


# TCSS 562: SOFTWARE ENGINEERING FOR CLOUD COMPUTING

## Cloud Computing Concepts and Models - II



Wes J. Lloyd  
 School of Engineering and Technology  
 University of Washington - Tacoma

1

## OFFICE HOURS - FALL 2024

- **Tuesdays:**
  - 2:30 to 3:30 pm - CP 229
- **Friday - THIS WEEK**
  - 1:00 pm to 2:00 pm - ONLINE via Zoom
- Or email for appointment

> Office Hours set based on Student Demographics survey feedback

2


## OBJECTIVES - 10/22

- **Questions from 10/17**
- Tutorials Questions
- Tutorial 5 - to be posted...
- From: Cloud Computing Concepts, Technology & Architecture: Chapter 4: Cloud Computing Concepts and Models:
  - Cloud computing delivery models
  - Cloud deployment models
- AWS Overview and demo
- **2<sup>nd</sup> hour:**
  - Activity 2 - Horizontal Scaling in the Cloud
  - Term Project Planning

3

## ONLINE DAILY FEEDBACK SURVEY

- Daily Feedback Quiz in Canvas - Take After Each Class
- Extra Credit for completing



4

### TCSS 562 - Online Daily Feedback Survey - 10/5

Started: Oct 7 at 1:13am

#### Quiz Instructions

Question 1 (0.5 pts)

On a scale of 1 to 10, please classify your perspective on material covered in today's class:

1 2 3 4 5 6 7 8 9 10

Mostly Review To Me      Equal New and Review      Mostly New To Me

Question 2 (0.5 pts)

Please rate the pace of today's class:

1 2 3 4 5 6 7 8 9 10

Slow      Just Right      Fast

5

## MATERIAL / PACE

- Please classify your perspective on material covered in today's class (**45** respondents):
  - ew, 10-mostly new
  - **Average - 6.01** (↓ - previous 6.50)
- Please rate the pace of today's class:
  - 1-slow, 5-just right, 10-fast
  - **Average - 5.24** (↓ - previous 5.59)
- **Response rates:**
  - TCSS 462: 29/42 - 69.05%
  - TCSS 562: 16/20 - 80.0%

6

### FEEDBACK FROM 10/17

- Will multi-tenancy slow down operations and increase costs?

October 22, 2024	TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma	L8.7
------------------	--	------

7

### IN CLASS QUIZZES

- Anticipated dates
- Designed for 1 hour (starting at 4:40pm)
- BHS 106 Room is available, so professor will stay late to allow additional time
- Open notes & books
- Closed laptop, smartphone, neighbor
- Quiz 1 - Tuesday November 5
- Quiz 2 - Tuesday November 26

October 22, 2024	TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma	L8.8
------------------	--	------

8

### AWS CLOUD CREDITS UPDATE

- AWS CLOUD CREDITS ARE NOW AVAILABLE FOR TCSS 462/562
- Credit codes must be securely exchanged
- Request codes by sending an email with the subject "AWS CREDIT REQUEST" to [wllloyd@uw.edu](mailto:wllloyd@uw.edu)
- Codes can also be obtained in person (or zoom), in the class, during the breaks, after class, during office hours, by appt
  - 41 credit requests fulfilled as of Oct 21 @ 11:59p
- To track credit code distribution, codes not shared via discord
- 46 of 62 students have completed AWS Cloud Credits Survey
  - 16 survey responses missing ???
- Are all students able to create AWS accounts ?**
- Tutorial 3 is due October 31st
  - OCT 31 is also a SOFT Deadline to request cloud computing credits
  - If you do not request by this date, and complete tutorial 3, you may experience cloud computing charges

October 10, 2023	TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma	L4.9
------------------	--	------

9

### OBJECTIVES - 10/22

- Questions from 10/17
- Tutorials Questions**
- Tutorial 5 - to be posted...
- From: Cloud Computing Concepts, Technology & Architecture: Chapter 4: Cloud Computing Concepts and Models:
  - Cloud computing delivery models
  - Cloud deployment models
- AWS Overview and demo
- 2<sup>nd</sup> hour:
  - Activity 2 - Horizontal Scaling in the Cloud
  - Term Project Planning

October 22, 2024	TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma	L8.10
------------------	--	-------

10

### TUTORIAL 0

- Getting Started with AWS
- [https://faculty.washington.edu/wllloyd/courses/tcss562/tutorials/TCSS462\\_562\\_f2024\\_tutorial\\_0.pdf](https://faculty.washington.edu/wllloyd/courses/tcss562/tutorials/TCSS462_562_f2024_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

October 22, 2024	TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma	L7.11
------------------	--	-------

11

### TUTORIAL 2 - DUE OCT 19 (CLOSES OCT 23)

- Introduction to Bash Scripting**
- [https://faculty.washington.edu/wllloyd/courses/tcss562/tutorials/TCSS462\\_562\\_f2024\\_tutorial\\_2.pdf](https://faculty.washington.edu/wllloyd/courses/tcss562/tutorials/TCSS462_562_f2024_tutorial_2.pdf)
- Review tutorial sections:
- Create a BASH webservice client
  - What is a BASH script?
  - Variables
  - Input
  - Arithmetic
  - If Statements
  - Loops
  - Functions
  - User Interface
- Call service to obtain IP address & lat/long of computer
- Call weatherbit.io API to obtain weather forecast for lat/long

October 11, 2022	TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma	L4.12
------------------	--	-------

12

### TUTORIAL 3 – DUE OCT 31

- Best Practices for Working with Virtual Machines on Amazon EC2
- [https://faculty.washington.edu/wlloyd/courses/tcss562/tutorials/TCSS462\\_562\\_f2024\\_tutorial\\_3.pdf](https://faculty.washington.edu/wlloyd/courses/tcss562/tutorials/TCSS462_562_f2024_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 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
 School of Engineering and Technology, University of Washington - Tacoma    L7.13

