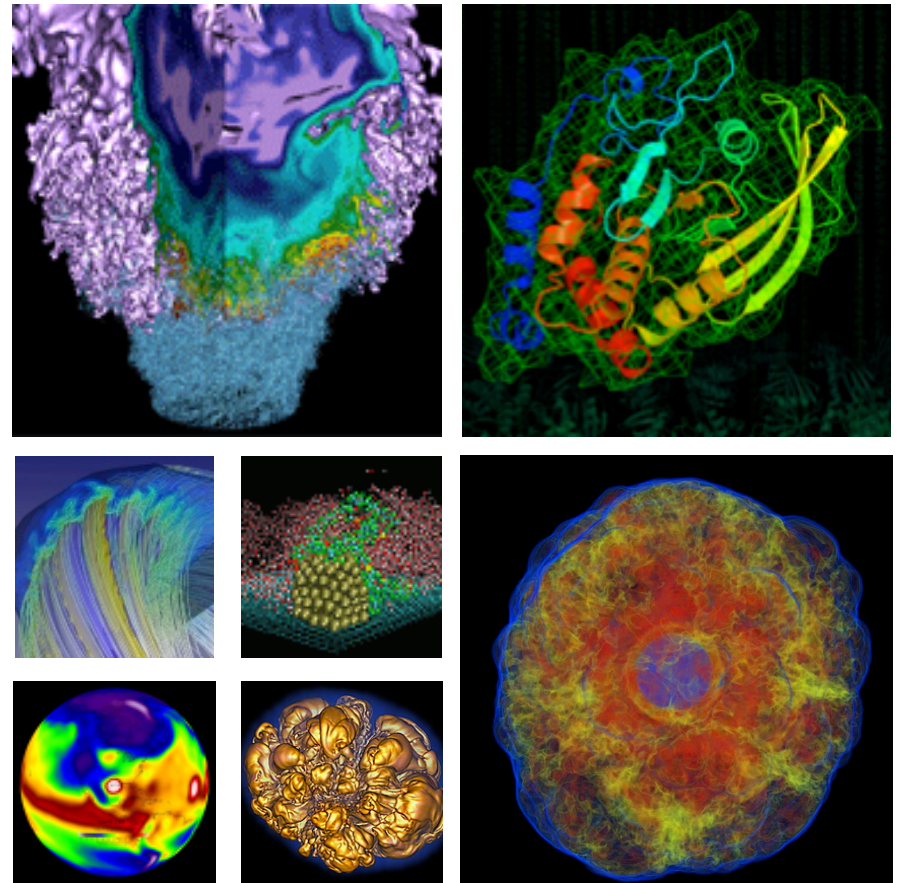


# Debugging and Optimization Tools



**Richard Gerber**  
NERSC  
User Services Group Lead

Thanks to Woo-Sun Yang and Helen He

- **Take-Aways**
- **Debugging**
- **Performance / Optimization**
- **NERSC “automatic” tools**

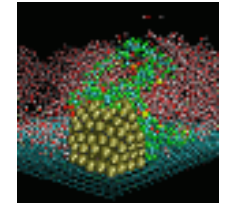
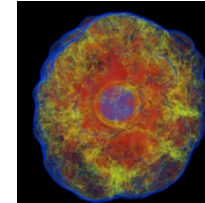
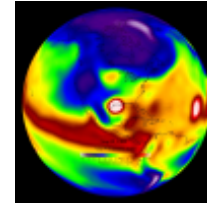
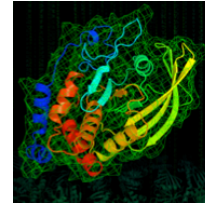
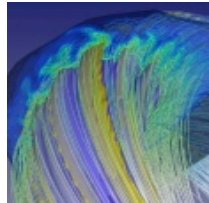
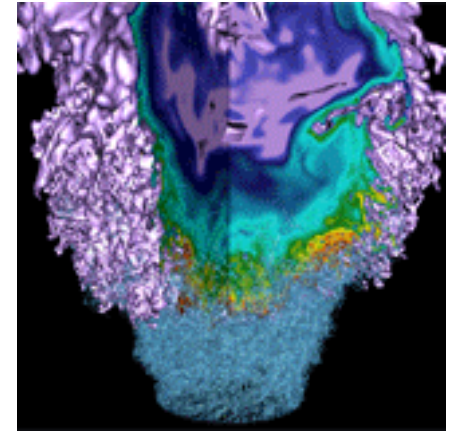
Videos, presentations, and references:

<http://www.nersc.gov/users/training/courses/CS267/>

Help at NERSC  
consult@nersc.gov  
<https://help.nersc.gov>

- **Tools can help you find errors in your program and locate performance bottlenecks**
- **In the world of HPC parallel computing, there are few widely adopted standard tools – try some and see what works for you**
  - Totalview and DDT debuggers
  - PAPI, Tau, various research tools, & vendor-specific performance tools
- **Common code problems**
- **How tools work in general**
- **Be suspicious of outliers among parallel tasks**
- **Where to get more information**

# Debugging



# What is a Bug?



- A bug is when your code

- crashes
- hangs (doesn't finish)
- gets inconsistent answers
- produces wrong answers
- behaves in any way you didn't want it to

1100 Started Cosine Tape (Sine check)  
1525 Started Mult+ Adder Test.



1545  
1630 Antenna started.  
1700 closed down.  
First actual case of bug being found.

The term "bug" was popularized by Grace Hopper (motivated by the removal of an actual moth from a computer relay in 1947)

# Common Causes of Bugs



- **“Serial” (Sequential might be a better word)**
  - Invalid memory references
  - Array reference out of bounds
  - Divide by zero
  - Use of uninitialized variables
- **Parallel** **Let’s concentrate on these**
  - Unmatched sends/receives
  - Blocking receive before corresponding send
  - Out of order collectives
  - Race conditions
  - Unintentionally modifying shared memory structures

# What to Do if You Have a Bug?



- **Find It**
  - You want to locate the part of your code that isn't doing what it's designed to do
- **Fix It**
  - Figure out how to solve it and implement a solution
- **Check It**
  - Run it to check for proper behavior



quickmeme.com

<http://www.geekherocomic.com/>



U.S. DEPARTMENT OF **ENERGY**

Office of Science





- **printf, write**
  - Versatile, sometimes useful
  - Doesn't scale well
  - Not interactive
  - Fishing expedition
- **Compiler / Runtime**
  - Bounds checking, exception handling
  - Dereferencing of NULL pointers
  - Function and subroutine interface checking
- **Serial gdb + friends**
  - GNU debugger, serial, command-line interface
  - See “man gdb”
- **Parallel debuggers**
  - DDT
  - Totalview
- **Memory debuggers**
  - MAP
  - Valgrind

See NERSC web site  
<https://www.nersc.gov/users/software/debugging-and-profiling/>

# Parallel Programming Bug



This code hangs because both Task 0 and Task N-1 are blocking on MPI\_Recv

```
if(task_no==0) {  
    ret = MPI_Recv(&herBuffer, 50, MPI_DOUBLE, totTasks-1, 0,  
MPI_COMM_WORLD, &status);  
    ret = MPI_Send(&myBuffer, 50, MPI_DOUBLE, totTasks-1, 0,  
MPI_COMM_WORLD);  
} else if (task_no==(totTasks-1)) {  
    ret = MPI_Recv(&herBuffer, 50, MPI_DOUBLE, 0, 0,  
MPI_COMM_WORLD, &status);  
    ret = MPI_Send(&myBuffer, 50, MPI_DOUBLE, 0, 0,  
MPI_COMM_WORLD);  
}
```

# NERSC NX – Accelerate Your X Connection



## FOR USERS

- » Live Status
- » My NERSC
- » Getting Started
- » Computational Systems
- » Data & File Systems
- » Network Connections
  - Connecting to NERSC
  - Using X Windows

### Connecting to NERSC with NX

#### NX FAQ

[Download Tested NX Player](#)

[NX Configuration File](#)

[Startup Tutorial](#)

[Transferring Data](#)

[Network Performance](#)

» [Queues and Scheduling](#)

» [Job Logs & Analytics](#)

» [Training & Tutorials](#)

» [Software](#)

» [Accounts & Allocations](#)

» [Policies](#)

» [Data Analytics & Visualization](#)

» [Data Management Policies](#)

Home » For Users » Network Connections » Connecting to NERSC with NX

## NERSC NX SERVICE - X-WINDOWS ACCELERATION AT NERSC

### Introduction

**NX** is a computer program that handles remote X Window System connections and it provides three benefits for NERSC users:

- **SPEED:** NX can greatly improve the performance of X Windows, allowing users with slow, high latency connections (e.g. on cell phone network, traveling in Africa) to use complex X Windows programs (such as rotating a plot in Matlab).
- **SESSION:** NX provides sessions that allow a user to disconnect from the session and reconnect to it at a later time while keeping the state of all running applications inside the session.
- **DESKTOP:** NX gives users a virtual desktop that's running at NERSC. You can customize the desktop according to your work requirement.

### TABLE OF CONTENTS

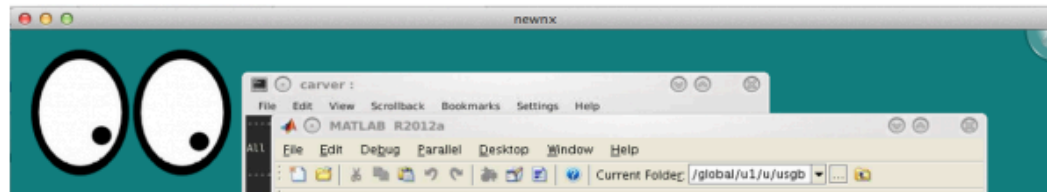
1. Introduction
2. New NX Service
3. Availability
4. Current Users (Live!)
5. Quick Start Up Tutorial
6. Got A Question?

### Related Information

[Download Tested NX Player](#)

[Download NX Configuration File](#)

[NX FAQ](#)



# Compile & Start DDT



## Compile for debugging

```
edison% make  
cc -c -g hello.c  
cc -o hello -g hello.o
```

## Set up the parallel run environment

```
edison% qsub -I -v -lmpwidth=24  
edison% cd $PBS_O_WORKDIR
```

## Start the DDT debugger

```
edison% module load ddt  
edison% ddt ./hello
```

# DDT Screen Shot



Press Go and then Pause when code appears hung.

Task 0 is at line 44

At hang, tasks are in 3 different places.

The screenshot shows the Allinea DDT v3.1-20 interface. The main window displays the following C code:

```
36 ret = MPI_Comm_size( MPI_COMM_WORLD, &totTasks );
37
38 printf("task_no is %6d of %6d total tasks\n",
39        task_no,
40        totTasks);
41
42
43 if(task_no==0) {
44     ret = MPI_Recv(&herBuffer, 50, MPI_DOUBLE, totTasks-1
45                 ret = MPI_Send(&myBuffer, 50, MPI_DOUBLE, totTasks-1,
46                 } else if (task_no==(totTasks-1)) {
47                 ret = MPI_Recv(&herBuffer, 50, MPI_DOUBLE, 0, 0, MPI_
48                 ret = MPI_Send(&myBuffer, 50, MPI_DOUBLE, 0, 0, MPI_C
49             }
50
51
52
53
```

The 'Stacks' panel at the bottom shows the following table:

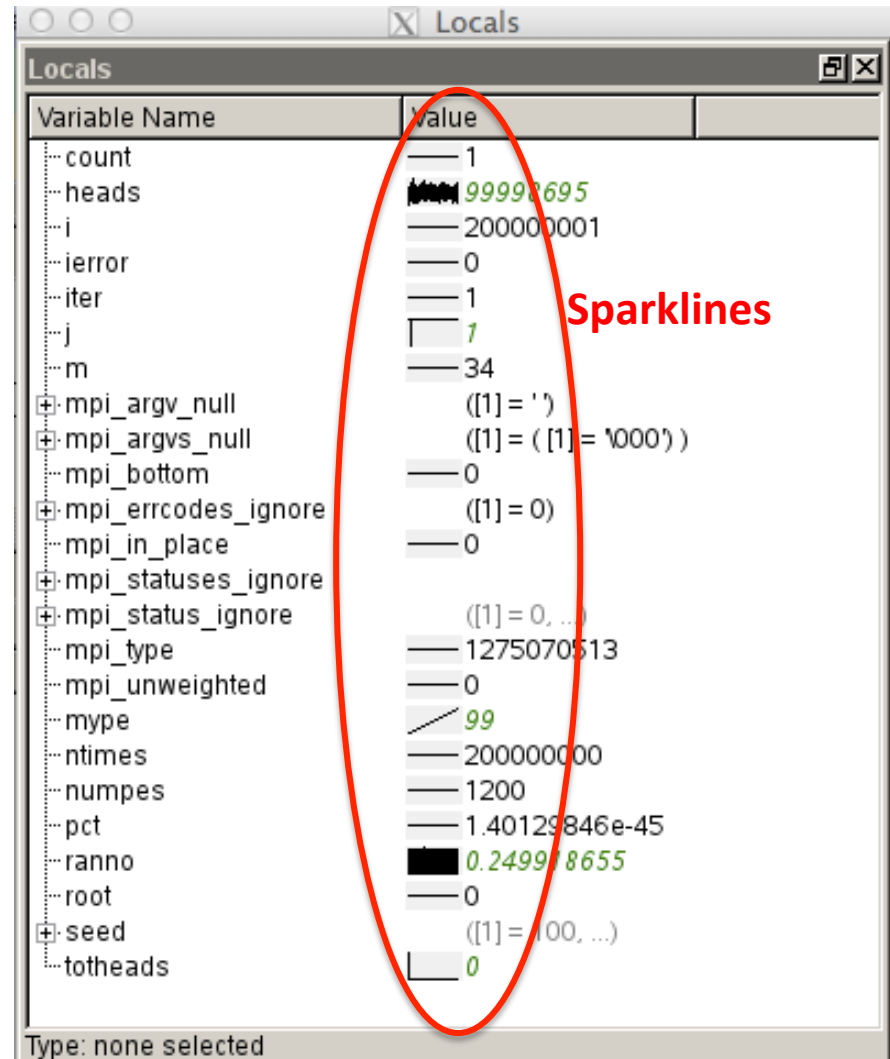
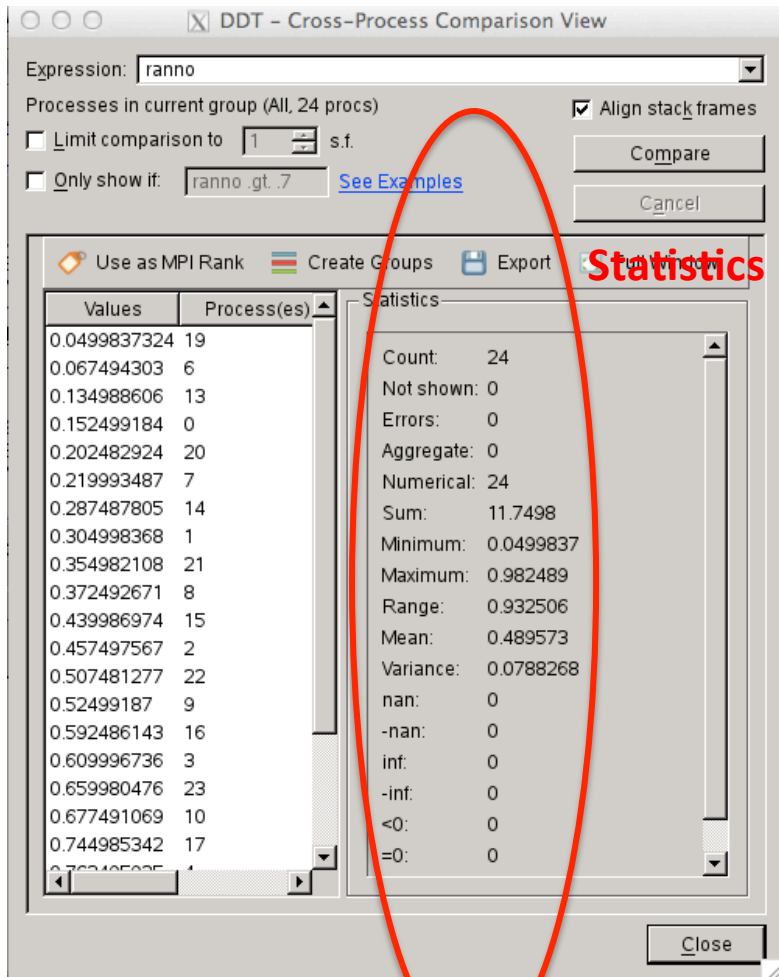
Processes	Function
1	main (hello.c:44)
1	main (hello.c:47)
2	main (hello.c:54)

# What About Massive Parallelism?



- **With 10K+ tasks/threads/streams it's impossible to examine every parallel instance in detail**
- **Make use of statistics and summaries**
- **Look for tasks that are doing something different**
  - Amount of memory used
  - Number of calculations performed (from counters)
  - Number of MPI calls
  - Wall time used
  - Time spent in I/O
  - One or a few tasks paused at a different line of code
- **We (NERSC) have been advocating for this statistical view for some time**

# Vendors are starting to listen (DDT)



# Debuggers on NERSC machines



- **Parallel debuggers with a graphical user interface**
  - DDT (Distributed Debugging Tool)
  - TotalView
- **Specialized debuggers on Hopper and Edison**
  - STAT (Stack Trace Analysis Tool)
    - Collect stack backtraces from all (MPI) tasks
  - ATP (Abnormal Termination Processing)
    - Collect stack backtraces from all (MPI) tasks when an application fails
  - CCDB (Cray Comparative Debugger)
    - Comparative debugging
- **Valgrind**
  - Suite of debugging and profiler tools

<https://www.nersc.gov/users/training/courses/CS267/> for links to recent training presentations

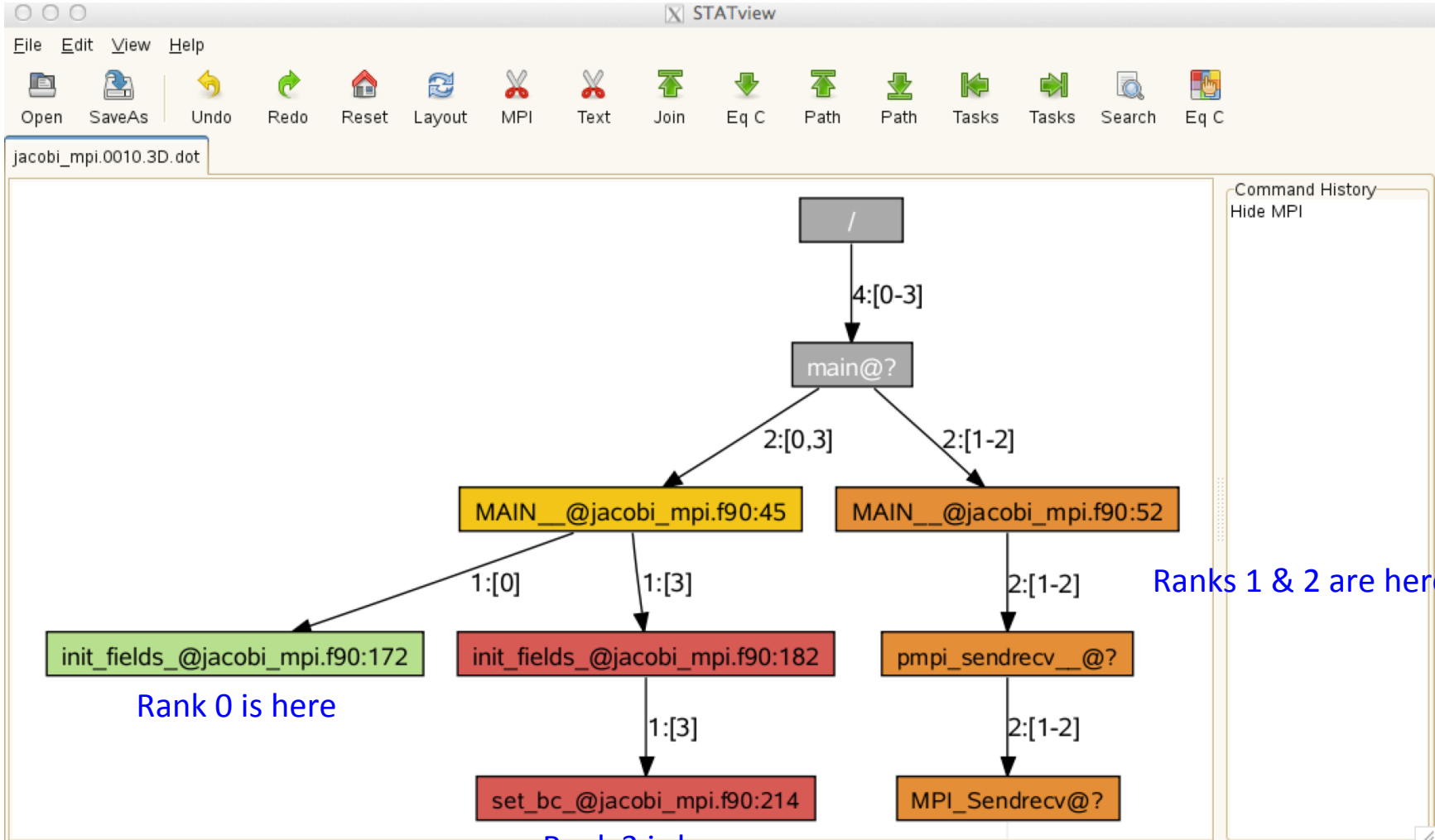


# STAT (Stack Trace Analysis Tool)



- **Gathers stack backtraces (showing the function calling sequences leading up to the ones in the current stack frames) from all (MPI) processes and merges them into a single file (\*.dot)**
  - Results displayed graphically as a call tree showing the location in the code that each process is executing and how it got there
  - Can be useful for debugging a hung application
  - With the info learned from STAT, can investigate further with DDT or TotalView
- **Works for MPI, CAF and UPC, but not OpenMP**
- **For more info:**
  - 'intro\_stat', 'STAT', 'statview' and 'statgui' man pages
  - [https://computing.llnl.gov/code/STAT/stat\\_userguide.pdf](https://computing.llnl.gov/code/STAT/stat_userguide.pdf)
  - <http://www.nersc.gov/users/software/debugging-and-profiling/stat-2/>

# Hung application with STAT (Cont'd)



# ATP (Abnormal Termination Processing)



- **ATP gathers stack backtraces from all processes of a failing application**
  - Invokes STAT underneath
  - Output in atpMergedBT.dot and atpMergedBT\_line.dot (which shows source code line numbers), which are to be viewed with statview
- **By default, the atp module is loaded on Hopper and Edison, but ATP is *not* enabled; to enable:**
  - `setenv ATP_ENABLED 1 # csh/tcsh`
  - `export ATP_ENABLED=1 # sh/bash/ksh`
- **For more info**
  - ‘intro\_atp’ man page
  - <http://www.nersc.gov/users/software/debugging-and-profiling/gdb-and-atp/>

# CCDB (Cray Comparative Debugger)



- Find a bug introduced in a version, by running two versions side by side and comparing data between them
- GUI
- Supports MPI; doesn't support threading
- For info:
  - ccdb man page and help pages
  - lgdb man page and help pages
  - *'Using the lgdb Comparative Debugging Feature'*,  
<http://docs.cray.com/books/S-0042-22/S-0042-22.pdf>
  - <http://www.nersc.gov/users/software/debugging-and-profiling/ccdb-lgdb/> (work in progress)

# Running CCDB



```
% qsub -IV -lmpwidth=48,walltime=30:00 -q debug
% cd $PBS_O_WORKDIR
% module load cray-ccdb
% ccdb
```

Request enough nodes to run two apps. simultaneously

PE set for 1<sup>st</sup> app

PE set for 2<sup>nd</sup> app

1<sup>st</sup> app

2<sup>nd</sup> app

- Run both Apps to here
- Run App-0 to here
- Restart and run both Apps to here
- Set breakpoint in both Apps
- Set breakpoint in App-0
- Clear breakpoint in both Apps
- Clear breakpoint in this App
- Build Assert

