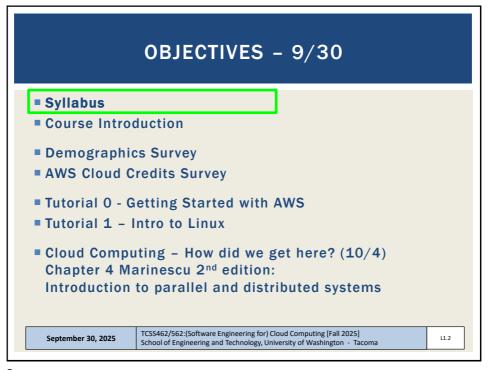


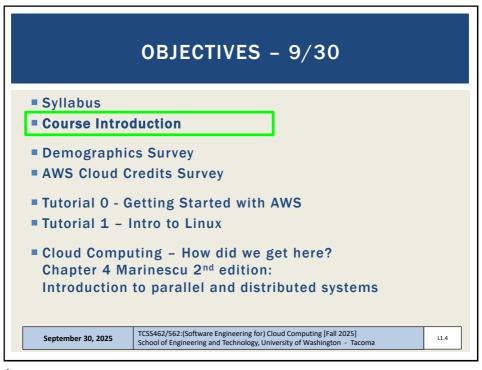
Τ



2

# TCSS562 - SOFTWARE ENGINEERING FOR CLOUD COMPUTING Course webpage is embedded into Canvas In CANVAS to access links: RIGHT-CLICK - Open in new window Syllabus online at: http://faculty.washington.edu/wlloyd/courses/tcss562/ Grading Schedule Assignments TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

3



4

### TCSS 462/562 - Fall 2025

- In-Person UWT JOY 215 Live Streamed on Zoom
- Class sessions are streamed LIVE and recorded for 24/7 availability
  - Recordings deleted after ~120 days
- ■18 class meetings
  - 2 Holidays: No Class on Nov 11, Nov 27
- This course will have 2 in-person quizzes
- ■This course can help with preparation for TCSS 558 Applied Distributed Computing

TCSS 462/ TCSS 562 FALL 2025

L1.5

5

### **DELIVERY FORMAT**

- Fall 2025 TCSS 462/562 :
  - In-person meetings JOY 215
  - Video live-stream of lectures + recordings by Zoom
  - Options to complete and submit most assignments remotely
- Please note: UWT does not provide professional video production services. Provided recordings/live-streams are provided with best-effort, but quality is not guaranteed.

September 30, 2025

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

1.6

6

# IMPROVING PERFORMANCE IN COLLEGE CLASSES

- From the DVD/Book: "Where there's a will there's an A"
- Three simple things the instructor remembers for improving grades in college classes:
  - 1. Attend every class
  - 2. Sit in the front row (or as close to the front as possible)
  - 3. Read the book (or assigned reading) all of it

• If not satisfied with recent grades, are you doing these things?

September 30, 2025

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

L1.7

7

### 462/562 RECOMMENDATIONS

- Attend in person as much as possible
  - Details of assignments can be easily explained with impromptu questions in-person which often do not occur online (Zoom)
- Stay current
  - Attend in-person or review lectures weekly
- Meet your project team members
  - Work out the best arrangements for the team (in-person, remote, etc.)
  - Expect some up-front in-person planning before remote work
- Quizzes and class presentations in person

September 30, 2025 TCSS462/56

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

L1.8

8

### INDUSTRY CHALLENGES

- Recently the job market has become increasingly competitive
- Many companies have returned to at/near 100% in-person
  - Amazon January 2025
  - Starbucks Corporate Offices
  - Commuter traffic has increased
- After an employee's job market (2021-23), we have now entered an employers job market (2024-)
  - Al is less the cause than the media says it is mostly economics
- Thie is a good time to asses your dedication and commitment to a CS degree, and set <u>GOALS.</u>.
- Many people are seeking to gain additional skills through graduate study and other specializations (i.e. data science certificates, etc.) to differentiate themselves in a competitive job market
- The job market is not impossible, but it is good time to assess your plans, set goals, and commit yourself to taking the best steps possible to be successful to meet job market challenges

September 30, 2025

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

L1.9

9

### **GRADUATE CREDIT OPPORTUNITY**

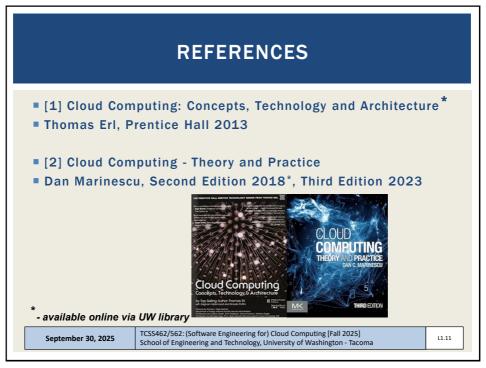
- Are you a BS CSS / BA CSS Student?
- Consider taking TCSS 562 this quarter instead of TCSS 462
- BS students can take 1 x 500-level course as a senior at UWT which can apply to both the BS and MS CSS degrees (double-dlp)
- On the fence about grad school? Taking one course as an undergrad reduces the total MS degree credits from:
  - 40 to 35 (capstone or thesis option)
  - 45 to 40 (coursework only option)
- Taking an MS CSS course while an undergrad saves money you pay the undergraduate tuition rate
  - Savings\*: \$3,086 (resident), \$1,857 (non-resident)
  - \* For registration in 5-credits, full-time savings is similar
- From: <a href="https://grad.uw.edu/policies/1-1-graduate-degree-requirements/">https://grad.uw.edu/policies/1-1-graduate-degree-requirements/</a>

