

TCSS 562: SOFTWARE ENGINEERING FOR CLOUD COMPUTING

Cloud Computing Concepts and Models

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
- **Fridays**
 - **<THIS WEEK>**
12:00 pm to 1:00 pm – ONLINE via Zoom
- **Or email for appointment**

> Office Hours set based on Student Demographics survey feedback

| | | |
|------------------|---|------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.2 |
|------------------|---|------|

2

OBJECTIVES – 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 | L7.3 |
|------------------|---|------|

3

ONLINE DAILY FEEDBACK SURVEY

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

Announcements

Assignments

Discussions

Zoom

Grades

People

Pages

Files

Quizzes

Collaborations

UW Libraries

UW Resources

▼ Upcoming Assignments

- 📄 **Class Activity 1 – Implicit vs. Explicit Parallelism**
Available until Oct 11 at 11:59pm | Due Oct 7 at 7:50pm | -/10 pts
- 📄 **Tutorial 1 - Linux**
Available until Oct 19 at 11:59pm | Due Oct 15 at 11:59pm | -/20 pts

▼ Past Assignments

- 📄 **TCSS 562 - Online Daily Feedback Survey - 10/5**
Available until Dec 18 at 11:59pm | Due Oct 6 at 8:59pm | -/1 pts
- 📄 **TCSS 562 - Online Daily Feedback Survey - 9/30**
Available until Dec 18 at 11:59pm | Due Oct 4 at 8:59pm | -/1 pts

| | | |
|------------------|---|------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.4 |
|------------------|---|------|

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

October 17, 2024 TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024]
School of Engineering and Technology, University of Washington - Tacoma L7.5

5

MATERIAL / PACE

- Please classify your perspective on material covered in today's class (44 respondents):
 - 1-mostly review, 5-equal new/review, 10-mostly new
 - **Average - 6.50 (↑ - previous 6.28)**
- Please rate the pace of today's class:
 - 1-slow, 5-just right, 10-fast
 - **Average - 5.59 (↑ - previous 5.51)**
- **Response rates:**
 - TCSS 462: 29/42 - 69.0%
 - TCSS 562: 15/20 - 75.0%

October 17, 2024 TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024]
School of Engineering and Technology, University of Washington - Tacoma L7.6

6

FEEDBACK FROM 10/15

- **How many layers of virtualization can a computer run?**
 - Nested virtualization – running a VM from within a VM
 - Too many layers will lead to high overhead

- **Are Lambda functions serverless ?**
 - Programmer does not manage servers
 - Programmer deploys function code, specifies memory setting
 - AWS Lambda is a function-as-a-service platform that is serverless
 - Servers are not specified or configured in using AWS Lambda
 - Servers are abstracted from the programmer's view

| | | |
|------------------|---|------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.7 |
|------------------|---|------|

7

FEEDBACK - 2

- **Still confused about exactly what AWS Lambda is used for**

AWS Compute Platforms:

- **Infrastructure-as-a-Service (IaaS):**
Elastic Compute Cloud (EC2)
 - Virtual machines on demand

- **Container-as-a-Service (CaaS):**
Elastic Container Service (ECS), Elastic Kubernetes Service (EKS)
 - Creates a container cluster consisting of VMs w/ Docker

- **Serverless containers: AWS Fargate**
 - Containers w/o managing a container cluster, or any VMs
 - Run containers on demand of short or long time
 - Only pay for resources used to run 1 container, which is less than a VM

| | | |
|------------------|---|------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.8 |
|------------------|---|------|

8

FEEDBACK - 3

- **Function-as-a-Service (FaaS):**
AWS Lambda
 - Hosts webservices/microservices
 - Specified as “functions”
 - No VMs or servers to manage
 - Developer specifies a function memory setting

| | | |
|------------------|---|------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.9 |
|------------------|---|------|

9

FEEDBACK - 4

- **What is the difference between asynchronous function calls (with AWS Lambda) and asynchronous programming in web/mobile app development (JavaScript with Promises)?**
- AWS Lambda clients and functions are prog. language agnostic
- Functions are invoked via HTTP REST interfaces (language agnostic)
- Promises provide a client-side calling construct (in Javascript) for performing asynchronous (non-blocking) function calls – the idea is usually to perform other work while waiting for a web response (promise fulfillment)

| | | |
|------------------|---|-------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.10 |
|------------------|---|-------|

10

AWS LAMBDA: VCPU SCALING W/ MEMORY

| Function Memory | CPU time share |
|-----------------|----------------------------|
| 1769 MB | 100 % = 1 vCPU |
| 2389 MB | 150 % = 1.5 vCPUs |
| 3008 MB | 200 % = 2 vCPUs |
| 4158 MB | 250 % = 2.5 vCPUs |
| 5307 MB | 300 % = 3 vCPUs |
| 6192 MB | 350 % = 3.5 vCPUs |
| 7076 MB | 400 % = 4 vCPUs (1 HT) |
| 7960 MB | 450 % = 4.5 vCPUs (1.5 HT) |
| 8845 MB | 500 % = 5 vCPUs (2 HT) |
| 9543 MB | 550 % = 5.5 vCPUs (2.5 HT) |
| 10240 MB | 600 % = 6 vCPUs (3 HT) |

Based on:
<https://stackoverflow.com/questions/66522916/aws-lambda-memory-vs-cpu-configuration>

| | | |
|------------------|---|-------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.11 |
|------------------|---|-------|

11

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
 - Include account number and email – used for tracking codes
- Codes can also be obtained in person (or zoom), in the class, during the breaks, after class, during office hours, by appt
 - All credit requests as of Oct 16 have been distributed
- To track credit code distribution, codes not shared via discord
- 40 credit requests fulfilled thus far

| | | |
|------------------|--|-------|
| October 10, 2023 | TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L4.12 |
|------------------|--|-------|

12

OBJECTIVES - 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 | L7.13 |
|------------------|---|-------|

13

TUTORIAL 0

- Getting Started with AWS
- https://faculty.washington.edu/wlloyd/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 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7/14 |
|------------------|---|-------|

14

TUTORIAL 2 – OCT 19

- **Introduction to Bash Scripting**
- https://faculty.washington.edu/wlloyd/courses/tcss562/tutorials/TCSS462_562_f2024_tutorial_2.pdf
- Review tutorial sections:
- Create a BASH webservice client
 1. What is a BASH script?
 2. Variables
 3. Input
 4. Arithmetic
 5. If Statements
 6. Loops
 7. Functions
 8. User Interface
- Call service to obtain IP address & lat/long of computer
- Call weatherbit.io API to obtain weather forecast for lat/long

October 11, 2022

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

L4.15

15

TUTORIAL 3 – 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
- 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 17, 2024

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

L7/16

16

OBJECTIVES – 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

L7.17

17

TUTORIAL 4 – TO BE POSTED

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

October 17, 2024

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

L7.18

18


OBJECTIVES – 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 | L7.19 |
|------------------|---|-------|

