



Intel Developer
FORUM



Hybrid Parallel Performance- Tuning for Multicore-based HPC Clusters

Wanqing He

SSG CRT\HPC Enabling Manager

SSGS004

Intel Developer
FORUM

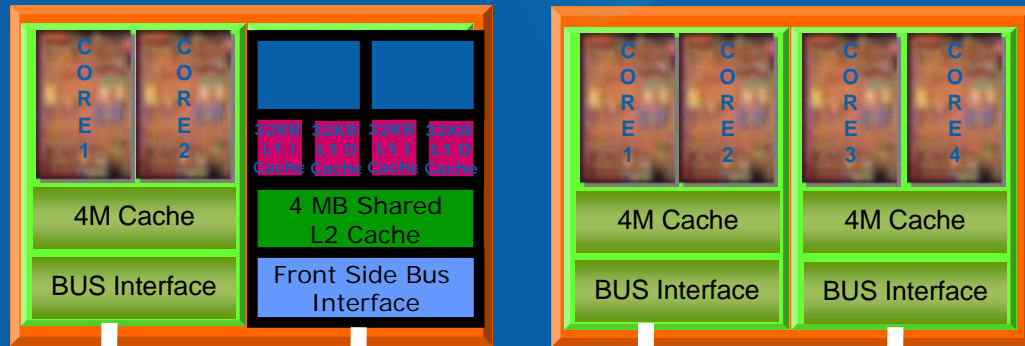
Disclaimer

- INFORMATION IN THIS DOCUMENT IS PROVIDED IN CONNECTION WITH INTEL® PRODUCTS. EXCEPT AS PROVIDED IN INTEL'S TERMS AND CONDITIONS OF SALE FOR SUCH PRODUCTS, INTEL ASSUMES NO LIABILITY WHATSOEVER, AND INTEL DISCLAIMS ANY EXPRESS OR IMPLIED WARRANTY RELATING TO SALE AND/OR USE OF INTEL PRODUCTS, INCLUDING LIABILITY OR WARRANTIES RELATING TO FITNESS FOR A PARTICULAR PURPOSE, MERCHANTABILITY, OR INFRINGEMENT OF ANY PATENT, COPYRIGHT, OR OTHER INTELLECTUAL PROPERTY RIGHT.
- Intel may make changes to specifications, product descriptions, and plans at any time, without notice.
- All dates provided are subject to change without notice.
- Intel and Itanium are trademarks or registered trademarks of Intel Corporation or its subsidiaries in the United States and other countries (regions).
- *Other names and brands may be claimed as the property of others.
- Copyright © 2007, Intel Corporation. All rights are protected.

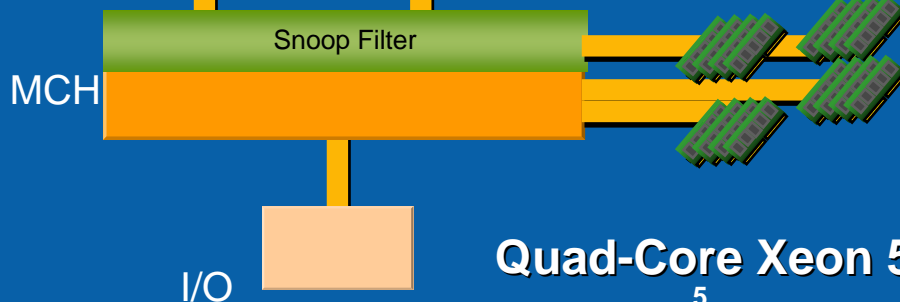
Agenda

- Define the hybrid parallel performance Tuning model
- Apply hybrid tuning to HPL: understand the threads dithering by ITAC
- Tune both MPI and OpenMP: apply hybrid tuning to real-world CFD application
- Try Cluster-OpenMP to fast deploy your Cluster app.
- Summary/Call to action

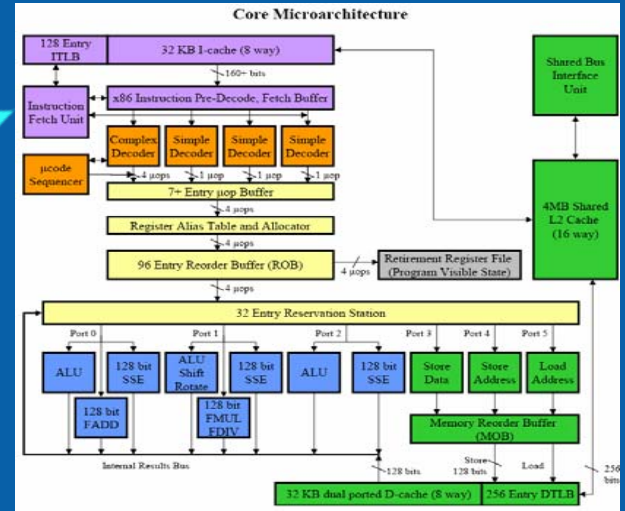
Intel's Xeon 5300 build SMP cluster base on Core's performance/watt



1066/1333 MT/s

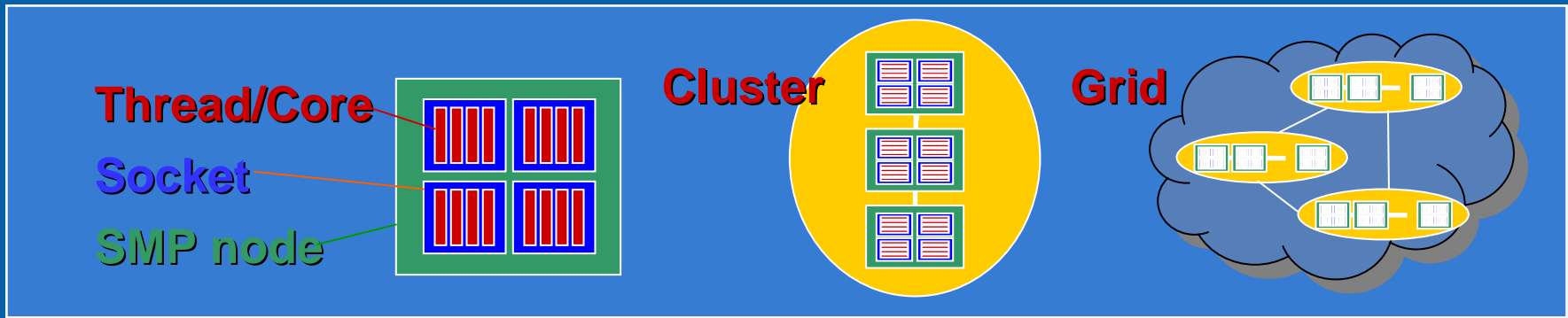


Quad-Core Xeon 5300



Hybrid parallel tuning:
Defined as MPI tuning between SMP nodes mixed with multithreading tuning inside SMP nodes to get both distributed mem. parallelism and share mem. parallelism.

Multicore building block for HPC



HPC System	Share memory/node	Distributed Memory /Hybrid Cluster	Compute Grid /
Program Model	Multithreading in share mem: OpenMP, pthread, winthread	MPI + OpenMP In distributed mem: Intel MPI, ClusterOpenMP	SOA, Virtualization
Develop Tools	<i>Intel® Thread Checker</i> <i>Intel® Thread Profiler</i> <i>VTune®</i> <i>Intel® MKL</i>	<i>Intel® MPI</i> <i>Intel® Trace Collector/Trace Analyser</i> <i>Intel® Cluster MKL</i> <i>Intel Message Checker</i> <i>Intel Threading tools</i>	<i>UNICORE, Globus*, etc.</i> <i>DRMAA*</i> <i>Grid Programming Environment (GPE)</i>

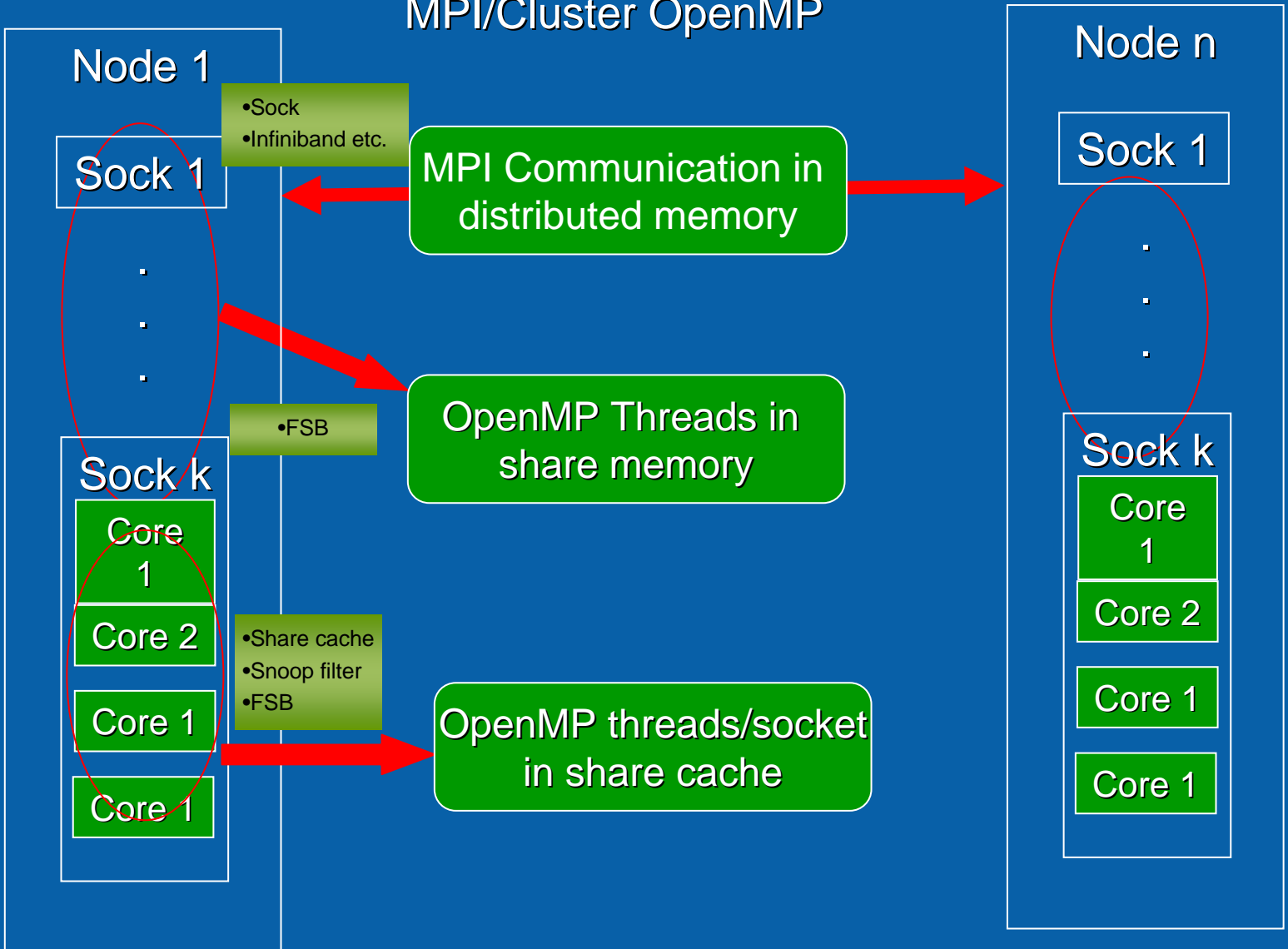
Multicore-SMPs Cluster Model

MPI/Cluster OpenMP

Level 1

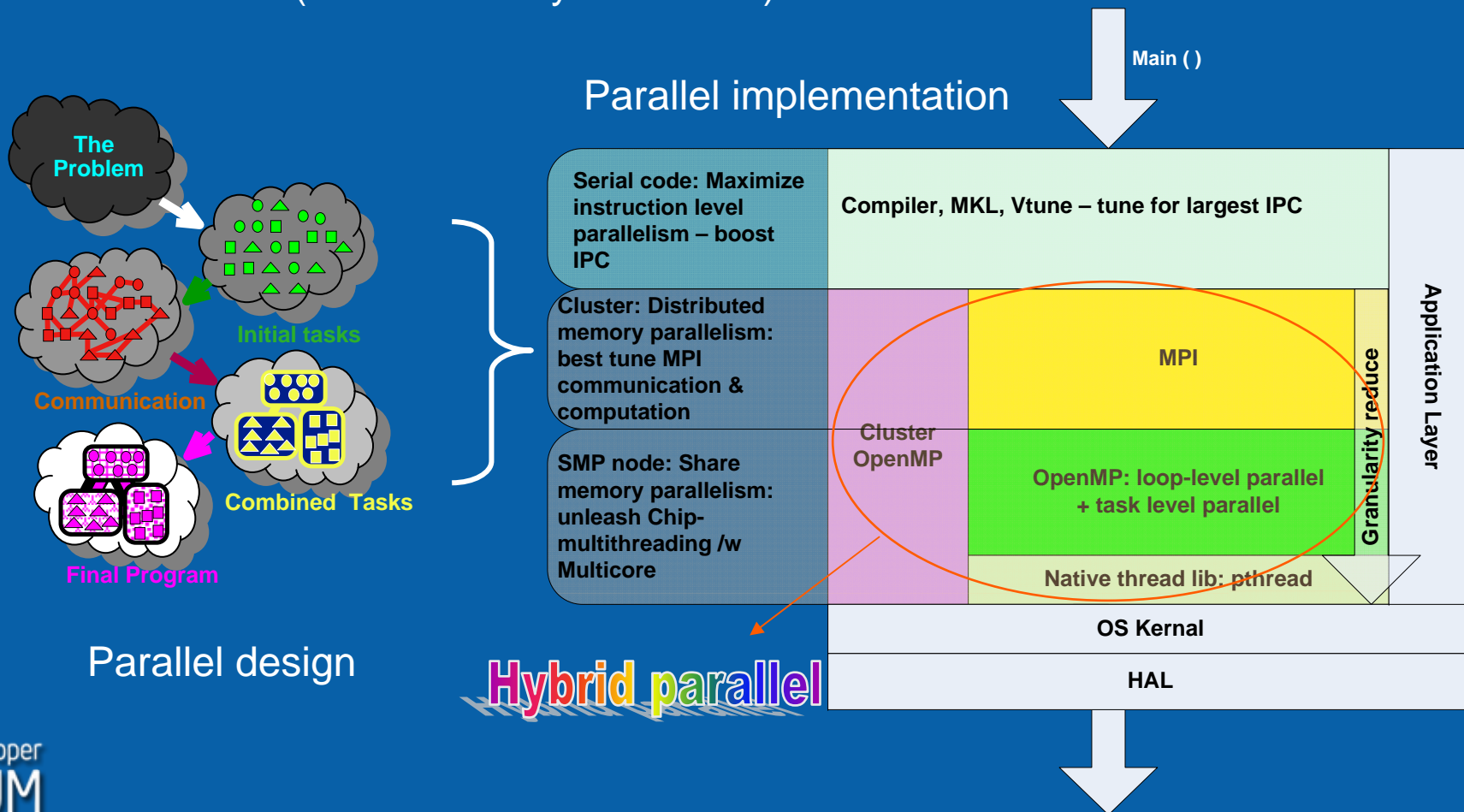
Level 2

Level 3



SMP Cluster parallel Model

Hybrid parallel should be implemented from coarse to fine granularity from distributed memory (MPI), share-memory (OpenMP/threading) and instruction level (serial code by IPC boost).



