# KAUST Supercomputing Laboratory

# Introduction to Performance Analysis tools on Shaheen II

George Markomanolis
Computational Scientist
April 17th, 2016
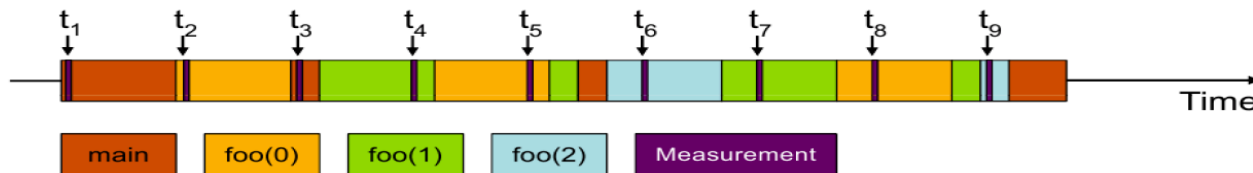
# Outline

❖ Introduction

❖ Test cases

❖ Cray tools

- Perftools

- Cray Apprentice 2

- Reveal

❖ Extrae/Paraver (briefly)

# Introduction

❖ **Why performance analysis?**

- Investigate the bottlenecks of an application
- Identify potential improvements
- Better usage of the hardware

❖ **Profiling**

- Sampling
  - Lightweight
  - Overhead depends on the sampling frequency
  - Can lack resolution if there are small function calls
- Event Tracing
  - Detailed information
  - Captures every event
  - Can capture communication events
  - Drawbacks, overhead and large amounts of data

# Sampling



```
int main()
{
  int i;

  for (i=0; i < 3; i++)
    foo(i);

  return 0;
}

void foo(int i)
{

  if (i > 0)
    foo(i - 1);

}
```
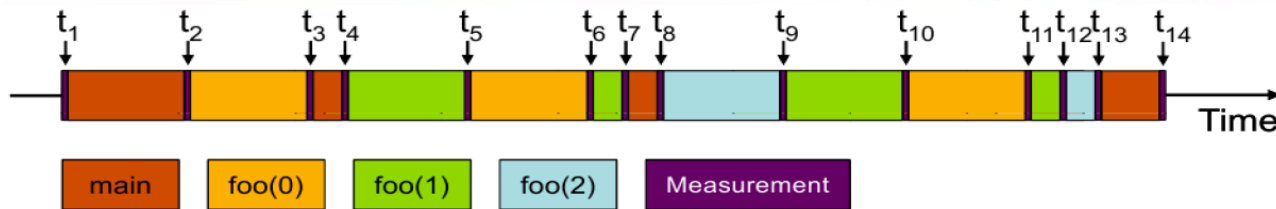
- Statistical inference of program behavior

- Not very detailed information

- Mainly for long-running applications

# Tracing



```
int main()
{
  int i;
  Enter("main");
  for (i=0; i < 3; i++)
    foo(i);
  Leave("main");
  return 0;
}

void foo(int i)
{
  Enter("foo");
  if (i > 0)
    foo(i - 1);
  Leave("foo");
}
```

- Every event is captured

- Detailed information

- Overhead (depends on many factors)

Technology

# Studying case

❖ NAS Parallel Benchmarks (NPB) consist of five kernels and three pseudo-applications, developed by NASA Advanced Supercomputing Division

❖ Why NPB/LU?

- LU stands for Lower-Upper Gauss-Seidel solver
- Simple application for testing purposes which combines computation and communication
- Compile with Cray, Intel, GNU compilers and fast

# CrayPat overview

❖ Assist the user with application performance analysis and optimization

- Provides concrete suggestions instead of just reporting

❖ Basic functionalities apply for all the compilers on the system

❖ Requires no source code or Makefile modification (for most of the cases)

# Components of CrayPat

❖ **Module perftools-base**
- pat_build – Instruments the program to be analyzed
- pat_report – Generates text reports from the performance data captured during program execution and exports data for use in other programs.
- Cray Apprentice2 – A graphical analysis tool that can be used to visualize and explore the performance data captured during program, execution
- Reveal – A graphical source code analysis tool that can be used to correlate performance analysis data with annotated source code listings, to identify key opportunities for optimization (it works only with Cray compiler)
- grid_order – Generates MPI rank order information that can be used with the MPICH_RANK_REORDER
- pat_help – Help system which provides extensive usage information

# Files generated during regular profiling

❖ A.out+pat+PID-node[s|t].xf: raw data files

- Depending on the profiling approach and conditions the execution of an instrumented application can create one or more .xf files where:
  - a.out is the name of the original program
  - PID is the process ID assigned to the instrumented program at runtime
  - Node is the physical node ID upon which the rank zero process executed
  - s|t is a letter code indicating the type of experiment performed, either **s** for sampling or **t** for tracing
- Pat_report tool dump the .xf file or export to another file format for use with other applications, i.e, *.ap2 files

❖ *.ap2 files: self contained compressed performance files

- Normally about 5 times smaller than the corresponding *.xf files
- Only one *.ap2 per experiment in comparison to potentially multiple *.xf files

# Prepare for the tutorial

- Connect to Shaheen II and copy the material:

  - ssh –X [username@shaheen.kaust.edu.sa](mailto:username@shaheen.kaust.edu.sa)

  - cp /scratch/tmp/performance_workshop.tgz **.**

  - tar zxvf performance_workshop.tgz

  - cd performance_workshop/NPB3.3-MPI
  - slides located in the folder performance_workshop/

❖ Load Perftools

- *module unload darshan*
- *module load perftools-base/6.3.2*
- *module load perftools/6.3.2*

❖ Compile the code

- *make clean*
- *make LU NPROCS=64 CLASS=C*
  - ▪ *"WARNING: PerfTools is saving object files from a temporary directory into directory…"*
- *cd bin*

