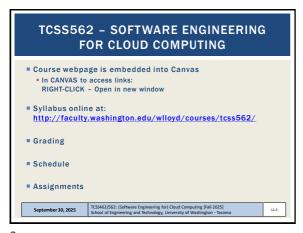


OBJECTIVES - 9/30 Course Introduction Demographics Survey AWS Cloud Credits Survey ■ Tutorial 0 - Getting Started with AWS ■ Tutorial 1 - Intro to Linux Cloud Computing - How did we get here? (10/4) Chapter 4 Marinescu 2nd edition: Introduction to parallel and distributed systems September 30, 2025



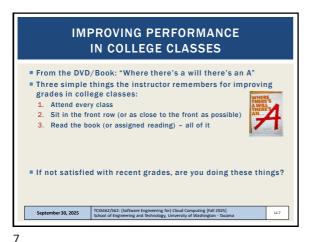
OBJECTIVES - 9/30 Syllabus Course Introduction ■ Demographics Survey AWS Cloud Credits Survey ■ Tutorial 0 - Getting Started with AWS ■ Tutorial 1 - Intro to Linux Cloud Computing - How did we get here? Chapter 4 Marinescu 2nd edition: Introduction to parallel and distributed systems September 30, 2025 L1.4

3



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

Slides by Wes J. Lloyd L1.1



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

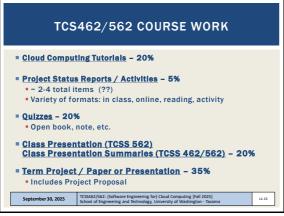
INDUSTRY CHALLENGES Recently - the job market has become increasingly competitive Many companies have returned to at/near 100% in-person Amazon - January 2025 Starbucks Corporate Offices Commuter traffic has increased After an employee's job market (2021-23), we have now entered an employers job market (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 GOALS... 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 mber 30, 2025

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: https://grad.uw.edu/policies/1-1-graduate-degreerequirements/ September 30, 2025 L1.10

11 12

Slides by Wes J. Lloyd

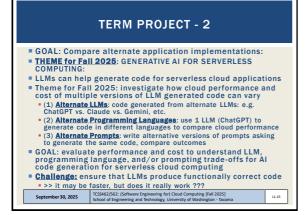


■ Project description to be posted
■ Teams of ~4, self formed, one project leader
■ Project scope can vary based on team size and background w/ instructor approval
■ Proposal due: Thursday October 16, 11:59pm (tentative)

■ 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 (\$)

| TCSS462/562: Eoftware Engineering for) Cloud Computing [Fall 2025]
| September 30, 2025 | TCSS462/562: Eoftware Engineering for) Cloud Computing [Fall 2025]

13 14



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

TSSS62/SS2: (Software Engineering for) Cloud Computing [fail 2025]
School of Engineering and Technology, University of Washington - Tacoma

15

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

School of Engineering and Enchnology, University of Washington - Taxoma

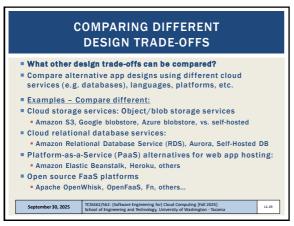
TERM PROJECT - 5

Primary goal for the term project is to implement a cloud-based 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)

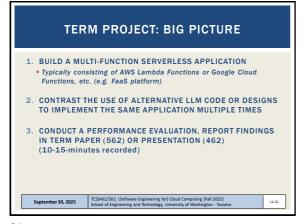
17 18

Slides by Wes J. Lloyd L1.3



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, 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 L1.20

19 20



1. Application should involve multiple processing steps
2. 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 fastest performance (runtime)?

**Lowest cost (\$)?

Best maintainability?

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

September 30, 2025

**TCSS462/S62: (Software Engineering for) Cloud Computing [Fall 2025]

**School of Engineering and Technology, University of Washington-Taccmu

**TCSS462/S62: (Software Engineering for) Cloud Computing [Fall 2025]

**School of Engineering and Technology, University of Washington-Taccmu

**TCSS462/S62: (Software Engineering for) Cloud Computing [Fall 2025]

**School of Engineering and Technology, University of Washington-Taccmu

**TCSS462/S62: (Software Engineering for) Cloud Computing [Fall 2025]

**School of Engineering and Technology, University of Washington-Taccmu

**TCSS462/S62: (Software Engineering for) Cloud Computing [Fall 2025]

21

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/wiloyd/research.html

September 30, 2025

TCSS462/S62: [Software Engineering for] Cloud Computing [Fall 2025]
School of Engineering and Technology, University of Washington-Tacoma

