# AWS Solutions Architect Associate

## Session 601

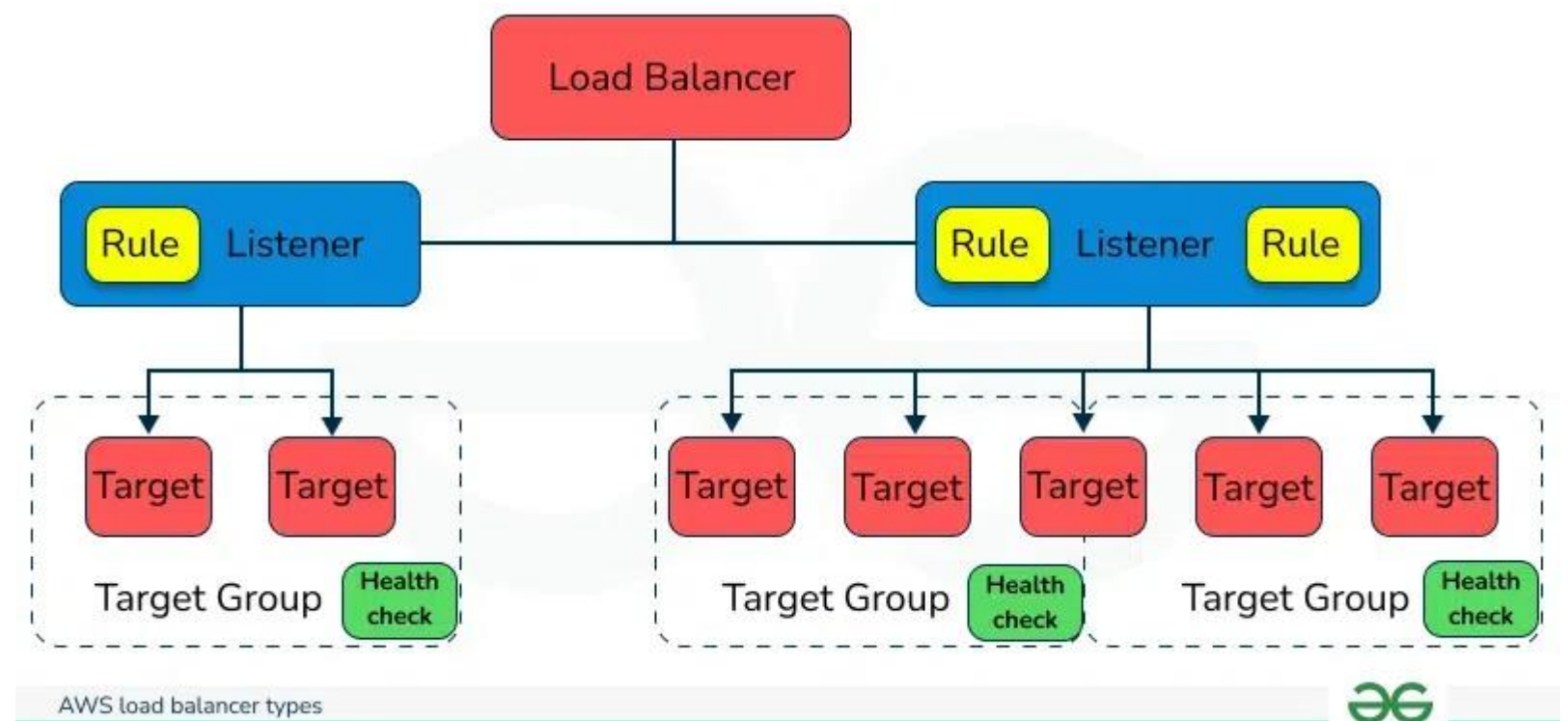## Compute: ELB and Mgmt &Gov: Cloudwatch

July/2024

Elastic Load
Balancing

A managed load balancing service that distributes incoming application traffic across multiple Amazon EC2 instances, containers, IP addresses and Lambda Functions.
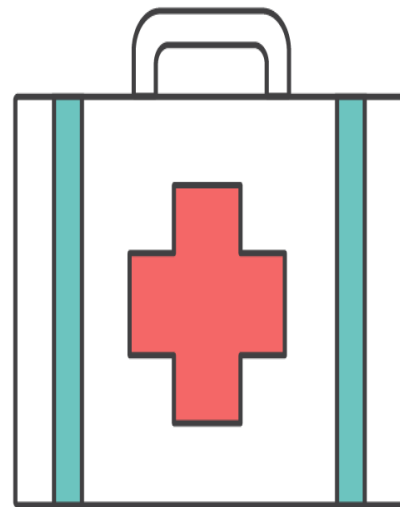


AWS load balancer types

Elastic Load Balancing

- Uses HTTP, HTTPS, TCP and SSL (secure TCP) protocols.

- Can be external or internal facing

- Each load balancer is given a DNS name

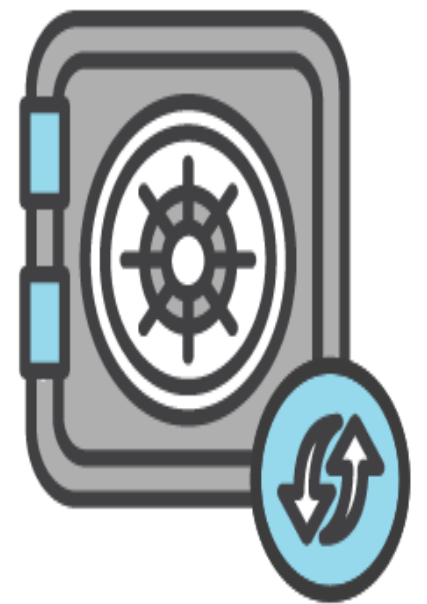- Recognizes and responds to healthy instances
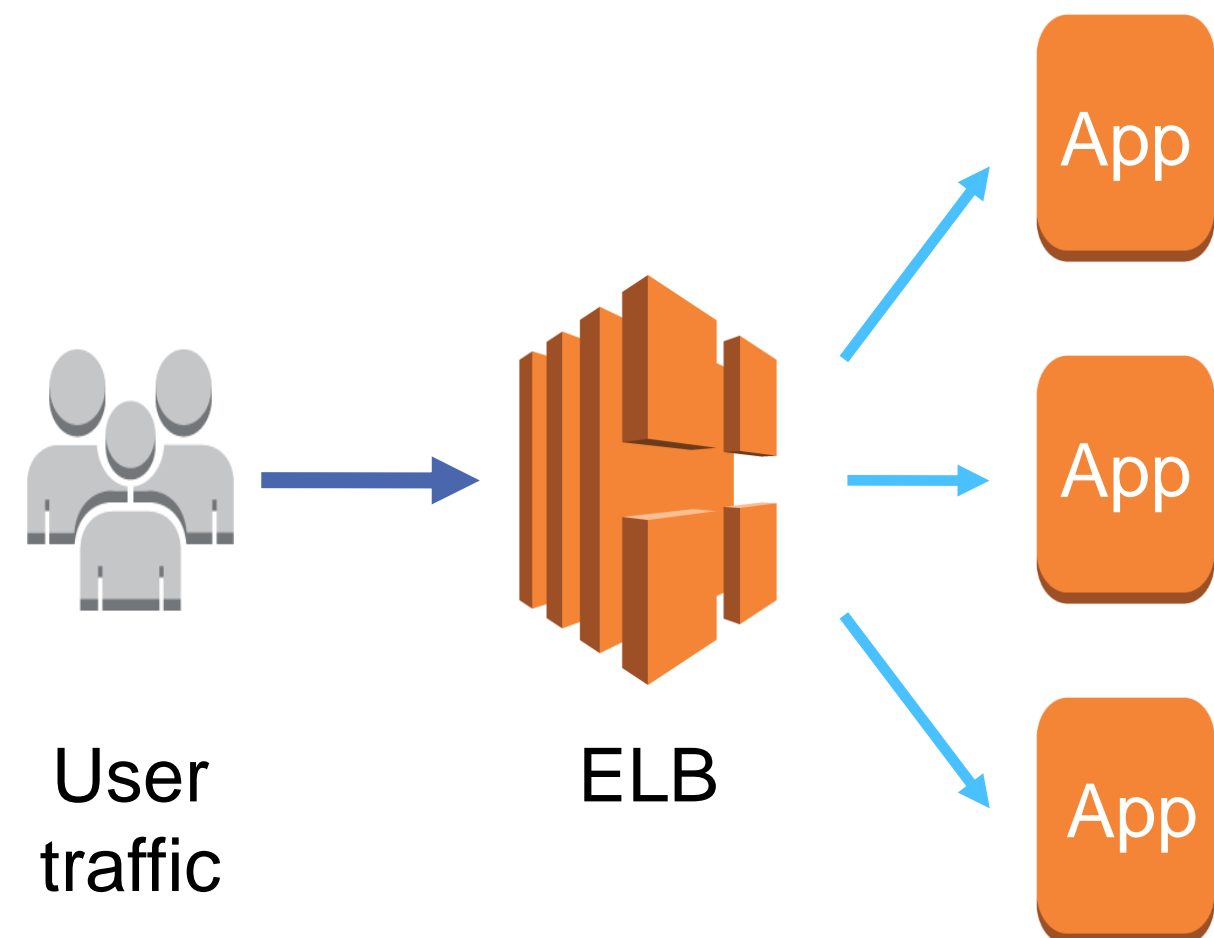
High
availability

Health
checks

Security
features

TLS
termination

Protocol,
port, path

If you need to remove an instance from your production fleet, but don't want to affect your users:
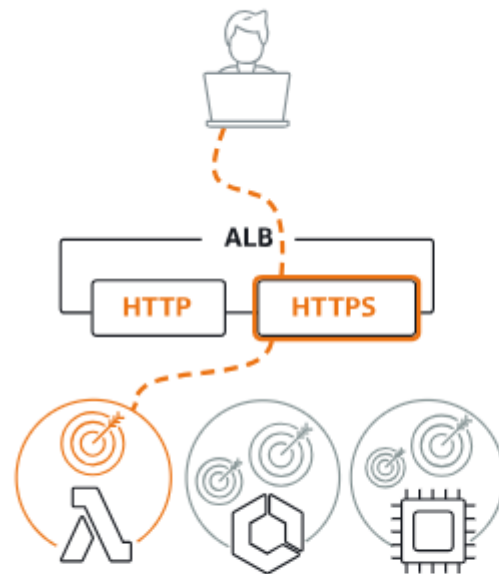
Affected backend instances will complete requests in progress before deregistration

User traffic
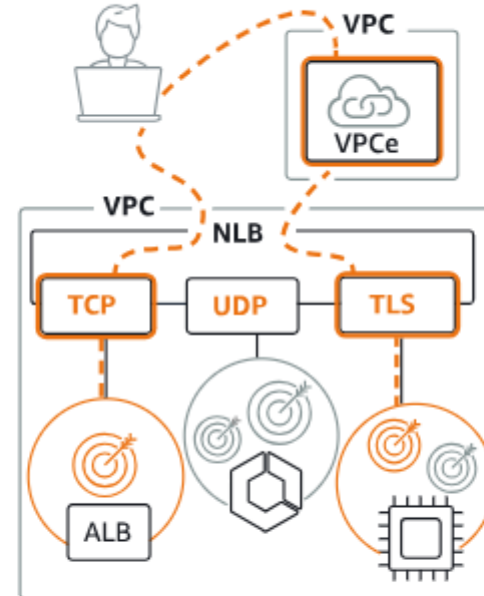
ELB

App

App

App

**Enable connection draining**

## Application Load Balancer Info



Choose an Application Load Balancer when you need a flexible feature set for your applications with HTTP and HTTPS traffic. Operating at the request level, Application Load Balancers provide advanced routing and visibility features targeted at application architectures, including microservices and containers.

