

# AWS Introduction

# AWS main components



Amazon EC2



Amazon ECR



Amazon ECS



AWS Elastic  
Beanstalk



AWS  
Lambda



Elastic Load  
Balancing



Amazon  
CloudFront



Amazon  
Kinesis



Amazon  
Route 53



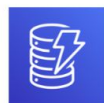
Amazon  
S3



Amazon  
RDS



Amazon  
Aurora



Amazon  
DynamoDB



Amazon  
ElastiCache



Amazon  
SQS



Amazon  
SNS



AWS Step Functions



Auto Scaling



Amazon API  
Gateway



Amazon  
SES



Amazon  
Cognito



IAM



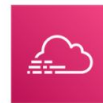
Amazon  
CloudWatch



Amazon EC2  
Systems Manager



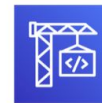
AWS  
CloudFormation



AWS  
CloudTrail



AWS  
CodeCommit



AWS  
CodeBuild



AWS  
CodeDeploy



AWS  
CodePipeline



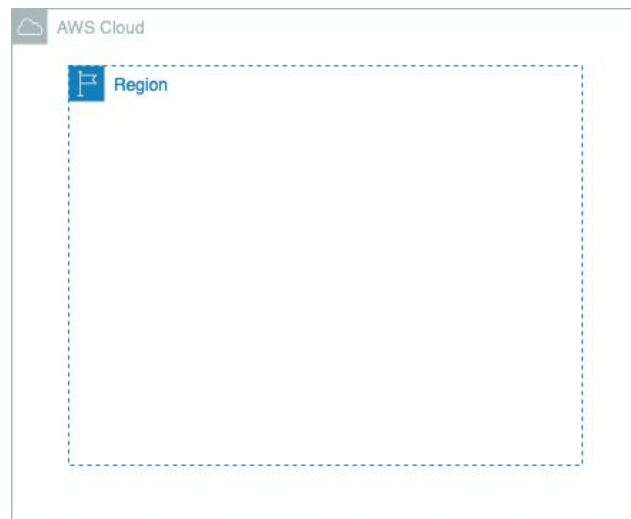
AWS  
X-Ray



AWS KMS

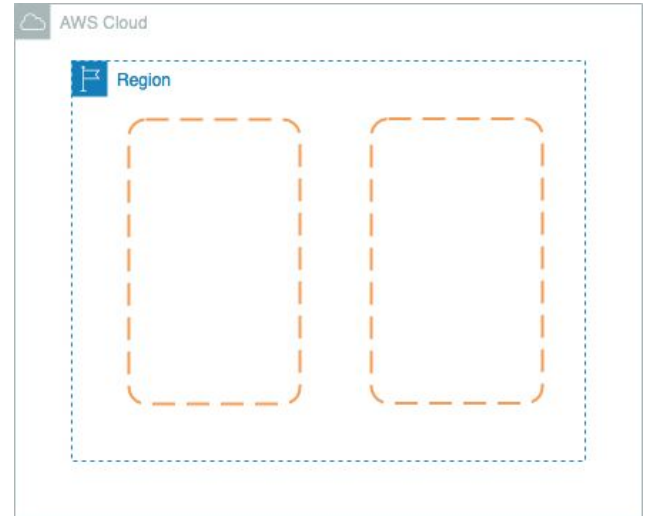
# AWS Regions

- Designed to be **isolated** from other Amazon Regions
- Achieve the greatest possible **fault tolerance** and **stability**
- Most AWS Resources are **tied** to the Regions except some Global Services like Identity and Access Management (IAM)
- For example, we may want to launch instances in the EU to be near European customers or to meet legal requirements

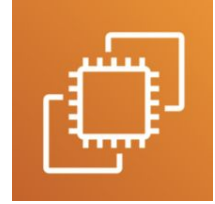


# AWS Availability Zones (AZ)

- Availability Zones are multiple, **isolated locations** within each **Region**
- Represented by a Region code followed by a letter identifier; for example, eu-central-1a
- Consist of one or more discrete **data centers**, each with redundant power, networking, and connectivity
- Offer the ability to operate applications that are more highly available, fault tolerant, and scalable




