

Τ

### **OFFICE HOURS - FALL 2025**

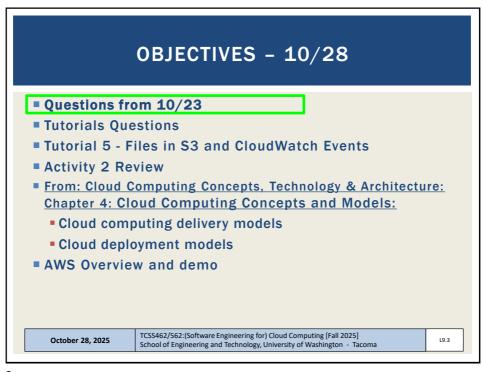
- Thursdays:
  - •6:00 to 7:00 pm CP 229 & Zoom
- <u>Friday \*\*\* THIS WEEK \*\*\*</u>
  - ■11:00 am to 12:00 pm ONLINE via Zoom
- Or email for appointment
- > Office Hours set based on Student Demographics survey feedback
- \* Friday office hours may be adjusted or canceled due meeting conflicts or other obligations. Adjustments will be announced via Canvas.

October 28, 2025

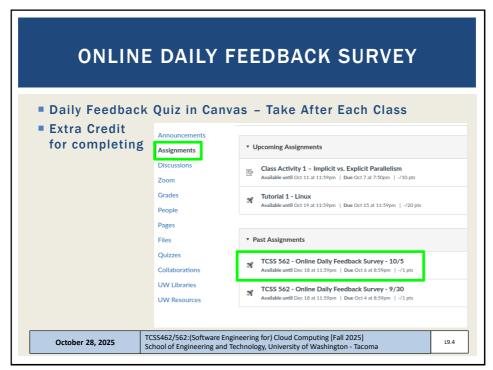
TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

9.2

2



3



4

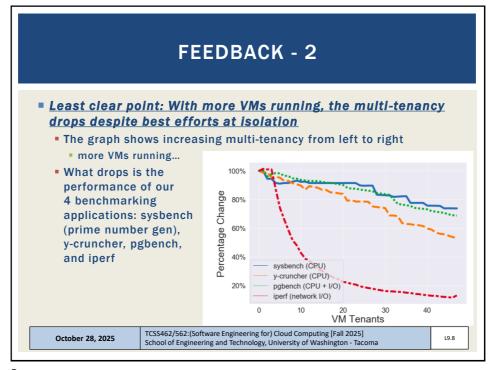
	Started	S 562 : Oct 7 at : z Instr	l:13am		Daily	Feedb	ack S	Surve	y - 10	/5			
		Questi	on 1								0.5 pts		
		On a so	cale of 1	l to 10, p	olease cl	assify yo	ur persp	ective o	n mater	ial cove	red in today's		
		1	2	3	4	5	6	7	8	9	10		
		Mostly Review			Ne	Equal w and Rev	view				Mostly New to Me		
		Questi	on 2								0.5 pts		
		Please	rate the	pace of	today's	class:							
		1	2	3	4	5	6	7	8	9	10		
		Slow			J	ust Right				F	ast		
October 2	28, 202	5	TC: Sch	SS462/5 nool of E	62:(Soft	ware Eng ng and T	gineering echnolog	g for) Clo gy, Unive	oud Compersity of V	puting [F Washing	Fall 2025] ton - Tacoma	L9.5	

5

# MATERIAL / PACE ■ Please classify your perspective on material covered in today's class (43 respondents, 25 in-person, 18 online): ■ 1-mostly review, 5-equal new/review, 10-mostly new ■ Average - 7.19 (↓ - previous 7.32) ■ Please rate the pace of today's class: ■ 1-slow, 5-just right, 10-fast ■ Average - 5.49 (↑ - previous 5.24) October 28, 2025 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] | School of Engineering and Technology, University of Washington - Tacoma

6

■ Why do	es the average	\$/hour pri	e within	a date ran	ge change
- wny ao	<u>es the average</u>	→/nour pri	e within	<u>a date ran</u>	ge cnange
OVON +h					g
	ough the same	_	•	-	
The a	verage price sho	ws is shown f	or each ava	ailability zon	ne within
	• .			•	
	gion over the tin	ie span of gr	ipii (3 nrs,	T day, T we	ek, I
mont	n, 3 months)				
	,				
ot Instance pricing history					
r instance type requirements, budget requ	rements, and application design will determine how to apply th		more, see Spot Instance Best Practices [		
	rements, and application design will determine how to apply the instance type    Schlarge	following best practices for your application. To le  Platform  Linux/i		Date range  Uniter range	
r instance type requirements, budget requ ph valiability Zones	Instance type	Platform		Date range	
r instance type requirements, budget requ ph valiability Zones ses	Instance type	Platform		Date range	
r instance type requirements, budget requ ph valiability Zones es	Instance type	Platform		Date range	
instance type requirements, budget requirements, bu	Instance type	Platform		Date range	
r instance type requirements, budget req. ph valiability Zones 255 100	Instance type	Platform		Date range  1 week	:-2a) \$0.0256 (0.0128 per vCPU)
r Instance type requirements, budget rep ph variability Zones es 100 Oct 21 0 0	Instance type	Flatform ▼ Lines/N		Date range	t-2a) \$0.0256 (0.0128 per vCPU) t-2b) \$0.0284 (0.0142 per vCPU)
Instance type requirements, budget requirements and provide requirement		Platform  ▼ Linux/1  Linux/1  On 14 On 15 000 11.00	X Del 25 Oct 25	Date range	t-2a) \$0.0256 (0.0128 per vCPU) t-2b) \$0.0284 (0.0142 per vCPU)
r instance type requirements, budget requirements,		Platform  ▼ Linux/1  Linux/1  On 14 On 15 000 11.00	X Del 25 Oct 25	Date range	t-2a) \$0.0256 (0.0128 per vCPU) t-2b) \$0.0284 (0.0142 per vCPU)
Instance type requirements, budget requirements, bu	### ##################################	Platform  ▼ Linux/1  Linux/1  On 14 On 15 000 11.00	X Del 25 Oct 25	Date range	t-2a) \$0.0256 (0.0128 per vCPU) t-2b) \$0.0284 (0.0142 per vCPU)
Instance type requirements, budget requirements, bu	### ##################################	Platform  ▼ Linux/1  Linux/1  On 14 On 15 000 11.00	X 0x25 0x25 0x25 0x00 1200	Date range	t-2a) \$0.0256 (0.0128 per vCPU) t-2b) \$0.0284 (0.0142 per vCPU)



8

### FEEDBACK - 3

- What is "session state" in a client/server application?
  - AWS Lambda functions are considered "stateless"
- What is "stateless"?
  - Each time a client calls an AWS Lambda function, requests are routed to a random function instance (worker) to process the call
  - If a client makes multiple calls, there is no guarantee it will run in the same function instance each time
  - AWS Lambda functions feature static global memory, but if client calls do not return to the same function instance, this memory can't be used to store session state
- Key Design Issue With Serverless Applications:
  - How do you persist session state on AWS Lambda?

October 28, 2025

TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L9.9

9

### **PRACTICE QUESTION 1**

- Which of the following can lead to performance problems for application hosting on cloud platforms?
- A. Resource sharing/contention
- B. Cloud consumer under-provisioning
- C. Heterogeneous hardware
- D. Cloud provider over-provisioning
- E. All of the above

October 28, 2025

TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L9.10

10

# PRACTICE QUESTION 2

- Which cloud computing delivery model often requires manual configuration to provide resource elasticity?
- A. Platform-as-a-Service
- B. Infrastructure-as-a-Service
- C. Serverless Database
- D. Function-as-a-Service
- E. All of the above

October 28, 2025

TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L9.11

11

### **PRACTICE QUESTION 3**

- On Amazon EC2, when using persistent spot requests, what occurs when you intentionally terminate the virtual machine?
- A. In addition to the virtual machine being deleted, the persistent spot request is also deleted
- B. VM termination is not supported using persistent spot requests
- C. Using the AWS management console, the user is prompted to enter a password prior to deletion of the virtual machine
- D. After a short delay, a replacement virtual machine is launched to satisfy the persistent spot request
- E. The virtual machine is stopped, not terminated, and can be later resumed without loss of data on the disk

October 28, 2025

TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L9.12

12

# OBJECTIVES - 10/28 Questions from 10/23 Tutorials Questions Tutorial 5 - Files in S3 and CloudWatch Events Activity 2 Review From: Cloud Computing Concepts, Technology & Architecture: Chapter 4: Cloud Computing Concepts and Models: Cloud computing delivery models Cloud deployment models AWS Overview and demo

13

# TUTORIAL 3 – OCT 30 (TEAMS OF 2) Best Practices for Working with Virtual Machines on Amazon EC2 https://faculty.washington.edu/wlloyd/courses/tcss562 /tutorials/TCSS462\_562\_f2025\_tutorial\_3.pdf Creating a spot VM Creating an image from a running VM Persistent spot request Stopping (pausing) VMs EBS volume types Ephemeral disks (local disks) Mounting and formatting a disk Disk performance testing with Bonnie++ Cost Saving Best Practices TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

14

## **TUTORIAL 4 - NOV 7**

Introduction to AWS Lambda with the Serverless Application Analytics Framework (SAAF)

https://faculty.washington.edu/wlloyd/courses/tcss562/tutorials/TCSS462\_562\_f2025\_tutorial\_4.pdf

- Set up Java development environment
- Introduction to Maven build files for Java
- Create and Deploy "hello" Java AWS Lambda Function
- Create API Gateway REST endpoint
- Sequential testing of "hello" AWS Lambda Function
  - API Gateway endpoint, AWS Lambda CLI Function invocation, AWS Function URL
- Profiling function performance with SAAF
- Concurrent function testing with faas\_runner
- Performance analysis using faas\_runner reports
- Two function pipeline development task: Caesar Cipher

