

High Performance Computing and Trends: Connecting Computational Requirements with Computing Resources

Jack Dongarra
Innovative Computing Laboratory
University of Tennessee

<http://www.cs.utk.edu/~dongarra/>

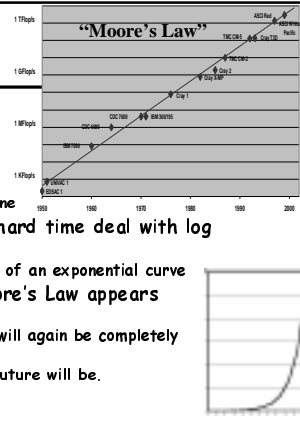
High Performance Computers

- ◆ ~ 20 years ago
 - 1×10^6 Floating Point Ops/sec (Mflop/s)
 - Scalar based
- ◆ ~ 10 years ago
 - 1×10^9 Floating Point Ops/sec (Gflop/s)
 - Vector & Shared memory computing, bandwidth aware
 - Block partitioned, latency tolerant
- ◆ ~ Today
 - 1×10^{12} Floating Point Ops/sec (Tflop/s)
 - Highly parallel, distributed processing, message passing, network based
 - data decomposition, communication/computation
- ◆ ~ 10 years away
 - 1×10^{15} Floating Point Ops/sec (Pflop/s)
 - Many more levels MH, combination/grids&HPC
 - More adaptive, LT and bandwidth aware, fault tolerant, extended precision, attention to SMP nodes

“Moore’s Wall”

— Horst Simon, NERSC

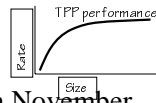
- ◆ Moore’s Law predicts exponential growth
 - Performance doubling every 18 months
 - Usually plotted on semi-log scale, appears as straight line
- ◆ Human experience has a hard time deal with log scale
 - We are sitting at the bend of an exponential curve
- ◆ From our perspective Moore’s Law appears as a “wall”
 - In a few years technology will again be completely different
 - Hard to predict what the future will be.



TOP500

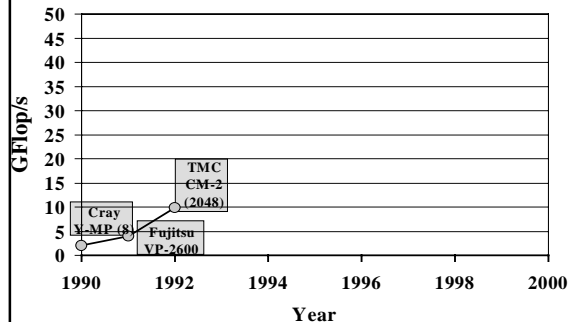
- Listing of the 500 most powerful Computers in the World
- Yardstick: Rmax from LINPACK MPP

$$Ax=b, \text{ dense problem}$$



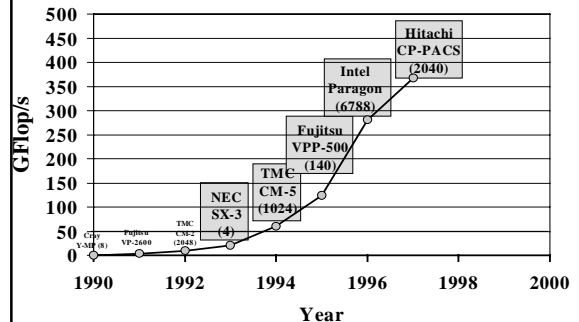
- Updated twice a year
- SC’xy in the States in November
- Meeting in Mannheim, Germany in June
- All data available from www.top500.org

Fastest Computer Over Time

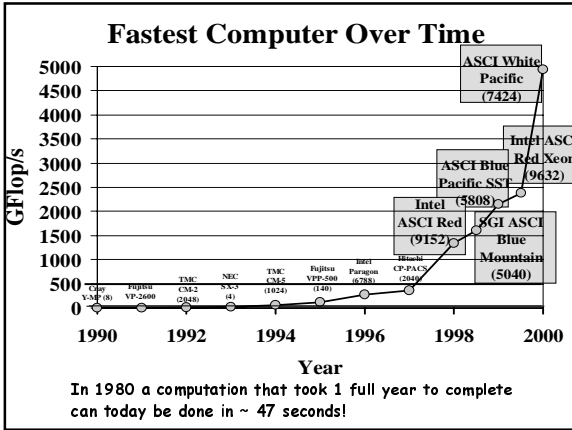


In 1980 a computation that took 1 full year to complete can now be done in ~ 10 hours!