❖ The new binary is called lu.C.64 is not instrumented yet

# Sampling instrumentation I

❖ Execute the application
  - *sbatch --reservation*=s1001_85 *submit.sh*
  - Check the output files (lu_C_64_out_...txt)

❖ Build the instrumented binary with sampling instrumentation
  - *pat_build –S lu.C.64*

❖ The instrumented binary is called *lu.C.64+pat*
❖ *Some results of the current presentation are acquired with 128 MPI processes.*

# Sampling instrumentation II

❖ Edit the submit.sh file, comment line 13 and uncomment line 16

- *sbatch --reservation=*s1001_85 *submit.sh*
- *The reservation of the nodes for this workshop is called s1001_85, you need to use it every time you submit jobs during this presentation.*

❖ The performance data are locate in a file called with the format
lu.C.64+PID-XXX**s**.xf  (PID and XXX are numbers)

# Create your first report with sampling instrumentation

❖ Execute the pat_report tool
  • *pat_report –o sampling_report_lu_C_64.txt lu.C.64+PID-XXX**s**.xf*
❖ Open the file *sampling_report_lu_C_64.txt*


❖ CrayPat/X:  Version 6.3.2 Revision rc1/6.3.2  02/25/16 18:26:21
  Number of PEs (MPI ranks):    64
  Numbers of PEs per Node:  32  PEs on each of  2  Nodes
  Numbers of Threads per PE:        1
  Number of Cores per Socket:    16
  Execution start time:  Wed Apr 13 16:57:06 2016
  System name and speed:  nid00035  2301 MHz (approx)
  Current path to data file:
   /scratch/markomg/NPB3.3.1/NPB3.3-MPI/bin/lu.C.64+pat+10974-35s.ap2
  (RTS)

# Create your first report with sampling instrumentation

❖ Table 1: Profile by Function

```
Samp% |  Samp | Imb.   | Imb.   |Group
      |       | Samp | Samp% | Function
      |       |      |       | PE=HIDE


100.0% | 1,039.1 |    -- |    -- |Total
|-------------------------------------------------------------
|  72.3% |  751.5 |    -- |    -- |USER
||------------------------------------------------------------
||  33.5% |  347.9 | 45.1 | 11.6% |rhs_
||   8.0% |   83.4 | 24.6 | 23.0% |blts_
||   7.9% |   82.1 | 18.9 | 18.9% |buts_
||   7.8% |   81.2 | 23.8 | 22.9% |jacld_
||   7.4% |   77.4 | 24.6 | 24.3% |jacu_
||   4.6% |   47.8 | 26.2 | 35.7% |exchange_3_
||   2.2% |   23.2 | 15.8 | 40.8% |ssor_
||============================================================
```

# Create your first report with sampling instrumentation (MPI with sampling is not helpful)

❖ Table 1: Profile by Function

```
Samp%  |  Samp | Imb.   | Imb.    |Group
       |       | Samp | Samp% | Function
       |       |       |         | PE=HIDE


  | 18.6% |  192.9 | -- |      -- |MPI
||-----------------------------------------------------------------
|| 6.3% |     65.1 | 154.9 | 70.9% |MPIDI_Cray_shared_mem_coll_bcast
|| 4.0% |     42.0 | 59.0 | 58.9% |MPIDI_CH3I_Progress
|| 2.2% |     22.5 | 99.5 | 82.2% |MPIDI_Cray_shared_mem_coll_barrier
|| 1.8% |     19.0 | 51.0 | 73.5% |MPID_nem_gni_poll
|| 1.5% |     15.2 | 39.8 | 72.9% |MPID_nem_gni_check_localCQ
||===============================================================
|  4.5% |     47.0 | -- |      -- |GNI
||-----------------------------------------------------------------
|| 4.3% |     44.8 | 118.2 | 73.1% |GNI_CqGetEvent
||===============================================================
|  4.4% |     45.8 | -- |      -- |ETC
||-----------------------------------------------------------------
|| 2.1% |     21.9 | 65.1 | 75.4% |GNII_DlaProgress
|| 1.0% |     10.4 | 8.6 | 45.8% |_cray_mpi_memcpy_snb
|===============================================================
```

# Profile by Group, Function, and Line

❖ Table 2: Profile by Group, Function, and Line

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

 100.0% | 1,039.1 |    -- |    -- |Total
|-------------------------------------------------------------------
|  72.3% |   751.5 |    -- |    -- |USER
||------------------------------------------------------------------
||  33.5% |   347.9 |    -- |    -- |rhs_
3|    |     |     |       | NPB3.3.1/NPB3.3-MPI/LU/rhs.f
||||-----------------------------------------------------------------
4|||   2.2% |22.5 |  13.5 | 37.8% |line.43
4|||   1.8% |18.2 |   9.8 | 35.2% |line.96
4|||   1.6% |16.4 |  10.6 | 39.6% |line.228
…
```

File rhs.f, line 43

```
do k = 1, nz
    do j = 1, ny
     do i = 1, nx
      do m = 1, 5
       rsd(m,i,j,k) = -
      frct(m,i,j,k)
      end do
     end do
    end do
end do
```

# More information from sampling

❖ Table 3:  Wall Clock Time, Memory High Water Mark (limited entries shown)