October 28, 2025

TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L9.15

15

## **OBJECTIVES - 10/28**

- Questions from 10/23
- Tutorials Questions
- Tutorial 5 Files in S3 and CloudWatch Events
- Activity 2 Review
- From: Cloud Computing Concepts, Technology & Architecture: Chapter 4: Cloud Computing Concepts and Models:
  - Cloud computing delivery models
  - Cloud deployment models
- AWS Overview and demo

October 28, 2025

TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L9.16

16

## **TUTORIAL 5 - TO BE POSTED**

- Introduction to Lambda II: Working with Files in S3 and CloudWatch Events
- Customize the Request object (add getters/setters)
  - Why do this instead of HashMap?
- Import dependencies (jar files) into project for AWS S3
- Create an S3 Bucket
- Give your Lambda function(s) permission to work with S3
- Write to the CloudWatch logs
- Use of CloudTrail to generate S3 events
- Creating CloudWatch rule to capture events from CloudTrail
- Have the CloudWatch rule trigger a target Lambda function with a static JSON input object (hard-coded filename)
- Optional: for the S3 PutObject event, dynamically extract the name of the file put to the S3 bucket for processing

October 27, 2022

TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2022] School of Engineering and Technology, University of Washington - Tacoma

L9.17

17

## **OBJECTIVES - 10/28**

- Questions from 10/23
- Tutorials Questions
- Tutorial 5 Files in S3 and CloudWatch Events

### Activity 2 Review

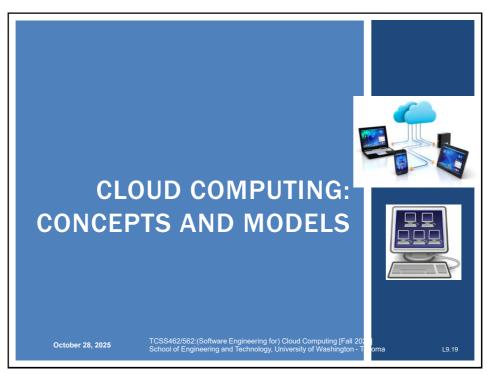
- From: Cloud Computing Concepts, Technology & Architecture: Chapter 4: Cloud Computing Concepts and Models:
  - Cloud computing delivery models
  - Cloud deployment models
- AWS Overview and demo

October 28, 2025

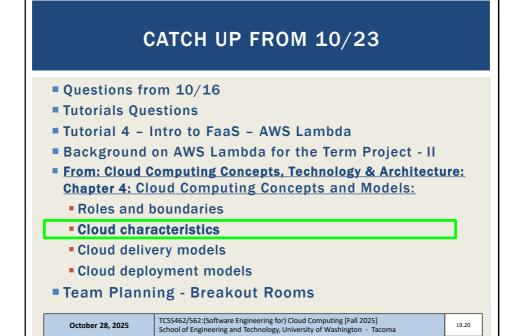
TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L9.18

18



19



20

# **CLOUD CHARACTERISTICS**

- On-demand usage
- Ubiquitous access
- Multitenancy (resource pooling)
- Elasticity
- Measured usage
- Resiliency
- Assessing these features helps measure the value offered by a given cloud service or platform

October 17, 2024

TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L7.21

21

### **MEASURED USAGE**

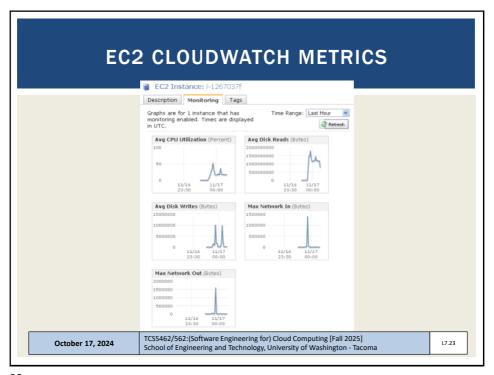
- Cloud platform tracks usage of IT resources
- For billing purposes
- Enables charging only for IT resources actually used
- Can be time-based (millisec, second, minute, hour, day)
  - Granularity is increasing...
- Can be throughput-based (data transfer: MB/sec, GB/sec)
- Can be resource/reservation based (vCPU/hr, GB/hr)
- Not all measurements are for billing
- Some measurements can support auto-scaling
- For example CPU utilization

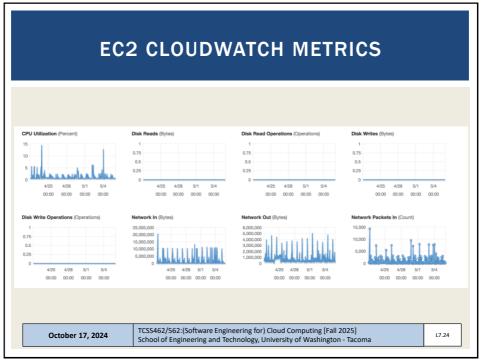
October 17, 2024

TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

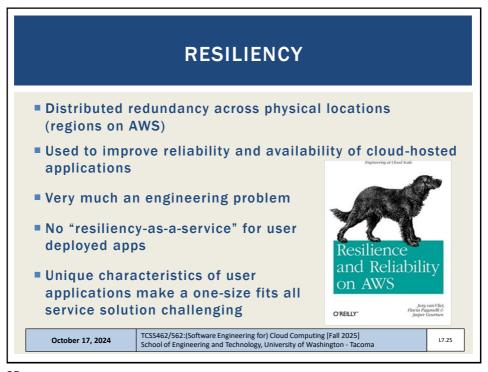
L7.22

22

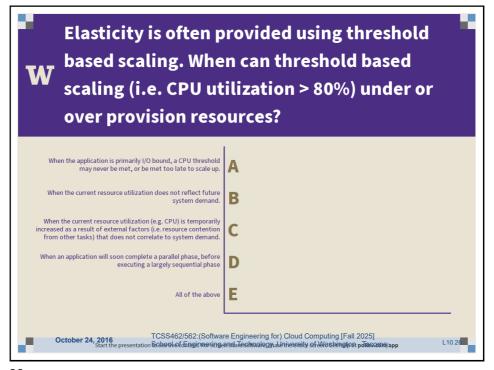




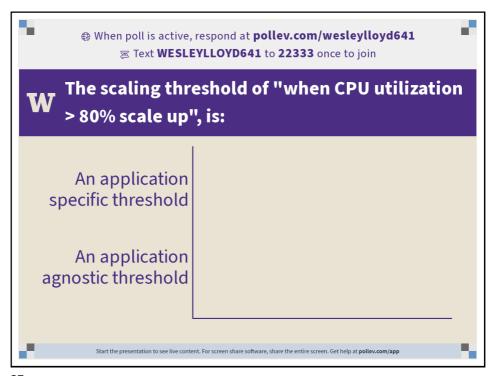
24



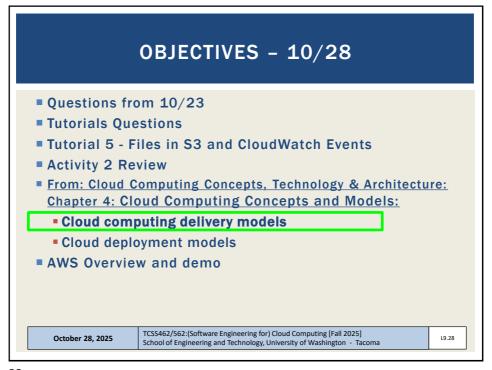
25



26



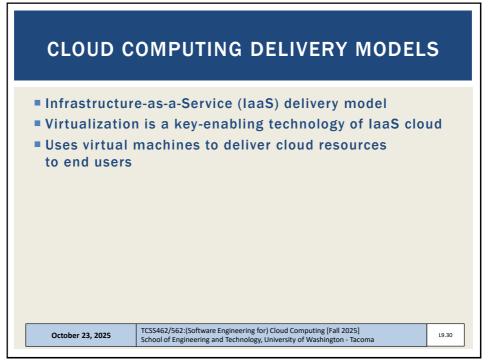
27



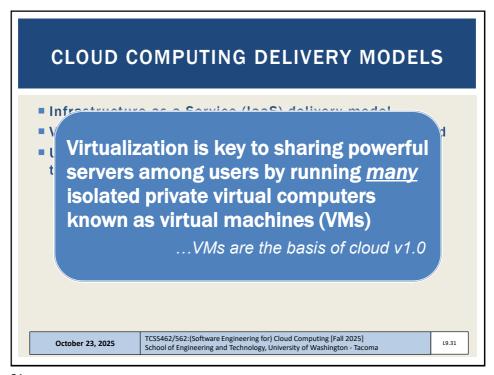
28

# CLOUD COMPUTING DELIVERY MODELS Infrastructure-as-a-Service (IaaS) Platform-as-a-Service (PaaS) Software-as-a-Service (SaaS) Serverless Computing: Function-as-a-Service (FaaS) Container-as-a-Service (CaaS) Other Delivery Models TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

