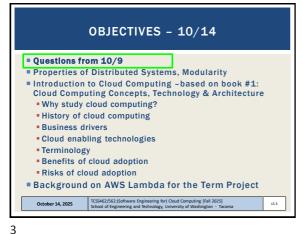


JOB RECRUITER IN CLASS - TODAY
(FIRST 10 MINUTES)

Tuesday October 14, 3:40pm
Fast Enterprises, LLC
Hires consultants that work with state and local governments, and also some federal agencies on an international scale.

Positions open to domestic and international graduates
They do not sponsor H1B visas
Professor has asked whether they can hire an OPT student on F-1

1



Thursdays:
6:00 to 7:00 pm - CP 229 and Zoom
Fridays
11:00 am to 12:00 pm - ONLINE via Zoom*
Or email for appointment

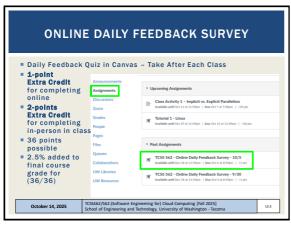
Office Hours set based on Student Demographics survey feedback

**- Friday office hours may be adjusted or canceled due meeting conflicts or other obligations. Adjustments will be announced via Canvas.

Cotober 14, 2025

TESS42/552/Software Engineering foil Cloud Computing [fail 2025]
School of Engineering and Technology University of Washington - Tacoma

3



5

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 *

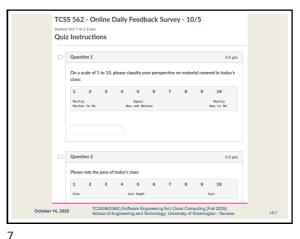
WE WILL NOT BE ABLE TO DUPLICATE CHECK SURVEYS FOR EACH CLASS SESSION AND MAKE CORRECTIONS

Cotober 14, 2025

CCODER 14, 2025

TCCS462/562: (Software Engineering for) Good Computing (fall 2025)
School of Engineering and Technology, University of Washington: Tacoma

Slides by Wes J. Lloyd L5.1



MATERIAL / PACE Please classify your perspective on material covered in today's class (48 respondents, 35 in-person, 13 online): 1-mostly review, 5-equal new/review, 10-mostly new - Average - 7.33 (1 - previous 7.14) ■ Please rate the pace of today's class: ■ 1-slow. 5-just right. 10-fast Average - 4.96 (↓ - previous 5.09) October 14, 2025 L5.8

FEEDBACK FROM 10/9 >Term project themes: If choosing LLM comparison or programming language comparison, just need to compare performance? Yes, but multiple metrics should be used to quantify performance: Average turnaround time, average runtime, average throughput, cost >If building the standard or an alternative application, need to select which design trade-offs to compare and to compare performance? • Two course themes for "design trade-off" comparison, choose: 1. LLM comparison: implement same app, with 2 or more LLMs 2. Prog Lang comparison: implement same app, in 2 or more LLMs Groups are free to compare other design trade-offs, but others are not the "standard" course themes for Fall 2025 TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tac October 14, 2025

FEEDBACK - 2 >What makes Amdahi's Law different from Gustafson's Law? Amdahl's law was conceived to consider the speed-up for a fixed problem size The amount of work is fixed, and some part is parallel, the rest sequential Gustafson's Law considers how much larger of a problem can be solved in the same amount of time when having access to more processors · Considers data parallelism, and the ability to process more data in less time with more compute resources The speed-up can grow linearly with the number of processors because the parallel portion (work) grows with the workload and dominates the run time KEY: Use Amdahi's for estimating speedup of a fixed amount of work, Gustafson for speedup of scalable amount of work TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacr L5.10

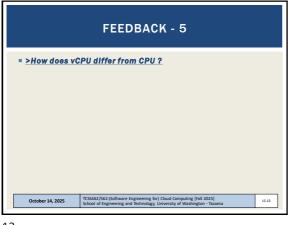
9

FEEDBACK - 3 Consequently, if computing the speed-up for a job where 50% is parallelizable, Amdahl's and Gustafson's disagree on speedup - Amdahl's Law (solid line) Amdahl's vs. Gustafson's Law (P = 50% r levels off near a max speedup of about 2×, showing diminishing returns once ½ the workload remains sequential. Gustafson's Law (dashed line) scales almost linearly up to about 64× with 128 processors, reflecting that if the workload grows the parallel half can dominate the runtime TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tac October 14, 2025 L5.11

