


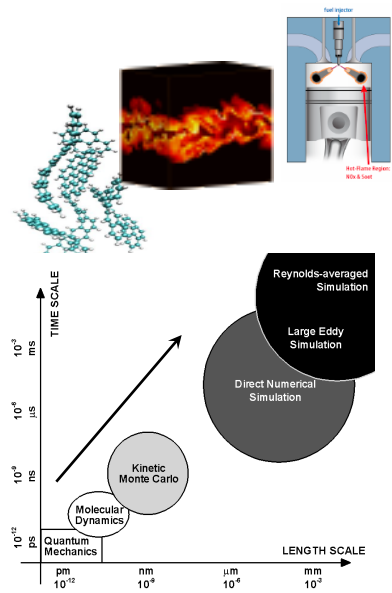

Introduction to parallel computing



1

Introduction

- As computers get faster, problems get more complex and computationally expensive.
- Most of engineering involves solving mathematical models of physical systems—
this means computing, either on paper, or by machine!
- Increase number of grid points
- Increase number of dimensions
 - x, y, z, time, state space, etc.
- Increase problem complexity
 - Add more chemical species
 - Resolve more timescales
 - Run in larger, more realistic domains
 - Do more realizations.
 - Relax simplifying assumptions
- Require faster turnaround time
 - How long are you willing to wait: 1 day to 2 weeks.
 - Often, real-time data processing is desired.

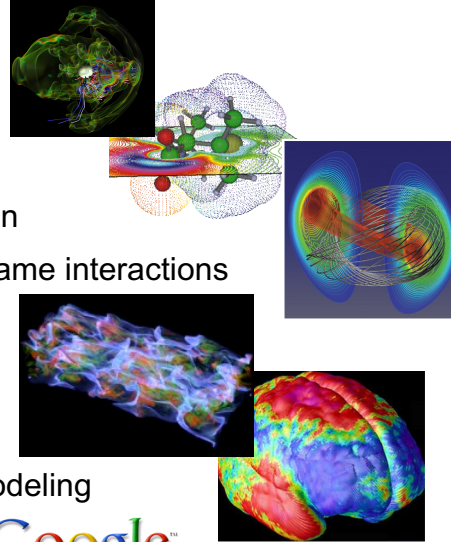



2

Parallel Computing Examples

3

- Astrophysics—supernovae
- Biology—protein dynamics
- Quantum Chemistry
- Fusion—tokamak simulation
- Combustion—turbulence flame interactions
- Climate—weather, CO₂
- Oil exploration
- Medical imaging
- Financial and economic modeling
- Web search engines



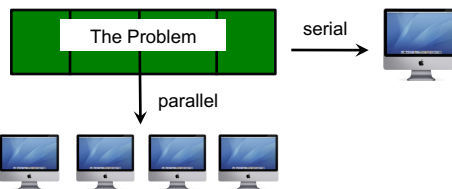
Google™



Parallel Computing

4

- Parallel computing
 - Reduce compute time
 - Increase problem size
 - Irrespective of time, some problems are just too big for one CPU.
- Split up problem to run on multiple CPUs



"What if five robbers were robbing for you and they each recruited five more robbers to rob for you?"



Classification—Flynn's Taxonomy

5

- **SISD**

- The usual workstation,
- *Nonparallel*



- **SIMD**

- All processors execute same instruction at the same time, but on different data streams.
- In many workstations: GPUs

- **MISD**

- Several processors work on same data.
- *Uncommon*
- Cryptography, data filtering



- **MIMD**

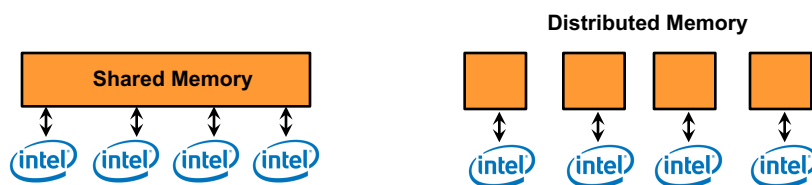
- Processors run independent code on independent data.
- *Most common parallel category*
- Often write one code, but each processor runs own version, with code execution depending on processor ID.