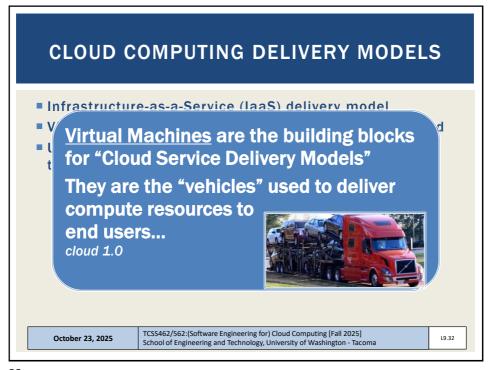
29



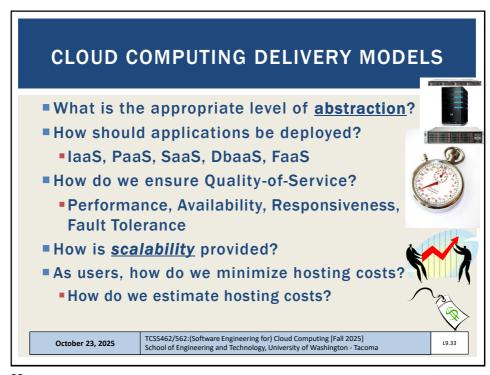
30



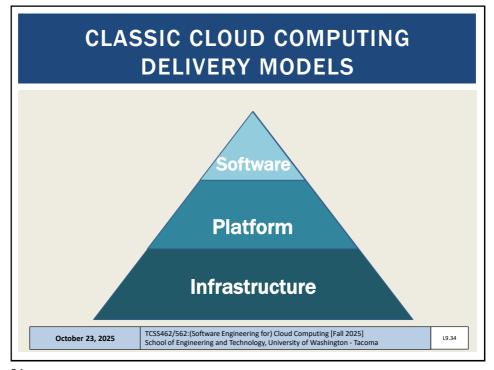
31



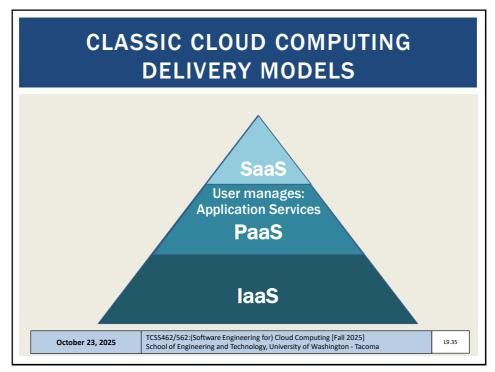
32

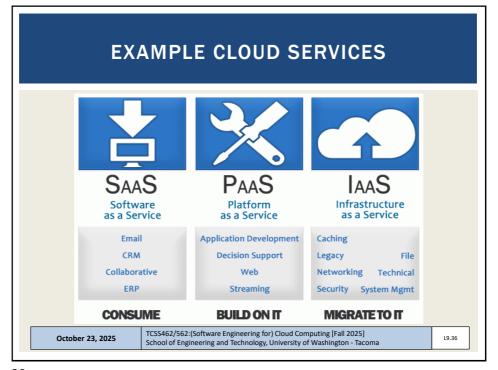


33

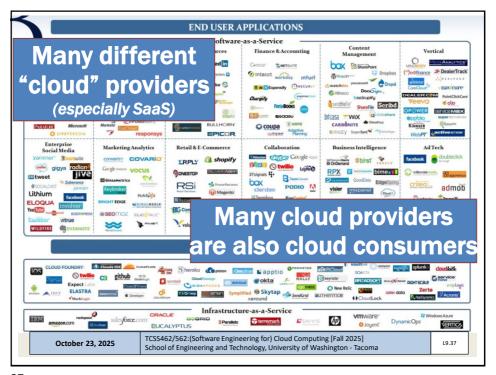


34





36

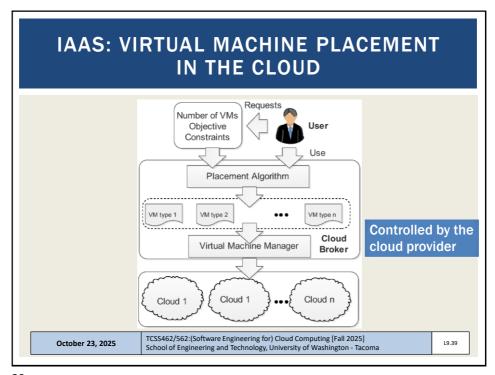


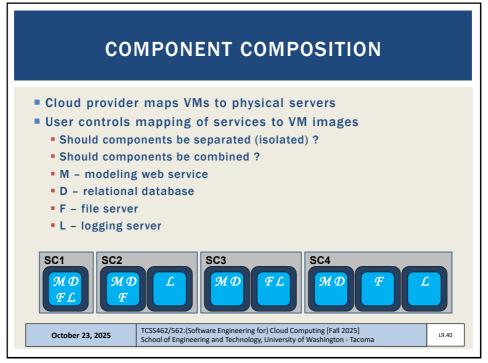


- Compute resources, on demand, as-a-service
  - Generally raw "IT" resources
  - Hardware, network, containers, operating systems
- Typically provided through virtualization
- Generally, not-preconfigured
- Administrative burden is owned by cloud consumer
- Best when high-level control over environment is needed
- Scaling is generally <u>not</u> automatic...
- Resources can be managed in bundles
- AWS CloudFormation: Scripts to specify creation of cloud infrastructures using JSON/YAML for app deployment

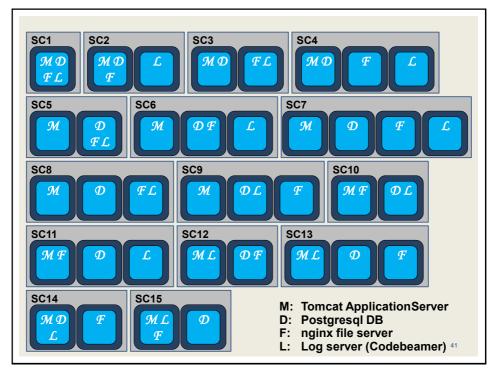
October 23, 2025 TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

