

DESY 03-073

June 2003

Lattice QCD Calculations on Commodity Clusters at DESY

A. Gellrich, D. Pop, P. Wegner, H. Wittig

*Deutsches Elektronen-Synchrotron DESY, 22603 Hamburg and 15735 Zeuthen, Germany
(e-mail: Andreas.Gellrich@desy.de, Peter.Wegner@ifh.de)*

M. Hasenbusch, K. Jansen

*John von Neumann Institut für Computing NIC and DESY, 15735 Zeuthen, Germany
(e-mail: Karl.Jansen@ifh.de)*

Lattice Gauge Theory is an integral part of particle physics that requires high performance computing in the multi-Tflops regime. These requirements are motivated by the rich research program and the physics milestones to be reached by the lattice community. Over the last years the enormous gains in processor performance, memory bandwidth, and external I/O bandwidth for parallel applications have made commodity clusters exploiting PCs or workstations also suitable for large Lattice Gauge Theory applications. For more than one year two clusters have been operated at the two DESY sites in Hamburg and Zeuthen, consisting of 32 resp. 16 dual-CPU PCs, equipped with Intel Pentium 4 Xeon processors. Interconnection of the nodes is done by way of Myrinet. Linux was chosen as the operating system. In the course of the projects benchmark programs for architectural studies were developed. The performance of the Wilson-Dirac Operator (also in an even-odd preconditioned version) as the inner loop of the Lattice QCD (LQCD) algorithms plays the most important role in classifying the hardware basis to be used. Using the SIMD Streaming Extensions (SSE/SSE2) on Intel's Pentium 4 Xeon CPUs give promising results for both the single CPU and the parallel version. The parallel performance, in addition to the CPU power and the memory throughput, is nevertheless strongly influenced by the behavior of hardware components like the PC chip-set and the communication interfaces. The paper starts by giving a short explanation about the physics background and the motivation for using PC clusters for Lattice QCD. Subsequently, the concept, implementation, and operating experiences of the two clusters are discussed. Finally, the paper presents benchmark results and discusses comparisons to systems with different hardware components including Myrinet-, GigaBit-Ethernet-, and Infiniband-based interconnects.

1. Introduction

Lattice field theory has established itself as an integral part of high energy physics by providing important, non-perturbatively obtained results for many physical observables. It complements standard approaches of theoretical particle physics such as perturbation theory and phenomenology, becoming an indispensable method to allow for first principles interpretation of experimentally obtained data. The aim of lattice field theory is to understand the structure of field theories, to test these theories against experiment and in this way to find physics behind the standard model of which we know that it has to be there, but not at what energy scale it should appear.

Another important aspect of lattice field theory is its high computational needs. A distinctive feature of the numerical computations in lattice field theory is the required performance of several Tflops and this performance is needed "in one piece". Such a requirement can only be fulfilled with massively parallel architectures having many processors that are connected via a very fast interconnecting network and work simultaneously on the same problem as a single machine. This distinguishes the computational needs in lattice field theory from the farming concepts usually employed in grid computing of experimental high energy physics.

The combination of high computational needs and the motivation for this in high energy physics as described above, makes a conference such as CHEP an ideal place to present the status and the perspectives of lattice field theory. What strengthens the connection even more is that, at least, a *data grid* is also needed in lattice field theory to exchange expensive data, the so-called configurations and propagators that are generated. An international lattice Data Grid initiative has been started [1].

The computational requirements with high performance in the multi-Tflops regime and fine grained communication can be realized on commercial supercomputers like Hitachi [2] or IBM [3], with specialized machines such as APE (Array Processor Experiment) [4] or QCDOC (QCD on Chip) [5], or with PC clusters as they are discussed in this article. Although commercial supercomputers can conveniently be used, since they are maintained by the corresponding computer centers, their price is normally very high, the efficiency of lattice field theory code is often not optimal and they have many users from different application fields leading consequently to a situation that a single user will not get much computer time.

Specialized, custom-made machines like APE and QCDOC are cost-effective and offer the best price/performance value for multi-Tflops installations. On the other hand, a lot of work has to be spent by

arXiv:physics/0306090 v2 16 Jun 2003

the physicists themselves in order to develop, build and maintain these machines. A compromise between supercomputers and custom-made machines might therefore be commercial PC clusters and in this article¹ we will concentrate on the experiences we have gained with these kind of machines at DESY. We also refer to [6] for a thorough discussion on the above point.

One observation in lattice field theory is that people gather in larger and larger collaborations. This results first of all from the huge computer resources required. It is aimed to use these resources wisely, not duplicating results and find the best strategies for solving the physics problems. One example of such an effort is SciDAC [7] in the US. Another one is the Lattice Forum (LATFOR) initiative in Germany [8]. Associated to LATFOR are also researchers in Austria and Switzerland. In this initiative groups from many universities and research institutes who work on lattice field theory combine their efforts to reach physics milestones in different areas of lattice field theory. They try to coordinate their physics program, to develop and share software and share configurations and propagators, which play the role of very expensive raw data for lattice field theory computations. The research areas of LATFOR are broad and cover

- ab initio calculations of QCD with *dynamical quarks*
 - Hadron spectrum and structure functions
 - fundamental parameters of QCD, i.e. the running strong coupling $\alpha_s(\mu)$ and the quark masses $\bar{m}(\mu)$
 - B-physics
- matter under extreme conditions
 - QCD thermodynamics
 - QCD at non-vanishing baryon density
- Non-QCD physics
 - Electroweak standard model
 - Supersymmetry
- Conceptual developments
 - exact chiral symmetry on the lattice
 - acceleration of the continuum limit
 - non-perturbative renormalization
 - finite size effects

¹This paper covers three talks [9, 10, 11] given during the parallel session on *Lattice Gauge Computing at Computing in High Energy Physics*, UCSD, La Jolla, USA, March 24-28, 2003.

- algorithm development

It would be too demanding (and it is not the purpose of this article) to discuss these topics in detail. The main target of lattice calculations is certainly QCD and we will therefore give in the next section a –presumably too short– introduction to QCD and discuss a few examples of the results that can be obtained.

2. Physics Motivation

2.1. Quantum Chromodynamics on the lattice

In lattice field theory [12], the continuum space-time is replaced by a 4-dimensional, euclidean grid with a lattice spacing a , measured in fm. The advantage of this procedure is that now the theory can be simulated on a computer. In order to obtain back the desired results in the continuum, the values of observables obtained at non-vanishing values of the lattice spacing have to be extrapolated in a continuum limit to zero lattice spacing. At this stage, the comparison to experiments becomes possible and a test of the validity of the model considered can be performed.

