



Module 6: Compute

AWS Academy Cloud Foundations

Section 1: Compute services overview

Module 6: Compute



AWS compute services

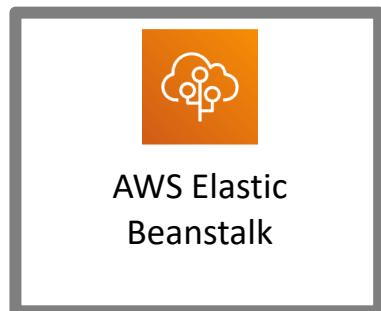
Amazon Web Services (AWS) offers many compute services. This module will discuss the highlighted services.



Amazon EC2
Auto Scaling



VMware Cloud
on AWS



Amazon Lightsail



AWS Batch



AWS Outposts



AWS Serverless
Application Repository

Categorizing compute services

Services	Key Concepts	Characteristics	Ease of Use
• Amazon EC2	<ul style="list-style-type: none">• Infrastructure as a service (IaaS)• Instance-based• Virtual machines	<ul style="list-style-type: none">• Provision virtual machines that you can manage as you choose	A familiar concept to many IT professionals.
• AWS Lambda	<ul style="list-style-type: none">• Serverless computing• Function-based• Low-cost	<ul style="list-style-type: none">• Write and deploy code that runs on a schedule or that can be triggered by events• Use when possible (architect for the cloud)	A relatively new concept for many IT staff members, but easy to use after you learn how.
• Amazon ECS • Amazon EKS • AWS Fargate • Amazon ECR	<ul style="list-style-type: none">• Container-based computing• Instance-based	<ul style="list-style-type: none">• Spin up and run jobs more quickly	AWS Fargate reduces administrative overhead, but you can use options that give you more control.
• AWS Elastic Beanstalk	<ul style="list-style-type: none">• Platform as a service (PaaS)• For web applications	<ul style="list-style-type: none">• Focus on your code (building your application)• Can easily tie into other services—databases, Domain Name System (DNS), etc.	Fast and easy to get started.

Choosing the optimal compute service

- The optimal compute service or services that you use will depend on your use case
- Some aspects to consider –
 - What is your application design?
 - What are your usage patterns?
 - Which configuration settings will you want to manage?
- Selecting the wrong compute solution for an architecture can lead to lower performance efficiency
- A good starting place—Understand the available compute options

Section 2: Amazon EC2

Module 6: Compute



Amazon Elastic Compute Cloud (Amazon EC2)

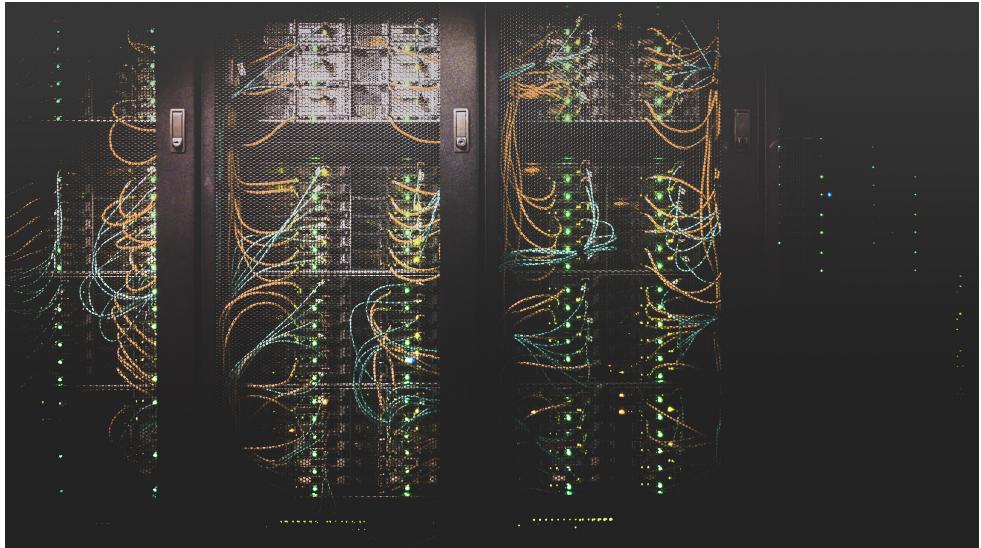


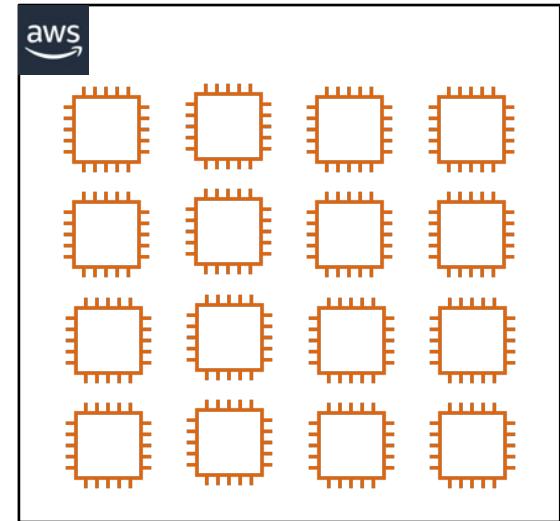
Photo by Taylor Vick on Unsplash

On-premises servers



Example uses of Amazon EC2 instances

- ✓ Application server
- ✓ Web server
- ✓ Database server
- ✓ Game server
- ✓ Mail server
- ✓ Media server
- ✓ Catalog server
- ✓ File server
- ✓ Computing server
- ✓ Proxy server



Amazon EC2 instances

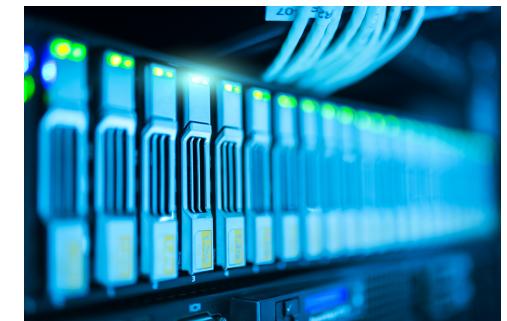
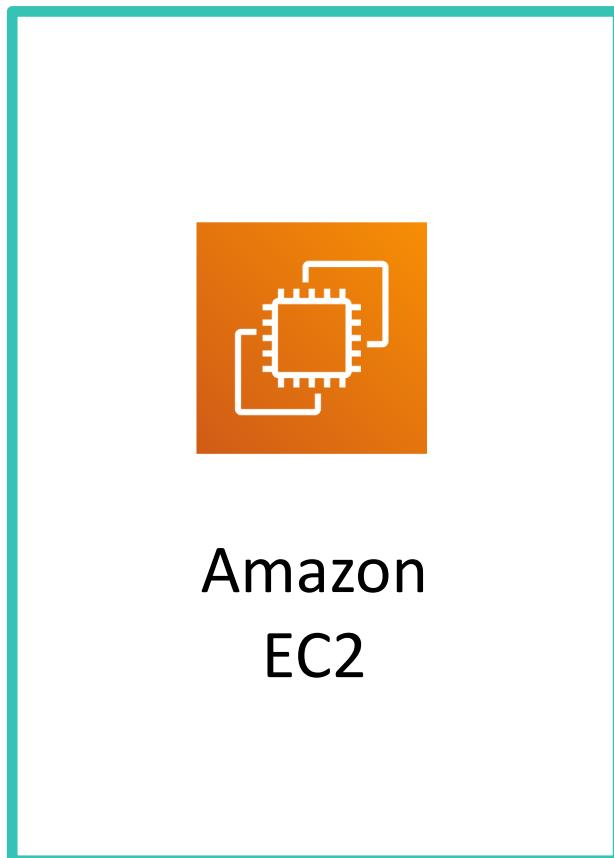


