# Parallelism is Everywhere

- And at all scales
- Students can begin working with parallel systems early
  - My son Chris ran on Jaguar at ORNL, then the Top500 #1, after his freshman year
  - More recently, ran on nearly 1.2M cores (more than I ever have!)
  - And he also found scalability bugs…



Concurrent vCPUs in US East (Northern Virginia)





NCSA

# Training students to think precisely and quantitatively about parallel computing

- Parallel programming is *hard*

- We need to be scientific about solving these problems

- We would all like parallel programming to be eaiser and more fun, but to accomplish that, we need to focus on the real problems

- And we must set a good example for our students.
  - What follows are some examples of fuzzy thinking that we, as a community, must strive to improve
  - We can do this by insisting that our students be rigorous and to follow the scientific method
  - We must remind each other to separate opinion from fact …

# Quotes from "Enabling Technologies for Petaflops Computing" (MIT Press 1995)

- **"The software for the current generation of 100 GF machines is not adequate to be scaled to a TF…"**

- "The Petaflops computer is achievable at reasonable cost with technology available in about 20 years [2014]."
  - (estimated clock speed in 2004 — 700MHz

- "Software technology for MPP's must ev~~ that is portable across a wide variety of can the small but important MPP sector leverage the massive investment that is being applied to commercial software for the business and commodity computer market."

- "To address the inadequate state of software productivity, there is a need to develop language systems able to integrate software components that use different paradigms and language dialects."

Except that software was adequate not only for TeraFLOP but for PetaFLOP and almost certainly ExaFLOP

# Quotes from "Enabling Technologies for Petaflops Computing" (MIT Press 1995)

- "The software for the current generation of 100 GF machines is not adequate to be scaled to a TF…"

- "**The Petaflops computer is achievable at reasonable cost with technology available in about 20 years [2014].**"
  - **(estimated clock speed in 2004 — 700MHz)**

- "Software technology for MPP's ~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~ e
  that is portable across a wide va        This clock speed estimate was        en
  can the small but important MPP         dead on target – IBM BG/L ran at       t
  leverage the massive investment         this speed, and this prediction was    ware
  for the business and commodity          based on quantitative projections

- "To address the inadequate state                                               l to
  develop language systems able t                                                ie
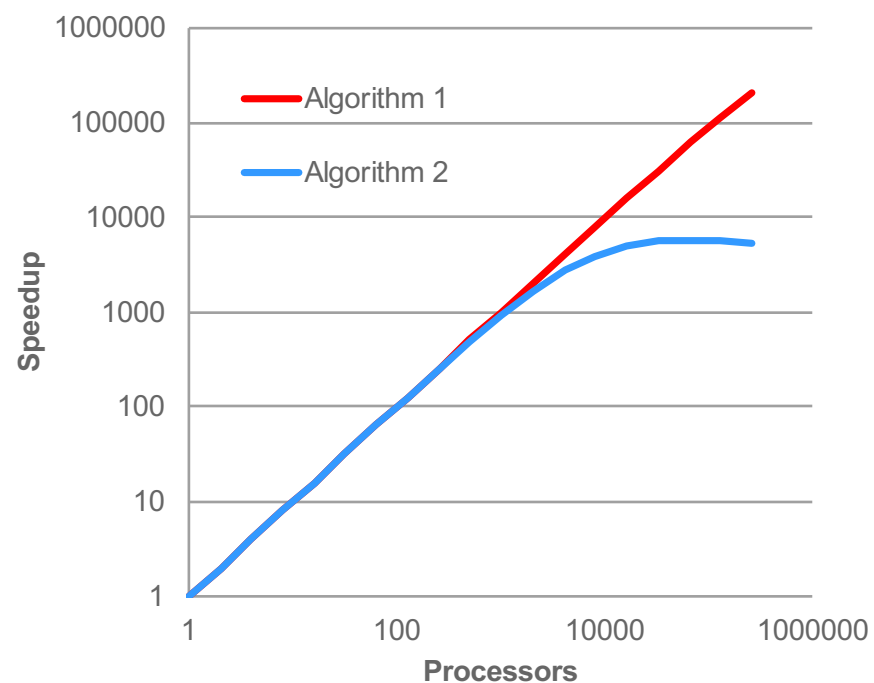  different paradigms and language dialects.

# What is Scalability?

- Is it
  - Using all cores on a chip efficiently
  - Using all nodes in a distributed memory system efficiently
  - Using many virtual CPUs in "the cloud" efficiently (and what is "many")
  - Ability to handle larger data
  - Ability to handle more users/requests
  - Running faster than a single core/chip/node?