On the pure computational side we are dealing with two kind of fields. The first one represents the quarks and are given by complex vectors

$$\Psi(x)_{\alpha,a,n_f}, \quad \begin{cases} \alpha = 1, 2, 3 & \text{color index} \\ a = 1, 2, 3, 4 & \text{Dirac index} \\ n_f = 1, \dots, 6 & \text{flavor index} \end{cases} \quad (1)$$

These fields live on the 4-dimensional space time point x which, on the computer, is represented by integer numbers, $x = (x_1, x_2, x_3, x_4) = a \cdot (i, j, k, l)$, $1 \leq i, j, k, l \leq N$. A second set of fields represents the gluons of QCD and are given by SU(3) matrices U , 3×3 complex matrix with unit norm. The fields $U(x, \mu)_{\alpha,\beta}$ carry again color indices through which they interact with the quark fields. The gluon fields live on the links of the lattice that connect points x and $x + \mu$ in direction $\mu = 1, 2, 3, 4$. The interaction is described by the action²

$$S = \bar{\Psi} M^{-1} \Psi. \quad (2)$$

The action of eq. (2) requires the inverse of the so-called fermion matrix M , or, to be more specific, the vector $X = M^{-1} \Psi$. Without giving the exact definition of the matrix M , the problem is that the fermion

²Of course, in QCD the quark fields are represented as Grassmann variables. We discuss here the *bosonized* form of the action as it is used in simulations.

matrix is high dimensional $O(10^6) \otimes O(10^6)$ and the numerical solution of the linear set of equations

$$M \cdot X = \Psi \quad (3)$$

employing such a matrix is clearly very demanding. It helps, however, that the matrix is sparse. For such a case a vast literature for solving eq. (3) exists [13]. What is important is that the algorithms that can be employed are self-correcting and many of them can be proven to converge in at most a number of steps that corresponds to the dimension of the matrix. Thus we are left with a well posed and regular numerical problem. Of course, in practice the number of iterations is much smaller and typical numbers of iterations to solve eq. (3) are 100–1000. Note that in each of these iterations the matrix M has to be applied to a vector of size $O(10^6)$.

In order to give a feeling about the computational demand, let us give an example. Let us consider a lattice of size $32^3 \cdot 64$ as is realistic in today's calculations in the quenched theory, where internal quark loops are neglected. Then we would need to solve eq. (3) twelve times per configuration for each color and Dirac component. Each solution needs $O(200)$ iterations and we want to perform this on typically $O(1000)$ configurations. Since one application of the matrix M on such a lattice needs 2.8 Gflop, we are left with approximately 6.6 Pflop for obtaining only one physical result for a single set of parameters, i.e. the bare coupling and the bare quark masses. On your standard PC which might run with 500 Mflops sustained for this problem, you would hence need about five months. In order to reach control over the finite size effects, the chiral extrapolation and the continuum limit, simulations at many values of the bare parameters have to be performed.

The numbers above hold for the quenched case, where the quark fields are left out as dynamical degrees of freedom in the simulation. If they are included, the cost of the simulations becomes at least a factor 100 more and a single physical result would need about 40 years on your PC. Clearly, better computers are needed such as the ones discussed in the introduction.

Of course, relying only on progress in the development of computers would be too risky and not enough. Lattice field theory has seen a number of conceptual improvements [14, 15] in the last years that allowed to accelerate the simulations themselves. In addition, many improvements in the algorithms used were found. Although each algorithmic improvement by itself was only a relatively small step [16], all in all a factor of 20 acceleration through algorithm improvement alone could be obtained in the last 15 years. Still, the development of machines were much faster in this period. The situation is illustrated in fig. 1 (taken from a LATFOR paper). Here we show the speedup obtained relative to the status in the year 1987. This year is

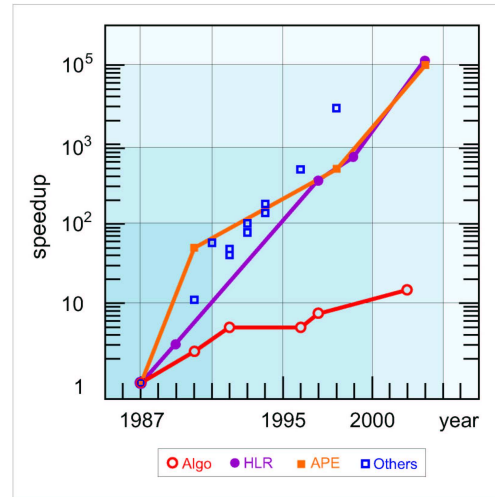


Figure 1: The performance gain of numerical simulations in lattice field theory in the last years relative to the –normalized to one– situation in 1987. Algorithm development (Algo) alone reached a gain of a factor 20. However, the performance gain through computer development appears to be orders of magnitude higher. We show this development at the example of the supercomputers exploited at the research center in Jülich (HLR) and of the APE computers (APE). As a comparison we also show other architectures (others) as used worldwide for lattice field theory computations. (taken from LATFOR).

special in that at that time the first exact algorithm for simulations of dynamical fermions was developed and used [17]. We see in the figure that since then the speedup resulting from better computer techniques, despite the impressive improvements by algorithmic developments, is orders of magnitude larger than the algorithm improvement. The figure also shows that special hardware, in this case the APE computer, shows the same scaling law as commercial supercomputers, in this case the CRAY. Other computers and their performance, relative to the status in 1987, used throughout the world for lattice field theory is also shown.

2.2. Some selected physics results

The machine, algorithmic and conceptual developments in lattice field theory allowed to compute a number of important physical quantities in the last years, despite the aforementioned very large computational requirements. An important aspect of the developments is that we understand not only the statistical errors of the numerical calculations, but also the systematic ones. For a number of quantities fully non-perturbatively renormalized results *in the continuum limit* could be obtained. The only restriction of these calculations is that they are still done in the

quenched approximation. However, there is no particular reason why the calculation should not be doable in a completely analogous way for the full theory. The advent of the next generation of machines in the multi-Tflops regime will then open the door to this exciting perspective.

Let us just discuss two examples of physics results to illustrate what we just said. One is the running coupling [18] and the other are moments of parton distribution functions in deep inelastic scattering. In a quantum field theory such as QCD, there is a steady generation of virtual particles that shield (or anti-shield) the charges of the elementary particles. This leads to renormalization effects. By changing the energy at which experiments are performed, the scale at which we look at the, say, color charge is altered, and, correspondingly the value of the charge itself depends on this energy scale.