13

### TUTORIAL 4 – DUE NOV 5

- Introduction to AWS Lambda with the Serverless Application Analytics Framework (SAAF)
- [https://faculty.washington.edu/wlloyd/courses/tcss562/tutorials/TCSS462\\_562\\_f2024\\_tutorial\\_4.pdf](https://faculty.washington.edu/wlloyd/courses/tcss562/tutorials/TCSS462_562_f2024_tutorial_4.pdf)
- Obtaining a Java development environment
- 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
- Performance analysis using faas\_runner reports
- Two function pipeline development task

October 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
 School of Engineering and Technology, University of Washington - Tacoma    L8.14

14

### OBJECTIVES – 10/22

- Questions from 10/17
- Tutorials Questions
- Tutorial 5 - to be posted...**
- From: Cloud Computing Concepts, Technology & Architecture: Chapter 4: Cloud Computing Concepts and Models:
  - Cloud computing delivery models
  - Cloud deployment models
- AWS Overview and demo
- 2<sup>nd</sup> hour:
  - Activity 2 - Horizontal Scaling in the Cloud
  - Term Project Planning

October 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
 School of Engineering and Technology, University of Washington - Tacoma    L8.15

15

### CATCH UP – 10/17

- Questions from 10/15
- 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 17, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
 School of Engineering and Technology, University of Washington - Tacoma    L7.16

16

### 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 2024]  
 School of Engineering and Technology, University of Washington - Tacoma    L7.17

17

### MULTITENANCY OF RESOURCES

- Where is the multitenancy?
  - >> What is shared? What is isolated?

October 17, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
 School of Engineering and Technology, University of Washington - Tacoma    L7.18

18

### RESOURCE CONTENTION FROM MUTLI-TENANCY

- Despite best efforts at isolation, co-resident VMs on a single cloud server running identical benchmarks simultaneously do not perform equally.

From Han, X., Schooley, R., Mackenzie, D., David, O., Lloyd, W., Characterizing Public Cloud Resource Contention to Support Virtual Machine Co-residency Prediction, 2020 8th IEEE International Conference on Cloud Engineering (IC2E 2020), Apr 21-24, 2020.

**Up to 48 VMs sharing same server!!**

October 17, 2024 | TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] | School of Engineering and Technology, University of Washington - Tacoma | L7.19

19

### RESOURCE CONTENTION FROM MUTLI-TENANCY - 2

- Performance variation from multi-tenancy is increasing as cloud servers add more CPU cores
- Running many idle operating system instances can impose significant overhead for some workloads

From Han, X., Schooley, R., Mackenzie, D., David, O., Lloyd, W., Characterizing Public Cloud Resource Contention to Support Virtual Machine Co-residency Prediction, 2020 8th IEEE International Conference on Cloud Engineering (IC2E 2020), Apr 21-24, 2020.

**Maximum potential resource contention (i.e. worst-case scenario)**

October 17, 2024 | TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] | School of Engineering and Technology, University of Washington - Tacoma | L7.20

20

### 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

October 17, 2024 | TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] | School of Engineering and Technology, University of Washington - Tacoma | L7.21

21

### PREDICTABLE DEMAND

- AWS EC2 Scaling Example:

**Auto-Scaling Example: Netflix**

From: Kawanishi, A. 2013. Month. Techniques for optimizing cloud footprint. In: 2013 IEEE Int. Conf. on Cloud Engineering (IC2E), pp. 258-268.

October 17, 2024 | TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] | School of Engineering and Technology, University of Washington - Tacoma | L7.22

22

### 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 2024] | School of Engineering and Technology, University of Washington - Tacoma | L7.23

23

### EC2 CLOUDWATCH METRICS

October 17, 2024 | TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] | School of Engineering and Technology, University of Washington - Tacoma | L7.24

24

### EC2 CLOUDWATCH METRICS

October 17, 2024 TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma L7.25

25

### RESILIENCY

- Distributed redundancy across physical locations (regions on AWS)
- Used to improve reliability and availability of cloud-hosted applications
- Very much an engineering problem
- No "resiliency-as-a-service" for user deployed apps
- Unique characteristics of user applications make a one-size fits all service solution challenging

October 17, 2024 TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma L7.26

26

### Elasticity is often provided using threshold based scaling. When can threshold based scaling (i.e. CPU utilization > 80%) under or over provision resources?

When the application is primarily I/O bound, a CPU threshold may never be met, or be met too late to scale up.

When the current resource utilization does not reflect future system demand.

When the current resource utilization (e.g. CPU) is temporarily increased as a result of external factors (i.e. resource contention from other tasks) that does not correlate to system demand.

When an application will soon complete a parallel phase, before executing a largely sequential phase.

All of the above

A B C D E

October 24, 2016 TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma L10.2

27

When poll is active, respond at [polllev.com/wesleyloyd641](https://polllev.com/wesleyloyd641)  
 Text WESLEYLOYD641 to 22333 once to join

### The scaling threshold of "when CPU utilization > 80% scale up", is:

An application specific threshold

An application agnostic threshold

Start the presentation to see live content. For screen share software, share the entire screen. Get help at [polllev.com/app](https://polllev.com/app)

28

## CLOUD COMPUTING: CONCEPTS AND MODELS

October 22, 2024 TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma L8.29

29

### OBJECTIVES - 10/22

- Questions from 10/17
- Tutorials Questions
- Tutorial 5 - to be posted...
- From: **Cloud Computing Concepts, Technology & Architecture: Chapter 4: Cloud Computing Concepts and Models:**
  - **Cloud computing delivery models**
  - Cloud deployment models
- AWS Overview and demo
- 2<sup>nd</sup> hour:
  - Activity 2 - Horizontal Scaling in the Cloud
  - Term Project Planning

October 22, 2024 TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma L8.30

30

### 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 22, 2024
TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2024]  
School of Engineering and Technology, University of Washington - Tacoma
L8.31

31

### CLOUD COMPUTING DELIVERY MODELS

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

October 22, 2024
TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2024]  
School of Engineering and Technology, University of Washington - Tacoma
L8.32

