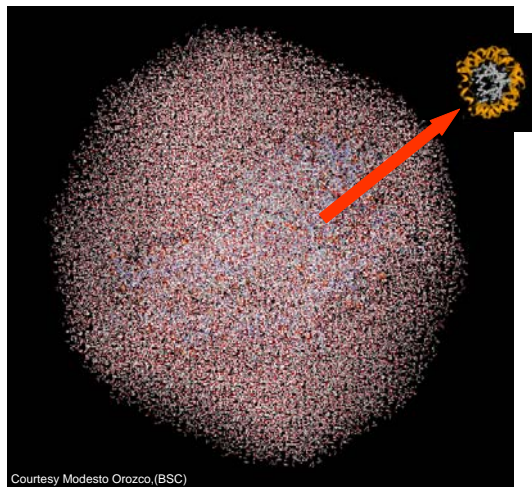# Advanced Profiling of GROMACS

**Jesus Labarta**
**Director Computer Sciences Research Dept.**
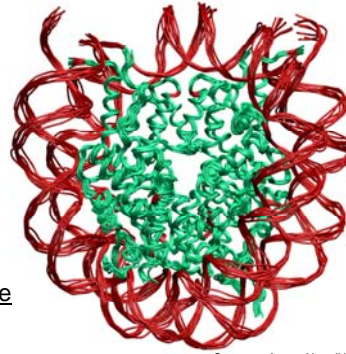**BSC**

---

## All I know about GROMACS

- A Molecular Dynamics application
- Heavily used @ BSC

- Not much ☹

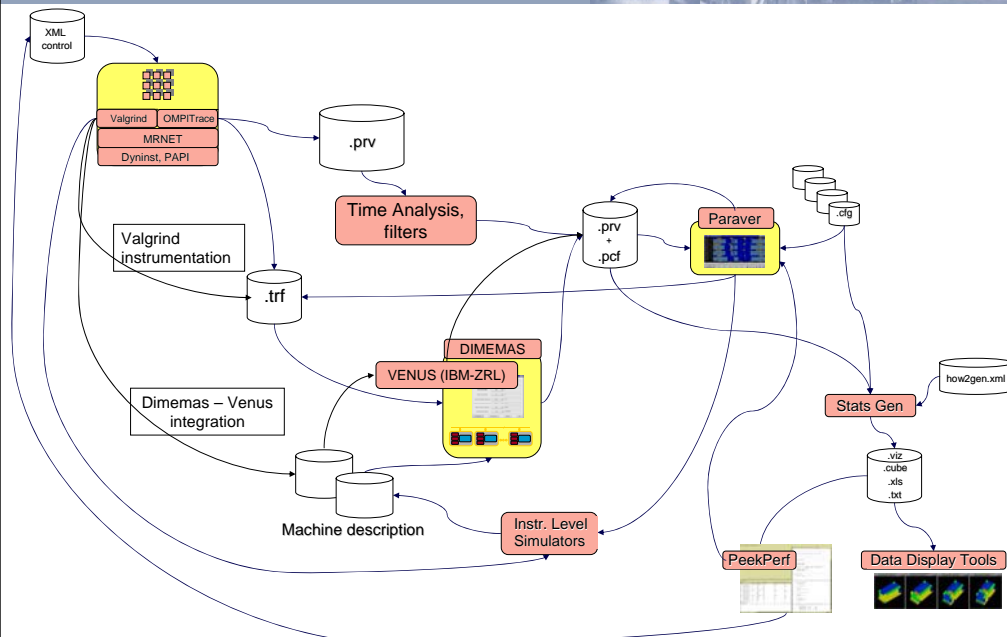Courtesy Modesto Orozco,(BSC)

1

# My interest

- Efficiency? Scalability?
- Can it be improved?

- Test case for our CEPBA-Tools environment
- Feedback/drive or research
  - Performance analysis tools
  - Programming models
  - Automatic Load balancing
  - Interconnects

Courtesy Agnes Noy (UB)

<u>Nucleosome test case</u>
145732 atoms
TIP3P water model
parmBSC0 force-field
Berendsen thermostat and barostat,
2 fs time step
constraints on Hydrogens

---

# Evaluation infrastructure



XML control

Valgrind | OMPITrace
MRNET
Dyninst, PAPI

.prv

Valgrind instrumentation

Time Analysis, filters

.prv + .pcf

Paraver

.cfg

.trf

DIMEMAS
VENUS (IBM-ZRL)

