

EECS 570

Lecture 1

Parallel Computer Architecture

Winter 2024

Prof. Ronald Dreslinski

<http://www.eecs.umich.edu/courses/eecs570/>

Slides developed in part by Profs. Austin, Adve, Falsafi, Martin, Narayanasamy, Nowatzky, Peh, and Wenisch of CMU, EPFL, MIT, UPenn, U-M, UIUC



Intel Paragon XP/S

EECS 570 Class Info

- Instructor: Prof. Ronald Dreslinski(rdreslin@umich.edu)
- Research interests:
 - Intersection of computer architecture/VLSI

EECS 570 Class Info

- GSIs:
 - ❑ Jonah Rosenblum(jonaher@umich.edu)
 - ❑ Joy Dong (joydong@umich.edu)
- Class info:
 - ❑ URL: <https://www.eecs.umich.edu/courses/eecs570/>
 - ❑ Canvas for reading quizzes and reporting grades
 - ❑ Piazza for technical discussions and project coordination
- Enrollment:
 - ❑ **Enrollment will not go over 85 under any circumstances**

Meeting Times

- Lecture:
 - ❑ MW 3:00-4:20pm
- Discussion
 - ❑ F 3:30-4:20pm
 - ❑ Used for talking about programming assignments and projects
 - ❑ Also for discussing class topics in small groups
- Office Hours:
 - ❑ Prof. Dreslinski: Mondays 2pm-2:50pm, Wednesday 5-6pm, 2637 BBB
 - ❑ GSI: TBD (will be posted on course webpage)
- Q&A:
 - ❑ Use Piazza for all technical questions
 - ❑ **Office hours for course policy feedback**
 - ❑ Use email very sparingly

Who Should Take 570?

- Graduate Students (and seniors interested in research)
 - ☐ Computer architects to be
 - ☐ Computer system designers
 - ☐ Those interested in computer systems
- Required Background
 - ☐ Computer Architecture (e.g., EECS 470)
 - ☐ C/C++ Programming
- If you do not have the required background (especially a good base in computer architecture), do NOT take this course.
 - ☐ Topics covered such as coherence and consistency are complex even for computer architects
 - ☐ This course requires prior knowledge of topics like out-of-order execution, speculation, rollback, caching, prefetching...

Grading (tentative)

- 2 Programming Assignments: 5% (PA1) and 10% (PA2)
- Reading Quizzes for (almost) every lecture: 10%
- Midterm Exam: 25%
- Final Exam: 25%
- Final Research Project: 25%
- **This course is a lot of work. Do NOT underestimate it.**
- The course grade is heavily weighted towards the end of the semester
 - ❑ Please don't ask what your final grade will be after the midterm; we honestly have no way of knowing

Grading (Contd.)

- Group studies are encouraged
- Group discussions are encouraged
- However,...
 - ❑ All programming assignments must be results of individual work
 - ❑ All reading quizzes must be done individually, questions/answers should not be posted publicly
- There is no tolerance for academic dishonesty. Please refer to the University Policy on cheating and plagiarism. Discussion and group studies are encouraged, but all submitted material must be the student's individual work (or in case of the project, individual group work).

Some Advice on Reading

- If you carefully read every paper from start to finish, it will take a very long time
- Learn to skim past details that are not critical to the paper's overall message

Reading Quizzes

- You must take an online quiz for *every* paper
- Quizzes must be completed by class start via Canvas
- There will be 2 multiple choice questions
 - ☐ The questions are chosen randomly from a list
 - ☐ You only have 5 minutes
 - Not enough time to find the answer if you haven't read the paper
 - ☐ You only get one attempt
- Some of the questions may be reused on the midterm/final
- 4 lowest quiz grades (of about 40) will be dropped over the course of the semester (e.g., can skip some if you fall ill)
 - ☐ Retakes/retries/reschedules will not be given for any reason
 - The quiz drops are intended to cover such cases

Reading Quizzes

- **Quizzes do not make the difference between you passing and failing**
 - ☐ They are only worth 10% of your grade
 - ☐ Don't worry too much if you get some quiz questions wrong
- The quiz questions have been used in 2 (possibly more) past semesters
 - ☐ They did not cause people to fail in those semesters
- **The quiz questions and quiz policy are not going to change**
- Please do not post on Piazza complaining about quizzes or asking for more to be dropped
 - ☐ Questions asking for changes in quiz policy will be summarily removed
 - ☐ If you have constructive feedback on the quizzes, please discuss with me in OH

Discussing Course Policy

- Please do not use Piazza to ask for changes in course policy
 - ❑ It just doesn't work well
- **Posts asking for course policy changes will be summarily removed**
 - ❑ I am happy to listen to constructive feedback during my OH
- Regular logistical questions are fine to ask via Piazza
 - ❑ **OK:** “The Canvas quizzes for Monday’s lecture are not available. Can the prof or GSI please enable them?”
 - ❑ **Not OK:** “Can we have 6 quizzes be dropped instead of 4?”
- This course’s structure has largely been the same for about a decade
 - ❑ Many iterations of students have made it through this course
 - ❑ Large changes are not necessary

Piazza Guidelines

- Make sure your question is clear and as precise as you can make it
 - The easier it is for us to understand your question, the better we will be able to help you
- Before posting, please search existing posts to see whether your question has already been answered
- Do not make multiple small posts when a single larger one would suffice
- You are responsible for reading and being aware of questions and posts on Piazza
- Piazza sign-up link will be posted in a Canvas announcement later today