```
Process |  Process   |PE=[mmm]
  Time   |   HiMem   |
         | (MBytes)  |

 20.455187 | 39.18 |Total
|----------------------------
| 23.922620 |39.74 |pe.34
| 19.638636 |39.57 |pe.107
| 16.558081 |39.66 |pe.68
|=============================
```

❖ ====================== Additional details =======================

Experiment:  samp_pc_time

Sampling interval:  10000 microsecs

# Automatic Profiling Analysis (APA)

❖ After the previous execution of the command *pat_report* two new files were created with extensions *apa* and *ap2*, the second one will be presented later.

❖ Open the file sampling_report_lu_C_64.apa

#      Collect the default PERFCTR group.

  -Drtenv=PAT_RT_PERFCTR=default

#      Alternatively, energy counters may be added to the default
#      list by commenting out the line above and enabling the
#      line below. Note that this may significantly increase the
#      runtime overhead for high trace counts. The parentheses
#      in the syntax below denote counters that are not available
#      on all platforms.

#  -Drtenv=PAT_RT_PERFCTR=default,(PM_ENERGY:NODE),(PM_ENERGY:ACC)

#      Libraries to trace.

  -g mpi

# Automatic Profiling Analysis (APA) II

\#      Local functions are listed for completeness, but cannot be traced.

 -w  \# Enable tracing of user-defined functions.

\# 33.49% 32799 bytes
     -T rhs_

\# 8.02% 3379 bytes
     -T blts_

\# 7.90% 3863 bytes
     -T buts_

\# 7.81% 14983 bytes
     -T jacld_
…
-o lu.C.128+apa \# New instrumented program.

# Automatic Profiling Analysis (APA) III

❖ In order to create the new binary with regard to APA, execute the following
- *pat_build -O sampling_report_lu_C_64.apa*

  ```
  WARNING: Tracing small, frequently called functions can add excessive overhead.
  WARNING: To set a minimum size, say 1200 bytes, for traced functions, use:
     -D trace-text-size=1200.
  INFO: A total of 7 selected non-group functions were traced.
  INFO: A maximum of 105 functions from group 'mpi' will be traced.
  ```

❖ The new instrumented binary is called *lu.C.64+apa*

❖ Edit the submit.sh file, comment line 16 and uncomment line 19
- *sbatch --reservation=s1001_85 submit.sh*

❖ The new performance file is called  *lu.C.64+apa+PID-XXXt.xf*

❖ Use the tool *pat_report*
- *pat_report -o report_apa_lu_C_64.txt  lu.C.64+apa+PID-XXXt.xf*

❖ Open the file report_apa_lu_C_64.txt

# Performance report I

Table 1:  Profile by Function Group and Function

```
 Time% |     Time |        Imb. |  Imb. |        Calls |Group
       |          |       Time | Time% |       | Function
       |          |       |        |       |  PE=HIDE


 100.0% | 12.081612 |      -- |       -- | 455,387.7 |Total
|------------------------------------------------------------------
|  73.8% |  8.922097 |      -- |       -- | 161,404.0 |USER
||-----------------------------------------------------------------
||  28.6% |  3.450003 | 0.416838 |  10.9% |      253.0 |rhs_
||  10.4% |  1.260820 | 0.153597 |  10.9% |  40,160.0 |buts_
||  10.4% |  1.259256 | 0.144344 |  10.4% |  40,160.0 |blts_
||   7.5% |  0.909228 | 0.122412 |  12.0% |  40,160.0 |jacld_
||   7.1% |  0.861425 | 0.130527 |  13.3% |  40,160.0 |jacu_
||   5.7% |  0.684862 | 0.139784 |  17.1% |       2.0 |ssor_
||   3.7% |  0.451014 | 0.295409 |  39.9% |      508.0 |exchange_3_
||=================================================================
```

# Performance report II

Table 1: Profile by Function Group and Function

```
 Time% |    Time |   Imb. |  Imb.   |  Calls |Group
       |         | Time   | Time% |          | Function
       |         |        |         |          | PE=HIDE

|  17.8% | 2.148878 |      -- |      -- | 293,958.7 |MPI
||--------------------------------------------------------------
|| 11.9% | 1.432456 | 3.029769 | 68.4% | 145,580.0 |MPI_RECV
||  3.8% | 0.465076 | 0.411500 | 47.3% | 146,502.9 |MPI_SEND
||  2.0% | 0.241474 | 1.003594 | 81.2% |     922.9 |mpi_wait
||==============================================================
|   8.4% | 1.010618 |      -- |      -- |     24.0 |MPI_SYNC
||--------------------------------------------------------------
||  8.2% | 0.991427 | 0.991319 | 100.0% |      1.0 |mpi_init_(sync)
|===============================================================
```

❖ If needed disable MPI Sync with
  • *export PAT_RT_MPI_SYNC=0*

# MPI topology

❖ MPI Grid Detection:
There appears to be point-to-point MPI communication in a 8 X 16
grid pattern. The 17.8% 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 Hilbert rank order
from the following table.

| Rank Order | On-Node Bytes/PE | On-Node Bytes/PE% of Total Bytes/PE | MPICH_RANK_REORDER_METHOD |
|---|---|---|---|
| Hilbert | 3.039e+10 | 87.40% | 3 |
| SMP | 2.947e+10 | 84.75% | 1 |
| Fold | 1.685e+10 | 48.46% | 2 |
| RoundRobin | 1.106e+10 | 31.82% | 0 |

❖ Example for 128 MPI processes

0,1,17,16,32,48...
68,84,85,69,70,71…

How to use the new MPI topology file:
1. **cp MPICH_RANK_ORDER.XXX MPICH_RANK_ORDER**
2. **export MPICH_RANK_REORDER_METHOD=3**

# Hardware counters

D1 cache utilization:

All instrumented functions with significant execution time had D1 cache hit ratios above the desirable minimum of 75.0%.

D1 + D2 cache utilization:

All instrumented functions with significant execution time had combined D1 and D2 cache hit ratios above the desirable minimum of 80.0%.

TLB utilization:

All instrumented functions with significant execution time had more than the desirable minimum of 200 data references per TLB miss.

Find more about hardware performance counters

❖ Execute:
  • *pat_help*
  • *counters haswell groups*

# Hardware counters