SMP Cluster Performance model

- Stem from Amdahl law:

Let T_1 be the execution time in serial mode; If p is the fraction of the code that is parallel, then:

$$T_N = (1-p) * T_1 + (p/N) * T_1$$

$$\text{Speedup} = \frac{1.0}{(1-p) + (p/N)}$$

- As we scale the problem for large N , the serial work remains fixed, while the parallelizable work grows with N .

$$p \rightarrow 1, \text{ Speedup} \rightarrow N$$

- *It is possible to realize linear speedup, if we also grow the problem size (Gustafson's law)*
- In hybrid mode, $N \approx N_{\text{MPI}} * N_{\text{thread}} * \text{TR}\%$ where N_{MPI} is the number of MPI processes and N_{thread} is the number of parallel threads running inside one MPI process, and $\text{TR}\%$ is the threading ratio that represent the percentage of Multi-threaded work in one process

Real world Cluster perf. Model : MPI

- Considering overhead time from MPI, we have:

$$T_N = (1-p) * T_1 + (p/N) * T_1 + T_{mpi}$$

When considering MPI Model, the communication effects limit the absolute performance and scalability:

- Latency or message initiation time (T_L)
- Bandwidth, a measure of transfer rate once it is initiated (B_w)
- Collective communication calls
- Synchronization, load imbalance

Let M = average message size

- T_c = collective communication time
- T_i = Load imbalance or synch time

$$T_{mpi} = A_1 * T_L + A_2 * M / B_w + A_3 * T_c + T_i$$

Real world Cluster perf. Model : hybrid

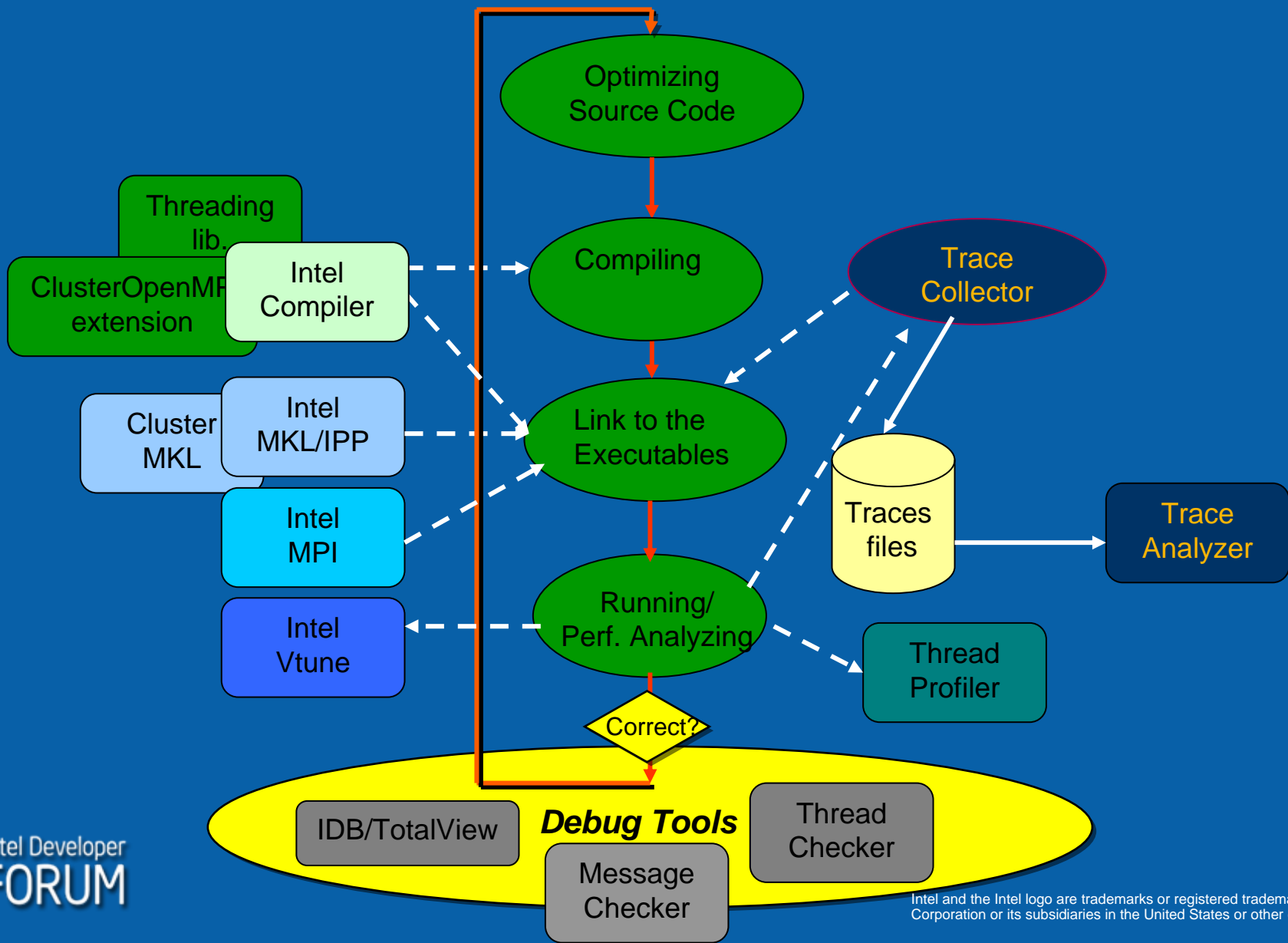
- When integrating Multithreaded code session into the MPI model, we got

$$T_N = (1-p) * T_1 + (p / (N'_{MPI} * N_{thread} * TR\%)) * T_1 + T'_{pi} + \sum T_{th}$$

The MPI communication time changed to T_{mpi} . Usually $T_{mpi} < T_{mpi}$ due to different partition. The parallelized work are roughly scaled out to $N'_{MPI} * N_{thread} * TR\%$ fork-join jobs. Here N'_{MPI} is new MPI processes number. The $\sum T_{th}$ stands for all overhead that threads bring into the system, which may including but not limited to: thread life cycle overhead, FSB latency/bandwidth overhead, share cache coherence overhead, Snoop filter overhead, threads sync. time etc.

- To be simplified, when $T_{mpi} + \sum T_{th} < T_{mpi}$, We can get overall benefit by introducing Multithread into SMP node. Or we have to fine-tune OpenMP to get more parallelism

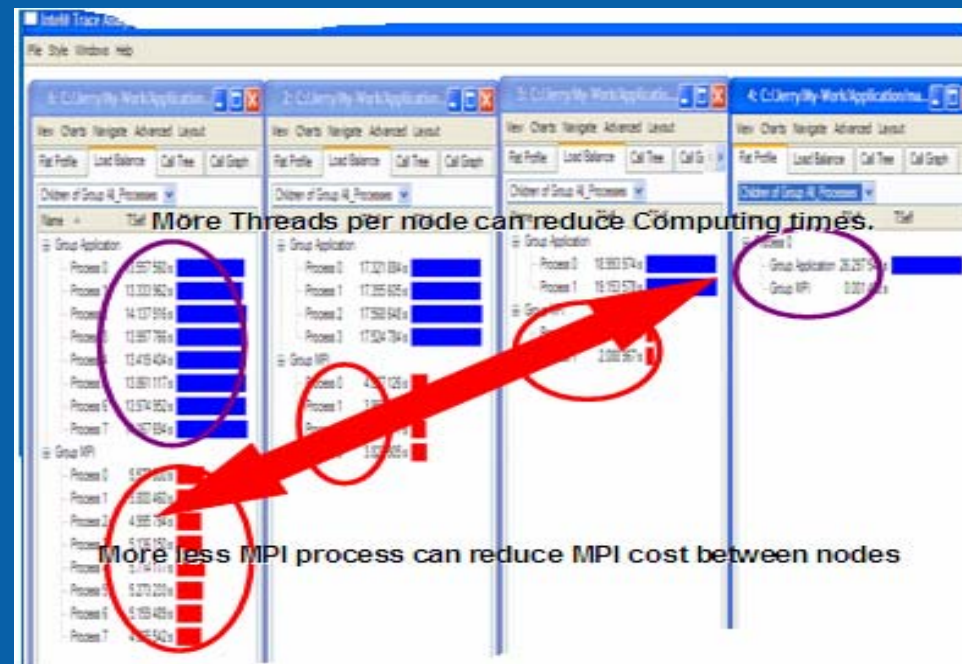
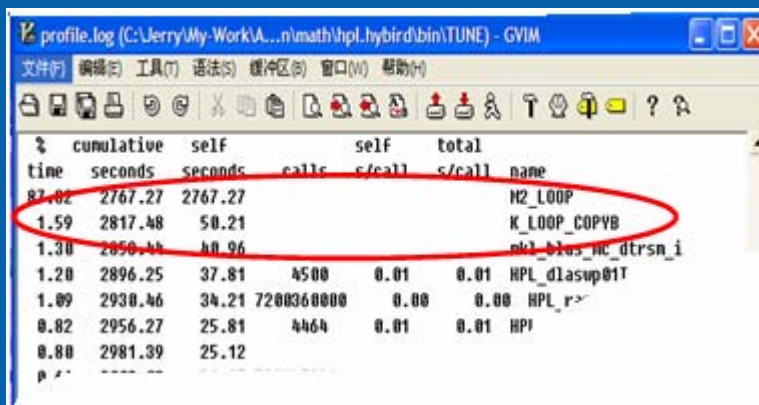
Hybrid programming tools: Compiler/ITAC/ITPT



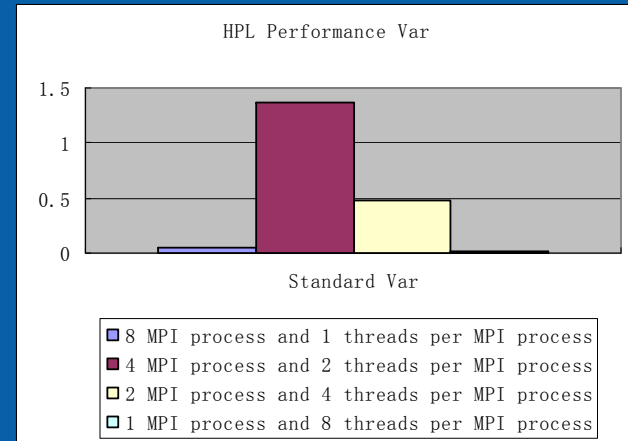
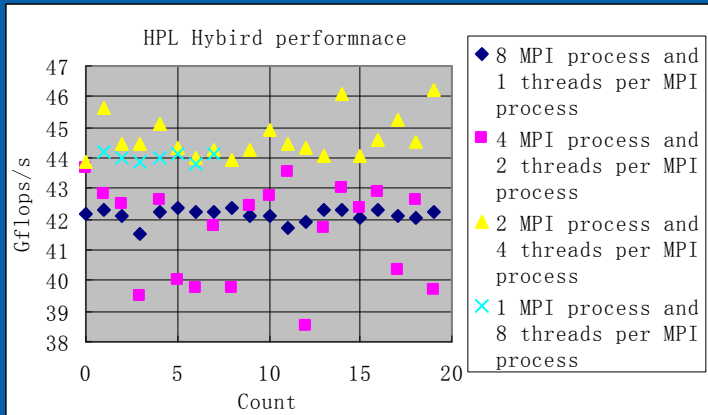
Apply Hybrid partition to HPL

- Compare pure MPI, pure OpenMP and MPI+OpenMP version, get slightly gain from hybrid model.
- OpenMP in SMP nodes can really reduce MPI communication cost, but if it's deserved to whole perf.?

