# CSE 291: Operating Systems in Datacenters

Amy Ousterhout

Sept. 22, 2022

# Agenda for Today

- Brief introductions
- Introduction to OS in Datacenters
- Course logistics
- Questions to ask when reading a paper

# Introductions

A bit about me:
- Amy Ousterhout

  "oh"-"stir"-"howt"

- Please call me Amy!
- Assistant Professor in CSE
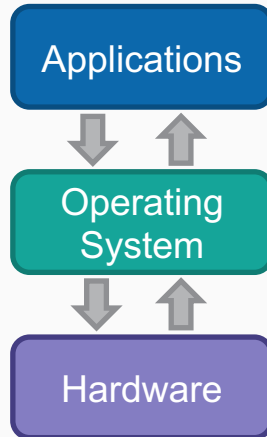- Research focus: resource efficiency in datacenters

Course TA:
- Anil Yelam

A bit about you!
- Your name
- Your background
- Why are you interested in operating systems and datacenters?

# Operating Systems in Datacenters

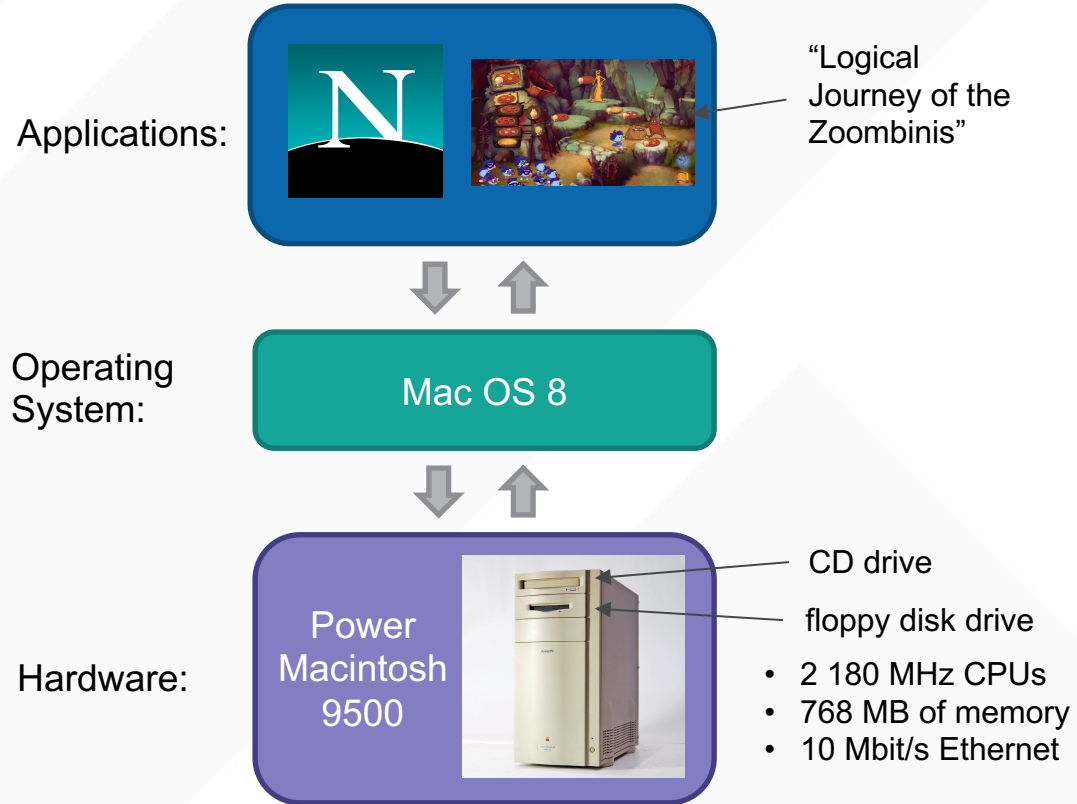# What is an Operating System?