Final Project

- Original research on a topic related to the course
 - ❑ Goal: conduct and present original research in computer architecture
 - ❑ 25% of overall grade
 - ❑ Done in groups of 4-5 (exceptions *may* be made for PhD students)
 - ❑ **Please don't join a final project group and then drop the course; it makes things difficult for the other group members**
- Available infrastructure
 - ❑ gem5 multiprocessor simulator
 - ❑ GPGPUsim
 - ❑ Pin
 - ❑ Xeon Phi accelerators
- Timeline (and further info) will be posted on course website

Late Assignment Policy

- No late submissions for PA1, PA2, or final project
- If you have extenuating circumstances (e.g., admitted to hospital), email us and we'll work something out
 - ❑ Having too much coursework to complete to be able to meet the deadline is not an extenuating circumstance
- You will have ample time to work on the assignments, so please start early!

Regrade Policy

- Applies only to PA1, PA2, and midterm
- Regrade requests must be submitted in writing within one week from the day the assignment/midterm is handed back
- Regrade requests must specify clearly what the grading issue is
 - ❑ There should be a clear technical grading mistake for a regrade to be justified
 - ❑ Do not use regrades to just try and squeeze out more points
- On a regrade, the **entire** assignment will be regraded, and the grade may go up or down

Announcements

No discussion this Friday.

- Online quizzes (Canvas) on 1st readings due Wednesday, 3:00pm
- Please sign up for Piazza (sign up link will be posted in announcement on Canvas)

Readings

For Wednesday 1/17 (**quizzes due by 3:00pm**)

- ❑ David Wood and Mark Hill. “Cost-Effective Parallel Computing,” *IEEE Computer*, 1995.
- ❑ Mark Hill et al. “21st Century Computer Architecture.” CCC White Paper, 2012.

For Monday 1/22:

Christina Delimitrou and Christos Kozyrakis. Amdahl's law for tail latency. *Commun. ACM* 61, July 2018

H Kim, R Vuduc, S Baghsorkhi, J Choi, Wenmei Hwu, Performance Analysis and Tuning for General Purpose Graphics Processing Units (GPGPU), Ch. 1

Parallel Computer Architecture

The Multicore Revolution
Why did it happen?

If you want to make your computer faster, there are only two options:

~~1. increase clock frequency →~~

2. execute two or more things in parallel

~~Instruction-Level Parallelism (ILP) →~~

Programmer specified explicit parallelism

The ILP Wall

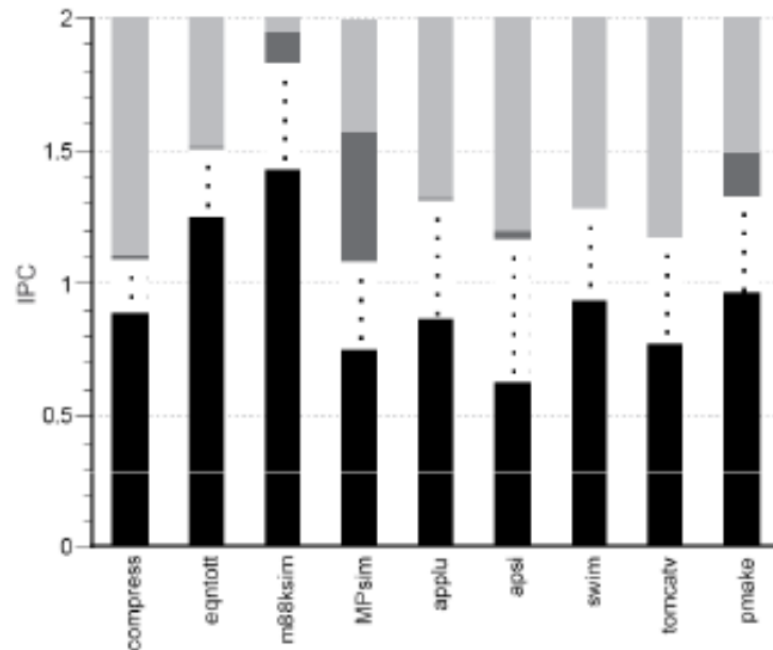


Figure 4. IPC Breakdown for a single 2-issue

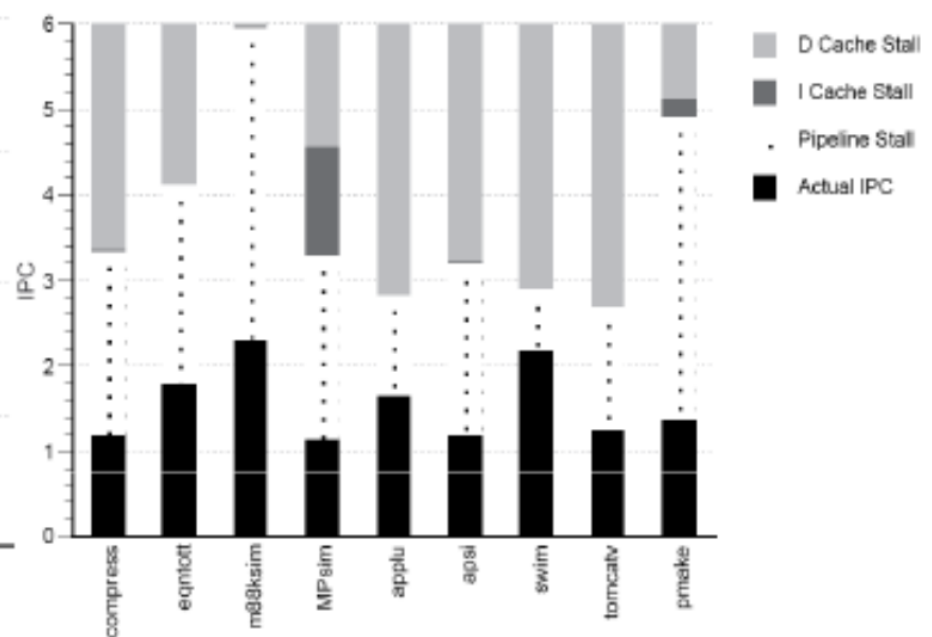
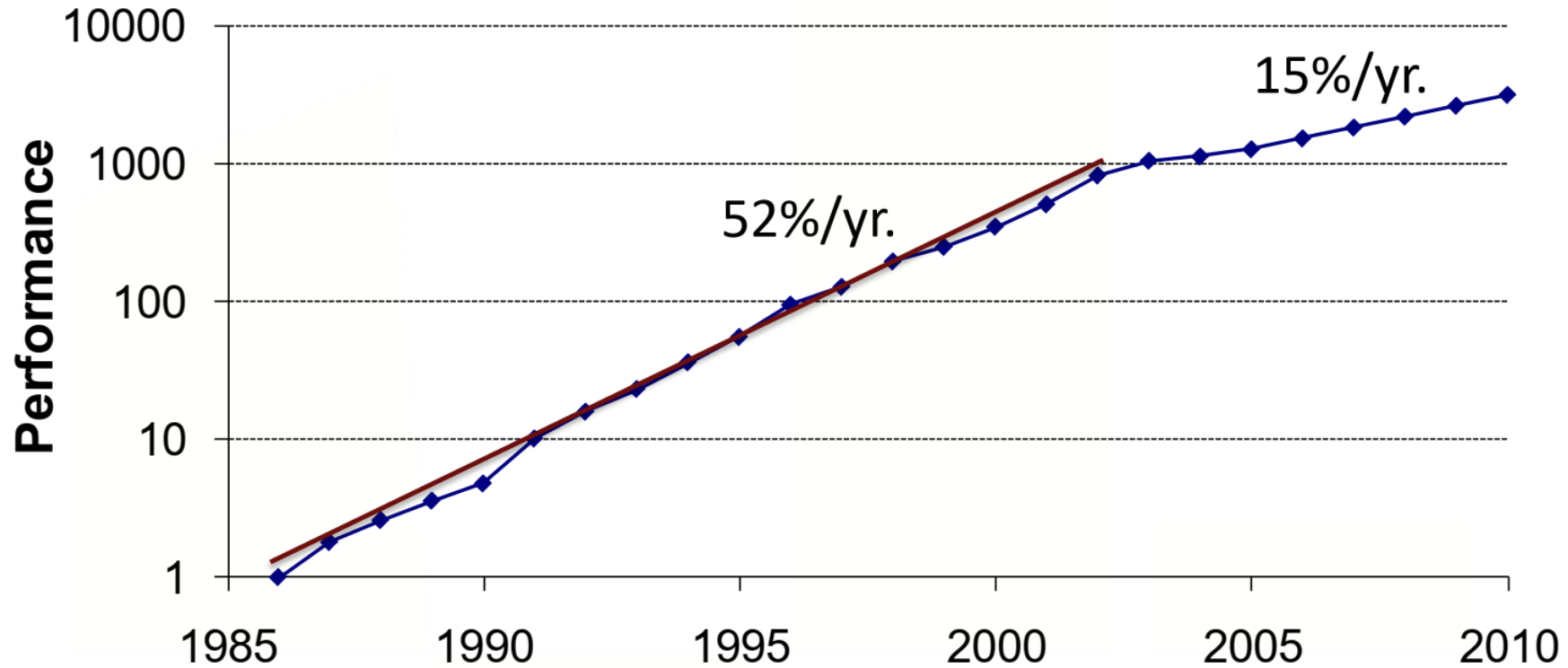


Figure 5. IPC Breakdown for the 6-issue processor.

Olukotun et al ASPLOS 96

- 6-issue has higher IPC than 2-issue, but not by 3x
 - Memory (I & D) and dependence (pipeline) stalls limit IPC