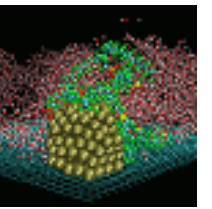
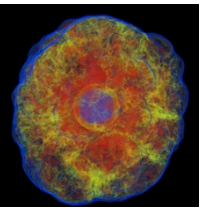
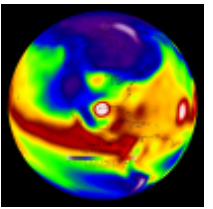
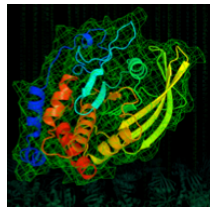
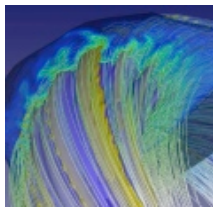
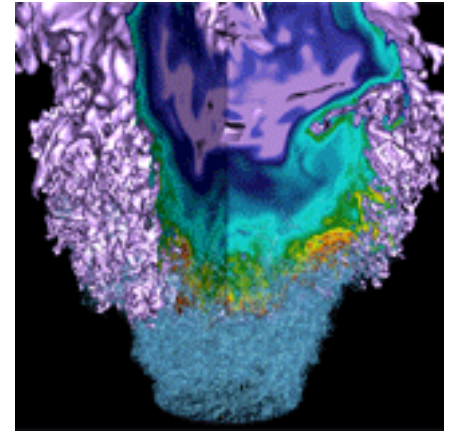
Application-0 Status: App0{0..3} Stopped hpc.c:20

Application-1 Status: App1{0..3} Stopped hpc.c:20

```
409: /*
410:  * Computes and displays norms, residuals ...
411:  */
412:  if( N <= 0 )
413:  {
414:      resid1 = HPL_rzero;
415:  }
416:  else
417:  {
418:      resid1 = resid0 / ( TEST->epsil * ( AnormI * Xnor
```

- Suite of debugging and profiler tools
- Tools include
  - **memcheck**: memory error and memory leaks detection
  - **cachegrind**: a cache and branch-prediction profiler
  - **callgrind**: a call-graph generating cache and branch prediction profiler
  - **massif, dhat (exp-dhat)**: heap profilers
  - **helgrind, drd**: pthreads error detectors
- For info:
  - <http://valgrind.org/docs/manual/manual.html>

# Performance / Optimization



- **How can we tell if a program is performing well? Or isn't? What is "good"?**
- **If performance is not "good," can we identify the causes?**
- **What can we do about it?**



# Is Your Code Performing Well?



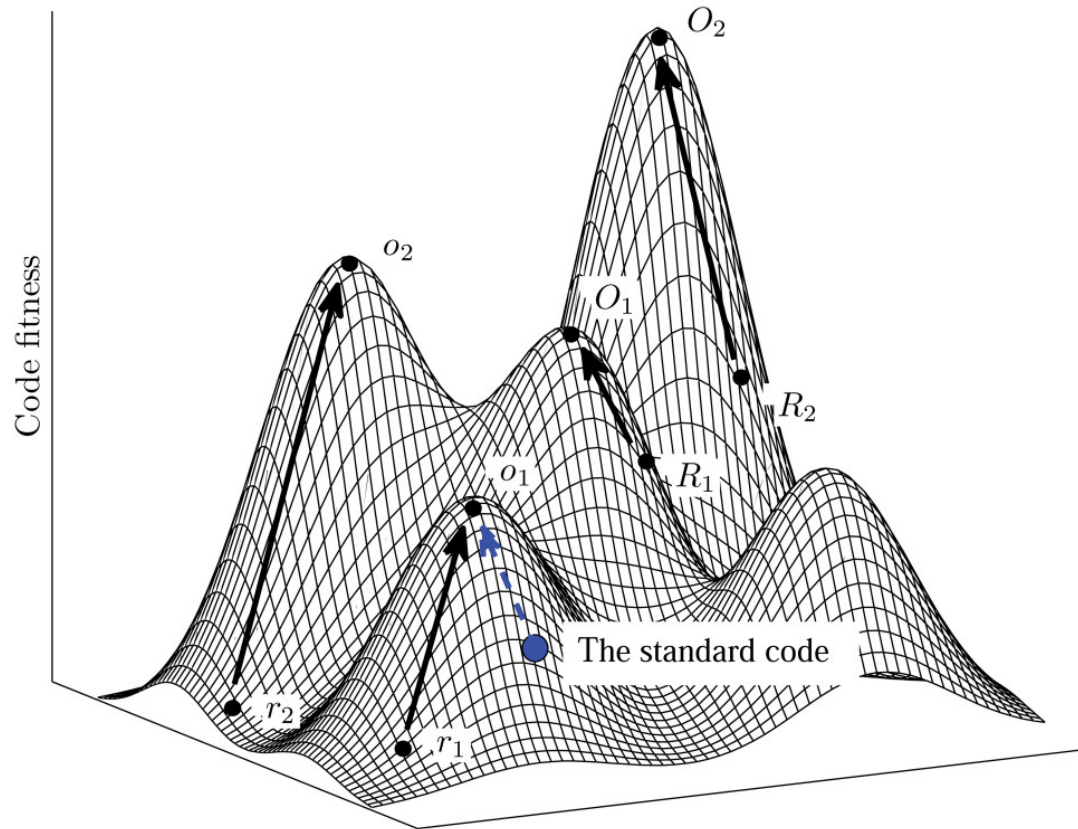
- **No single answer, but**
  - Does it scale well?
  - Is MPI time <20% of total run time?
  - Is I/O time <10% of total run time?
  - Is it load balanced?
  - If GPU code, does GPU+Processor perform better than 2 Processors?
- **“Theoretical” CPU performance vs. “Real World” performance in a highly parallel environment**
  - Cache-based x86 processors: >10% of theoretical is pretty good
  - GPUs, Xeon Phi: >few% in today’s real full HPC applications pretty good? **This your challenge!**

# What can we do about it



- Minimize latency effects (aggregate messages)
- Maximize work vs. communication
- Minimize data movement (recalculate vs. send)
- Use the “most local” memory
- Take advantage of vector (SIMD) capabilities
- Use large-block I/O
- Use a balanced strategy for I/O
  - Avoid “too many” tasks accessing a single file, but “too many” files performs poorly ~1000s
  - Use “enough” I/O tasks to maximum I/O bandwidth, but “too many” causes contention 1/node

# Performance Landscape



# Can We Identify the Causes? Use Tools



- **Vendor Tools:**
  - CrayPat on Crays
  - INTEL VTune
- **Community Tools :**
  - TAU (U. Oregon via ACTS)
  - PAPI (Performance API)
  - gprof
  - HPC Toolkit
  - Scalasca
- **NERSC “automatic” and/or easy-to-use tools**
  - e.g. IPM, Darshan

See NERSC web site  
<https://www.nersc.gov/users/software/debugging-and-profiling/>

# Example: CrayPat



- **Suite of tools that provides a wide range of performance-related information**
- **Can be used for both sampling and tracing**
  - with or without hardware or network performance counters
  - Built on PAPI
- **Supports Fortran, C, C++, UPC, MPI, Coarray Fortran, OpenMP, Pthreads, SHMEM**
- **Man pages**
  - `intro_craypat(1)`, `intro_app2(1)`, `intro_papi(1)`

- 1. Access the tools**
  - module load perftools
- 2. Build your application; keep .o files**
  - make clean
  - make
- 3. Instrument application**
  - `pat_build ... a.out`
  - Result is a new file, `a.out+pat`
- 4. Run instrumented application to get top time consuming routines**
  - `aprun ... a.out+pat`
  - Result is a new file `XXXXX.xf` (or a directory containing `.xf` files)
- 5. Run `pat_report` on that new file; view results**
  - `pat_report XXXXX.xf > my_profile`
  - view `my_profile`
  - Also produces a new file: `XXXXX.ap2` that can be viewed with apprentice GUI application

# Cray perftools and perftools-lite



- **Reports:**

- execution time
- memory high water mark
- aggregate FLOPS rate
- top time consuming user function
- MPI information
- IO information
- hardware performance counters
- load balance ...

- **Start with perftools-lite**

- **Available on Hopper and Edison.**

- **Documentation:**

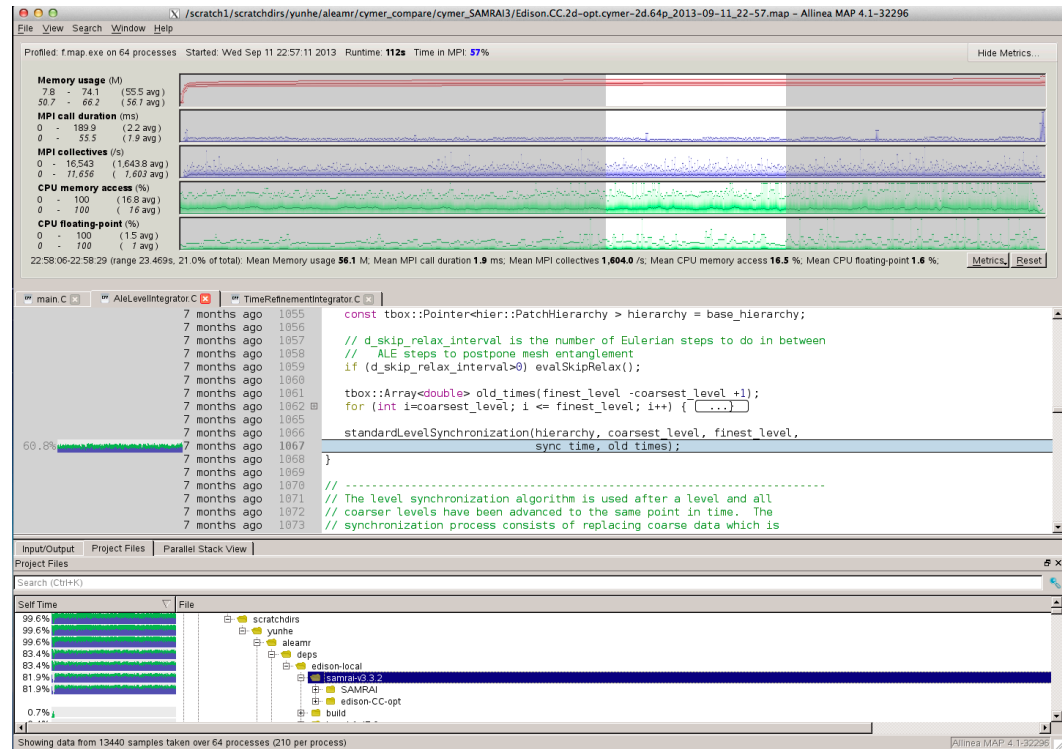
- <https://www.nersc.gov/users/software/debugging-and-profiling/craypat/>
- <http://www.nersc.gov/assets/Uploads/UsingCrayPat.pdf>
- <http://www.nersc.gov/assets/Training-Materials/UsingApprentice2012.pdf>
- <http://www.nersc.gov/assets/Uploads/Perftools-lite-2013.pdf>

```
Number of PEs (MPI ranks): 240
Numbers of PEs per Node: 24 PEs on each of 10 Nodes
Numbers of Threads per PE: 1
Number of Cores per Socket: 12
Execution start time: Sun Feb 2 13:38:33 2014
System name and speed: nid01665 2401 MHz
Wall Clock Time: 290.822940 secs
High Memory: 243.36 MBytes
MFLOPS (aggregate): Not supported (see observation below)
I/O Read Rate: 46.30 MBytes/Sec
I/O Write Rate: 5.91 MBytes/Sec
Table 1: Profile by Function Group and Function (top 10 functions shown)
100.0% | 28484.6 | -- | -- | Total
|-----|
| 61.8% | 17598.4 | -- | -- | USER
| |-----|
| | 36.3% | 10328.2 | 58.8 | 0.6% | decompmod_initdecomp_
...
| |=====|
| 29.6% | 8432.1 | -- | -- | MPI
| |-----|
| | 9.0% | 2571.0 | 129.0 | 4.8% | MPI_GATHERV
```

# Allinea MAP



- Allinea MAP is a parallel MPI profiler with GUI, small overhead.
- Reports: Memory usage, MPI usage, CPU time, CPU instructions, I/O, etc. as a function of time.
- Available on Hopper, Edison, and Carver.



