

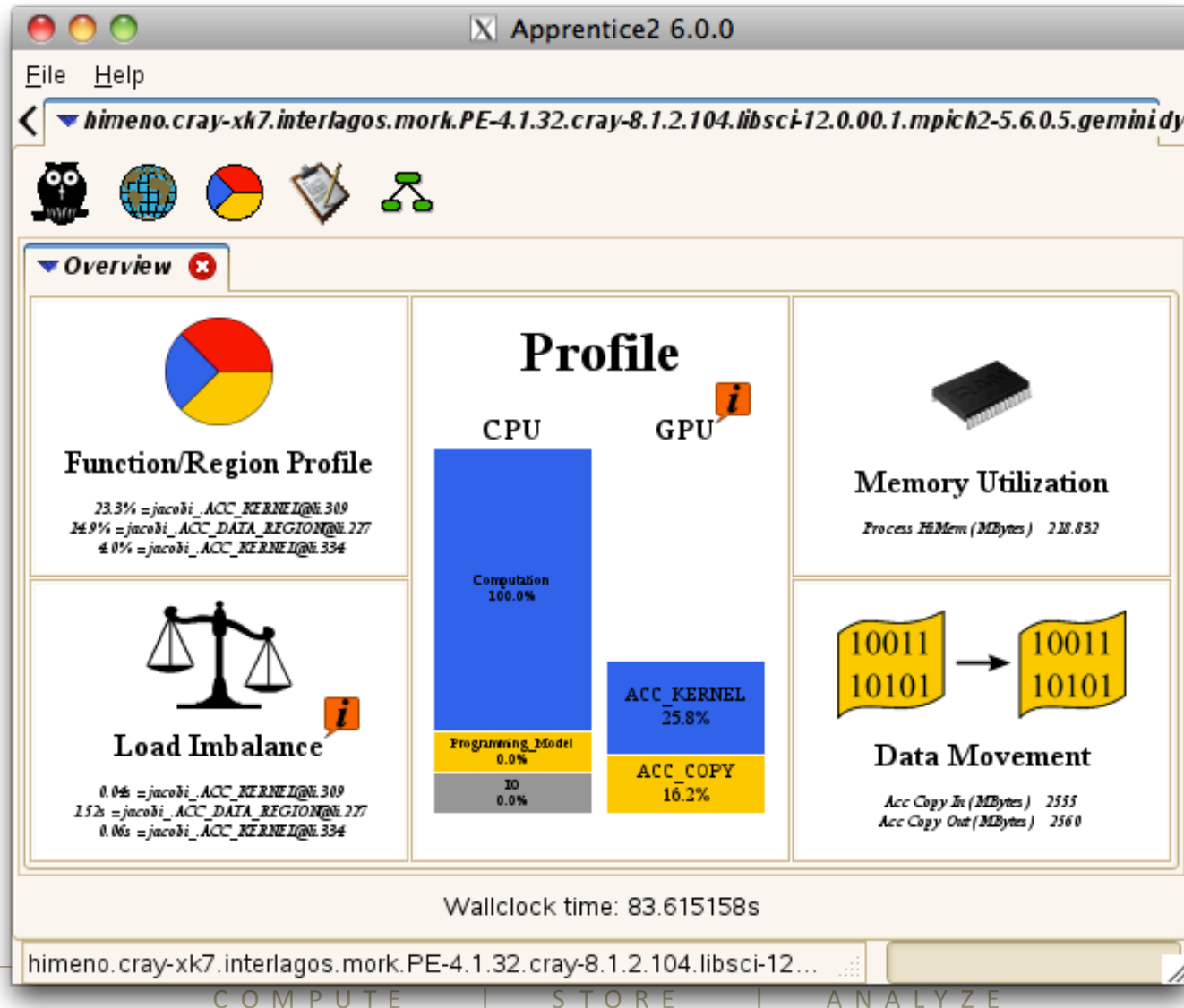


# **Cray Performance Measurement, Analysis and Porting Tools**

**Luiz DeRose**  
Sr. Principal Engineer,  
Programming  
Environments Director

**Heidi Poxon**  
Technical Lead &  
Sr. Manager, Performance  
Tools

# Cray Performance Tools Refresher





# Safe Harbor Statement

This presentation may contain forward-looking statements that are based on our current expectations. Forward looking statements may include statements about our financial guidance and expected operating results, our opportunities and future potential, our product development and new product introduction plans, our ability to expand and penetrate our addressable markets and other statements that are not historical facts. These statements are only predictions and actual results may materially vary from those projected. Please refer to Cray's documents filed with the SEC from time to time concerning factors that could affect the Company and these forward-looking statements.

