

**Paul's Norwegian Vacation
(or "Experiences with Cluster
Computing")**

Paul Sack

20 September, 2002

sack@stud.ntnu.no

www.stud.ntnu.no/~sack/

Outline

- Background information
- Work on clusters
- Profiling tools
- Porting Fortran 90 code
- Performance analysis
- MPI Broadcast w/ multicast

IDI Cluster: Clustis

- Description of clustis: 1.46 GHz AMD processors, 2 GB RAM, 40 GB drives
- Use OpenPBS batch system (PBS implementation)
- Classification of nodes:
 - `clustis`: master node
 - `node01-node08`: interactive use
 - `node09-node40`: batch use
- Several queues of different priorities:
 - privileged users
 - guest users
 - `sif80bi` users

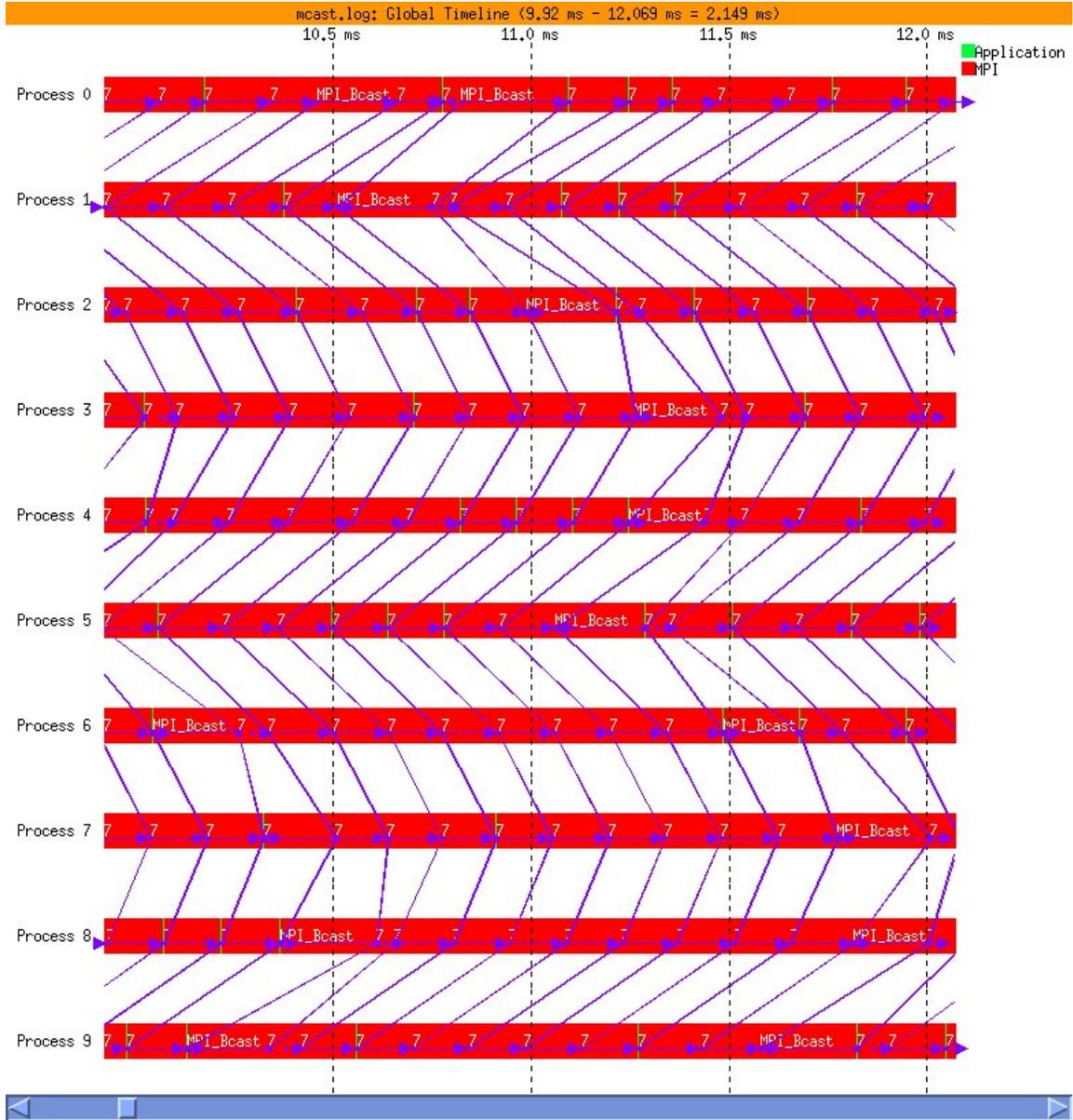
- Documentation
- `qmpirun` script:
 - Usual MPI invocation:
`mpirun -np 4 program`
 - To use batch system:
`qmpirun -np 4 program`
 - Output saved in `program.out` and `program.err`
 - Creates and submits shell script for batch system
 - Cleans itself up when finished
- `mpirun` modified to support OpenPBS

