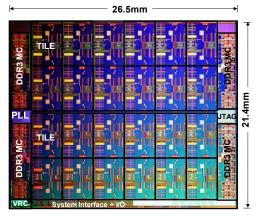
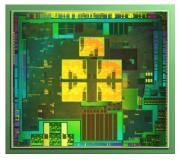


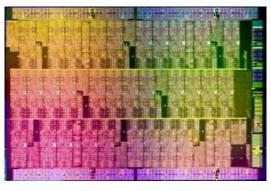
NVIDIA GTX 480 processor



Intel labs 48 core SCC processor



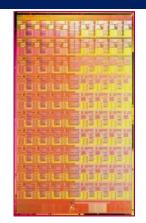
NVIDIA Tegra 3 (quad Arm Corex A9 cores + GPU)



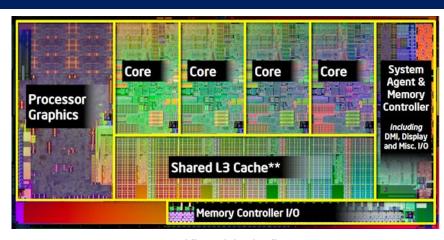
An Intel MIC processor

# **Introduction to Parallel Computing**

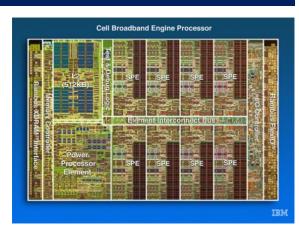
#### Tim Mattson (Intel Labs)



Intel Labs 80 core Research processor



Intel "Sandybridge" processor



IBM Cell Broadband engine processor

Third party names are the property of their owners

Other than the Intel lab's research processors. Die photos from UC Berkeley CS194 lecture notes

# **Disclaimer**READ THIS ... its very important



- The views expressed in this talk are those of the speakers and not their employer.
- This is an academic style talk and does not address details of any particular Intel product. You will learn nothing about Intel products from this presentation.
- This was a team effort, but if I say anything really stupid, it's my fault ... don't blame my collaborators.



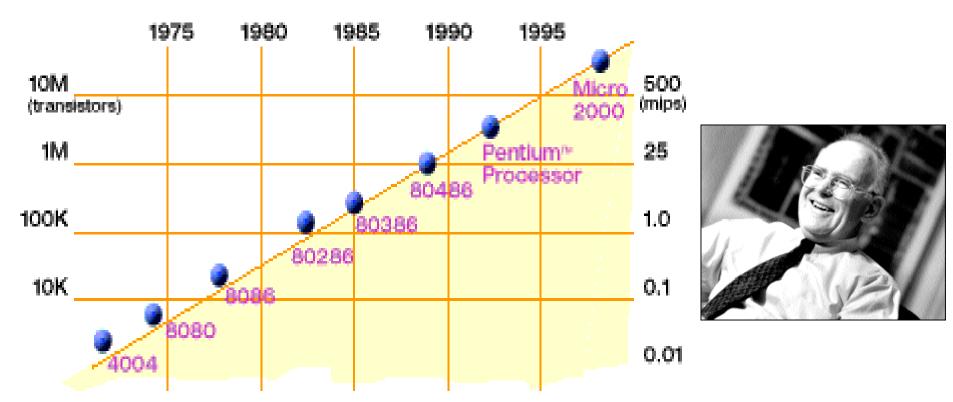
Slides marked with this symbol were produced-with Kurt Keutzer and his team for CS194 ... A UC Berkeley course on Architecting parallel applications with Design Patterns.

#### **Outline**

- Motivation: We all must be parallel programmers
  - Key concepts in parallel Computing
  - An introduction to parallel hardware
  - Software for parallel systems: key design patterns
  - Closing comments

# **Moore's Law**

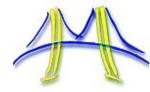


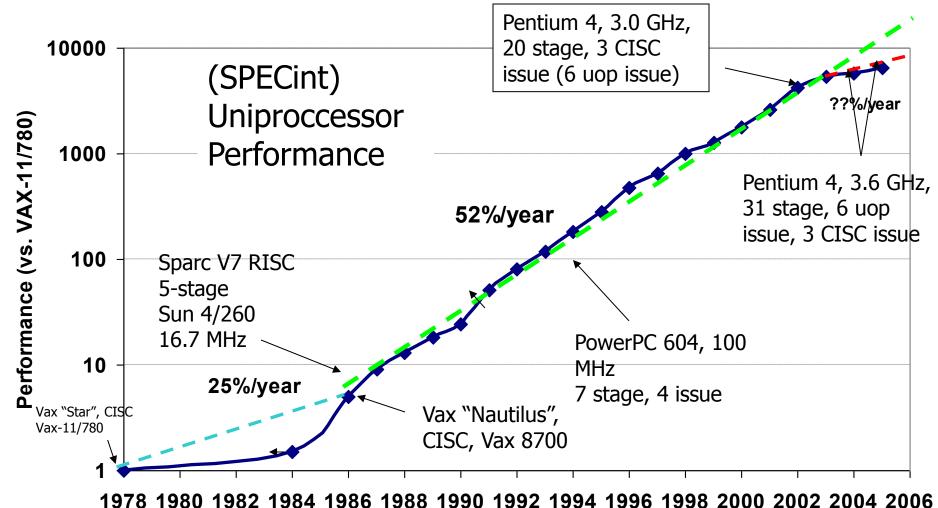


- In 1965, Intel co-founder Gordon Moore predicted (from just 3 data points!) that semiconductor density would double every 18 months.
  - He was right! Transistors are still shrinking at the same rate

Slide source: UCB CS 194 Fall'2010

### The good old days ...

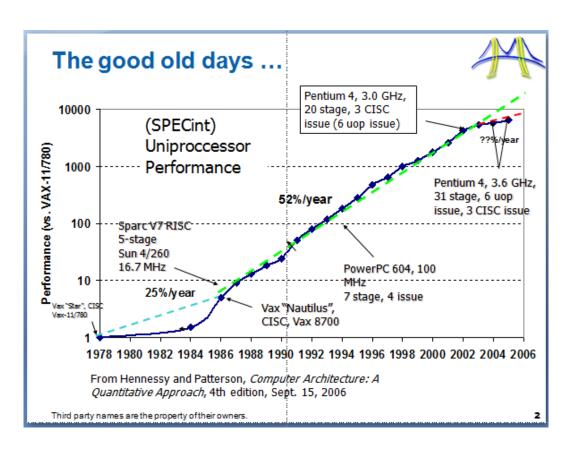




From Hennessy and Patterson, *Computer Architecture: A Quantitative Approach*, 4th edition, Sept. 15, 2006

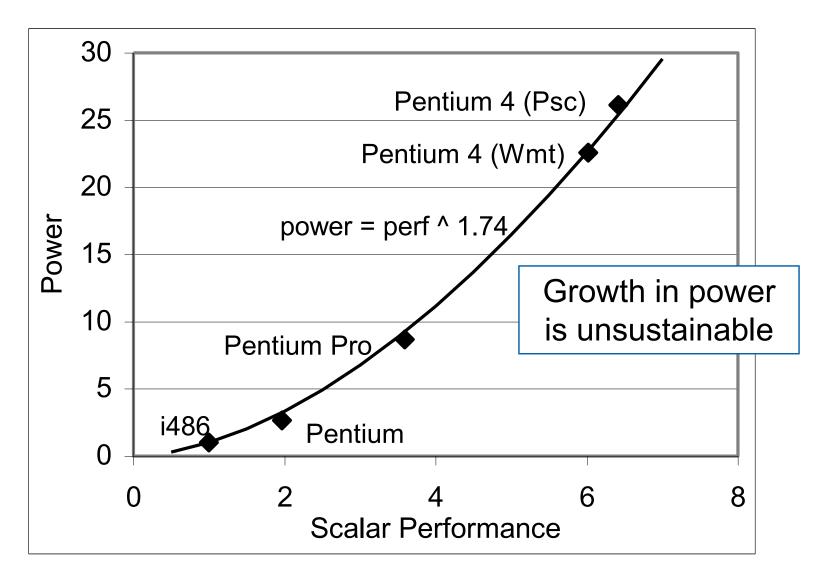
#### The Hardware/Software contract

 Write your software as you choose and we HW-geniuses will take care of performance.



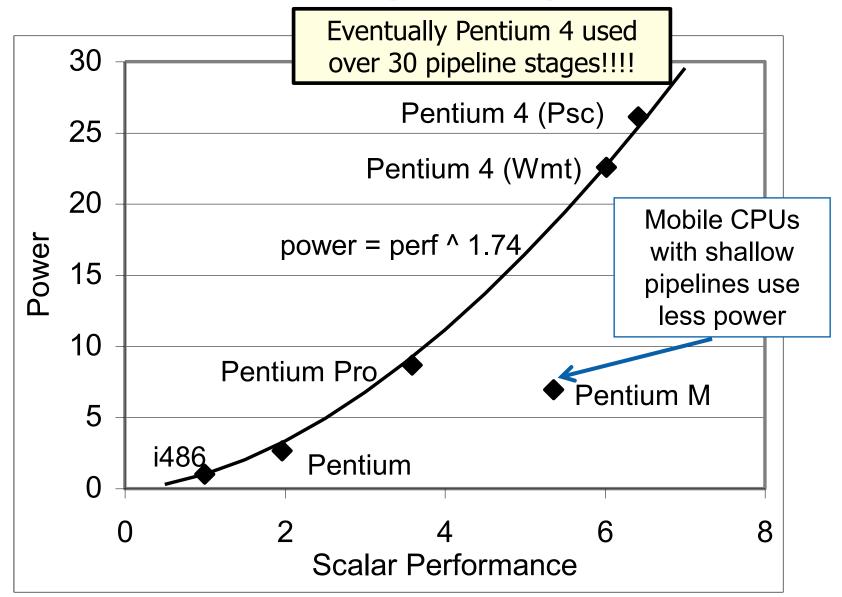
• The result: Generations of performance ignorant software engineers using performance-handicapped languages (such as Java) ... which was OK since performance was a HW job.

#### ... Computer architecture and the power wall



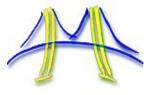
Source: E. Grochowski of Intel

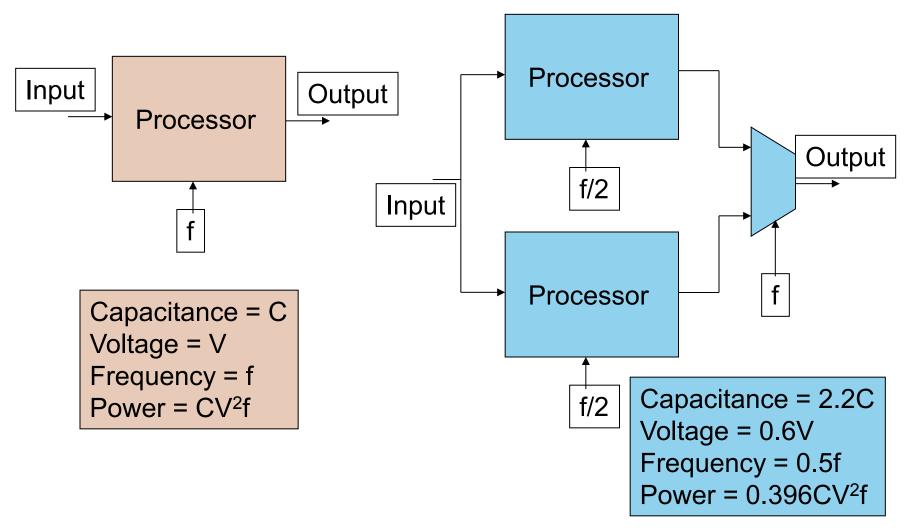
#### ... partial solution: simple low power cores



Source: E. Grochowski of Intel

#### ... The rest of the solution add cores



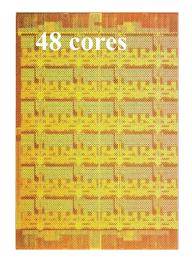


Chandrakasan, A.P.; Potkonjak, M.; Mehra, R.; Rabaey, J.; Brodersen, R.W., "Optimizing power using transformations," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*,, vol.14, no.1, pp.12-31, Jan 1995

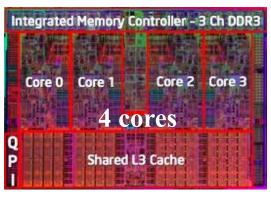
Source: Vishwani Agrawal

### Microprocessor trends

Individual processors are many core (and often heterogeneous) processors.



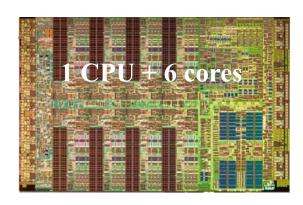
**Intel SCC Processor** 



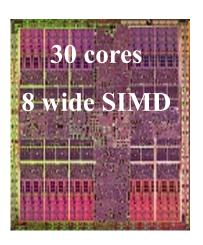
Intel Nehalem



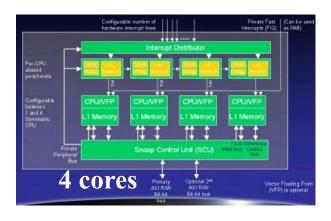
**ATI RV770** 



IBM Cell

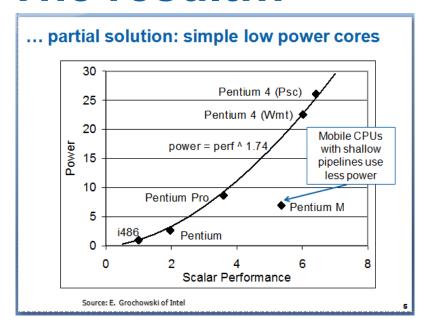


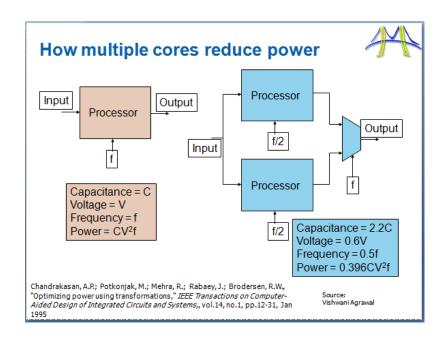
NVIDIA Tesla C1060



ARM MPCORE

#### The result...





A new contract ... HW people will do what's natural for them (lots of simple cores) and SW people will have to adapt (rewrite everything)

The problem is this was presented as an ultimatum ... nobody asked us if we were OK with this new contract ... which is kind of rude.

### The many core challenge

- A harsh assessment ...
  - We have turned to multi-core chips <u>not</u> because of the success of our parallel software but because of <u>our failure</u> to continually increase CPU frequency.
- Result: a fundamental and dangerous (for the computer industry) mismatch
  - □ Parallel hardware is ubiquitous.
  - □ Parallel software is rare
- The Many Core challenge ...
  - Parallel software must become as common as parallel hardware

Fortunately, we don't have to start over "from scratch".

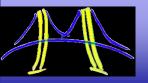
We can draw from past experience with parallelism
from high performance computing

#### **Outline**

- Motivation: We all must be parallel programmers
- Key concepts in parallel Computing
  - An introduction to parallel hardware
  - Software for parallel systems: key design patterns
  - Closing comments

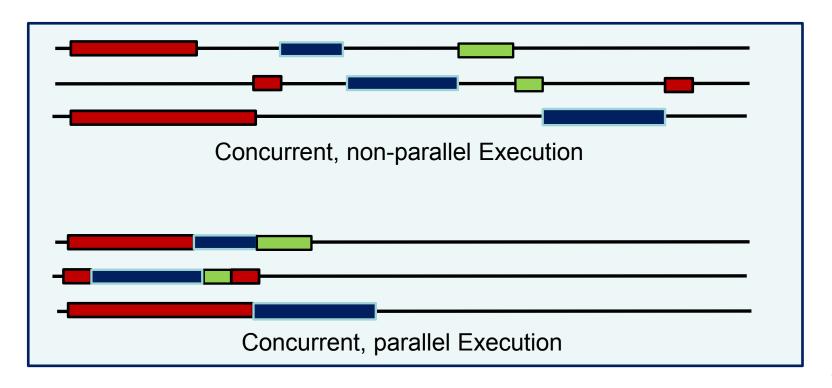
#### **Outline**

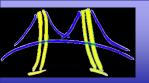
- Motivation: We all must be parallel programmers
- Key concepts in parallel Computing
- Basic definitions: Parallelism and Concurrency
  - Notions of parallel performance
  - The limits of scalability
  - Sources of parallel overhead
- An introduction to parallel hardware
- Software for parallel systems: key design patterns
- Closing comments



## Concurrency vs. Parallelism

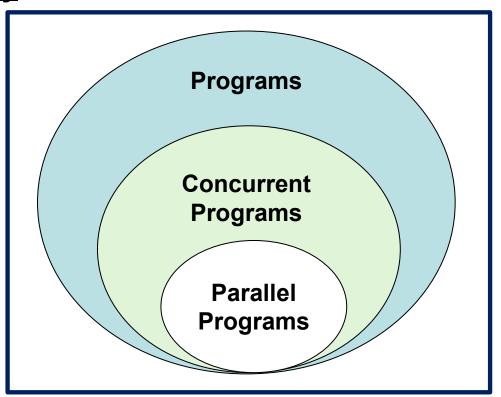
- Two important definitions:
  - Concurrency: A condition of a system in which multiple tasks are logically active at one time.
  - Parallelism: A condition of a system in which multiple tasks are <u>actually</u> active at one time.

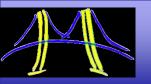




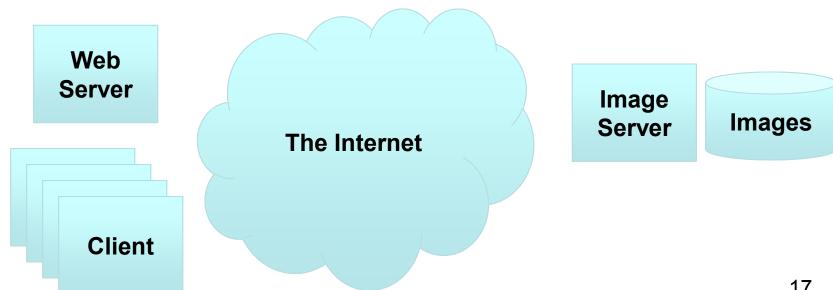
### Concurrency vs. Parallelism

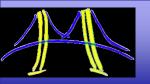
- Two important definitions:
  - Concurrency: A condition of a system in which multiple tasks are logically active at one time.
  - Parallelism: A condition of a system in which multiple tasks are <u>actually</u> active at one time.