Photo by panumas nikhomkhai from Pexels

Amazon EC2 overview



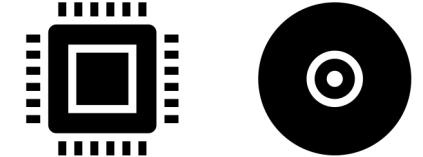
- **Amazon Elastic Compute Cloud (Amazon EC2)**
 - Provides **virtual machines**—referred to as **EC2 instances**—in the cloud.
 - Gives you *full control* over the guest operating system (Windows or Linux) on each instance.
 - You can launch instances of any size into an **Availability Zone** anywhere in the world.
 - Launch instances from **Amazon Machine Images (AMIs)**.
 - Launch instances with a few clicks or a line of code, and they are ready in minutes.
 - You can control traffic to and from instances.

2. Select an instance type

Choices made using the Launch Instance Wizard:

1. AMI
2. **Instance Type**
3. Network settings
4. IAM role
5. User data
6. Storage options
7. Tags
8. Security group
9. Key pair

- Consider your use case
 - How will the EC2 instance you create be used?
- The **instance type** that you choose determines –
 - Memory (RAM)
 - Processing power (CPU)
 - Disk space and disk type (Storage)
 - Network performance
- Instance type categories –
 - General purpose
 - Compute optimized
 - Memory optimized
 - Storage optimized
 - Accelerated computing
- Instance types offer *family, generation, and size*



EC2 instance type naming and sizes

Instance type naming

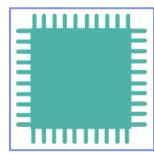
- Example: **t3.large**
 - T is the family name
 - 3 is the generation number
 - Large is the size

Example instance sizes

Instance Name	vCPU	Memory (GB)	Storage
t3.nano	2	0.5	EBS-Only
t3.micro	2	1	EBS-Only
t3.small	2	2	EBS-Only
t3.medium	2	4	EBS-Only
t3.large	2	8	EBS-Only
t3.xlarge	4	16	EBS-Only
t3.2xlarge	8	32	EBS-Only



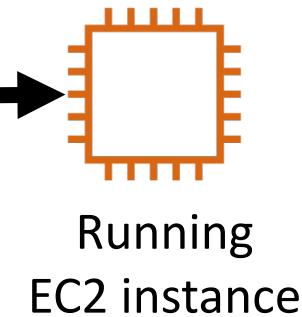
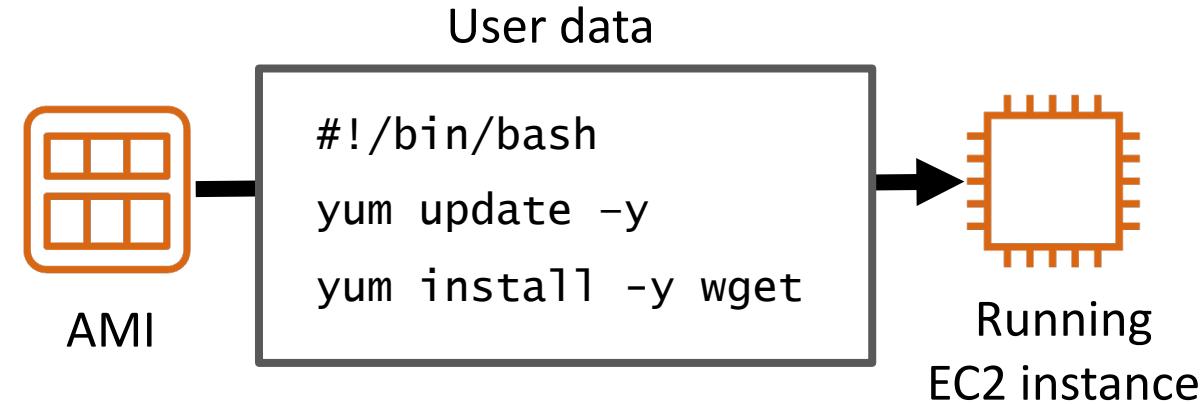
Select instance type: Based on use case

	 General Purpose	 Compute Optimized	 Memory Optimized	 Accelerated Computing	 Storage Optimized
Instance Types	a1, m4, m5, t2, t3	c4, c5	r4, r5, x1, z1	f1, g3, g4, p2, p3	d2, h1, i3
Use Case	Broad	High performance	In-memory databases	Machine learning	Distributed file systems

5. User data script (optional)

Choices made by using the Launch Instance Wizard:

1. AMI
2. Instance Type
3. Network settings
4. IAM role
5. User data
6. Storage options
7. Tags
8. Security group
9. Key pair



Running
EC2 instance

- Optionally specify a user data script at instance launch
- Use **user data** scripts to customize the runtime environment of your instance
 - Script runs the first time the instance starts
 - Can be used strategically
 - For example, reduce the number of custom AMIs that you build and maintain

6. Specify storage

Choices made by using the Launch Instance Wizard:

1. AMI
2. Instance Type
3. Network settings
4. IAM role
5. User data
6. Storage options
7. Tags
8. Security group
9. Key pair

- Configure the **root volume**
 - Where the guest operating system is installed
- Attach **additional storage volumes (optional)**
 - AMI might already include more than one volume
- For each volume, specify:
 - The **size** of the disk (in GB)
 - The **volume type**
 - Different types of solid state drives (SSDs) and hard disk drives (HDDs) are available
 - If the volume will be deleted when the instance is terminated
 - If **encryption** should be used



Amazon EC2 storage options

- **Amazon Elastic Block Store (Amazon EBS) –**
 - Durable, block-level storage volumes.
 - You can stop the instance and start it again, and the data will still be there.
- **Amazon EC2 Instance Store –**
 - Ephemeral storage is provided on disks that are attached to the host computer where the EC2 instance is running.
 - If the instance stops, data stored here is deleted.
- Other options for storage (not for the root volume) –
 - Mount an **Amazon Elastic File System (Amazon EFS)** file system.
 - Connect to **Amazon Simple Storage Service (Amazon S3)**.