# Amazon EC2



- EC2 = Elastic Compute Cloud = Infrastructure as a Service
- You can use Amazon EC2 to **launch** as many or as few virtual servers as you need, configure security and networking, and manage storage
- Knowing EC2 is **fundamental** to understand how the Cloud works
- Operating System (**OS**): Linux, Windows or Mac OS
- How much compute power & cores (**CPU**)
- How much random-access memory (**RAM**)
- How much **storage** space
- **Network** card: speed of the card, Public IP address

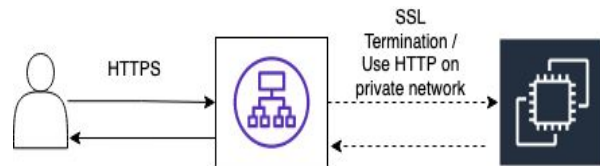
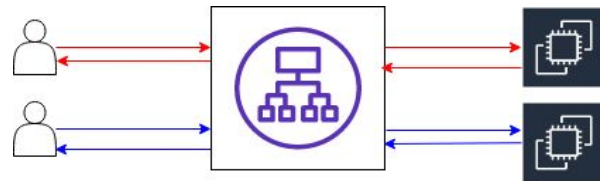
# EC2 Types

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

# AWS ELB



- An ELB (EC2 Load Balancer) is a **managed load balancer**
- AWS takes care of upgrades, maintenance
- **Spreads load** across multiple downstream instances
- Exposes a **single point of access** (DNS) to your application
- Does regular **health checks** to your instances
- **High availability** across zones
- **Separates** public traffic from private traffic
- Provide **SSL termination** (HTTPS) for your websites



# Types of load balancer on AWS

- **Classic Load Balancer** (v1 - old generation) – HTTP, HTTPS, TCP
- **Application Load Balancer** (v2 - new generation) – HTTP, HTTPS, WebSocket
- **Network Load Balancer** (v2 - new generation) – TCP, TLS & UDP

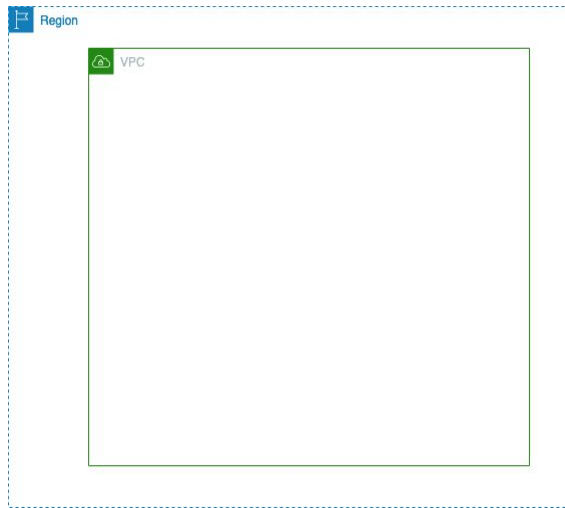
You can setup **internal** (private) or **external** (public) ELBs





# AWS VPC

- VPC = **Virtual Private Cloud** to hold all of our AWS resources
- Restricts what sort of traffic, IP addresses and also the users that can access our instances
- VPC is **private**, only the Private IP ranges are allowed (10.0.0.0–10.255.255.255 / 172.16.0.0–172.31.255.255 / 192.168.0.0–192.168.255.255)
- Up to **5 per region** – soft limit
- A VPC's CIDR (Classless Inter-Domain Routing) should **not overlap** with your other networks

