

*The High Performance Cluster for Lattice QCD  
Calculations:  
System Monitoring and Benchmarking*

Lucas Fernandez Seivane

[quevedin@mail.desy.de](mailto:quevedin@mail.desy.de)

Michal Kapalka

[kapalka@icslab.agh.edu.pl](mailto:kapalka@icslab.agh.edu.pl)

Supervisor: Andreas Gellrich

**IT Division**

September 2002



# A cluster?!?

- Lattice QCD calculations are **complicated**
- We need a lot of **power**, so what do we do?  
Buy a **parallel computer**?  
Or maybe a **cluster** → good solution, but we have to **MANAGE** it!
- Manage = install software + test + **monitor** + **benchmark** + **analyse** + repair + etc.

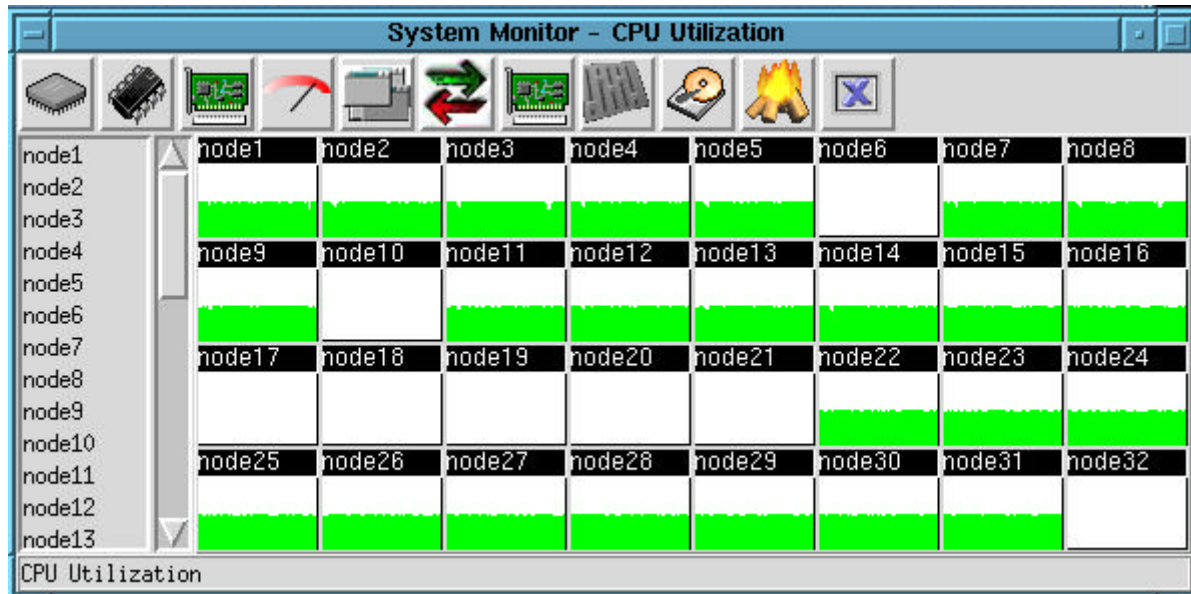


# Monitoring

- **Observing**, how things are going on
- **Measuring** CPU/network/whatever load
- Making **logs** (so nobody plays Quake on 500,000 \$ cluster)
- And calling system administrator at 4 am when a node refuses working



# Monitoring 'lattice' – before



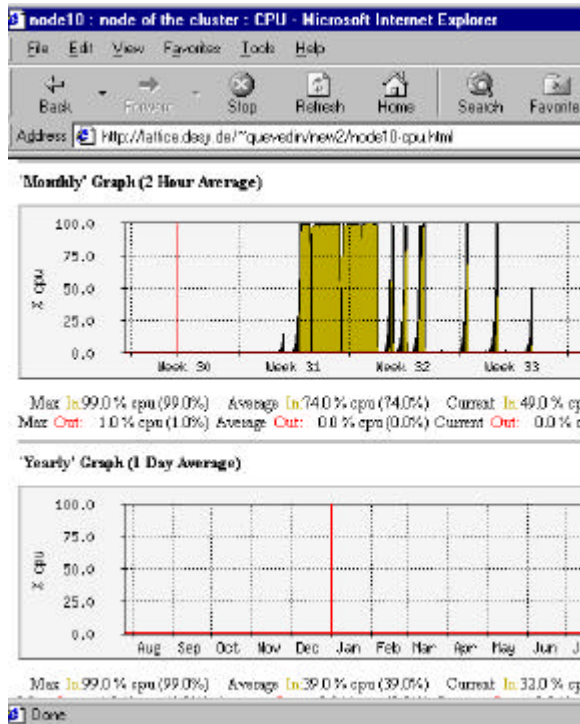
The **lattice** cluster = master + 32 nodes

Each node = dual Xeon CPU (16 x 1.7 GHz + 16 x 2.0 GHz), 1 GB RAM, 18 GB SCSI HDD

Network: Fast Ethernet + **Myrinet**



# Monitoring – after



Lattice Cluster Monitor - Microsoft Internet Explorer

Adresse: [http://lattice.desy.de/cgi-bin/quevedin/cgi\\_clumon2.pl](http://lattice.desy.de/cgi-bin/quevedin/cgi_clumon2.pl)