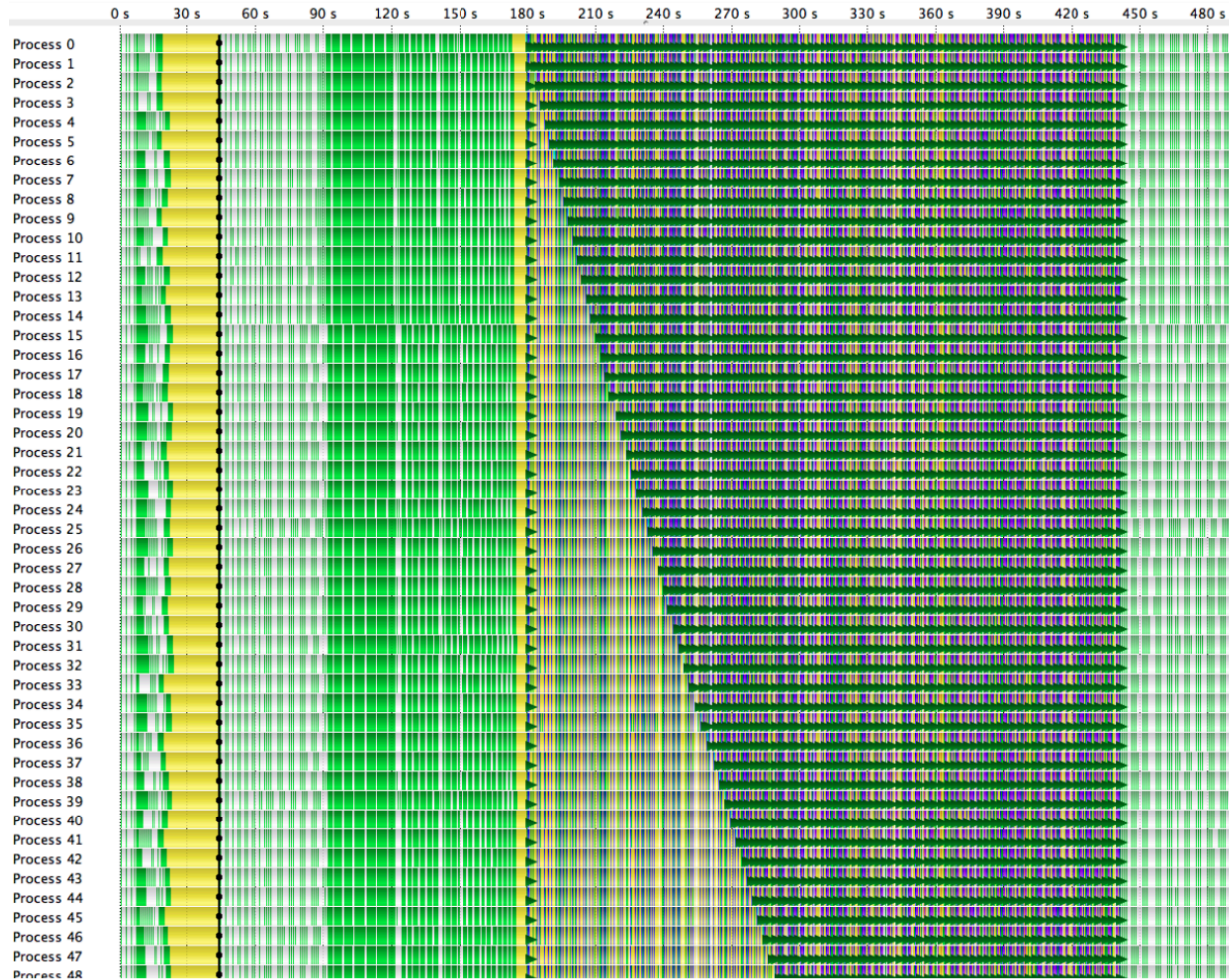
## Documentation:

<http://www.nersc.gov/users/software/debugging-and-profiling/MAP/>

<http://www.allinea.com/products/map/>



# Vampir and Vampirtrace



- **Using even the best tools can be tedious**
  - “Follow these 10 steps to perform the basic analysis of your program” – from a supercomputer center web site for a well-known tool
- **NERSC wants to enable easy access to information that can help you improve your parallel code**
  - **automatic** data collection
  - provide useful tools through the web
- **Efforts**
  - Work with vendors (e.g., CRAY ARU, Allinea Perf. Report)
  - IPM (MPI profiling, chip HW counters, memory used)
  - Accounting & UNIX resource usage
  - System-level I/O monitoring
  - User-level I/O profiling (Darshan)

# NERSC Completed Jobs



Show  entries Search:

#	Host	JobID	Job Name	User	Nds	Complete	Wall hrs
0	edison	<a href="#">782557</a>	my_job	hdixit	11	02/17/14 13:14	0.160
1	edison	<a href="#">782431</a>	cp_ref_big	hsinyu	13	02/17/14 13:14	0.849
2	edison	<a href="#">782591</a>	tvsoi	wangyu	2	02/17/14 13:14	0.005
3	edison	<a href="#">782560</a>	my_job	hdixit	11	02/17/14 13:13	0.131
4	edison	<a href="#">719618</a>	nacl_big	carnevav	16	02/17/14 13:12	11.453
5	edison	<a href="#">782403</a>	slab-Ni-pa... <a href="#">[Full Name]</a>	hdixit	11	02/17/14 13:12	1.006
6	edison	<a href="#">782590</a>	NL_test	khkim	86	02/17/14 13:12	0.008

APID	Command	Nodes	Tasks	Threads per Task	Tasks per Node	Max Task Mem (MB)	Run Time (secs)	St
3140812	<code>gribmeanp.x</code>	3	64	1	24	235	5	
	<b>Command line:</b>	<code>/opt/cray/alps/default/bin/aprun -n 64 -N 24            /scratch1/scratchdirs/whitaker/edison/newstuff/bin/gribme            /scratch1/scratchdirs/whitaker/gfsenkf_t126_1999iau_biasc</code>						
	<b>Node list:</b>	2628-2630						

# IPM: An Easy to Use Performance Tool



Just load the module, relink, and run.

```
# host      : s05601/006035314C00_AIX
# start    : 11/30/04/14:35:34
# stop     : 11/30/04/14:36:00
# gbytes   : 6.65863e-01 total

#                               [total]
# wallclock      953.272
# user           837.25
# system         60.6
# mpi            264.267
# %comm          27.7234
# gflop/sec      2.33478
# gbytes         0.665863
# PM_FPU0_CMPL   2.28827e+10
# PM_FPU1_CMPL   1.70657e+10
# PM_FPU_FMA     3.00371e+10
# PM_INST_CMPL   2.78819e+11
# PM_LD_CMPL     1.25478e+11
# PM_ST_CMPL     7.45961e+10
# PM_TLB_MISS    2.45894e+08
# PM_CYC         3.0575e+11

#                               [time]
# MPI_Send       188.386
# MPI_Wait       69.5032
# MPI_Irecv      6.34936
# MPI_Barrier    0.0177442
# MPI_Reduce     0.00540609
# MPI_Comm_rank 0.00465156
# MPI_Comm_size 0.000145341

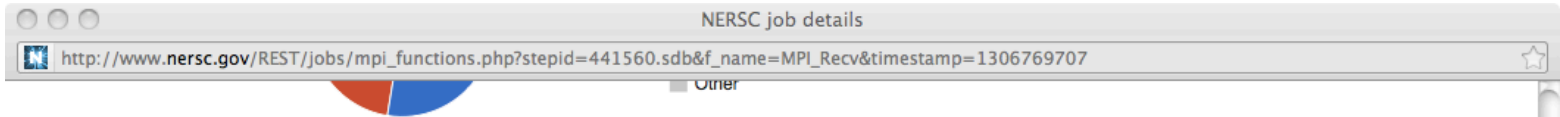
mpi_tasks : 32 on 2 nodes
wallclock : 29.97514 sec
%comm     : 27.72
gflop/sec : 2.33478e+00 total

<avg>          min          max
29.7897        29.6092        29.9752
26.1641        25.71          26.92
1.89375        1.52           2.59
8.25834        7.73025        8.70985
27.7234        25.8873        29.3705
0.0729619     0.072204       0.0745817
0.0208082     0.0195503      0.0237541
7.15084e+08   7.07373e+08    7.30171e+08
5.33304e+08   5.28487e+08    5.42882e+08
9.3866e+08    9.27762e+08    9.62547e+08
8.71309e+09   8.20981e+09    9.21761e+09
3.92118e+09   3.74541e+09    4.11658e+09
2.33113e+09   2.21164e+09    2.46327e+09
7.68418e+06   6.98733e+06    2.05724e+07
9.55467e+09   9.36585e+09    9.62227e+09

#                               [calls]
#                               <%mpi>
#                               <%wall>
# MPI_Send       639616          71.29          19.76
# MPI_Wait       639616          26.30          7.29
# MPI_Irecv      639616          2.40           0.67
# MPI_Barrier    32             0.01           0.00
# MPI_Reduce     32             0.00           0.00
# MPI_Comm_rank 32             0.00           0.00
# MPI_Comm_size 32             0.00           0.00
```



# IPM Examples



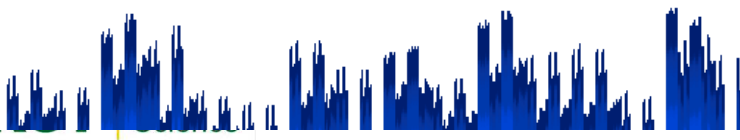
Time spent by each task in *MPI\_Recv* as a percentage of the maximum value

The MPI rank represented by each cell in the table is the sum of the cell's column and row indices.

Table Columns:

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31
0	74	76	69	70	77	77	70	71	60	62	60	61	64	65	65	64	78	79	72	72	78	79	73	73	63	64	63	64	66	65	66	66
32	70	71	63	65	72	72	65	66	52	54	53	55	56	57	58	58	71	73	67	68	72	73	67	69	55	55	54	55	58	58	57	57
64	91	92	88	90	91	93	90	91	76	77	74	76	79	81	79	79	95	98	92	96	98	98	96	96	78	81	77	80	84	84	83	82
96	86	87	80	80	86	89	81	82	62	63	61	61	65	67	65	65	91	94	83	86	94	94	87	86	64	66	63	65	70	69	69	68
128	69	71	64	65	70	72	65	66	54	56	53	55	60	60	59	59	69	70	64	65	69	70	65	66	56	57	55	56	60	59	58	58
160	65	67	57	59	67	68	60	61	42	43	43	44	47	46	47	47	66	67	58	60	67	67	60	60	43	44	45	46	46	46	48	48
192	86	87	80	81	88	89	82	83	67	68	66	67	71	72	70	70	85	88	81	84	89	88	84	84	68	71	66	69	74	73	72	71
224	79	79	71	71	80	80	72	72	49	49	49	50	52	52	53	52	76	79	70	73	79	79	73	74	50	52	48	52	55	55	54	54
256	83	84	81	83	85	85	84	85	69	71	69	70	74	75	74	74	85	86	85	87	86	87	86	87	73	73	71	71	76	75	74	73
288	72	73	66	68	73	74	69	70	60	62	59	60	64	65	63	63	75	76	69	71	76	76	71	71	63	63	61	61	65	65	64	64
320	94	95	93	94	93	96	94	95	77	78	75	75	80	81	78	78	99	99	96	96	99	99	98	97	79	80	77	78	83	82	81	80
352	82	83	74	75	83	84	75	76	61	62	66	66	65	65	69	69	84	85	76	77	85	85	77	78	63	64	66	67	66	66	69	69
384	84	86	76	78	86	88	79	80	64	64	63	64	68	68	67	67	85	86	77	78	85	86	78	79	65	65	66	66	68	67	69	68
416	69	70	60	62	71	72	63	64	45	46	46	47	49	50	50	51	68	69	61	62	69	70	62	62	47	47	46	47	49	49	47	49
448	98	98	99	100	99	98	100	100	85	85	78	77	89	89	81	80	94	94	95	96	94	94	97	97	87	87	78	79	91	90	82	81
480	82	82	72	73	82	84	74	74	54	54	52	52	56	57	54	55	83	83	72	73	84	84	71	73	54	55	53	54	57	55	56	56
512	88	89	89	91	90	91	92	93	78	80	78	80	83	84	84	84	93	94	93	95	94	94	95	94	80	81	82	83	85	84	86	85
544	90	91	79	81	91	93	82	83	63	65	63	64	67	69	68	68	96	97	84	84	98	97	85	85	66	66	66	67	69	68	70	69
576	76	77	74	75	76	77	76	76	65	65	66	67	67	67	70	70	79	82	77	80	83	82	81	81	67	69	69	71	72	72	74	73
608	76	77	69	69	76	78	70	71	57	57	57	58	60	59	60	60	79	81	72	75	82	81	75	75	60	62	58	61	64	63	63	62
640	92	94	85	87	96	96	90	89	75	77	75	77	81	81	81	81	93	94	87	88	94	94	88	88	77	78	76	77	81	80	80	79
672	86	87	78	78	88	88	80	80	57	58	58	59	62	62	63	63	87	87	78	79	88	87	79	79	59	59	60	60	62	61	63	63
704	78	79	74	75	80	80	76	75	63	64	63	63	68	67	67	66	78	81	74	77	81	81	77	77	64	67	63	66	70	69	69	67
736	70	71	65	65	71	71	66	67	52	52	50	51	55	54	53	53	71	73	64	67	74	73	68	68	53	55	52	54	57	57	57	56
768	94	96	87	88	96	98	90	90	71	73	70	71	75	77	75	75	95	95	87	87	96	96	87	88	74	75	73	74	79	78	77	77
800	79	80	70	71	81	83	72	73	57	58	62	63	62	62	67	67	80	79	71	71	80	80	71	72	60	60	64	64	63	63	66	67
832	82	83	77	77	82	84	78	78	63	64	64	64	66	66	67	66	86	86	79	80	87	86	81	80	65	66	64	65	69	68	68	67
864	69	70	65	65	70	70	65	66	56	56	56	56	59	59	58	58	73	73	67	67	73	73	67	67	58	59	58	59	61	59	61	59
896	92	92	85	86	93	94	88	88	71	72	69	70	76	75	74	74	92	92	87	88	94	93	88	89	73	73	70	71	76	76	74	74
928	73	74	64	65	75	76	66	67	43	44	48	49	47	48	53	54	74	73	64	65	73	74	64	66	45	45	51	52	47	48	53	54
960	74	74	70	70	75	75	71	71	63	64	61	61	66	66	64	63	76	77	72	73	77	77	73	73	63	64	62	62	66	66	66	64
992	63	63	57	57	63	64	58	58	41	41	42	42	43	42	44	43	61	61	56	57	62	62	56	57	41	41	42	42	43	42	43	43