PROJECT SUPPORT

Project cloud infrastructure support:

AWS Account - Pald 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

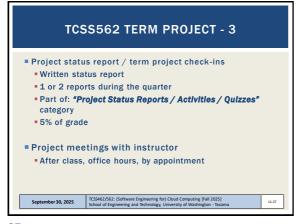
23 24

Slides by Wes J. Lloyd L1.4



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 L1.26

25 26



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

TCSS42/S62: (Software Engineering for) Cloud Computing [fail 2025] school of Engineering and Technology, University of Washington - Tacoma

27

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 will not factor into the grading of class TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacc September 30, 2025 L1.29 CLASS PRESENTATION PEER REVIEWS – 2

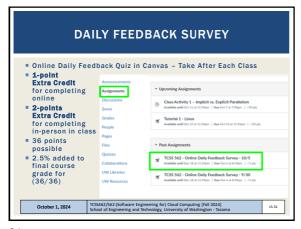
For TCSS 462 – the 4 required peer reviews will count for the entire presentation score
For TCSS 562 – the peer reviews will count as ~20% of the presentation score

September 30, 2025

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

29 30

Slides by Wes J. Lloyd L1.5



TCSS 562 - Online Daily Feedback Survey - 10/5

Stanted Cot 7 at 1:32m

Quiz Instructions

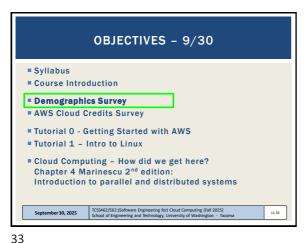
Question 1

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

Maxily Maxily Stanted Stante Stant

31 32



DEMOGRAPHICS SURVEY

Please complete the ONLINE demographics survey:

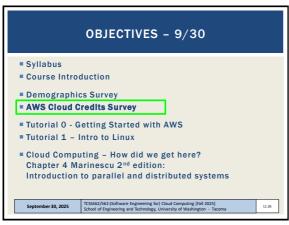
https://forms.gle/QNUW2hUV7fR7BDmv7

Linked from course webpage in Canvas:

http://faculty.washington.edu/willoyd/courses/tcss562/announcements.html

TCS462/56: [Software Engineering for] Cloud Computing [Fall 2025]
School of Engineering and Technology, University of Mashington - Tacoma

33



AWS CLOUD CREDITS SURVEY

Please complete the AWS CLOUD CREDITS survey:

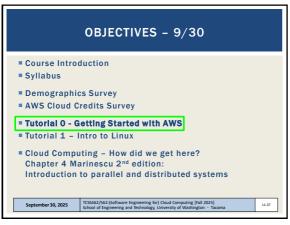
https://forms.gle/Y4IWvBRFVLRPnPX37

Linked from course webpage in Canvas:

http://faculty.washington.edu/willoyd/courses/tcss562/announcements.html

35 36

Slides by Wes J. Lloyd L1.6



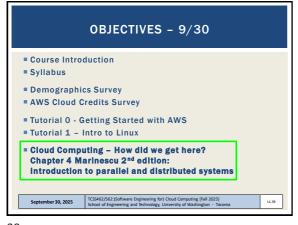
OBJECTIVES - 9/30

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 2nd edition:
Introduction to parallel and distributed systems

September 30, 2025

CASSAZ/SSZ/SGTWARE Engineering for Cloud Computing [Fall 2025]
School of Engineering and Technology, University of Washington - Tacoma

37



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

| ICSS462/SGZ: (Software Engineering for) Cloud Computing [fall 2025]
| School of Engineering and Technology, University of Washington - Tacoma

| 11.40|

39



Each Year We Get More Processors

18,000,000

Intel CPU Trends
(toource: kind, Walquedu, K. Ohkotun)

Pentium

Pentium

Rectium A

September 30, 2025

Tocsse2/562/(Soft) (Software Engineering for) Cloud Computing (Fall 2025)
School of Engineering and Technology, University of Washington-Tacons

Historically:
Boost single-stream performance via more complex chips.

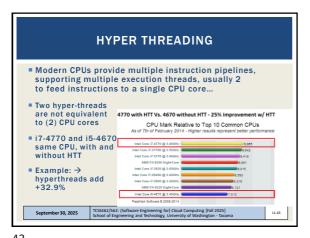
Now:
Deliver more cores per chip (+ GPU, NIC, SoC).