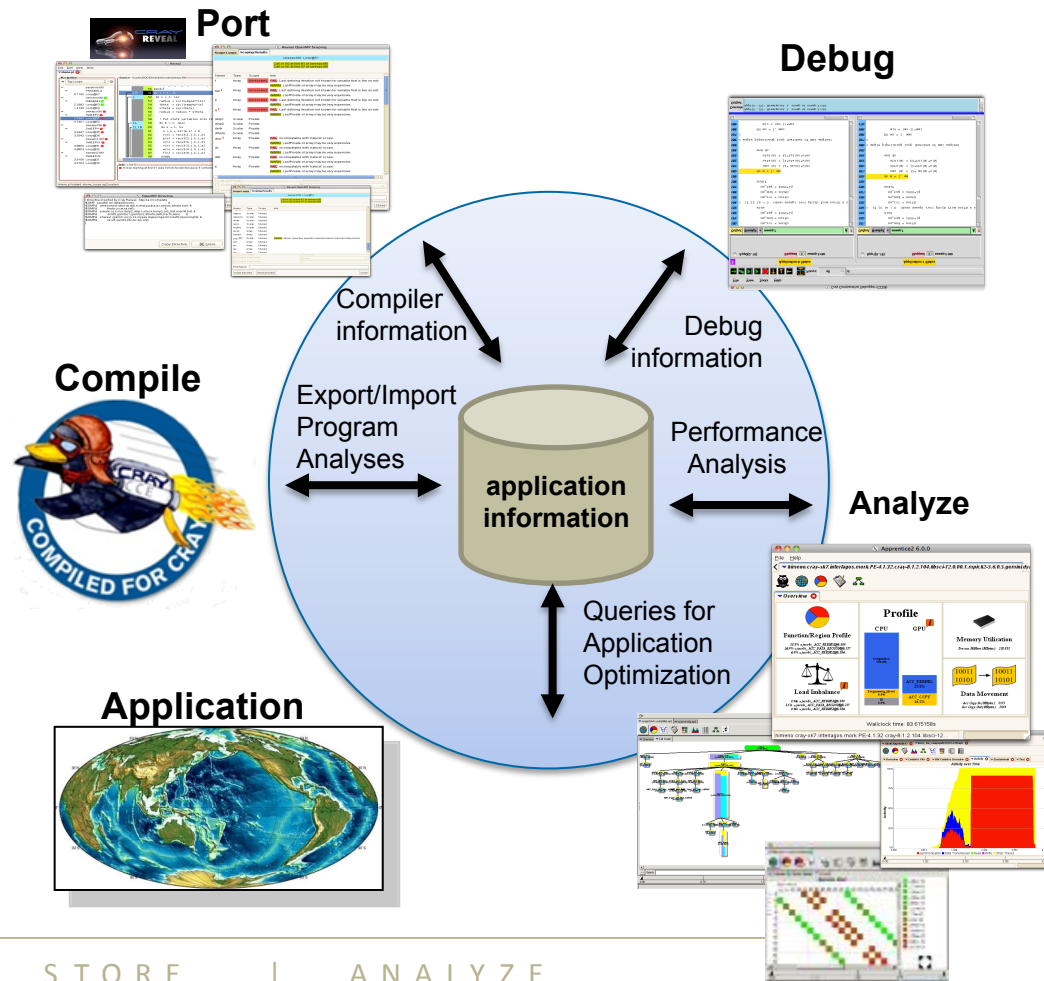


# Agenda

- 10:20 – 10:30 Overview of Cray PE topics for the day
- 10:30 – 11:15 Cray Performance Tools refresher
- **11:15 – 12:00 lab (CrayPat-lite, loop stats, Rank reorder)**
- 12:00 – 13:00 Lunch
- 13:00 – 14:00 Reveal
- 14:00 – 14:30 CCE and OpenACC update
- **14:30 – 15:00 lab with Reveal**
- 15:00 – 15:20 Break
- 15:20 – 16:00 CCDB
- **16:00 – 17:00 find a bug using CCDB, general lab time**

# The Programming Environment Mission

- It is the role of the Programming Environment to **close the gap** between observed performance and achievable performance
- Support the **application development life cycle** by providing a **tightly coupled** environment with compilers, libraries, and tools that will **hide the complexity** of the system
  - Address issues of scale and complexity of HPC systems
  - Target **ease of use** with extended **functionality** and increased **automation**
  - Close **interaction with users**
    - For feedback targeting functionality enhancements



COMPUTE | STORE | ANALYZE

# Cray Programming Environment Focus

- **Performance**

- Help users **maximize the cycles to the application**
  - Address issues of scale and complexity of HPC systems

- **Programmability**

- How do you get intuitive behavior and best performance with the least amount of effort
  - Provide the **best environment to develop, debug, analyze, and optimize** applications for production supercomputing
  - Provide programming environment **consistency across Cray platforms**



# Recent Enhancements

- **Cray Apprentice2 available for Linux, Mac and Windows with .dmg and .exe installers**
  - Available in `$CRAYPAT_ROOT/share/desktop_installers`
- **Router-aware MPI rank placement**
- **Simplified Interface to performance counters**
- **Program timeline**
- **I/O statistics in CrayPat-lite**
- **pat\_build default is now pat\_build -O apa**



# Two Interfaces to the Performance Tools

- **CrayPat** offers a wealth of performance measurement, analysis and presentation options for in-depth performance investigation and tuning assistance
- **CrayPat-lite** offers easy access to an application performance summary for users not familiar with the Cray performance tools or who may not be familiar with performance analysis
- **CrayPat classic and CrayPat-lite are designed to compliment each other**
  - Produce files with same format
  - Users familiar with CrayPat can easily switch back and forth between the two interfaces
  - CrayPat-lite users become familiar with reporting style also used with CrayPat





# CrayPat-lite

- **Produces application performance statistics at the end of a job**
  - Focus is to offer a simplified interface to basic application performance information for users not familiar with the Cray performance tools and perhaps new to application performance analysis
  - Gives sites the option to enable/disable application performance data collection for all users for a period of time
- **Compliments “classic” perftools**
- **Provides a simple way to transition from perftools-lite to perftools to encourage further tool use for more in-depth performance analysis**

# Using CrayPat-lite

## Access light version of performance tools software

```
> module load perftools-lite
```

## Build program

```
> make
```



```
a.out (instrumented program)
```

## Run program (no modification to batch script)

```
aprun a.out
```



```
Condensed report to stdout  
a.out*.rpt (same as stdout)  
a.out*.ap2  
MPICH_RANK_XXX files
```



# Performance Statistics Available

- **Set of predefined experiments, enabled with the CRAYPAT\_LITE environment variable**
  - Sample\_profile
  - Event\_profile
  - GPU
- **Job information**
  - Number of MPI ranks, ranks per node, number of threads
  - Wallclock
  - High memory water mark
  - Aggregate MFLOPS (CPU only)
  - I/O
- **Profile of top time consuming routines with load balance**
- **Observations**
- **Instructions on how to get more information**

# CrayPat-lite Output Example

```

CrayPat/X: Version 6.1.4.12457 Revision 12457 (xf 12277) 02/26/14 13:58:24
Experiment:          lite lite/sample_profile
Number of PEs (MPI ranks): 8164
Numbers of PEs per Node: 16 PEs on each of 510 Nodes
                    4 PEs on 1 Node
Numbers of Threads per PE: 1
Number of Cores per Socket: 8
Execution start time: Fri Feb 28 23:06:31 2014
System name and speed: hera2 2100 MHz

```

```

Wall Clock Time: 999.595275 secs
High Memory: 475.52 MBytes
MFLOPS (aggregate): 806112.33 M/sec
I/O Read Rate: 33.57 MBytes/Sec
I/O Write Rate: 215.40 MBytes/Sec

```

Table 1: Profile by Function Group and Function (top 7 functions shown)

| Time%  | Time       | Imb. Time | Imb. Time% | Calls     | Group                                | Function |
|--------|------------|-----------|------------|-----------|--------------------------------------|----------|
| 100.0% | 101.961423 | --        | --         | 5315211.9 | Total                                | PE=HIDE  |
| -----  |            |           |            |           |                                      |          |
| 92.5%  | 94.267451  | --        | --         | 5272245.9 | USER                                 |          |
| -----  |            |           |            |           |                                      |          |
| 75.8%  | 77.248585  | 2.356249  | 3.0%       | 1001.0    | LAMMPS_NS::PairLJCut::compute        |          |
| 6.5%   | 6.644545   | 0.105246  | 1.6%       | 51.0      | LAMMPS_NS::Neighbor::half_bin_newton |          |
| 4.1%   | 4.131842   | 0.634032  | 13.5%      | 1.0       | LAMMPS_NS::Verlet::run               |          |
| 3.8%   | 3.841349   | 1.241434  | 24.8%      | 5262868.9 | LAMMPS_NS::Pair::ev_tally            |          |
| 1.3%   | 1.288463   | 0.181268  | 12.5%      | 1000.0    | LAMMPS_NS::FixNVE::final_integrate   |          |
| =====  |            |           |            |           |                                      |          |
| 7.0%   | 7.110931   | --        | --         | 42637.0   | MPI                                  |          |
| -----  |            |           |            |           |                                      |          |
| 4.8%   | 4.851309   | 3.371093  | 41.6%      | 12267.0   | MPI_Send                             |          |
| 1.5%   | 1.536106   | 2.592504  | 63.8%      | 12267.0   | MPI_Wait                             |          |
| =====  |            |           |            |           |                                      |          |