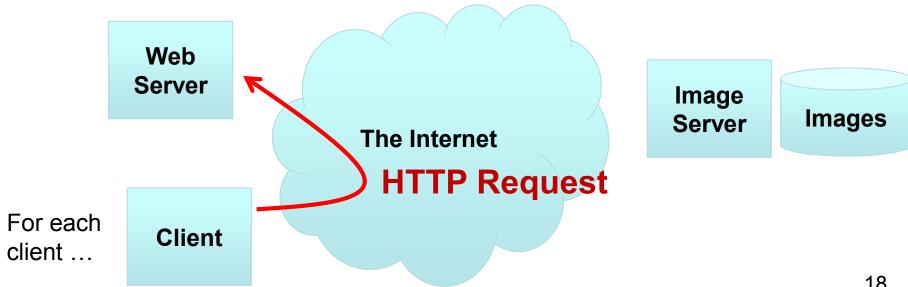


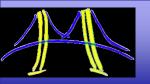
- An Web Server is a Concurrent Application (the problem is fundamentally defined in terms of concurrent tasks):
  - An arbitrary, large number of clients make requests which reference per-client persistent state
- Consider an Image Server, which relieves load on primary web servers by storing, processing, and serving only images



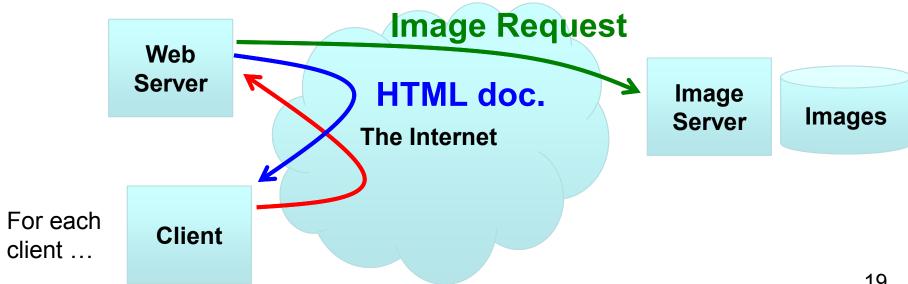


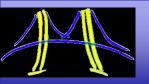
- A **Web Server** is a Concurrent Application (the problem is fundamentally defined in terms of concurrent tasks):
  - An arbitrary, large number of clients make requests which reference per-client persistent state
- Consider an Image Server, which relieves load on primary web servers by storing, processing, and serving only images



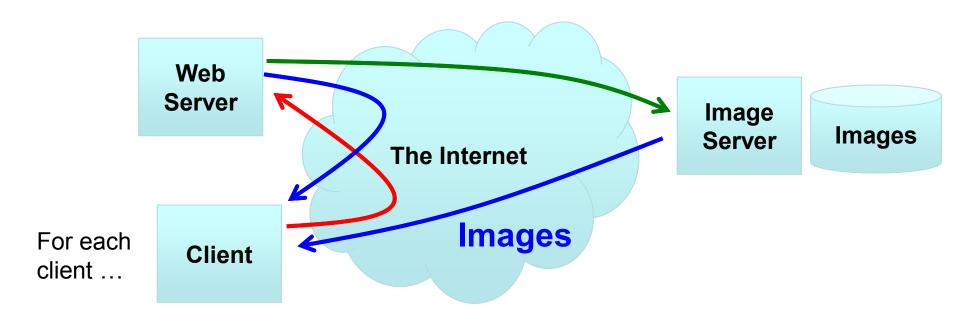


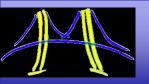
- A **Web Server** is a Concurrent Application (the problem is fundamentally defined in terms of concurrent tasks):
  - An arbitrary, large number of clients make requests which reference per-client persistent state
- Consider an Image Server, which relieves load on primary web servers by storing, processing, and serving only images



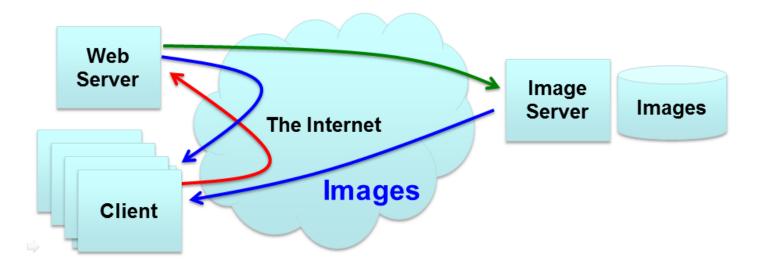


- A Web Server is a Concurrent Application (the problem is fundamentally defined in terms of concurrent tasks):
  - An arbitrary, large number of clients make requests which reference per-client persistent state
- Consider an Image Server, which relieves load on primary web servers by storing, processing, and serving only images



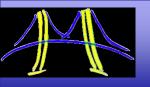


The Web server, image server, and clients (you have to plan on having many clients) all execute at the same time



The problem of one or more clients interacting with a web server not only contains concurrency, the problem is fundamentally current. It doesn't exist as a serial problem.

**Concurrent application**: An application for which the problem definition is fundamentally concurrent.

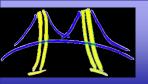


The Mandelbrot set: An iterative map in the complex plane

 $z_{n+1} = z_n^2 + c$ 

Color each point in the complex plain of C values based on convergence or divergence of the iterative map.

 $z_0 = 0$ , c is constant Cimaginary



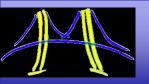
```
int mandel (complex C) {
  int n;
 double a = C.real();
 double b = C.imag();
 double zr = 0.0, zi = 0.0;
 double tzr, tzi;
  n = 0:
 while (n < max_iters && sqrt (zr*zr + zi*zi) < t) {
   tzr = (zr*zr - zi*zi) + a;
    tzi = (zr*zi + zr*zi) + b;
   zr = tzr;
   zi = tzi;
   n = n+1;
  return n;
```

Function to compute the iterative map for a single point C where

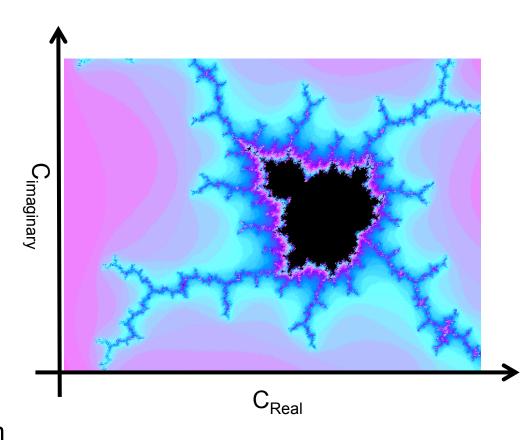
$$C = a + b * i$$

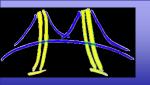
Where i is the square root of (-1)

"t" is a constant that defines a threshold beyond which we consider the iterative map to diverge.



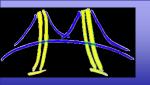
- To generate the famous Mandelbrot set image, we use the function mandel(C) where C comes from the points in the complex plane.
- At each point C, use n=mandel(C) to determine if:
  - The map converges (n=max\_iters), assign the color black
  - The map diverges (n<max\_iters), assign the color based on the value of n
- The computation for each point is independent of all the other points ... a so-called embarrassingly parallel problem





The following is simplified code for the serial Mandelbrot program.

```
for (i=0; i<N; i++){
    for (j=0; j<N; j++) {
        complex c = get_const_at_pixel(i,j);
        complex image[i][j] = mandel( c);
    }
}</pre>
```



- The following is simplified code for the serial Mandelbrot program.
- Loop iterations are independent, so we can create a parallel version of this program as follows ...

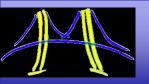
```
for (i=0; i<N; i++)

Combine the two loops into one big loop and execute them in parallel

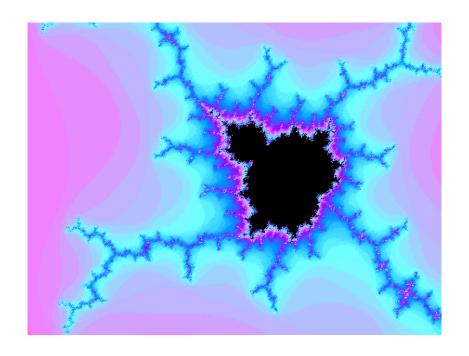
complex c = get_const_at_pixel(i,j);

complex image[i][j] = mandel( c);

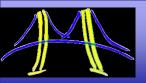
}
```



- The problem of generating an image of the Mandelbrot set can be viewed serially.
- We choose to exploit the concurrency contained in this problem so we can generate the image in less time



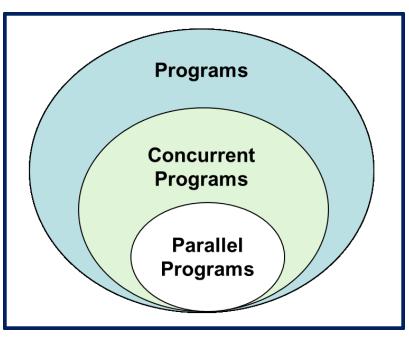
**Parallel application**: An application composed of tasks that actually execute concurrently in order to (1) consider larger problems in fixed time or (2) complete in less time for a fixed size problem.



### Concurrency vs. Parallelism: wrap up

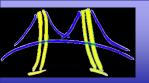
#### Key points:

- A web server had concurrency in its problem definition ... it doesn't make sense to even think of writing a "serial web server".
- The Mandelbrot program didn't have concurrency in its problem definition. It would take a long time, but it could be serial

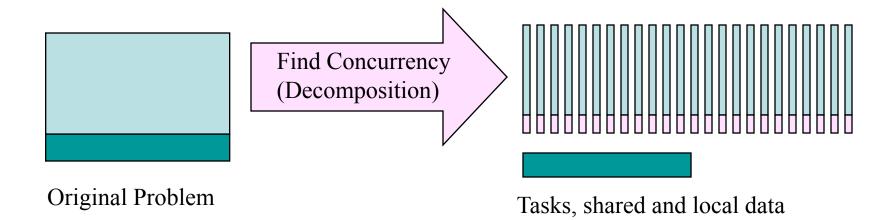


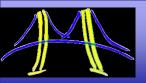
#### Both cases use concurrency:

- A concurrent application is concurrent by definition.
- A parallel application solves a problem that could be serial, but it is run in parallel by ...
  - 1. find concurrency in the problem
  - 2. expose the concurrency in the source code.
  - 3. exploit the exposed concurrency to complete a job in less time.

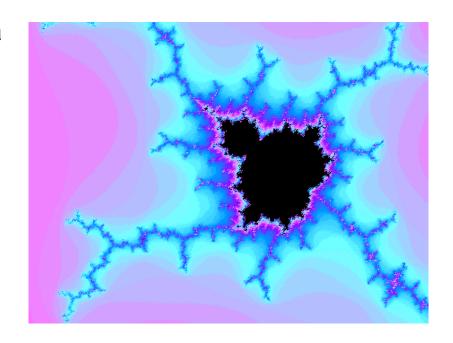


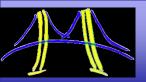
# The Parallel programming process:





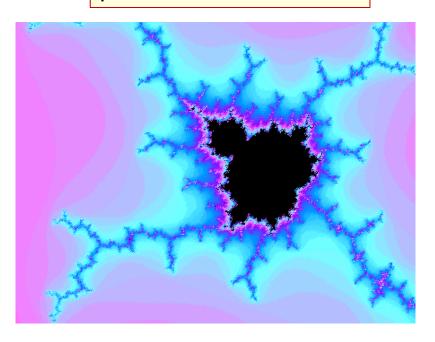
- Every parallel program is based on concurrency ... i.e. tasks defined by an application that can run at the same time.
- EVERY parallel program requires a task decomposition and a data decomposition:
  - Task decomposition: break the application down into a set of tasks that can execute concurrently..
  - Data decomposition: How must the data be broken down into chunks and associated with threads/processes to make the parallel program run efficiently.

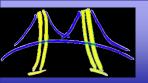




- Every parallel program is based on concurrency ... i.e. tasks defined by an application that can run at the same time.
- EVERY parallel program requires a task decomposition and a data decomposition:
  - Task decomposition: break the application down into a set of tasks that can execute concurrently..
  - Data decomposition: How must the data be broken down into chunks and associated with threads/processes to make the parallel program run efficiently.

What's a task decomposition for this problem?





- Every parallel program is based on concurrency ... i.e. tasks defined by an application that can run at the same time.
- EVERY parallel program requires a task decomposition and a data decomposition:
  - Task decomposition: break the application down into a set of tasks that can execute

```
for (i=0; i<N; i++){

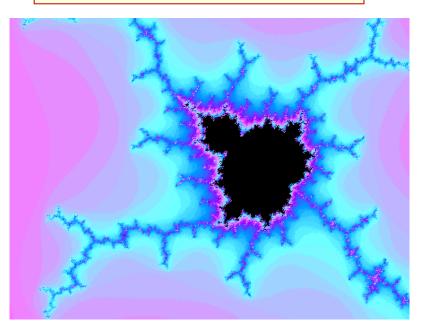
for (j=0; j<N; j++) {

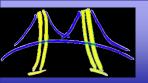
complex c = get_const_at_pixel(i,j);

complex image[i][j] = mandel( c);

}
```

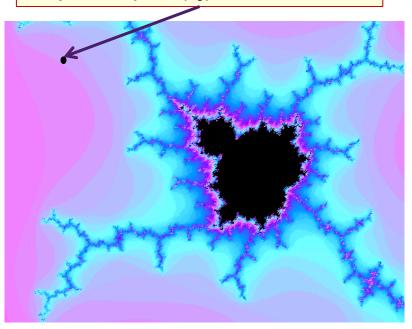
Hint: Think of the source code and work that is compute-intensive that can execute independently

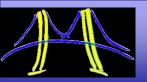




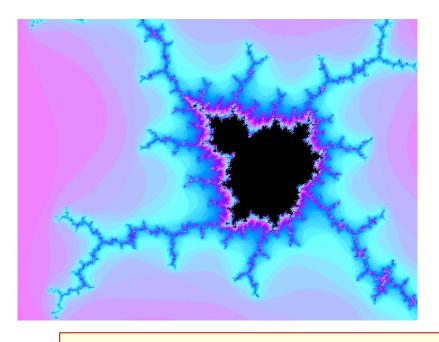
- Every parallel program is based on concurrency ... i.e. tasks defined by an application that can run at the same time.
- EVERY parallel program requires a task decomposition and a data decomposition:
  - Task decomposition: break the application down into a set of tasks that can execute concurrently..
  - Data decomposition: How must the data be broken down into chunks and associated with threads/processes to make the parallel program run efficiently.

Task: the computation required for each pixel ... the body of the loop for a pair (i,j).

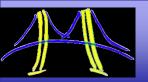




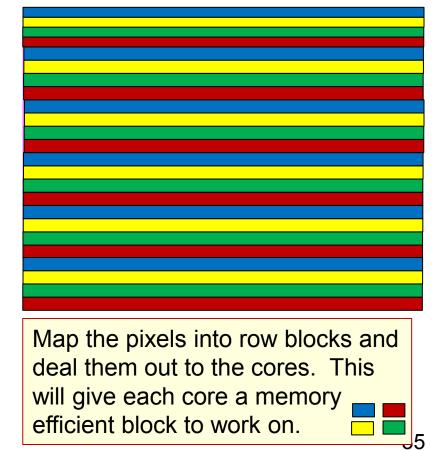
- Every parallel program is based on concurrency ... i.e. tasks defined by an application that can run at the same time.
- EVERY parallel program requires a task decomposition and a data decomposition:
  - Task decomposition: break the application down into a set of tasks that can execute concurrently..
  - Data decomposition: How must the data (the complex plain, C) be broken down into chunks and associated with threads/processes to make the parallel program run efficiently.

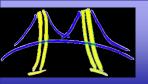


Suggest a data decomposition for this problem ... assume a quad core shared memory PC.



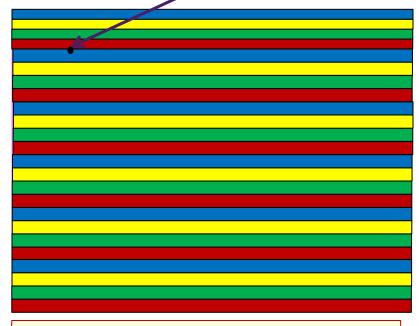
- Every parallel program is based on concurrency ... i.e. tasks defined by an application that can run at the same time.
- EVERY parallel program requires a task decomposition and a data decomposition:
  - Task decomposition: break the application down into a set of tasks that can execute concurrently..
  - Data decomposition: How must the data (the complex plain, C) be broken down into chunks and associated with threads/processes to make the parallel program run efficiently.



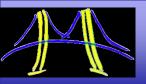


- Every parallel program is based on concurrency ... i.e. tasks defined by an application that can run at the same time.
- **EVERY parallel program** requires a <u>task decomposition</u> and a <u>data</u> <u>decomposition</u>:
  - Task decomposition: break the application down into a set of tasks that can execute concurrently..
  - Data decomposition: How must the data (the complex plain, C) be broken down into chunks and associated with threads/processes to make the parallel program run efficiently.

But given this data decomposition, it is effective to think of a task as the update to a pixel? Should we update our task definition given the data decomposition?



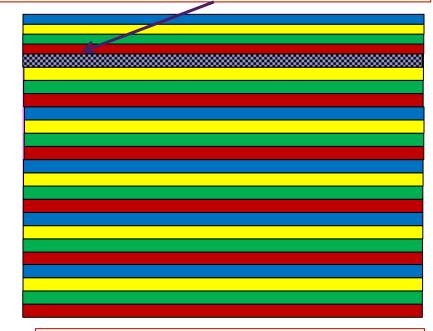
Map the pixels into row blocks and deal them out to the cores. This will give each core a memory efficient block to work on.



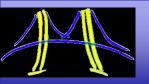
### Decomposition in parallel programs

- Every parallel program is based on concurrency ... i.e. tasks defined by an application that can run at the same time.
- **EVERY parallel program** requires a <u>task decomposition</u> and a <u>data</u> <u>decomposition</u>:
  - Task decomposition: break the application down into a set of tasks that can execute concurrently..
  - Data decomposition: How must the data (the complex plain, C) be broken down into chunks and associated with threads/processes to make the parallel program run efficiently.

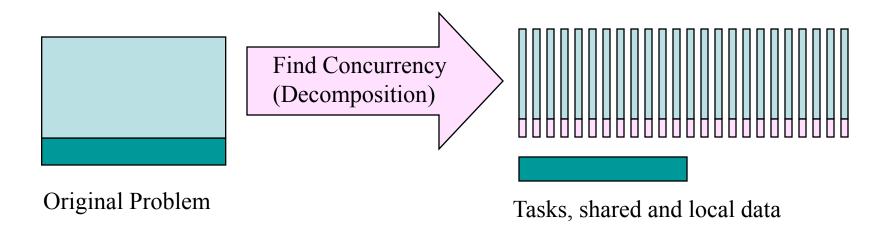
Yes. You go back and forth between task and data decomposition until you have a pair that work well together. In this case, let's define a task as the update to a row-block

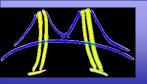


Map the pixels into row blocks and deal them out to the cores. This will give each core a memory efficient block to work on.