Single-thread performance

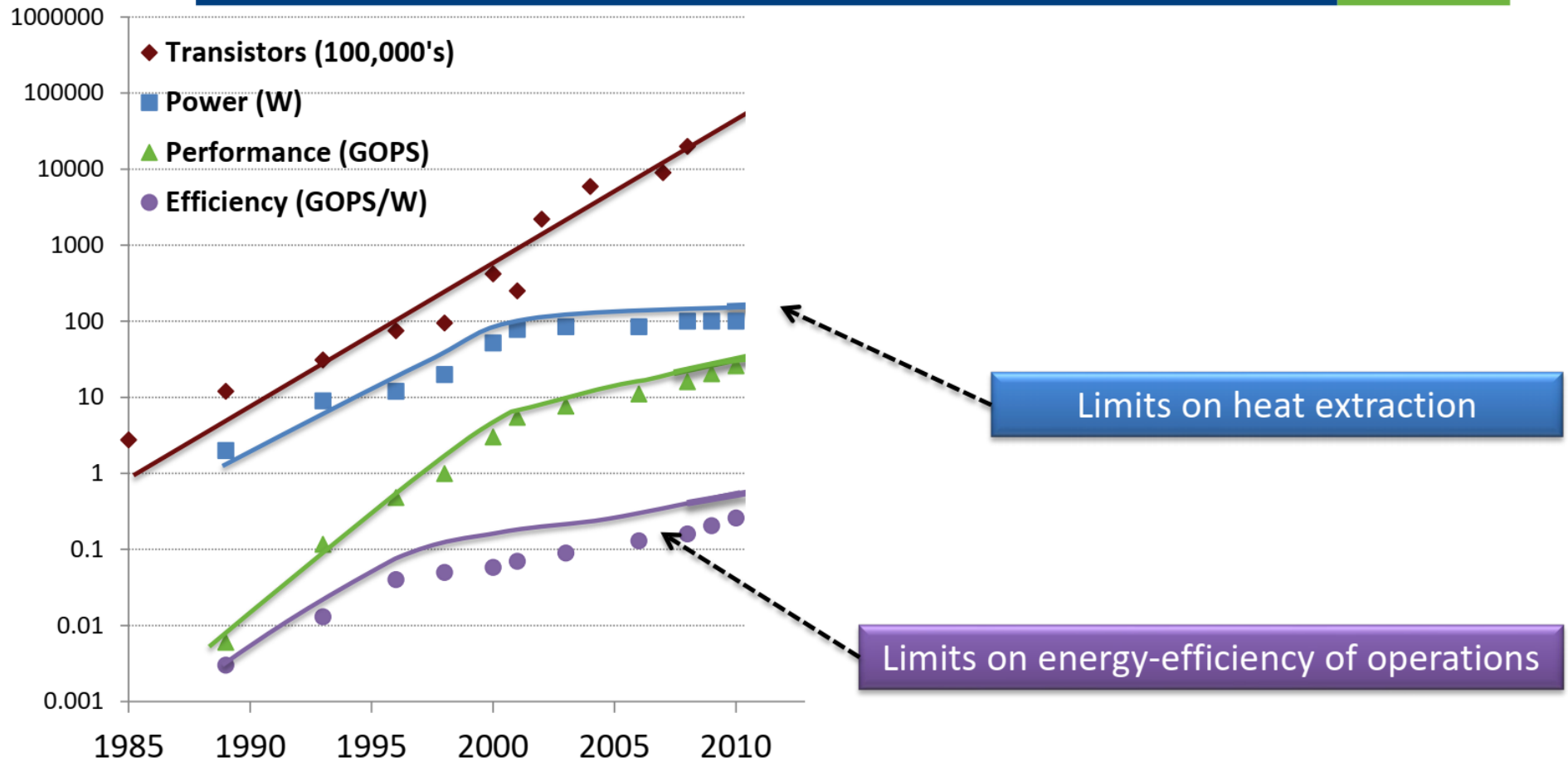


Source: Hennessy & Patterson, *Computer Architecture: A Quantitative Approach*, 4th ed.

Conclusion: Can't scale MHz or issue width to keep selling chips
Hence, multicore!

