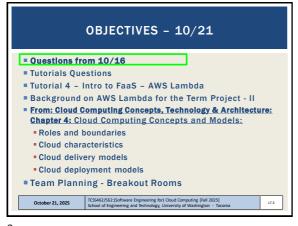
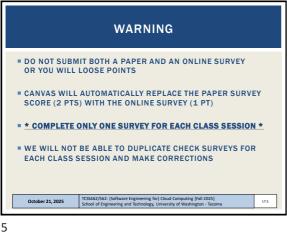


OFFICE HOURS - FALL 2025 Thursdays: 6:00 to 7:00 pm - CP 229 & Zoom Fridays ■11:00 am to 12:00 pm - ONLINE via Zoom* ■Or email for appointment Office Hours set based on Student Demographics survey feedback * - Friday office hours may be adjusted or canceled due meeting conflicts or other obligations. Adjustments will be announced via Canvas. October 21, 2025



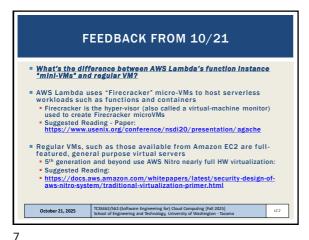
ONLINE DAILY FEEDBACK SURVEY Daily Feedback Quiz in Canvas - Take After Each Class 1-point
 Extra Credit for completing online Class Activity 1 - Implicit vs. Explicit Parallelism
Available until Oct 11 at 11:59pm | Due Oct 7 at 7:50pm | -2-points Extra Credit for completing in-person in class 36 points possible 2 5% added to final course grade for TCSS 562 - Online Daily Feedback Survey - 9/30 (36/36) TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacor October 21, 2025 L7.4

3



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 (1 - 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 L7.6

6



MICROVMS VS. VMS Firecracker Micro VMs on AWS EC2 VM No direct connections allowed. The guest OS can only be accessed via a serverless function or container code Users can directly connect using SSH for console sessio to interact with the guest OS VMs run any OS selected by the Micro VMs run only Amazon Linux 2023 user 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) ■ 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 Public/private IP addresses Internal hidden IP addresses TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Taco

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

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 (\$3) is most commonly used to persist data, but any database service can be used. Common alternatives: DynamoDB (No SQL DB) Amazon RDS & Aurora (managed relational database service) SQS - Simple Queuing Service SNS - Simple Notification Service 5 Elasticache **Document DB** Amazon MO October 21, 2025 L7.10

9

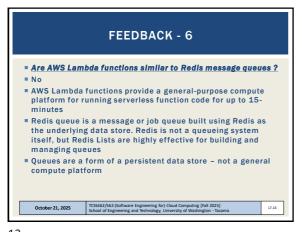
FEEDBACK - 4 Do asynchronous client calls to a server, save the client time moreso than synchronous calls? A synchronous call blocks the calling thread to wait for a result from the server If the programmer has not designed the client to be multithreaded, then the client essentially is frozen while waiting 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 TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tac October 21, 2025 L7.11

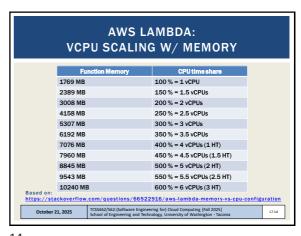
FEEDBACK - 5 When should synchronous vs. asynchronous client calls to a server be used ? - Asynchronous calls are best for long operations Maintaining a network connection for more than 30-seconds is error Example: mobile device traveling down I-5 switching cell towers Synchronous calls block the client 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 TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Taco October 21, 2025 L7.12

11 12

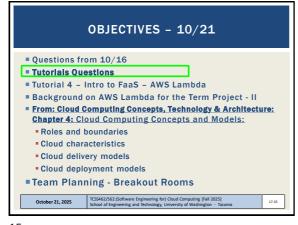
Slides by Wes J. Lloyd L7.2

8





13 14



15

TUTORIAL 2 - OCT 21	
	b Bash Scripting y.washington.edu/wlloyd/courses/tcss562/tutorials/T 12025_tutorial_2.pdf
Review tutoriaCreate a BASH1. What is a B	l webservice client
 Variables Input Arithmetic 	
5. If Statemen 6. Loops 7. Functions 8. User Interfa	
Call service to	obtain IP address & lat/long of computer it.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

TUTORIAL 3 – OCT 30 (TEAMS OF 2)

