King Abdullah University of Science and Technology | Core Labs and Research Infrastructure

# Programming Environment, Performance & Debugging tools

## Bilel Hadri

bilel.hadri@kaust.edu.sa

**Computational Scientist**

**KAUST Supercomputing Lab**

@mnoukhiya   @KAUST_HPC

King Abdullah University of Science and Technology | SHAHEEN SUPERCOMPUTING LABORATORY

# Shaheen 2 Cray XC40

- Edit and Compile only your code on login. To run, submit jobs.

- **Compiler available Cray CCE (default), Intel and GNU supported**

  - Compiler wrappers for serial and parallel
    - **ftn for Fortran code**
    - **cc for C code**
    - **CC for C++ code**

  - Do not purge.



```
Use of this system is limited to users who have been properly authorised by
the KAUST Supercomputing Laboratory. Unauthorised users must disconnect
immediately.

For support, see http://www.hpc.kaust.edu.sa/
or email help@hpc.kaust.edu.sa
Last login: Mon Sep 12 18:36:13 2022 from 10.200.0.112
```



```
Shaheen is a 36 rack Cray XC40 system. The front-end environment is
running SUSE Linux Enterprise Server 15.
                        ---------------------------------
 ?cd12:~> module list
Currently Loaded Modulefiles:
  1) modules/3.2.11.4                 8) pmi/5.0.17                       15) rca/2.2.20-7.0.3.1_3.18__g8e3fb5b.ari    22) eproxy/2.0.24-7.0.3.1_3.9__g8e04b33.ari
  2) craype-network-aries             9) dmapp/7.1.1-7.0.3.1_3.21__g93a7e9f.ari    16) atp/3.14.5                      23) craype-haswell
  3) cce/12.0.3                      10) gni-headers/5.0.12.0-7.0.3.1_3.9__gd0d73fe.ari  17) perftools-base/21.09.0     24) xalt/1.1.2
  4) craype/2.7.10                   11) xpmem/2.2.27-7.0.3.1_3.10__gada73ac.ari   18) PrgEnv-cray/6.0.10              25) darshan/3.3.1
  5) cray-libsci/20.09.1             12) job/2.2.4-7.0.3.1_3.17__g36b56f4.ari       19) cray-mpich/7.7.18              26) ksl/ksl
  6) udreg/2.3.2-7.0.3.1_3.16__g5f0d670.ari    13) dvs/2.12_2.2.224-7.0.3.1_3.14__gc77db2af   20) slurm/slurm
  7) ugni/6.0.14.0-7.0.3.1_6.4__g8101a58.ari   14) alps/6.6.67-7.0.3.1_3.21__gb91cd181.ari    21) dws/3.0.36-7.0.3.1_3.19__g6985c90.ari
```

# Compiler Driver Wrappers

Use them exactly like you would use the original compiler, e.g. To compile.

```
#to use Cray  compilers
ftn -o myexe myprog.f90 # Fortran
cc -o myexe myprog.c # for C
CC -o myexe myC++code.C # for C++
```

```
#to use Intel compilers


module swap PrgEnv-cray PrgEnv-intel
ftn -o myexe myprog.f90 # Fortran
cc -o myexe myprog.c # for C
CC -o myexe myC++code.C # for C++
```

```
#to use GNU compilers


module swap PrgEnv-intel PrgEnv-gnu
ftn -o myexe myprog.f90 # Fortran
cc -o myexe myprog.c # for C
CC -o myexe myC++code.C # for C++
```

ftn, cc, and CC, are not Cray compilers; they invoke the Intel, GNU, or Cray compilers under

the hood, depending on the loaded programming environment module (PrgEnv-xxx)

No need to call for mpicc/mpif90... only ftn/cc/CC

# Compilers

- Intel - better chance of getting processor specific optimizations

- Cray compiler – many new features and optimizations, especially with Fortran; useful tools like reveal work with Cray compiler only

- GNU - widely used by open software

- More information from compilers options on the man page

| PrgEnv | Description | Real Compilers |
|---|---|---|
| PrgEnv-cray | Cray Compilation Environment | crayftn, craycc, crayCC |
| PrgEnv-intel | Intel Composer Suite | ifort, icc, icpc |
| PrgEnv-gnu | GNU Compiler Collection | gfortran, gcc, g++ |

# Compilers

- Use ftn, cc, and CC to compile instead of the underlying native compilers (ifort, icc, icpc, gfortran, gcc, g++..)

