

TCSS 562: SOFTWARE ENGINEERING FOR CLOUD COMPUTING

Cloud Computing:
How did we get here? – cont'd

Introduction to Cloud Computing

Wes J. Lloyd
School of Engineering and Technology
University of Washington - Tacoma



FEEDBACK FROM 10/1

- Parallel threads accessing the same data
- Threads and how they affect the running time of a program
 - When 2 threads run, don't they have to take turns?
 - How is runtime cut in half with 2 threads if they have to take turns?

October 3, 2018

TCSS562: Software Engineering for Cloud Computing [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L3.2

FEEDBACK - 2

- **Roofline model - relates to “arithmetic intensity”**

- **Arithmetic Intensity:**

Ratio of work vs. memory traffic (RW)

$$I = \frac{W}{Q}$$

I=Arithmetic intensity; W=Work; Q=memory traffic

- **Roofline model:**

When performance bottleneck changes from
memory to GPU/CPU

- Threshold: when arithmetic intensity of code is high
 - Number of operations outweighs memory traffic

October 3, 2018

TCSS562: Software Engineering for Cloud Computing [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L3.3

ARITHMETIC INTENSITY

- **Arithmetic Intensity:** Ratio of work (W) to
memory traffic r/w (Q)

$$I = \frac{W}{Q}$$

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!)

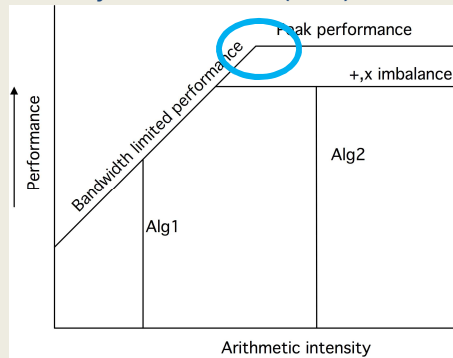
October 3, 2018

TCSS562: Software Engineering for Cloud Computing [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L3.4

ROOFLINE MODEL

- When program reaches a given arithmetic intensity performance of code running on CPU hits a “roof”
- CPU performance bottleneck changes from:
memory bandwidth (left) → floating point performance (right)



Key take-aways:

When a program's has low Arithmetic Intensity, memory bandwidth limits performance..

With high Arithmetic intensity, the system has peak parallel performance...

October 3, 2018

TCSS562: Software Engineering for Cloud Computing [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L3.5

FEEDBACK - 3

- Recommended material to refresh and prep for class
 - Textbooks, read chapters listed on Schedule Page
 - Book #1: Cloud Computing: Concepts, Technology, and Architecture
 - Book #2: Cloud Computing: Theory and Practice
 - 1st Edition is online as a PDF

October 3, 2018

TCSS562: Software Engineering for Cloud Computing [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L3.6

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

- Summary of the key points from 10/1:
 - 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

October 3, 2018

TCCS562: Software Engineering for Cloud Computing [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L3.7

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

October 3, 2018

TCCS562: Software Engineering for Cloud Computing [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L3.8

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

October 3, 2018

TCSS562: Software Engineering for Cloud Computing [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L3.9

ARITHMETIC INTENSITY

- **Arithmetic intensity:** Ratio of work (W) to memory traffic r/w (Q) $I = \frac{W}{Q}$

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!)

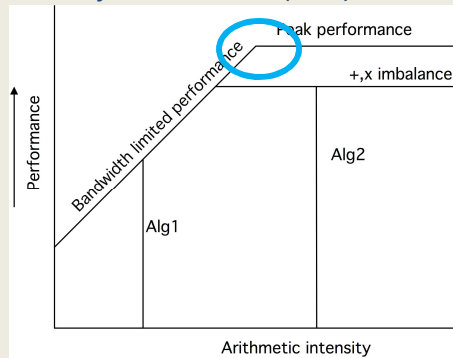
October 3, 2018

TCSS562: Software Engineering for Cloud Computing [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L3.10

ROOFLINE MODEL

- When program reaches a given arithmetic intensity performance of code running on CPU hits a “roof”
- CPU performance bottleneck changes from: memory bandwidth (left) → floating point performance (right)



Key take-aways:

When a program's has low Arithmetic Intensity, memory bandwidth limits performance..

With high Arithmetic intensity, the system has peak parallel performance...

October 3, 2018

TCSS562: Software Engineering for Cloud Computing [Fall 2018]
 School of Engineering and Technology, University of Washington - Tacoma

L3.11

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

October 3, 2018

TCSS562: Software Engineering for Cloud Computing [Fall 2018]
 School of Engineering and Technology, University of Washington - Tacoma

L3.12

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:

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

$T(1)$ → execution time of total sequential computation

$T(N)$ → execution time for performing N parallel computations in parallel

October 3, 2018

TCSS562: Software Engineering for Cloud Computing [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L3.13

SPEED-UP EXAMPLE

- Consider embarrassingly parallel image processing
- Eight images (multiple data)
- Apply image transformation (greyscale) in parallel
- 8-core CPU, 16 hyperthreads
- 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: $S(N) = 24 / 3 = 8\times$ speedup
- Called “**perfect scaling**”
- Must consider data transfer and computation setup time

October 3, 2018

TCSS562: Software Engineering for Cloud Computing [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L3.14

AMDAHL'S LAW

- Portion of computation which cannot be parallelized determines the overall speedup
 - For an embarrassingly parallel job of fixed size
 - Assuming no overhead for distributing the work, and a perfectly even work distribution
- α : fraction of program run time which can't be parallelized (e.g. must run sequentially)
- Maximum speedup is:
$$S = 1 / \alpha$$
 - **Example:**
Consider a program where 25% cannot be parallelized
Q: What is the maximum possible speedup of the program?

October 3, 2018

TCSS562: Software Engineering for Cloud Computing [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L3.15

GUSTAFSON'S LAW

- Calculates the scaled speed-up using “N” processors
$$S(N) = N + (1 - \alpha) N$$
- 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 except for 25% that cannot be parallelized. $\alpha=.25$
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, 2018

TCSS562: Software Engineering for Cloud Computing [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L3.16

GUSTAFSON'S EXAMPLE

- **QUESTION:**

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

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

$$N=2, \alpha=.25$$

$$S(N) = 2 + (1 - 2) .25$$

$$S(N) = ?$$

- What is the maximum theoretical speed-up on a **4-core CPU**?

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

$$N=4, \alpha=.25$$

$$S(N) = 4 + (1 - 4) .25$$

$$S(N) = ?$$

October 3, 2018

TCSS562: Software Engineering for Cloud Computing [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L3.17

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
 - Transition from: increasing clock rates → 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
- **Dynamic core processor** – processing resources and speed can be dynamically configured among cores
- **Observation: asymmetric processors offer a higher speedup**

October 3, 2018

TCSS562: Software Engineering for Cloud Computing [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L3.18

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

October 3, 2018

TCCS562: Software Engineering for Cloud Computing [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L3.19

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

October 3, 2018

TCCS562: Software Engineering for Cloud Computing [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L3.20

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

October 3, 2018

TCSS562: Software Engineering for Cloud Computing [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L3.21

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 3, 2018

TCSS562: Software Engineering for Cloud Computing [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L3.22

QUESTIONS



October 3, 2018

TCSS562: Software Engineering for Cloud Computing [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L3.66