Introduction to HPC Computing







Agenda

- HPC Definition, Customers and Motivations
- HPC Principles and Common Architectures
- HPC Challenges







Experimental methods and quantification

of Newton's laws. Maxwell's equations ... phenomena

theory, experiment and simulation with massive data sets from multiple sources and disciplines

© 2014 NVIDIA





THE HIV-1

CAPSID

Atomic structure of the AIDS pathogen's protein coat PAGE 643

ACCELERATING discoveries

USING A SUPERCOMPUTER POWERED BY 3,000 TESLA PROCESSORS, UNIVERSITY OF ILLINOIS SCIENTISTS PERFORMED THE FIRST ALL-ATOM SIMULATION OF THE HIV VIRUS AND DISCOVERED THE CHEMICAL STRUCTURE OF ITS CAPSID – "THE PERFECT TARGET FOR FIGHTING THE INFECTION."

WITHOUT GPU, THE SUPERCOMPUTER WOULD NEED TO BE 5X LARGER FOR SIMILAR PERFORMANCE.



High Performance Computing

- High Performance Computing refers to the use of high advanced computational capabilities to solve cutting edge problems in science, engineering and business.
- The complex problems HPC has to deal with require capabilities beyond those of standalone computers
 - Aircraft design utilizing composite materials
 - Vehicle fuel efficiency and safety improvements
 - Simulations of enzyme catalysis, protein folding
 - Reservoir Simulation & Seismic Processing
 - Targeted material and drug design
 - Financial portfolio risk modeling
 - High Performance Computing enables scientists and engineers to perform simulations to:
 - Investigate phenomena where economics and constraints preclude experimentation
 - Evaluate complex models and manage massive data volumes
 - Transform business and engineering practices



Challenges

Simulation is key for Scientists and Engineers

- Previously: Theory Experiments
- Now: Theory Simulations Experiments

Agenda

Goals and Benefits



Principles

Challenges

Who is using High Performance Computing?



Digital Media

Digital content creation, management and distribution

Petroleum

Oil and gas exploration and production



CAE, EDA, CAD/PDM for electronics, automotive, and aerospace



Life Sciences

Research, drug discovery, diagnostics, information-based medicine



Business Intelligence Data warehousing and data

warehousing and data mining

Financial Services

Optimizing IT infrastructure, risk management and compliance, analytics

Government & Higher Education

Scientific research, classified/defense, weather/environmental sciences





5 Agenda



Challenges

CFD Contribution to the 767 in 1975





Challenges

CFD Contribution to the 767 in 2005







HPC Units

- High Performance Computing usually try to solve data- and/or compute-intensive problems
- The number of floating point operations per second – called Flops – is the unit used to evaluate compute intensive workloads and the computational performance of HPC computers.
 - Latest processors are able to issue 4 Floating Point Operations per clock cycle per core.
 - Ex: Dual Core processor @ 2.8GHz 4 x 2 x 2.8 = 22.4 GFlops = 22.6 10⁹ Flops
 - Ex: RoadRunner ~ 1 PFlops = 10¹⁵ Flops
- The amount of data being used during HPC workloads is measure is Bytes.

We also measure the amount of time needed for the data to be transferred within the different components of a HPC system in seconds.

- Mega = $2^{20} \approx 10^6$
- Giga = $2^{30} \approx 10^9$
- Tera = $2^{40} \approx 10^{12}$
- Peta = $2^{50} \approx 10^{15}$

Let's say you can write 1000 numbers on 1 sheet of paper 0.1mm thick,

10¹⁵ numbers represents 100.000km of stacked sheets = ¼ distance to the moon (Source: Jack Dongarra)



© 2014 NVIDIA

Agenda

22



HPC Units(Continued)