Itanium cluster

- Two IA-64 nodes and master
- Setup OpenPBS and `qmpirun` script as for clustis

Profiling Tools: Vampir

- Produced by Pallas GmbH in Cologne and Dortmund
- Easiest to use: wrappers for `mpicc`, `mpiCC`, `mpif90`
- Simple, effective design: static libraries
- Shows individual MPI calls during a timeslice
- Shows aggregate statistics in varying levels of detail
- Expensive: 75 000 Kroner for clustis



Profiling Tools: Paradyn

- Written at University of Wisconsin
- Difficult to use: non-intuitive GUI interface
- Complex design: run-time code-patching
- Shows even more data than Vampir
- Does not work in batch environments
- Impossible to compile
- Inexpensive: Free

Profiling Tools: Jumpshot

- Written at various universities and research centres in USA
- Comes with MPICH
- Simple design: static libraries
- Has years of testing and works well
- Limited data: no aggregate statistics
- Also free



Porting Fortran 90 Code

- Compilers not standard-compliant
- Most compilers don't fully support standard **and** exceed the standard
- This step took the longest in the project
- SGI compiler
 - SGI supercomputers used for scientific computing
 - Scientists and mathematicians like Fortran
 - Hence, SGI writes a Fortran 90 compiler that produces very optimized code
 - Code originally written on SGI

Intel compilers

- NAGWare (NTNU site-licensed)
 - Compiler produces faulty code in array reshaping code
 - Debugger segfaults
 - NAGWare acknowledges bug, recommends upgrading (£££!)
- Intel
 - Intel also produces supercomputers
 - Good compiler, but strictly standard-conformant
 - My code has some obsolete Fortran 77 syntax
 - Free!

- Portland Group (Bergen and Linköping)
 - Compiler mishandles 3D-array reshape calls
 - Successfully rewrote this as a series of 2D-array reshape calls
 - Sometimes segfaults
- No GNU Fortran 90 compiler yet

Getting MPICH to work with Fortran 90

- Simple mapping in C between function/variable names and symbol names in object files (`foo()` = `foo`)
- Fortran 90 compilers sometimes add one underscore and sometimes add two (`foo()` = `foo_` or `foo__`)
- MPICH configure script is supposed to handle this, but doesn't work
- Since Fortran 90 is mostly used for scientific computing, the Portland Group provides a custom configure script, which does work

- Porting
 - Used Portland Group compiler
 - Joakim Hove and Knut Petter assisted me in correcting many errors
 - Most of the errors were in extraneous portions of the code and could be commented out.
 - By comparing the code with standard-compliant code and a bit of guesswork, I managed to get all of the necessary modules to compile without error