32

### CLOUD COMPUTING DELIVERY MODELS

**Virtualization is key to sharing powerful servers among users by running *many* isolated private virtual computers known as virtual machines (VMs)**

*...VMs are the basis of cloud v1.0*

October 22, 2024
TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2024]  
School of Engineering and Technology, University of Washington - Tacoma
L8.33


33

### CLOUD COMPUTING DELIVERY MODELS

**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 22, 2024
TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2024]  
School of Engineering and Technology, University of Washington - Tacoma
L8.34

34

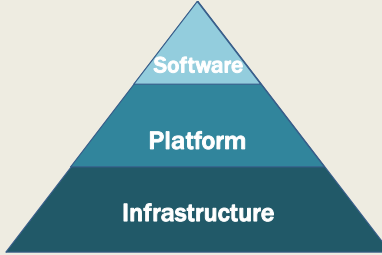
### CLOUD COMPUTING DELIVERY MODELS

- What is the appropriate level of **abstraction**?
- How should applications be deployed?
  - IaaS, PaaS, SaaS, DbaaS, FaaS
- How do we ensure Quality-of-Service?
  - Performance, Availability, Responsiveness, Fault Tolerance
- How is **scalability** provided?
- As users, how do we minimize hosting costs?
  - How do we estimate hosting costs?

October 22, 2024
TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2024]  
School of Engineering and Technology, University of Washington - Tacoma
L8.35

35

### CLASSIC CLOUD COMPUTING DELIVERY MODELS



October 22, 2024
TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2024]  
School of Engineering and Technology, University of Washington - Tacoma
L8.36

36

### CLASSIC CLOUD COMPUTING DELIVERY MODELS

**SaaS**  
 User manages:  
 Application Services  
**PaaS**  
**IaaS**

October 22, 2024 TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma L8.37

37

### EXAMPLE CLOUD SERVICES

SaaS Software as a Service	PaaS Platform as a Service	IaaS Infrastructure as a Service
Email CRM Collaborative ERP	Application Development Decision Support Web Streaming	Caching Legacy File Networking Technical Security System Mgmt
<b>CONSUME</b>	<b>BUILD ON IT</b>	<b>MIGRATE TO IT</b>

October 22, 2024 TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma L8.38

38

### END USER APPLICATIONS

Many different "cloud" providers (especially SaaS)

Many cloud providers are also cloud consumers

October 22, 2024 TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma L8.39

39

### 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**: Scripts to specify creation of cloud infrastructures using JSON/YAML for app deployment

October 22, 2024 TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma L8.40

40

### IAAS: VIRTUAL MACHINE PLACEMENT IN THE CLOUD

October 22, 2024 TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma L8.41

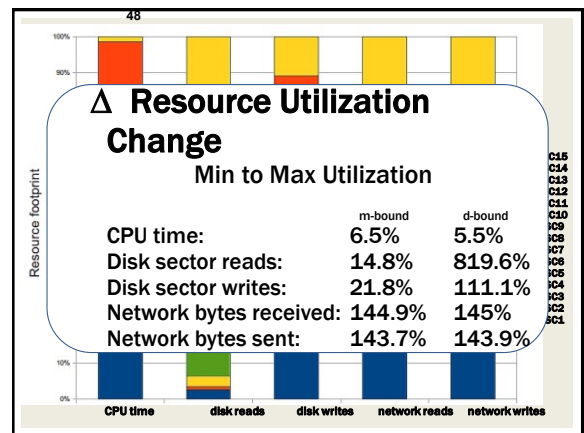
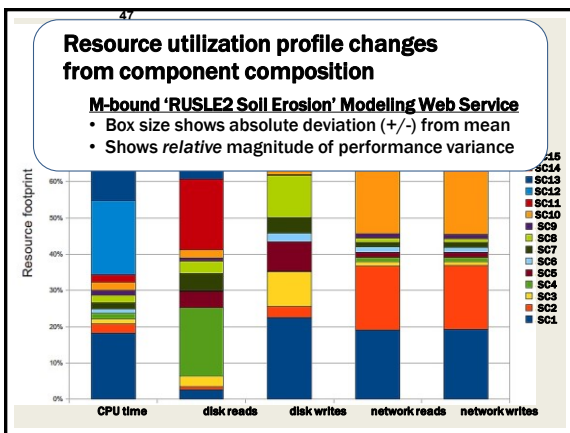
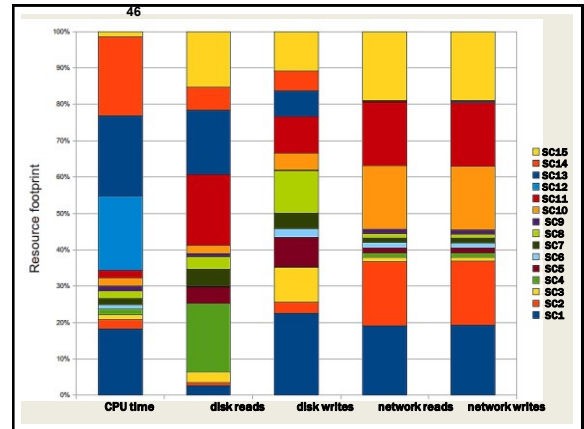
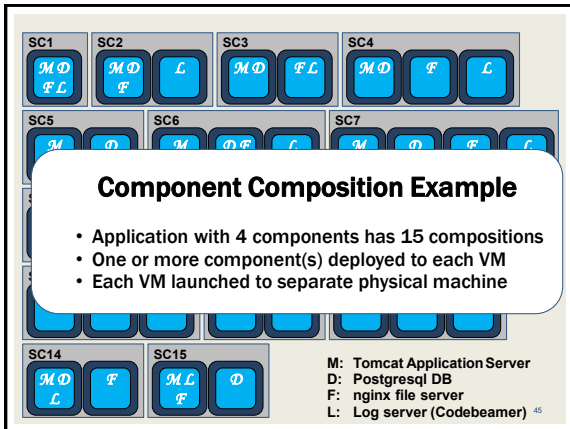
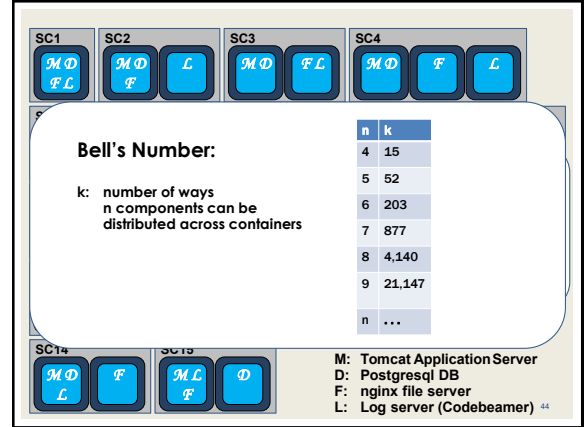
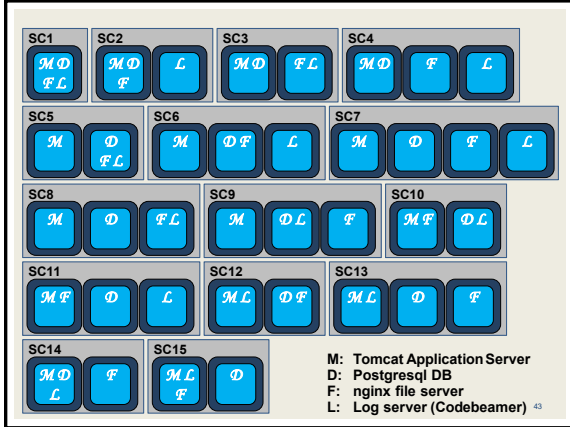
41