- High Performance Computing (HPC) units are:
 - Flop: floating point operation, usually double precision unless noted
 - Flop/s: floating point operations per second
 - Bytes: size of data (a double precision floating point number is 8 bytes)
- Typical sizes are millions, billions, trillions...

```
Mega Mflop/s = 106 flop/sec Mbyte = 2^{20} = 1048576 ~ 106 bytesGiga Gflop/s = 109 flop/sec Gbyte = 2^{30} ~ 109 bytesTera Tflop/s = 10^{12} flop/secTbyte = 2^{40} ~ 10^{12} bytesPeta Pflop/s = 10^{15} flop/secPbyte = 2^{50} ~ 10^{15} bytesExaEflop/s = 10^{18} flop/secEbyte = 2^{60} ~ 10^{18} bytesZettaZflop/s = 10^{21} flop/secZbyte = 2^{70} ~ 10^{21} bytesYottaYflop/s = 10^{24} flop/secYbyte = 2^{80} ~ 10^{24} bytes
```

- Current fastest (public) machine ~ 55 Pflop/s, 3.1M cores
 - Up-to-date list at www.top500.org

© 2014 NVIDIA

Goals and Benefits



Principles

Challenges

What Drives HPC? – "The Need for Speed..." Computational Needs exceed the Petaflops Range



CFD Wing Simulation
 512x64x256 Grid
 (8.3 x 10⁶ mesh points)
 5000 Flops per mesh points
 5000 time steps / cycle
 2.15 x 10¹⁴ Flop



Source: A Jameson, et al



- CFD Full Plane Simulation
 512x64x256 Grid
 (3.5 x 10¹⁷ mesh points)
 5000 Flops per mesh points
 5000 time steps / cycle
 8.17 x 10²⁴ Flop
 - Digital Movies and Special Effects

~1E14 Flops per frame 50 frames / sec 90 mins movie **2.7 x 10¹⁹ Flop** ~150 days on 2000 1 GFlops CPU Spare Parts Inventory Planning

Modeling the optimized deployment of 10,000 part numbers across 100 parts depots and requires: **2.14 x 10¹⁵ Flops**

(1 hour turnaround time)



Source: D Bailey, NERSC

Industry trend for rapid, frequent modeling for timely business decision support drives higher sustained performance

Materials Science



Magnetic Materials: Future: HDD Simulation

Source: B Dietrich, IBM

Future: HDD Simulation – **30TF, 2TB** Current: 2000 atoms – **2.64TF, 512GB**

Electronic Structures:

Current: 300 atoms – **0.5TF, 100GB** Future: 3000 atoms – **50TF, 2TB**

ACCELERATING INSIGHTS

"Now You Can Build Google's \$1M Artificial Brain on the Cheap"

WIRED

GOOGLE DATACENTER



STANFORD AI LAB



Deep learning with COTS HPC systems, A. Coates, B. Huval, T. Wang, D. Wu, A. Ng, B. Catanzaro ICML 2013

World's Fastest Enterplace Supercomputer



IBM.

3 Petaflops Linpack Performance

Most Energy-Efficient Petascale System in the World

3,000 NVIDIA Tesla K20X GPU Accelerators

Maximizing Opportunity for Oil Discovery with GPU-powered Supercomputer



Challenges

Goals and Benefits

Principles

History of fastest Architectures



Super Scalar/Special Purpose/Parallel

35 Agenda

History of fastest Architectures Old times

Basic ideas

- A super computer is like a Ferrari \rightarrow use specific components
- Super computing is like F1 or WRG \rightarrow adapt to the application scenario
- Vector Computer (Cray/ETA/Convex)
- Database Computer
- An alternative idea: Super SIMD (Connection Machines)
- Incredible creativity in architecture and network design

© 2014 NVIDIA



Challenges







What Accelerators

- There exists different type of processor accelerators...*
- An additional computation unit within a core – ex: vectorization unit
- SOC "System On a Chip" with lightweight generalist cores and integrated GPU – ex: NVIDIA Tegra K1
- Accelerator cards which can be plugged within the same server – ex: GPGPU (General Purpose Graphic Processing Unit)



Challenges



⊣ ⊢ High speed Link

Accelerator Card

Supercomputing... A Quick Look at the Web

- Top500.org
 - performance development
 - logarithmic progression! (x10 in 3years)
 - clusters, clusters (84%)!
 - 54% in industry
 - max power efficiency: 2.9 Gflops/W
 - #500: 96 TFlops! Total : 223 Pflops
 - poster Top500
- Graph500.org
 - BlueGene++
- Green500.org
 - GPU++
 - max: 4+ Gflops/W (Immersive cooling)
 - #1 green500 = #467 top500 (1 T00flops)
 - #1 top500 = #32 green500

