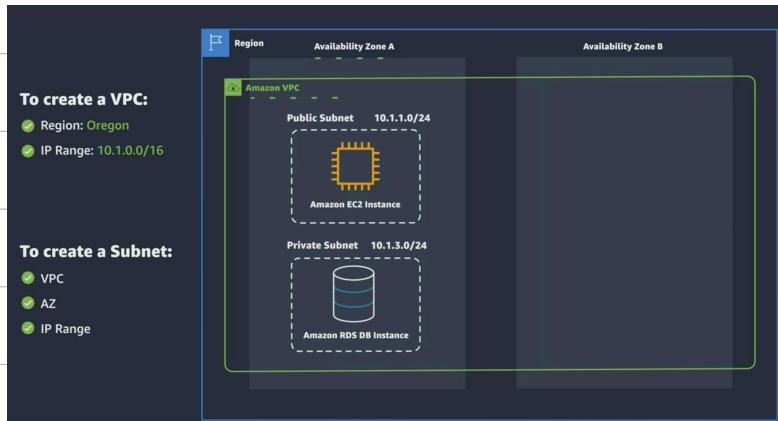
- **Networking**: how you connect computers around the world and allow them to communicate with one another.

### ◦ Amazon VPC (Virtual Private Cloud)

- VPC create walls/boundaries where application and resources is isolated from outside.

◦ To create VPC we need: ① region ② IP Range in CIDR notation

To create subnet . need: ① VPC ② AZ ③ IP Range



◦ **CIDR** (Classless Inter-Domain Routing) notation.

E.g. 192.168.1.0/24

fixed      flexible      the number of 24 bits are fixed.

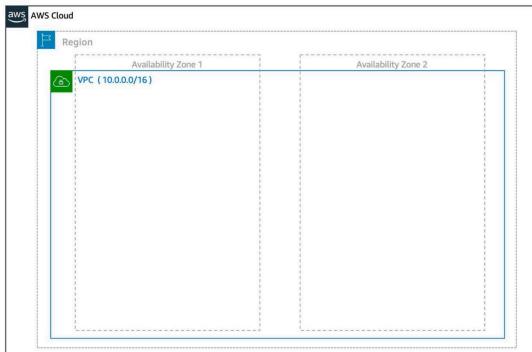
即  $32 - 24 = 8$  bits flexible,  $2^8 = 256$

IP 可以用.

.. 0/24 in range to ... 0/16 range in. ( $2^{32-16} = 65,536$ )

- Eg. 実操 Amazon VPC

### Step 1: Create VPC



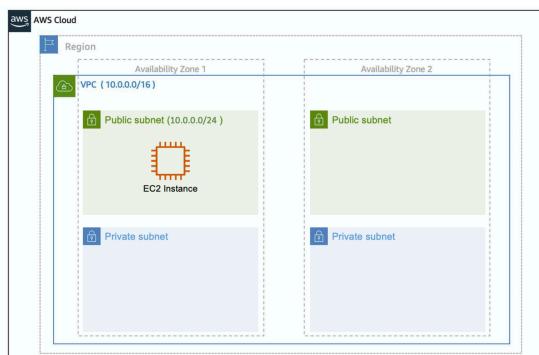
要:

① Name of VPC

② Region of VPC to live on

③ IP Range for VPC in CIDR

### Step 2: Create Subnet



- Subnets = Smaller networks inside base network or a virtual area networks (VLANs)
- Create subnets with availability and providing diff connectivity options.

要:

- ① VPC want to live in
- ② AZ want to live in
- ③ CIDR block of subnet

- High Availability with a VPC



COPY PASTE していきなさい。

### · Reserved IP

- ② DNS (Domain Name System)
- 這些 Reserved IP 用于 ① routing ③ networking management.

## Gateway

### Internet Gateway (IGW)

To ensure internet connectivity for VPC.

### Virtual Private Gateway (VGW)

Allows to connect AWS VPC to another private network.