# The Cray Performance Analysis Framework

- **Supports traditional post-mortem performance analysis**
  - Automatic identification of performance problems
    - Indication of causes of problems
    - Suggestions of modifications for performance improvement
  - `pat_build`: provides automatic instrumentation
  - **CrayPat run-time library** collects measurements (transparent to the user)
  - `pat_report` performs analysis and generates text reports
  - `pat_help`: online help utility
  - **Cray Apprentice2**: graphical visualization tool

# Steps to Using CrayPat “classic”

## Access performance tools software

```
> module load perftools
```

## Build program, retaining .o files

```
> make
```



```
a.out
```

## Instrument binary

```
> pat_build -O apa a.out
```



```
a.out+pat
```

## Modify batch script and run program

```
aprun a.out+pat
```



```
a.out+pat*.xf
```

## Process raw performance data and create report

```
> pat_report a.out+pat*.xf
```



```
a.out+pat*.ap2  
Text report to stdout  
a.out+pat*.apa  
MPICH_RANK_XXX
```



# Sampling with Line Number information

```

heidi@limited: /h/heidi — ssh — 81x26
Table 2: Profile by Group, Function, and Line

Samp% | Samp | Imb. | Imb. | Group
      |      | Samp | Samp% | Function
      |      |      |      | Source
      |      |      |      | Line
      |      |      |      | PE=HIDE

100.0% | 8376.9 | -- | -- | Total
-----
| 93.2% | 7804.0 | -- | -- | USER
-----
|| 51.7% | 4328.7 | -- | -- | calc3_
3|      |      |      |      | heidi/DARPA/cache_util/calc3.do300-ijswap.F
||||-----
4||| 15.7% | 1314.4 | 93.6 | 6.8% | line.78
4||| 13.9% | 1167.7 | 98.3 | 7.9% | line.79
4||| 14.5% | 1211.6 | 97.4 | 7.6% | line.80
4||| 1.2% | 103.1 | 26.9 | 21.2% | line.93
4||| 1.1% | 88.4 | 22.6 | 20.8% | line.94
4||| 1.0% | 84.5 | 17.5 | 17.6% | line.95
4||| 1.0% | 86.8 | 33.2 | 28.2% | line.96
4||| 1.3% | 105.0 | 23.0 | 18.4% | line.97
4||| 1.4% | 116.5 | 24.5 | 17.7% | line.98
||||-----
|||| 144,1 38%

```



# APA File Example

```
# You can edit this file, if desired, and use it
# to reinstrument the program for tracing like this:
#
# pat build -O standard.cray-
# xt.PE-2.1.56HD.pgi-8.0.amd64.pat-5.0.0.2-Oapa.512.quad.cores.seal.
# 090405.1154.mpi.pat_rt_exp=default.pat_rt_hwpc=none.14999.xf.xf.apa
#
# These suggested trace options are based on data from:
#
# /home/users/malice/pat/Runs/Runs.seal.pat5001.2009Apr04/./pat.quad/
# homme/standard.cray-xt.PE-2.1.56HD.pgi-8.0.amd64.pat-5.0.0.2-Oapa.
# 512.quad.cores.seal.
# 090405.1154.mpi.pat_rt_exp=default.pat_rt_hwpc=none.14999.xf.xf.cdb
# -----
# HWPC group to collect by default.
#
-Drtenv=PAT_RT_HWPC=1 # Summary with TLB metrics.
# -----
# Libraries to trace.
#
-g mpi
# -----
# User-defined functions to trace, sorted by % of samples.
#
# The way these functions are filtered can be controlled with
# pat_report options (values used for this file are shown):
#
# -s apa_max_count=200 No more than 200 functions are listed.
# -s apa_min_size=800 Commented out if text size < 800 bytes.
# -s apa_min_pct=1 Commented out if it had < 1% of samples.
# -s apa_max_cum_pct=90 Commented out after cumulative 90%.
#
# Local functions are listed for completeness, but cannot be traced.
#
-w # Enable tracing of user-defined functions.
# Note: -u should NOT be specified as an additional option.
```

```
# 31.29% 38517 bytes
# -T prim_advance_mod_preq_advance_exp_
#
# 15.07% 14158 bytes
# -T prim_si_mod_prim_diffusion_
#
# 9.76% 5474 bytes
# -T derivative_mod_gradient_str_nonstag_
#
. . .
# 2.95% 3067 bytes
# -T forcing_mod_apply_forcing_
#
# 2.93% 118585 bytes
# -T column_model_mod_applycolumnmodel_
#
# Functions below this point account for less than 10% of samples.
#
# 0.66% 4575 bytes
# -T bndry_mod_bndry_exchangev_thsave_time_
#
# 0.10% 46797 bytes
# -T baroclinic_inst_mod_binst_init_state_
#
# 0.04% 62214 bytes
# -T prim_state_mod_prim_printstate_
#
. . .
# 0.00% 118 bytes
# -T time_mod_timelevel_update_
# -----
-o preqx.cray-xt.PE-2.1.56HD.pgi-8.0.amd64.pat-5.0.0.2.x+apa
# New instrumented program.
#
./AUTO/cray/css.pe_tools/malice/craypat/build/pat/2009Apr03/2.1.56HD/
amd64/homme/pgi/pat-5.0.0.2/homme/2005Dec08/build.Linux/preqx.cray-
xt.PE-2.1.56HD.pgi-8.0.amd64.pat-5.0.0.2.x # Original program.
```





# CrayPat Runtime Options

- Runtime controlled through PAT\_RT\_XXX environment variables
- See [intro\\_craypat\(1\)](#) man page
- **Examples of control**
  - Enable full trace
  - Change number of data files created
  - Enable collection of HW counters
  - Enable collection of network counters
  - Enable tracing filters to control trace file size (max threads, max call stack depth, etc.)



# Example Runtime Environment Variables

- **Optional timeline view of program available**
  - export `PAT_RT_SUMMARY=0`
  - View trace file with Cray Apprentice<sup>2</sup>
- **Write 1 file per node:**
  - export `PAT_RT_EXPFIL_MAX=0`
- **Request hardware performance counter information:**
  - export `PAT_RT_PERFCTR=<HW counter group or event(s)>`
  - Can specify individual events or predefined groups

# Generating Profile from APA

- Instrument application for further analysis (a.out+apa)

```
% pat_build -O <apafilename>.apa
```

- Run application

```
% aprun ... a.out+apa (or qsub <apa script>)
```

- Generate text report and visualization file (.ap2)

```
% pat_report -o my_text_report.txt [<datafile>.xf |  
  <datadir>]
```

- View report in text and/or with Cray Apprentice<sup>2</sup>

```
% app2 <datafile>.ap2
```

# Files Generated and the Naming Convention

| File Suffix             | Description   |
|-------------------------|---|
| a.out+pat               | Program instrumented for data collection  |
| a.out...s.xf            | Raw data for sampling experiment, available after application execution                               |
| a.out...t.xf            | Raw data for trace (summarized or full) experiment, available after application execution             |
| a.out...st.ap2          | Processed data, generated by pat_report, contains application symbol information                      |
| a.out...s.apa           | Automatic profiling analysis template, generated by pat_report (based on pat_build -O apa experiment) |
| a.out+apa               | Program instrumented using .apa file  |
| MPICH_RANK_ORDER.Custom | Rank reorder file generated by pat_report from automatic grid detection and reorder suggestions       |