Create

## Network Load Balancer Info



Choose a Network Load Balancer when you need ultra-high performance, TLS offloading at scale, centralized certificate deployment, support for UDP, and static IP addresses for your applications. Operating at the connection level, Network Load Balancers are capable of handling millions of requests per second securely while maintaining ultra-low latencies.
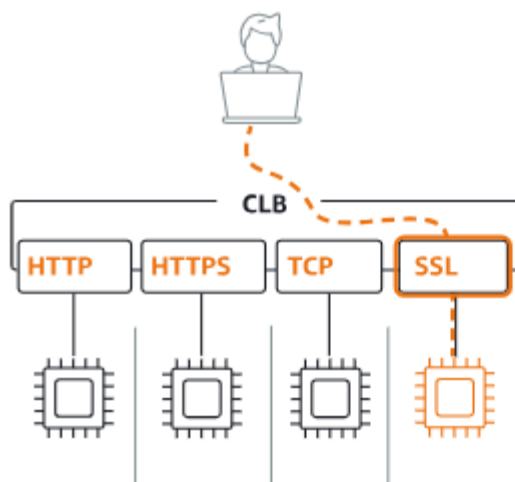
## Gateway Load Balancer Info



Choose a Gateway Load Balancer when you need to deploy and manage a fleet of third-party virtual appliances that support GENEVE. These appliances enable you to improve security, compliance, and policy controls.

Create

▼ Classic Load Balancer - *previous generation*

## Classic Load Balancer Info



Choose a Classic Load Balancer when you have an existing application running in the EC2-Classic network.

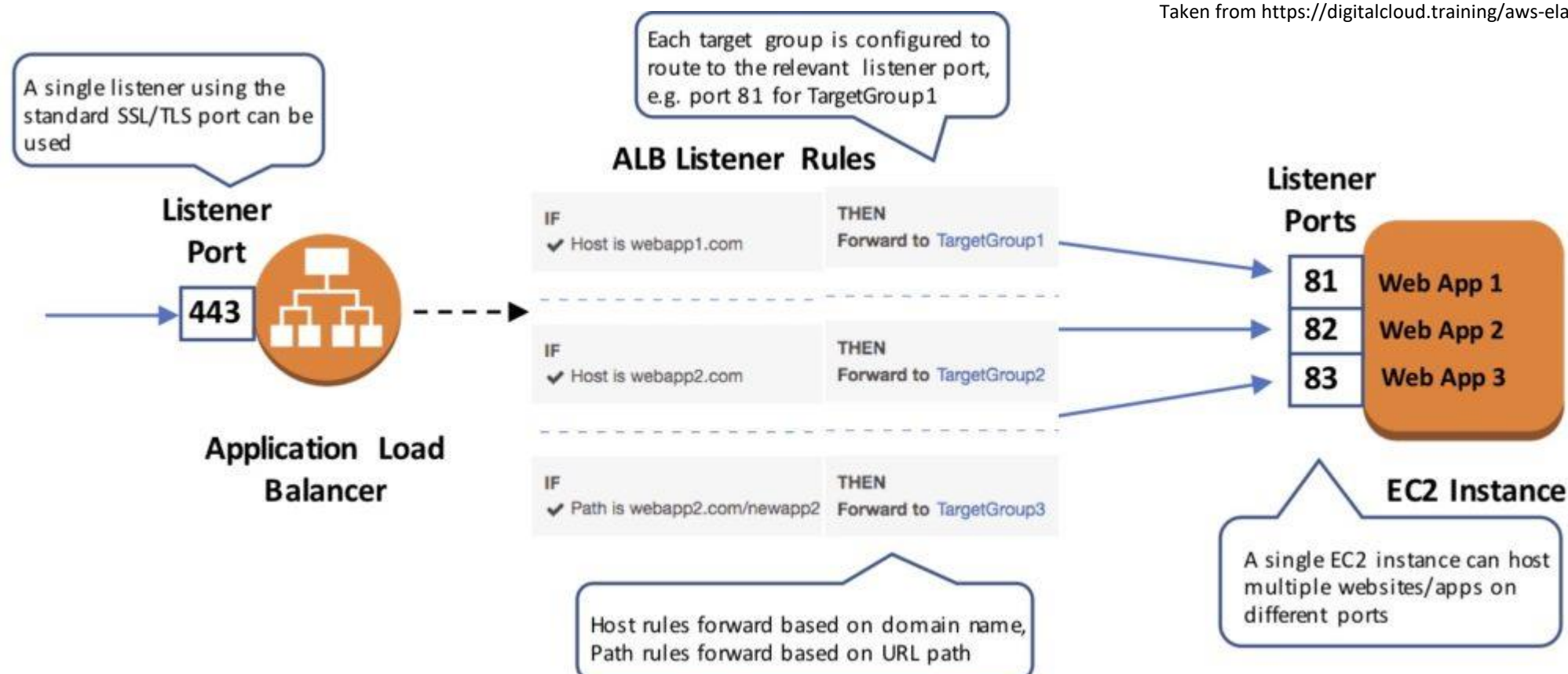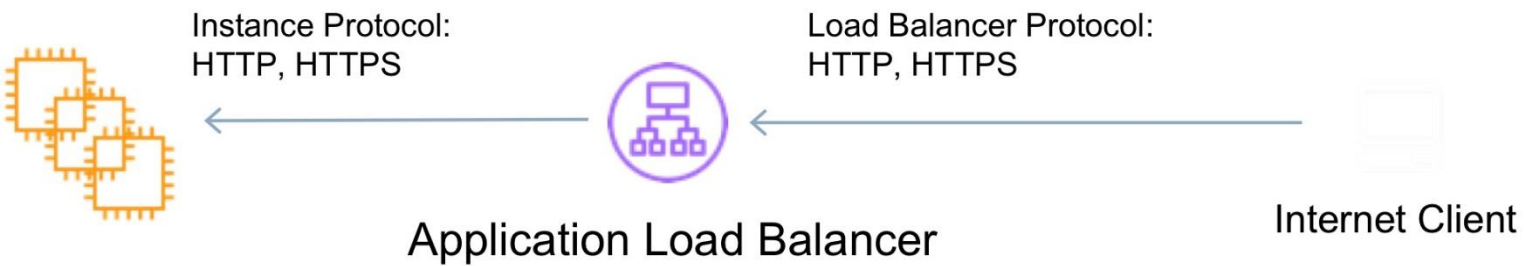Create

Taken from AWS Console (18/07/2024)

I: Target Groups: Type (Instance/IP), Port, Target Registration

II: Load Balancer: Type and Listener.

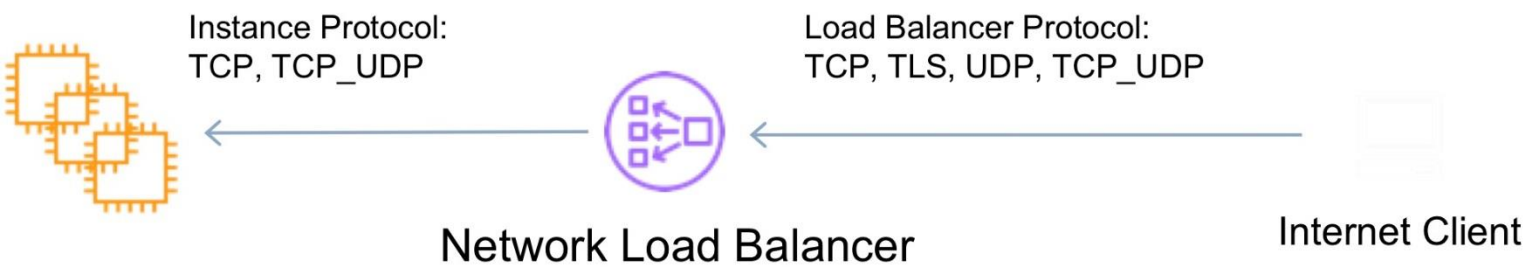**Note:** Can be done using Infrastructure or Kubernetes Commands.

Taken from https://digitalcloud.training/aws-elastic-load-balancing-aws-elb/  (18/07/2024)

Instance Protocol:
HTTP, HTTPS

Load Balancer Protocol:
HTTP, HTTPS

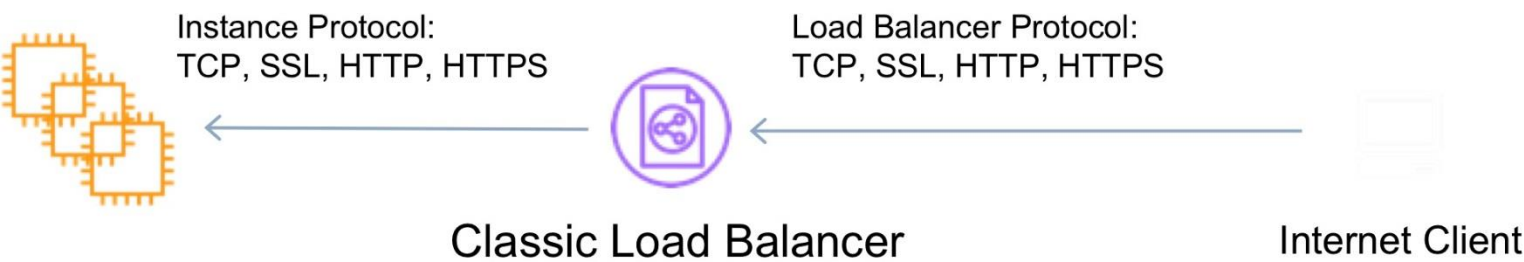Application Load Balancer

Internet Client

## Application Load Balancer

- Operates at the request level
- Routes based on the content of the request (layer 7)
- Supports path-based routing, host-based routing, query string parameter-based routing, and source IP address-based routing
- Supports IP addresses, Lambda Functions and containers as targets

Instance Protocol:
TCP, TCP_UDP

Load Balancer Protocol:
TCP, TLS, UDP, TCP_UDP

Network Load Balancer

Internet Client

## Network Load Balancer

- Operates at the connection level
- Routes connections based on IP protocol data (layer 4)
- Offers ultra high performance, low latency and TLS offloading at scale
- Can have static IP / Elastic IP
- Supports UDP and static IP addresses as targets