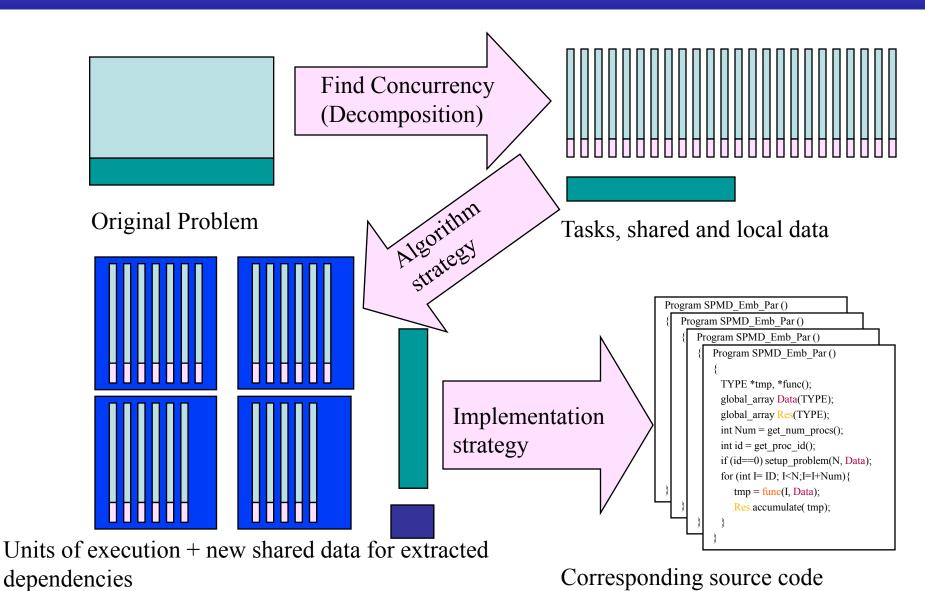


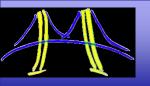
# The Parallel programming process:



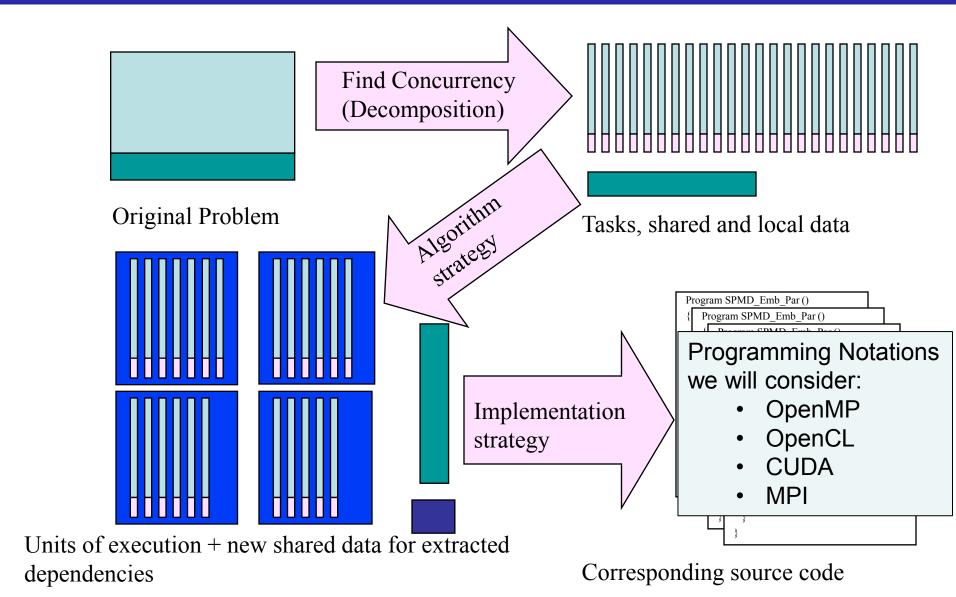


## The Parallel programming process:



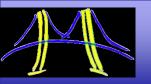


## The Parallel programming process:



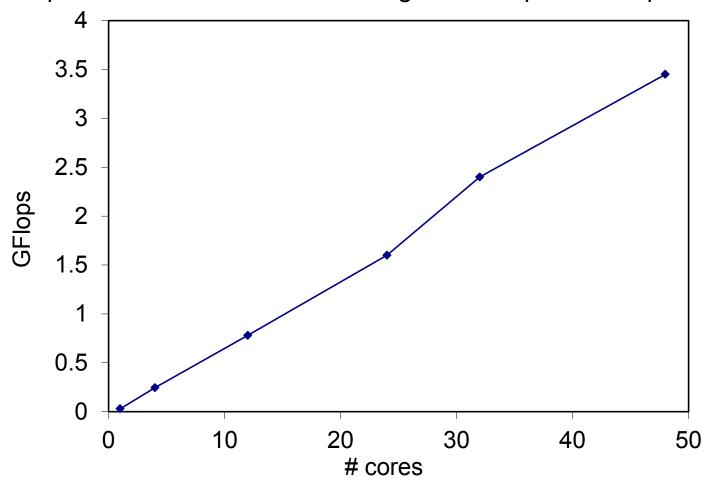
#### **Outline**

- Motivation: We all must be parallel programmers
- Key concepts in parallel Computing
  - Basic definitions: Parallelism and Concurrency
- Notions of parallel performance
  - The limits of scalability
  - Sources of parallel overhead
- An introduction to parallel hardware
- Software for parallel systems: key design patterns
- Closing comments

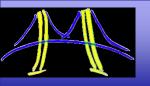


#### **Parallel Performance**

MP Linpack benchmark, order 1000 matrix (solve a dense system of linear equations ... the dense linear algebra computational pattern).



Intel SCC 48 processor, 500 MHz core, 1 GHz router, DDR3 at 800 MHz.



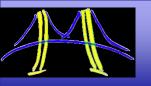
#### Talking about performance

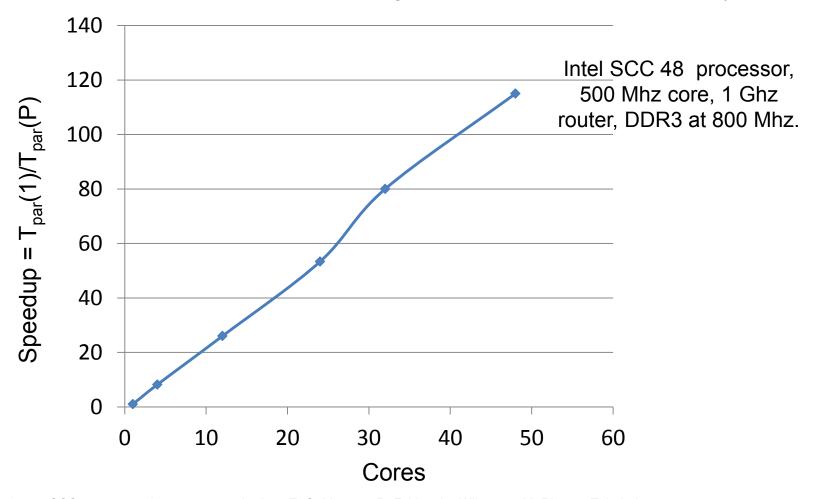
Speedup: the increased performance from running on P processors

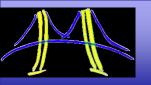
Perfect Linear Speedup: happens when no parallel overhead and algorithm is 100% parallel.

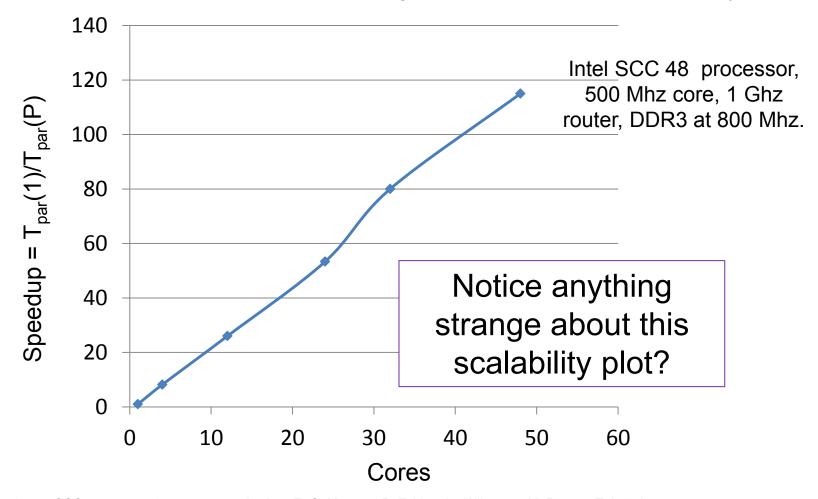
$$S(P) = \frac{Time_{seq}(1)}{Time_{par}(P)}$$

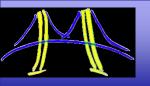
$$S(P) = P$$

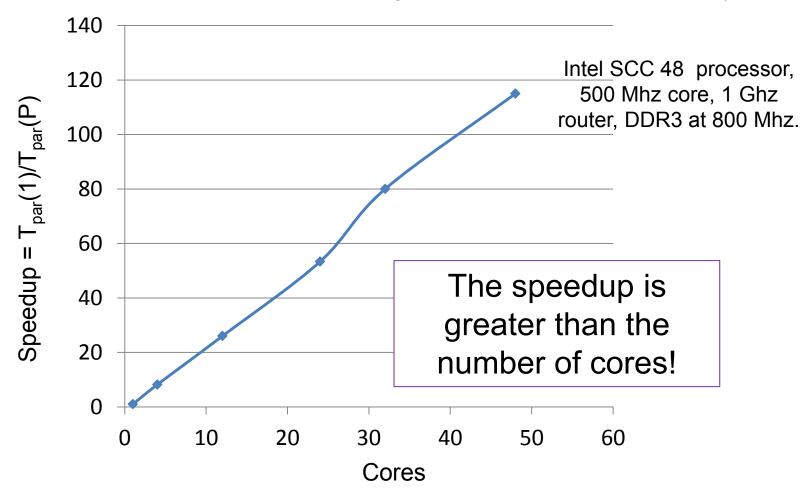


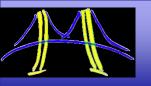












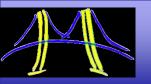
#### Talking about performance

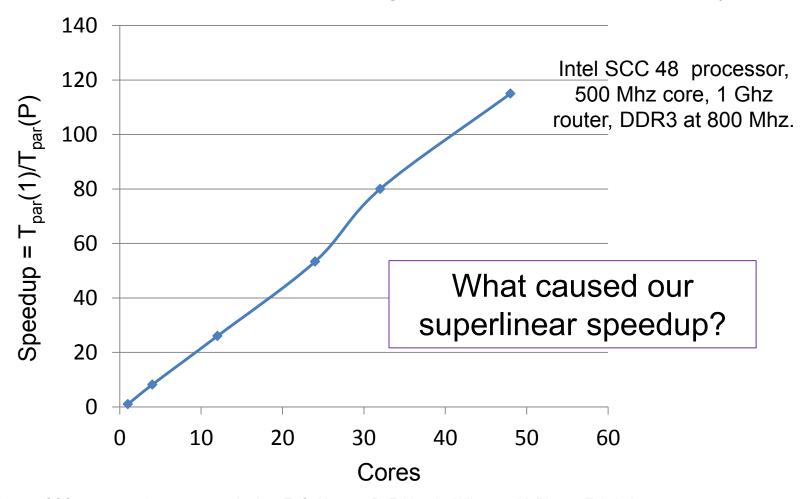
Speedup: the increased performance from running on P processors

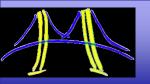
- Perfect Linear Speedup: happens when no parallel overhead and algorithm is 100% parallel.
- Super-linear Speedup: Speed grows faster than the number of processing elements

$$S(P) = \frac{Time_{seq}(1)}{Time_{par}(P)}$$

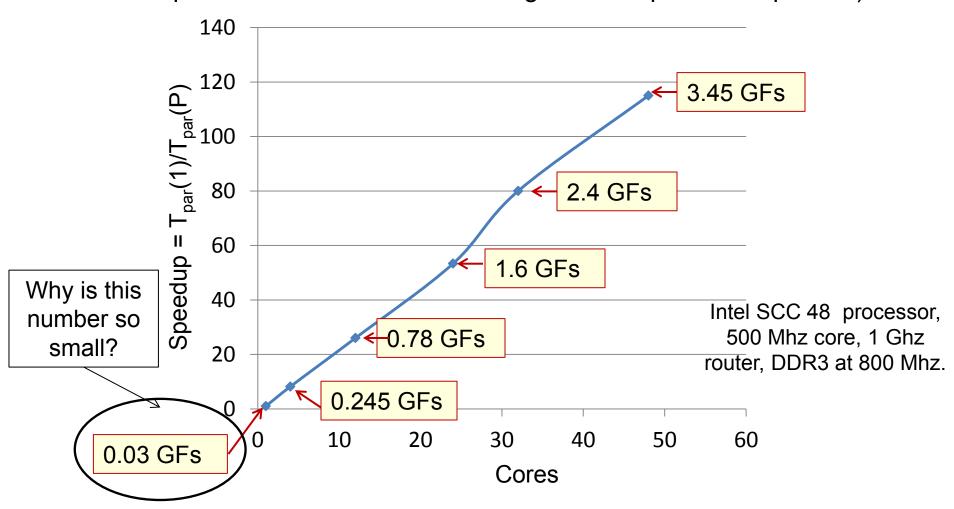
$$S(P) = P$$



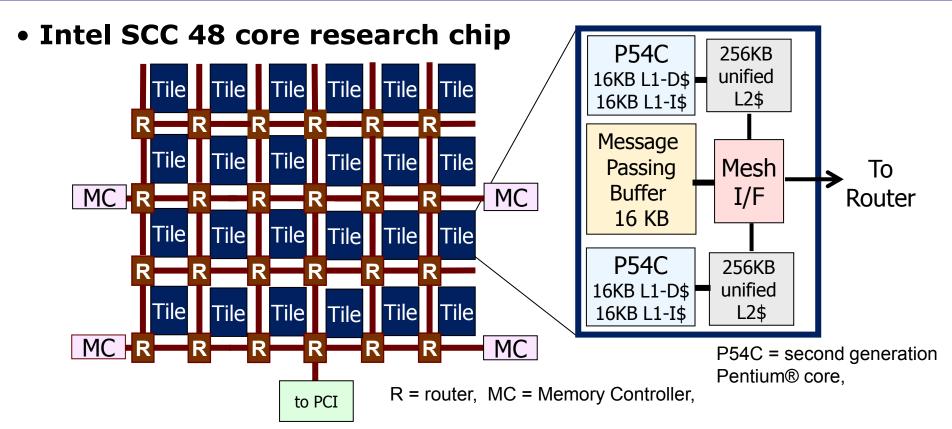




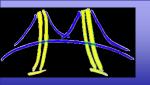
#### SuperLinear Speedup



# Why the Superlinear speedup?

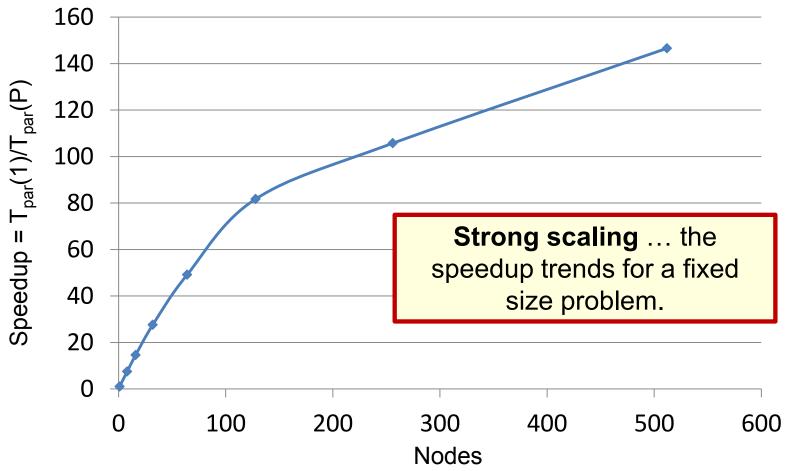


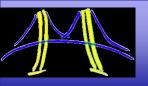
- SCC caches are so small, even a small portion of our O(1000) matrices won't fit.
  - Hence the single node performance measures memory overhead.
- As you add more cores, the aggregate cache size grows.
  - ➤ Eventually the tiles of the matrices being processed fits in the caches and performance sharply increases → superlinear speedup.



## A more typical speedup plot

■ CHARMM molecular dynamics program running the myoglobin benchmark on an Intel Paragon XP/S supercomputer with 32 Mbyte nodes running OSF R 1.2. (The nbody computational pattern). Speedup relative to running the parallel program on one node.



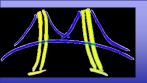


#### **Efficiency**

Efficiency measures how well the parallel system's resources are being utilized.

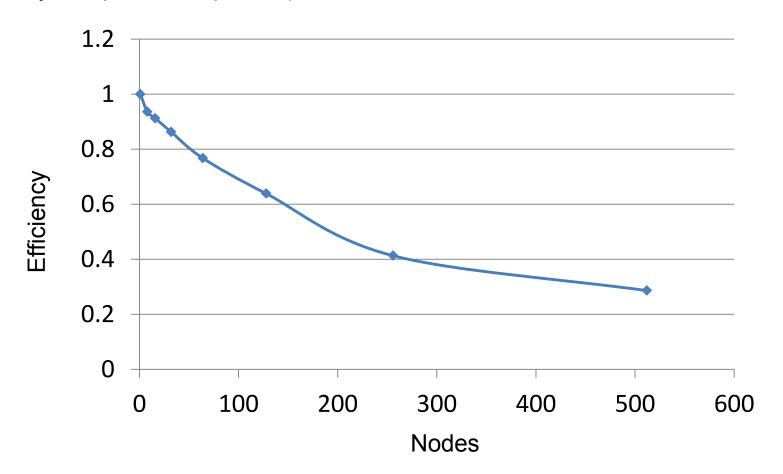
$$\varepsilon = \frac{Time_{seq}}{P*Time_{par}(P)} = \frac{S(P)}{P}$$

Where P is the number of nodes and T is the elapsed runtime.



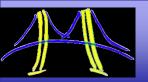
### **Efficiency**

■ CHARMM molecular dynamics program running the myoglobin benchmark on an Intel Paragon XP/S supercomputer with 32 Mbyte nodes running OSF R 1.2. (The nbody computational pattern). Speedup relative to running the parallel program on one node.



#### **Outline**

- Motivation: We all must be parallel programmers
- Key concepts in parallel Computing
  - Basic definitions: Parallelism and Concurrency
  - Notions of parallel performance
- The limits of scalability
  - Sources of parallel overhead
- An introduction to parallel hardware
- Software for parallel systems: key design patterns
- Closing comments



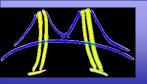
#### **Amdahl's Law: History**

- Gene Amdahl was a computer architect in the 1960's at IBM
- In 1967, refuted the idea that parallel computing was a practical path to improving program performance.
- Example: Compare these two systems



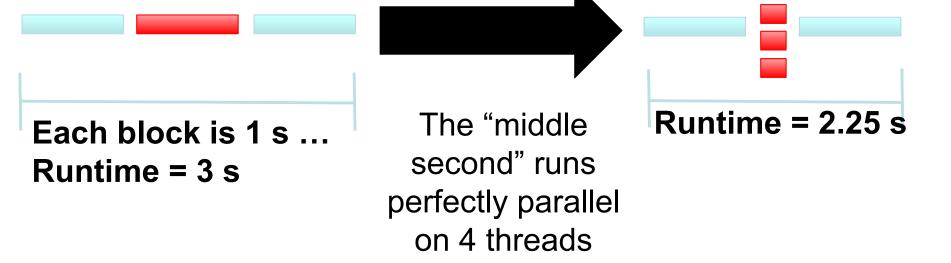
IBM System 360, ca. 1964

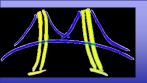
- The IBM System 360:
  - A single-processor machine, running at 16 MHz.
  - 1 FP addition per 60 ns cycle, and 1 FP mul in ~10 60 ns cycles, and execute multiple instructions simultaneously
- ILLIAC IV:
  - "The first Supercomputer" ... installed at NASA Ames in 1975.
  - 256 processors ... could perform 256 FP adds in 240 ns.



