



Bio sketch

Education

undergrad: B. Tech, IIT Kharagpur, India, 1980
 grad: Ph.D. University of Utah, 1986

Professional:

- □ Assistant Prof, CS University of Oregon (86-91)
- □ Assistant Prof, ECE Oregon State University (91-92)
- □ Researcher/Prof Associé IRISA, Rennes (1993-2001)
- □ Colorado State University (2001-...)

Research contributions/interests:

- $\hfill\square$ High-performance computing, embedded systems, VLSI, algorithms, compilers, programming languages, ...
- □ Polyhedral model: a mathematical framework for describing, transforming and "compiling" massively parallel computations



Topics

- Answer: Why parallel computing?
- Overview the types of parallel computing
- Overview of speedup and efficiency
- Review the plan for this class
- Interactive/Illustrative example

Why Parallel Programming

Need for speed

- Many applications require orders of magnitude more compute power than we have right now.
 Speed came for a long time from technology improvements, mainly increased clock speed and chip density.
- The technology improvements have slowed down, and the only way to get more speed is to exploit parallelism. All new computers are now parallel computers.





Corollary of exponential growth

When two quantities grow exponentially, but at different rates, their ratio also grows exponentially.

$$y_1 = a^x$$
, and $y_2 = b^x$ for $a \ge b \ge 1$, therefore $\left(\frac{a}{b}\right)^x = r^x$, $r \ge 1$

$$1.1^{n} \neq O(2^{n})$$

 \Box or 2^n grows a lot faster than $(1.1)^n$

 Consequence for computer architecture: growth rate for e.g. memory is not as high as for processors, therefore, memory gets slower and slower (in terms of clock cycles) as compared to processors.

This gives rise to so called gaps or walls







Implicit Parallel Programming

- Let the compiler / runtime system detect parallelism, do task and data allocation, and scheduling.
- This has been a "holy grail" of compiler and parallel computing research for a long time and, after >40 years, still requires research on automatic parallelization, languages, OS, fault tolerance, etc.
- We still don't have a general purpose high performance implicit parallel programming language that compiles to a general purpose parallel machine.

Explicit Parallel Programming

Let the programmer express parallelism, task and data partitioning, allocation, synchronization, and scheduling, using programming languages extended with explicit parallel programming constructs.

- This makes the programming task harder, but a lot of fun and a large source of income
- (very large actually ⁽²⁾)













Clarifying Reading assignments

- 1. Inside Front Cover
- 2. Chapter 3 Parallel Algorithm Design









Amdahl's Law (from your reading)





Class Format

- Hands on: write efficient parallel code. This often starts with writing efficient sequential code.
- four versions of a code: naïve sequential, naïve parallel, smart sequential, smart parallel
- Lab covers many details, discussions follow up on what's learned in labs.
- PAs: Write effective and scalable parallel programs:

 report and interpret their performance.



Locality

A process (running on a processor) needs to have its data as close to it as possible. Processors have non uniform memory access (NUMA):

- registers 0 (extra) cycles
- cache 1,2,3 cycles there are multiple caches on a multicore chip
- local memory 50s to 100s of cycles
- global memory, processor-processor interconnect, disk ...

Process

- COARSE in OS world: a program in execution
- FINER GRAIN in HPC world: (a number of) loop iterations, or a function invocation

Exercise: summing numbers

10x10 grid of 100 processors. Add up 4000 numbers initially at the corner processor[0,0], processors can only communicate with NEWS neighbors

Different cost models:

- □ Each communication "event" may have an arbitrarily large set of numbers (cost independent of volume).
 - Alternate: Communication cost is proportional to volume
- Cost: Communication takes 100x the time of adding
 - Alternate: communication and computation take the same time







Weekly Roundup

- Moore's law (AKA why parallel?)
- SIMD vs MIMD
- 3 major types of parallel computing
- Explicit versus implicit parallelism
- Locality
- Calculation speedup
 - \Box Measured
 - □ Theoretical ideal
 - \Box Cost of communication