Instance Protocol:
TCP, SSL, HTTP, HTTPS

Load Balancer Protocol:
TCP, SSL, HTTP, HTTPS

Classic Load Balancer

Internet Client

## Classic Load Balancer

- Old generation; not recommended for new applications
- Performs routing at Layer 4 and Layer 7
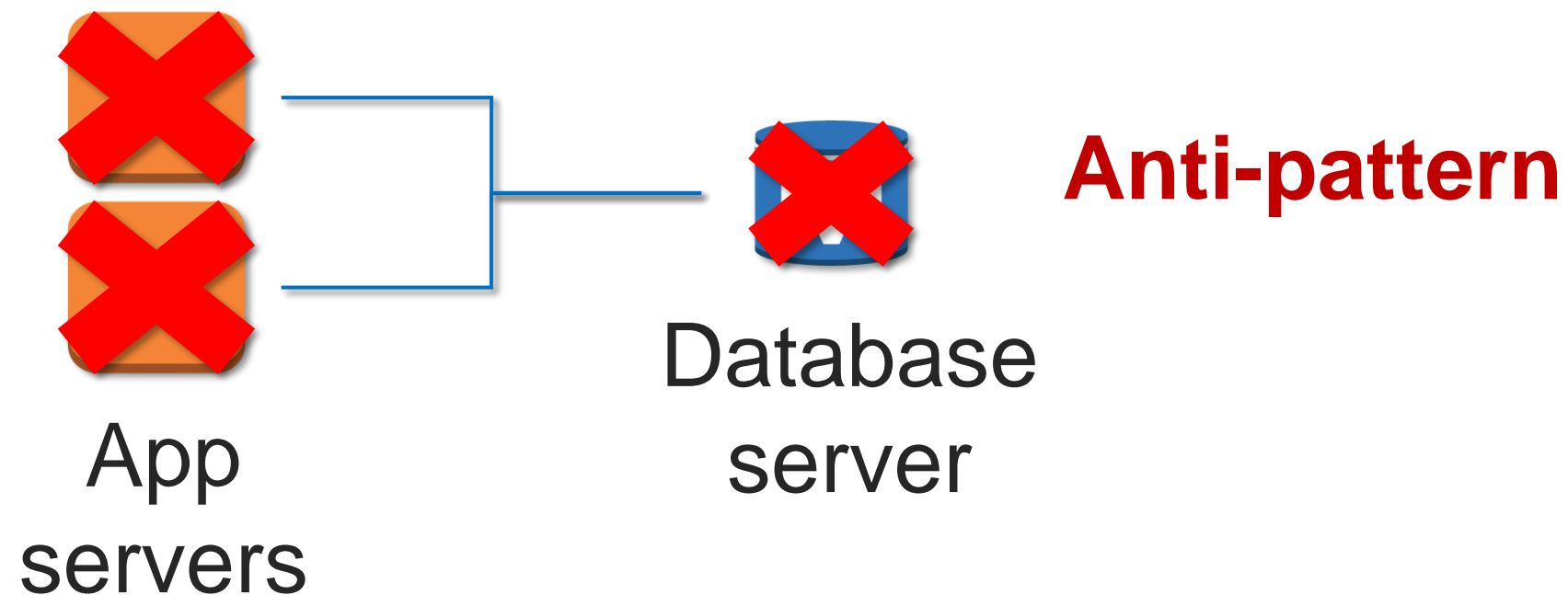- Use for existing applications running in EC2-Classic

Taken from https://digitalcloud.training/certification-training/aws-solutions-architect-associate/compute/elastic-load-balancing/ (06/07/2020)

Your application can recover from a failure or roll over to a secondary source within an acceptable amount of degraded performance time.

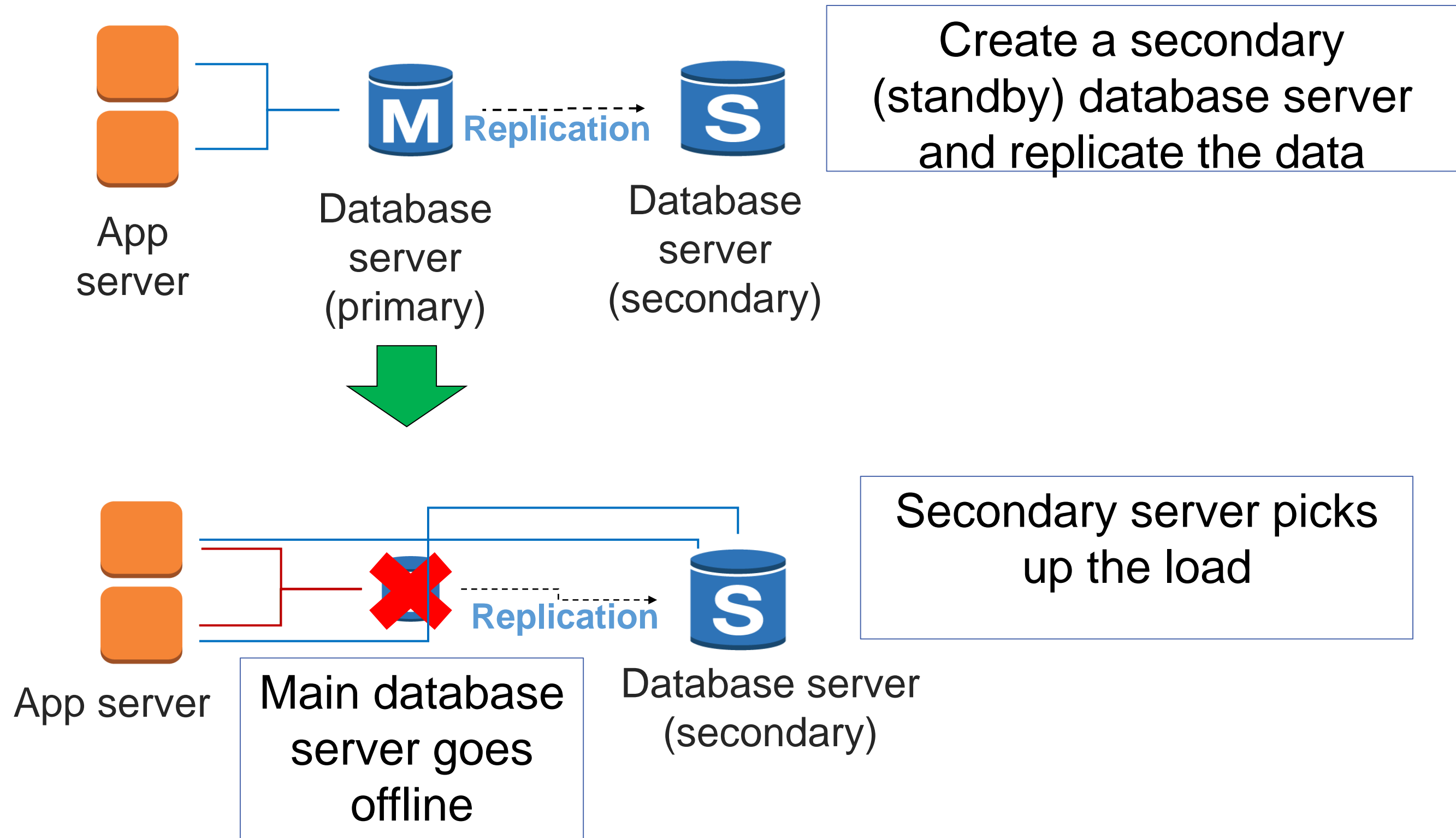| Percent of Uptime | Max Downtime per Year | Equivalent Downtime per Day |
| --- | --- | --- |
| 90% | 36.5 days | 2.4 hrs |
| 99% | 3.65 days | 14 min |
| 99.9% | 8.76 hrs | 86 sec |
| 99.99% | 52.6 min | 8.6 sec |
| 99.999% | 5.25 min | .86 sec |

*Assume everything fails, and design backward*

Implement redundancy where possible in order to prevent single failures from bringing down an entire system.



**Anti-pattern**

App servers

Database server

**Best practice**

App
server

Database
server
(primary)

**Replication**

Database
server
(secondary)

Create a secondary
(standby) database server
and replicate the data

App server

Main database
server goes
offline

**Replication**

Database server
(secondary)

