



TCSS 462/562: (SOFTWARE ENGINEERING FOR) CLOUD COMPUTING

Introduction to Cloud Computing

Wes J. Lloyd
School of Engineering and Technology
University of Washington – Tacoma
TR 5:50-7:50 PM



1

OBJECTIVES – 10/18

- **Questions from 10/13**
 - Properties of Distributed Systems, Modularity
 - Introduction to Cloud Computing –From book #1 - Chapter 3: Understanding Cloud Computing
Cloud Computing Concepts, Technology & Architecture
 - Why study cloud computing?
 - History of cloud computing
 - Business drivers
 - Cloud enabling technologies
 - Terminology
 - Benefits of cloud adoption
 - Risks of cloud adoption
 - From Book #1:
Chapter 4: Cloud Computing Concepts and Models
 - At the end: Questions on the Term Project

October 18, 2022	TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2022] School of Engineering and Technology, University of Washington - Tacoma	L6.2
------------------	---	------

2

OFFICE HOURS – FALL 2022

■ **Tuesdays:**

■ **4:20 to 5:20 pm - CP 229**

■ **Fridays**

■ **12:00 to 1:00 pm – ONLINE via Zoom**

■ **Or email for appointment**

> Office Hours set based on Student Demographics survey feedback

October 18, 2022

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

L6.3

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 18, 2022

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

L6.4

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:

12345678910

Mostly Review To MeEqual New and ReviewMostly New to Me

Question 2

0.5 pts

Please rate the pace of today's class:

12345678910

SlowJust RightFast

October 18, 2022

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

L6.5

5

MATERIAL / PACE

Please classify your perspective on material covered in today's class (49 respondents):

1-mostly review, 5-equal new/review, 10-mostly new

Average - 6.61 (↓ - previous 7.43)

Please rate the pace of today's class:

1-slow, 5-just right, 10-fast

Average - 5.53 (↓ - previous 5.83)

Response rates:

TCSS 462: 25/33 - 75.8%

TCSS 562: 24/26 - 92.3%

October 18, 2022

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

L6.6

6

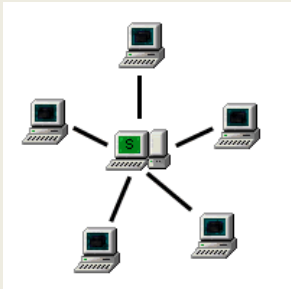
FEEDBACK FROM 10/13

■ Could you please explain more about the multiple points of control and failure?

■ How many nodes can the system suffer the loss of?

■ Depends on which node fails

■ What is the role of each node?



Centralized Architecture
Single-point-of-failure

October 18, 2022

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

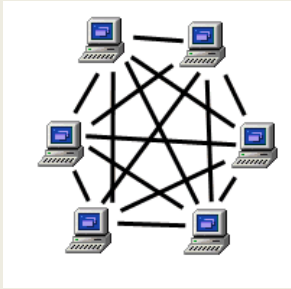
L6.7

7

FEEDBACK FROM 10/13

■ Could you please explain more about the multiple points of control and failure?

■ How many nodes can the system suffer the loss of?



Distributed Architecture
Multiple-points-of-failure

October 18, 2022

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

L6.8

8

FEEDBACK - 2

- *I find most acronyms unfamiliar, and it would be better to see specific examples of how these concepts are utilized in practice.*
 - I am happy to elaborate on specific examples, but would need to know which one(s)..<
- *Many concepts for the project have been discussed and made me feel overwhelmed as most of these things are new to me.*

October 18, 2022

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

L6.9

9

AWS CLOUD CREDITS

- IAM User Accounts Create – please let me know of any issues with these accounts
- If you did not provide your AWS account number on the AWS CLOUD CREDITS SURVEY to request AWS cloud credits and you would like credits this quarter, please contact the professor

October 11, 2022

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

L4.10

10

TUTORIAL 1

- **Introduction to Linux & the Command Line**
- https://faculty.washington.edu/wlloyd/courses/tcss562/tutorials/TCSS462_562_f2022_tutorial_1.pdf
- **Tutorial Sections:**
 1. The Command Line
 2. Basic Navigation
 3. More About Files
 4. Manual Pages
 5. File Manipulation
 6. VI – Text Editor
 7. Wildcards
 8. Permissions
 9. Filters
 10. Grep and regular expressions
 11. Piping and Redirection
 12. Process Management

October 11, 2022

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

L4.11

11

TUTORIAL 2

- **Introduction to Bash Scripting**
- https://faculty.washington.edu/wlloyd/courses/tcss562/tutorials/TCSS462_562_f2022_tutorial_2.pdf
- Review tutorial sections:
 1. What is a BASH script?
 2. Variables
 3. Input
 4. Arithmetic
 5. If Statements
 6. Loops
 7. Functions
 8. User Interface
- Create BASH webservice client
- Call service to obtain IP address & lat/long of computer
- Call weatherbit service to obtain weather forecast for lat/long
 - ➔ *** WEATHERBIT now limited to 7 days ***

October 11, 2022

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

L4.12

12

TUTORIAL 0

- Getting Started with AWS
- http://faculty.washington.edu/wlloyd/courses/tcss562/tutorials/TCSS462_562_f2022_tutorial_0.pdf
- Create an account
- Create account credentials for working with the CLI
- Install awsconfig package
- Setup awsconfig for working with the AWS CLI

October 18, 2022	TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2022] School of Engineering and Technology, University of Washington - Tacoma	L6.13
------------------	---	-------

13

TUTORIAL 3

- Best Practices for Working with Virtual Machines on Amazon EC2
- http://faculty.washington.edu/wlloyd/courses/tcss562/tutorials/TCSS462_562_f2022_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 18, 2022	TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2022] School of Engineering and Technology, University of Washington - Tacoma	L6.14
------------------	---	-------

14

OBJECTIVES – 10/18

- **Questions from 10/13**
- Properties of Distributed Systems, **Modularity**
- Introduction to Cloud Computing – From book #1 - Chapter 3: Understanding Cloud Computing
Cloud Computing Concepts, Technology & Architecture
 - Why study cloud computing?
 - History of cloud computing
 - Business drivers
 - Cloud enabling technologies
 - Terminology
 - Benefits of cloud adoption
 - Risks of cloud adoption
- From Book #1:
Chapter 4: Cloud Computing Concepts and Models
- At the end: Questions on the Term Project

October 18, 2022

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

L6.15

15

TYPES OF MODULARITY

- **Soft modularity:** TRADITIONAL
 - Divide a program into modules (classes) that call each other and communicate with shared-memory
 - A procedure calling convention is used (or method invocation)
 - Examples: object-oriented programming, modularity, etc.
- **Enforced modularity:** CLOUD COMPUTING
 - Program is divided into modules that communicate only through **message passing**
 - The ubiquitous **client-server** paradigm
 - Clients and servers are independent decoupled modules
 - System is more robust if servers are stateless
 - May be scaled and deployed separately
 - May also **FAIL** separately!

