

Τ

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

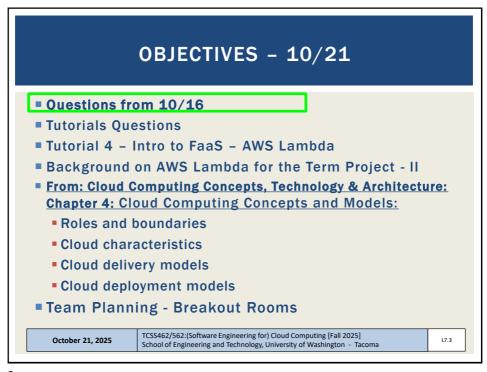
- Thursdays:
  - •6:00 to 7:00 pm CP 229 & Zoom
- **■**Fridays
  - ■11:00 am to 12:00 pm ONLINE via Zoom\*
- Or email for appointment
- > Office Hours set based on Student Demographics survey feedback
- \* Friday office hours may be adjusted or canceled due meeting conflicts or other obligations. Adjustments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

October 21, 2025

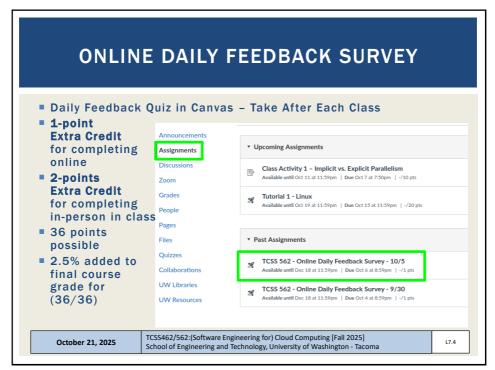
TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L7.2

2



3



4

### **WARNING**

- DO NOT SUBMIT BOTH A PAPER AND AN ONLINE SURVEY OR YOU WILL LOOSE POINTS
- CANVAS WILL AUTOMATICALLY REPLACE THE PAPER SURVEY SCORE (2 PTS) WITH THE ONLINE SURVEY (1 PT)
- \* COMPLETE ONLY ONE SURVEY FOR EACH CLASS SESSION \*
- WE WILL NOT BE ABLE TO DUPLICATE CHECK SURVEYS FOR EACH CLASS SESSION AND MAKE CORRECTIONS

October 21, 2025

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

L7.5

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 7.00 (↑ previous 6.91)
- Please rate the pace of today's class:
- 1-slow, 5-just right, 10-fast
- Average 5.44 (↑ previous 5.14)

October 21, 2025

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

L7.6

6

### FEEDBACK FROM 10/21

- What's the difference between AWS Lambda's function instance "mini-VMs" and regular VM?
- AWS Lambda uses "Firecracker" micro-VMs to host serverless workloads such as functions and containers
  - Firecracker is the hyper-visor (also called a virtual-machine monitor) used to create Firecracker microVMs
  - Suggested Reading Paper: https://www.usenix.org/conference/nsdi20/presentation/agache
- Regular VMs, such as those available from Amazon EC2 are fullfeatured, general purpose virtual servers
  - 5<sup>th</sup> generation and beyond use AWS Nitro nearly full HW virtualization:
  - Suggested Reading:
  - https://docs.aws.amazon.com/whitepapers/latest/security-design-of-aws-nitro-system/traditional-virtualization-primer.html

October 21, 2025

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

L7.7

7

### MICROVMS VS. VMS

### Firecracker Micro VMs on AWS

- No direct connections allowed. The guest OS can only be accessed via a serverless function or container code
- Micro VMs run only Amazon Linux 2023
- Boot time ~125 ms
- Based on Firecracker which is based on Google crosvm which is based on KVM
- Massively stripped down virtual computer w/ simplified device model: no BIOS, no PCI
- Internal hidden IP addresses

### EC2 VM

- Users can directly connect using SSH for console sessions to interact with the guest OS
- VMs run any OS selected by the user
- Boot time 3-8+ sec, (more w/ special OS or software)
- Full-featured VMs based on AWS Nitro hypervisor based on KVM (Linux Kernel Virtual Machine)
- Public/private IP addresses

October 21, 2025

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

L7.8

8

### FEEDBACK - 2

- If making a function call (to AWS Lambda) and my internet disconnects, will the function keep running on the instance and keep the result for me or does it simply terminate?
- If a client connection to an AWS Lambda function fails, the function should continue to run
- The issue is how is the result provided to the user
- If the result is not saved (persisted) somewhere, and only returned as a REST response object, the result is lost
- If the result is persisted in a data store, such as the simple storage service (S3), then the client can retrieve the result later on

October 21, 2025

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

L7.9

9

### FEEDBACK - 3

- If we want to get the response from an asynchronous call, do we need to make a synchronous call? (b/c async call has no response)
- No. To obtain the response from an asynchronous call, the result must be fetched from a data store
- The simple storage service (S3) is most commonly used to persist data, but any database service can be used.
- Common alternatives:
- 1. DynamoDB (No SQL DB)
- 2. Amazon RDS & Aurora (managed relational database service)
- 3. SQS Simple Queuing Service
- 4. SNS Simple Notification Service
- 5. Elasticache
- 6. Document DB
- 7. Amazon MQ

October 21, 2025

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

L7.10

10

### FEEDBACK - 4

- Do asynchronous client calls to a server, save the client time moreso than synchronous calls?
- YES
- A synchronous call blocks the calling thread to wait for a result from the server
- If the programmer has not designed the client to be multithreaded, then the client essentially is frozen while waiting for a results from the server – it can do nothing but wait
- The programmer can "spawn" a thread for the synchronous call while the parent thread performs other work
- Multi-threaded programming is more complex and resource intensive, however
- An asynchronous call frees the main thread immediately to do other work, and does not block and wait

October 21, 2025

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

L7.11

11

### FEEDBACK - 5

- When should synchronous vs. asynchronous client calls to a server be used ?
- Asynchronous calls are best for long operations
  - Maintaining a network connection for more than 30-seconds is error prone
  - Example: mobile device traveling down I-5 switching cell towers
- Synchronous calls block the client program unless it is a multi-threaded client
  - No other work can happen while waiting
  - Good for short calls that are expected to quickly return a result (within a few seconds)
- Clients and servers can run out of "connections" if too many synchronous sessions occur simultaneously
  - Asynchronous calls close connections and lower the burden

October 21, 2025

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

L7.12

12

### FEEDBACK - 6

- Are AWS Lambda functions similar to Redis message queues ?
- No
- AWS Lambda functions provide a general-purpose compute platform for running serverless function code for up to 15minutes
- Redis queue is a message or job queue built using Redis as the underlying data store. Redis is not a queueing system itself, but Redis Lists are highly effective for building and managing queues
- Queues are a form of a persistent data store not a general compute platform

October 21, 2025

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

L7.13

L7.14

13

## AWS LAMBDA: VCPU SCALING W/ MEMORY

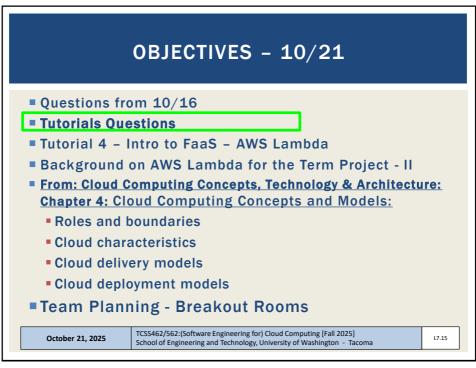
|   | Function Memory | CPU time share             |
|---|-----------------|----------------------------|
|   | 1769 MB         | 100 % = 1 vCPU             |
|   | 2389 MB         | 150 % = 1.5 vCPUs          |
|   | 3008 MB         | 200 % = 2 vCPUs            |
|   | 4158 MB         | 250 % = 2.5 vCPUs          |
|   | 5307 MB         | 300 % = 3 vCPUs            |
|   | 6192 MB         | 350 % = 3.5 vCPUs          |
|   | 7076 MB         | 400 % = 4 vCPUs (1 HT)     |
|   | 7960 MB         | 450 % = 4.5 vCPUs (1.5 HT) |
|   | 8845 MB         | 500 % = 5 vCPUs (2 HT)     |
|   | 9543 MB         | 550 % = 5.5 vCPUs (2.5 HT) |
| Boood on.   | 10240 MB        | 600 % = 6 vCPUs (3 HT)     |
| Based on: <a href="https://stackoverflow.com/questions/66522916/aws-lambda-memory-vs-cpu-con/">https://stackoverflow.com/questions/66522916/aws-lambda-memory-vs-cpu-con/</a> |                 |                            |

14

October 21, 2025

Slides by Wes J. Lloyd L7.7

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



15

## TUTORIAL O Getting Started with AWS https://faculty.washington.edu/wlloyd/courses/tcss562/tutorials/TCSS462\_562\_f2025\_tutorial\_0.pdf Create an AWS account Create account credentials for working with the CLI Install awsconfig package Setup awsconfig for working with the AWS CLI TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

16

### TUTORIAL 2 - OCT 21

- Introduction to Bash Scripting
- https://faculty.washington.edu/wlloyd/courses/tcss562/tutorials/T CSS462\_562\_f2025\_tutorial\_2.pdf
- Review tutorial sections:
- Create a BASH webservice client
  - 1. What is a BASH script?
  - 2. Variables
  - 3. Input
  - 4. Arithmetic
  - 5. If Statements
  - 6. Loops
  - 7. Functions
  - 8. User Interface
- Call service to obtain IP address & lat/long of computer
- Call weatherbit.io API to obtain weather forecast for lat/long