See top500.org for current supercomputer info.

Hardware

6



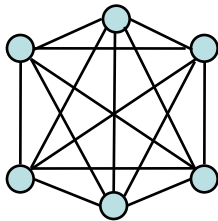
- Two main machine types: shared memory and distributed memory
- Shared: each processor has direct access to all memory.
 - Fast data transfers, simpler programming
 - Limited number of processors (low scalability)
 - Program with *POSIX threads, openMP*
- Distributed: each processor has own memory.
 - Highly scalable
 - More difficult to program.
 - Program with message passing: *MPI*



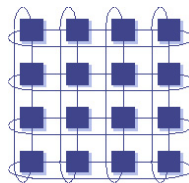
Interconnection Networks

7

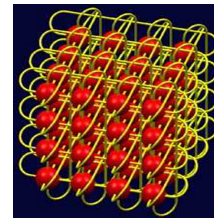
Direct Connection



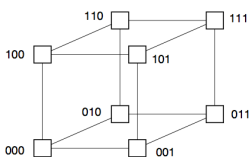
2D Torus



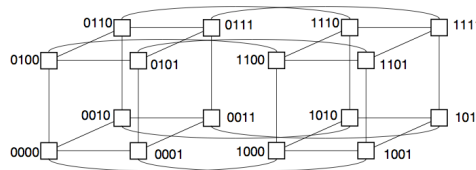
3D Torus: Jaguar



3D Hypercube



4D Hypercube



Other: Switch, Mesh, Omega network

Parallel Programming

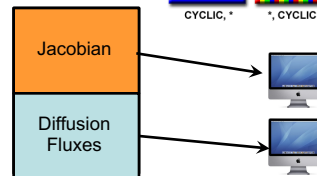
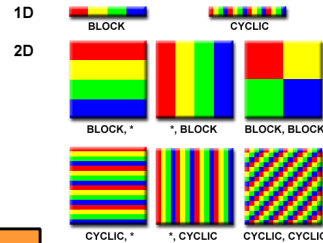
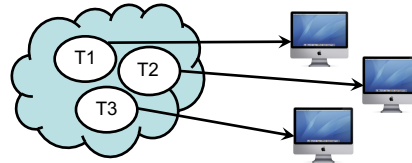
8

- Threads Model
 - Usually shared memory
 - Library of routines
 - Compiler directives
 - Think of threads as running different subroutines at once (*MULTI-TASKING*)
 - POSIX threads or P-threads
 - C language.
 - Very explicit, hard to code.
 - OpenMP
 - Compiler directives
 - Modify your existing code
 - Fortran and C/C++
 - Example: wrapping do loops



How to parallelize your problem?

- **Embarassingly parallel**
 - Little or no effort to separate the problem into parallel tasks.
 - Farm out obvious, independent tasks
 - Graphics processing
 - Each pixel is independent
 - Post processing multiple files in the same way.
- **Domain decomposition**
 - Perfect for spatial problems, numerical grids.
 - Particle flows.
 - Data based
- **Functional decomposition**
 - Algorithm-based.
 - Parallelize on tasks

