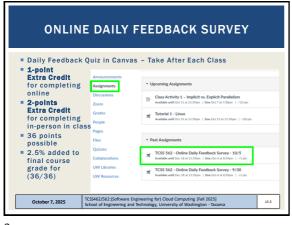


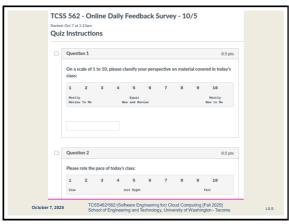
OBJECTIVES - 10/7 Questions from 10/2 ■ Tutorial 0, Tutorial 1, Tutorial 2 Cloud Computing - How did we get here? (Marinescu Ch. 2 - 1st edition, Ch. 4 - 2nd edition) ■ Class Activity 1 - Implicit vs Explicit Parallelism SIMD architectures, vector processing, multimedia Graphics processing units Speed-up, Amdahl's Law, Scaled Speedup ■ Properties of distributed systems Modularity October 7, 2025 L3.2

2



WARNING ■ DO NOT SUBMIT BOTH A PAPER AND AN ONLINE SURVEY OR YOU WILL LOOSE POINTS **CANVAS WILL AUTOMATICALLY REPLACE THE PAPER SURVEY** SCORE (2 PTS) WITH THE ONLINE SURVEY (1 PT) * * COMPLETE ONLY ONE SURVEY FOR EACH CLASS SESSION * ■ WE WILL NOT BE ABLE TO DUPLICATE CHECK SURVEYS FOR EACH CLASS SESSION AND MAKE CORRECTIONS October 7, 2025 L3.4

3

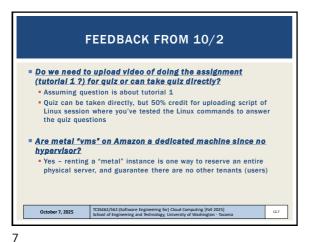


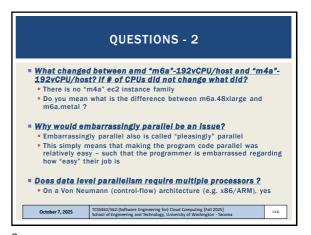
5

MATERIAL / PACE Please classify your perspective on material covered in today's class (46 respondents, 38 in-person, 8 online): ■ 1-mostly review, 5-equal new/review, 10-mostly new **Average - 7.24** (**1-** previous 7.20) Please rate the pace of today's class: ■ 1-slow, 5-just right, 10-fast ■ <u>Average - 5.20</u> (<u>↑ - previous 5.16</u>) October 7, 2025 L3.6

6

Slides by Wes J. Lloyd L3.1





9

DEMOGRAPHICS SURVEY

Please complete the ONLINE demographics survey:

We have received 41 responses so far.
We are waiting on ~12 responses.

https://forms.gle/QNUW2hUV7fR7BDmv7

Linked from course webpage in Canvas:

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

105462/862:(Software Engineering for) Choud Computing [fail 2025]
School of Engineering and Technology, University of Washington-Tacoma

AWS CLOUD CREDITS SURVEY

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

https://forms.gle/Y4IWvBRFVLRPnPX37.

Linked from course webpage in Canvas:

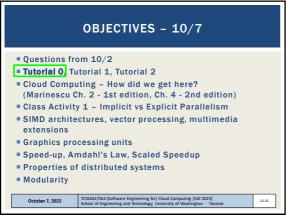
http://faculty.washington.edu/wlloyd/courses/tcss562/announcements.html

October 7, 2025

| TC5562/562: (Software Engineering for) Cloud Computing [Fall 2025]
| School of Engineering and Technology, University of Washington - Tacoma

11 12

Slides by Wes J. Lloyd L3.2



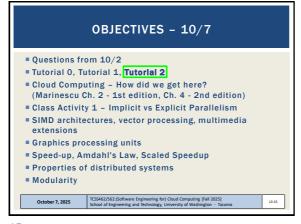
OBJECTIVES - 10/7

Questions from 10/2
Tutorial 0 Tutorial 1, Tutorial 2
Cloud Computing - How did we get here?
(Marinescu Ch. 2 - 1st edition, Ch. 4 - 2nd edition)
Class Activity 1 - Implicit vs Explicit Parallelism
SIMD architectures, vector processing, multimedia extensions
Graphics processing units
Speed-up, Amdahl's Law, Scaled Speedup
Properties of distributed systems
Modularity

