

### CSS LIFE AFTER GRADUATION SEMINAR EXTRA CREDIT OPPORTUNITY

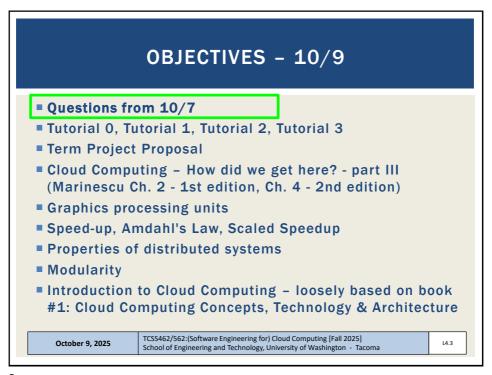
- When: Wednesday, October 9, 12:30pm-1:20pm
- Where:Milgard Hall (MLG) 110
- Is there life after graduating from CSS in SET?
- Yes! Dr. Donald Chinn and Andrew Fry will discuss the two main career paths after getting your bachelors degree: graduate school and industry.
- Whether you are a senior dreading the prospect of looking for a job or doing more school, or a junior who wonders what courses to take and how to get an internship, it is never too early (or late) to learn about what you can do now to prepare yourself for your life after graduation.
- This session is also open to CSS MS students
- Extra credit: 3 additional points in the daily feedback category
  - Enables 3 surveys to be missed, or for 23 out of 20 points for 2.3% grade bonus

October 9, 2025

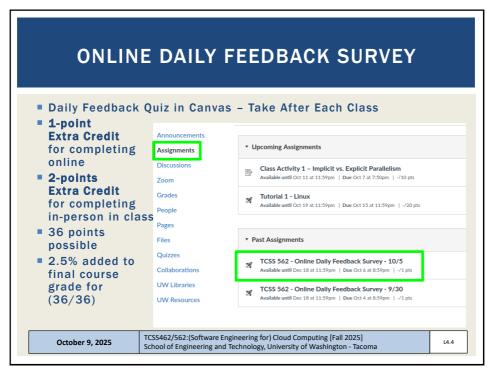
TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

4.2

2



3



4

### **WARNING**

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

October 9, 2025

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

L4.5

5

### MATERIAL / PACE

- Please classify your perspective on material covered in today's class (43 responses, 32 in-person, 11 online):
- 1-mostly review, 5-equal new/review, 10-mostly new
- Average 7.14 ( $\downarrow$  previous 7.24)
- Please rate the pace of today's class:
- 1-slow, 5-just right, 10-fast
- Average 5.09 ( $\downarrow$  previous 5.20)

October 9, 2025

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

L4.6

6

### FEEDBACK FROM 10/9

- Still confused on what makes specific ec2 c series instance good for specific avg or peak TLP
- For q9, can choose c7a over c5a? diff btwn c5a and c7a assuming same # of vCPUs?
  - Yes, in tutorial 3 we explain that "c5" refers to the 5<sup>th</sup> generation, while "c7" refers to the 7<sup>th</sup> generation of VMs
  - The higher the generation the new the CPU & hardware
  - Any generation of 5 or greater is considered still relevant
    - AWS Lambda currently use 5<sup>th</sup> and 6<sup>th</sup> generation CPUs from EC2
  - 4<sup>th</sup> generation or earlier is considered legacy

October 9, 2025

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

L4.7

7

### FEEDBACK - 2

- Homogeneous and heterogeneous systems: how do these systems differ from each other in terms of real life usage?
  - On the large public clouds, cloud services/platforms based on homogeneous hardware may be quite rare
  - A service or platform implemented using identical hardware should have less performance variability
  - AWS wants users to accept on average 10% performance variation:
    - if avg runtime is 60 sec, acceptable runtime is 54 to 66 sec.
- What is message passing?
- Is Karatsuba's algorithm (a more efficient multiplication algorithm) considered data-level parallelism?
  - No this algorithm is considered task parallelism as the operations performed on the data are not the same

October 9, 2025

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

L4.8

8

### ■ Is there a benefit to having a low arithmetic intensity? ■ A program having low arithmetic intensity has: ■ Low work (W) – few computations ■ High r/w memory traffic (Q) – lots of fetching ■ In modern computers, performing computations is significantly faster than fetching data from main memory (RAM). ■ Processors execute billions of calculations per second ■ Time to retrieve data from RAM is a bottleneck that can cause the CPU to sit idle. ■ Computers use a tiered memory hierarchy with multiple levels of extremely fast, small, and expensive cache memory located between the CPU and the slower, larger, and cheaper main memory.

9

October 9, 2025

# OBJECTIVES - 10/9 - Questions from 10/7 - Tutorial 0, Tutorial 1, Tutorial 2, Tutorial 3 - Term Project Proposal - Cloud Computing - How did we get here? - part III (Marinescu Ch. 2 - 1st edition, Ch. 4 - 2nd edition) - Graphics processing units - Speed-up, Amdahl's Law, Scaled Speedup - Properties of distributed systems - Modularity - Introduction to Cloud Computing - loosely based on book #1: Cloud Computing Concepts, Technology & Architecture

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

School of Engineering and Technology, University of Washington - Tacoma

10

### **DEMOGRAPHICS SURVEY**

Please complete the ONLINE demographics survey:

We have received 42 responses so far. >>> Random Drawing

- https://forms.gle/QNUW2hUV7fR7BDmv7
- Linked from course webpage in Canvas:
- http://faculty.washington.edu/wlloyd/courses/tcss562/ announcements.html