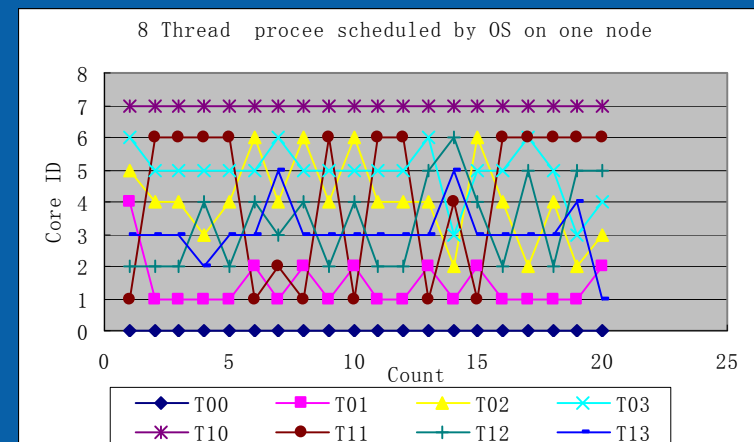
	Core number (MPI process x Threads)	Performance (Gflops/s)
MPI performance	8 x 1	42.48
OpenMP Performance	1 x 8	44.21
Hybird Performance	2 x 4	46.43



Performance statistic dithering due to threads drifting among processors

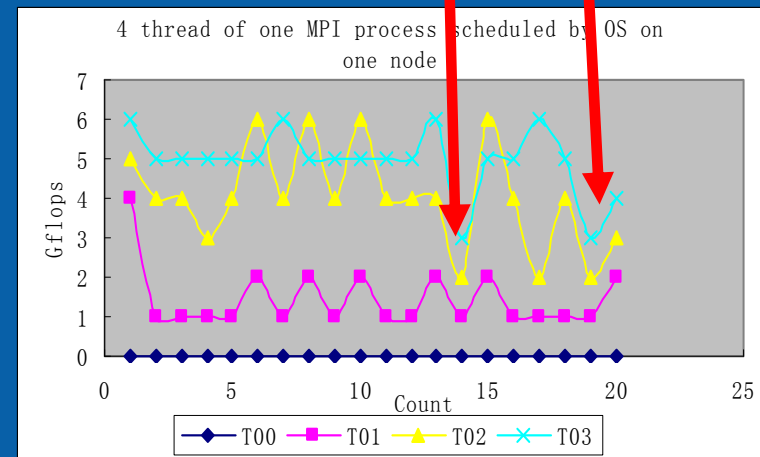
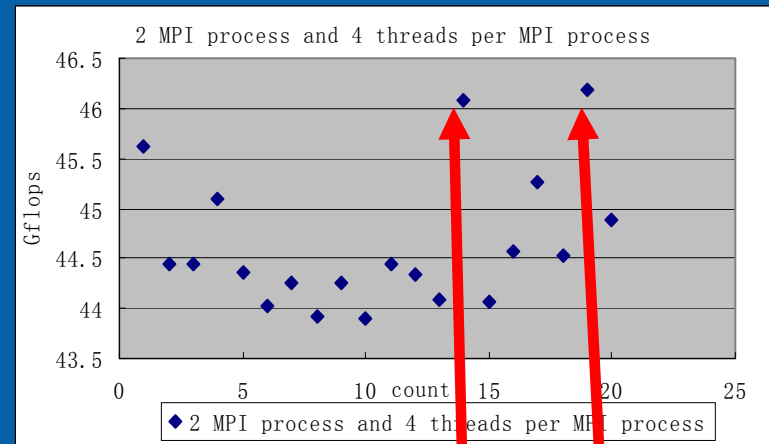


- Performance has statistic dithering in each threading application.
- Unless bundle thread to processor core), dithering is inevitable
- Threads jump between processors in run time

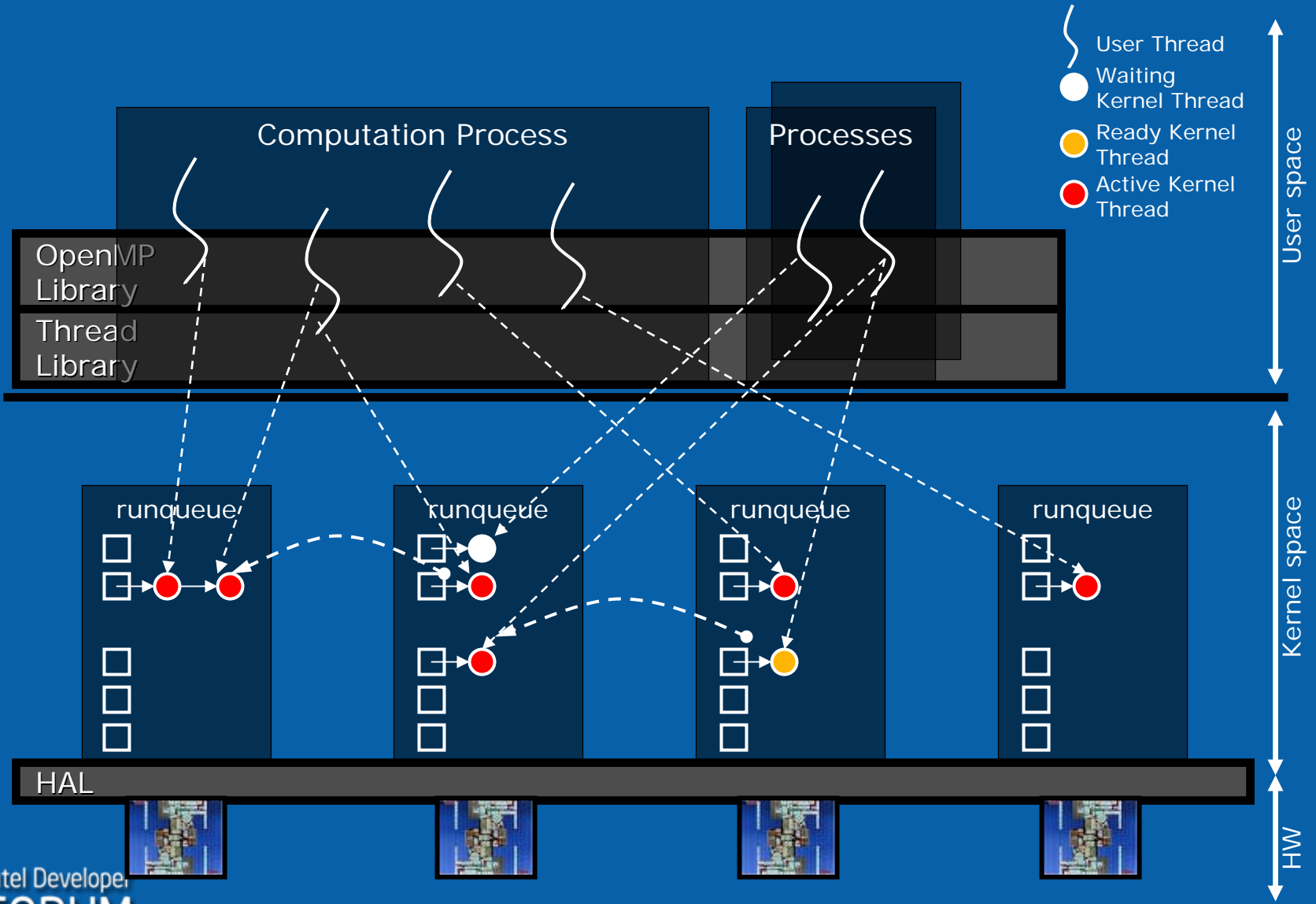


Threads drifting detected by OS thread scheduling

- When bundling MPI process on the Core 0 and Core 7. The other threads will be scheduled by OS.
- From T00 to T03 is created by the No.1 MPI process. T00 is main thread.
- When thread per MPI process was created on one socket, get best performance on Clover-town platform



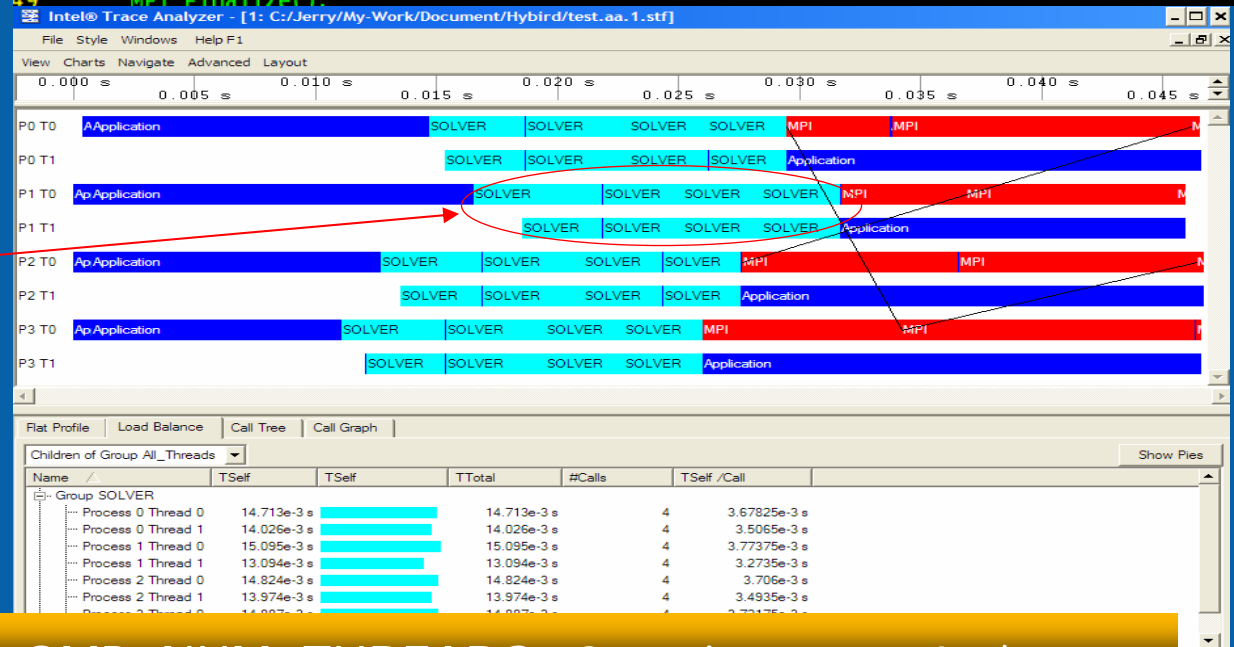
OpenMP threads inside Linux 2.6



ITAC support OpenMP threads tracing by manual adding probe code

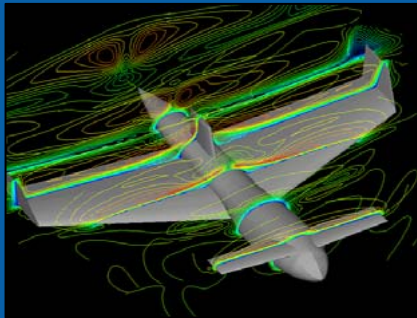
- Step 1: Insert function VT_funcdef/VT_enter/VT_leave on source code of the functions which were then executed inside threads
- Step 2: Link Intel Trace Analyzer Tools lib when compiler source code
- Step 3: Run application.
- Step 4: See the trace data with Intel Trace Analyzer Tools after application finish to run.
- Notice the threading parts
- Pay attention to overhead brought by probing code inserted

```
20 // Set VT enter
21 VT_funcdef("SOLVER:execute", VT_NOCLASS, &statehandle);
22 VT_enter(statehandle, VT_NOSCL);
23 +-- 13 行: for (k=0; k<size; k++) {-----
36 }
37 if ( myid%2 == 0 ) {
38 MPI_Send(x, 1000000, MPI_FLOAT, dest, 100, MPI_COMM_WORLD);
39 MPI_Recv(x, 1000000, MPI_FLOAT, source, 100, MPI_COMM_WORLD, &stat);
40 } else {
41 MPI_Recv(x, 1000000, MPI_FLOAT, source, 100, MPI_COMM_WORLD, &stat);
42 MPI_Send(x, 1000000, MPI_FLOAT, dest, 100, MPI_COMM_WORLD);
43 }
44 }
45 // Clear VT enter
46 VT_leave(VT_NOSCL);
47 //
48
49 MPI_Finalize();
```

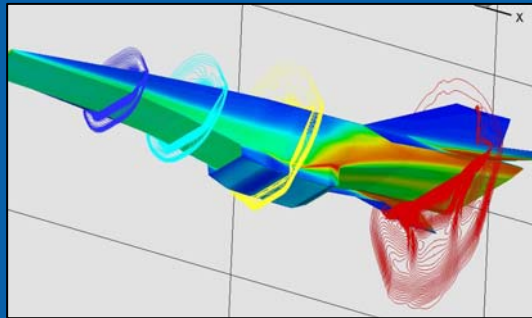


NAPA Hybrid parallel: Tune both MPI and OpenMP

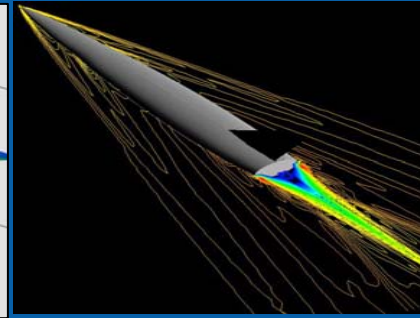
- Napa* is a NUAA self-developed CFD application, The code has been fully MPI-tuned by Intel
- MPI version has lots of point-to-point communication to update the data of ghost grids every iteration, and has no much collective communication
- MPI version Only get no more than 60 speedup on Ethernet, low scalability. This is a good candidate to apply hybrid parallel model.



MAV



X-43



NASA Lift body

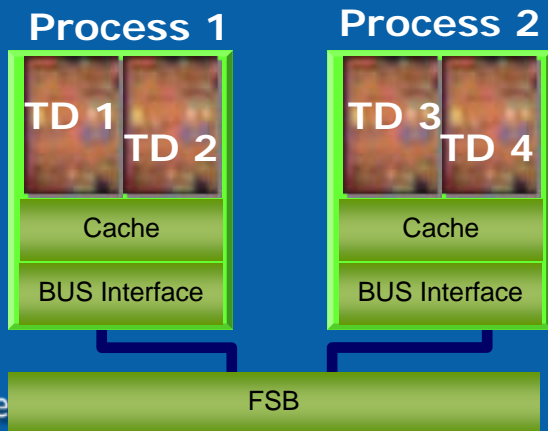
MPI tuning: Consider both Computation and Communication

- Overlapping Computation and Communication?
 - For general MPI implementation (Mpich1.2, Intel mpi 1.0) in the Share memory, Ethernet, Infiniband, Numa-link, the answer is No.
 - Some specific Message Passing Structure, the answer is Yes.
 - For example: QCDOC Supercomputer and its LQCD Message Passing API (QMP)
 - NAPA parallelization complete for cluster, validate with benchmark results.