Fastest Computer Over Time

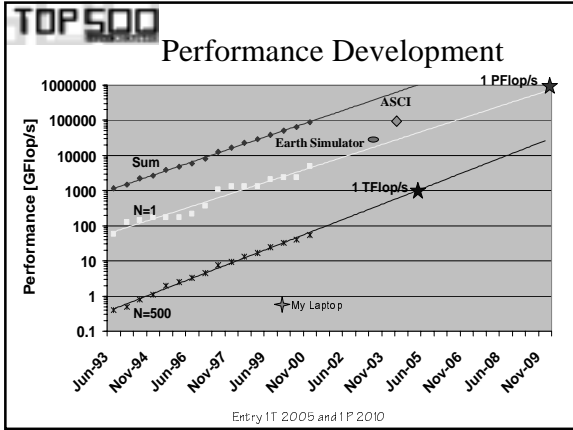
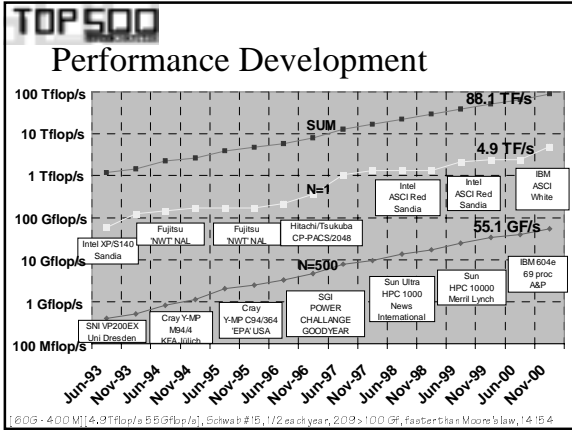


In 1980 a computation that took 1 full year to complete can now be done in ~ 16 minutes!

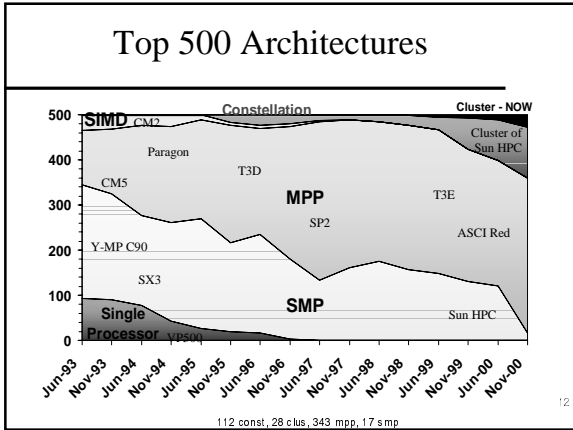


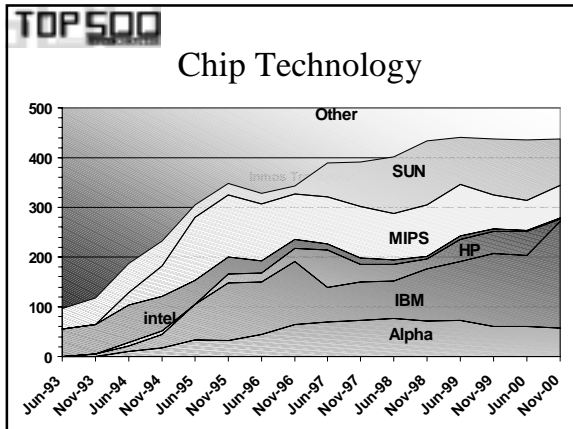
TOP500 Top 10 Machines (Nov 2000)