October 7, 2025

TESS462/S62/Software Engineering for Cloud Computing [Fall 2025]
School of Engineering and Technology, University of Washington - Tacoma

13



OBJECTIVES - 10/7

Questions from 10/2
Tutorial 0, Tutorial 1, Tutorial 2

Cloud Computing - How did we get here?
(Marlnescu Ch. 2 - 1st edition, Ch. 4 - 2nd edition)

Class Activity 1 - Implicit vs Explicit Parallelism

SIMD architectures, vector processing, multimedia extensions

Graphics processing units

Speed-up, Amdahl's Law, Scaled Speedup

Properties of distributed systems

Modularity

October 7, 2025

INSEED SCALE SCHOMANIE Engineering for) Cloud Computing [fail 2025]
School of Engineering and Technology, University of Washington - Tacoma

15

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
Autonatic management and maintenance- largely with little human intervention
Self organizing
User requested resources organize themselves to satisfy requests on-demand

Cotober 7, 2025

Lattrian American Cotober 1 (Note Computing (Mail 2025) (Stobal of Engineering and Technology, University of Washington - Taxonsa)

Lattrian Cotober 7, 2025

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?

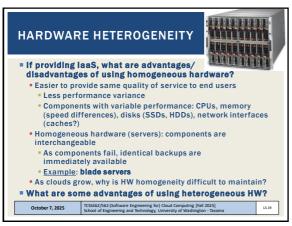
Ctober 7, 2025

TCSS462/S621/Software Engineering for Cloud Computing [Fall 2025]
School of Engineering and Technology, University of Washington - Tacoma

17 18

Slides by Wes J. Lloyd L3.3

14



OBJECTIVES - 10/7

Questions from 10/2

Tutorial 0, Tutorial 1, Tutorial 2

Cloud Computing - How did we get here?
(Marinescu Ch. 2 - 1st edition, Ch. 4 - 2nd edition)

Class Activity 1 - Implicit vs Explicit Parallelism

SIMD architectures, vector processing, multimedia extensions

Graphics processing units

Speed-up, Amdahl's Law, Scaled Speedup

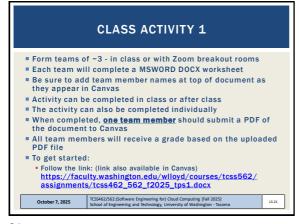
Properties of distributed systems

Modularity

October 7, 2025

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

19

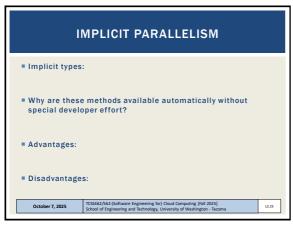


CLASS ACTIVITY 1

Solutions to be discussed..

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

21



EXPLICIT PARALLELISM

Explicit types:

Advantages:

Disadvantages:

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

23 24

Slides by Wes J. Lloyd L3.4

20

L3.26

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

PARALLELISM QUESTIONS 7. For bit-level parallelism, should a developer be concerned with the available number of virtual CPU processing cores when choosing a cloud-based virtual machine if wanting to obtain the best possible speed-up? (Yes / No) ■ 8. For instruction-level parallelism, should a developer be concerned with the physical CPU's architecture used to host a cloud-based virtual machine if wanting to obtain the best possible speed-up? (Yes / No) October 7, 2025 L3.25

PARALLELISM QUESTIONS - 2 9. An application developer measures the average and peak thread level parallelism (TLP) of an application prior to deployment on the AWS EC2. The developer measures an average TLP of 2.3, and a peak TLP of 7.3. The application is to be deployed using a compute-optimized (c-series) ec2 instance. Using resources online, such as the websites below, , propose a good virtual machine (ec2 type) that satisfies average TLP, and a second for satisfying peak TLP that does not under-provision or over-provision vCPUs for the TLP goal, in order to control costs. https://docs.aws.amazon.com/ec2/latest/instancetypes/ co.html https://instances.vantage.sh/ October 7, 2025

25 26

PARALLELISM QUESTIONS - 3 What is a good ec2 c-series instance for average TLP? Why is this instance good/sufficient for satisfying average What is a good ec2 c-series instance for peak TLP? Why is this instance good/sufficient for satisfying peak TLP? TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tar October 7, 2025 L3.27

OBJECTIVES - 10/7 Questions from 10/2 ■ Tutorial 0, Tutorial 1, Tutorial 2 Cloud Computing - How did we get here? (Marinescu Ch. 2 - 1st edition, Ch. 4 - 2nd edition) Class Activity 1 - Implicit vs Explicit Parallelism SIMD architectures, vector processing, multimedia extensions Graphics processing units Speed-up, Amdahl's Law, Scaled Speedup Properties of distributed systems October 7, 2025 L3.28