### COMPONENT COMPOSITION

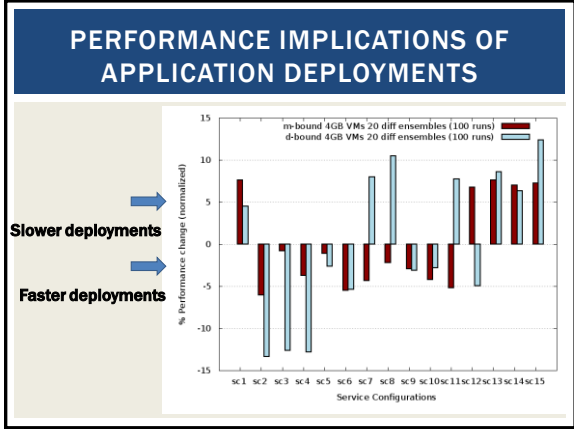
- Cloud provider maps VMs to physical servers
- User controls mapping of services to VM images
  - Should components be separated (isolated) ?
  - Should components be combined ?
  - M - modeling web service
  - D - relational database
  - F - file server
  - L - logging server

October 22, 2024 TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma L8.42

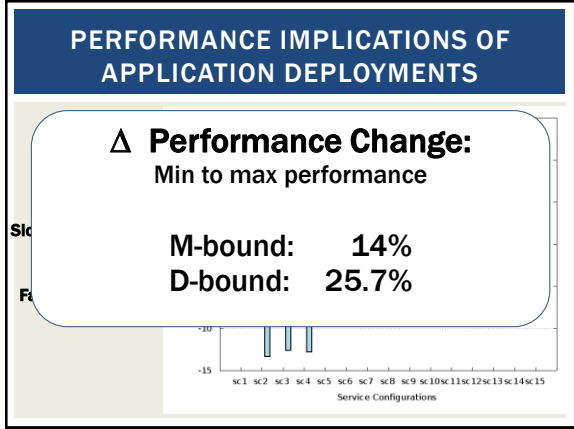
42







49



50

### 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 22, 2024 | TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] | School of Engineering and Technology, University of Washington - Tacoma | LR.51

51

### 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 22, 2024 | TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] | School of Engineering and Technology, University of Washington - Tacoma | LR.52

52

### 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 IaaS

October 22, 2024 | TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] | School of Engineering and Technology, University of Washington - Tacoma | LR.53

53

### 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 22, 2024 | TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] | School of Engineering and Technology, University of Washington - Tacoma | LR.54

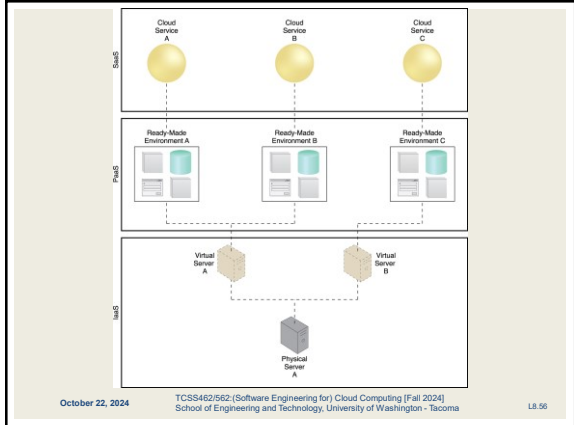
54

## SOFTWARE-AS-A-SERVICE

- Software applications as shared cloud service
- Nearly all server infrastructure management is abstracted away from the user
- Software is generally configurable
- SaaS can be a complete GUI/UI based environment
- Or UI-free (database-as-a-service)
- SaaS offerings
  - Google Docs
  - Office 365
  - Cloud9 Integrated Development Environment
  - Salesforce

October 22, 2024
TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
School of Engineering and Technology, University of Washington - Tacoma
LR.55

55



56

## 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 22, 2024
TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
School of Engineering and Technology, University of Washington - Tacoma
LR.57

57

## SERVERLESS COMPUTING

Introducing Cloud 2.0

### Serverless Computing

Deploy Applications Without Fiddling With Servers

October 22, 2024
TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
School of Engineering and Technology, University of Washington - Tacoma
LR.58

58

## SERVERLESS COMPUTING

Servers

(AHHHHHHHH!!!)