FEEDBACK - 4 >What is the grading rubric for the term project? The full body of work produced is considered. The final presentation or term paper is graded. Code and other artifacts can be considered Points are assigned to each of the components: . Description of the case study Description of the application Description of the experiments (*) Project results presentation including graphs, tables, statistics, Results analysis and discussion in paper/presentation Results conclusions in paper/presentation Formatting 4.0 = Projects that produce high quality reports that clearly communicate the implications and key outcomes of design trade-offs of their respective case studies TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tace October 14, 2025

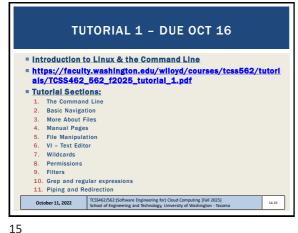
11 12

Slides by Wes J. Lloyd L5.2



TUTORIAL 0 Getting Started with AWS https://faculty.washington.edu/wlloyd/courses/tcss562/tutori als/TCSS462_562_f2025_tutorial_0.pdf ■ Create an AWS account Create account credentials for working with the CLI Install awsconfig package Setup awsconfig for working with the AWS CLI October 14, 2025 L5.14

13 14



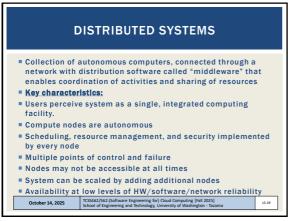
TUTORIAL 2 - DUE OCT 21 Introduction to Bash Scripting https://faculty.washington.edu/wlloyd/courses/tcss562/tutorials/TCSS462_562_f2025_tutorial_2.pdf Review tutorial sections: Create a BASH webservice client 1. What is a BASH script? 2. Variables Input
 Arithmetic If Statements Loops Functions 8. User Interface Call service to obtain IP address & lat/long of computer Call weatherbit.io API to obtain weather forecast for lat/long October 11, 2022 L4.16

TUTORIAL 3 - DUE OCT 30 (TEAMS OF 2) ■ Best Practices for Working with Virtual Machines on Amazon EC2 https://faculty.washington.edu/wlloyd/courses/tcss562/tutori als/TCSS462_562_f2025_tutorial_3.pdf ■ Creating a spot VM Creating an image from a running VM ■ Persistent spot request Stopping (pausing) VMs ■ EBS volume types Ephemeral disks (local disks) ■ Mounting and formatting a disk ■ Disk performance testing with Bonnie++ ■ Cost Saving Best Practices TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tac October 14, 2025 L5.17

OBJECTIVES - 10/14 Questions from 10/9 ■ Properties of Distributed Systems Modularity ■ Introduction to Cloud Computing – based on book #1: 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 Background on AWS Lambda for the Term Project TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Ta October 14, 2025

17 18

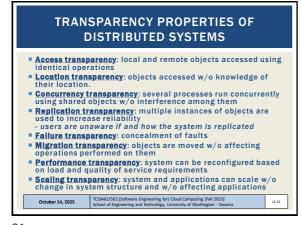
Slides by Wes J. Lloyd L5.3



| 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

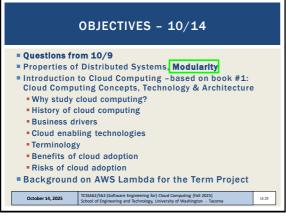
| October 14, 2025 | TCSS462/562; Software Engineering for Cloud Computing [Fall 2025] | School of Engineering and Technology, University of Washington - Tacoma | 15.20

19 20



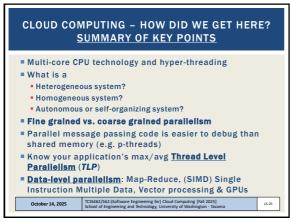
WE WILL RETURN AT 5:02PM

21



23 24

Slides by Wes J. Lloyd L5.4



CLOUD COMPUTING - HOW DID WE GET HERE?