October 9, 2025

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

L4.11

11

### **AWS CLOUD CREDITS SURVEY**

Please complete the AWS Cloud Credits survey:

Please only complete survey after setting up AWS account or if requiring an IAM user (no-credit card option)

- https://forms.gle/Y4iWvBRFVLRPnPX37
- Linked from course webpage in Canvas:
- http://faculty.washington.edu/wlloyd/courses/tcss562/ announcements.html

October 9, 2025

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

L4.12

12

### **TUTORIAL 1** Introduction to Linux & the Command Line https://faculty.washington.edu/wlloyd/courses/tcss562/tutori als/TCSS462\_562\_f2025\_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 TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma L4.13 October 9, 2025

13

### **TUTORIAL 2** Introduction to Bash Scripting https://faculty.washington.edu/wlloyd/courses/tcss562/tutorials/T CSS462\_562\_f2025\_tutorial\_2.pdf Review tutorial sections: Create a BASH webservice client 1. What is a BASH script? 2. Variables 3. Input 4. Arithmetic 5. If Statements 6. Loops 7. Functions 8. User Interface Call service to obtain IP address & lat/long of computer Call service to obtain weather forecast for lat/long TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] October 9, 2025 L4.14 School of Engineering and Technology, University of Washington - Tacoma

14

### OBJECTIVES - 10/9 Questions from 10/7 Tutorial 0, Tutorial 1, Tutorial 2, Tutorial 3 Term Project Proposal Cloud Computing - How did we get here? - part III (Marinescu Ch. 2 - 1st edition, Ch. 4 - 2nd edition) Graphics processing units Speed-up, Amdahl's Law, Scaled Speedup Properties of distributed systems

15

Modularity

October 9, 2025

# OBJECTIVES - 10/9 Questions from 10/7 Tutorial 0, Tutorial 1, Tutorial 2, Tutorial 3 Term Project Proposal Cloud Computing - How did we get here? - part III (Marinescu Ch. 2 - 1st edition, Ch. 4 - 2nd edition) Graphics processing units Speed-up, Amdahl's Law, Scaled Speedup Properties of distributed systems Modularity Introduction to Cloud Computing - loosely based on book #1: Cloud Computing Concepts, Technology & Architecture October 9, 2025 TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

Introduction to Cloud Computing – loosely based on book
 #1: Cloud Computing Concepts, Technology & Architecture

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

School of Engineering and Technology, University of Washington - Tacoma

16

### OBJECTIVES - 10/9

- Questions from 10/7
- Tutorial 0, Tutorial 1, Tutorial 2, Tutorial 3
- Term Project Proposal
- Cloud Computing How did we get here? part III (Marinescu Ch. 2 - 1st edition, Ch. 4 - 2nd edition)
- Graphics processing units
- Speed-up, Amdahl's Law, Scaled Speedup
- Properties of distributed systems
- Modularity
- Introduction to Cloud Computing loosely based on book #1: Cloud Computing Concepts, Technology & Architecture

October 9, 2025

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

L4.17

17

### **ARITHMETIC INTENSITY**

Arithmetic intensity: Ratio of work (W) to memory traffic r/w (Q)  $I = \frac{W}{Q}$ 

Example: # of floating point ops per byte of data read

- Characterizes application scalability with SIMD support
  - SIMD can perform many fast matrix operations in parallel
- High arithmetic Intensity:

**Programs** with dense matrix operations scale up nicely (many calcs vs memory RW, supports lots of parallelism)

Low arithmetic intensity:

Programs with sparse matrix operations do not scale well with problem size

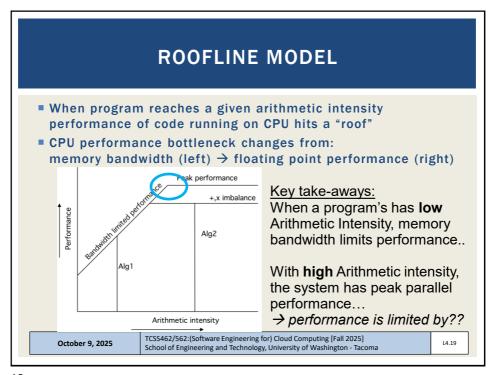
(memory RW becomes bottleneck, not enough ops!)

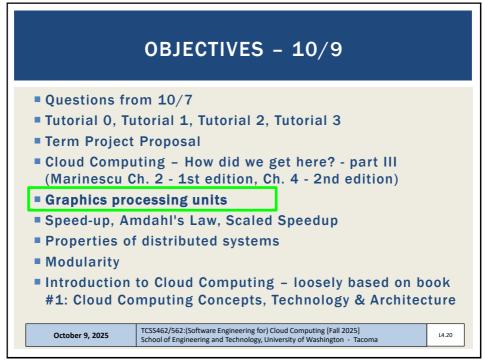
October 9, 2025

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

4.18

18





20

### **GRAPHICAL PROCESSING UNITS (GPUs)**

- GPU provides multiple SIMD processors
- Typically 7 to 15 SIMD processors each
- 32,768 total registers, divided into 16 lanes (2048 registers each)
- GPU programming model: single instruction, multiple thread
- Programmed using CUDA- C like programming language by NVIDIA for GPUs
- CUDA threads single thread associated with each data element (e.g. vector or matrix)
- Thousands of threads run concurrently

October 9, 2025

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

L4.21

21

### **OBJECTIVES - 10/9**

- Questions from 10/7
- Tutorial 0, Tutorial 1, Tutorial 2, Tutorial 3
- Term Project Proposal
- Cloud Computing How did we get here? part III (Marinescu Ch. 2 - 1st edition, Ch. 4 - 2nd edition)
- Graphics processing units
- Speed-up, Amdahl's Law, Scaled Speedup
- Properties of distributed systems
- Modularity
- Introduction to Cloud Computing loosely based on book
   #1: Cloud Computing Concepts, Technology & Architecture