```
Total
-------------------------------------------------------------------------
  Time%                                100.0%
  Time                                 12.081612 secs
  Imb. Time                                     -- secs
  Imb. Time%                                    --
  Calls                    0.038M/sec               455,387.7 calls
  CPU_CLK_THREAD_UNHALTED:THREAD_P          47,351,574,846
  CPU_CLK_THREAD_UNHALTED:REF_XCLK           2,124,810,371
  DTLB_LOAD_MISSES:MISS_CAUSES_A_WALK            6,686,929
  DTLB_STORE_MISSES:MISS_CAUSES_A_WALK           2,823,391
  L1D:REPLACEMENT                            1,404,754,113
  L2_RQSTS:ALL_DEMAND_DATA_RD                  515,418,048
  L2_RQSTS:DEMAND_DATA_RD_HIT                  197,719,491
  MEM_UOPS_RETIRED:ALL_LOADS                20,512,449,601
  CPU_CLK                  2.23GHz
  TLB utilization          2,156.86 refs/miss        4.21 avg uses
  D1 cache hit,miss ratios         93.2% hits         6.8% misses
  D1 cache utilization (misses)    14.60 refs/miss    1.83 avg hits
  D2 cache hit,miss ratio     77.4% hits          22.6% misses
  D1+D2 cache hit,miss ratio       98.5% hits         1.5% misses
  D1+D2 cache utilization          64.57 refs/miss    8.07 avg hits
  D2 to D1 bandwidth       2,603.843MiB/sec  32,986,755,044 bytes
  Average Time per Call                      0.000027 secs
  CrayPat Overhead : Time          8.0%
```

# Hardware Counters - Description

Hardware performance counter events:

  CPU_CLK_THREAD_UNHALTED:REF_XCLK        Count core clock cycles whenever the clock signal on the specificcore is running (not halted):Cases when the core is unhalted at 100Mhz

  CPU_CLK_THREAD_UNHALTED:THREAD_P        Count core clock cycles whenever the clock signal on the specificcore is running (not halted):Cycles when thread is not halted

  DTLB_LOAD_MISSES:MISS_CAUSES_A_WALK   Data TLB load misses:Misses in all DTLB levels that cause page walks

  DTLB_STORE_MISSES:MISS_CAUSES_A_WALK  Data TLB store misses:Misses in all DTLB levels that cause page walks

  L1D:REPLACEMENT                                       L1D cache:L1D Data line replacements

  L2_RQSTS:ALL_DEMAND_DATA_RD          L2 requests:Any data read request to L2 cache

  L2_RQSTS:DEMAND_DATA_RD_HIT          L2 requests:Demand Data Read requests that hit L2 cache

  MEM_UOPS_RETIRED:ALL_LOADS             Memory uops retired (Precise Event):All load uops retired
  PM_ENERGY:NODE                               Compute node accumulated energy
  CYCLES_RTC                      User Cycles (approx, from rtc)

# Load Balance with MPI Message stats

Table 3:  Load Balance with MPI Message Stats (limited entries shown)

```
Time% |      Time |   MPI Msg | MPI Msg Bytes |  Avg MPI |Group
      |           | Count     |               | Msg Size | PE=[mmm]


 100.0% | 12.081612 | 146,522.9 | 271,667,585.0 | 1,854.10 |Total
|------------------------------------------------------------------
|  73.8% |  8.922097 |       0.0 |           0.0 |     -- |USER
||-----------------------------------------------------------------
||  80.6% |  9.739499 |       0.0 |           0.0 |     -- |pe.26
||  75.8% |  9.160217 |       0.0 |           0.0 |     -- |pe.61
||  45.1% |  5.442844 |       0.0 |           0.0 |     -- |pe.127
||=================================================================
|  17.8% |  2.148878 | 146,522.9 | 271,667,585.0 | 1,854.10 |MPI
||-----------------------------------------------------------------
||  48.8% |  5.891394 |  80,852.0 | 143,737,828.0 | 1,777.79 |pe.127
||  15.5% |  1.874838 | 161,678.0 | 293,895,236.0 | 1,817.78 |pe.43
||  10.5% |  1.263484 | 161,678.0 | 303,691,732.0 | 1,878.37 |pe.26
||=================================================================
|   8.4% |  1.010618 |       0.0 |           0.0 |     -- |MPI_SYNC
||-----------------------------------------------------------------
||  22.0% |  2.653814 |       0.0 |           0.0 |     -- |pe.103
||   7.4% |  0.895974 |       0.0 |           0.0 |     -- |pe.123
||   0.1% |  0.012597 |       0.0 |           0.0 |     -- |pe.0
```

# Load Balance with MPI message stats by caller

Table 4:  MPI Message Stats by Caller (limited entries shown)

```
          MPI | MPI Msg Bytes |  MPI Msg | MsgSz |  16<=  | 256<= | 64KiB<= |Function
          Msg |               |   Count  |  <16   | MsgSz | MsgSz |  MsgSz | Caller
      Bytes%  |               |          | Count  |  <256 | <4KiB |  <1MiB |PE=[mmm]
              |               |          |        | Count |  Count |  Count |
       100.0% | 271,667,585.0     | 146,522.9 | 14.0   |  6.9  | 145,581.3 |   920.8 |Total
      |--------------------------------------------------------------------------
      | 100.0% | 271,667,261.0 | 146,502.9 |   0.0 |   0.9 | 145,581.3 |   920.8 |MPI_SEND
      ||-------------------------------------------------------------------------
      ||  67.5% | 183,314,340.0 |      920.8 |   0.0 |   0.0 |      0.0 |   920.8 |exchange_3_
      3|  67.2% | 182,592,630.0 |      917.1 |   0.0 |   0.0 |      0.0 |   917.1 | rhs_
      4|      |               |       |      |      |      |          | ssor_
      5|      |               |       |      |      |      |          | applu_
      ||||||-----------------------------------------------------------------------
      6|||||  77.2% | 209,848,320.0 |  1,012.0 |   0.0 |   0.0 |    0.0 | 1,012.0 |pe.17
      6|||||  72.4% | 196,732,800.0 |  1,012.0 |   0.0 |   0.0 |    0.0 | 1,012.0 |pe.88
      6|||||  36.2% |  98,366,400.0 |    506.0 |   0.0 |   0.0 |      0.0 |   506.0 |pe.127
```

❖  In order to adjust the size of the MPI eager mode (default 8KB, max value 128KB) according to the MPI message stats, use the following command in your job script, where
   - export MPICH_GNI_MAX_EAGER_MSG_SIZE=131072
   - export MPICH_ENV_DISPLAY=1

# Wall clock and memory high water mark

Table 5:  Wall Clock Time, Memory High Water Mark (limited entries shown)