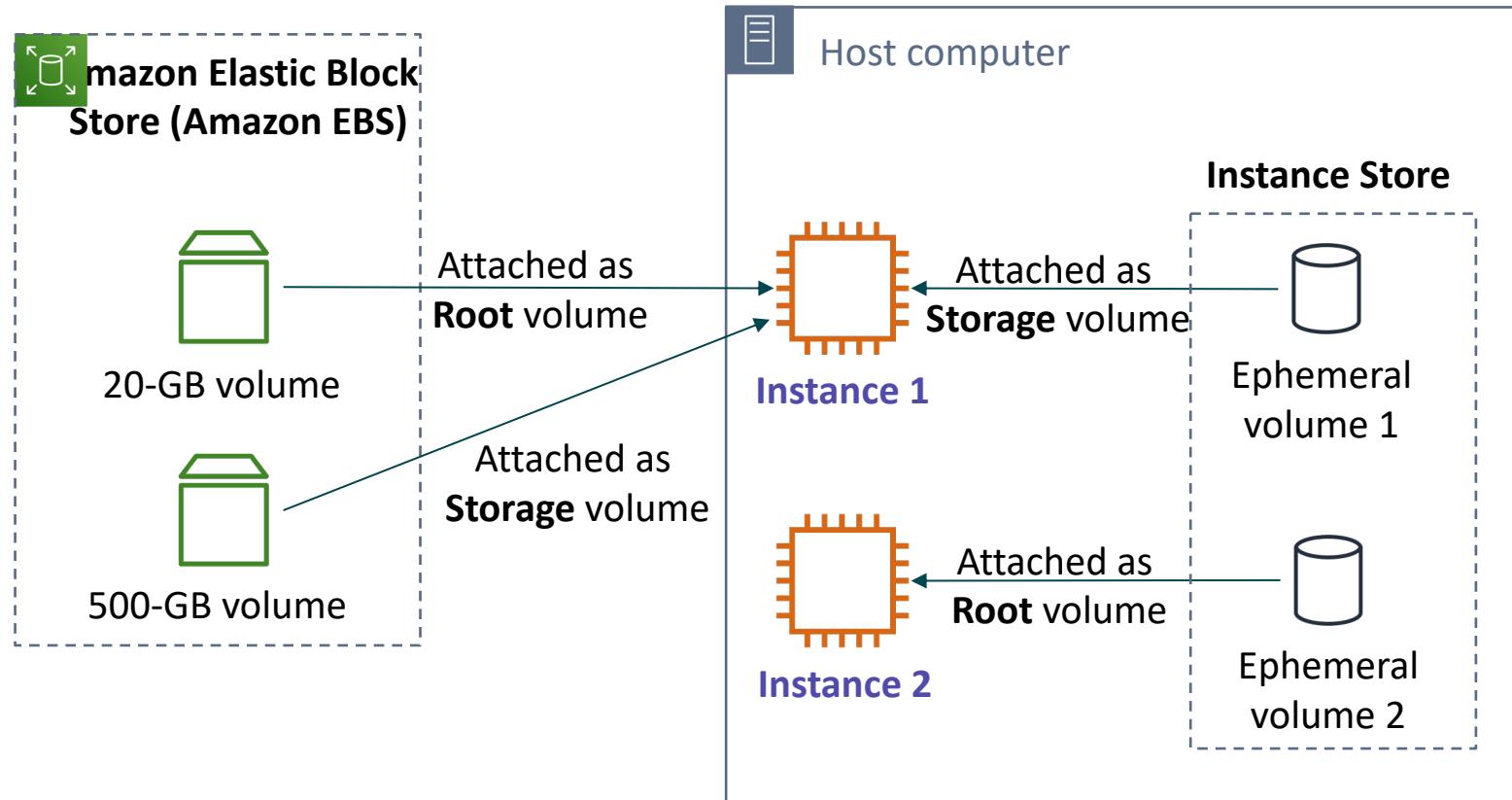
Example storage options

- **Instance 1 characteristics –**

- It has an **Amazon EBS root volume** type for the operating system.
- What will happen if the instance is stopped and then started again?

- **Instance 2 characteristics –**

- It has an **Instance Store root volume** type for the operating system.
- What will happen if the instance stops (because of user error or a system malfunction)?



Amazon EC2 console view of a running EC2 instance

The screenshot shows the AWS EC2 Management Console interface. On the left, a sidebar menu lists various services under 'Instances' (Instances, Launch Templates, Spot Requests, Reserved Instances, Dedicated Hosts, Scheduled Instances, Capacity Reservations), 'Images' (AMIs, Bundle Tasks), and 'Elastic Block Store' (Volumes, Snapshots). The main content area displays a table of instances. A search bar at the top of the table results section shows a result for 'i-092b6f3efba959a53'. The table has columns for Name, Instance ID, Instance Type, Instance State, Status Checks, Public DNS (IPv4), and IPv4 Public IP. One row is highlighted for the instance with ID 'i-092b6f3efba959a53', which is a 't2.micro' type in the 'running' state. Below the table, detailed information for this instance is shown in a card format. The card includes tabs for 'Description', 'Status Checks', 'Monitoring', and 'Tags'. The 'Description' tab is selected and displays the following details:

Attribute	Value
Instance ID	i-092b6f3efba959a53
Public DNS (IPv4)	ec2-54-159-171-63.compute-1.amazonaws.com
IPv4 Public IP	54.159.171.63
IPv6 IPs	-
Private DNS	ip-172-31-82-44.ec2.internal
Private IPs	172.31.82.44
Secondary private IPs	
VPC ID	vpc-e4e9859e
Subnet ID	subnet-d22779fc
Network interfaces	eth0

Other tabs in the card include 'Status Checks' (which shows 'Initializing'), 'Monitoring', and 'Tags'.

Another option: Launch an EC2 instance with the AWS Command Line Interface

- EC2 instances can also be created programmatically.
- This example shows how simple the command can be.
 - This command assumes that the key pair and security group already exist.
 - More options could be specified. See the [AWS CLI Command Reference](#) for details.