This scale dependence (the running) of the coupling can be computed in lattice field, starting from the QCD Lagrangian alone, without any further assumptions. The trick is to use a suitable lattice renormalization scheme, the so-called Schrödinger functional (SF) scheme, that is defined in a finite volume and hence ideally suited for numerical simulations. By going to very high energies, contact to perturbation theory can safely be established and renormalization group invariant quantities can be extracted. In the case of the running strong coupling this corresponds to the Λ -parameter of QCD. The advantage of the knowledge of renormalization group invariant quantities is that they can be translated to any preferred renormalization scheme. In this way it becomes possible to translate results for the running coupling from lattice simulations to continuum results in the, say, $\overline{\text{MS}}$ -scheme as it is conventionally used in perturbation theory.

In fig. 2 we show an example of such a calculation. The results are already in the continuum. They cover a broad energy range and are very precise (the error bars of the simulation points are well below the size of the symbols). In the plot, also a comparison to perturbation theory is shown and a good agreement down to surprisingly small energy scales is found. Another example is a moment of a parton distribution function as they can be extracted from global analyses of experimental data. Such moments can be expressed as expectation values of local operators and are hence computable in lattice simulations. The renormalization procedure of such moments follow the general strategy of using the finite volume SF renormalization scheme discussed above for the running coupling. We show in fig. 3 an example of the continuum limit of the first moment $\langle x \rangle$ of a twist-2, non-singlet operator in a pion [19]. In the plot two different lattice formulations of QCD were used. It is reassuring that in the limit that the lattice spacing is sent to zero both unphysical lattice versions of QCD extrapolate to the same num-

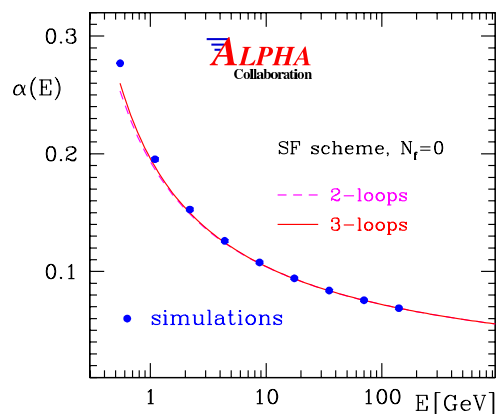


Figure 2: The running strong coupling constant in the continuum as function of the energy scale. The quenched approximation is used.

ber. As a final result for this study it is found that the lattice gives a preliminary value of $\langle x \rangle = 0.30(3)$ while the experimental number reads $\langle x \rangle = 0.23(3)$. This lattice number has to be taken with care since it is obtained in the quenched approximation. But, again, there is no other reason than missing computer power to repeat this calculation also for the full theory. If in this case the results just mentioned were stable, a really interesting situation would have emerged.

The results that we have just discussed were actually obtained on APE machines. There are, however, a number of physics problems, where PC clusters were used extensively. The simulations concern in particular the recently discovered chirally invariant formulations of QCD on the lattice. Examples for results on the PC clusters that are installed at DESY are given in [20]. The next sections are devoted to a discussion on PC cluster systems that are installed at DESY.

3. Commodity Clusters at DESY

3.1. Conceptual considerations

Sufficient computing power to perform *Lattice QCD* (LQCD) calculations as described above can obviously not be drawn from a single processor. The solution is to parallelize the physical problem in order to concurrently deploy many CPUs. As a consequence massive

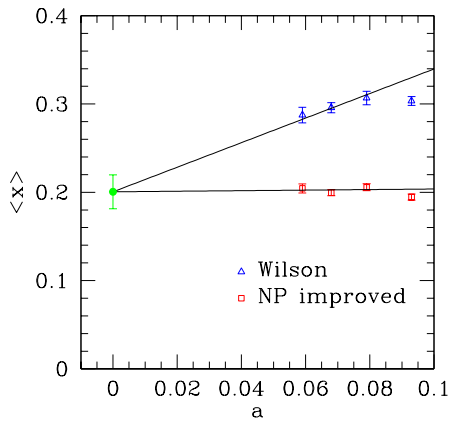


Figure 3: Continuum limit of the second moment of a twist-2 operator in a pion. Two versions of lattice QCD are used, ordinary Wilson fermions (Wilson) and $O(a)$ -improved fermions (NP improved).

intercommunication between the CPUs is needed. Typical *High Performance Computing* (HPC) or *Supercomputing* applications are characterized by high demands on:

- CPU (especially *Floating Point Unit* (FPU)) performance,
- memory throughput,
- interconnectivity (bandwidth and latency) between nodes.

One example for so-called *supercomputers* are *Symmetric Multi-Processing* (SMP) machines with up to hundreds of processors. They appeared as single machines and are optimized for memory access and interconnectivity between the CPUs.

In the LQCD area custom-made special purpose machines such as APE (Array Processor Experiment) [4] or QCDOC (QCD on Chip) [5] have been developed.

In the last years PCs, exploiting processors with competitive computing power, hit the commodity market. Concurrently, modern network technologies have achieved performances, which allow for high-speed low latency interconnects between PCs. These developments paved the way to build PC *clusters* [21]. Commodity clusters draw computing power from up to hundreds or even thousands of in principle independent PCs with one or two CPUs, called *nodes*. Those clusters benefit from their scalability and the possibility to deploy components of the commodity market

with good price/performance ratios. Interconnectivity is provided by exploiting modern network technologies. Such a cluster must:

- deliver sufficient CPU (FPU) performance and memory throughput,
- provide good connectivity between CPUs (bandwidth as well as latency),
- be scalable,
- be reliable,
- incorporate tools for easy installation and administration,
- provide a usable software environment for the applications,
- be connected to backup and archiving facilities,
- fit boundary conditions such as space, cooling and power supply capacities.

In clusters, a set of main building blocks can be identified:

- The **computing nodes** which actually provide the computing power, optionally with local disk space,
- a **high-speed, low latency network** for the parallelized physics application,
- an **auxiliary network** to remotely control and administer the nodes,
- a **host system** for login, compiling, linking, batch job submission, and central disk space,
- optionally, a **slow control network**, e.g. based on a field bus.

A schematic view is shown in fig. 4. The high speed network can either be organized as a switched network (e.g. the DESY clusters using Myrinet-switches) or by a *n*-dimensional *mesh* to allow for nearest-neighbor communication, see fig. 5. In [22] the network is organized as a 2-dimensional GigaBit-Ethernet mesh.