Dimemas – Venus integration

how2gen.xml

Stats Gen

Machine description

Instr. Level Simulators

.viz
.cube
.xls
.txt
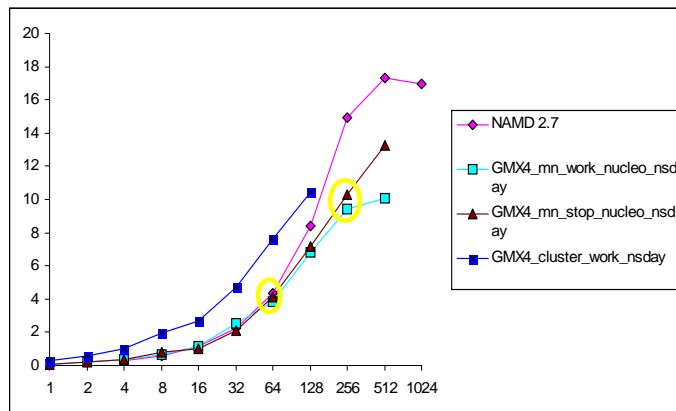
PeekPerf

Data Display Tools

# Index

- A story on
  - Performance
    - for (;;)
      - Parallel efficiency
      - Sequential performance (VMX,…)
    - Hybrid programming

---

# Performance

- Raw user metric (ns/day)

3

## Problems?

| | Procs | phase | eff |
|---|---|---|---|
| Measured | 64 | all | 0,55 |
| | 64 | FFTs | 0,55 |
| | 64 | Particles | 0,54 |
| | 256 | all | 0.32 |
| | 256 | Particles | 0.35 |
| | 256 | FFTs | 0.29 |

Not really good

Bad !!!

- Usual suspects
  - Lot of communication, large messages,…
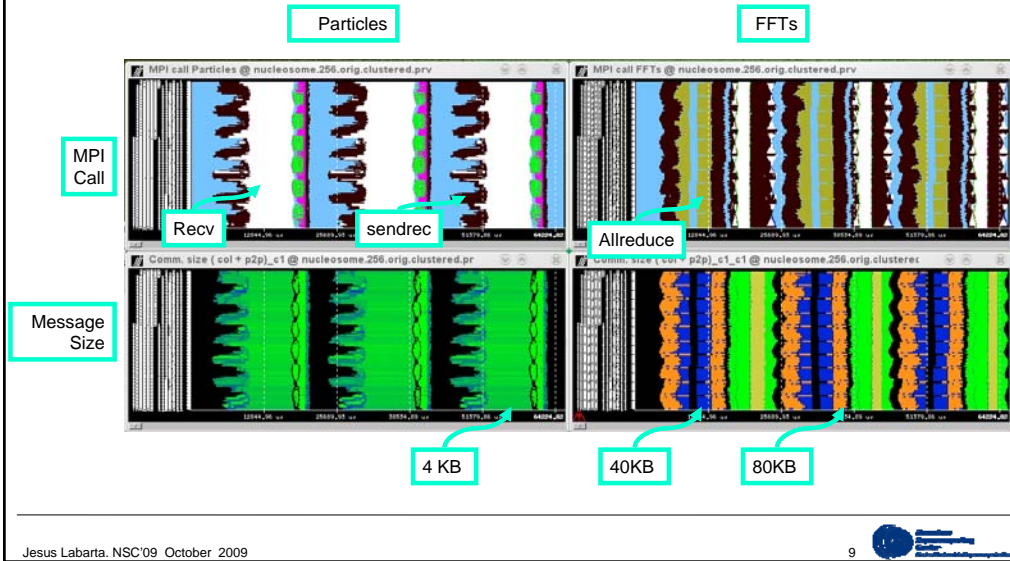
- VMX

---

# Parallel efficiency

4

Lots of communications, message sizes, …

- Suspect → convicted?

9

---

# Convicted

- Who?
  - Machine? Program?

- Average Bandwidth per core ?
  - 12 -15 MB/s

- Ideal interconnect

| | Procs | phase | eff |
|---|---|---|---|
| Measured | 64 | all | 0,55 |
| | 64 | FFTs | 0,55 |
| | 64 | Particles | 0,54 |
| | 256 | all | 0.32 |
| | 256 | Particles | 0.35 |
| | 256 | FFTs | 0.29 |
| Prediction ideal | 256 | all | 0.59 |
| | 256 | Particles | 0.63 |
| | 256 | FFTs | 0.52 |

Better !!!!!!

Not that good !!!!!!