- Use same wrapper even for MPI codes. Do not use mpicc/mpif90....


- Default compiling is dynamic on Shaheen
  - just add the -static flag to the command and link lines,
  - or set CRAYPE_LINK_TYPE=static in the environment


- Compiler wrappers do cross compilation
  - Compiling on login nodes to run on compute nodes
  - One may run into trouble with GNU  automake  or  cmake.
  - Add the  specifier –host=x86_64-unknown-linux-gnu  for the  configure tool .


- By default, Cray C/C++ is using CLANG

# OpenMP

OpenMP is supported by all of the PrgEnvs.

| PrgEnv | Enable OpenMP |
|---|---|
| PrgEnv-cray | C/C++: -fopenmp<br>Fortran: -h omp |
| PrgEnv-intel | -qopenmp |
| PrgEnv-gnu | -fopenmp |

# Cray Scientific Libraries

- Compiler wrappers takes care of not only the compiler but also libs like BLAS, LAPACK, SCALAPACK, MPI,..

- Cray Scientific Libraries package, LibSci, is a collection of numerical routines optimized for best performance on Cray systems.
  - LibSci is loaded by default and this is for all programming environment
  - No user flags or options are required for compiling or linking.
  - LibSci library collection contains; BLAS, BLACS, LAPACK, ScaLAPACK, IRT, CRAFFT, CASE, FFT, FFTW2, FFTW3

- Both cray-python and cray-R call the OpenMP threaded version of cray-libscic calls the OpenMP threaded version of cray-libsci.
  - It is recommended to set the number of desired threads with the OMP_NUM_THREADS environment variable.

# Cray Scientific Libraries

- FFTW: Cray's main FFT library is FFTW from MIT with some  additional optimizations for Cray hardware

- Cray PETSc , Cray Trilinos.

- Just need to module load and compile your code

- No need to put the whole path of the libraries

- Cray-python, cray-R .
  - Just load the module and use the tools

# Cray PE DL Plugin

- craype-dl-plugin - introduces the Cray PE DL Plugin for accelerating distributed deep learning DESCRIPTION The Cray PE DL Plugin provides a highly tuned communication layer that can be easily added to any deep learning framework.

- Plugin has both a C and Python 3 API and supports multiple DL datatypes

- Compatible with TensorFlow and PyTorch frameworks

- Can be used with popular DL frameworks or integrated into a project via its API


- module load craype-dl-plugin

- man intro_dl_plugin

# Cray Scientific Libraries

- Cray TPSL (Third Party Scientific Libraries) contains a collection of outside mathematical libraries that can be used with PETSc and Trilinos
  - The TPSL increase the flexibility of PETSc and Trilinos by providing users with multiple options for solving problems in dense and sparse linear algebra
  - The cray-tpsl module is automatically loaded when PETSc or Trilinos is loaded. The libraries included are MUMPs, SuperLU, SuperLU_dist, ParMetis, Hypre, Sundials, and Scotch.

- Intel MKL: The Intel Math Kernel libraries is an alternative to LibSci
  - Features tuned performance for Intel CPUs as well
  - Linking is quite complicated but with Intel compilers (PrgEnv-intel) is usually straightforward using the Intel Link advisor
  - http://software.intel.com/en-us/articles/intel-mkl-link-line-advisor

# Modules available on Shaheen

# Shaheen 2 Cray XC40 flyer

- **Get started with the flyer**

https://hpc.kaust.link/Shaheen_flyer

# CDT: updates

- In order, KSL will be updating the Cray Programming Environment, namely the Cray Developer Toolkit (CDT), to provide a predictable, stable and consistent programming environment while still making necessary software updates.

- Using more recent packages may result in faster execution of the code.

- CDT consists of compilers, MPI, scientific and I/O libraries, profiling and debugging tools, etc.

- New CDT software will be installed at least twice a year. The new versions will not be made the defaults when installed. You need to load them.