- System software that:
    - Manages computer hardware and software
    - Provides services to computer programs
- Acts as the interface between hardware and applications

# Example Operating System: Mac OS 8

- Released in 1997
- Key OS features:
  - Processes and threads
  - Storage (disks)
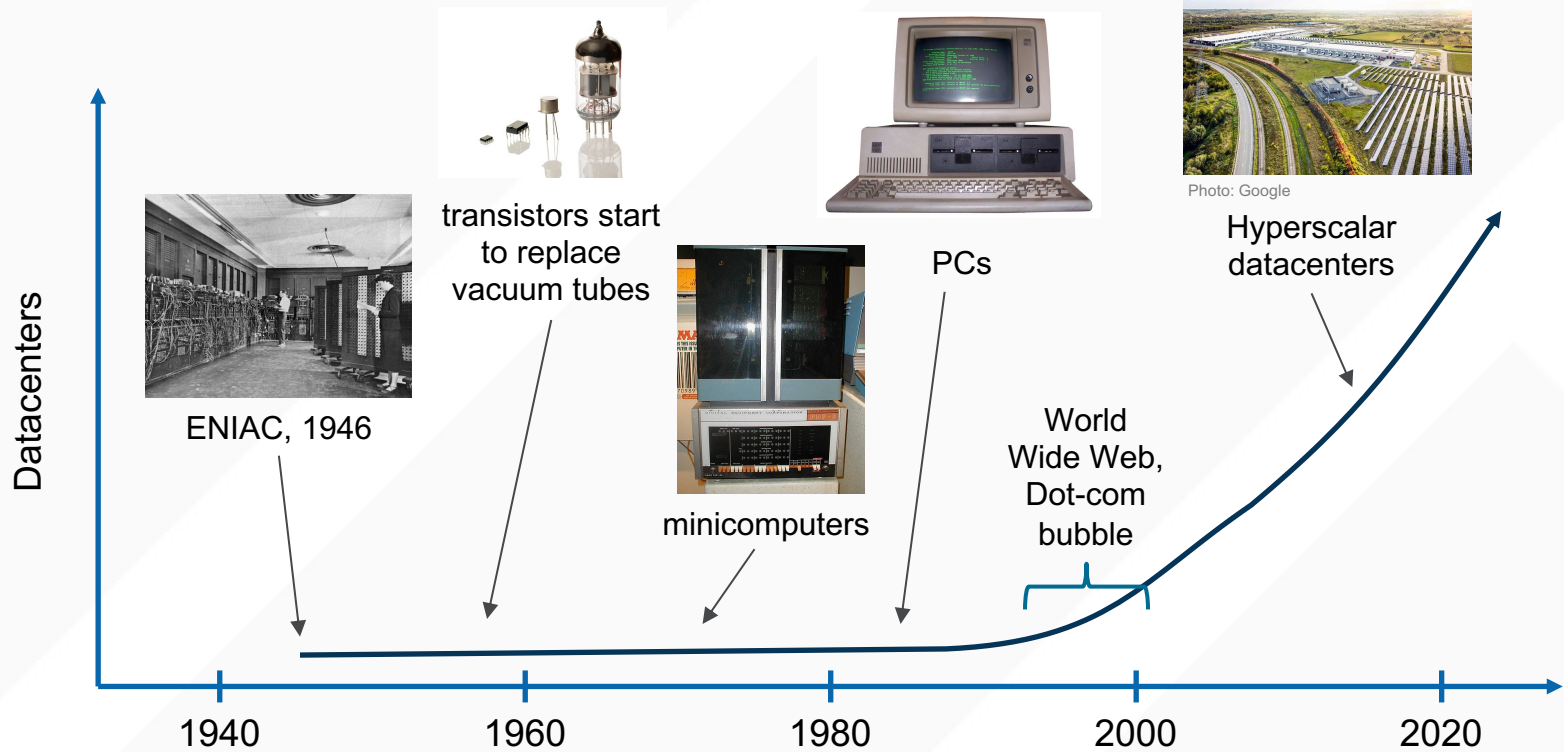  - Virtual memory
  - File system
  - Network stack