Secondary server picks
up the load
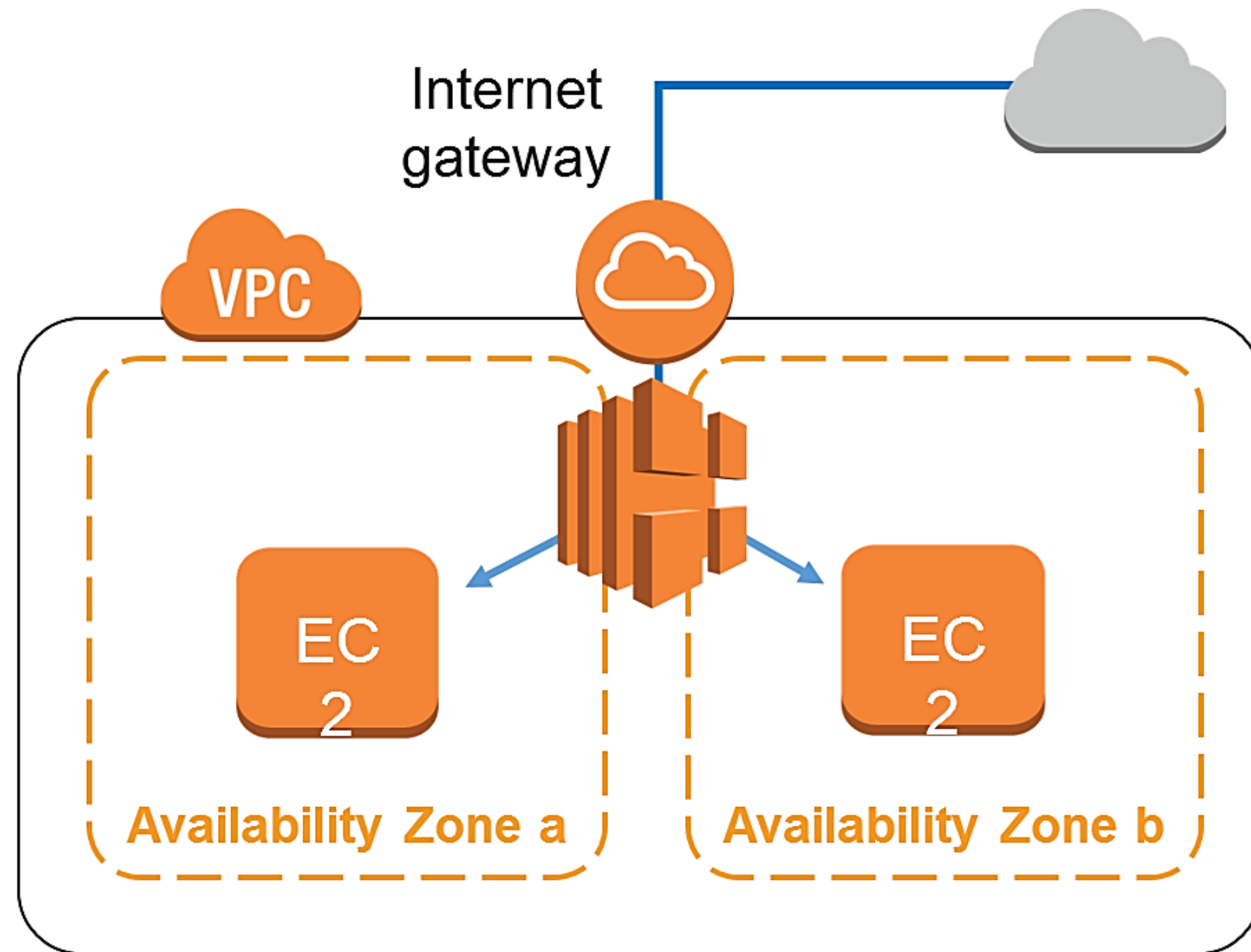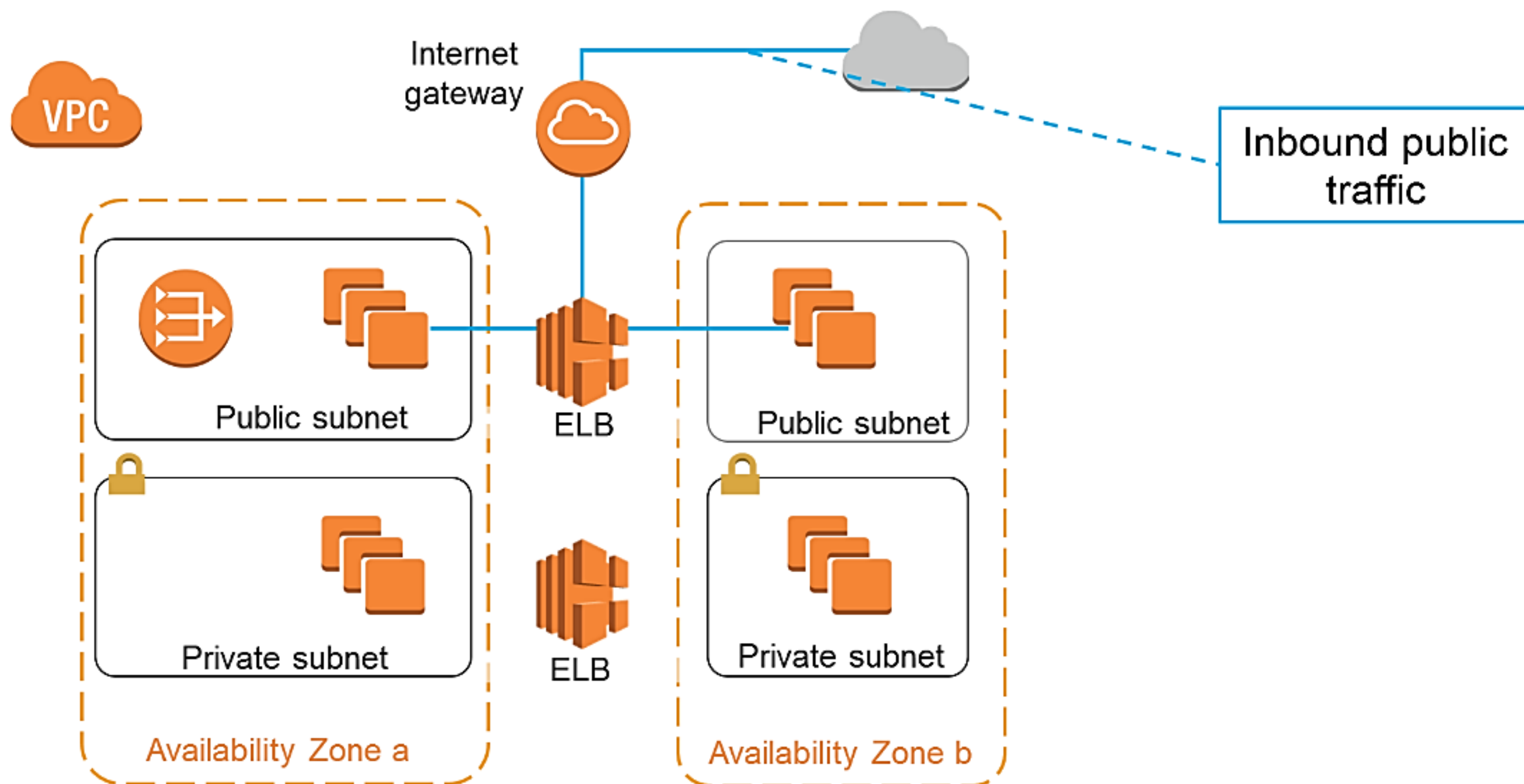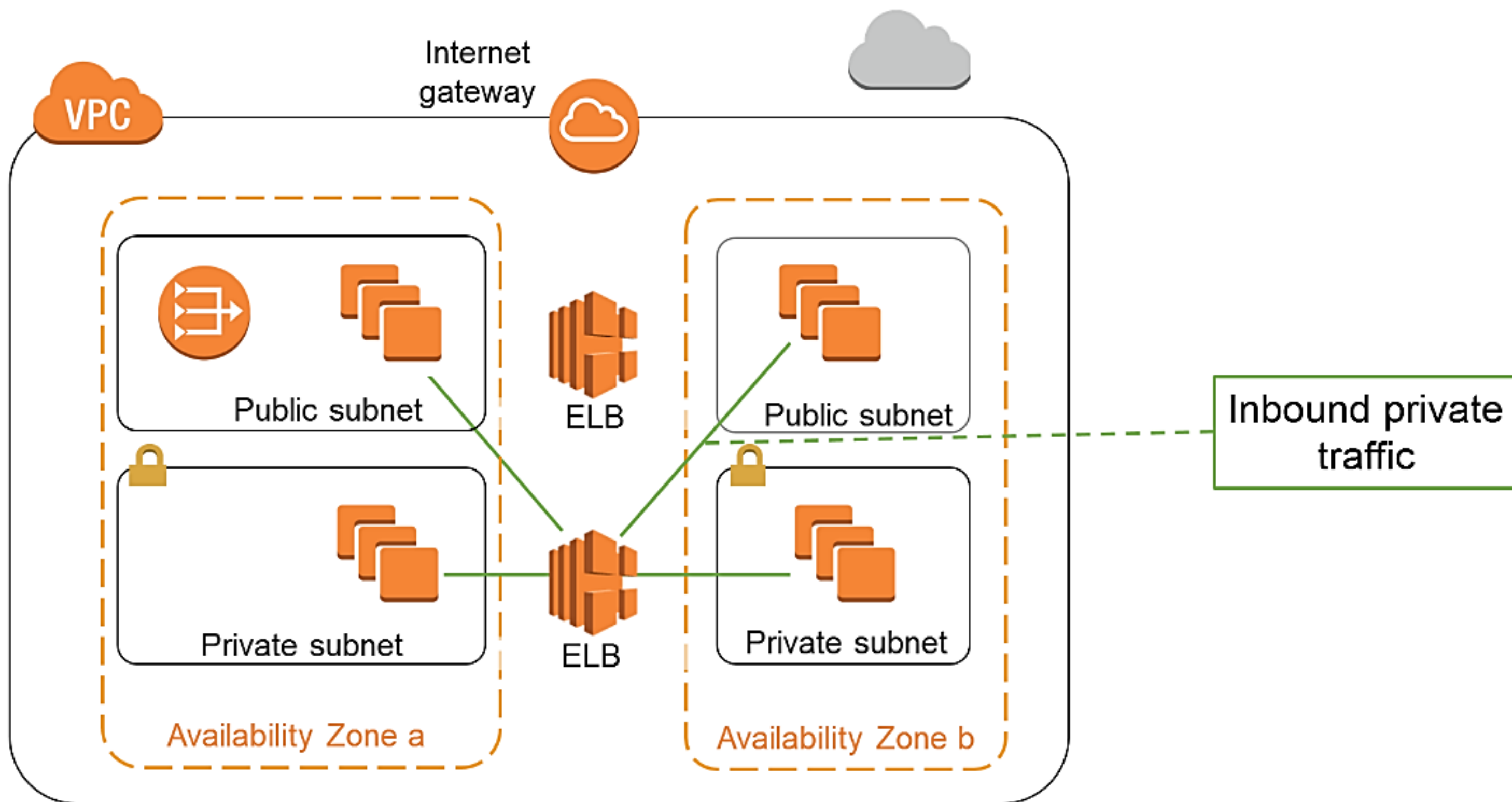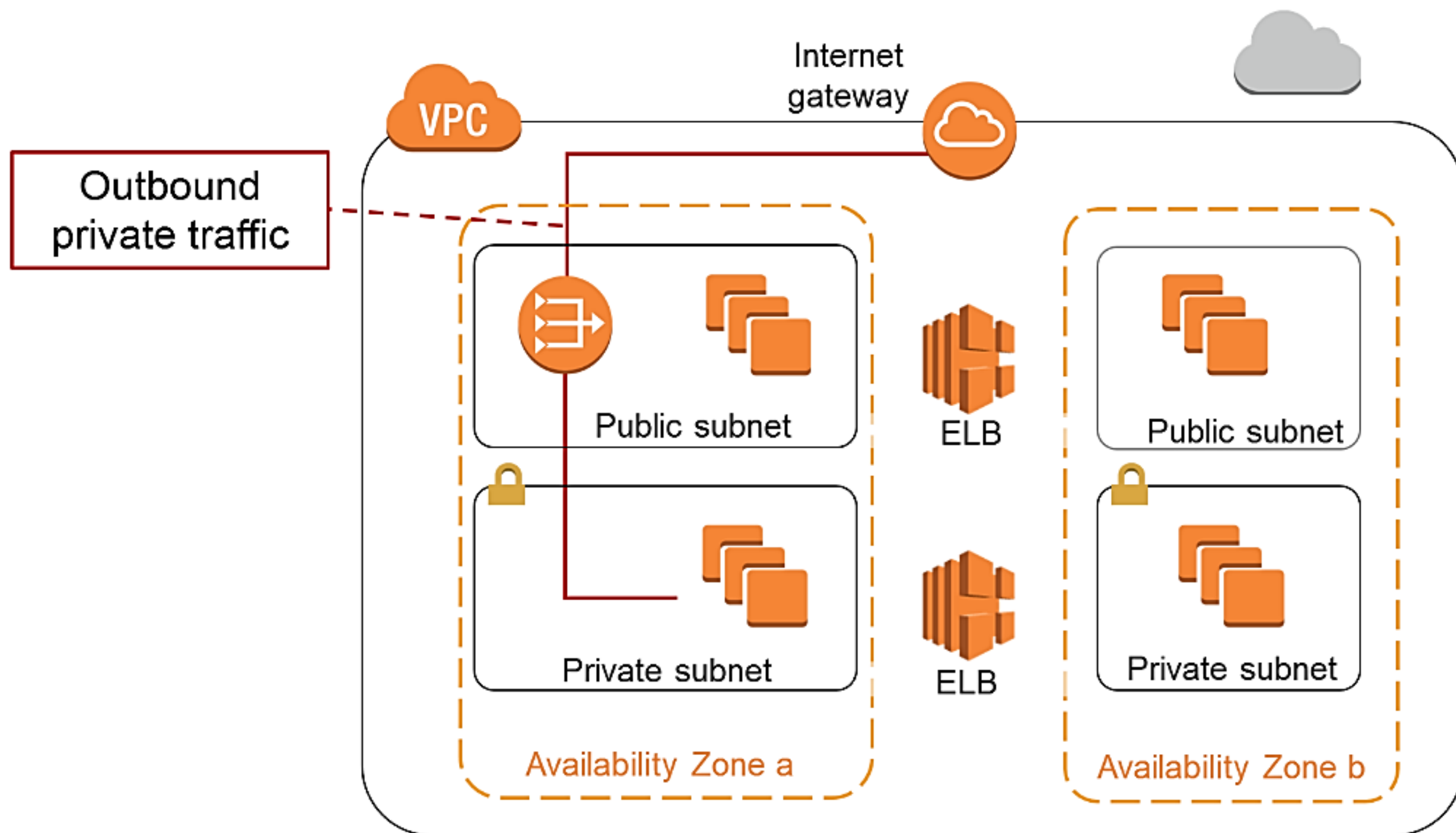
Start with two Availability Zones per AWS Region.