© 2014 NVIDIA



Challenges



Challenges



History of fastest Architectures 2013... The Tianhe-2 (Milky Way-2)

- Ranked 1st in the top500 list of the most "powerful" (computing intensive) computers (June 2013)
- Ranked 6th in the graph500 list of the most "powerful" (data intensive processing) computers (June 2013)
- Ranked 32nd in the green500 list of the most energy efficient computer (June 2013)
- China (National University of Defense Technology)







2013 #2... The Titan (Cray XK7)

- Ranked 2nd in the top500 list (1st in Nov. 2012)
- 299008 cores Memory: 710 TB Cray Gemini Interconnect
- 18688 Opteron 6274 16 cores 2.200GHz + 18,688 Nvidia Tesla K20X GPUs
- Rmax = 17590 Rpeak = 27112 (computing efficiency : 65 %)
- Power: 8.2 MW... only!



Principles

Challenges



CSCS Piz Daint among Top 10 in Top 500 and Green 500

| Top500 Rank | TFLOPS/s | | | | | |
|--|----------|--|--|--|--|--|
| 1 | 33,862.7 | National Super Computer Centre Guangzhou | | | | |
| 2 | 17,590.0 | Oak Ridge National Lab#1 USA | | | | |
| 3 | 17,173.2 | DOE, United States | | | | |
| 4 | 10,510.0 | RIKEN Advanced Institute for Computational Science | | | | |
| 5 | 8,586.6 | Argonne National Lab | | | | |
| 6 | 6,271.0 | Swiss National Supercomputing Centre (CSCS) #1 Europe | | | | |
| 7 | 5,168.1 | University of Texas | | | | |
| 8 | 5,008.9 | Forschungszentrum Juelich | | | | |
| 9 | 4,293.3 | DOE, United States | | | | |
| 10 | 3,143.5 | Government | | | | |
| Piz Daint: 5,272 nodes with 1 x Xeon E5-2670 (SB) CPU, 1 x NVIDIA K20X GPU, Cray XC <u>30 at 2.0 MW</u> | | | | | | |

| reen500 Rank | MFLOPS/W | Site Site |
|-----------------|----------|--|
| 1 | 4,389.82 | GSIC Center, Tokyo Tech |
| 2 | 3,631.70 | Cambridge University |
| 3 | 3,517.84 | University of Tsukuba |
| 4 | 3,459.46 | SURFsara |
| 5 | 3,185.91 | Swiss National Supercomputing (CSCS) |
| 6 | 3,131.06 | ROMEO HPC Center |
| 7 | 3,019.72 | CSIRO |
| 8 | 2,951.95 | GSIC Center, Tokyo Tech |
| 9 | 2,813.14 | Eni |
| 10 | 2,629.10 | (Financial Institution) |
| 16 | 2,495.12 | Mississippi State (top non- NVIDIA) |
| 59 | 1,226.60 | ICHEC (top X86 cluster) |

Variety of node architectures



Cray XE6: dual-socket x 2-die x 6-core, 24 cores



Cray XC30: dual-socket x 8-core, 16 cores



Cray XK7: 16-core AMD + K20X GPU



Intel MIC: 16-core host + 60⁺ cores co-processor



Challenges



Performance History of HPC Systems





2010 2011 2012 2013

Intersect360 HPC User Site Census: Systems, July 2013 IDC HPC End-User MSC Study, 2013 **NVIDIA GPU is**

Accelerator of Choice

© 2014 NVIDIA





Microprocessor Transistors Trend

- Moore's Law
 - Gordon E Moore, Co-Founder of Intel
 - 1965 paper called "Cramming more components onto integrated circuits"
 - The number of transistors per chip doubles every 18 months.
- 1.0**EH0**0 The law has been true for ~45 years! 1.0E**1 0** Number of Transistors Transistors Transistors Table 1 Transistors Table 1 Transistors Thanks to -50% CAGR manufacturing progress 45 nm, 32 nm, 22 nm 1.0E105 ۲ 1980 1985 1990 1995 2000 2005 2010