September 30, 2025

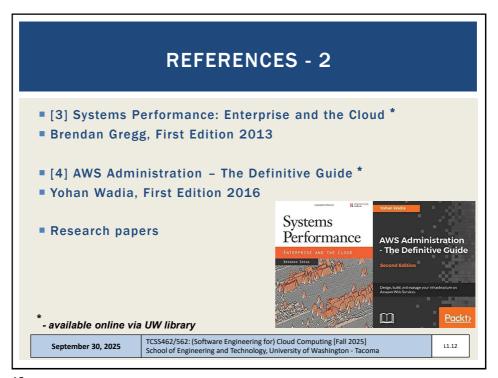
TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L1.10

10



11



12

## TCS462/562 COURSE WORK

- Cloud Computing Tutorials 20%
- Project Status Reports / Activities 5%
  - ~ 2-4 total items (??)
  - Variety of formats: in class, online, reading, activity
- Quizzes 20%
  - Open book, note, etc.
- Class Presentation (TCSS 562)
   Class Presentation Summaries (TCSS 462/562) 20%
- Term Project / Paper or Presentation 35%
  - Includes Project Proposal

September 30, 2025

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

L1.13

13

### **TERM PROJECT**

- Project description to be posted
- Teams of ~4, self formed, one project leader
- Project scope can vary based on team size and background w/ instructor approval
- Proposal due: Thursday October 16, 11:59pm (tentative)
- Approach:
  - Build a "cloud native" web services application
    - Using serverless computing, containerization, or other
    - App will consist of multiple services (FaaS functions)
    - Objective is to compare alternate implementations / designs
      - Performance (runtime)
      - Cost (\$)

September 30, 2025

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

L1.14

14

### TERM PROJECT - 2

- GOAL: Compare alternate application implementations:
- THEME for Fall 2025: GENERATIVE AI FOR SERVERLESS COMPUTING:
- LLMs can help generate code for serverless cloud applications
- Theme for Fall 2025: investigate how cloud performance and cost of multiple versions of LLM generated code can vary
  - (1) <u>Alternate LLMs</u>: code generated from alternate LLMs: e.g. ChatGPT vs. Claude vs. Gemini, etc.
  - (2) <u>Alternate Programming Languages</u>: use 1 LLM (ChatGPT) to generate code in different languages to compare cloud performance
  - (3) <u>Alternate Prompts</u>: write alternative versions of prompts asking to generate the same code, compare outcomes
- GOAL: evaluate performance and cost to understand LLM, programming language, and/or prompting trade-offs for Al code generation for serverless cloud computing
- Challenge: ensure that LLMs produce functionally correct code

>> it may be faster, but does it really work ???

September 30, 2025

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

L1.15

15

### **TERM PROJECT - 3**

- A & B Testing
  - Compare performance of approaches: language A vs. B
  - Use statistical methods to infer which performs better
    - t-tests: student t-test, Welch's t-test (unequal sample sizes or variances), Mann-Whitney U test (non-normal data)
  - Specify and test specific performance goals
  - Performance: runtime (ms), throughput (requests/sec), network latency (ms), data throughput (MB/sec), others...
- Focus is on performance & cost
- Other quality aspects, we assume the cloud provides for us: high availability, accessibility, resilience to failure, usability

September 30, 2025

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

L1.16

16

### **TERM PROJECT - 4**

- **■** Deliverables:
  - TCSS 562: Project paper (4-6 pgs IEEE format, template provided)
  - TCSS 462: Comprehensive recorded video presentation (12-15 minutes), project paper option
  - GitHub (project source)
  - How-To document describing how to test the system (via GitHub markdown)
- Suggested application:
  - Implement a multi-function data processing pipeline:
     Extract-Transform-Load (ETL) data processing pipeline combing AWS Lambda, S3, and Amazon Aurora DB

September 30, 2025

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

L1.17

17

### **TERM PROJECT - 5**

- Primary goal for the term project is to implement a cloudbased application and investigate 1 or more design trade-offs
  - Fall 2025 focus LLM code generation performance comparison
- Teams evaluate the impact of different designs (implementations) on performance and cost objectives and report on the results
- Creative projects encouraged!
- Groups do not have to follow the Fall 2025 THEME
- Groups can propose and implement any project that analyzes other design trade-offs (besides LLM code generation)

September 30, 2025

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

L1.18

18

# COMPARING DIFFERENT DESIGN TRADE-OFFS

- What other design trade-offs can be compared?
- Compare alternative app designs using different cloud services (e.g. databases), languages, platforms, etc.
- Examples Compare different:
- Cloud storage services: Object/blob storage services
  - Amazon S3, Google blobstore, Azure blobstore, vs. self-hosted
- Cloud relational database services:
  - Amazon Relational Database Service (RDS), Aurora, Self-Hosted DB
- Platform-as-a-Service (PaaS) alternatives for web app hosting:
  - Amazon Elastic Beanstalk, Heroku, others
- Open source FaaS platforms
  - Apache OpenWhisk, OpenFaaS, Fn, others...

September 30, 2025

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

L1.19

19

# COMPARING DIFFERENT DESIGN TRADE-OFFS - 2

- Serverless storage alternatives
  - On AWS: Amazon EFS, S3, Containers, others
- Container platforms
  - Amazon ECS/Fargate, AKS, Azure Kubernetes, Self-hosted Kubernetes cluster on cloud VMs
- Contrasting queueing service alternatives
  - Amazon SQS, Amazon MQ, Apache Kafka, RabbitMQ, Omq, others
- NoSQL database services
  - DynamoDB, Google BigTable, MongoDB, Cassandra
- CPU architectures
  - Intel (x86\_64), AMD (x86\_64), ARM (Graviton), MAC (M1)
- Service designs or compositions