# Performance Counters

- **Cray supports raw counters, derived metrics and thresholds for:**
  - Processor
  - Network
  - Accelerator
  - Power
- **Predefined groups**
  - Groups together suggested counters for experiments
- **Single interface to access counters (PAT\_RT\_PERFCTR environment variable)**
- **PAPI with Cray custom components for network, uncore, power (available to 3<sup>rd</sup> party tool developers)**

# Example: HW counter data and Derived Metrics

```

PAPI_TLB_DM  Data translation lookaside buffer misses
PAPI_L1_DCA  Level 1 data cache accesses
PAPI_FP_OPS  Floating point operations
DC_MISS      Data Cache Miss
User_Cycles  Virtual Cycles
  
```

=====

USER

-----

```

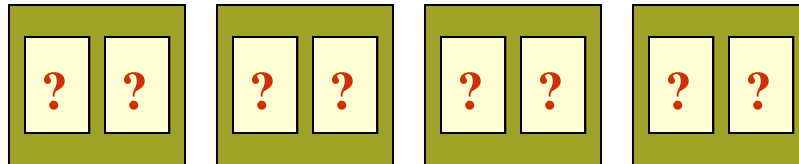
Time%                98.3%
Time                 4.434402 secs
Imb.Time             -- secs
Imb.Time%            --
Calls                0.001M/sec    4500.0 calls
PAPI_L1_DCM          14.820M/sec    65712197 misses
PAPI_TLB_DM          0.902M/sec    3998928 misses
PAPI_L1_DCA          333.331M/sec  1477996162 refs
PAPI_FP_OPS          445.571M/sec  1975672594 ops
User time (approx)   4.434 secs    11971868993 cycles  100.0%Time
Average Time per Call 0.000985 sec
CrayPat Overhead : Time 0.1%
HW FP Ops / User time 445.571M/sec  1975672594 ops  4.1%peak (DP)
HW FP Ops / WCT      445.533M/sec
Computational intensity 0.17 ops/cycle    1.34 ops/ref
MFLOPS (aggregate)    1782.28M/sec
TLB utilization        369.60 refs/miss  0.722 avg uses
D1 cache hit,miss ratios 95.6% hits    4.4% misses
D1 cache utilization (misses) 22.49 refs/miss  2.811 avg hits
  
```

=====

**PAT\_RT\_HWPC=1**  
**Flat profile data**  
**Raw counts**  
**Derived metrics**

# MPI Rank Reorder

- MPI rank placement with environment variable



- Distributed placement
- SMP style placement
- Folded rank placement
- User provided rank file



# MPI Rank Reorder (cont'd)

- **Non-default MPI rank placements are useful when point-to-point communication consumes significant fraction of the program time, and there is a significant load imbalance**
- **Performance report contains a prioritized list of placement scenarios and includes instructions on how to choose one of the placements for subsequent program execution**
- **Custom placement files automatically generated, user just chooses one**
- **Utilities available to create MPI rank placements for applications with grid or lattice topologies.**





# Automatic Communication Grid Detection

- **Cray performance tools produce a custom rank order if it's beneficial based on grid size, grid order and cost metric**
- **Heuristics available for:**
  - MPI sent message statistics
  - User time (time spent in user functions) – can be used for PGAS codes
  - Hybrid of sent message and user time)
- **Summarized findings in report**
- **Available with sampling or tracing**
- **Describe how to re-run with custom rank order**

# MPI Rank Order Observations

Table 1: Profile by Function Group and Function

| Time%        | Time       | Imb.<br>Time | Imb.<br>Time% | Calls   | Group<br>Function<br>PE=HIDE |
|--------------|------------|--------------|---------------|---------|------------------------------|
| 100.0%       | 463.147240 | --           | --            | 21621.0 | Total                        |
| <b>52.0%</b> | 240.974379 | --           | --            | 21523.0 | <b>MPI</b>                   |
| 47.7%        | 221.142266 | 36.214468    | 14.1%         | 10740.0 | mpi_recv                     |
| 4.3%         | 19.829001  | 25.849906    | 56.7%         | 10740.0 | MPI_SEND                     |
| 43.3%        | 200.474690 | --           | --            | 32.0    | USER                         |
| 41.0%        | 189.897060 | 58.716197    | 23.6%         | 12.0    | sweep_                       |
| 1.6%         | 7.579876   | 1.899097     | 20.1%         | 12.0    | source_                      |
| 4.7%         | 21.698147  | --           | --            | 39.0    | MPI_SYNC                     |
| 4.3%         | 20.091165  | 20.005424    | 99.6%         | 32.0    | mpi_allreduce_(sync)         |
| 0.0%         | 0.000024   | --           | --            | 27.0    | SYSCALL                      |

COMPUTE | STORE | ANALYZE

# MPI Rank Order Observations (2)

## MPI Grid Detection:

There appears to be point-to-point MPI communication in a **96 X 8 grid pattern**. The **52% of the total execution time spent in MPI functions** might be reduced with a rank order that maximizes communication between ranks on the same node. The effect of several rank orders is estimated below.

A file named `MPICH_RANK_ORDER.Grid` was generated along with this report and contains usage instructions and the Custom rank order from the following table.

| Rank Order | On-Node Bytes/PE | On-Node Bytes/PE% of Total Bytes/PE | MPICH_RANK_REORDER_METHOD |
|------------|------------------|-------------------------------------|---------------------------|
| Custom     | 2.385e+09        | 95.55%                              | 3                         |
| SMP        | 1.880e+09        | 75.30%                              | 1                         |
| Fold       | 1.373e+06        | 0.06%                               | 2                         |
| RoundRobin | 0.000e+00        | 0.00%                               | 0                         |

# MPICH\_RANK\_ORDER File

```

# The 'Custom' rank order in this file targets nodes with multi-core
# processors, based on Sent Msg Total Bytes collected for:
#
# Program:    /lus/nid00030/heidi/sweep3d/mod/sweep3d.mpi
# Ap2 File:   sweep3d.mpi+pat+27054-89t.ap2
# Number PEs: 48
# Max PEs/Node: 4
#
# To use this file, make a copy named MPICH_RANK_ORDER, and set the
# environment variable MPICH_RANK_REORDER_METHOD to 3 prior to
# executing the program.
#
# The following table lists rank order alternatives and the grid_order
# command-line options that can be used to generate a new order.
...