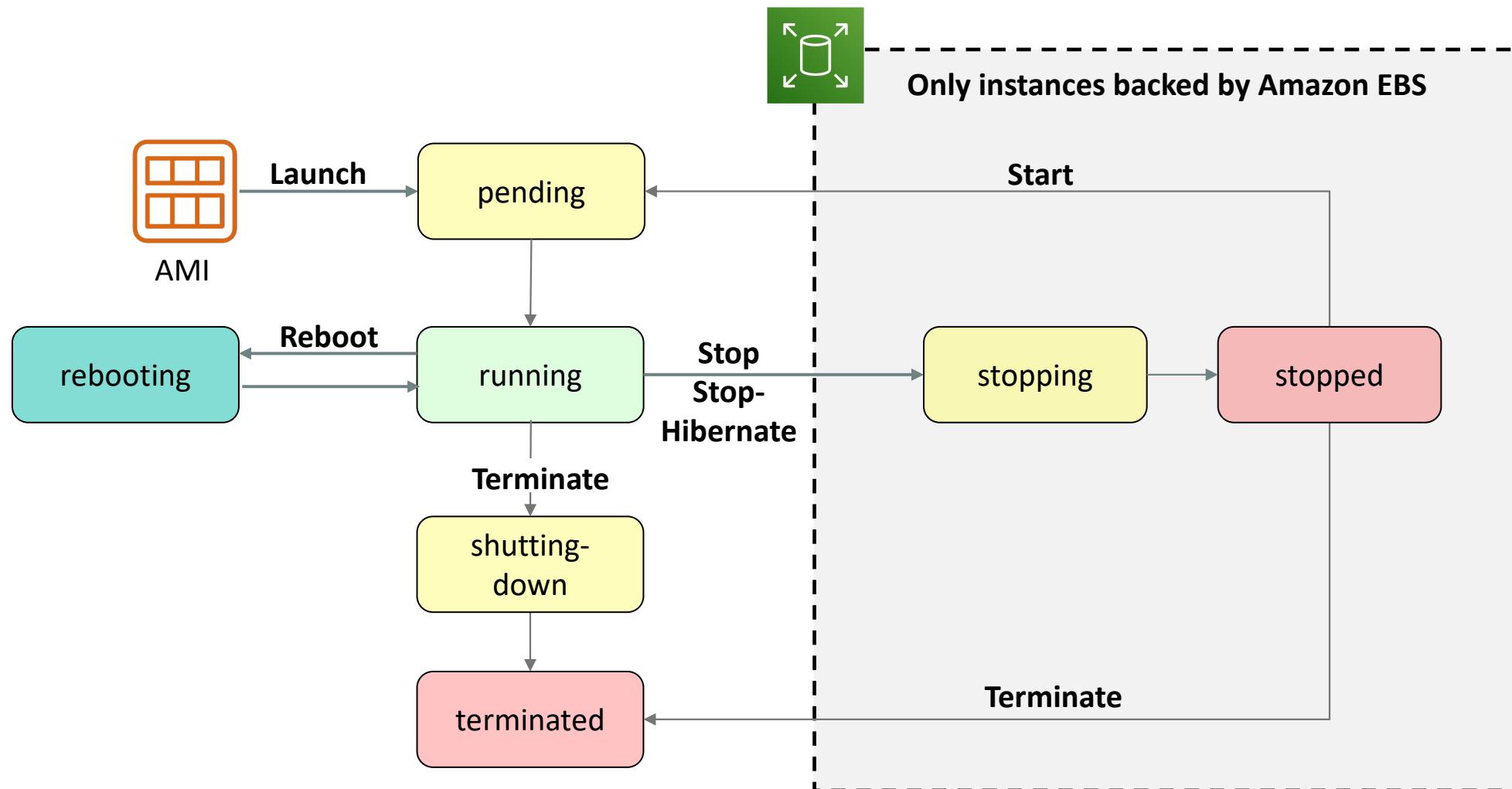


AWS Command Line Interface (AWS CLI)

Example command:

```
aws ec2 run-instances \
--image-id ami-1a2b3c4d \
--count 1 \
--instance-type c3.large \
--key-name MyKeyPair \
--security-groups MySecurityGroup \
--region us-east-1
```

Amazon EC2 instance lifecycle



Activity: Check your understanding

1. Between Amazon EC2 or Amazon RDS, which provides a managed service? What does *managed service* mean?
 - **ANSWER:** Amazon RDS provides a managed service. Amazon RDS handles provisioning, installation and patching, automated backups, restoring snapshots from points in time, high availability, and monitoring.
2. Name at least one advantage of deploying Microsoft SQL Server on Amazon EC2 instead of Amazon RDS.
 - **ANSWER:** Amazon EC2 offers complete control over every configuration, the OS, and the software stack.
3. What advantage does the Quick Start provide over a manual installation on Amazon EC2?
 - **ANSWER:** The Quick Start is a reference architecture with proven best practices built into the design.
4. Which deployment option offers the best approach for all use cases?
 - **ANSWER:** Neither. The correct deployment option depends on your specific needs.
5. Which approach costs more: using Amazon EC2 or using Amazon RDS?
 - **ANSWER:** It depends. Managing the database deployment on Amazon EC2 requires more customer oversight and time. If time is your priority, then Amazon RDS might be less expensive. If you have in-house expertise, Amazon EC2 might be more cost-effective.

Section 3: Amazon EC2 cost optimization

Module 6: Compute



Amazon EC2 pricing models

On-Demand Instances

- Pay by the hour
- No long-term commitments.
- Eligible for the [AWS Free Tier](#).

Dedicated Hosts

- A physical server with EC2 instance capacity fully dedicated to your use.

Dedicated Instances

- Instances that run in a VPC on hardware that is dedicated to a single customer.

Reserved Instances

- Full, partial, or no upfront payment for instance you reserve.
- Discount on hourly charge for that instance.
- 1-year or 3-year term.

Scheduled Reserved Instances

- Purchase a capacity reservation that is always available on a recurring schedule you specify.
- 1-year term.

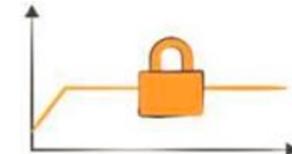
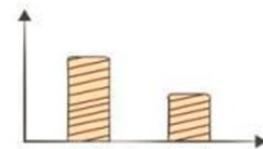
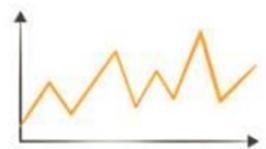
Spot Instances

- Instances run as long as they are available and your bid is above the Spot Instance price.
- They can be interrupted by AWS with a 2-minute notification.
- Interruption options include terminated, stopped or hibernated.
- Prices can be significantly less expensive compared to On-Demand Instances
- Good choice when you have flexibility in when your applications can run.

Per second billing available for On-Demand Instances, Reserved Instances, and Spot Instances that run Amazon Linux or Ubuntu.

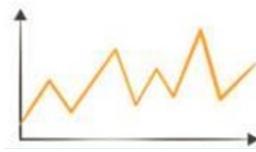


Amazon EC2 pricing models: Benefits

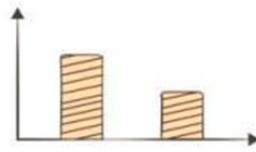


On-Demand Instances	Spot Instances	Reserved Instances	Dedicated Hosts
<ul style="list-style-type: none">Low cost and flexibility	<ul style="list-style-type: none">Large scale, dynamic workload	<ul style="list-style-type: none">Predictability ensures compute capacity is available when needed	<ul style="list-style-type: none">Save money on licensing costsHelp meet compliance and regulatory requirements