September 30, 2025

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

L1.20

20

### TERM PROJECT: BIG PICTURE

- 1. BUILD A MULTI-FUNCTION SERVERLESS APPLICATION
  - Typically consisting of AWS Lambda Functions or Google Cloud Functions, etc. (e.g. FaaS platform)
- 2. CONTRAST THE USE OF ALTERNATIVE LLM CODE OR DESIGNS TO IMPLEMENT THE SAME APPLICATION MULTIPLE TIMES
- 3. CONDUCT A PERFORMANCE EVALUATION, REPORT FINDINGS IN TERM PAPER (562) OR PRESENTATION (462) (10-15-minutes recorded)

September 30, 2025

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

L1.21

21

### **TERM PROJECT - KEY REQUIREMENTS**

- 1. Application should involve multiple processing steps
- Implementation does not have to be Function-as-a-Service (FaaS)
- 3. Implementation leverages multiple cloud services (e.g. databases, object stores, queues)
- 4. Projects will contrast alternate designs/code
- 5. Define your comparison metrics:
- Which designs offer the <u>fastest performance (runtime)</u>?
- Lowest cost (\$)?
- Best maintainability?

Consider size: lines of code (LOC), smaller programs are generally considered to be easier to maintain

September 30, 2025

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

L1.22

22

### TERM PROJECT: RESEARCH

- Alternative: Conduct a cloud-related research project on any topic focused on specific research goals / questions
  - Can help spur MS Capstone/Thesis or BS honors thesis projects
  - Identify and investigate 1 2 research questions
  - Implement a novel solution to an open problem
  - Complete initial research towards publishing a conference or workshop paper
  - If you're interested in this option, please talk with the instructor
- Instructor will help guide projects throughout the quarter
- Explore our growing body of cloud research publications at: http://faculty.washington.edu/wlloyd/research.html

September 30, 2025

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

L1.23

23

### **PROJECT SUPPORT**

- Project cloud infrastructure support:
- AWS Account Paid Plan
  - Create standard AWS account with UW email
  - Credit card required
  - Provides access to all AWS services
  - Initial \$100 free credit, second \$100 free credit 6 mo expiration
  - Additional credits available from Instructor throughout Fall quarter
- AWS Account Free Plan
  - No Credit Card required
  - Provides access to a subset of AWS services
  - Initial \$100 free credit, second \$100 free credit 6 mo expiration
  - Only free tier services accessible after credits exhausted

September 30, 2025

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

L1.24

24

### PROJECT SUPPORT - 2

- Other Support :
- Github Student Developer Pack:
  - https://education.github.com/pack
  - Formerly offered AWS credits, but Microsoft bought GitHub
  - Includes up to \$200 in Digital Ocean Credits
  - Includes up to \$100 in Microsoft Azure Credits
  - Unlimited private git repositories
  - Several other benefits
- Microsoft Azure for Students
  - \$100 free credit per account valid for 1 year no credit card (?)
  - https://azure.microsoft.com/en-us/free/students/
- Google Cloud
  - \$300 free credit for 1 year
  - https://cloud.google.com/free/
- Chameleon / CloudLab
  - Bare metal NSF cloud free

September 30, 2025

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

L1.25

25

# TERM PROJECT RESEARCH OPPORTUNITIES

- Projects can lead to papers or posters presented at ACM/IEEE/USENIX conferences, workshops
  - Networking and research opportunity
    - ... travel ???
  - Conference participation (posters, papers) helps differentiate your resume/CV from others
- Project can support preliminary work for:
   UWT BS honors, MS capstone/thesis projects
- Research projects provide valuable practicum experience with cloud systems analysis, prototyping
- Publications are key for building your resume/CV, Also very important for applying to PhD programs

September 30, 2025

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

L1.26

26

### TCSS562 TERM PROJECT - 3

- Project status report / term project check-ins
  - Written status report
  - 1 or 2 reports during the quarter
  - Part of: "Project Status Reports / Activities / Quizzes" category
  - 5% of grade
- Project meetings with instructor
  - After class, office hours, by appointment

September 30, 2025

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

L1.27

27

### TCSS 562: CLASS PRESENTATION

- TCSS 562 students will give a team presentation teams of ~3
- Cloud Service Review Presentation
  - PPT Slides, demonstration
  - Present a cloud service not covered in class
  - Present overview of features, performance, etc.
- Cloud Research Paper Review Presentation
  - PPT slides, identify research contributions, strengths and weaknesses of paper, possible areas for future work

September 30, 2025

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

L1.28

28

### **CLASS PRESENTATION PEER REVIEWS**

- ALL Students will submit reviews of class presentations using rubric worksheet (~ 1-page)
- Students will review a minimum of one presentation for each presentation day, for a minimum of 4 reviews
- In addition to the reviews, students will write two questions about content in the presentation. These can be questions to help clarify content from the presentation that was not clear, or any related questions inspired by the presentation.
- To ensure intellectual depth of questions, questions should not have yes-no answers.
- Peer reviews will be shared with presentation groups to provide feedback but <u>will not</u> factor into the grading of class presentations

September 30, 2025

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

L1.29

29

### **CLASS PRESENTATION PEER REVIEWS - 2**

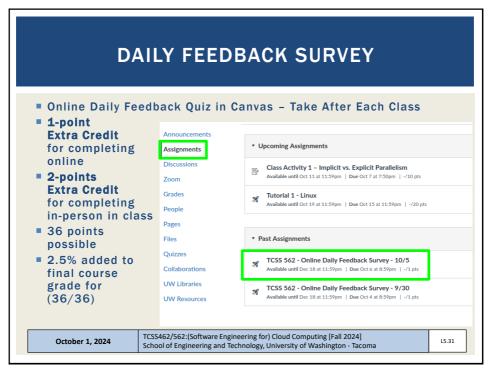
- For TCSS 462 the 4 required peer reviews will count for the entire presentation score
- For TCSS 562 the peer reviews will count as ~20% of the presentation score