In scientific computing Unix-like operating systems have always played the dominant role. The development of Linux along with the triumphal procession of PCs into the scientific world made Linux-PCs the systems of choice in most universities, physics institutes, and laboratories. In order to actually operate Linux-PC clusters, further system aspects must be taken into account:

- Installation and administration of the **operating system** Linux,
- **security** issues (login, open ports, private network),

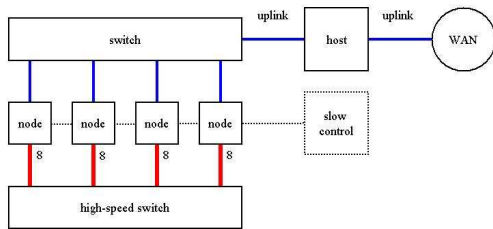
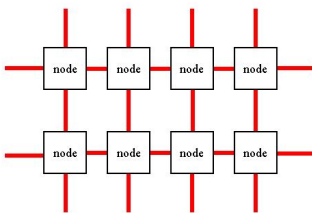


Figure 4: Schematic view of a cluster.

Figure 5: Schematic view of a *mesh* cluster.

- user **administration**,
- application software **installation**,
- **backup** and **archiving** issues,
- **monitoring** and **alarming**.

3.2. Implementation

The recently improved PC architectures are well suitable for Supercomputing by exploiting:

- Increasing CPU clock rates following *Moore's Law* now extending the 2 GHz border,
- larger caches at full processor speed,
- vector units (SSE, SSE2),
- cache pre-fetch,
- fast memory interfaces,
- PCI-Bus at 66 MHz/64-bit,
- high external bandwidth (Myrinet, GigaBit-Ethernet), partly with very low latency (Myrinet).

3.2.1. Hardware

At both DESY sites in Hamburg and Zeuthen commodity clusters are operated since January 2002 and December 2001 respectively. See tab. I for the set-ups.

Table I DESY's PC clusters.

Item	Hamburg	Zeuthen
Nodes	32	16
CPUs/node	2	2
CPUs/1.7 GHZ	2 × 16	2 × 16
CPUs/2.0 GHZ	2 × 16	

The cluster nodes as well as the host system are equipped with high-end commodity components (see tab. II).

Table II Cluster Hardware.

Item	Implementation
Computing Nodes	
Chassis	rack-mounted 4U module
Main-board	SuperMicro P4DC6
Processors	2 Intel Pentium 4 Xeon 1.7/2.0 GHz
Chip-set	Intel i860
Memory	4 × 256 MB RDRAM
Disk	18 GB SCSI IBM IC35L018UWD210-0
Host System	
Chassis	rack-mounted 4U module
Main-board	SuperMicro P4DC6
Processors	2 Intel Pentium 4 Xeon 1.7 GHz
Chip-set	Intel i860
Memory	Rambus 4 × 256 MB RDRAM
Disk	36 GB SCSI IBM DDYS-T36950N
Uplink	Intel EtherExpress PRO 1000 F
Downlink	Intel EtherExpress PRO 1000 T
High-speed Network	
Interface cards	Myrinet M3F-PCI64B-2
Chassis	Myrinet M3-E32 5 slot
Line cards	Myrinet M3-SW16-8F
Mngmnt card	Myrinet M3-M
Auxiliary Network	
Interface Card	on-board Intel 82557 100Base T
Switch	Compu-Shack GIGALine 2024M 48-port 100Base T
Uplink	Module 1000Base T

Intel Pentium 4 Xeon processors [23], which became available at the end of 2001, showed much enhanced performance compared to Pentium III Tualatine CPUs due to features such as vector units, and

larger caches (see the section 4 on benchmarking results). The SuperMicro motherboard P4DC6 with the Intel i860 chip-set was the only possible combination until the end of 2002. It provides PCI-Bus support. For the memory Rambus modules were chosen. The communication between the nodes in the physics application is done by means of Myrinet [24]. It provides bandwidths up to 240 MB/s with very low latencies in the order of a few μ s. For administration purposes and to actually submit jobs and copy data the nodes are interconnected via Fast-Ethernet in a private subnet. The host system is connected to the central switch by a copper GigaBit-Ethernet link and has a separate fiber GigaBit-Ethernet link to the outside. Each node consists of a PC with two CPUs and is housed by a 4U rack-mounted chassis. Up to 9 nodes (18 CPUs) are installed in a cabinet.

3.2.2. Software

The basic installation plus most of the software support was purchased with the hardware from the German company MEGWARE [25].

For all software related aspects the host system is used as a server for the nodes.

For the Linux operating system a S.u.S.E. 7.2 distribution [26] was chosen, which contains the kernel version 2.4.17 with SMP capabilities. Temperature and fan sensors are read out by the kernel module *lmsensors* [27].

The nodes are booted via *DHCP*, *TFTP*, and Intel's [30] *Pre-boot Execution Environment* (PXE).

The nodes are operated in a private network (192.168.1.0) behind the server. For security reasons, external user login is only possible to the host system via *ssh*. Individual login from the host system to the nodes can be done over the auxiliary network by means of *rsh*. Users are registered on the host system which exports the home directories (*/home*) and a data partition (*/data*) to the nodes. For the user administration standard Unix tools are used (*useradd*). The necessary files (*/etc/groups*, */etc/passwd*, */etc/shadow*) are manually distributed to the nodes. Source code is compiled and linked on the host system. Software is mostly written in C/C++ but also compilers for Fortran77/90 are required. On the host system in addition to GNU compilers of the Portland Group [28], KAI [29], and Intel [30] are available.

The parallelization of the computation is done within the application by means of the *Message Passing Interface* (MPI). Since MPI is running over Myrinet, a special library which uses low level Myricom communications is installed (MPICH-GM [31]). For the Myrinet network a static routing table is used.

Zeuthen uses the open source *Portable Batch System* (OpenPBS) for job submission. In Hamburg nodes are manually distributed to users on good-will basis. Time synchronization is done via *XNTP*. The nodes

synchronize with respect to the host system which gets the correct time from DESY's central server.

Backup and archiving are basically different items: Regular backups are done to provide security against loss of system data and home directories. Under normal conditions backups will never be retrieved. For the Hamburg cluster DESY's standard backup environment based on IBM's *Tivoli Storage Manager* (TSM) is used. It automatically creates incremental backups of the disk of the host system and a regular basis. At DESY Zeuthen a copy of the host system's disk is stored on a second disk.

Archiving tools allow users to arbitrarily store and retrieve large amounts of data. DESY uses *dCache* which provides a simplified and unified tool to access the tertiary storage [32]. It provides a unique view into the storage repository, hiding the physical location of the file data, cached or tape only. Dataset staging and disk space management is performed invisibly to the data clients. Currently around 1 TB of data are stored. The archiving system is mainly used to store temporary *check-points* and final results of long time computing jobs, so-called *configurations*, which can be used for further analyses.