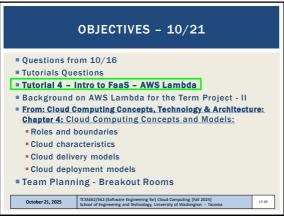
Best Practices for Working with Virtual Machines on Amazon EC2
https://faculty.washington.edu/wlloyd/courses/tcss562
//tutorials/TCSS462_562_f2025_tutorial_3.pdf
Creating a spot VM
Creating an image from a running VM
Persistent spot request
Stopping (pausing) VMs
EBS volume types
Ephemeral disks (local 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

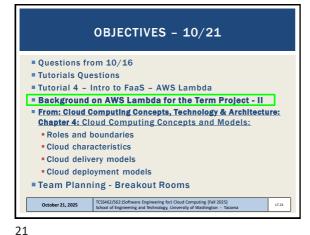
17 18

Slides by Wes J. Lloyd L7.3



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 TCSS462/562:(School of Engir L7.20

19 20



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 Cent05 9 Stream
Suggested Reading:
https://docs.aws.amazon.com/lambda/latest/dg/gettingstarted-limits.html

Cotober 21, 2025

TCSS62/S62/Software Engineering for Cloud Computing [Fail 2025] school of Engineering and Technology, University of Washington - Tacoma

21

```
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
   3. 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 8th column returned
       Metric can be read using SAAF in tutorial #4
                TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025]
School of Engineering and Technology, University of Washington - Tacoma
  October 21, 2025
                                                                        L7.23
```

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",

"vmTD": "tiUCnA",

"platform": "AWS Lambda",

"newcontainer": 1,

"cpuUsrDelta": "904",

"cpuNiceDelta": "904",

"cpuNiceDelta": "885",

"cpuIdeDelta": "882428",

"cpuIdeDelta": "882428",

"cpuIdeDelta": "7",

"vmcpuScoftIrqDelta": "7",

"vmcpuscablelta": "7",

"trameworkRuntime": 35.72,

"message": "Hello Fred Smith!",

"runtime": 38.94

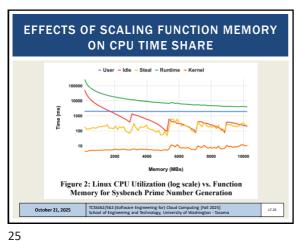
}

October 21, 2025

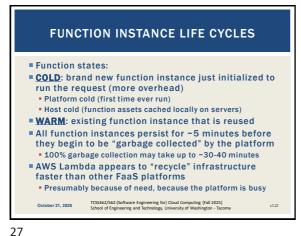
**Coctober 21, 202

23 24

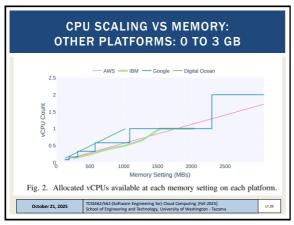
Slides by Wes J. Lloyd L7.4



EFFECTS OF SCALING FUNCTION MEMORY ON CPU TIME SHARE - User - Idle - Steal - Runtime - Kernel Key observations: Runtime decreases as vCPUs and CPU time share increase CPU user time remains constant for the prime number generation task – work doesn't change CPU idle time gradually decreases as memory and vCPUs increase (the idle time is becoming active time) When the 4th vCPU is added, cpuSteal tracks closely with cpuldle time (hyperthreading effect) There is more cpu Kernel time after the 4th vCPU is added



WARM VS COLD FUNCTION INSTANCES New F Figure 3: AWS Lambda Function Instance Replacement vs. Function Call Interval over 24-hours October 21, 2025

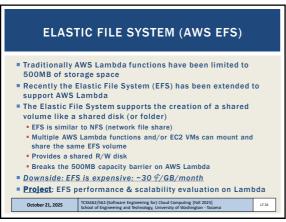


CPU SCALING VS MEMORY: OTHER PLATFORMS: 0 TO 3 GB **Key observations:** Google only supports strict memory steps AWS gradually increases the CPU time share as memory is increased IBM is similar but slope is not constant Digital Ocean only scales up to 1 GB Fig. 2. Allocated vCPUs available at each memory setting on each platform. October 21, 2025