38

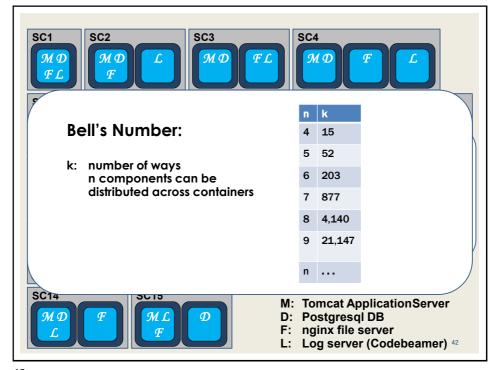




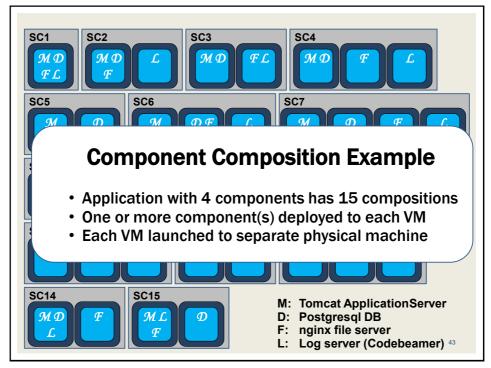
40

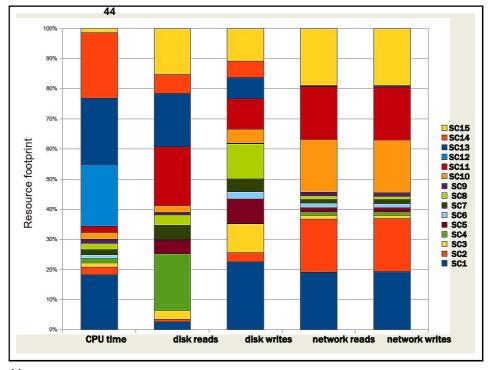


41

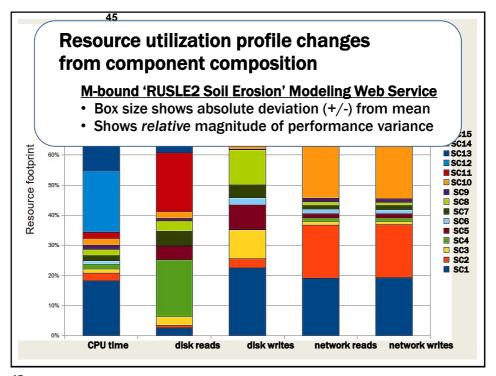


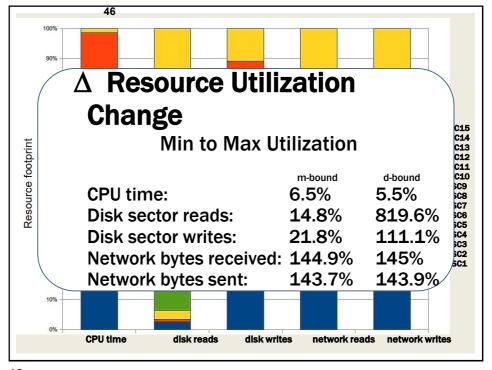
42



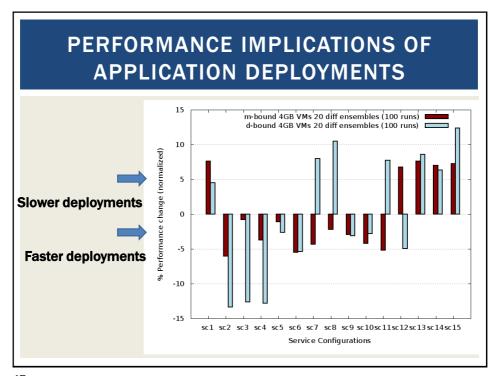


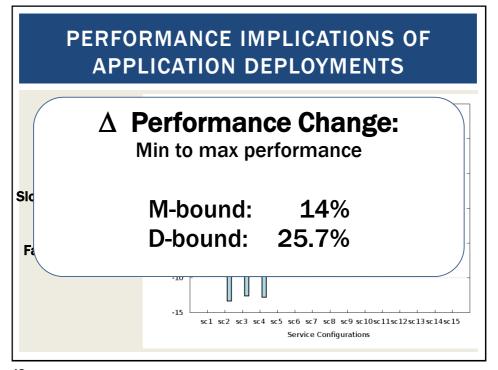
44



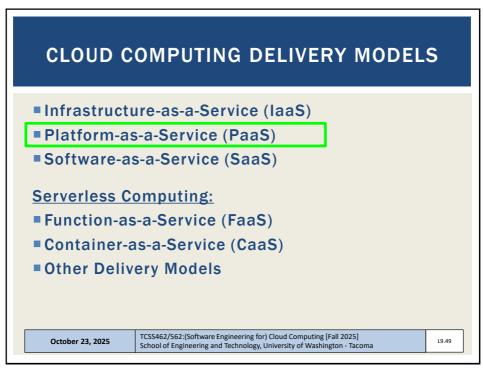


46

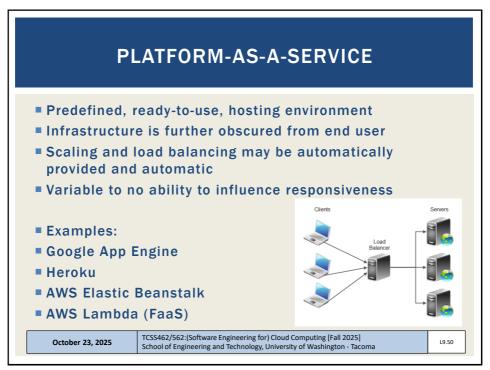




48



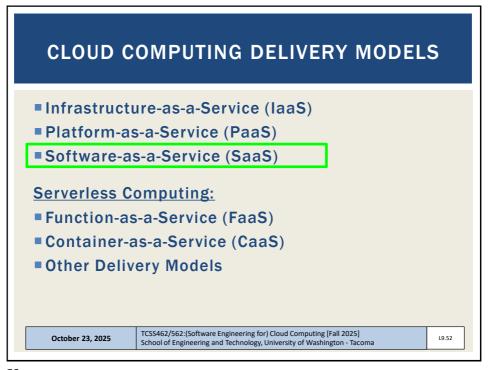
49



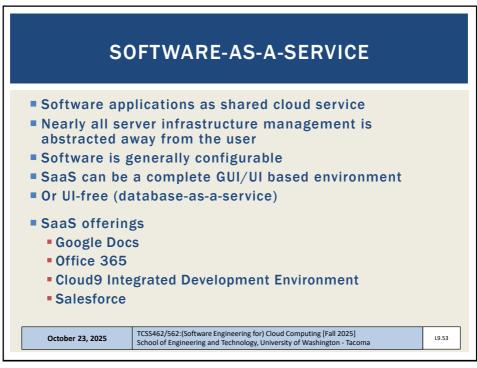
50

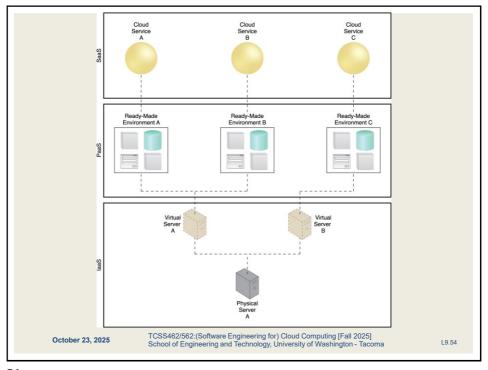
# USES FOR PAAS Cloud consumer Wants to extend on-premise environments into the cloud for "web app" hosting Wants to entirely substitute an on-premise hosting environment Cloud consumer wants to become a cloud provider and deploy its own cloud services to external users PaaS spares IT administrative burden compared to laaS TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

51

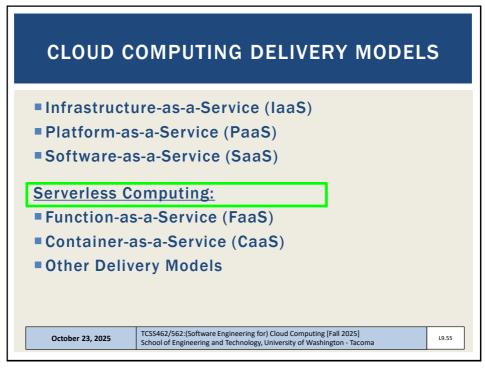


52





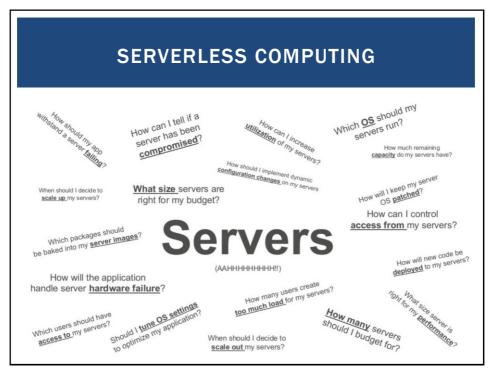
54

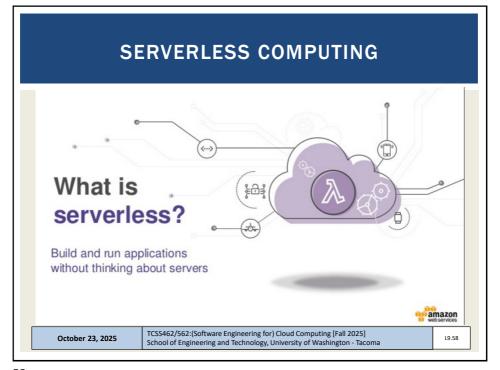


55

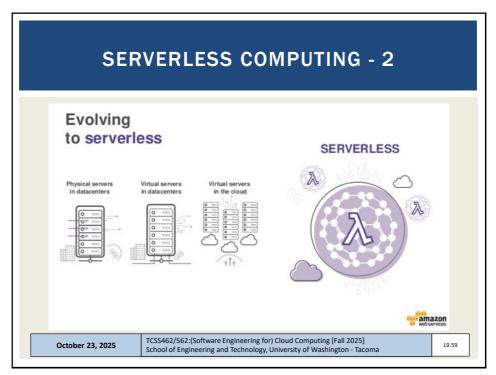


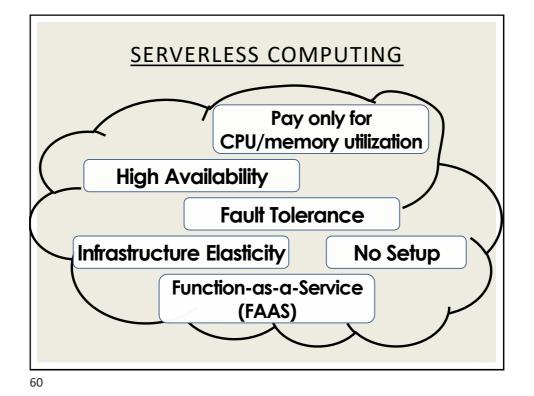
56





58





# SERVERLESS COMPUTING

# **Why Serverless Computing?**

Many features of distributed systems, that are challenging to deliver, are provided automatically

...they are built into the platform

61

### **CLOUD COMPUTING DELIVERY MODELS**

- ■Infrastructure-as-a-Service (laaS)
- Platform-as-a-Service (PaaS)
- Software-as-a-Service (SaaS)

## **Serverless Computing:**

- Function-as-a-Service (FaaS)
- Container-as-a-Service (CaaS)
- Other Delivery Models

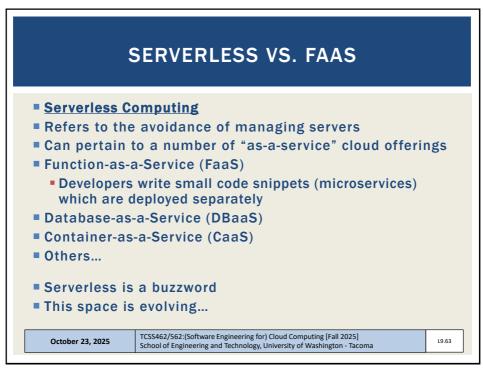
October 23, 2025

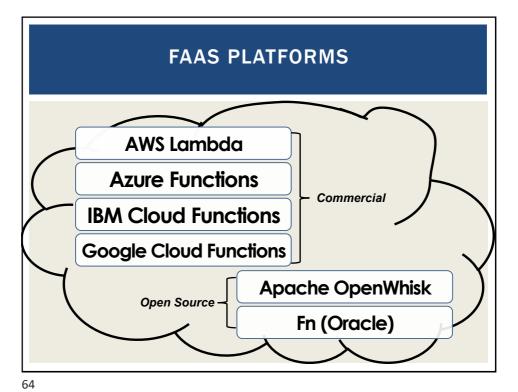
TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

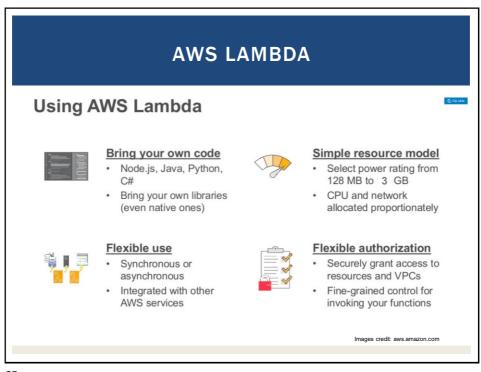
-

Slides by Wes J. Lloyd L9.31

62







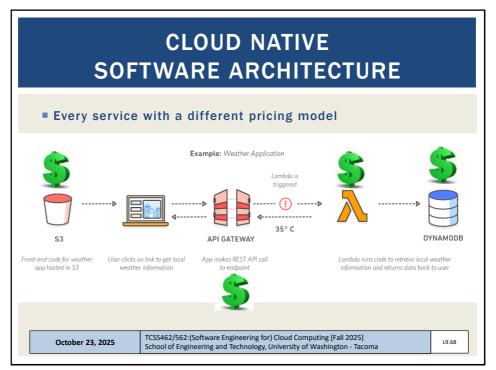
65

# FAAS PLATFORMS - 2 New cloud platform for hosting application code Every cloud vendor provides their own: AWS Lambda, Azure Functions, Google Cloud Functions, IBM OpenWhisk Similar to platform-as-a-service Replace opensource web container (e.g. Apache Tomcat) with abstracted vendor-provided black-box environment October 23, 2025 TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