26

Agenda



Challenges

Microprocessor Clock Speed Trend

- Clock Speed and Power Consumption used to increase with the numbers of transistors and was thus associated with Moore's Law
- Clock Speed is not increasing any more because it involves a power consumption and dissipation increase that
 - manufacturers are not able to follow.
- Power evolves like frequency³







Microprocessor Computational Power Trend

- The divergence between transistor density and frequency scaling drives toward multicore design
- Multicores trend makes the computational power of processors and systems to continue to grow.
- SpecInt rate is a standard benchmark measuring the performance of a full system
 - www.spec.org



Parallel Computing

- The complex problems HPC has to deal with require capabilities beyond those of standalone computers
- Parallel Computing enables scientists and engineers to compute faster and/or to solve bigger, more realistic problems
- The parallelism can be:
 - Internal inside the processor architecture – ex: Vectorization
 - By adding more processors
 - Shared Memory: SMP
 - Distributed Memory: MPP
 - Mixed Mode: both SMP and MPP



Challenges

SMP

 Multiple processors can operate independently but share the same memory resources



MPP

Mixed

- Processors have their own local memory
- Nodes are interconnected using a network, used by the processes to communicate



- MPP system with SMP Nodes
- Currently 90% of HPC Systems

Massively parallel

- The massively parallel system is a system with many interconnected cores
- The network is designed to be very efficient for communication between process
- The idea is that each core taken individually is not very good but taken globally it provides very good performances
- University of Juelich's supercomputer is running around 300 000 cores running at 0.85 GHz
- The interconnection was a 3D Torus based on standard 10 Gigabit Ethernet
- The performances is around 800 TFlops



Grids

- The grids consist of the interconnection of many computers
- This is an heterogeneous cluster: the computers and networks configurations differ from one to an other
- To be able to use the grid as a mean for computation, the problem must be split in many small problem that can be solved independently
- The servers then collect the results to build the global solution
- Example of project using grids are seti@home and folding@home
- <u>Seti@home</u> annouce a performance of 364.8 TFlops. The average is 40 TFlops



Goals and Benefits

Power Consumption in the Data Center

- Power and Cooling costs are increasing faster than acquisition costs
 - Today, each dollar of new servers cost more than \$0.50 to power and cool
- The current Data Centers are not able to sustain the electrical power needs of big HPC systems
- Data Center Power Consumption limits for the size of the HPC systems
- If we scale the current technology, each next generation exascale computer will require ~125MW of Total Power
 - For information, this is roughly equivalent to 120 Wind Turbines or 10% of a nuclear reactor

51 Agenda







Challenges

Performance vs Power Consumption

- The ratio to optimize for HPC systems is the <u>Flops/Watt ratio</u>
 - Power Consumption issues are currently addressed by the fastest HPC systems either by:
 - Using a large number of low power processors with limited single thread performance - Massively Parallel Computers
 - Using high single thread performance processors with heterogeneous cores or accelerators
 - With both solutions, <u>the Flops/Watt ratio is improved</u>

Principles

Challenges



Power is the Problem



Jaguar (Nov. '11) 2.3 petaflops @ 7 megawatts 7,000 Homes = 7 megawatts Small city

© 2014 NVIDIA

Power is the Problem

120 petaflops | 376 megawatts Enough to power all of San Francisco

1

ZNZNZNZNZNZNZNZNZNZNZNZNZNZNZNZ



Concepts and Terminology

© 2014 NVIDIA

Flynn's Taxonomy of Parallel Arch.



Distinguishes parallel architecture by instruction and data streams

SISD: classical uniprocessor architecture

| <u>SISD</u> | <u>S I M D</u> | | |
|--|---|--|--|
| Single Instruction, | Single Instruction, | | |
| Single Data | Multiple Data | | |
| <u>M I S D</u> Multiple Instruction, Single Data | <u>M I M D</u> Multiple Instruction, Multiple Data | | |

(Introduction to Parallel Computing, Blaise Barney) 68

Single Instruction, Multiple Data (SIMD)