Rank	Company	Machine	Procs	Gflop/s	Place	Country	Year
1	IBM	ASCI White	8192	4938	Livermore National Laboratory	Livermore	2000
2	Intel	ASCI Red	9632	2380	Sandia National Labs	Albuquerque	USA
3	IBM	ASCI Blue-Pacific SST, IBM SP 604e	5808	2144	Lawrence Livermore National Laboratory	Livermore	USA
4	SGI	ASCI Blue Mountain	6144	1608	Los Alamos National Laboratory	Los Alamos	USA
5	IBM	SP Power3 375 MHz	1336	1417	Naval Oceanographic Office (NAVOCEANO)		USA
6	IBM	SP Power3 375 MHz	1104	1179	National Center for Environmental Protection		USA
7	Hitachi	SR8000-F1/112	112	1035	Leibniz Rechenzentrum	Muenchen	Germany
8	IBM	SP Power3 375 MHz, 8 way	1152	929	UCSD/San Diego Supercomputer Center		USA
9	Hitachi	SR8000-F1/100	100	917	High Energy Accelerator Research Organization /KEK	Tsukuba	Japan
10	Cray Inc.	T3E1200	1084	892	Government		USA



- ### Petaflop Computers Within the Next Decade
- Five basis design points:
 - Conventional technologies
 - 4.8 GHz processor, 8000 nodes, each w/16 processors
 - Processing-in-memory (PIM) designs
 - Reduce memory access bottleneck
 - Superconducting processor technologies
 - Digital superconductor technology, Rapid Single-Flux-Quantum (RSFQ) logic & hybrid technology multi-threaded (HTMT)
 - Special-purpose hardware designs
 - Specific applications e.g. GRAPE Project in Japan for gravitational force computations
 - Schemes utilizing the aggregate computing power of processors distributed on the web
 - SETI@home





High-Performance Computing Directions: Beowulf-class PC Clusters

Definition:

- ♦ **COTS PC Nodes**
 - Pentium, Alpha, PowerPC, SMP
- ♦ **COTS LAN/SAN Interconnect**
 - Ethernet, Myrinet, Gigaset, ATM
- ♦ **Open Source Unix Computing**
 - Linux, BSD
- ♦ **Message Passing**
 - MPI, PVM
 - HPF

Advantages:

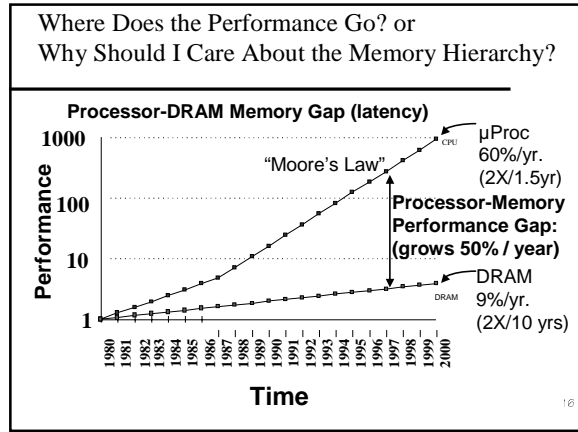
- ♦ Best price-performance
- ♦ Low entry-level cost
- ♦ Just-in-place configuration
- ♦ Vendor invulnerable
- ♦ Scalable
- ♦ Rapid technology tracking

Enabled by PC hardware, networks and operating system achieving capabilities of scientific workstations at a fraction of the cost and availability of industry standard message passing libraries. However, much more of a contact sport.

Clusters + TOP500

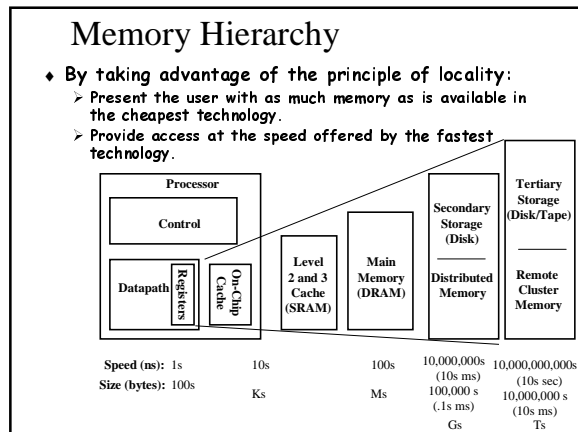
Rank	Site	System	Nodes	Peak Perf. (MFLOP/s)	Interconnect
1	LANL	IBM SP2	1000	1000	Myrinet
2	ORNL	IBM SP2	1000	1000	Myrinet
3	LANL	IBM SP2	1000	1000	Myrinet
4	LANL	IBM SP2	1000	1000	Myrinet
5	LANL	IBM SP2	1000	1000	Myrinet

