



Intro to Amazon Cloud

EC2 overview

Larry Babarinde

Solutions Architect, WWPS

11/12/21

Agenda

- Introduction to AWS Cloud
 - Global Reach
 - EC2 Overview
 - EC2 Details
-
- Visit <https://coatl.awsimmersionday.net> for the full day agenda

What is AWS?





AWS provides a highly reliable, scalable, low-cost infrastructure platform in the cloud that powers millions of businesses in over 190 countries around the world.

Benefits

- Low Cost
- Elasticity & Agility
- Open & Flexible
- Secure
- Global Reach



What sets AWS apart?

-  Security
Fine-grained control
-  Service Breadth & Depth; pace of innovation
175+ services to support any cloud workload; rapid customer driven releases
-  Experience: 1M+ customers
Building and managing cloud since 2006
-  Global Footprint
80+ Availability Zones within 25 geographic Regions, 1 Local Zone, 216 Points of Presence (205 Edge Locations and 11 Regional Edge Caches) in 84 cities across 42 countries.
-  Machine Learning
More machine learning happens on AWS than anywhere else. Machine learning in the hands of every developer and data scientist
-  Ecosystem
Tens of thousands of APN partners. The AWS Marketplace offers 50 categories, and more than 8,000 software listings
-  Enterprise leader
AWS positioned as a Leader in the Gartner Magic Quadrant for Cloud Infrastructure as a Service, Worldwide

Pricing Philosophy

High volume / low margin businesses are in our core DNA

Trade CapEX for
variable expense

Pay for what
you use

Our economies of
scale provide us
with lower costs

85 price
reductions
since 2006

Pricing model
choice to support
variable and stable
workloads

On-demand
Reserved Instances
Spot

Save more money as
you grow bigger

Tiered pricing
Volume discounts
Custom pricing

Customer obsessed



90%

of roadmap originates with customer requests and are designed to meet specific needs



“Performance, reliability, and responsiveness are fundamental to our customer experience, and T3 instances help us to deliver on that customer promise while also controlling our costs.”

—Heroku

It's greener in the cloud.

AWS's infrastructure is

3.6x more energy efficient

than the median of the surveyed U.S. enterprise data centers

AWS performs the same task with an

88% lower carbon footprint

Source: 451 Research, 2019, all rights reserved

Reducing water used for cooling

AWS has multiple initiatives to improve our water use efficiency for cooling data centers:

- **Evaporative cooling**
- **Reduce potable water usage**
- **Recycled cooling water**
- **Invest in reclaimed water infrastructure**