## AWS VPC Routing and Security

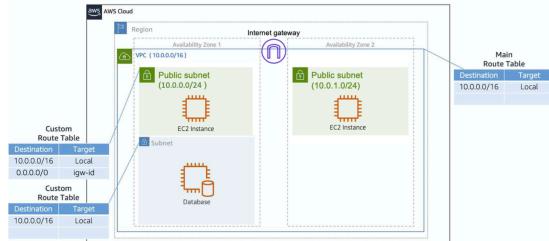
### Main Route Table

当创建 VPC 时，AWS 自动生成 main route table，用来  
determine where network traffic is directed. AWS 会自动  
当创建 VPC with subnets 时，你想要 traffic to follow  
between them.

### Custom Route Table

want to be more granular about how you route your  
traffic for specific subnets.

当有 custom route table 时，  
不会使用 main route table。  
By default, custom  
use local table.

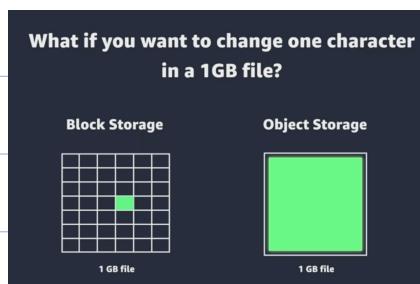
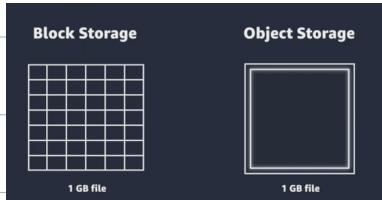


# Week 3: Storage & Databases on AWS

## • Storage on AWS

### • Storage Types on AWS

- ① Block
- ② Object
- ③ File



→ object storage often follows a WORM model

### - File Storage

- ideal when require a centralized access to files that need to be easily shared and managed by multiple host computers.
- typically, this storage is mounted onto multiple hosts and requires **file locking and integration** with existing file system communication protocols.
- common use:**
  - ① large content repositories
  - ② development environment
  - ③ user home directories

### - Block Storage

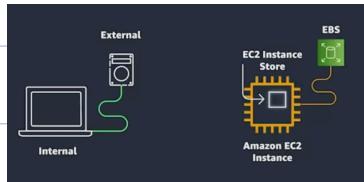
- Split files into file-size chunks of data cell (block) that have their own address, so each block can be retrieved effectively
- is optimized for low latency operations**, it's a typical

storage for high-performance enterprise workloads, such as databases or enterprise resource planning (ERP) systems.

- Object Storage.

- Stored in a flat structure instead of hierarchy. Each obj is a file with unique identifier. Identifier, along with any additional metadata, is bundled with the data and stored.
- Can store almost every type of data, and no limit to the number of obj stored, making it easy to scale.
- **常用于** storing large data sets, unstructured files.

- Amazon EC2 Instance Storage and Amazon Elastic Block Storage.



### EC2 Instance Storage:

①: Fast

ephemeral storage

②: life cycle bounded with instance, so it's often referred as <sup>✓</sup>.

### EBS :

①: 像个硬盘. life cycle 不受限 能从一个 AZ 中在其他 instance 附上.

②: EBS 与 instance 可以解绑, EBS Multi-Attach.

## • EBS.

- Block-level storage can attach to EC2 instance.

• Storage devices are called Amazon EBS volumes. (亚马逊大王)

### • Scale EBS Throughput:

① Increase volume size, <sup>超过</sup> size limit (16 TB).

② Attach multiple volumes to single Amazon EC2 instance.

### • Use case:

① DSS ② DB ③ enterprise app ④ throughput-intensive application.