19

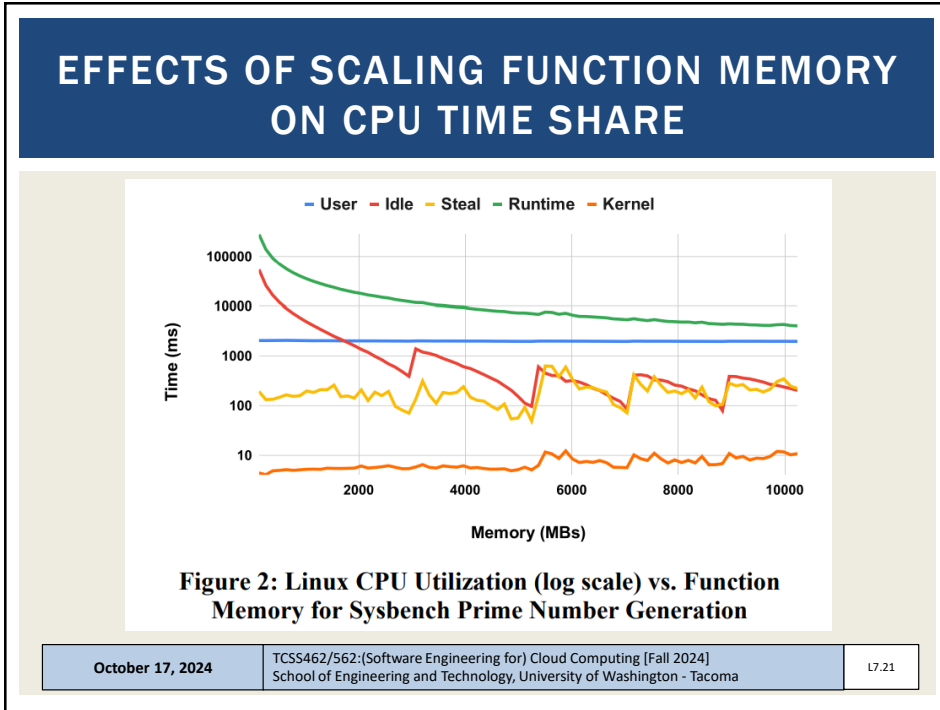
CPUSTEAL



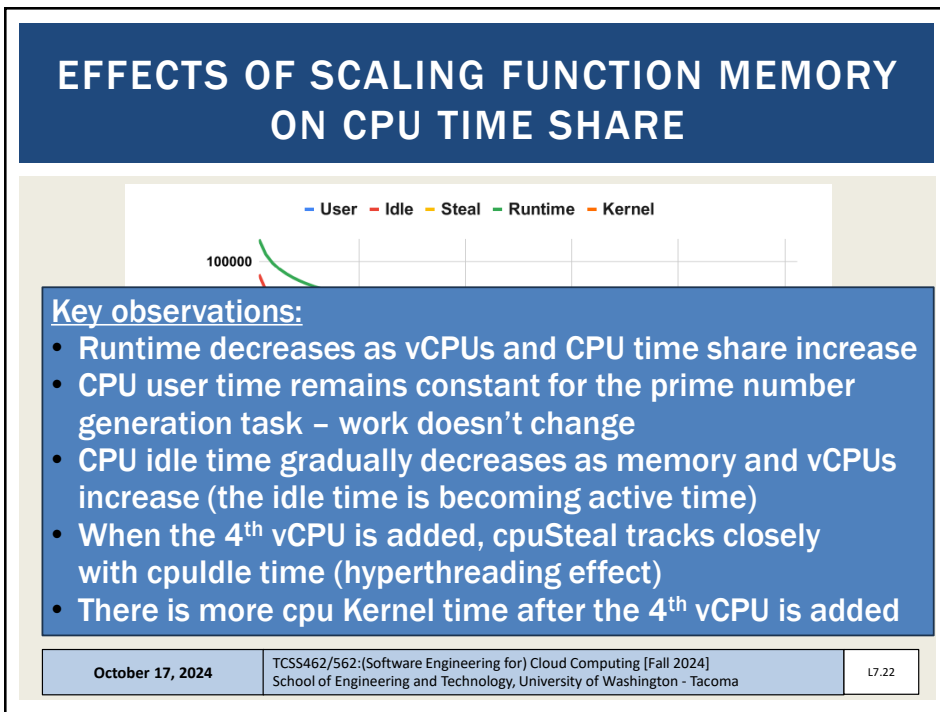
- *CpuSteal*: Metric that measures when a CPU core is ready to execute but the physical CPU core is busy and unavailable
- Symptom of over provisioning physical servers in the cloud
- Factors which cause *CpuSteal*: (x86 hyperthreading)
 1. Physical CPU is shared by too many busy VMs
 2. Hypervisor kernel 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

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

20



21



22

FUNCTION INSTANCE LIFE CYCLES

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

October 17, 2024

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

L7.23

23

WARM VS COLD FUNCTION INSTANCES

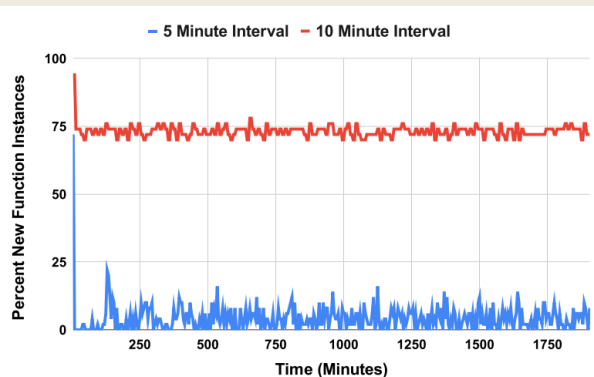


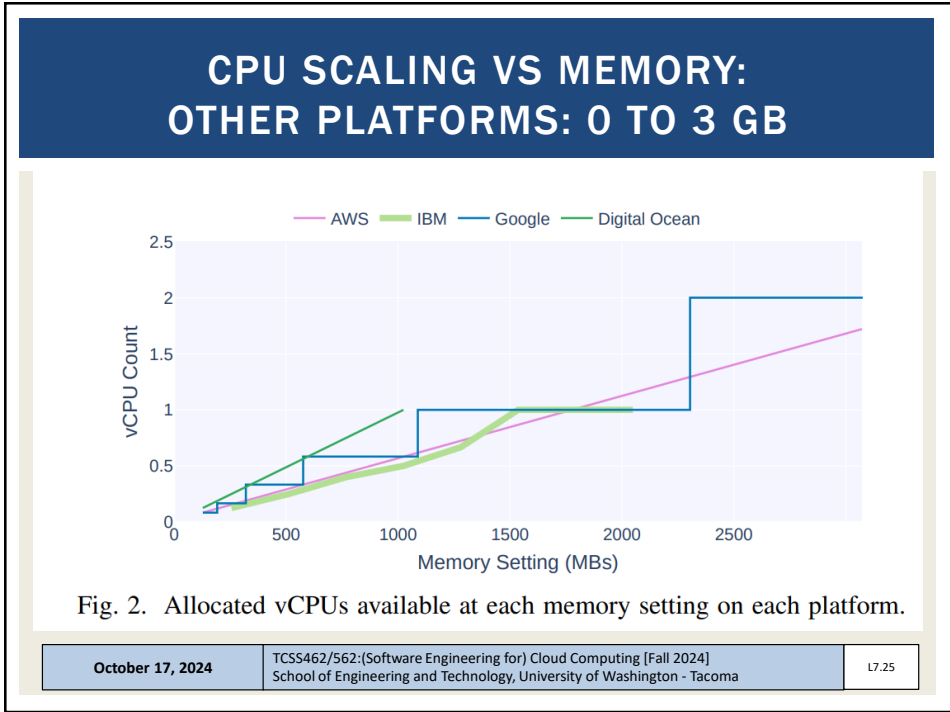
Figure 3: AWS Lambda Function Instance Replacement vs. Function Call Interval over 24-hours

October 17, 2024

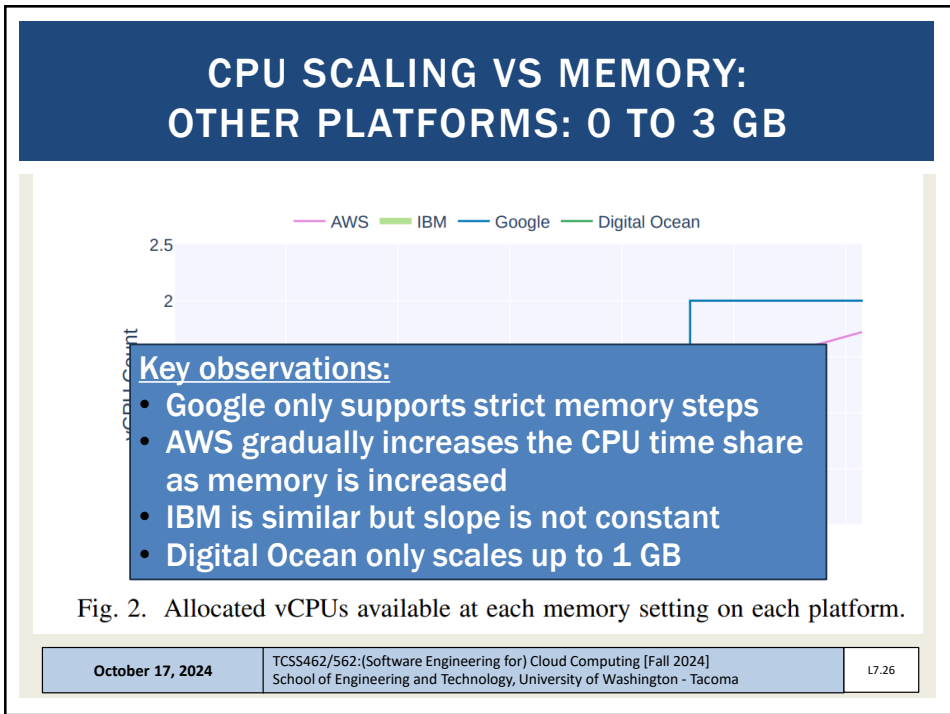
TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024]
School of Engineering and Technology, University of Washington - Tacoma

L9.24

24



25



26

ELASTIC FILE SYSTEM (AWS EFS)

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

October 17, 2024

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

L7.27

27

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 17, 2024

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

L7.28

28

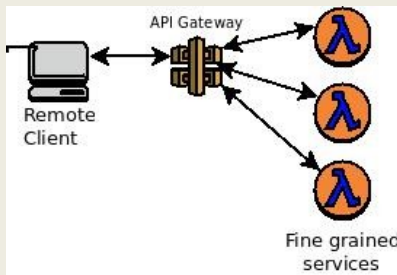
SERVERLESS FILE STORAGE COMPARISON PROJECT

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

| | | |
|------------------|---|-------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.29 |
|------------------|---|-------|

29

SERVICE COMPOSITION



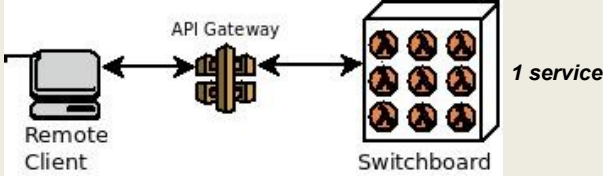
| | | | |
|-------------|---|---|--|
| A | B | C | 3 services Full Service Isolation |
| A | B | C | 2 services |
| A | B | C | 2 services |
| A B C | | | 1 service Full Service Aggregation |

Other possible compositions: group by library, functional cohesion, etc.

| | | |
|------------------|---|-------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.30 |
|------------------|---|-------|

30

SWITCH-BOARD ARCHITECTURE



Single deployment package with consolidated codebase (Java: one JAR file)

Entry method contains “switchboard” logic
Case statement that route calls to proper service