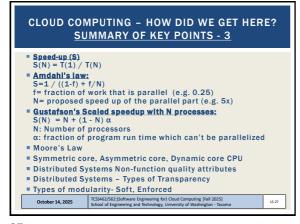
SUMMARY OF KEY POINTS - 2

Bit-level parallelism
Instruction-level parallelism (CPU pipelining)
Fiynn'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 14, 2025

TCSS42/S62/Software Engineering for) Cloud Computing [Fall 2025]
School of Engineering and Technology, University of Washington - Taxoma

25



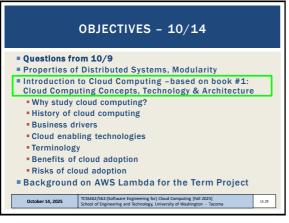
INTRODUCTION TO CLOUD COMPUTING

October 14, 2025

TCSS4C2/SC2 (Software Engineering for) Cloud Computing Fiel 202
School of Engineering and Technology, University of Westlengton Tomas

LS 28

27



OBJECTIVES - 10/14

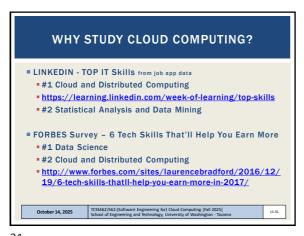
- Questions from 10/9
- Properties of Distributed Systems, Modularity
- Introduction to Cloud Computing - based on book #1:
Cloud Computing Concepts, Technology & Architecture
- Why study cloud computing?
- History of cloud computing
- Business drivers
- Cloud enabling technologies
- Terminology
- Benefits of cloud adoption
- Risks of cloud adoption
- Risks of cloud adoption
- Background on AWS Lambda for the Term Project

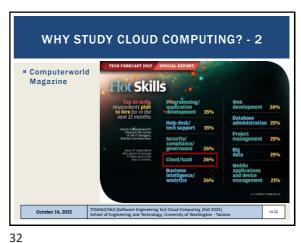
October 14, 2025

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

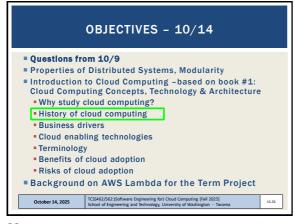
29 30

Slides by Wes J. Lloyd L5.5





31 3



A BRIEF HISTORY OF CLOUD COMPUTING

John McCarthy, 1961
 Turing award winner for contributions to Al

"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 14, 2025

| TCSS482/SG2/Software Engineering for) Cloud Computing [Fall 2025] | School of Engineering and Technology, University of Wisblington -Tacoma

| 13.34

33



CLOUD HISTORY: SERVICES - 1

Late 1990s - Early Software-as-a-Service (SaaS)
Salesforce: Remotely provisioned services for the enterprise

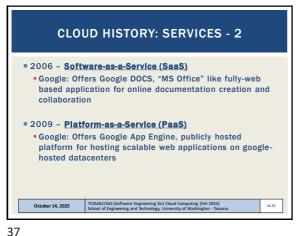
2002 Amazon Web Services (AWS) platform: Enterprise oriented services for remotely provisioned storage, computing resources, and business functionality

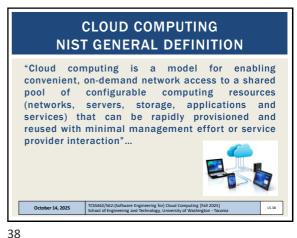
2006 - Infrastructure-as-a-Service (laaS)
Amazon launches Elastic Compute Cloud (EC2) service
Organization can "lease" computing capacity and processing power to host enterprise applications
Infrastructure

(Ctober 14, 2025 | Ctober 14, 202

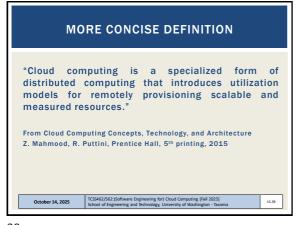
35 36

Slides by Wes J. Lloyd L5.6





7



OBJECTIVES - 10/14

• Questions from 10/9

• Properties of Distributed Systems, Modularity

• Introduction to Cloud Computing - based on book #1:
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