How should my app withstand a server failure?

How can I tell if a server has been compromised?

How can I increase utilization of my servers?

Which OS should my servers run?

When should I decide to scale up my servers?

What size servers are right for my budget?

How should I implement dynamic configuration changes on my servers?

How much remaining capacity do my servers have?

Which packages should be baked into my server images?

How will the application handle server hardware failure?

How will I keep my server OS patched?

How can I control access from my servers?

How will new code be deployed to my servers?

Which users should have access to my servers?

Should I tune OS settings to optimize my application?

How many users create too much load for my servers?

How many servers should I budget for?

When should I decide to scale out my servers?

How will new code be deployed to my servers?

What size server is right for my performance?

October 22, 2024
TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
School of Engineering and Technology, University of Washington - Tacoma
LR.60

59

## SERVERLESS COMPUTING

### What is serverless?

Build and run applications without thinking about servers

October 22, 2024
TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
School of Engineering and Technology, University of Washington - Tacoma
LR.60

60

## SERVERLESS COMPUTING - 2

**Evolving to serverless**

Physical servers in datacenters | Virtual servers in datacenters | Virtual servers in the cloud

**SERVERLESS**

October 22, 2024 | TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L8.61

61

## SERVERLESS COMPUTING

**Pay only for CPU/memory utilization**

**High Availability**

**Fault Tolerance**

**Infrastructure Elasticity** | **No Setup**

**Function-as-a-Service (FAAS)**

62

## SERVERLESS COMPUTING

**Why Serverless Computing?**

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

*...they are built into the platform*

63

## 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 22, 2024 | TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L8.64

64

## SERVERLESS VS. FAAS

- **Serverless Computing**
- Refers to the avoidance of managing servers
- Can pertain to a number of "as-a-service" cloud offerings
- **Function-as-a-Service (FaaS)**
  - Developers write small code snippets (microservices) which are deployed separately
- Database-as-a-Service (DBaaS)
- Container-as-a-Service (CaaS)
- Others...
- Serverless is a buzzword
- This space is evolving...

October 22, 2024 | TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L8.65

65

## FAAS PLATFORMS

**Commercial**

- AWS Lambda
- Azure Functions
- IBM Cloud Functions
- Google Cloud Functions

**Open Source**

- Apache OpenWhisk
- Fn (Oracle)

66

## AWS LAMBDA

### Using AWS Lambda

**Bring your own code**

- Node.js, Java, Python, C#,
- Bring your own libraries (even native ones)

**Simple resource model**

- Select power rating from 128 MB to 3 GB
- CPU and network allocated proportionately

**Flexible use**

- Synchronous or asynchronous
- Integrated with other AWS services

**Flexible authorization**

- Securely grant access to resources and VPCs
- Fine-grained control for invoking your functions

Images credit: aws.amazon.com

67

## 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 22, 2024
TCCS462/562: (Software Engineering for) Cloud Computing [Fall 2024]  
School of Engineering and Technology, University of Washington - Tacoma
L8.68

68

## FAAS PLATFORMS - 3

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

October 22, 2024
TCCS462/562: (Software Engineering for) Cloud Computing [Fall 2024]  
School of Engineering and Technology, University of Washington - Tacoma
L8.69

69

## CLOUD NATIVE SOFTWARE ARCHITECTURE

▪ Every service with a different pricing model

October 22, 2024
TCCS462/562: (Software Engineering for) Cloud Computing [Fall 2024]  
School of Engineering and Technology, University of Washington - Tacoma
L8.70

70

## IAAS BILLING MODELS

- Virtual machines as-a-service at \$ per hour
- No premium to scale:

$$= \frac{1000 \text{ computers}}{1 \text{ computer}} @ \frac{1 \text{ hour}}{1000 \text{ hours}}$$

- Illusion of infinite scalability to cloud user
- As many computers as you can afford
- Billing models are becoming increasingly granular
  - By the minute, second, 1/10th sec
- Auction-based instances: Spot instances →

October 22, 2024
TCCS462/562: (Software Engineering for) Cloud Computing [Fall 2024]  
School of Engineering and Technology, University of Washington - Tacoma
L8.71

71

## 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):**
  - \$0.0000002 per request
  - \$0.00000208 to rent 128MB / 100-ms
  - \$0.0001667 GB / second

October 22, 2024
TCCS462/562: (Software Engineering for) Cloud Computing [Fall 2024]  
School of Engineering and Technology, University of Washington - Tacoma
L8.72

72

### 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)

- **ON AWS EC2:** Amazon EC2 c5.large 2-vCPU VM x 4GB
- 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?**

October 22, 2024 | TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L8.73

73

### PRICING OBFUSCATION

Assume 1 month = 30.41667 days (365d / 12)

Workload: (4 GB) 10,512,000 GB-sec

Worst-case FaaS scenario = ~2.72x !

AWS EC2: \$62.05

AWS Lambda: \$168.91

Break Even: 3,702,459 GB-sec @4GB ~10.71 days

BREAK-EVEN POINT: \$62.05 - \$0.33 (calls) = \$61.72

\$61.72 / .00001667 GB-sec = ~3,702,459 GB-sec-mon/4GB/call = ~925,614 sec or ~10.71 days

Point at which using FaaS costs the same as IaaS

October 22, 2024 | TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L8.74

74

### 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 22, 2024 | TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L8.75

75

### 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?

October 22, 2024 | TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L8.76

76

### VENDOR ARCHITECTURAL LOCK-IN

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

Example: Weather Application

Front-end code for weather app hosted in S3 | User clicks on link to get local weather information | App makes REST API call to endpoint | Lambda is triggered | Lambda runs code to retrieve local weather information and returns data back to user

Images credit: aws.amazon.com

▪ **Increased dependencies → increased hosting costs**

October 22, 2024 | TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L8.77

77

### PRICING OBFUSCATION

- **VM pricing:** hourly rental pricing, billed to nearest second is intuitive...
- **FaaS pricing:**

AWS Lambda Pricing

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

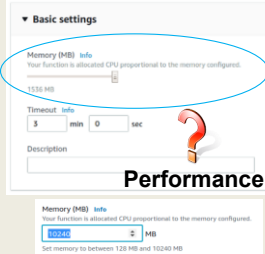
- Afterwards: \$0.0000002 per request  
 \$0.000000208 to rent 128MB / 100-ms

October 22, 2024 | TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L8.78

78

## MEMORY RESERVATION QUESTION...

- Lambda memory reserved for functions
- UI provides text box formerly "slider bar" to set function's memory
- Resource capacity (CPU, disk, network) coupled to slider bar: "every doubling of memory, doubles CPU..."
- But how much memory do FaaS functions require?