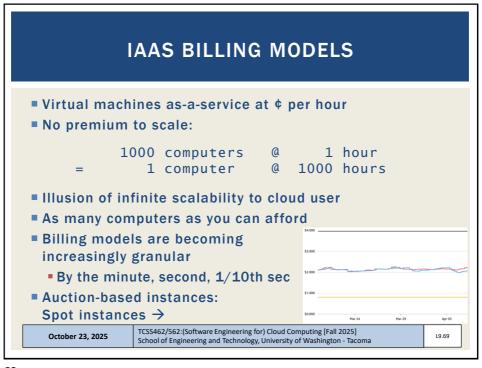
66

# FAAS PLATFORMS - 3 Many challenging features of distributed systems are provided automatically Built into the platform: Highly availability (24/7) Scalability Fault tolerance

67



68



69

# PRICING OBFUSCATION ■ VM pricing: hourly rental pricing, billed to nearest second is intuitive... non-intuitive pricing policies FaaS pricing: • FREE TIER: first 1,000,000 function calls/month $\rightarrow$ FREE first 400,000 GB-sec/month → FREE Afterwards: obfuscated pricing (AWS Lambda): \$0.0000002 per request \$0.00000208 to rent 128MB / 100-ms \$0.00001667 GB /second TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma October 23, 2025

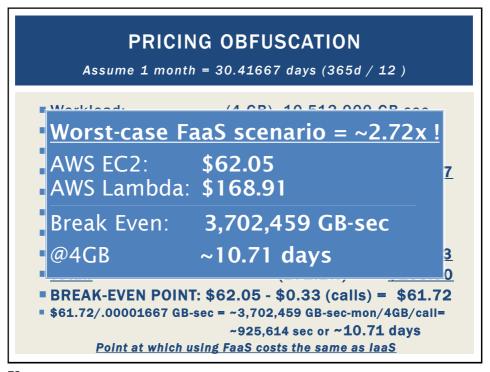
70

# CLOUD COMPUTING DELIVERY MODELS Infrastructure-as-a-Service (IaaS) Platform-as-a-Service (PaaS) Software-as-a-Service (SaaS) Serverless Computing: Function-as-a-Service (FaaS) Container-as-a-Service (CaaS) Other Delivery Models ICSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

71

### WEBSERVICE HOSTING EXAMPLE ON AWS Lambda ■ Each service call: 100% of 2 CPU-cores 100% of 4GB of memory Workload: uses 2 continuous threads Duration: 1 month (30.41667 days) Amazon EC2 c5.large 2-vCPU VM x 4GB ON AWS EC2: ■ c5.large: 8.5¢/hour, 24 hrs/day x 30.41667 days Hosting cost: \$62.05/month How much would hosting this workload cost on AWS Lambda? TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma October 22, 2024

72

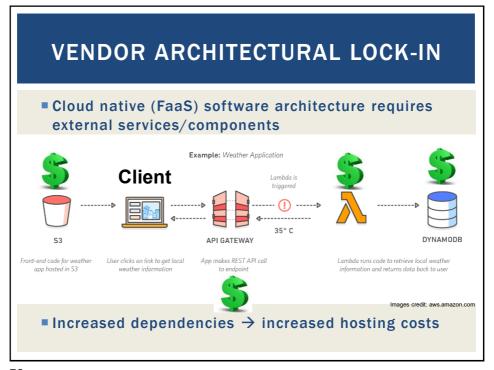


## FAAS PRICING Break-even point is the point where renting VMs or deploying to a serverless platform (e.g. Lambda) is exactly the same. Our example is for one month Could also consider one day, one hour, one minute What factors influence the break-even point for an application running on AWS Lambda? October 28, 2025 TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

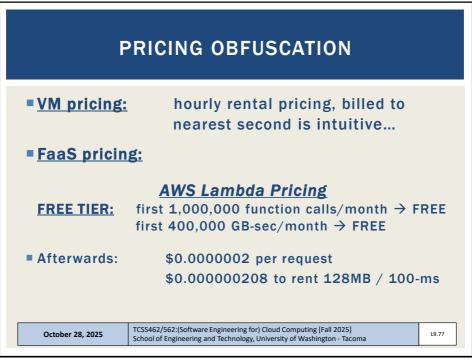
74

## FAAS CHALLENGES Vendor architectural lock-in – how to migrate? Pricing obfuscation – is it cost effective? Memory reservation – how much to reserve? Service composition – how to compose software? Infrastructure freeze/thaw cycle – how to avoid? Performance – what will it be? Ccober 28, 2025 TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

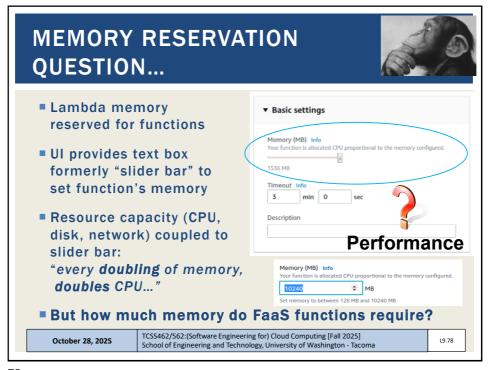
75



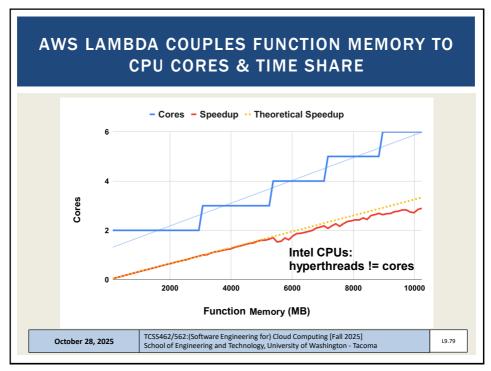
76

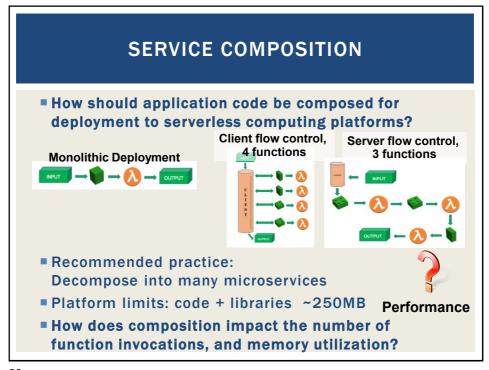


77

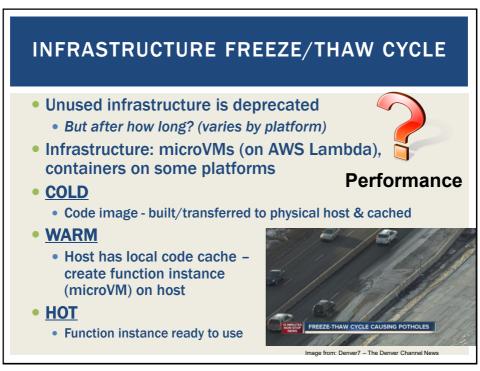


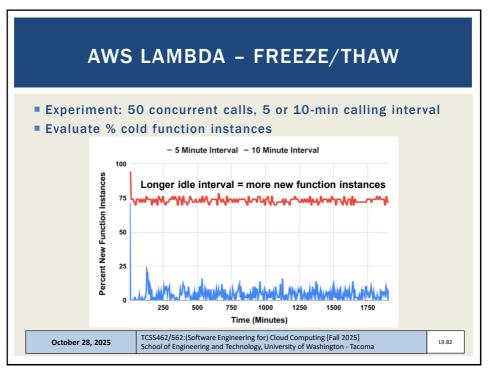
78





80





82

### FACTORS IMPACTING PERFORMANCE OF FAAS COMPUTING PLATFORMS

- Infrastructure scaling/elasticity
- Resource contention (CPU, network, memory caches)
- Hardware heterogeneity (CPU types, hyperthread, etc)
- Load balancing / provisioning variation
- Infrastructure retention: COLD vs. WARM
  - Infrastructure freeze/thaw cycle
- Function memory reservation size
- Application service composition

October 28, 2025

TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

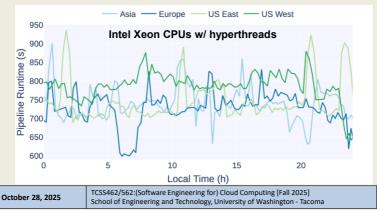
L9.83

L9.84

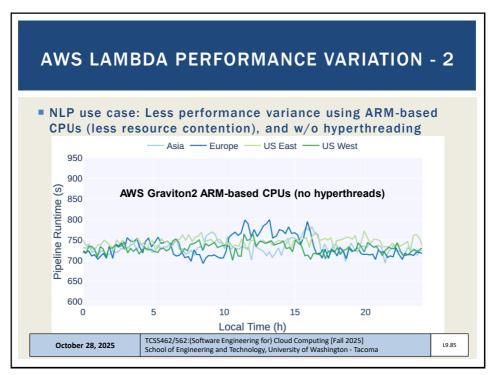
83

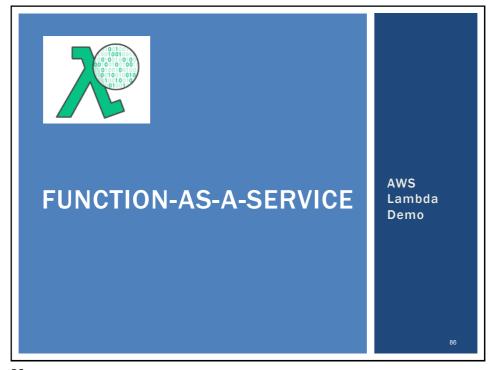
### AWS LAMBDA PERFORMANCE VARIATION

- NLP processing pipeline use case
- Performance variance from: diurnal changes in load (e.g. resource contention), Intel hyperthreading