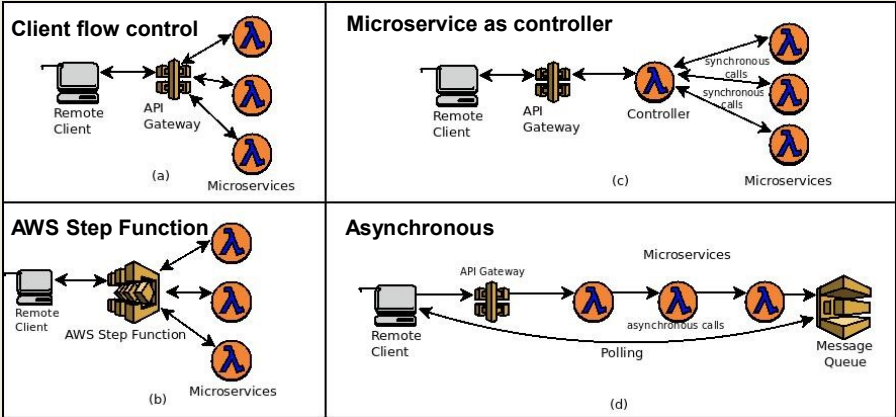
Routing is based on data payload
Check if specific parameters exist, route call accordingly

Goal: reduce # of COLD starts to improve performance

| | | |
|------------------|---|-------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.31 |
|------------------|---|-------|

31

APPLICATION FLOW CONTROL - 3



(a) Client flow control: Remote Client → API Gateway → Microservices

(b) AWS Step Function: Remote Client → AWS Step Function → Microservices

(c) Microservice as controller: Remote Client → API Gateway → Controller → Microservices (synchronous calls)

(d) Asynchronous: Remote Client → API Gateway → Microservices (asynchronous calls) → Message Queue → Polling → Remote Client

| | | |
|------------------|---|-------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.32 |
|------------------|---|-------|

32

PROGRAMMING LANGUAGE COMPARISON

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

| | | |
|------------------|---|-------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.33 |
|------------------|---|-------|

33

FAAS PLATFORMS

- Many commercial and open source FaaS platforms exist
- TCSS562 projects can choose to compare performance and cost implications of alternate platforms.

- Supported by SAAF:
 - AWS Lambda
 - Google Cloud Functions
 - Azure Functions
 - IBM Cloud Functions
 - Apache OpenWhisk (*open source, deploy your own FaaS*)
 - Open FaaS (*open source, deploy your own FaaS*)

| | | |
|------------------|---|-------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.34 |
|------------------|---|-------|

34

DATA PROVISIONING

- Consider performance and cost implications of the data-tier design for the serverless application
- Use different tools as the relational datastore to support service #2 (LOAD) and service #3 (EXTRACT)

- **SQL / Relational:**
- Amazon Aurora (serverless cloud DB), Amazon RDS (cloud DB), DB on a VM (MySQL), DB inside Lambda function (SQLite, Derby)

- **NO SQL / Key/Value Store:**
- Dynamo DB, MongoDB, S3

| | | |
|------------------|---|-------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.35 |
|------------------|---|-------|

35

PERFORMANCE VARIABILITY

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

| | | |
|------------------|---|-------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.36 |
|------------------|---|-------|

36

CPU STEAL CASE STUDY

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

October 17, 2024

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

L7.37

37

CPU ARCHITECTURE & PERFORMANCE

- X86_64 Intel
 - Intel Xeon Platinum 8259 CL @ 2.5 GHz
- ARM64 Graviton2

October 17, 2024

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

L7.38

38


OBJECTIVES - 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 | L7.39 |
|------------------|---|-------|

39

CLOUD COMPUTING: CONCEPTS AND MODELS



| | | |
|------------------|---|-------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.40 |
|------------------|---|-------|

40

ROLES

- **Cloud provider**
 - Organization that provides cloud-based resources
 - Responsible for fulfilling SLAs for cloud services
 - Some cloud providers “resell” IT resources from other cloud providers
 - Example: Heroku sells PaaS services running atop of Amazon EC2
- **Cloud consumers**
 - Cloud users that consume cloud services
- **Cloud service owner**
 - Both cloud providers and cloud consumers can own cloud services
 - A cloud service owner may use a cloud provider to provide a cloud service (e.g. Heroku)

| | | |
|------------------|---|-------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.41 |
|------------------|---|-------|

41

ROLES - 2

- **Cloud resource administrator**
 - Administrators provide and maintain cloud services
 - Both cloud providers and cloud consumers have administrators
- **Cloud auditor**
 - Third-party which conducts independent assessments of cloud environments to ensure security, privacy, and performance.
 - Provides unbiased assessments
- **Cloud brokers**
 - An intermediary between cloud consumers and cloud providers
 - Provides service aggregation
- **Cloud carriers**
 - Network and telecommunication providers which provide network connectivity between cloud consumers and providers

| | | |
|------------------|---|-------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.42 |
|------------------|---|-------|

42

ORGANIZATION BOUNDARY

October 17, 2024 TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024]
School of Engineering and Technology, University of Washington - Tacoma L7.43

43

TRUST BOUNDARY

October 17, 2024 TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024]
School of Engineering and Technology, University of Washington - Tacoma L7.44

44

OBJECTIVES - 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 | L7.45 |
|------------------|---|-------|

45

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 | L7.46 |
|------------------|---|-------|

46

ON-DEMAND USAGE

- The freedom to self-provision IT resources
- Generally, with automated support
- Automated support requires no human involvement
- Automation through software services interface

Internet Data Centre National Informatics Centre
Data Centre and Web Services Division
Virtual Machine Request Form

You are requested to please go through the IDC security policies before filling up this form.

1. Name of the NIC Group / Division

2. Name of the Project / Service
(Please describe the Architecture on a separate sheet)

3. Category: Web | Database | Other

Others if any specify:


4. Virtual Machine Specifications

- Name of the Virtual Machine
- Hard Disk Required (GB) (Please specify the size)
- CPU Required
- RAM Required

5. Software Environment

- Operating System (with version)
- Software & Tools
- Software Licenses Details (including IP)

Instructions provide user ID for self-provisioning the application.



| | | |
|------------------|---|-------|
| October 17, 2024 | TCCS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.47 |
|------------------|---|-------|

47

UBIQUITOUS ACCESS

- Cloud services are widely accessible
- Public cloud: internet accessible
- Private cloud: throughout segments of a company's intranet
- 24/7 availability

| | | |
|------------------|---|-------|
| October 17, 2024 | TCCS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.48 |
|------------------|---|-------|

48

MULTITENANCY

- Cloud providers pool resources together to share them with many users
- Serve multiple cloud service consumers
- IT resources can be dynamically assigned, reassigned based on demand
- Multitenancy can lead to performance variation

| | | |
|------------------|---|-------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.49 |
|------------------|---|-------|

49

SINGLE TENANT MODEL

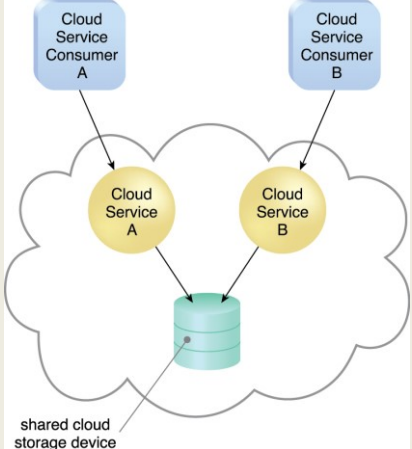
The diagram illustrates the Single Tenant Model. It features a central cloud shape containing two distinct service paths. On the left, 'Cloud Service Consumer A' (blue box) connects to 'Cloud Service A' (yellow circle), which in turn connects to 'Cloud Storage Device A' (teal cylinder). On the right, 'Cloud Service Consumer B' (blue box) connects to 'Cloud Service B' (yellow circle), which connects to 'Cloud Storage Device B' (teal cylinder). A large blue double-headed arrow labeled '> Isolation <' spans the gap between the two service paths, indicating that each consumer's resources are isolated from the other's.

| | | |
|------------------|---|-------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.50 |
|------------------|---|-------|

50

MULTITENANT MODEL

- Resource is “multiplexed” and share amongst multiple users
- Goal is to increase utilization
- Often server resources are underutilized
- There are many “sunk costs” whether usage is 0% or 100%
- Cloud computing tries to maximize “sunk cost” investments through **multi-tenancy**




The diagram illustrates a multitenant model. At the top, two boxes labeled 'Cloud Service Consumer A' and 'Cloud Service Consumer B' have arrows pointing to two yellow circles labeled 'Cloud Service A' and 'Cloud Service B' respectively. These services are contained within a cloud shape. Both 'Cloud Service A' and 'Cloud Service B' have arrows pointing to a central green cylinder representing a 'shared cloud storage device'.

| | | |
|------------------|---|-------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.51 |
|------------------|---|-------|

51

MULTITENANT DATABASE

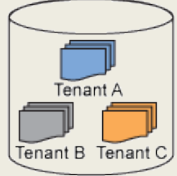
Isolated



Separate database

E1

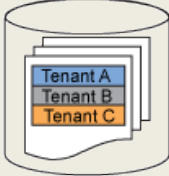
Semi-shared



**Shared database
Separate schema**

E2

Shared



**Shared database
Shared schema**

E3

- Many users on a single database instance
- ***What issues may occur when sharing a single database instance?***

| | | |
|------------------|---|-------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.52 |
|------------------|---|-------|