Profiling: Systems

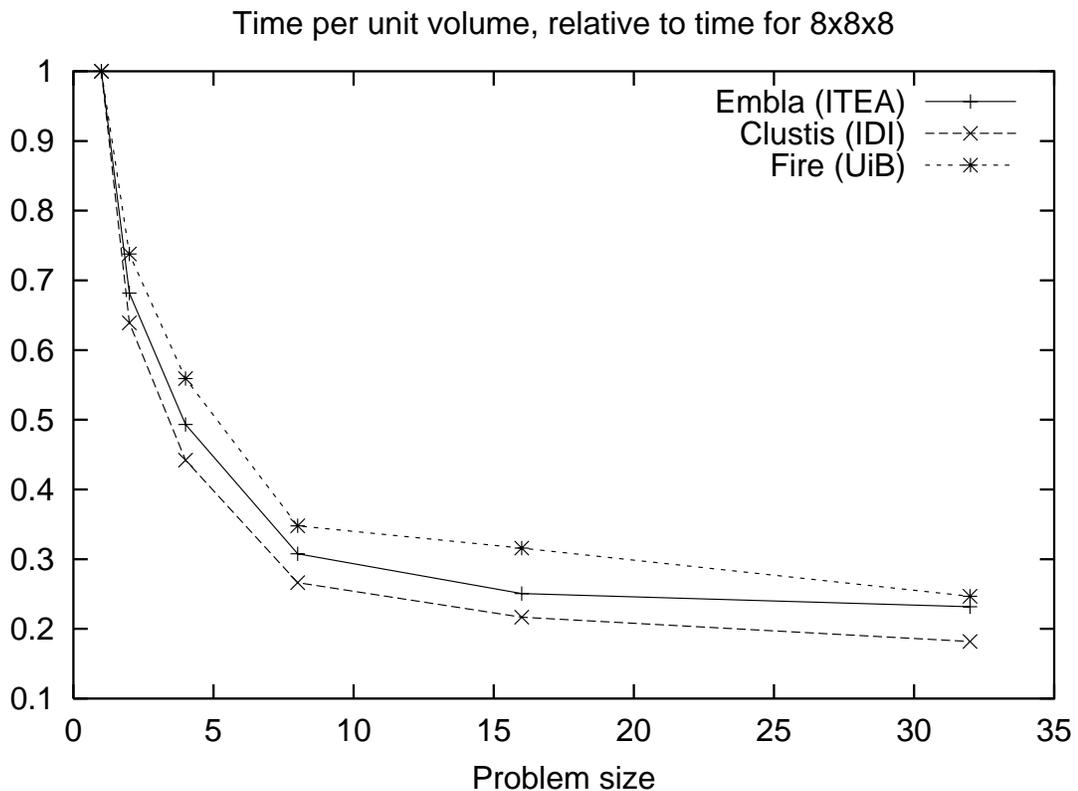
- Cluster in Bergen: 32 × dual 1.26 GHz Pentium IIIs
100 Megabit network
- Cluster in Linköping: 33 × 900 MHz Pentium IIIs
100 Megabit and SCI (full-duplex Gigabit) networks
- Supercomputer (embla): 512 × 600 MHz MIPS processors
- Did not use Clustis

Description of Program

- Simulation of behaviour of subatomic particles in near absolute-zero temperatures under an electromagnetic field
- Three-dimensional parallelepiped divided between processes
- One thousand iterations. Each iteration has a separate communication and a computation phase
- Initial hypothesis:
 - Cluster better for small number of processes (faster processors)
 - Supercomputer better for large number of processes (faster interconnect network)