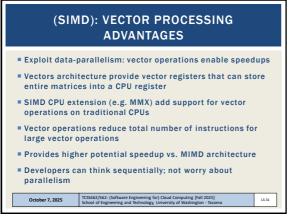
27

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: (Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacor October 7, 2025

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 October 7, 2025 L3.30

29 30

Slides by Wes J. Lloyd L3.5



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

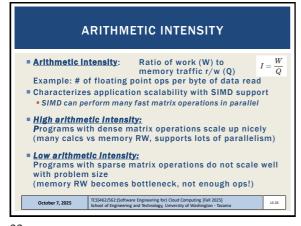
MIMD systems have different methods of sharing memory

Uniform Memory Access (UMA)

Cache Only Memory Access (COMA)

Non-Uniform Memory Access (NUMA)

31 32



33

OBJECTIVES - 10/7

Questions from 10/2
Tutorial 0, Tutorial 1, Tutorial 2
Cloud Computing - How did we get here?
(Marinescu Ch. 2 - 1st edition, Ch. 4 - 2nd edition)
Class Activity 1 - Implicit vs Explicit Parallelism
SIMD architectures, vector processing, multimedia extensions
Graphics processing units
Speed-up, Amdahl's Law, Scaled Speedup
Properties of distributed systems
Modularity

October 7, 2025

INCS462/S62/Scritwure Engineering for) Cloud Computing [fall 2025]
School of Engineering and Technology, University of Washington - Tacoma

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

Ctober 7, 2025

Ctober 7, 2025

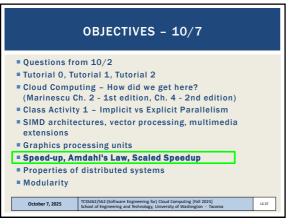
Ctober 7, 2025

Ctober 7, 2025

1236

35 36

Slides by Wes J. Lloyd L3.6



PARALLEL COMPUTING

■ Parallel hardware and software systems allow:

■ Solving problems needing resources not available on a single system.

■ Reduced time required to obtain solution

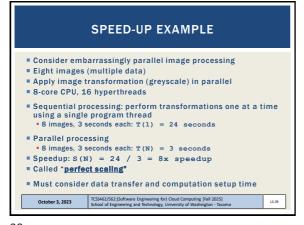
■ The speed-up (S) measures effectiveness of parallelization:

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

T(1) → execution time of total sequential computation T(N) → execution time for performing N parallel computations in parallel
computations in parallel

October 3, 2023 | TCSSE2/262/56thware Engineering for) Cloud Computing [Fall 2025] | School of Engineering and Technology, University of Wishington - Tacoms

37 38



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

39

 $Speed-up formula \Rightarrow S = \frac{1}{(1-f) + \frac{f}{N}}$ = S = theoretical speedup of the whole task $= f = \text{fraction of work that is parallel} \qquad (ex. 25\% \text{ or } 0.25)$ $= N = \text{proposed speed up of the parallel part} \qquad (ex. 5 \text{ times speedup})$ $= \% \text{ improvement} \qquad \text{of task execution} \qquad = 100 * (1 - (1/S))$ = Using Amdahl's law, we can find the maximum possible speed-up (S) for a given scenario (e.g. ~8X)... $\text{October 3, 2023} \qquad \text{ICSS62/562 isoftware Engineering fool Cloud Computing [Fall 2025]} \qquad \text{School of Engineering and Technology, University of Washington - Tacoma}$

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

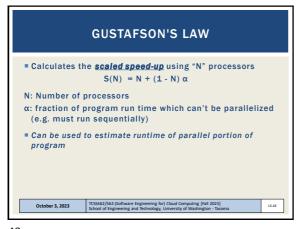
N = 5 (speedup of part B)
F= .25 (only 25% of the whole job (A+B) will be sped-up)
S=1 / ((1-f) + f/S)
S=1 / ((1-f) + .25/5)
S=1.25 (speed up is 1.25x faster)
with improvement = 100 * (1 - 1/1.25) = 20%

October 3, 2023

TCSS462/S621/Sotrouse Engineering for/ Cloud Computing [Fall 2025]
School of Engineering and Technology, University of Washington-Tacoma