29 30

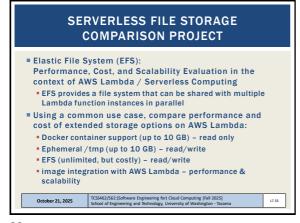
Slides by Wes J. Lloyd L7.5

26



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:(School of Engir L7.32

31 32



SERVICE COMPOSITION

API Gateway

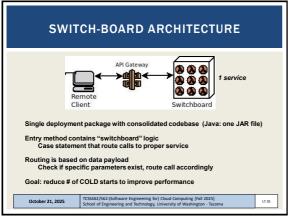
API Gateway

A B C 3 services
Full Service Isolation

A B C 2 services

Cotober 21, 2025 INCS482/502 [Software Engineering for) Cloud Computing [Fail 2025] School of Engineering and Technology, University of Washington - Tacoma 12 34

33



APPLICATION FLOW CONTROL - 3

Client flow control

APPLICATION FLOW CONTROL - 3

Client flow control

APPLICATION FLOW CONTROL - 3

Microservices

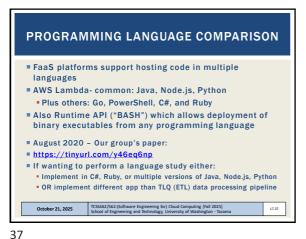
AWS Step Function

Asynchronous

Asyn

35 36

Slides by Wes J. Lloyd L7.6



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)
October 21, 2025
TCSS62/S62-Software Engineering for/ Cloud Computing [Fall 2025]
School of Engineering and Technology, University of Washington - Tacoma

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

| 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 11, 2025 | School of Engineering for Loud Computing [Fall 2025] | School of Engineering and Technology, University of Washington - Tacoms | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-40 | 12-

39

L7.39

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

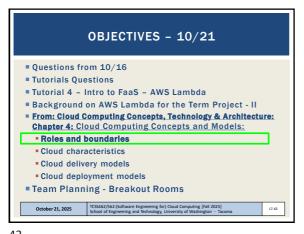
CPU STEAL CASE STUDY

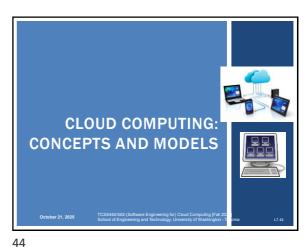
On AWS Lambda (or other FaaS platforms), when we run functions, how much CpuSteal do we observe?
How does CpuSteal vary for different workloads? (e.g., functions that have different resource requirements)
How does CpuSteal vary over time hour, day, week, location?
How does CpuSteal relate to function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

41 42

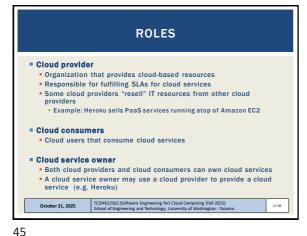
Slides by Wes J. Lloyd L7.7

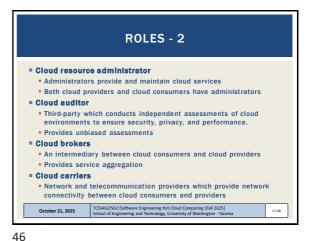
38



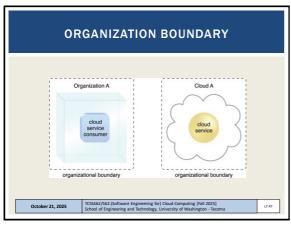


43 44





+3



TRUST BOUNDARY

trust boundary

Cloud A

Cloud A

Cloud A

Cloud A

organizational boundary

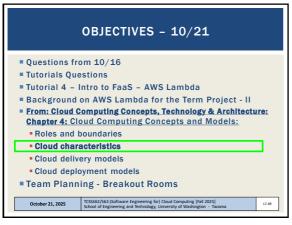
organizational boundary

October 21, 2025

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

17.48

47 48

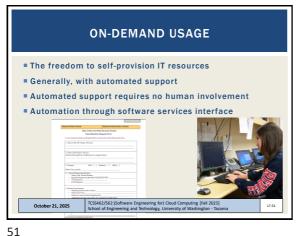


CLOUD CHARACTERISTICS

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

49



UBIQUITOUS ACCESS

Cloud services are widely accessible

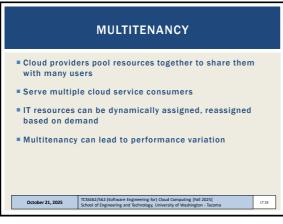
Public cloud: internet accessible

Private cloud: throughout segments of a company's intranet

24/7 availability

TCSS42/562/50t/wave Engineering for/ Cloud Computing [Fall 2025] school of Engineering and Technology, University of Washington - Tacoma