Simple OpenMP vs. MPI

- benchmark result on Xeon 5100 1 nodes (with 4 cores)
 - mpich 1.2.7p1 compiled with Intel compiler
 - Intel Fortran compiler 9.1.033

cores \ speedup	1 cores	2 cores	4 cores
Serial	1	-	-
OpenMP	-	1.77	2.8
MPI	-	1.76	3.14 (4 MPI process)
Hybrid	-	-	3.10 (2 threads/process)

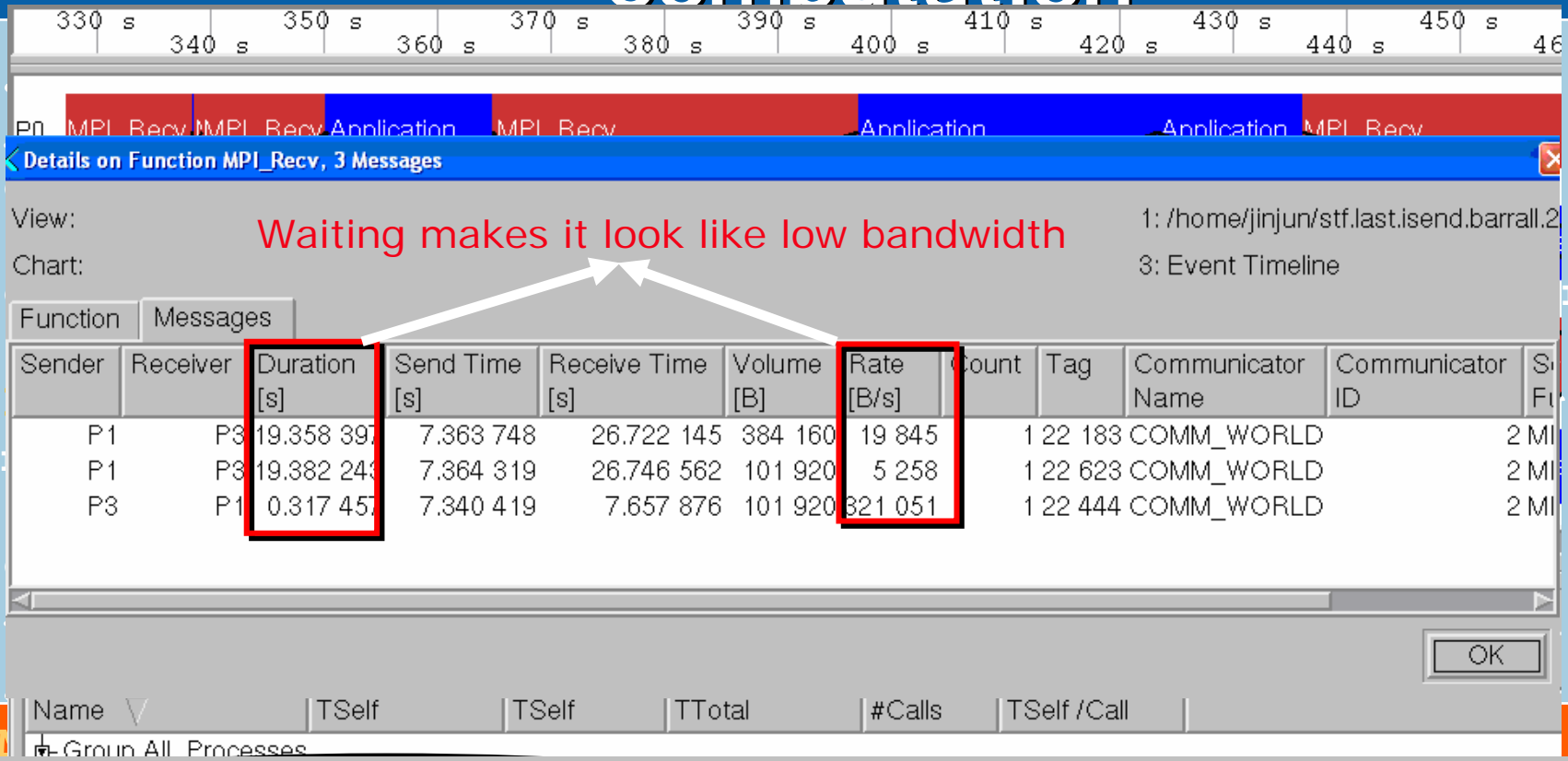


OpenMP scale not so good from 2cores to 4 cores, the overhead $T_{th} > T_{mpi}$

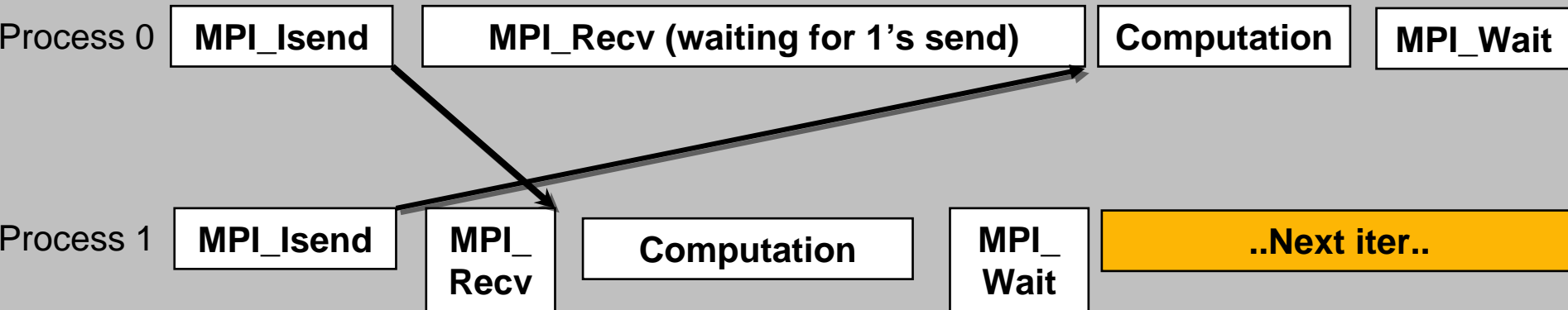
MPI has better performance on 4 cores

NAPA: MPI Communication and Computation

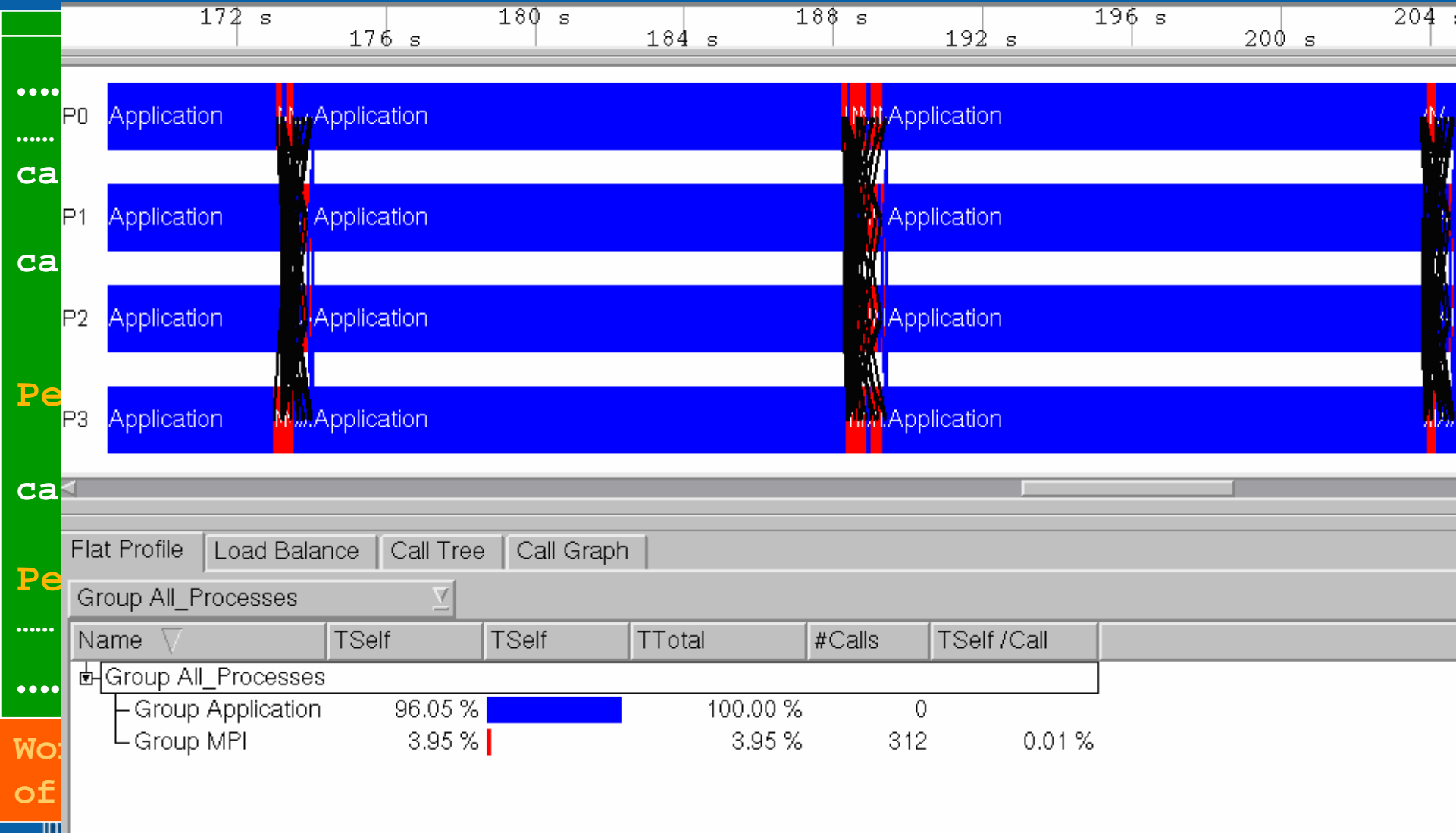
Grid:97*97*97*24



Waiting makes it look like low bandwidth

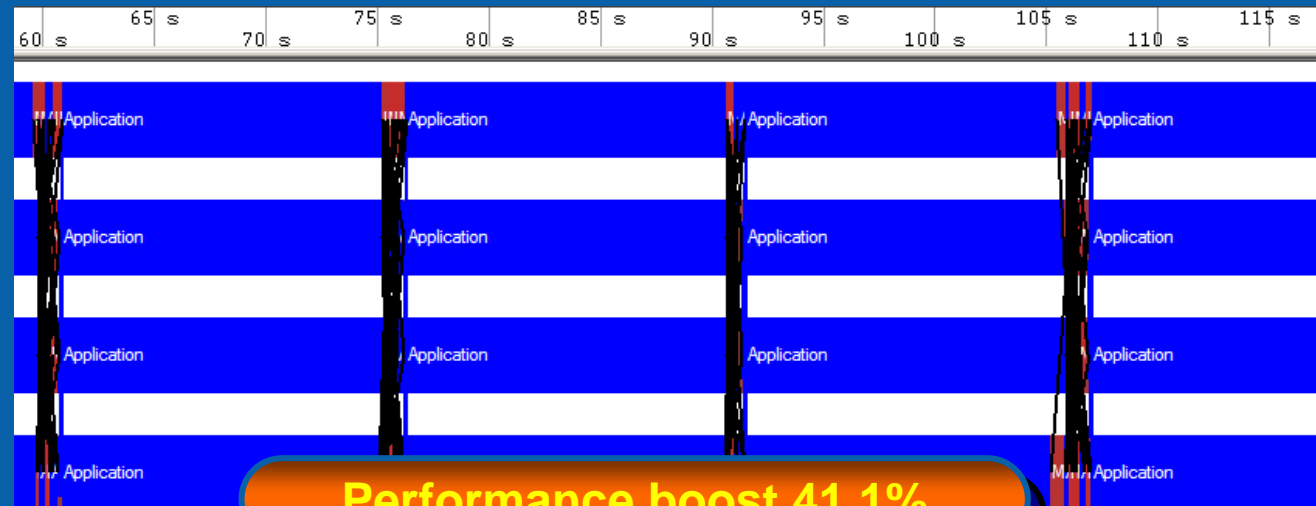
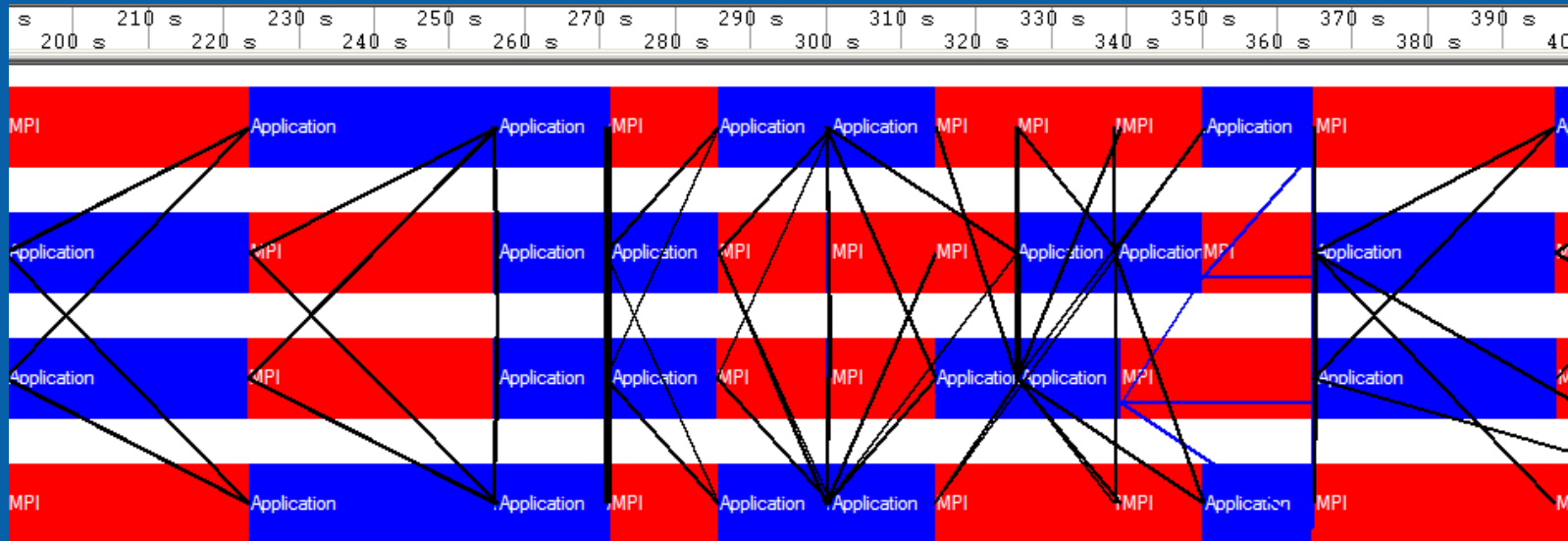


MPI Communication and Computation (cont.)