September 30, 2025

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

L1.30

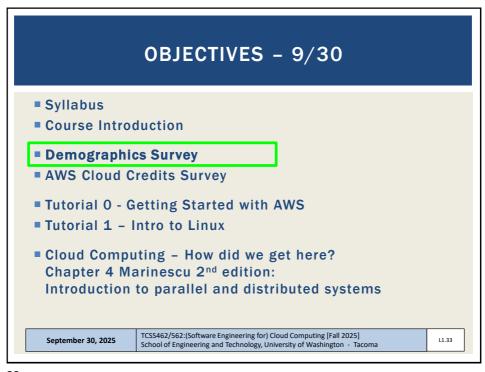
30



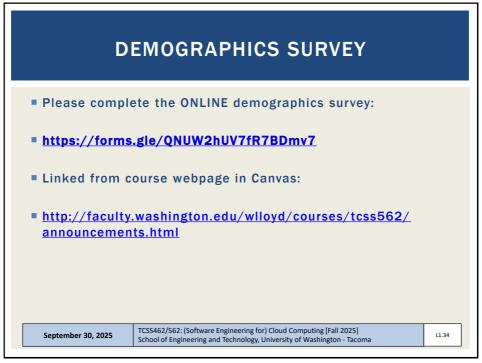
31

TCS	SS 562 - Online Daily Feedback Survey - 10/5
	d: Oct 7 at 1:13am
Qui	z Instructions
	Question 1 0.5 pts
	On a scale of 1 to 10, please classify your perspective on material covered in today's class:
	1 2 3 4 5 6 7 8 9 10
	Mostly Equal Mostly Review To Me New and Review New to Me
D	Question 2 0.5 pts
	Please rate the pace of today's class:
	1 2 3 4 5 6 7 8 9 10
	Slow Just Right Fast
October 1, 202	TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2024] School of Engineering and Technology, University of Washington - Tacoma  L5.32

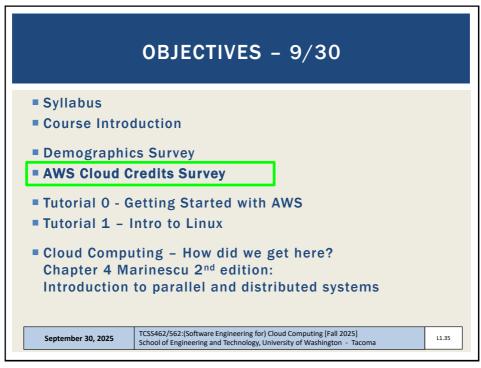
32



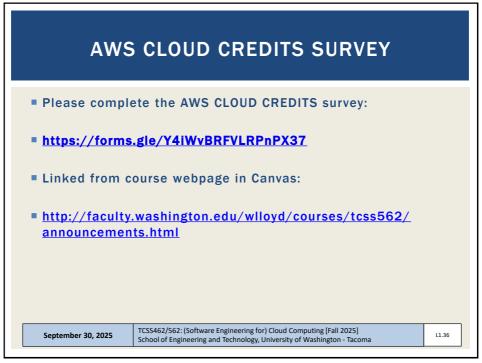
33



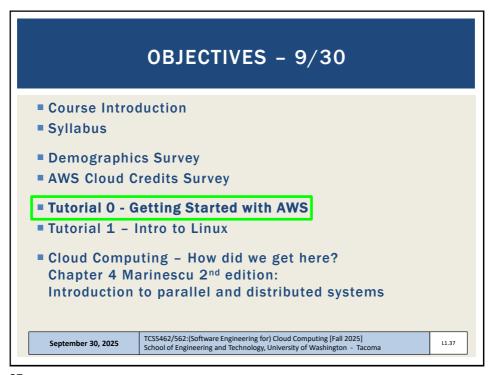
34



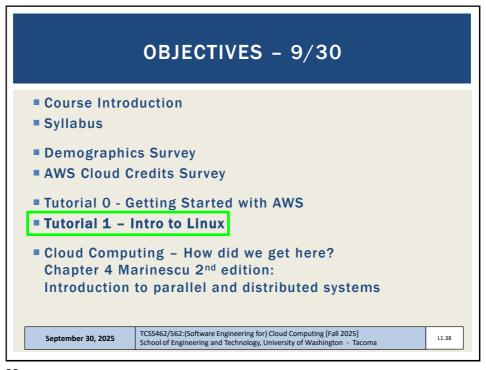
35



36



37



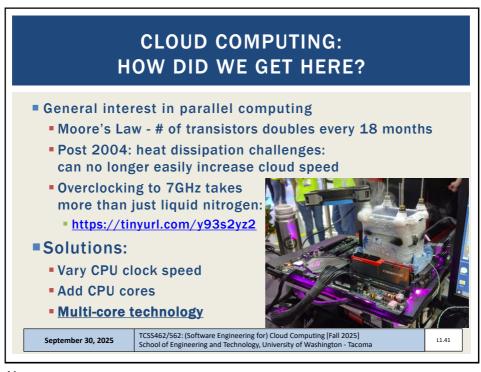
38

# OBJECTIVES - 9/30 Course Introduction Syllabus Demographics Survey AWS Cloud Credits Survey Tutorial 0 - Getting Started with AWS Tutorial 1 - Intro to Linux Cloud Computing - How did we get here? Chapter 4 Marinescu 2<sup>nd</sup> edition: Introduction to parallel and distributed systems CS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