52

MULTITENANCY OF RESOURCES

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

Traditional On Premise

Single Tenant (Hosted)

Multi-Tenant

Virtual Appliance

October 17, 2024
TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024]
School of Engineering and Technology, University of Washington - Tacoma
L7.53

53

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.

| VM Tenants | sysbench (CPU) | y-cruncher (CPU) | pgbench (CPU + I/O) | iperf (network I/O) |
|------------|----------------|------------------|---------------------|---------------------|
| 0 | 100% | 100% | 100% | 100% |
| 10 | 95% | 90% | 95% | 40% |
| 20 | 90% | 80% | 90% | 25% |
| 30 | 85% | 70% | 85% | 20% |
| 40 | 80% | 60% | 80% | 15% |
| 48 | 75% | 55% | 75% | 10% |

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
L7.54

54

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) →

| Instance Family | iperf (network) | pgbench (CPU + I/O) | sysbench (CPU) | y-cruncher (CPU) | Total Variance (%) |
|-----------------|-----------------|---------------------|----------------|------------------|--------------------|
| c3 | 19.2% | 0.3% | 24.5% | 0.3% | 44.3% |
| c4 | 42.1% | 0.2% | 0.1% | 0.1% | 42.5% |
| z1d | 84.6% | 11.2% | 0.2% | 3.3% | 99.3% |
| m5d (t) | 94.6% | 33.0% | 20.8% | 48.0% | 196.4% |

October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] | L7.55
 School of Engineering and Technology, University of Washington - Tacoma

55

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] | L7.56
 School of Engineering and Technology, University of Washington - Tacoma

56

PREDICTABLE DEMAND

- AWS EC2 Scaling Example:

Auto-Scaling Example: Netflix

From: Kejarawal, A., 2013, March. 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 | L7.57 |
|------------------|---|-------|

57

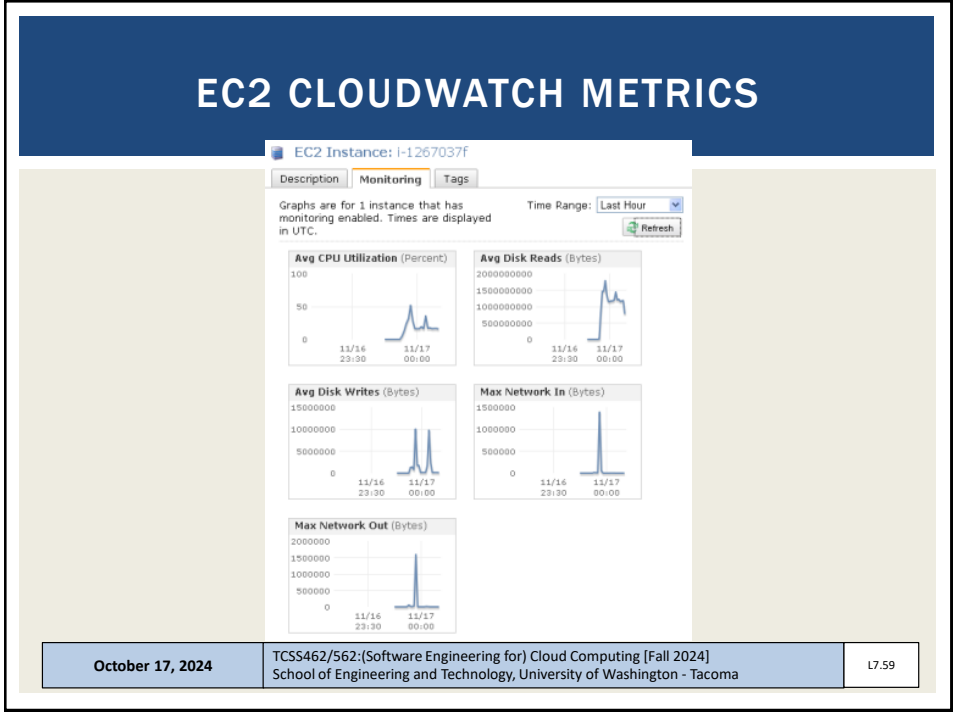
MEASURED USAGE

- 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)

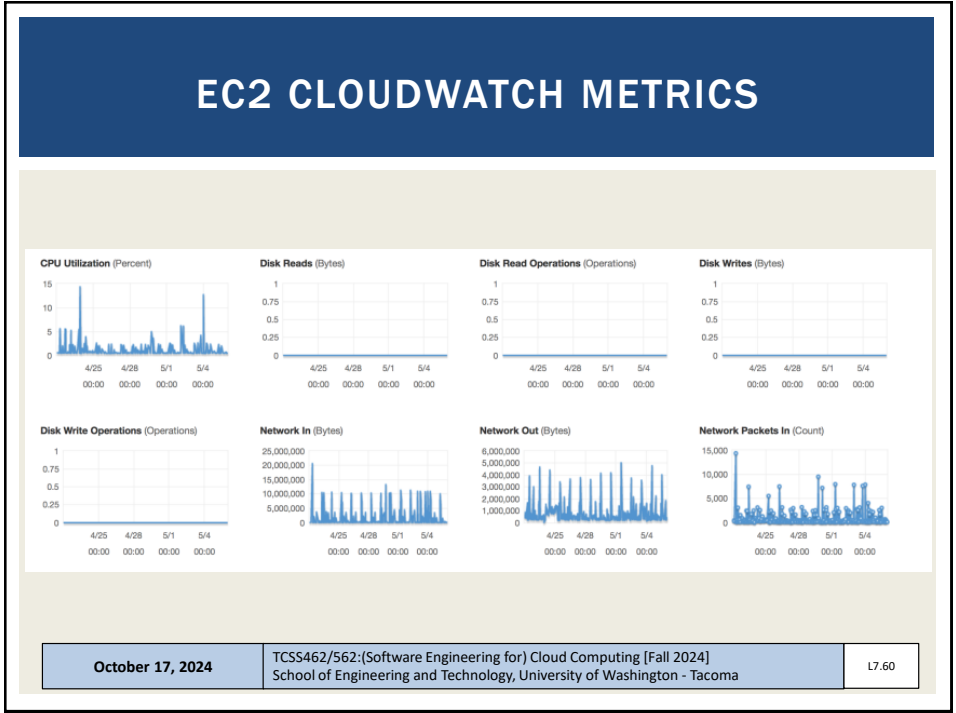
- Not all measurements are for billing
- Some measurements can support auto-scaling
- For example CPU utilization

| | | |
|------------------|---|-------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.58 |
|------------------|---|-------|

58



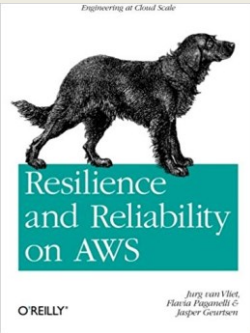
59



60

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 | L7.61 |
|------------------|---|-------|

61

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

- A When the application is primarily I/O bound, a CPU threshold may never be met, or be met too late to scale up.
- B When the current resource utilization does not reflect future system demand.
- C 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.
- D When an application will soon complete a parallel phase, before executing a largely sequential phase
- E All of the above

| | | |
|------------------|---|--------|
| October 24, 2016 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L10.62 |
|------------------|---|--------|

62

When poll is active, respond at pollev.com/wesleylloyd641
Text **WESLEYLLOYD641** to **22333** once to join

W 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 pollev.com/app

63

OBJECTIVES – 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 | L7.64 |
|------------------|---|-------|

64

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 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.65 |
|------------------|---|-------|

65

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 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.66 |
|------------------|---|-------|

66

CLOUD COMPUTING DELIVERY MODELS

- Infrastructure-as-a-Service (IaaS) delivery model
- Virtual Machines
- Users

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 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.67 |
|------------------|---|-------|

67


CLOUD COMPUTING DELIVERY MODELS

- Infrastructure-as-a-Service (IaaS) delivery model
- Virtual Machines
- Users

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 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.68 |
|------------------|---|-------|

68

CLOUD DELIVERY MODELS

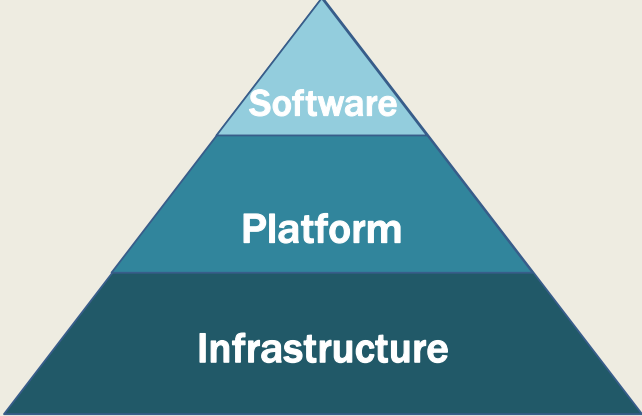
- What is the appropriate level of **abstraction**?
- How should applications be deployed?
 - IaaS, PaaS, SaaS, DBaaS, FaaS
- How do we ensure Quality-of-Service?
 - Performance, Availability, Responsiveness, Fault Tolerance
- How is **scalability** provided?
- As users, how do we minimize hosting costs?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



| | | |
|------------------|---|-------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.69 |
|------------------|---|-------|