October 21, 2025

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

L7.17

17

### 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

October 21, 2025

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

L7.18

18

### **OBJECTIVES - 10/21**

- Questions from 10/16
- Tutorials Questions
- Tutorial 4 Intro to FaaS AWS Lambda
- Background on AWS Lambda for the Term Project II
- From: Cloud Computing Concepts, Technology & Architecture: Chapter 4: Cloud Computing Concepts and Models:
  - Roles and boundaries
  - Cloud characteristics
  - Cloud delivery models
  - Cloud deployment models
- Team Planning Breakout Rooms

October 21, 2025

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

L7.19

19

### **TUTORIAL 4 - TO BE POSTED**

- Introduction to AWS Lambda with the Serverless Application Analytics Framework (SAAF)
- (link to be posted)
- Setting up a Java development environment (IDE)
- Introduction to Maven build files for Java
- Create and Deploy "hello" Java AWS Lambda Function
  - Creation of API Gateway REST endpoint
- Sequential testing of "hello" AWS Lambda Function
  - API Gateway endpoint
  - AWS CLI Function invocation
- Observing SAAF profiling output
- Parallel testing of "hello" AWS Lambda Function with faas\_runner tool
- Performance analysis using faas\_runner reports
- Two function pipeline development task: Caesar Cipher

October 21, 2025

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

L7.20

20

### OBJECTIVES - 10/21

- Questions from 10/16
- Tutorials Questions
- Tutorial 4 Intro to FaaS AWS Lambda
- Background on AWS Lambda for the Term Project II
- From: Cloud Computing Concepts, Technology & Architecture: Chapter 4: Cloud Computing Concepts and Models:
  - Roles and boundaries
  - Cloud characteristics
  - Cloud delivery models
  - Cloud deployment models
- Team Planning Breakout Rooms

October 21, 2025

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

L7.21

21

### **AWS LAMBDA PLATFORM LIMITATIONS - 2**

- 10 concurrent function executions inside account (default)
- Function payload: 6MB (synchronous), 256KB (asynchronous)
- Deployment package: 50MB (compressed), 250MB (unzipped)
- Container image size: 10 GB
- Processes/threads: 1024
- File descriptors: 1024
- Function instances run Amazon Linux 2023
  - Based on a combination of Red Hat open-source Linux distributions:
     Fedora (versions 34, 35, 36) and CentOS 9 Stream
- Suggested Reading:
- https://docs.aws.amazon.com/lambda/latest/dg/gettingstart ed-limits.html

October 21, 2025

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

L7.22

22

### **CPUSTEAL**



- CpuSteal: Metric that measures when a CPU core is ready to execute but the physical CPU core is busy and unavailable
- Symptom of over provisioning physical servers in the cloud
- Factors which cause *CpuSteal*: (x86 hyperthreading)
  - 1. Physical CPU is shared by too many busy VMs
  - 2. Hypervisor kernel is using the CPU
    - On AWS Lambda this would be the Firecracker MicroVM which is derived from the KVM hypervisor
  - VM's CPU time share <100% for 1 or more cores, and 100% is needed for a CPU intensive workload.
- Man procfs press "/" type "proc/stat"
  - CpuSteal is the 8<sup>th</sup> column returned
  - Metric can be read using SAAF in tutorial #4

October 21, 2025

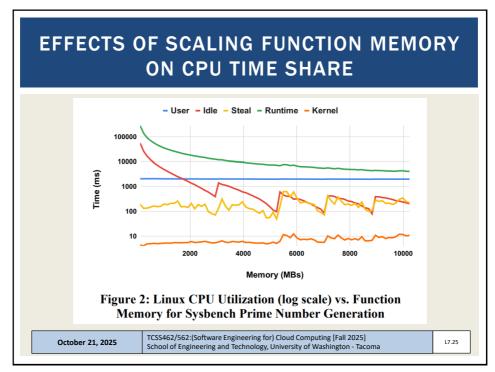
TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

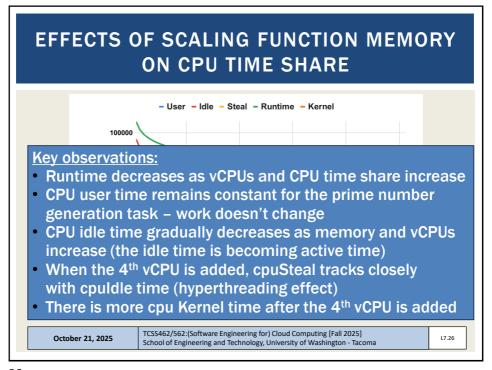
L7.23

23

```
Snippet of sample output returned by SAAF (Tutorial 4)
{
   "version": 0.2,
   "lang": "python",
    "cpuType": "Intel(R) Xeon(R) Processor @ 2.50GHz",
    "cpuModel": 63,
    "vmuptime": 1551727835,
    "uuid": "d241c618-78d8-48e2-9736-997dc1a931d4",
    "vmID": "tiUCnA",
   "platform": "AWS Lambda",
   "newcontainer": 1,
   "cpuUsrDelta": "904",
   "cpuNiceDelta": "0",
    "cpuKrnDelta": "585"
    "cpuIdleDelta": "82428",
    "cpuIowaitDelta": "226",
   "cpuIrqDelta": "0",
   "cpuSoftIrqDelta": "7",
   "vmcpustealDelta": "1594",
   "frameworkRuntime": 35.72,
   "message": "Hello Fred Smith!",
   "runtime": 38.94
                TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma
October 21, 2025
```

24





26

### **FUNCTION INSTANCE LIFE CYCLES**

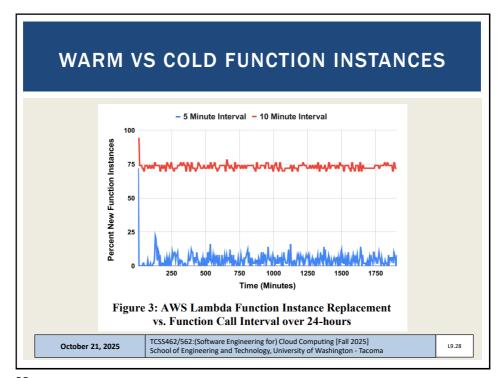
- Function states:
- COLD: brand new function instance just initialized to run the request (more overhead)
  - Platform cold (first time ever run)
  - Host cold (function assets cached locally on servers)
- WARM: existing function instance that is reused
- All function instances persist for ~5 minutes before they begin to be "garbage collected" by the platform
  - 100% garbage collection may take up to ~30-40 minutes
- AWS Lambda appears to "recycle" infrastructure faster than other FaaS platforms
  - Presumably because of need, because the platform is busy

October 21, 2025

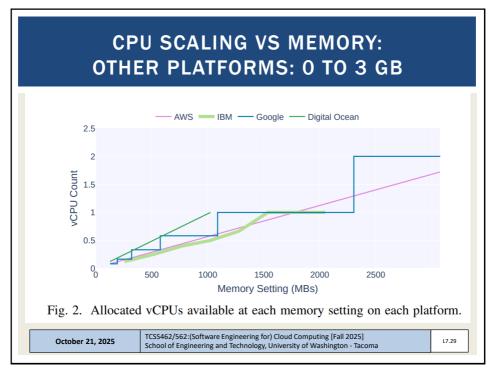
TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

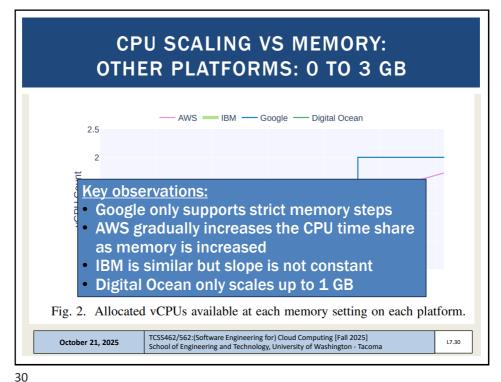
L7.27

27



28





30

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

- Traditionally AWS Lambda functions have been limited to 500MB of storage space
- Recently the Elastic File System (EFS) has been extended to support AWS Lambda
- The Elastic File System supports the creation of a shared volume like a shared disk (or folder)
  - EFS is similar to NFS (network file share)
  - Multiple AWS Lambda functions and/or EC2 VMs can mount and share the same EFS volume
  - Provides a shared R/W disk
  - Breaks the 500MB capacity barrier on AWS Lambda
- Downside: EFS is expensive: ~30 \$\psi/GB/month\$
- **Project**: EFS performance & scalability evaluation on Lambda

October 21, 2025

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

L7.31

31

## SERVERLESS APPLICATION - DESIGN TRADEOFFS

- Serverless file systems: EFS, docker container, extended /tmp
- Service/function composition / decomposition
- Switchboard architecture
- Application control flow
- Programming language comparison (course theme w/ LLMs)
- FaaS platforms: AWS, Azure, Google, etc.
- Alternate data services/backends for application state, large data transfer, short to long term data persistence
- Performance variability
  - Temporal: 24 hour, 7 days, etc. (diurnal patterns?)
  - Geospatial: By Region, availability zone
  - From HW heterogeneity (alternate CPUs)