Amazon EBS Volume Types

	EBS Provisioned IOPS SSD	EBS General Purpose SSD	Throughput Optimized HDD	Cold HDD
Description	Highest performance SSD designed for latency-sensitive transactional workloads	General purpose SSD that balances price and performance for a wide variety of transactional workloads	Low-cost HDD designed for frequently accessed, throughput intensive workloads	Lowest cost HDD designed for less frequently accessed workloads
Use Cases	I/O-intensive NoSQL and relational databases	Boot volumes, low-latency interactive apps, development, and test	Big data, data warehouses, log processing	Colder data requiring fewer scans per day
Volume Size	4 GB-16 TB	1 GB-16 TB	500 GB-16 TB	500 GB-16 TB
Max IOPS/Volume	64,000	16,000	500	250
Max Throughput/Volume	1,000 MB/s	250 MB/s	500 MB/s	250 MB/s

## • Obj Storage with Amazon S3

### • Amazon S3 (an object storage service)

• it's a standalone storage that isn't tied to compute.

• It enables retrieving data from anywhere else on the web.

• In S3 there is no state container called buckets.

Steps: ① Region ② bucket name, must be unique across all AWS accounts.

Bucket

Object/Key

<http://doc.s3.amazonaws.com/2006-03-01/AmazonS3.html>

• Use case :

① backup & storage ② media sharing ③ SW delivery

④ data lakes ⑤ static website ⑥ static content.

• S3 bucket policy

```
{  
    "Version": "2012-10-17",  
    "Statement": [  
        {  
            "Sid": "PublicRead",  
            "Effect": "Allow",  
            "Principal": "*",  
            "Action": "s3:GetObject",  
            "Resource": ["arn:aws:s3:::employebucket/*"]  
        }  
    ]  
}
```

• Encrypt S3:

encrypt when saving

- Server side encryption : decrypt when downloading.

- Client-side encryption : 自己解密。

• Use Versioning → Preserve Objects.

• Amazon S3 Storage Classes.

1. S3 Standard

2. S3 Intelligent-Tiering

3. S3 Standard-Infrequent Access

;

- Choose the right storage service.

- EC2 Instance Storage

- EBS

- S3

- EFS (Elastic File System) & FSx

## ● Database on AWS

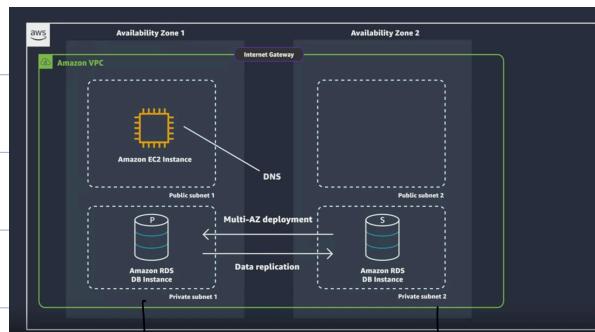
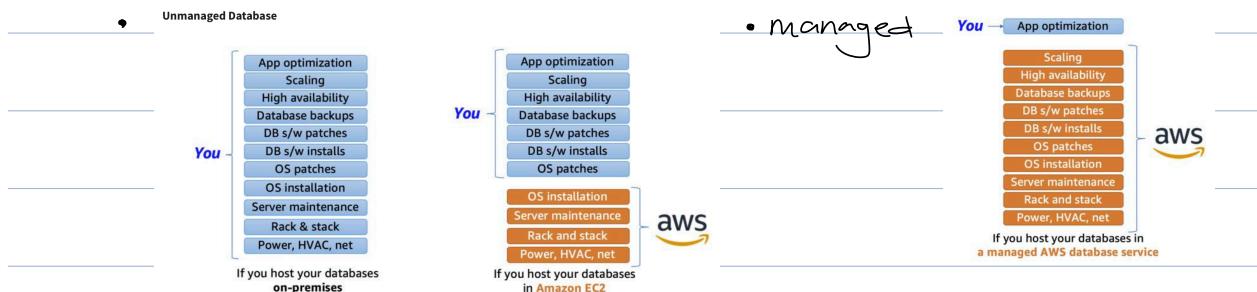
### • Database on AWS

- Relational DB (MySQL, Amazon Aurora)

→ super fast & large volume

- ① joins   ② reduced redundancy

- ③ familiarity   ④ accuracy.



↓ Primary DB      ↓ Secondary DB

当两个 AZ 中的 DB 连不上时会连 Secondary DB.