main components of an undergrad OS course

Applications:
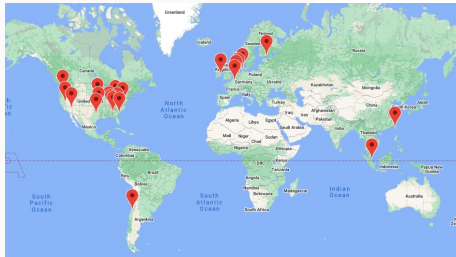
"Logical Journey of the Zoombinis"

Operating System:

Mac OS 8

Hardware:

Power Macintosh 9500

CD drive

floppy disk drive

- 2 180 MHz CPUs
- 768 MB of memory
- 10 Mbit/s Ethernet

# What is a Datacenter?

- Dedicated space that contains:
  - Computers
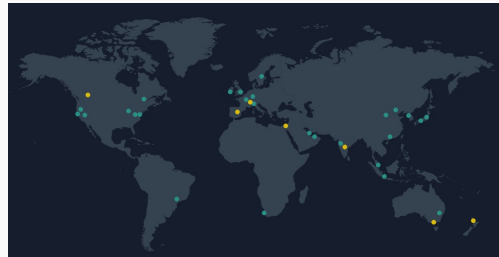  - Communication systems
  - Storage systems

# History of Datacenters



ENIAC, 1946

transistors start
to replace
vacuum tubes

minicomputers

PCs

World
Wide Web,
Dot-com
bubble

Photo: Google

Hyperscalar
datacenters

Datacenters

1940     1960     1980     2000     2020

# Datacenters Today

- Over 8,000 datacenters globally
- Over 2,600 datacenters in the U.S.
- Huge energy consumers – almost 2% of global energy use
  - Usually built near energy sources (hydroelectric, wind, solar)



Google's Datacenter
Locations



Amazon AWS Locations



Photo: Google

Google datacenter

# Inside a Datacenter

- Servers arranged into racks
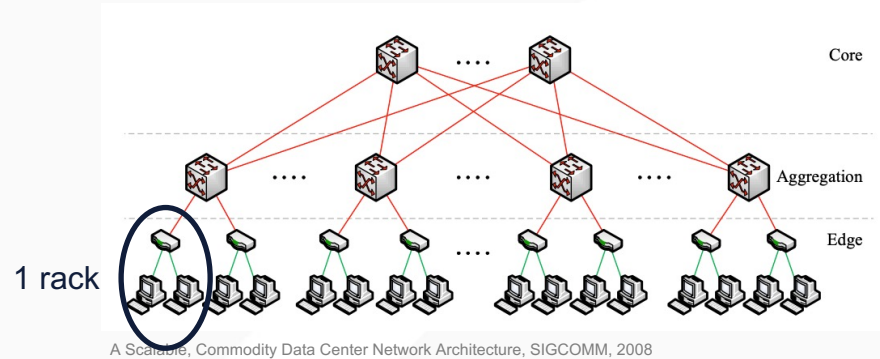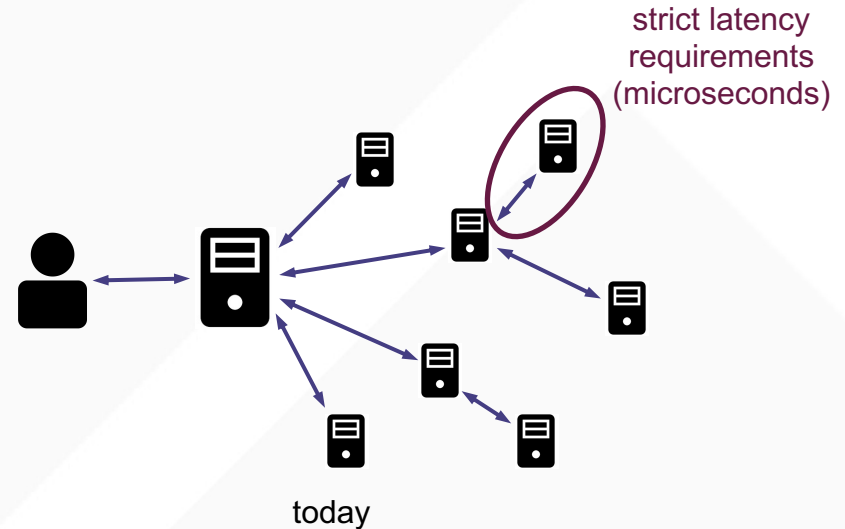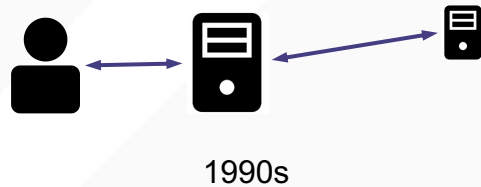- Racks connected by a hierarchical network topology