#### **Amdahl's Law**

- Clearly, the ILLIAC will run programs much faster than the S/360: It has 60x higher instruction throughput!
  - ... if you always have 256 independent instructions
- Amdahl argued that large portions of many programs are not parallelizable. Parallel hardware does not help serial code:





#### **Amdahl's Law**

- What is the maximum speedup you can expect from a parallel program?
- Consider a sequential program with runtime:

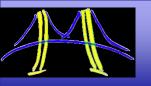
$$Time_{Seq}$$

- We can think of this program as consisting of two parts ... one that can benefit from multiple processing elements (parallel) and a second part that is fundamentally serial.
- The runtime is therefore:

$$Time_{seq} = Time_{Serial} + Time_{parallelizable}$$

We can express this in terms of a fraction of the program that is serial and a fraction of the program that is parallel or

 $Time_{seq} = serial\_fraction * Timeseq + parallel\_fraction * Time_{seq}$ 



#### Amdahl's Law

If we run the program on P processing elements and assume linear speedup, then our time for the parallel program becomes:

$$Time_{par}(P) = (serial\_fraction + \frac{parallel\_fraction}{P})*Time_{seq}$$

If the serial\_fraction is  $\alpha$  and the parallel\_fraction is (1- $\alpha$ ), the speedup is:

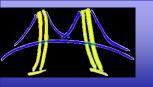
$$S(P) = \frac{Time_{seq}}{Time_{par}(P)} = \frac{Time_{seq}}{\left(\alpha + \frac{1 - \alpha}{P}\right) * Time_{seq}} = \frac{1}{(\alpha + \frac{1 - \alpha}{P})}$$

If you had an unlimited number of processors:

$$\lim_{P \to \infty} \frac{1 - \alpha}{P} = 0$$

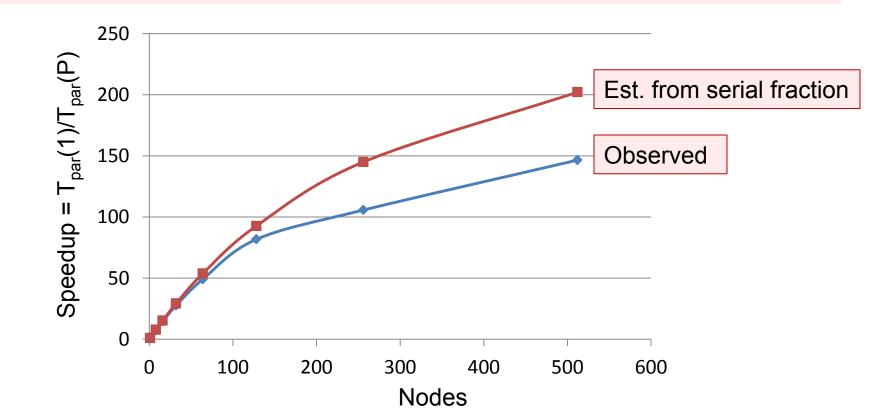
The maximum possible speedup is:

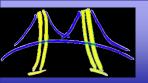
$$S = \frac{1}{\alpha}$$
 Amdahl's Law



# Amdahl's Law and the CHARMM MD program

- We Profiled CHARMM running on the Paragon XPS to find the time spent in code that was not parallelized ... concluded that CHARMM has a serial fraction of ~0.003.
  - The maximum possible speedup is: S = 1/0.003 = 333



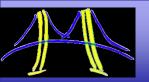


#### What if the problem size grows

- Consider the dense linear algebra
- A key feature is that operations between matrices (such as LU factorization or matrix multiplication) scale as the cube of the order of the matrix.
- Assume we can parallelize the linear algebra operation (O(N³)) but not the loading of the matrices from memory (O(N²)). How does the serial fraction vary with matrix order (assume loading from memory is much slower than a floating point op).

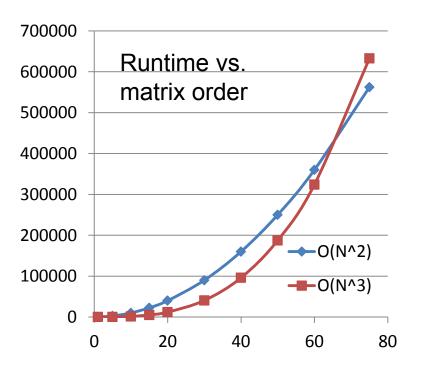
What would plots of runtime vs. problem size look like for the N squared and N cubed terms?

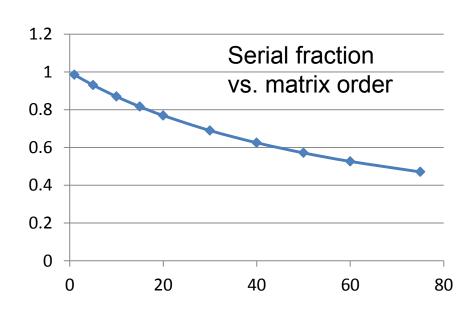
What would plots of serial fraction vs. problem size look like for the N squared and N cubed terms?

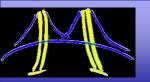


#### What if the problem size grows

- Consider dense linear algebra
- A key feature is that operations between matrices (such as LU factorization or matrix multiplication) scale as the cube of the order of the matrix.
- Assume we can parallelize the linear algebra operation (O(N³)) but not the loading of the matrices from memory (O(N²)). How does the serial fraction vary with matrix order (assume loading from memory is much slower than a floating point op).



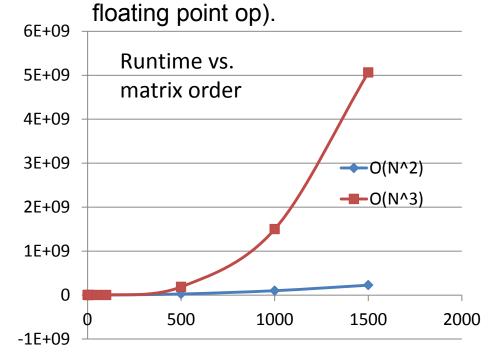


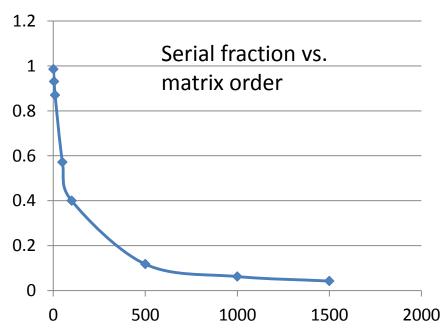


# What if the problem size grows

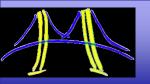
- Consider the dense linear algebra design pattern (which we will cover in much more detail later).
- A key feature is that operations between matrices (such as LU factorization or matrix multiplication) scale as the cube of the order of the matrix.

Assume we can parallelize the linear algebra operation (O(N³)) but not the loading of the matrices from memory (O(N²)). How does the serial fraction vary with matrix order (assume loading from memory is much slower than a





For much larger Matrix orders ...



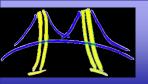
### Weak Scaling: a response to Amdahl

Gary Montry and John Gustafson (1988, Sandia National Laboratories) observed that for many problems the serial fraction of a function of the problem size (N) decreases:

$$S(P,N) = \frac{T_{seq}(1)}{(\alpha(N) + \frac{1 - \alpha(N)}{P}) * T_{seq}(1)} \lim_{N \to N_{large}} \alpha(N) = 0$$

$$S(P,N_{large}) \to P$$

- In other words ... if parallelizable computations asymptotically dominate the runtime, then you can increase a problem size until limitations due to Amdahl's law can be ignored. This is an easier form of scalability for a programmer to meet ... so its called "weak scaling":
  - Weak Scaling: Performance of an application when the problem size increases with the number of processors (fixed size problem per node)

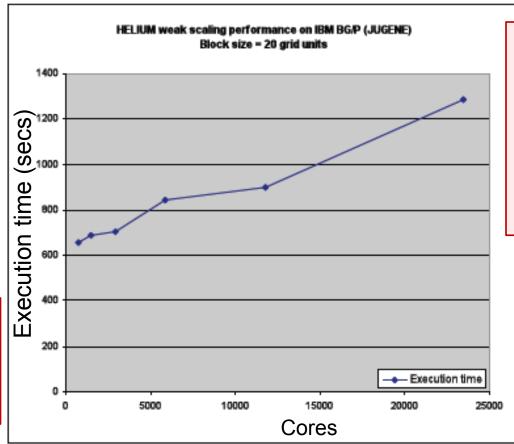


#### Example of weak scaling

#### HELIUM Weak Scaling Performance on BG/P

epcc

Local block size fixed to 20 grid units



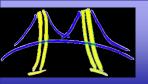
A time dependent Quantum simulation of helium atoms with 20 grid units per processing element.

IBM Blue Gene P, 0.85 GHz, PowerPC 450, 4way processors

May 13, 10

NAMD & HELIUM Enabling Work on the PRACE IBM Prototypes

34

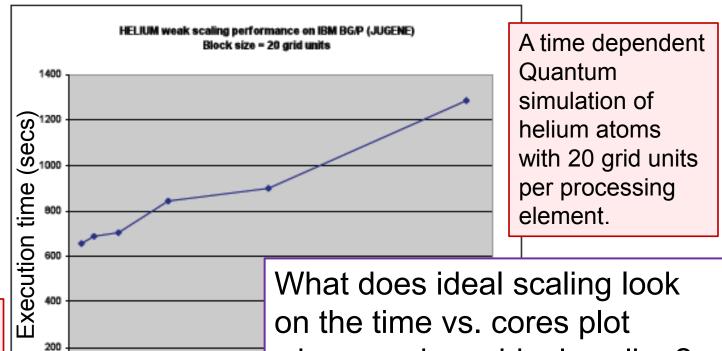


## Example of weak scaling

#### HELIUM Weak Scaling Performance on BG/P

epcc

Local block size fixed to 20 grid units



IBM Blue Gene P, 0.85 GHz, PowerPC 450, 4way processors

NAMD & HELIUM Enabling Work on the PRACE IBM Prototypes

15000

Cores

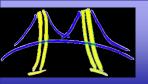
20000

10000

34

when you have ideal scaling?

5000

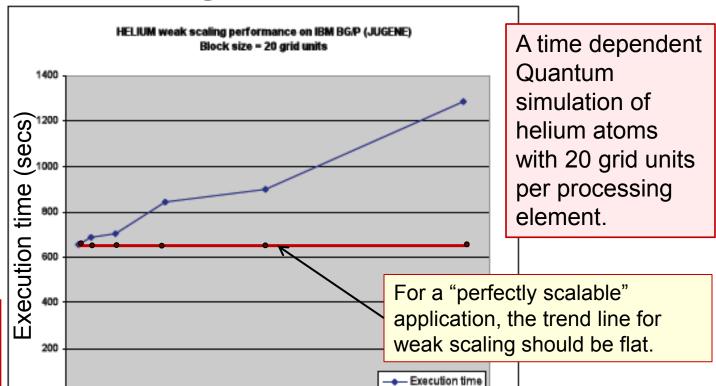


### Example of weak scaling

#### HELIUM Weak Scaling Performance on BG/P

epcc

Local block size fixed to 20 grid units



15000

Cores

20000

IBM Blue Gene P, 0.85 GHz, PowerPC 450, 4way processors

May 13, 10 NAMD & HELIUM Enabling Work on the PRACE IBM Prototypes

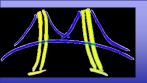
10000

34

5000

#### **Outline**

- Motivation: We all must be parallel programmers
- Key concepts in parallel Computing
  - Basic definitions: Parallelism and Concurrency
  - Notions of parallel performance
  - The limits of scalability
- Sources of parallel overhead
- An introduction to parallel hardware
- Software for parallel systems: key design patterns
- Closing comments



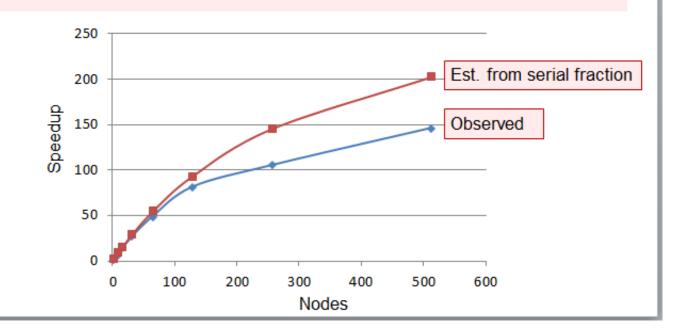
#### Limitations to scalability

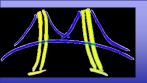
Remember the speedup plot we discussed earilier?

# 1

# Amdahl's Law and the CHARMM MD program

- We Profiled CHARMM running on the Paragon XPS to find the time spent in code that was not parallelized ... concluded that CHARMM has a serial fraction of ~0.003.
  - The maximum possible speedup is: S= 1/0.003 = 333





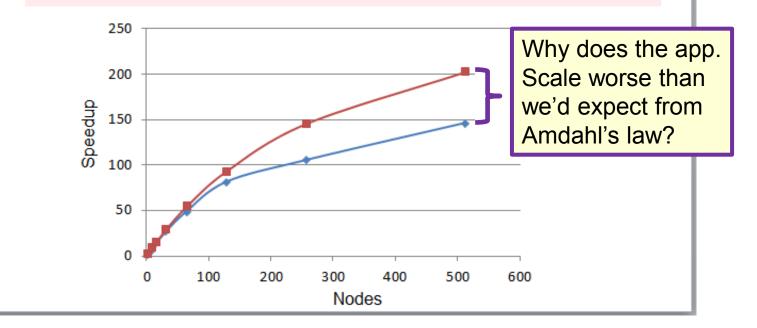
#### Limitations to scalability

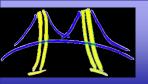
Remember the speedup plot we discussed from last time?

# **M**

# Amdahl's Law and the CHARMM MD program

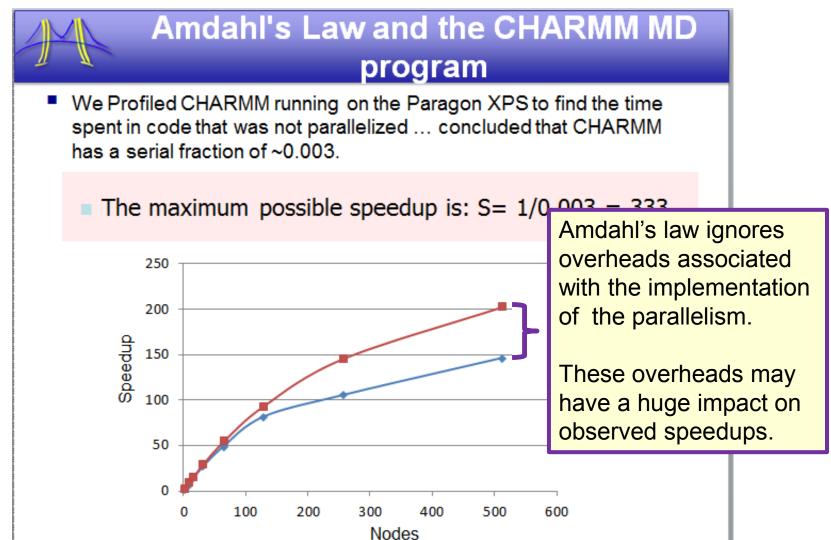
- We Profiled CHARMM running on the Paragon XPS to find the time spent in code that was not parallelized ... concluded that CHARMM has a serial fraction of ~0.003.
  - The maximum possible speedup is: S= 1/0.003 = 333

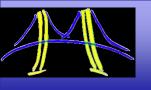




#### Limitations to scalability

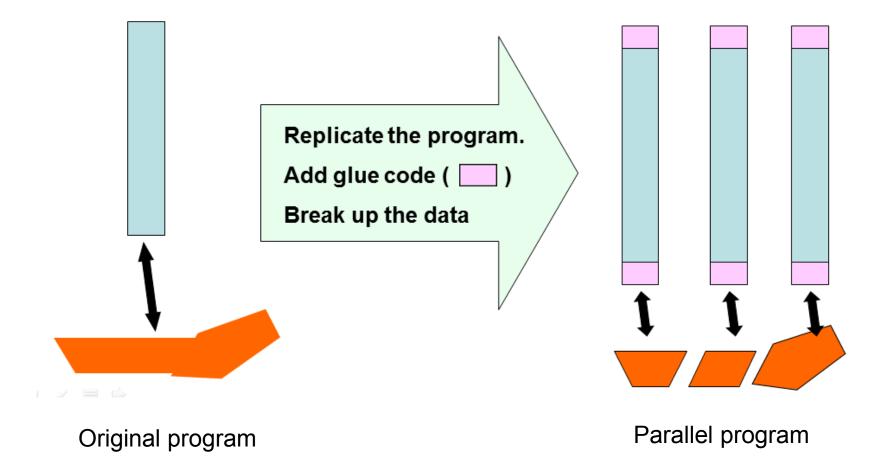
Remember the speedup plot we discussed from last time?

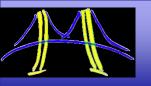




# Parallel overheads: The algorithmic structure of many HPC codes (part 1)

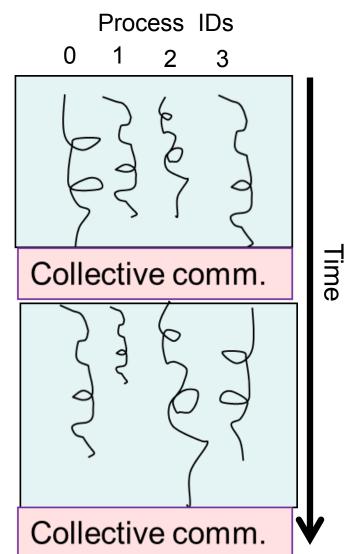
A large fraction of HPC applications (such as CHARMM) use a message passing notation with the Single Program Multiple Data or SPMD design pattern.

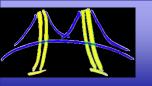




# Parallel overheads: The algorithmic structure of many HPC codes (part 2)

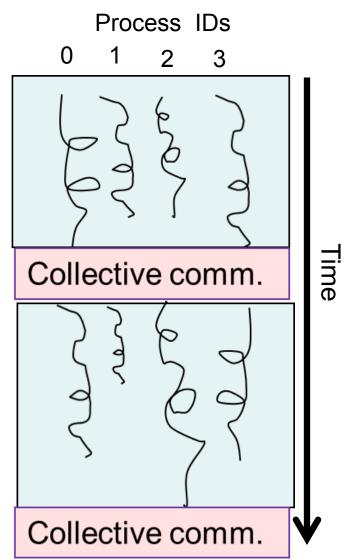
- And many SPMD programs use an additional simplification ... "Bulk Synchronous Processing".
  - Each process maintains a local view of the global data
  - A problem is broken down into phases each composed of two subphases:
    - Compute on local view of data (the "squiggles" in the figure)
    - Communicate to update global view on all processes (collective communication).
  - Continue phases until complete

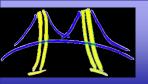




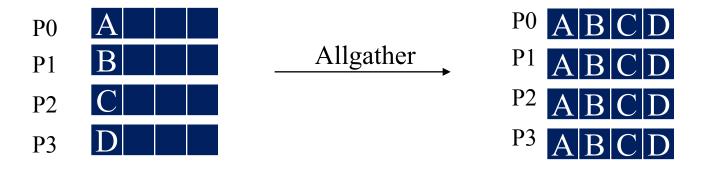
# Parallel overheads with the Bulk Synchronous Processing pattern

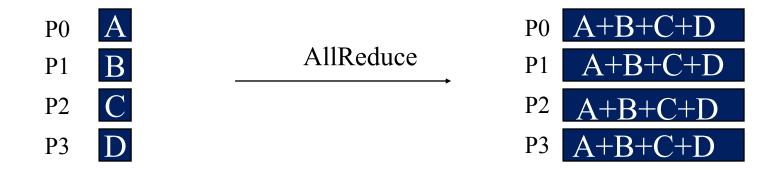
- Two major sources of parallel overhead:
- 1. Load imbalance: the slowest process determines when everyone is done. Time waiting for other processes to finish (i.e. unequal lengths of the "squiggles" in the figure ) is time wasted.
- 2. Communication overhead: A cost only incurred by the parallel program. Grows with the number of processes for collective comm.