The backup and archiving scheme is shown in fig. 6.

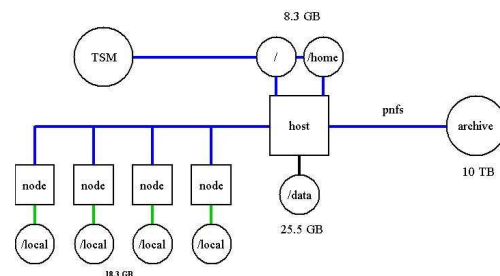


Figure 6: The backup and archiving scheme.

LQCD calculations require the availability of **all** nodes dedicated to the problem for the **entire** runtime of the job. The failure of only one node spoils the entire job. Therefore, stability of the cluster in terms of availability and sustained performance is a crucial credential of the clusters. Usually not all nodes are used for one single job. Many calculations are done on a limited number (2 – 8) nodes. This allows to (manually) restart jobs using the check-points on different nodes in case of failures. In order to use the resources of the clusters most efficiently a well-defined monitoring and alarming scheme is needed.

The company MEGWARE delivered with the software installation a monitoring package called *clustware*. It provides a snapshot of all relevant properties of the cluster nodes in a graphical user interface, including CPU usage, load, I/O, and temperature. A long term

history ($> 1min$) is not shown.

Alternatively, a DESY in-house development called *ClusterMonitor* (CluMon) [33] is in use. Every node runs a simple Perl-written daemon which periodically dumps status information of relevant node properties such as uptime, load, CPU usage, memory usage, swap usage, and all temperatures into a node-specific file on the host system. The host system runs an *Apache* web server, which allows to remotely access status information of the cluster with any web browser. History of the quantities is kept by means of the *MRTG* package [34] and is available from the web page. In addition, *CluMon* provides alarming by e-mail based on the time period since the last update of the status file.

3.3. Operational experiences

The two clusters in Hamburg and Zeuthen are located in the computer centers. The clusters are operated and administrated in cooperation by members of the computer centers and the DESY theory and NIC groups.

Of around twenty registered users per cluster just a handful can be classified as *Power Users*, running regularly resource consuming jobs (see tab. III).

Table III Accounts and users.

Site	Accounts	Power Users	Strategy
Hamburg	20	7	good-will basis
Zeuthen	18	6	batch system

The two clusters in Hamburg and Zeuthen deploy in total 52 dual-CPU PCs: 32 nodes (Hamburg), plus 16 nodes (Zeuthen), plus 2 spare nodes, plus 2 servers. During the almost 17 month of operation quite a number of problems occurred, so usually not **all** nodes have been available at **all** times. Tab. IV lists all major hardware problems.

Taking into account that the clusters exploit commodity hardware components and offer considerably better price/performance ratios than big mainframe SMP-machines, failures of certain components such as disks and power supplies were expected. The stability of the PC hardware after replacing the obviously systematically misbehaving IBM disks was reasonable. Some annoyance was caused the repeating failures of the Myrinet interface and line cards, which was also seen at Fermilab [35] and traced back to broken optical receivers. Nevertheless, the general opinion of the users on the performance and the stability of the clusters is very positive.

Table IV List of failures.

Component	Faults	Total
PC		
Motherboards	1	52
CPUs	0	104
Memory Modules	0	208
Power Supplies	1	52
Disks	35	52
Ethernet Chips	0	52
CPU Fans	0	104
Chassis Fans	0	52
Myrinet		
Fibers	1	48
Slot Chassis	0	2
Line Cards	4	10
Interface Cards	6	48
Infrastructure		
Cabinet Fans	3	24

3.4. Future developments

The considerations in section 2 require cluster sizes of $O(1000)$ nodes to approach the Tflops regime. Even more, in order to actually deliver a few Tflops sustained for hours, days or even weeks, **all** nodes would need to run at the same time. As discussed earlier the failure of just one node would spoil the **entire** calculation. Accepting the experiences so far, this seems not be possible within the current concept³.

Recent tests have shown that GigaBit-Ethernet might be an interesting alternative to Myrinet. Benchmark showed bandwidths of 2×1 Gbit/s bidirectional with special switches which would imply a much better price/performance ratio than Myrinet. Since GigaBit-Ethernet is widely used now, one might also expect more stability and reliability compared to the niche product Myrinet (see section 4).

Disks –even SCSI– are the most likely components to break in PCs. Though the replacement of disks is easy, the affected node is down and needs to be re-installed completely afterwards. This could be avoided by running the nodes disk-less. Such a concept would require a stable and reliable server which could be achieved by setting up one or more RAID-systems in a tree-like architecture to distribute load. The server could also provide the boot image. In another scenario booting could be done from a local block device such as an EPROM or a memory stick.

The current set-up relies on one single server which

³1 broken node per week in a 50 node cluster is equivalent to a maximal lifetime of a complete 1000 node cluster of 8 hours.

exports `/home` and `/data` directories to the nodes. It also serves as a login host for the users and is used for code development, compilation, linking, and job submission. This machine is clearly a single point of failure. In a bigger system one would opt for redundancy in the server arrangement by distributing different functionalities to different machines. In particular a separate file server for the exported directories is needed.

Space, power consumption, and cooling will become a major issue when planning for thousands of nodes. Recent developments of so-called *blades* place the motherboards vertically to improve the air-flow for cooling in order to increase CPU densities.

Software installation, administration, and monitoring of thousands of nodes is a challenge which requires a very careful choice of appropriate tools. Remote administration could be enabled by exploiting the serial consoles of the PCs. They could be connected to a dedicated terminal server or—in a much cheaper scenario—subsequently from node to node.

4. Benchmarks

The hardware of commodity PCs has been extremely improved over the last years. The performance increase inside the CPU is due to higher clock rates and enhanced building blocks (e.g. SSE1/SSE2 instructions) following Moore's Law which predicts a performance doubling all 18 month. Moore's Law gives a technology estimation of mainly the CMOS density or number of transistors which can be integrated on a chip of a given size. It does not work well for the other interacting PC components like the memory interface and external busses. On the other hand also a big step forward in the development of fast memory architectures like Rambus and DDR RAM and a series of high bandwidth PCI bus based interconnects like Myrinet and QSNNet is going on. Therefore PC clusters are becoming more and more attractive for classical high performance parallel computing and therefore also as a hardware basis for LQCD applications.