```



# Auto-Generated MPI Rank Order File

```
# The rank order in this file targets nodes with multi-core processors, based on Sent Msg Total Bytes collected for:
# Program: /lus/nid00023/malice/craypat/WORKSHOP/bh2o-demo/Rank/sweep3d/src/sweep3d
# Ap2 File: sweep3d.gmpi-u.ap2
# Number PEs: 768
# Max PEs/Node: 16
#
# To use this file, make a copy named MPICH_RANK_ORDER, and set the environment variable MPICH_RANK_REORDER_MET HOD to 3 prior to executing the program.
#
0,532,64,564,32,572,96,540,8,596,72,524,40,604,24,588,104,556,16,628,80,636,56,620,48,516,112,580,88,548,120,612
```

|                        |                         |                        |                         |
|------------------------|-------------------------|------------------------|-------------------------|
| 1,403,65,435,33,411,97 | 5,439,37,407,69,447,10  | 3,440,35,432,67,400,99 | 257,345,265,313,281,30  |
| ,443,9,467,25,499,105  | ,1,415,13,471,45,503,29 | ,408,11,464,43,496,27  | ,5,273,337,609,369,577, |
| 507,41,475             | ,479,77,511             | 472,51,504             | 377,617,329,513,529     |
| 73,395,81,427,57,459,1 | 53,399,85,431,21,463,6  | 19,392,75,424,59,456,8 | 545,297,633,361,625,32  |
| 7,419,113,491,49,387,8 | 1,391,109,423,93,455,1  | 3,384,107,416,91,488,1 | 1,585,537,601,289,553,  |
| 9,451,121,483          | 17,495,125,487          | 15,448,123,480         | 353,593,521,569,561     |
| 6,436,102,468,70,404,3 | 2,530,34,562,66,538,98  | 132,401,196,441,164,40 | 256,373,261,341,264,34  |
| 8,412,14,444,46,476,11 | ,522,10,570,42,554,26,  | 9,228,433,236,465,204, | 9,280,317,272,381,269,  |
| 0,508,78,500           | 594,50,602              | 473,244,393,188,497    | 309,285,333,277,365     |
| 86,396,30,428,62,460,5 | 18,514,74,586,58,626,8  | 252,505,140,425,212,45 | 352,301,320,325,288,35  |
| 4,492,118,420,22,452,9 | 2,546,106,634,90,578,1  | 7,156,385,172,417,180, | 7,328,304,360,312,376,  |
| 4,388,126,484          | 14,618,122,610          | 449,148,489,220,481    | 293,296,368,336,344     |
| 129,563,193,531,161,57 | 135,315,167,339,199,34  | 131,534,195,542,163,56 | 258,338,266,346,282,31  |
| 1,225,539,241,595,233, | 7,259,307,231,371,239,  | 6,227,526,235,574,203, | 4,274,370,766,306,710,  |
| 523,249,603,185,555    | 379,191,331,247,299     | 598,243,558,187,606    | 378,742,330,678,362     |
| 153,587,169,627,137,63 | 175,363,159,323,143,35  | 251,590,211,630,179,63 | 646,298,750,322,718,35  |
| 5,201,619,177,515,145, | 5,255,291,207,275,183,  | 8,139,622,155,550,171, | 4,758,290,734,662,686,  |
| 579,209,547,217,611    | 283,151,267,215,223     | 518,219,582,147,614    | 670,726,702,694,654     |
| 7,405,71,469,39,437,10 | 133,406,197,438,165,47  | 761,660,737,652,705,66 | 262,375,263,343,270,31  |
| 3,413,47,445,15,509,79 | 0,229,414,245,446,141,  | 8,745,692,673,700,641, | 1,271,351,286,319,278,  |
| ,477,31,501            | 478,237,502,253,398     | 684,713,644,753,724    | 342,287,350,279,374     |
| 111,397,63,461,55,429, | 157,510,189,462,173,43  | 729,732,681,756,721,71 | 294,318,358,383,359,31  |
| 87,421,23,493,119,389, | 0,205,390,149,422,213,  | 6,764,676,697,748,689, | 0,295,382,326,303,327,  |
| 95,453,127,485         | 454,181,494,221,486     | 657,740,665,649,708    | 367,366,335,302,334     |
| 134,402,198,434,166,41 | 130,316,260,340,194,37  | 760,528,736,536,704,56 | 765,661,709,663,741,65  |
| 0,230,442,238,466,174, | 2,162,348,226,308,234,  | 0,744,520,672,568,712, | 3,711,669,767,655,743,  |
| 506,158,394,246,474    | 380,242,332,250,300     | 592,752,552,640,600    | 671,749,695,679,703     |
| 190,498,254,426,142,45 | 202,364,186,324,154,35  | 728,584,680,624,720,51 | 677,727,751,693,647,70  |
| 8,150,386,182,418,206, | 6,138,292,170,276,178,  | 2,696,632,688,616,664, | 1,717,687,757,685,733,  |
| 490,214,450,222,482    | 284,210,218,268,146     | 544,608,656,648,576    | 725,719,735,645,759     |
| 128,533,192,541,160,56 | 4,535,36,543,68,567,10  | 762,659,738,651,706,66 |                         |
| 5,232,525,224,573,240, | 0,527,12,599,44,575,28  | 7,746,643,714,691,674, |                         |
| 597,184,557,248,605    | ,559,76,607             | 699,754,683,730,723    |                         |
| 168,589,200,517,152,62 | 52,591,20,631,60,639,8  | 722,731,763,658,642,75 |                         |
| 9,136,549,176,637,144, | 4,519,108,623,92,551,1  | 5,739,675,707,650,682, |                         |
| 621,208,581,216,613    | 16,583,124,615          | 715,698,666,690,747    |                         |

COMPUTE | STORE | ANALYZE



# grid\_order Utility

- Can use `grid_order` utility without first running the application with the Cray performance tools if you know a program's data movement pattern
- Originally designed for MPI programs, but since reordering is done by PMI, it can be used by other programming models (since PMI is used by MPI, SHMEM and PGAS programming models)
- Utility available if `perftools` modulefile is loaded
- See [grid\\_order\(1\)](#) man page or run `grid_order` with no arguments to see usage information



# Reorder Example for Bisection Bandwidth

- Assume 32 ranks

- Decide on row or column ordering:

- `$ grid_order -R -g 2,16`

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

- `$ grid_order -C -g 2,16`

0, 2, 4, 6, 8, 10, 12, 14, 16, 18, 20, 22, 24, 26, 28, 30

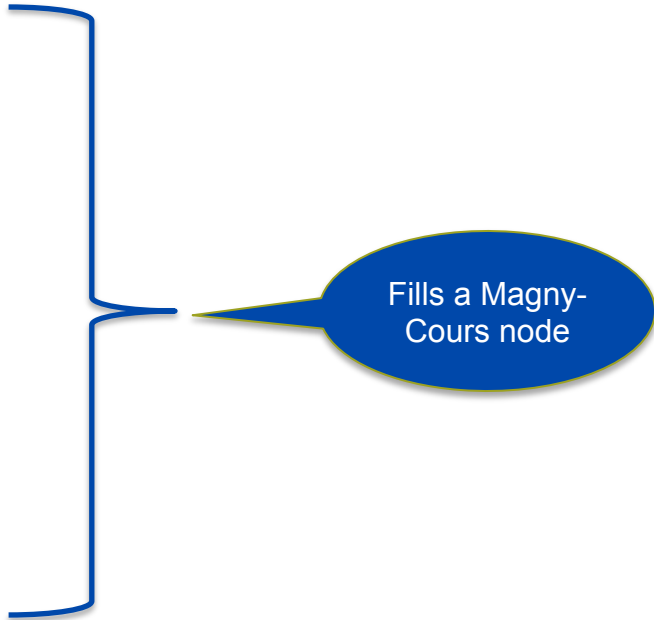
1, 3, 5, 7, 9, 11, 13, 15, 17, 19, 21, 23, 25, 27, 29, 31

