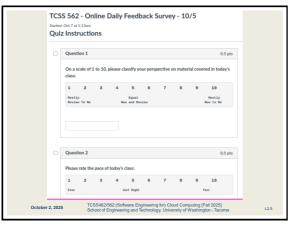


3

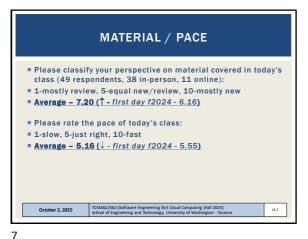


5

OBJECTIVES - 10/2 Daily Feedback Surveys Questions from 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 TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tao October 2, 2025 L2.6

6

Slides by Wes J. Lloyd L2.1



FEEDBACK FROM 9/30

What is the difference between the TCSS 462 and TCSS 562 term project?

TCSS 562 teams and hybrid-teams (462+562) submit a 4 to 6 page term paper
TCSS 462-only teams submit a 10-15 minute presentation recording with the option of submitting a term paper instead
The content of the paper and presentation are the same
A template is provided for both the term paper and presentation
The sections within are the same

Is TCSS 562 credit transferrable?
Yes - if taking TCSS 562 as an undergraduate senior, the credit will transfer to the UWT MS CSS program. There is a good chance the credit can transfer other MS Computer Science programs besides UW. This can save -\$8,000+.

FEEDBACK - 2 How can the term project contribute to a thesis or capstone project? If pursuing a thesis/capstone related to cloud computing or distributed systems, then the skills learned in the course will be important. If having an idea for a capstone/thesis project, it is possible to explore the idea by proposing and conducting preliminary research in the term project. Performance investigation projects may be extendable into a capstone project Thesis is more difficult as there needs to be an original (novel) research contribution TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2025 School of Engineering and Technology, University of Washington - Tac October 2, 2025 12.9

FEEDBACK - 3 I am unclear regarding the workload in the class Expect 1 tutorial per week throughout the quarter, 2 quizzes, plus the term project. The term project become a heavy focus after week 6. In general, this course is easier for students who have strong computing skills including the ability to learn how to use new services, troubleshoot issues/challenges, and solve problems. 462/562 is not a programming course, but a systems course. Is the project 100% from the slides, or will I need to do additional research? The term project involves applying and synthesizing what is learned across the entire course into a final project TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Taco October 2, 2025 L2.10

9

FEEDBACK - 4

What is the difference in quizzes for TCSS 462 vs. 562?
Each section (462 or 562) is graded using a separate curve
The 462 curve is usually more generous

Term Project? Jargon?
It is normal for the term project to be confusing on day 1, as we are just getting started

October 2, 2025

TCSM62/562: (Software Engineering for) Cloud Computing [Fall 2025] school of Engineering and Technology, University of Washington - Tacoma

OBJECTIVES - 10/2

Daily Feedback Surveys
Questions from Course Introduction

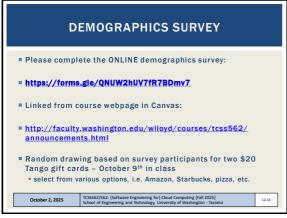
Demographics Survey
AWS Cloud Credits Survey
Tutorial 0 - Getting Started with AWS
Tutorial 1 - Intro to Linux
Cloud Computing - How did we get here?
Chapter 4 Marinescu 2nd edition:
Introduction to parallel and distributed systems

October 2, 2025

CSC462/562/Software Engineering for) Cloud Computing [Fall 2025]
School of Engineering and Technology, University of Washington - Tacoma

11 12

Slides by Wes J. Lloyd L2.2



OBJECTIVES - 10/2

Daily Feedback Surveys
Questions from Course Introduction

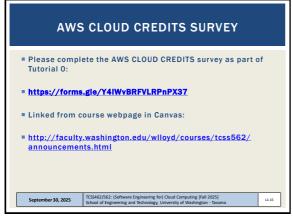
Demographics Survey

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

13



CLOUD CREDITS are being dispersed on request

AWS bills monthly, with charges applied to the credit card (or credit balance) on the last day of the month, for the month's charges

END OF MONTH:

CHECK YOUR CLOUD BILL
at least a few days before the end of the month

Billing Alarms – can be configured to generate email when there is a charge – can generate email if charges exceed \$0.01