4.1. Benchmark systems

Apart from the DESY clusters described above, the following systems were used in order to test the ability of PC clusters for LQCD applications:

Mellanox: Blade dual Pentium 4 Xeon cluster connected via Infiniband, running MPICH for VIA/Infiniband with patch from Ohio State University [37],

ParTec: Dual Pentium 4 Xeon cluster connected via Myrinet running ParaStation MPI [36],

MEGWARE: Dual Pentium 4 Xeon cluster connected via Myrinet running MPICH-GM from Myricom [31],

Leibniz-Rechenzentrum Munich: (single CPU tests) Pentium 4 and dual Xeon PCs with CPUs with clock rates between 2.4 and 3.06 GHz,

University of Erlangen: GigaBit-Ethernet dual Pentium 4 Xeon cluster.

4.2. Benchmarks and results

Representative benchmarks for the evaluation of different PC systems have been developed. Already in the year 2000 a first benchmark of M. Lüscher (CERN) has shown the potential in using the SSE1 and SSE2 instructions for the Wilson-Dirac operator [38]. This program takes heavily advantage of the Pentium 4 memory-to-cache pre-fetch capabilities and the SSE registers and instructions which are implemented by using assembly in-line code, compatible to the gcc and Intel compilers. Fig. 7 shows on the left hand side the performance gain of the highly optimized 32-bit and 64-bit Dirac-Operator kernel which linearly follows the evolution of the CPU performance expressed by their clock rate. The value of 1.5 Gflops for the 32-bit implementation or 0.8 Gflops for the 64-bit implementation respectively was unexpected high for a PC in the year 2000 and encouraged groups working on LQCD algorithms on PC clusters also to use the Pentium 4 capabilities to improve their algorithms on PC clusters.

The Wilson-Dirac operator Benchmarks are accompanied by two tests called `add_assign_field` (similar to the BLAS `daxpy`) and `square_norm` which are representing the linear algebra part of the benchmark. Both parts are strongly memory bound which means that they cannot benefit from the SSE-environment. This results in a relative small improvement shown on the right hand side of fig. 7 which also gives an impression of the slowly evolving memory interface architectures since the introduction of the dual channel 800 MHz Rambus.

Another version of such a single node benchmark was developed by M. Hasenbusch (DESY) for the even-odd preconditioned Wilson-Dirac operator [39]. Meanwhile (using recent FSB800 based PCs equipped with a 3.06 MHz Pentium 4) a performance of about 2.6 Gflops for the 32-bit implementation and about 1.4 Gflops for the 64-bit implementation can be observed.

To evaluate the behavior of PC cluster interconnects a 1-dimensional parallel even-odd preconditioned Dirac Operator Benchmark on a 2×16^3 lattice (also written by M. Hasenbusch) was used (see Appendix). The aim of the parallel benchmark was to compare different parallel PC based architectures against each

other rather than achieving the best performance for a given system. Fig. 8 shows the results on clusters with different numbers of nodes. In addition to the CPU power and the memory interface the throughput of the external PCI-Bus depending on the given chip-set and the interconnecting interface itself are dominating the entire performance. Both early Intel Pentium 4 Xeon based clusters at DESY are using the i860 chip-set which came with a relative poor 33 MHz/64-bit PCI-Bus performance. A bus-read (send) of 227 MB/s and a bus-write (recv) of 315 MB/s of maximal 528 MB/s and expected 450 – 460 MB/s was measured. This ends up in an external unidirectional bandwidth of about 160 MB/s of maximal 240 MB/s. In the bidirectional case we measured (90 + 90) MB/s. Meanwhile more advanced chip-sets like the E7500 were available for the Pentium 4 Xeon which are providing a PCI throughput close to the expected numbers. The influence of the chip-set is dominating the results shown in fig. 8, whereas in the left hand side on CPU per node and on the right hand side two CPUs per node communicating via shared memory MPI are used.

Infiniband is a new promising communication technology especially designed for cluster interconnections fig. 9. Beside the high throughput shown in fig. 10 the latency at relative small buffer size in the order of 2 kB is significantly higher then using Myrinet which could be an advantage for applications using small local lattices. In fig. 11 the performance of a 4 node partition of the DESY Myrinet Pentium 4 Xeon cluster is compared to the performance of a corresponding Infiniband cluster from Mellanox using a 2-dim parallel Wilson-Dirac Operator Benchmark developed at DESY Hamburg by M. Lüscher. Due to some problems using the assembly in-lines within the 2.96 gcc compiler coming with the RedHat Linux system on the Mellanox cluster during the short time in which the cluster was available for testing a code without the SSE optimizations was used. The Infiniband cluster performed approximately 1.8 times better in the 32-bit case and approximately 1.6 times better in the 64-bit case as the Myrinet cluster.

The 1-dimensional (non-optimized) parallel even-odd preconditioned Dirac Operator Benchmark was modified to test whether an improvement using non-blocking MPI communication functions can be achieved. No effect was seen in using MPICH over GM provided by Myricom. Tests on a 4 node cluster which runs the MPI version of the ParaStation software results in an performance increase of about 20% (see tab. V).

Fig. 12 gives a summary of the communication behavior of the different architectures resulting from the parallel benchmarks. The efficiency number is the ratio between the performance of the same benchmark with and without communication. Included is also a test on a 4 node GigaBit-Ethernet cluster connected

Table V Parastation3 non-blocking I/O support (non-SSE).

MPI blocking I/O	MPI non-blocking I/O
308 Mflops	367 Mflops

via a non-blocking GigaBit-Ethernet switch. The efficiency number for GigaBit-Ethernet in the rightmost column of fig. is 12 even better than one can achieve using Myrinet on a system with a chip-set which provides a slow PCI-bus throughput. The positive influence of non-blocking communication support (ParaStation) and fast communication support (Infiniband) is shown by the efficiency numbers.

Compared to the single node benchmarks the results of the parallel benchmarks imply that that the capacity of bandwidth is crucial for the efficient use of PC clusters as a scalable platform for LQCD applications.

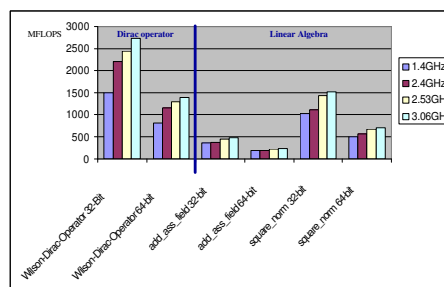


Figure 7: Single node Wilson-Dirac operator and linear algebra benchmark.