October 21, 2025

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

L7.32

32

## SERVERLESS FILE STORAGE COMPARISON PROJECT

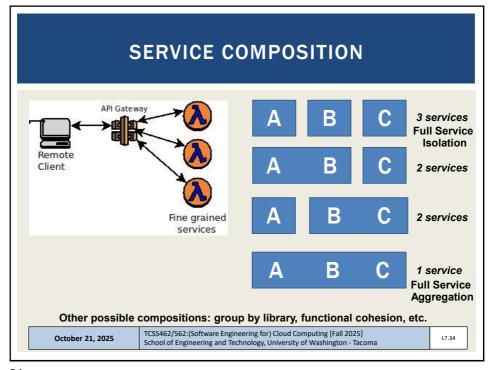
- Elastic File System (EFS):
   Performance, Cost, and Scalability Evaluation in the context of AWS Lambda / Serverless Computing
  - EFS provides a file system that can be shared with multiple Lambda function instances in parallel
- Using a common use case, compare performance and cost of extended storage options on AWS Lambda:
  - Docker container support (up to 10 GB) read only
  - Ephemeral /tmp (up to 10 GB) read/write
  - EFS (unlimited, but costly) read/write
  - image integration with AWS Lambda performance & scalability

October 21, 2025

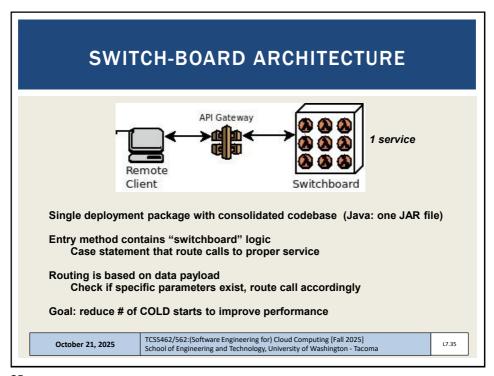
TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

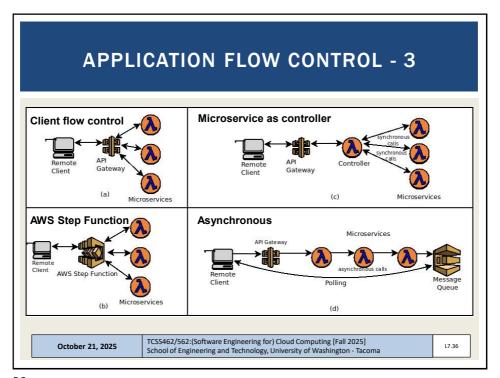
L7.33

33



34





36

### PROGRAMMING LANGUAGE COMPARISON

- FaaS platforms support hosting code in multiple languages
- AWS Lambda- common: Java, Node.js, Python
  - Plus others: Go, PowerShell, C#, and Ruby
- Also Runtime API ("BASH") which allows deployment of binary executables from any programming language
- August 2020 Our group's paper:
- https://tinyurl.com/y46eq6np
- If wanting to perform a language study either:
  - Implement in C#, Ruby, or multiple versions of Java, Node.js, Python
  - OR implement different app than TLQ (ETL) data processing pipeline

October 21, 2025

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

L7.37

37

### **FAAS PLATFORMS**

- Many commercial and open source FaaS platforms exist
- TCSS562 projects can choose to compare performance and cost implications of alternate platforms.
- Supported by SAAF:
- AWS Lambda
- Google Cloud Functions
- Azure Functions
- IBM Cloud Functions
- Apache OpenWhisk (open source, deploy your own FaaS)
- Open FaaS (open source, deploy your own FaaS)

October 21, 2025

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

L7.38

38

### **DATA PROVISIONING**

- Consider performance and cost implications of the data-tier design for the serverless application
- Use different tools as the relational datastore to support service #2 (LOAD) and service #3 (EXTRACT)
- SQL / Relational:
- Amazon Aurora (serverless cloud DB), Amazon RDS (cloud DB), DB on a VM (MySQL), DB inside Lambda function (SQLite, Derby)
- NO SQL / Key/Value Store:
- Dynamo DB, MongoDB, S3

October 21, 2025

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

L7.39

39

### PERFORMANCE VARIABILITY

- Cloud platforms exhibit performance variability which varies over time
- Goal of this case study is to measure performance variability (i.e. extent) for AWS Lambda services by hour, day, week to look for common patterns
- Can also examine performance variability by availability zone and region
  - Do some regions provide more stable performance?
  - Can services be switched to different regions during different times to leverage better performance?
- Remember that performance = cost
- If we make it faster, we make it cheaper...

October 21, 2025

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

L7.40

40

### CPU STEAL CASE STUDY

- On AWS Lambda (or other FaaS platforms), when we run functions, how much CpuSteal do we observe?
- How does CpuSteal vary for different workloads? (e.g. functions that have different resource requirements)
- How does CpuSteal vary over time hour, day, week, location?
- How does CpuSteal relate to function performance?

October 21, 2025

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

L7.41

41

### **CPU ARCHITECTURE & PERFORMANCE**

- **X86\_64** Intel
  - Intel Xeon Platinum 8259 CL @ 2.5 GHz
- ARM64 Graviton2

October 21, 2025

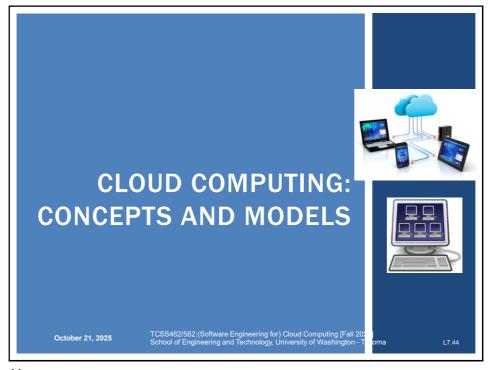
TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L7.42

42

# OBJECTIVES - 10/21 Questions from 10/16 Tutorials Questions Tutorial 4 - Intro to FaaS - AWS Lambda Background on AWS Lambda for the Term Project - II From: Cloud Computing Concepts, Technology & Architecture: Chapter 4: Cloud Computing Concepts and Models: Roles and boundaries Cloud characteristics Cloud delivery models Cloud deployment models Team Planning - Breakout Rooms October 21, 2025 TCSS462/562:(Software Engineering for) Cloud Computing [Fail 2025] School of Engineering and Technology, University of Washington - Tacoma

43



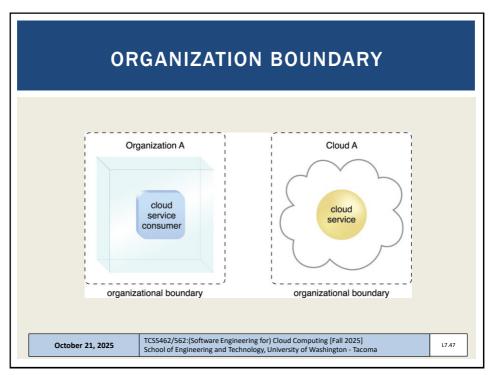
44

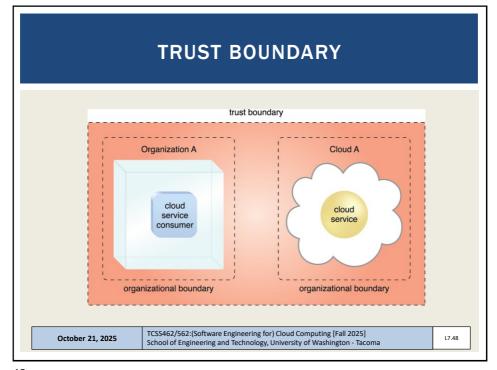
## Cloud provider Organization that provides cloud-based resources Responsible for fulfilling SLAs for cloud services Some cloud providers "resell" IT resources from other cloud providers Example: Heroku sells PaaS services running atop of Amazon EC2 Cloud consumers Cloud users that consume cloud services Cloud service owner Both cloud providers and cloud consumers can own cloud services A cloud service owner may use a cloud provider to provide a cloud service (e.g. Heroku) October 21, 2025 TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

45

### ROLES - 2 Cloud resource administrator Administrators provide and maintain cloud services Both cloud providers and cloud consumers have administrators Cloud auditor Third-party which conducts independent assessments of cloud environments to ensure security, privacy, and performance. Provides unbiased assessments Cloud brokers An intermediary between cloud consumers and cloud providers Provides service aggregation Cloud carriers Network and telecommunication providers which provide network connectivity between cloud consumers and providers TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma October 21, 2025 17.46

46





48

### **OBJECTIVES - 10/21** Questions from 10/16 ■ Tutorials Questions ■ Tutorial 4 - Intro to FaaS - AWS Lambda Background on AWS Lambda for the Term Project - II ■ From: Cloud Computing Concepts, Technology & Architecture: **Chapter 4: Cloud Computing Concepts and Models:** Roles and boundaries Cloud characteristics Cloud delivery models Cloud deployment models ■ Team Planning - Breakout Rooms TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] October 21, 2025 School of Engineering and Technology, University of Washington - Tacoma