• Risks of cloud adoption

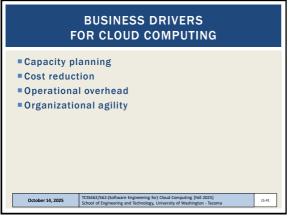
• Background on AWS Lambda for the Term Project

October 14, 2025

TESSAG/ISC2 (Software Engineering for) Cloud Computing [Fall 2025]

School of Engineering and Technology University of Washington - Tacoma

39



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 14, 2025

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

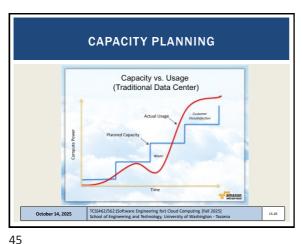
41 42

Slides by Wes J. Lloyd L5.7



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 Load prediction Capacity planning helps anticipate demand flucations L5.44

44



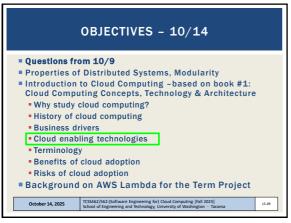
CAPACITY PLANNING - 2 ■ Ca October 14, 2025

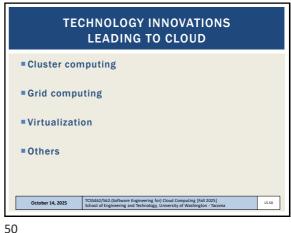
BUSINESS DRIVERS FOR CLOUD - 3 Cost reduction • IT Infrastructure acquisition IT Infrastructure maintenance Operational overhead Technical personnel to maintain physical IT infrastructure System upgrades, patches that add testing to deployment cycles Utility bills, capital investments for power and cooling Security and access control measures for server rooms Admin and accounting staff to track licenses, support agreements, purchases TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tac October 14, 2025