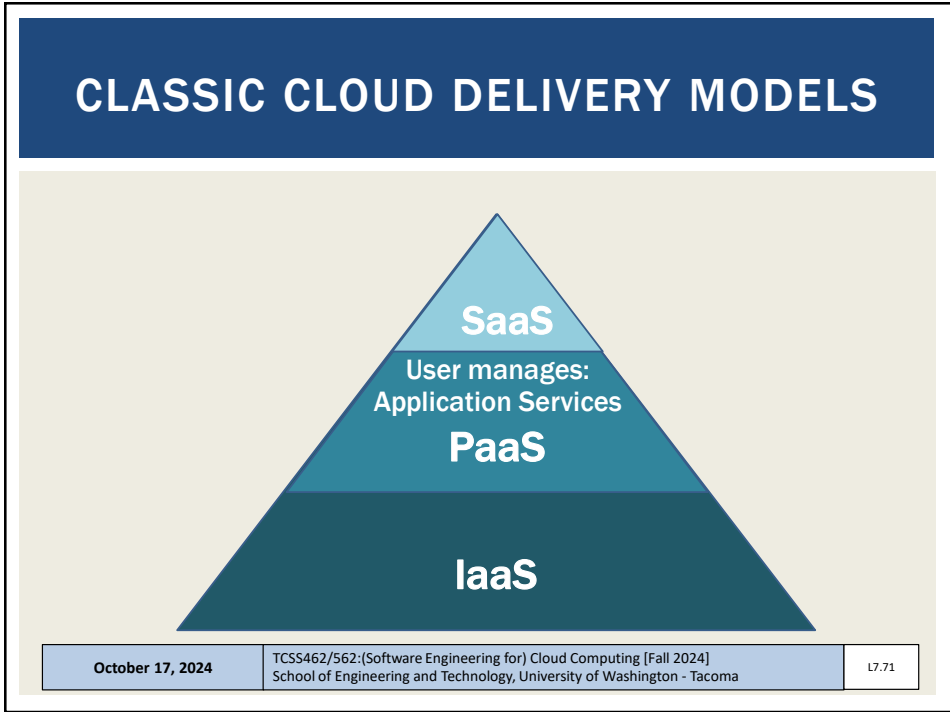
69

CLASSIC CLOUD DELIVERY MODELS

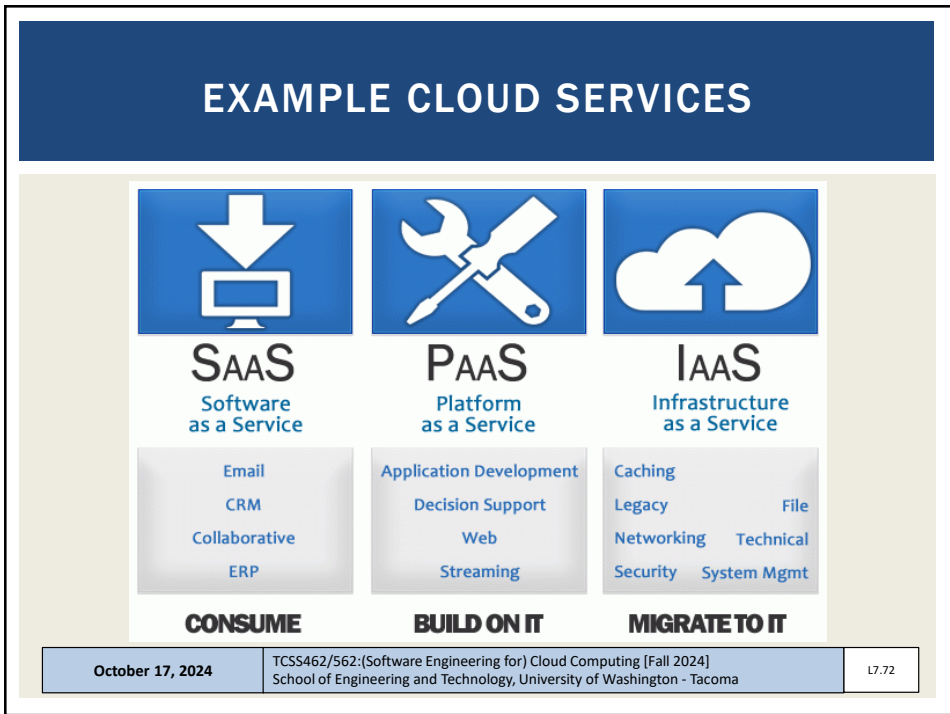


| | | |
|------------------|---|-------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.70 |
|------------------|---|-------|

70



71



72

END USER APPLICATIONS

Many different "cloud" providers (especially SaaS)

Many cloud providers are also cloud consumers

| | | |
|------------------|---|-------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.73 |
|------------------|---|-------|

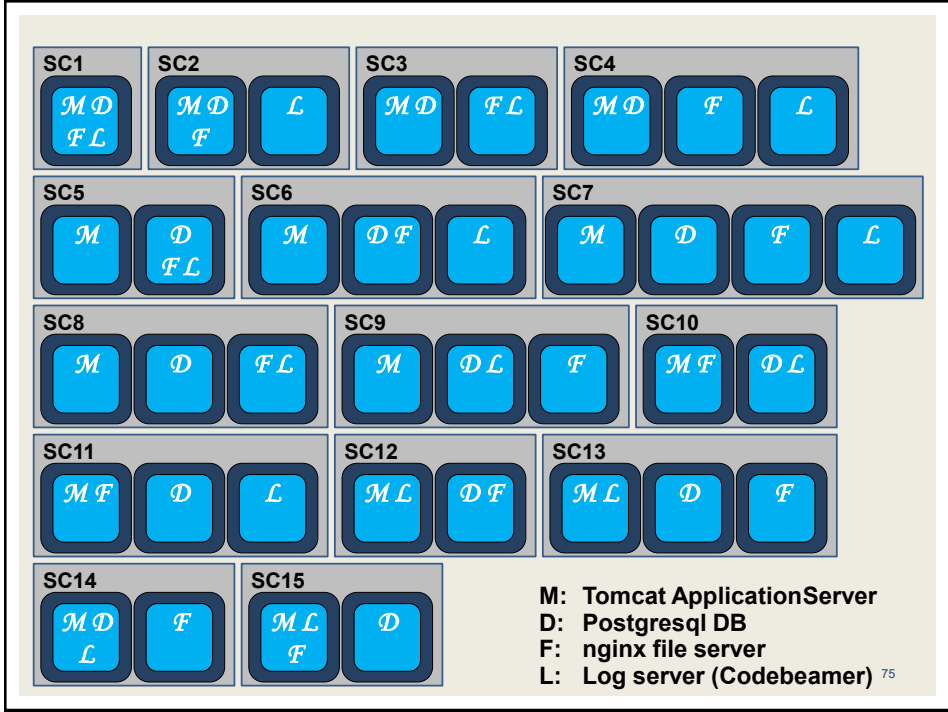
73

INFRASTRUCTURE-AS-A-SERVICE

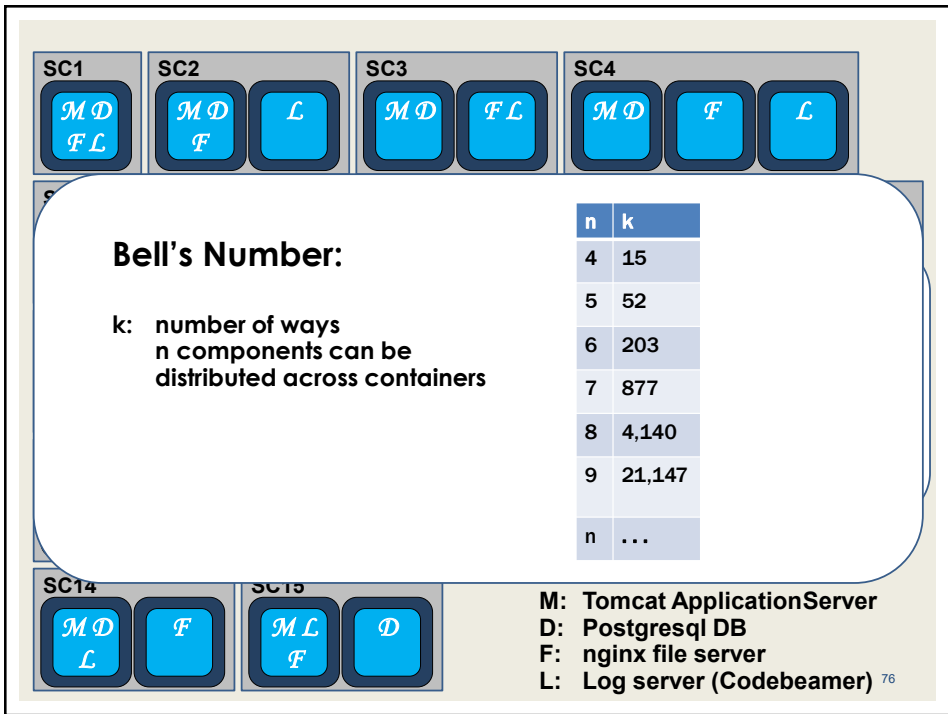
- Compute resources, on demand, as-a-service
 - Generally raw "IT" resources
 - Hardware, network, containers, operating systems
- Typically provided through virtualization
- Generally, not-preconfigured
- Administrative burden is owned by cloud consumer
- Best when high-level control over environment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

| | | |
|------------------|---|-------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.74 |
|------------------|---|-------|