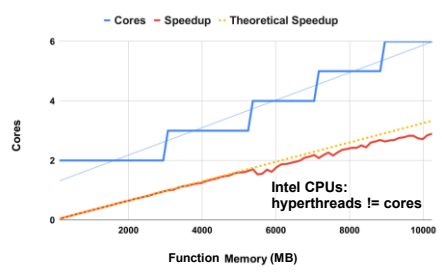


**Performance**

October 22, 2024 | TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] | School of Engineering and Technology, University of Washington - Tacoma | L8.79

79

## AWS LAMBDA COUPLES FUNCTION MEMORY TO CPU CORES & TIME SHARE



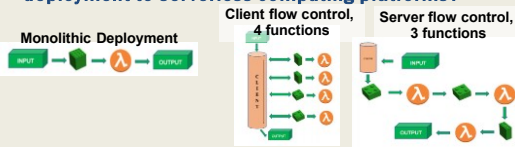
**Performance**

October 22, 2024 | TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] | School of Engineering and Technology, University of Washington - Tacoma | L8.80

80

## SERVICE COMPOSITION

- How should application code be composed for deployment to serverless computing platforms?



- Recommended practice: Decompose into many microservices
- Platform limits: code + libraries ~250MB
- How does composition impact the number of function invocations, and memory utilization?


**Performance**

October 22, 2024 | TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] | School of Engineering and Technology, University of Washington - Tacoma | L8.81

81

## INFRASTRUCTURE FREEZE/THAW CYCLE

- Unused infrastructure is deprecated
  - But after how long? (varies by platform)
- Infrastructure: microVMs (on AWS Lambda), containers on some platforms
- COLD**
  - Code image - built/transferred to physical host & cached
- WARM**
  - Host has local code cache - create function instance (microVM) on host
- HOT**
  - Function instance ready to use



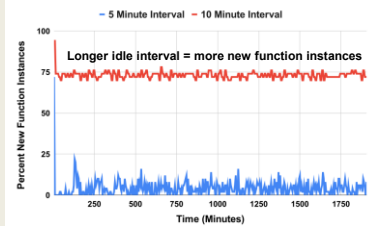
**Performance**

October 22, 2024 | TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] | School of Engineering and Technology, University of Washington - Tacoma | L8.82

82

## AWS LAMBDA - FREEZE/THAW

- Experiment: 50 concurrent calls, 5 or 10-min calling interval
- Evaluate % cold function instances



**Performance**

October 22, 2024 | TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] | School of Engineering and Technology, University of Washington - Tacoma | L8.83

83

## 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

**Performance**

October 22, 2024 | TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] | School of Engineering and Technology, University of Washington - Tacoma | L8.84

84

### AWS LAMBDA PERFORMANCE VARIATION

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

October 22, 2024 | TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] | School of Engineering and Technology, University of Washington - Tacoma | L8.85

85

### AWS LAMBDA PERFORMANCE VARIATION - 2

- NLP use case: Less performance variance using ARM-based CPUs (less resource contention), and w/o hyperthreading

October 22, 2024 | TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] | School of Engineering and Technology, University of Washington - Tacoma | L8.86

86

## FUNCTION-AS-A-SERVICE

AWS Lambda Demo

87

### 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

88

### 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
  - 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

89

### 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

90

### OTHER CLOUD SERVICE MODELS

- IaaS
  - Storage-as-a-Service
- PaaS
  - Integration-as-a-Service
- SaaS
  - Database-as-a-Service
  - Testing-as-a-Service
  - Model-as-a-Service
- ?
  - Security-as-a-Service
  - Integration-as-a-Service

October 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
 School of Engineering and Technology, University of Washington - Tacoma    L8.91

91

### OBJECTIVES - 10/22

- Questions from 10/17
- Tutorials Questions
- Tutorial 5 - to be posted...
- From: **Cloud Computing Concepts, Technology & Architecture: Chapter 4: Cloud Computing Concepts and Models:**
  - Cloud computing delivery models
  - **Cloud deployment models**
- AWS Overview and demo
- 2<sup>nd</sup> hour:
  - Activity 2 - Horizontal Scaling in the Cloud
  - Term Project Planning

October 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
 School of Engineering and Technology, University of Washington - Tacoma    L8.92

92

### CLOUD DEPLOYMENT MODELS

- Distinguished by ownership, size, access
- Four common models
  - Public cloud
  - Community cloud
  - Hybrid cloud
  - Private cloud

October 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
 School of Engineering and Technology, University of Washington - Tacoma    L8.93

93

### PUBLIC CLOUDS

October 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
 School of Engineering and Technology, University of Washington - Tacoma    L8.94

94

### COMMUNITY CLOUD

- Specialized cloud built and shared by a particular community
- Leverage economies of scale within a community
- Research oriented clouds
- Examples:
  - Bionimbus - bioinformatics
  - Chameleon
  - CloudLab

October 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
 School of Engineering and Technology, University of Washington - Tacoma    L8.95

95

### PRIVATE CLOUD

- Compute clusters configured as IaaS cloud
- Open source software
  - Eucalyptus
  - Openstack
  - Apache Cloudstack
  - Nimbus
- Virtualization: XEN, KVM, ...