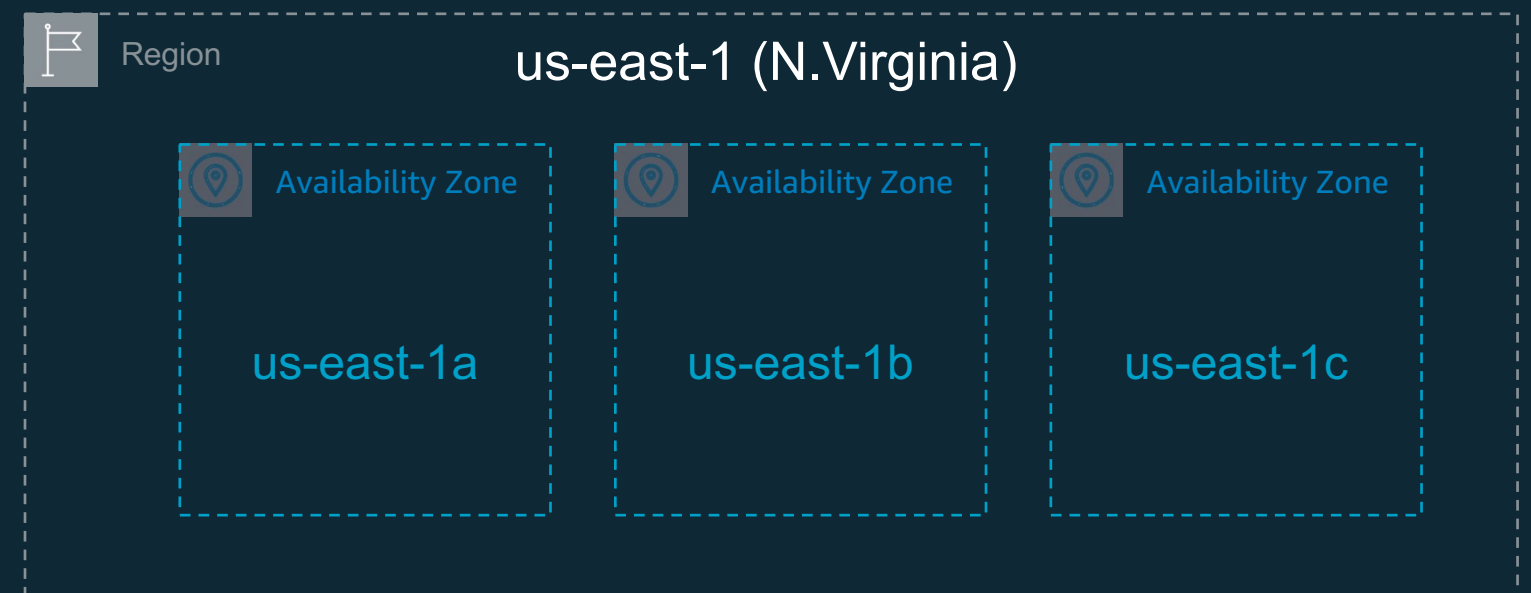
AWS Global Reach

25
Regions



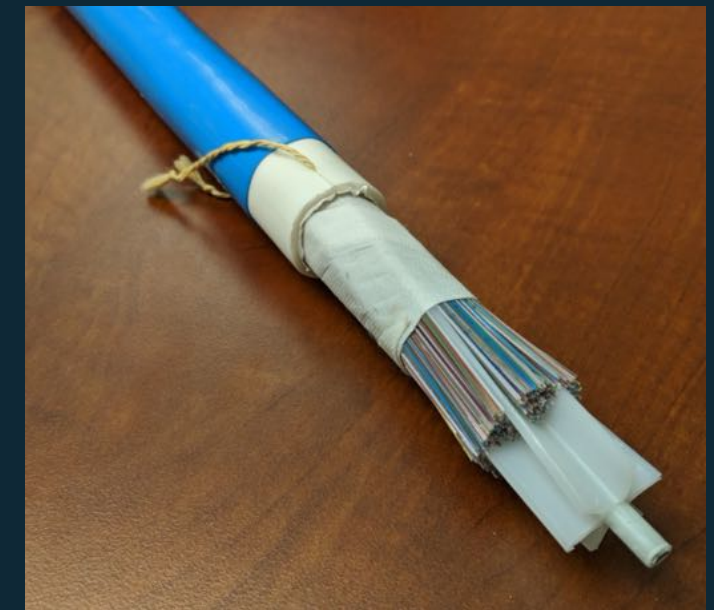
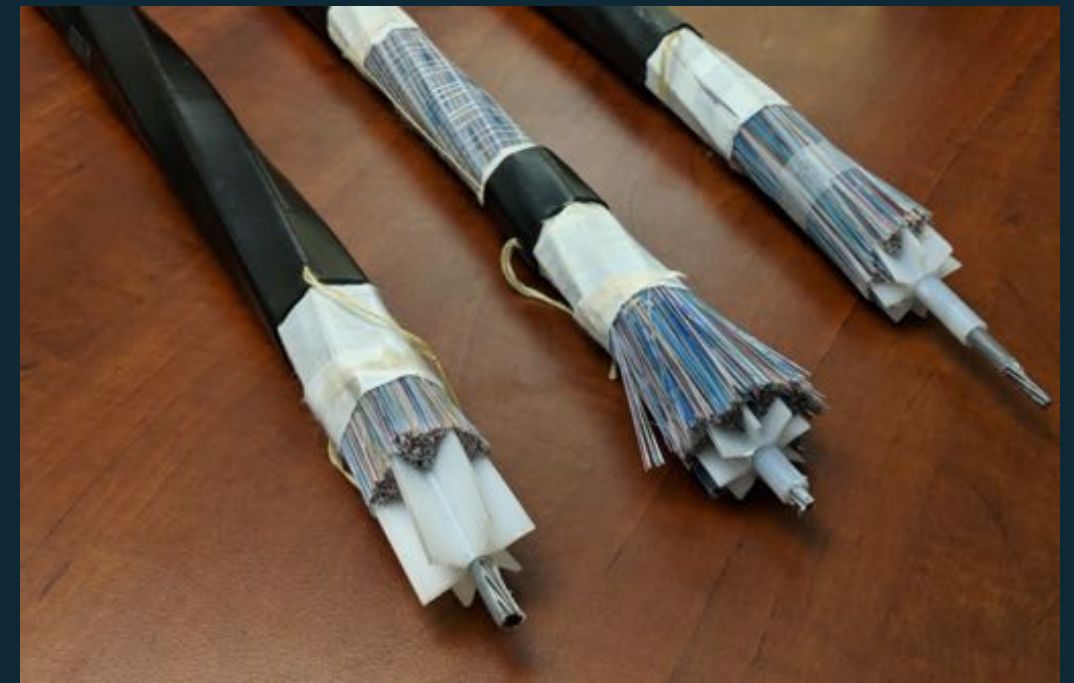
Availability Zones

- A region is comprised of multiple Availability Zones (typically 3)
- An Availability Zone (AZ) is one or more discrete data centers with redundant power, networking, and connectivity in an AWS Region
- High throughput, low latency (<10mS) network between Availability Zones
- All traffic between AZ's is encrypted
- Physical Separation < 100km