October 18, 2022

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

L6.16

16

CLOUD COMPUTING – HOW DID WE GET HERE? SUMMARY OF KEY POINTS

- Multi-core CPU technology and hyper-threading
- What is a
 - Heterogeneous system?
 - Homogeneous system?
 - Autonomous or self-organizing system?
- **Fine grained vs. coarse grained parallelism**
- Parallel message passing code is easier to debug than shared memory (e.g. p-threads)
- Know your application's max/avg **Thread Level Parallelism (TLP)**
- **Data-level parallelism:** Map-Reduce, (SIMD) Single Instruction Multiple Data, Vector processing & GPUs

October 18, 2022

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

L6.17

17

CLOUD COMPUTING – HOW DID WE GET HERE? SUMMARY OF KEY POINTS - 2

- **Bit-level parallelism**
- **Instruction-level parallelism** (CPU pipelining)
- **Flynn's taxonomy:** computer system architecture classification
 - **SISD** – Single Instruction, Single Data (modern core of a CPU)
 - **SIMD** – Single Instruction, Multiple Data (Data parallelism)
 - **MIMD** – Multiple Instruction, Multiple Data
 - MISD is RARE; application for fault tolerance...
- **Arithmetic Intensity:** ratio of calculations vs memory RW
- **Roofline model:**
Memory bottleneck with low arithmetic intensity
- **GPUs:** ideal for programs with high arithmetic intensity
 - SIMD and Vector processing supported by many large registers

October 18, 2022

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

L6.18

18

CLOUD COMPUTING – HOW DID WE GET HERE?
SUMMARY OF KEY POINTS - 3

- **Speed-up (S)**
 $S(N) = T(1) / T(N)$
- **Amdahl's law:**
 $S = 1 / ((1-f) + f/N)$
f= fraction of work that is parallel (e.g. 0.25)
N= proposed speed up of the parallel part (e.g. 5x)
- **Gustafson's Scaled speedup with N processes:**
 $S(N) = N + (1 - \alpha) \alpha$
N: Number of processors
 α : fraction of program run time which can't be parallelized
- Moore's Law
- Symmetric core, Asymmetric core, Dynamic core CPU
- Distributed Systems Non-function quality attributes
- Distributed Systems – Types of Transparency
- Types of modularity- Soft, Enforced


October 18, 2022

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

L6.19

19

INTRODUCTION TO
CLOUD COMPUTING



October 18, 2022

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

L6.20

20

OBJECTIVES – 10/18

- Questions from 10/13
- Properties of Distributed Systems, Modularity
- Introduction to Cloud Computing –From book #1
 - Chapter 3: Understanding Cloud Computing: Cloud Computing Concepts, Technology & Architecture
 - Why study cloud computing?
 - History of cloud computing
 - Business drivers
 - Cloud enabling technologies
 - Terminology
 - Benefits of cloud adoption
 - Risks of cloud adoption
- From Book #1:
Chapter 4: Cloud Computing Concepts and Models
- At the end: Questions on the Term Project

October 18, 2022

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

L6.21

21

OBJECTIVES – 10/18

- Questions from 10/13
- Properties of Distributed Systems, Modularity
- Introduction to Cloud Computing –From book #1 - Chapter 3: Understanding Cloud Computing Cloud Computing Concepts, Technology & Architecture
 - Why study cloud computing?
 - History of cloud computing
 - Business drivers
 - Cloud enabling technologies
 - Terminology
 - Benefits of cloud adoption
 - Risks of cloud adoption
- From Book #1:
Chapter 4: Cloud Computing Concepts and Models
- At the end: Questions on the Term Project

October 18, 2022

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

L6.22

22

WHY STUDY CLOUD COMPUTING?

- **LINKEDIN - TOP IT Skills** from job app data
 - #1 Cloud and Distributed Computing
 - <https://learning.linkedin.com/week-of-learning/top-skills>
 - #2 Statistical Analysis and Data Mining
- **FORBES Survey – 6 Tech Skills That’ll Help You Earn More**
 - #1 Data Science
 - #2 Cloud and Distributed Computing
 - <http://www.forbes.com/sites/laurencebradford/2016/12/19/6-tech-skills-thatll-help-you-earn-more-in-2017/>

October 18, 2022

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

L6.23

23

WHY STUDY CLOUD COMPUTING? - 2

■ **Computerworld Magazine**

TECH FORECAST 2017 SPECIAL REPORT

Hot Skills

Top 10 skills respondents plan to hire for in the next 12 months:

Source: Computerworld's Forecast 2017 survey of 196 IT managers, directors and executives.

Base: 57 respondents who expect to increase IT head count in the next 12 months.

Programming/application development	35%
Help desk/tech support	35%
Security/compliance/governance	26%
Cloud/SaaS	26%
Business intelligence/analytics	26%
Web development	26%
Database administration	25%
Project management	25%
Big data	25%
Mobile applications and device management	21%

© COMPUTERWORLD

October 18, 2022

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

L6.24

24

OBJECTIVES – 10/18

- Questions from 10/13
- Properties of Distributed Systems, Modularity
- Introduction to Cloud Computing –From book #1 - Chapter 3: Understanding Cloud Computing
Cloud Computing Concepts, Technology & Architecture
 - Why study cloud computing?
 - History of cloud computing
 - Business drivers
 - Cloud enabling technologies
 - Terminology
 - Benefits of cloud adoption
 - Risks of cloud adoption
- From Book #1:
Chapter 4: Cloud Computing Concepts and Models
- At the end: Questions on the Term Project

October 18, 2022


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

L6.25

25

A BRIEF HISTORY OF CLOUD COMPUTING

- John McCarthy, 1961
 - Turing award winner for contributions to AI
- “If computers of the kind I have advocated become the computers of the future, then computing may someday be organized as a public utility just as the telephone system is a public utility... The computer utility could become the basis of a new and important industry...”



October 18, 2022

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

L6.26

26

CLOUD HISTORY - 2

- Internet based computer utilities
 - Since the mid-1990s
 - Search engines: Yahoo!, Google, Bing
 - Email: Hotmail, Gmail
- 2000s
 - Social networking platforms: MySpace, Facebook, LinkedIn
 - Social media: Twitter, YouTube
- Popularized core concepts
- Formed basis of cloud computing

October 18, 2022

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

L6.27

27

CLOUD HISTORY: SERVICES - 1

- Late 1990s – Early Software-as-a-Service (SaaS)
 - Salesforce: Remotely provisioned services for the enterprise
- 2002 -
 - Amazon Web Services (AWS) platform: Enterprise oriented services for remotely provisioned storage, computing resources, and business functionality
- 2006 – Infrastructure-as-a-Service (IaaS)
 - Amazon launches Elastic Compute Cloud (EC2) service
 - Organization can “lease” computing capacity and processing power to host enterprise applications
 - Infrastructure

October 18, 2022

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

L6.28

28

Cloud History: Services - 2

- 2006 – **Software-as-a-Service (SaaS)**
 - Google: Offers Google DOCS, “MS Office” like fully-web based application for online documentation creation and collaboration
- 2009 – **Platform-as-a-Service (PaaS)**
 - Google: Offers Google App Engine, publicly hosted platform for hosting scalable web applications on google-hosted datacenters

October 18, 2022


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

L6.29

29

Cloud Computing
NIST General Definition

“Cloud computing is a model for enabling convenient, on-demand network access to a shared pool of configurable computing resources (networks, servers, storage, applications and services) that can be rapidly provisioned and reused with minimal management effort or service provider interaction”...



October 18, 2022

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

L6.30

30

MORE CONCISE DEFINITION

“Cloud computing is a specialized form of distributed computing that introduces utilization models for remotely provisioning scalable and measured resources.”