October 9, 2025

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

1.22

22

### PARALLEL COMPUTING

- Parallel hardware and software systems allow:
  - Solve problems demanding resources not available on single system.
  - Reduce time required to obtain solution
- The speed-up (S) measures effectiveness of parallelization:

$$S(N) = T(1) / T(N)$$

 $T(1) \rightarrow$  execution time of total sequential computation

T(N) → execution time for performing N parallel computations in parallel

October 9, 2025

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

L4.23

23

### **SPEED-UP EXAMPLE**

- Consider embarrassingly parallel image processing
- Eight images (multiple data)
- Apply image transformation (greyscale) in parallel
- 8-core CPU, 16 hyperthreads
- Sequential processing: perform transformations one at a time using a single program thread
  - 8 images, 3 seconds each: T(1) = 24 seconds
- Parallel processing
  - 8 images, 3 seconds each: T(N) = 3 seconds
- Speedup: S(N) = 24 / 3 = 8x speedup
- Called "perfect scaling"
- Must consider data transfer and computation setup time

October 9, 2025

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

L4.24

24

### **AMDAHL'S LAW**

- Amdahl's law is used to estimate the speed-up of a job using parallel computing
- 1. Divide job into two parts
- 2. Part A that will still be sequential
- 3. Part B that will be sped-up with parallel computing
- Portion of computation which cannot be parallelized will determine (i.e. limit) the overall speedup
- Amdahl's law assumes jobs are of a fixed size
- Also, Amdahl's assumes no overhead for distributing the work, and a perfectly even work distribution

October 9, 2025

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

L4.25

25

### AMDAHL'S LAW

Speed-up formula
→

$$S = \frac{1}{(1-f) + \frac{f}{N}}$$

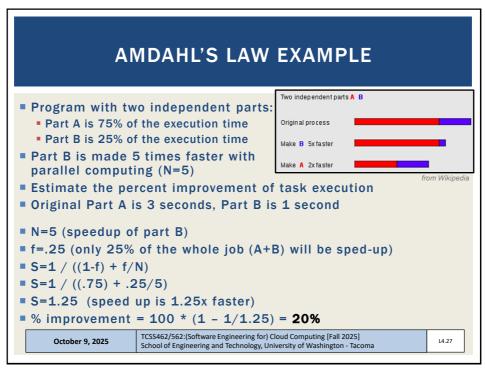
- S = theoretical speedup of the whole task
- f= fraction of work that is parallel (ex. 25% or 0.25)
- N= proposed speed up of the parallel part (ex. 5 times speedup)
- % improvement of task execution = 100 \* (1 - (1 / S))
- Using Amdahl's law, we can find the maximum possible speed-up (S) for a given scenario (e.g. ~8x) ...

October 9, 2025

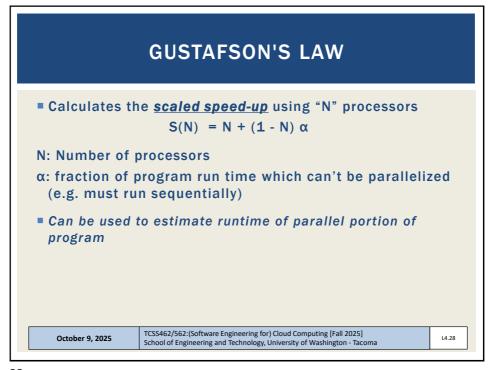
TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L4.26

26



27



28

### **GUSTAFSON'S LAW**

■ Calculates the <u>scaled speed-up</u> using "N" processors

 $S(N) = N + (1 - N) \alpha$ 

N: Number of processors

- α: fraction of program run time which can't be parallelized (e.g. must run sequentially)
- Can be used to estimate runtime of parallel portion of program
- Where  $\alpha = \sigma / (\pi + \sigma)$
- Where  $\sigma$ = sequential time,  $\pi$  =parallel time
- Our Amdahl's example:  $\sigma$ = 3s,  $\pi$  =1s,  $\alpha$  =.75

October 9, 2025

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

L4.29

29

### **GUSTAFSON'S LAW**

Calculates the <u>scaled speed-up</u> using "N" processors

$$S(N) = N + (1 - N) \alpha$$

N: Number of processors

α: fraction of program run time which can't be parallelized (e.g. must run sequentially)

**Example:** 

Consider a program that is embarrassingly parallel, but 75% cannot be parallelized.  $\alpha$ =.75

QUESTION: If deploying the job on a 2-core CPU, what scaled speedup is possible assuming the use of two processes that run in parallel?

October 9, 2025

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

L4.30

30

### GUSTAFSON'S EXAMPLE QUESTION: What is the maximum theoretical speed-up on a 2-core CPU? $S(N) = N + (1 - N) \alpha$ $N=2, \alpha=.75$ S(N) = 2 + (1 - 2) .75 S(N) = ?What is the maximum theoretical speed-up on a 16-core CPU? $S(N) = N + (1 - N) \alpha$ $N=16, \alpha=.75$ S(N) = 16 + (1 - 16) .75 S(N) = ?