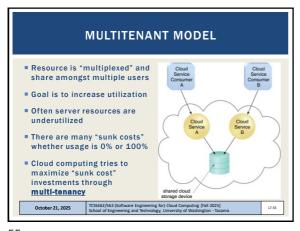
, 1

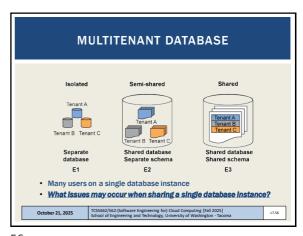


53 54

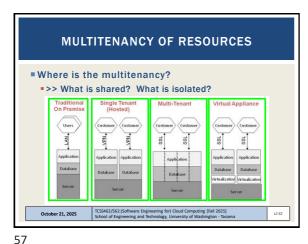
Slides by Wes J. Lloyd L7.9

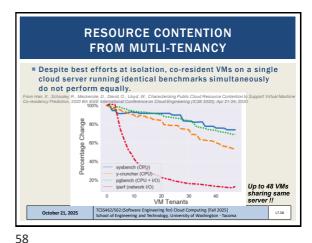
50



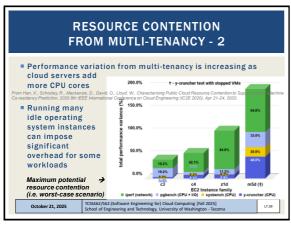


55 56





5



ELASTICITY

■ Automated ability of cloud to transparently scale resources

■ Scaling based on runtime conditions or pre-determined by cloud consumer or cloud provider

■ Threshold based scaling

• CPU-utilization > threshold_A, Response_time > 100ms

• Application agnostic vs. application specific thresholds

• Why might an application agnostic threshold be non-ideal?

■ Load prediction

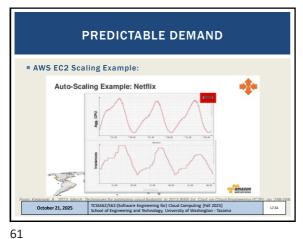
• Historical models

• Real-time trends

| TCSS42/562/Sc/Software Engineering for Cloud Computing [fail 2025]
| School of Engineering and Technology, University of Wisblington -Tacoma

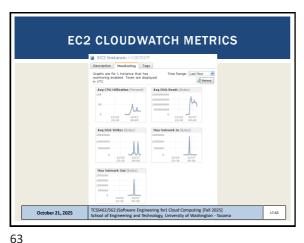
| U 500

59 60

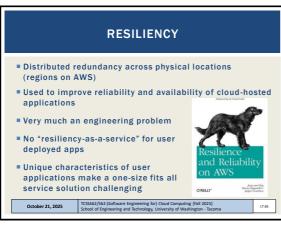


MEASURED USAGE Cloud platform tracks usage of IT resources ■ For billing purposes ■ Enables charging only for IT resources actually used Can be time-based (millisec, second, minute, hour, day) Granularity is increasing... Can be throughput-based (data transfer: MB/sec, GB/sec) Can be resource/reservation based (vCPU/hr, GB/hr) Not all measurements are for billing Some measurements can support auto-scaling ■ For example CPU utilization October 21, 2025 L7.62

62



EC2 CLOUDWATCH METRICS TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Taco October 21, 2025 L7.64



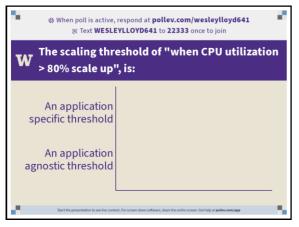
Elasticity is often provided using threshold based scaling. When can threshold based scaling (i.e. CPU utilization > 80%) under or over provision resources? TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2025]

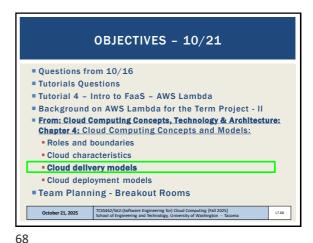
October 24, 2016

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

65 66

Slides by Wes J. Lloyd L7.11





67

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

ICSS42/S62/Software Engineering for/ Cloud Computing (Fall 2025)
School of Engineering and Technology, University of Visibilington-Taxoma