From Cloud Computing Concepts, Technology, and Architecture
Z. Mahmood, R. Puttini, Prentice Hall, 5th printing, 2015

October 18, 2022	TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2022] School of Engineering and Technology, University of Washington - Tacoma	L6.31
------------------	---	-------

31

OBJECTIVES – 10/18

- Questions from 10/13
- Properties of Distributed Systems, Modularity
- Introduction to Cloud Computing –From book #1 - Chapter 3: Understanding Cloud Computing
Cloud Computing Concepts, Technology & Architecture
 - Why study cloud computing?
 - History of cloud computing
 - Business drivers
 - Cloud enabling technologies
 - Terminology
 - Benefits of cloud adoption
 - Risks of cloud adoption
- From Book #1:
Chapter 4: Cloud Computing Concepts and Models
- At the end: Questions on the Term Project

October 18, 2022	TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2022] School of Engineering and Technology, University of Washington - Tacoma	L6.32
------------------	---	-------

32

**BUSINESS DRIVERS
FOR CLOUD COMPUTING**

- Capacity planning
- Cost reduction
- Operational overhead
- Organizational agility

October 18, 2022

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

L6.33

33

**BUSINESS DRIVERS
FOR CLOUD COMPUTING**

- Capacity planning
 - Process of determining and fulfilling future demand for IT resources
- Capacity vs. demand
 - Discrepancy between capacity of IT resources and actual demand
- Over-provisioning: resource capacity exceeds demand
- Under-provisioning: demand exceeds resource capacity
- Capacity planning aims to minimize the discrepancy of available resources vs. demand

October 18, 2022

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

L6.34

34



Dwight, The Office TV sitcom

October 18, 2022

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

L6.35

35

BUSINESS DRIVERS FOR CLOUD - 2

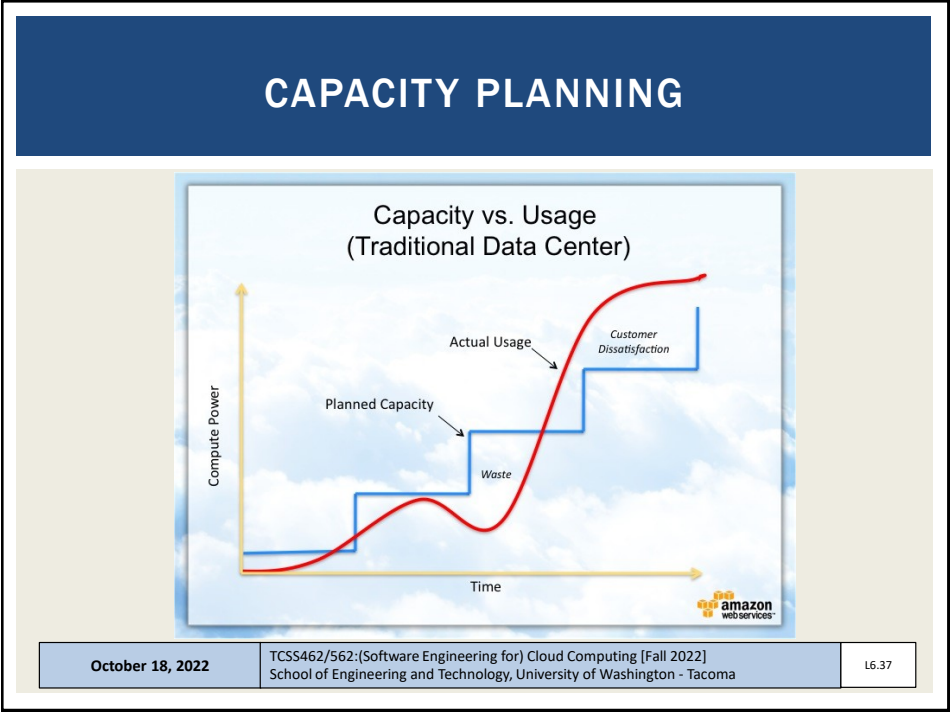
- Capacity planning
 - Over-provisioning: is costly due to too much infrastructure
 - Under-provisioning: is costly due to potential for business loss from poor quality of service
- Capacity planning strategies
 - Lead strategy: add capacity in anticipation of demand (pre-provisioning)
 - Lag strategy: add capacity when capacity is fully leveraged
 - Match strategy: add capacity in small increments as demand increases
- Load prediction
 - Capacity planning helps anticipate demand fluctuations

October 18, 2022

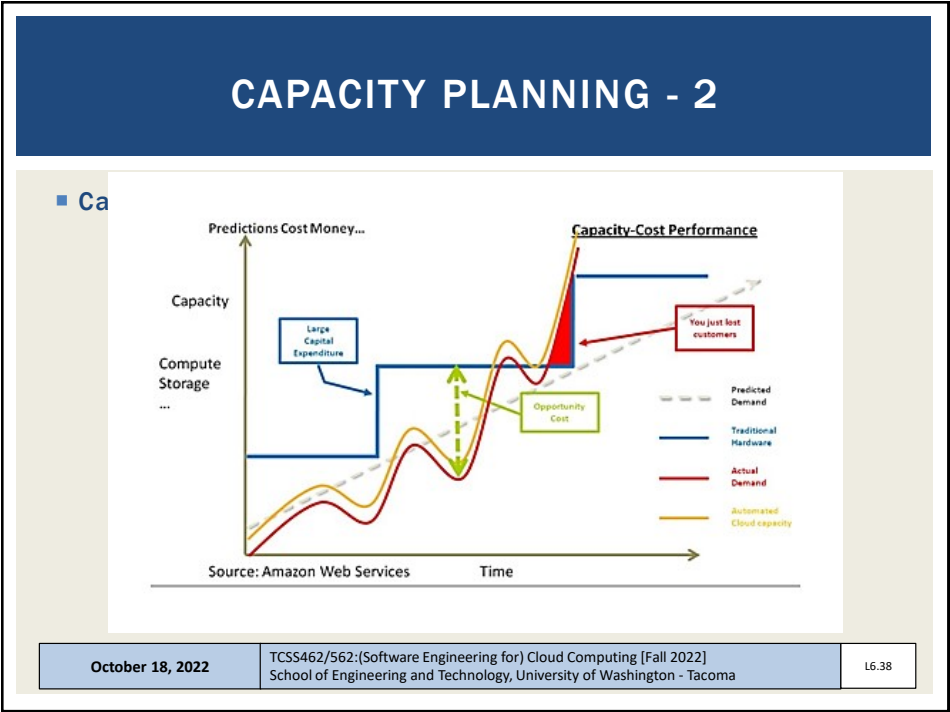
TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2022]
School of Engineering and Technology, University of Washington - Tacoma

L6.36

36



37



38

BUSINESS DRIVERS FOR CLOUD - 3

- **Cost reduction**
 - IT Infrastructure acquisition
 - IT Infrastructure maintenance
- **Operational overhead**
 - Technical personnel to maintain physical IT infrastructure
 - System upgrades, patches that add testing to deployment cycles
 - Utility bills, capital investments for power and cooling
 - Security and access control measures for server rooms
 - Admin and accounting staff to track licenses, support agreements, purchases

October 18, 2022

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

L6.39

39

BUSINESS DRIVERS FOR CLOUD - 4

- **Organizational agility**
 - Ability to adapt and evolve infrastructure to face change from internal and external business factors
 - Funding constraints can lead to insufficient on premise IT
 - Cloud computing enables IT resources to scale with a lower financial commitment

October 18, 2022

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

L6.40

40

OBJECTIVES – 10/18

- Questions from 10/13
- Properties of Distributed Systems, Modularity
- Introduction to Cloud Computing –From book #1 - Chapter 3: Understanding Cloud Computing
Cloud Computing Concepts, Technology & Architecture
 - Why study cloud computing?
 - History of cloud computing
 - Business drivers
 - Cloud enabling technologies
 - Terminology
 - Benefits of cloud adoption
 - Risks of cloud adoption
- From Book #1:
Chapter 4: Cloud Computing Concepts and Models
- At the end: Questions on the Term Project