## Lattice Cluster Monitor

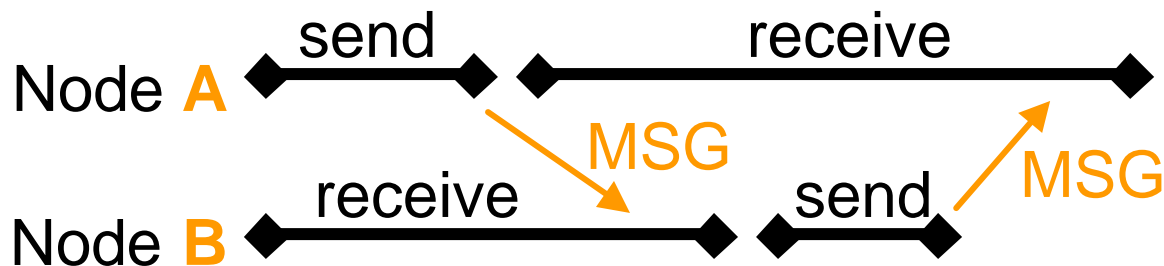
Status: Mon Sep 2 17:55:00 2002 (Version: 2.0 of 31-08-2002)

Host	Date	Load	CPU	Mem [MB]	Swap [MB]	Fan	T [C]
<b>Rack 1</b>							
<a href="#">node1</a>	Mon Sep 2 17:54:15 2002	45	1.00	<a href="#">0.99</a>	<a href="#">952</a>	<a href="#">1</a>	<a href="#">5357</a> <a href="#">33.0</a>
<a href="#">node2</a>	Mon Sep 2 17:50:36 2002	264	1.04	<a href="#">0.99</a>	<a href="#">377</a>	<a href="#">0</a>	<a href="#">5192</a> <a href="#">35.5</a>
<a href="#">node3</a>	Mon Sep 2 17:53:37 2002	83	1.00	<a href="#">0.99</a>	<a href="#">365</a>	<a href="#">0</a>	<a href="#">5113</a> <a href="#">35.0</a>
<a href="#">node4</a>	Mon Sep 2 17:50:16 2002	284	1.01	<a href="#">0.99</a>	<a href="#">364</a>	<a href="#">0</a>	<a href="#">5000</a> <a href="#">33.5</a>
<a href="#">node5</a>	Mon Sep 2 17:52:17 2002	163	1.07	<a href="#">1.00</a>	<a href="#">363</a>	<a href="#">0</a>	<a href="#">5113</a> <a href="#">32.5</a>
<a href="#">node6</a>	Mon Sep 2 17:53:02 2002	118	0.00	<a href="#">0.00</a>	<a href="#">608</a>	<a href="#">0</a>	<a href="#">5152</a> <a href="#">24.0</a>
<a href="#">node7</a>	Mon Sep 2 17:53:00 2002	120	1.08	<a href="#">0.99</a>	<a href="#">366</a>	<a href="#">0</a>	<a href="#">5113</a> <a href="#">35.0</a>
<a href="#">node8</a>	Mon Sep 2 17:51:33 2002	207	1.02	<a href="#">0.99</a>	<a href="#">360</a>	<a href="#">0</a>	<a href="#">5113</a> <a href="#">30.5</a>
<a href="#">node9</a>	Mon Sep 2 17:50:16 2002	284	1.00	<a href="#">0.99</a>	<a href="#">361</a>	<a href="#">0</a>	<a href="#">5113</a> <a href="#">35.5</a>
<b>Rack 2</b>							
<a href="#">node10</a>	Mon Sep 2 17:53:00 2002	120	0.00	<a href="#">0.00</a>	<a href="#">600</a>	<a href="#">0</a>	<a href="#">5113</a> <a href="#">28.0</a>



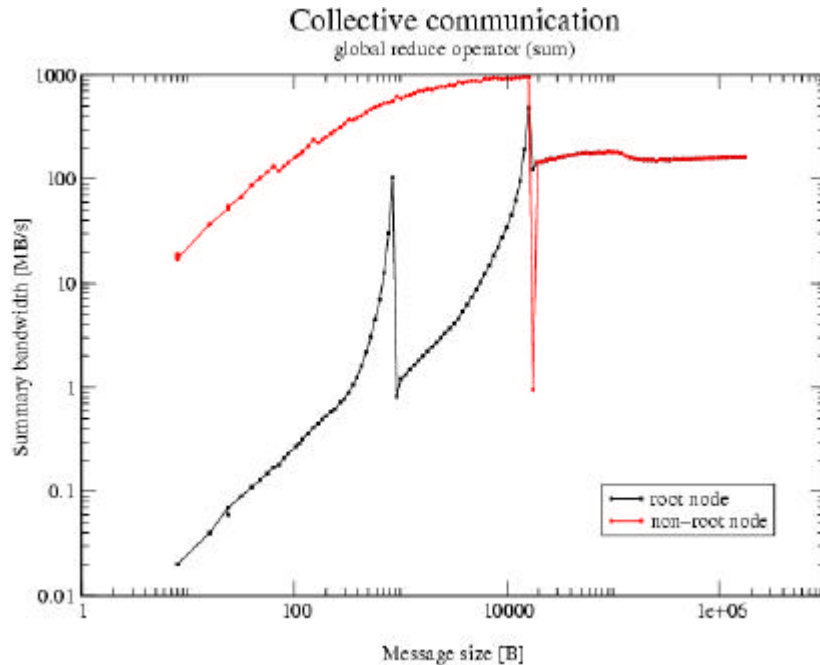
# Benchmarking

- **Benchmark** = compare or **test** (in order to improve things)
- On cluster we test:
  - single** nodes (CPU, memory,
  - communication** between them (here: MPI))
- **MPI** a standard → sending data = passing **messages** between processes





# Benchmarking – results



... NOT trivial ☹️

Interpretation  
is usually ...