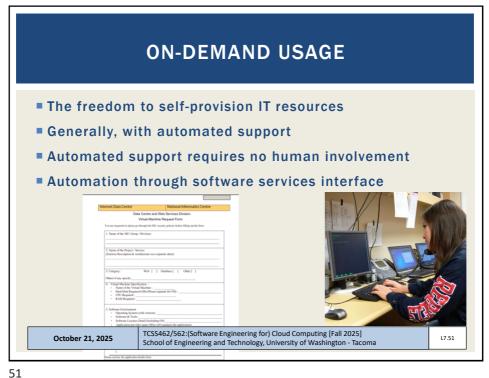
### **CLOUD CHARACTERISTICS** On-demand usage Ubiquitous access Multitenancy (resource pooling) Elasticity Measured usage Resiliency Assessing these features helps measure the value offered by a given cloud service or platform TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025]

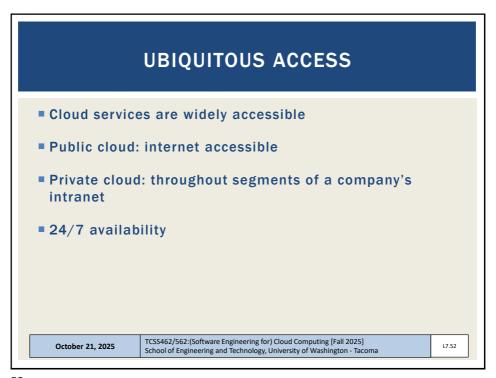
50

October 21, 2025

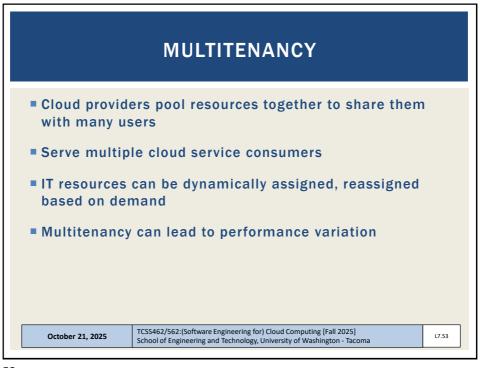
Slides by Wes J. Lloyd L7.25

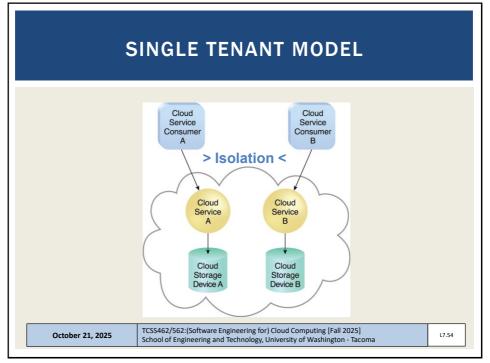
School of Engineering and Technology, University of Washington - Tacoma