Performance has been improved from 6m52s to 4m03s, 1.45X speedup

NAPA: MPI Parallel optimization apply to real workload

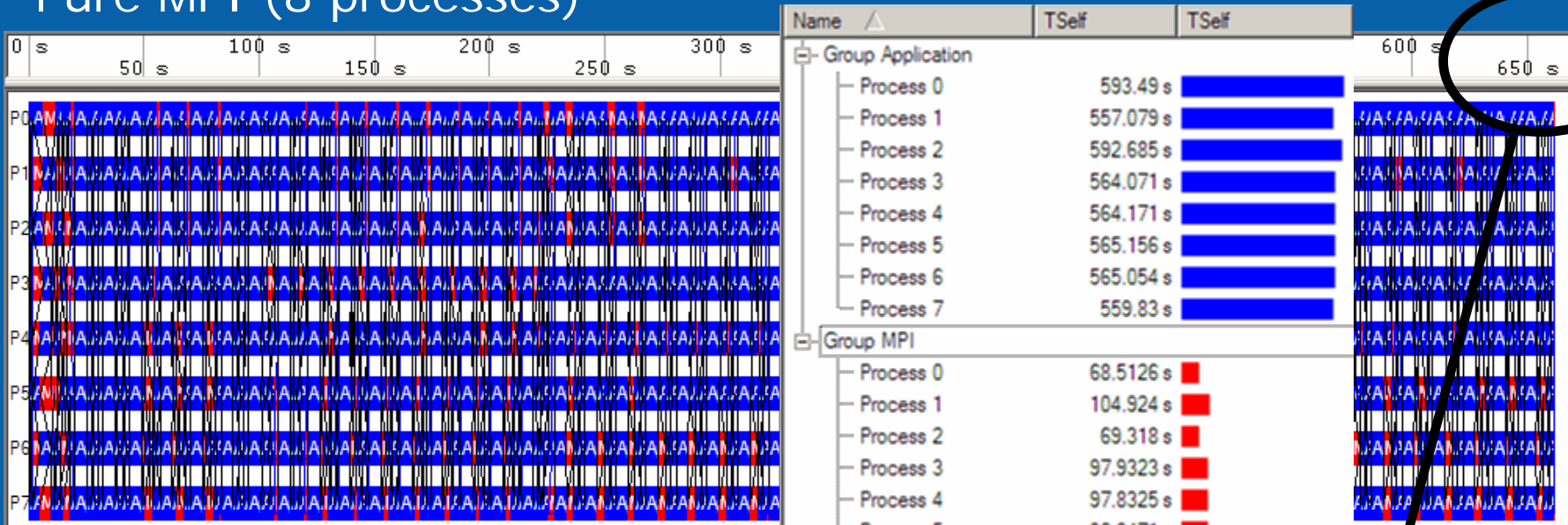


Performance boost 41.1%

Workload: Ultra-sonic in-let
4 processes

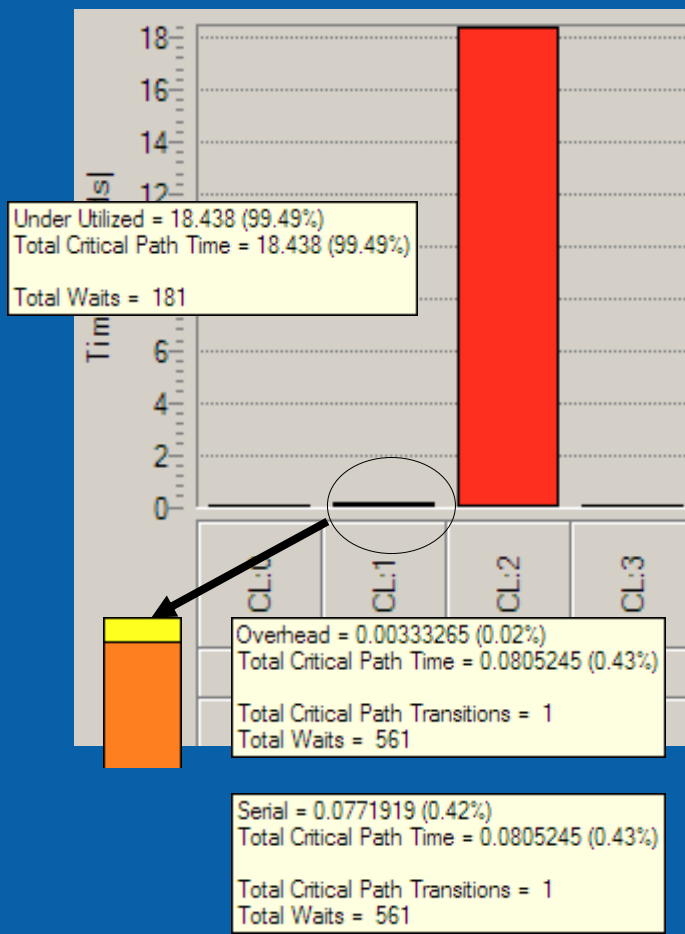
Threads reduce the computation

Pure MPI (8 processes)

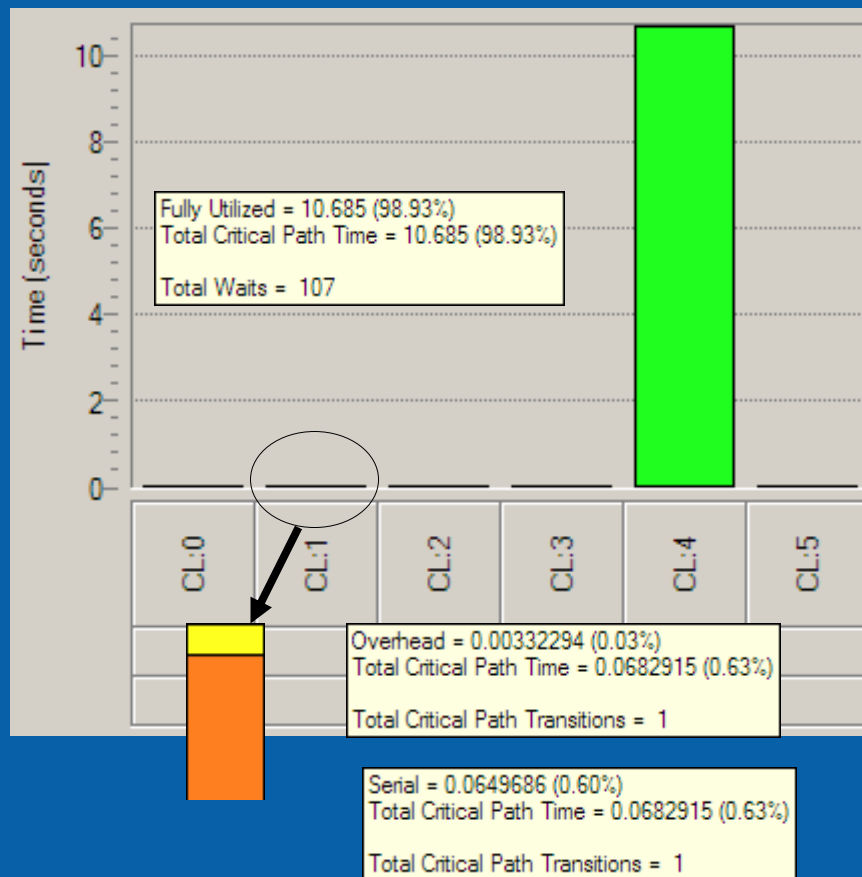


Application time per process is nearly the same in Hybrid model, the communication is greatly reduced from 13.9% (PureMPI) to 3.84%(Hybrid)

Profile OpenMP thread by Intel threading profiler



2 threads



4 threads

Thread State

- Active
- Spin
- Wait

Critical Path ...

- Concurrency
- Behavior

- No Thread Active
- Serial
- Under Utilized
- Fully Utilized
- Over Utilized
- Overhead

After MPI tuning, avoid naive OpenMP – fine tune it.

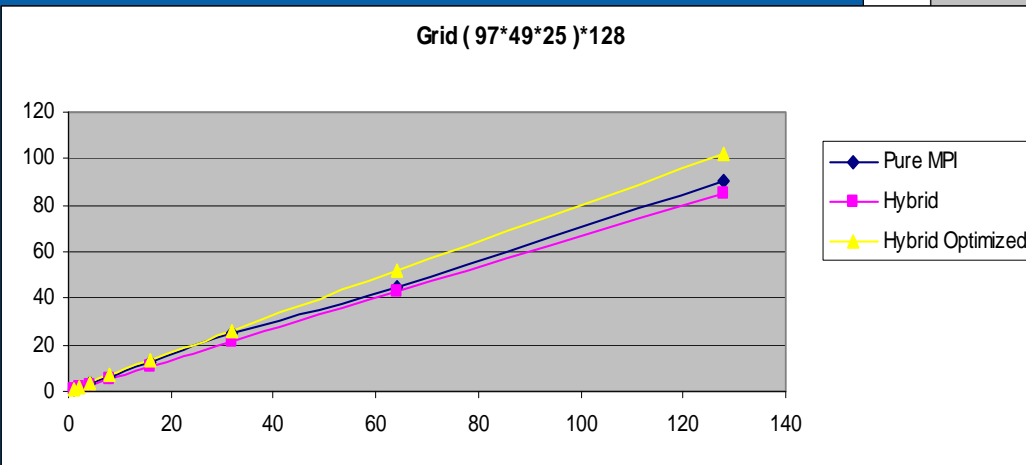
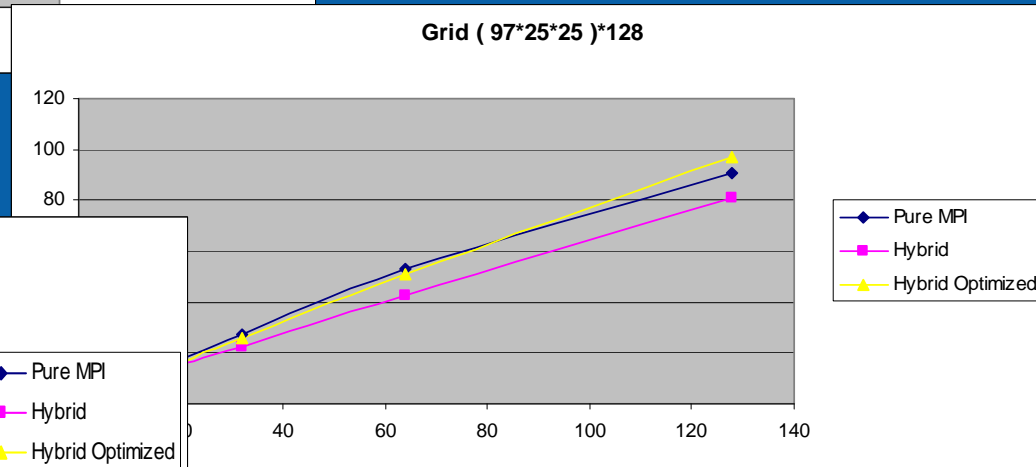
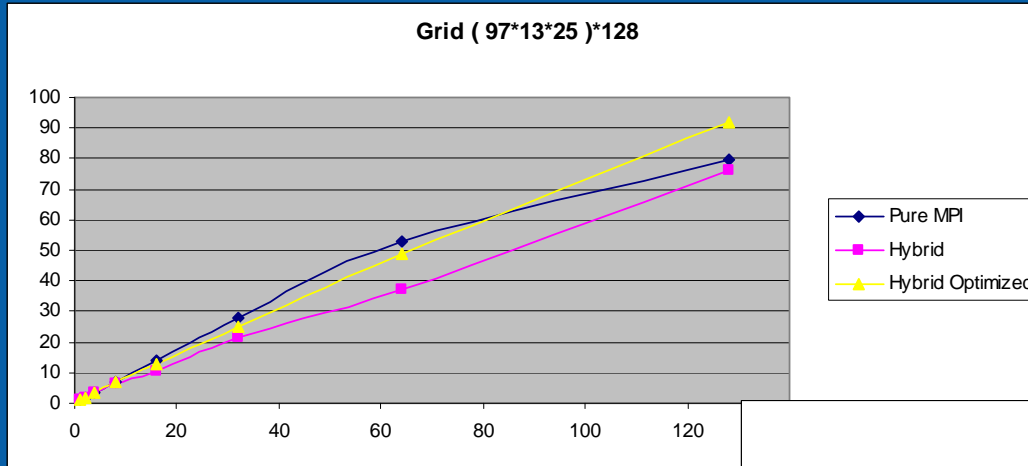
<pre> 1 2 C----- 3 C NEXT PLANES 4 C----- 5 DO 1300 I = 3, IL-1 6 7 DO 1200 K = 2, KL 8 DO 1200 J = 2, JL 9 AM(J,K) = AD(J,K) 10 AD(J,K) = AP(J,K) 11 12 DISI = DISS(I+1, J, K, 1) 13 CFLI = CFL*DISI/(RADI(I+1, J, K)+MAX(RADJ(I+1, J, K), RADK(I+1, J, K))) 14 CFLA = 0.25 * (CFLI*CFLI/(CFL*CFLL)-1.) 15 AP(J, K) = MAX(CFLA, SMOOP1) 16 17 18 DO 1300 K = 2, KL 19 DO 1300 J = 2, JL 20 T = 1./(1.+2.*AD(I,K)-AM(I,K)*DI(I-1,J,K,1)) 21 DI(I, J, K, 1) = T*AP(J, K) 22 23 DO 1300 N = 1, NVAR 24 W(I, J, K, N) = T*(W(I, J, K, N) + AM(J, K)*W(I-1, J, K, N)) 25 1300 CONTINUE </pre>	<pre> 1 !\$omp parallel do private(disi,cfli,cfla) 2 C----- 3 C NEXT PLANES 4 C----- 5 6 DO 1200 K = 2, KL 7 DO 1200 J = 2, JL 8 DO 1200 I = 3, IL-1 9 10 AM(I-1, J, K) = AD(I-1, J, K) 11 AD(I, J, K) = AP(I, J, K) 12 13 DISI = DISS(I+1, J, K, 1) 14 CFLI = CFL*DISI/(RADI(I+1, J, K)+MAX(RADJ(I+1, J, K), RADK(I+1, J, K))) 15 CFLA = 0.25 * (CFLI*CFLI/(CFL*CFLL)-1.) 16 AP(I+1, J, K) = MAX(CFLA, SMOOP1) 17 18 TT(I, J, K) = 1./(1.+2.*AD(I, J, K)-AM(I-1, J, K)*DI(I-1, J, K, 1)) 19 DI(I, J, K, 1) = TT(I, J, K)*AP(I, J, K) 20 1200 CONTINUE 21 22 DO 1300 N = 1, NVAR 23 !\$omp parallel do 24 DO 1300 K = 2, KL 25 DO 1300 J = 2, JL 26 DO 1300 I = 3, IL-1 27 W(I, J, K, N) = TT(I, J, K)*(W(I, J, K, N) + AM(I-1, J, K)*W(I-1, J, K, N)) 28 1300 CONTINUE </pre>
--	---