Photo: Google



A Scalable, Commodity Data Center Network Architecture, SIGCOMM, 2008
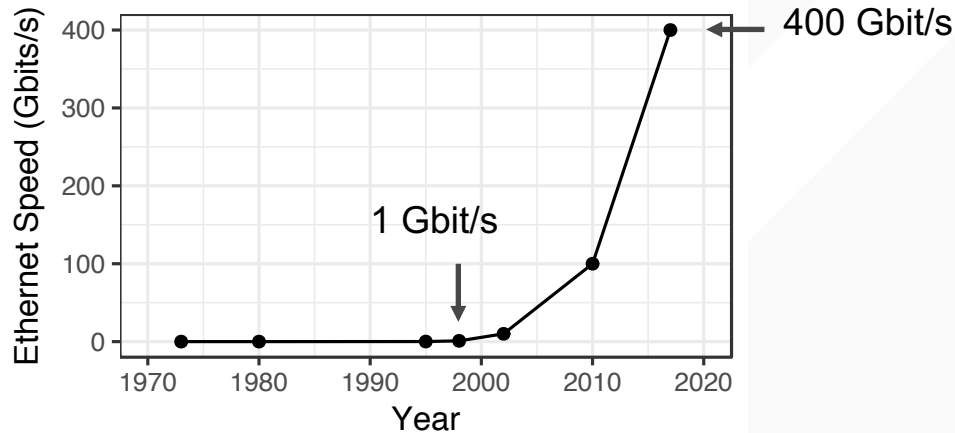
# Trend #1: Increasingly Complex Applications

- 1990s: static web pages served by a single server
- 2010: tens to hundreds of servers involved
  - Web search, social networks, etc.
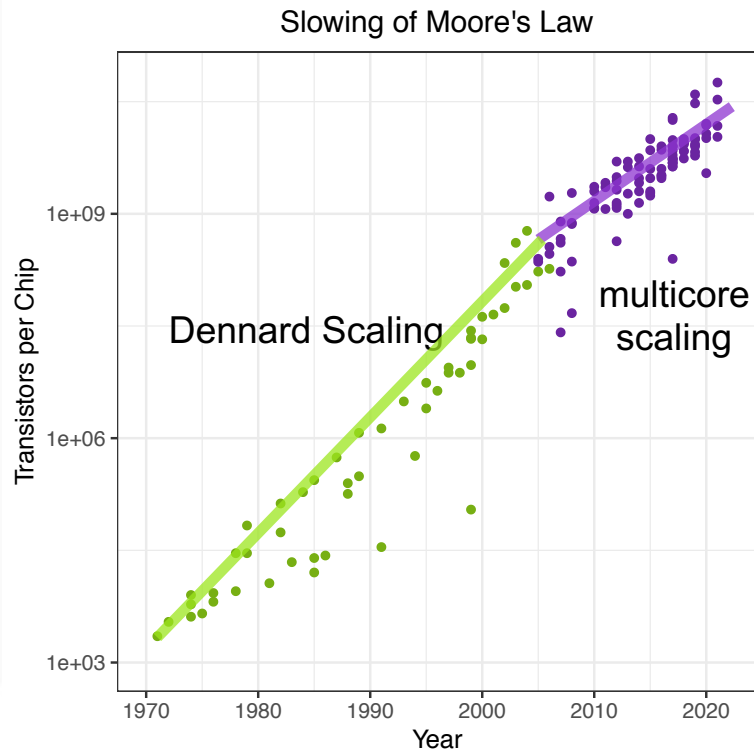- 2020: hundreds to thousands of servers involved