## \* Amazon Relational DB Service (RDS).

· Types:

① commercial: Oracle. SQL Server

② open source: MySQL. MariaDB. PostgreSQL

③ cloud native: Amazon Aurora.

· DB instance:

① Standard ② memory optimized ③ burstable performance.

· EBS (Elastic Block Storage) volume storage:

① General Purpose (SSD)

② Provisioned IOPS (SSD)

③ Magnetic storage (not recommended)

· Work with Amazon RDS in Amazon Virtual Private Cloud.

· To create DB subnet group:

① At include the subnet want to added

② subnet in that At where DB instances are placed.

· # AWS Identity and Access Management (IAM) # 身份認證安全

· Backup data: ① automatic backup

② manual snapshots

· Purpose Built DB on AWS

· # DB 來源於 purpose data. → 這 DB, 不該 data.



Amazon DynamoDB  
NoSQL database



Amazon DocumentDB  
with MongoDB compatibility



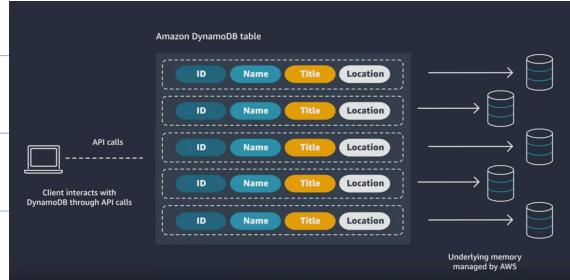
Amazon Neptune  
graph database



Amazon QLDB  
ledger database

## Amazon DynamoDB

不用自己建 table schema, 只需要提 data 定义



不用自己 maintenance

Amazon DynamoDB 处理能力

10 billion+ request

Database Type	Use Cases	AWS Service
Relational	Traditional applications, ERP, CRM, e-commerce	Amazon RDS, Amazon Aurora, Amazon Redshift
Key-value	High-traffic web apps, e-commerce systems, gaming applications	Amazon DynamoDB
In-memory	Caching, session management, gaming leaderboards, geospatial applications	Amazon ElastiCache for Memcached, Amazon ElastiCache for Redis
Document	Content management, catalogs, user profiles	Amazon DocumentDB (with MongoDB compatibility)
Wide column	High-scale industrial apps for equipment maintenance, fleet management, and route optimization	Amazon Keyspaces (for Apache Cassandra)
Graph	Fraud detection, social networking, recommendation engines	Amazon Neptune
Time series	IoT applications, DevOps, industrial telemetry	Amazon Timestream
Ledger	Systems of record, supply chain, registrations, banking transactions	Amazon QLDB

# Week 4: Monitoring, Optimizing and Curing Service on AWS

## • Monitoring on AWS

- Monitor before user notice

- U of Monitoring:

① Respond to operational issue proactively before you end

user are aware of them.

- ② Improve the performance and reliability of your resources.
- ③ Recognize security threats and events
- ④ Create more cost-effective solutions.

- Gain visibility.

- Detect anomalous behavior in your environments.
- Set alarms to alert you when something's not right.
- Visualize logs and metrics with the AWS Management Console.
- Take automated actions like scaling.
- Troubleshoot issues.
- Discover insights to keep your applications healthy.

## • Amazon CloudWatch.

- For Basic Monitoring & Detailed Monitoring.

- Break down metrics

namespace, dimensions.

### Learn the CloudWatch Logs Terminology

Log data sent to CloudWatch Logs can come from different sources, so it's important you understand how they're organized and the terminology used to describe your logs.

**Log event:** A log event is a record of activity recorded by the application or resource being monitored, and it has a timestamp and an event message.

**Log stream:** Log events are then grouped into log streams, which are sequences of log events that all belong to the same resource being monitored. For example, logs for an EC2 instance are grouped together into a log stream that you can then filter or query for insights.

**Log groups:** Log streams are then organized into log groups. A log group is composed of log streams that all share the same retention and permissions settings. For example, if you have multiple EC2 instances hosting your application and you are sending application log data to CloudWatch Logs, you can group the log streams from each instance into one log group. This helps keep your logs organized.