84





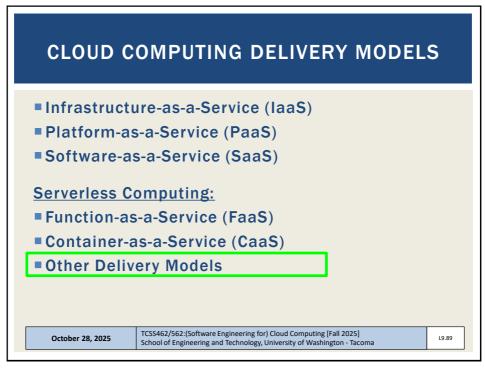
86

## CLOUD COMPUTING DELIVERY MODELS Infrastructure-as-a-Service (IaaS) Platform-as-a-Service (PaaS) Software-as-a-Service (SaaS) Serverless Computing: Function-as-a-Service (FaaS) Container-as-a-Service (CaaS) Other Delivery Models Ccober 28, 2025 CCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

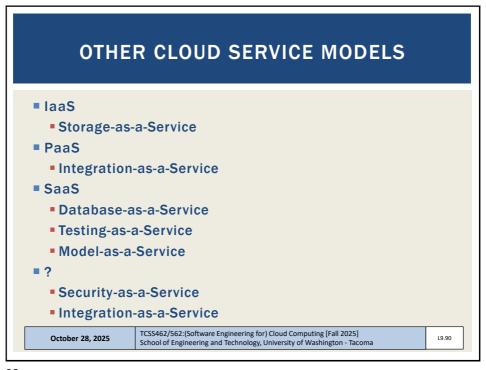
87

### **CONTAINER-AS-A-SERVICE** Cloud service model for deploying application containers (e.g. Docker containers) to the cloud Deploy containers without worrying about managing infrastructure: Servers (virtual machines) Or container orchestration platforms Container platform examples: Kubernetes, Docker swarm, Apache Mesos/Marathon, Amazon Elastic Container Service Container platforms support creation of container clusters on the using cloud hosted VMs CaaS Examples: AWS Fargate Google Cloud Run Azure Container Instances TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma October 28, 2025 19 88

88



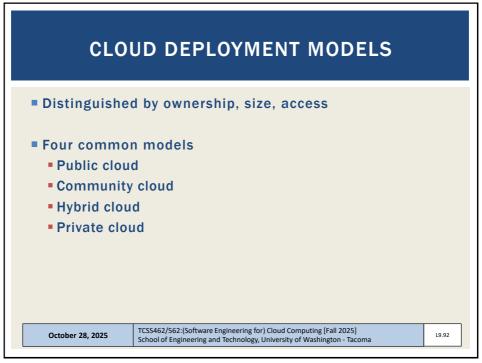
89



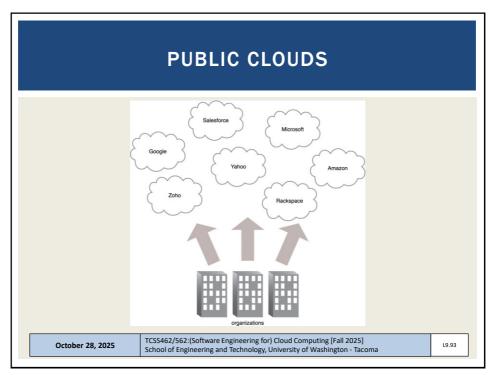
90

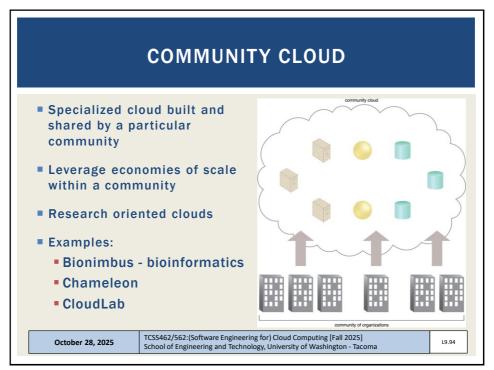
# OBJECTIVES - 10/28 Questions from 10/23 Tutorials Questions Tutorial 5 - Files in S3 and CloudWatch Events Activity 2 Review From: Cloud Computing Concepts, Technology & Architecture: Chapter 4: Cloud Computing Concepts and Models: Cloud computing delivery models Cloud deployment models AWS Overview and demo Ctober 28, 2025 TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

91

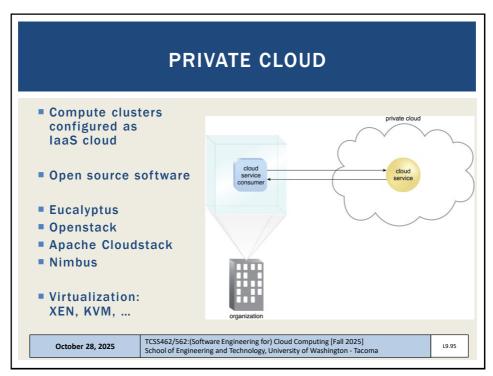


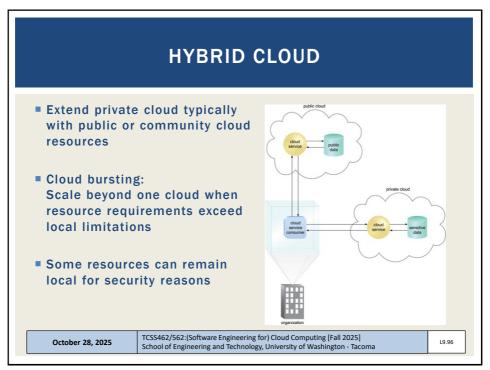
92





94





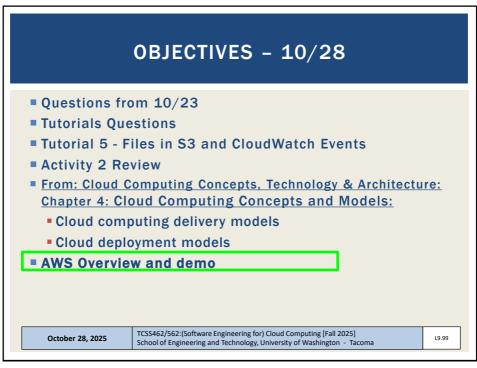
96

## Federated cloud Simply means to aggregate two or more clouds together Hybrid is typically private-public Federated can be public-public, private-private, etc. Also called inter-cloud Virtual private cloud Google and Microsoft simply call these virtual networks Ability to interconnect multiple independent subnets of cloud resources together Resources allocated private IPs from individual network subnets can communicate with each other (10.0.1.0/24) and (10.0.2.0/24) Subnets can span multiple availability zones within an AWS region Cotober 28, 2025 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] | School of Engineering and Technology, University of Washington - Tacoma

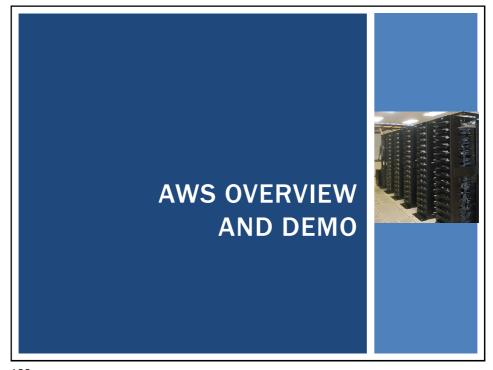
97



98



99



100

### **ONLINE CLOUD TUTORIALS**

- From the eScience Institute @ UW Seattle: https://escience.washington.edu/
- Online cloud workshops
- Introduction to AWS, Azure, and Google Cloud
- Task: Deploying a Python DJANGO web application
- Self-guided workshop materials available online:
- https://cloudmaven.github.io/documentation/
- AWS Educate provides access to many online tutorials / learning resources:
- https://aws.amazon.com/education/awseducate/

October 28, 2025

TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L9.101

101

### **LIST OF TOPICS**

- AWS Management Console
- Elastic Compute Cloud (EC2)
- Instance Storage: Virtual Disks on VMs
- Elastic Block Store: Virtual Disks on VMs
- Elastic File System (EFS)
- Amazon Machine Images (AMIs)
- EC2 Paravirtualization
- EC2 Full Virtualization (hvm)
- EC2 Virtualization Evolution

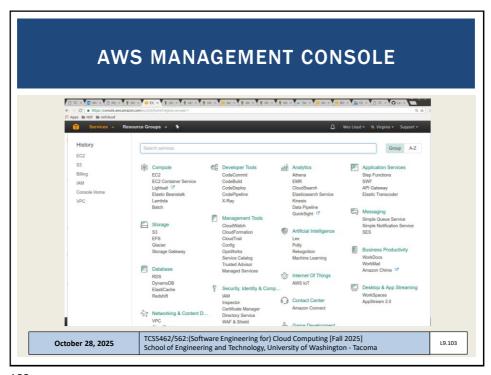
- (VM) Instance Actions
- EC2 Networking
- EC2 Instance Metadata Service
- Simple Storage Service (S3)
- AWS Command Line Interface (CLI)
- Legacy / Service Specific CLIs
- AMI Tools
- Signing Certificates
- Backing up live disks
- Cost Savings Measures