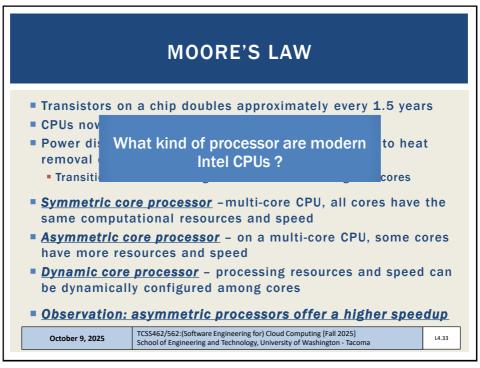
TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

31

October 9, 2025

### **GUSTAFSON'S EXAMPLE** What is the maximum theoretical speed-up on a 2-core CPU? $S(N) = N + (1 - N) \alpha$ $N=2, \alpha=$ For 2 CPUs, speed up is 1.25x S(N) =S(N) = ?For 16 CPUs, speed up is 4.75x ■ What is the maximum theoretical speed-up on a 16-core CPU? $S(N) = N + (1 - N) \alpha$ $N=16, \alpha=.75$ S(N) = 16 + (1 - 16).75S(N) = ?TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma October 9, 2025 L4.32

32



33

# OBJECTIVES - 10/9 Questions from 10/7 Tutorial 0, Tutorial 1, Tutorial 2, Tutorial 3 Term Project Proposal Cloud Computing - How did we get here? - part III (Marinescu Ch. 2 - 1st edition, Ch. 4 - 2nd edition) Graphics processing units Speed-up, Amdahl's Law, Scaled Speedup Properties of distributed systems Modularity Introduction to Cloud Computing - loosely based on book #1: Cloud Computing Concepts, Technology & Architecture

34

### **DISTRIBUTED SYSTEMS**

- Collection of autonomous computers, connected through a network with distribution software called "middleware" that enables coordination of activities and sharing of resources
- Key characteristics:
- Users perceive system as a single, integrated computing facility.
- Compute nodes are autonomous
- Scheduling, resource management, and security implemented by every node
- Multiple points of control and failure
- Nodes may not be accessible at all times
- System can be scaled by adding additional nodes
- Availability at low levels of HW/software/network reliability

October 9, 2025

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

L4.35

35

### **DISTRIBUTED SYSTEMS - 2**

- Key non-functional attributes
  - Known as "ilities" in software engineering
- Availability 24/7 access?
- Reliability Fault tolerance
- Accessibility reachable?
- Usability user friendly
- Understandability can under
- Scalability responds to variable demand
- Extensibility can be easily modified, extended
- Maintainability can be easily fixed
- Consistency data is replicated correctly in timely manner

October 9, 2025

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

L4.36

36

### TRANSPARENCY PROPERTIES OF DISTRIBUTED SYSTEMS

- Access transparency: local and remote objects accessed using identical operations
- Location transparency: objects accessed w/o knowledge of their location.
- Concurrency transparency: several processes run concurrently using shared objects w/o interference among them
- Replication transparency: multiple instances of objects are used to increase reliability
  - users are unaware if and how the system is replicated
- Failure transparency: concealment of faults
- Migration transparency: objects are moved w/o affecting operations performed on them
- Performance transparency: system can be reconfigured based on load and quality of service requirements
- Scaling transparency: system and applications can scale w/o change in system structure and w/o affecting applications

October 9, 2025

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

L4.37

37

### **OBJECTIVES - 10/9**

- Questions from 10/7
- Tutorial 0, Tutorial 1, Tutorial 2, Tutorial 3
- Term Project Proposal
- Cloud Computing How did we get here? part III (Marinescu Ch. 2 - 1st edition, Ch. 4 - 2nd edition)
- Graphics processing units
- Speed-up, Amdahl's Law, Scaled Speedup
- Properties of distributed systems
- Modularity
- Introduction to Cloud Computing loosely based on book
   #1: Cloud Computing Concepts, Technology & Architecture

October 9, 2025

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

4.38

38

### 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)
- 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 9, 2025

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

L4.39

39

### CLOUD COMPUTING - HOW DID WE GET HERE? PART III 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 9, 2025

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

L4.40

40

### CLOUD COMPUTING - HOW DID WE GET HERE? PART III 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 9, 2025

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

L4.41

41

### CLOUD COMPUTING - HOW DID WE GET HERE? PART III SUMMARY OF KEY POINTS - 3

- Speed-up (S)
  S(N) = T(1) / T(N)
- Amdahl's law:

 $S = 1/\alpha$ 

 $\alpha$  = percent of program that must be sequential

Scaled speedup with N processes:

 $S(N) = N - \alpha(N-1)$ 

Moore's Law

October 9, 2025

- Symmetric core, Asymmetric core, Dynamic core CPU
- Distributed Systems Non-function quality attributes
- Distributed Systems Types of Transparency
- Types of modularity- Soft, Enforced

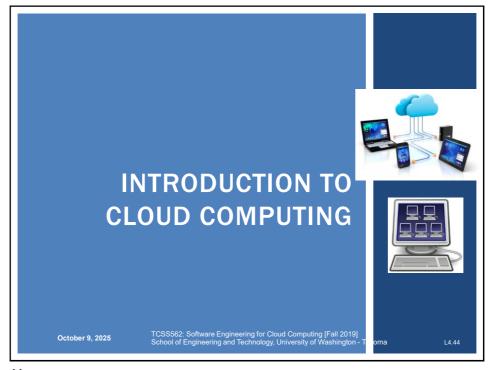
TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L4.42