strict latency requirements (microseconds)

1990s

today

# Trend #2: Faster Networks

- Network bandwidth has increased 400x
- Network latencies have decreased too
    - Network latency = transmission + propagation + switching
    - Transmit 1500 bytes at 1 Gbit/s: 12 μs
    - Transmit 1500 bytes at 400 Gbit/s: 30 ns

servers are expected
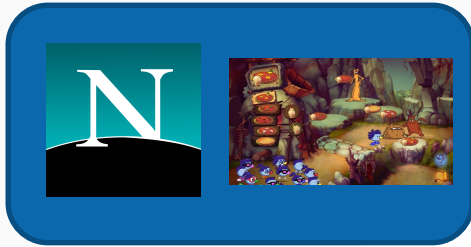to process large
amounts of network
traffic very quickly

# Trend #3: End of Moore's Law

- Increasing demand for compute
- Faster CPUs every few years!
- But, Moore's Law is ending
- Consequences:
  - More cores per server (multicore)
  - Move tasks to hardware with custom accelerators



Slowing of Moore's Law

# Operating Systems Requirements in Datacenters

Applications:

Applications tolerate ms-scale overheads

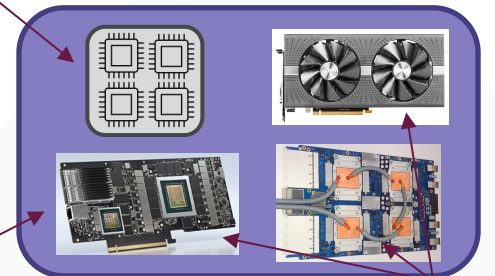Applications demand μs-scale overheads

Operating System:

Mac OS 8

This course!

Hardware:

Power Macintosh 9500

1-2 CPU cores

CD drive

Low I/O bandwidth

floppy disk drive

- 2 180 MHz CPUs
- 768 MB of memory
- 10 Mbit/s Ethernet

Multicore CPUs

I/O latency: 10s of ms

1990s

I/O: hundreds of Gbits/s, μs latency

today

GPUs, TPUs, SmartNICs

# Challenges Imposed on OS by Datacenters

- New hardware
    - Multicore
    - Fast networks
    - Heterogeneity (GPUs, TPUs, SmartNICs)
- New applications
    - Large and complex
    - Expect extremely low latency
- This course: how should Operating Systems adapt?

# Course Logistics

# Overview

- Graduate-level research-focused course
- The goals of the course are:
  - To learn about recent OS techniques to address datacenter challenges
  - To practice reading and discussing research papers
  - To conduct a research project

# Readings and Reviews

- We will read 1-2 papers per class
    - Everyone should read the papers ahead of time
    - Submit a short review by 11:59 pm the evening before
    - Details about reviews will be on Canvas    coming soon!
- Come prepared to discuss!
- Class format
    - Brief overview of each topic
    - Paper discussions

# Leading a Discussion

- Each student will lead 1 paper discussion
- Preparation:
    - Outline your discussion
    - Share with instructor at least 2 days before discussion
- No need for slides