Fine tune the OpenMP part, make its grain more coarse and decouple the data dependence in loop

Benchmark Result: Hybrid boost speedup based on MPI

Hybrid: 2 threads get 1.6~ speedup

Hybrid opt: 2 threads get 1.77 speedup

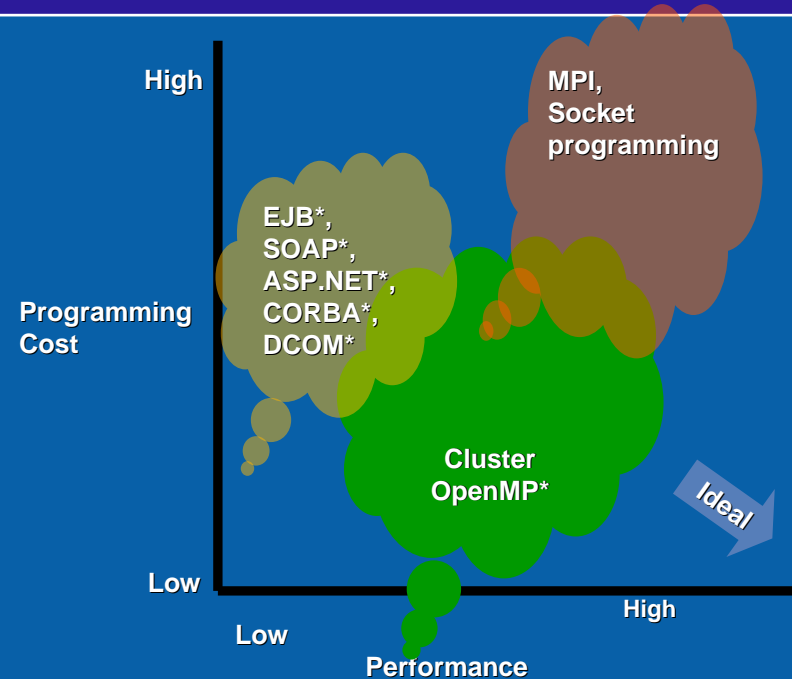


- For hybrid, every process contain 2 threads
- MPI benchmark implemented Share-Memory on same Node

Intel Cluster OpenMP as “Hybrid” model: A tool to avoid MPI complexity

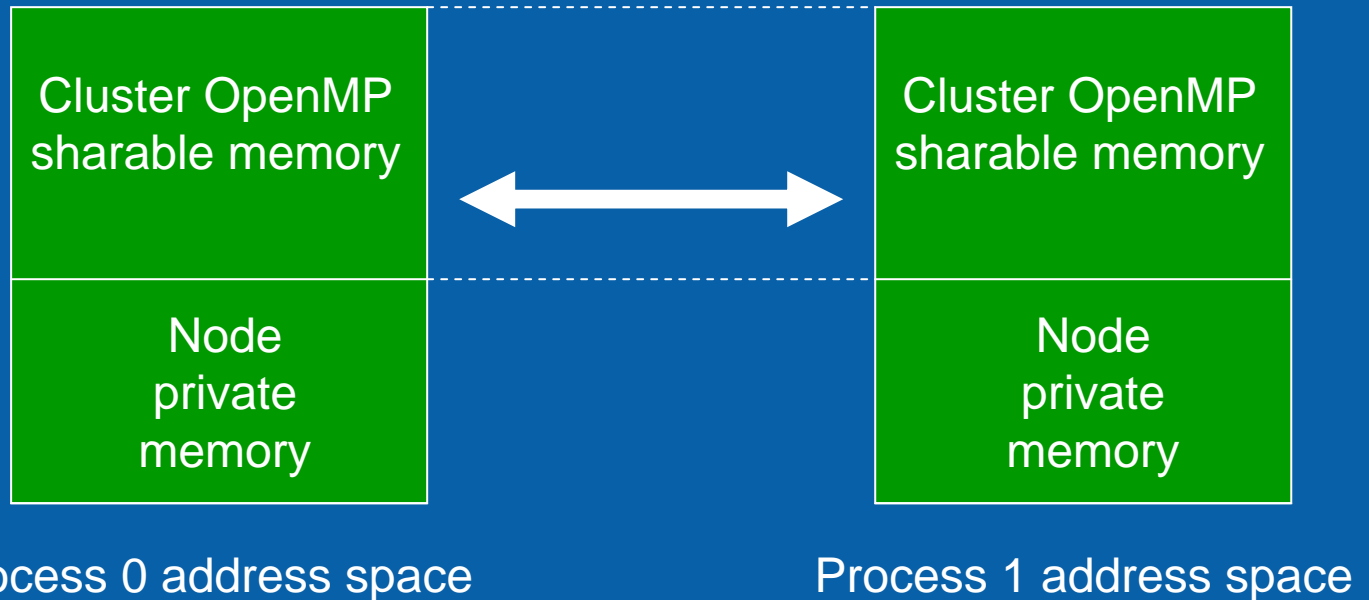
- Intel Cluster OpenMP is an OpenMP extension to apply OpenMP threads across Cluster nodes.
- Cluster OpenMP* is an easier way to program clusters. Certain codes can achieve good performance and scaling with Cluster OpenMP. The process of porting a code to Cluster OpenMP is systematic and much less work than the porting process to a message-passing form. Performance tools help tune the program for best performance on a cluster.

Enables the execution of OpenMP applications on Intel based 64 bit clusters.



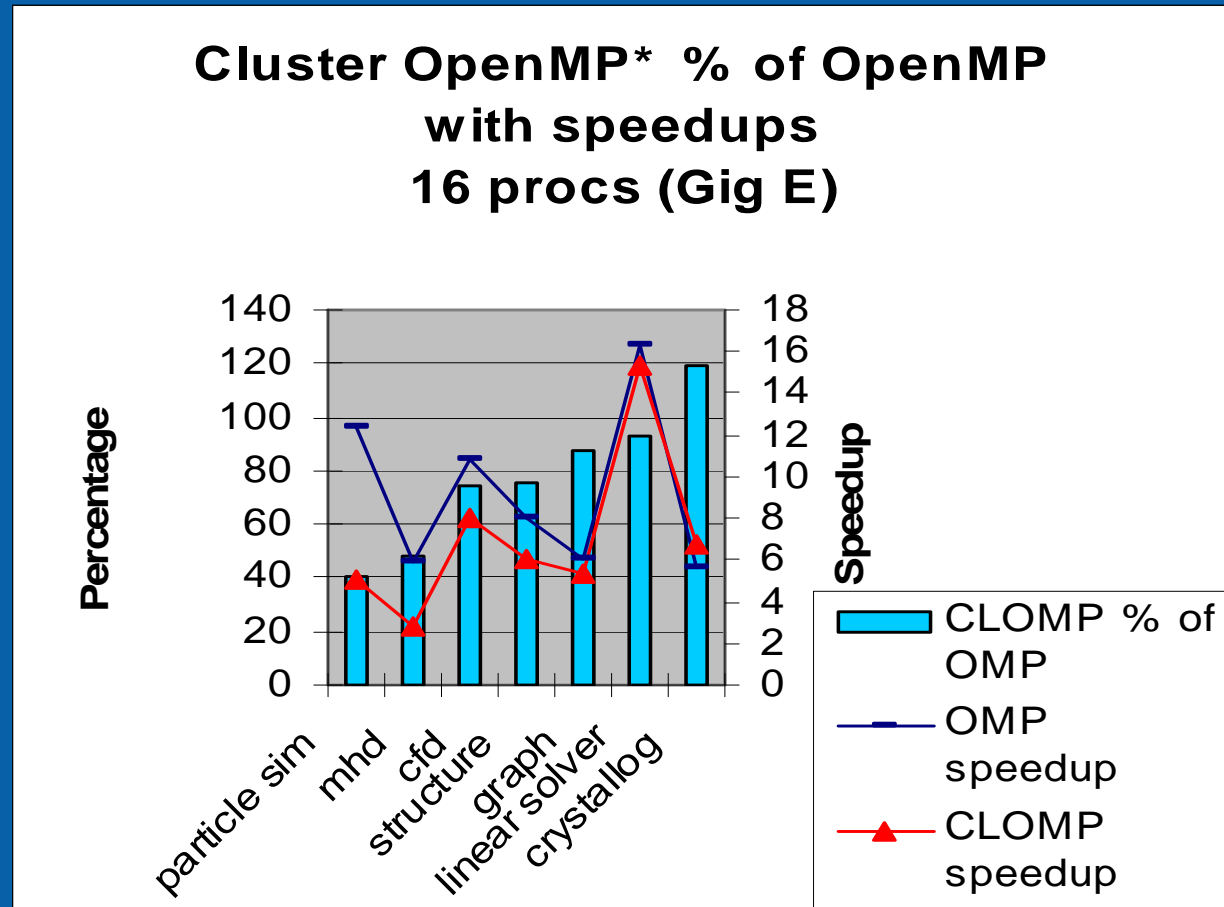
How does Cluster OpenMP work?

Cluster OpenMP* Memory Model



Cluster OpenMP is implemented by compiler extensions and a run-time library that supports running an OpenMP program on a cluster. The ClusterOpenMP memory model create a virtual share-memory region across the distributed nodes, which make it easy for users to avoid MPI programming while take advantage of distributed memory.

Scalability of Cluster OpenMP in Hybrid Architecture

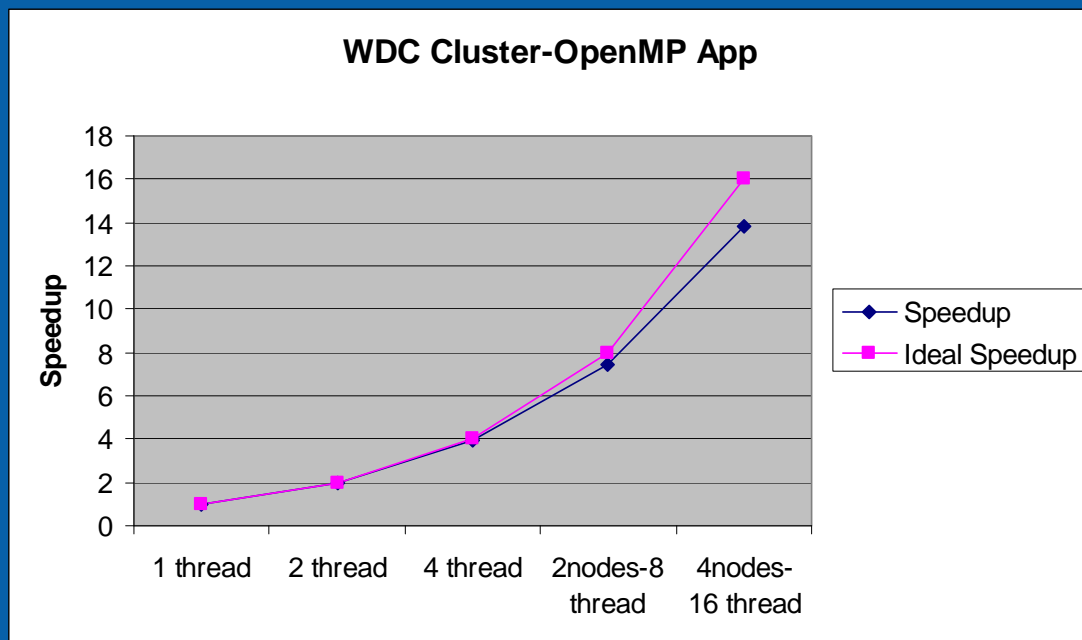


* Other names and brands may be claimed as the property of others.

