



This work is licensed under a Creative Commons Attribution-Noncommercial-Share Alike 3.0 United States See http://creativecommons.org/licenses/by-nc-sa/3.0/us/ for details



What is Cloud Computing?

- 1. Web-scale problems
- 2. Large data centers
- 3. Different models of computing
- 4. Highly-interactive Web applications

1. Web-Scale Problems

- Characteristics:
 - Definitely data-intensive
 - May also be processing intensive
- Examples:
 - Crawling, indexing, searching, mining the Web
 - "Post-genomics" life sciences research
 - Other scientific data (physics, astronomers, etc.)
 - Sensor networks
 - Web 2.0 applications
 - ...



How much data?

- Wayback Machine has 2 PB + 20 TB/month (2006)
- Google processes 20 PB a day (2008)
- "all words ever spoken by human beings" ~ 5 EB
- NOAA has ~1 PB climate data (2007)
- CERN's LHC will generate 15 PB a year (2008)











2. Large Data Centers

- Web-scale problems? Throw more machines at it!
- Clear trend: centralization of computing resources in large data centers
 - Necessary ingredients: fiber, juice, and space
 - What do Oregon, Iceland, and abandoned mines have in common?

Important Issues:

- Redundancy
- Efficiency
- Utilization
- Management











4. Web Applications

 A mistake on top of a hack built on sand held together by duct tape?

- What is the nature of software applications?
 - From the desktop to the browser
 - SaaS == Web-based applications
 - Examples: Google Maps, Facebook
- How do we deliver highly-interactive Web-based applications?
 - AJAX (asynchronous JavaScript and XML)
 - For better, or for worse...



- MapReduce: the "back-end" of cloud computing
 - Batch-oriented processing of large datasets
- Ajax: the "front-end" of cloud computing
 - Highly-interactive Web-based applications
- Computing "in the clouds"
 - Amazon's EC2/S3 as an example of utility computing



<section-header><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item>

This course is not for you...

- If you're not genuinely interested in the topic
- If you're not ready to do a lot of programming
- o If you're not open to thinking about computing in new ways

University of Maryland

- If you can't cope with uncertainly, unpredictability, poor documentation, and immature software
- If you can't put in the time

Otherwise, this will be a richly rewarding course!





Cloud Computing Zen

- Don't get frustrated (take a deep breath)...
 - This is bleeding edge technology
 - Those W\$*#T@F! moments
- Be patient...
 - This is the second first time I've taught this course
- Be flexible...
 - There will be unanticipated issues along the way
- Be constructive...
 - Tell me how I can make everyone's experience better











Things to go over...

- Course schedule
- Assignments and deliverables
- Amazon EC2/S3

The iSchool University of Maryland

Web-Scale Problems?

- Don't hold your breath:
 - Biocomputing
 - Nanocomputing
 - Quantum computing
 - ...
- It all boils down to...
 - Divide-and-conquer
 - Throwing more hardware at the problem

Simple to understand... a lifetime to master...



Different Workers

- Different threads in the same core
- Different cores in the same CPU
- Different CPUs in a multi-processor system
- Different machines in a distributed system



Choices, Choices, Choices

- Commodity vs. "exotic" hardware
- Number of machines vs. processor vs. cores
- Bandwidth of memory vs. disk vs. network
- Different programming models















Parallelization Problems

- How do we assign work units to workers?
- What if we have more work units than workers?
- What if workers need to share partial results?
- How do we aggregate partial results?
- How do we know all the workers have finished?
- What if workers die?

What is the common theme of all of these problems?



General Theme?

- Parallelization problems arise from:
 - Communication between workers
 - Access to shared resources (e.g., data)
- Thus, we need a synchronization system!
- This is tricky:
 - Finding bugs is hard
 - Solving bugs is even harder



Managing Multiple Workers

Difficult because

- (Often) don't know the order in which workers run
- (Often) don't know where the workers are running
- (Often) don't know when workers interrupt each other

• Thus, we need:

- Semaphores (lock, unlock)
- Conditional variables (wait, notify, broadcast)
- Barriers
- Still, lots of problems:
 - Deadlock, livelock, race conditions, ...
- Moral of the story: be careful!
 - Even trickier if the workers are on different machines



- Parallel computing has been around for decades
- Here are some "design patterns" ...









Rubber Meets Road

- From patterns to implementation:
 - pthreads, OpenMP for multi-threaded programming
 - MPI for clustering computing
 - ...

• The reality:

- Lots of one-off solutions, custom code
- Write you own dedicated library, then program with it
- Burden on the programmer to explicitly manage everything
- MapReduce to the rescue!
 - (for next time)