Time vs. MPI Rank for *MPI\_Recv*



# Cori – NERSC's next supercomputer system



- NERSC will be installing a 30-petaflop/s Intel KNL-based Cray system in the 2016 time frame, named after American Biochemist Gerty Cori
- Over 9,300 single-socket nodes in the system with each node > 3TeraFLOPS/s theoretical peak performance
- Cray Aries high speed "dragonfly" topology interconnect
- Cabinets, Liquid cooling
- Lustre filesystem with > 430 GB/sec I/O bandwidth and 28 PB of disk capacity; 1-2 TB/sec Burst Buffer bandwidth



Gerty Cori

# Intel Xeon Phi Knights Landing Processor overview



- **Single socket node with NUMA**
- **Greater than 60 cores per node with up to 4 hardware threads each**
- **AVX512 Vector pipelines with a hardware vector length of 512 bits (eight double-precision elements)**
- **Better performance per watt than previous generation Xeon Phi™ systems and 3X single-thread performance**
- **64-128 GB of DRAM memory per node**
- **On-package, high-bandwidth memory, up to 16GB capacity with bandwidth projected to be 5X that of DDR4**



# NESAP Tools Partners



- HPCToolkit
  - <http://hpctoolkit.org/>
- TAU
  - <http://www.cs.uoregon.edu/research/tau/home.php>
- Open | Speedshop
  - <http://www.openspeedshop.org/wp/>
- PAPI
  - <http://icl.cs.utk.edu/papi/>
- MAP (Allinea)
  - <http://www.allinea.com/products/map>
- Vampir
  - <https://www.vampir.eu/>
- PerfExpert
  - <https://www.tacc.utexas.edu/research-development/tacc-projects/perfexpert>
- Scalasca
  - <http://www.scalasca.org/>
- VTUNE
  - <https://software.intel.com/en-us/intel-vtune-amplifier-xe>
- TotalView
  - <http://www.roguewave.com/products-services/totalview>
- DDT (Allinea)
  - <http://www.allinea.com/products/ddt>

# Challenges with debugging and profiling on Cori



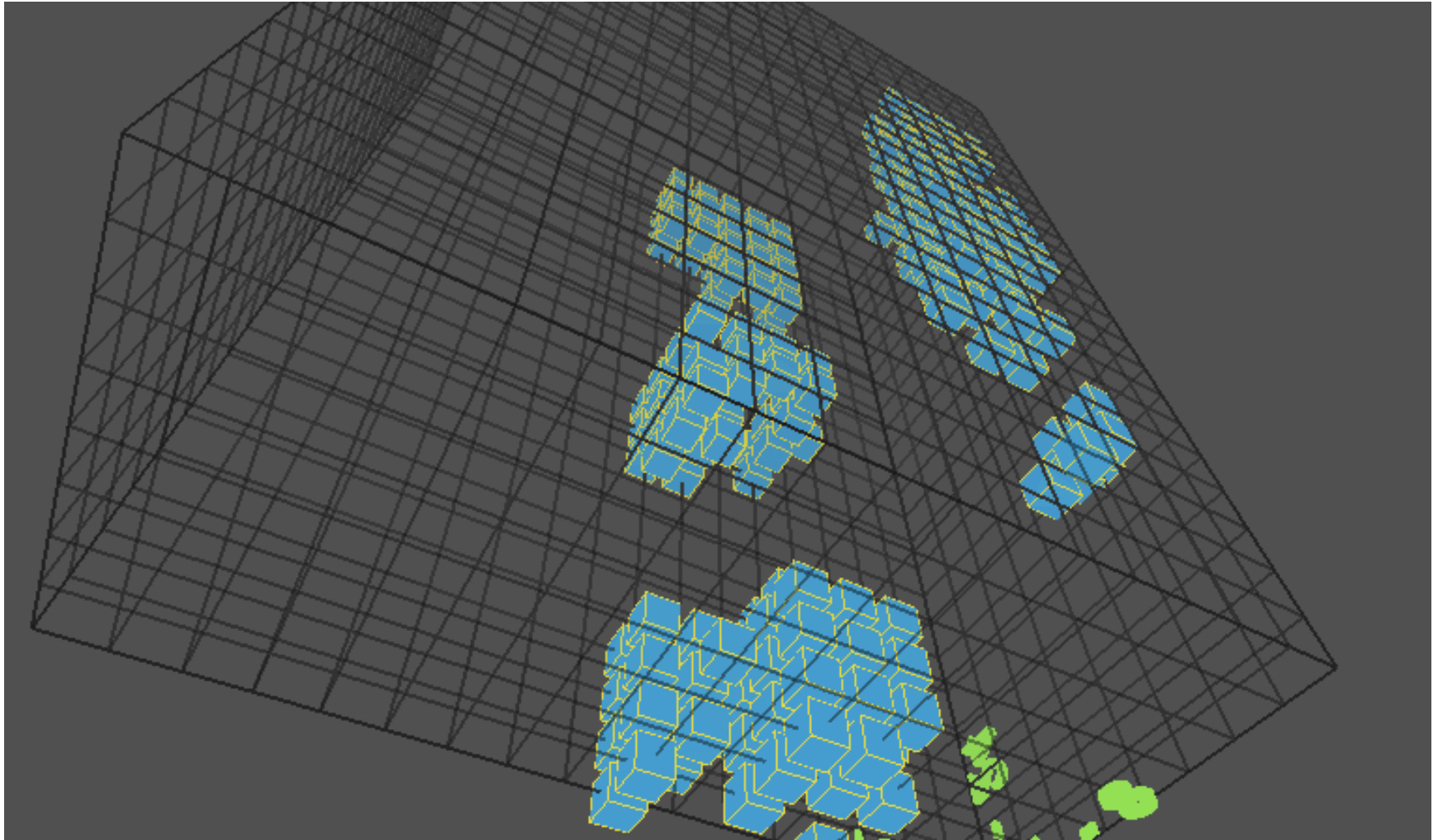
- **Debugging and performance tools need to scale up to thousands (or more?) of threads and tasks**
- **NESAP participants will want to know**
  - Hotspots
  - Memory usage, bandwidth and errors with on-package memory
  - Cache utilization
  - MPI communication performance - same interconnect fabrics for more powerful nodes
  - Threading performance
  - I/O statistics
  - ...

- **Debugging and Parallel Code Optimization can be hard**
- **Tools can help**
  - See NERSC web pages for recommendations
  - Use the ones that work for you
- **Be aware of some of the more common errors and best practices**
- **Look for outliers in parallel programs**
- **Refer to NERSC web pages for details**
  - <http://www.nersc.gov/users/software/debugging-and-profiling/>
  - <http://www.nersc.gov/users/training/courses/CS267/>

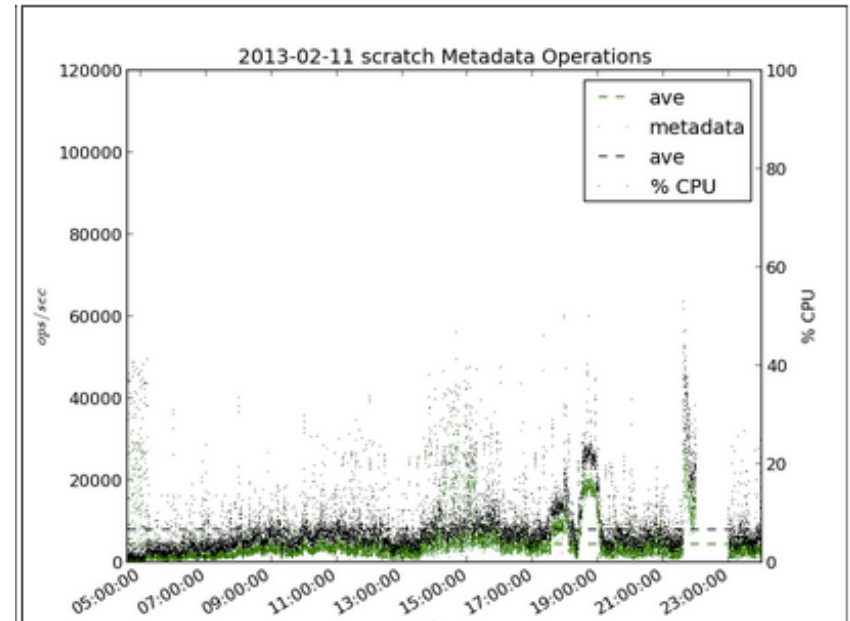
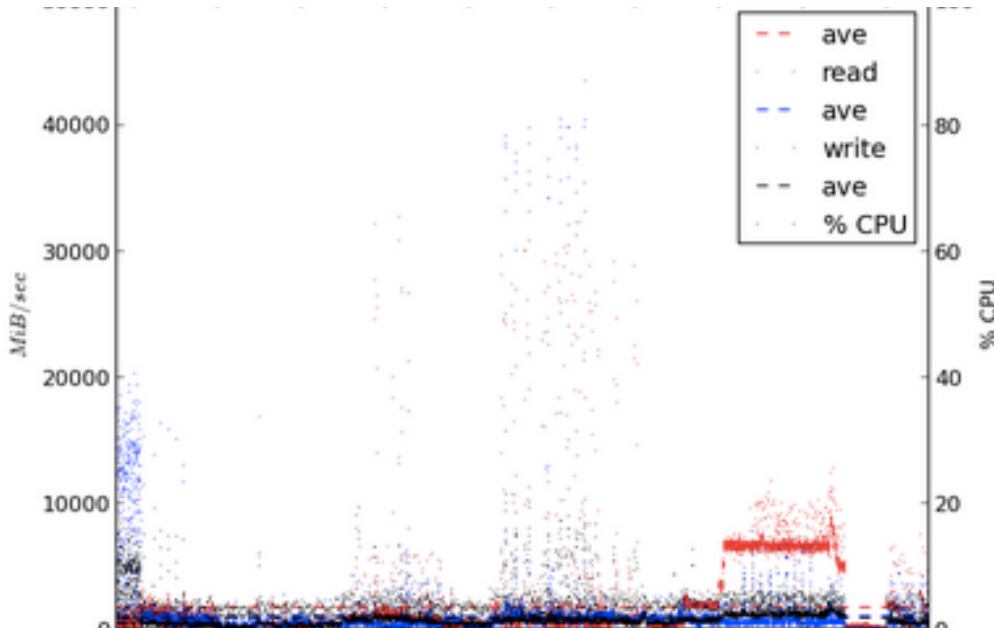


**National Energy Research Scientific Computing Center**

# Job Physical Topology



# System-Level I/O Monitoring



Users can see the system-wide I/O activity while their job ran to look for contention.