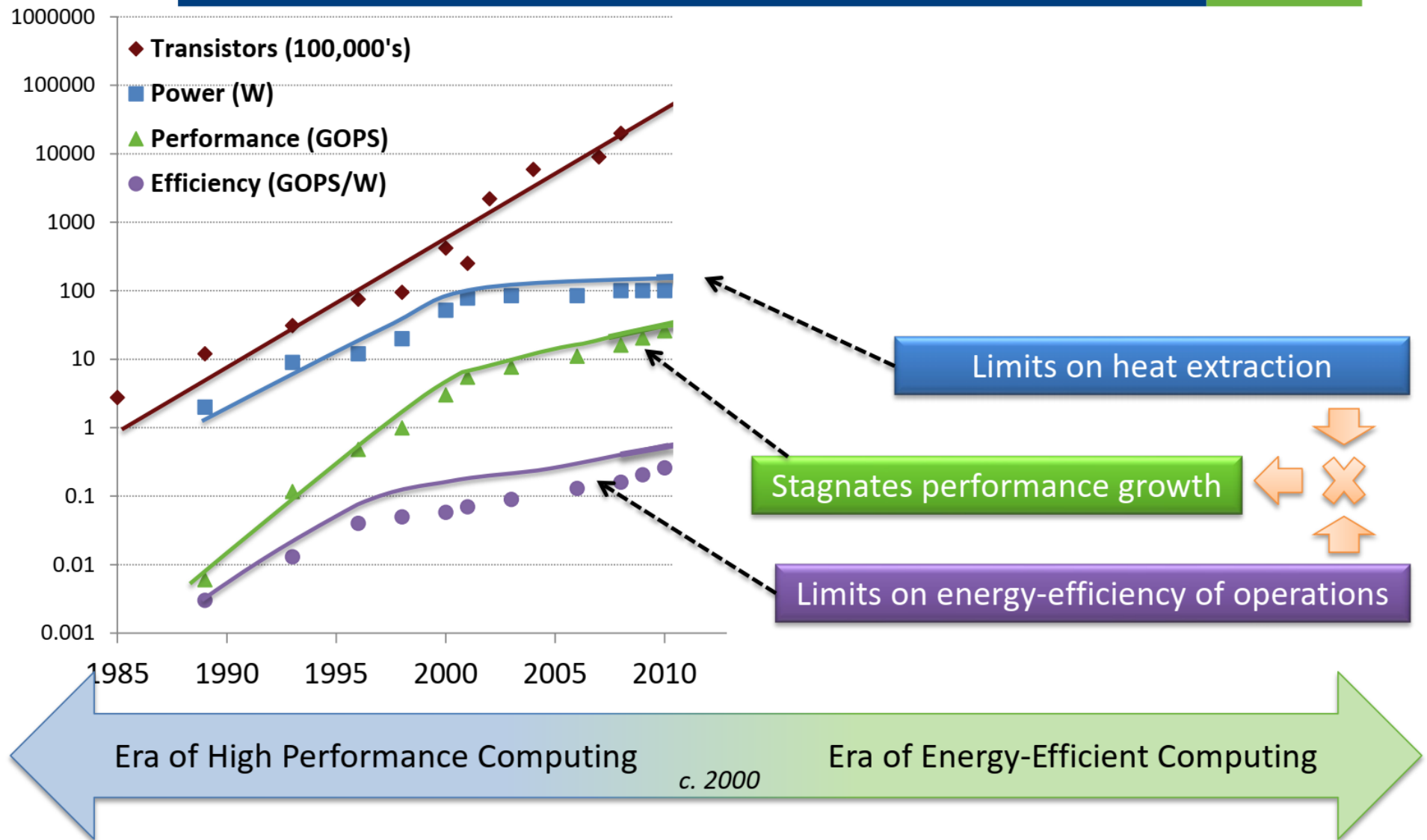


The Power Wall





The Power Wall



Classic CMOS Dennard Scaling: the Science behind Moore's Law

Scaling:

Voltage: V/α

Oxide: t_{ox}/α

Wire width: W/α

Gate width: L/α

Diffusion: x_d/α

Substrate: αN_A

Results:

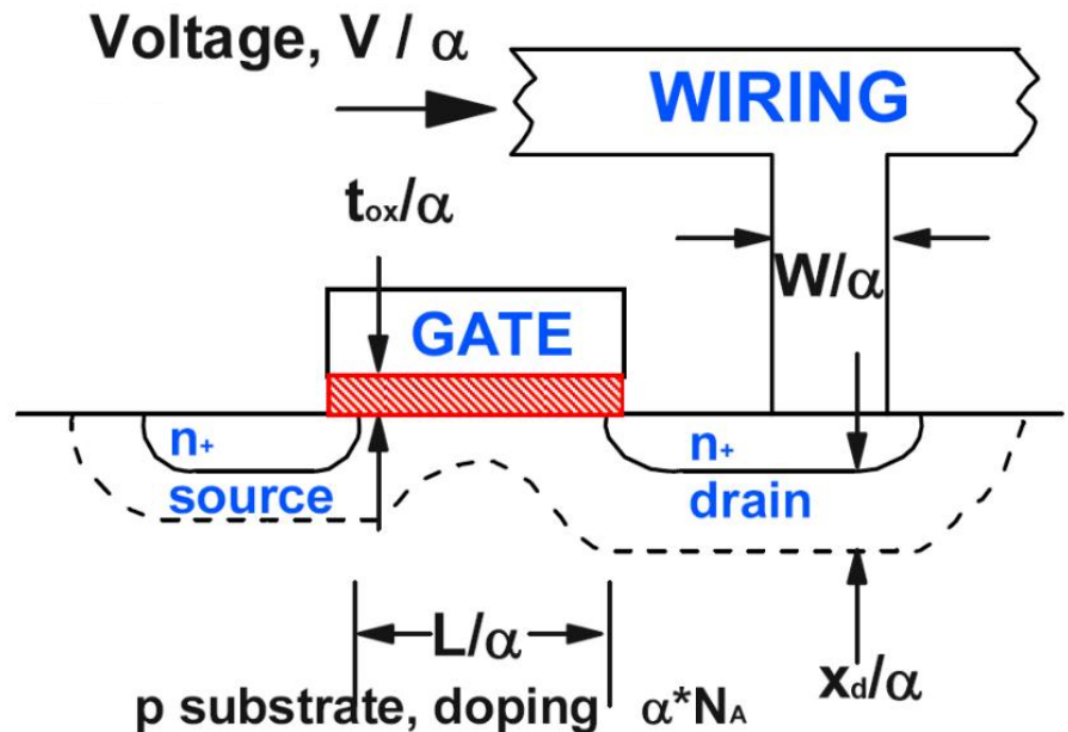
Higher Density: $\sim \alpha^2$

Higher Speed: $\sim \alpha$

Power/ckt: $1/\alpha^2$

Power Density: $\sim \text{Constant}$

Source: Future of Computing Performance:
Game Over or Next Level?,
National Academy Press, 2011



$$P = C V^2 f$$

R. H. Dennard et al.,
IEEE J. Solid State Circuits, (1974).

Post-classic CMOS Dennard Scaling

Post Dennard CMOS Scaling Rule

TODO:

Chips w/ higher power (no), smaller (☹), dark silicon (☺), or other (?)

