**Profiling and Analyzing Program Performance Using Cray Tools**

Heidi Poxon

CUG 2017 CAFFEINATED COMPUTING

Redmond, Washington May 7-11, 2017

# 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, and URIKA. The following are trademarks of Cray Inc.:  APPRENTICE2, CHAPEL, CLUSTER CONNECT, CRAYPAT, CRAYPORT, ECOPHLEX, LIBSCI, NODEKARE, REVEAL, 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.*

COMPUTE | STORE | ANALYZE

# Agenda

- **Recent Cray performance tools enhancements**

- **What's coming next?**

- **Managing MCDRAM usage on Intel® Xeon Phi™ 7250 (hereafter referred to as "KNL")**

# Cray Performance Tools

- **Cray tools offer functionality that reduces the time investment associated with porting and tuning applications on Cray systems**

- **Whole program performance analysis across many nodes to help you find critical performance bottlenecks within a program**

- **Novice and advanced user interfaces for ease of use**

# Highlights Since Last CUG

- **Switch to perftools-base + instrumentation modules (6.4.0)**
  - perftools-base provides access Reveal, Apprentice2, pat_report and man pages without modification to applications

  - **perftools-base** loaded by default starting with **cdt-prgenv 6.0.4** (May 2017)

  - Load an instrumentation module to collect performance data

```
Examples:

$ module load perftools-lite
$ module load perftools
```

# Highlights Since Last CUG (continued)

- **Memory high water mark per NUMA domain (6.4.4)**

- **Charm++ support (6.4.4)**
  - Build and Run Charm++ and NAMD on the Cray XC with CrayPat
    - http://docs.cray.com/books/S-2802-10//S-2802-10.pdf

- **Address job termination issues for OpenMP programs built with Intel compiler (6.4.4)**

- **MCDRAM configuration statistics included in job summary (6.4.4)**

- **HBM memory analysis tool (6.4.4)**

COMPUTE | STORE | ANALYZE

# Reveal



- **Reduce effort associated with adding OpenMP to MPI programs**

- **Get insight into optimizations performed by the Cray compiler**

- **Use to add OpenMP or as a first step to parallelize loops that will target GPUs**

- **Track requests to memory and evaluate the bandwidth contribution of objects within a program**

# Reveal Enhancements

- **Reveal auto-parallelization (6.4.0)**
  - With one-click, build experimental binary that includes automatic runtime-assisted parallelization

- **Reveal client for Mac OS X (6.4.0)**
  - Improved tool response time with client that executes locally on laptop