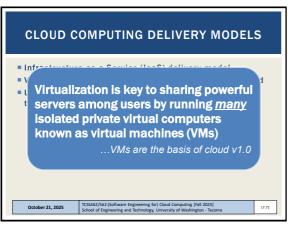
CLOUD COMPUTING DELIVERY MODELS

Infrastructure-as-a-Service (laaS) delivery model
Virtualization is a key-enabling technology of laaS cloud
Uses virtual machines to deliver cloud resources to end users

October 21, 2025

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

69



CLOUD COMPUTING DELIVERY MODELS

Infrastructure-as-a-Service (IaaS) delivery model

Virtual Machines are the building blocks for "Cloud Service Delivery Models"
They are the "vehicles" used to deliver compute resources to end users...

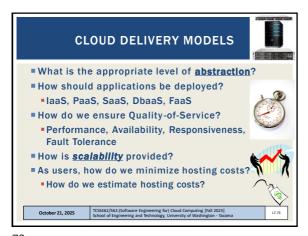
cloud 1.0

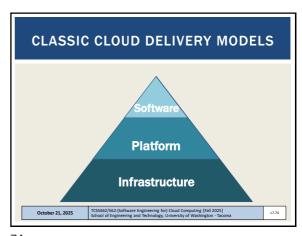
October 21, 2025

TCSS62/562/Software Engineering for) Cloud Computing [Fall 2025]
School of Engineering and Technology, University of Washington-Taxoma

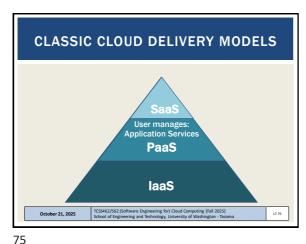
71 72

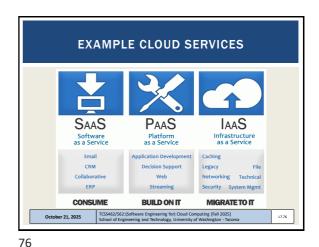
Slides by Wes J. Lloyd L7.12



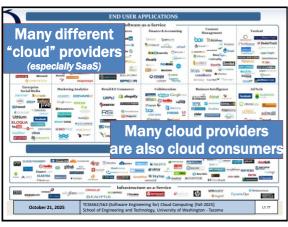


73 74





J



INFRASTRUCTURE-AS-A-SERVICE

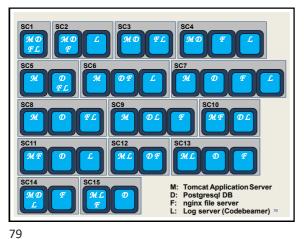
Compute resources, on demand, as-a-service
Generally raw "IT" resources
Hardware, network, containers, operating systems

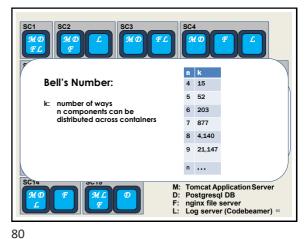
Typically provided through virtualization
Generally, not-preconfigured
Administrative burden is owned by cloud consumer
Best when high-level control over environment is needed
Scaling is generally not automatic...
Resources can be managed in bundles
AWS CloudFormation: Allows specification in JSON/YAML of cloud infrastructures

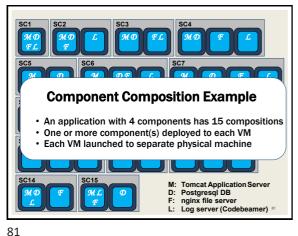
October 21, 2025

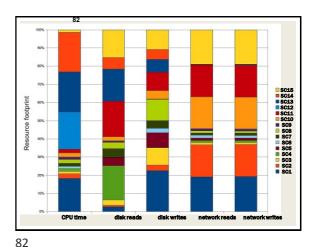
TCSS402/GG/Software Engineering for) Cloud Computing [Fall 2025]
School of Engineering and Technology, University of Washington - Tacorna