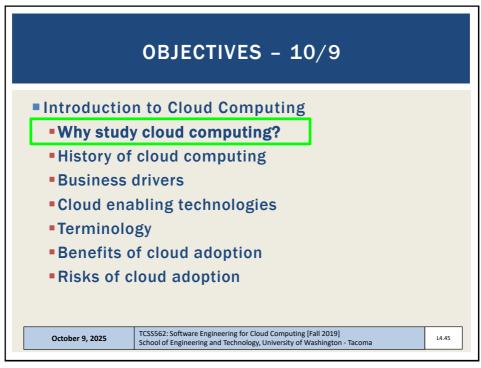
42

# OBJECTIVES - 10/9 Questions from 10/7 Tutorial 0, Tutorial 1, Tutorial 2, Tutorial 3 Term Project Proposal Cloud Computing - How did we get here? - part III (Marinescu Ch. 2 - 1st edition, Ch. 4 - 2nd edition) Graphics processing units Speed-up, Amdahl's Law, Scaled Speedup Properties of distributed systems Modularity Introduction to Cloud Computing - loosely based on book #1: Cloud Computing Concepts, Technology & Architecture

43

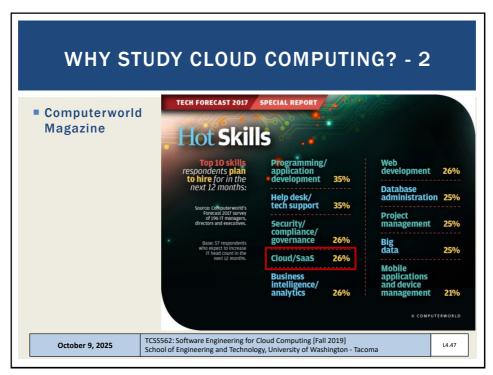


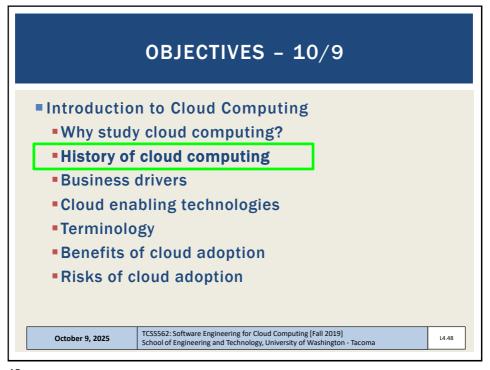
44



# 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/

46





48

### 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 9, 2025

TCSS562: Software Engineering for Cloud Computing [Fall 2019] School of Engineering and Technology, University of Washington - Tacoma

L4.49

49

### **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 9, 2025

TCSS562: Software Engineering for Cloud Computing [Fall 2019] School of Engineering and Technology, University of Washington - Tacoma

L4.50

50

### 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 (laaS)
  - Amazon launches Elastic Compute Cloud (EC2) service
  - Organization can "lease" computing capacity and processing power to host enterprise applications
  - Infrastructure

October 9, 2025

TCSS562: Software Engineering for Cloud Computing [Fall 2019] School of Engineering and Technology, University of Washington - Tacoma

L4.51

51

### **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 googlehosted datacenters

October 9, 2025

TCSS562: Software Engineering for Cloud Computing [Fall 2019] School of Engineering and Technology, University of Washington - Tacoma

L4.52

52

### 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 9, 2025

TCSS562: Software Engineering for Cloud Computing [Fall 2019] School of Engineering and Technology, University of Washington - Tacoma

L4.53

53

### MORE CONCISE DEFINITION

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

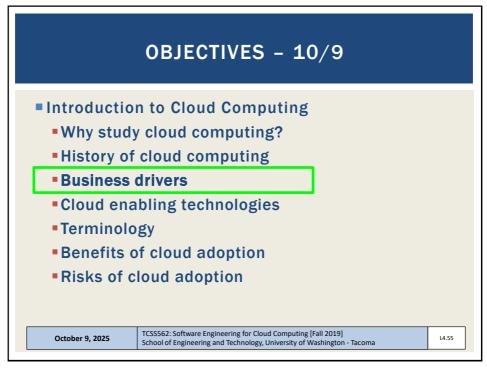
From Cloud Computing Concepts, Technology, and Architecture Z. Mahmood, R. Puttini, Prentice Hall, 5<sup>th</sup> printing, 2015

October 9, 2025

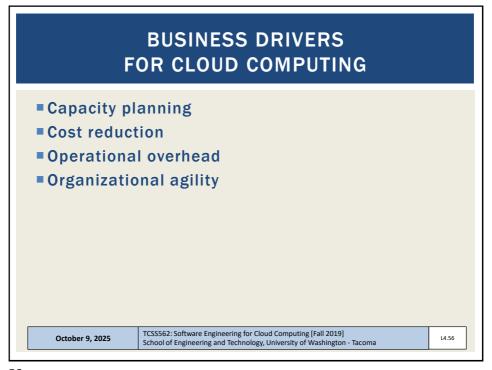
TCSS562: Software Engineering for Cloud Computing [Fall 2019] School of Engineering and Technology, University of Washington - Tacoma

L4.54

54



55



56

### 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 9, 2025

TCSS562: Software Engineering for Cloud Computing [Fall 2019] School of Engineering and Technology, University of Washington - Tacoma