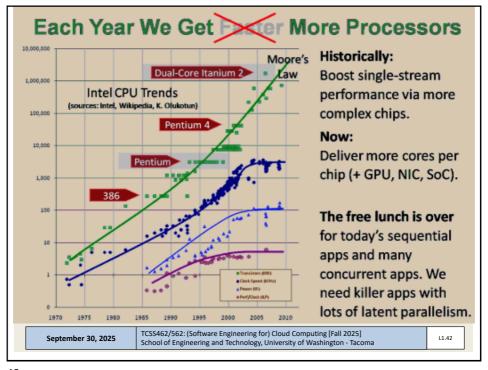
39

# Cloud Computing: How did we get here? Parallel and distributed systems (Marinescu Ch. 2 - 1st edition, Ch. 4 - 2nd edition) Data, thread-level, task-level parallelism Parallel architectures SIMD architectures, vector processing, multimedia extensions Graphics processing units Speed-up, Amdahl's Law, Scaled Speedup Properties of distributed systems Modularity TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

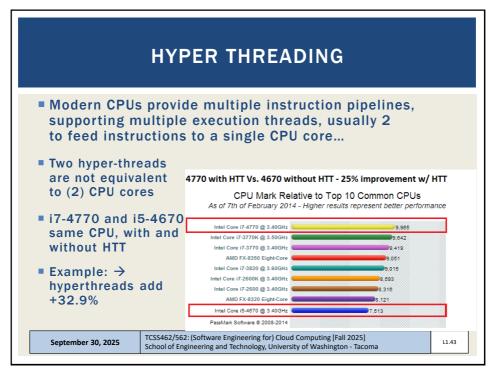
40



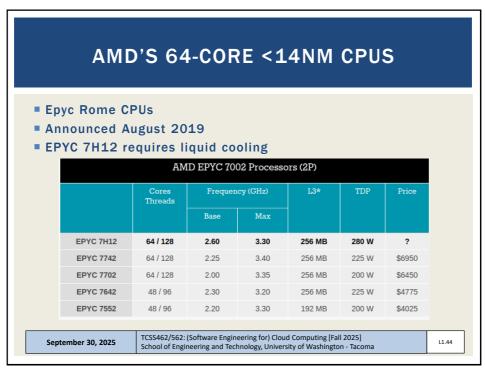
41



42



43



44

### AMD'S 64-CORE < 14NM CPUS

### AMD EPYC 9654/9654P/9684X/9R14 (AWS OEM):

- June 2023: <u>96 cores</u>, 192 hyper-threads CPUs
- Mixes 4nm:APU (combines CPUs+GPU), 5nm:L3 cache (8 CPU-chiplet), and 6nm:I/O dies, 2.25 to 3.7 burst GHz, up to 400 watts
- \$10,625 to \$14,756

AMD EPYC 9754: 128 cores, 256 hyperthreads!

- 2.25 to 3.1 burst GHz, 360 watts
- \$11,900

AMD EPYC 9005: 192 cores, 384 threads, 3nm (in dev)

September 30, 2025

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

L1.45

45

## CLOUD COMPUTING: HOW DID WE GET HERE? - 2

- To make computing faster, we must go "parallel"
- Difficult to expose parallelism in scientific applications
- Not every problem solution has a parallel algorithm
  - Chicken and egg problem...
- Many commercial efforts promoting pure parallel programming efforts have failed
- Enterprise computing world has been skeptical and less involved in parallel programming

September 30, 2025

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

L1.46

46

# CLOUD COMPUTING: HOW DID WE GET HERE? - 3

- Cloud computing provides access to "infinite" scalable compute infrastructure on demand
- Infrastructure availability is key to exploiting parallelism
- Cloud applications
  - Based on client-server paradigm
  - Thin clients leverage compute hosted on the cloud
  - Applications run many web service instances
  - Employ load balancing

September 30, 2025

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

L1.47

47

## CLOUD COMPUTING: HOW DID WE GET HERE? - 4

- Big Data requires massive amounts of compute resources
- MAP REDUCE
  - Single instruction, multiple data (SIMD)
  - Exploit data level parallelism
- Bioinformatics example

September 30, 2025

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

1.48

48

# SMITH WATERMAN USE CASE



- Applies dynamic programming to find best local alignment of two protein sequences
  - Embarrassingly parallel, each task can run in isolation
  - Use case for GPU acceleration
- AWS Lambda Serverless Computing Use Case: Goal: Pair-wise comparison of all unique human protein sequences (20,336)
  - Python client as scheduler
  - C Striped Smith-Waterman (SSW) execution engine

From: Zhao M, Lee WP, Garrison EP, Marth GT: SSW library: an SIMD Smith-Waterman C/C++ library for use in genomic applications. PLoS One 2013, 8:e82138

September 30, 2025

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

L1.49

49

### **SMITH WATERMAN RUNTIME**

- Laptop server and client (2-core, 4-HT): 8.7 hours
- AWS Lambda FaaS, laptop as client: 2.2 minutes
  - Partitions 20,336 sequences into 41 sets
  - Execution cost: ~ 82¢ (~237x speed-up)
- AWS Lambda server, EC2 instance as client: 1.28 minutes
  - Execution cost: ~ 87¢ (~408x speed-up)
- Hardware
  - Laptop client: Intel i5-7200U 2.5 GHz :4 HT, 2 CPU
  - Cloud client: EC2 Virtual Machine m5.24xlarge: 96 vCPUs
  - Cloud server: Lambda ~1000 Intel E5-2666v3 2.9GHz CPUs

September 30, 2025

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

L1.50

50

# CLOUD COMPUTING: HOW DID WE GET HERE? - 3

- Compute clouds are large-scale distributed systems
  - Heterogeneous systems
  - Homogeneous systems
  - Autonomous
  - Self organizing

September 30, 2025

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

L1.51

51

### **OBJECTIVES**

- Cloud Computing: How did we get here?
  - Parallel and distributed systems
     (Marinescu Ch. 2 1<sup>st</sup> edition, Ch. 4 2<sup>nd</sup> edition)
  - Data, thread-level, task-level parallelism
  - Parallel architectures
  - SIMD architectures, vector processing, multimedia extensions
  - Graphics processing units
  - Speed-up, Amdahl's Law, Scaled Speedup
  - Properties of distributed systems
  - Modularity