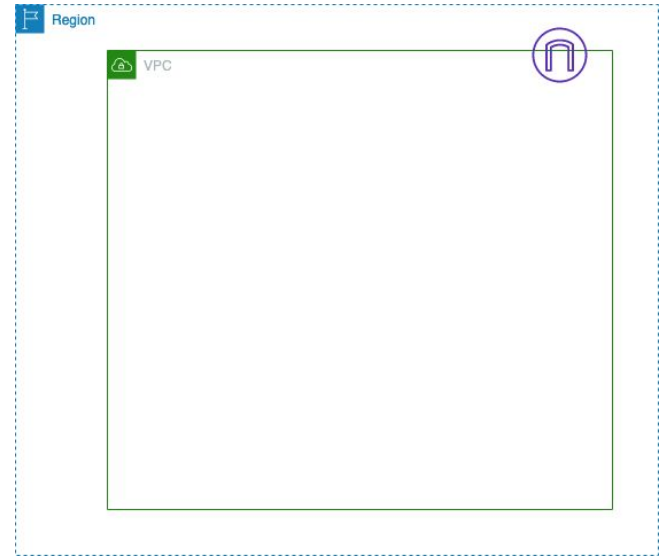


# AWS VPC Components

- **Subnet:** A segment of a VPC's IP address range where you can place groups of isolated resources
- **Internet Gateway:** The Amazon VPC side of a connection to the public Internet
- **NAT Gateway:** Highly available, managed service for resources in a private subnet to access the Internet
- **Virtual private gateway:** The Amazon VPC side of a VPN connection
- **Peering Connection:** Route traffic via private IP addresses between two peered VPCs
- **VPC Endpoints:** Enables private connectivity to services hosted in AWS, from within your VPC
- **Egress-only Internet Gateway:** A stateful gateway to provide egress only access for IPv6 traffic from the VPC to the Internet

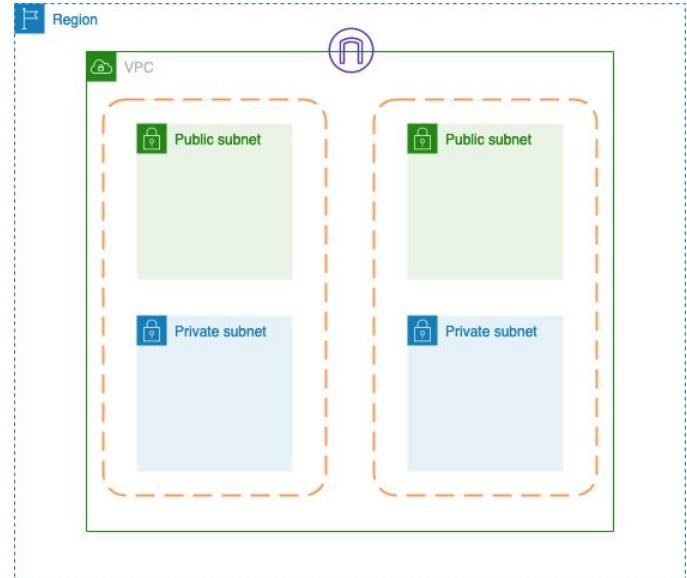
# AWS Internet Gateways (IG)

- VPC are in a private network -> Can **not reach** internet
- IG **helps** our VPC instances connect with the **internet**
- Managed by AWS, scales **horizontally** and is **HA**
- One VPC can only be **attached** to **one** IGW and vice versa



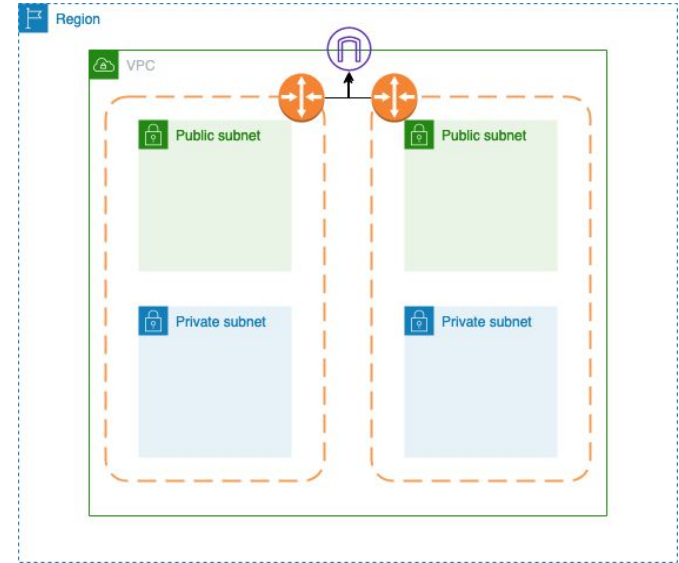
# AWS Subnets

- Are containers within VPC that **segment** off a slice of the **CIDR block** you define in your VPC
- Subnets allow you to give **different access rules** and place resources in different containers where those rules should apply
- Is a Availability Zone resource
- Can be **public** (accessible from the internet) or **private** (**not** accessible from the internet)