```
 Process |  Process |PE=[mmm]
     Time | HiMem |
          | (MBytes) |


 20.166938 |       48.25 |Total
|----------------------------
| 23.813279 |       48.70 |pe.98
| 20.039177 |       49.79 |pe.82
| 17.694283 |       49.70 |pe.0
|============================
```

❖ In order to extract the profling information for all the processes and not aggregate data, the pat_report tool can be used as following:
  - *pat_report -s pe=ALL -o sampling_results_all.txt txt lu.C.64+apa+PID-XXXt.xf*
  - *pat_report -s filter_input='pe<=5' ...*
  - *pat_report -s filter_input='pe%2==0' ...*

# Apprentice2

A GUI for the raw data

# How to start with Apprentice2

❖ The *pat_report* tool has created one file with extension *ap2*
  - ls –ltr *.ap2

❖ In order to visualize the performance data
  - Connect to Shaheen II with  "ssh –X …"
  - module load perftools-base/6.3.2
  - app2 lu.C.64+**apa**+PID-XX**t**.ap2

❖ The example of the presentation is for lu.C.128

# Apprentice2 – Generic view

# Apprentice2 – Generic view

# Apprentice2 – Generic view

# Apprentice2 – Profile I

# Apprentice2 – Profile II

# Apprentice2 – Load Balance I

# Apprentice2 – Load Balance II

# Apprentice2 – Load Balance III

KAUST   King Abdullah University of Science and Technology

# Apprentice2 – Activity



KAUST   King Abdullah University of Science and Technology

# Apprentice2 – Call Tree

KAUST   King Abdullah University of Science and Technology

# Apprentice2 – Mosaic I

# Apprentice2 – Mosaic II

# Apprentice2 – Mosaic IV

# Apprentice2 – Mosaic V

# Apprentice2 – Mosaic VI

KAUST   King Abdullah University of Science and Technology

# Apprentice2 – Profile comparison (v6.3.2)

# Apprentice2 – Profile comparison

KAUST   King Abdullah University of Science and Technology