## ◆ Optimization

### ⇒ Optimizing Solutions on AWS

#### - Availability -

Availability (%)	Downtime (per year)
90% ("one nine")	36.53 days
99% ("two nines")	3.65 days
99.9% ("three nines")	8.77 hours
99.95% ("three and a half nines")	4.38 hours
99.99% ("four nines")	52.60 minutes
99.995% ("four and a half nines")	26.30 minutes
99.999% ("five nines")	5.26 minutes

#### ◦ Improve Application Availability.

#### ◦ Use 2<sup>nd</sup> AZ

#### ◦ Manage Replication, Redirection, and High Availability.

#### ◎ Amazon EC2 Auto Scaling.

#### • Route Traffic with Amazon Elastic Load Balancing.

##### • Load Balancer

###### - FEATURES OF ELB

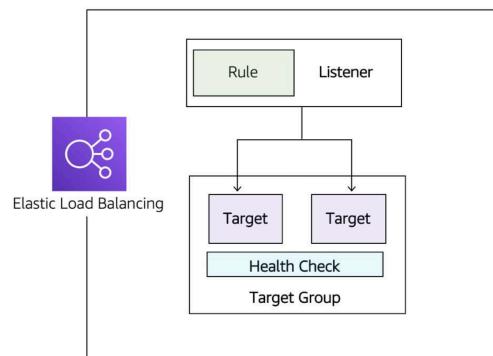
The ELB service provides a major advantage over using your own solution to do load balancing, in that you don't need to manage or operate it. It can distribute incoming application traffic across EC2 instances as well as containers, IP addresses, and AWS Lambda functions.

- The fact that ELB can load balance to IP addresses means that it can work in a hybrid mode as well, where it also load balances to on-premises servers.
- ELB is highly available. The only option you have to ensure is that the load balancer is deployed across multiple Availability Zones.
- In terms of scalability, ELB automatically scales to meet the demand of the incoming traffic. It handles the incoming traffic and sends it to your backend application.

##### • Health check

###### - ELB COMPONENTS

The ELB service is made up of three main components.



- **Listeners:** The client connects to the listener. This is often referred to as client-side. To define a listener, a port must be provided as well as the protocol, depending on the load balancer type. There can be many listeners for a single load balancer.
- **Target groups:** The backend servers, or server-side, is defined in one or more target groups. This is where you define the type of backend you want to direct traffic to, such as EC2 Instances, AWS Lambda functions, or IP addresses. Also, a health check needs to be defined for each target group.
- **Rules:** To associate a target group to a listener, a rule must be used. Rules are made up of a condition that can be the source IP address of the client and a condition to decide which target group to send the traffic to.

## Network Load Balancer (NLB)

**NLB uses a flow hash routing algorithm.** The algorithm is based on:

- The protocol
- The source IP address and source port
- The destination IP address and destination port
- The TCP sequence number

If all of these parameters are the same, then the packets are sent to the exact same target. If any of them are different in the next packets, then the request may be sent to a different target.

**NLB has sticky sessions.** Different from ALB, these sessions are based on the source IP address of the client instead of a cookie.

**NLB supports TLS offloading.** NLB understands the TLS protocol. It can also offload TLS from the backend servers similar to how ALB works.

**NLB handles millions of requests per second.** While ALB can also support this number of requests, it needs to scale to reach that number. This takes time. NLB can instantly handle this amount of requests.

**NLB supports static and elastic IP addresses.** There are some situations where the application client needs to send requests directly to the load balancer IP address instead of using DNS. For example, this is useful if your application can't use DNS or if the connecting clients require firewall rules based on IP addresses. In this case, NLB is the right type of load balancer to use.

**NLB preserves source IP address.** NLB preserves the source IP address of the client when sending the traffic to the backend. With ALB, if you look at the source IP address of the requests, you will find the IP address of the load balancer. While with NLB, you would see the real IP address of the client, which is required by the backend application in some cases.