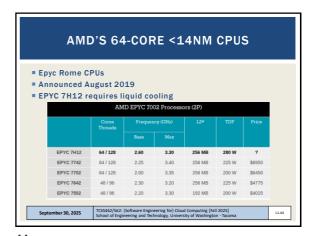
The free lunch is over for today's sequential apps and many concurrent apps. We need killer apps with lots of latent parallelism.

41 42

Slides by Wes J. Lloyd L1.7

38





43 44

AMD EPYC 9654/9654P/9684X/9R14 (AWS 0EM):

June 2023: <u>96 cores</u>, 192 hyper-threads CPUs

Mixes 4nm:APU (combines CPUs+GPU), 5nm:L3 cache
(8 CPU-chiplet), and 6nm:l/0 dies, 2.25 to 3.7 burst
GHz, up to 400 watts

\$10,625 to \$14,756

AMD EPYC 9754: <u>128 cores</u>, 256 hyperthreads!

2.25 to 3.1 burst GHz, 360 watts

\$11,900

AMD EPYC 9005: <u>192 cores</u>, 384 threads, 3nm (in dev)

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

45

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 Visabington - Tacoma

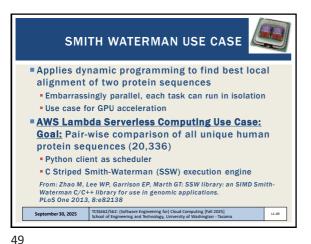
CLOUD COMPUTING:
HOW DID WE GET HERE? - 4

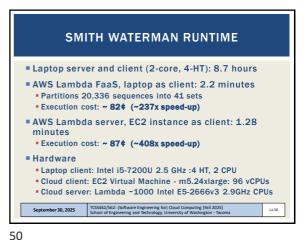
BIg Data requires massive amounts of compute resources

MAP - REDUCE
Single instruction, multiple data (SIMD)
Exploit data level parallelism
Bioinformatics example

47 48

Slides by Wes J. Lloyd L1.8





CLOUD COMPUTING:
HOW DID WE GET HERE? - 3

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

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

September 30, 2025

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

1.132

51

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 TCSS462/562: (Software Engineering for) Cloud Computing (Fall 2025) School of Engineering and Technology, University of Washington - Tac September 30, 2025 L1.53 PARALLELISM - 2

Coordination of nodes
Requires message passing or shared memory
Debugging parallel message passing code is easier than parallel shared memory 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!!

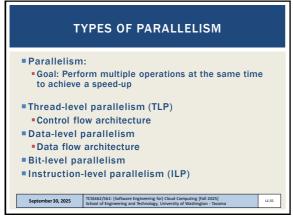
Processing speed is orders of magnitude faster than communication speed (CPU > memory bus speed)
Avoiding coordination achieves the best speed-up

Tissid2/j62: Icoflower Engineering for Cloud Computing [Iril 2025]
September 30, 2025

Tissid2/j62: Icoflower Engineering for Cloud Computing [Iril 2025]
September 30, 2025

53 54

Slides by Wes J. Lloyd L1.9



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 average, maximum (peak)

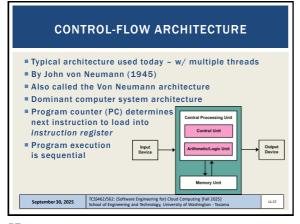
QUESTION: What are the consequences of average (TLP) for scheduling an application to run on a computer with a fixed number of CPU cores and hyperthreads?
Let's say there are 4 cores, or 8 hyper-threads...

Key to avoiding waste of computing resources is knowing your application's TLP....

September 30, 2025

TEMSAG/SSZ: Esoftware Engineering for Cloud Computing [Tail 2025]
School of Engineering and Technology (Windship) of Washington: Tacons

55



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

1CSS462/562: (Software Engineering for) Cloud Computing [Fall 2025]
School of Engineering and Technology, University of Washington - Tacoma

57

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

INSSECTIONS (Software Engineering for Obous Computing Fall 2025)

School of Engineering and Enchnology, University of Washington - Taxoma

11.59

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

Execution window size typically 32 to 200 instructions

Utility of data flow architectures has been much less than envisioned

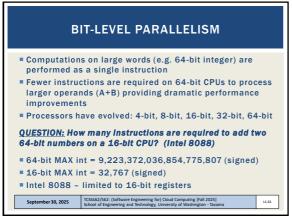
September 30, 2025

| TOSAGI/SGZ: Coltware Engineering for Cloud Computing [rail 2025] Shool of Engineering and Technology University of Versibington - Tecono

59 60

Slides by Wes J. Lloyd L1.10

56



INSTRUCTION-LEVEL PARALLELISM (ILP)