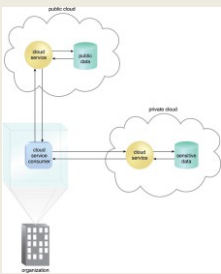
October 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
 School of Engineering and Technology, University of Washington - Tacoma    L8.96

96



## HYBRID CLOUD

- Extend private cloud typically with public or community cloud resources
- Cloud bursting:  
Scale beyond one cloud when resource requirements exceed local limitations
- Some resources can remain local for security reasons



October 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
School of Engineering and Technology, University of Washington - Tacoma    L8.97

97

## OTHER CLOUDS

- 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

October 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
School of Engineering and Technology, University of Washington - Tacoma    L8.98

98

# WE WILL RETURN AT 5:50 PM



99

## OBJECTIVES - 10/22

- Questions from 10/17
- Tutorials Questions
- Tutorial 5 - to be posted...
- From: Cloud Computing Concepts, Technology & Architecture: Chapter 4: Cloud Computing Concepts and Models:
  - Cloud computing delivery models
  - Cloud deployment models
  - **AWS Overview and demo**
- 2<sup>nd</sup> hour:
  - Activity 2 - Horizontal Scaling in the Cloud
  - Term Project Planning

October 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
School of Engineering and Technology, University of Washington - Tacoma    L8.100

100

# AWS OVERVIEW AND DEMO



101

## 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 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
School of Engineering and Technology, University of Washington - Tacoma    L8.102

102

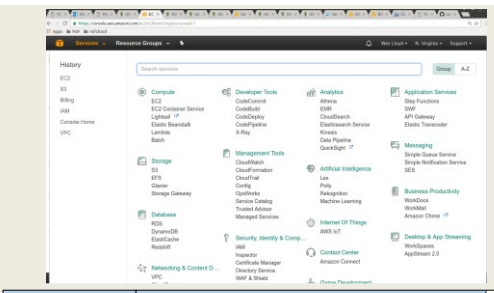
## LIST OF TOPICS

- 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 22, 2024      TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
 School of Engineering and Technology, University of Washington - Tacoma      LB.10

103

## AWS MANAGEMENT CONSOLE



October 22, 2024      TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
 School of Engineering and Technology, University of Washington - Tacoma      LB.104

104

## AWS EC2

- Elastic Compute Cloud
- Instance types: <https://ec2instances.info>
  - On demand Instance - full price
  - Reserved Instance - contract based where customer guarantees VM rental for a fixed period of time (e.g. 1 year, 3 years, etc.)  
 Deeper discounts with longer term commitments
  - Spot Instance - portion of cloud capacity reserved for low cost instances, when demand exceeds supply instances are randomly terminated with 2 minute warning
    - Users can make diverse VM requests using different types, zones, regions, etc. to minimize instance terminations
    - Developers can design for failure because often only 1 or 2 VMs in a cluster fail at any given time. They then need to be replaced.
  - Dedicated host - reserved private HW (server)
  - Instance families -  
 General, compute-optimized, memory-optimized, GPU, etc.

October 22, 2024      TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
 School of Engineering and Technology, University of Washington - Tacoma      LB.105

105

## 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 22, 2024      TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
 School of Engineering and Technology, University of Washington - Tacoma      LB.106

106

## INSTANCE STORAGE

- Also called ephemeral storage
- Persisted using images saved to S3 (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 S3
  - 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 22, 2024      TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
 School of Engineering and Technology, University of Washington - Tacoma      LB.107

107

## 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
  - **Data volumes:** 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 22, 2024      TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
 School of Engineering and Technology, University of Washington - Tacoma      LB.108

108

### 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
  - 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 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma    LB.109

109

### 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 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma    LB.110

110

### 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 per GB (for volumes sized 10GB to 1024 GB)
  - Credits – .72 minutes/day per GB

October 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma    LB.111

111

### ELASTIC FILE SYSTEM (EFS) - 2

*Information subject to revision*

- **Burstable Throughput Rates**
  - Throughput rates: baseline vs burst
  - Credit model for bursting: maximum burst per day

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%

October 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma    LB.112

112

### ELASTIC FILE SYSTEM (EFS) - 3

*Information subject to revision*

- **Throughput Models**
- **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

October 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma    LB.113

113

### ELASTIC FILE SYSTEM (EFS) - 4

*Information subject to revision*

Performance Comparison, Amazon EFS and Amazon EBS

	Amazon EFS	Amazon EBS Provisioned IOPS
Per-operation latency	Low, consistent latency.	Lowest, consistent latency.
Throughput scale	10+ GB per second.	Up to 2 GB per second.

Storage Characteristics Comparison, Amazon EFS and Amazon EBS

	Amazon EFS	Amazon EBS Provisioned IOPS
Availability and durability	Data is stored redundantly across multiple AZs.	Data is stored redundantly in a single AZ.
Access	Up to thousands of Amazon EC2 instances, from multiple AZs, can connect concurrently to a file system.	A single Amazon EC2 instance in a single AZ can connect to a file system.
Use cases	Big data and analytics, media processing workflows, content management, web serving, and home directories.	Boot volumes, transactional and NoSQL databases, data warehousing, and ETL.