Intra & inter-AZ connectivity

- Dark fiber “spans”
 - Optimized for low-latency & physical diversity
- Amazon controlled infrastructure
- Geospatial coordinates
- Dense wavelength division multiplexing (DWDM)



275

Amazon
CloudFront
Points of
Presence



108

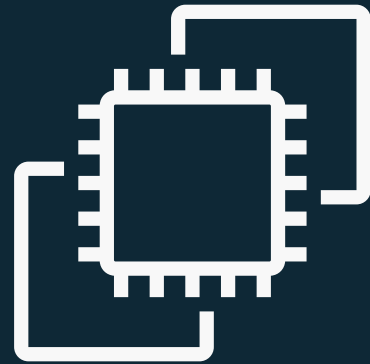
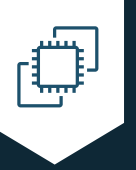
AWS Direct
Connect
locations



2

EC2 Overview

Choices for Compute



Amazon EC2

Virtual server instances
in the cloud



Amazon ECS, EKS, and Fargate

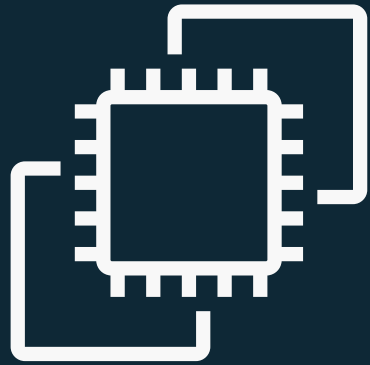
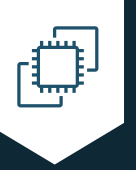
Container management service
for running
Docker on a managed
cluster of EC2



AWS Lambda

Serverless compute
for stateless code execution in
response to triggers

Amazon EC2



Amazon EC2

Linux | Windows | Mac

Arm and x86 architectures

General purpose and workload optimized

Bare metal, disk, networking capabilities

Packaged | Custom | Community AMIs

Multiple purchase options: On-demand, RI, Spot

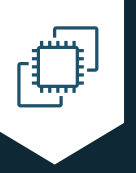


What's a virtual CPU? (vCPU)

- A vCPU is typically a hyper-threaded physical core*
 - Divide vCPU count by 2 to get core count
 - On Linux, "A" threads enumerated before "B" threads
 - On Windows, threads are interleaved
-
- Cores by Amazon EC2 & RDS DB Instance type:
<https://aws.amazon.com/ec2/virtualcores/>

** CPU Optimizing options allow disabling hyperthreading and reduce number of cores*

Memory and Storage



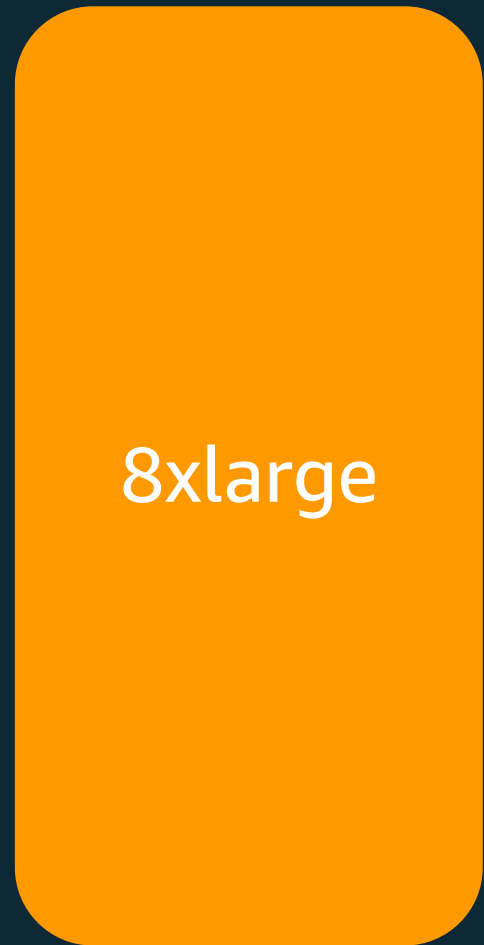
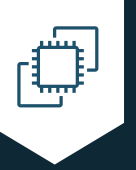
What's a GiB?

- Memory is presented as GibiBytes (GiB) and not Gigabytes (GB)
- 256 GiB = 275 GB

What about storage?