September 30, 2025

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

1.52

52

### **PARALLELISM**

- Discovering parallelism and development of parallel algorithms requires considerable effort
- Example: numerical analysis problems, such as solving large systems of linear equations or solving systems of Partial Differential Equations (PDEs), require algorithms based on domain decomposition methods.
- How can problems be split into independent chunks?
- Fine-grained parallelism
  - Only small bits of code can run in parallel without coordination
  - Communication is required to synchronize state across nodes
- Coarse-grained parallelism
  - Large blocks of code can run without coordination

September 30, 2025

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

L1.53

53

### **PARALLELISM - 2**

- Coordination of nodes
- Requires <u>message passing</u> or <u>shared memory</u>
- Debugging parallel <u>message passing</u> code is easier than parallel <u>shared memory</u> code
- Message passing: all of the interactions are clear
  - Coordination via specific programming API (MPI)
- **Shared memory**: interactions can be implicit must read the code!!
- Processing speed is orders of magnitude faster than communication speed (CPU > memory bus speed)
- Avoiding coordination achieves the best speed-up

September 30, 2025

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

L1.54

54

### TYPES OF PARALLELISM

- Parallelism:
  - Goal: Perform multiple operations at the same time to achieve a speed-up
- ■Thread-level parallelism (TLP)
  - Control flow architecture
- Data-level parallelism
  - Data flow architecture
- Bit-level parallelism
- ■Instruction-level parallelism (ILP)

September 30, 2025

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

L1.55

55

### THREAD LEVEL PARALLELISM (TLP)

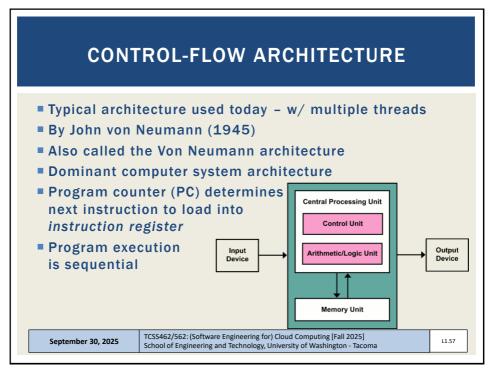
- Number of threads an application runs at any one time
- Varies throughout program execution
- As a metric:
- Minimum: 1 thread
- Can measure <u>average</u>, <u>maximum (peak)</u>
- QUESTION: What are the consequences of <u>average</u> (TLP) for scheduling an application to run on a computer with a fixed number of CPU cores and hyperthreads?
- Let's say there are 4 cores, or 8 hyper-threads...
- **Key to avoiding waste of computing resources** is knowing your application's TLP...

September 30, 2025

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

L1.56

56



57



- Partition data into big chunks, run separate copies of the program on them with little or no communication
- Problems are considered to be embarrassingly parallel
- Also perfectly parallel or pleasingly parallel...
- Little or no effort needed to separate problem into a number of parallel tasks
- MapReduce programming model is an example

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

School of Engineering and Technology, University of Washington - Tacoma

58

September 30, 2025

### DATA FLOW ARCHITECTURE

- Alternate architecture used by network routers, digital signal processors, special purpose systems
- Operations performed when input (data) becomes available
- Envisioned to provide much higher parallelism
- Multiple problems has prevented wide-scale adoption
  - Efficiently broadcasting data tokens in a massively parallel system
  - Efficiently dispatching instruction tokens in a massively parallel system
  - Building content addressable memory large enough to hold all of the dependencies of a real program

September 30, 2025

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

L1.59

59

### **DATA FLOW ARCHITECTURE - 2**

- Architecture not as popular as control-flow
- Modern CPUs emulate data flow architecture for dynamic instruction scheduling since the 1990s
  - Out-of-order execution reduces CPU idle time by not blocking for instructions requiring data by defining execution windows
  - Execution windows: identify instructions that can be run by data dependency
  - Instructions are completed in data dependency order within execution window
    - Execution window size typically 32 to 200 instructions

<u>Utility of data flow architectures has been</u> <u>much less than envisioned</u>

September 30, 2025

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

L1.60

60

### **BIT-LEVEL PARALLELISM**

- Computations on large words (e.g. 64-bit integer) are performed as a single instruction
- Fewer instructions are required on 64-bit CPUs to process larger operands (A+B) providing dramatic performance improvements
- Processors have evolved: 4-bit, 8-bit, 16-bit, 32-bit, 64-bit

QUESTION: How many instructions are required to add two 64-bit numbers on a 16-bit CPU? (Intel 8088)

- 64-bit MAX int = 9,223,372,036,854,775,807 (signed)
- 16-bit MAX int = 32,767 (signed)
- Intel 8088 limited to 16-bit registers

September 30, 2025

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

L1.61

61

### **INSTRUCTION-LEVEL PARALLELISM (ILP)**

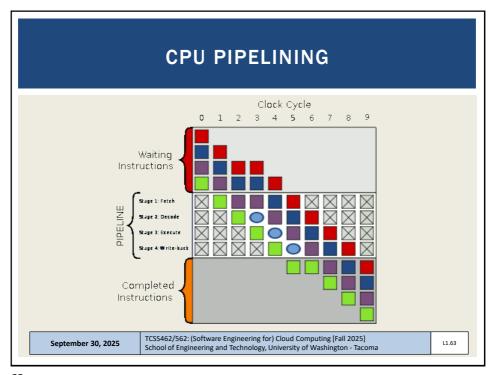
- CPU pipelining architectures enable ILP
- CPUs have multi-stage processing pipelines
- Pipelining: split instructions into sequence of steps that can execute concurrently on different CPU circuitry
- Basic RISC CPU Each instruction has 5 pipeline stages:
- IF instruction fetch
- ID- instruction decode
- EX instruction execution
- MEM memory access
- WB write back