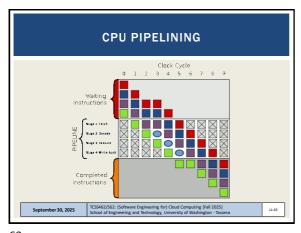
CPU pipelining architectures enable ILP
CPUs have multi-stage processing pipelines
Pipelining: split instructions into sequence of steps that can execute concurrently on different CPU circuitry

Basic RISC CPU - Each instruction has 5 pipeline stages:
IF - instruction fetch
De instruction decode
EX - instruction execution
MEM - memory access
WB - write back

September 30, 2025

TCSS462/SS2: (Software Engineering for) Cloud Computing [Fall 2025)
School of Engineering and Technology, University of Vocabington - Taconsu

61



INSTRUCTION LEVEL PARALLELISM - 2

RISC CPU:
After 5 clock cycles, all 5 stages of an instruction are loaded
Starting with 6th clock cycle, one full instruction completes each cycle
The CPU performs 5 tasks per clock cycle!
Fetch, decode, execute, memory read, memory write back
Pentium 4 (CISC CPU) - processing pipeline w/ 35 stages!

63

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

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

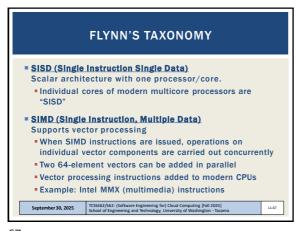
For fault tolerance, may want to execute same instructions redundantly to detect and mask errors – for task replication

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

65 66

Slides by Wes J. Lloyd

62



(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

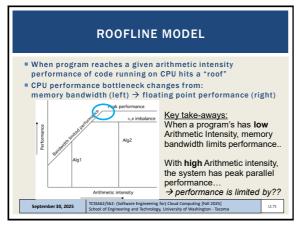
Vector operations reduce total number of instructions for large vector operations
Provides higher potential speedup vs. MIMD architecture
Developers can think sequentially; not worry about parallelism

ICSS462/562: (Software Engineering for) Gloud Compouting [Fall 2025]
School of Engineering and Technology, University of Washington - Tacoma

67

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 CCS462/562: (Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacomas

69



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

September 30, 2025

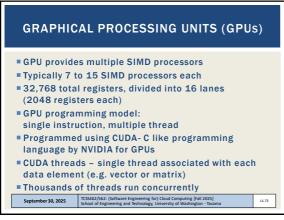
TICSS62/S62: (Software Engineering for) Cloud Computing [fall 2025]
School of Engineering and Technology, University of Washington - Tacoma

1.172

71 72

Slides by Wes J. Lloyd L1.12

68



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

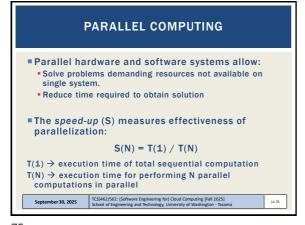
Speed-up, Amdahl's Law, Scaled Speedup

Properties of distributed systems

Modularity

TCSS402/562: [Software Engineering for) Cloud Computing [Fall 2025]
School of Engineering and Technology, University of Washington - Tacoma

73



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 TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacco September 30, 2025 L1.76

75

AMDAHL'S LAW

Amdahl's law is used to estimate the speed-up of a job using parallel computing

Divide job into two parts
Part A that will still be sequential
Portion of computation which cannot be parallelized will determine (i.e. limit) the overall speedup

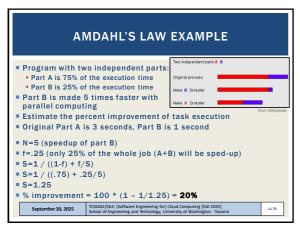
Amdahl's law assumes jobs are of a fixed size
Also, Amdahl's assumes no overhead for distributing the work, and a perfectly even work distribution

Amount of the strict of the s

77 78

Slides by Wes J. Lloyd L1.13

74



GUSTAFSON'S LAW

Calculates the scaled speed-up 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

TCSS62: Software Engineering for Cloud Computing [Fall 2020] School of Engineering and Technology, University of Washington - Tacoma

79

GUSTAFSON'S LAW

Calculates the scaled speed-up 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 σ = sequential time, π =parallel time

Our Amdahl's example: σ = 3s, π =1s, α =.75

GUSTAFSON'S LAW

Calculates the scaled speed-up 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

Ticss42/562: Coffwave Engineering for J Cloud Computing [fall 2025] school of Engineering and Technology, University of Washington - Tacoma