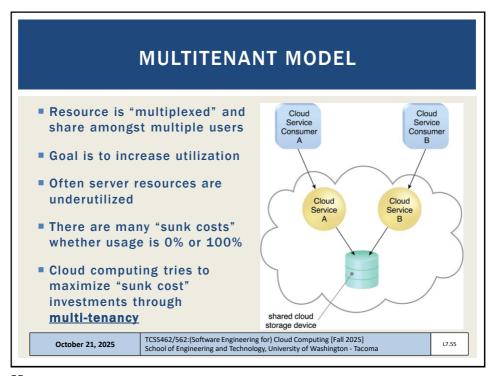


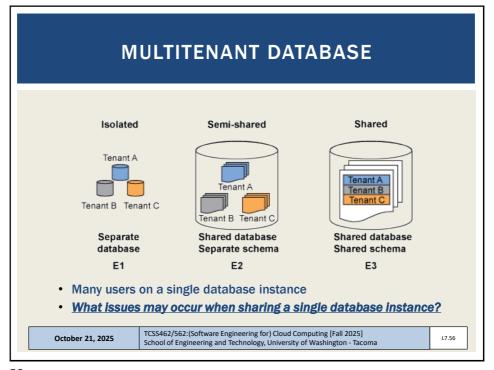
52



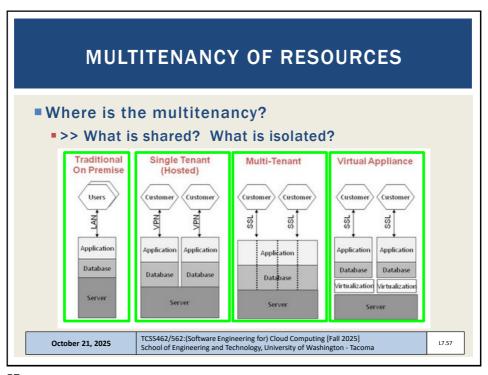


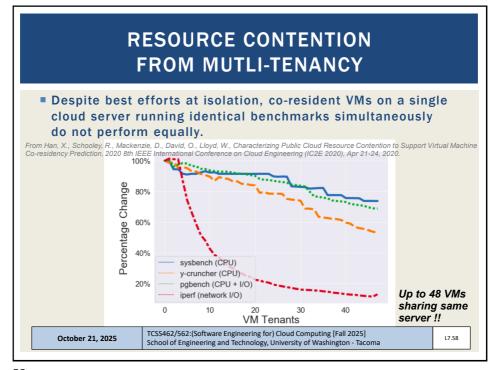
54



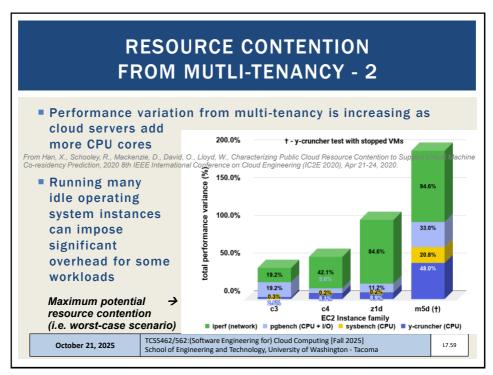


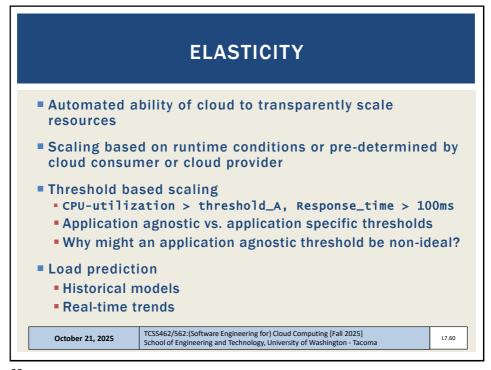
56



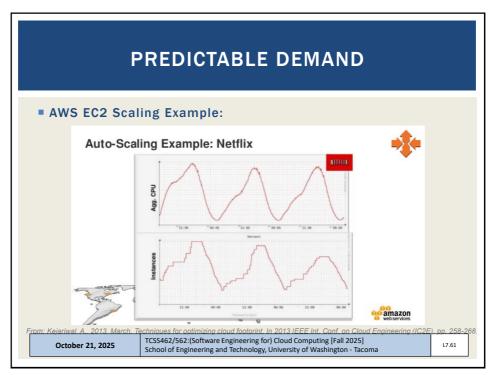


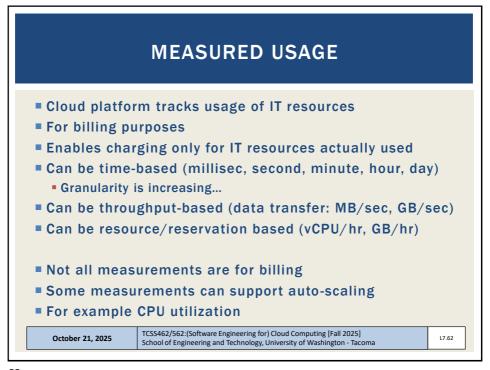
58



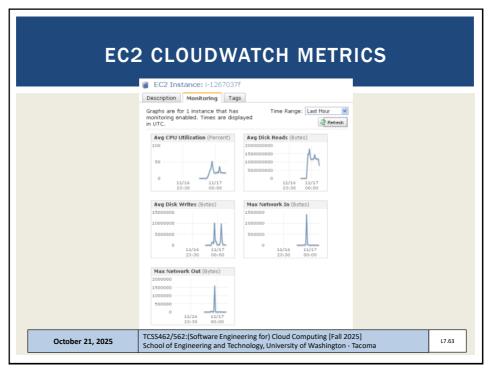


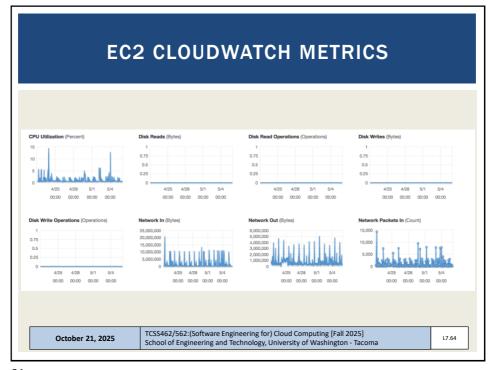
60



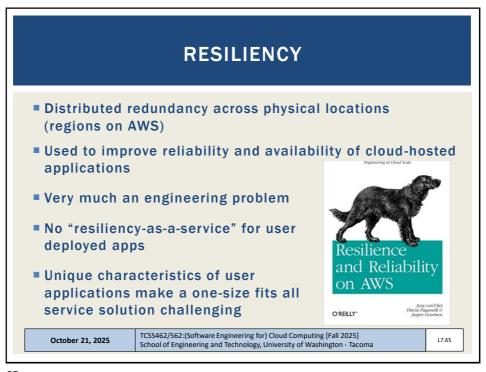


62

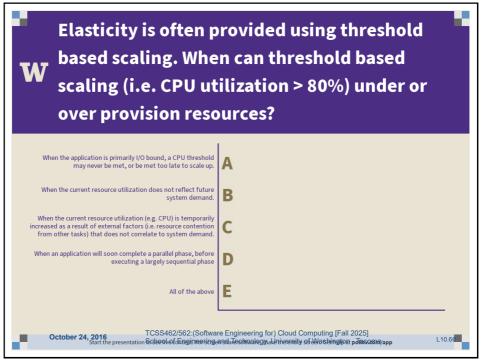




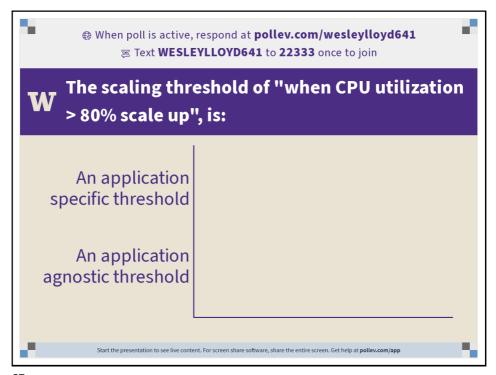
64



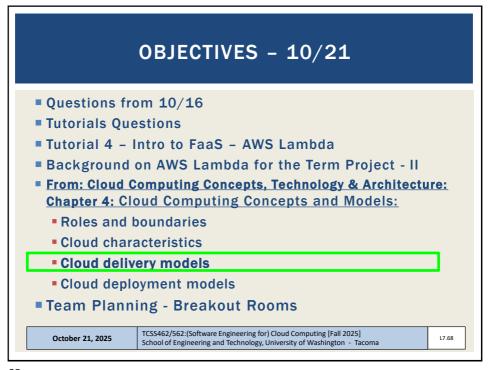
65



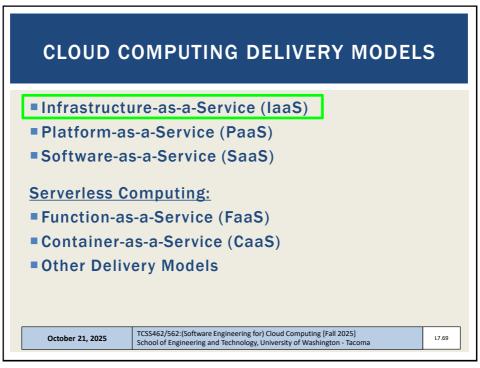
66

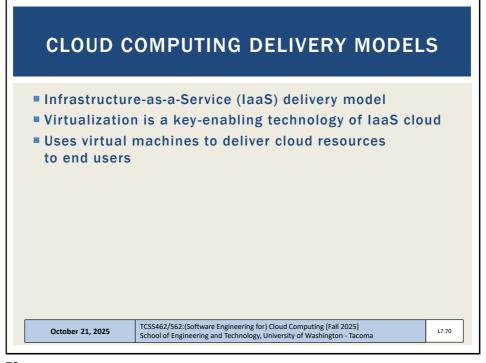


67

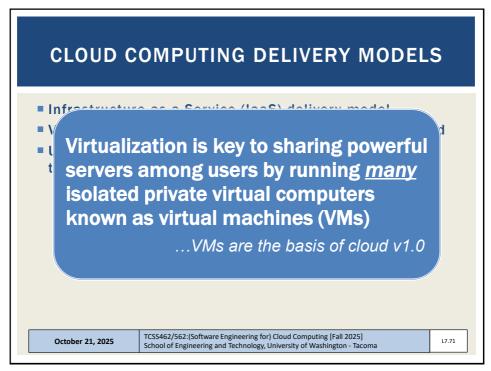


68

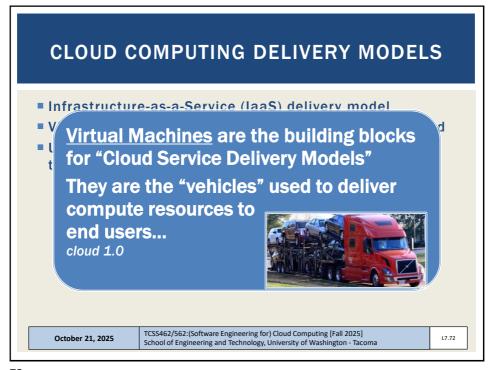




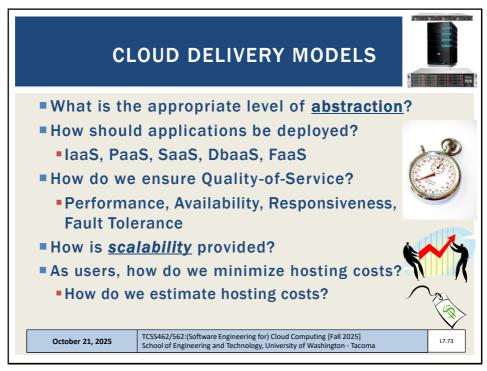
70



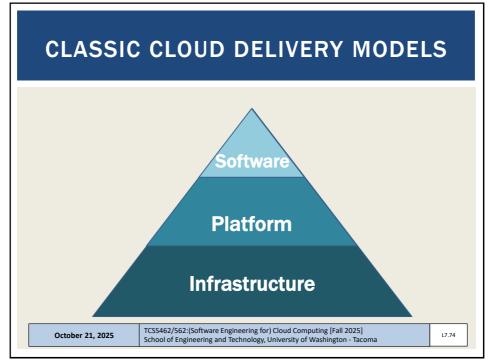
71



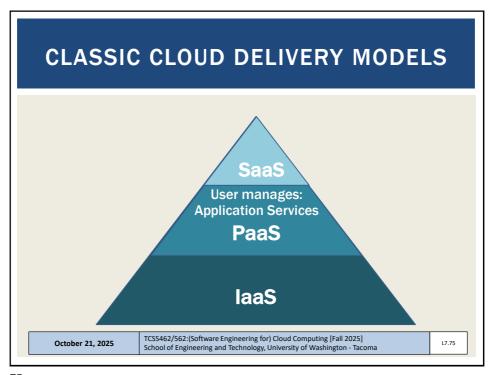
72

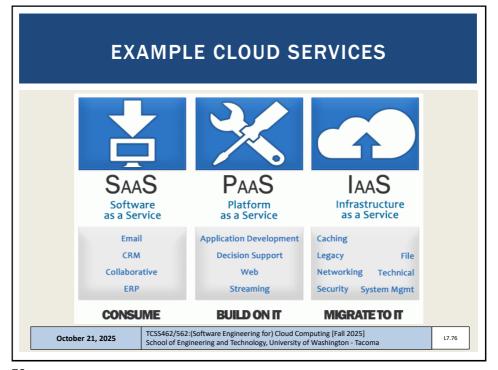


73

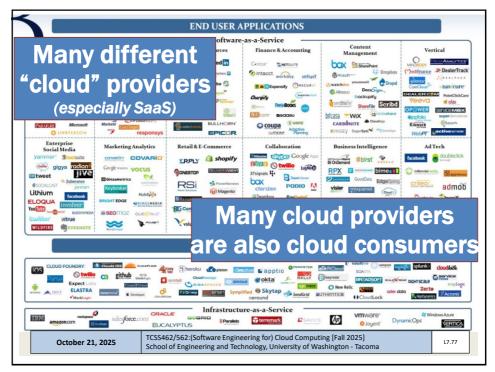


74





76

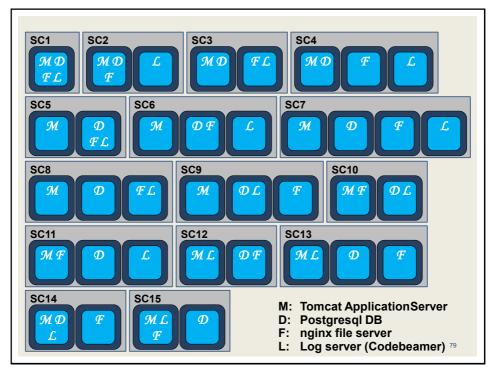




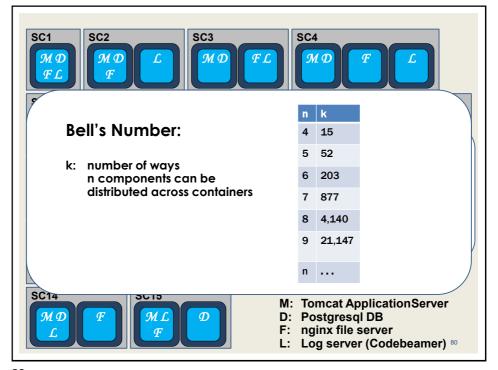
- Compute resources, on demand, as-a-service
  - Generally raw "IT" resources
  - Hardware, network, containers, operating systems
- Typically provided through virtualization
- Generally, not-preconfigured
- Administrative burden is owned by cloud consumer
- Best when high-level control over environment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