If resources in one Availability Zone are unreachable, your application shouldn't fail.

Internet gateway

VPC

EC 2

EC 2

Availability Zone a
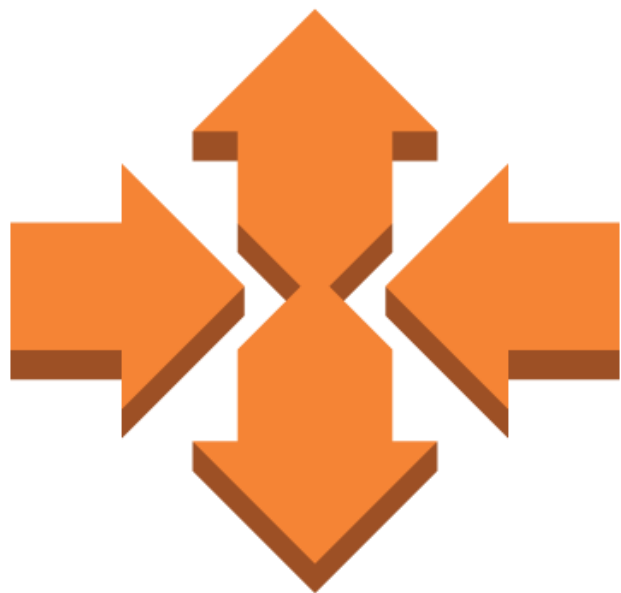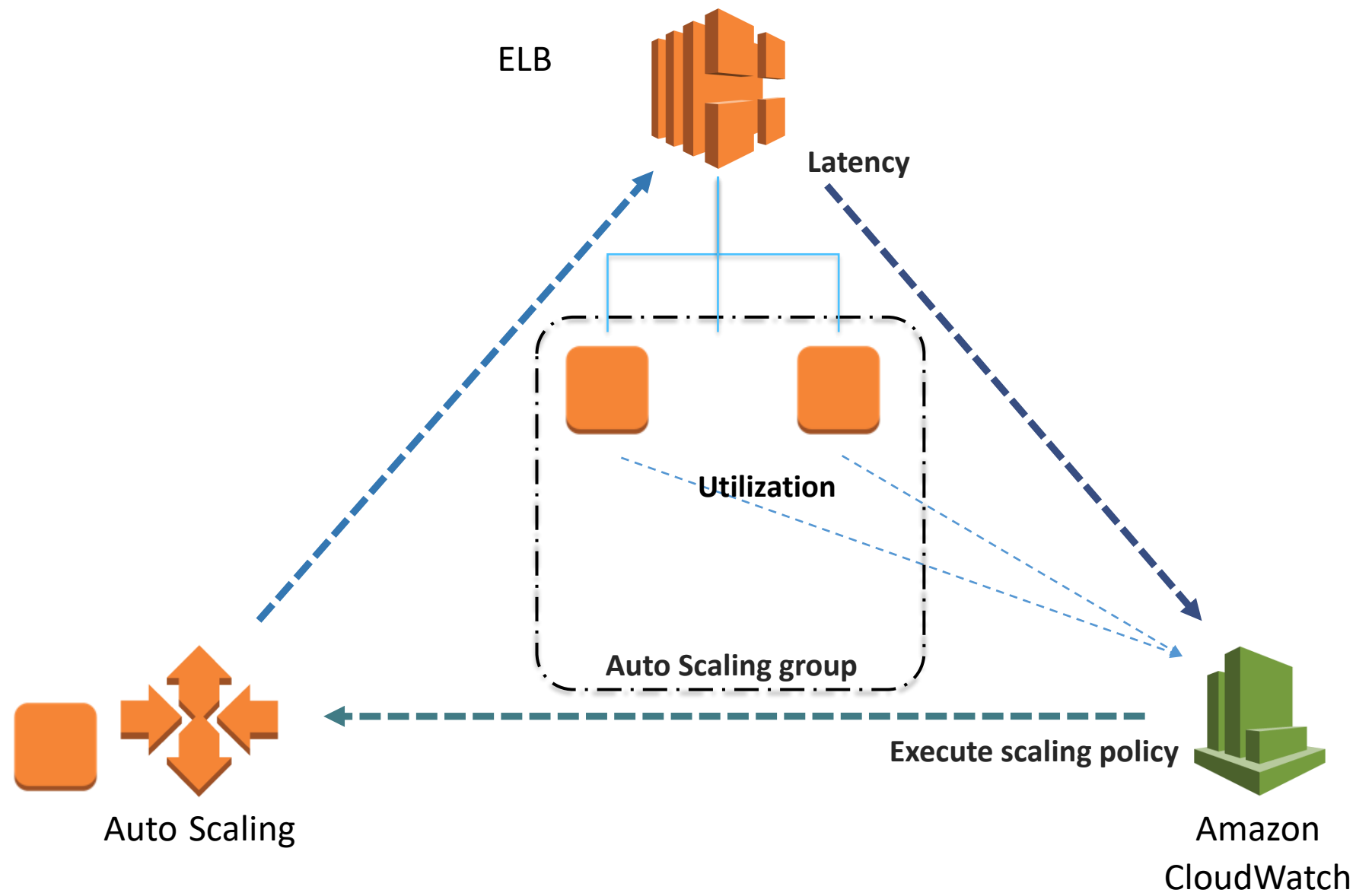
Availability Zone b

Auto Scaling

- Launches or terminates instances based on specified conditions

- Automatically registers new instances with load balancers when specified

- Can launch across Availability Zones

- Additional AutoScaling instead of EC2: Application AS (ECS, EMR, DynamoDB), AWS AutoScaling (Managed Both).
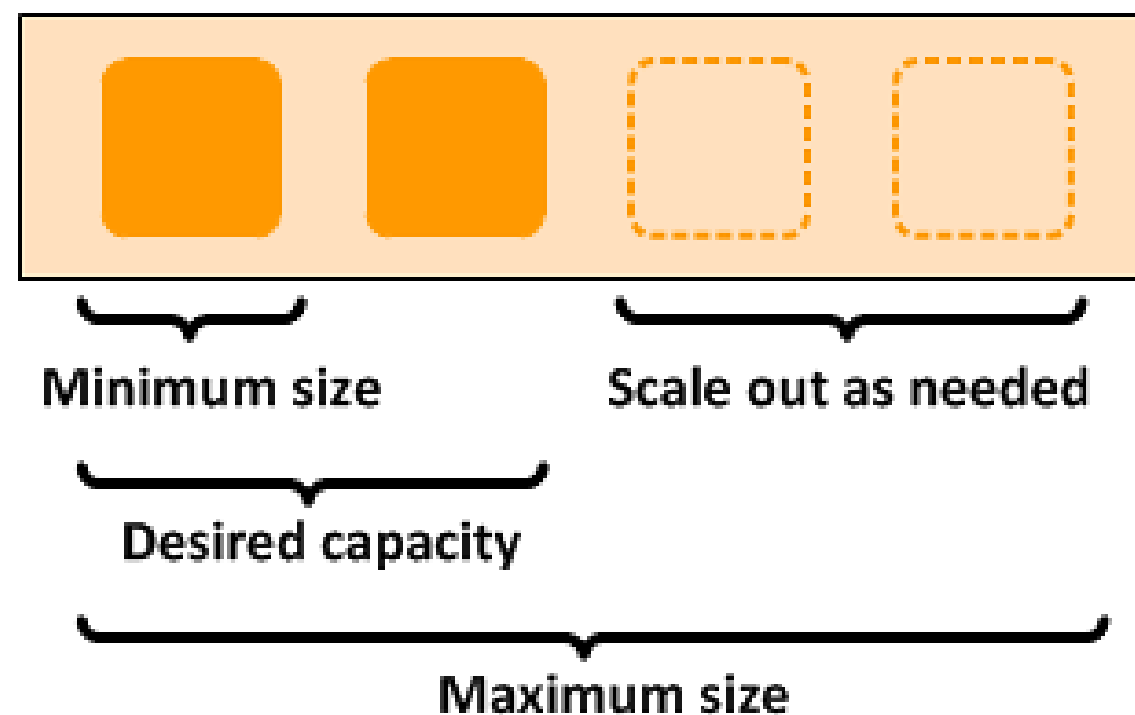
- Desired capacity

- Minimum capacity

- Maximum capacity

What would be a good **minimum** capacity to set it to?
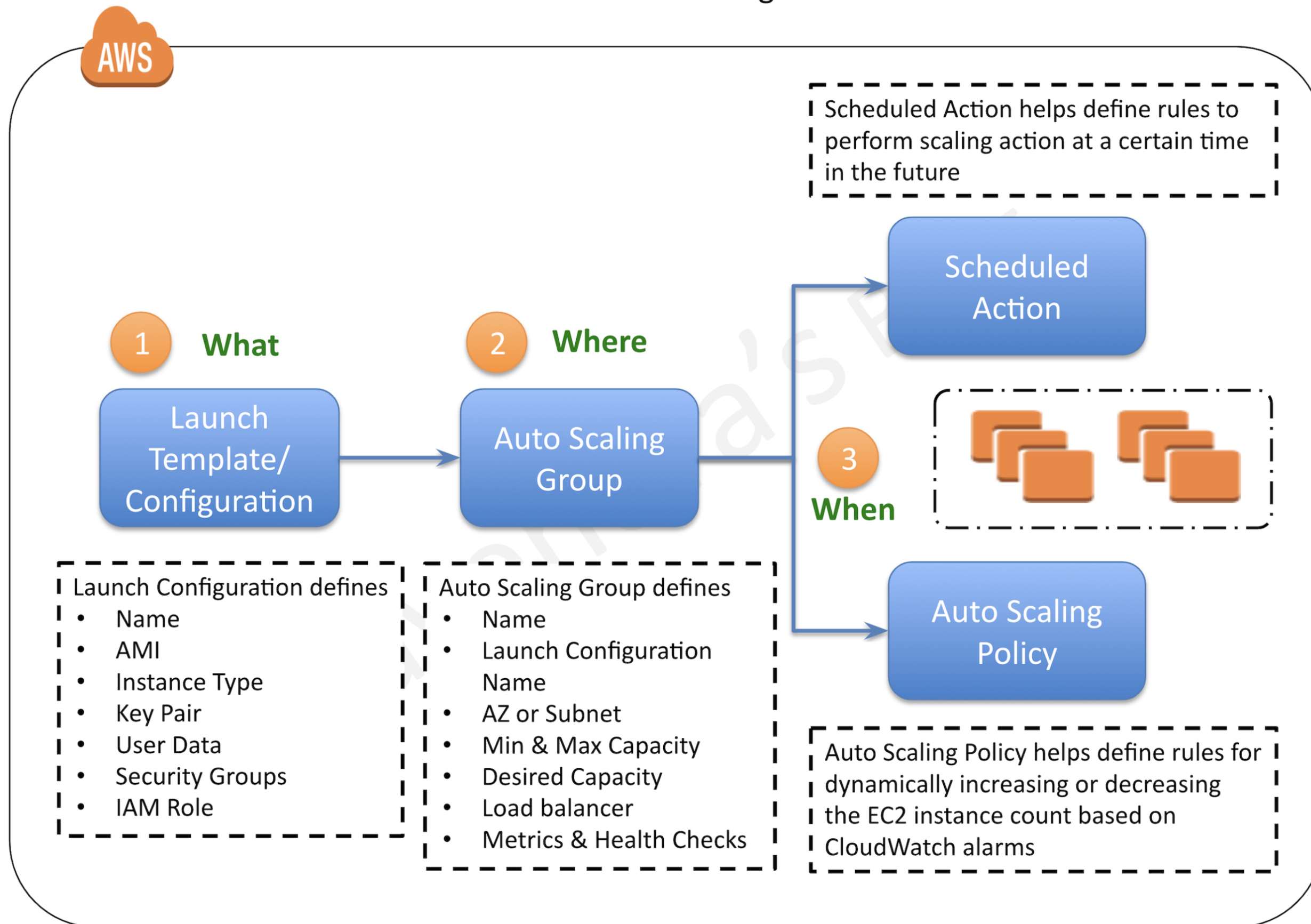
What would be a good **maximum** number?