- Storage is independent of compute
- You allocate drives known as EBS volumes
- Max 16 TiB per volume
- Some instance types provide physically attached (ephemeral) storage

Instance sizing



8xlarge

c4.8xlarge

≈



4xlarge

2 - c4.4xlarge

≈



4xlarge

≈



2xlarge



2xlarge



2xlarge



2xlarge

4 - c4.2xlarge



xlarge



xlarge



xlarge



xlarge



xlarge



xlarge



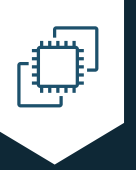
xlarge



xlarge

8 - c4.xlarge

EC2 Naming Explained



Instance generation

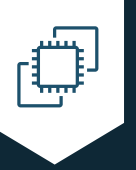
c5n.xlarge

Instance
family

Attribute

Instance size

Instance Types



General Purpose

Compute Optimized

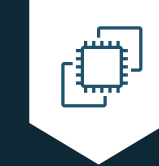
Memory Optimized

Accelerated Computing

Storage Optimized

	General Purpose		Compute Optimized		Memory Optimized				Accelerated Computing			Storage Optimized		
	Burstable performance	General Purpose	Compute Intensive	Compute +memory up to 100 Gbps	Memory Optimized	In-memory	Memory Intensive	Compute and Memory Intensive	Graphics Intensive	General Purpose GPU	FPGA	High I/O	Dense Storage	Big Data Optimized
intel	T3	M5	C5	C5n	R5	X1	X1e		G3	P3	F1		D2	H1
Local storage (NVMe SSD)		M5d	C5d		R5d			Z1d				I3		
AMD	T3a	M5a			R5a									
metal		M5m	c5m		R5m		u-12tb1	Z1dm				I3m		
others	A1	M6g	C6g		R6g					P3dn		I3en		
	arm													

Choose your processor and architecture



Intel® Xeon® Scalable
(Skylake) processor



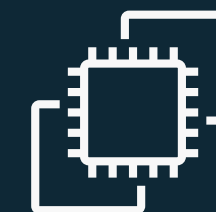
NVIDIA V100
Tensor Core GPUs



AMD EPYC processor



AWS Graviton
Processor (arm)



FPGAs for custom
hardware acceleration

Right compute for the right application and workload

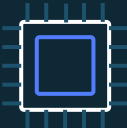
AWS Graviton2 Processor

Enabling the best price/performance for your cloud workloads

Graviton Processor



First Arm-based processor available in major cloud



Built on 64-bit Arm Neoverse cores with AWS-designed silicon using 16 nm manufacturing technology



Up to 16 vCPUs, 10 Gbps enhanced networking, 3.5 Gbps EBS bandwidth

Graviton2 Processor



7x performance, 4x compute cores, and 5x faster memory

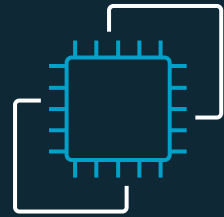


Built with 64-bit Arm Neoverse cores with AWS-designed silicon using 7 nm manufacturing technology



Up to 64 vCPUs, 25 Gbps enhanced networking, 18 Gbps EBS bandwidth

Choice of accelerators for specialized workloads



Elastic Graphics

Easily add graphics acceleration to your EC2 instance

Configure right amount of graphics acceleration for your workload

Accelerate application for fraction of cost of standalone graphics instances



Elastic Inference

Reduce deep learning inference costs by up to 75%

Easily attach fractional sizes of a full GPU instance to EC2 or SageMaker instances

Scale inference acceleration up or down as needed with EC2 Auto Scaling

What is an Amazon Machine Image (AMI)?