# Reveal Auto-Parallelization Feedback
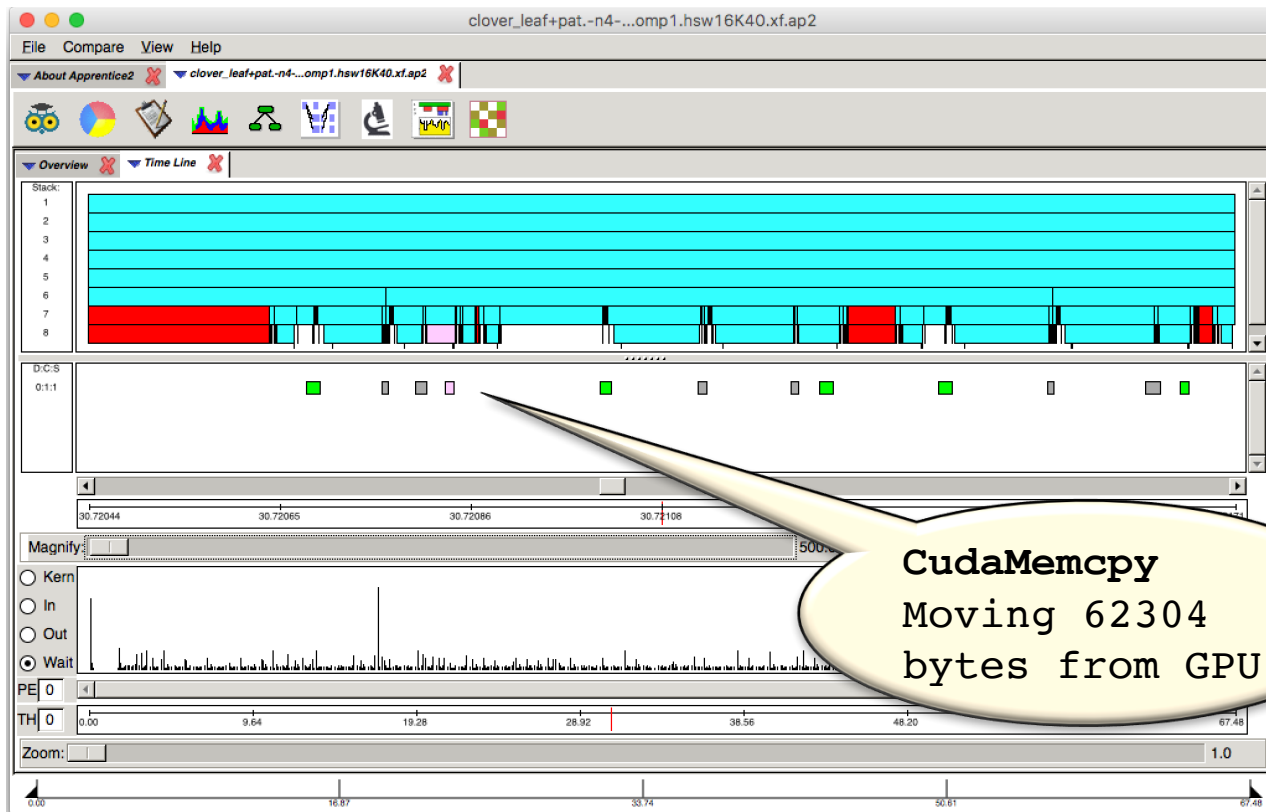
COMPUTE | STORE | ANALYZE

# What's Coming Next..

# Functionality Coming Next

- **CrayPat lite mode enhancement to disable instrumentation of test programs built through build systems (CMake/GNU Autotools) as part of application build**
  - export CRAYPAT_LITE_BLACKLIST=test1,test2

- **CUDA support in Apprentice2 (overview, GPU timeline)**

# CUDA Support in Apprentice2

# Functionality Coming Next (Continued)

- **Consolidation of Cray performance tools results files into single directory**
  - **MPICH_RANK_ORDER.Grid,**
  - **MPICH_RANK_ORDER.USER_Time**
  - **MPICH_RANK_ORDER.USER_Time_hybrid**
  - **stencil_order+pat+49144-225t.xf**
  - **stencil_order+pat+49144-225t.ap2**
  - **stencil_order+pat+49144-225t.rpt**
  - **stencil_order+pat+49144-225t.apa**

- **Same prefix naming scheme as used with current xf files**
  - **stencil_order+pat+49144-225t/**

COMPUTE | STORE | ANALYZE

# HBM Memory Analysis Assistance

# Managing Multi-tiered Memory

- **Future systems are likely to include a high performance tier (for bandwidth or latency) and a high capacity tier at a lower cost**

- **Our goals to assist users with using more than one type of explicitly addressable on-node memory:**
  - Provide easy-to-use interface for user to allocate data into HBM

  - Provide assistance with making best use of limited capacity

# Identifying Arrays Is Difficult

```
subroutine ax_e()

do i=1,n
    wr = g(1,i)*ur(i) + g(2,i)*us(i) + g(3,i)*ut(i)
    ws = g(2,i)*ur(i) + g(4,i)*us(i) + g(5,i)*ut(i)
    wt = g(3,i)*ur(i) + g(5,i)*us(i) + g(6,i)*ut(i)
    ur(i) = wr
    us(i) = ws
    ut(i) = wt
enddo
```

```
subroutine glsc3()

do i=1,n
    tmp = tmp + a(i)*b(i)*mult(i)
continue
```

Arrays a, b, and `mult` have a higher bandwidth sensitivity than array g

# MCDRAM Usage Assistance

- **Combination of CCE, CrayPat and Reveal are used to identify arrays that contribute most to memory bandwidth**