Scaling:

Voltage: ~~V/α~~ V

Oxide: t_{ox}/α

Wire width: W/α

Gate width: L/α

Diffusion: x_d/α

Substrate: αN_A

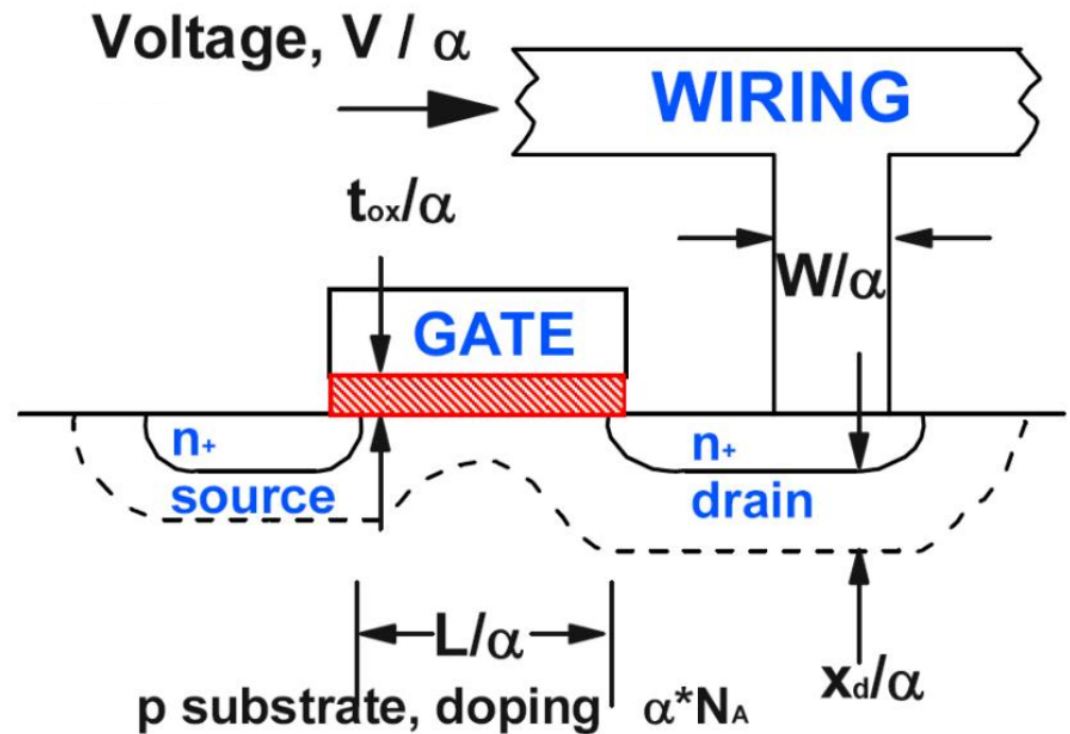
Results:

Higher Density: $\sim \alpha^2$

Higher Speed: $\sim \alpha$

Power/ckt: ~~$1/\alpha^2$~~ 1

Power Density: ~~$\sim \text{Constant}$~~ α^2



$$P = C V^2 f$$

R. H. Dennard et al.,
IEEE J. Solid State Circuits, (1974).

Lecture 1
Slide 20

Leakage Killed Dennard Scaling

Leakage:

- Exponential in inverse of V_{th}
- Exponential in temperature
- Linear in device count

To switch well

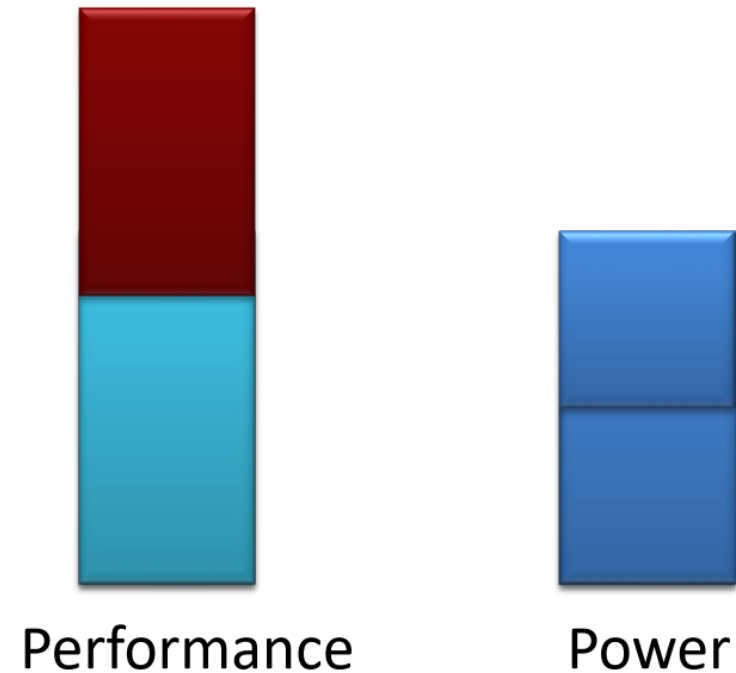
- must keep $V_{dd}/V_{th} > 3$

→ V_{dd} can't go down

Multicore: Solution to Power-constrained design?