# Research Project

- Open-ended research project
- Can work alone or in groups of 2-3 students
- You choose the topic
  - Broadly related to OS in datacenters
  - Implementation, experimental, algorithmic, theoretical
- How to pick a topic?
  - Propose your own idea
  - I will suggest some ideas for how to find a topic
- Computing platform
  - I recommend CloudLab

# Research Project Components

- ~1-page proposal, due 10/20
- Project presentations, in-class 11/29 and 12/1
- ~6-page project write-up, due 12/8
- We will meet with you throughout the quarter to check-in
- More details will be posted on Canvas

# Warm-Up Assignment

- Goals:
  - Show you how to use CloudLab
  - Give you some experience with RDMA and DPDK

# Grading

- There are no exams
- 15% paper reviews
- 15% class participation
- 10% discussion lead
- 10% warm-up assignment
- 50% project

# Course website

https://amyousterhout.com/cse291-fall22

| Date | Topics | Papers |
|------|--------|--------|
| Th 9/22 | Course overview | |
| Tu 9/27 | Multicore, intro to CloudLab | Multikernel (SOSP '09), CloudLab (ATC '19) - only first 2 sections |
| Th 9/29 | Network stacks | IX (OSDI '14), XDP (CoNEXT '18) |
| Tu 10/4 | RDMA and RPCs | FaRM (NSDI '14) |
| Th 10/6 | RDMA and RPCs | eRPC (NSDI '19), PRISM (SOSP '21) |
| Tu 10/11 | Congestion control | Homa (SIGCOMM '18), Swift (SIGCOMM '20) |
| Th 10/13 | CPU scheduling | Killer Microseconds (CACM '17), Shenango (NSDI '19) |
| Tu 10/18 | CPU scheduling | ghOSt (SOSP '21) |
| Th 10/20 | Performance diagnosis | NSight (NSDI '22), Collie (NSDI '22) |
| Tu 10/25 | Datacenter tax | Warehouse-scale computers (ISCA '15) |
| Th 10/27 | SmartNICs | AccelNet (NSDI '18) |
| Tu 11/1 | SmartNICs | iPipe (SIGCOMM '19), nanoPU (OSDI '21) |
| Th 11/3 | GPUs | PTask (SOSP '11) |
| Tu 11/8 | TPUs | TensorFlow (OSDI '16) |
| Th 11/10 | FPGAs | AmorphOS (OSDI '18), Coyote (OSDI '20) |
| Tu 11/15 | Disaggregation | LegoOS (OSDI '18) |
| Th 11/17 | Memory management | Llama (ASPLOS '20) |
| Tu 11/22 | Memory management | TLB shootdowns (EuroSys '20) |

# For Tuesday

**Multikernel**
- Regular paper discussion
- Do submit a review

**CloudLab**
- To learn more about what CloudLab is
- First 2 sections only
- No need to submit a review

# Questions to Ask When Reading a Paper

# High-Level Questions

- What is the problem?
  - Why is it important?
- What is the solution?
  - What is new about the solution?
- Which parts did you not understand?

# More Detailed Questions

- Solution
  - What is their approach?
  - What are the key components and how important is each one?
  - Did the paper solve the problem?
  - Are there limitations? How fundamental are they?
- Evaluation
  - How did they evaluate their work?
  - Are the experiments realistic (testbed, workloads, etc.)?
  - Do they demonstrate that the solution works?
- Impact
  - Do you think this work will be impactful? Why?
  - What kind of impact do you think it will have?

# More Detailed Questions

- Authors
  - Who are they and why did they write this paper now?
- Extensions
  - Useful for you to think about as a researcher!
  - What weaknesses does the paper have/how could it be improved?
  - Could you apply these ideas to other problems or in other domains?

# Summary

- CSE 291: Operating Systems in Datacenters
- Focus on learning about new OS techniques in datacenters
- Read and discuss research papers
- Undertake a class project