# Detailed instrumentation

❖ Do <span style="color:red">not</span> follow these instructions during the hands-on session

❖ Disable the summary of the performance data and create one file per node

- *export PAT_RT_SUMMARY=0*

- *export PAT_RT_EXPFILE_MAX=0*

- *sbatch --reservation=s001_85 submit.sh*

❖ Expect more overhead, the trace file size can increase from some MB to GB

❖ Create the *ap2* file

- *pat_report –o detailed_report_lu_C_64.txt lu.C.64+apa+PID-XXt*

❖ *Use Apprentice2*

- *app2 lu.C.64+apa+PID-XXt.ap2*

# Detailed instrumentation – Profile

# Detailed instrumentation – Activity over time

KAUST   King Abdullah University of Science and Technology

# Detailed instrumentation – Traffic Report

KAUST   King Abdullah University of Science and Technology

# Detailed instrumentation – Traffic Report with links

# Detailed instrumentation – Plots

# Detailed instrumentation – Counters Plot

# Reveal

A tool to port your application to OpenMP

# Reveal

❖ Reveal is Cray's next-generation integrated performance analysis and code optimization tool.

- Source code navigation using whole program analysis (data provided by the Cray compilation environment only)
- Coupling with performance data collected during execution by CrayPAT. Understand which high level serial loops could benefit from parallelism.
- Enhanced loop mark listing functionality.
- Dependency information for targeted loops
- Assist users optimize code by providing variable scoping feedback and suggested compile directives.

# Prepare for Reveal

## ❖ Load Perftools

- *module unload darshan*
- *module load perftools-base/6.3.2*
- *module load perftools/6.3.2*

❖ Compile the code
  - *cd performance_workshop/NPB3.3-MPI_reveal*