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 14, 2025

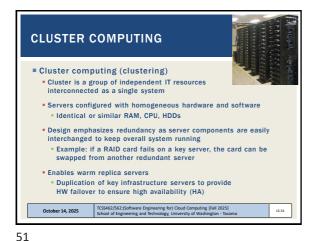
47 48

Slides by Wes J. Lloyd L5.8



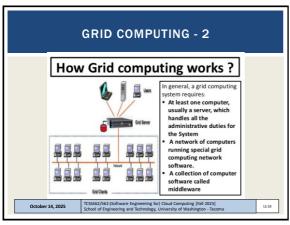


49





31



VIRTUALIZATION

Virtual Machine

OS Kernel

Threads

Processes

Drivers

Hypervisor

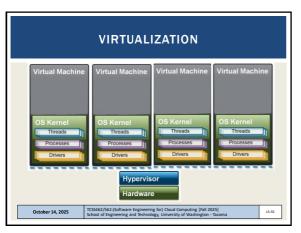
Hardware

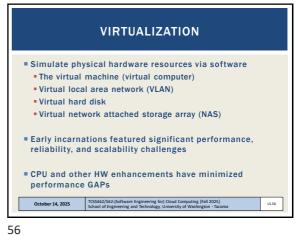
Actober 14, 2025

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

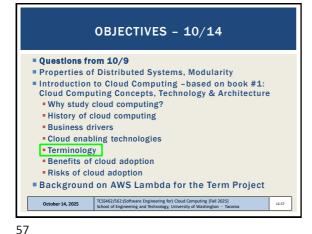
53 54

Slides by Wes J. Lloyd L5.9



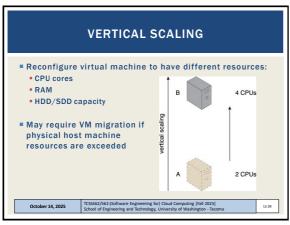


55



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 14, 2025 L5.58

0/



HORIZONTAL SCALING

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

pooled physical servers

virtual servers

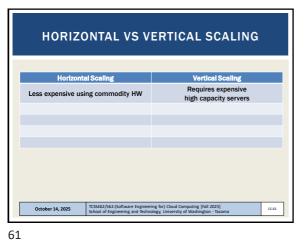
Virtual Servers

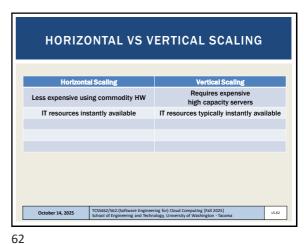
Virtual Servers

TCSS462/562:[software Engineering for) Cloud Computing [fall 2025]
School of Engineering and Technology, University of Visinington - Tacoma

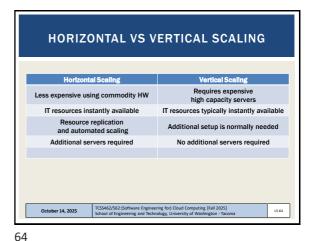
59 60

Slides by Wes J. Lloyd L5.10



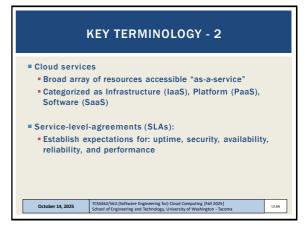


HORIZONTAL VS VERTICAL SCALING		
Horizontal Scaling	Vertical Scaling	
Less expensive using commodity HW	Requires expensive high capacity servers	
IT resources instantly available	IT resources typically instantly available	
Resource replication and automated scaling	Additional setup is normally needed	
	TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma	



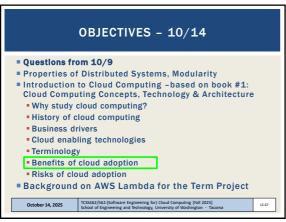
63

Horizontal S	caling	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
		ring for) Cloud Computing [Fall 2025] ology, University of Washington - Tacoma



65 66

Slides by Wes J. Lloyd L5.11



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 14, 2025

Cotober 14, 2025

67 68



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)

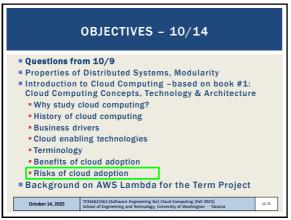
69



CLOUD BENEFITS Increased scalability Example demand over a 24-hour day → 10.000 9.000 8,000 Increased availability 7,000 6,000 5.000 ■ Increased reliability 4,000 3.000 2,000 4 6 8 10 12 14 16 18 20 22 24 time (h) October 14, 2025

71 72

Slides by Wes J. Lloyd L5.12



CLOUD ADOPTION RISKS

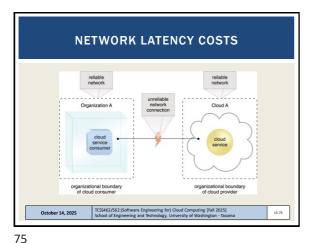
Increased security vulnerabilities
Expansion of trust boundaries now include the external cloud
Security responsibility shared with cloud provider

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 14, 2025

ICCS462/562/Software Engineering for/ Cloud Computing [Tail 2025]
School of Engineering and Technology, University of Wischington - Tacoma

73



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

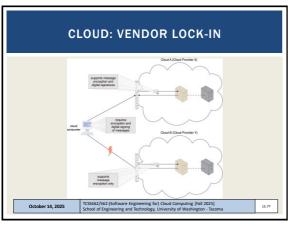
October 14, 2025

Cotober 14, 2025

Cotober 14, 2025

Listo of Engineering and Technology University of Washington - Tacoma

73



OBJECTIVES - 10/14

- Questions from 10/9
- Properties of Distributed Systems, Modularity
- Introduction to Cloud Computing - based on book #1:
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
- Risks of cloud adoption
- Background on AWS Lambda for the Term Project

October 14, 2025

| TCSS462/562/Software Engineering forly Cloud Computing [Fail 2025]
| School of Engineering and Technology, University of Washington - Tacoma

77 78

Slides by Wes J. Lloyd L5.13

74



SERVERLESS – KEY CONCEPTS

**Function-as-a-Service (FaaS) platform

• A platform where developers deploy "functions" written in common languages (e.g. Java, Python, Go, Node.js) that run as microservices

• AWS Lambda is a FaaS platform

• We will discuss platform limitations

**Function Instances*

• This is an instantiation of a running function

• A function instance is created when a client (user) calls the serverless function

• Each concurrent (parallel) call to AWS Lambda to the same function will create a unique function instance to handle the request

• The default maximum number of concurrently running function instances in your account is 10.