L4.57

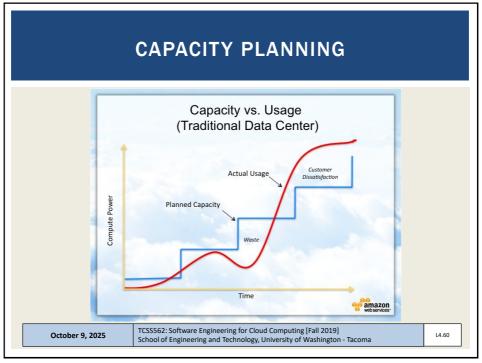
57



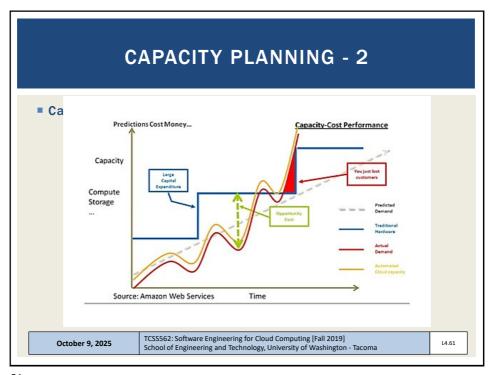
58

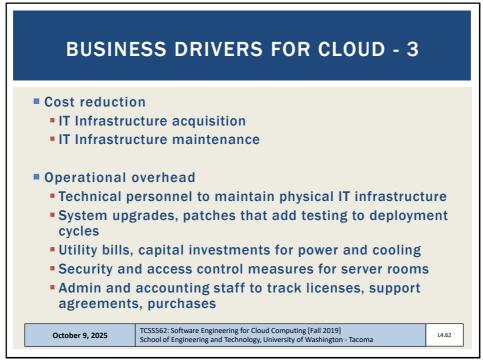
### **BUSINESS DRIVERS FOR CLOUD - 2** Capacity planning Over-provisioning: is costly due to too much infrastructure Under-provisioning: is costly due to potential for business loss from poor quality of service Capacity planning strategies Lead strategy: add capacity in anticipation of demand (preprovisioning) 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 flucations TCSS562: Software Engineering for Cloud Computing [Fall 2019] October 9, 2025 L4.59 School of Engineering and Technology, University of Washington - Tacoma

59



60

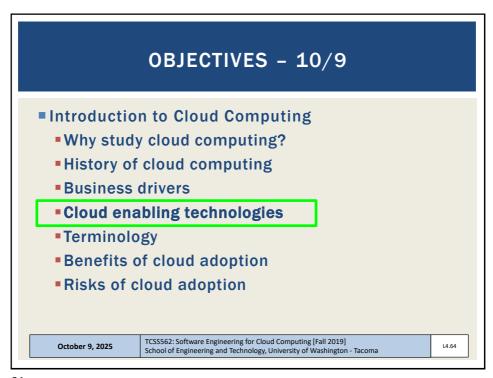




62

# 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 Ctober 9, 2025 | TCSSS62: Software Engineering for Cloud Computing [Fall 2019] | School of Engineering and Technology, University of Washington - Tacoma

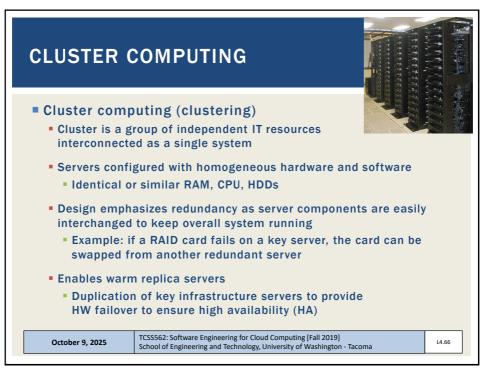
63



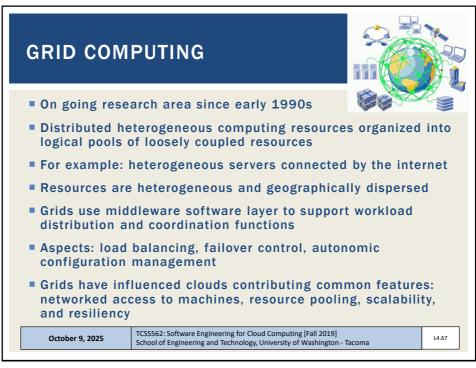
64

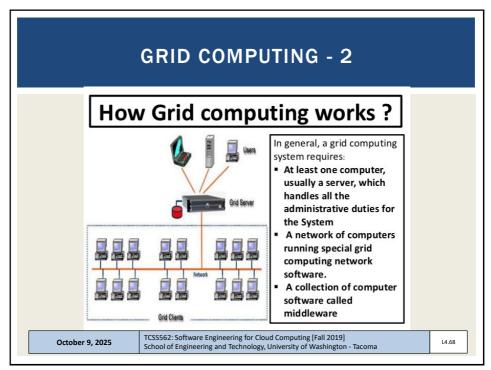
## TECHNOLOGY INNOVATIONS LEADING TO CLOUD Cluster computing Grid computing Virtualization October 9, 2025 TCSSS62: Software Engineering for Cloud Computing [Fall 2019] School of Engineering and Technology, University of Washington - Tacoma