Power = CV^2F $F \propto V$
Scale clock frequency to 80%

Now add a second core



Same power budget, but 1.6x performance!

But:

- ❑ Must parallelize application
- ❑ Remember Amdahl's Law!

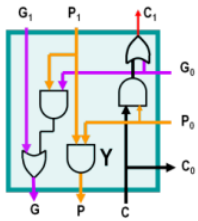
What Is a Parallel Computer?

“A collection of processing elements that communicate and cooperate to solve large problems fast.”

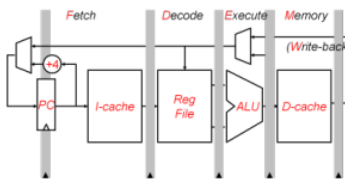
Almasi & Gottlieb, 1989

Spectrum of Parallelism

Bit-level

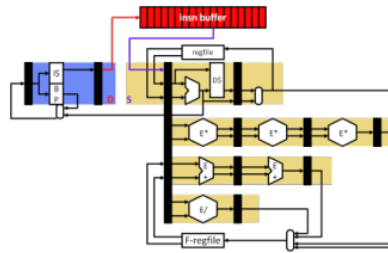


Pipelining



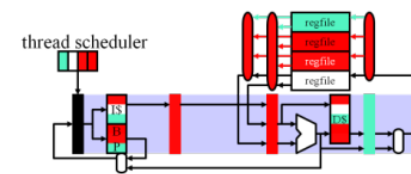
EECS 370

ILP



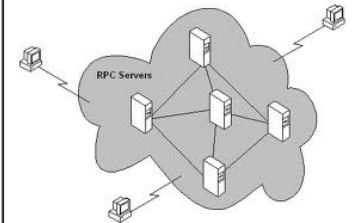
EECS 470

Multithreading
Multiprocessing



EECS 570

Distributed



EECS 591

Why multiprocessing?

- Desire for performance
- Techniques from 370/470 difficult to scale further

Why Parallelism Now?

- These arguments are no longer theoretical
- All major processor vendors are producing multicore chips
 - ❑ Most machines today are already parallel machines
 - ❑ All programmers will be parallel programmers???
- New software model
 - ❑ Want a new feature? Hide the “cost” by speeding up the code first
 - ❑ All programmers will be performance programmers???
- Some may eventually be hidden in libraries, compilers, and high level languages
 - ❑ But a lot of work is needed to get there
- Big open questions:
 - ❑ How should the chips, languages, OS be designed to make it easier for us to develop parallel programs?

Multicore in Products

- “We are dedicating all of our future product development to multicore designs. ... This is a sea change in computing”

Paul Otellini, President, Intel (2005)

- All microprocessor companies switch to MP (2X cores / 2 yrs)

	Intel's Nehalem-EX	Azul's Vega	nVidia's Tesla
Processors/System	4	16	4
Cores/Processor	8	48	448
Threads/Processor	2	1	
Threads/System	64	768	1792

Revolution Continues..



Azul's Vega 3 7300
54-core chip
864 cores
768 GB Memory
May 2008



Blue Gene/Q Sequoia
16-core chip
1.6 million cores
1.6 PB
2012



Sun's Modular DataCenter '08
8-core chip, 8-thread/core
816 cores / 160 sq.feet

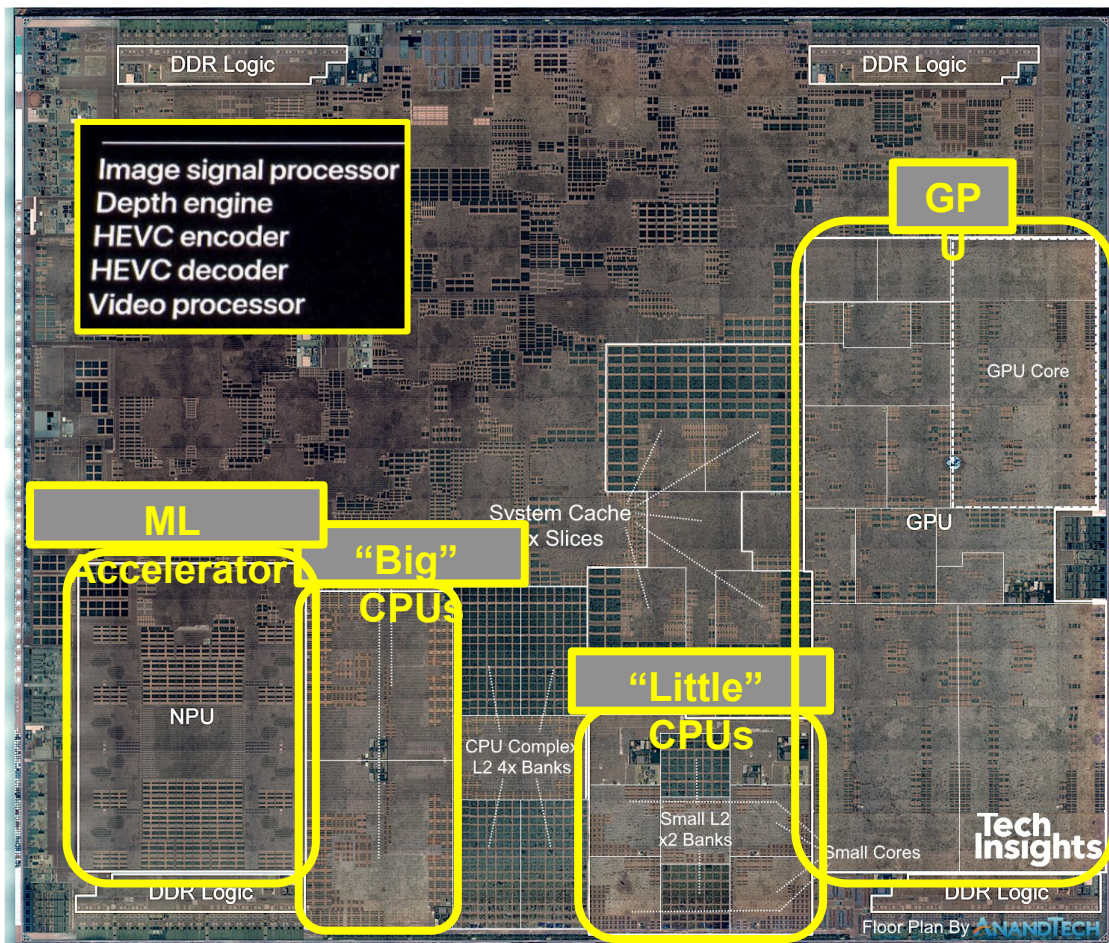
Lakeside Datacenter (Chicago)
1.1 milion sq.feet
~45 million threads