October 21, 2025 TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

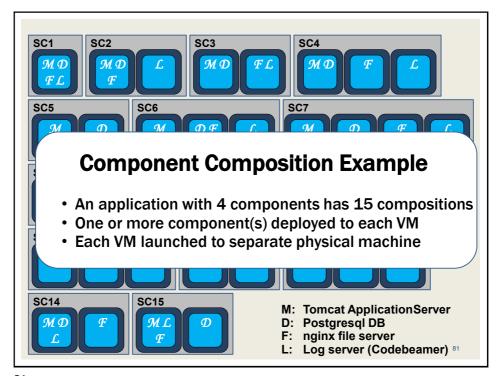
78

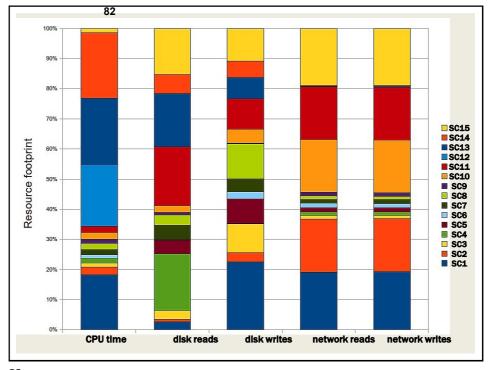


79

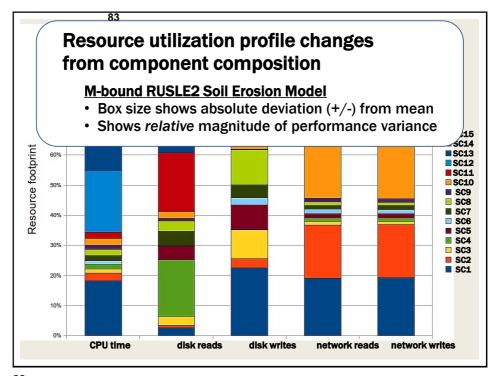


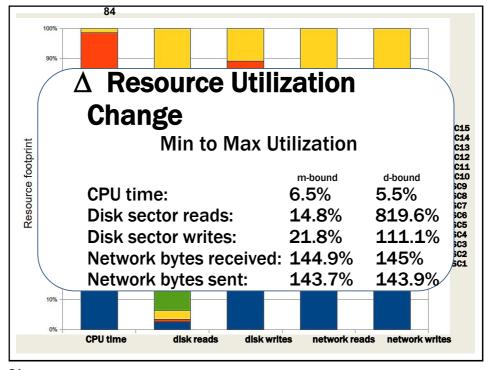
80



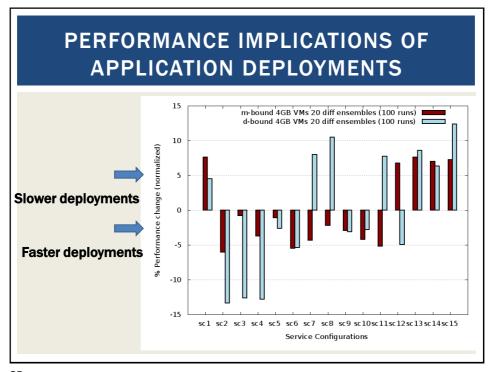


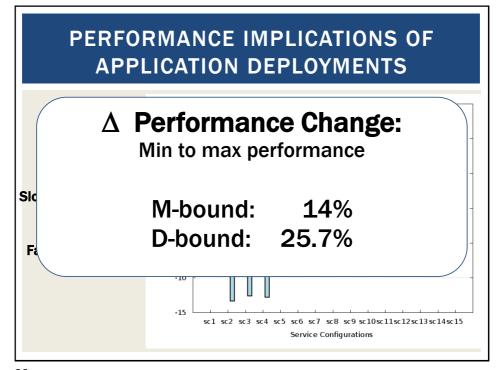
82



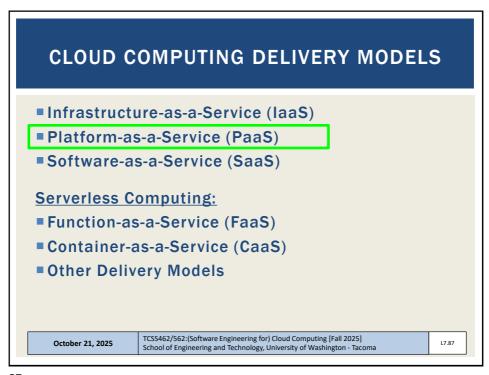


84

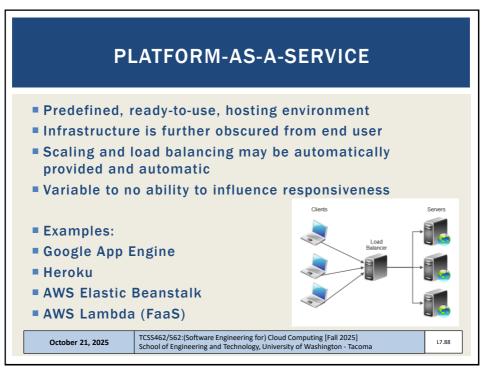




86



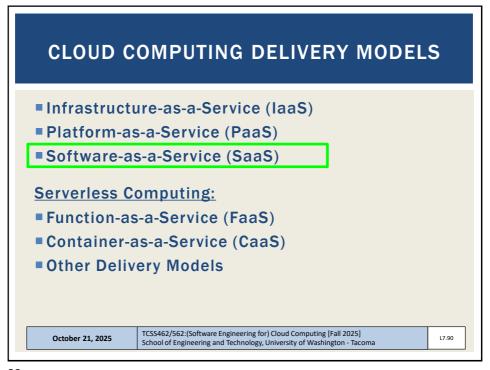
87



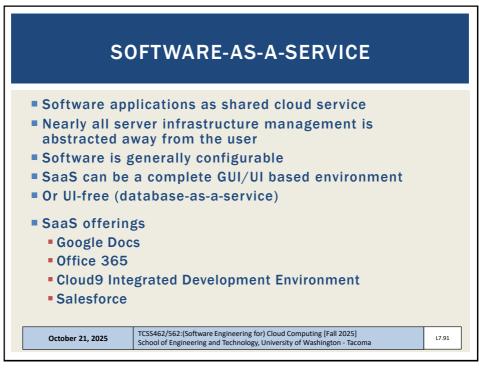
88

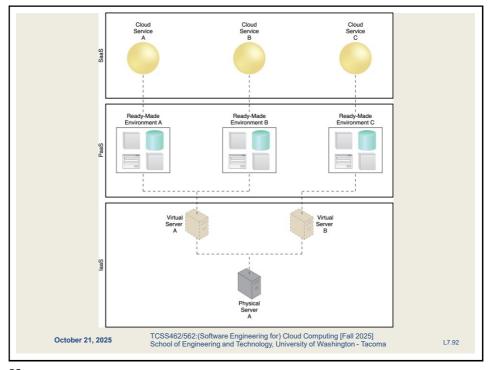
## USES FOR PAAS Cloud consumer Wants to extend on-premise environments into the cloud for "web app" hosting Wants to entirely substitute an on-premise hosting environment Cloud consumer wants to become a cloud provider and deploy its own cloud services to external users PaaS spares IT administrative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