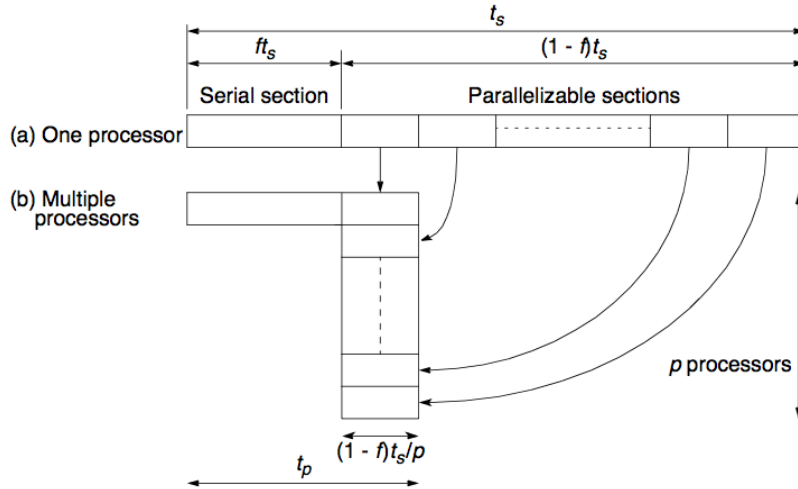


Some concepts and definitions

- **Speedup Factor**
 - $S(p) = t_{\text{serial}} / t_{\text{parallel}}$
 - $S(p) = (\text{best serial time}) / (\text{parallel time on } p \text{ procs})$
 - Serial and parallel algorithms are usually different.
 - Maximum speedup is p
 - Unless machines are different, or there is a nondeterministic algorithm.
 - This is rare, but happens.
- **Efficiency**
 - $E = S/p = t_{\text{serial}} / p * t_{\text{parallel}}$
- **Tradeoff between communication and computation.**
 - Communication takes time, reduces efficiency
 - Minimize communication between processors
 - Algorithm choice
 - Parallelization scheme
 - Communication network
 - As p increases, time decreases, but E decreases.

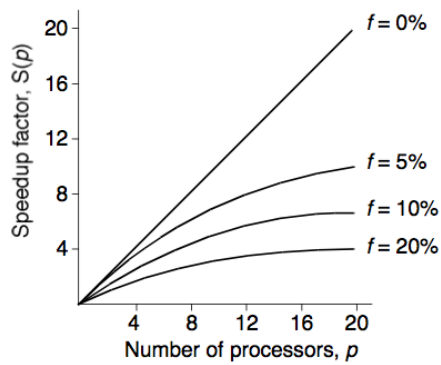


Amdahl's Law



Amdahl's Law

$$S(p) = \frac{t_s}{f t_s + (1 - f) t_s / p} = \frac{p}{1 + (p - 1) f}$$



Linear speedup is ideal

Max Speedup limited to $1/f$ as p gets large!

So if $f=5\%$, max speedup = 20!

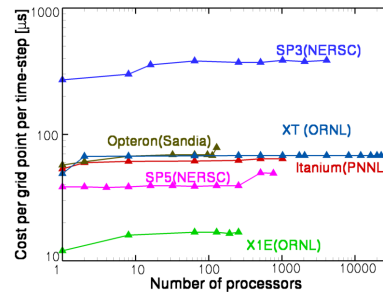
What is "wrong" with this analysis?

Assumes that f is constant as p increases. Usually as p increases the problem size increases so f decreases. (That was close!)



Example, S3D DNS Solver

- Cartesian domain decomposition
 - Only nearest neighbor communication.
 - Gives ~linear weak scaling since communication to computation ratio remains fixed
 - Communication scales with cube surface area, but computation goes with cube volume
- Explicit finite difference solver
 - 8th order central differences
 - 4th order, low storage RK integrator
- Solve reacting, compressible, Navier-Stokes equations.



Strong scaling: fix problem size, vary processors

Weak scaling: vary problem size with processors to keep a fixed load.

Which is more relevant is problem dependent (weak scaling will give "nicer" results).