Multiprocessors Are Here To Stay

- Moore's law is making the multiprocessor a commodity part
 - ❑ 1B transistors on a chip, what to do with all of them?
 - ❑ Not enough ILP to justify a huge uniprocessor
 - ❑ Really big caches? t_{hit} increases, diminishing $\%_{miss}$ returns
- **Chip multiprocessors (CMPs)**
 - ❑ Every computing device (even your cell phone) is now a multiprocessor

Accelerator-Level Parallelism [Reddi and Hill 2019]

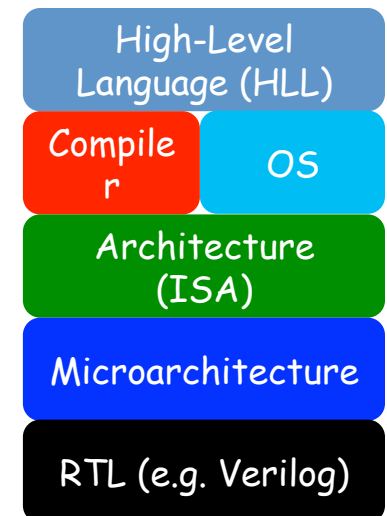
- Accelerators are specialized processing elements for different types of tasks
 - ❑ Machine learning, crypto, graphics (GPUs), etc



- Apple A12 System-on-Chip (SoC)
 - ❑ More than 40 accelerators [Wang and Shao 2019]

Heterogeneous Parallelism Across the Stack

- Parallelism has percolated up to high-level languages and compilers
 - Java, C/C++11, C#, etc. all have threads
- OS modified to best utilise parallel hardware
- New software toolchains for hardware accelerators
 - e.g. TensorFlow, PyTorch, TVM
- New landscape!



Course Outline

- Unit I – Parallel Programming Models
 - ❑ Message passing, shared memory (pthreads and GPU)
- Unit II – Synchronization
 - ❑ Synchronization, Locks, Lock-free structures
 - ❑ Transactional Memory
- Unit III – Coherency and Consistency
 - ❑ Snooping bus-based systems
 - ❑ Directory-based distributed shared memory
 - ❑ Memory Consistency Models
- Unit IV – Interconnection Networks
 - ❑ On-chip and off-chip networks
- Unit V – Applications
 - ❑ ML, health, and data-center applications

Parallel Programming Intro

Motivation for MP Systems

- Classical reason for multiprocessing:
More performance by using multiple processors in parallel
 - ❑ Divide computation among processors and allow them to work concurrently
 - ❑ Assumption 1: There is parallelism in the application
 - ❑ Assumption 2: We can exploit this parallelism

Finding Parallelism

1. Functional parallelism
 - ❑ Car: {engine, brakes, entertain, nav, ...}
 - ❑ Game: {physics, logic, UI, render, ...}
 - ❑ Signal processing: {transform, filter, scaling, ...}
2. Data parallelism
 - ❑ Vector, matrix, db table, pixels, ...
3. Request parallelism
 - ❑ Web, shared database, telephony, ...

Computational Complexity of (Sequential) Algorithms

- Model: Each step takes a unit time
- Determine the time (/space) required by the algorithm as a function of input size



Sequential Sorting Example

- Given an array of size n
- MergeSort takes $O(n \log n)$ time
- BubbleSort takes $O(n^2)$ time
- But, a BubbleSort implementation can sometimes be faster than a MergeSort implementation
- Why?

Sequential Sorting Example

- Given an array of size n
- MergeSort takes $O(n \log n)$ time
- BubbleSort takes $O(n^2)$ time
- But, a BubbleSort implementation can sometimes be faster than a MergeSort implementation
- The model is still useful
 - Indicates the scalability of the algorithm for large inputs
 - Lets us prove things like a sorting algorithm requires at least $O(n \log n)$ comparisons

We need a similar model for parallel algorithms