# AWS Route Tables

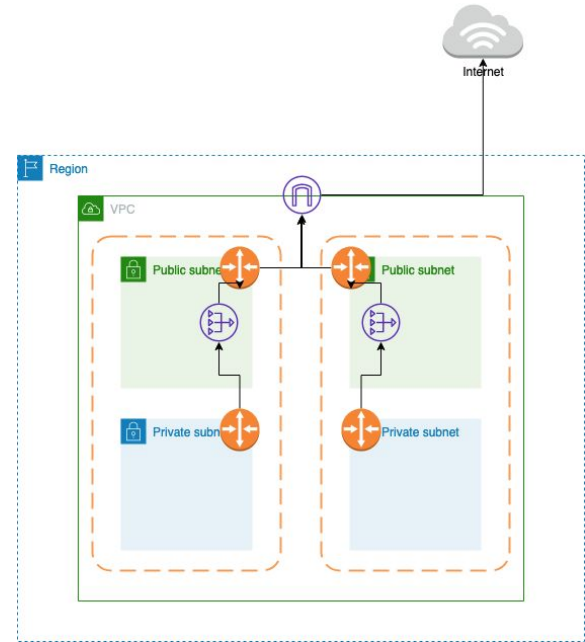
- Contains a **set of rules**, called routes, that are used to **determine** where network **traffic** from your subnet or gateway is directed
- Each **subnet** in your VPC must be **associated** with a route table, which controls the routing for the subnet (subnet route table)
- Each route in a table specifies a **destination** and a **target**
- For example, to enable a subnet to access the internet through an internet gateway, we can use the route table entry from the second image



Destination	▼	Target
192.168.0.0/16		local
0.0.0.0/0		<a href="#">igw-087fc142cabecf6d8</a>

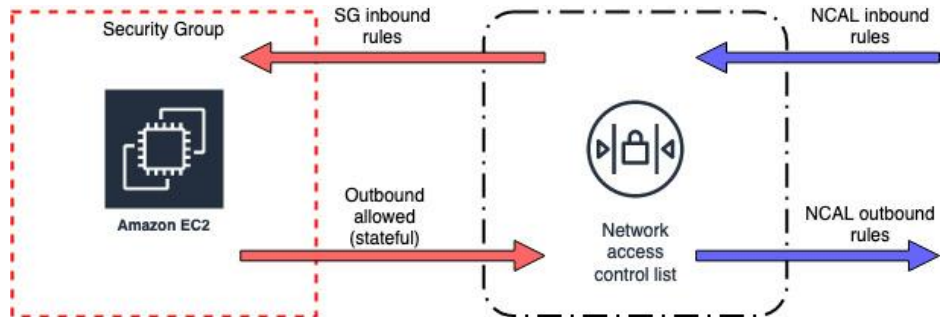
# AWS NAT Gateway

- Allows instances in the **private** subnets to **connect** to the internet.
- Must be launched in a **public** subnet.
- **Managed** by AWS
- NAT is created in a **specific** AZ, uses an EIP
- 1 NAT per AZ to have **fault-tolerance** and **HA** (High Availability)
- Requires an IGW (**Private Subnet => NAT => IGW**)



# Network ACLs

- NACL are like a **firewall** which control **traffic from and to subnet**
- Are placed on **subnet level**
- **Default** NACL allows everything outbound and everything inbound
- **One** NACL per Subnet
- **Deny** and **Allow** rules
- **Stateless**