Provides the information required to launch an instance

Launch multiple instances from a single AMI

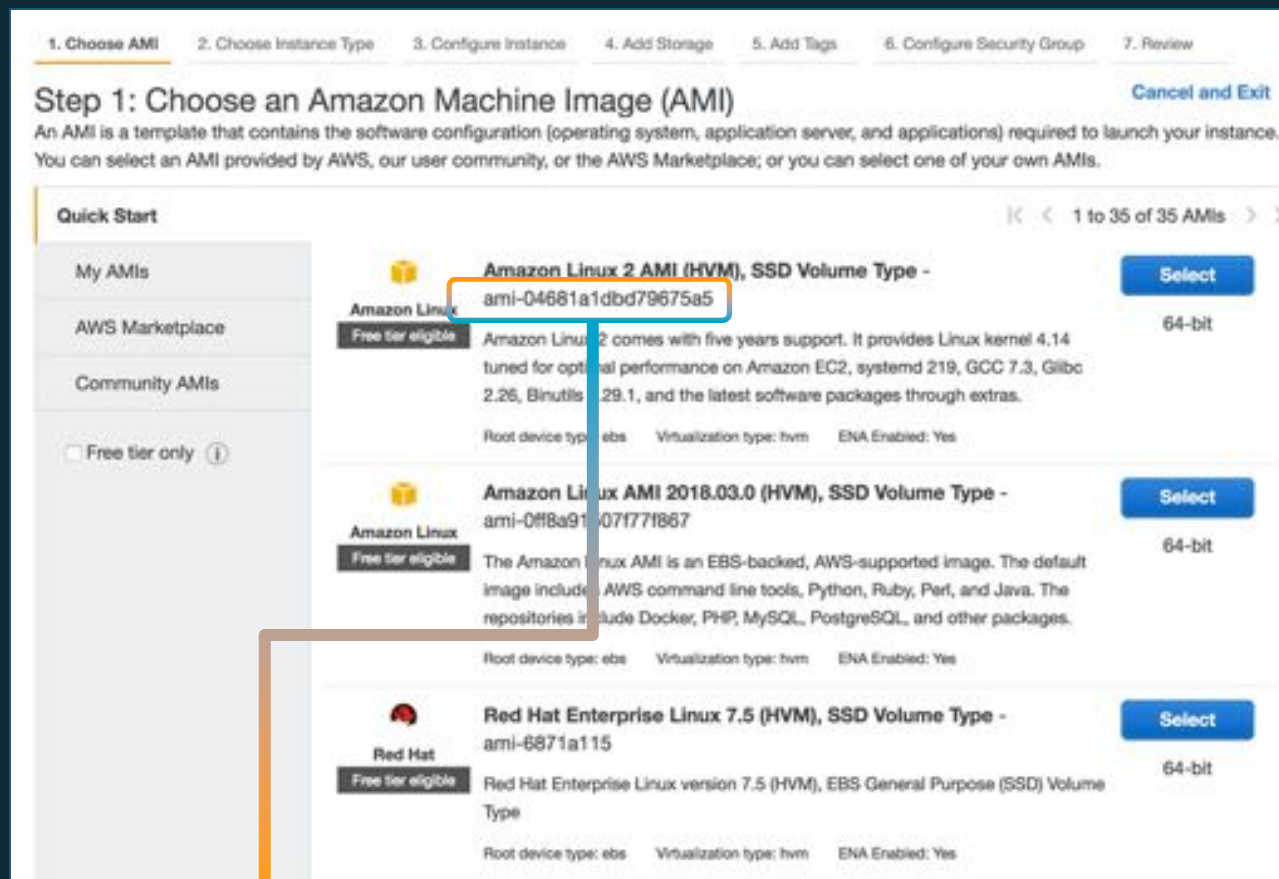
An AMI includes the following

- A template for the root volume (for example, operating system, applications)
- Launch permissions that control which AWS accounts can use the AMI
- Block device mapping that specifies volumes to attach to the instance