By default, **Desired** (# instances all the time) is equal to **minimum**.

**?**

**Auto Scaling group**

**Availability Zone 1**

Minimum size

Desired capacity

Scale out as needed

Maximum size

## AWS Auto Scaling

**AWS**

Scheduled Action helps define rules to perform scaling action at a certain time in the future

**1** **What**

**2** **Where**

Launch Template/ Configuration → Auto Scaling Group → Scheduled Action

**3**

**When**

Auto Scaling Policy

Launch Configuration defines
- Name
- AMI
- Instance Type
- Key Pair
- User Data
- Security Groups
- IAM Role

Auto Scaling Group defines
- Name
- Launch Configuration Name
- AZ or Subnet
- Min & Max Capacity
- Desired Capacity
- Load balancer
- Metrics & Health Checks

Auto Scaling Policy helps define rules for dynamically increasing or decreasing the EC2 instance count based on CloudWatch alarms

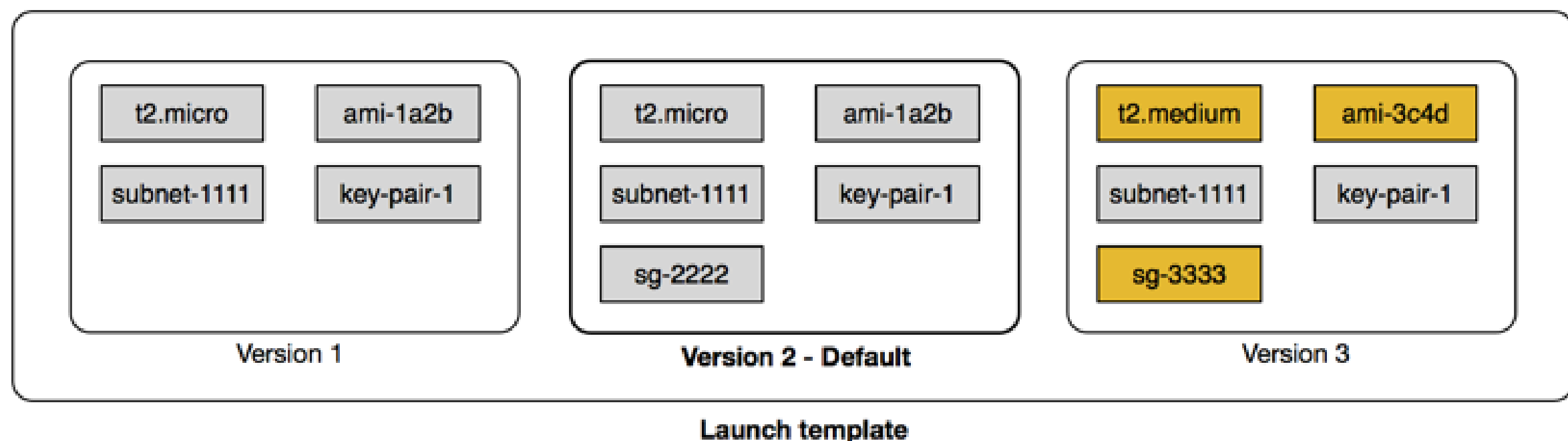Taken from https://jayendrapatil.com/aws-auto-scaling/ (18/07/2024)

Launch template is a better option tan Launch Configuration because it allows **reused** and add configuration to new templates.
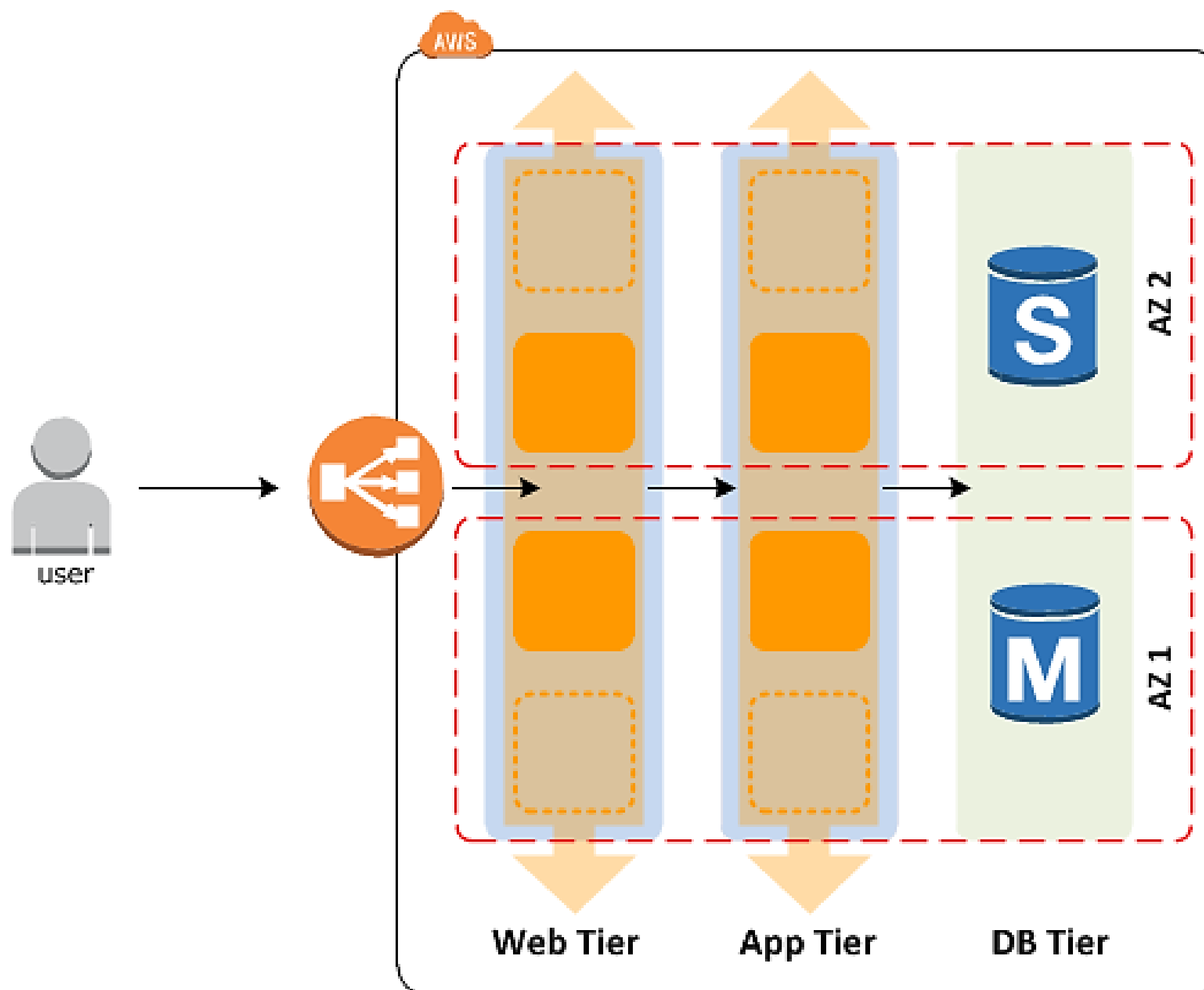
In addition, launch templates allows all configuration from AutoScaling groups, for instance mix OnDemand and Spot Instance on the same group or launch on Dedicated Hosts.
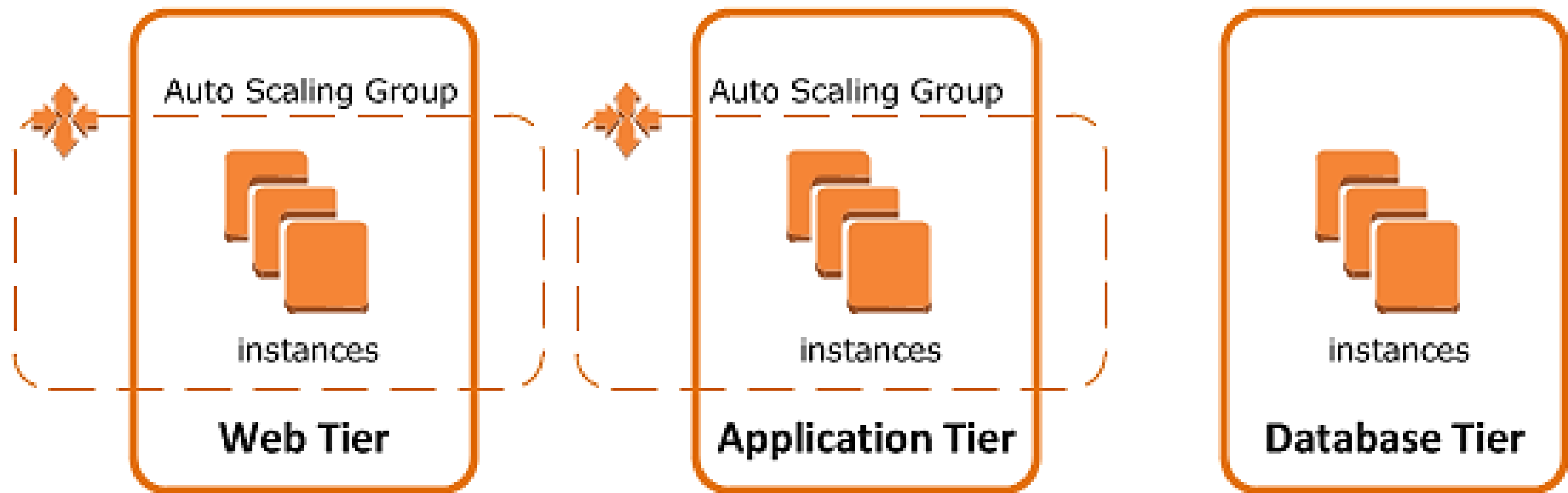
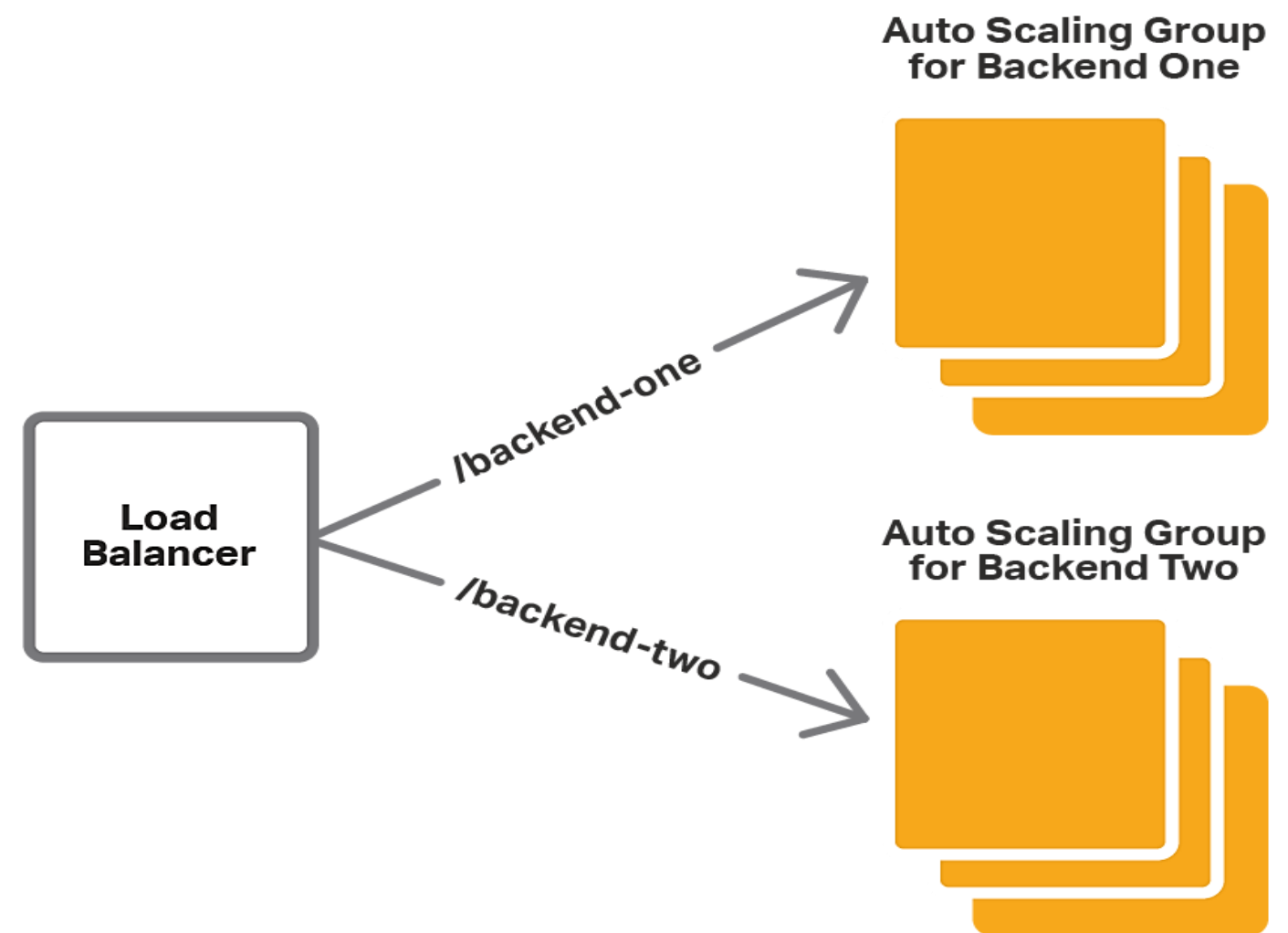Both options contains the same information about instances: type, AMI, Keypair, Security Group, storage, etc.