With cloud credits, there should be no monthly charges

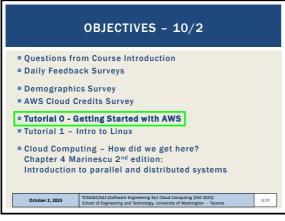
PROBLEM: when accounts have a credit balance, they do not generate billing alarms for high service usage – it is necessary to manually inspect service usage

October 2, 2025

| TOSSIGNITIES of Contracting (not Cound Companing [call 2025] school of Engineering and Enchnology University of Washington - Taxoma

| Tostal Contracting | Tostal 2025 | Tost

15



OBJECTIVES - 10/2

Questions from Course Introduction
Daily Feedback Surveys
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

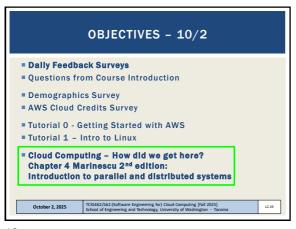
October 2, 2025

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

17 18

Slides by Wes J. Lloyd L2.3

14



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
SIMD architectures, vector processing, multimedia extensions
Graphics processing units
Speed-up, Amdahl's Law, Scaled Speedup
Properties of distributed systems
Modularity

Cotober 2, 2025

CSSIGS/SSZ: (Software Engineering for) Chood Computing [fall 2025]
School of Engineering and Technology, University of Washington - Tacoma

12.20

19 20



.1

• • • • • • • • • • • • • • • • • • • •	PER THREADIN	G
supporting multip	vide multiple instructio ble execution threads, t ns to a single CPU core	usually 2
■ Two hyper-threads		
are not equivalent	4770 with HTT Vs. 4670 without HTT	F - 25% improvement w/ HTT
are not equivalent to (2) CPU cores	4770 with HTT Vs. 4670 without HTT CPU Mark Relative to Tr As of 7th of February 2014 - Higher re	op 10 Common CPUs
to (2) CPU cores	CPU Mark Relative to To As of 7th of February 2014 - Higher re	op 10 Common CPUs esuits represent better performance
to (2) CPU cores • i7-4770 and i5-476	CPU Mark Relative to To As of 7th of February 2014 - Higher re	op 10 Common CPUs esuits represent better performance
to (2) CPU cores i7-4770 and i5-476 same CPU, with an	CPU Mark Relative to T As of 7th of February 2014 - Higher re	op 10 Common CPUs esuits represent better performance
to (2) CPU cores • i7-4770 and i5-476	CPU Mark Relative to T As of 7th of February 2014 - Higher re d bits Core (7-3770 @ 3.400fe; list Core (7-3770 @ 3.400fe;	op 10 Common CPUs esuits represent better performance
to (2) CPU cores i7-4770 and i5-476 same CPU, with an without HTT	CPU Mark Relative to T As of 7th of February 2014 - Higher re that Core 17-4770 @ 3-400Hz that Core 17-4770 @ 3-500Hz that Core 17-3770 @ 3-400Hz AMD 17-4390 @ 19-400Hz	op 10 Common CPUs esuits represent better performance
to (2) CPU cores i i7-4770 and i5-476 same CPU, with an without HTT Example: →	CPU Mark Relative to T As of 7th of February 2014 - Higher re d bits Core (7-3770 @ 3.400fe; list Core (7-3770 @ 3.400fe;	op 10 Common CPUs esuits represent better performance
to (2) CPU cores i7-4770 and i5-476 same CPU, with an without HTT	CPU Mark Relative to T As of 7th of February 2014 - Higher re black Cover 47.479 (3.4000+ black Cover 47.479 (3.4000+ black Cover 47.479 (3.4000+ black Cover 47.479 (3.4000+ AMD 74.800 (3.4000+ black Cover 47.479 (3.4000+ blac	op 10 Common CPUs suits represent better performance 9 565 9 542 9 419 9 551 9 201
to (2) CPU cores ■ i7-4770 and i5-476 same CPU, with an without HTT ■ Example: →	As of 710 of February 2014 - Higher re See 1 Se	op 10 Common CPUs ssuits represent better performance 9 505 90 442 92 413 99 501 90 505
to (2) CPU cores ■ i7-4770 and i5-476 same CPU, with an without HTT ■ Example: → hyperthreads add	CPU Mark Relative to T As of 7th of February 2014 - Higher re that Cave 17-17th 0 1-400ic bits Cave 17-17th 0 1-400ic	op 10 Common CPUs ssuts represent better performance #2 865 #3 542 #3 419 #3 055 #4 015 #5 0.59 #3 316