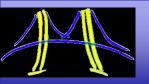


#### **Collective Data Movement**





74



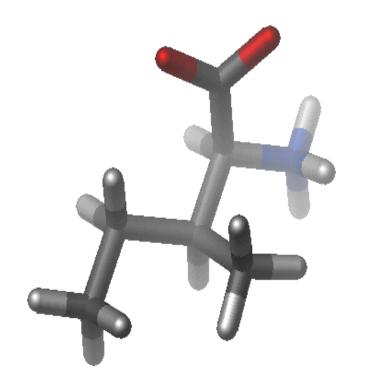
### Molecular dynamics



 Models motion of atoms in molecular systems by solving Newton's equations of motion:

$$\vec{M}\frac{d^2\vec{r}}{dt^2} = -\nabla U(\vec{r}) = F(\vec{r})$$

- The potential energy, U(r), is divided into two parts:
  - Bonded terms Groups of atoms connected by chemical bonds.
  - Non-bonded terms longer range forces (e.g. electrostatic).
    - An N-body problem ... i.e. every atom depends on every other atom, an O(N<sup>2</sup>) problem.



Bonds, angles and torsions

# Molecular dynamics simulation



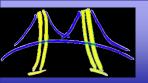
```
We used a cutoff method ... the
real atoms (3, N)
                             potential energy drops off quickly so
                             atoms beyond a neighborhood can be
real force (3, N)
                             ignored in the nonbonded force calc.
int neighbors (MX, N)
// Every PE has a copy of atoms and force
loop over time steps
    parallel loop over atoms
      Compute neighbor list (for my atoms)
      Compute nonbonded forces (my atoms and neighbors)
      Barrier
      All reduce (Sum force arrays, each PE gets a copy)
      Compute bonded forces (for my atoms)
       Integrate to Update position (for my atoms)
      All gather (update atoms array)
    end loop over atoms
```

end loop

# Molecular dynamics simulation



```
real atoms (3, N)
real force (3, N)
int neighbors (MX, N) //MX = max neighbors an atom may have
// Every PE has a copy of atoms and force
loop over time steps
    parallel loop over atoms
      Compute neighbor list (for my atoms)
      Compute long range for synchronization nd neighbors)
      Barrier
      All reduce (Sum force arrays,
                                                     copy)
      Compute bonded forces (for my
      Integrate to Update position Communication
      All gather (update atoms array)
    end loop over atoms
end loop
```



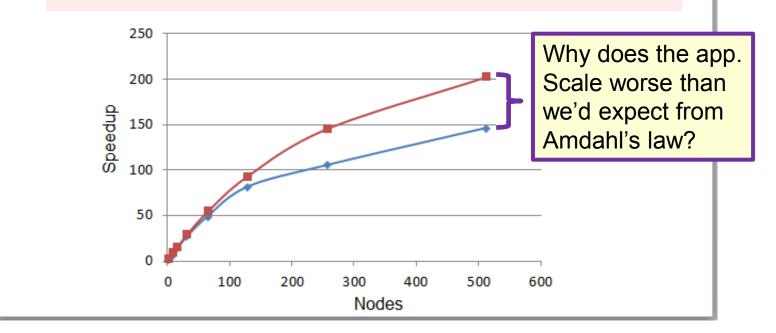
### Limitations to scalability

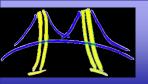
Remember the speedup plot we discussed from last time?

# **M**

# Amdahl's Law and the CHARMM MD program

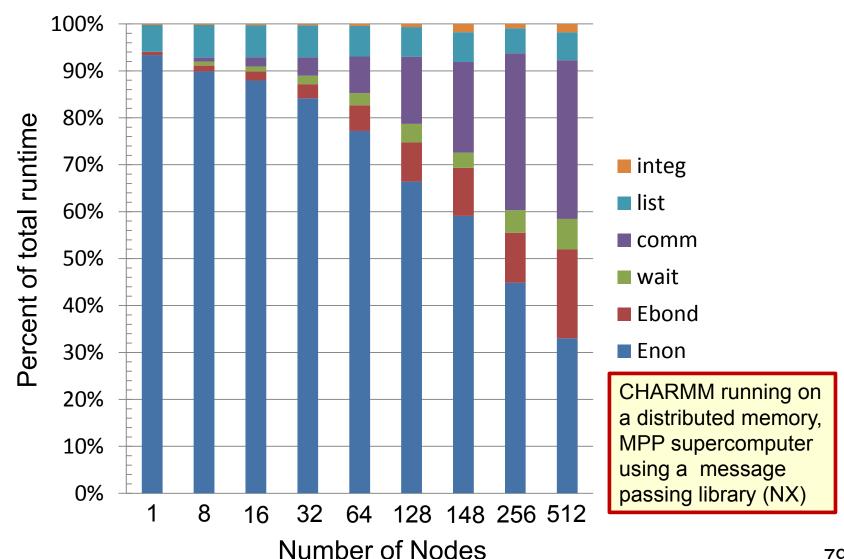
- We Profiled CHARMM running on the Paragon XPS to find the time spent in code that was not parallelized ... concluded that CHARMM has a serial fraction of ~0.003.
  - The maximum possible speedup is: S= 1/0.003 = 333



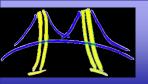


# **CHARMM Myoglobin Benchmark**

Percent of runtime for the different phases of the computation

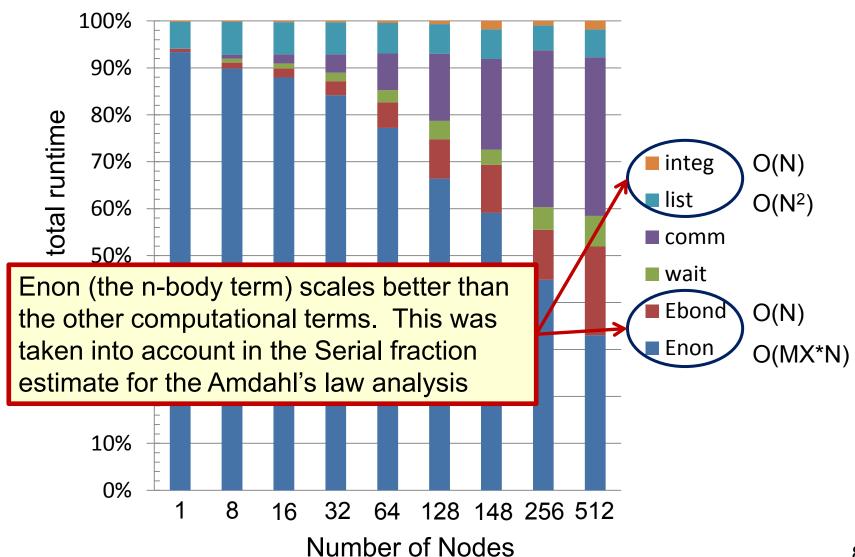


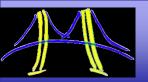
79



# **CHARMM Myoglobin Benchmark**

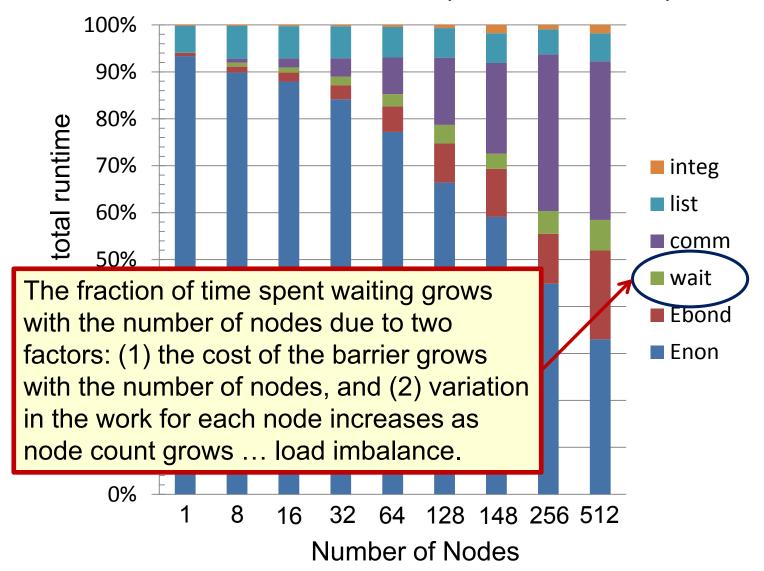
Percent of runtime for the different phases of the computation

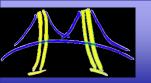




# Charm Myoglobin Benchmark

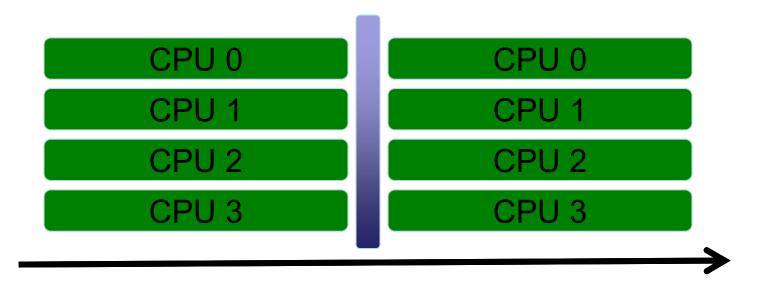
Percent of runtime for the different phases of the computation

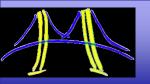




# Synchronization overhead

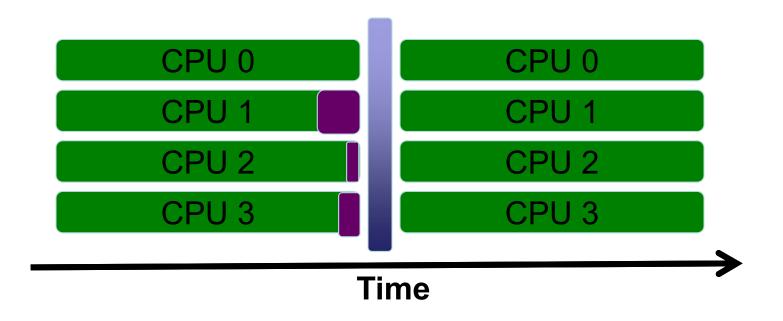
- Processes finish their work and must assure that all processes are finished before the results are combined into the global force array.
  - This is parallel overhead since this doesn't occur in a serial program.
  - The synchronization construct itself takes time and in some cases (such as a barrier) the cost grows with the number of nodes.

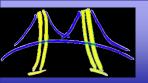




### Load imbalance

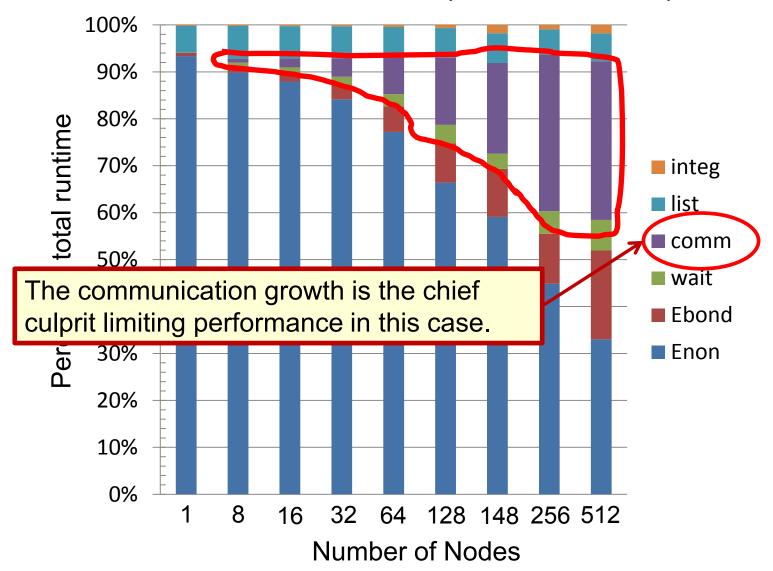
- If some processes finish their share of the computation early, the time spent waiting for other processors is wasted.
  - This is an example of Load Imbalance

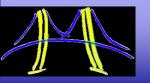




# Charm Myoglobin Benchmark

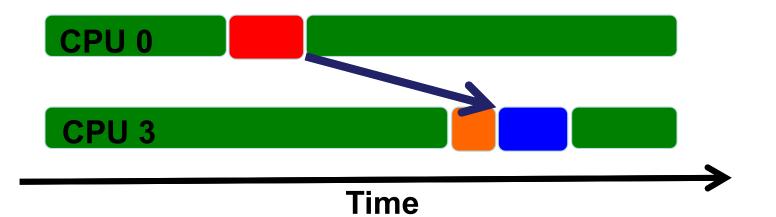
Percent of runtime for the different phases of the computation

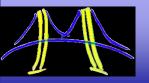




#### Communication

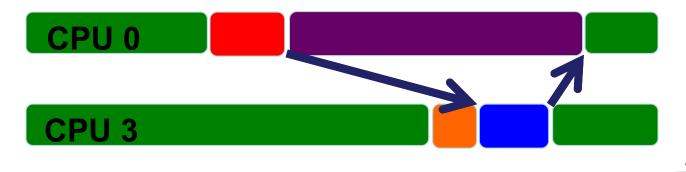
- On distributed-memory machines (e.g. a cluster), communication can only occur by sending discrete messages over a network
  - The sending processor marshals the shared data from the application's data structures into a message buffer
  - The receiving processor must wait for the message to arrive ...
  - ... and un-pack the data back into data structures

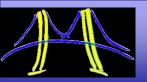




#### Communication

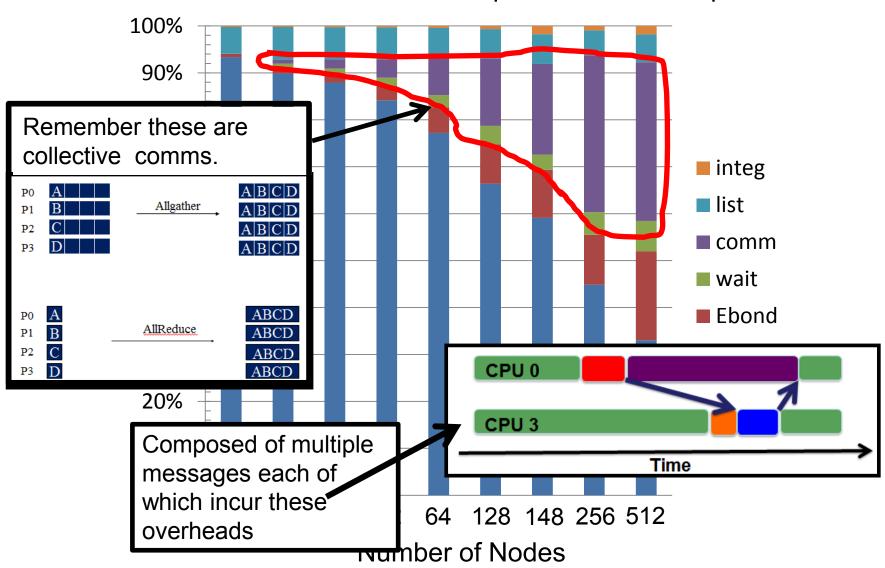
- On distributed-memory machines (e.g. a cluster), communication can only occur by sending discrete messages over a network
  - The sending processor marshals the shared data from the application's data structures into a message buffer
  - The receiving processor must wait for the message to arrive ...
  - ... and un-pack the data back into data structures
- If the communication protocol is synchronous, then the sending processor must wait for acknowledgement that the message was received

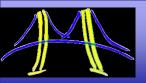




# Charm Myoglobin Benchmark

Percent of runtime for the different phases of the computation





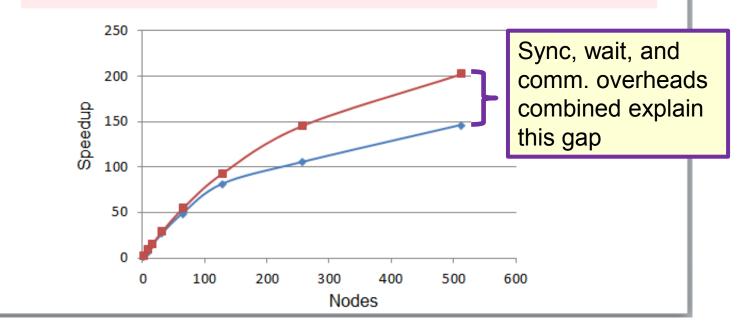
# Limitations to scalability

Remember the speedup plot we discussed from last time?