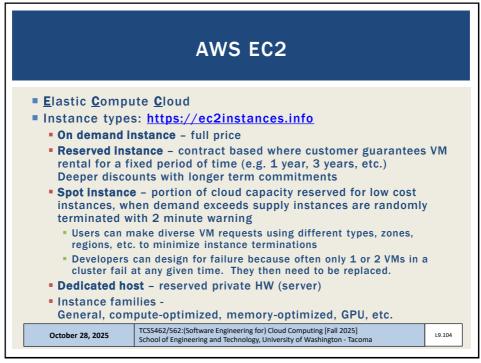
October 28, 2025

TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L9.10

102





104

### **AWS EC2 - 2**

- Storage types
  - Instance storage ephemeral storage
    - Temporary disk volumes stored on disks local to the VM
    - Evolution: physical hard disk drives (HDDs)
    - Solid state drives (SSDs)
    - Non-volatile memory express (NVMe) drives (closer to DRAM speed)
  - EBS Elastic block store
    - Remotely hosted disk volumes
  - EFS Elastic file system
    - Shared file system based on network file system
    - VMs, Lambdas, Containers mount/interact with shared file system
    - Somewhat expensive

October 28, 2025

TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L9.105

105

### **INSTANCE STORAGE**

- Also called ephemeral storage
- Persisted using images saved to \$3 (simple storage service)
  - ~2.3¢ per GB/month on S3
  - 5GB of free tier storage space on S3
- Requires "burning" an image
- Multi-step process:
  - Create image files
  - Upload chunks to \$3
  - Register image
- Launching a VM
  - Requires downloading image components from S3, reassembling them...
     is potentially slow
- VMs with instance store backed root volumes not pause-able
- Historically root volume limited to 10-GB max- faster imaging...

October 28, 2025

TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L9.106

106

### **ELASTIC BLOCK STORE**

- EBS provides 1 drive to 1 virtual machine (1:1) (not shared)
- EBS cost model is different than instance storage (uses S3)
  - ~10¢ per GB/month for General Purpose Storage (GP2)
  - ~8¢ per GB/month for General Purpose Storage (GP3)
  - 30GB of free tier storage space
- EBS provides "live" mountable volumes
  - Listed under volumes
  - <u>Data volumes</u>: can be mounted/unmounted to any VM, dynamically at any time
  - Root volumes: hosts OS files and acts as a boot device for VM
  - In Linux drives are linked to a mount point "directory"
- Snapshots back up EBS volume data to S3
  - Enables replication (required for horizontal scaling)
  - EBS volumes not actively used should be snapshotted, and deleted to save EBS costs...

October 28, 2025

TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L9.107

107

### **EBS VOLUME TYPES - 2**

- Metric: I/O Operations per Second (IOPS)
- General Purpose 2 (GP2)
  - 3 IOPS per GB, min 100 IOPS (<34GB), max of 16,000 IOPS</li>
  - 250MB/sec throughput per volume
- General Purpose 3 (GP3 new Dec 2020)
  - Max 16,000 IOPS, Default 3,000 IOPS
  - GP2 requires creating a 1TB volume to obtain 3,000 IOPS
  - GP3 all volumes start at 3000 IOPS and 125 MB/s throughput
  - 1000 additional IOPS beyond 3000 is \$5/month up to 16000 IOPS
  - 125 MB/s additional throughput is \$5/month up to 1000 MB/s throughput

October 28, 2025

TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L9.108

108

### **EBS VOLUME TYPES - 3**

- Provisioned IOPS (IO1)
  - Legacy, associated with GP2
  - Allows user to create custom disk volumes where they pay for a specified IOPS and throughput
  - 32,000 IOPS, and 500 MB/sec throughput per volume MAX
- Throughput Optimized HDD (ST1)
  - Up to 500 MB/sec throughput
  - 4.5 ¢ per GB/month
- Cold HDD (SC1)
  - Up to 250 MB/sec throughput
  - 2.5 ¢ per GB/month
- Magnetic
  - Up to 90 MB/sec throughput per volume
  - 5 ¢ per GB/month

October 28, 2025

TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L9.109

109

### **ELASTIC FILE SYSTEM (EFS)**

- EFS provides 1 volume to many client (1:n) shared storage
- Network file system (based on NFSv4 protocol)
- Shared file system for EC2, Fargate/ECS, Lambda
- Enables mounting (sharing) the same disk "volume" for R/W access across multiple instances at the same time
- Different performance and limitations vs. EBS/Instance store
- Implementation uses abstracted EC2 instances
- ~ 30 ¢ per GB/month storage default burstable throughput
- Throughput modes:
- Can modify modes only once every 24 hours
- Burstable Throughput Model:
  - Baseline 50kb/sec per GB
  - Burst 100MB/sec pet GB (for volumes sized 10GB to 1024 GB)
  - Credits .72 minutes/day per GB

October 28, 2025 TCSS462/562:

TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L9.110

110

ELASTIC FILE SYSTEM (EFS) - 2				
<ul> <li>Burstable Throughput Rates</li> <li>Throughput rates: baseline vs burst</li> <li>Credit model for bursting: maximum burst per day</li> </ul>				
File System Size (GiB)	Baseline Aggregate Throughput (MiB/s)	Burst Aggregate Throughput (MiB/s)	Maximum Burst Duration (Min/Day)	% of Time File System Can Burst (Per Day)
10	0.5	100	7.2	0.5%
256	12.5	100	180	12.5%
512	25.0	100	360	25.0%
1024	50.0	100	720	50.0%
1536	75.0	150	720	50.0%
2048	100.0	200	720	50.0%
3072	150.0	300	720	50.0%
4096	200.0	400	720	50.0%

### **ELASTIC FILE SYSTEM (EFS) - 3**

■ Throughput Models

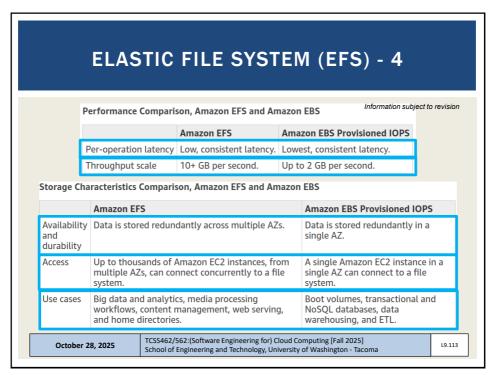
Information subject to revision

- Provisioned Throughput Model
- For applications with:
- high performance requirements, but low storage requirements
- Get high levels of performance w/o overprovisioning capacity
- \$6 MB/s-Month (Virginia Region)
  - Default is 50kb/sec for 1 GB, .05 MB/s = 30 ¢ per GB/month
- If file system metered size has higher baseline rate based on size, file system follows default Amazon EFS Bursting Throughput model
  - No charges for Provisioned Throughput below file system's entitlement in Bursting Throughput mode
  - Throughput entitlement = 50kb/sec per GB

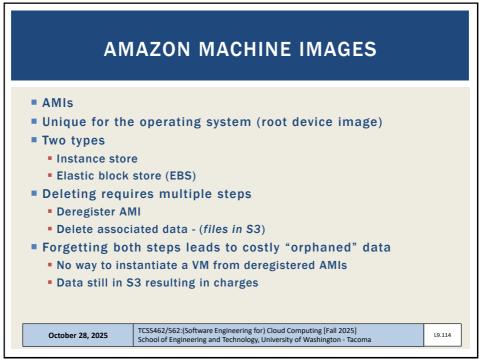
TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma October 28, 2025

L9.112

112



113



114

### **EC2 VIRTUALIZATION - PARAVIRTUAL**

- 1<sup>st</sup>, 2<sup>nd</sup>, 3<sup>rd</sup>, 4<sup>th</sup> generation → XEN-based
- 5<sup>th</sup> generation instances → AWS Nitro virtualization
- XEN two virtualization modes
- XEN Paravirtualization "paravirtual"
  - 10GB Amazon Machine Image base image size limit
  - Addressed poor performance of old XEN HVM mode
  - I/O performed using special XEN kernel with XEN paravirtual mode optimizations for better performance
  - Requires OS to have an available paravirtual kernel
  - PV VMs: will use common <u>AKI</u> files on AWS Amazon kernel image(s)
    - Look for common identifiers

October 28, 2025

TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L9.115

115

### **EC2 VIRTUALIZATION - HVM**

- XEN HVM mode
  - Full virtualization no special OS kernel required
  - Computer entirely simulated
  - MS Windows runs in "hvm" mode
  - Allows work around: 10GB instance store root volume limit
  - Kernel is on the root volume (under /boot)
  - No AKIs (kernel images)
  - Commonly used today (EBS-backed instances)

October 28, 2025

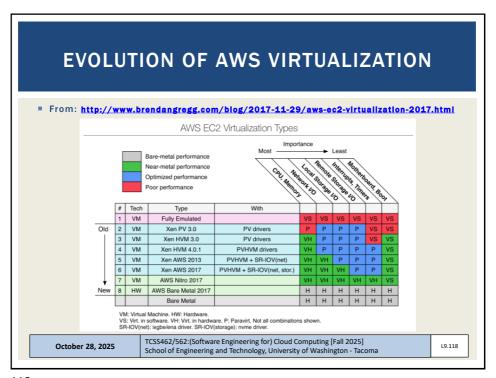
TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L9.116