# AWS Security Groups

- They **control** how **traffic** is allowed into or out of our EC2 Instances.
- Security groups **only** contain **rules**
- Security groups rules can reference by **IP** or by **security group**
- **Stateful**: Changes in **incoming rules** applied to **outgoing rules**

Security group ID ▾	Security group na... ▾	VPC ID ▾	Description ▾
<a href="#">sg-01c6514886e02738b</a>	alb-allow-http	<a href="#">vpc-03222a142e126c781</a> <a href="#">🔗</a>	Allow HTTP access to ALB
<a href="#">sg-0a5f80739ed40184b</a>	ec2-sg	<a href="#">vpc-03222a142e126c781</a> <a href="#">🔗</a>	Allow access to ec2 only from ALB
<a href="#">sg-0d964290f4b09a593</a>	default	<a href="#">vpc-03222a142e126c781</a> <a href="#">🔗</a>	default VPC security group

Type	Protocol	Port range	Source	Description - optional
HTTP	TCP	80	0.0.0.0/0	Allow access http on port 80
HTTP	TCP	80	::/0	Allow access http on port 80

Type	Protocol	Port range	Source	Description - optional
All TCP	TCP	0 - 65535	sg-01c6514886e02738b / alb-allow-http	–



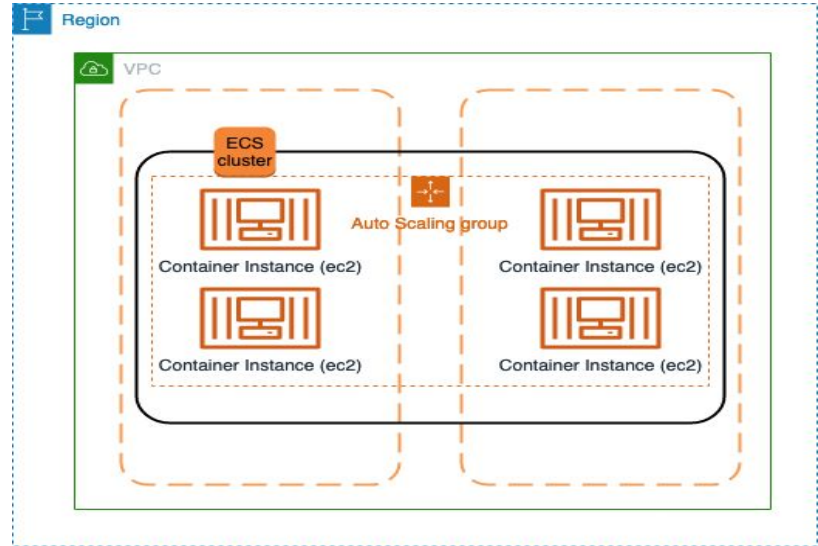
# AWS NACLs vs SG

Security Group	NACL
Instance level	Subnet level
Stateful	Stateless
Allow rules only	Allow and Deny rules
All rules are evaluated before traffic is allowed	Rules are evaluated in the order specified
First layer of defense for egress traffic	First layer of defense for ingress traffic

# AWS ECS



- ECS = Elastic Container Service
- Launch **Docker** containers on AWS
- Simplifies running containers in a HA manner across multiple Availability Zones within a Region
- Serverless with AWS Fargate



# ECS Cluster

- Is Region specific
- Is a logical grouping of tasks and services
- Uses one or more EC2 Instances to run tasks
- EC2 instances of the cluster run the ECS agent
- The ECS agent registers the instance to the Cluster
- Serverless using AWS Fargate

# ECS Task Definition

- A **JSON** file that **describes** one or more containers for ECS to run
- Can be thought of as a **blueprint for your application**
- **Docker image** to use with each container in your task
- **CPU** and **memory** to use with each task
- Which **ports** should be opened for your application
- What **data volumes** should be used with the containers in the task

# ECS Services

- Allows to **run and maintain** a specified number of tasks
- If any of the tasks **fails**, ECS launches another task in order to maintain the desired number of tasks in the service
- **Task placement strategies** and constraints to customise task placement decisions
- Three **deployment types**: rolling update, blue/green, and external
- Can be linked to an **ELB** (Load Balancer)