- A type of parallel computer
- Single instruction: All processing units execute the same instruction at any given clock cycle
- Multiple data: Each processing unit can operate on a different data element
- This type of machine typically has an instruction dispatcher, a very highbandwidth internal network, and a very large array of very small-capacity instruction units.
- Best suited for specialized problems characterized by a high degree of regularity, such as image processing.
- Synchronous (lockstep) and deterministic execution
 - Two varieties: Processor Arrays and Vector Pipelines Examples:
 - Processor Arrays: Connection Machine CM-2, Maspar MP
 - Vector Pipelines: IBM 9000, Cray C90, Fujitsu VP, NEC SX-2, Hitachi S820



Multiple Instruction, Multiple Data (MIMD)



- Currently, the most common type of parallel computer. Most modern computers fall into this category.
- Multiple Instruction: every processor may be executing a different instruction stream
- Multiple Data: every processor may be working with a different data stream
- Execution can be synchronous or asynchronous, deterministic or nondeterministic
- Examples: most current supercomputers, networked parallel computer "grids" and multi-processor SMP computers - including some types of PCs.



Parallel Control Mechanisms



| Name | Meaning | Examples |
|--|--|---|
| Single Instruction, Multiple Data (SIMD) | A single thread of control, same computation applied across "vector" elts | Array notation as in Fortran 90: A[1:n] = A[1:n] + B[1:n] |
| Multiple Instruction, Multiple Data (MIMD) | Multiple threads of control, processors periodically synch | Parallel loop: forall (i=0; i <n; i++)<="" td=""></n;> |
| Single Program, Multiple Data (SPMD) | Multiple threads of control, but each processor executes same code | Processor-specific code: if (\$myid == 0) { } |

Some General Parallel Terminology



Like everything else, parallel computing has its own "jargon". Some of the more commonly used terms associated with parallel computing are listed below. Most of these will be discussed in more detail later.

Task

A logically discrete section of computational work. A task is typically a program or program-like set of instructions that is executed by a processor.

Parallel Task

A task that can be executed by multiple processors safely (yields correct results)

Serial Execution

Execution of a program sequentially, one statement at a time. In the simplest sense, this is what happens on a one processor machine. However, virtually all parallel tasks will have sections of a parallel program that must be executed serially.

Parallel Execution



Execution of a program by more than one task, with each task being able to execute the same or different statement at the same moment in time.

Shared Memory

From a strictly hardware point of view, describes a computer architecture where all processors have direct (usually bus based) access to common physical memory. In a programming sense, it describes a model where parallel tasks all have the same "picture" of memory and can directly address and access the same logical memory locations regardless of where the physical memory actually exists.

Distributed Memory

In hardware, refers to network based memory access for physical memory that is not common. As a programming model, tasks can only logically "see" local machine memory and must use communications to access memory on other machines where other tasks are executing.



Communications

Parallel tasks typically need to exchange data. There are several ways this can be accomplished, such as through a shared memory bus or over a network, however the actual event of data exchange is commonly referred to as communications regardless of the method employed.

Synchronization

- The coordination of parallel tasks in real time, very often associated with communications. Often implemented by establishing a synchronization point within an application where a task may not proceed further until another task(s) reaches the same or logically equivalent point.
 - Synchronization usually involves waiting by at least one task, and can therefore cause a parallel application's wall clock execution time to increase.



Granularity

- In parallel computing, granularity is a qualitative measure of the ratio of computation to communication.
- Coarse: relatively large amounts of computational work are done between communication events
- Fine: relatively small amounts of computational work are done between communication events
- **Observed Speedup**
 - Observed speedup of a code which has been parallelized, defined as:

wall-clock time of serial execution

wall-clock time of parallel execution

One of the simplest and most widely used indicators for a parallel program's performance.



Parallel Overhead

- The amount of time required to coordinate parallel tasks, as opposed to doing useful work. Parallel overhead can include factors such as:
 - Task start-up time
 - Synchronizations
 - Data communications
 - Software overhead imposed by parallel compilers, libraries, tools, operating system, etc.
 - Task termination time

Massively Parallel

Refers to the hardware that comprises a given parallel system having many processors. The meaning of many keeps increasing, but currently BG/L pushes this number to 6 digits.