October 18, 2022

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

L6.41

41

TECHNOLOGY INNOVATIONS
LEADING TO CLOUD

- Cluster computing
- Grid computing
- Virtualization
- Others


October 18, 2022

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

L6.42

42

CLUSTER COMPUTING




- Cluster computing (clustering)
 - Cluster is a group of independent IT resources interconnected as a single system
 - Servers configured with homogeneous hardware and software
 - Identical or similar RAM, CPU, HDDs
 - Design emphasizes redundancy as server components are easily interchanged to keep overall system running
 - Example: if a RAID card fails on a key server, the card can be swapped from another redundant server
 - Enables warm replica servers
 - Duplication of key infrastructure servers to provide HW failover to ensure high availability (HA)

October 18, 2022	TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2022] School of Engineering and Technology, University of Washington - Tacoma	L6.43
------------------	---	-------

43

GRID COMPUTING



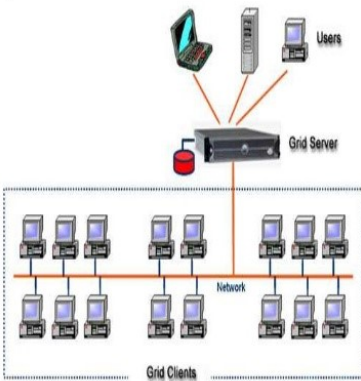
- On going research area since early 1990s
- Distributed heterogeneous computing resources organized into logical pools of loosely coupled resources
- For example: heterogeneous servers connected by the internet
- Resources are heterogeneous and geographically dispersed
- Grids use middleware software layer to support workload distribution and coordination functions
- Aspects: load balancing, failover control, autonomic configuration management
- Grids have influenced clouds contributing common features: networked access to machines, resource pooling, scalability, and resiliency

October 18, 2022	TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2022] School of Engineering and Technology, University of Washington - Tacoma	L6.44
------------------	---	-------

44

GRID COMPUTING - 2

How Grid computing works ?



The diagram illustrates the architecture of a grid computing system. At the top, three user devices (laptop, PDA, and desktop) are labeled 'Users'. Arrows from these users point to a central 'Grid Server' represented by a server rack icon. Below the server, a 'Network' is shown as a horizontal line with multiple vertical connections leading to a group of computer icons labeled 'Grid Clients'. The entire system is enclosed in a dashed-line box.

In general, a grid computing system requires:

- At least one computer, usually a server, which handles all the administrative duties for the System
- A network of computers running special grid computing network software.
- A collection of computer software called middleware

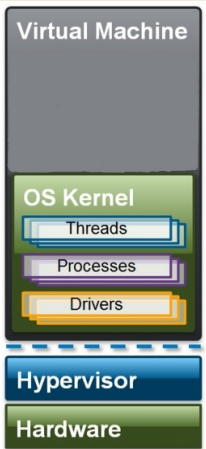
October 18, 2022

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

L6.45

45

VIRTUALIZATION



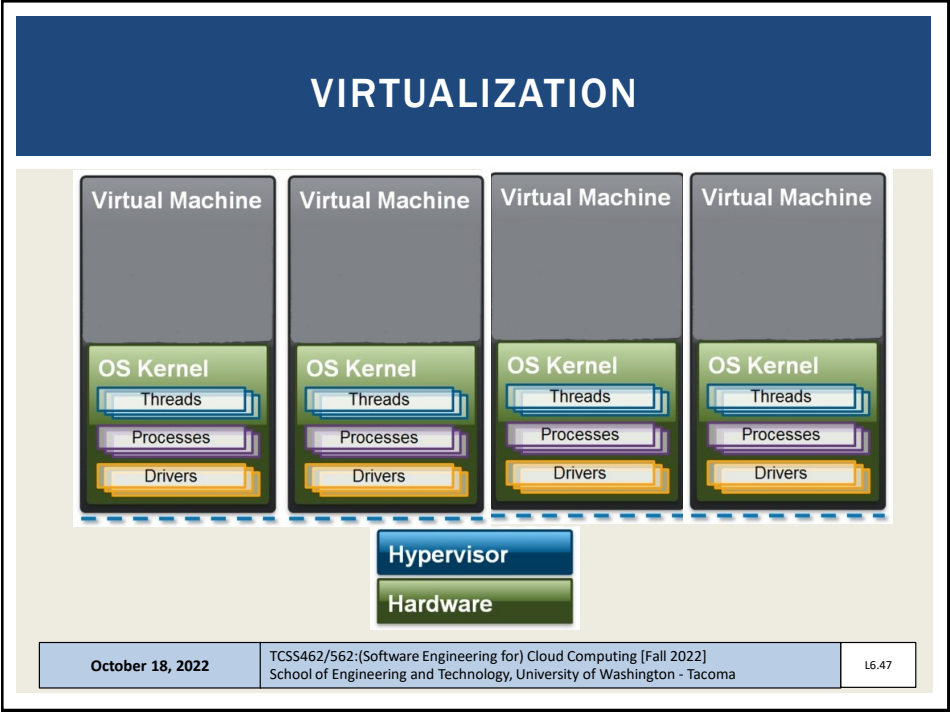
The diagram shows a vertical stack of components representing the virtualization layer. From top to bottom, the layers are: 'Virtual Machine' (grey box), 'OS Kernel' (green box) which contains sub-components 'Threads' (blue), 'Processes' (purple), and 'Drivers' (yellow); 'Hypervisor' (blue box); and 'Hardware' (green box). A dashed line separates the OS Kernel from the Hypervisor.

October 18, 2022

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

L6.46

46



47

VIRTUALIZATION

- Simulate physical hardware resources via software
 - The virtual machine (virtual computer)
 - Virtual local area network (VLAN)
 - Virtual hard disk
 - Virtual network attached storage array (NAS)
- Early incarnations featured significant performance, reliability, and scalability challenges
- CPU and other HW enhancements have minimized performance GAPS

October 18, 2022	TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2022] School of Engineering and Technology, University of Washington - Tacoma	L6.48
------------------	---	-------

48

OBJECTIVES – 10/18

- **Questions from 10/13**
- Properties of Distributed Systems, Modularity
- Introduction to Cloud Computing –From book #1 - Chapter 3: Understanding Cloud Computing
Cloud Computing Concepts, Technology & Architecture
 - Why study cloud computing?
 - History of cloud computing
 - Business drivers
 - Cloud enabling technologies
 - Terminology
 - Benefits of cloud adoption
 - Risks of cloud adoption
- At the end: Questions on the Term Project

October 18, 2022

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

L6.49

49

KEY TERMINOLOGY

- On-Premise Infrastructure
 - Local server infrastructure not configured as a cloud
- Cloud Provider
 - Corporation or private organization responsible for maintaining cloud
- Cloud Consumer
 - User of cloud services
- Scaling
 - Vertical scaling
 - Scale up: increase resources of a single virtual server
 - Scale down: decrease resources of a single virtual server
 - Horizontal scaling
 - Scale out: increase number of virtual servers
 - Scale in: decrease number of virtual servers

October 18, 2022

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

L6.50

50

VERTICAL SCALING

■ Reconfigure virtual machine to have different resources:

- CPU cores
- RAM
- HDD/SDD capacity

■ May require VM migration if physical host machine resources are exceeded

The diagram illustrates vertical scaling. It shows two virtual machines, A and B, represented as 3D blocks. Machine A is at the bottom, labeled '2 CPUs'. Machine B is above it, labeled '4 CPUs'. A vertical arrow points from A to B, with the text 'vertical scaling' written vertically along the arrow.