65

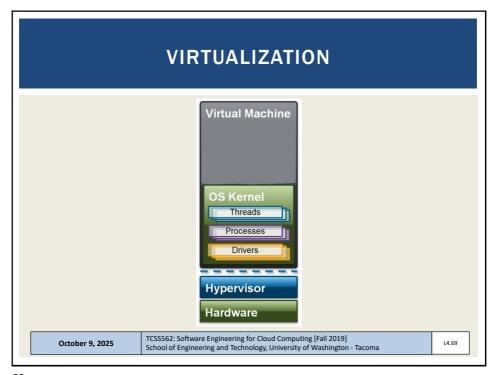


66

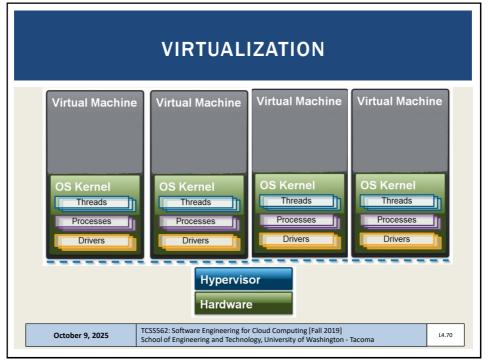




68



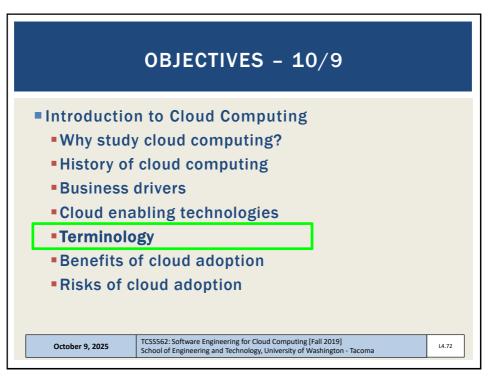
69



70

# 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 TCSSS62: Software Engineering for Cloud Computing [Fall 2019] School of Engineering and Technology, University of Washington - Tacoma

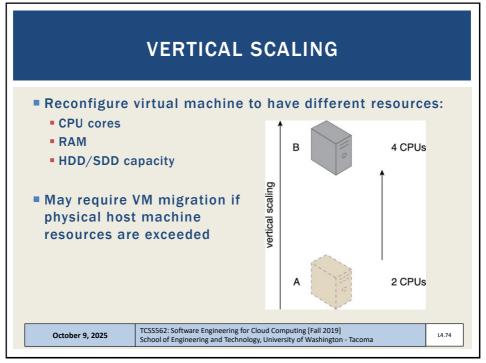
71



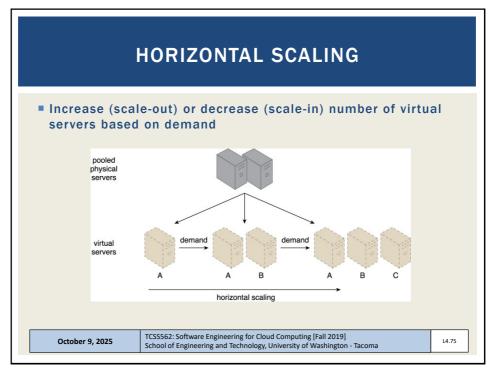
72

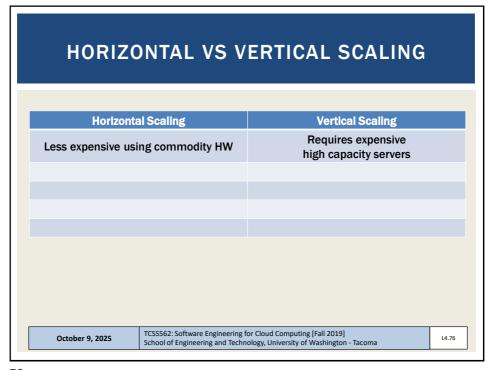
### **KEY TERMINOLOGY** ■ On-Premise Infrastructure Local server infrastructure not configured as a cloud Cloud Provider Corporation or private organization responsible for maintaining cloud Cloud Consumer User of cloud services Scaling Vertical scaling Scale up: increase resources of a single virtual 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 TCSS562: Software Engineering for Cloud Computing [Fall 2019] October 9, 2025 School of Engineering and Technology, University of Washington - Tacoma