  - *make clean*
  - In *the config.make.def* file
    - *MPIF77 = ftn -h profile_generate -hpl=npb_lu.pl -h noomp -h noacc*
    - *FMPI_LIB  = -h profile_generate -hpl=npb_lu.pl -h noomp -h noacc*
  - *make LU NPROCS=64 CLASS=C*
    - *"WARNING: PerfTools is saving object files from a temporary directory into directory…"*
  - *cd bin*
❖ The new binary is called lu.C.64 is not instrumented yet

# Prepare and load Reveal

❖ **Prepare the binary for tracing**
- *pat_build –w lu.C.64*

❖ **Uncomment the line 16 in file submit.sh (the one with lu.C.64+pat)**

❖ sbatch --reservation=s1001_85 submit.sh

❖ pat_report -o reveal.txt lu.C.64+pat+PID-XXt.xf

❖ reveal ../LU/npb_lu.pl ./lu.C.64+pat+PID-XXt.ap2

# Reveal – Loop Performance

# Reveal – Loop performance – Potential Speedup

# Reveal – Scoping

# Reveal – Scoping Results

# Reveal – OpenMP Directives

# Reveal – Compiler messages

# Summary

- ❖ Craypat seems easy to use
- ❖ The user should be careful though
- ❖ Studying in detail the communication with Craypat is difficult
- ❖ Reveal tool could be really helpful
- ❖ Probably other tool(s) could be used for more detailed analysis
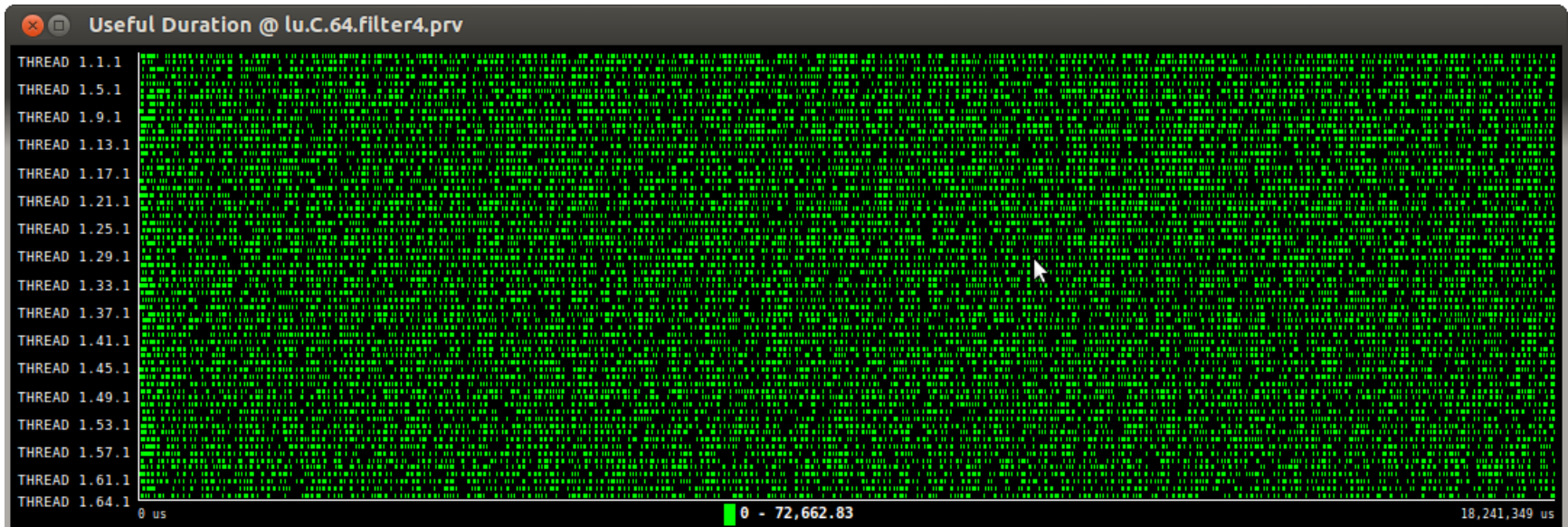
# Extrae/Paraver

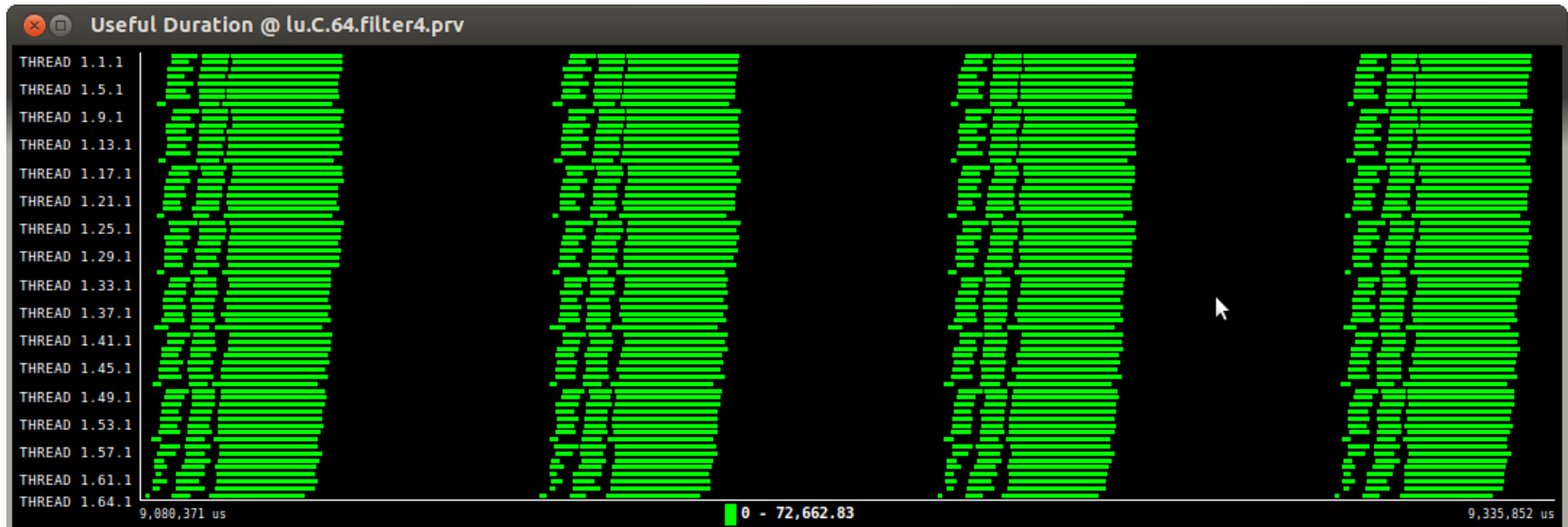A profiling tool from Barcelona Supercomputing Center

# Extrae/Paraver (briefly)

❖ Instrumentation tool from Barcelona Supercomputing Center

❖ The main details are defined in an XML file

❖ For dynamic compilation a wrapper and LD_PRELOAD is enough

❖ For static compilation, linking is necessary

❖ Need to compile with at least -g option and -finstrument-functions for functions instrumentation with Intel and GNU compilers

❖ The trace for LU.C.64 is around to 5 GB