# **M**

# Amdahl's Law and the CHARMM MD program

- We Profiled CHARMM running on the Paragon XPS to find the time spent in code that was not parallelized ... concluded that CHARMM has a serial fraction of ~0.003.
  - The maximum possible speedup is: S = 1/0.003 = 333

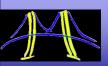


#### **Outline**

- Motivation: We all must be parallel programmers
- Key concepts in parallel Computing
- An introduction to parallel hardware
  - Software for parallel systems: key design patterns
  - Closing Comments

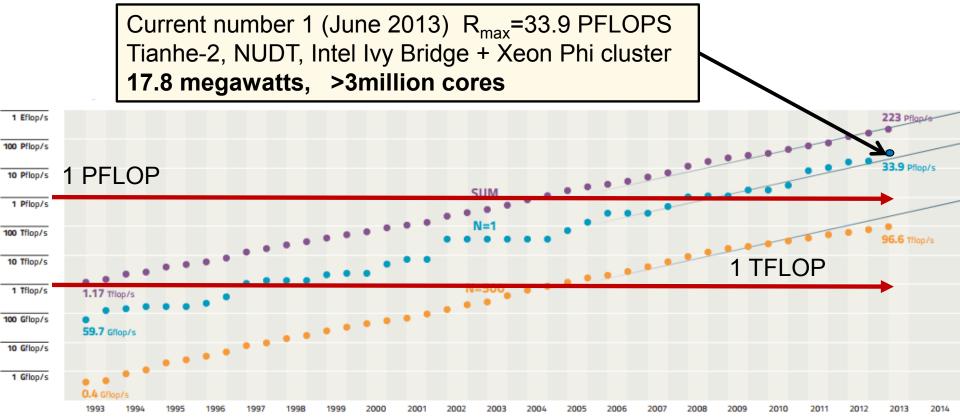
#### **Outline**

- Motivation: We all must be parallel programmers
- Key concepts in parallel Computing
- An introduction to parallel hardware
- History of parallel hardware
  - The major building blocks of modern parallel systems
    - Multicore processors
    - The GPU
- Software for parallel systems: key design patterns
- Closing Comments



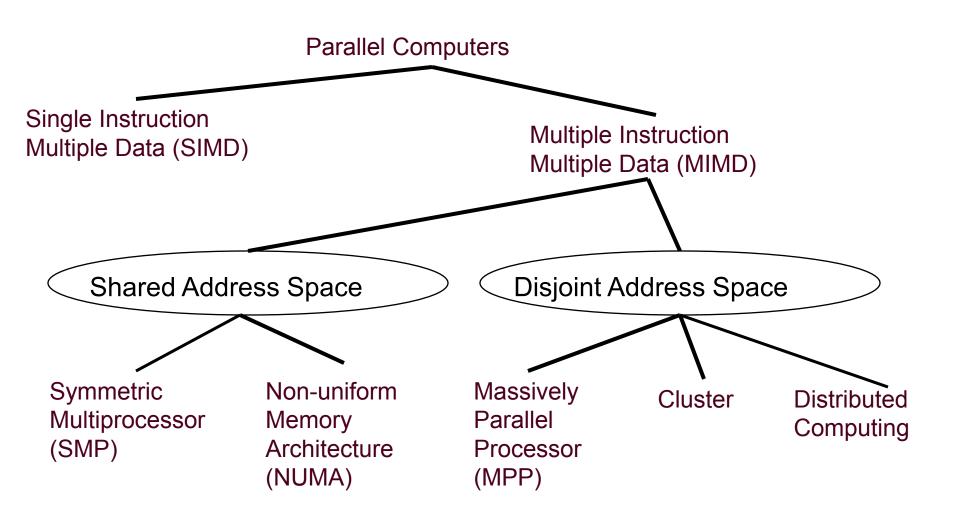
# Tracking Supercomputers: Top500

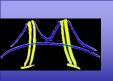
- Top500: a list of the 500 fastest computers in the world (www.top500.org)
- Computers ranked by solution to the MPLinpack benchmark:
  - Solve Ax=b problem for any order of A
- List released twice per year: in June and November



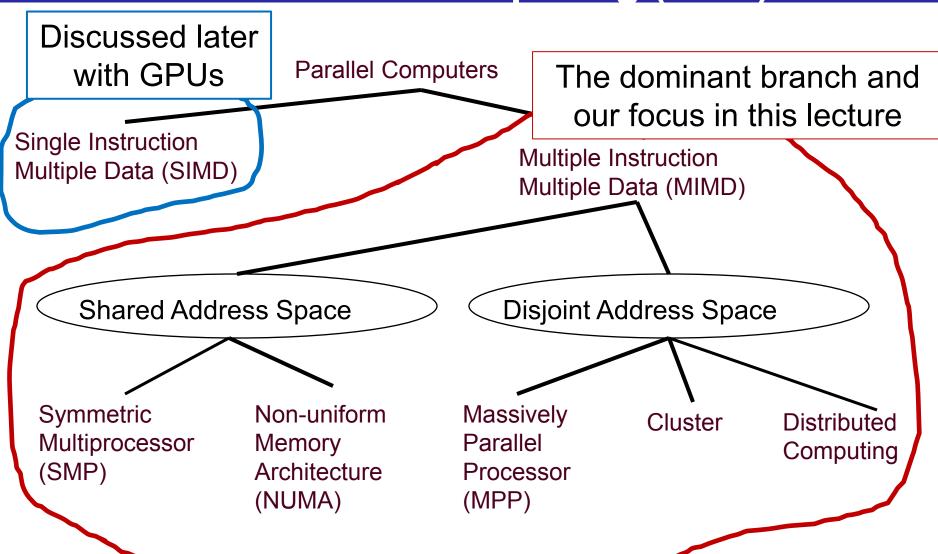


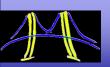
# Hardware Architectures for High Performance Computing (HPC)





# Hardware Architectures for High Performance Computing (HPC)





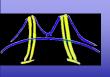
# The birth of Supercomputing



On July 11, 1977, the CRAY-1A, serial number 3, was delivered to NCAR. The system cost was \$8.86 million (\$7.9 million plus \$1 million for the disks).

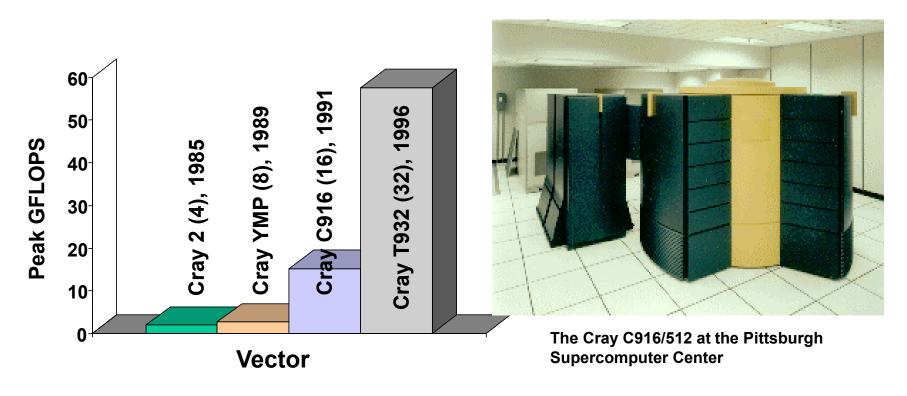
http://www.cisl.ucar.edu/computers/gallery/cray/cray1.jsp

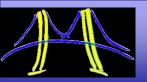
- The CRAY-1A:
  - 2.5-nanosecond clock,
  - 64 vector registers,
  - 1 million 64-bit words of highspeed memory.
  - Peak speed:
    - 80 MFLOPS scalar.
    - 250 MFLOPS vector (but this was VERY hard to achieve)
- Cray software ... by 1978
  - Cray Operating System (COS),
  - the first automatically vectorizing Fortran compiler (CFT),
  - Cray Assembler Language (CAL) were introduceg₄



# History of Supercomputing:

- Large mainframes that operated on vectors of data
- Custom built, highly specialized hardware and software
- Multiple processors in an shared memory configuration
- Required modest changes to software (vectorization)





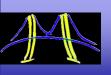
#### The attack of the killer micros



- The Caltech Cosmic Cube developed by Charles Seitz and Geoffrey Fox in1981
- 64 Intel 8086/8087 processors
- 128kB of memory per processor
- 6-dimensional hypercube network

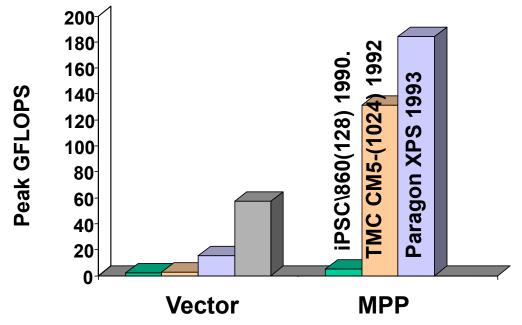
The cosmic cube, Charles Seitz Communications of the ACM, Vol 28, number 1 January 1985, p. 22 Launched the "attack of the killer micros" Eugene Brooks, SC'90

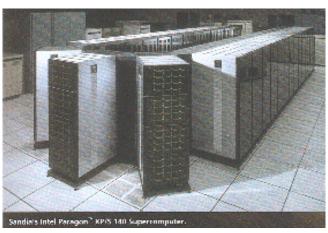
http://calteches.library.caltech.edu/3419/1/Cubism.pdf



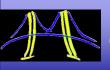
# It took a while, but MPPs came to dominate supercomputing

- Parallel computers with large numbers of microprocessors
- High speed, low latency, scalable interconnection networks
- Lots of custom hardware to support scalability
- Required massive changes to software (parallelization)



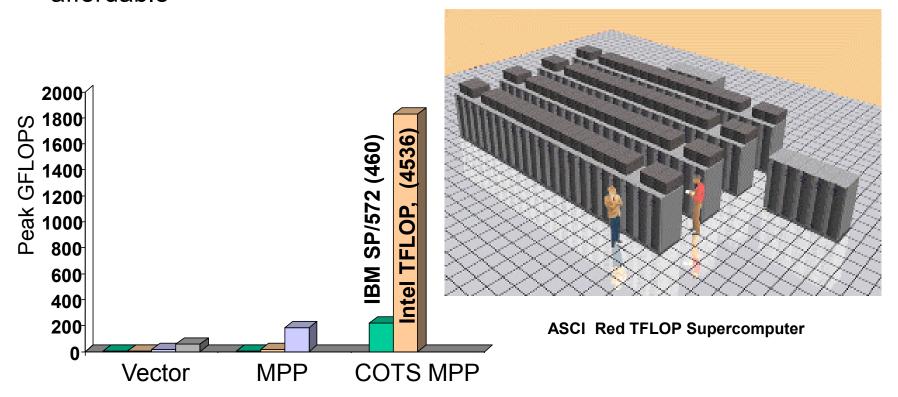


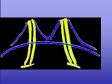
Paragon XPS-140 at Sandia National labs in Albuquerque NM



# The cost advantage of mass market COTS

- MPPs using <u>Mass market Commercial</u> off the shelf (COTS) microprocessors and standard memory and I/O components
- Decreased hardware and software costs makes huge systems affordable

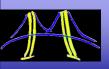




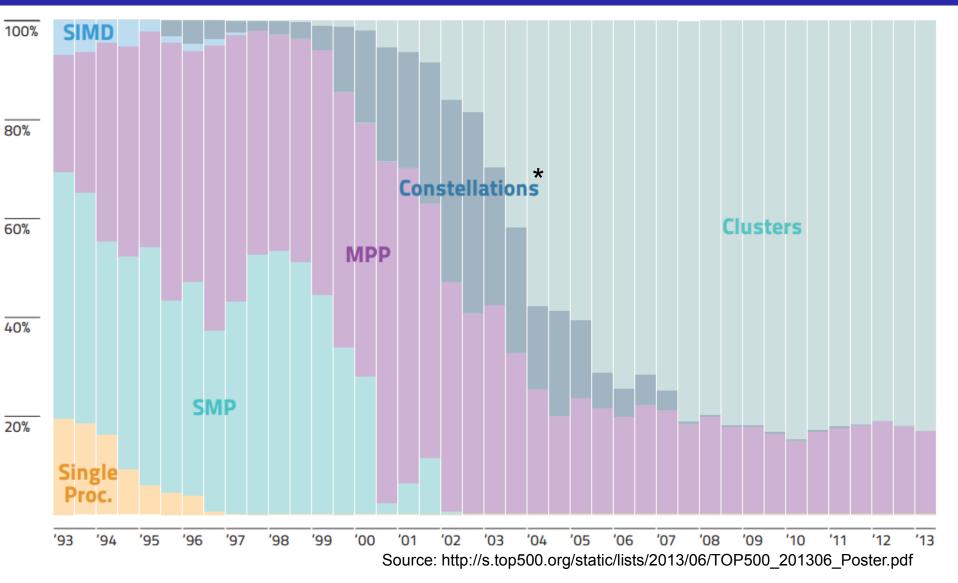
# The MPP future looked bright ... but then clusters took over

- A cluster is a collection of connected, independent computers that work in unison to solve a problem.
- Nothing is custom ... motivated users could build cluster on their own
- First clusters appeared in the late 80's (Stacks of "SPARC pizza boxes")
- The Intel Pentium Pro in 1995 coupled with Linux made them competitive.
  - NASA Goddard's Beowulf cluster demonstrated publically that high visibility science could be done on clusters.
- Clusters made it easier to bring the benefits due to Moores's law into working supercomputers





# **Top 500 list: System Architecture**



\*Constellation: A cluster for which the number of processors on a node is greater than the number of nodes in the cluster. I've never seen anyone use this term outside of the top500 list.

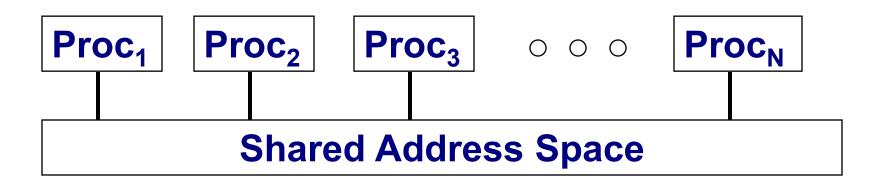
#### **Outline**

- Motivation: We all must be parallel programmers
- Key concepts in parallel Computing
- An introduction to parallel hardware
  - History of parallel hardware
  - The major building blocks of modern parallel systems
  - Multicore processors
    - The GPU
- Software for parallel systems: key design patterns



# How do we connect cores together?

- A symmetric multiprocessor (SMP) consists of a collection of processors that share a single address space:
  - Multiple processing elements.
  - A shared address space with "equal-time" access for each processor.
  - The OS treats every processor the same

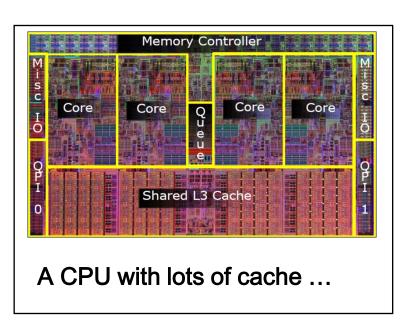






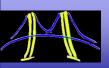
#### How realistic is this model?

 Some of the old supercomputer mainframes followed this model,



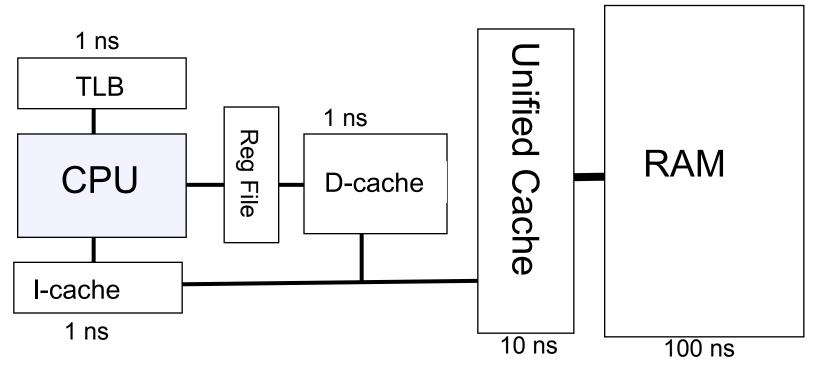


- But as soon as we added caches to CPUs, the SMP model fell apart.
  - Caches ... all memory is equal, but some memory is more equal than others.

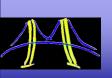


# **Memory Hierarchies**

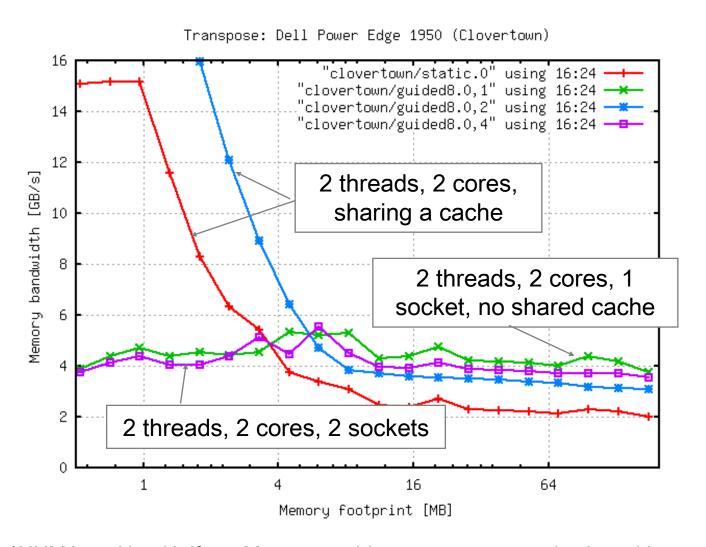
A typical microprocessor memory hierarchy



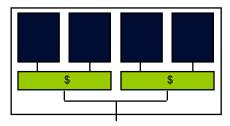
- Instruction cache and data cache pull data from a unified cache that maps onto RAM.
- TLB implements virtual memory and brings in pages to support large memory foot prints.



# NUMA\* issues on a Multicore Machine 2-socket Clovertown Dell PE1950



A single quadcore chip is a NUMA machine!

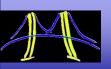


Xeon® 5300 Processor block diagram

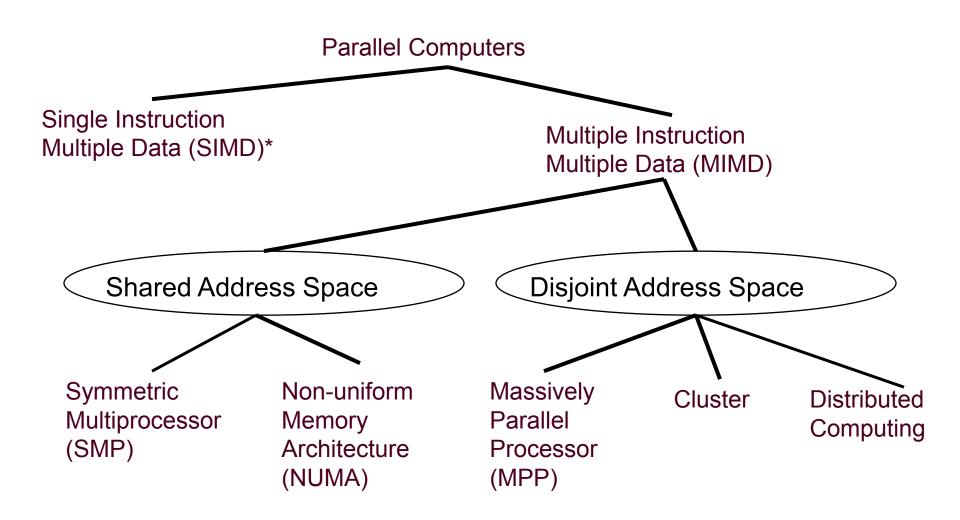
\*NUMA == Non Uniform Memory architecture ... memory is shared but access times vary.

#### **Outline**

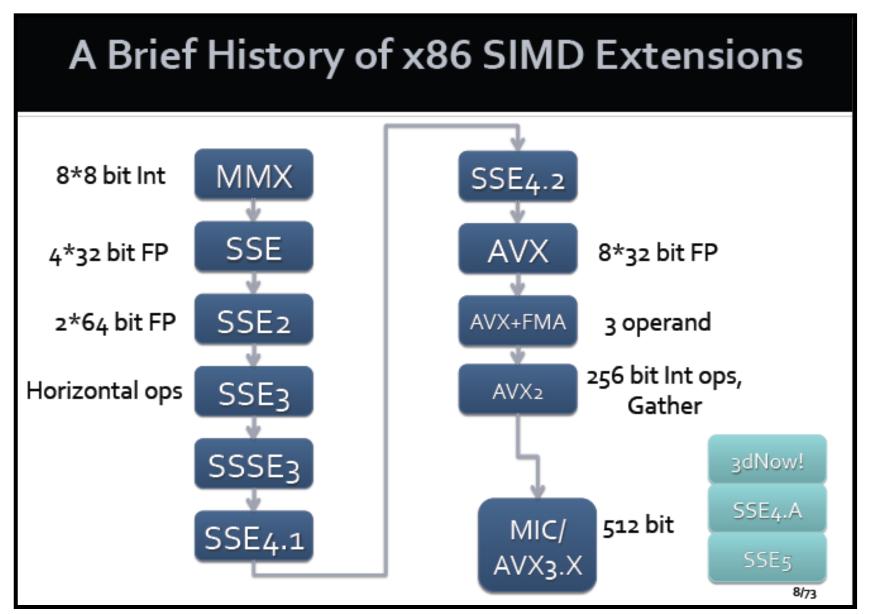
- Motivation: We all must be parallel programmers
- Key concepts in parallel Computing
- An introduction to parallel hardware
  - History of parallel hardware
  - The major building blocks of modern parallel systems
    - Multicore processors
  - The GPU
- Software for parallel systems: key design patterns
- Closing comments



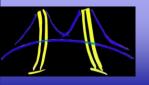
### What happened to SIMD?