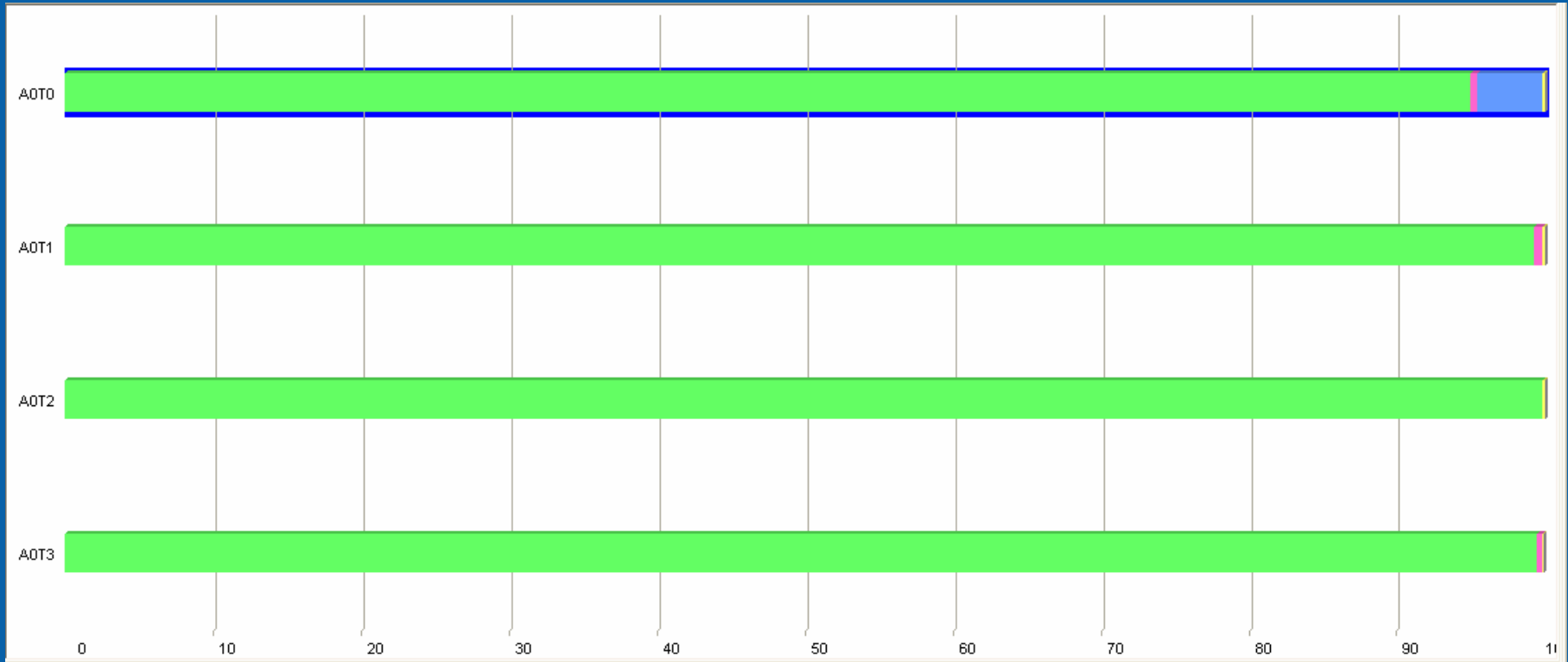
Case Study: GalRen in 2-way Xeon 5100 Hybrid Cluster

- Main function: rendering a movie of a flight through space.
- It uses a database of stars (brightness, color, position, etc) and does ray-tracing through them.
- Each thread takes a set of scan-lines of the frame buffer and computes the pixels.

14X Speedup for 16 Threads



Case Study: GaRen:4threads in one node



Label	Total	% Parallel	% Sequential	% Imbalance	% Barrier	% Locks	% Synchronized	% Parallel overheads	% Sequential overheads
A0T2	11...	99.815	0.000	0.000	0.085	0.000	0.000	0.100	0.000
A0T0	11...	95.028	4.471	0.000	0.407	0.000	0.000	0.094	0.000
A0T1	11...	99.305	0.000	0.000	0.596	0.000	0.000	0.099	0.000
A0T3	11...	99.476	0.000	0.000	0.426	0.000	0.000	0.098	0.000

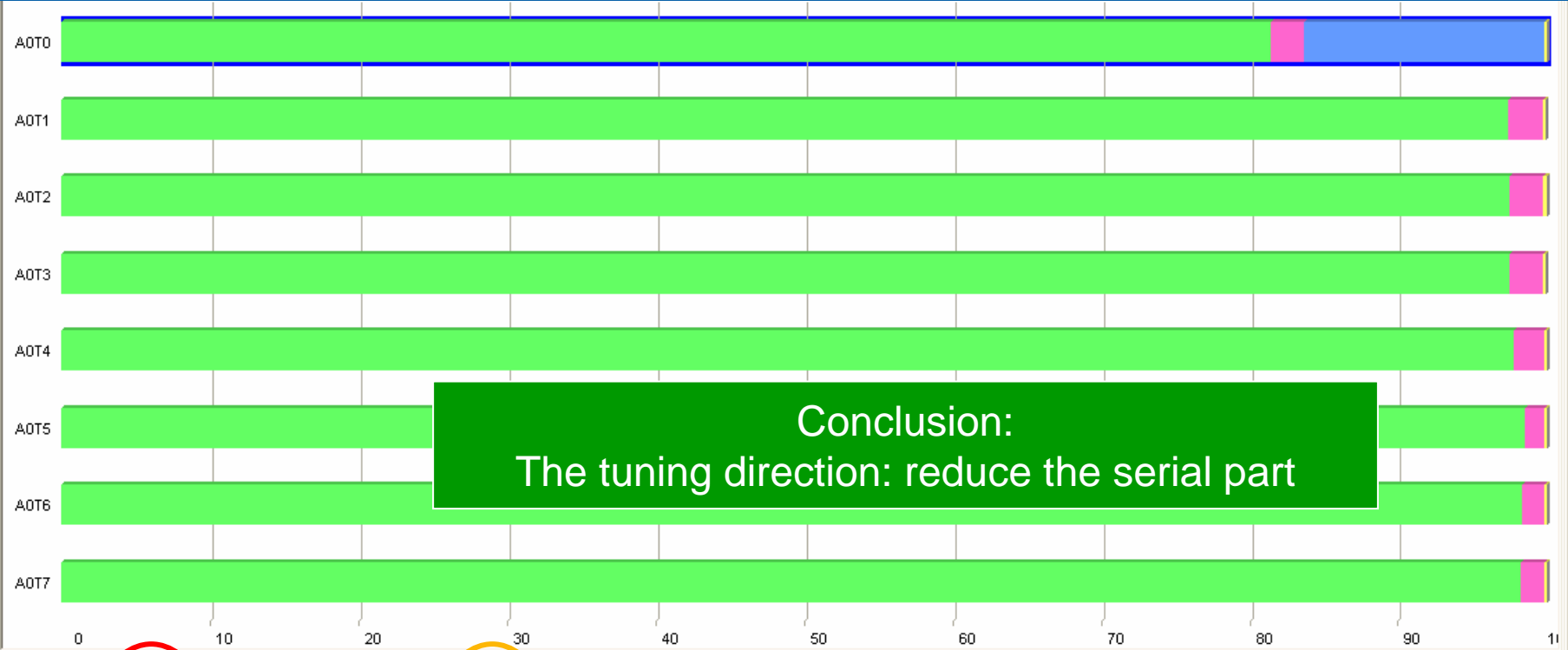
95%+ parallel

4% cost serial part

0.5% for Barrier, Synchron

1 runs 1 showing, 7 regions 7 showing

8threads in 2 nodes



Label	Total	% Parallel	% Sequential	% Imbalance	% Barrier	% Locks	% Synchronized	% Parallel overheads	% Sequential overheads
A0T0	8.420	81.418	16.152	0.000	2.257	0.000	0.000	0.174	0.000
A0T3	7.050	97.515	0.000	0.000	2.270	0.000	0.000	0.216	0.000
A0T1	7.050	97.392	0.000	0.000	2.408	0.000	0.000	0.200	0.000
A0T4	7.030	97.835	0.000	0.000	1.991	0.000	0.000	0.173	0.000
A0T2	7.050	97.518	0.000	0.000	2.266	0.000	0.000	0.215	0.000
A0T6	7.010	98.404	0.000	0.000	1.422	0.000	0.000	0.174	0.000
A0T5	7.030	98.543	0.000	0.000	1.280	0.000	0.000	0.177	0.000
A0T7	7.050	98.263	0.000	0.000	1.560	0.000	0.000	0.175	0.000

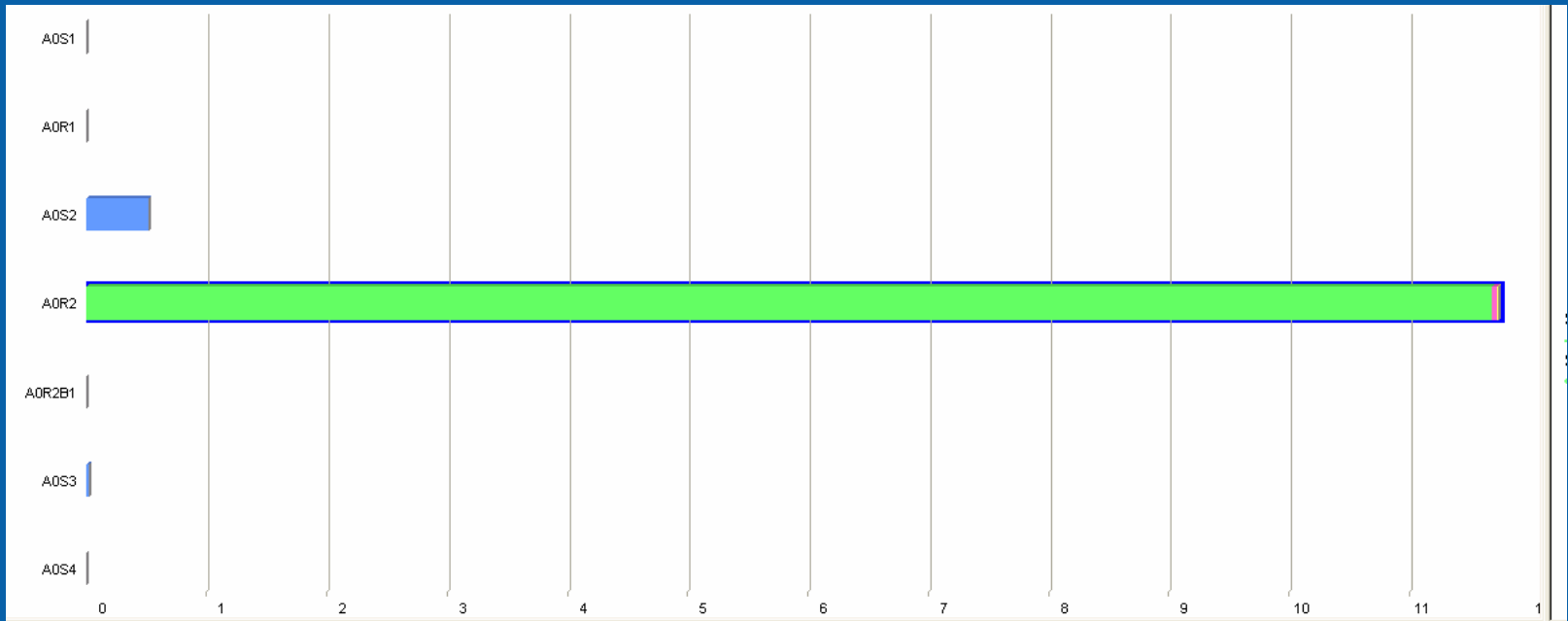
81%+ parallel

16% cost for serial part

2.5% for Barrier, Synch Part

How to find out the source code: Location vs Threads

4threads in one node



Label	Num. threads	Total	Parallel	Sequential	I...	B...	L...	S...	P...	S...	Location
AOS2	1	0.520	0.000	0.520	0.0...	0.0...	0.0...	0.0...	0.0...	0.0...	/home/it/clusteromp/galren/Galren/galaxy/galren.c (main) after par. region @73, before par. region /home/it/clusteromp/galren/Galren/galaxy/render.c@
AOS3	1	0.030	0.000	0.030	0.0...	0.0...	0.0...	0.0...	0.0...	0.0...	/home/it/clusteromp/galren/Galren/oalaxw/render.c (render_oalaxw) after par. region @753, before par. region @741
AOR2	4	11.750	11.693	0.000	0.0...	0.0...	0.0...	0.0...	0.0...	0.0...	/home/it/clusteromp/galren/Galren/galaxy/render.c (render_galaxy) @ line 741

Parallel part location

serial code location

```
#pragma omp parallel private(i, j, R, G, B, A)
```

```
{
```

```
#if defined(OMP)
```

```
int thread_num = omp_get_thread_num();  
int my_slice = xsize / omp_get_num_threads();  
int my_offset = (thread_num * my_slice);  
int my_bytes = ( sizeof( uByte8 ) * my_slice );  
uByte8 *my_pixels;
```

```
my_pixels = malloc( my_bytes );  
memset( my_pixels, 0x55, my_bytes );
```

```
#endif
```

```
#pragma omp for schedule(runtime)
```

```
for (j=0; j < ysize; j++) {  
    GraphicsDrawRow( my_offset, ysize - j - 1, my_slice, my_pixels );  
    for (i=0; i < xsize; i++) {  
        ComputeRay( g, s, image, i, j);  
    } /* end for index */  
    GraphicsDrawRow( 0, ysize - j - 1, xsize, &image[xsize*j]);  
}
```

```
#if defined(OMP)
```

```
free( my_pixels );
```