October 18, 2022

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

L6.51

51

HORIZONTAL SCALING

■ Increase (scale-out) or decrease (scale-in) number of virtual servers based on demand

The diagram illustrates horizontal scaling. At the top, two physical servers are labeled 'pooled physical servers'. Arrows point from these servers to a row of virtual servers. The virtual servers are labeled A, A, B, A, B, and C. Arrows labeled 'demand' point from the first A to the second A, and from the second A to the B. A long arrow at the bottom points to the right and is labeled 'horizontal scaling'.

October 18, 2022

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

L6.52

52

HORIZONTAL VS VERTICAL SCALING

Horizontal Scaling	Vertical Scaling
Less expensive using commodity HW	Requires expensive high capacity servers

October 18, 2022

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

L6.53

53

HORIZONTAL VS VERTICAL SCALING

Horizontal Scaling	Vertical Scaling
Less expensive using commodity HW	Requires expensive high capacity servers
IT resources instantly available	IT resources typically instantly available

October 18, 2022

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

L6.54

54

HORIZONTAL VS VERTICAL SCALING

Horizontal Scaling	Vertical Scaling
Less expensive using commodity HW	Requires expensive high capacity servers
IT resources instantly available	IT resources typically instantly available
Resource replication and automated scaling	Additional setup is normally needed

October 18, 2022

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

L6.55

55

HORIZONTAL VS VERTICAL SCALING

Horizontal Scaling	Vertical Scaling
Less expensive using commodity HW	Requires expensive high capacity servers
IT resources instantly available	IT resources typically instantly available
Resource replication and automated scaling	Additional setup is normally needed
Additional servers required	No additional servers required

October 18, 2022

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

L6.56

56

HORIZONTAL VS VERTICAL SCALING

Horizontal Scaling	Vertical Scaling
Less expensive using commodity HW	Requires expensive high capacity servers
IT resources instantly available	IT resources typically instantly available
Resource replication and automated scaling	Additional setup is normally needed
Additional servers required	No additional servers required
Not limited by individual server capacity	Limited by individual server capacity

October 18, 2022

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

L6.57

57

KEY TERMINOLOGY - 2

- Cloud services
 - Broad array of resources accessible “as-a-service”
 - Categorized as Infrastructure (IaaS), Platform (PaaS), Software (SaaS)
- Service-level-agreements (SLAs):
 - Establish expectations for: uptime, security, availability, reliability, and performance

October 18, 2022

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

L6.58

58

OBJECTIVES – 10/18

- Questions from 10/13
 - Properties of Distributed Systems, Modularity
 - Introduction to Cloud Computing –From book #1 - Chapter 3: Understanding Cloud Computing
 - Cloud Computing Concepts, Technology & Architecture
 - Why study cloud computing?
 - History of cloud computing
 - Business drivers
 - Cloud enabling technologies
 - Terminology
 - Benefits of cloud adoption
 - Risks of cloud adoption
 - From Book #1:
 - Chapter 4: Cloud Computing Concepts and Models
 - At the end: Questions on the Term Project

October 18, 2022

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

L6.59

59

GOALS AND BENEFITS

- Cloud providers
 - Leverage economies of scale through mass-acquisition and management of large-scale IT resources
 - Locate datacenters to optimize costs where electricity is low
- Cloud consumers
 - Key business/accounting difference:
 - Cloud computing enables anticipated capital expenditures to be replaced with operational expenditures
 - Operational expenditures always scale with the business
 - Eliminates need to invest in server infrastructure based on anticipated business needs
 - Businesses become more agile and lower their financial risks by eliminating large capital investments in physical infrastructure

October 18, 2022

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


L6.60

60

CLOUD BENEFITS - 2

- On demand access to pay-as-you-go resources on a short-term basis (less commitment)
- Ability to acquire “unlimited” computing resources on demand when required for business needs
- Ability to add/remove IT resources at a fine-grained level
- Abstraction of server infrastructure so applications deployments are not dependent on specific locations, hardware, etc.
 - The cloud has made our software deployments more agile...

Before Cloud Computing?



October 18, 2022	TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2022] School of Engineering and Technology, University of Washington - Tacoma	L6.61
------------------	---	-------

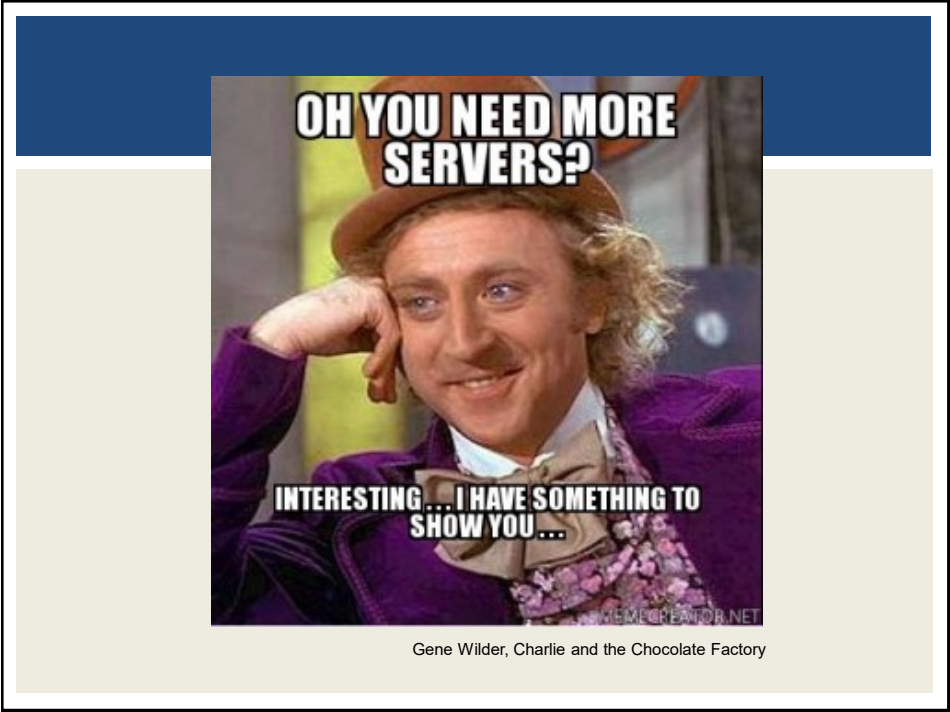
61

CLOUD BENEFITS - 3

- Example: Using 100 servers for 1 hour costs the same as using 1 server for 100 hours
- Rosetta Protein Folding Use Case: Working with a UW-Tacoma graduate student, we deployed this science model across 5,900 compute cores on Amazon for 2-days...
- *What is the cost to purchase 5,900 compute cores?*
- Recent Dell Server purchase example:
20 cores on 2 servers for \$4,478...
- Using this ratio 5,900 cores costs \$1.3 million (purchase only)

October 18, 2022	TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2022] School of Engineering and Technology, University of Washington - Tacoma	L6.62
------------------	---	-------

62



63

CLOUD BENEFITS

- Increased scalability
 - Example demand over a 24-hour day →
- Increased availability
- Increased reliability

A line graph showing the number of concurrent users over a 24-hour period. The y-axis is labeled "concurrent users" and ranges from 1,000 to 10,000 in increments of 1,000. The x-axis is labeled "time (h)" and ranges from 2 to 24 in increments of 2. The graph shows a smooth curve that starts at approximately 2,000 users at 2 AM, dips to a minimum of about 1,500 users at 6 AM, then rises sharply to a peak of about 9,500 users at 16:00 (4 PM), and finally declines to about 2,500 users by 24:00 (midnight). The area under the curve is shaded in a light orange color.