AMD'S 64-CORE 7NM CPUS ■ Epyc Rome CPUs Announced August 2019 ■ EPYC 7H12 requires liquid cooling 64 / 128 EPYC 7702 64 / 128 3.35 EPYC 7642 48 / 96 2.30 3.20 256 MB 225 W 48/96 192 MB EPYC 7552 2.20 3.30 200 W \$4025 October 2, 2025

23 24

Slides by Wes J. Lloyd L2.4

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

June 2023: 96 cores, 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)

X86_64 HOST SERVER VCPUS - AMAZON EC2 **INFRASTRUCTURE-AS-A-SERVICE CLOUD** Cloud server virtual CPUs/host (x86_64) Growth since 2006 - Amazon Compute Cloud (EC2) ■ 1st generation Intel: m1 - 8 vCPUs / host (Aug 2006) ■ 2nd generation Intel: m2 - 16 vCPUs / host (Oct 2009) ■ 3rd generation Intel: m3 - 32 vCPUs / host (Oct 2012) ■ 4th generation Intel: m4 - 48 vCPUs / host (June 2015) ■ 5th generation Intel: m5 - 96 vCPUs / host (Nov 2017) ■ 6th generation Intel: m6i - 128 vCPUs / host (Aug 2021) ■ 6th generation AMD: m6a - 192 vCPUs / host (Nov 2021) (Aug 2023) ■ 7th generation Intel: m7i - 192 vCPUs / host ■ 7th generation AMD: m7a - 192 vCPUs / host (Aug 2023) October 2, 2025

25

ARM64 HOST SERVER VCPUS - AMAZON EC2 INFRASTRUCTURE-AS-A-SERVICE CLOUD Cloud server virtual CPUs/host (ARM64) Launched in 2018 on the Amazon Compute Cloud (EC2) • 64-bit ARM CPUs designed by AWS subsidiary Annapurna Labs Lower energy consumption compared to x86-64 Fixed (non-variable) clock rates, No hyperthreading ■ Each new release - performance boost of ~ 30% Cost savings of ~20% less for ARM resources on AWS ■ 1st generation Graviton: a1 - 16 vCPUs / host (Nov 2018) ■ 2nd generation Graviton2: m6g-64 vCPUs/host (Dec 2019) AWS Lambda limited to Graviton2 ■ 3rd generation Graviton3: m7g-64 vCPUs/host (May 2022) 4th generation Graviton4: m8g-192 vCPUs/host(Sept 2024) TCSS462/562: (Software Engineering for) Cloud Computing [Fall 202 School of Engineering and Technology, University of Washington - Ta October 2, 2025

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

27

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

- October 2, 2025
- ICCS662/562: (Software Engineering for) Cloud Computing [Fall 2025]
- School of Engineering and Technology, University of Washington - Tacoma

CLOUD COMPUTING:
HOW DID WE GET HERE? - 4

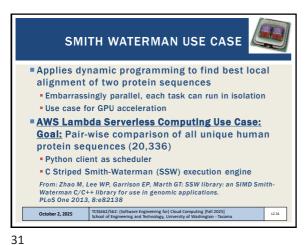
BIg Data requires massive amounts of compute resources

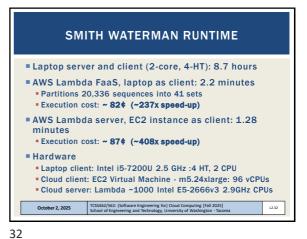
MAP - REDUCE
Single instruction, multiple data (SIMD)
Exploit data level parallelism
Bioinformatics example

29 30

Slides by Wes J. Lloyd L2.5

26





CLOUD COMPUTING:
HOW DID WE GET HERE? - 5

Compute clouds are large-scale distributed systems
Heterogeneous systems
Many services/platforms w/ diverse hw + capabilities
Homogeneous systems
Within a platform - illusion of identical hardware
Autonomous
Autonomous
Automatic management and maintenance- largely with little human intervention
Self organizing
User requested resources organize themselves to satisfy requests on-demand

33

CLOUD COMPUTING:
HOW DID WE GET HERE? - 6

Compute clouds are large-scale distributed systems
Infrastructure-as-a-Service (IaaS) Cloud
Provide VMs on demand to users
ec2Instances.Info (AWS EC2)

Clouds can consist of
Homogeneous hardware (servers, etc.)
Heterogeneous hardware (servers, etc.)
Which is preferable?

