

# Lattice Gauge Computing Summary

---

Don Holmgren  
Fermilab

# Outline

- n Communities / Collaborations
- n Existing Facilities
- n LQCD on Purpose-Built Machines
- n LQCD on Clusters

# Lattice Gauge Computing Requirements

- n Very large floating point requirements:
  - u Need 100's to 1000's of Gflop/sec-years for results
  - u Primary work is inverting large 4-D sparse matrices
  - u Sustained performance is often low (20% of peak)
  - u Multiprocessors are essential – clusters or massively parallel
- n Inter-processor communications requirements:
  - u High bandwidth (100's Mbyte/sec)
  - u Low latency -  $O(\mu\text{sec})$

# Communities / Collaborations

- n LATFOR (Lattice Forum) (K. Jansen)
  - u Forum of German lattice physicists, plus groups in Austria (Graz, Wien) and Switzerland (Bern)
  - u Universities: Berlin, Bielefeld, Darmstadt, Leipzig, Münster, Regensburg, Wuppertal
  - u Research Labs: DESY, GSI, NIC
- n Efforts:
  - u Coordinate physics program
  - u Share software
  - u Share configurations/propagators (ILDG = International Lattice Data Grid)

# Communities / Collaborations

## n LATFOR - Requirements

### u Typical application profile:

Lattice Size	$32^3 \cdot 64$
Memory	100 Gbyte
I/O request	0.1 Gbyte/Gflop
Runtime	15 Teraflops-years

### u Need **12.5 Teraflops** sustained

# Communities / Collaborations

## n SciDAC Lattice Gauge Computing

- u Majority of US community
- u Labs: Brookhaven, Fermilab, Jefferson Lab
- u Universities: Indiana, Utah, UCSB, UCSD, MIT, BU, ASU, Illinois, OSU, Columbia...
- u <http://www.lqcd.org/>

## n Efforts

- u QCDOC (Columbia/BNL)
- u Cluster Prototypes (Fermilab, JLAB)
- u Software (QMP, QLA, QIO, QDP) to allow applications to run on either type of hardware (Chulwoo Jung)
- u O(10 Tflops) in 2004 (QCDOC), 2005/6 (clusters)

# Facilities

## n SciDAC:

### u Jefferson Lab

Now: 128 nodes / 128 processors (2.0 GHz Xeon), Myrinet

Soon: 256 nodes / 256? Processors (Xeon), GigE mesh

### u Fermilab

Now: 176 nodes / 352 processors (2.0/2.4 GHz Xeon), Myrinet

Autumn: O(256) processor expansion (Xeon? Itanium2?)

Myrinet? GigE?

# Facilities

- n DESY Clusters:
  - u DESY Hamburg
    - 16 Dual 1.7 GHz Xeon
    - 16 Dual 2.0 GHz Xeon
    - Myrinet
  - u DESY Zeuthen
    - 16 Dual 1.7 GHz Xeon
    - Myrinet

# APEmille installations

n	Bielefeld	130 GF	(2 crates)
n	Zeuthen	520 GF	(8 crates)
n	Milan	130 GF	(2 crates)
n	Bari	65 GF	(1 crates)
n	Trento	65 GF	(1 crates)
n	Pisa	325 GF	(5 crates)
n	Rome 1	520 GF	(8 crates)
n	Rome 2	130 GF	(2 crates)
n	Orsay	16 GF	(1/4 crates)
n	Swansea	65 GF	(1 crates)
n	Gr. Total	~1966 GF	



# LQCD On Purpose-Built Hardware

## n QCDOC: (T. Wettig)

### u “QCD on a Chip” ASIC (0.18 $\mu$ ):

- ◀ Partner with IBM
- ◀ Based on PPC940 + 64-bit FPU, 1 Gflops peak
- ◀ 4 MB EDRAM + controller for external DDR
- ◀ 6 bidirectional LVDS channels, each 2 X 500 Mbps
- ◀ Ethernet
- ◀ ~ 5 Watts
- ◀ \$1/MF assuming 50% of peak

### u Packaging:

- ◀ 2 ASICs/daughterboard, 32 daughterboards per motherboard
- ◀ 8 motherboards per backplane (512 processors)

# QCDOC

## n Schedule:

- u Funded: 10 Tflops each in late 2003 for Edinburgh, RIKEN-BNL  
5 Tflops for Columbia (SciDAC) [pending]
- u 10 – 20 Tflops for US @ BNL in 2004
- u ASIC design is done, 1<sup>st</sup> chips in May