Launch template

user

Web Tier

App Tier

DB Tier

AWS

AZ 2

AZ 1

Auto Scaling Group — Web Tier — instances

Auto Scaling Group — Application Tier — instances

Database Tier — instances

**Policies:**
- Manual
- Scheduled
- Dynamic
- Dynamic: Target Tracking Policy
- Dynamic: Simple Tracking Policy → Wait until health check and cooldown before evaluate
- Dynamic: Step Tracking Policy

- Metrics on CloudWatch apply on Dynamic.



Auto Scaling Group for Backend One

Auto Scaling Group for Backend Two

Load Balancer

/backend-one

/backend-two

1. Choose LB (i.e. SSL Termination, Globl Balancing with R53)
2. Choose the appropriate instance families and sizing based on the workload in that ASG
3. Consider placement groups in ASG –if applicable-
4. Use Launch Templates
5. Group instances by purpose (i.e. Domains)
6. Set up health checks
7. Utilize target tracking scaling policies
8. Implement cooldown periods
9. Configure notifications
10. Implement proper security (NACLs? Sec Groups)
11. Implement instance termination policies
12. Distribute instances across Availability Zones
13. Use Auto Scaling lifecycle hooks (Launch/Terminate)

Taken from https://www.nops.io/blog/aws-ec2-autoscaling/ (18/07/2024)

Amazon CloudWatch    Alarm    Rule    Event (event-based)    Event (time-based)

- Monitoring and Observability Service to cover AWS Resources, especially computing services such Applications and Containers, in addition Logs, Databases. It collects and track metrics from those resources.
- You can create: Dashboards, Insights, Alarms and Events. You collect: Metrics and Logs.
- Services associated: Cloud Trail, IAM
- Services to trigger actions: SNS – Simple Notification Service, AutoScaling Groups, Lambda; however for Events you can trigger more services.

**Metrics:**

Metrics: Timestamp (UTC), Retention

High-Resolution Monitoring: - Retention: 3 h.

Detailed monitoring: 1 mins (Paid) – R: 15 d.

Basic monitoring: 5 mins – R: 63 d.
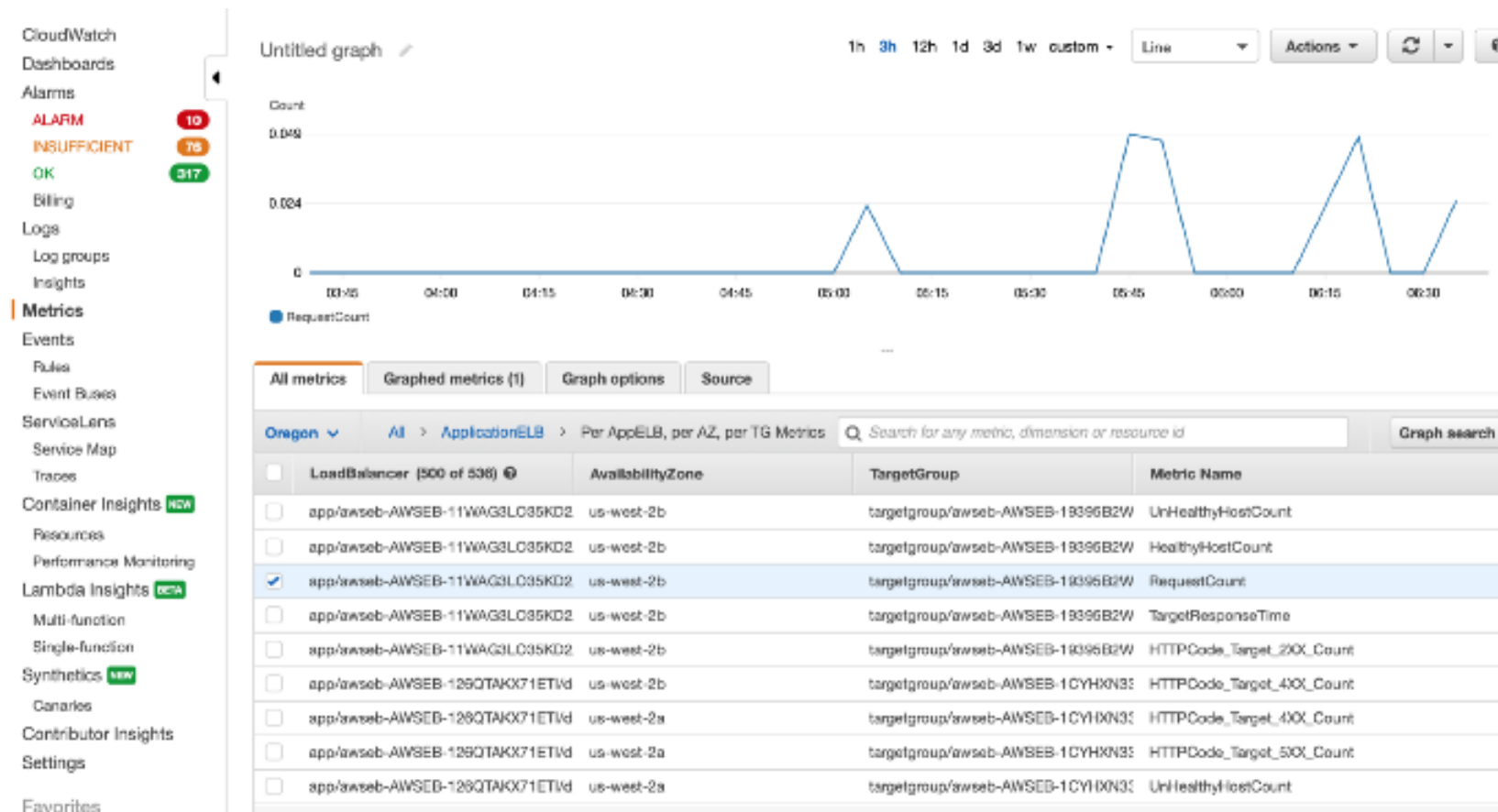
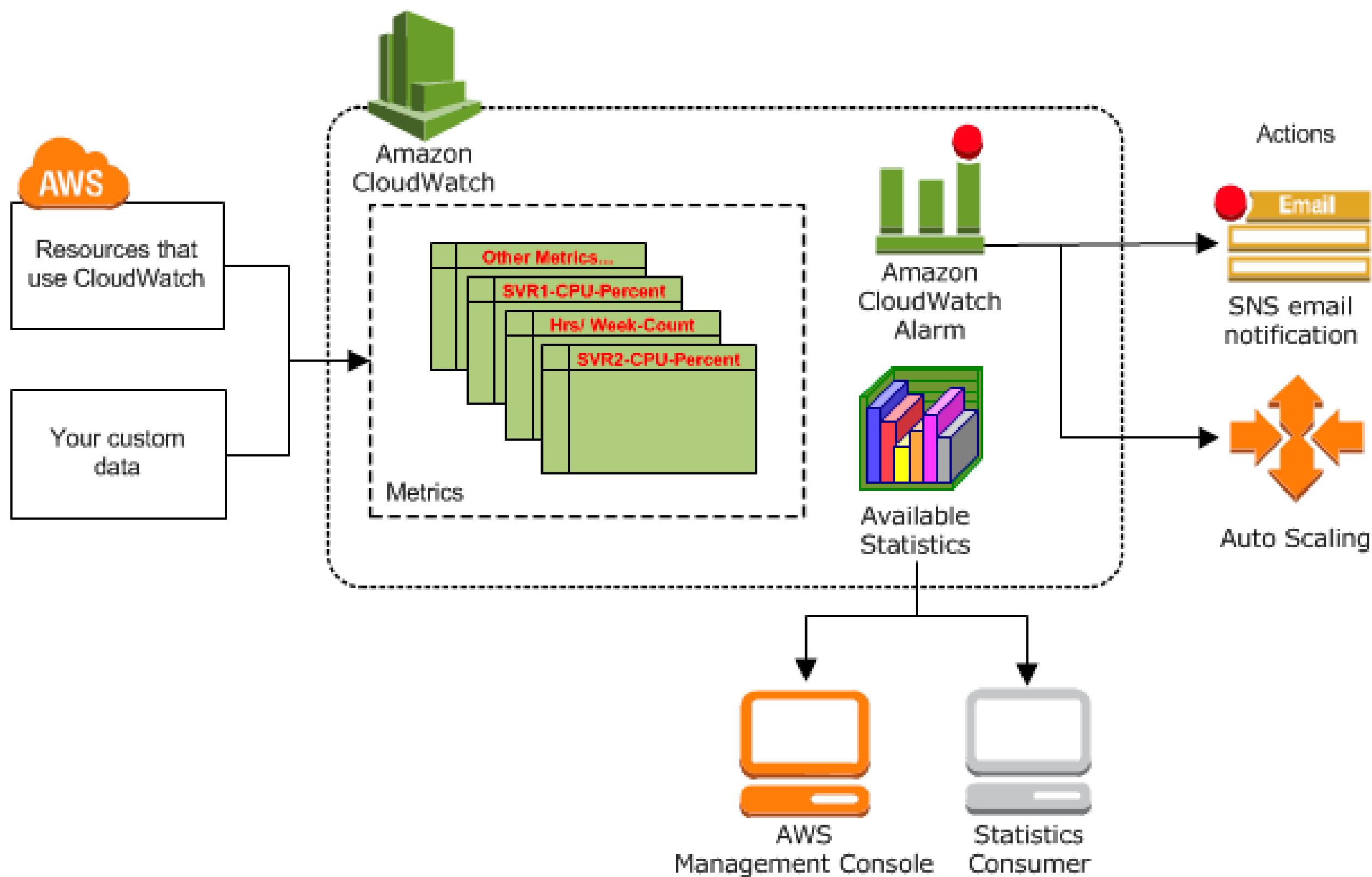Statistics: Max, Min, Average, Sum, Count, Percentiles