89

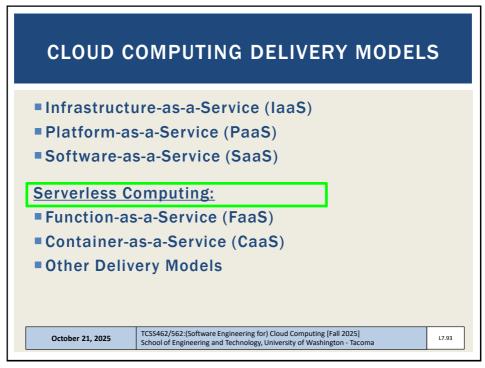


90





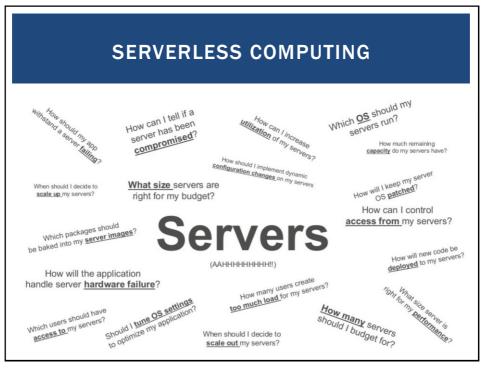
92

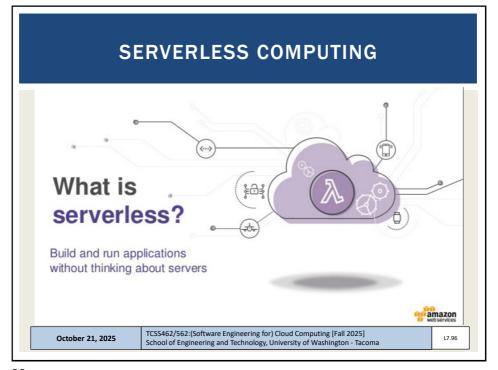


93

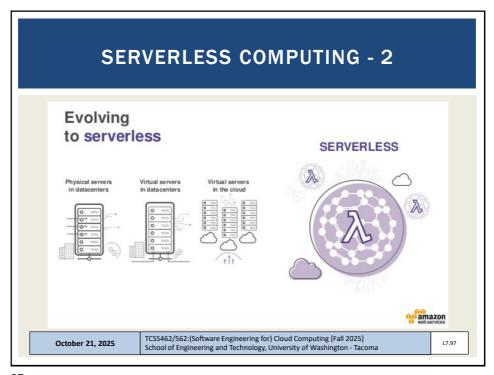


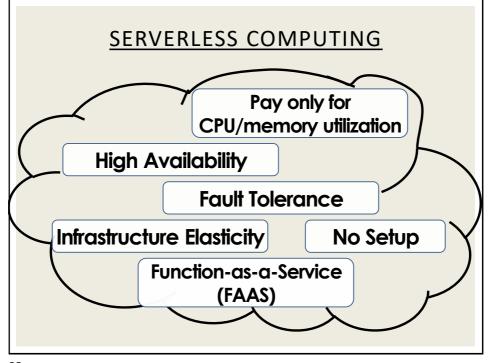
94





96

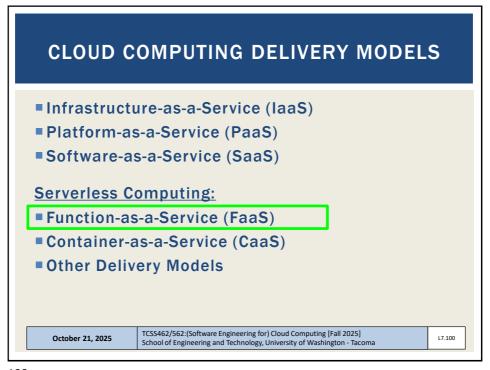




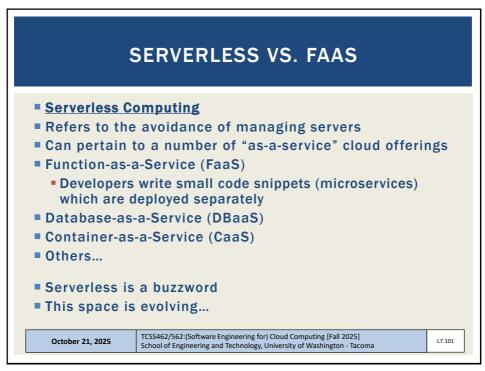
98

### Why Serverless Computing? Many features of distributed systems, that are challenging to deliver, are provided automatically ...they are built into the platform

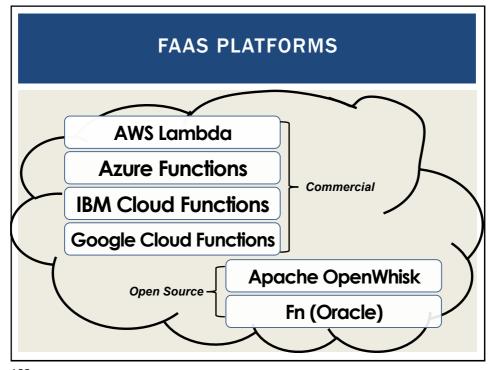
99



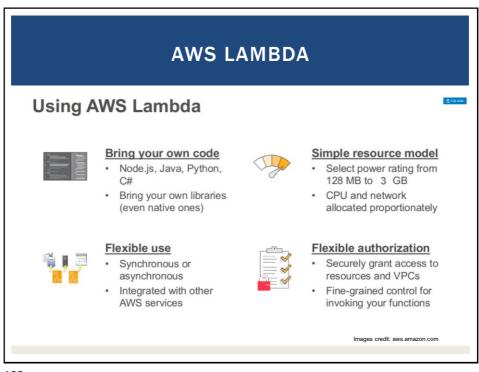
100



101



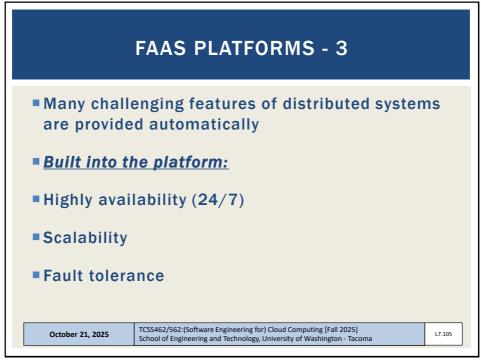
102

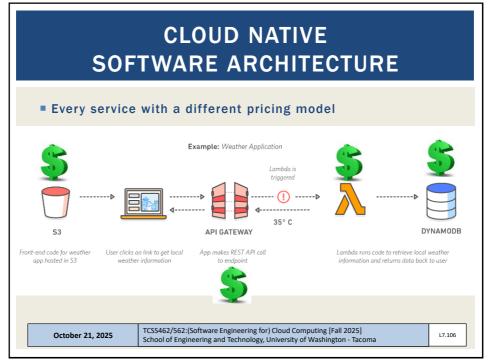


103

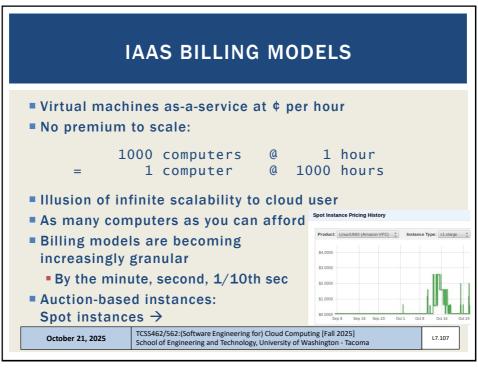
## FAAS PLATFORMS - 2 New cloud platform for hosting application code Every cloud vendor provides their own: AWS Lambda, Azure Functions, Google Cloud Functions, IBM OpenWhisk Similar to platform-as-a-service Replace opensource web container (e.g. Apache Tomcat) with abstracted vendor-provided black-box environment Cotober 21, 2025 TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

104



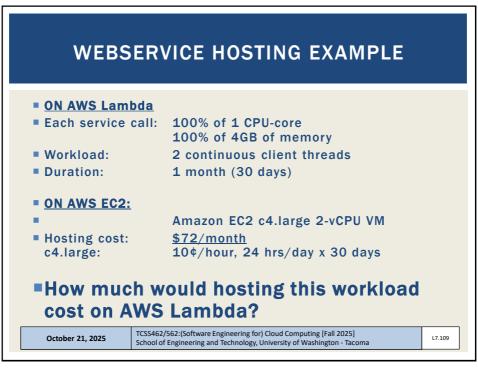


106

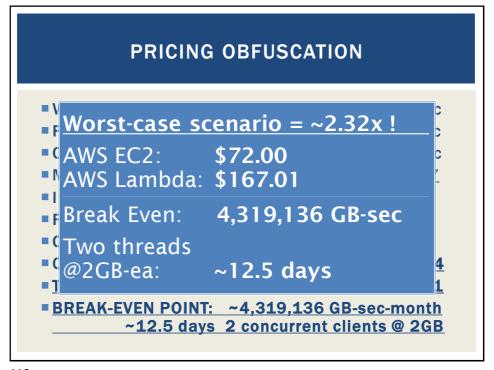


### PRICING OBFUSCATION ■ VM pricing: hourly rental pricing, billed to nearest second is intuitive... non-intuitive pricing policies FaaS pricing: • FREE TIER: first 1,000,000 function calls/month $\rightarrow$ FREE first 400,000 GB-sec/month → FREE Afterwards: obfuscated pricing (AWS Lambda): \$0.0000002 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 21, 2025 17 108

108



109



110

### **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 21, 2025

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

L7.111

111

### FACTORS IMPACTING PERFORMANCE OF FAAS COMPUTING PLATFORMS

- Infrastructure elasticity
- Load balancing
- Provisioning variation
- Infrastructure retention: COLD vs. WARM
  - Infrastructure freeze/thaw cycle
- Memory reservation
- Service composition

October 21, 2025

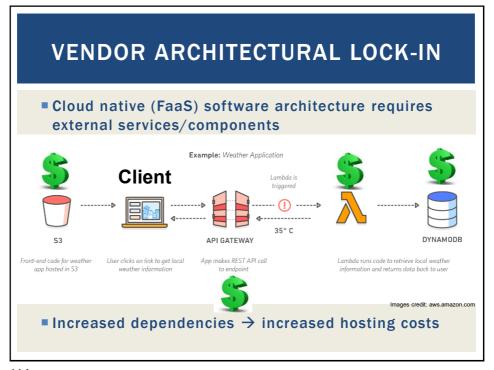
TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L7.112