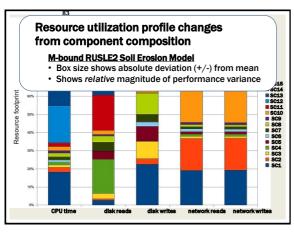
77 78





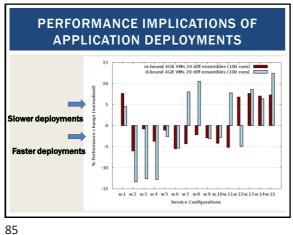




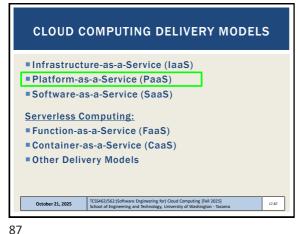


△ Resource Utilization Change Min to Max Utilization Resource footprint d-bound CPU time: 6.5% 5.5% Disk sector reads: 14.8% 819.6% Disk sector writes: 21.8% 111.1% Network bytes received: 144.9% 145% Network bytes sent: 143.7% 143.9% CPU time

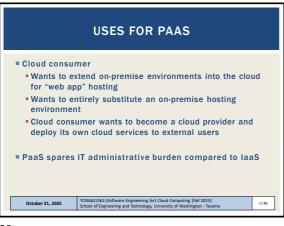
83 84



PERFORMANCE IMPLICATIONS OF APPLICATION DEPLOYMENTS **△** Performance Change: Min to max performance Sid M-bound: 14% 25.7% D-bound: F sc1 sc2 sc3 sc4 sc5 sc6 sc7 sc8 sc9 sc10sc11sc12sc13sc14sc15 Service Configurations



PLATFORM-AS-A-SERVICE ■ Predefined, ready-to-use, hosting environment Infrastructure is further obscured from end user. Scaling and load balancing may be automatically provided and automatic Variable to no ability to influence responsiveness ■ Examples: ■ Google App Engine ■ Heroku AWS Elastic Beanstalk AWS Lambda (FaaS) October 21, 2025



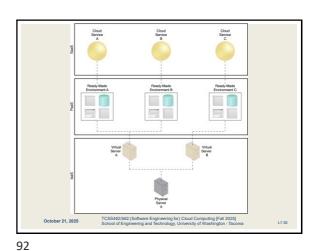
CLOUD COMPUTING DELIVERY MODELS ■Infrastructure-as-a-Service (IaaS) ■ Platform-as-a-Service (PaaS) Software-as-a-Service (SaaS) **Serverless Computing:** ■ Function-as-a-Service (FaaS) ■ Container-as-a-Service (CaaS) Other Delivery Models October 21, 2025