Jesus Labarta. NSC'09  October  2009

10

5

# Sensitivity to interconnect parameters

| | Procs | phase | eff |
|---|---|---|---|
| Measured | 64 | all | 0,55 |
| | 64 | FFTs | 0,55 |
| | 64 | Particles | 0,54 |
| | 256 | all | 0.32 |
| | 256 | Particles | 0.35 |
| | 256 | FFTs | 0.29 |
| Prediction ideal | 256 | all | 0.59 |
| | 256 | Particles | 0.63 |
| | 256 | FFTs | 0.52 |



latency impact for different bandwidths (MB/s)



Bandwidth impact for different latencies (us)
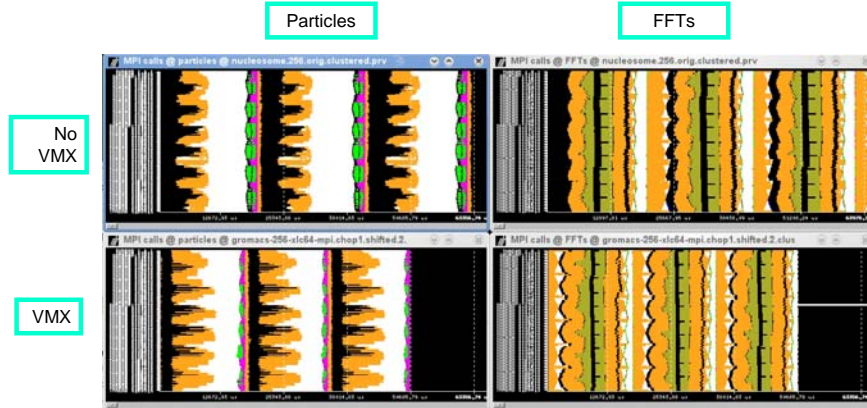
- Latency. Need to group messages?
- Sensitivity to network bandwidth?

---

# VMX

6

# VMX

- Not as much gain as expected ???

Particles

FFTs

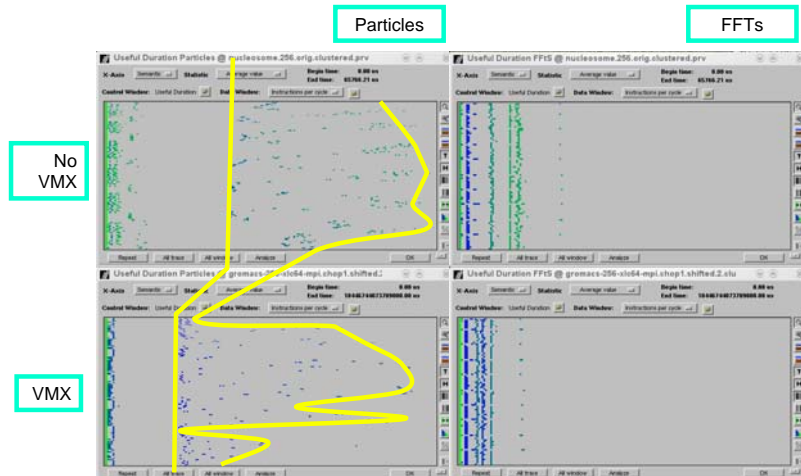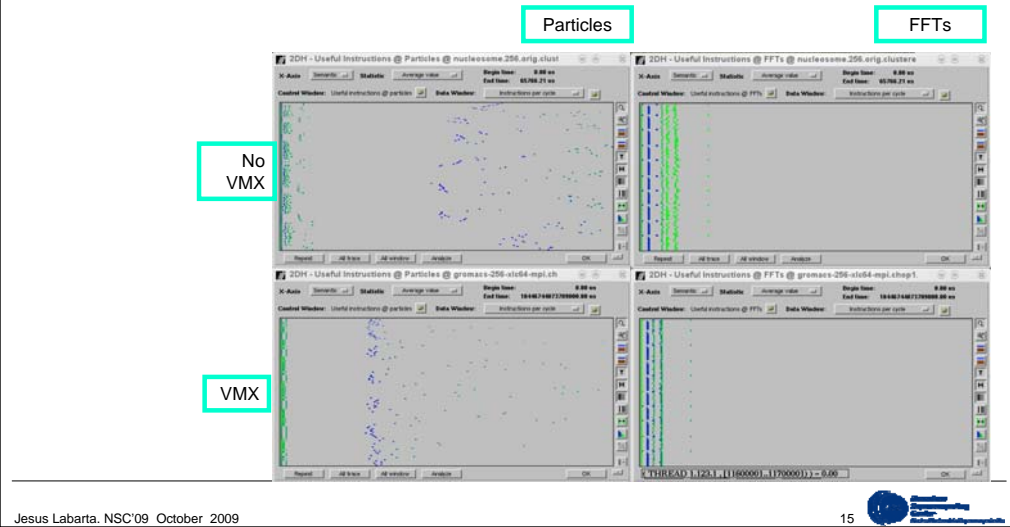No VMX