Amazon EC2 pricing models: Use cases



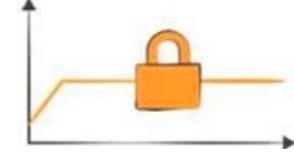
Spiky Workloads



Time-Insensitive Workloads



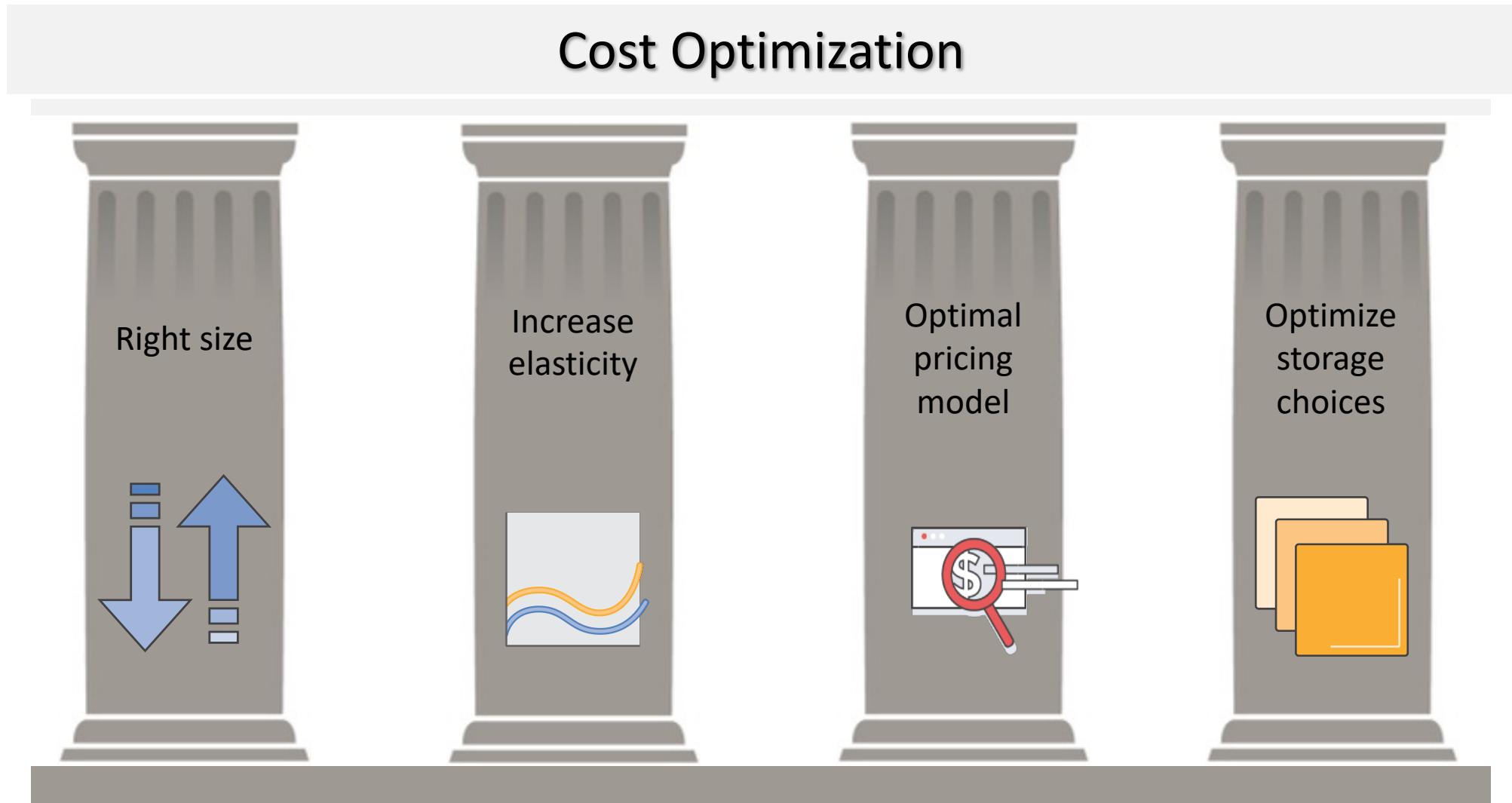
Steady-State Workloads



Highly Sensitive Workloads

On-Demand Instances	Spot Instances	Reserved Instances	Dedicated Hosts
<ul style="list-style-type: none">• Short-term, spiky, or unpredictable workloads• Application development or testing	<ul style="list-style-type: none">• Applications with flexible start and end times• Applications only feasible at very low compute prices• Users with urgent computing needs for large amounts of additional capacity	<ul style="list-style-type: none">• Steady state or predictable usage workloads• Applications that require reserved capacity, including disaster recovery• Users able to make upfront payments to reduce total computing costs even further	<ul style="list-style-type: none">• Bring your own license (BYOL)• Compliance and regulatory restrictions• Usage and licensing tracking• Control instance placement

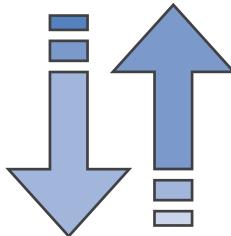
The four pillars of cost optimization



Pillar 1: Right size

Pillars:

1. Right size
2. Increase elasticity
3. Optimal pricing model
4. Optimize storage choices



✓ Provision instances to match the need

- CPU, memory, storage, and network throughput
- Select appropriate [instance types](#) for your use

✓ Use Amazon CloudWatch metrics

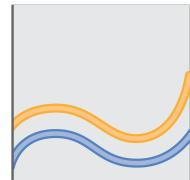
- How idle are instances? When?
- Downsize instances

✓ Best practice: Right size, then reserve

Pillar 2: Increase elasticity

Pillars:

1. Right-Size
2. **Increase Elasticity**
3. Optimal pricing model
4. Optimize storage choices

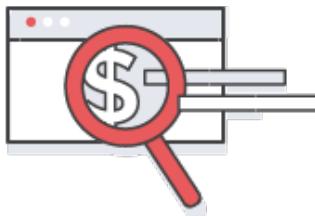


- ✓ Stop or hibernate Amazon EBS-backed instances that are not actively in use
 - Example: non-production development or test instances
- ✓ Use automatic scaling to match needs based on usage
 - Automated and time-based elasticity

Pillar 3: Optimal pricing model

Pillars:

1. Right-Size
2. Increase Elasticity
- 3. Optimal pricing model**
4. Optimize storage choices

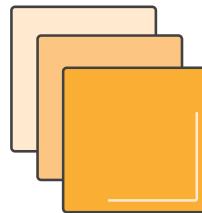


- ✓ Leverage the right pricing model for your use case
 - Consider your usage patterns
- ✓ Optimize and *combine* purchase types
- ✓ Examples:
 - Use **On-Demand Instances** and **Spot Instances** for variable workloads
 - Use **Reserved Instances** for predictable workloads
- ✓ Consider serverless solutions (**AWS Lambda**)

Pillar 4: Optimize storage choices

Pillars:

1. Right-Size
2. Increase Elasticity
3. Optimal pricing model
4. **Optimize storage choices**



- ✓ Reduce costs while maintaining storage performance and availability
- ✓ Resize EBS volumes
- ✓ Change EBS volume types
 - ✓ Can you meet performance requirements with less expensive storage?
 - ✓ Example: **Amazon EBS Throughput Optimized HDD (st1)** storage typically costs half as much as the default **General Purpose SSD (gp2)** storage option.
- ✓ Delete EBS snapshots that are no longer needed
- ✓ Identify the most appropriate destination for specific types of data
 - ✓ Does the application need the instance to reside on Amazon EBS?
 - ✓ Amazon S3 storage options with lifecycle policies can reduce costs

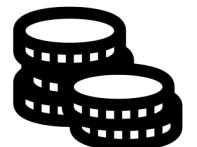
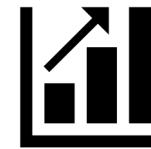
Measure, monitor, and improve

- Cost optimization is an ongoing process.



- Recommendations –

- Define and enforce **cost allocation tagging**.
- Define metrics, set targets, and review regularly.
- Encourage teams to **architect for cost**.
- Assign the responsibility of optimization to an individual or to a team.