- Since rank 0 talks to rank 16, and not with rank 1, we choose Row ordering

# Reorder Example for Bisection Bandwidth (2)

- Specify cell (or chunk) to make sure rank pairs live on same node (but don't care how many pairs live on a node)
- \$ grid\_order -R -g 2,16 -c 2,1

0,16  
 1,17  
 2,18  
 3,19  
 4,20  
 5,21  
 6,22  
 7,23  
 8,24  
 9,25  
 10,26  
 11,27  
 12,28  
 13,29  
 14,30  
 15,31





# Using New Rank Order

- Save grid\_order output to file called **MPICH\_RANK\_ORDER**
- Export **MPICH\_RANK\_REORDER\_METHOD=3**
- Run non-instrumented binary with and without new rank order to check overall wallclock time for improvements

# Collecting Loop Work Estimates

- Load PrgEnv-cray module (must use CCE)
- Load perftools module
- Compile **AND** link with `-h profile_generate`
  - `cc -h profile_generate -o my_program my_program.c`
- Instrument binary for tracing
  - `pat_build -w my_program`
- Run application
- Create report with loop statistics
  - `pat_report my_program.xf > loops_report`

pat\_report produces report plus .ap2 file that can be used with Reveal

# Example Report – Inclusive Loop Time

**Table 2:** Loop Stats by Function (from `-hprofile_generate`)

| Loop<br>Incl<br>Time<br>Total | Loop<br>Hit | Loop<br>Trips<br>Avg | Loop<br>Trips<br>Min | Loop<br>Trips<br>Max | Function=/.LOOP[.]<br>PE=HIDE |
|-------------------------------|-------------|----------------------|----------------------|----------------------|-------------------------------|
| 8.995914                      | 100         | 25                   | 0                    | 25                   | sweepy_.LOOP.1.li.33          |
| 8.995604                      | 2500        | 25                   | 0                    | 25                   | sweepy_.LOOP.2.li.34          |
| 8.894750                      | 50          | 25                   | 0                    | 25                   | sweepz_.LOOP.05.li.49         |
| 8.894637                      | 1250        | 25                   | 0                    | 25                   | sweepz_.LOOP.06.li.50         |
| 4.420629                      | 50          | 25                   | 0                    | 25                   | sweepx2_.LOOP.1.li.29         |
| 4.420536                      | 1250        | 25                   | 0                    | 25                   | sweepx2_.LOOP.2.li.30         |
| 4.387534                      | 50          | 25                   | 0                    | 25                   | sweepx1_.LOOP.1.li.29         |
| 4.387457                      | 1250        | 25                   | 0                    | 25                   | sweepx1_.LOOP.2.li.30         |
| 2.523214                      | 187500      | 107                  | 0                    | 107                  | riemann_.LOOP.2.li.63         |
| 1.541299                      | 20062500    | 12                   | 0                    | 12                   | riemann_.LOOP.3.li.64         |
| 0.863656                      | 1687500     | 104                  | 0                    | 108                  | parabola_.LOOP.6.li.67        |

# Cray Apprentice2 Overview

File Help
About Apprentice2 ✕ sweep3d.gmpi-u.ap2 ✕

Overview ✕

### Function/Region Profile

63.9% = *mpi\_recv*  
29.2% = *sweep\_*  
3.2% = *mpi\_allreduce\_(sync)*

# Profile

## CPU

Observations and suggestions

**NPI Grid Detection:**

There appears to be point-to-point MPI communication in a 96 X 8 grid pattern. The execution time spent in MPI functions might be reduced with a rank order that maximises communication between ranks on the same node. The effect of several rank orders is estimated below.

A file named `MPICH_RANK_ORDER.Grid` was generated along with this report and contains usage instructions and the Custom rank order from the following table.

| Rank Order | On-Node Bytes/FE | On-Node Bytes/FE <sup>4</sup> of Total Bytes/FE | MPICH_RANK_REORDER_METHOD |
|------------|------------------|---|---------------------------|
| Custom     | 1.851e+012       | 96.56%  | 3                         |
| SMP        | 1.459e+012       | 76.08%  | 1                         |
| Fold       | 1.056e+009       | 0.06%   | 2                         |
| RoundRobin | 0.000e+000       | 0.00%   | 0                         |

**Metric-Based Rank Order:**

When the use of a shared resource like memory bandwidth is unbalanced across nodes, total execution time may be reduced with a rank order that improves the balance. The metric used here for resource usage is: **USER Time**

For each node, the metric values for the ranks on that node are summed. The maximum and average value of those sums are shown below for both the current rank order and a custom rank order that seeks to reduce the maximum value.

A file named `MPICH_RANK_ORDER.USER_Time` was generated along with this report and contains usage instructions and the Custom rank order from the following table.

| Rank Order | Node Metric Imb. | Reduction in Max Value | Maximum Value | Average Value |
|------------|------------------|------------------------|---------------|---------------|
| Current    | 29.45%           |                        | 4.509e+003    | 3.181e+003    |
| Custom     | 0.34%            | 29.211%                | 3.192e+003    | 3.181e+003    |

### Memory Utilization

Process HiMem (MBytes) 658.835

### Load Imbalance

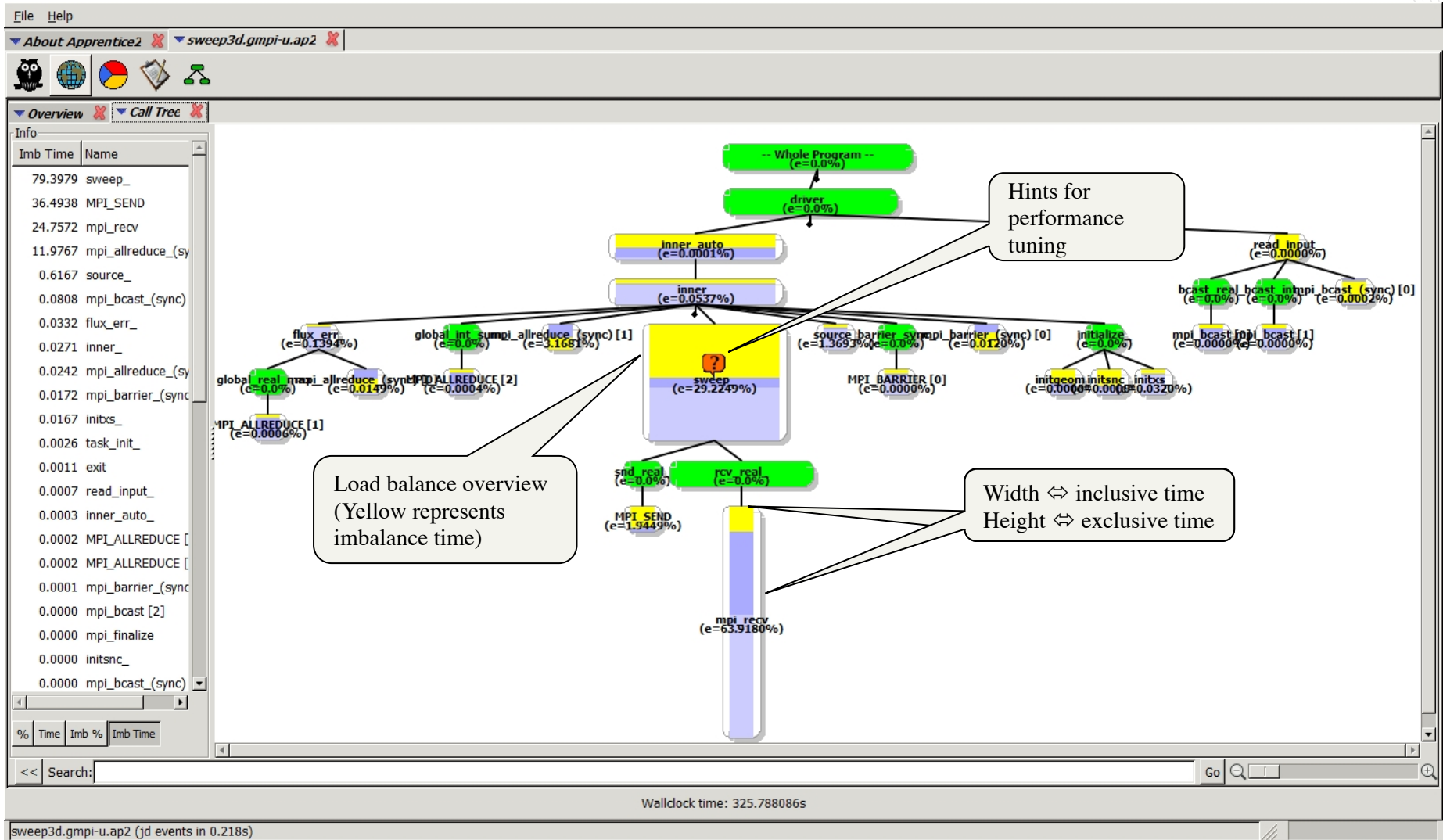
24.76s = *mpi\_recv*  
79.40s = *sweep\_*  
10.20s = *mpi\_allreduce\_(sync)*