- At least strong and weak scaling (when discussed together) are usually interpreted as #1 or #2
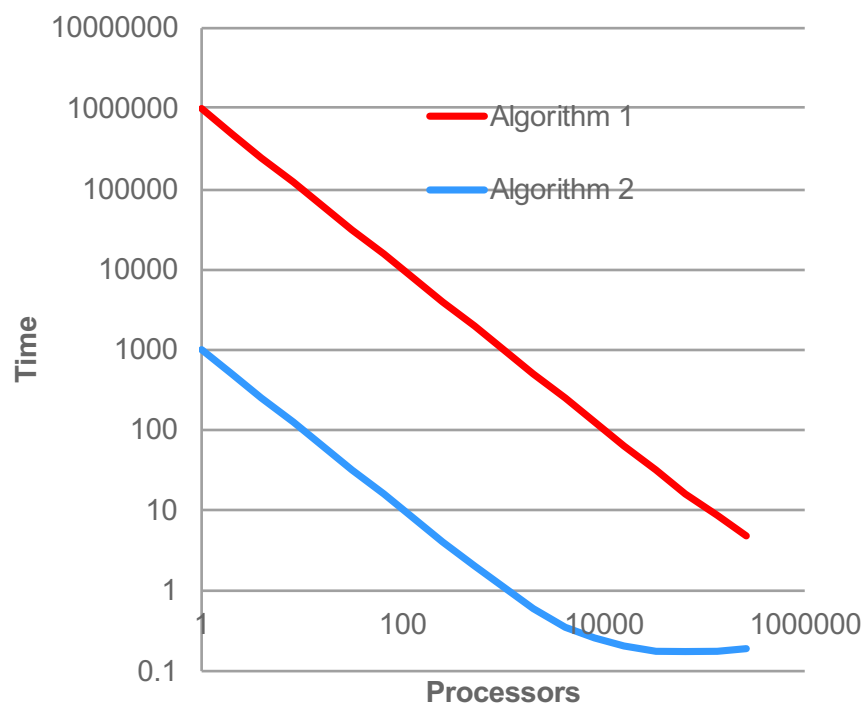
- Preciseness in language is important!

# Scalability – What does it tell you?

- Should I use algorithm 1 or algorithm 2?

# Scalability – What does it tell you?

- Should I use algorithm 1 or algorithm 2?

- Here's the same data, only showing time instead of speedup

- Its time to solution that is important

- You can always improve scalability by:
  - Decreasing per-core performance
  - Using a more computationally intensive algorithm

# How should you evaluate speedup?

- When measured, you need to
  - Understand the algorithm and options
    - Some are not obvious – accurate n-body force calculations can be done in $O(n)$, not $O(n^2)$ time; Broadcast of n words takes $O(n)$, not $O(n \log p)$
  - Evaluate implementation quality
    - Can use simple quantitative performance model
    - For many apps, STREAM is appropriate



W. K. Anderson, William D. Gropp, D. K. Kaushik, D. E. Keyes, and B. F. Smith. *Achieving high sustained performance in an unstructured mesh CFD application*, SC99 (Gordon Bell Prize)

# What Parameters Are Important?

- Raw rates (clock speed, FLOPS, memory bus bandwidth) are rarely useful by themselves
  - I once saw a paper that claimed memory bandwidth wasn't a useful indicator of performance
    - What they were looking at was the memory bus bandwidth – which was *far* higher than the sustained memory performance seen by applications or STREAM
    - Memory bus bandwidth *was* useless, but *sustained bandwidth* as measured by STREAM was useful

- Complex interactions of latency, bandwidth, instruction hazards, etc. make simple comparisons suspect
  - Students often latch onto the hype around raw performance numbers.

# How do you compare two programming models?

- First – are you comparing programming *models*, programming *systems*, or *implementations* of programming systems?
  - Answer – Almost always implementations
  - Implication – No paper should be accepted that claims to compare X to Y when all it does is compare an implementation of X on Z to an implementation of Y on Z

- Second – what conclusions can you draw from the comparison?
  - Example (drawn from a real vendor document)
  - Intended message: use our non-standard method – its faster!
  - Actual message: We have no clue how to implement MPI correctly
  - Good Test Question: Why do I make this claim?

### We're Better than MPI!



Chart axes:
- Y-axis: Bandwidth (1.00E+04 to 1.00E+09)
- X-axis: Data Size (1.E+00, 1.E+02, 1.E+04, 1.E+06)
- Legend: MPI, VendorX
- Annotation: 3X!

# A comparison of implementations

Collaborative Filtering (Weak scaling, 250 M edges/node)



Some MPI implementation

Factor of 100!

**Navigating the Maze of Graph Analytics Frameworks using Massive Graph Datasets**
Nadathur Satish, Narayanan Sundaram, Md. Mostofa Ali Patwary, Jiwon Seo, Jongsoo Park, M. Amber Hassaan, Shubho Sengupta, Zhaoming Yin, and Pradeep Dubey

# Reproducibility and significant digits

- A keystone of science is the reproducible experiment
  - Part of this is measuring and reporting only what is reproducible
- There are many challenges in making experiments reproducible
  - Complete description of the apparatus (e.g., system, compiler version/options, source code, …)
- But some are (relatively easy)
  - Don't report overly precise measurements
    - Or, don't use output from %e in your paper; even %.2e is often better
  - Describe your test environment
  - Make the code and data available