81

 $\begin{array}{c} \textbf{GUSTAFSON'S EXAMPLE} \\ \hline \\ \textbf{* QUESTION:} \\ \textbf{What is the maximum theoretical speed-up on a 2-core CPU$?} \\ \textbf{S}(N) &= N + (1 - N) \ \alpha \\ \textbf{N} &= 2, \ \alpha = .75 \\ \textbf{S}(N) &= 2 + (1 - 2) .75 \\ \textbf{S}(N) &= ? \\ \hline \\ \textbf{* What is the maximum theoretical speed-up on a 16-core CPU$?} \\ \textbf{S}(N) &= N + (1 - N) \ \alpha \\ \textbf{N} &= 16, \ \alpha = .75 \\ \textbf{S}(N) &= 16 + (1 - 16) .75 \\ \textbf{S}(N) &= ? \\ \hline \\ \textbf{September 30, 2025} \\ \hline \\ \textbf{TCSS462/562: Goftware Engineering for Cloud Computing [Fall 2025]} \\ \textbf{School of Engineering and Technology, University of Washington - Taccina} \\ \hline \\ \textbf{1.13} \\ \hline \end{array}$

GUSTAFSON'S EXAMPLE

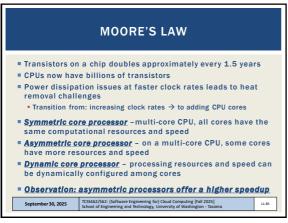
**QUESTION:* What is the maximum theoretical speed-up on a 2-core CPU? $S(N) = N + (1 - N) \alpha$ $N = 2, \alpha = S(N) = 2$ S(N) = 2 S(N) = 7For 16 CPUs, speed up is 1.25x S(N) = 7For 16 CPUs, speed up is 4.75x

What is the maximum theoretical speed-up on a 10-core CPU? $S(N) = N + (1 - N) \alpha$ $N = 16, \alpha = .75$ S(N) = 16 + (1 - 16) .75 S(N) = ?September 30, 2025 | TCSM62/562: (Software Engineering for) Cloud Computing [Fall 2025] | School of Engineering and Technology, University of Washington - Taxoma | LLM|

83 84

Slides by Wes J. Lloyd L1.14

80



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

September 30, 2025

TCSG62/S62: (Software Ingineering for) Coud Computing [Fall 2025]
School of Engineering and Technology, University of Washington - Tacoma

85 86

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 September 30, 2025 Installability at low levels of HW/software/network reliability

87

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

**Fallure 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

**Scaling transparency: System and applications can scale w/o change in system structure and w/o affecting applications

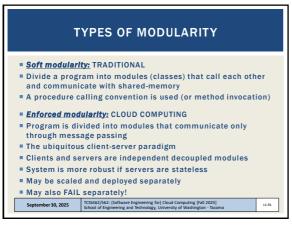
**Scaling transparency: System and applications can scale w/o change in system structure and w/o affecting applications

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

| September 30, 2025 | TCSS62/S62: [Software Engineering for) Cloud Computing [Fall 2025] | School of Engineering and Technology, University of Washington - Tacoma

89 90

Slides by Wes J. Lloyd L1.15



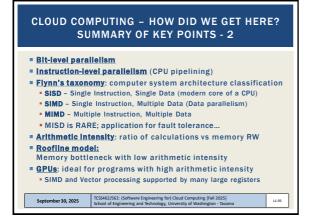
CLOUD COMPUTING - HOW DID WE GET HERE?
SUMMARY OF KEY POINTS

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

TCS462/582: Softwave Engineering (poil Could Computing [Fail 2025]
September 30, 2025

91



CLOUD COMPUTING – HOW DID WE GET HERE?

SUMMARY OF KEY POINTS - 3

Speed-up (S)
S(N) = T(1) / T(N)

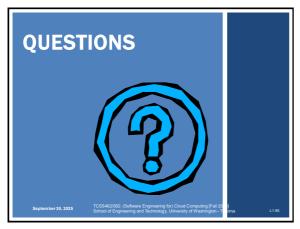
Amdahl's law:
S=1 / ((1-f) + f/N),s=latency, f=parallel fraction, N=speed-up
α = percent of program that must be sequential
Scaled speedup with N processes:
S(N) = N - α(N-1)

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

TCS:GE/SE2: (Schlware Engineering for) Cloud Computing [fail 2025]
School of Engineering and Technology, University of Washington - Tacoms

93



95

Slides by Wes J. Lloyd L1.16

92