116

### Nitro based on Kernel-based-virtual-machines Stripped down version of Linux KVM hypervisor Uses KVM core kernel module I/O access has a direct path to the device Goal: provide indistinguishable performance from bare metal | October 28, 2025 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] | School of Engineering and Technology, University of Washington - Tacoma | 19,117 | School of Engineering and Technology, University of Washington - Tacoma | 19,117 | School of Engineering and Technology, University of Washington - Tacoma | 19,117 | School of Engineering and Technology, University of Washington - Tacoma | 19,117 | School of Engineering and Technology, University of Washington - Tacoma | 19,117 | School of Engineering and Technology, University of Washington - Tacoma | 19,117 | School of Engineering and Technology, University of Washington - Tacoma | 19,117 | School of Engineering and Technology, University of Washington - Tacoma | 19,117 | School of Engineering and Technology, University of Washington - Tacoma | 19,117 | School of Engineering and Technology, University of Washington - Tacoma | 19,117 | School of Engineering and Technology, University of Washington - Tacoma | 19,117 | School of Engineering and Technology, University of Washington - Tacoma | 19,117 | School of Engineering and Technology, University of Washington - Tacoma | 19,117 | School of Engineering and Technology, University of Washington - Tacoma | 19,117 | School of Engineering and Technology, University of Washington - Tacoma | 19,117 | School of Engineering and Technology, University of Washington - Tacoma | 19,117 | School of Engineering and Technology, University of Washington - Tacoma | 19,117 | School of Engineering and Technology, University of Washington - Tacoma | 19,117 | School of Engineering and Technology, University of Washington - Tacoma | 19,117 | School of Engineering and Technology, University of Washington - Tacoma | 19,117 | School of Engineering and Technology | 19,117 | Scho

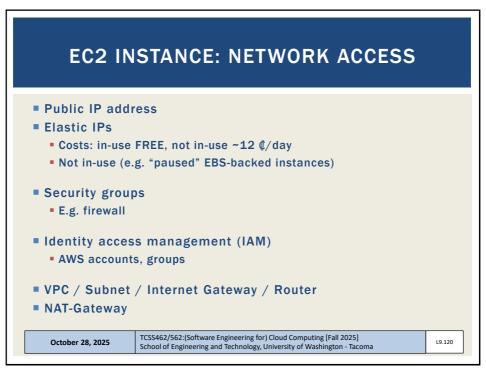
117



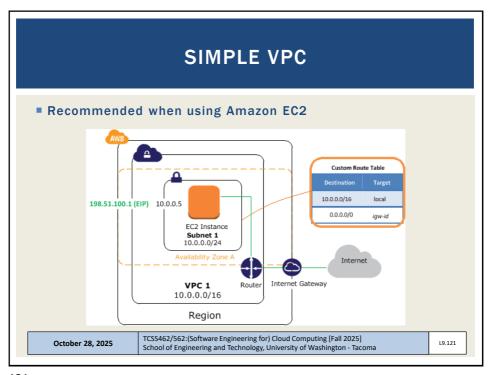
118

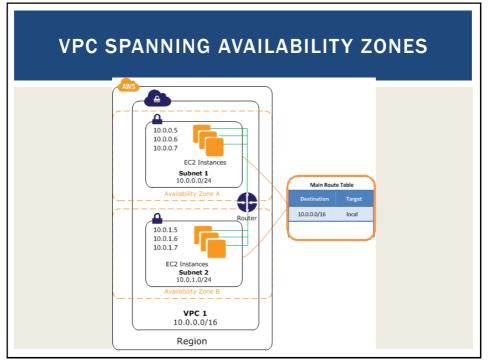
# INSTANCE ACTIONS Stop Costs of "pausing" an instance Terminate Reboot Image management Creating an image EBS (snapshot) Bundle image Instance-store October 28, 2025 TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

119



120





122

### INSPECTING INSTANCE INFORMATION

- EC2 VMs run a local metadata service
- Can query instance metadata to self discover cloud configuration attributes
- Find your instance ID:

```
curl http://169.254.169.254/
curl http://169.254.169.254/latest/
curl http://169.254.169.254/latest/meta-data/
curl http://169.254.169.254/latest/meta-data/instance-id
```

- ec2-get-info command
- Python API that provides easy/formatted access to metadata

October 28, 2025

; echo

TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L9.123

123

### SIMPLE STORAGE SERVICE (S3)

- Key-value blob storage
- What is the difference vs. key-value stores (NoSQL DB)?
- Can mount an S3 bucket as a volume in Linux
  - Supports common file-system operations
- Provides eventual consistency
- Can store Lambda function state for life of container.

October 28, 2025

TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L9.124

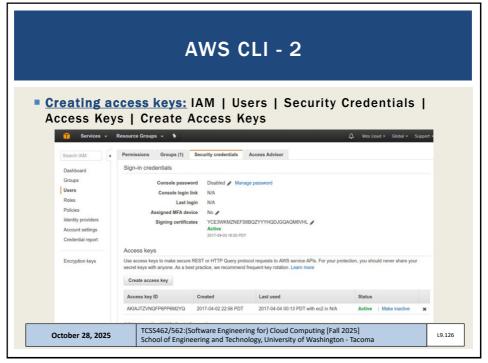
124

```
AWS CLI

Launch Ubuntu 16.04 VM
Instances | Launch Instance

Install the general AWS CLI
sudo apt install awscli

Create config file
[default]
aws_access_key_id = <access key id>
aws_access_key_id = <secret access key>
region = us-east-1
```



126

# AWS CLI - 3 Export the config file Add to /home/ubuntu/.bashrc export AWS\_CONFIG\_FILE=\$HOME/.aws/config Try some commands: aws help aws command help aws ec2 help aws ec2 describes-instances --output text aws ec2 describe-instances --output json aws s3 ls aws s3 ls aws s3 ls vmscaleruw October 28, 2025 TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

127

### LEGACY / SERVICE SPECIFIC CLI(S) sudo apt install ec2-api-tools Provides more concise output Additional functionality Define variables in .bashrc or another sourced script: export AWS\_ACCESS\_KEY={your access key} export AWS\_SECRET\_KEY={your secret key} ec2-describe-instances ec2-run-instances ec2-request-spot-instances EC2 management from Java: http://docs.aws.amazon.com/AWSJavaSDK/latest/javad oc/index.html Some AWS services have separate CLI installable by package TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] October 28, 2025 L9.128 School of Engineering and Technology, University of Washington - Tacoma

128

### **AMI TOOLS**

- Amazon Machine Images tools
- For working with disk volumes
- Can create live copies of any disk volume
  - Your local laptop, ec2 root volume (EBS), ec2 ephemeral disk
- Installation:

 $\frac{https://docs.aws.amazon.com/AWSEC2/latest/UserGuide/ami}{-tools-commands.html}$ 

- AMI tools reference:
- https://docs.aws.amazon.com/AWSEC2/latest/UserGuide/ami--tools-commands.html
- Some functions may require private key & certificate files

October 28, 2025

TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L9.129

129

### PRIVATE KEY AND CERTIFICATE FILE

- Install openssl package on VM
- # generate private key file

\$openssl genrsa 2048 > mykey.pk

# generate signing certificate file

\$openssl req -new -x509 -nodes -sha256 -days 36500 -key mykey.pk -outform PEM -out signing.cert

- Add signing.cert to IAM | Users | Security Credentials | -- new signing certificate --
- From: http://docs.aws.amazon.com/AWSEC2/latest/UserGuide/set-up-ami-tools.html?icmpid=docs\_iam\_console#ami-tools-create-certificate

October 28, 2025

TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L9.130

130

### PRIVATE KEY, CERTIFICATE FILE

- These files, combined with your AWS\_ACCESS\_KEY and AWS\_SECRET\_KEY and AWS\_ACCOUNT\_ID enable you to publish new images from the CLI
- Objective:
- 1. Configure VM with software stack
- 2. Burn new image for VM replication (horizontal scaling)
- An alternative to bundling volumes and storing in S3 is to use a containerization tool such as Docker. . .
- Create image script . . .

October 28, 2025

TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L9.131

131

### SCRIPT: CREATE A NEW INSTANCE STORE IMAGE FROM LIVE DISK VOLUME

```
image=$1
echo "Burn image $image"
echo "$image" > image.id
mkdir /mnt/tmp
AWS_KEY_DIR=/home/ubuntu/.aws
export EC2_URL=http://ec2.amazonaws.com
export S3_URL=https://s3.amazonaws.com
export EC2_PRIVATE_KEY=${AWS_KEY_DIR}/mykey.pk
export EC2_CERT=${AWS_KEY_DIR}/signing.cert
export AWS_USER_ID={your account id}
export AWS_ACCESS_KEY={your aws access key}
export AWS_SECRET_KEY={your aws secret key}
ec2-bundle-vol -s 5000 -u ${AWS_USER_ID} -c ${EC2_CERT} -k ${EC2_PRIVATE_KEY} --ec2cert /etc/ec2/amitools/cert-ec2.pem --no-inherit -r x86_64 -p $image -i
/etc/ec2/amitools/cert-ec2.pem
cd /tmp
ec2-upload-bundle -b tcss562 -m $image.manifest.xml -a ${AWS_ACCESS_KEY} -s ${AWS_SECRET_KEY} --url http://s3.amazonaws.com --location US
ec2-register tcss562/$image.manifest.xml --region us-east-1 --kernel aki-
88aa75e1
                          TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025]
      October 28, 2025
                          School of Engineering and Technology, University of Washington - Tacoma
```

132

### **COST SAVINGS MEASURES**

- From Tutorial 3:
- #1: ALWAYS USE SPOT INSTANCES FOR COURSE/RESEARCH RELATED PROJECTS
- #2: NEVER LEAVE AN EBS VOLUME IN YOUR ACCOUNT THAT IS NOT ATTACHED TO A RUNNING VM
- #3: BE CAREFUL USING PERSISTENT REQUESTS FOR SPOT INSTANCES
- #4: TO SAVE/PERSIST DATA, USE EBS SNAPSHOTS AND THEN
- #5: DELETE EBS VOLUMES FOR TERMINATED EC2 INSTANCES.
- #6: UNUSED SNAPSHOTS AND UNUSED EBS VOLUMES SHOULD BE PROMPTLY DELETED !!
- #7: USE PERSISTENT SPOT REQUESTS AND THE "STOP" FEATURE TO PAUSE VMS DURING SHORT BREAKS

October 28, 2025

TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L9.133

133



134