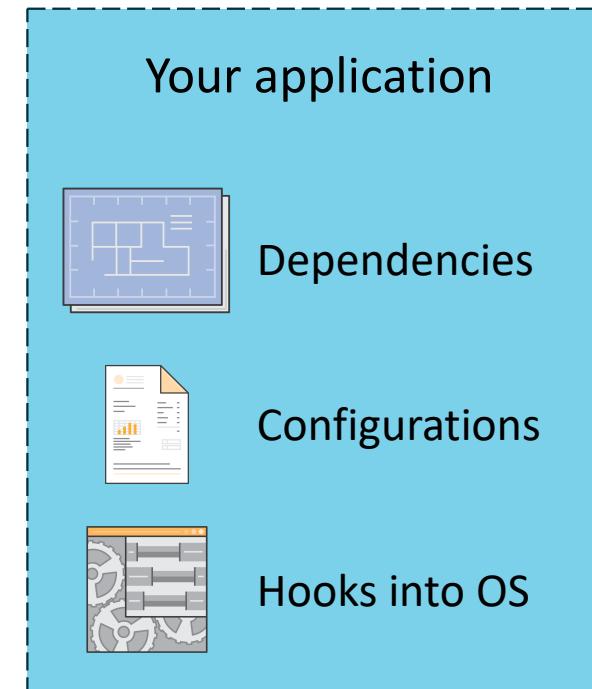
Container basics

- **Containers** are a method of operating system virtualization.

- Benefits –

- Repeatable.
- Self-contained environments.
- Software runs the same in different environments.
 - Developer's laptop, test, production.
- Faster to launch and stop or terminate than virtual machines

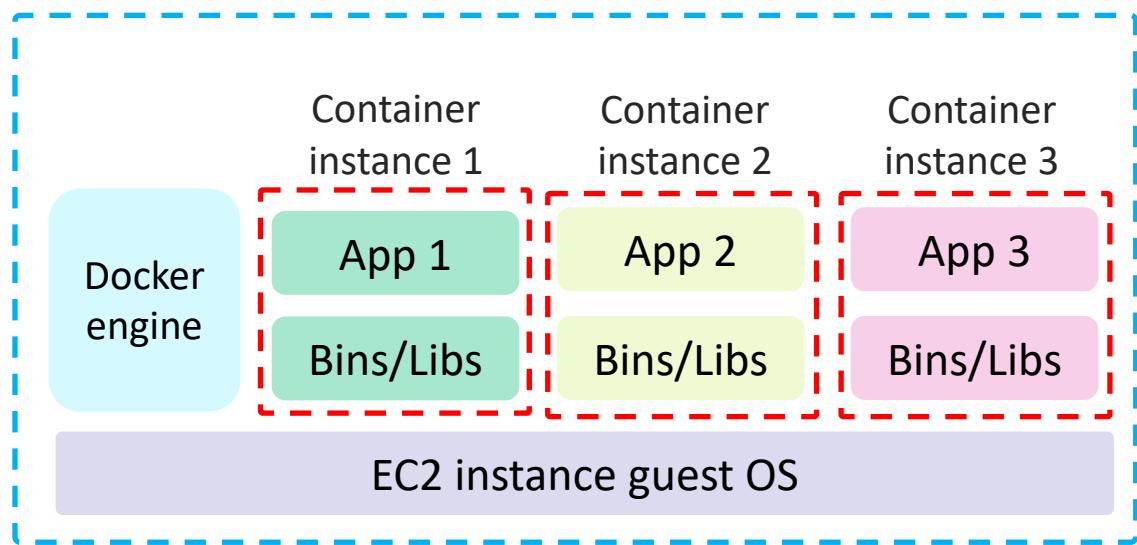
Your Container



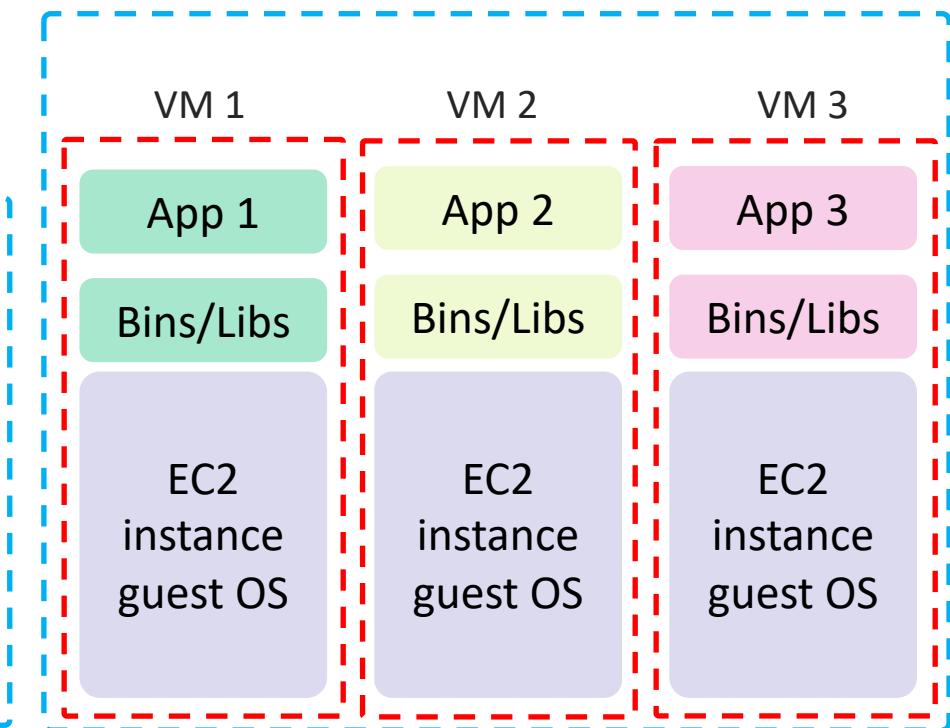
Containers versus virtual machines

Example

Three containers on one EC2 instance



Three virtual machines on three EC2 instances



Hypervisor

Host operating system

Physical server

Part of
AWS Global
Infrastructure

Amazon Elastic Container Service (Amazon ECS)