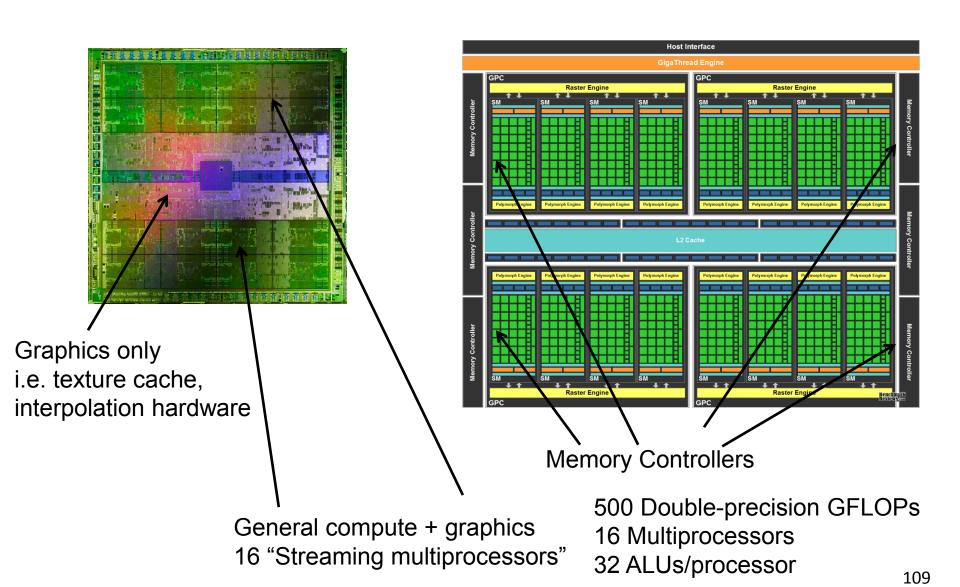
#### SIMD and sx86 multimedia extensions.



Source: Bryan Catanzaro, NVIDIA, UCB Parlab Bootcamp, 2013

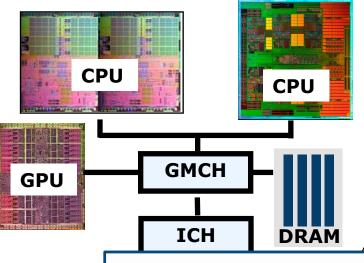


## **NVIDIA GTX 480**



## The end of the discrete GPU

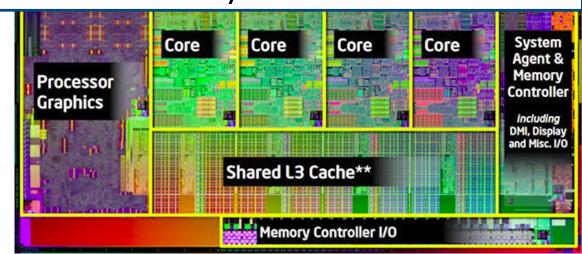




- A modern platform has:
  - CPU(s)
  - GPU(s)
  - DSP processors
  - ... other?

Absorption into CPU (remove "off chip" penalty) but uncertain standards story → success unclear

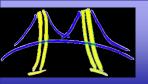
• Current this functionality onto a single chip ... mitigates the PCIe bottleneck in GPGPU computing!



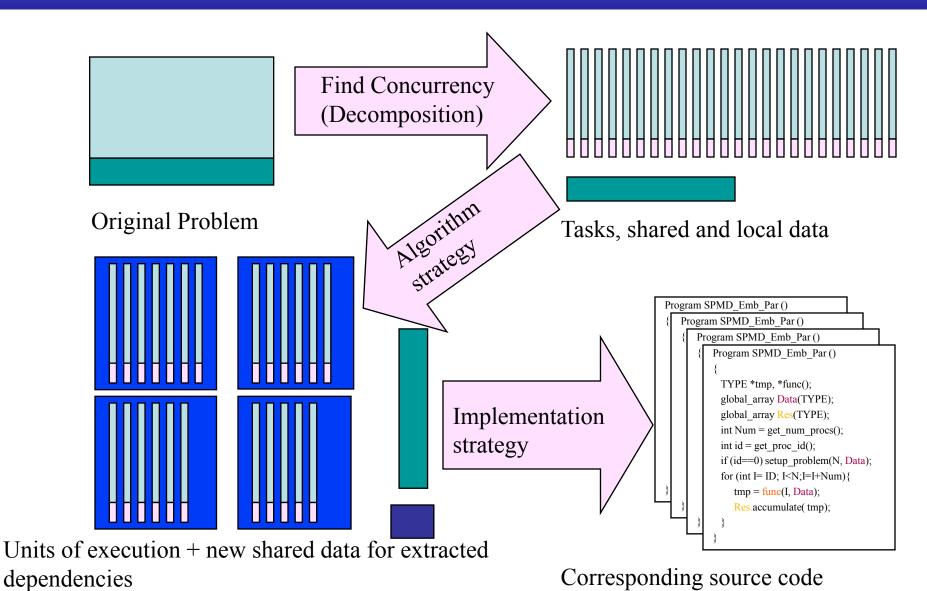
Intel® Core™ i5-2500K Desktop Processor (Sandy Bridge) Intel HD Graphics 3000 (2011)

### **Outline**

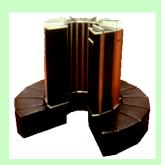
- Motivation: We all must be parallel programmers
- Key concepts in parallel Computing
- An introduction to parallel hardware
- Software for parallel systems: key design patterns
  - Closing comments



## The Parallel programming process:



## Parallel computing: It's old

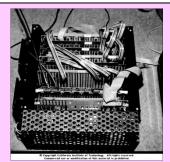






Cray C-90 (1991)

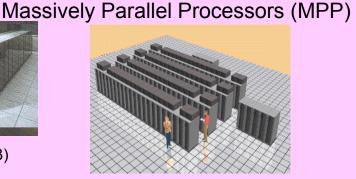
Vector Computers



Cosmic cube (1983)



Paragon (1993)



**ASCI Red (1997)** 

**Cluster Computers** 



Clusters (late 80's)

Linux PC Clusters (~1995)

Late 80's Late 90's

## We tried to solve the parallel programming problem by searching for the right programming environment

### Parallel programming environments in the 90's

ABCPL	CORRELATE	GLU	Mentat	Parafrase2	pC++
ACE	CPS	GUARD	Legion	Paralation	SCHEDULE
ACT++	CRL	HAsL.	Meta Chaos	Parallel-C++	SciTL
Active messages	CSP	Haskell	Midway	Parallaxis	POET
Adl	Cthreads	HPC++	Millipede	ParC	SDDA.
Adsmith	CUMULVS	JAVAR.	CparPar	ParLib++	SHMEM
ADDAP	DAGGER	HORUS	Mirage	ParLin	SIMPLE
AFAPI	DAPPLE	HPC	MpC	Parmacs	Sina
ALWAN	Data Parallel C	IMPACT	MOSIX	Parti	SISAL.
AM	DC++	ISIS.	Modula-P	pC	distributed smalltalk
AMDC	DCE++	JAVAR	Modula-2*	pC++	SMI.
AppLeS	DDD	JADE	Multipol	PCN	
Amoeba	DICE.	Java RMI	MPI	PCP:	SONIC Salit C
ARTS	DIPC	javaPG	MPC++	PH	Split-C.
Athapascan-0b	DOLIB	JavaSpace	Munin	PEACE	SR
Aurora	DOME	JIDL	Nano-Threads	PCU	Sthreads
Automap	DOSMOS.	Joyce	NESL	PET	Strand.
bb_threads	DRL	Khoros	NetClasses++	PETSc	SUIF.
Blaze	DSM-Threads	Karma	Nexus	PENNY	Synergy
BSP	Ease .	KOAN/Fortran-S	Nimrod	Phosphorus	Telegrphos
BlockComm	ECO	LAM	NOW	POET.	SuperPascal
C*.	Eiffel	Lilac	Objective Linda	Polaris	TCGMSG.
"C* in C	Eilean	Linda	Occam	POOMA	Threads.h++.
C**	Emerald	JADA	Omega	POOL-T	TreadMarks
CarlOS	EPL	WWWinda	OpenMP	PRESTO	TRAPPER
Cashmere	Excalibur	ISETL-Linda	Orca	P-RIO	uC++
C4	Express	ParLin	OOF90	Prospero	UNITY
CC++	Falcon	Eilean	P++	Proteus	UC
Chu	Filaments	P4-Linda	P3L	QPC++	V
Charlotte	FM	Glenda	p4-Linda	PVM	ViC*
Charm	FLASH	POSYBL	Pablo	PSI	Visifold V-NUS
Charm++	The FORCE	Objective-Linda	PADE	PSDM	VPE
Cid	Fork	LiPS	PADRE	Quake	Win32 threads
Cilk	Fortran-M	Locust	Panda	Quark	WinPar
CM-Fortran	FX	Lparx	Papers	Quick Threads	WWWinda
Converse	GA	Lucid	AFAPI.	Sage++	XENOOPS
Code	GAMMA	Maisie	Para++	SCANDAL	XPC
COOL	Glenda	Manifold	Paradigm	SAM	Zounds
COOL	Civilan		- uruurgiii	071111	ZPL

## Throwing new languages at the problem didn't work: the "Dead Architecture Society"



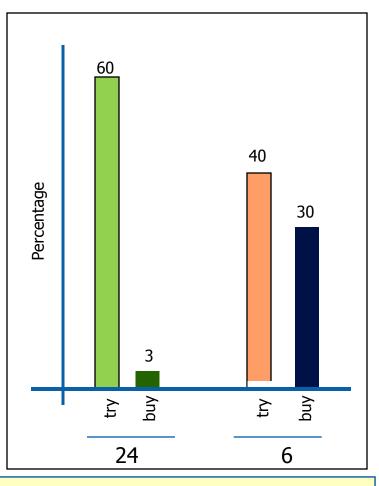
1980 1990 2000

## Language obsessions: More isn't always

better

 The Draeger Grocery Store experiment consumer choice :

- Two Jam-displays with coupon's for purchase discount.
  - 24 different Jam's
  - 6 different Jam's
- How many stopped by to try samples at the display?
- Of those who "tried", how many bought jam?



The findings from this study show that an extensive array of options can at first seem highly appealing to consumers, yet can reduce their subsequent motivation to purchase the product.

Iyengar, Sheena S., & Lepper, Mark (2000). When choice is demotivating: Can one desire too much of a good thing? *Journal of Personality and Social Psychology*, 76, 995-1006.

## My optimistic view from 2005 ...

## Parallel Programming API's today

- Thread Libraries
  - Win32 API
  - POSIX threads.
- Compiler Directives
  - OpenMP portable shared memory parallelism.
- Message Passing Libraries
  - MPI message passing
- Coming soon ... a parallel language for managed runtimes? Java or X10?

We don't want to scare away the programmers ... Only add a new API/language if we can't get the job done by fixing an existing approach.

Third party names are the property of their owners.

We've learned our lesson ... we emphasize a small number of industry standards

## But we didn't learn our lesson History is repeating itself!



A small sampling of Programming environments from the **NEW golden age of parallel programming** (from the literature 2010-2012)

AM++**ISPC** OpenACC Scala Copperhead ArBB SIAL CUDA Java PAMI Parallel Haskell BSP STAPL DryadOpt Liszt C++11 Erlang MapReduce ParalleX STM C++AMP **Fortress** MATE-CG **PATUS SWARM** Charm++ GA MCAPI PLINQ TBB Chapel GO MPI PPL **UPC** Cilk++ NESL Pthreads Win32 Gossamer CnC threads **GPars** OoOJava **PXIF** coArray Fortran GRAMPS OpenMP PyPar X10 Codelets Hadoop OpenCL Plan42 **XMT HMPP** OpenSHMEM **RCCE** ZPL

Note: I'm not criticizing these technologies. I'm criticizing our collective urge to create so many of them.

## Maybe its time to try something different?

## But we didn't learn our lesson History is repeating itself!



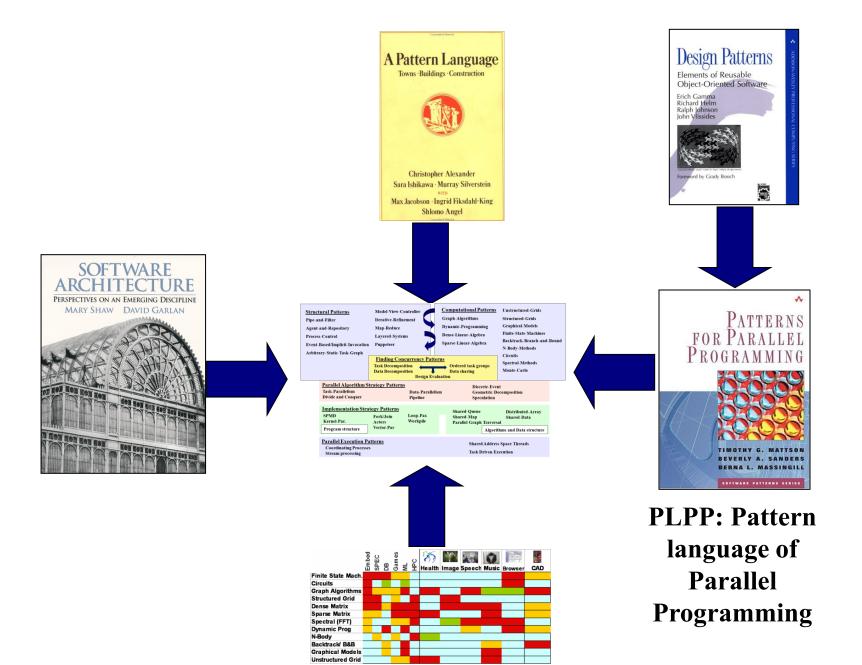
A small sampling of Programming environments from the NEW golden age of parallel programming (from the literature 2010-2012)

AM++	Copperhea	ISPC	OpenACC	Scala
ArBB	d	Java	PAMI	SIAL
BSP	CUDA	Liszt	Parallel Haskell	STAPL
C++11	DryadOpt	MapReduce	ParalleX	STM
C++AMP	Erlang	MATE-CG	PATUS	SWARM
Charm++	Fortress	MCAPI	PLINQ	TBB
Chapel	GA	MPI	PPL	UPC
Cilk++	GO	NESL	Pthreads	Win32
CnC	Gossamer	OoOJava	PXIF	threads
coArray Fortran	GPars	OpenMP	PyPar	X10
Codelets	GRAMPS	OpenCL	Plan42	XMT
	Hadoop	OpenSHME	RCCE	ZPL
	HMMP	M		

Note: I'm not criticizing these technologies. I'm criticizing our collective urge to create so many of them.

2

Third party names are the property of their owners.



13 dwarves

## OPL Pattern Language (Keutzer & Mattson 2010)

### **Applications**

**Structural Patterns** 

**Agent-and-Repository** 

**Model-View-Controller** 

**Iterative-Refinement** 

**Graph-Algorithms** 

Pipe-and-Filter

**Layered-Systems** 

**Map-Reduce** 

**Puppeteer** 

**Dynamic-Programming** 

**Computational Patterns** 

**Process-Control** 

Dense-Linear-Algebra

**Event-Based/Implicit-Invocation** 

Sparse-Linear-Algebra

**Unstructured-Grids** 

Structured-Grids

**Graphical-Models** 

**Finite-State-Machines** 

Backtrack-Branch-and-Bound

**N-Body-Methods** 

Circuits

**Spectral-Methods** 

Monte-Carlo

Arbitrary-Static-Task-Graph

**Finding Concurrency Patterns** 

**Task Decomposition Data Decomposition** 

Ordered task groups **Data sharing Design Evaluation** 

**Parallel Algorithm Strategy Patterns** 

Task-Parallelism **Divide and Conquer**  **Data-Parallelism Pipeline** 

Discrete-Event

**Geometric-Decomposition** 

**Speculation** 

**Implementation Strategy Patterns** 

**SPMD** 

Kernel-Par.

Fork/Join Actors

Loop-Par. Workpile

**Shared-Queue** Shared-Map

**Distributed-Array** Shared-Data

**Parallel Graph Traversal** 

**Program structure** 

Vector-Par

**Algorithms and Data structure** 

**Parallel Execution Patterns** 

**Coordinating Processes** Stream processing

**Shared Address Space Threads** 

**Task Driven Execution** 

**Concurrency Foundation constructs (not expressed as patterns)** 

Thread/proc management

Communication

**Synchronization** 

Source: Keutzer and Mattson Intel Technology Journal, 2010

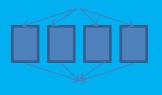
## Pattern examples





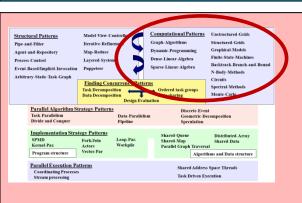


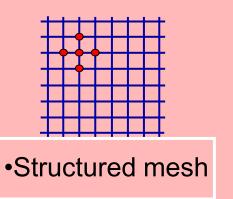


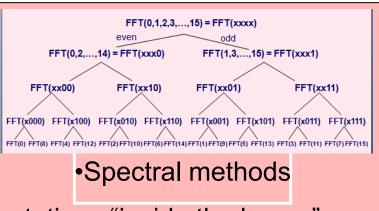


Pipe-and-FilterIterative refinementMapReduce

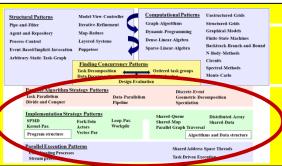
Structural Patterns: Define the software structure .. Not what is computed







Computational Patterns: Define the computations "inside the boxes"

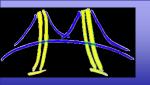


•Fork-join

•SPMD

Data parallel

Parallel Patterns: Defines parallel algorithms

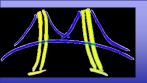


## Seven strategies for parallelizing software

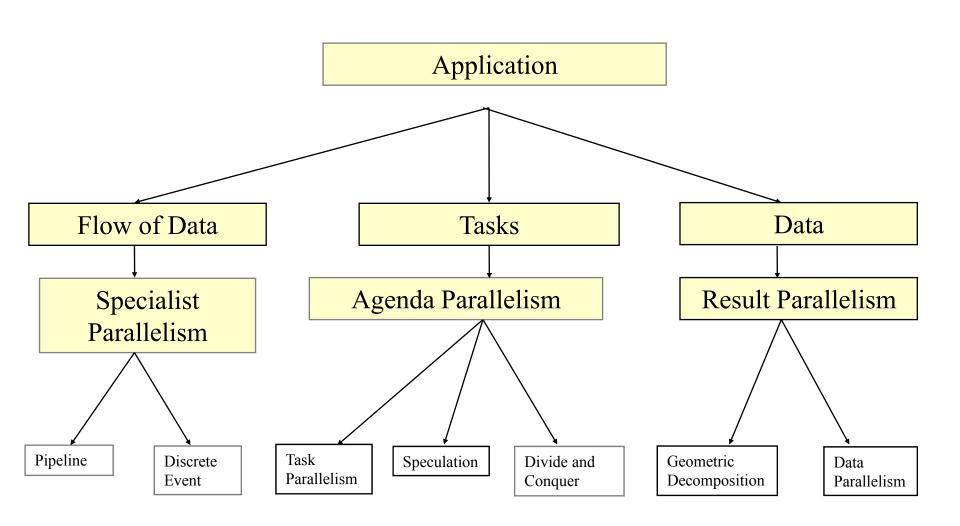
- These seven strategies for parallelizing software give us:
  - Names: so we can communicate better
  - Categories: so we can gather and share information
  - A palette (like an artist's palette) of approaches that is:
    - Necessary: we should consider them all and
    - Sufficient: once we have considered them all then we don't' have to worry that we forgot something

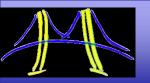
#### Parallel Algorithm Strategy Patterns

Task-Parallelism Divide and Conquer Data-Parallelism Pipeline Discrete-Event Geometric-Decomposition Speculation



## **Parallel Algorithmic Strategies**





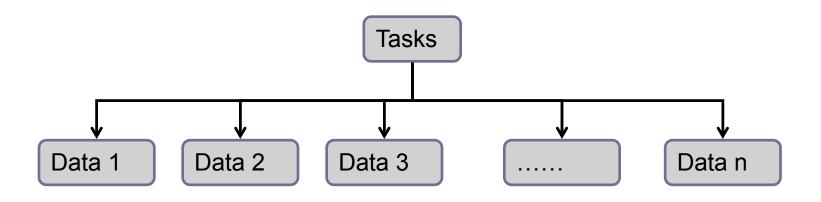
## Data Parallelism Pattern

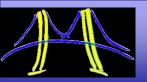
#### Use when:

 Your problem is defined in terms of collections of data elements operated on by a similar (if not identical) sequence of instructions; i.e. the concurrency is in the data.