VMX

# VMX

- Histogran of duration: Not uniform impact !!!!
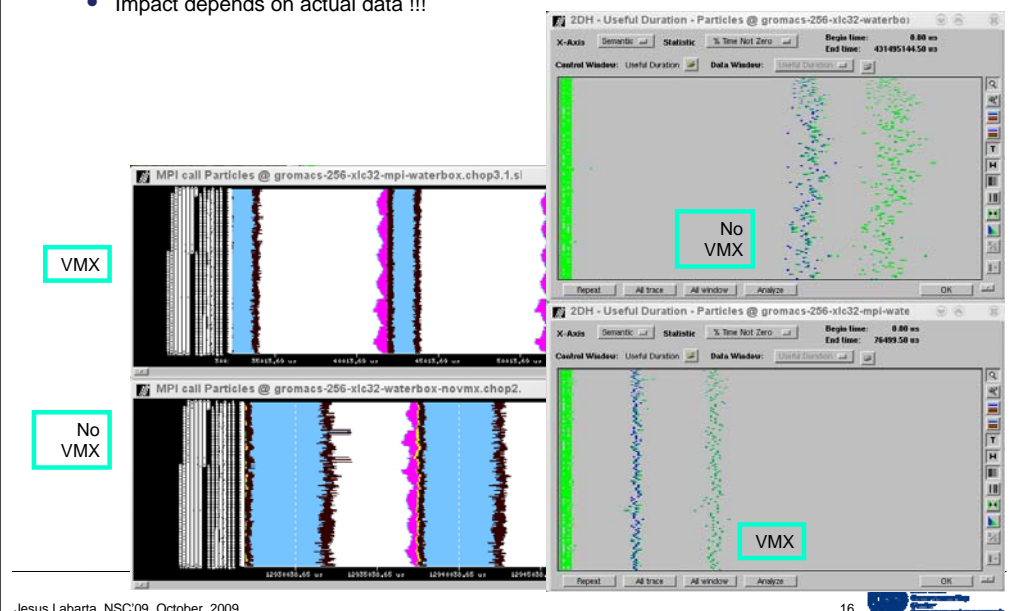  - Across Particle processes
  - Between FFTs and Particles FFTs

Particles

FFTs

No VMX

VMX

7

# VMX

- Histograms of Instructions

Particles    FFTs

No VMX

VMX

---

# VMX @ waterbox

- Impact depends on actual data !!!

No VMX

VMX

No VMX

VMX

# Back to parallel performance

---

# Performance factors

$$Parallel\_\eta = LB * \overbrace{microLB * Transfer}^{TotalLB}$$

$$\eta = Parallel\_\eta * relIPC * Comp$$

$$\underbrace{\phantom{LB * microLB * Transfer}}_{Comm}$$

**Performance Factors**



Legend:
- Parallel_Eff
- LB
- Comm
- rel IPC
- Comp
- uLB
- Transfer
- Total LB
- Efficiency