Scalability

- Refers to a parallel system's (hardware and/or software) ability to demonstrate a proportionate increase in parallel speedup with the addition of more processors. Factors that contribute to scalability include:
 - Hardware particularly memory-cpu bandwidths and network communications
 - Application algorithm
 - Parallel overhead related
 - Characteristics of your specific application and coding



Parallel machines & programming model (hardware & software)

Parallelism within single processor - pipelining



- Like assembly line in manufacturing
- Instruction pipeline allows overlapping execution of multiple instructions with the same circuitry

| Instr. No. | Pipeline Stage | | | | | | |
|----------------|----------------|----|----|-----|-----|-----|-----|
| 1 | IF | ID | EX | мем | WB | | |
| 2 | | IF | ID | EX | мем | WB | |
| 3 | | | IF | ID | EX | мем | WB |
| 4 | | | | IF | ID | EX | мем |
| 5 | | | | | IF | ID | EX |
| Clock Cycle | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

IF = Instruction Fetch ID = Instruction Decode EX = Execute MEM = Memory access WB = Register write back

- Sequential execution: 5 (cycles) * 5 (inst.) = 25 cycles
- Pipelined execution: 5 (cycles to fill the pipe, latency) + 5 (cycles, 1 cycle/inst. throughput) = 10 cycles
- Arithmetic unit pipeline: A FP multiply may have latency 10 cycles, but throughput of 1/cycle
- Pipeline helps throughput/bandwidth, but not latency

Parallelism within single processor - SIMD
 SIMD: Single Instruction, Multiple Data

Scalar processing

X + Y

- traditional mode
- one operation produces
 one result

- SIMD processing
 - with SSE / SSE2
 - SSE = streaming SIMD extensions









SSE / SSE2 SIMD on Intel

• SSE2 data types: anything that fits into 16 bytes, e.g.,





- Instructions perform add, multiply etc. on all the data in this 16-byte register in parallel
- Challenges:
 - Need to be contiguous in memory and aligned
 - Some instructions to move data around from one part of register to another
- Similar on GPUs, vector processors (but many more simultaneous operations)

85

© 2014 NVIDIA

Memory Hierarchy ... Flops is not everything



- Most programs have a high degree of locality in their accesses
 - spatial locality: accessing things nearby previous accesses
 - temporal locality: reusing an item that was previously accessed
- Memory hierarchy tries to exploit locality to improve average



| Speed | 1ns | 10ns | 100ns | 10ms | 10sec |
|-------|-----|------|-------|------|-------|
| Size | KB | MB | GB | ТВ | PB |



Amdahl's Law – Limitations to Parallelism

<u>Are the HPC applications</u> <u>able to use a high number</u> <u>of processors?</u>



number of processors P

<u>Speedup</u>

Theoritical Maximum S(P)=P We can decompose the execution of a program on a parallel computer into:

 $T(1) = T_p + T_s$ and $T(P) = T_p/P + T_s$ Where P is the number of processors T_p is the part of the program which can be made parallel

 T_s is the part of the program which cannot be made parallel

We measure scalability of an application with the speedup: Speed Up : $S(P) = T(1) / T(P) = P / (f_p + f_s P)$

For example: $f_p=0.8$ and $f_s=0.2$ leads to $S(P) < 1/T_s < 5$

Amdahl's Law defines that the non-parallelizable fraction of a program limits performance and poses an upper limit on the scalability

46

Agenda

Adenda

Goals and Benefits

Gustafson's Law – Re-evaluating Amdahl's Law

- One of Amdahl's law hypothesis is that the problem size remains constant with the number of processors.
- Gustafson refuted this assumption and reevaluated Amdahl's law stating that:
 - when given a more powerful system, the problem generally expands to make use of the increased facilities.
 - the fraction of the program which can me made parallel increases with the problem size.

Speed Up : $S(P) = T(1) / T(P) = P / (f_p(N) + f_s(N).P)$

where P is the number of processors N is the problem size fp is the fraction of the program which can be made parallel fs is the fraction of the program which cannot be made parallel

 f_s(N) diminishes with the problem size then speedup approaches P as N approaches infinity



efficiently parallelized.





Δ



Summary

- Latest innovations have enable us to build reliable and programmable Petascale systems
- But new architectures will be needed to reach the ExaFlops barrier
 - The main challenges next HPC architectures will have to deal with are:
 - Power
 - Resiliency
 - Memory Capacity
 - Efficiency
 - Programmability