- ♦ Peak performance
- ♦ Interconnection
- ♦ <http://clusters.top500.org>
- ♦ Benchmark results to follow in the coming months



Optimizing Computation and Memory Use

- ♦ **Computational optimizations**
 - Theoretical peak: $(\# \text{ fpus}) * (\text{flops/cycle}) * \text{Mhz}$
 - PIII: $(1 \text{ fpus}) * (1 \text{ flop/cycle}) * (850 \text{ Mhz}) = 850 \text{ MFLOP/s}$
 - Athlon: $(2 \text{ fpus}) * (1 \text{ flop/cycle}) * (600 \text{ Mhz}) = 1200 \text{ MFLOP/s}$
 - Power3: $(2 \text{ fpus}) * (2 \text{ flops/cycle}) * (375 \text{ Mhz}) = 1500 \text{ MFLOP/s}$
- ♦ **Operations like:**
 - $\alpha = x^T y$: 2 operands (16 Bytes) needed for 2 flops; at 850 Mflop/s will requires 1700 MW/s bandwidth
 - $y = \alpha x + y$: 3 operands (24 Bytes) needed for 2 flops; at 850 Mflop/s will requires 2550 MW/s bandwidth
- ♦ **Memory optimization**
 - Theoretical peak: $(\text{bus width}) * (\text{bus speed})$
 - PIII: $(32 \text{ bits}) * (133 \text{ Mhz}) = 532 \text{ MB/s} = 66.5 \text{ MW/s}$
 - Athlon: $(64 \text{ bits}) * (133 \text{ Mhz}) = 1064 \text{ MB/s} = 133 \text{ MW/s}$
 - Power3: $(128 \text{ bits}) * (100 \text{ Mhz}) = 1600 \text{ MB/s} = 200 \text{ MW/s}$

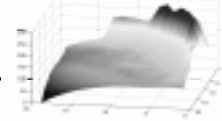


Self-Adapting Numerical Software (SANS)

- ◆ Today's processors can achieve high-performance, but this requires extensive machine-specific hand tuning.
- ◆ Operations like the BLAS require many man-hours / platform
 - Software lags far behind hardware introduction
 - Only done if financial incentive is there
- ◆ Hardware, compilers, and software have a large design space w/many parameters
 - Blocking sizes, loop nesting permutations, loop unrolling depths, software pipelining strategies, register allocations, and instruction schedules.
 - Complicated interactions with the increasingly sophisticated micro-architectures of new microprocessors.
- ◆ Need for quick/dynamic deployment of optimized routines.
- ◆ ATLAS - Automatic Tuned Linear Algebra Software

19

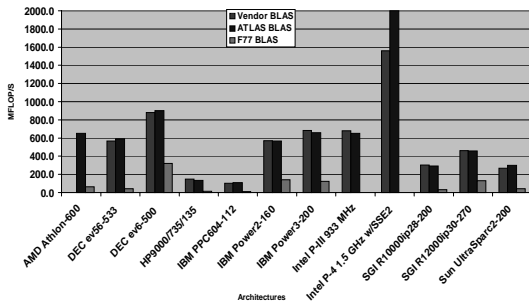
Software Generation Strategy



- ◆ Level 1 cache multiply optimizes for:
 - TLB access
 - L1 cache reuse
 - FP unit usage
 - Memory fetch
 - Register reuse
 - Loop overhead minimization
- ◆ Takes about 30 minutes to run.
- ◆ "New" model of high performance programming where critical code is machine generated using parameter optimization.
- ◆ Code is iteratively generated & timed until optimal case is found. We try:
 - Differing NBs
 - Breaking false dependencies
 - M, N and K loop unrolling
- ◆ Designed for RISC arch
 - Super Scalar
 - Need reasonable C compiler
- ◆ Today ATLAS in use by Matlab, Mathematica, Octave, Maple, Debian, Scyld Beowulf, SuSE, ...