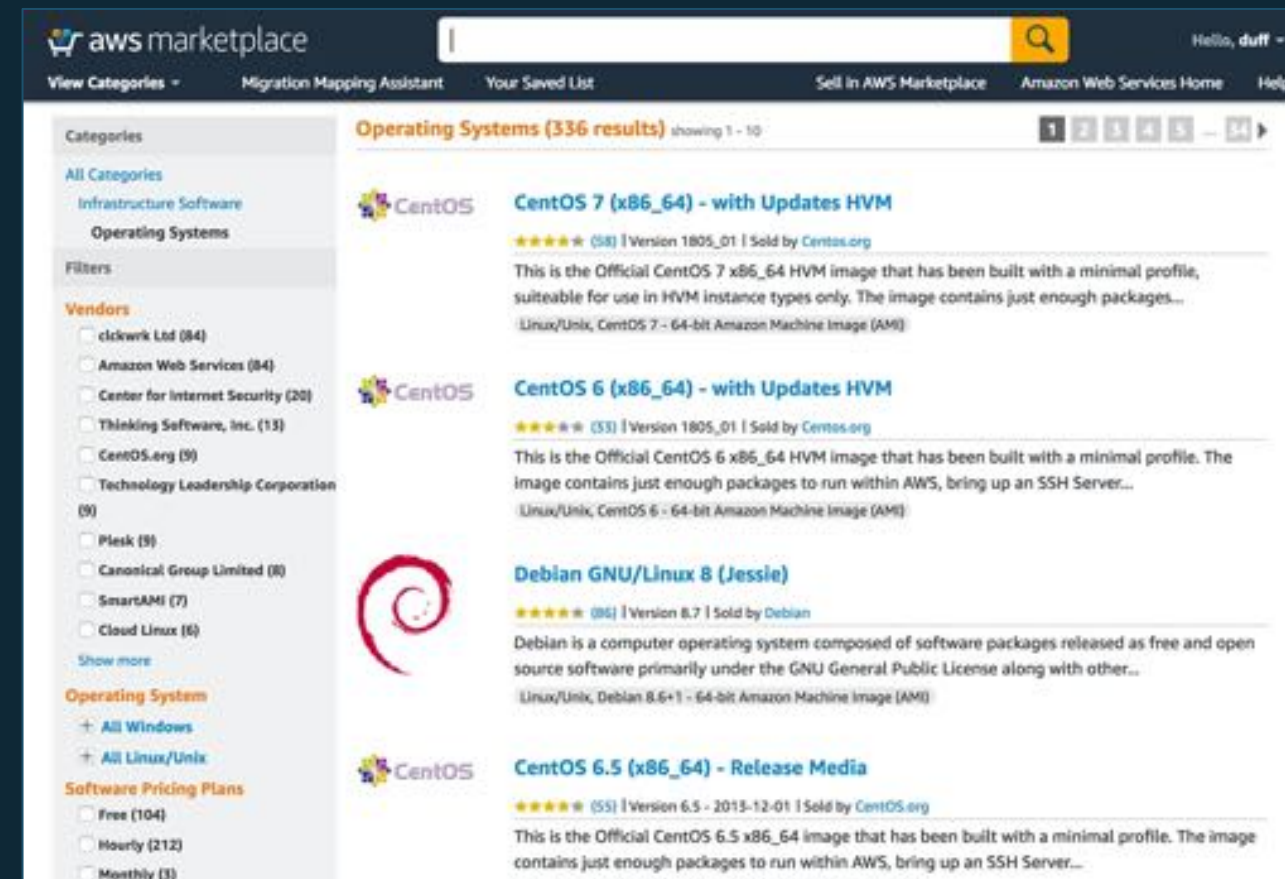
Choosing an AMI



AWS Console



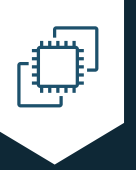
AWS Marketplace



Use the AMI ID to launch through the API or AWS Command Line Interface (AWS CLI)

```
aws ec2 run-instances --image-id ami-04681a1dbd79675a5 --instance-type c4.8xlarge --count 10 --key-name MyKey
```





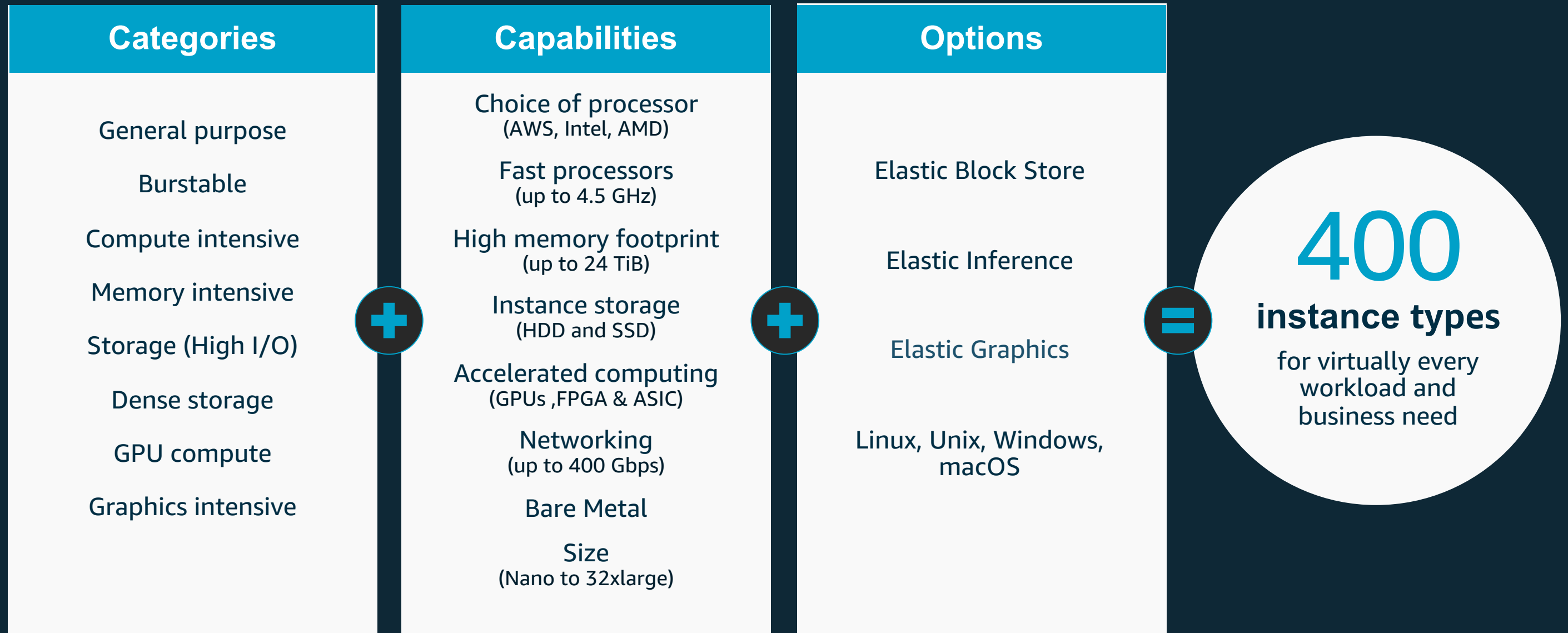
EC2 Operating Systems Supported

- Windows 2003R2*/2008*/2008R2*/2012/2012R2/2016/2019
- Amazon Linux
- Debian
- Suse
- CentOS
- Red Hat Enterprise Linux
- Ubuntu
- Mac