#### Solution

- Define collections of data elements that can be updated in parallel.
- Define computation as a sequence of collective operations applied together to each data element.





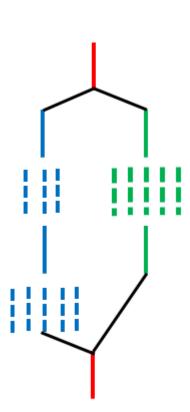
## **Task Parallelism Pattern**

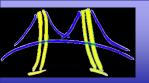
#### Use when:

 The problem naturally decomposes into a distinct collection of tasks

#### Solution

- Define the set of tasks and a way to detect when the computation is done.
- Manage (or "remove") dependencies so the correct answer is produced regardless of the details of how the tasks execute.
- Schedule the tasks for execution in a way that keeps the work balanced between the processing elements of the parallel computer and

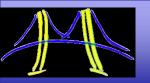




## Task Parallelism in practice

- Embarrassingly parallel:
  - The tasks are independent, so the parallelism is "so easy to exploit it's embarrassing".
- Separable dependencies:
  - Turn a problem with dependent tasks into an "embarrassingly parallel" by "replicating data between tasks, doing the work, then recombining data (often a reduction) to restore global state.
- Functional Decomposition
  - A task is associated with a functional decomposition of the problem to produce a coarse grained parallel program

Its becoming common to associate this case as the prototypical "task parallel" approach ... but to us old-timers, the previous two cases are overwhelming more common.



## **Divide and Conquer Pattern**

#### Use when:

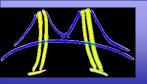
A problem includes a method to divide into subproblems and a way to recombine solutions of subproblems into a global solution.

#### Solution

- Define a split operation
- Continue to split the problem until subproblems are small enough to solve directly.
- Recombine solutions to subproblems to solve original global problem.

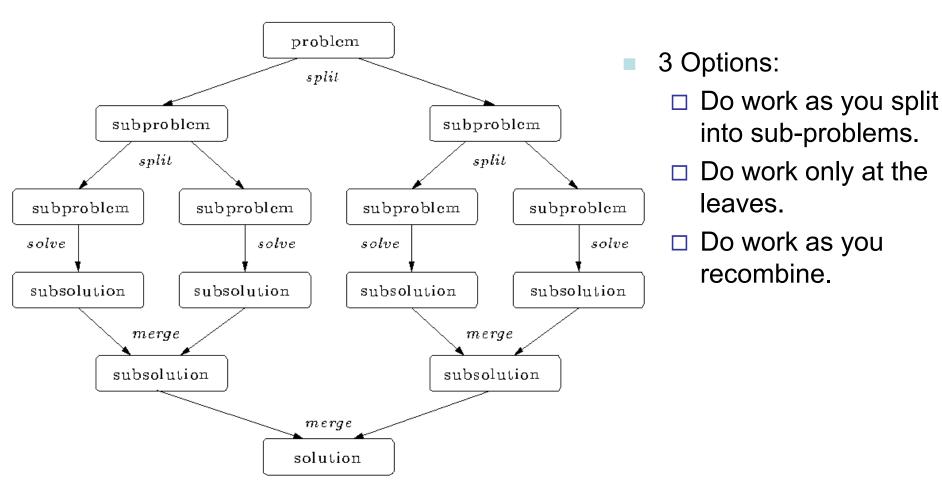
#### Note:

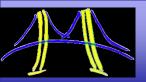
Computing may occur at each phase (split, leaves, recombine).



## Divide and conquer

Split the problem into smaller sub-problems. Continue until the subproblems can be solve directly.





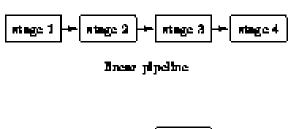
## **Pipeline Pattern**

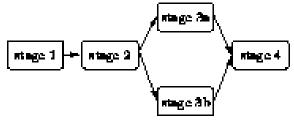
#### Use when:

 Your problem can be described as data flowing through a sequence of computational stages

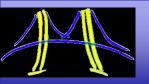
#### Solution

- □ Define a set of stages setup with data-flow connections between them.
- Set up input/output channels to support data driven execution.
- Parallelism comes from multiple stages acrive at one time.





nonlinear pipeline



## **Geometric Decomposition**

#### Use when:

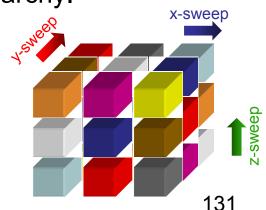
 The problem is organized around a central data structure that can be decomposed into smaller segments (chunks) that can be updated concurrently.

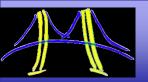
#### Solution

- Typically, the data structure is updated iteratively where a new value for one chunk depends on neighboring chunks.
- The computation breaks down into three components: (1) exchange boundary data, (2) update the interiors or each chunk, and (3) update boundary regions. The optimal size of the chunks is dictated by the properties of the memory hierarchy.

#### Note:

 This pattern is often used with the Structured grid and linear algebra computational strategy pattern.





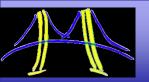
## **Speculation**

#### Use when:

 Suppose that the computation has been decomposed into a number of tasks that are not completely independent, but where conflicts are expected to only infrequently occur when the computation is actually executed. Solution

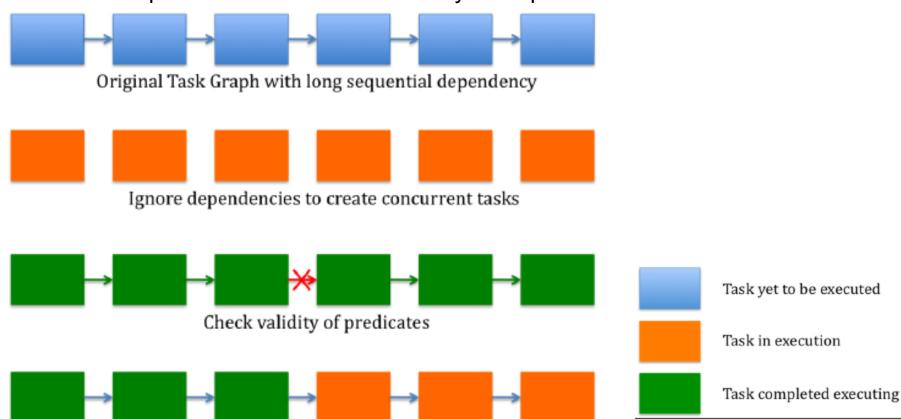
#### Solution:

- An effective solution may be to just run the tasks independently, that is speculate that no conflicts will occur, and then clean up after the fact and retry in the rare situations where a conflict does occur. Two essential element of this solution are:
  - Have an easily identifiable safety check to determine whether the computation ran without conflicts and can thus be committed
  - 2. The ability to rollback and re-compute the cases where conflicts occur.



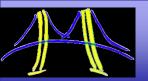
## **Speculative Parallelism**

- Speculative Parallelism:
  - Speculate on state of dependencies
  - Check validities of speculations
  - Recompute as needed to correct any mis-speculations



Recompute tasks (and its children) for which predicates are invalid

133



### **Discrete-Event**

#### Use when:

 The computation has been structured as loosely connected sequence of tasks that interact at unpredictable points in time.

#### Solution

- Setup an event handler infrastructure
- Launch a collection of tasks whose interaction is handled through the event handler. The handler is an intermediary between tasks, and in many cases the tasks do not need to know the source or destination for the events.

#### Note:

Discrete event is often used with problems, such as GUIs and discrete event simulations, that are handled with the Eventbased implicit invocation, model-view-controller, or process control patterns.

## OPL Pattern Language (Keutzer & Mattson 2010)

#### **Applications**

**Structural Patterns** 

Pipe-and-Filter

**Map-Reduce** 

**Puppeteer** 

Fork/Join

**Vector-Par** 

Actors

**Iterative-Refinement** 

**Model-View-Controller** 

**Agent-and-Repository** 

**Process-Control Layered-Systems** 

**Event-Based/Implicit-Invocation** 

Arbitrary-Static-Task-Graph

**Computational Patterns** 

**Graph-Algorithms** 

**Dynamic-Programming** 

Dense-Linear-Algebra

Sparse-Linear-Algebra

**Unstructured-Grids** 

Structured-Grids

**Graphical-Models** 

**Finite-State-Machines** 

Backtrack-Branch-and-Bound

**N-Body-Methods** 

Circuits

ecomposition

**Spectral-Methods** 

Monte-Carlo

**Finding Concurrency Patterns** 

**Task Decomposition Data Decomposition** 

**Data sharing Design Evaluation** 

**Parallel Algorithm Strategy Patterns** 

Implementation Strategy Patterns

Task-Parallelism

**SPMD** 

Kernel-Par.

**Divide and Conquer** 

Program structure

7 patterns to turn algorithms into code

Loop-Par.

Workpile

SHALEU-QUEUE

Ordered task groups

Shared-Map **Parallel Graph Traversal**  **Distributed-Array** Shared-Data

**Algorithms and Data structure** 

**Parallel Execution Patterns** 

**Coordinating Processes** Stream processing

**Shared Address Space Threads** 

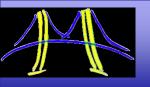
**Task Driven Execution** 

**Concurrency Foundation constructs (not expressed as patterns)** 

Thread/proc management

Communication Source: Keutzer and Mattson Intel Technology Journal, 2010

**Synchronization** 



## Seven strategies for implementing our algorithms as software

- These seven strategies for implementing our parallel algorithms give us:
  - Names: so we can communicate better
  - Categories: so we can gather and share information
  - A palette (like an artist's palette) of approaches that is:
    - Necessary: we should consider them all and
    - Sufficient: once we have considered them all then we don't' have to worry that we forgot something

### **Implementation Strategy Patterns**

**SPMD** 

Actors

Fork/Join

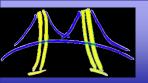
Workpile

**Program structure** 

Loop-Parallel

**Kernel-Parallel** 

**Vector-Parallel** 



## Implementation Strategy patterns

The most commonly used implementation strategy patterns:

SPMD	One program replicated, specialized by ID and NumProcs	
Fork-Join	Single thread forks a team as needed and later joins	
Work-pile	Create a pile of tasks for a set of workers to process	
Loop-Parallel	Make expensive loops independent and use a "parallel for"	
Vector-Parallel	Unroll loops to expose blocks, vector ops process blocks	
Kernel-Parallel	Fine-Grained SPMD kernels . Large numbers to address little's law.	

- Programming models are often optimized around the needs of these patterns. For "our" programming models:
  - MPI: SPMD, work-pile
  - OpenMP: Loop-parallel, fork-join ... SPMD on large NUMA systems.
  - OpenCL and CUDA: Kernel-parallelism
  - OpenACC: Loop-parallel and Kernel Parallel

### **Outline**

- Motivation: We all must be parallel programmers
- Key concepts in parallel Computing
- An introduction to parallel hardware
- Software for parallel systems: key design patterns
- Closing comments

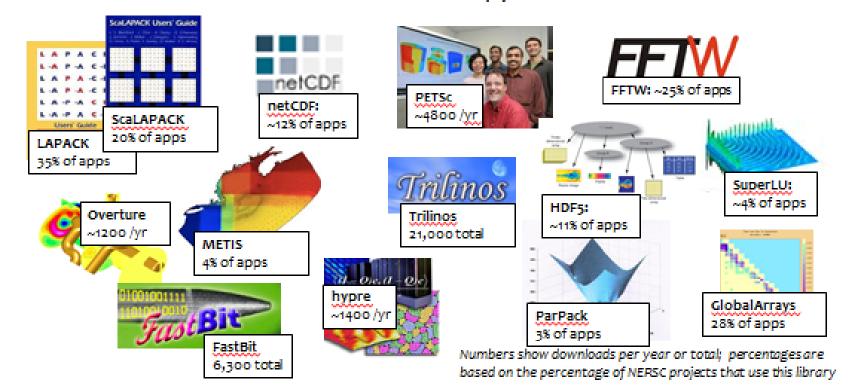
## Parallel programming is really hard

- Programming is hard whether you write serial or parallel code.
  - Parallel programming is just a new wrinkle added to the already tough problem of writing high quality, robust and efficient code.
- Why does Parallel programming seems so complex?
  - The literature overwhelms with hundreds of languages/APIs and a countless assortment of algorithms.
  - Experienced parallel programmers love to tell "war stories" of Herculean efforts to make applications scale ... which can scare people away.
  - It's new: synchronization, scalable algorithms, distributed data structures, concurrency bugs, memory models ... hard or not it's a bunch of new stuff to learn.

## But it's really not that bad (part 1): parallel libraries

## Programming Challenges and NITRD Solutions

- Application complexity grew due to parallelism and more ambitious science problems (e.g., multiphysics, multiscale)
- Scientific libraries enable these applications



2

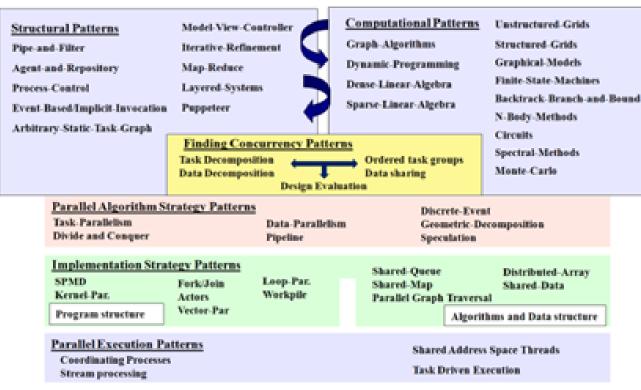
Source: Kathy Yelick

## But its really not that bad: part 2

- Don't let the glut of parallel programming languages confuse you.
- Leave research languages to C.S. researchers and stick to the small number of broadly used languages/APIs:
  - Industry standards:
    - Pthreads (eventually, C++'11 threads)
    - OpenMP
    - MPI
    - OpenCL
    - TBB (For C++ ... might be replaced by parallelism in C++ standard?)
  - or a broadly deployed solutions tied to your platform of choice
    - CUDA and OpenACC (for NVIDIA platforms and PGI compilers)
    - .NET and C++ AMP (Microsoft)

## But its really not that bad: part 3

 Most algorithms are based on a modest number of recurring patterns.

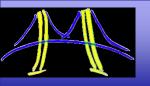


- Almost every parallel program is written in terms of just 7 basic patterns:
  - SPMD
  - Kernel Parallelism
  - Fork/join
  - Actors

- Vector Parallelism
- Loop Parallelism
- Work Pile

## Parallel programming is easy

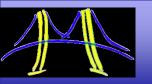
- So all you need to do is:
  - Pick your language.
    - I suggest sticking to industry standards and open source so you can move around between hardware platforms:
  - pthreads– OpenMP– OpenCL– MPI– TBB
  - Learn the key 7 patterns
    - SPMDVector Parallelism
    - Kernel ParallelismLoop Parallelism
    - Fork/joinWork Pile
    - Actors
    - Master the few patterns common to your platform and application domain ... for example, most application programmers just use these three patterns
    - SPMDKernel ParallelismLoop Parallelism



# Comparing parallel programming languages/APIs

To compare programming languages and APIs at a high level, we can think in terms of four key elements

Units of Execution	A distinct executable agent that carries out the work of a program. Examples include the threads managed by an OS, processes running on the node of a cluster, or work-items in an OpenCL program
Tasks/mapping	Tasks are a logically related set of operations used to organize the computations in a program. A key aspect of a parallel program is how these tasks are associated (or mapped) onto the units of execution.
Coordination	Mechanisms to manage units of execution (e.g. create, destroy, suspend) and how they interact (e.g. synchronization and communication).
Hardware targets	Most programming models were designed with a particular class of parallel hardware in mind.



# Comparing parallel programming languages/APIs

	Units of execution	Tasks/mapping	Coordination	Hardware targets
Pthreads	threads	Fork join	Shared variables and <b>explicit</b> synchronization constructs	Shared address space computers
OpenMP	threads	Teams of threads with worksharing (loops and tasks)	Shared variables and synchronization constructs	Shared address space computers
MPI	processes	SPMD*	Message passing	Any MIMD* computer
OpenCL	Work-items	Kernel parallelism*		Heterogeneous computers*
CUDA	CUDA-threads	Kernel parallelism*		NVIDIA GPUs

<sup>\*</sup> MIMD (multiple instruction multiple data) and heterogeneous computers will be covered in a latter lecture on parallel hardware. The SPMD (single Program Multiple Data) and kernel parallelism patterns will be covered in our parallel design patterns lecture.

## If you become overwhelmed during this course ...

 Come back to this slide and remind yourself ... things are not as bad as they seem

#### Parallel programming is easy So all you need to do is: Pick your language. - I suggest sticking to industry standards and open source so you can move around between hardware platforms: - TBB - MPI pthreads OpenMP OpenCL Learn the key 7 patterns SPMD Vector Parallelism Kernel Parallelism Loop Parallelism Fork/join - Work Pile Actors - Master the few patterns common to your platform and application domain ... for example, most application programmers just use these three patterns Kernel Parallelism - SPMD Loop Parallelism