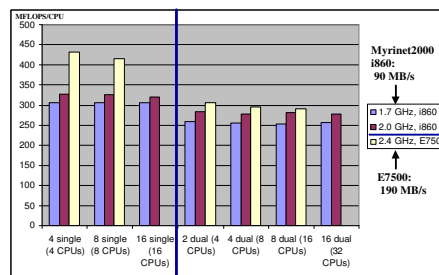


Figure 8: Parallel (1-dim) Wilson-Dirac Operator Benchmark (SSE), even-odd preconditioned, 2×16^3 lattice, Xeon CPUs, single CPU performance.

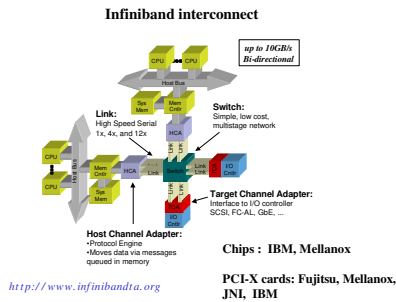


Figure 9: Infiniband interconnect.

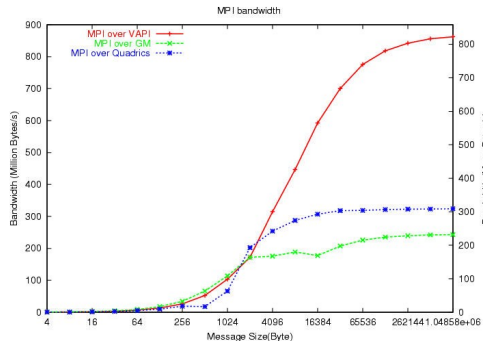


Figure 10: Infiniband bandwidth compared to Myrinet and QSNNet (source Mellanox).

5. Conclusions

We have discussed the usage of commodity clusters based on PCs for LQCD calculations. Such installations would not have been considered competitive a few years ago. However, our experience with such kind of machines at NIC and DESY adds further evidence that for problems in QCD that require below, say, 1 Tflops computer power, PC clusters are a valuable and cots-effective tool for computing physics results in LQCD. In Hamburg and Zeuthen clusters with 64 and 32 CPUs are successfully in operation for more than one year.

Investigations using representative benchmarks on the DESY clusters and also other architectures were carried out with promising results.

Applying the SSE/SSE2 (SIMD+pre-fetch) instructions on Pentium 4 like CPUs, the single node performance of the Wilson-Dirac operator is increasing according to the clock rate improvements of those CPUs used in commodity PCs.

The performance of memory bounded parts of the LQCD algorithms, especially the linear algebra routines, depends strongly on the throughput of the memory interfaces. Those interfaces did not show the same level of enhancements as it has been observed and will

Infiniband vs Myrinet performance (MFLOPS):

	XEON 1.7 GHz Myrinet, i860 chipset		XEON 2.2 GHz Infiniband, E7500 chipset	
	32-Bit	64-Bit	32-Bit	64-Bit
8x8 ³ lattice, 2x2 processor grid	370	281	697	477
16x16 ³ lattice, 2x4 processor grid	338	299	609	480

Figure 11: Infiniband and Myrinet performance comparisons using a parallel (2-dim) Wilson-Dirac Operator Benchmark on 4 node Pentium 4 Xeon clusters, single CPU performance, without SSE optimization, the local lattice size is 4²x8³ for the 8x8³, and 8x4x16² for the global 16x16³ lattice.

Maximal Efficiency of external I/O

	MFLOPs (without communication)	MFLOPs (with communication)	Maximal Bandwidth	Efficiency
Myrinet (i860), SSE	579	307	90 + 90	0.53
Myrinet/GM (E7500), SSE	631	432	190 + 190	0.68
Myrinet/Parastation (E7500), SSE	675	446	181 + 181	0.66
Myrinet/Parastation (E7500), non-blocking, non-SSE	406	368	hidden	0.91
Gigabit, Ethernet, non-SSE	390	228	100 + 100	0.58
Infiniband non-SSE	370	297	210 + 210	0.80

Figure 12: I/O Efficiency.

be expected further in the case of the CPUs.

The performance of the 1-dimensional and 2-dimensional parallel implementations of the Dirac-Wilson operator depends on the behavior of the external interconnects, i.e. is mainly dependent on the PCI-bus throughput coming given by the chip-set and the interface card itself. Results coming from PC clusters consisting of different components have shown an enhancement in both the quality of the chip-sets (e.g. E7500) and the throughput of the communication interfaces (e.g. Infiniband).

Non-blocking MPI communication can improve the performance by using adequate MPI implementations (e.g. ParaStation).

In summary, it might be envisaged, as done by e.g. LATFOR, that heterogeneous computer landscapes will be available to the user with centers that host machines in the multi-Tflops regime, still enabled by specialized machines, and many smaller installations at universities as well as research centers in the few hundred to 1 Tflops range realized by PC clusters.

Appendix: Discussion of the even-odd benchmark

The benchmark program applies the even-odd preconditioned Wilson-Dirac matrix that is defined on a $L^3 \times T$ lattice to a spinor-field. The program is implemented in C plus some in-lined SSE2 extensions. For parallelization we have used the MPI message passing library. The code is derived from Martin Lüscher's benchmark code presented at the lattice conference 2001 [38]. Even with communication switched off, the present code performs worse than the one of Martin Lüscher (579 Mflops vs. 880 Mflops on 1.7 GHz Pentium 4). The reason is twofold:

- Less variables that reside in the cache can be reused than in the standard case.
- We have skipped the cache-optimized order of the lattice-points to simplify the parallelization.

Strategy of the parallelization

The Wilson-Dirac matrix is a sparse matrix. The hopping part of the matrix only connects nearest neighbor sites on the lattice. Therefore, for parallelization, it is natural to divide the lattice in sub-blocks of size $t \times l_x \times l_y \times l_z$. Each of the MPI-processes takes one such sub-block. For simplicity we have done the parallelization only in one direction: $l_x = l_y = l_z = L$ and $tn_p = T$, where n_p is the number of processes.

The hopping part of the Wilson-Dirac matrix connects nearest neighbor sites on the lattice. Therefore each application of H_{eo} or H_{oe} (H_{oe} connects even with odd sites and H_{eo} vice versa.) the spinor-fields at the right boundary of the left neighbor and the spinor-fields at the left boundary of the right neighbor of each of the processes has to be sent and received.