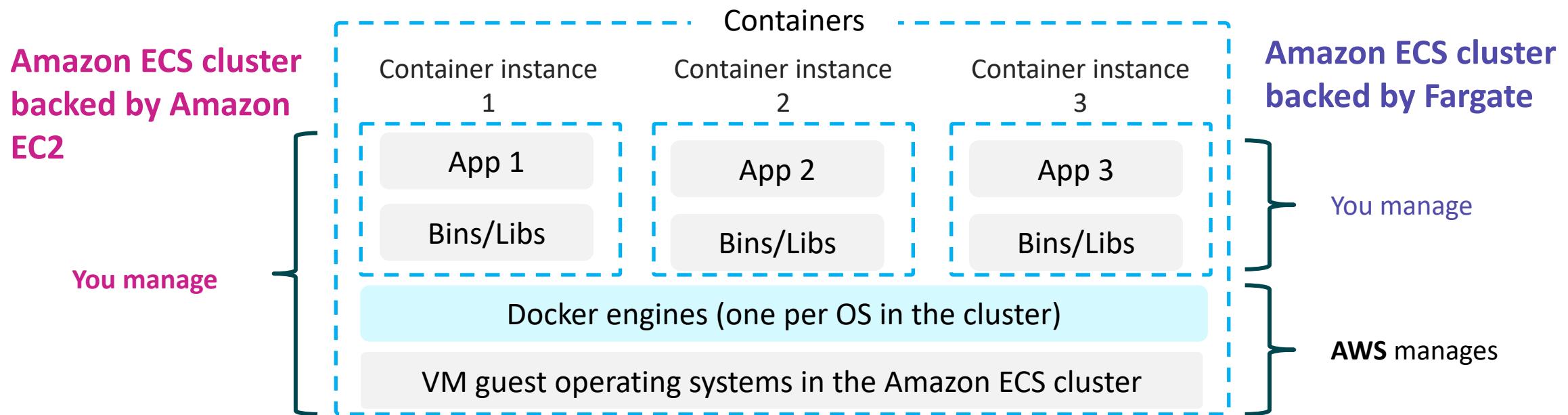
- Amazon Elastic Container Service ([Amazon ECS](#)) –
 - A highly scalable, fast, [container management service](#)
- Key benefits –
 - Orchestrates the running of Docker containers
 - Maintains and scales the fleet of nodes that run your containers
 - Removes the complexity of standing up the infrastructure
- Integrated with features that are familiar to Amazon EC2 service users –
 - Elastic Load Balancing
 - Amazon EC2 security groups
 - Amazon EBS volumes
 - IAM roles



Amazon Elastic
Container Service

Amazon ECS cluster options

- Key question: Do *you* want to manage the Amazon ECS cluster that runs the containers?
 - If yes, create an **Amazon ECS cluster backed by Amazon EC2** (provides more granular control over infrastructure)
 - If no, create an **Amazon ECS cluster backed by AWS Fargate** (easier to maintain, focus on your applications)



Amazon Elastic Kubernetes Service (Amazon EKS)

- Amazon Elastic Kubernetes Service (**Amazon EKS**)

- Enables you to run Kubernetes on AWS
- Certified Kubernetes conformant (supports easy migration)
- Supports Linux and Windows containers
- Compatible with Kubernetes community tools and supports popular Kubernetes add-ons



Amazon Elastic
Kubernetes Service

- Use Amazon EKS to –

- Manage clusters of Amazon EC2 compute instances
- Run containers that are orchestrated by Kubernetes on those instances

Amazon Elastic Container Registry (Amazon ECR)

Amazon ECR is a fully managed Docker [container registry](#) that makes it easy for developers to store, manage, and deploy Docker container images.



