

## MPI Network Performance

- Scalability
- Latency
- Bandwidth



1

## Example

$$\begin{array}{|c|c|} \hline a & b \\ \hline c & d \\ \hline \end{array} \begin{array}{|c|} \hline e \\ \hline f \\ \hline \end{array} = \begin{array}{|c|} \hline a \cdot e + b \cdot f \\ \hline c \cdot e + d \cdot f \\ \hline \end{array}$$

CPU 1

CPU 2

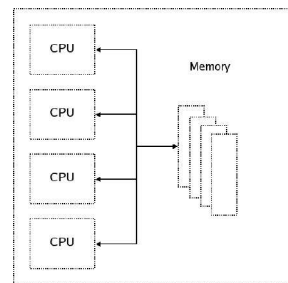
$$M\vec{b} = \vec{c}$$

Many numerical methods use matrix calculation and can be parallelized.  
•BLAS -> ATLAS -> Pthread

2

## Parallel Approaches

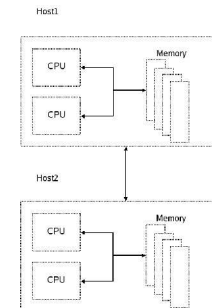
- Posix Threads
  - Well understood
  - Shared Memory
  - Simple Mutexs
  - Not Cheap



3

## Parallel Approaches

- MPI (Message Passing Interface)
  - Shared or distributed Memory
  - Well supported
  - Portable
  - Explicit Data Passing



4

## The Networks

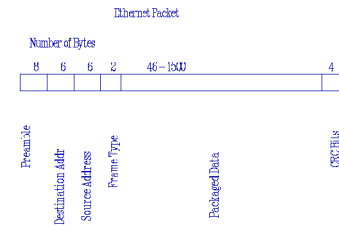
- Myrinet 2000
  - 2Gb/s
  - Uses GM driver
- Ethernet
  - 1Gb/s
  - Jumbo Frames



5

## Ethernet

- Cheap
- Reliable
- Jumbo Frames
- Slow
- TCP/IP



6

## Myrinet

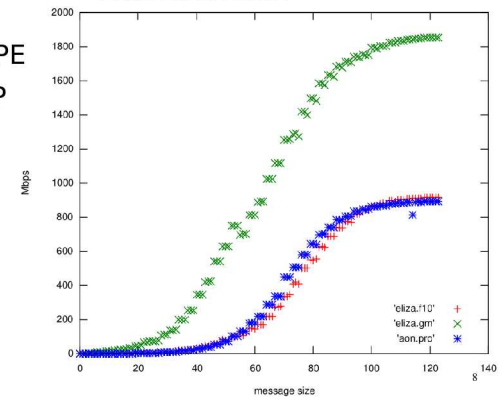
- Fast (For Now)
- No TCP/IP
- Well Supported



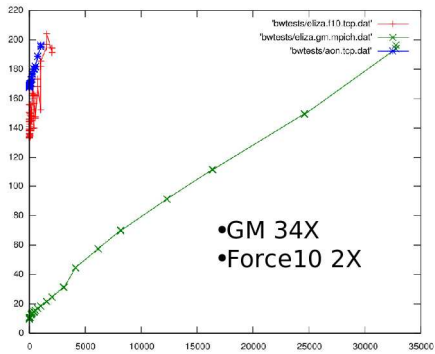
7

## Bandwidth

- NetPIPE
- TCP/IP

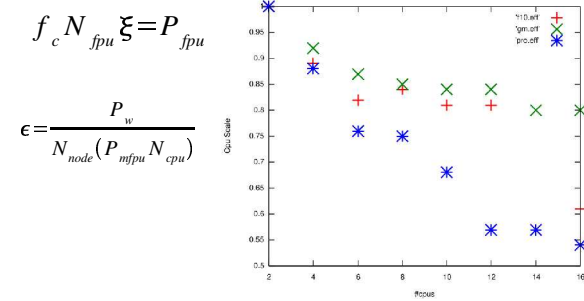


### Latency



9

### Cpu Scaling



10

### Recommendations

- Embarrassingly Parallel
  - MCNP5
  - Seti@home
- Tightly Coupled
  - Boundary Condition
  - HPL



11

### Checklist

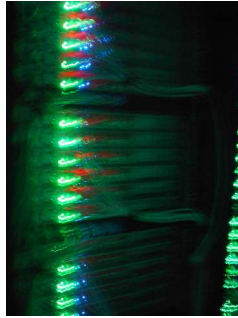
- Problem Run Time
- Problem Nature
- Cost
- Shared System
- Do you *NEED* Shared Memory?



12

## Who are we?

- 584 Nodes (1,168 CPU's)
- 1,244 GB RAM
- 11 TB Shared Disk
- 30 TB Scratch
- 0.58 Tb/s Network
- 4 Clusters 3 Platforms 2 OS's



13