October 18, 2022

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

L6.64

64

OBJECTIVES – 10/18

- **Questions from 10/13**
- Properties of Distributed Systems, Modularity
- Introduction to Cloud Computing –From book #1 - Chapter 3: Understanding Cloud Computing
Cloud Computing Concepts, Technology & Architecture
 - Why study cloud computing?
 - History of cloud computing
 - Business drivers
 - Cloud enabling technologies
 - Terminology
 - Benefits of cloud adoption
 - Risks of cloud adoption
- From Book #1:
Chapter 4: Cloud Computing Concepts and Models
- At the end: Questions on the Term Project

October 18, 2022	TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2022] School of Engineering and Technology, University of Washington - Tacoma	L6.65
------------------	---	-------

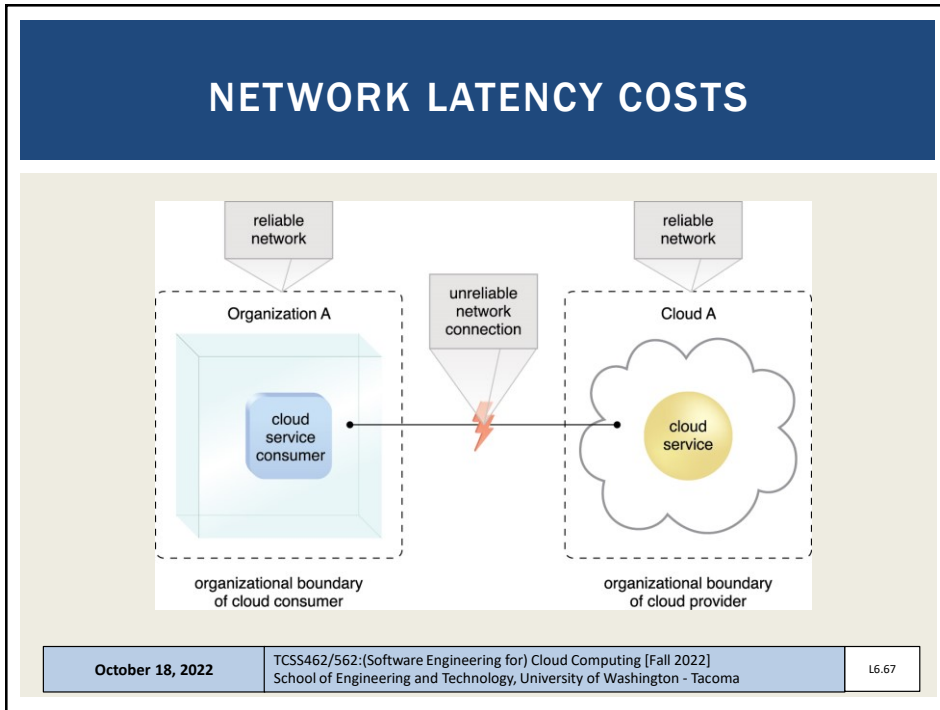
65

CLOUD ADOPTION RISKS

- **Increased security vulnerabilities**
 - Expansion of trust boundaries now include the external cloud
 - Security responsibility shared with cloud provider
- **Reduced operational governance / control**
 - Users have less control of physical hardware
 - Cloud user does not directly control resources to ensure quality-of-service
 - Infrastructure management is abstracted
 - Quality and stability of resources can vary
 - Network latency costs and variability

October 18, 2022	TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2022] School of Engineering and Technology, University of Washington - Tacoma	L6.66
------------------	---	-------

66



67

CLOUD RISKS - 2

- **Performance monitoring of cloud applications**
 - Cloud metrics (AWS cloudwatch) support monitoring cloud infrastructure (network load, CPU utilization, I/O)
 - Performance of cloud applications depends on the health of aggregated cloud resources working together
 - User must monitor this aggregate performance
- **Limited portability among clouds**
 - Early cloud systems have significant “vendor” lock-in
 - Common APIs and deployment models are slow to evolve
 - Operating system containers help make applications more portable, but containers still must be deployed
- **Geographical issues**
 - Abstraction of cloud location leads to legal challenges with respect to laws for data privacy and storage

October 18, 2022	TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2022] School of Engineering and Technology, University of Washington - Tacoma	L6.68
------------------	---	-------

68

CLOUD: VENDOR LOCK-IN

Cloud A (Cloud Provider X)

supports message encryption and digital signatures

cloud consumer

requires encryption and digital signing of messages

Cloud B (Cloud Provider Y)

supports message encryption only

October 18, 2022

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

L6.69

69

CLOUD COMPUTING:
CONCEPTS AND MODELS

October 18, 2022

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

L6.70

70

OBJECTIVES – 10/18

- Questions from 10/13
- From: Cloud Computing Concepts, Technology & Architecture:
Chapter 3: Understanding Cloud Computing
- From: Cloud Computing Concepts, Technology & Architecture:
Chapter 4: Cloud Computing Concepts and Models:
 - Roles and boundaries
 - Cloud characteristics
- At the end: Questions on the Term Project
 - TCSS 462/562 Term Project
 - Team Planning

October 18, 2022

TCSS562:Software Engineering for Cloud Computing [Fall 2021]
School of Engineering and Technology, University of Washington - Tacoma

L6.71

71

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 18, 2022

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

L6.72

72

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 18, 2022

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

L6.73

73

ORGANIZATION BOUNDARY

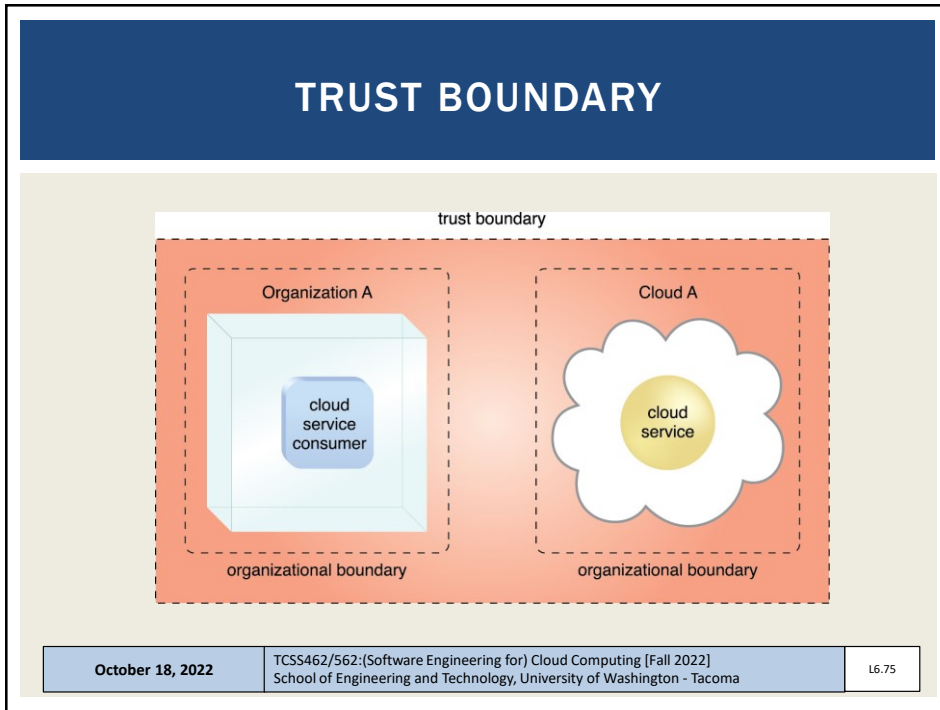
The diagram illustrates the concept of an organization boundary in cloud computing. It consists of two side-by-side boxes, each labeled "Organization A" and "Cloud A" at the top. The left box contains a blue cube labeled "cloud service consumer" and is labeled "organizational boundary" at the bottom. The right box contains a yellow circle labeled "cloud service" and is also labeled "organizational boundary" at the bottom. Both boxes are enclosed in dashed lines.