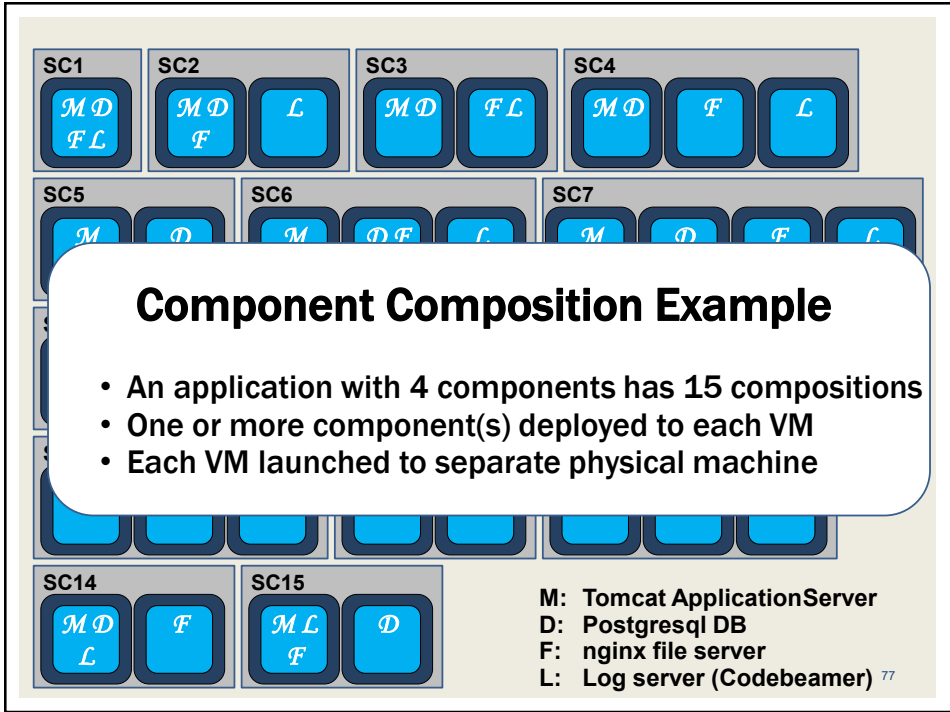
74



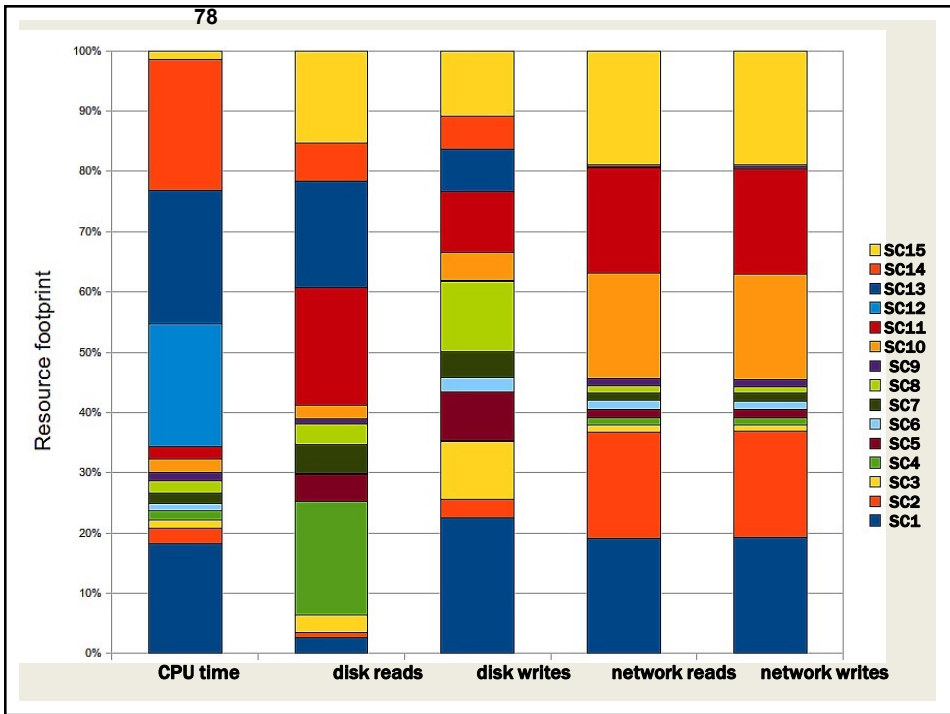
75



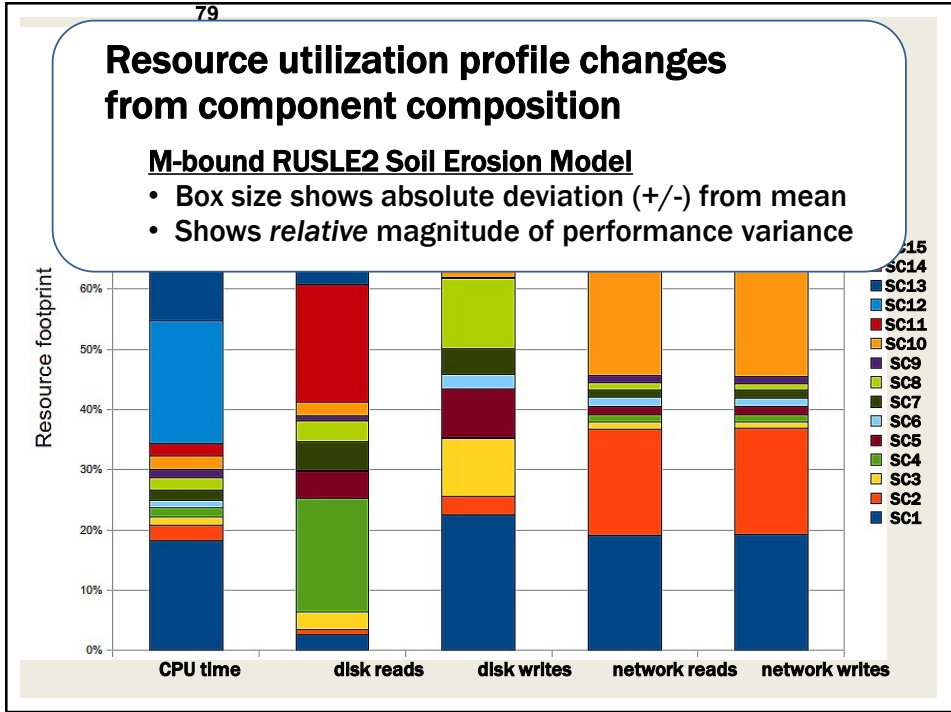
76



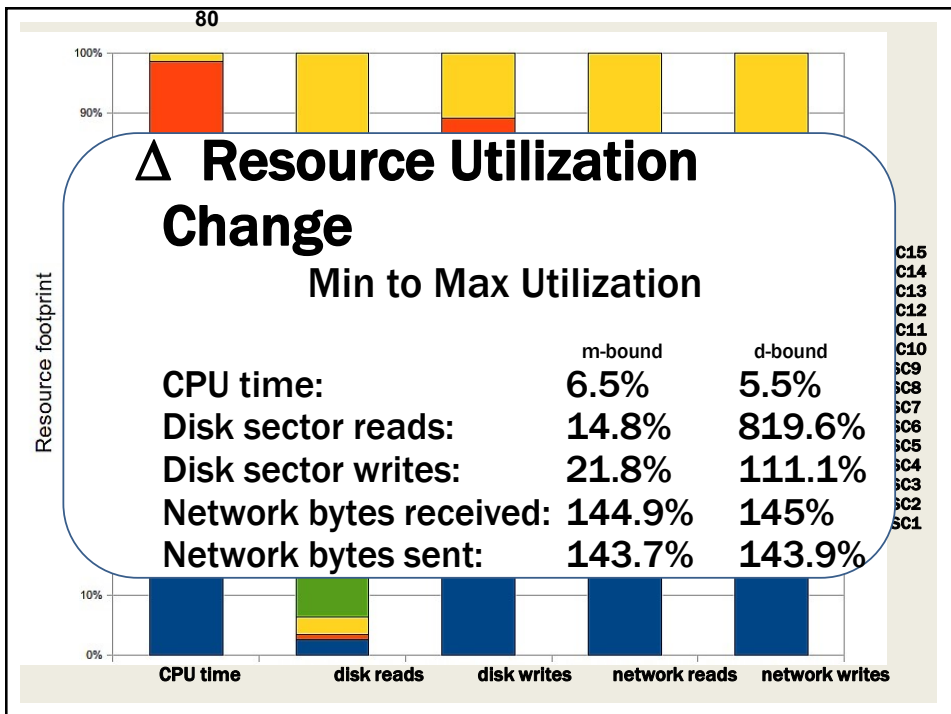
77



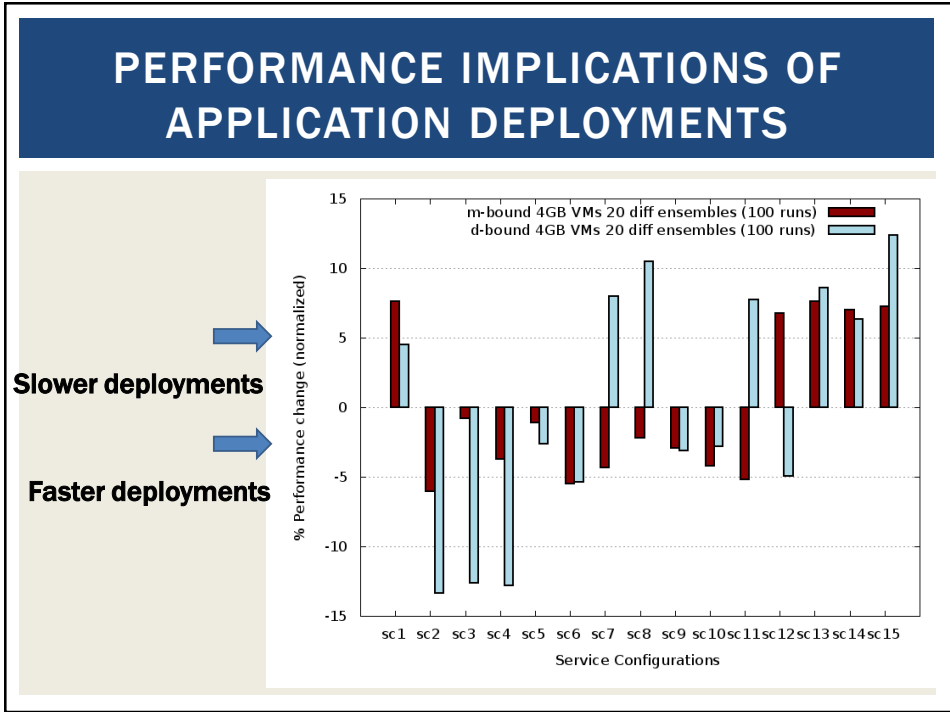
78



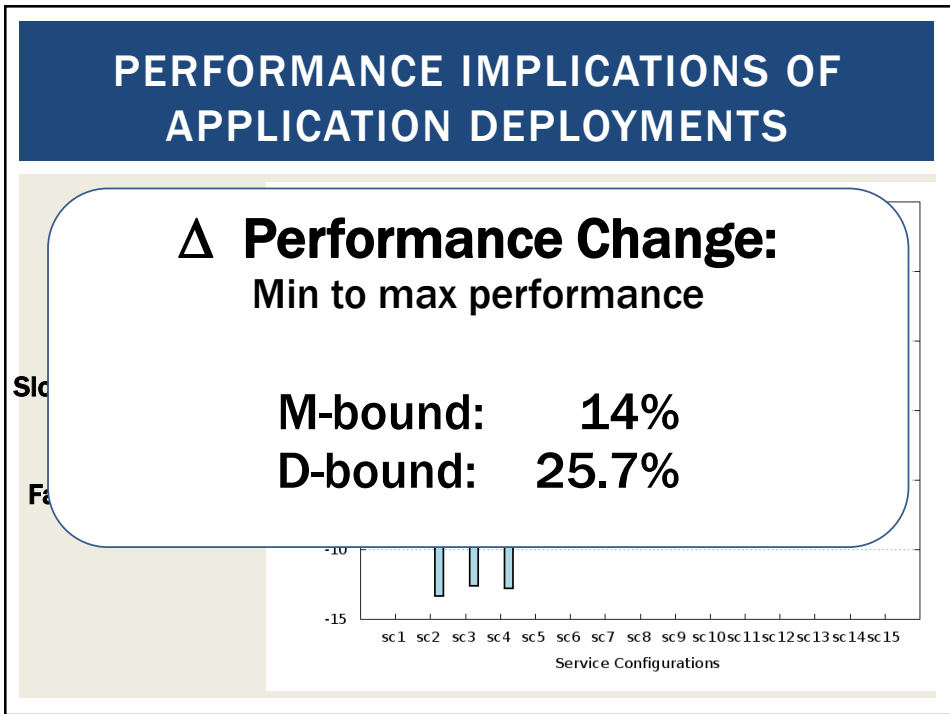
79



80



81



82

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 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.83 |
|------------------|---|-------|

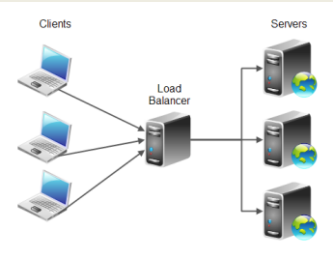
83

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)



The diagram illustrates the Platform-as-a-Service (PaaS) architecture. On the left, three laptop icons labeled 'Clients' have arrows pointing to a central server icon labeled 'Load Balancer'. From the 'Load Balancer', three arrows point to three server rack icons labeled 'Servers', representing the distribution of traffic across multiple instances.

| | | |
|------------------|---|-------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.84 |
|------------------|---|-------|

84

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 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.85 |
|------------------|---|-------|

85

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 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.86 |
|------------------|---|-------|

86

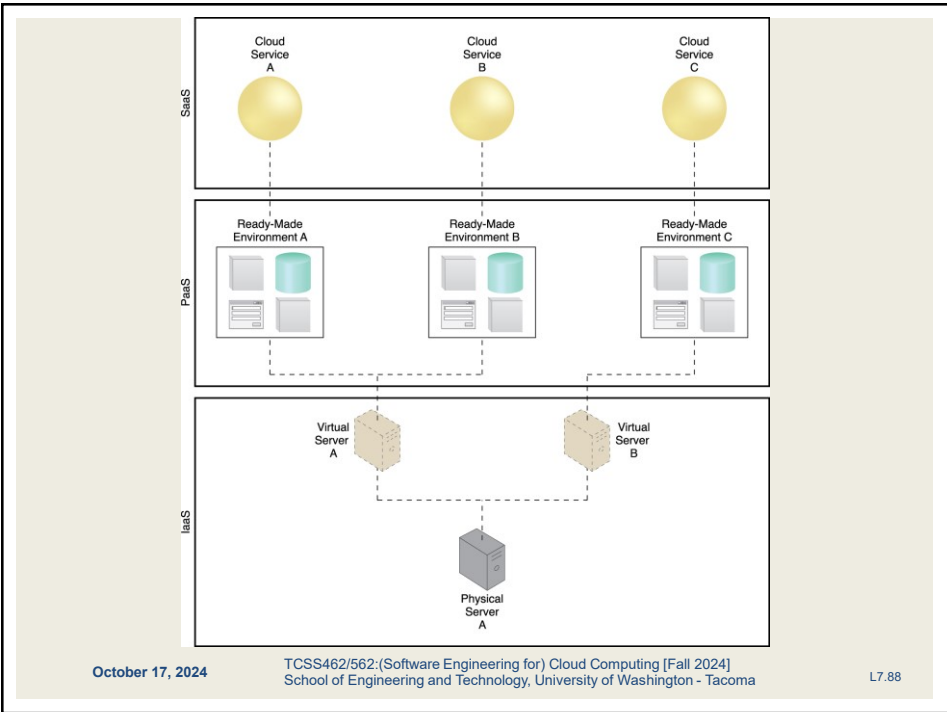
SOFTWARE-AS-A-SERVICE

- 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 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.87 |
|------------------|---|-------|

87



88

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 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.89 |
|------------------|---|-------|

89

SERVERLESS COMPUTING

Introducing Cloud 2.0

Serverless Computing

Deploy Applications Without
Fiddling With Servers



Image from: <https://mobisoftinfotech.com/resources/blog/serverless-computing-deploy-applications-without-fiddling-with-servers/>

90

SERVERLESS COMPUTING

Servers

(AAHHHHHHHHH!!!)

How should my app withstand a server falling?

How can I tell if a server has been compromised?