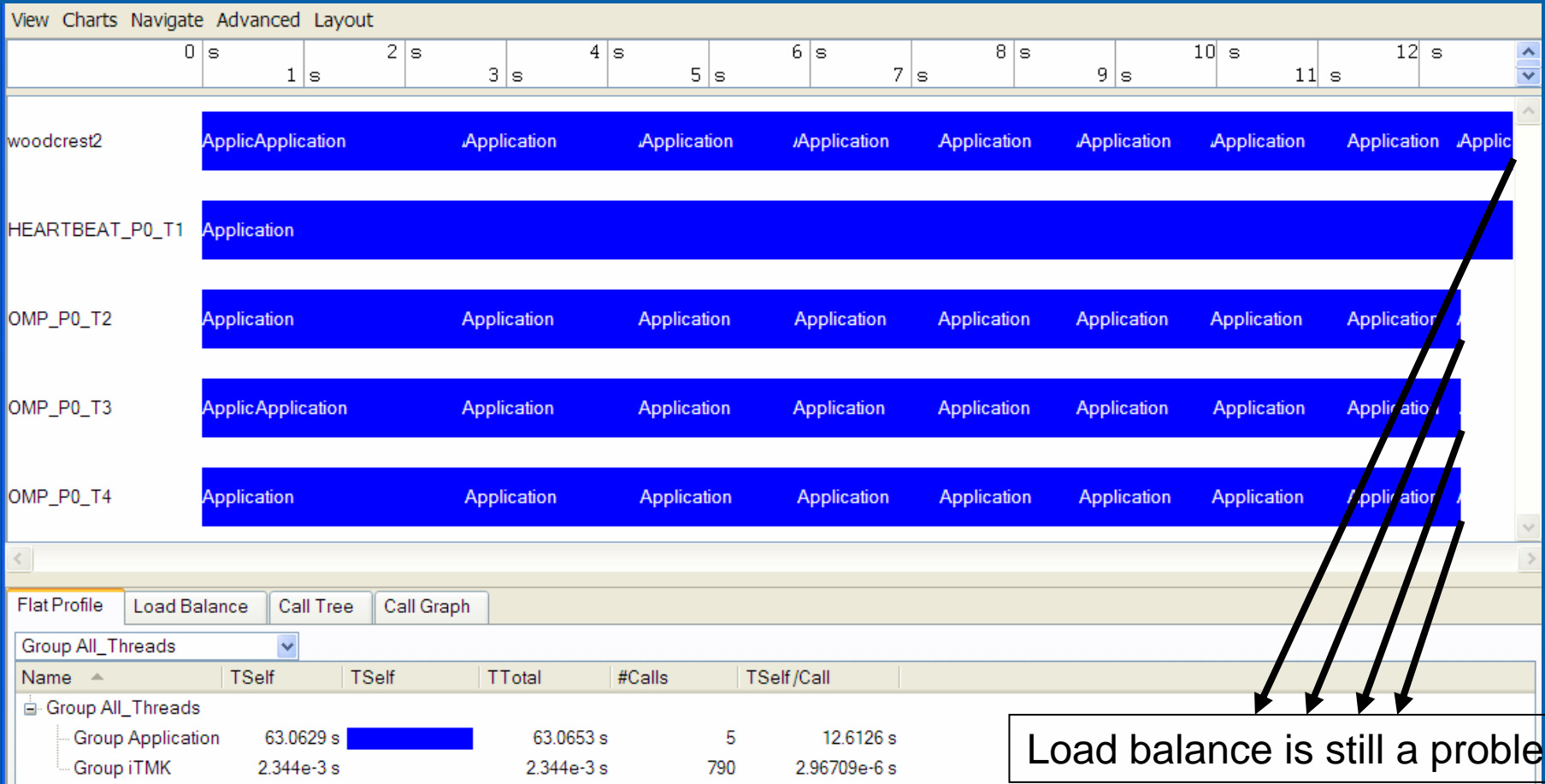
```
#endif
```

```
}
```

*One Part of Cluster OpenMP
Pseudo Code*

Very Little Change needed

ITC/ITA(4 threads* 1 nodes)



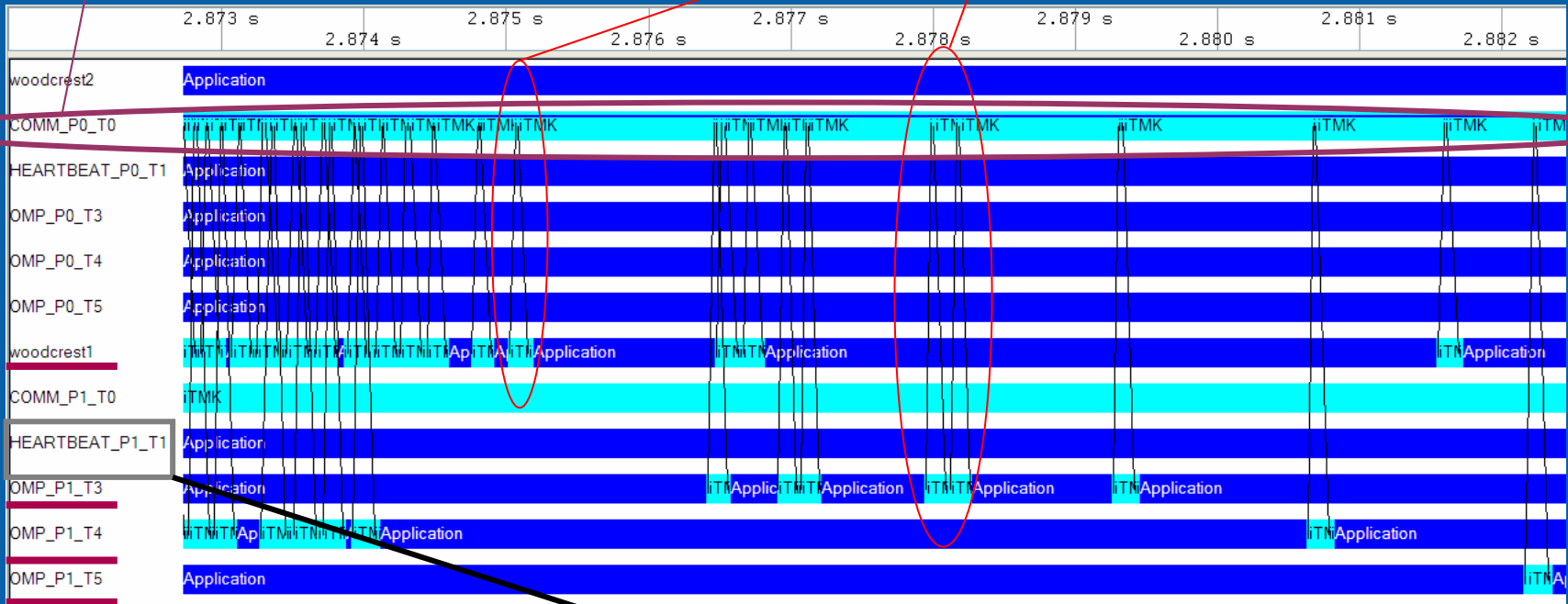
Load balance is still a problem

Almost Non-Communication part in single node

ITC/ITA(4 threads* 2 nodes)

COMM_P*_T0 Thread deals communication with threads in other nodes

Message only between nodes

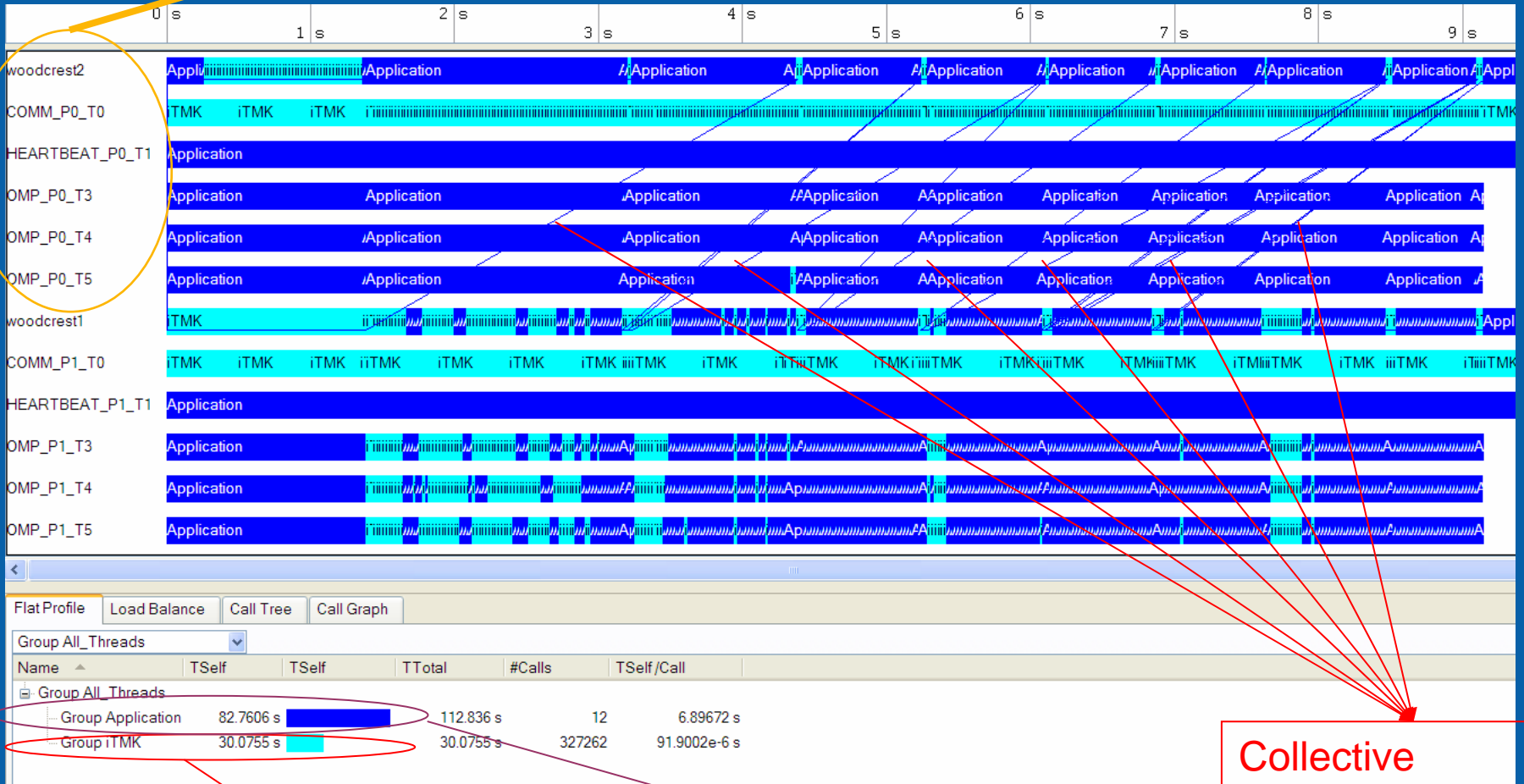


Main thread and the other **3 slave threads** in one node **for computation**

Heart beat Thread each node

ITC/ITA(4 threads * 2 nodes)

More threads get better overlap for communication and computation



Collective Communication

Communication part

Computation part

Summary/Call to Action

- Consider both threading and MPI parallelism to maximize scalability in current Multicore SMP Cluster.
- Tradeoff between MPI and OpenMP, use Intel ITAC and threading tools to Analyze the threads' behavior and overhead.
- Considering parallel granularity, apply MPI tuning (large granularity) before introducing fine-tuned OpenMP threading.
- Use Cluster OpenMP can fast deploy/debug Multithreading in SMP cluster, get reasonable speedup to certain type of applications without MPI/PVM programming.

Please fill out the Session Evaluation Form

**Thank You for your input, we use it
to improve future Intel Developer
Forum events**

**Save your date for IDFs this Fall:
San Francisco, USA 2007 September 18 -20
Taipei, Taiwan, 2007 October 15 -16**

Thanks!

- Thanks for ICSC CRT/HPC enabling team, especially for Wang Zhe, Qiao Nan, Jin Jun and Xu Jin.
- Also thanks for Hoeflinger, Jay P, Meadows, Lawrence F and Ohlay, Patrick for their help in Cluster OpenMP and ITAC threading tracing problem

Backup

Advantages of Cluster OpenMP

- To the user
 - Expand problem size
 - Reduce the cost of computing.
 - Eliminate capacity bottlenecks
- To the developer
 - Reduce development time
 - Improve maintainability
 - Could take advantage of existing OpenMP code
 - The implementation is easy going, no need to master MPI or PVM knowledge.
 - Code is easy to read, keeping consistent with the serial code more than MPI/PVM code.
 - Debug is simple as it's small step based

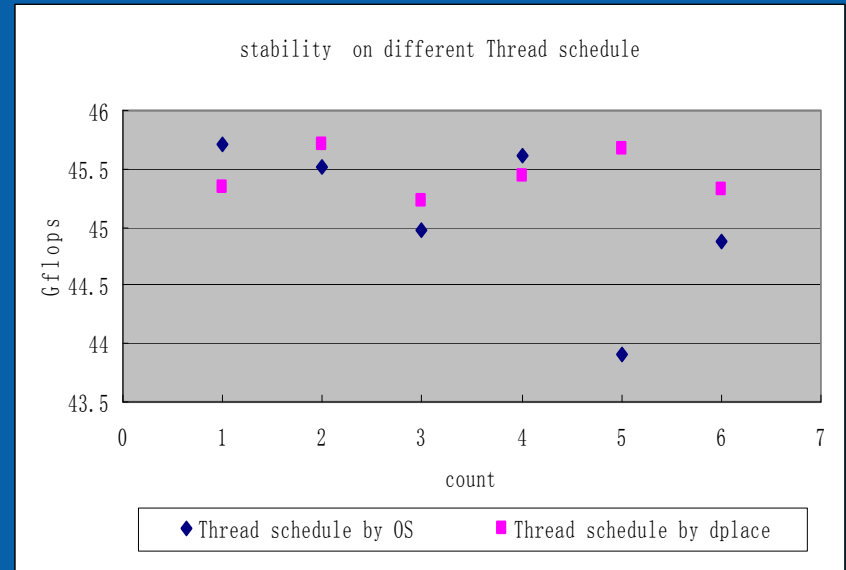
Solution:

Method #1:

Some OS can supply tools to permit to schedule Thread|Process on different core|cpu. Such as SGI or BULL's NUMA tools.

Method #2:

```
Insert OpenMP function sched_setffinity on OpenMP region. Such as :  
#pragma omp parallel  
{  
    unsigned long mask = (1 << omp_get_thread_num()) * 2 );  
  
    sched_setaffinity(0, sizeof(mask), &mask);  
    source code  
}
```



Benchmark on SGI Altix. Use dplace to schedule threads