HARDWARE HETEROGENEITY If providing laaS, what are advantages/ disadvantages of using homogeneous hardware? Easier to provide same quality of service to end users Less performance variance Components with variable performance: CPUs, memory (speed differences), disks (SSDs, HDDs), network interfaces (caches?) Homogeneous hardware (servers): components are interchangeable As components fail, identical backups are immediately available Example: blade servers As clouds grow, why is HW homogeneity difficult to maintain? What are some advantages of using heterogeneous HW? TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tac October 2, 2025

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

October 2, 2025

| ICSS462/S62: (Software Engineering for) Cloud Computing [fail 2025]
| School of Engineering and Technology, University of Washington - Tacoma

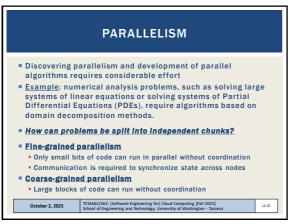
| ICSS462/S62: (Software Engineering for) Cloud Computing [fail 2025]
| School of Engineering and Technology, University of Washington - Tacoma

| ICSS462/S62: (Software Engineering for) Cloud Computing [fail 2025]
| School of Engineering and Technology, University of Washington - Tacoma

| ICSS462/S62: (Software Engineering for) Cloud Computing [fail 2025]

35 36

Slides by Wes J. Lloyd L2.6



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)

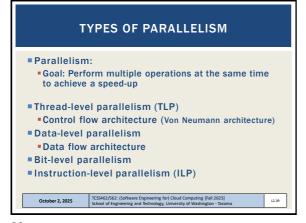
Shared memory: interactions can be implicit – must read the code!!

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

October 2, 2025

TCSM62/502: Esoftware Engineering from Cloud Computing [Fill 2025]
Shood of Engineering and Technology University of Volkheighein - Taccoma

37 38

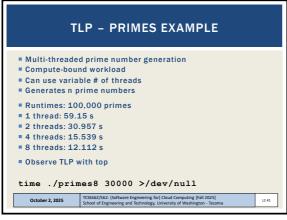


THREAD LEVEL PARALLELISM (TLP)

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...

October 2, 2025
School of Engineering Indi Gloud Computing [Fill 2025]
School of Engineering and Technology, University of Washington - Tacoms

39



Typical architecture used today - w/ multiple threads
Each thread runs a sequential program sequence
By John von Neumann (1945), also called the Von Neumann architecture
Dominant computer system architecture
Program counter (PC) determines next instruction to load into instruction register
Program execution is sequential

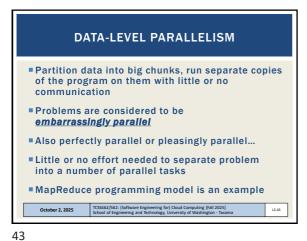
October 2, 2025

TCSS402/562: (Software Engineering for) Cloud Computing [Fall 2025]
Shool of Engineering and Technology, University of Versibington: Tecono

12.42

41 42

Slides by Wes J. Lloyd L2.7



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

October 2, 2025

October 2, 2025

School of Engineering and Technology, University of Washington - Tacoma

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

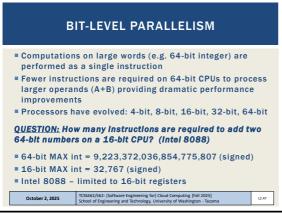
Utility of data flow architectures has been much less than envisioned

October 2, 2025

WE WILL RETURN AT 4:50PM

45 46

L2.45



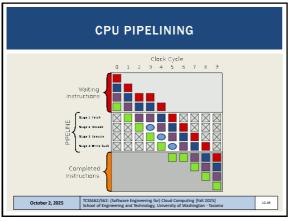
INSTRUCTION-LEVEL PARALLELISM (ILP)

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

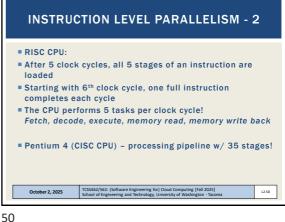
47 48

Slides by Wes J. Lloyd L2.8



49

51



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 Cotober 2, 2025 School of Engineering and Technology, University of Washington - Tacoma L234