# IPM Examples



NERSC job details

[http://www.nersc.gov/REST/jobs/mpi\\_functions.php?stepid=619349.sdb&f\\_name=MPI\\_Allreduce&timestamp=1310766809](http://www.nersc.gov/REST/jobs/mpi_functions.php?stepid=619349.sdb&f_name=MPI_Allreduce&timestamp=1310766809)

Time spent by each task in *MPI\_Allreduce* as a percentage of the maximum value

The MPI rank represented by each cell in the table is the sum of the cell's column and row indices.

Table Columns: 48

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47		
0	100	100	100	99	99	98	98	98	98	97	97	97	97	97	97	96	94	95	94	94	94	93	93	93	92	92	92	92	92	92	92	92	89	89	88	88	87	87	87	87	87	86	86	86	86	86	86			
48	82	82	81	81	80	81	80	80	79	79	79	79	79	80	79	79	79	78	78	77	76	77	76	76	76	76	77	76	76	76	76	73	72	73	71	71	71	71	71	70	70	70	70	70	69	71	70	69		
96	76	77	76	75	75	75	75	74	74	74	75	74	74	75	74	74	75	74	73	73	74	73	72	72	72	72	72	72	73	72	71	69	69	69	68	68	68	68	67	67	68	67	67	67	67	67	67	67		
144	67	66	66	65	64	65	64	64	64	64	64	64	64	64	64	64	63	55	54	54	53	53	54	53	53	53	54	54	53	53	53	52	52	51	53	51	51	51	50	50	50	50	51	52	51	50	50			
192	50	50	50	49	49	50	49	49	49	50	49	49	50	50	50	50	49	49	48	48	48	48	47	47	47	47	47	49	48	50	51	51	50	50	50	50	50	50	50	50	50	51	52	52	51	52	51			
240	51	51	50	50	50	51	50	49	50	49	49	50	51	51	51	58	58	57	57	57	57	57	57	57	57	57	57	57	59	58	56	56	55	56	56	56	56	56	55	56	56	57	57	58	58	58	58			
288	58	58	58	58	57	58	58	57	58	58	57	58	58	60	59	59	56	56	56	56	55	56	57	56	56	58	57	57	58	59	59	60	58	58	59	58	58	58	58	58	58	58	58	58	59	60	60	61		
336	62	61	62	62	62	63	62	63	63	64	63	64	65	65	66	65	65	65	65	65	65	65	65	65	65	65	65	65	66	67	68	61	61	62	61	61	62	62	62	62	63	63	63	63	65	65	66			
384	57	56	56	56	57	57	57	58	58	59	60	61	60	59	60	59	60	60	60	60	60	60	60	60	60	61	61	61	61	63	62	63	69	69	69	69	69	69	69	69	69	69	69	69	69	69	69	69		
432	80	80	80	79	80	80	81	80	80	81	81	81	81	82	82	83	77	76	76	76	76	76	76	77	77	77	77	78	78	79	80	78	78	78	78	78	78	78	78	79	79	79	79	80	80	81	82	82		
480	75	75	75	75	75	75	76	75	76	76	76	76	77	78	78	78	79	79	79	79	79	79	79	80	80	80	80	81	82	83	85	92	92	92	92	91	91	91	91	91	90	90	90	90	90	90	90	90		
528	86	87	86	86	85	85	85	84	84	84	84	83	83	83	83	85	84	84	83	83	82	82	82	82	82	82	82	82	81	78	77	77	76	76	76	76	75	75	75	75	75	75	75	75	75	75	75	75		
576	73	73	73	71	71	72	71	71	71	71	71	71	71	71	71	68	67	67	66	66	66	66	65	65	65	64	65	64	66	65	64	68	68	68	68	66	66	66	66	66	66	66	66	66	65	65	66	65	65	
624	65	65	65	64	64	64	64	63	63	63	63	62	62	55	55	55	54	54	53	53	53	53	53	53	53	53	53	53	53	51	50	51	50	49	50	49	49	49	49	49	48	48	49	50	49	48	48			
672	44	44	44	43	43	43	43	43	43	43	43	43	43	43	43	42	42	42	42	42	41	41	41	41	41	41	41	41	41	41	41	41	41	41	41	41	41	41	41	41	41	41	41	41	41	41	41	41		
720	41	40	41	40	39	40	40	40	38	39	39	39	39	40	40	39	44	44	45	43	44	44	44	44	44	44	44	44	44	45	45	46	45	44	44	44	44	44	44	44	43	43	44	43	43	44	45	44		
768	54	54	53	53	53	54	54	54	54	54	54	55	55	54	51	52	52	51	51	51	52	51	52	52	52	52	52	52	52	53	53	53	53	53	53	53	53	53	53	52	53	53	53	52	53	53	55	54	54	
816	51	52	52	52	52	52	52	54	53	54	55	54	55	54	54	54	54	54	54	54	54	54	54	54	54	54	54	55	55	55	55	56	56	57	58	58	58	58	59	59	59	60	60	61	60	61	61	61	62	
864	61	60	61	61	61	61	61	61	61	61	61	61	61	62	63	63	64	59	58	59	59	59	59	60	60	60	61	61	61	61	63	62	63	55	55	56	56	56	56	56	56	56	56	56	57	57	57	58	58	59
912	58	58	59	59	58	59	59	60	60	60	61	62	62	62	68	68	68	68	68	68	68	68	68	68	68	69	69	69	69	70	71	71	72	79	80	80	79	80	80	81	80	80	81	81	81	82	82	83		
960	75	74	75	75	75	75	75	76	76	76	76	77	77	78	78	77	77	77	77	77	78	78	78	78	78	78	78	79	79	80	81	81	74	74	74	74	74	74	74	74	75	74	75	75	75	75	76	77	77	
1008	78	78	79	79	79	79	79	80	80	80	80	81	83	84	86	92	92	91	91	91	91	91	91	91	90	90	90	89	90	89	84	85	84	83	82	81	81	81	81	80	80	80	80	80	80	80	80	80		
1056	82	82	81	81	81	80	80	80	79	79	80	80	80	80	80	75	75	74	74	73	74	73	73	72	72	72	72	72	72	71	71	71	70	70	70	69	69	69	69	69	69	69	69	69	69	69	69	69		
1104	67	66	67	65	65	65	64	64	64	64	64	64	65	64	64	69	69	69	68	67	68	67	67	67	67	67	67	66	66	67	66	66	64	64	64	63	63	63	62	62	62	62	62	62	62	62	62	62	62	
1152	55	55	55	54	53	53	53	52	53	53	52	53	53	53	51	50	51	50	49	49	49	49	49	49	49	49	48	48	48	49	49	48	43	43	43	42	42	42	42	42	42	42	42	42	42	42	42	42	42	
1200	43	43	42	43	41	41	41	41	41	41	41	41	41	41	41	42	42	42	42	42	42	42	42	42	42	42	42	42	42	42	42	42	41	41	41	41	41	41	41	41	41	41	41	41	41	41	41	41	41	
1248	43	43	43	42	42	43	42	42	42	43	43	42	43	44	44	43	43	43	42	42	43	42	42	42	42	42	42	42	42	42	43	43	43	53	54	53	53	53	53	53	53	53	53	54	54	54	53	55	55	54
1296	51	51	51	50	51	50	51	50	51	50	51	50	51	52	52	50	51	50	50	50	50	50	50	50	50	50	51	51	50	51	53	52	52	51	52	52	51	51	52	52	52	51	53	52	54	55	54	55	54	55
1344	54	54	55	54	55	55	55	55	55	55	57	57	57	53	53	53	53	53	53	53	53	53	54	54	55	54	55	56	56	57	56	55	55	55	55	55	55	55	55	55	55	55	55	55	56	56	56	58	58	59
1392	54	54	55	54	55	55	55	55	55	56	56	56	56	58	58	59	56	56	56	56	56	57	57	57	57	57	58	57	58	59	60	62	61	62	62	63	63	63	63	63	63	64	64	64	65	65	66	66		
1440	70	70	70	70	70	71	70	71	71	71	73	73	73	73	73	80	81	81	80	81	81	82	81	82	83	83	83	84	85	78	77	77	77	77	77	77	77	77	78	78	78	78	78	78	78	79	80	80	80	
1488	79	79	79	79	79	80	80	80	81	81	81	83	83	83	77	77	78	77	77	78	77	78	77	78	78	79	79	79	80	80	79	81	81	81	82	82	82	82	82	82	83	83	83	84	84	84	85	86	88	88
1536	98	97	97	96	96	95	95	95	95	94	94	94	94	94	94	94	88	90	89	88	87	87	87	86	86	86	86	85	85	85	85	85	85	85	84	84	84	83	83	83	83	83	82							

# IPM Examples



NERSC job details  
[http://www.nersc.gov/REST/jobs/ipm\\_summary.php?stepid=627129.sdb&name=memory&timestamp=1311046935](http://www.nersc.gov/REST/jobs/ipm_summary.php?stepid=627129.sdb&name=memory&timestamp=1311046935)

## Task distribution of *Maximum Memory Usage (GBytes)* - as a percentage of maximum

The MPI rank is the sum of the column and row indices in the table.