20

ATLAS (DGEMM n = 500)



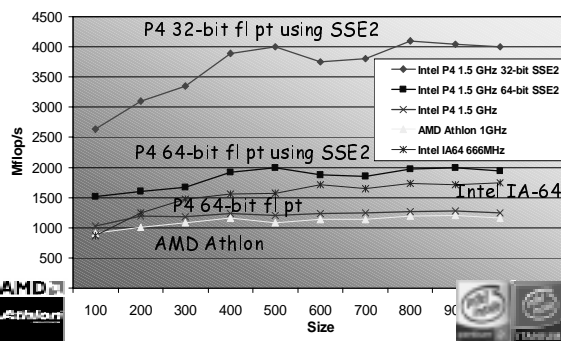
- ◆ ATLAS is faster than all other portable BLAS implementations and it is comparable with machine-specific libraries provided by the vendor.

Related Tuning Projects

- ◆ PHIPAC
 - Portable High Performance ANSI C
 - www.icsi.berkeley.edu/~bilmes/hipac initial automatic GEMM generation project
- ◆ FFTW Fastest Fourier Transform in the West
 - www.fftw.org
- ◆ UHFFT
 - tuning parallel FFT algorithms
 - rodin.cs.uh.edu/~mirkovic/fft/parfft.htm
- ◆ SPIRAL
 - Signal Processing Algorithms Implementation Research for Adaptable Libraries maps DSP algorithms to architectures
- ◆ Sparsity
 - Sparse-matrix-vector and Sparse-matrix-matrix multiplication
 - www.cs.berkeley.edu/~ejim/publication/ tunes code to sparsity structure of matrix more later in this tutorial

22

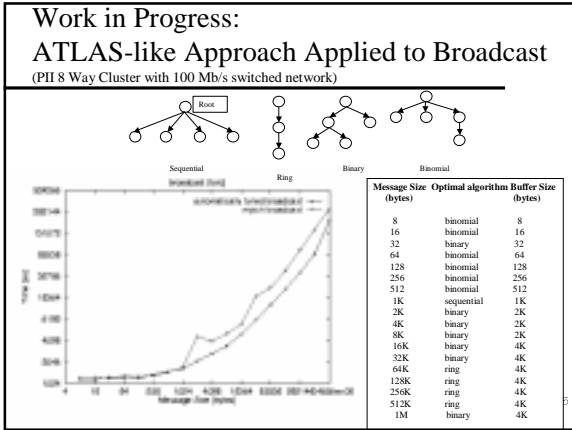
ATLAS Matrix Multiply (64 & 32 bit floating point results)



24

Machine-Assisted Application Development and Adaptation

- ◆ Communication libraries
 - Optimize for the specifics of one's configuration.
- ◆ Algorithm layout and implementation
 - Look at the different ways to express implementation



CG Variants by Dynamic Selection at Run Time

- Variants combine inner products to reduce communication bottleneck at the expense of more scalar ops.
- Same number of iterations, no advantage on a sequential processor
- With a large number of processor and a high-latency network may be advantages.
- Improvements can range from 15% to 50% depending on size.

CG Variants by Dynamic Selection at Run Time

- Variants combine inner products to reduce communication bottleneck at the expense of more scalar ops.
- Same number of iterations, no advantage on a sequential processor
- With a large number of processor and a high-latency network may be advantages.
- Improvements can range from 15% to 50% depending on size.

SETI@home

- Use thousands of Internet-connected PCs to help in the search for extraterrestrial intelligence.
- Uses data collected with the Arecibo Radio Telescope, in Puerto Rico
- When their computer is idle or being wasted this software will download a 300 kilobyte chunk of data for analysis.
- The results of this analysis are sent back to the SETI team, combined with thousands of other participants.