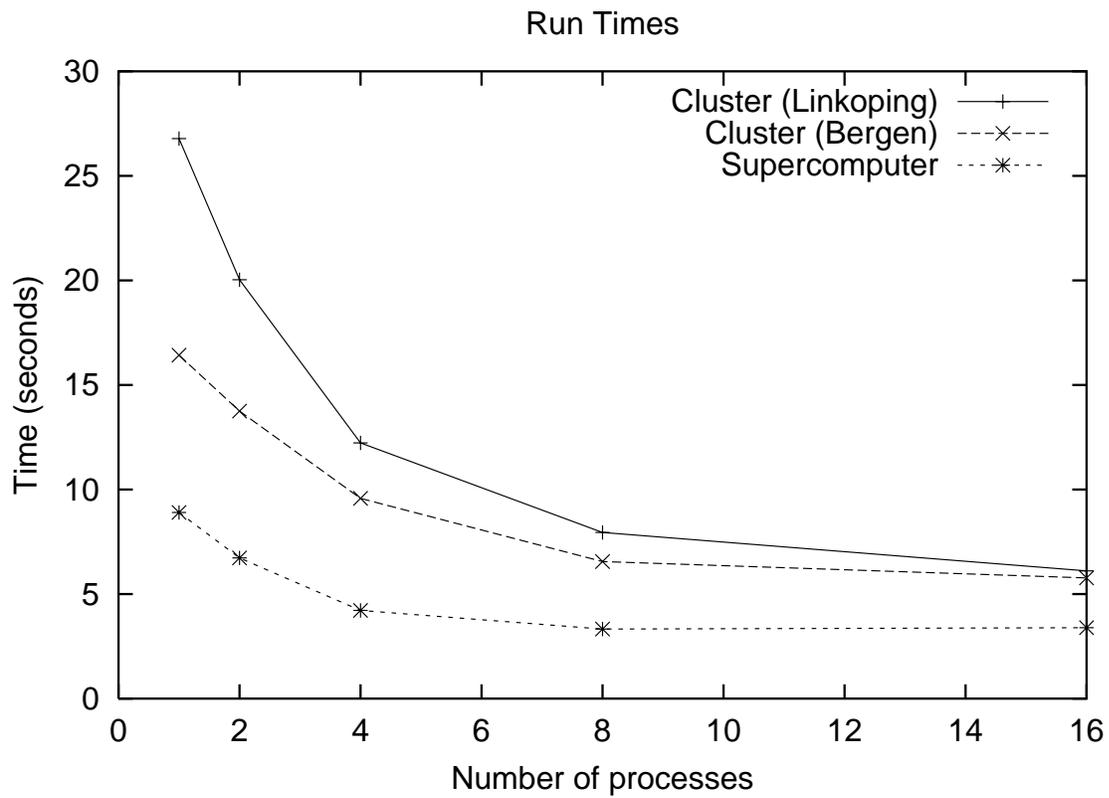
Analysis of Problem Size

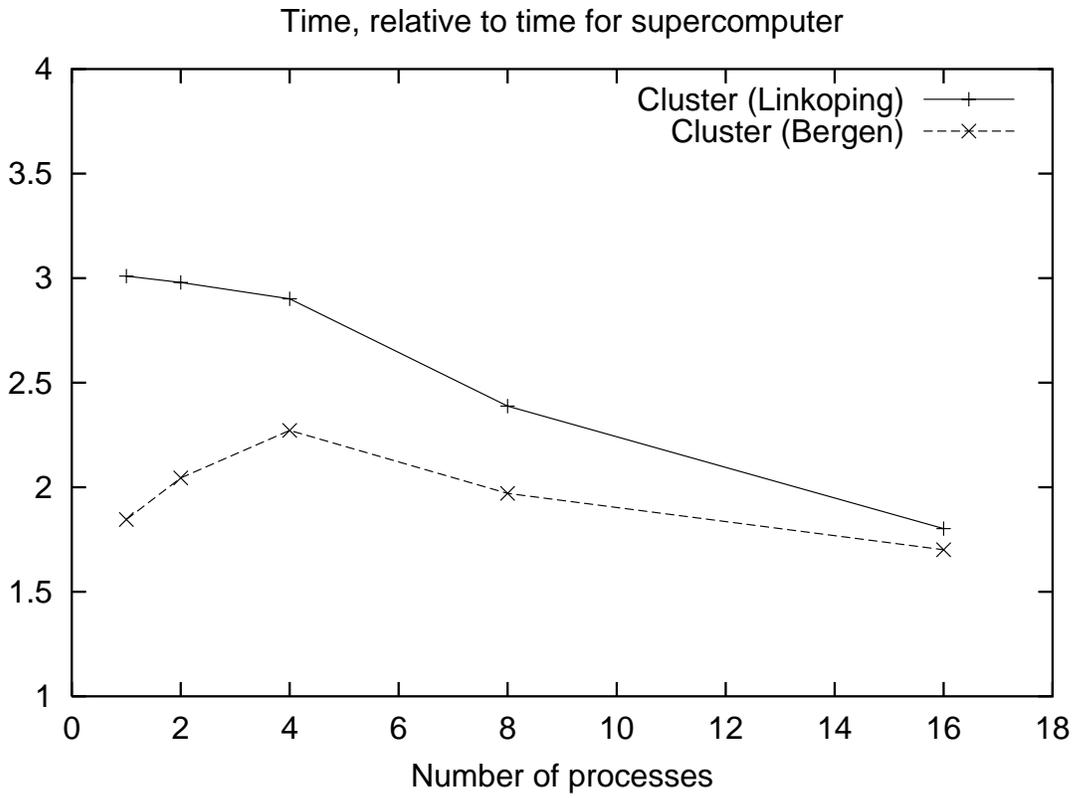
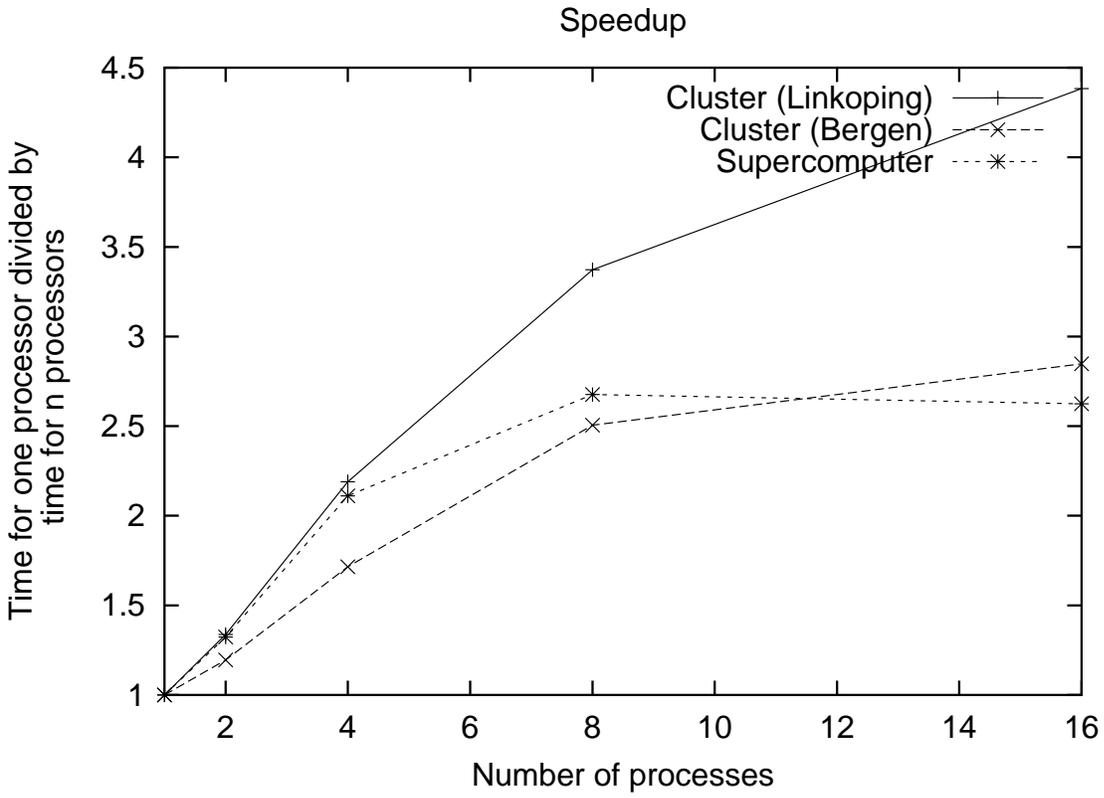
- Number of messages passed is proportional to the internal surface areas
- Amount of computation proportional to the volume
- I expected the program to perform better *per-unit-volume* with larger problem sizes