for more OSes see: <https://aws.amazon.com/marketplace/b/2649367011>

Broadest and deepest platform choice



Amazon EC2 purchase options

On-Demand

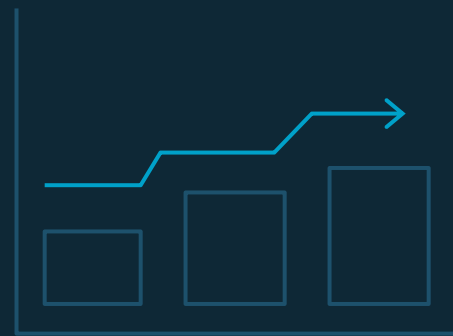
Pay for compute capacity by **the second** with no long-term commitments



Spiky workloads, to define needs

Reserved Instances

Make a 1 or 3 year commitment and receive a **significant discount** off On-Demand prices



Committed and steady-state usage

Savings Plan

Same great discounts as Amazon EC2 RIs with **more flexibility**



Committed flexible access to compute

Spot Instances

Spare Amazon EC2 capacity at **savings of up to 90%** off On-Demand prices



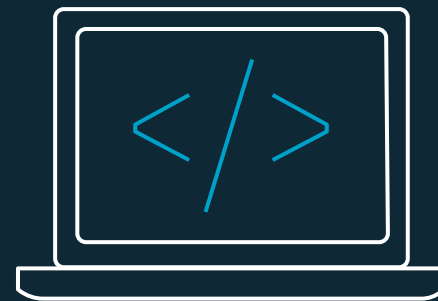
Fault-tolerant, flexible, stateless workloads

Hibernate Amazon EC2 Instances

Maintain a fleet of pre-warmed instances to quickly get to a productive state



Available with Amazon EBS-backed instances



Use familiar Stop and Start APIs



Memory data saved in EBS root volume



RAM contents are encrypted on EBS

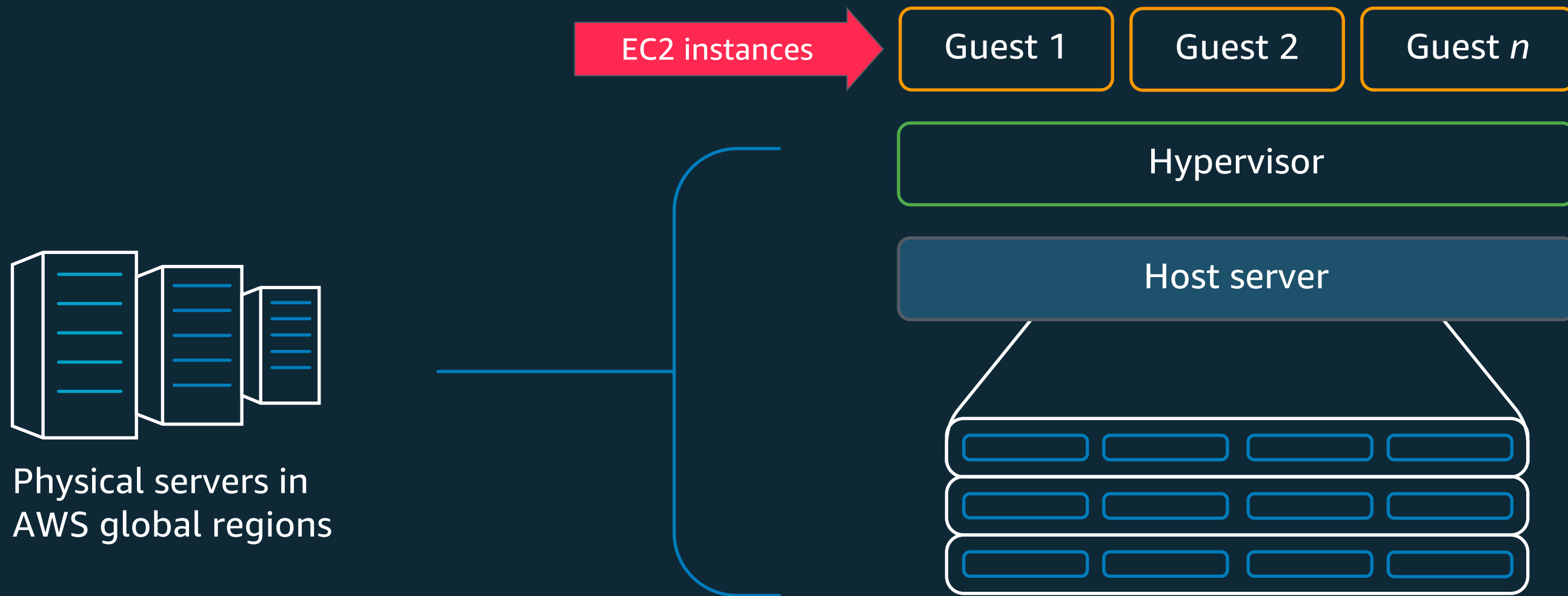
Its just like closing and opening your laptop!