73

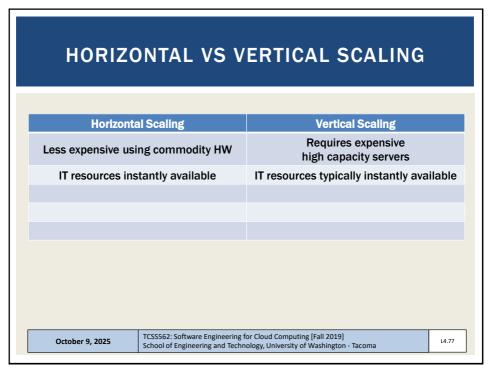


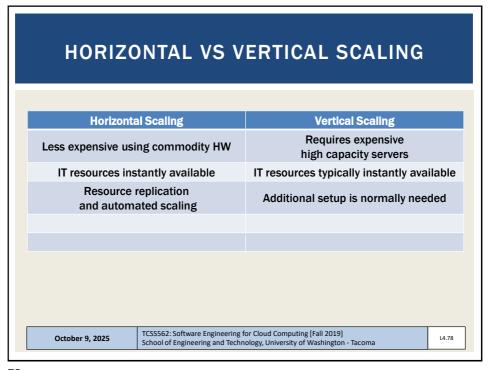
74



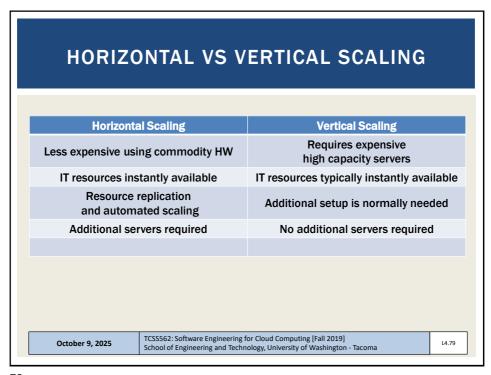


76





78

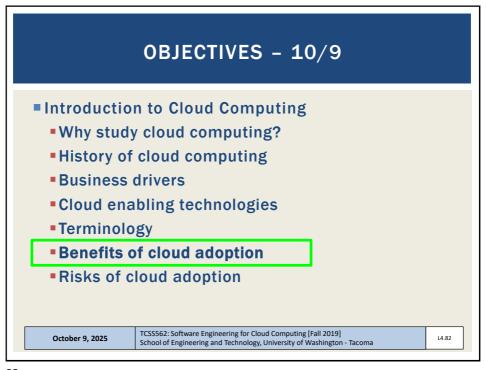


HORIZONTAL VS VERTICAL SCALING	
Horizontal Scaling	Vertical Scaling
Less expensive using commodity HW	Requires expensive high capacity servers
IT resources instantly available	IT resources typically instantly available
Resource replication and automated scaling	Additional setup is normally needed
Additional servers required	No additional servers required
Not limited by individual server capacity	Limited by individual server capacity
October 9, 2025 TCSS562: Software Engineering School of Engineering and Technology	for Cloud Computing [Fall 2019] hology, University of Washington - Tacoma

80

# KEY TERMINOLOGY - 2 Cloud services Broad array of resources accessible "as-a-service" Categorized as Infrastructure (laaS), Platform (PaaS), Software (SaaS) Service-level-agreements (SLAs): Establish expectations for: uptime, security, availability, reliability, and performance October 9, 2025 TCSSS62: Software Engineering for Cloud Computing [Fall 2019] School of Engineering and Technology, University of Washington - Tacoma

81



82

### **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 9, 2025

TCSS562: Software Engineering for Cloud Computing [Fall 2019] School of Engineering and Technology, University of Washington - Tacoma

L4.83

83

### **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...

October 9, 2025

TCSS562: Software Engineering for Cloud Computing [Fall 2019] School of Engineering and Technology, University of Washington - Tacoma

L4.84

84

### **CLOUD BENEFITS - 3**

- Example: Using 100 servers for 1 hour costs the same as using 1 server for 100 hours
- Rosetta Protein Folding: Working with a UW-Tacoma graduate student, we recently 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 9, 2025

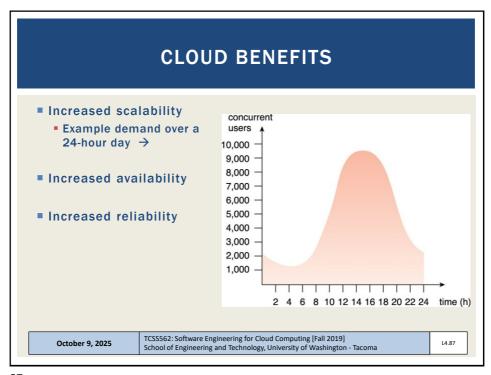
TCSS562: Software Engineering for Cloud Computing [Fall 2019] School of Engineering and Technology, University of Washington - Tacoma

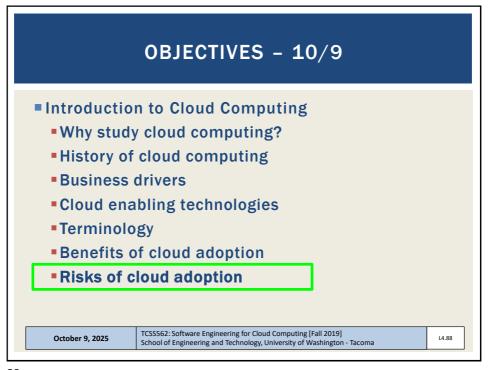
L4.85

85



86

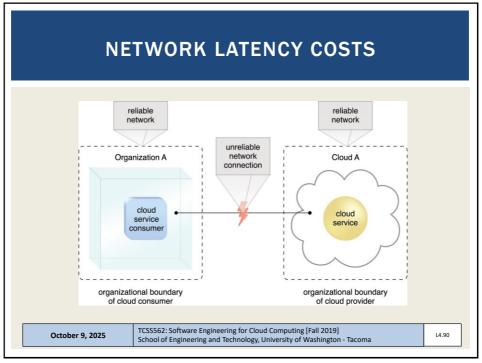




88

### **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 TCSS562: Software Engineering for Cloud Computing [Fall 2019] October 9, 2025 L4.89 School of Engineering and Technology, University of Washington - Tacoma

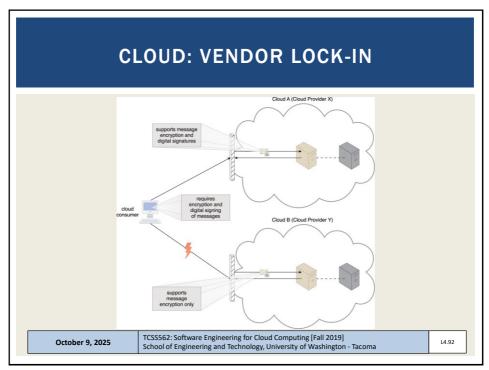
89



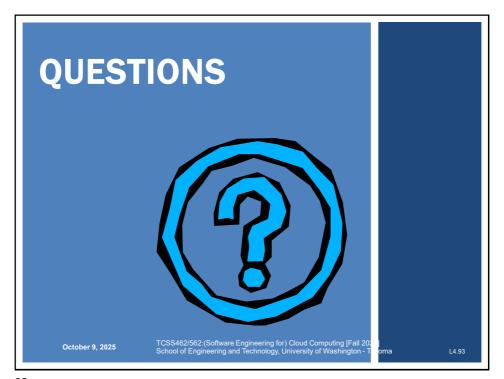
90

### **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 TCSS562: Software Engineering for Cloud Computing [Fall 2019] October 9, 2025 School of Engineering and Technology, University of Washington - Tacoma

91



92



93