Profiling Results: Speedup

In this test, the problem size is constant and the number of processes used varies.





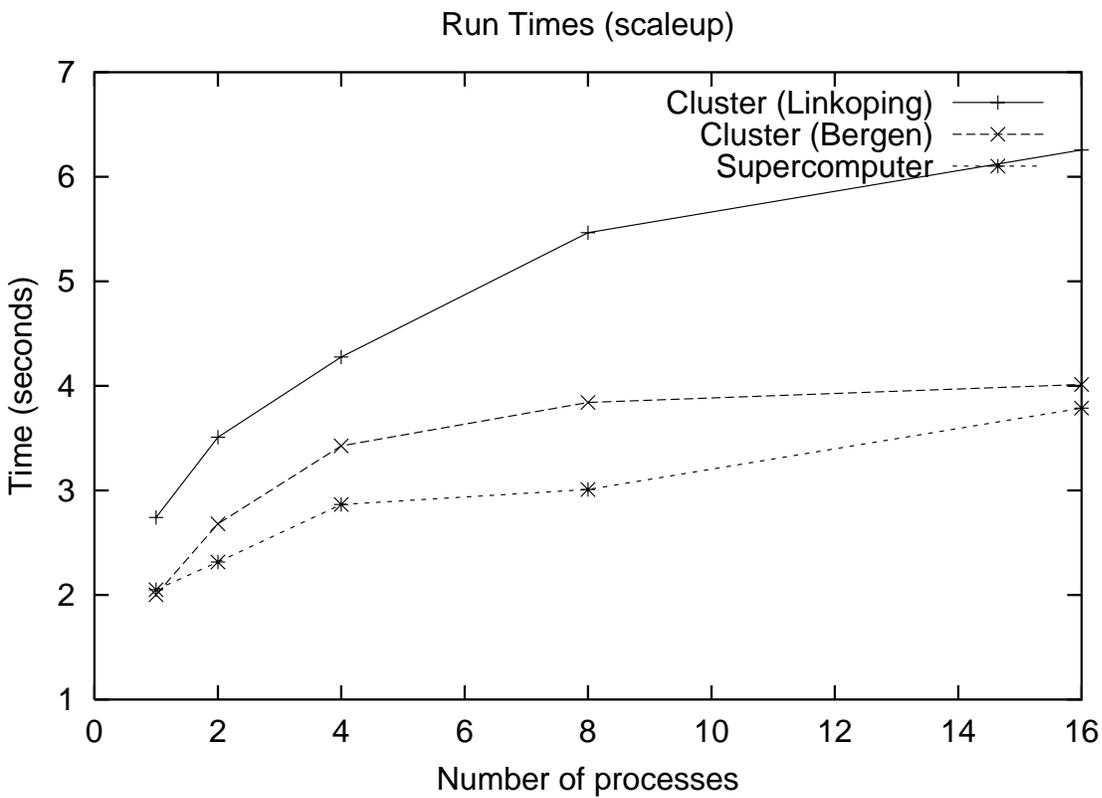
Profiling Results: Speedup

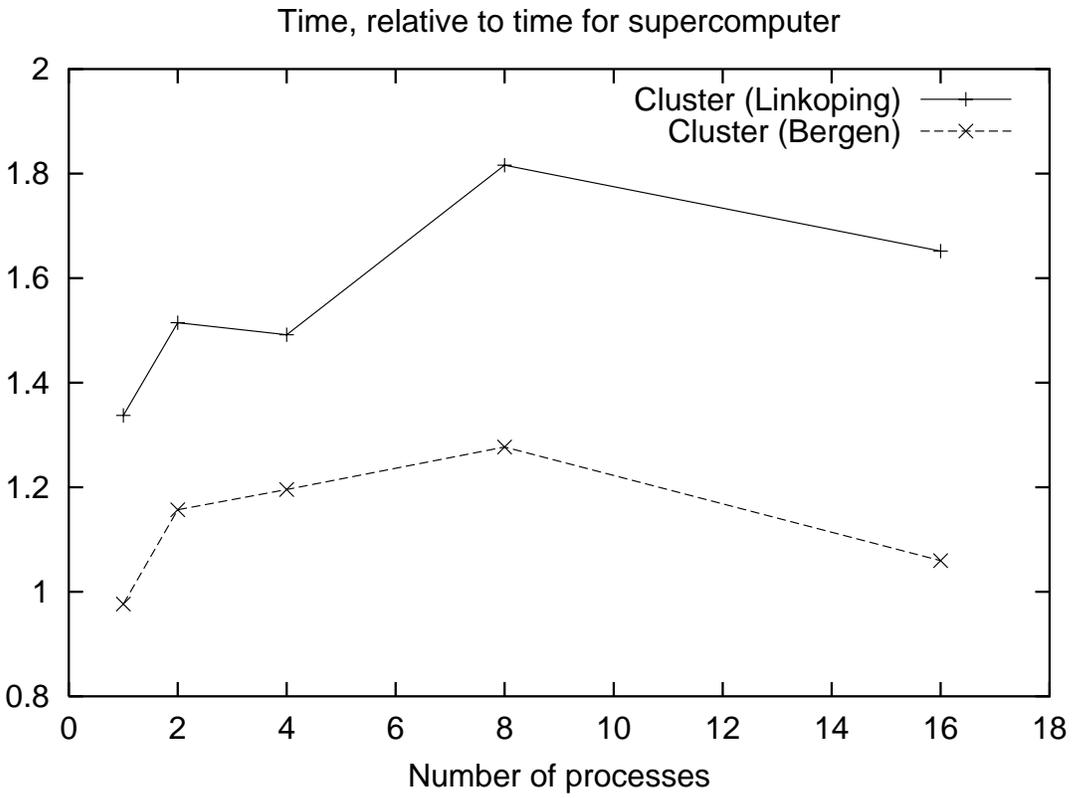
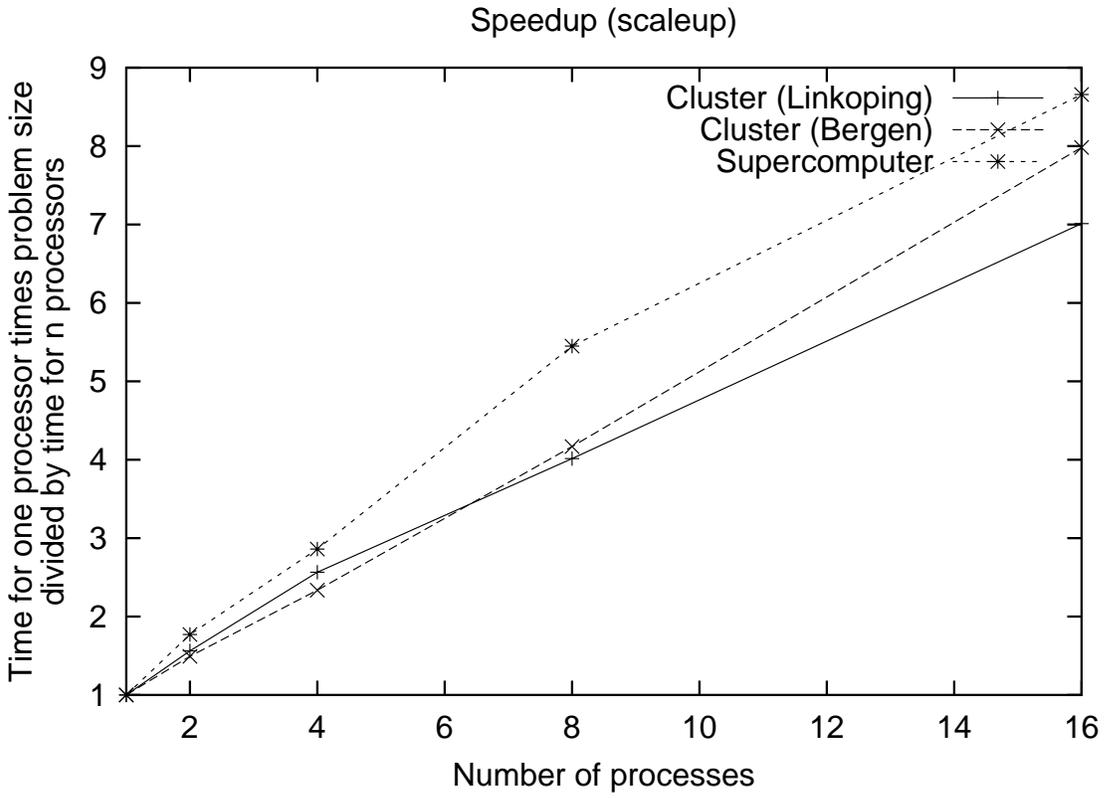
Observations:

- Unexpectedly, the supercomputer outperforms the clusters for a small number of processes
 - SGI compiler produces more optimized code
 - MIPS CPUs have larger caches
- With a large number of processors, the two clusters' performance is approximately equal

Profiling Results: Scaleup

In this test, the problem size is grown in proportion to the number of processes. In other words, the problem size per process is kept constant.





Profiling Results: Scaleup

Observations:

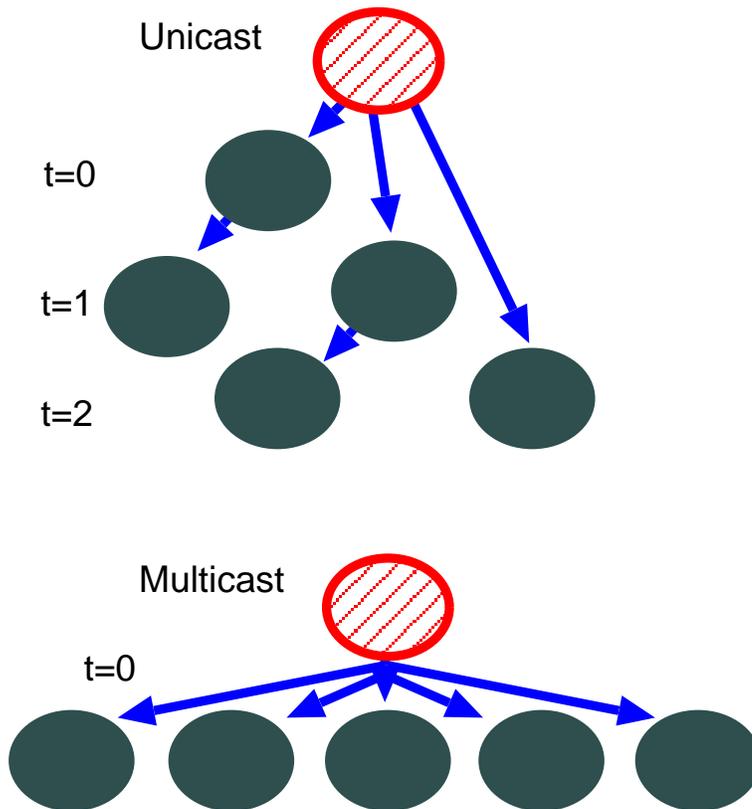
- The clusters fare much better in this test
- The Bergen cluster's performance is within 20% of the supercomputer's
- The Linköping cluster's performance is within 80% of the supercomputer's
- This is a more real-world test, since the amount of work per process is a reasonable amount.
- All the systems exhibit better speedup than in previous test