How can I increase utilization of my servers?

Which OS should my servers run?

How much remaining capacity do my servers have?

How should I implement dynamic configuration changes on my servers?

When should I decide to scale up my servers?

What size servers are right for my budget?

How will I keep my server OS patched?

How can I control access from my servers?

Which packages should be baked into my server images?

How will the application handle server hardware failure?

How will new code be deployed to my servers?

What size server is right for my performance?

How many users create too much load for my servers?

How many servers should I budget for?

Which users should have access to my servers?

Should I tune OS settings to optimize my application?


When should I decide to scale out my servers?


91

SERVERLESS COMPUTING

What is serverless?

Build and run applications without thinking about servers





| | | |
|------------------|---|-------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.92 |
|------------------|---|-------|

92

SERVERLESS COMPUTING - 2

Evolving to serverless

The diagram illustrates the evolution of computing through four stages: 1. Physical servers in datacenters (represented by server racks), 2. Virtual servers in datacenters (represented by server racks with cloud icons), 3. Virtual servers in the cloud (represented by server racks with cloud icons and arrows pointing to a cloud), and 4. SERVERLESS (represented by the AWS Lambda logo inside a cloud). The Amazon Web Services logo is visible in the bottom right corner of the diagram area.

| | | |
|------------------|---|-------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.93 |
|------------------|---|-------|

93

SERVERLESS COMPUTING

The diagram features a large cloud shape containing six text boxes, each representing a characteristic of serverless computing: **Pay only for CPU/memory utilization**, **High Availability**, **Fault Tolerance**, **Infrastructure Elasticity**, **No Setup**, and **Function-as-a-Service (FAAS)**.

94

SERVERLESS COMPUTING

Why Serverless Computing?