9

## Structure: A few colors



MPI call FFTs

recv
sendrec
allreduce

User functions @ FFTSs

make_bsplines
gather_f_bsplines
solve_pme

MPI call Particles

recv
sendrec

User functions @ Particles

ewald_LRcorrection
do_nonbonded

## Scaling structure



User function @ Particles      User function @ FFTSs

32

64

128

256

10

# PRACE Test case

Useful Duration

MPI calls
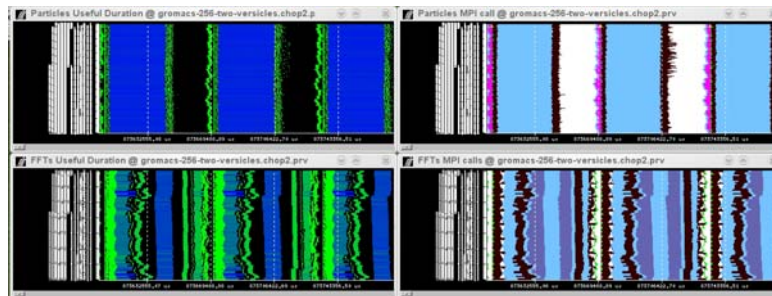
Particles

FFTs

- Average point to point bandwidth: 13 MB/s
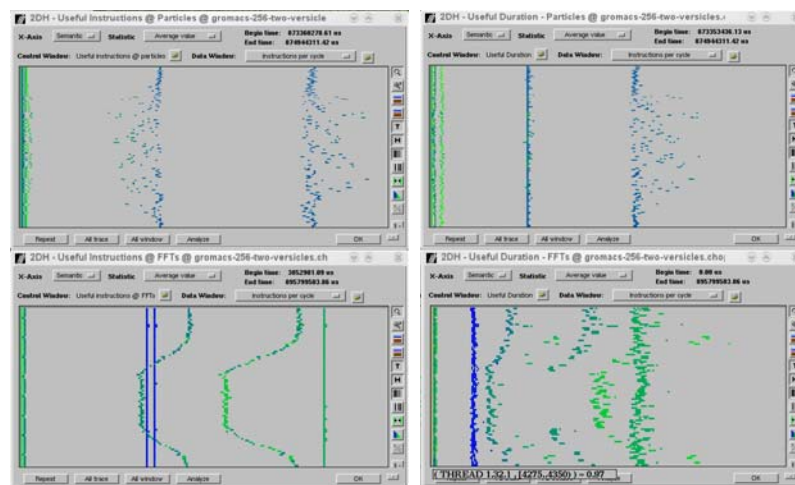
# PRACE Test case

- Effects → opportunities

Instructions

Duration

Particles

FFTs

11

# Back to sequential performance
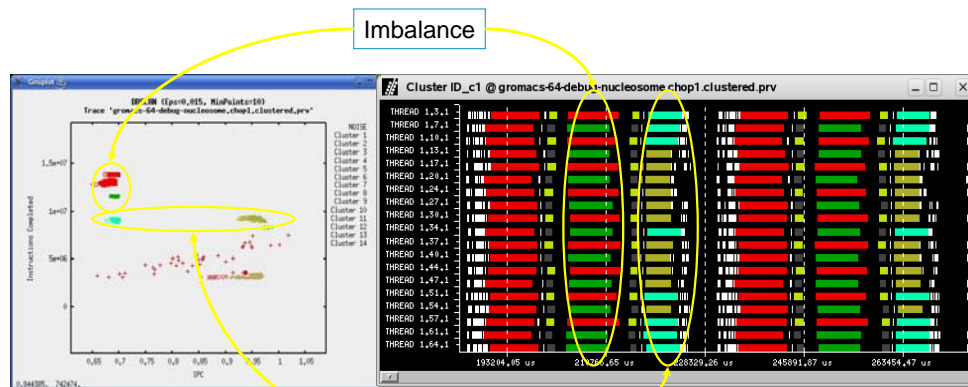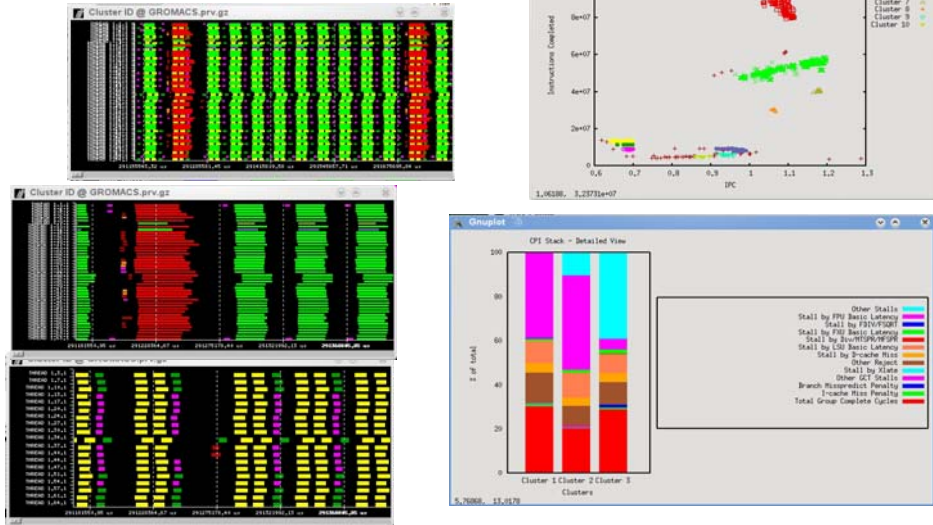
## 64 processes: FFTs
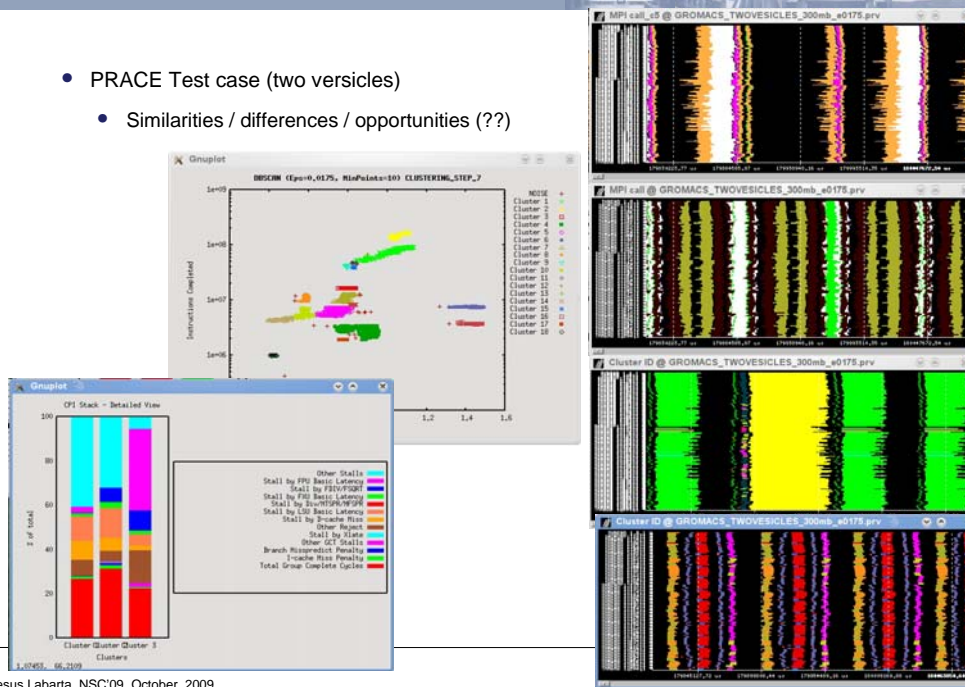
- Globally balanced but some local imbalance