Profiling Results: Conclusions

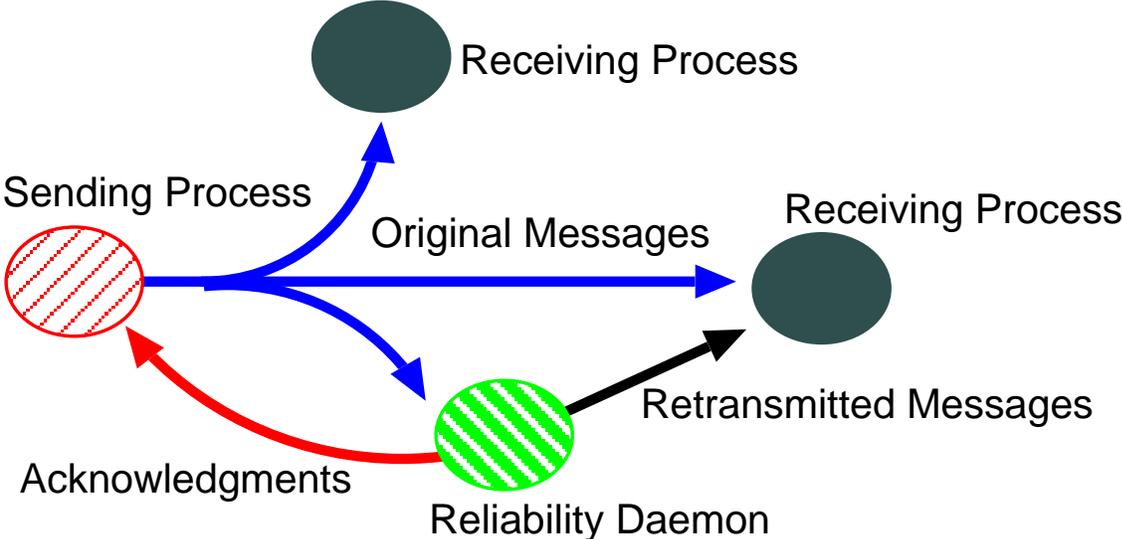
- Clusters perform quite well, considering their cost
- Middle ground: clusters with proprietary Gigabit networks
- Could not get SCI (Gigabit) network to work

MPI Broadcast w/ Multicast

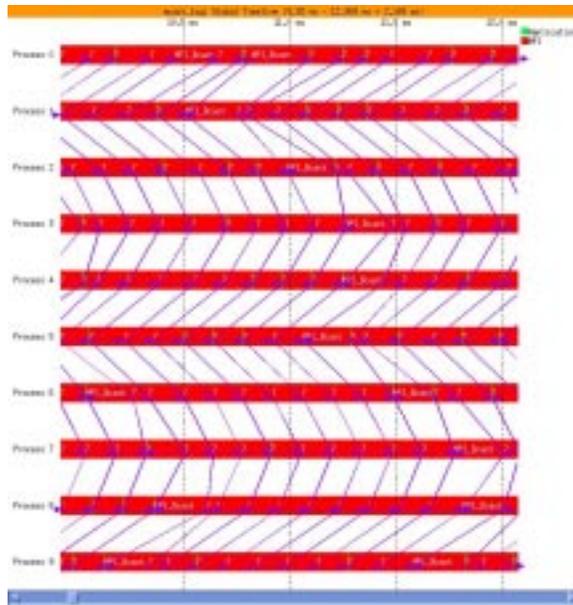
- Broadcast
- With Unicast (TCP/IP) is $O(\log_2 n)$
- With Multicast is $O(1)$
- But Multicast uses UDP, which is unreliable



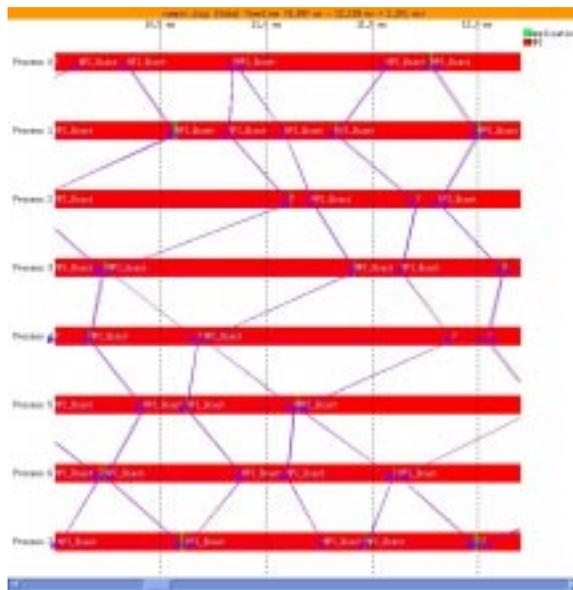
Reliability protocol



Multicast:



Unicast:



Results