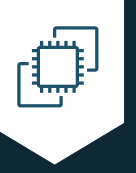
Applications can pick up right where it left off

3

EC2 Design

EC2 Host Virtualization

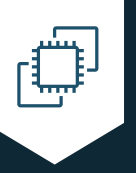




Resource allocation

- All resources assigned to you are dedicated to your instance with no over commitment*
 - All vCPUs are dedicated to you
 - Memory allocated is assigned only to your instance
 - Network resources are partitioned to avoid “noisy neighbors”
- Curious about the number of instances per host?
 - See “Dedicated Hosts Configuration Table” for a guide.

*the “T” family is special



Which hypervisor do we use?

Original host architecture: **Xen-based**

- Hypervisor consumed resources from the underlying host
- Limited optimization

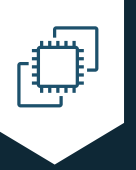
AWS Nitro Hypervisor: **Custom KVM based hypervisor**

- AWS Nitro System (launched on Nov 2017)
- Less server resources used, more resources for the customer
- AWS optimized

Bare metal: **Direct access to processor and memory resources**

- Built on the AWS Nitro system
- Enables custom hypervisors and micro-VM runtimes

AWS Nitro System

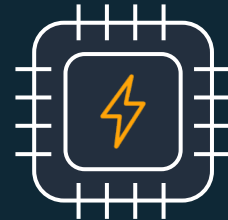


Nitro Card



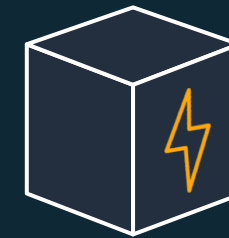
Local NVMe storage
Elastic Block Storage
Networking, monitoring,
and security

Nitro Security Chip



Integrated into motherboard
Protects hardware resources

Nitro Hypervisor



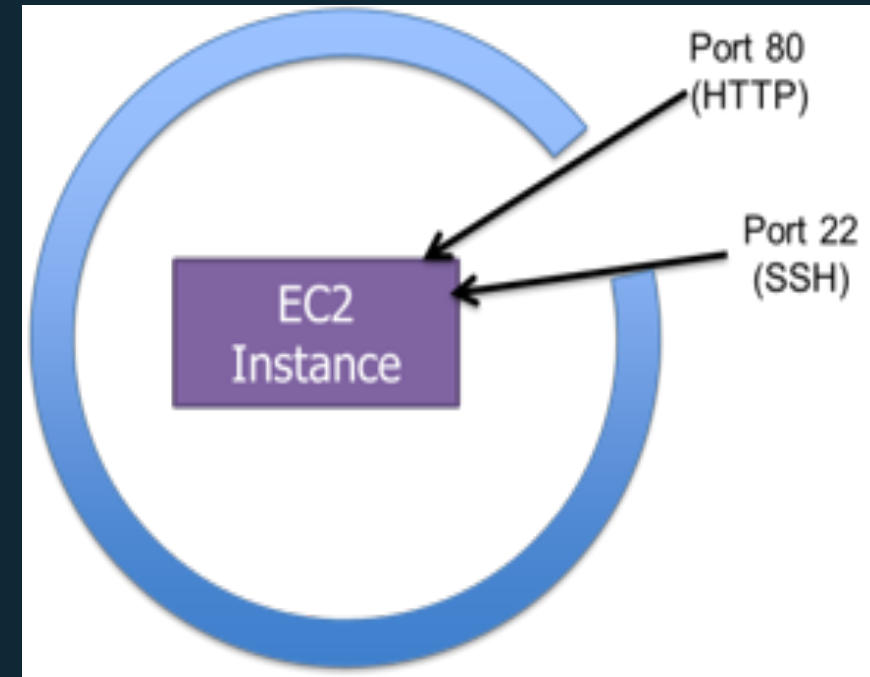
Lightweight hypervisor
Memory and CPU allocation
Bare metal-like performance

Modular building blocks for rapid design and delivery of **EC2** instances

EC2 Security Groups

Security Group Rules

- Name
- Description
- Protocol
- Port range
- IP address, IP range, Security Group name



EC2-Specific Credentials

EC2 key pairs

- Linux – SSH key pair for first-time host login
- Windows – Retrieve Administrator password

Standard SSH RSA key pair

- Public/Private Keys
- Private keys are not stored by AWS

AWS approach for providing initial access to a generic OS

- Secure
- Personalized
- Non-generic (NIST, PCI DSS)



“Public Half” inserted by Amazon into each EC2 instance that you launch



“Private Half” downloaded to your desktop

Any Questions?

Up Next: EC2 Hands on Lab