In the case of even-odd pre-conditioning the spinor-field only resides on the even (or odd) sites. Therefore $L^3/2$ spinors have to be sent and received. A single data package has the size

$$24 \times 8 \times L^3/2 \text{ Byte} = 96 \times L^3 \text{ Byte} \quad (4)$$

In our blocking implementation, the communication of the data and the computation is performed in a consecutive way. First the spinor-fields are exchanged using the MPLSendrecv function. This is followed by the application of H_{eo} or H_{oe} on the single nodes.

Both times required for communication t_{comm} and calculation t_{calc} are measured separately. The effective bandwidth is computed as:

$$\text{Bandwidth} = \frac{96 \times L^3/2 \text{ Byte}}{t_{comm}} \quad (5)$$

The performance per node without communication is computed as:

$$P_0 = \frac{t \times L^3/2 \times 1392 \text{ flop}}{t_{calc}} \quad (6)$$

Correspondingly, the performance including the communication of the data is given by:

$$P = \frac{t \times L^3/2 \times 1392 \text{ flop}}{t_{comm} + t_{calc}} \quad (7)$$

In the case of the non-blocking communication, we had to divide H_{oe} (or H_{eo}) into a part that only acts on spinors that reside on the local lattice and a part that acts on spinors that reside on the neighbors:

- Initialize send and receive (MPI_Isend and MPI_Irecv),
- perform the calculation for the local part of H_{oe} ,
- Wait for the communication to finish (MPI_Wait),
- do the rest of H_{oe} .

Acknowledgments

The authors would like to thank Martin Lüscher (CERN) for the benchmark codes and the fruitful discussions about PCs for LQCD, and Isabel Campos Plasencia (Leibnitz-Rechenzentrum Munich), Gerhard Wellein (Uni Erlangen), Holger Müller (MEGWARE), Norbert Eicker (ParTec), Chris Eddington (Mellanox) for the opportunity to run the benchmarks on their clusters.

The authors also wish to thank the computer centers of DESY Hamburg and Zeuthen.

References

- [1] <http://www.lqcd.org/>.
- [2] <http://www.lrz-muenchen.de/services/compute/hlr/b/>.
- [3] <http://www.ibm.com/redbooks/>.
- [4] R. Alfieri et al. (apeNEXT-collaboration), hep-lat/0102011; R. Ammendola et al. (apeNEXT-collaboration), hep-lat/0211031.
- [5] D. Chen et al., hep-lat/0011004, P.A. Boyle et al., hep-lat/0110124, P.A. Boyle et al., hep-lat/0210034.
- [6] M. Hasenbusch, K. Jansen, T. Lippert, H. Stüben, P. Wegner, T. Wettig, and H. Wittig, *Evaluating Supercomputer platforms for lattice QCD applications*, in preparation.
- [7] <http://www.osti.gov/scidac/henp/index.htm>.
- [8] <http://www-zeuthen.desy.de/latfor/>.

- [9] K. Jansen, “Lattice Gauge Theory and High Performance Computing: The LATFOR initiative in Germany”, Talk given at CHEP03.
- [10] A. Gellrich, “A Commodity Cluster for Lattice QCD Calculations at DESY”, Talk given at CHEP03.
- [11] P. Wegner, “LQCD Benchmarks on Cluster Architectures”, Talk given at CHEP03.
- [12] H.J. Rothe, *Lattice Gauge Theories*, World Scientific Lecture Notes in Physics, Vol. 43, (World Scientific, Singapore, 1992); Montvay and G. Münster, *Quantum Fields on a Lattice*, Cambridge Univ. Press, 1994; J. Smit, *Introduction to Quantum Fields on a Lattice*, Cambridge Lecture Notes in Physics, Cambridge Univ. Press, 2002.
- [13] Y. Saad, *Iterative Methods for Sparse Linear Systems*, PWS Publishing Company, Boston, 1996.
- [14] G. Rossi, Nucl.Phys.Proc.Suppl. 53 (1997) 3; R. Sommer, Schladming Lectures 1997, hep-ph/9711243; M. Lüscher, Lectures given at Les Houches Summer School, 1997, hep-lat/9802029.
- [15] P.H. Ginsparg and K.G. Wilson, Phys. Rev. **D25** (1982) 2649; D.B. Kaplan, Phys.Lett. **B288** (1992) 342; P. Hasenfratz, Nucl.Phys. B (Proc.Suppl.) **63A-C** (1998) 53; P. Hasenfratz, V. Laliena, and F. Niedermayer, hep-lat/9801021; M. Lüscher, Phys.Lett. **B428** (1998) 342.
- [16] K. Jansen, Nucl.Phys.Proc.Suppl. 53 (1997) 127; M. Peardon, Nucl.Phys.Proc.Suppl. 106 (2002) 3, hep-lat/0201003.
- [17] S. Duane, A.D. Kennedy, B.J. Pendleton, and D. Roweth, Phys.Lett. **B195** (1987) 216.
- [18] R. Sommer and H. Wittig, physics/0204015.
- [19] M. Guagnelli, K. Jansen, F. Palombi, R. Petronzio, A. Shindler, and I. Wetzorke, hep-lat/0303012; K. Jansen, hep-lat/0010038.
- [20] M. Hasenbusch, Phys.Lett. **B519** (2001) 177; M. Hasenbusch and K. Jansen, hep-lat/0210036, hep-lat/0211042; K. Jansen and C. Urbach, in preparation; L. Giusti, Ch. Hoelbling, M. Lüscher, and H. Wittig, hep-lat/0212012.
- [21] G.F. Pfister, *In Search of Clusters*, Prentice Hall PTR, 2nd Edition, 1998.
- [22] Z. Fodor, S.D. Katz, and G. Papp, hep-lat/0202030.
- [23] <http://www.intel.com/>.
- [24] <http://www.myricom.com/>.
- [25] <http://www.megware.com/>.
- [26] <http://www.suse.de/>.
- [27] <http://secure.netroedge.com/lm78/>.
- [28] <http://www.pgroup.com/>.
- [29] <http://www.kai.com/>.
- [30] <http://www.intel.com/software/products/>.
- [31] <http://www.myri.com/>.
- [32] <http://dcache.desy.de/>.
- [33] <http://www.desy.de/~gellrich/clumon/>.
- [34] <http://people.ee.ethz.ch/~oetiker/webtools/mrtg/>.
- [35] D. Holmgren, private communications, CHEP03.
- [36] <http://www.par-tec.com/>.
- [37] <http://www.mellanox.com/>.
- [38] M. Lüscher, hep-lat/0110007, Nucl.Phys.Proc.Suppl. 106 (2002) 21.
- [39] M. Hasenbusch, <http://www.theorie.physik.uni-wuppertal.de/Cluster2002/Talks/hasenbusch.ps>.