September 30, 2025

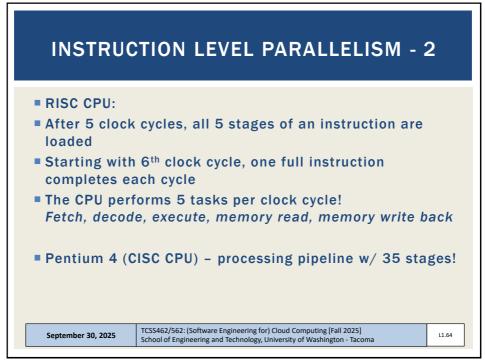
TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

1.62

62



63



64

### **OBJECTIVES**

- Cloud Computing: How did we get here?
  - Parallel and distributed systems
     (Marinescu Ch. 2 1<sup>st</sup> edition, Ch. 4 2<sup>nd</sup> edition)
  - Data, thread-level, task-level parallelism
  - Parallel architectures
  - SIMD architectures, vector processing, multimedia extensions
  - Graphics processing units
  - Speed-up, Amdahl's Law, Scaled Speedup
  - Properties of distributed systems
  - Modularity

September 30, 2025

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

L1.65

65

# MICHAEL FLYNN'S COMPUTER ARCHITECTURE TAXONOMY

- Michael Flynn's proposed taxonomy of computer architectures based on concurrent instructions and number of data streams (1966)
- SISD (Single Instruction Single Data)
- SIMD (Single Instruction, Multiple Data)
- MIMD (Multiple Instructions, Multiple Data)
- LESS COMMON: MISD (Multiple Instructions, Single Data)
- Pipeline architectures: functional units perform different operations on the same data
- For fault tolerance, may want to execute same instructions redundantly to detect and mask errors – for task replication

September 30, 2025

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

L1.66

66

### FLYNN'S TAXONOMY

- SISD (Single Instruction Single Data)
  - Scalar architecture with one processor/core.
  - Individual cores of modern multicore processors are "SISD"
- SIMD (Single Instruction, Multiple Data)

Supports vector processing

- When SIMD instructions are issued, operations on individual vector components are carried out concurrently
- Two 64-element vectors can be added in parallel
- Vector processing instructions added to modern CPUs
- Example: Intel MMX (multimedia) instructions

September 30, 2025

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

L1.67

67

# (SIMD): VECTOR PROCESSING ADVANTAGES

- Exploit data-parallelism: vector operations enable speedups
- Vectors architecture provide vector registers that can store entire matrices into a CPU register
- SIMD CPU extension (e.g. MMX) add support for vector operations on traditional CPUs
- Vector operations reduce total number of instructions for large vector operations
- Provides higher potential speedup vs. MIMD architecture
- Developers can think sequentially; not worry about parallelism

September 30, 2025

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

L1.68

68

### FLYNN'S TAXONOMY - 2

- MIMD (Multiple Instructions, Multiple Data) system with several processors and/or cores that function asynchronously and independently
- At any time, different processors/cores may execute different instructions on different data
- Multi-core CPUs are MIMD
- Processors share memory via interconnection networks
  - Hypercube, 2D torus, 3D torus, omega network, other topologies
- MIMD systems have different methods of sharing memory
  - Uniform Memory Access (UMA)
  - Cache Only Memory Access (COMA)
  - Non-Uniform Memory Access (NUMA)

September 30, 2025

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

L1.69

69

### **ARITHMETIC INTENSITY**

**Arithmetic intensity**: Ratio of work (W) to memory traffic r/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:

**P**rograms 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

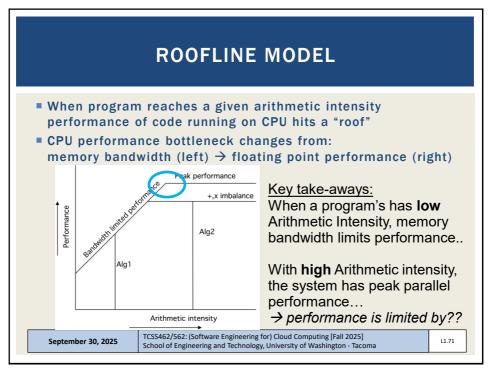
(memory RW becomes bottleneck, not enough ops!)

September 30, 2025

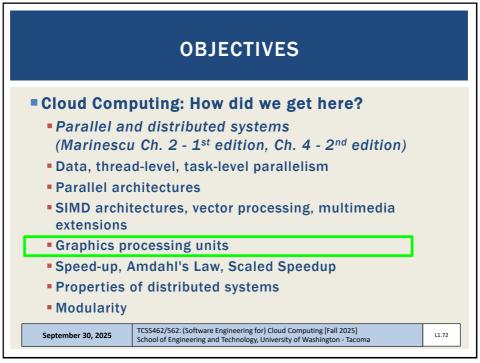
TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

L1.70

70



71



72

# 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

September 30, 2025

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

L1.73

73

#### **OBJECTIVES**

- Cloud Computing: How did we get here?
  - Parallel and distributed systems
     (Marinescu Ch. 2 1<sup>st</sup> edition, Ch. 4 2<sup>nd</sup> edition)
  - Data, thread-level, task-level parallelism
  - Parallel architectures
  - SIMD architectures, vector processing, multimedia extensions
  - Graphics processing units
  - Speed-up, Amdahl's Law, Scaled Speedup
  - Properties of distributed systems
  - Modularity

September 30, 2025

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

L1.74

74

#### 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

September 30, 2025

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

L1.75

75

# SPEED-UP EXAMPLE

- Consider embarrassingly parallel image processing
- Eight images (multiple data)
- Apply image transformation (greyscale) in parallel
- 8-core CPU, 16 hyper threads
- 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