Table Columns: 32

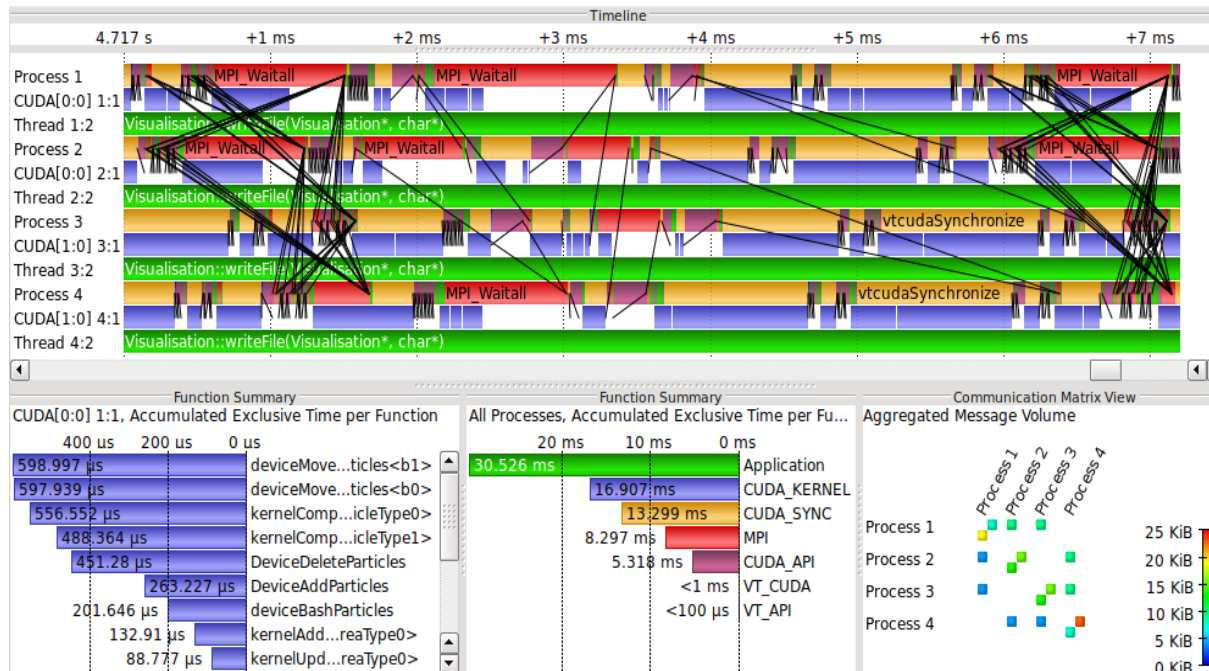
	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	
0	67	68	67	69	67	69	68	70	68	69	69	70	69	70	71	71	70	72	69	71	71	73	70	72	71	72	71	72	72	74	72	73	
32	69	71	70	72	70	72	70	72	71	72	72	73	71	72	72	73	73	75	72	74	73	75	72	74	75	76	74	75	74	75	74	75	
64	71	71	72	73	72	73	73	74	72	72	73	73	73	73	74	74	74	75	74	75	76	77	75	76	76	76	75	75	77	77	77	76	
96	74	75	75	76	74	75	75	76	75	75	76	75	75	76	75	75	76	75	78	79	77	79	78	79	77	78	79	79	78	78	79	78	
128	86	88	86	87	87	89	86	89	88	89	87	88	88	89	88	89	86	87	85	87	87	89	86	88	87	88	86	87	88	89	88	89	
160	89	91	89	91	89	91	88	90	91	92	90	91	90	91	90	91	89	90	88	90	88	90	88	90	90	91	90	91	90	91	90	91	
192	95	96	95	96	96	97	96	96	95	95	97	97	97	96	95	95	94	95	96	97	95	96	95	95	95	95	95	95	97	97	96	96	
224	98	99	98	98	99	97	98	100	99	98	98	99	99	98	97	98	99	97	98	99	97	98	98	99	97	98	99	99	98	97	99	98	97
256	70	72	71	73	70	72	71	73	72	73	72	73	71	72	73	73	75	73	74	73	75	72	74	75	76	74	75	74	76	74	75		
288	67	69	68	70	66	68	66	68	69	70	70	71	67	68	68	69	71	73	70	72	69	71	68	70	72	73	72	73	70	71	70	71	
320	74	75	76	76	74	75	75	77	75	75	76	76	75	75	76	76	78	79	77	78	78	79	77	79	80	79	79	78	79	79	78	78	
352	72	73	73	74	71	71	72	72	72	73	73	70	70	72	71	75	76	75	76	74	74	73	74	76	76	76	75	74	74	73	73	73	
384	90	92	89	91	90	91	89	91	91	93	91	92	91	92	91	91	89	91	89	90	89	91	88	90	91	92	90	91	90	91	90	91	
416	87	89	87	89	85	86	89	90	88	89	87	88	86	87	86	88	86	88	85	87	84	86	88	89	87	88	86	87	86	87	86	87	
448	99	100	98	99	99	100	98	99	100	100	99	99	100	100	99	99	99	99	98	99	98	99	98	99	99	99	99	99	99	99	99	98	
480	96	97	95	96	94	95	93	94	97	97	96	96	95	95	94	94	96	97	95	96	94	95	93	94	97	96	96	95	95	94	94	94	
512	71	72	72	73	73	74	71	72	71	72	72	73	73	74	73	75	75	74	75	76	77	75	76	74	75	73	74	75	76	75	76	75	
544	74	75	75	76	74	75	75	76	74	75	75	76	74	75	75	76	78	79	77	78	78	79	77	78	77	78	77	78	77	78	77	78	
576	68	69	70	70	69	70	72	68	67	69	68	69	68	71	69	72	73	71	72	73	74	72	74	72	70	71	70	73	72	73	71		
608	72	73	73	74	72	72	72	74	71	70	72	71	71	70	72	71	76	77	75	76	76	76	74	76	75	74	74	73	75	74	74	73	
640	94	95	94	94	95	96	94	95	94	94	93	93	94	95	94	94	94	93	94	95	96	94	95	93	94	92	93	94	95	93	95		
672	98	98	97	97	98	96	97	97	98	96	97	96	97	96	97	96	97	97	98	96	97	97	98	96	97	97	96	97	96	97	96	97	
704	88	89	87	88	89	90	89	90	88	86	86	86	89	88	88	87	88	89	87	88	89	90	89	90	88	86	87	86	89	88	88	87	
736	91	92	91	91	91	92	90	91	91	90	89	91	89	91	89	90	88	91	92	91	92	91	92	90	91	91	90	90	89	91	89	90	
768	75	76	76	77	75	75	76	76	74	75	76	74	75	76	74	75	76	78	79	78	78	78	79	77	78	78	79	77	78	77	78	77	
800	72	73	73	74	71	71	71	72	71	73	72	74	70	70	71	72	76	76	75	76	73	74	73	73	75	76	74	75	73	73	72	73	
832	72	72	73	74	71	72	72	74	71	70	72	71	71	70	72	71	75	76	75	76	75	76	75	76	75	74	75	73	75	74	74	73	
864	69	70	70	71	67	68	68	69	69	67	69	68	66	65	68	66	73	74	72	73	71	72	70	71	72	71	71	71	70	69	69	68	
896	99	99	98	98	99	97	98	98	99	97	97	97	97	98	96	97	98	99	97	98	98	99	97	98	99	97	98	97	96	98	96	97	
928	95	96	95	95	93	94	92	93	95	96	94	95	92	93	92	93	95	96	94	95	93	94	92	93	94	95	94	95	92	93	92	93	
960	92	93	91	92	91	92	91	92	92	90	91	90	92	90	91	89	92	93	91	92	91	92	91	92	91	92	90	91	90	91	90	89	
992	89	90	89	89	87	88	87	88	89	87	88	87	87	86	86	85	89	90	88	90	88	88	87	88	89	87	88	87	87	85	86	84	
1024	87	88	86	87	87	90	87	89	88	89	87	88	89	90	88	89	86	88	85	87	87	89	86	89	87	88	86	87	88	90	88	89	
1056	90	92	89	91	89	91	89	90	91	92	91	92	91	92	90	91	89	91	88	90	89	91	88	90	91	92	90	91	90	91	90	91	
1088	95	96	95	95	96	97	96	97	96	96	95	95	97	97	97	96	95	96	94	95	96	97	96	96	95	96	95	95	97	97	96	96	
1120	98	99	98	99	98	99	97	98	100	99	99	99	99	99	98	98	98	99	97	98	98	99	97	98	98	99	97	98	99	98	98	97	
1152	71	73	71	73	73	74	72	74	73	74	72	74	74	75	74	74	69	71	69	70	70	72	69	71	71	72	70	71	72	73	71	72	
1184	75	77	74	76	75	76	74	76	76	77	75	77	76	77	76	76	72	74	71	73	72	74	71	73	74	75	73	74	74	75	73	74	
1216	76	77	76	76	78	79	76	78	77	77	76	76	78	78	77	74	75	73	74	75	76	75	75	75	75	74	73	76	76	75	74		
1248	80	80	79	80	79	81	79	80	81	81	80	80	81	81	79	80	77	78	77	78	77	78	77	78	79	79	77	77	78	78	77	76	
1280	90	92	90	92	90	92	89	91	92	93	91	92	91	92	91	91	89	92	89	91	89	92	89	91	91	92	90	92	91	92	90	91	
1312	87	89	87	89	86	87	85	86	89	90	88	90	87	88	86	87	89	86	88	85	87	85	86	88	89	87	89	87	88	86	87		
1344	99	100	98	100	99	100	98	99	100	100	99	99	100	100	99	99	99	99	98	99	98	99	98	99	98	99	100	100	99	99	99	98	
1376	96	97	95	96	94	95	93	94	97	97	96	96	95	95	94	94	96	97	95	96	94	95	93	94	97	97	96	96	95	95	94	94	
1408	75	77	74	77	75	77	74	76	77	78	76	77	76	77	76	77	73	75	72	74	73	74	72	74	75	76	74	75	74	75	74	74	
1440	72	75	72	74	71	72	70	72	74	75	74	74	72	73	72	72	70	72	69	71	68	70	67	70	72	72	71	72	70	71	69	70	
1472	80	81	80	80	81	79	80	81	81	79	80	81	80	80	80	79	78	79	77	78	77	79	77	78	79	77	77	78	79	77	77	77	
1504	77	78	77	78	75	76	75	76	79	78	77	77	77	76	76	75	75	76	75	75	73	74	72	73	76	76	75	75	74	74	73	72	
1536	94	95	94	94	95	96	94	95	94	94	93	94	95	96	94	95	94	95	93	94	95	96	94	95	93	94	92	93	94	95	94	95	



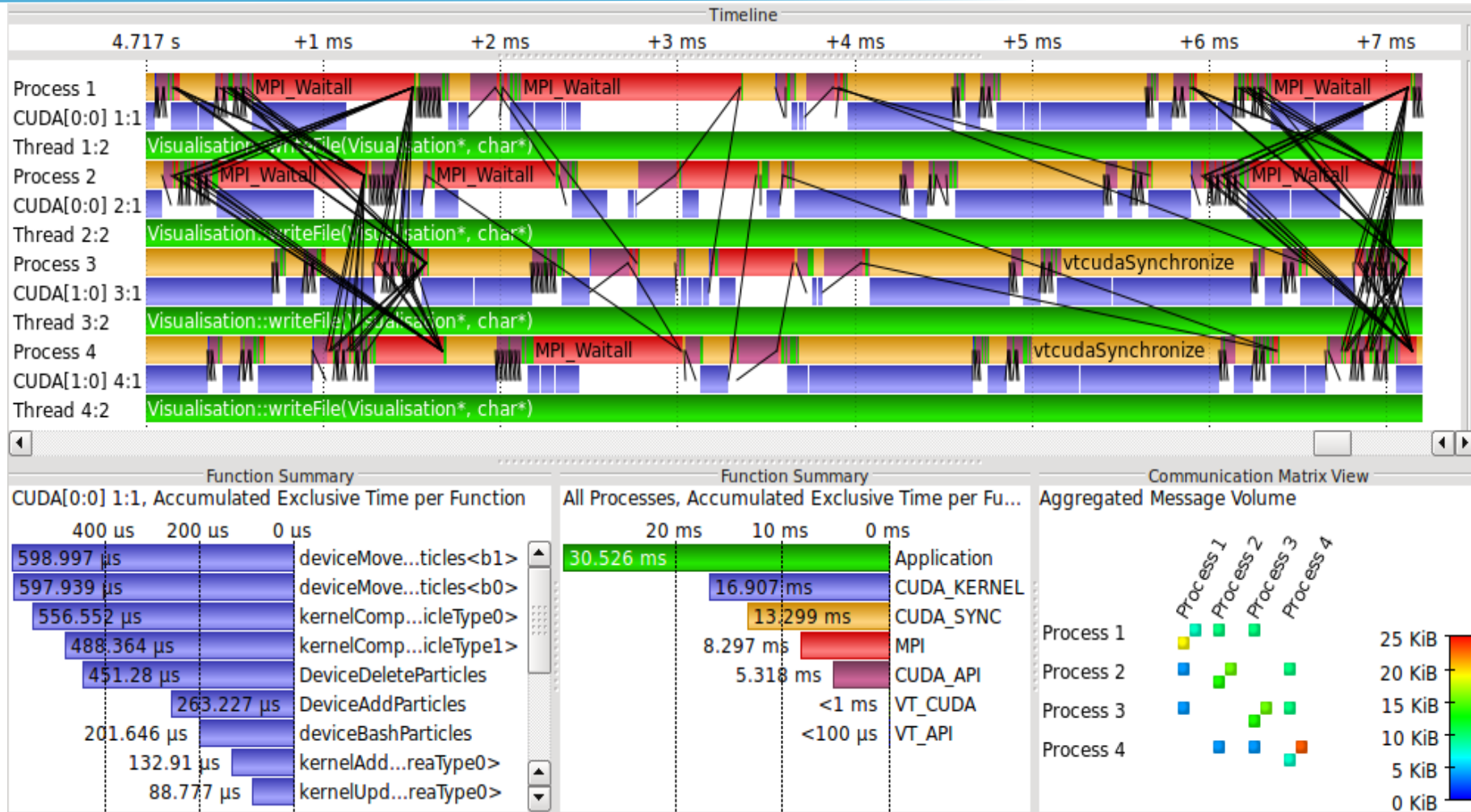
# Parallel Tools for the Masses



- Using even the best tools can be tedious
  - “Follow these 10 steps to perform the basic analysis of your program”