### Data Movement

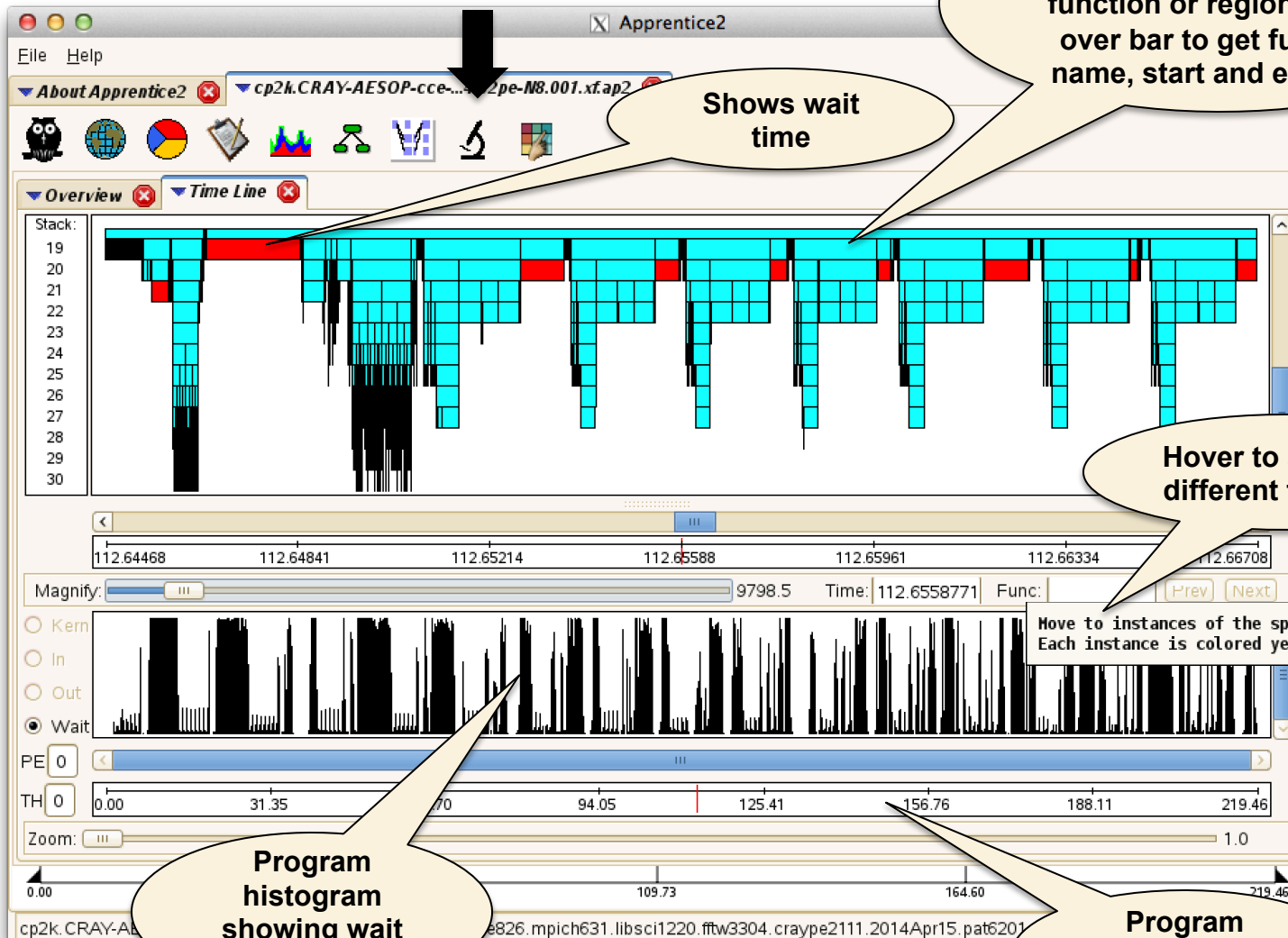
MPI Msg MBytes 2380.775

sweep3d.gmpi-u.ap2 (jd events in 0.218s)

# Call Tree View with Load Imbalance Information



# CPU Program Timeline: 36GB CP2K Full Trace



CPU call stack:  
Bar represents CPU function or region: Hover over bar to get function name, start and end time