Imbalance



Imbalance

# Online analysis

- Nucleosome

# Online analysis

- PRACE Test case (two versicles)
  - Similarities / differences / opportunities (??)

## Detailed metrics

- Application characterization

| Metric Description | Cluster 1 | Cluster 2 | Cluster 3 | Cluster 4 |
|---|---|---|---|---|
| **FMA ops per floating point instruction**<br>PM_FPU_FMA / PM_FPU0_FIN + PM_FPU1_FIN | 0,681 | 0,768 | 0,700 | 0,436 |
| **Instructions per Load/Store**<br>PM_INST_CMPL / (PM_LD_REF_L1 + PM_ST_REF_L1) | 1,470 | 1,465 | 1,659 | 1,306 |
| **HW floating point instructions (flips)**<br>PM_FPU0_FIN + PM_FPU1_FIN | 1.174.346 | 2.896.665 | 3.549.551 | 1.465.935 |
| **Total floating point operations (flops)**<br>PM_FPU0_FIN + PM_FPU1_FIN + PM_FPU_FMA - PM_FPU_STF | 1.436.155 | 3.979.120 | 6.459.096 | 3.391.547 |
| **Total FP Load&Store operations (fp_tot_ls)**<br>PM_LSU_LDF + PM_FPU_STF | 1.348.179 | 3.006.794 | 3.656.377 | 2.014.815 |
| **FMA %**<br>2 * PM_FPU_FMA / flops | 36,51% | 44,12% | 37,19% | 81,21% |
| **Memory Mix**<br>(PM_LD_REF_L1+ PM_ST_REF_L1) / PM_INST_DISP | 32,61% | 24,70% | 30,26% | 25,20% |
| **Load Mix**<br>PM_LD_REF_L1 / PM_INST_DISP | 25,83% | 19,39% | 24,93% | 24,76% |
| **Store Mix**<br>PM_ST_REF_L1 / PM_INST_DISP | 6,77% | 5,32% | 5,32% | 0,44% |
| **FPU Mix**<br>flips / PM_INST_DISP | 1,06% | 8,22% | 2,00% | 5,32% |
| **FXU Mix**<br>PM_FXU_FIN / PM_INST_DISP | 20,85% | 16,19% | 24,38% | 16,07% |