October 18, 2022

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

L6.74

74



75

OBJECTIVES – 10/18

- Questions from 10/13
- **From: Cloud Computing Concepts, Technology & Architecture:**
Chapter 3: Understanding Cloud Computing
- **From: Cloud Computing Concepts, Technology & Architecture:**
Chapter 4: Cloud Computing Concepts and Models:
 - Roles and boundaries
 - **Cloud characteristics**
- At the end: Questions on the Term Project
 - TCSS 562 Term Project
 - Team Planning

October 18, 2022	TCSS562:Software Engineering for Cloud Computing [Fall 2021] School of Engineering and Technology, University of Washington - Tacoma	L6.76
------------------	---	-------

76

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 18, 2022

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

L6.77

77

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 Center

National Informatics Centre

Data Center and Web Services Division

Virtual Machine Request Form

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

1. Name of the VMC Group / Division

2. Name of the Project / Service
(If Machine Description & Architecture are a separate sheet)

3. Category: Web | Database | Other |

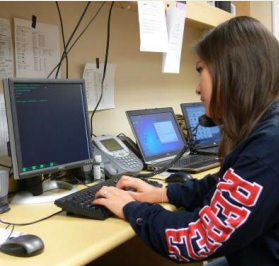
Others if any specify:

4. Virtual Machine Specification

- Name of the Virtual Machine
- Operating System (OS) (Please specify the VM)
- CPU Required
- RAM Required

5. Software Environment

- Operating System (with version)
- Software & Tools
- Software Licenses (Detail including VM)
- Application provide access VM to will maintain the application



October 18, 2022

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

L6.78

78

UBIQUITOUS ACCESS

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

October 18, 2022

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

L6.79

79

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 18, 2022

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

L6.80

80

SINGLE TENANT MODEL

> Isolation <

October 18, 2022

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

L6.81

81

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

shared cloud storage device

October 18, 2022


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

L6.82

82

MULTITENANT DATABASE

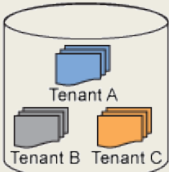
Isolated



Separate database

E1

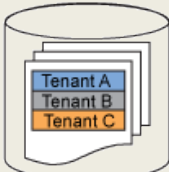
Semi-shared



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 18, 2022

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

L6.83

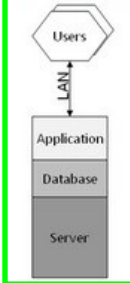
83

MULTITENANCY OF RESOURCES

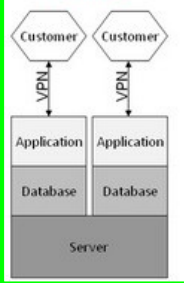
■ Where is the multitenancy?

■ >> What is shared? What is isolated?

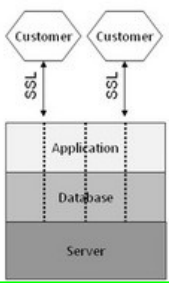
Traditional On Premise



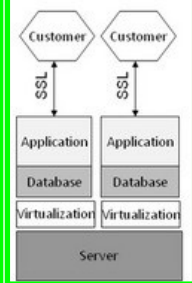
Single Tenant (Hosted)



Multi-Tenant



Virtual Appliance



October 18, 2022

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

L6.84

84

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%	85%	90%	25%
30	85%	75%	85%	15%
40	80%	65%	80%	10%

Up to 48 VMs sharing same server !!

October 18, 2022

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

L6.85

85

RESOURCE CONTENTION FROM MUTLI-TENANCY - 2

■ Performance variation from multi-tenancy is increasing as cloud servers add more CPU cores

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.

EC2 Instance family	iperf (network)	pgbench (CPU + I/O)	sysbench (CPU)	y-cruncher (CPU)
c3	19.2%	19.2%	0.3%	24.6%
c4	42.1%	5.6%	0.2%	0.1%
z1d	84.6%	11.2%	0.2%	0.0%
m5d (t)	94.6%	33.0%	20.8%	48.0%

■ Running many idle operating system instances can impose significant overhead for some workloads

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

October 18, 2022

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

L6.86

86

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 18, 2022

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

L6.87

87

PREDICTABLE DEMAND

- AWS EC2 Scaling Example:

Auto-Scaling Example: Netflix

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

October 18, 2022

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

L6.88

88

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 18, 2022

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

L6.89

89

EC2 CLOUDWATCH METRICS

EC2 Instance: i-1267037f

Description Monitoring Tags

Graphs are for 1 instance that has monitoring enabled. Times are displayed in UTC.

Time Range: Last Hour Refresh

Avg CPU Utilization (Percent)

Avg Disk Reads (Bytes)

Avg Disk Writes (Bytes)

Max Network In (Bytes)

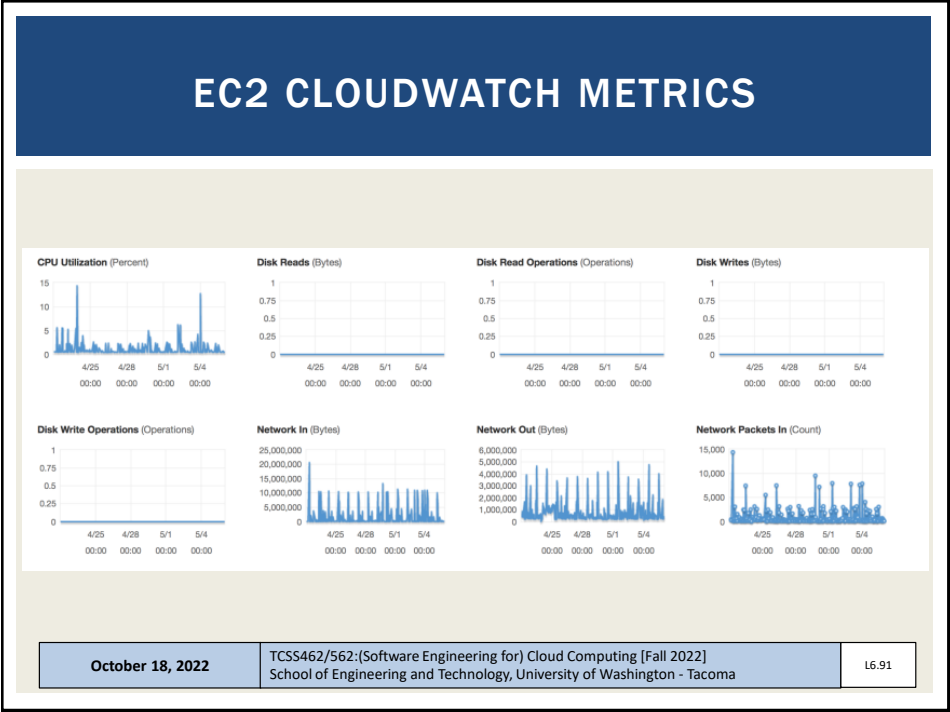
Max Network Out (Bytes)

October 18, 2022

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

L6.90

90



91

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

Engineering at Cloud Scale

Resilience and Reliability on AWS

Jurg van Vleet, Flavio Paegemöller & Jasper Geurtsen

O'REILLY®

October 18, 2022	TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2022] School of Engineering and Technology, University of Washington - Tacoma	L6.92
------------------	---	-------