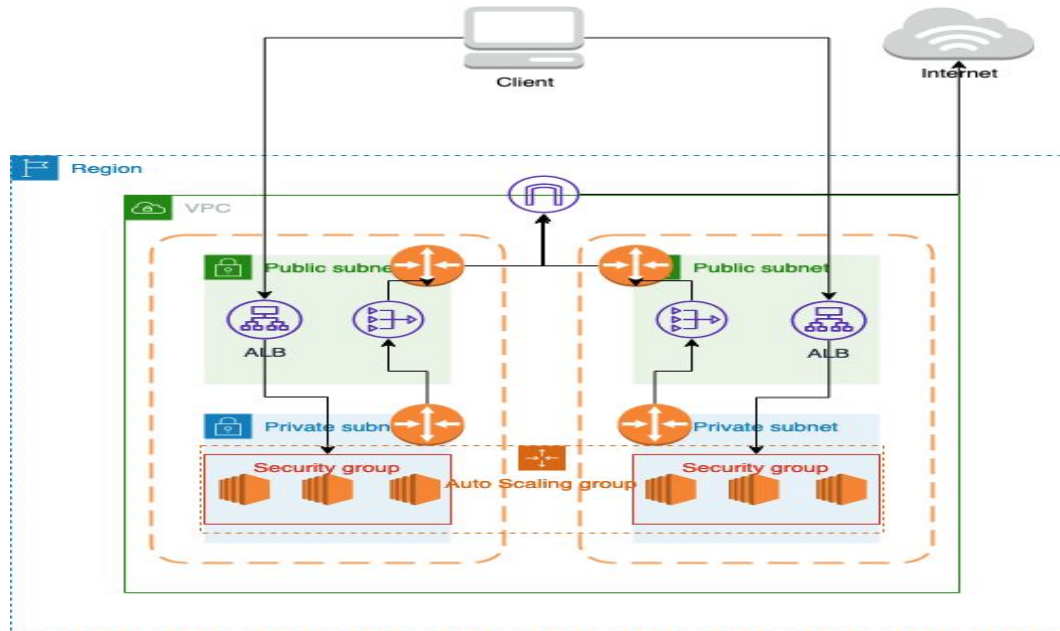
# Terraform

- **Infrastructure as Code** (described using a high-level configuration syntax)
- Is a tool for **building, changing, and versioning** infrastructure safely and efficiently
- Configuration files describe to Terraform the components needed to run
- Generates an **execution plan** describing what it will do to reach the desired state
- Executes the plan to build the described infrastructure
- Determines what changed and creates **incremental** execution plans
- Can manage low-level components (compute instances, networking), as well as high-level components (DNS entries, SaaS features)

# Terraform hands-on

```
main.tf x
devops > aws > ecs_fargate > main.tf
1 provider "aws" {
2   shared_credentials_file = "SHOME/.aws/credentials"
3   profile                 = "default"
4   region                 = var.aws_region
5 }
6
7 module "networking" {
8   source           = "../modules/network"
9   create_vpc       = var.create_vpc
10  create_igw        = var.create_igw
11  single_nat_gateway = var.single_nat_gateway
12  enable_nat_gateway = var.enable_nat_gateway
13  region            = var.aws_region
14  vpc_name           = var.vpc_name
15  cidr_block         = var.cidr_block
16  availability_zones = var.availability_zones
17  public_subnet_cidrs = var.public_subnet_cidrs
18  private_subnet_cidrs = var.private_subnet_cidrs
19 }
20
21 #####
22 # ECS Tasks Execution IAM
23 #####
24 # ECS task execution role data
25 data "aws_iam_policy_document" "ecs_task_execution_role" {
26   version = "2012-10-17"
27   statement {
28     sid     = ""
29     effect  = "Allow"
30     actions = ["sts:AssumeRole"]
31
32     principals {
33       type       = "Service"
34       identifiers = ["ecs-tasks.amazonaws.com"]
35     }
36   }
37 }
```

# AWS Monolith API design with fault-tolerance and HA





Thank you  
Alexandros and Dimosthenis



arconsis