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

Cotober 2, 2025

Cotober 2, 2025

Cotober 2, 2025

LESS COMPUTER

COTOBER 2015

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

(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

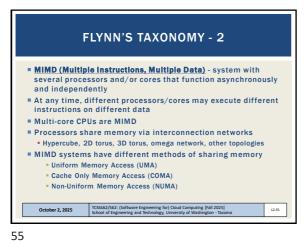
Developers can think sequentially; not worry about parallelism

October 2, 2025

TISSIGS/552: (Software Engineering for) Cloud Computing [Fall 2025]
Shool of Engineering and Technology, University of Washington - Taccoma

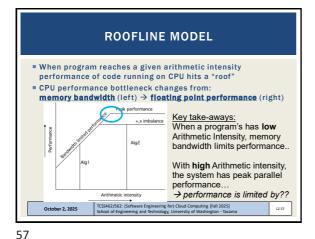
53 54

Slides by Wes J. Lloyd L2.9

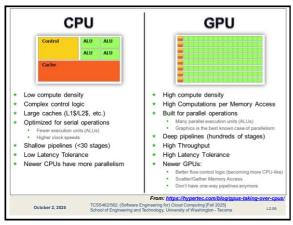


ARITHMETIC INTENSITY Arithmetic intensity: Ratio of work (W) to memory traffic r/w (Q) Example: # of floating-point ops per byte of data read Characterizes application scalability with SIMD support SIMD can perform many fast matrix operations in parallel High arithmetic intensity: Programs with dense matrix operations scale up nicely (many calcs vs memory RW, supports lots of parallelism) Low arithmetic intensity: Programs with sparse matrix operations do not scale well with problem size (memory RW becomes bottleneck, not enough ops!) October 2, 2025 L2.56

56



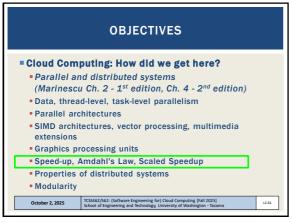
OBJECTIVES Cloud Computing: How did we get here? Parallel and distributed systems (Marinescu Ch. 2 - 1st edition, Ch. 4 - 2nd edition) Data, thread-level, task-level parallelism Parallel architectures SIMD architectures, vector processing, multimedia extensions Graphics processing units Speed-up, Amdahl's 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 - Taco October 2, 2025 L2.58



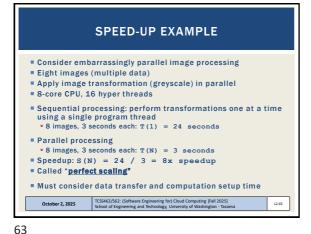
GRAPHICAL PROCESSING UNITS (GPUs) ■ GPU provides multiple SIMD processors ■ Typically 7 to 15 SIMD processors each 32,768 total registers, divided into 16 lanes (2048 registers each) ■ GPU programming model: single instruction, multiple thread Programmed using CUDA- C like programming language by NVIDIA for GPUs CUDA threads - single thread associated with each data element (e.g. vector or matrix) Thousands of threads run concurrently TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Taco October 2, 2025 L2.60

59 60

Slides by Wes J. Lloyd L2.10



61



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
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

Cotober 2, 2025

Cotober 2, 2025

CSSAGO/SOZ: (Software Engineering (et) Cloud Computing [fall 2025]
School of Engineering and Technology, University of Washington - Tacoma

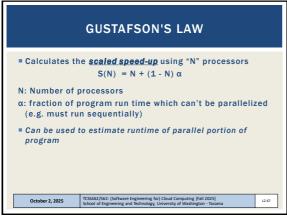
03

AMDAHL'S LAW EXAMPLE Program with two independent parts Part A is 75% of the execution time Part B is 25% of the execution time ■ Part B is made 5 times faster with parallel computing Estimate the percent improvement of task execution Original Part A is 3 seconds, Part B is 1 second ■ N=5 (speedup of part B) f=.25 (only 25% of the whole job (A+B) will be sped-up) ■ S=1 / ((1-f) + f/S) ■ S=1 / ((.75) + .25/5) ■ S=1.25 ■ % improvement = 100 * (1 - 1/1.25) = 20% TCSS462/562: (Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Taco October 2, 2025