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

...they are built into the platform

95

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 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.96 |
|------------------|---|-------|

96

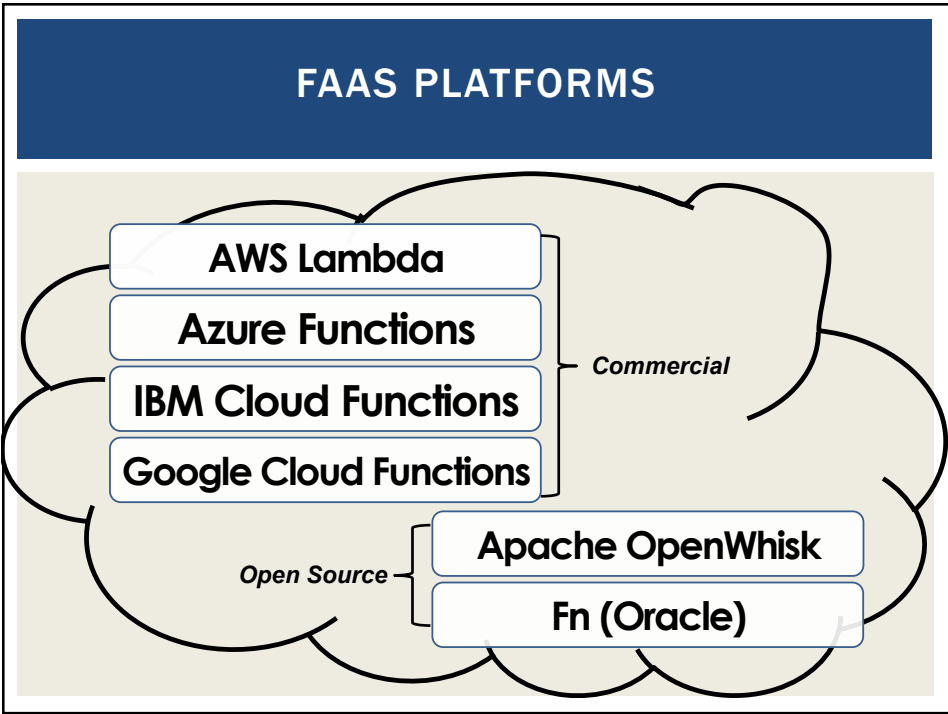
SERVERLESS VS. FAAS

- **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 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.97 |
|------------------|---|-------|


97



98


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

99

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 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.100 |
|------------------|---|--------|

100

FAAS PLATFORMS - 3

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

| | | |
|------------------|---|--------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.101 |
|------------------|---|--------|

101

CLOUD NATIVE SOFTWARE ARCHITECTURE

- Every service with a different pricing model

Example: Weather Application

The diagram illustrates the flow of a weather application. It starts with S3 (Front-end code for weather app hosted in S3), which is accessed by a user (User clicks on link to get local weather information). The user's request goes to an API Gateway (App makes REST API call to endpoint). The API Gateway triggers Lambda (Lambda is triggered), which then interacts with DynamoDB (Lambda runs code to retrieve local weather information and returns data back to user). A temperature of 35° C is shown as the result of the Lambda function.

| | | |
|------------------|---|--------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.102 |
|------------------|---|--------|

102

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 17, 2024
TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024]
School of Engineering and Technology, University of Washington - Tacoma
L7.103

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):*
 - \$0.0000002 per request
 - \$0.000000208 to rent 128MB / 100-ms
 - \$0.00001667 GB /second

October 17, 2024
TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024]
School of Engineering and Technology, University of Washington - Tacoma
L7.104

104

WEBSERVICE HOSTING EXAMPLE

- **ON AWS Lambda**
- Each service call: 100% of 1 CPU-core
100% of 4GB of memory
- Workload: 2 continuous client threads
- Duration: 1 month (30 days)

- **ON AWS EC2:**
- Amazon EC2 c4.large 2-vCPU VM
- Hosting cost: \$72/month
c4.large: 10¢/hour, 24 hrs/day x 30 days

- **How much would hosting this workload cost on AWS Lambda?**

| | | |
|------------------|---|--------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.105 |
|------------------|---|--------|

105

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

106

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 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.107 |
|------------------|---|--------|

107

FACTORS IMPACTING PERFORMANCE OF FAAS COMPUTING PLATFORMS

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

| | | |
|------------------|---|--------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.108 |
|------------------|---|--------|

108

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?

| | | |
|------------------|---|--------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.109 |
|------------------|---|--------|

109

VENDOR ARCHITECTURAL LOCK-IN

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

Example: Weather Application

The diagram illustrates a weather application architecture with the following components and flow:

- S3:** Front-end code for weather app hosted in S3.
- Client:** User clicks on link to get local weather information.
- API GATEWAY:** App makes REST API call to endpoint.
- Lambda:** Lambda is triggered (35° C) and runs code to retrieve local weather information and returns data back to user.
- DYNAMODB:** Database for weather information.

Images credit: aws.amazon.com

- Increased dependencies → increased hosting costs

110


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 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.111 |
|------------------|---|--------|

111

MEMORY RESERVATION QUESTION...



- Lambda memory reserved for functions
- UI provides “slider bar” to set function’s memory allocation
- Resource capacity (CPU, disk, network) coupled to slider bar:
“every doubling of memory, doubles CPU...”
- **But how much memory do model services require?**

▼ Basic settings

Memory (MB) [Info](#)
Your function is allocated CPU proportional to the memory configured.

1536 MB

Timeout [Info](#)
3 min 0 sec

Description

Performance

| | | |
|------------------|---|--------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.112 |
|------------------|---|--------|

112

SERVICE COMPOSITION

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

Monolithic Deployment

Client flow control, 4 functions

Server flow control, 3 functions

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

113


INFRASTRUCTURE FREEZE/THAW CYCLE

- Unused infrastructure is deprecated
 - But after how long?
- Infrastructure: VMs, “containers”
- Provider-COLD / VM-COLD**
 - “Container” images - built/transferred to VMs
- Container-COLD**
 - Image cached on VM
- Container-WARM**
 - “Container” running on VM

Performance

Image from: Denver7 - The Denver Channel News

114



FUNCTION-AS-A-SERVICE

AWS
Lambda
Demo

115

115

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 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.116 |
|------------------|---|--------|

116

CONTAINER-AS-A-SERVICE

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

| | | |
|------------------|---|--------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.117 |
|------------------|---|--------|

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 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.118 |
|------------------|---|--------|

118

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 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L10.119 |
|------------------|---|---------|

119

OBJECTIVES – 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 | L7.120 |
|------------------|---|--------|

120

CLOUD DEPLOYMENT MODELS

- Distinguished by ownership, size, access

- Four common models
 - Public cloud
 - Community cloud
 - Hybrid cloud
 - Private cloud

| | | |
|------------------|---|--------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.121 |
|------------------|---|--------|

121

PUBLIC CLOUDS

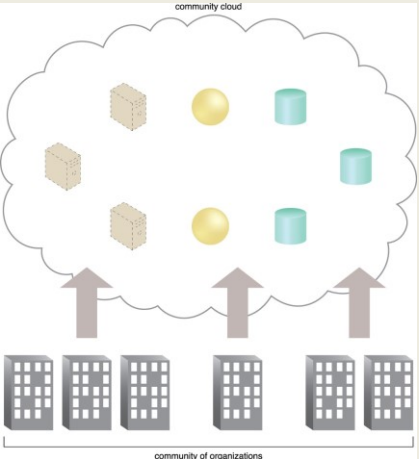
The diagram illustrates the concept of public clouds. At the bottom, three server rack icons are labeled 'organizations'. Three arrows point upwards from these racks to a collection of seven cloud icons. Each cloud icon contains the name of a major public cloud provider: Salesforce, Microsoft, Amazon, Yahoo, Google, Zoho, and Rackspace. This visualizes how organizations utilize services from these external cloud providers.

| | | |
|------------------|---|--------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.122 |
|------------------|---|--------|

122

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

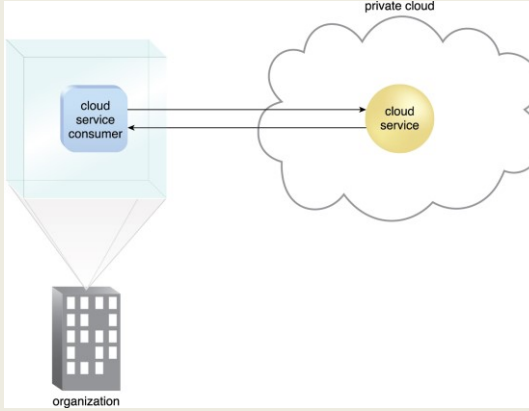


| | | |
|------------------|---|--------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.123 |
|------------------|---|--------|

123

PRIVATE CLOUD

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



| | | |
|------------------|---|--------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.124 |
|------------------|---|--------|

124

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

The diagram illustrates a hybrid cloud architecture. At the bottom, an 'organization' (represented by a server rack icon) is connected to a 'cloud service consumer' (represented by a blue box). This consumer is linked to two cloud environments: a 'public cloud' (top left) and a 'private cloud' (top right). The public cloud contains a 'cloud service' (yellow circle) and 'public data' (green box). The private cloud contains a 'cloud service' (yellow circle) and 'sensitive data' (green box). Bidirectional arrows indicate communication between the consumer and both clouds, and between the two clouds.

| | | |
|------------------|---|--------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.125 |
|------------------|---|--------|

125

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 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.126 |
|------------------|---|--------|

126



OBJECTIVES - 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 | L7.127 |
|------------------|---|--------|

127

TCSS 462/562 TERM PROJECT



| | | |
|------------------|---|--------|
| October 17, 2024 | TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.128 |
|------------------|---|--------|

128

QUESTIONS

October 17, 2024

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

L7.129

129