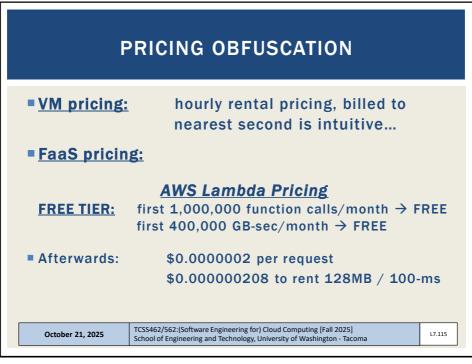
112

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

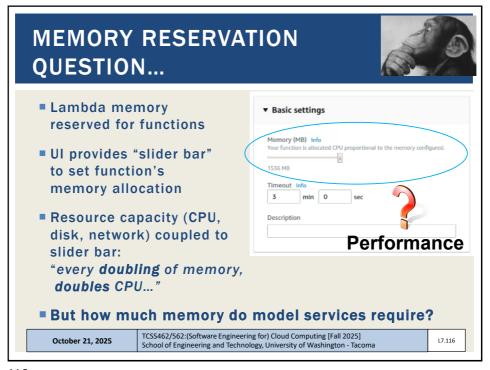
113



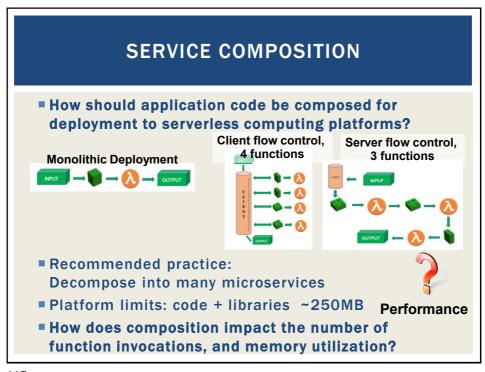
114

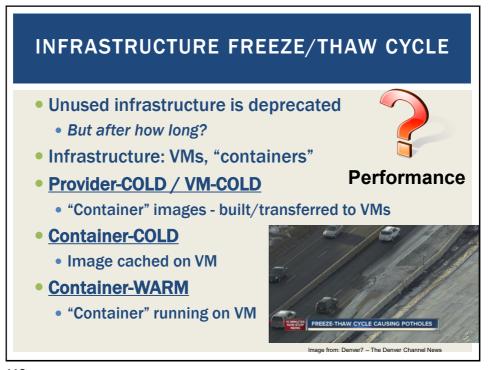


115

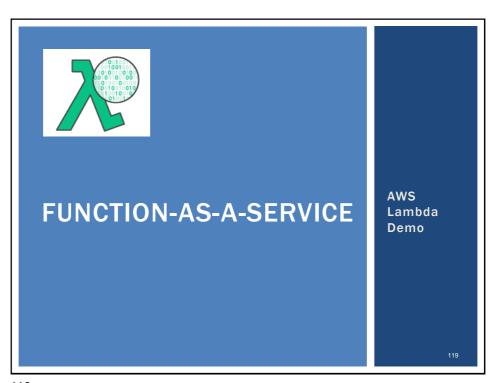


116

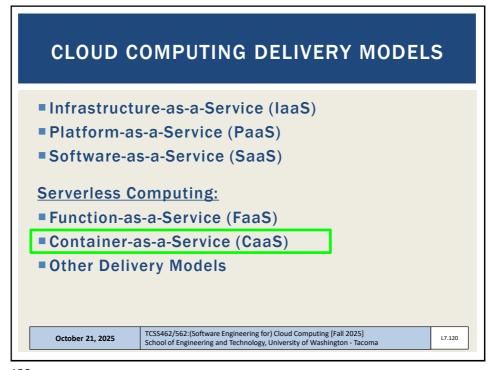




118



119



120

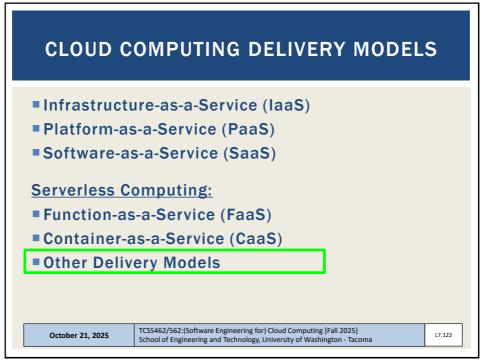
### CONTAINER-AS-A-SERVICE Cloud service model for deploying application containers (e.g. Docker) to the cloud Deploy containers without worrying about managing infrastructure: Servers Or container orchestration platforms Container platform examples: Kubernetes, Docker swarm, Apache Mesos/Marathon, Amazon Elastic Container Service Container platforms support creation of container clusters on the using cloud hosted VMs CaaS Examples: AWS Fargate Azure Container Instances Google KNative

TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025]

School of Engineering and Technology, University of Washington - Tacoma

121

October 21, 2025

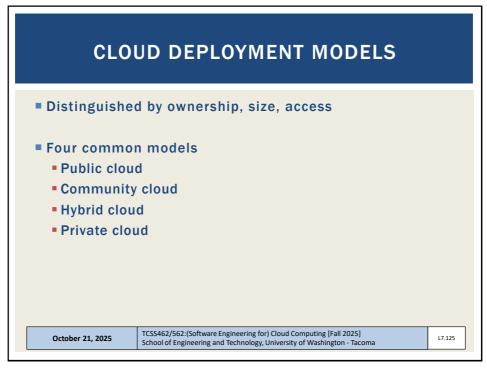


122

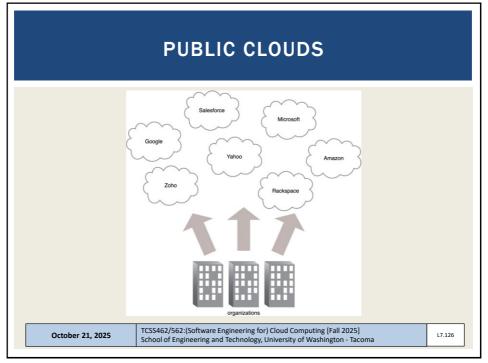
123

# OBJECTIVES - 10/21 Questions from 10/16 Tutorials Questions Tutorial 4 - Intro to FaaS - AWS Lambda Background on AWS Lambda for the Term Project - II From: Cloud Computing Concepts, Technology & Architecture: Chapter 4: Cloud Computing Concepts and Models: Roles and boundaries Cloud characteristics Cloud delivery models Cloud deployment models Team Planning - Breakout Rooms October 21, 2025 TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

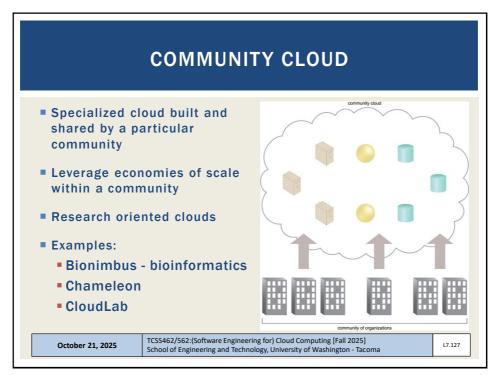
124

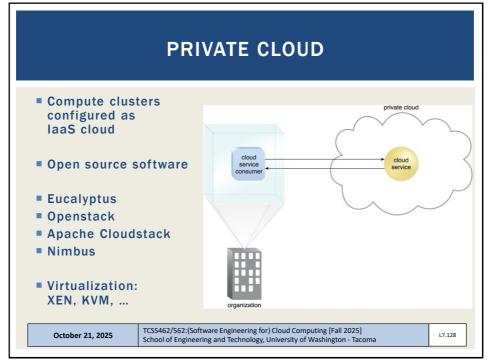


125

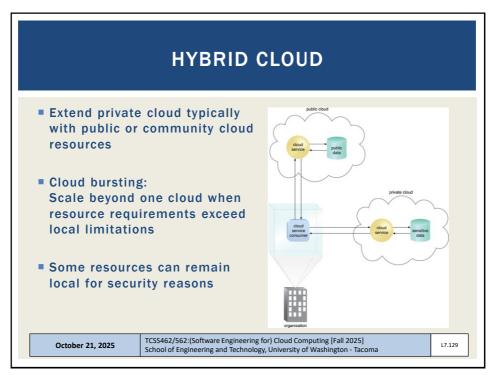


126

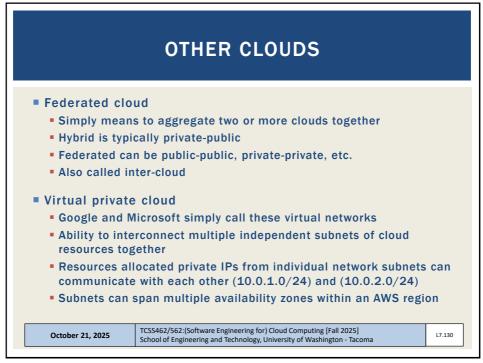




128



129



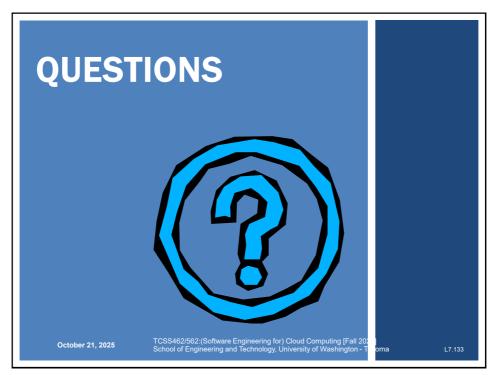
130

# OBJECTIVES - 10/21 Questions from 10/16 Tutorials Questions Tutorial 4 - Intro to FaaS - AWS Lambda Background on AWS Lambda for the Term Project - II From: Cloud Computing Concepts, Technology & Architecture: Chapter 4: Cloud Computing Concepts and Models: Roles and boundaries Cloud characteristics Cloud delivery models Cloud deployment models Team Planning - Breakout Rooms October 21, 2025 TCSS462/562:(Software Engineering for) Cloud Computing [Fail 2025] School of Engineering and Technology, University of Washington - Tacoma

131



132



133