My first MPI Program

- Open MPI
- Compile:
 - mpicxx myfirstMPI.cc
- Run:
 - mpirun -np 128 a.out
- Include the mpi library
- MPI_Init: startup MPI, the first function called
- MPI_Comm_rank: get rank of processor making the call.
 - Relative to communicator MPI_COMM_WORLD (which just means all the processors). A communicator is a group of processors.
 - Processors ordered 0 to N-1.
- MPI_Comm_size: get # of procs, N
- MPI_Finalize: shutdown MPI (last MPI call).
- Note syntax.

```
#include <iostream>
#include <mpi.h>
using namespace std;
int main(int argc, char* argv[]) {
    int myid, nprocs, ierr;
    MPI_Status status;

    MPI_Init(&argc, &argv);
    MPI_Comm_rank(MPI_COMM_WORLD, &myid);
    MPI_Comm_size(MPI_COMM_WORLD, &nprocs);

    if(myid==0)
        cout << endl << "Hello from process " << myid << endl;
    if(myid==1)
        cout << endl << "Hello from process " << myid << endl;
    if(myid==2)
        cout << endl << "Hello from process " << myid << endl;
    if(myid==3)
        cout << endl << "Hello from process " << myid << endl;

    MPI_Finalize();

    return 0;
}
```

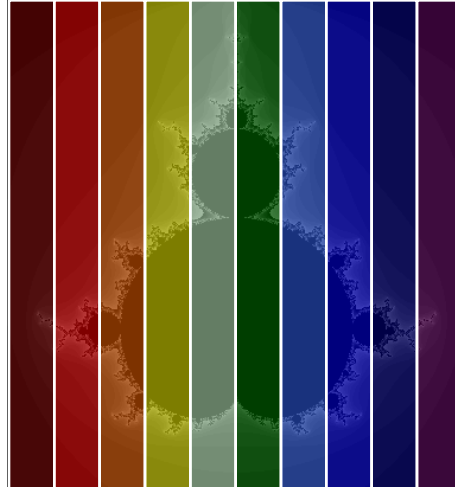
This works, but its silly, why?

Just replace ALL if statements with 1 cout statement



Parallelize the Mandelbrot Set

- Embarrassingly parallel, (but not as embarrassing as some).
- Each pixel is computed independently.
- How to divide up the domain?
 - Random? (why, too hard)
 - Pixel by pixel? (way more pixels than processors, so too much assignment (scalar) and gathering (communication).
 - Square grid? (its *embarrassing*, why make it hard—2D)
 - Go line-by-line, or chunk the lines.



Code Structure

Initialize MPI

Master Collects Data From Slaves

Slaves Compute and Send to Master

```
#include <mpi.h>
int X_RESN = 600, Y_RESN = 600; // Resolution
int main (int argc, char* argv[]) {
    MPI_Status status;
    int myid, nproc, ierr;

    MPI_Init(&argc, &argv);
    MPI_Comm_rank(MPI_COMM_WORLD, &myid);
    MPI_Comm_size(MPI_COMM_WORLD, &nproc);

    int nipp = X_RESN/(nproc-1); // lines per processor
    int pdata[nipp][600]; // processor chunk data

    if (myid == 0) { //===== MASTER
        for(int npdone=0; npdone < nproc-1; npdone++) {
            MPI_Recv(&pdata[0][0], 600*nipp, MPI_INT, MPI_ANY_SOURCE,
                    MPI_ANY_TAG, MPI_COMM_WORLD, &status);

            // GRAPH DATA
        }
    }
    else { //===== SLAVES
        // Compute the pdata on the processor
        MPI_Send(&pdata[0][0], 600*nipp, MPI_INT, 0, 0, MPI_COMM_WORLD);
    }

    MPI_Finalize();
    return 0;
}
```



Send and Receive Messages

```

MPI_Send(void*      message,      MPI_Send(&data[0][0],
        int          count,        600*nlp,
        MPI_Datatype datatype,    MPI_INT,
        int          dest,        0,
        int          tag,        0,
        MPI_COMM    comm)        MPI_COMM_WORLD)
  
```

```

MPI_Recv(void*      message,      MPI_Recv(&data[0][0],
        int          count,        600*nlp,
        MPI_Datatype datatype,    MPI_INT,
        int          source,      MPI_ANY_SOURCE,
        int          tag,        MPI_ANY_TAG,
        MPI_COMM    comm,        MPI_COMM_WORLD,
        MPI_Status   status)
  
```

MPI_ANY_SOURCE is a built-in variable, but could be just an integer corresponding to a given processor.