Units, Period

Aggregation

**Concepts:**  Namespaces (i.e. Service) / Dimensions (i.e. Group Metrics)

**Alarms:** Trigger SNS and/or AutoScaling Groups
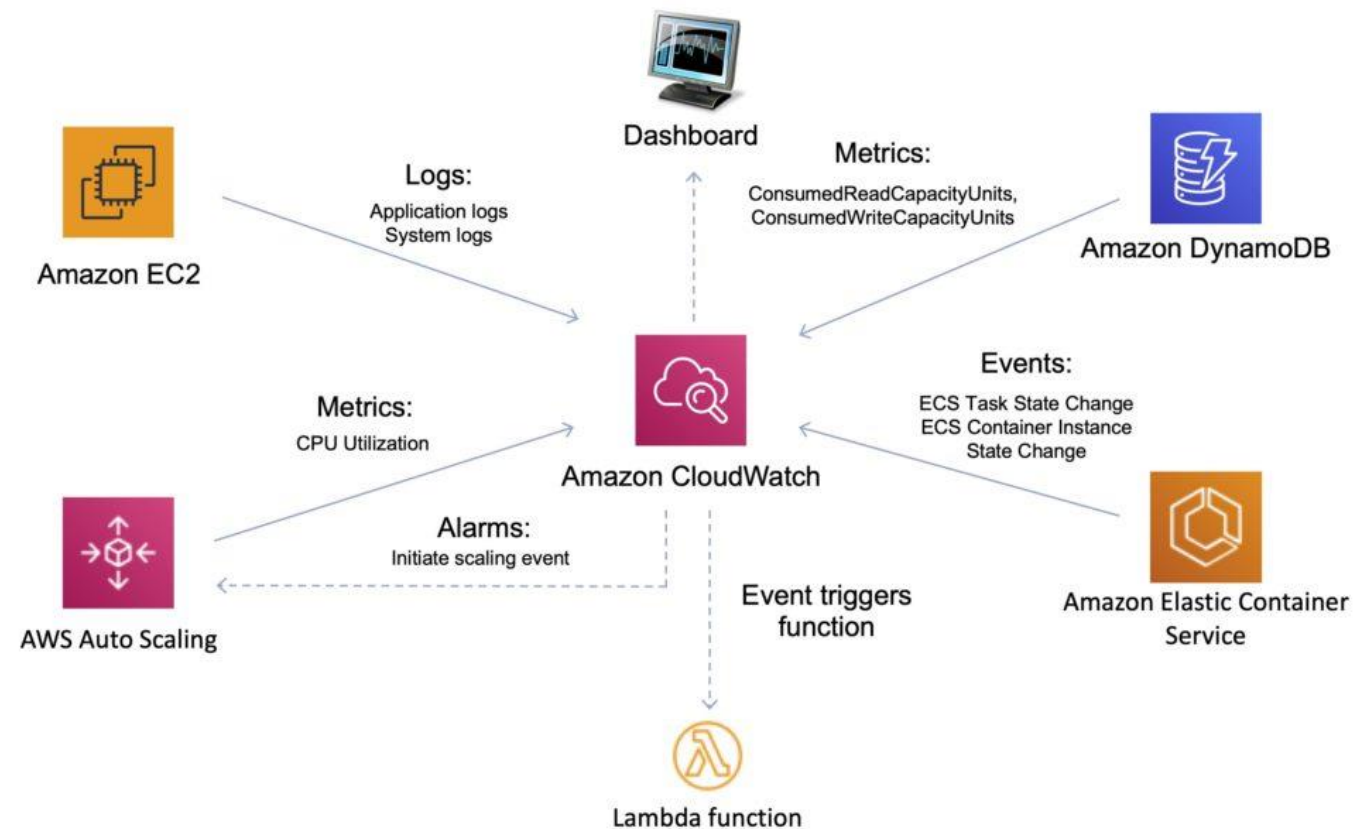
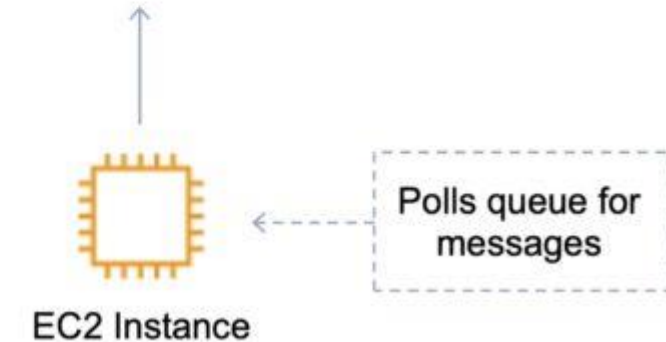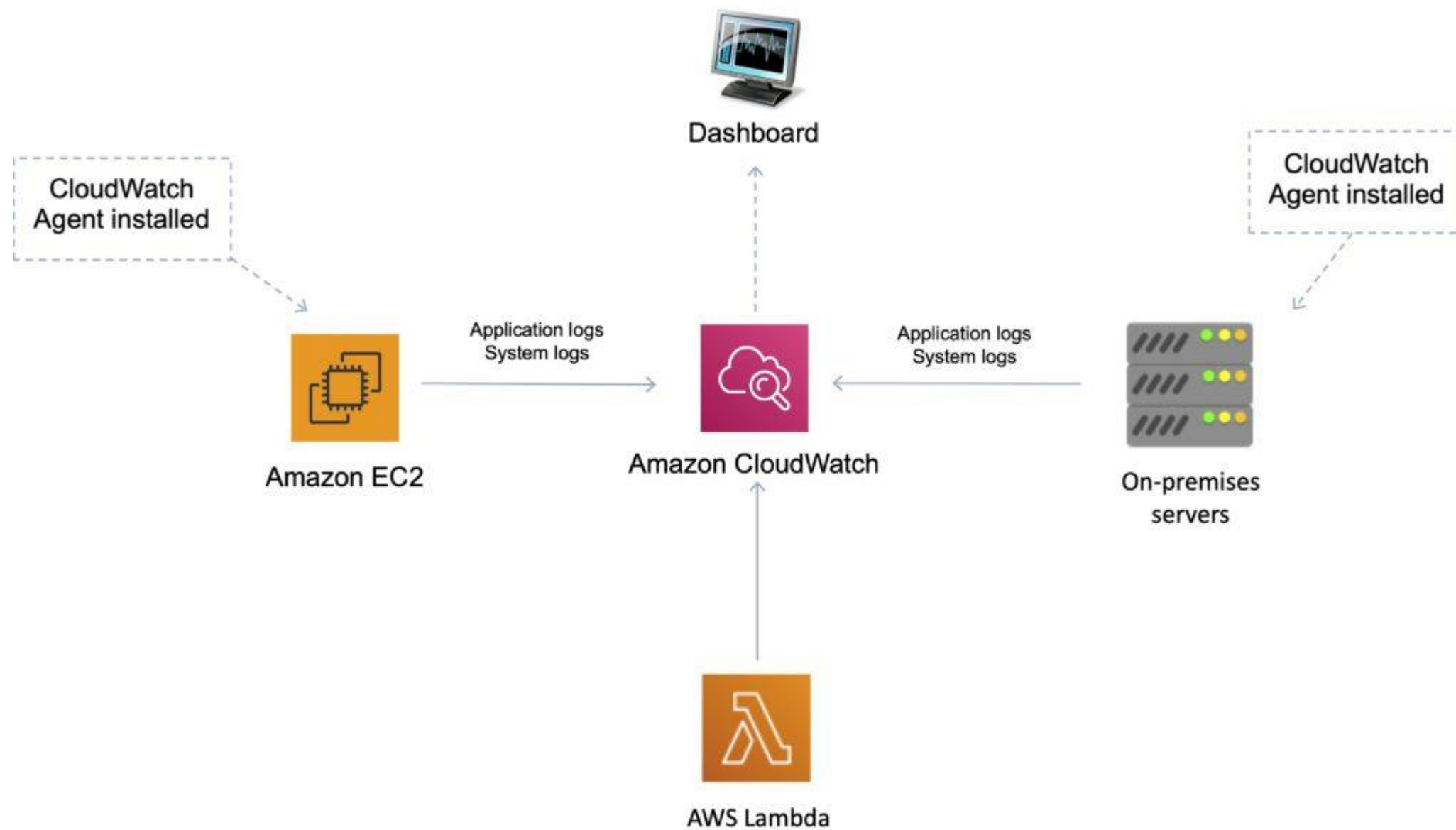**Events:** Trigger something by Event or Schedule

**CloudWatch Alarms**

use <u>conditions</u> to trigger notifications/actions



**CloudWatch Events**

use <u>rules</u> to trigger <u>targets</u>



| Alarms | Events |
|---|---|
| Can be trigger by events/sustained events | Can be trigger by events or timing |
| Limited number of actions: SNS, AutoScaling Group, EC2 Actions. | Many services to trigger, which called Targets: Lambda, SNS, etc. |
| Watch Single Metric (Defined by a period) | Respond to Actions (Near real-time) |
| Can be added to Dashboard | No can be added to Dashboard |

It receives Application and System Logs, and act as centralized repository and then, it can be feed by a Lambda or ElasticSearch Cluster.

Taken from https://digitalcloud.training/amazon-cloudwatch/ (18/07/2024)