92

OBJECTIVES – 10/18

- Questions from 10/13
- Properties of Distributed Systems, Modularity
- Introduction to Cloud Computing –From book #1 - Chapter 3: Understanding Cloud Computing
 - Cloud Computing Concepts, Technology & Architecture
 - Why study cloud computing?
 - History of cloud computing
 - Business drivers
 - Cloud enabling technologies
 - Terminology
 - Benefits of cloud adoption
 - Risks of cloud adoption
- From Book #1:
 - Chapter 4: Cloud Computing Concepts and Models
 - At the end: Questions on the Term Project



October 18, 2022

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

L6.93

93

TCSS 462/562
TERM PROJECT



October 18, 2022

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

L6.94

94

TCSS 462/562 TERM PROJECT

- Build a serverless cloud native application
- Application provides case study to investigate architecture/design trade-offs
 - Application provides a vehicle to compare and contrast one or more trade-offs
- Alternate 1: Cloud Computing Related Research Project
- Alternate 2: Literature Survey/Gap Analysis
 - *- as an individual project*

October 18, 2022

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

L6.95

95

DESIGN TRADE-OFFS

- Service composition
 - Switchboard architecture:
 - compose services in single package
 - Address COLD Starts
 - Infrastructure Freeze/Thaw cycle of AWS Lambda (FaaS)
 - Full service isolation (each service is deployed separately)
- Application flow control
 - client-side, step functions, server-side controller, asynchronous hand-off
- Programming Languages
- Alternate FaaS Platforms

October 18, 2022

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

L6.96

96

DESIGN TRADE-OFFS - 2

- **Alternate Cloud Services (e.g. databases, queues, etc.)**
 - Compare alternate data backends for data processing pipeline
- **Performance variability (by hour, day, week, and host location)**
 - Deployments (to different zones, regions)
- **Service abstraction**
 - Abstract one or more services with cloud abstraction middleware: Apache libcloud, apache jcloud; make code cross-cloud; measure overhead

October 18, 2022

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

L6.97

97

OTHER PROJECT IDEAS

- **Elastic File System (EFS)**
Performance & Scalability Evaluation
- **Docker container image integration with AWS Lambda – performance & scalability**
- **Resource contention study using CpuSteal metric**
 - Investigate the degree of CpuSteal on FaaS platforms
 - What is the extent? Min, max, average
 - When does it occur?
 - Does it correlate with performance outcomes?
 - Is contention self-inflicted?
- **& others**

October 18, 2022

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

L6.98

98

SERVERLESS APPLICATIONS

- **Extract Transform Load Data Processing Pipeline**
 - * >>>This is the STANDARD project<<< *
 - Batch-oriented data
 - Stream-oriented data
- **Image Processing Pipeline**
 - Apply series of filters to images
- **Stream Processing Pipeline**
 - Data conversion, filtering, aggregation, archival storage
 - What throughput (records/sec) can Lambda ingest directly?
 - Comparison with AWS Kinesis Data Streams and DB backend:
 - <https://aws.amazon.com/getting-started/hands-on/build-serverless-real-time-data-processing-app-lambda-kinesis-s3-dynamodb-cognito-athena/>
 - Kinesis data streams claims multiple GB/sec throughput
 - What is the cost difference?

October 18, 2022

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

L6.99

99

SERVERLESS APPLICATIONS - 2

- **Map-Reduce Style Application**
 - Function 1: split data into chunks, usually sequentially
 - Function 2: process individual chunks concurrently (in parallel)
 - Data process is considered to be Embarrassingly Parallel
 - Function 3: aggregate and summarize results
- **Image Classification Pipeline**
 - Deploy pretrained image classifiers in a multi-stage pipeline
- **Machine Learning**
 - Multi-stage inferencing pipelines
 - Natural Language Processing (NLP) pipelines
 - Training (?)

October 18, 2022

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

L6.100

100

AWS LAMBDA PLATFORM LIMITATIONS

- Maximum 10 GB memory per function instance
- Maximum 15-minutes execution per function instance
- 500 MB of temporary disk space for local I/O (default)
- 10 GB ephemeral storage (for additional charge)
 - <https://aws.amazon.com/blogs/aws/aws-lambda-now-supports-up-to-10-gb-ephemeral-storage/>
- Access up to 6 vCPUs depending on memory reservation size
- 1,000 concurrent function executions inside account (default)
- Function payload: 6MB (synchronous), 256KB (asynchronous)
- Deployment package: 50MB (compressed), 250MB (unzipped)
- Container image size: 10 GB
- Processes/threads: 1024
- File descriptors: 1024
- See: <https://docs.aws.amazon.com/lambda/latest/dg/gettingstarted-limits.html>

October 18, 2022	TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2022] School of Engineering and Technology, University of Washington - Tacoma	L6.101
------------------	---	--------

101

EXTRACT TRANSFORM LOAD DATA PIPELINE

- Service 1: **TRANSFORM**
 - Read CSV file, perform some transformations
 - Write out new CSV file
- Service 2: **LOAD**
 - Read CSV file, load data into relational database
 - Cloud DB (AWS Aurora), or local DB (Derby/SQLite)
 - Derby DB and/or SQLite code examples to be provided in Java

October 18, 2022	TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2022] School of Engineering and Technology, University of Washington - Tacoma	L6.102
------------------	---	--------

102

EXTRACT TRANSFORM LOAD
DATA PIPELINE - 2

- Service 3: **QUERY**
- Using relational database, apply filter(s) and/or functions to aggregate data to produce sums, totals, averages
- Output aggregations as JSON

October 18, 2022

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

L6.103

103

SERVICE COMPOSITION

Remote Client

API Gateway

Fine grained services

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.

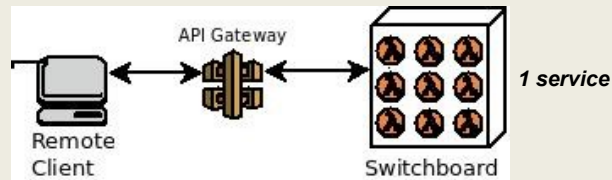
October 18, 2022

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

L6.104

104

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 18, 2022

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

L6.105

105

APPLICATION FLOW CONTROL

- **Serverless Computing:**
 - AWS Lambda (FAAS: [Function-as-a-Service](#))
 - Provides HTTP/REST like web services
 - Client/Server paradigm
- **Synchronous web service:**
 - Client calls service
 - Client blocks (freezes) and waits for server to complete call
 - Connection is maintained in the “OPEN” state
 - Problematic if service runtime is long!
 - Connections are notoriously dropped
 - System timeouts reached
 - Client can't do anything while waiting unless using threads

October 18, 2022

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

L6.106

106

APPLICATION FLOW CONTROL - 2

- **Asynchronous web service**
- Client calls service
- Server responds to client with OK message
- Client closes connection
- Server performs the work associated with the service
- Server posts service result in an external data store
 - AWS: S3, SQS (queueing service), SNS (notification service)

October 18, 2022

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

L6.107

107

APPLICATION FLOW CONTROL - 3

Client flow control

(a)

Microservice as controller

(c)

AWS Step Function

(b)

Asynchronous

(d)

October 18, 2022

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

L6.108

108

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 18, 2022

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

L6.109

109

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

October 18, 2022

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

L6.110

110

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 18, 2022

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

L6.111

111

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 18, 2022

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

L6.112

112

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 18, 2022

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

L6.113

113

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**:
 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 18, 2022

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

L6.114

114

CPUSTEAL 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 18, 2022

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

L6.115

115

QUESTIONS



October 18, 2022

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

L6.116

116