- **First introduced in December 2016**
  - cce/8.5.6
  - perftools-base/6.4.4

- See http://docs.cray.com/books/S-2803-10//S-2803-10.pdf

# Perform Analysis on Xeon or Phi

- **Haswell, Broadwell or KNL processors supported**
  - Memory traffic due to the hardware prefetcher on KNL is untracked
  - HSW is not KNL (differences with page table walks, L3 cache, etc.)
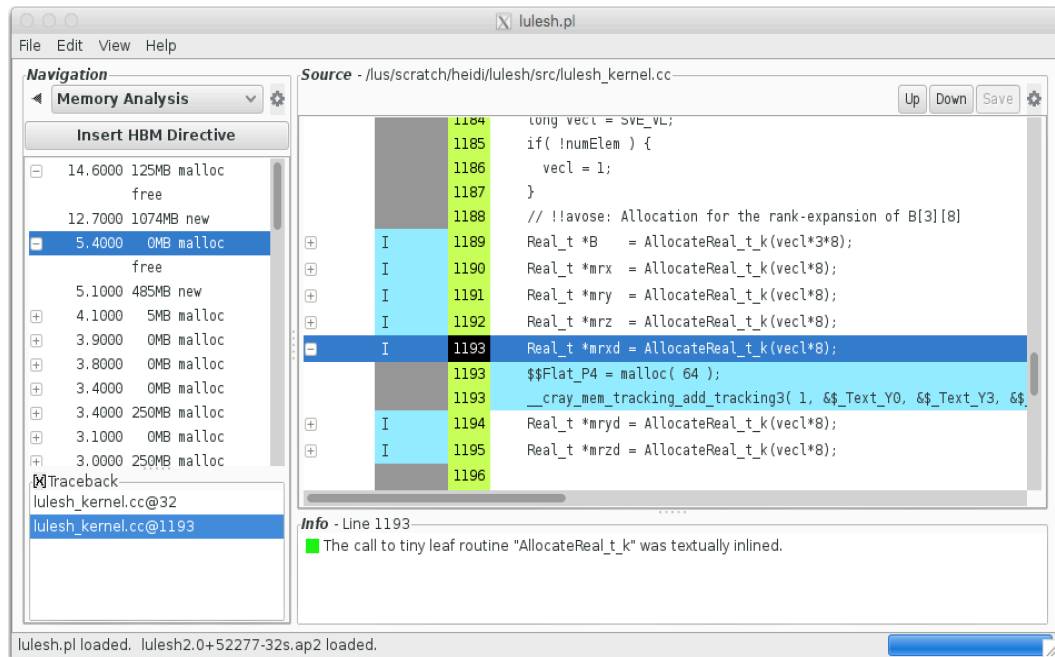  - Prefer HSW for applications with streaming accesses

```
$ module load PrgEnv-cray, craype-haswell
$ module load perftools-lite-hbm

$ cc -h pl=/path/my_program.pl my_program.c

Run program (no batch script modifications required) to create .ap2  file

$ reveal my_program.pl traffic_results.ap2
```

# Ranked Arrays and Allocation Sites



- **Objects sorted by memory bandwidth contribution**

- **Match free()s to mallocs()**

- **For C++, Reveal shows how STL objects rank with other arrays**
  - User must find declaration site, modify declaration to point to an hbw-aware allocator

# Build CCE Memory Allocation Directive

# Memory Analysis Results for Nekbone



- **Low data collection overhead**
  - ~1% - few %
  - Can analyze large jobs
  - Functions that are entered and exited often and allocate arrays big enough to be tracked can increase runtime

COMPUTE | STORE | ANALYZE

# Use Combination of Reveal and Text Report

- **Check program memory footprint**
  - Available at top of report with job summary information
  - If less than 16GB, use `numactl -- membind=1` to force all allocations into MCDRAM

- **Get ranking of which objects are referenced the most**

- **Pinpoint memory activity**

- **Review memory high water mark per NUMA domain**

```
Process |     HiMem |     HiMem |Numanode
  HiMem |      Numa |      Numa | PE=ALL
(MBytes) |    Node 0 |    Node 1 |
         |  (MBytes) |  (MBytes) |


  786.9 |     534.5 |     252.4 |Total
|-------------------------------------------
|    786.9 |     534.5 |     252.4 |numanode.0
||------------------------------------------
||    794.3 |     538.9 |     255.4 |pe.0
||    791.3 |     537.6 |     253.8 |pe.64
||    791.2 |     537.3 |     253.9 |pe.128
||    791.1 |     537.3 |     253.8 |pe.192
...
```

# Where to Find Memory Activity in Reports

- **High water mark by allocation site**
  - Number of all objects active at any time

- **Profile by Function table**
  - Shows where most of the memory traffic is happening within program

- **Profile by Group, Function, and Line table**
  - Identifies memory traffic hot spots within a function

# MCDRAM Allocation Assistance Recap

- **Cray Tools track requests to memory and evaluate the bandwidth contribution of objects within a program**

- **Helpful for memory-intensive programs that cannot fit within MCDRAM**

- **Reduces time investment associated with selectively allocating data into KNL's MCDRAM**

- **The result is performance portable code**
  - CCE memory allocation directives are ignored on X86 processors