89 90

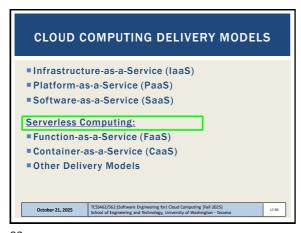
Slides by Wes J. Lloyd L7.15

86



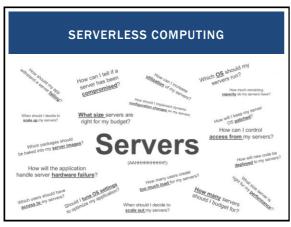


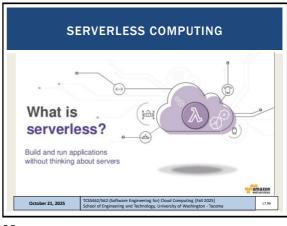
31



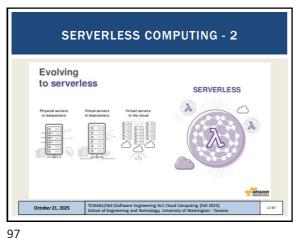


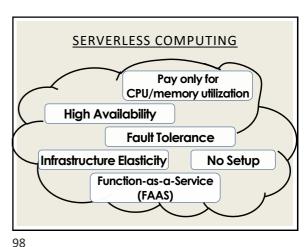
93

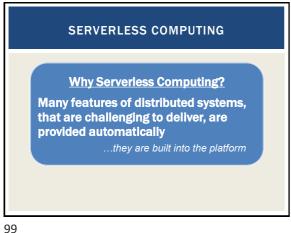




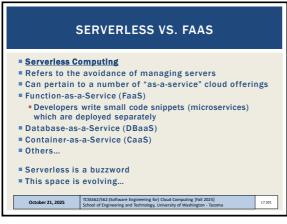
95 96

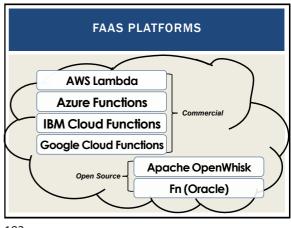






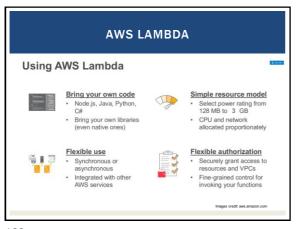
CLOUD COMPUTING DELIVERY MODELS ■ Infrastructure-as-a-Service (IaaS) ■ Platform-as-a-Service (PaaS) ■ Software-as-a-Service (SaaS) **Serverless Computing:** ■ Function-as-a-Service (FaaS) ■ Container-as-a-Service (CaaS) Other Delivery Models TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Taco October 21, 2025

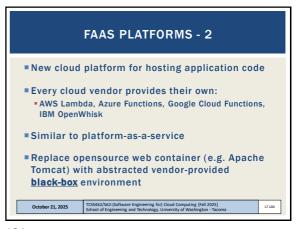




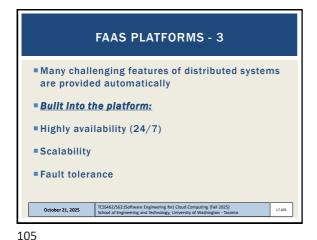
101 102

Slides by Wes J. Lloyd L7.17





103 104



CLOUD NATIVE SOFTWARE ARCHITECTURE

Every service with a different pricing model

Example: Weather Application

Lambda in Engaged

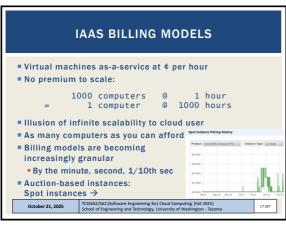
APPLOATEWAY

TCSS462/562/25/25/56/tware Engineering for) Cloud Computing [Fall 2025]

School of Engineering and Technology, University of Washington - Taxoma

Lambda vans code to refuse food washer eight and refuse dad book to user

103



PRICING OBFUSCATION

■ VM pricing: hourly rental pricing, billed to nearest second is intuitive...

■ FaaS pricing: non-intuitive pricing policies

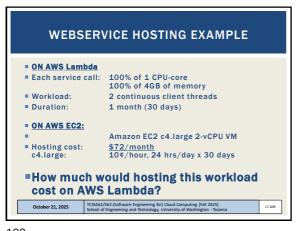
■ FREE TIER: first 1,000,000 function calls/month → FREE first 400,000 GB-sec/month → FREE

■ Afterwards: obfuscated pricing (AWS Lambda): \$0.000002 per request \$0.0000028 to rent 128MB / 100-ms \$0.00001667 GB / second

■ CCSS462/562/Software Engineering for) Cloud Computing [Fail 2025] School of Engineering and Technology, University of Washington - Taccma

107 108

Slides by Wes J. Lloyd L7.18



PRICING OBFUSCATION

Worst-case scenario = ~2.32x!

AWS EC2: \$72.00

AWS Lambda: \$167.01

Break Even: 4,319,136 GB-sec

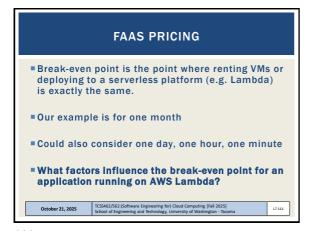
Two threads

@2GB-ea: ~12.5 days

BREAK-EVEN POINT: ~4,319,136 GB-sec-month

~12.5 days 2 concurrent clients @ 2GB

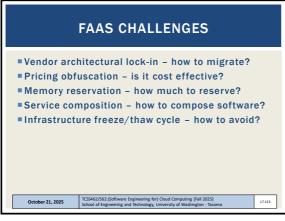
109



FACTORS IMPACTING PERFORMANCE OF FAAS COMPUTING PLATFORMS

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

111



VENDOR ARCHITECTURAL LOCK-IN

■ Cloud native (FaaS) software architecture requires external services/components