September 30, 2025

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

L1.76

76

#### 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 14, 2020

TCSS562: Software Engineering for Cloud Computing [Fall 2020] School of Engineering and Technology, University of Washington - Tacoma

L5.77

77

#### AMDAHL'S LAW

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

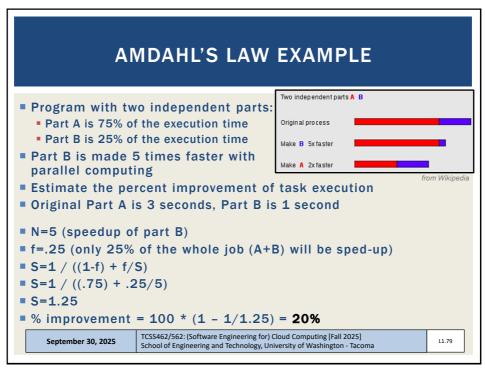
- S = theoretical speedup of the whole task
- f= fraction of work that is parallel (ex. 25% or 0.25)
- N= proposed speed up of the parallel part (ex. 5 times speedup)
- " % improvement
   of task execution = 100 \* (1 (1 / S))
- Using Amdahl's law, what is the maximum possible speed-up?

October 14, 2020

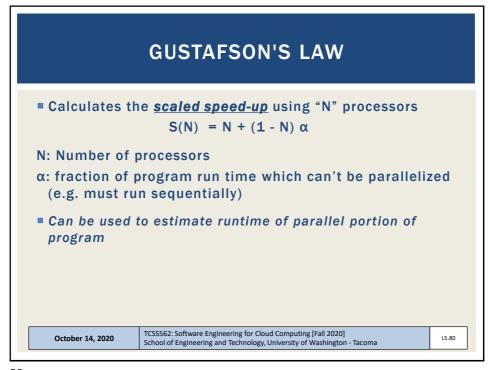
TCSS562: Software Engineering for Cloud Computing [Fall 2020] School of Engineering and Technology, University of Washington - Tacoma

L5.78

78



79



80

#### **GUSTAFSON'S LAW**

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

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

N: Number of processors

- a: 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 14, 2020

TCSS562: Software Engineering for Cloud Computing [Fall 2020] School of Engineering and Technology, University of Washington - Tacoma

L5.81

81

#### **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?

September 30, 2025

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

L1.82

82

# **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?

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

$$N=16, \alpha=.75$$

$$S(N) = 16 + (1 - 16).75$$

$$S(N) = ?$$

September 30, 2025

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

L1.83

83

# **GUSTAFSON'S EXAMPLE**

- QUESTION:

What is the maximum theoretical speed-up on a 2-core CPU?

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

 $N=2, \alpha=$ 

S(N) = 2

For 2 CPUs, speed up is 1.25x

S(N) = ?

For 16 CPUs, speed up is 4.75x

■ What is the maximum theoretical speed-up on a <u>ro-core CPU</u>?

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

 $N=16, \alpha=.75$ 

S(N) = 16 + (1 - 16).75

S(N) = ?

September 30, 2025

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

L1.84

84

#### **MOORE'S LAW**

- Transistors on a chip doubles approximately every 1.5 years
- CPUs now have billions of transistors
- Power dissipation issues at faster clock rates leads to heat removal challenges
  - ullet Transition from: increasing clock rates ullet to adding CPU cores
- Symmetric core processor multi-core CPU, all cores have the same computational resources and speed
- Asymmetric core processor on a multi-core CPU, some cores have more resources and speed
- <u>Dynamic core processor</u> processing resources and speed can be dynamically configured among cores
- Observation: asymmetric processors offer a higher speedup

September 30, 2025

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

L1.85

85

#### **OBJECTIVES**

- Cloud Computing: How did we get here?
  - Parallel and distributed systems
     (Marinescu Ch. 2 1<sup>st</sup> edition, Ch. 4 2<sup>nd</sup> edition)
  - Data, thread-level, task-level parallelism
  - Parallel architectures
  - SIMD architectures, vector processing, multimedia extensions
  - Graphics processing units
  - Speed-up, Amdahl's Law, Scaled Speedup
  - Properties of distributed systems
  - Modularity

September 30, 2025

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

L1.86

86

## **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

September 30, 2025

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

L1.87

87

#### **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

September 30, 2025

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

L1.88

88

# 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

September 30, 2025

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

L1.89

89

#### **OBJECTIVES**

- Cloud Computing: How did we get here?
  - Parallel and distributed systems
     (Marinescu Ch. 2 1<sup>st</sup> edition, Ch. 4 2<sup>nd</sup> edition)
  - Data, thread-level, task-level parallelism
  - Parallel architectures
  - SIMD architectures, vector processing, multimedia extensions
  - Graphics processing units
  - Speed-up, Amdahl's Law, Scaled Speedup
  - Properties of distributed systems
  - Modularity

September 30, 2025

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

L1.90

90

#### 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!

September 30, 2025

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

L1.91

91

# 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

September 30, 2025

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

L1.92

92

## 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

September 30, 2025

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

93

#### **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), s=latency, f=parallel fraction, N=speed-up

- Scaled speedup with N processes:

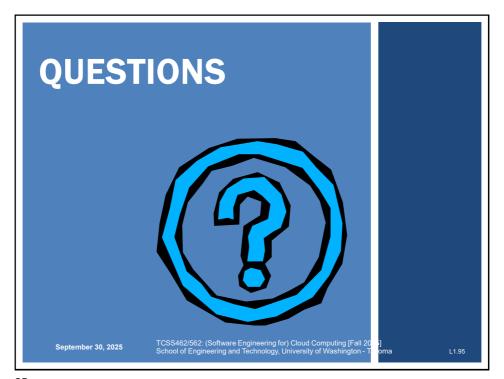
 $S(N) = N - \alpha(N-1)$ 

- 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

September 30, 2025

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

94



95