41 42

Slides by Wes J. Lloyd L3.7



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

Cotober 3, 2023

CSS602/562/Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma

43 44

```
GUSTAFSON'S LAW

Calculates the scaled speed-up using "N" processors S(N) = N + (1 - N) \alpha

N: Number of processors \alpha: fraction of program run time which can't be parallelized (e.g. must run sequentially)

Example:
Consider a program that is embarrassingly parallel, but 75% cannot be parallelized. \alpha=.75
QUESTION: If deploying the job on a 2-core CPU, what scaled speedup is possible assuming the use of two processes that run in parallel?

October 3, 2023

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

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 3, 2023 | TCSS42/S62/Software Engineering for Cloud Computing [Fail 2025] | School of Engineering and Technology, University of Washington - Tacoma

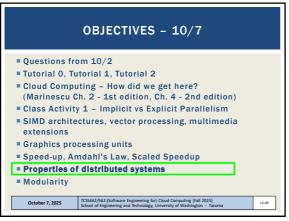
45

```
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) =
 S(N) = ?
              For 16 CPUs, speed up is 4.75x
What is the maximum theoretical speed-up on a 16-core CPU?
 S(N) = N + (1 - N) \alpha
 N=16, \alpha=.75
 S(N) = 16 + (1 - 16).75
 S(N) = ?
  October 3, 2023
                                                           L3.47
```

MOORE'S LAW Transistors on a chip doubles approximately every 1.5 years ■ CPUs nov ■ Power di What kind of processor are modern to heat removal Intel CPUs? Transiti **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 Dynamic core processor - processing resources and speed can be dynamically configured among cores Observation: asymmetric processors offer a higher speedup TCSS462/562:(Software Engineering for) Cloud Computing [Fall 2025] School of Engineering and Technology, University of Washington - Tacoma October 3, 2023

47 48

Slides by Wes J. Lloyd L3.8



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 7, 2025

Totoska/JSc2/Johnwe Engineering for Courd Computing [Isla 2025]

Totoska/JSc2/Johnwe Engineering for Courd Computing Isla 2025]

49 50



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 7, 2025

October 7, 2025

Costober 7, 202

51

OBJECTIVES - 10/7

Questions from 10/2

Tutorial 0, Tutorial 1, Tutorial 2

Cloud Computing - How did we get here?
(Marinescu Ch. 2 - 1st edition, Ch. 4 - 2nd edition)

Class Activity 1 - Implicit vs Explicit Parallelism

SIMD architectures, vector processing, multimedia extensions

Graphics processing units

Speed-up, Amdahl's Law, Scaled Speedup

Properties of distributed systems

Modularity

October 7, 2025

CSS652/652/Software Engineering for) Good Computing (Tail 2025)
School of Engineering and Technology, University of Washington - Tacoma

TYPES OF MODULARITY

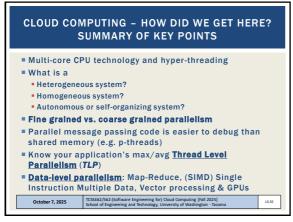
Soft modularity: TRADITIONAL
Divide a program into modules (classes) that call each other and communicate with shared-memory
A procedure calling convention is used (or method invocation)

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!

TCSS462/5621;Software Engineering for Cloud Computing [Fail 2025]
School of Engineering and Technology, University of Washington - Tacoma

53 54

Slides by Wes J. Lloyd L3.9



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

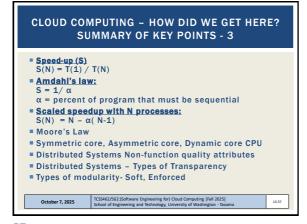
Bit-level parallelism
Instruction-level parallelism (CPU pipelining)
Flynn's taxonomy: computer system architecture classification
SISD - Single Instruction, Single Data (modern core of a CPU)
SIMD - Single Instruction, Multiple Data parallelism)
MIMD - Multiple Instruction, Multiple Data
MISD is RARE; application for fault tolerance...
Arlthmetic Intensity: ratio of calculations vs memory RW
Roofline model:
Memory bottleneck with low arithmetic intensity
SIMD and Vector processing supported by many large registers

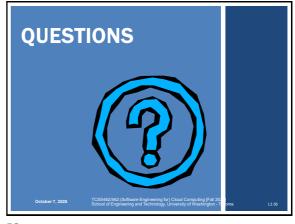
Cotaber 7, 2025

Cotaber 7, 2025

Interior Count Computing [fall 2025]
School of Engineering and Technology, University of Washington - Taxoma

55 56





57 58

Slides by Wes J. Lloyd L3.10