# What's in a name?

- In science, rigor and precision are important
- Example:
  - "Programming exascale systems requires moving beyond the BSP programming model; we need to replace MPI with <some new thing>."
- What is wrong with this statement?
  - "…exascale systems requires moving beyond BSP"
  - This is a hypothesis. How would you test or demonstrate it?
  - But there is a more subtle error…

# MPI is not a BSP system

- BSP = Bulk Synchronous Parallel
  - Programmers **like** the BSP model, adopting it even when not necessary (see FIB)
  - Unlike most programming models, *designed* with a performance model to encourage *quantitative* design in programs
- MPI makes it easy to emulate a BSP system
  - Rich set of collectives, barriers, blocking operations
- MPI (even MPI-1) sufficient for dynamic adaptive programming
  - The main issues are performance and "progress"
  - Improving implementations and better HW support for integrated CPU/NIC coordination more effective at supporting need
  - (This is not to say that there aren't issues – see Marc Snir's talk: "MPI is too High Level/MPI is too Low Level" at http://www.mcs.anl.gov/mpi-symposium/slides/marc_snir_25yrsmpi.pdf )

Leslie G. Valiant, A bridging model for parallel computation, Communications of the ACM, Volume 33 Issue 8, Aug. 1990

NCSA

# Understand what you ask for

- Current I/O performance is often appallingly poor
  - Even relative to what current systems can achieve
  - Part of the problem is the I/O interface semantics
- Many applications need to rethink their approach to I/O
  - Not sufficient to "fix" current I/O implementations
- HPC Centers have been complicit in causing this problem
  - By asking users the wrong question
  - By using their response as an excuse to keep doing the same thing
  - What is that question that causes so much trouble?
    - Do you want/need POSIX I/O?

# Just how bad Is current I/O performance?



*"A Multiplatform Study of I/O Behavior on Petascale Supercomputers,"* Huong Luu, Marianne Winslett, William Gropp, Robert Ross, Philip Carns, Kevin Harms, Prabhat, Suren Byna, and Yushu Yao, proceedings of HPDC'15.

# One Example in Parallel I/O

- Original code used a simple I/O strategy for reading input grid file
- Typical science run was many hours (6-24). I/O less than 10-20% of time
- But debugging runs are a few seconds – after input loaded
- Understanding the achievable performance, and using the appropriate I/O semantics (not POSIX), gives more than **1000X** speedup (MeshIO library of Ed Karrels)



Legacy I/O Times - 1M Grid Points on Vulcan



MPI I/O Times - 1M Grid Points on Vulcan

# What are some of the problems?

- POSIX I/O has a strong consistency model
  - Extremely hard to cache effectively
  - Applications need to transfer block-aligned and sized data to achieve performance
  - Complexity adds to fragility of file system, the major cause of failures on large scale HPC systems

- Files as I/O objects add metadata "choke points"
  - Serialize operations, even with "independent" files

- Burst buffers will not fix <u>these</u> problems – must change the semantics of the operations

NCSA

# Other communities have matched their data systems to their needs

- "Big Data" file systems have very different consistency models and metadata structures, designed for their application needs
  - Why doesn't HPC?
    - There have been some efforts, such as PVFS, but the requirement for POSIX has held up progress

# Why is the question "wrong"?

- To almost all application developers, asking if they need POSIX I/O means (to them)
  - Do you need open/seek/read/write/close?
- And if they answer "no", the implication is
  - They will need to rewrite their application and reformat their files
- With this interpretation, **no-one** would ever answer "no" to "Do you need POSIX I/O?"
  - Questions to which only one answer is reasonable don't tell you anything

# No science application code needs POSIX I/O semantics

- Many are single reader or single writer
  - Eventual consistency is fine
- Some are disjoint reader or writer
  - Eventual consistency is fine, but must handle non-block-aligned writes
- Some applications use the file system as a simple data base
  - Use a data base – we know how to make these fast and reliable
- Some applications use the file system to implement interprocess mutex
  - Use a mutex service – even MPI point-to-point
- A few use the file system as a bulletin board
  - Most likely better off using RDMA
  - Only need release or eventual consistency
- *Correct* Fortran codes do not require POSIX in any form
  - Standard requires unique open, enabling correct and aggressive client and/or server-side caching
- MPI-IO would be better off without POSIX

# Where's the hope?

- I firmly believe that a rigorous, scientific approach is the only way to solve the great challenges facing us in making better, more productive use of parallel computing

- And look where it has gotten us already
  - Computing power in a single system has increased one *billion* fold over my career. Data capacity and networking bandwidths have seen similar increases
  - We have been able to program these systems
    - Often the real challenge is in single core or single node performance

- We can accelerate this progress by (training our students in) challenging conventional wisdom and being rigorous in applying the scientific method

# Thanks!

- For funding from
  - National Science Foundation
  - Department of Energy
  - ExxonMobil and JumpLabs
  - State of Illinois
- My co-workers at Yale, Argonne, and Illinois
- My colleagues around the world
- My students
- And especially Ken Kennedy, for his contributions and for the example he set as a scholar and gentleman

**NCSA**