October 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma    LB.114

114

### AMAZON MACHINE IMAGES

- AMIs
- Unique for the operating system (root device image)
- Two types
  - Instance store
  - Elastic block store (EBS)
- Deleting requires multiple steps
  - Deregister AMI
  - Delete associated data - (files in S3)
- Forgetting both steps leads to costly "orphaned" data
  - No way to instantiate a VM from deregistered AMIs
  - Data still in S3 resulting in charges

October 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
 School of Engineering and Technology, University of Washington - Tacoma    L8.115

115

### 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 **AKI** files on AWS – **Amazon kernel Image(s)**
    - Look for common identifiers

October 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
 School of Engineering and Technology, University of Washington - Tacoma    L8.116

116

### 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 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
 School of Engineering and Technology, University of Washington - Tacoma    L8.117

117

### EC2 VIRTUALIZATION - NITRO

- 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 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
 School of Engineering and Technology, University of Washington - Tacoma    L8.118

118

### EVOLUTION OF AWS VIRTUALIZATION

From: <http://www.brendanregg.com/blog/2017-11-29/aws-ec2-virtualization-2017.html>

AWS EC2 Virtualization Types

#	Tech	Type	With	Importance → Least				
				Local Storage I/O	Network I/O	Interrupts	Memory Access	Direct Access
1	VM	Fully Emulated		V/S	V/S	V/S	V/S	V/S
2	VM	Xen PV 3.0	PV drivers	P	P	P	V/S	V/S
3	VM	Xen HVM 3.0	PV drivers	V/S	P	P	P	V/S
4	VM	Xen HVM 4.0.1	PVHVM drivers	V/S	P	P	P	V/S
5	VM	Xen AWS 2013	PVHVM + SR-IOV(net)	V/S	V/S	P	P	V/S
6	VM	Xen AWS 2017	PVHVM + SR-IOV(net, stor)	V/S	V/S	V/S	P	V/S
7	VM	AWS Nitro 2017		H	H	H	H	V/S
8	HW	AWS Bare Metal 2017		H	H	H	H	H

VM: Virtual Machine; HW: Hardware.  
 V/S: VM in software; H: VM in hardware; P: Paravirt; Net all combinations shown.  
 SR-IOV(net): network driver; SR-IOV(stor): storage driver.

October 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
 School of Engineering and Technology, University of Washington - Tacoma    L8.119

119

### INSTANCE ACTIONS

- Stop
  - Costs of "pausing" an instance
- Terminate
- Reboot
- Image management
  - Creating an image
    - EBS (snapshot)
  - Bundle image
  - Instance-store

October 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024]  
 School of Engineering and Technology, University of Washington - Tacoma    L8.120

120

### EC2 INSTANCE: NETWORK ACCESS

- Public IP address
- Elastic IPs
  - Costs: in-use FREE, not in-use ~12 €/day
  - Not in-use (e.g. "paused" EBS-backed instances)
- Security groups
  - E.g. firewall
- Identity access management (IAM)
  - AWS accounts, groups
- VPC / Subnet / Internet Gateway / Router
- NAT-Gateway

October 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma    L8.121

121

### SIMPLE VPC

- Recommended when using Amazon EC2

Destination	Target
10.0.0/16	local
0.0.0.0/0	igw-af

October 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma    L8.122

122

### VPC SPANNING AVAILABILITY ZONES

Destination	Target
10.0.0/16	local

October 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma    L8.124

123

### 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 ; echo
```
- ec2-get-info command
- Python API that provides easy/formatted access to metadata

October 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma    L8.124

124

### 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 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma    L8.125

125

### AWS CLI

- Launch Ubuntu 16.04 VM
  - Instances | Launch Instance
- Install the general AWS CLI
  - sudo apt install awscli
- Create config file
 

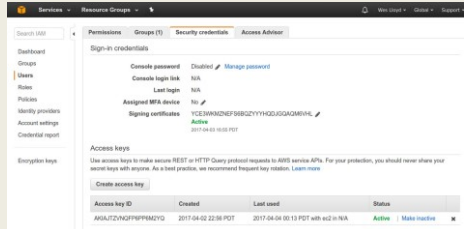
```
[default]
aws_access_key_id = <access key id>
aws_secret_access_key = <secret access key>
region = us-east-1
```

October 22, 2024    TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma    L8.126

126

## AWS CLI - 2

- **Creating access keys:** IAM | Users | Security Credentials | Access Keys | Create Access Keys



October 22, 2024 TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma L8.127

127

## 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 vmscaleruw`

October 22, 2024 TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma L8.128

128

## 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/javadoc/index.html>
- Some AWS services have separate CLI installable by package

October 22, 2024 TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma L8.129

129

## 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:
  - <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 22, 2024 TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma L8.130

130

## 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](http://docs.aws.amazon.com/AWSEC2/latest/UserGuide/set-up-ami-tools.html?icmpid=docs_iam_console#ami-tools-create-certificate)

October 22, 2024 TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma L8.131

131

## 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 22, 2024 TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma L8.132

132

### 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 -l
/etc/ec2/amitools/cert-ec2.pem
cd /tmp
ec2-upload-bundle -b tc562 -n $image.manifest.xml -a ${AWS_ACCESS_KEY} -s
${AWS_SECRET_KEY} --url http://s3.amazonaws.com --location US
ec2-register tc562/$image.manifest.xml --region us-east-1 --kernel ak1-
88aa75e1
```

October 22, 2024 TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma L8.133

133

### 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 22, 2024 TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma L8.134

134

### OBJECTIVES - 10/22

- Questions from 10/17
- Tutorials Questions
- Tutorial 5 - to be posted...
- From: Cloud Computing Concepts, Technology & Architecture: Chapter 4: Cloud Computing Concepts and Models:
  - Cloud computing delivery models
  - Cloud deployment models
- AWS Overview and demo
- 2<sup>nd</sup> hour:
  - Activity 2 - Horizontal Scaling in the Cloud
  - Term Project Planning

October 22, 2024 TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma L8.135

135


### OBJECTIVES - 10/22

- Questions from 10/17
- Tutorials Questions
- Tutorial 5 - to be posted...
- From: Cloud Computing Concepts, Technology & Architecture: Chapter 4: Cloud Computing Concepts and Models:
  - Cloud computing delivery models
  - Cloud deployment models
- AWS Overview and demo
- 2<sup>nd</sup> hour:
  - Activity 2 - Horizontal Scaling in the Cloud
  - Term Project Planning

October 22, 2024 TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma L8.136

136


## TCSS 462/562 TERM PROJECT



October 22, 2024 TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma L8.137

137

## QUESTIONS



October 22, 2024 TCSS462/562: Software Engineering for Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma L8.138

138