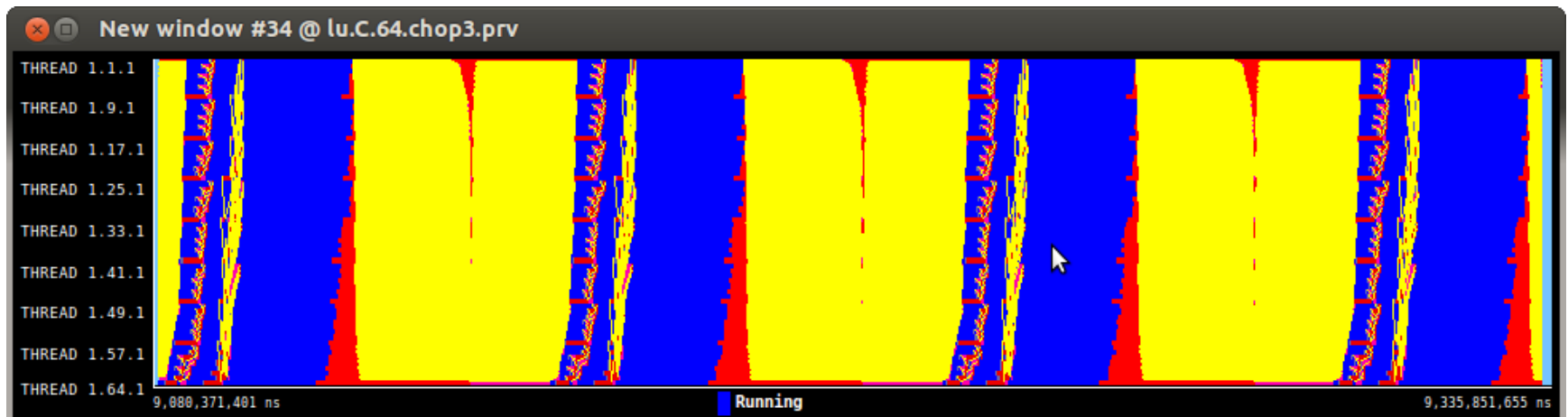❖ Paraver is the tool to visualize and handle the traces from Extrae

# Paraver – Useful duration I
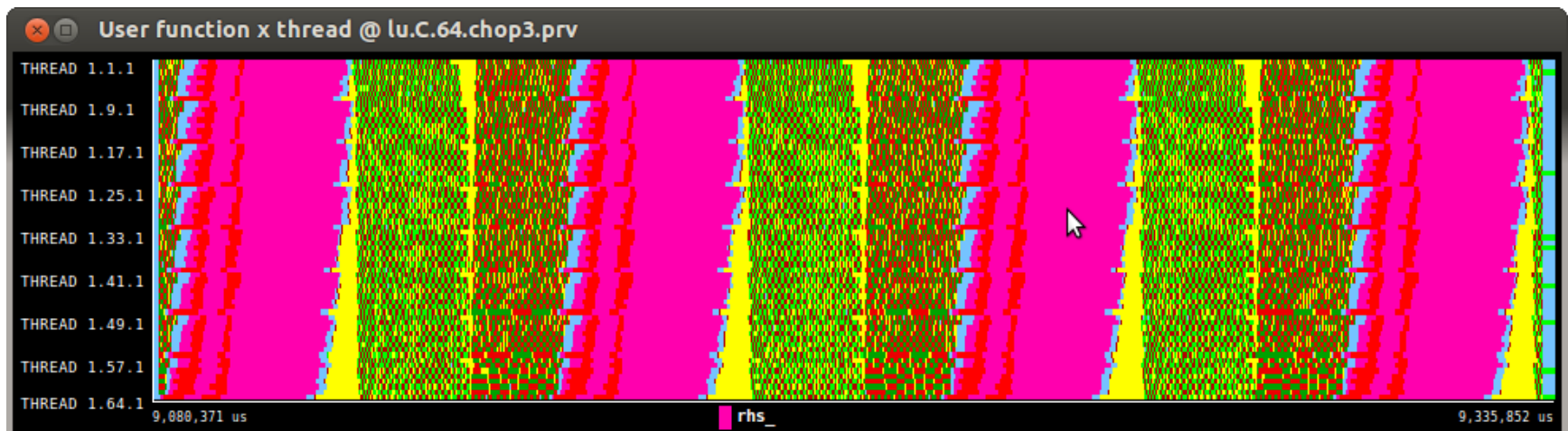
# Paraver – Useful duration II - zoom
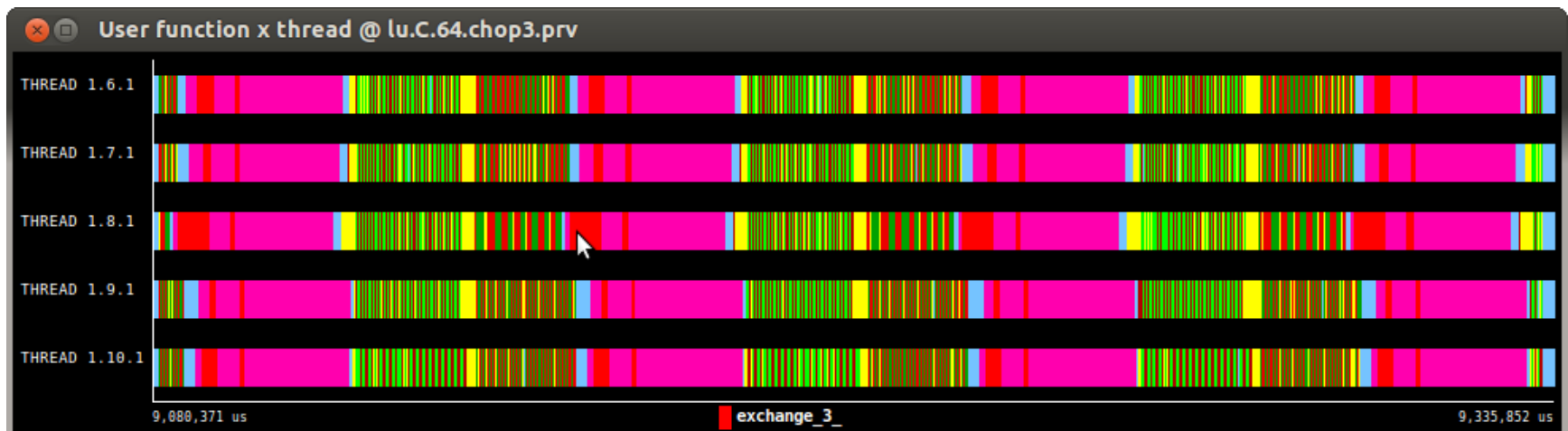
# Paraver – Visualize events

# Paraver – User functions
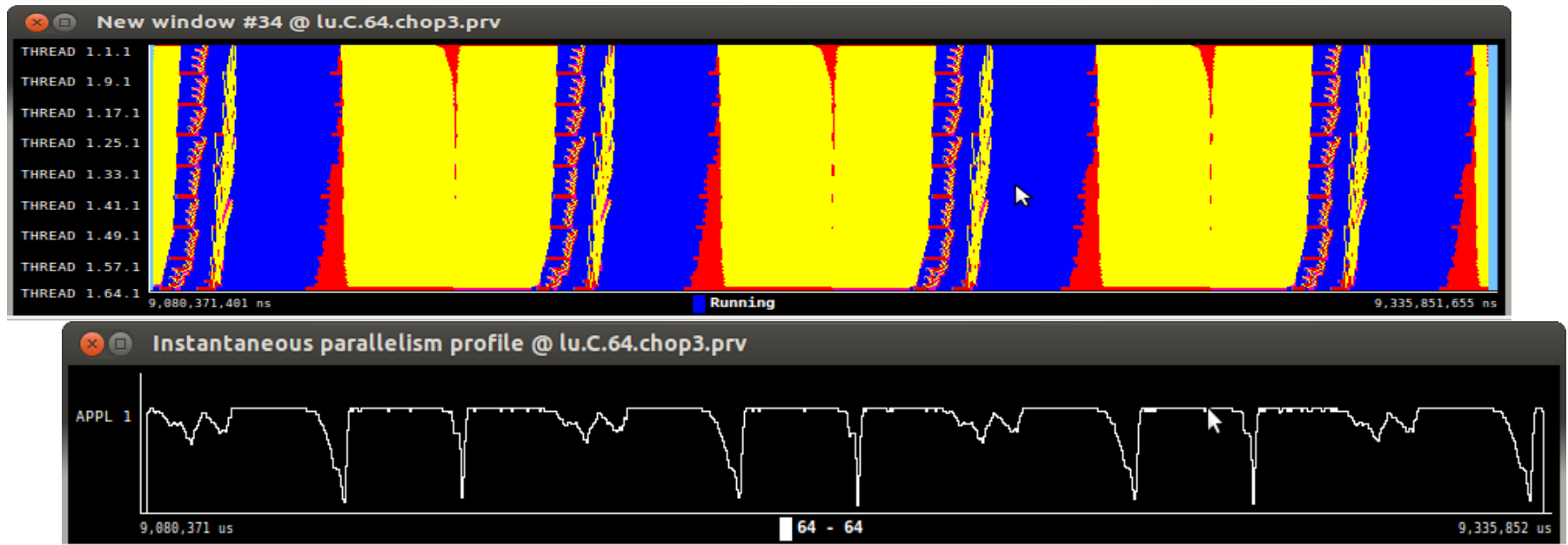
# Paraver – User Functions Profile



| | exchange_3_ | rhs_ | jacld_ | blts_ | exchange_1_ | jacu_ | buts_ | |
|---|---|---|---|---|---|---|---|---|
| **THREAD 1.1.1** | 8 | 12 | 503 | 1,508 | 2,044 | 519 | 1,557 | |
| **THREAD 1.2.1** | 8 | 12 | 502 | 1,503 | 2,040 | 518 | 1,557 | |
| **THREAD 1.3.1** | 8 | 12 | 500 | 1,500 | 2,027 | 513 | 1,539 | |
| **THREAD 1.4.1** | 8 | 12 | 500 | 1,497 | 2,023 | 512 | 1,536 | |
| **THREAD 1.5.1** | 8 | 12 | 499 | 1,494 | 2,020 | 511 | 1,536 | |
| **THREAD 1.6.1** | 8 | 12 | 498 | 1,493 | 2,018 | 511 | 1,533 | |
| **THREAD 1.7.1** | 8 | 12 | 497 | 1,490 | 2,015 | 510 | 1,533 | |
| **THREAD 1.8.1** | 8 | 12 | 497 | 1,488 | 1,992 | 499 | 1,500 | |
| **THREAD 1.9.1** | 8 | 12 | 499 | 1,496 | 2,035 | 518 | 1,557 | |
| **THREAD 1.10.1** | 8 | 12 | 498 | 1,491 | 2,031 | 518 | 1,554 | |
| **THREAD 1.11.1** | 8 | 12 | 497 | 1,490 | 2,019 | 512 | 1,539 | |
| **THREAD 1.12.1** | 8 | 12 | 497 | 1,488 | 2,016 | 511 | 1,536 | |
| **THREAD 1.13.1** | 8 | 12 | 496 | 1,487 | 2,014 | 511 | 1,533 | |
| **THREAD 1.14.1** | 8 | 12 | 495 | 1,484 | 2,011 | 510 | 1,533 | |
| **THREAD 1.15.1** | 8 | 12 | 495 | 1,482 | 2,009 | 510 | 1,530 | |
| **THREAD 1.16.1** | 8 | 12 | 494 | 1,481 | 1,986 | 499 | 1,497 | |

# Paraver – Timeline selecting specific MPI processes

# Paraver – Instantaneous parallelism profile

**help@hpc.kaust.edu.sa**

**www.hpc.kaust.edu.sa**

**Twitter:
KAUST_HPC**

**Thank You!**