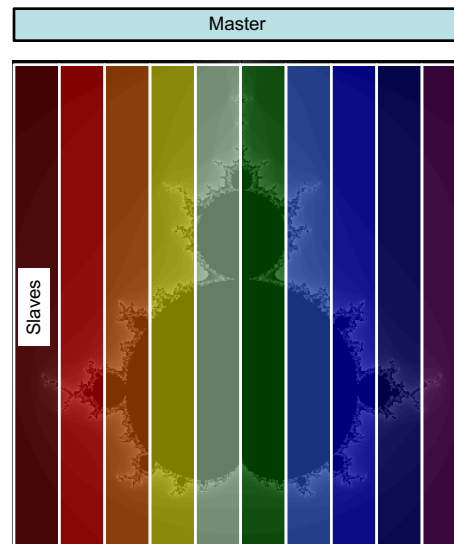
Same with MPI_ANY_TAG.

The status variable is declared as MPI_Status status.



Load Balancing

- Master-slave arrangement.
- One processor coordinates, gathers, directs, the others compute.
- Cost is not uniform over domain
 - Some processors will finish before others, and waste time.
- How to fix?
 - Divide into smaller chunks and farm them out as processors become available.
 - One could even dynamically determine workload and adjust optimal chunk size on the fly.



Resources

The screenshot displays three browser windows:

- Left Window:** "Ira and MaryLou Fulton Supercomputing Laboratory - Resources". It lists various system configurations:
 - manylinux:** 2482 processor cores, 4884 GB total memory, benchmarked at approximately 12.5 TFlops.
 - manylinux_Quad-core nodes:** 623 compute nodes, 2 x Quad-core Intel EM64T processors 2.6 GHz (4 cores) each, 8 GB memory per node.
 - manylinux_Quad-core nodes:** 256 processor cores, 512 GB total memory.
 - quad:** 32 compute nodes, 2 x Quad-core Intel Harpertown EM64T processors 2.5 GHz (8 cores) each, 16 GB memory per node.
 - Large memory (128 GB) nodes:** 32 processor cores, 256 GB total memory.
 - memf22:** 2 compute nodes, 4 x Quad-core AMD Opteron 8356 processors (16 cores) each, 128 GB memory per node.
- Middle Window:** "Introduction to Parallel Computing" by Blake Barney, Lawrence Livermore National Laboratory. It includes a "Table of Contents" with sections like:
 - 1. Abstract
 - 2. Overview
 - 3. Concepts and Terminology
 - 4. Parallel Computer Memory Architectures
 - 5. Shared Memory
 - 6. Distributed Memory
 - 7. Shared Memory Model
 - 8. Message Passing Model
 - 9. Data Parallel Model
 - 10. Limits and Costs of Parallel Programming
 - 11. Performance Analysis and Tuning
 - 7. Parallel Examples
- Right Window:** "Home | TOP500 Supercomputer Sites". It features a "TOP 10 Systems - 11/2008" list:
 - Roadrunner - BladeCenter Q320/S321 Cluster, PowerCell R1 3.2 GHz / Opteron OC 1.8 GHz - Volume Infiniband
 - Jaguar - Cray XT5 QC 2.3 GHz
 - Pleiades - SGI Altix ICE 820EX, Xeon QC 3.0/2.8 GHz
 - BlueGene/L - eServer Blue Gene Solution
 - Blue Gene/P Solution
 - Ranger - SunBlade x6420, Opteron QC 2.3 GHz, Infiniband
 - Franklin - Cray XT4 QuadCore 2.1 GHz
 - Jaguar - Cray XT4 QuadCore 2.1 GHz
 - Red Storm - Sandia Cray Red Storm, XT3A, 2.4/2.2 GHz dual/quad core
 - Dawning 2000A - Dawning 2000A, GC Opteron 1.3 GHz, Infiniband, Windows HPC 2008



- ✓ Many online tutorials for MPI and parallel programming.
- ✓ Pacheco "Parallel Programming with MPI" is a good introductory text