65 66

Slides by Wes J. Lloyd L2.11

62



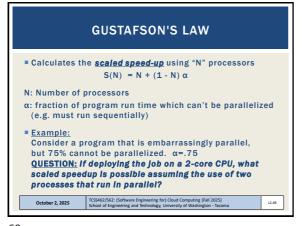
GUSTAFSON'S LAW

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 $\sigma = \text{sequential time}$, $\pi = \text{parallel time}$ Our Amdahl's example: $\sigma = 3s$, $\pi = 1s$, $\alpha = .75$

67



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) = ?**October 2, 2025

**TCSS462/S62: (Software Engineering for) Good Computing (Fall 2025)
School of Engineering and Technology, University of Washington -Taxoma

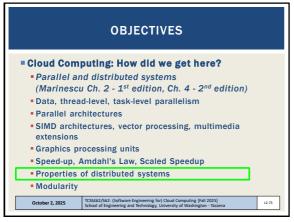
69

GUSTAFSON'S EXAMPLE OUESTION: What is the maximum theoretical speed-up on a 2-core CPU? $S(N) = N + (1 - N) \alpha$ N=2, α= For 2 CPUs, speed up is 1.25x S(N) = 3S(N) = ?For 16 CPUs, speed up is 4.75x ■ What is the ma $S(N) = N + (1 - N) \alpha$ $N=16. \alpha = .75$ S(N) = 16 + (1 - 16).75S(N) = ?October 2, 2025 L2.71

71 72

Slides by Wes J. Lloyd L2.12

68



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

Cotober 2, 2025

Cotober 3, 2025

Cotober 2, 2025

Cotober 3, 2025

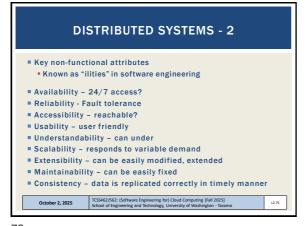
Cotober 2, 2025

Cotober 3, 2025

Cotober 4, 2025

Coto

73 74



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

October 2, 2005

October 2, 2005

October 2, 2005

**TOSHAGI/SEZ: (Software Engineering for Good Computing [Fail 2025])

October 2, 2005

**TOSHAGI/SEZ: (Software Engineering for Good Computing [Fail 2025])

October 2, 2005

**TOSHAGI/SEZ: (Software Engineering for Good Computing [Fail 2025])

**TOSHAGI/SEZ: (Software Engineering for Good Computing [Fail 2025])

75

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

October 2, 2025

| ICCSM62/S62: (Software Engineering for) Cloud Computing [Fall 2025] | Chool of Engineering and Technology, University of Washington - Taxoma | 12.77

TYPES OF MODULARITY

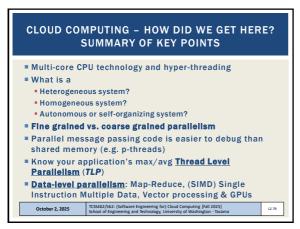
Soft modularity: TRADITIONAL
Divide a program into modules (classes) that call each other and communicate with shared-memory
A procedure calling convention is used (or method invocation)
Enforced modularity: CLOUD COMPUTING
Program is divided into modules that communicate only through message passing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

October 2, 2025

TCSS62/56: Clothware Engineering for Cloud Computing [Fall 2025]
TCSS62/56: Clothware Engineering for Cloud Computing [Fall 2025]

77 78

Slides by Wes J. Lloyd L2.13



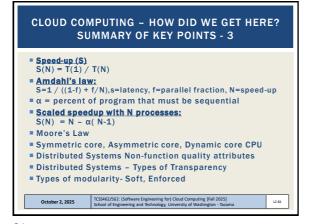
CLOUD COMPUTING - HOW DID WE GET HERE?
SUMMARY OF KEY POINTS - 2

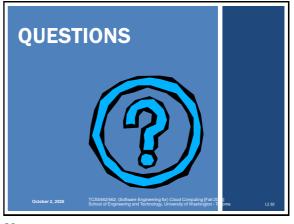
Bit-level parallelism
Instruction-level parallelism (CPU pipelining)
Fiynn's taxonomy: computer system architecture classification
SISD - Single Instruction, Single Data (modern core of a CPU)
SIMD - Single Instruction, Multiple Data (Data parallelism)
MIMD - Multiple Instruction, Multiple Data
MISD 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 2, 2025

TCSS62/SG2: (Software Engineering for) Cloud Computing [Fall 2025]
School of Engineering and Technology, University of Washington - Tacoma

79 80





81 82

Slides by Wes J. Lloyd L2.14