## Detailed metrics

- Behavior on architecture

| Metric Description | Cluster 1 | Cluster 2 | Cluster 3 | Cluster 4 |
|---|---|---|---|---|
| **% Instr. Completed**<br>PM_INST_CMPL / PM_INST_DISP | 47,93% | 36,20% | 50,21% | 32,90% |
| **L1 misses per Kinstr.**<br>(PM_LD_MISS_L1 + PM_ST_MISS_L1 / PM_INST_CMPL) * 1000 | 10,18 | 37,04 | 10,38 | 63,81 |
| **L2 misses per Kinstr.**<br>(PM_DATA_FROM_MEM / PM_INST_CMPL) * 1000 | 0,093 | 0,103 | 0,057 | 1,008 |
| **Bytes from maim memory per floating point instruction finished**<br>(PM_DATA_FROM_MEM * mem_line_size) / PM_FPU_FIN | 0,084 | 0,090 | 0,068 | 9,712 |
| **Number of Loads per Load miss**<br>PM_LD_REF_L1 / PM_LD_MISS_L1 | 121,81 | 74,99 | 102,29 | 11,84 |
| **Number of Stores per Store miss**<br>PM_ST_REF_L1 / PM_ST_MISS_L1 | 24,53 | 4,91 | 19,19 | 50,31 |
| **Number of Loads&Stores per L1 miss**<br>(PM_LD_REF_L1 + PM_ST_REF_L1) / (PM_LD_MISS_L1 + PM_ST_MISS_L1) | 66,80 | 18,42 | 58,05 | 12,00 |
| **L1 cache hit rate**<br>1 - (PM_LD_MISS_L1 + PM_ST_MISS_L1) / (PM_LD_REF_L1 + PM_ST_REF_L1) | 98,50% | 94,57% | 98,28% | 91,67% |
| **Number of Loads per (D)TLB miss**<br>PM_LD_REF_L1 / PM_DTLB_MISS | 149.034,25 | 112.018,98 | 128.501,34 | 49.827,87 |
| **Number of Loads&Stores per (D)TLB miss**<br>(PM_LD_REF_L1 + PM_ST_REF_L1) / PM_DTLB_MISS | 188.100,04 | 142.754,08 | 155.941,21 | 50.707,36 |
| **% TLB misses per cycle**<br>PM_DTLB_MISS / PM_CYC | 0,000% | 0,000% | 0,000% | 0,001% |
| **Total Loads from local L2 (M) (total_ld_l_L2)**<br>PM_DATA_FROM_L2 / 1024*1024 | 0,173 | 0,060 | #VALUE! | 0,381 |
| **Local L2 load traffic (MB)**<br>L1_cache_line_size * total_ld_l_L2 | 22,126 | 7,635 | #VALUE! | 48,781 |

14

## Detailed metrics

- Real Performance

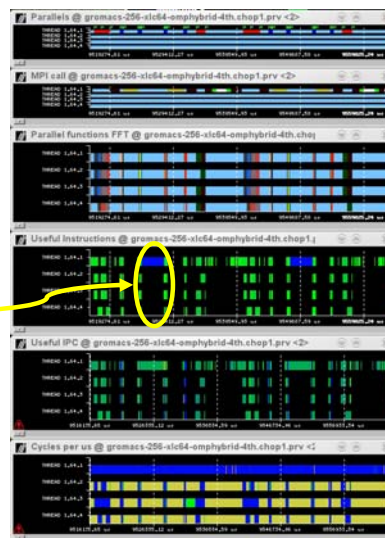| Metric Description | Cluster 1 | Cluster 2 | Cluster 3 | Cluster 4 |
|---|---|---|---|---|
| **% Duration** | 0.60741 | 0.20625 | 0.09719 | 0.04445 |
| **Avg. Burst Duration (µs)** | 21.260,58 | 8.243,17 | 36.604,90 | 5.790,58 |
| **Total preemted time (µs)** <br> Avg. Burst Duration – (PM_CYC * processor_cycle_time) | 366,81 | 123,88 | 576,02 | 75,37 |
| **% preempted time** <br> Preempted Time / Avg. Burst Time | 1,725% | 1,503% | 1,574% | 1,302% |
| **IPC** <br> PM_INST_CMPL / PM_CYC | 1,11 | 0,68 | 1,08 | 0,69 |
| **CPI** <br> 1/IPC | 0,90 | 1,46 | 0,93 | 1,45 |
| **MIPs** <br> PM_INST_CMPL / Avg. Burst Duration | 2496,91 | 1547,98 | 2431,90 | 1566,34 |
| **Mem.BW (MB/s)** <br> (PM_DATA_FROM_MEM * mem_line_size ) / Avg. Burst Duration | 29,74 | 20,43 | 17,60 | 202,08 |
| **Memory instructions per second** <br> (PM_LD_REF_L1 + PM_ST_REF_L1) / Avg. Burst Duration | 1698,69 | 1056,39 | 1465,48 | 1199,69 |
| **HW floating point instructions per cycle** <br> flips / PM_CYC | 0,024 | 0,155 | 0,043 | 0,112 |
| **Flop rate (MFLOPs)** <br> flops / Avg. Burst Duration | 67,55 | 482,72 | 176,45 | 585,70 |
| **HW floating point instructions rate** <br> flips / Avg. Burst Duration | 55,24 | 351,40 | 96,97 | 253,16 |
| **Computation intensity** <br> flops / fp_tot_ls | 1,065 | 1,323 | 1,767 | 1,683 |
| **Local L2 load bandwidth per processor (MB/s)** <br> L1_cache_line_size * total_ld_l_L2 / Avg. Burst Duration | 1.040,72 | 926,19 | #VALUE! | 8.424,28 |
| **% Loads from local L2 per cycle** <br> PM_DATA_FROM_L2 / PM_CYC | 0,378% | 0,335% | #VALUE! | 3,043% |