- **Largest distributed computation project in existence**
 - ~ 400,000 machines
 - Averaging 26 Tflop/s
- Today many companies trying this for profit.

Distributed and Parallel Systems

- Gather (unused) resources
- Steal cycles
- System SW manages resources
- System SW adds value
- 10% - 20% overhead is OK
- Resources drive applications
- Time to completion is not critical
- Time-shared

- Bounded set of resources
- Apps grow to consume all cycles
- Application manages resources
- System SW gets in the way
- 5% overhead is maximum
- Apps drive purchase of equipment
- Real-time constraints
- Space-shared

The Grid

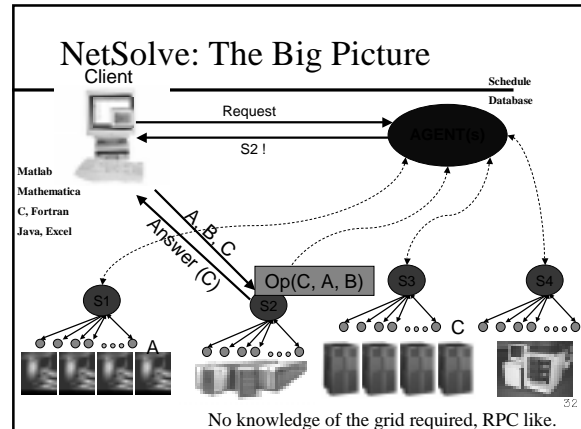
- To treat CPU cycles and software like commodities.
- on steroids.
- Enable the coordinated use of geographically distributed resources - in the absence of central control and existing trust relationships.
- Computing power is produced much like utilities such as power and water are produced for consumers.
- Users will have access to "power" on demand

NetSolve

Network Enabled Server

- ◆ NetSolve is an example of a grid based hardware/software server.
- ◆ Easy-of-use paramount
- ◆ Based on a RPC model but with ...
 - resource discovery, dynamic problem solving capabilities, load balancing, fault tolerance
 - asynchronicity, security, ...
- ◆ Other examples are NEOS from Argonne and NINF Japan.
- ◆ Use a resource, not tie together geographically distributed resources for a single application.

31



Basic Usage Scenarios



- ◆ Grid based numerical library routines
 - User doesn't have to have software library on their machine, LAPACK, SuperLU, ScaLAPACK, PETSc, AZTEC, ARPACK
- ◆ Task farming applications
 - "Pleasantly parallel" execution
 - eg Parameter studies
- ◆ Remote application execution
 - Complete applications with user specifying input parameters and receiving output
- ◆ "Blue Collar" Grid Based Computing
 - Does not require deep knowledge of network programming
 - Level of expressiveness right for many users
 - User can set things up, no "su" required
 - In use today, up to 200 servers in 9 countries

33

Futures for Numerical Algorithms and Software

- ◆ Numerical software will be adaptive, exploratory, and intelligent
- ◆ Determinism in numerical computing will be gone.
 - After all, its not reasonable to ask for exactness in numerical computations.
 - Auditability of the computation, reproducibility at a cost
- ◆ Importance of floating point arithmetic will be undiminished.
 - 16, 32, 64, 128 bits and beyond.
- ◆ Reproducibility, fault tolerance, and auditability
- ◆ Adaptivity is a key so applications can function appropriately

34

Contributors to These Ideas

- ◆ Top500
 - Erich Strohmaier, LBL
 - Hans Meuer, Mannheim U
- ◆ Linear Algebra
 - Victor Eijkhout, UTK
 - Piotr Luszczyk, UTK
 - Antoine Petit, UTK
 - Clint Whaley, UTK
- ◆ NetSolve
 - Dorian Arnold, UTK
 - Susan Blackford, UTK
 - Henri Casanova, UCSB
 - Michelle Miller, UTK
 - Sathish Vadhiyar, UTK

For additional information see...
www.netlib.org/top500/
www.netlib.org/atlas/
www.netlib.org/netsolve/
www.cs.utk.edu/~dongarra/

Many opportunities within the group at Tennessee

35