• The default was originally 1,000 when the platform was introduced, and was dropped to 100, then 50, and is now just 10 in response to the growing popularity of AWS Lambda (they are running out of servers??)

• You will want to request an increase in your AWS account's default concurrency. A minimum of 100 is recommended

79 80

TYPES OF FUNCTION CALLS:
SYNCHRONOUS

Serverless Computing:

AWS Lambda supports synchronous and asynchronous function calls
Clients typically orchestrate synchronous calls and pipelines
Asynchronous calls are often made via events

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 14, 2025
School of Engineering Into Cound Computing [Fall 2025]
School of Engineering and Technology, University of Washington - Tacoms

81

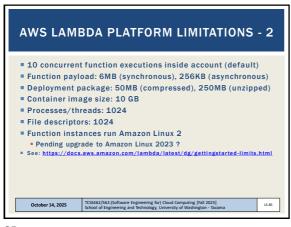
TYPES OF FUNCTION CALLS:
ASYNCHRONOUS

- 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)

AWS LAMBDA PLATFORM LIMITATIONS Maximum 10 GB memory per function instance Maximum 15-minutes execution per function instance 500 MB of /tmp disk space for local I/O (default) Up to 10 GB /tmp ephemeral storage (for additional charge) https://aws.amazon.com/ blogs/aws/aws-lambda-now-supports-up-to-10gb-ephemeral-storage/ Access up to 6 vCPUs depending on memory reservation size Figure 1: AWS Lambda Performance Speedup for Sysbench Prime Number Generation vs. Function Me October 14, 2025

83 84

Slides by Wes J. Lloyd L5.14



CPUSTEAL

CPUSTEAL

CPUSteal: Metric that measures when a CPU core is ready to execute but the physical CPU core is busy and unavailable

Symptom of over provisioning physical servers in the cloud

Factors which cause CpuSteal: (x86 hyperthreading)

Physical CPU is shared by too many busy VMs

Hypervisor kernel is using the CPU

On AWS Lambda this would be the Firecracker MicroVM which is derived from the KVM hypervisor

Wis 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

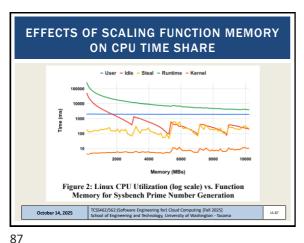
Metric can be read using SAAF in tutorial #4

Cotober 14, 2025

TSSS462/S621/Schaue Engineering for Cloud Computing [Fall 2025]
School of Engineering and Technology, University of Weshington - Taccoma

Lago

85 86



EFFECTS OF SCALING FUNCTION MEMORY
ON CPU TIME SHARE

-User - Idle - Steal - Runtime - Kernel

1000000

Key observations:
Runtime decreases as vCPUs and CPU time share increase
CPU user time remains constant for the prime number generation task - work doesn't change
CPU idle time gradually decreases as memory and vCPUs increase (the idle time is becoming active time)
When the 4th vCPU is added, cpuSteal tracks closely with cpuIdle time (hyperthreading effect)
There is more cpu Kernel time after the 4th vCPU is added

October 14, 2025

| TCSAEZ/SGZ/Software Engineering for) Cloud Computing [fiell 2025]
| School of Engineering and Technology, University of Washington - Tacoma

0/

FUNCTION INSTANCE LIFE CYCLES

Function states:

GOLD: brand new function instance just initialized to run the request (more overhead)
Platform cold (first time ever run)
Host cold (function assets cached locally on servers)

WARM: existing function instance that is reused
All function instances persist for ~5 minutes before they begin to be "garbage collected" by the platform
100% garbage collection may take up to ~30-40 minutes