# QCDOC

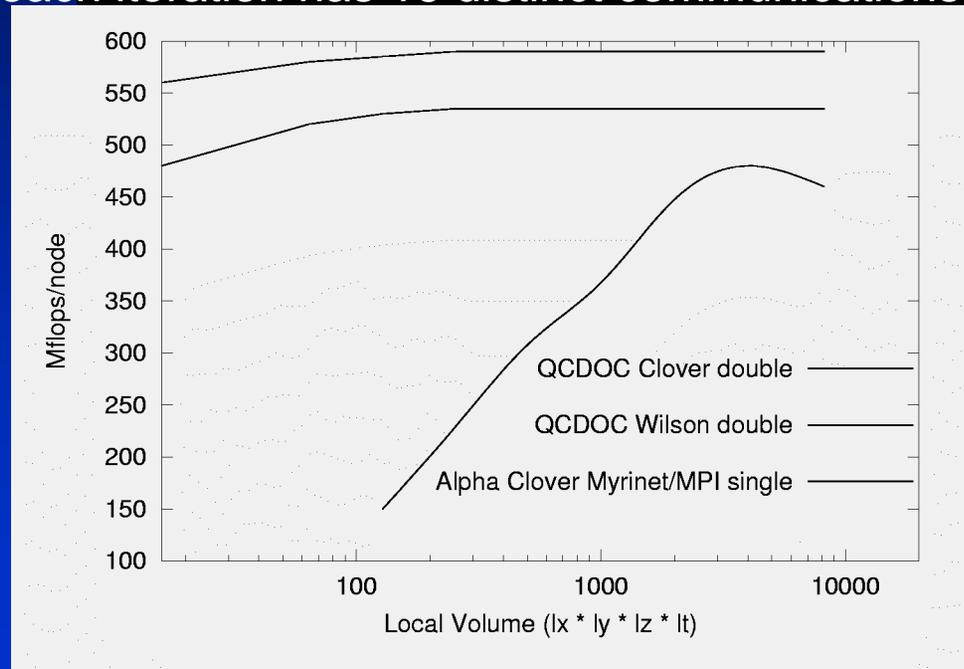
## Performance: (P. Boyle)

u Single node kernels, based on simulations:

- ◀ SU3 x SU3 800 Mflops/node
- ◀ SU3 x 2spinor 780 Mflops/node

u Multinode:

- ◀ Example – Wilson  $2^4$ , 470 Mflops/node, 22 $\mu$ sec/iteration, each iteration has 16 distinct communications



# QCDOC

- n Fast global sums via network passthru:

## Estimated Scalability

Based on cycle accurate simulation.

**Wilson** CG performance on a  $32^3 \times 64$  Lattice.

Clover and Domain Wall Fermions **even more scalable**.

Nodes	$M^\dagger M$	Gsum	Sust. Tflops
4096	2620 $\mu$ s	10 $\mu$ s	2.15
8192	1310 $\mu$ s	11.5 $\mu$ s	4.2
16384	680 $\mu$ s	13 $\mu$ s	8.1
32768	340 $\mu$ s	15 $\mu$ s	15.6

# APEnext (R. De Pietri)

- n In 0.18 $\mu$  ASIC, native implementation of “fused multiply-add”:  
 $a \times b + c$  (complex numbers)
- n 256 registers, each holding a complex number (double precision)
- n 200 MHz, ~ 5 watts, 1.6 Gflops/chip, \$1/Mflop @ 50% of peak
- n 6+1 channels of LVDS communications, each 200 Mbyte/sec
  - u 3D torus (x, y, z)
  - u Fast I/O link to UNIX host
- n Schedule:
  - u 300 to 600 ASIC's to be delivered in June `03
  - u From these, assemble and test 256-processor crate
  - u Mass production to start in Sept `03

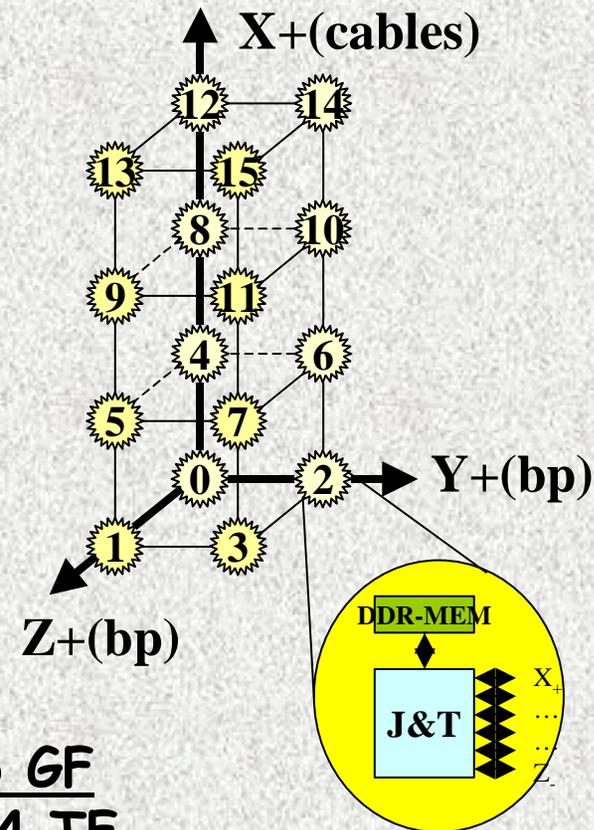
# The apeNEXT architecture (2)

n Two directions (Y,Z) on the backplane

n Direction X through front panel cables

n System topologies:

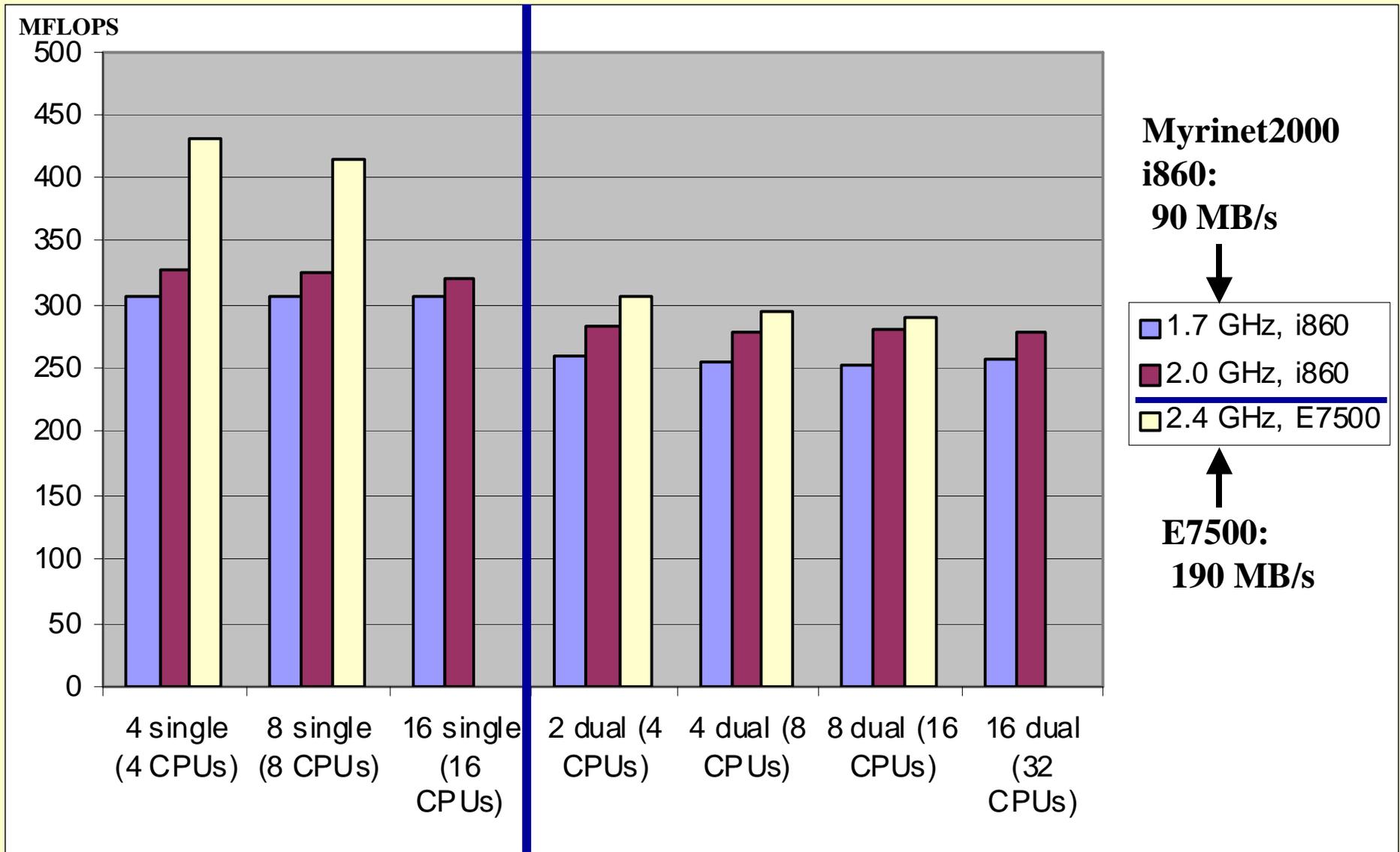
q	Processing Board	$4 \times 2 \times 2 \sim \underline{26 \text{ GF}}$
q	subCrate (16 PB)	$4 \times 8 \times 8 \sim \underline{0.4 \text{ TF}}$
q	Crate (32 PB)	$8 \times 8 \times 8 \sim \underline{0.8 \text{ TF}}$
q	Large systems	$(8*n) \times 8 \times 8$