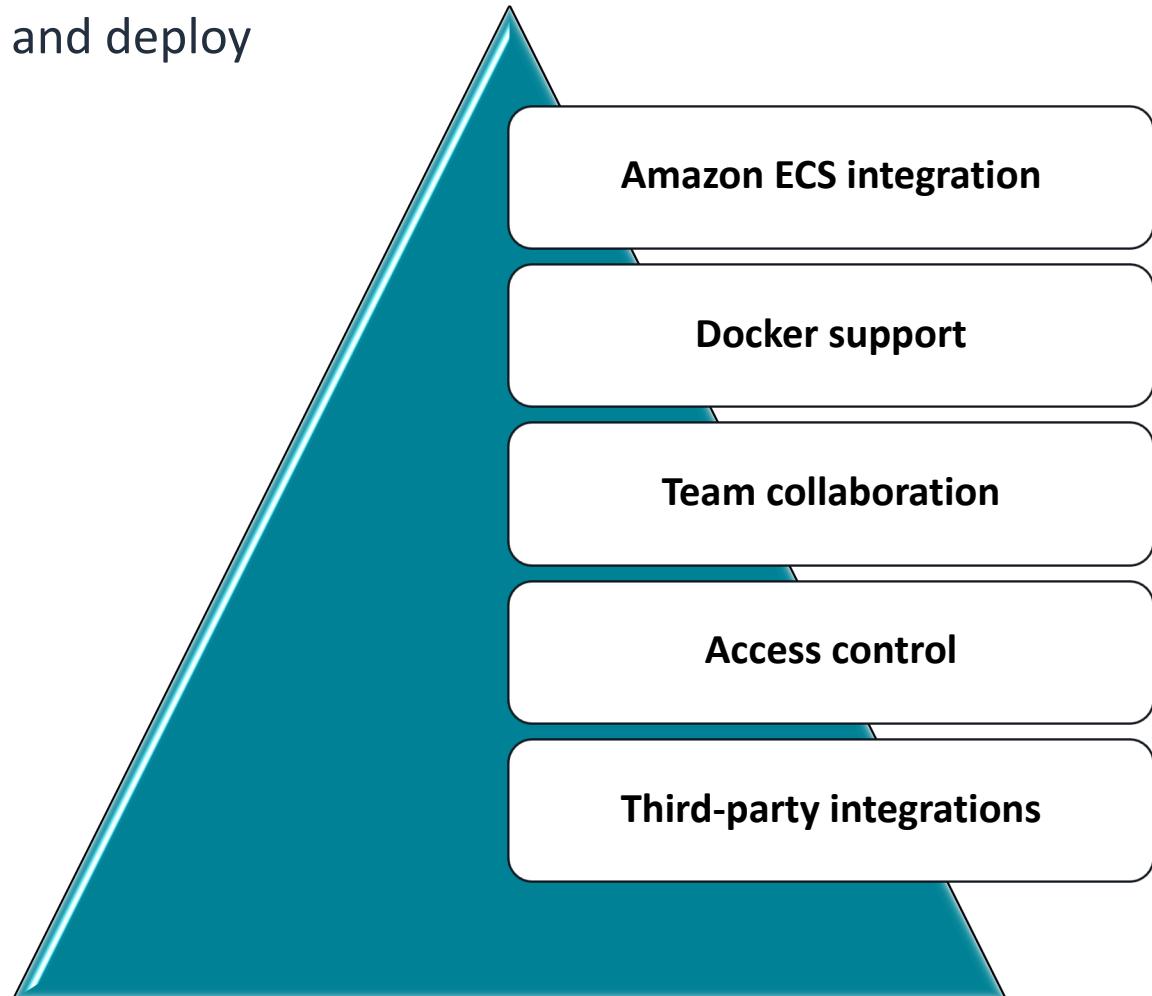
Amazon Elastic
Container Registry



Image

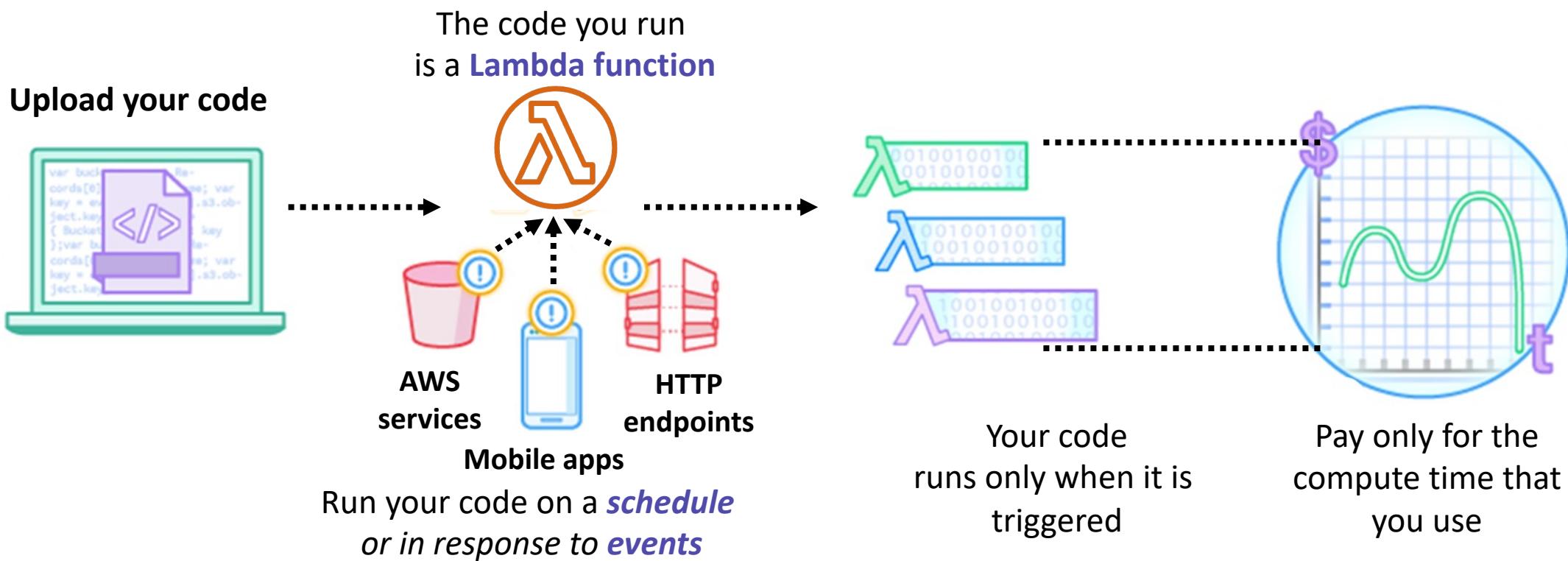


Registry



AWS Lambda: Run code without servers

AWS Lambda is a **serverless** compute service.



Benefits of Lambda



AWS
Lambda



It supports multiple programming languages



Completely automated administration



Built-in fault tolerance

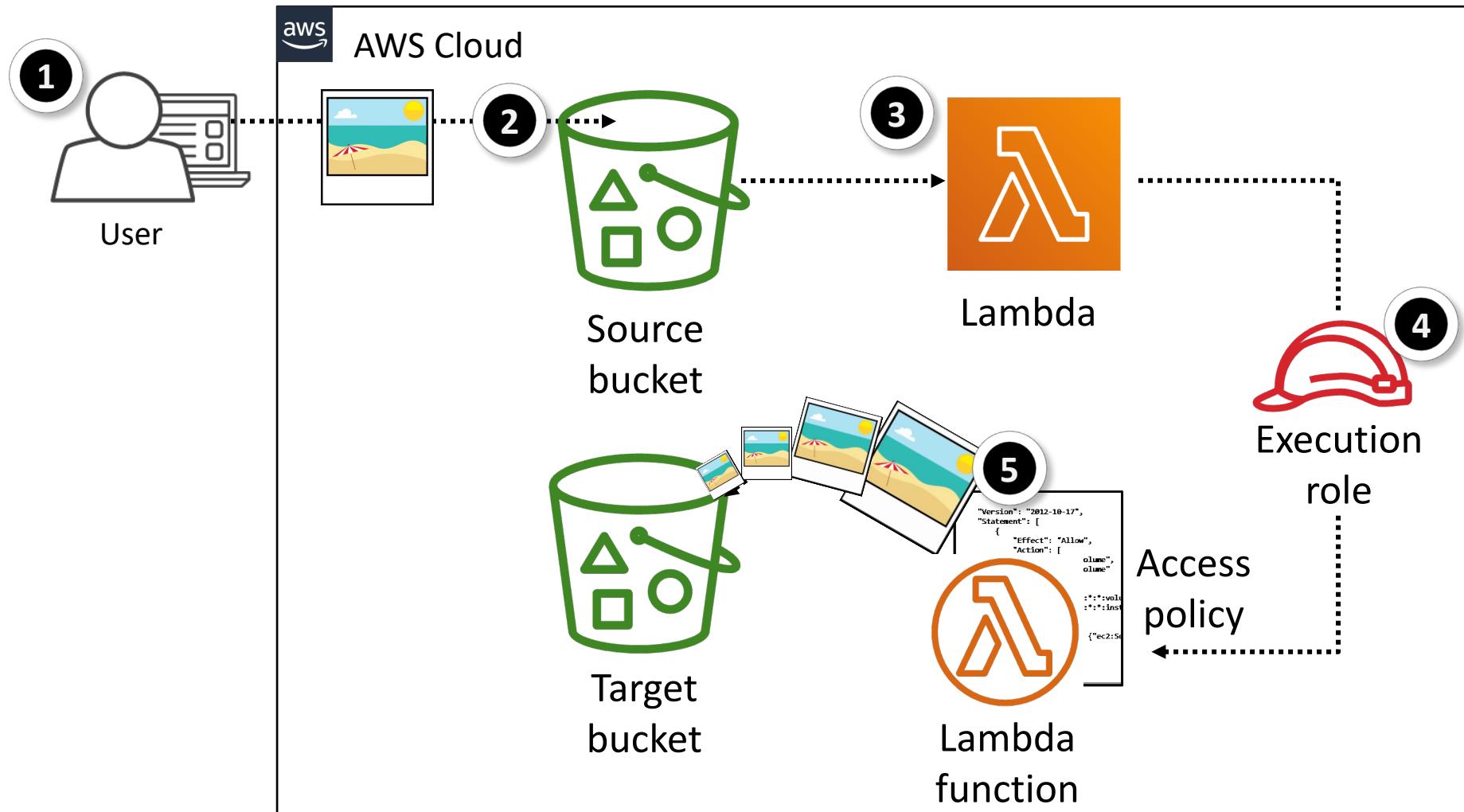


It supports the orchestration of multiple functions



Pay-per-use pricing

Event-based Lambda function example: Create thumbnail images



AWS Elastic Beanstalk

- An easy way to get **web applications** up and running

- A **managed service** that automatically handles –



**AWS Elastic
Beanstalk**

- Infrastructure provisioning and configuration
- Deployment
- Load balancing
- Automatic scaling
- Health monitoring
- Analysis and debugging
- Logging

- No additional charge for Elastic Beanstalk
- Pay only for the underlying resources that are used

AWS Elastic Beanstalk deployments

- It supports web applications written for common platforms
 - Java, .NET, PHP, Node.js, Python, Ruby, Go, and Docker

- You upload your code
 - Elastic Beanstalk automatically handles the deployment
 - Deploys on servers such as Apache, NGINX, Passenger, Puma, and Microsoft Internet Information Services (IIS)