# Vampir w/ CUDA



# Compiler runtime bounds checking



Out of bounds reference  
in source code for  
program “flip”

```
...  
  
allocate(put_seed(random_size)  
)  
  
...  
  
bad_index = random_size+1  
put_seed(bad_index) = 67
```

Intel compiler:

```
ftn -c -g -Ktrap=fp -check bounds flip.f90  
ftn -c -g -Ktrap=fp -check bounds printit.f90  
ftn -o flip flip.o printit.o -g
```

```
% qsub -I -qdebug -lmpwidth=16  
% cd $PBS_O_WORKDIR  
%  
% aprun -n 16./flip
```

```
forrtl: severe (408): fort: (2): Subscript #1  
of the array SEED has value 3 which is greater  
than the upper bound of 2
```

- **For a list of compiler options, see the man pages for the individual compilers**
  - man pgcc | pgCC | pgf90
  - man icc | icpc | ifort
  - man gcc | g++ | gfortran
- **Use your favorite search engine to find vendor manuals on line**

- **Try different compilers**
  - Diagnostic messages and language spec compliances differ
- **Look for memory corruption**
  - Bad memory reference in one place (array out of bounds) can make code crash elsewhere
  - It might appear that you're crashing on a perfectly valid line of code
- **Check the arguments to your MPI calls**
- **Call the NERSC Consultants (800-66-NERSC or 510 486-8600)**

- **Hardware Event Counters**
  - Special registers count events on processor
  - E.g. number of floating point instructions
  - Many possible events
  - Only a few can be recorded at a time (~4 counters)
  - Can give you an idea of how efficiently you are using the processor hardware

# Apprentice Basic View



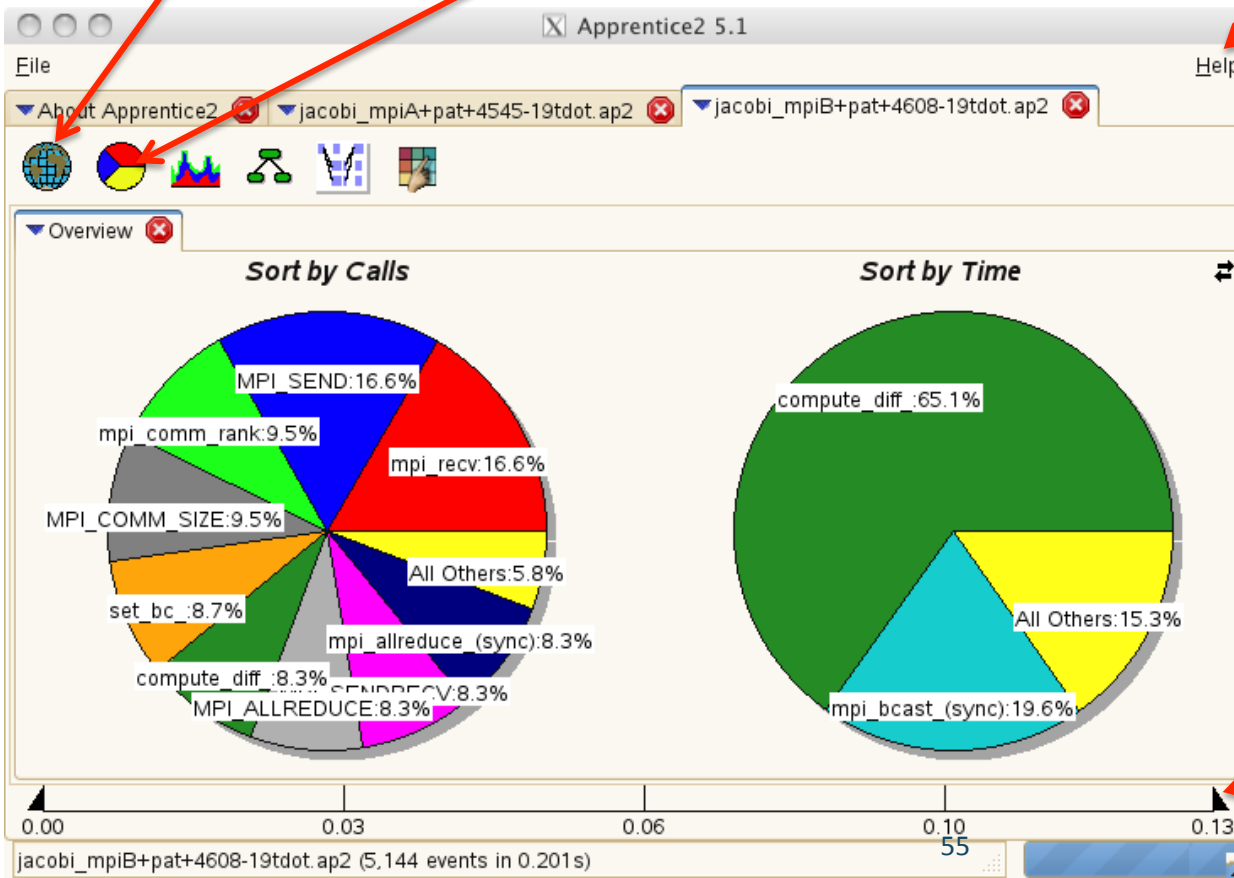
Can select new (additional) data file and do a screen dump

Worthless

Useful

Can select other views of the data

Can drag the "calipers" to focus the view on portions of the run



- **You will have a homework assignment using TAU**
  - `%module load tau`
  - Define paths in Makefile
  - Modify header file to define TAU macros
  - Add macro calls to the code
  - Compile and submit to batch queue
  - Use pprof to produce readable output
- **Good reference**
  - <http://acts.nersc.gov/events/Workshop2011/Talks/TAU.pdf>



# Users Want (Need?) Tools



- **Users are asking for tools because HPC systems and programming models are changing**
- **More and more components to worry about**
  - CPU (caches, FPUs, pipelining, ...)
  - Data movement to main memory, GPU memory, levels of cache
  - I/O
  - Network (message passing)
  - CPU Threads (OpenMP)
  - GPU performance

# Questions to You

---



- **What tools do you use?**
- **What tools do you want?**
- **What would you like centers to support?**
- **Can you get to exascale without tools?**

# What I Want in a Tool



- **Let the users help themselves**
- **Work for everyone all (most of?) the time**
- **Easy to use**
- **Useful**
- **Easy to interpret the results**
- **Affordable & ubiquitous**
- **Simple, supplement existing complex tools**
  - Point the way for a “deeper dive” in problem areas

- **Integrated Performance Monitoring**
- **Developed by David Skinner at NERSC**
- **MPI profiling, hardware counter metrics, IO profiling (?)**
- **IPM requires no code modification & no instrumented binary**
- **IPM uses hooks already in the MPI library to intercept your MPI calls and wrap them with timers and counters**

- **How it works (user perspective)**
    - % module load IPM\*
    - Run program as normal
    - Look at results on the web
  - **It's that easy!**
    - And extremely low overhead, so IPM is examining your production code
- \* (As long as your system supports dynamic load libs)

# IPM Examples



Click on the metric you are want.

NERSC job details  
[http://www.nersc.gov/REST/jobs/job\\_details.php?stepid=732423.sdb&timestamp=1313679078&completion=1313679081](http://www.nersc.gov/REST/jobs/job_details.php?stepid=732423.sdb&timestamp=1313679078&completion=1313679081)

## IPM Summary

Executable					./wrf.exe
Number of tasks	512	Aggregate GFlop/sec	0.1482	Average GFlop/sec/task	0.0003
Average wall secs	8.861e+01	Aggregate memory (GB)	32.2626	Average memory/task (GB)	0.0630
Average MPI secs/task	7.898e+01	MPI time %	89.14	Aggregate MPI calls made	7.027e+07

## IPM Summary Statistics - 512 tasks

Metric	Sum over all tasks	Average (per task)	Task CV (%)	Task Minimum	Task Maximum
<a href="#">Aggregate Floating Point Operations (Flop x 10**9)</a>	1.313e+01	2.565e-02	6.10	1.713e-02	2.758e-02
<a href="#">GFlop/sec</a>	1.482e-01	2.895e-04	6.10	1.934e-04	3.114e-04
<a href="#">Maximum Memory Usage (GBytes)</a>	3.226e+01	6.301e-02	10.12	5.701e-02	1.947e-01
<a href="#">Time Spent in MPI Routines (sec)</a>	4.044e+04	7.898e+01	4.05	9.801e+00	8.359e+01
<a href="#">Wallclock Time (sec)</a>	4.537e+04	8.861e+01	0.10	8.848e+01	8.895e+01

Memory in units of gigabytes; time in seconds.

## Hardware counter statistics - 512 tasks

Counter Name	Sum over all tasks	Average (per task)	Task CV (%)	Task Minimum	Task Maximum
<a href="#">PAPI FP OPS</a>	1.161799e+12	2.269139e+09	6.09	1.515023e+09	2.439529e+09

## MPI Time Statistics - 512 tasks

Call	Sum over all tasks	Average (per task)	Task CV (%)	Task Minimum	Task Maximum	% of MPI	% of wall
<a href="#">MPI Bcast</a>	3.517e+04	6.869e+01	4.48	4.342e-01	7.269e+01	86.969	77.520
<a href="#">MPI Scatterv</a>	2.589e+03	5.057e+00	5.79	1.059e+00	5.540e+00	6.403	5.707
<a href="#">MPI Wait</a>	2.176e+03	4.249e+00	17.82	1.250e+00	4.968e+00	5.380	4.795
<a href="#">MPI Gatherv</a>	4.312e+02	8.422e-01	36.44	3.552e-03	2.271e+00	1.066	0.950
<a href="#">MPI Isend</a>	5.250e+01	1.025e-01	11.96	7.182e-02	1.259e-01	0.130	0.116
<a href="#">MPI Irecv</a>	1.033e+01	2.017e-02	10.21	1.217e-02	2.613e-02	0.026	0.023
<a href="#">MPI Gather</a>	1.021e+01	1.995e-02	502.07	1.391e-03	1.434e+00	0.025	0.023
<a href="#">MPI Comm rank</a>	4.563e-01	8.913e-04	4.74	7.799e-04	1.404e-03	0.001	0.001
<a href="#">MPI Comm size</a>	9.629e-02	1.881e-04	10.65	1.462e-04	4.859e-04	0.000	0.000
<a href="#">MPI Init</a>	0.000e+00	0.000e+00	0.000e+00	0.000e+00	0.000	0.000	0.000
<a href="#">MPI Finalize</a>	0.000e+00	0.000e+00	0.000e+00	0.000e+00	0.000	0.000	0.000

Average MPI Time per Task



- MPI\_Bcast
- MPI\_Scatterv
- MPI\_Wait
- MPI\_Gatherv

- **Sampling**
  - Regularly interrupt the program and record where it is
  - Build up a statistical profile of time spent in various routines
  - Concentrate first on longest running sections or routines
- **Tracing**
  - Insert hooks into program to record and time program events (logging)
  - Reasonable for sequential programs
  - Unwieldy for large parallel programs (too much data!)

- Optional visualization tool for Cray's perftools data
- Use it in a X Windows environment
- Uses a data file as input (**XXX.ap2**) that is prepared by **pat\_report**

```
app2 [--limit_per_pe tags] XXX.ap2
```



- **PAPI (Performance API) provides a standard interface for use of the performance counters in major microprocessors**
- **Predefined actual and derived counters supported on the system**
  - To see the list, run ‘papi\_avail’ on compute node via aprun:

```
qsub -I -lmppwidth=24
module load perftools
aprun -n 1 papi_avail
```
- **AMD native events also provided; use ‘papi\_native\_avail’:**

```
aprun -n 1 papi_native_avail
```

- **Tuning and Analysis Utilities**
- **Fortran, C, C++, Java performance tool**
- **Procedure**
  - Insert macros
  - Run the program
  - View results with pprof
- **More info than gprof**
  - E.g. per process, per thread info; supports pthreads
- **<http://acts.nersc.gov/tau/index.html>**

- **IPM “only” gives a high-level, entire-program-centric view**
- **Still, very valuable guidance**
  - Shows whole-run info per MPI task, OpenMP thread, (CUDA under development)
  - Many pieces of data in one place
- **Reveals what many users don’t know about their code**
  - High-water memory usage (per task)
  - Load balance
  - Call imbalance
  - MPI time
  - I/O time