AWS Lambda appears to "recycle" infrastructure faster than other FaaS platforms
Presumably because of need, because the platform is busy

WARM VS COLD FUNCTION INSTANCES

-5 Minute interval - 10 Minute interval

-5 Minute interval - 10 Minute interval

-5 Minute interval - 10 Minute interval

-5 Minute interval

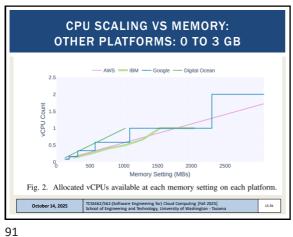
-6 Minute interval

-6 Minute interval

-7 Minute

89 90

Slides by Wes J. Lloyd L5.15

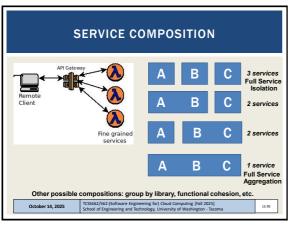


CPU SCALING VS MEMORY: OTHER PLATFORMS: 0 TO 3 GB Key observations: Google only supports strict memory steps AWS gradually increases the CPU time share as memory is increased IBM is similar but slope is not constant Digital Ocean only scales up to 1 GB Fig. 2. Allocated vCPUs available at each memory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

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 \$\tilde{\pi}\$/GB/month • Project: EFS performance & scalability evaluation on Lambda

SERVERLESS FILE STORAGE COMPARISON PROJECT Elastic File System (EFS): Performance, Cost, and Scalability Evaluation in the context of AWS Lambda / Serverless Computing EFS provides a file system that can be shared with multiple Lambda function instances in parallel Using a common use case, compare performance and cost of extended storage options on AWS Lambda: Docker container support (up to 10 GB) - read only Emphemeral /tmp (up to 10 GB) - read/write EFS (unlimited, but costly) - read/write image integration with AWS Lambda - performance & scalability TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Taco October 14, 2025 L5.94

93

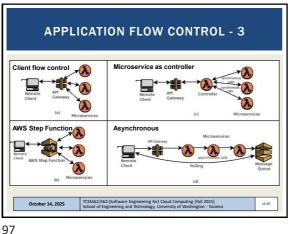


SWITCH-BOARD ARCHITECTURE (A) (A) (A) 000 **3 3 3** Remote Switchboard Single deployment package with consolidated codebase (Java: one JAR file) Entry method contains "switchboard" logic
Case statement that route calls to proper service Routing is based on data payload Check if specific parameters exist, route call accordingly Goal: reduce # of COLD starts to improve performance TCSS462/562:(Software Engineering for) Cloud Computing (Fall 2025 School of Engineering and Technology, University of Washington - Tar October 14, 2025

95 96

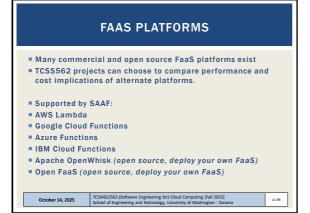
Slides by Wes J. Lloyd L5.16

92



98





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, Mongo DB, S3

October 14, 2025

STSS462/582:Schware Engineering for Coud Computing [Fall 2025]

99

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 14, 2025

TCSS462/S621/Software Engineering for) Cloud Computing [fail 2025] school of Engineering and Technology, University of Washington - Taccoma

CPU STEAL CASE STUDY

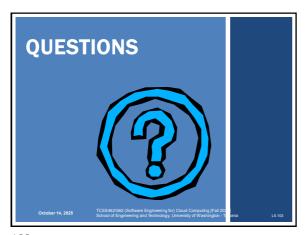
On AWS Lambda (or other FaaS platforms), when we run functions, how much CpuSteal do we observe?
How does CpuSteal vary for different workloads? (e.g. functions that have different resource requirements)
How does CpuSteal vary over time hour, day, week, location?
How does CpuSteal relate to function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

101 102

Slides by Wes J. Lloyd L5.17

CSS 462: Cloud Computing [Fall 2025]

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



103

Slides by Wes J. Lloyd L5.18