Example: Weather Application

Client

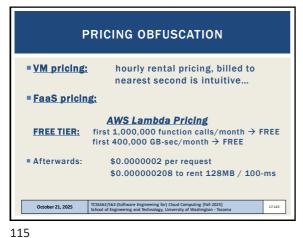
Lambda is triggered

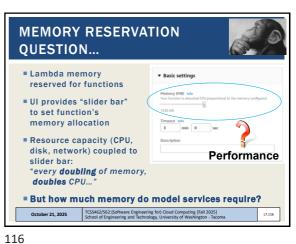
Lambda is reference and the reference did both to same section reference did both to same section reference and processed to the reference did both to same images credit two amazons come

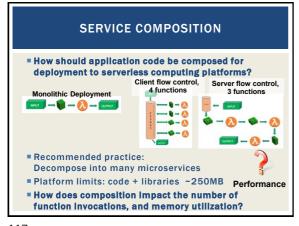
■ Increased dependencies → increased hosting costs

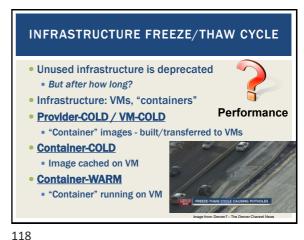
113 114

Slides by Wes J. Lloyd L7.19

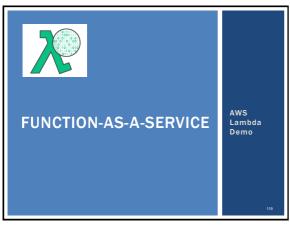








117



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 October 21, 2025

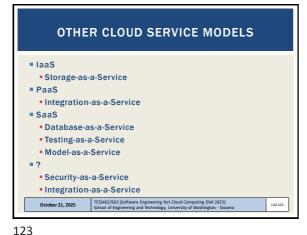
119 120



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

121 122



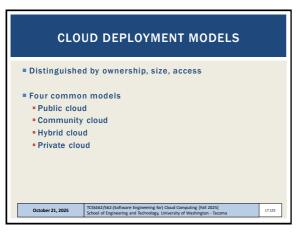
OBJECTIVES - 10/21

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 delivery models
Team Planning - Breakout Rooms

October 21, 2025

TCS462/S62: Software Engineering for Cloud Computing [Fall 2025]
School of Engineering and Technology, University of Washington - Tacoma

123



PUBLIC CLOUDS

Salestons

Mercent

Manager

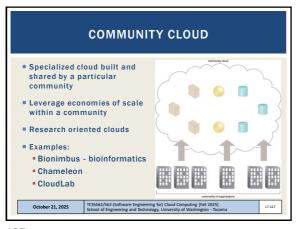
Analogo

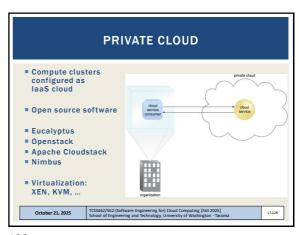
Particle (Code Computative Fig. 12025)
School of Engineering and Technology, University of Washington - Tacoma

12.126

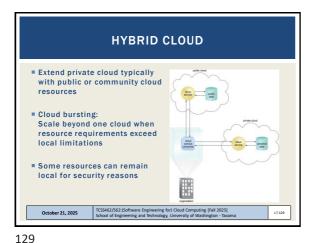
125 126

Slides by Wes J. Lloyd L7.21





127 128



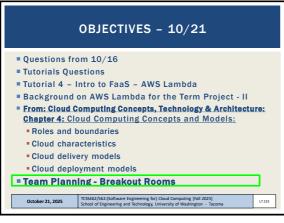
Pederated cloud
Simply means to aggregate two or more clouds together
Hybrid is typically private-public
Federated can be public-public, private-private, etc.
Also called inter-cloud

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

October 21, 2025

TCSS462/562/Scholwave Engineering for) Cloud Computing [fail 2025]
School of Engineering and Technology, University of Washington-Tacoma

.29 130



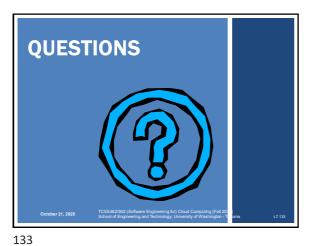
TCSS 462/562
TERM PROJECT

TCSS462862 (Software Engineering for) Cloud Computing [Fall 202 school of Engineering and Technology, University of Washington: To cona 17:132

131 132

[Fall 2025]

TCSS 462: Cloud Computing TCSS 562: Software Engineering for Cloud Computing School of Engineering and Technology, UW-Tacoma



L7.23 Slides by Wes J. Lloyd