Shows wait time

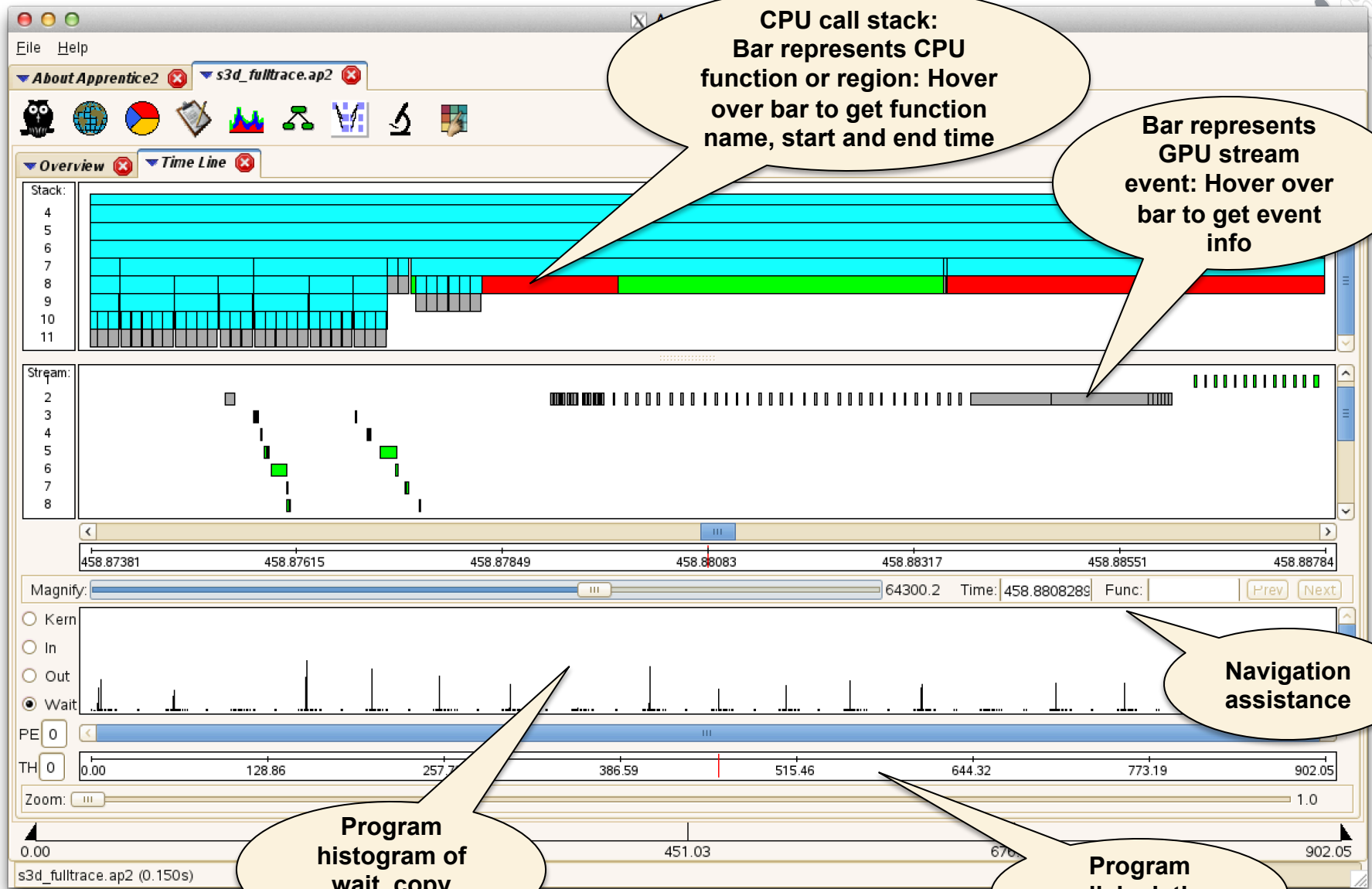
Hover to see what different filters do

Move to instances of the specified function. Each instance is colored yellow.

Program histogram showing wait time

Program wallclock time line

# GPU Program Timeline



Program histogram of wait, copy kernel time

Navigation assistance

Program wallclock time line

# Example Default Accelerator Statistics

Table 1: Time and Bytes Transferred for Accelerator Regions

| Host   | Host  | Acc   | Acc Copy | Acc Copy | Calls  | Calltree  |
|--------|-------|-------|----------|----------|--------|---|
| Time%  | Time  | Time  | In       | Out      |        | PE=HIDE   |
|        |       |       | (MBytes) | (MBytes) |        |   |
| 100.0% | 2.750 | 2.015 | 2812.760 | 13.568   | 103    | Total   |
| -----  |       |       |          |          |        |   |
| 100.0% | 2.750 | 2.015 | 2812.760 | 13.568   | 103    | lbm3d2p_d_  |
|        |       |       |          |          |        | lbm3d2p_d_.ACC_DATA_REGION@li.104                         |
| -----  |       |       |          |          |        |   |
| 3      | 63.5% | 1.747 | 1.747    | 2799.192 | --     | 1  lbm3d2p_d_.ACC_COPY@li.104                             |
| 3      | 22.1% | 0.609 | 0.088    | 12.304   | 12.304 | 36  streaming_  |
| -----  |       |       |          |          |        |   |
| 4      | 20.6% | 0.566 | 0.046    | 12.304   | 12.304 | 27  streaming_exchange_                                   |
| 5      |       |       |          |          |        | streaming_exchange_.ACC_DATA_REGION@li.526                |
| 6      | 18.8% | 0.517 | --       | --       | --     | 1   streaming_exchange_.ACC_DATA_REGION@li.526(exclusive) |
| 4      | 1.6%  | 0.043 | 0.042    | --       | --     | 9  streaming_.ACC_DATA_REGION@li.907                      |
| 5      | 1.1%  | 0.031 | 0.031    | --       | --     | 4   streaming_.ACC_REGION@li.909                          |
| 6      | 1.1%  | 0.031 | --       | --       | --     | 1   streaming_.ACC_REGION@li.909(exclusive)               |
| =====  |       |       |          |          |        |   |

...





# Lab Time

- **bw.ncsa.illinois.edu:~heidi/scratch/lab/**
  - CrayPat-lite
    - pr01
  - Loop statistics
    - pr05
  - MPI rank reorder
    - pr03
  - Apprentice2
    - pr04
  - Reveal
    - pr05, pr01
- **CCDB**



# Questions ?

# Legal Disclaimer

*Information in this document is provided in connection with Cray Inc. products. No license, express or implied, to any intellectual property rights is granted by this document.*

*Cray Inc. may make changes to specifications and product descriptions at any time, without notice.*

*All products, dates and figures specified are preliminary based on current expectations, and are subject to change without notice.*

*Cray hardware and software products may contain design defects or errors known as errata, which may cause the product to deviate from published specifications. Current characterized errata are available on request.*

*Cray uses codenames internally to identify products that are in development and not yet publically announced for release. Customers and other third parties are not authorized by Cray Inc. to use codenames in advertising, promotion or marketing and any use of Cray Inc. internal codenames is at the sole risk of the user.*

*Performance tests and ratings are measured using specific systems and/or components and reflect the approximate performance of Cray Inc. products as measured by those tests. Any difference in system hardware or software design or configuration may affect actual performance.*

*The following are trademarks of Cray Inc. and are registered in the United States and other countries: CRAY and design, SONEXION, URIKA, and YARCDATA. The following are trademarks of Cray Inc.: ACE, APPRENTICE2, CHAPEL, CLUSTER CONNECT, CRAYPAT, CRAYPORT, ECOPHLEX, LIBSCI, NODEKARE, THREADSTORM. The following system family marks, and associated model number marks, are trademarks of Cray Inc.: CS, CX, XC, XE, XK, XMT, and XT. The registered trademark LINUX is used pursuant to a sublicense from LMI, the exclusive licensee of Linus Torvalds, owner of the mark on a worldwide basis. Other trademarks used in this document are the property of their respective owners.*

*Copyright 2014 Cray Inc.*