---

# Hybrid MPI + OpenMP

# Hybrid parallelization

- OpenMP parallelization
  - Particles by Christian Simarro (BSC)
  - FFTs by Sebastian von Alfthan (CSC)

- 256 MPI x 4 openMP

130 procs!!!!

# Hybrid parallelization

- FFTs
  - Partial parallelization
  - Too synchronous
    - MPI – OpenMP –MPI – OpenMP ....
  - OS scheduling issues

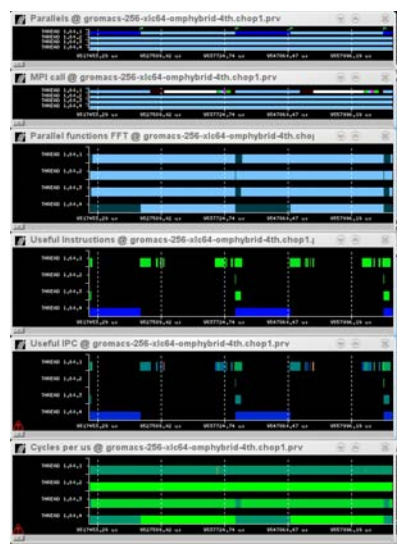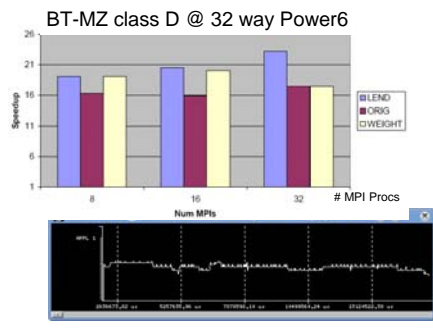Not Parallelized

# Hybrid parallelization

- Particles

```
for(n=n0; (n<n1);
#pragma omp pa...  f private(tabledata, nrnb_ind, nlist, outeriter,
                  inneriter, tabletype, kernelptr) shared(fr,nblists)
  for(i=i...   1); i++) {
    …
    if (something)
       compute();
}
```

Right place?

Very imbalanced

Few iterations
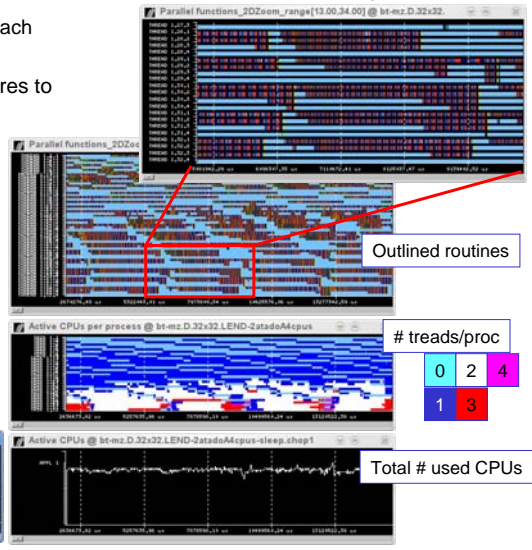
Big difference across processes

---

# Dynamic Load Balancing run time

- MPI + OpenMP  (StarSs)
  - Job is started with P MPI processes each with 1 OpenMP thread
- Processes blocking at MPI calls lend cores to other processes
- GADGET @ 800 cores: Sup=2.5

BT-MZ class D @ 32 way Power6

BT-MZ class D @ 32 way Power6



Outlined routines

# treads/proc

| 0 | 2 | 4 |
|---|---|---|
| 1 | 3 |   |

Total # used CPUs

# Conclusion

- Applications: Nothing IS, everything CAN BE
  - Can always be different, input data dependent
  - Important to get insight into actual behavior.
  - Need dynamic run time optimization

- GROMACS: an interesting case
  - Really heterogeneous application
  - Serialization and load balancing issues
  - A target for MPI + SMPSs
    - Automatic Load Balance
    - Overlap communication/computation

- The need for cooperative work

18