# LQCD on Clusters

- n Performance (DESY: P. Wegner, Fermilab: D. Holmgren)
- n Common themes:
  - u SSE/SSE2 very important for optimizing performance
  - u Large memory bandwidth increase in Pentium 4 critical to LQCD
  - u Building clusters with Myrinet, investigating other networks
  - u Pay attention to PCI bus performance, varies with chipset

# Parallel (1-dim) Dirac Operator Benchmark (SSE), even-odd preconditioned, $2 \times 16^3$ , XEON CPUs, single CPU performance

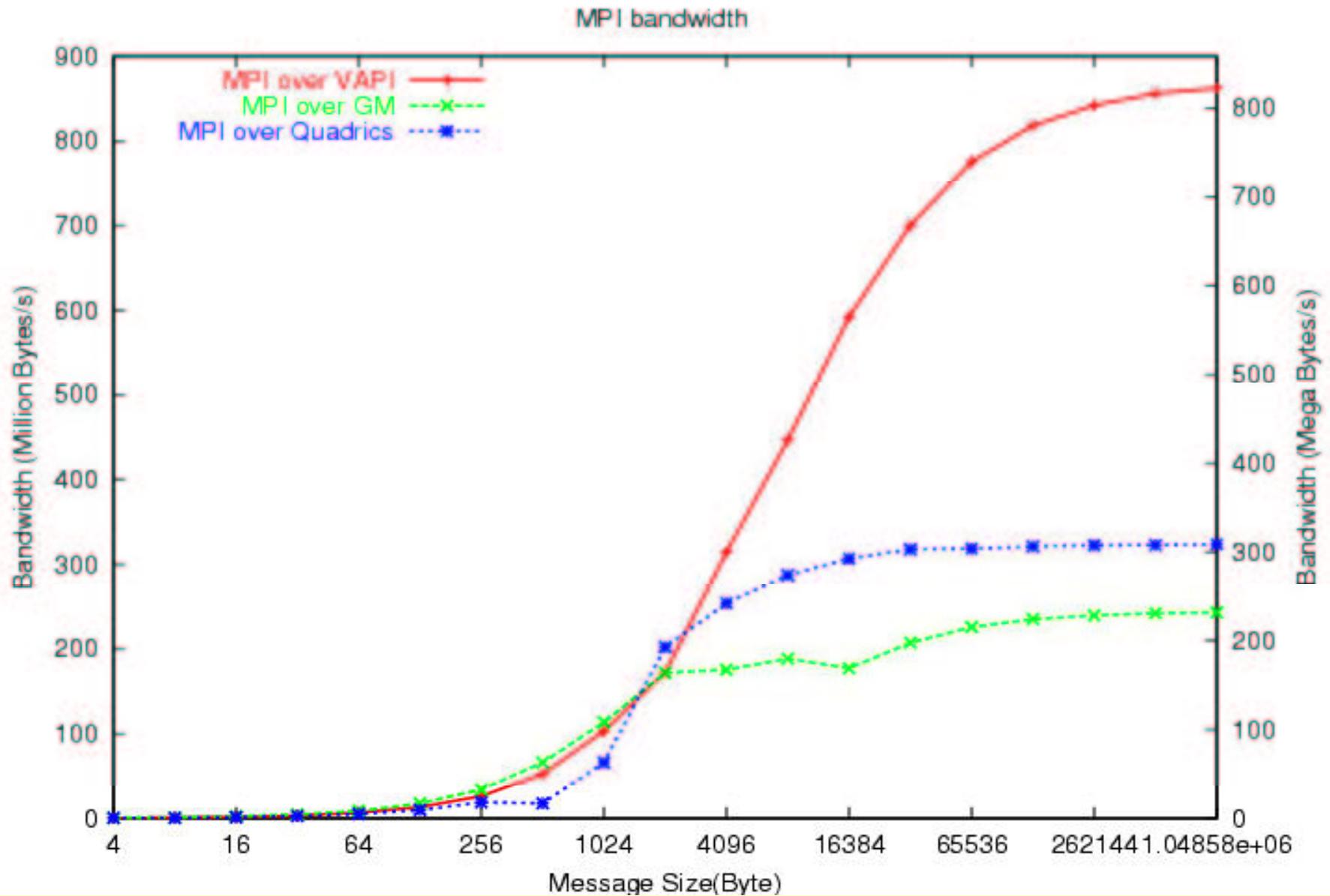


## Maximal Efficiency of external I/O



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

# Infiniband interconnect



# LQCD on Clusters - Management

- n See talks by [A. Gellrich](#), [A. Singh](#)
- n Common themes:
  - u Linux, PBS, Myrinet, MPI, private networks
  - u Web based monitoring and alarms
  - u MRTG history plots (e.g. temperatures, fans)
  - u Some Myrinet hardware reliability issues