- Module load cdt/21.09 is the default.
  - cdt/22.09 is available

# Performance

# Why Performance Analysis ?

- You want to get the best expected performance.
  - Ex: Internet Bandwidth, RPM vehicles
  - Need to identity the issue

- Economic: TIME is MONEY
  - Lifetime of HPC systems is short (4/5 years)
  - Large HPC machines cost in O($10M)

- Qualitative: Do more science
  - Get codes run faster
  - Perform more time steps
  - Simulation higher resolutions

- **Must strive to evaluate how your code is running.**

- **Learn to think of performance during the entire cycle of**

250 km        100 km        28 km

*© Karlsruhe Institute of Technology*

# Typical Performance Analysis Procedure

- Measuring the wallclock time is not enough.

- Need to know what's really happening under the hood.



- Do I have a performance problem at all?
  - Time / speedup / scalability

- What is the key bottleneck ?
  - computation / communication

- Where is the key bottleneck?
  - Detailed profiling

- Why does the code have scalability problems?
  - Load imbalance analysis, compare profiles at various sizes function-by - function, performance modeling

# Performance measurement

- No single solution is sufficient
  - Timing manually… Not efficient and accurate
  - Don't reinvent the wheel


- Need to use a combination of different methods, tools and techniques is needed!
  - Measurement Sampling and profiling
  - Analysis Statistics, visualization, automatic analysis, data mining, …

# Performance/Monitoring tools

- Many tools are available on HPC systems:
    - Gprof
    - PAPI
    - VTUNE
    - Allinea/ARM  Tools
    - VAMPIR
    - TAU
    - Scalasca
    - Likwid
    - VAMPIR
    - HPCToolkit
    - Paraver/Extrae
    - Darshan
    - Perftools ( Cray systems)

It doesn't matter how many resources you have...

If you don't know how to use them, it will never be enough.

# Profile a Python code

- Just type:
  - python -m cProfile myscript.py

- For call graphs
  - pycallgraph graphviz -- ./myscript.py
  - Display pycallgraph.png

# Overview on performance of code over time

# ARM/DDT

- ARM performance Tools

- Provides quick overview of performance issues:
  - The time spent in various categories of instruction: memory access, numeric operations, floating point operations
  - Overview on I/O, Memory, Communication, Threads , Energy usage
  - Energy Saves data in HTML, CVS or text form

- To get the report in html or txt
  - Load arm-reports module
  - make-profiler-libraries
  - Relink dynamically your code as shown in the output
  - perf-report srun –n 2 ./mycode

# ARM/DDT general Overview

**arm**
PERFORMANCE
REPORTS

Command:      srun wave.exe
Resources:    4 nodes (32 physical, 64 logical cores per node)
Memory:       126 GiB per node
Tasks:        4 processes
Machine:      nid00024
Start time:   Fri Feb 23 2018 08:29:34 (UTC+03)
Total time:   121 seconds (about 2 minutes)
Full path:    /lustre/project/k01/hadrib/allinea_workshop/
              1_reporting/f90

## Summary: wave.exe is Compute-bound in this configuration

**Compute**  93.6%   Time spent running application code. High values are usually good.
This is **very high**; check the CPU performance section for advice

**MPI**  6.4%   Time spent in MPI calls. High values are usually bad.
This is **very low**; this code may benefit from a higher process count

**I/O**  0.0%   Time spent in filesystem I/O. High values are usually bad.
This is **negligible**; there's no need to investigate I/O performance

This application run was Compute-bound. A breakdown of this time and advice for investigating further is in the CPU section below.

As very little time is spent in MPI calls, this code may also benefit from running at larger scales.

# ARM/DDT Detailed

## CPU

A breakdown of the 93.6% CPU time:

| | | |
|---|---|---|
| Scalar numeric ops | 28.6% | |
| Vector numeric ops | 0.0% | |
| Memory accesses | 71.4% | |

The per-core performance is memory-bound. Use a profiler to identify time-consuming loops and check their cache performance.

No time is spent in vectorized instructions. Check the compiler's vectorization advice to see why key loops could not be vectorized.

## MPI

A breakdown of the 6.4% MPI time:

| | | |
|---|---|---|
| Time in collective calls | 0.8% | |
| Time in point-to-point calls | 99.2% | |
| Effective process collective rate | 470 kB/s | |
| Effective process point-to-point rate | 2.34 MB/s | |

Most of the time is spent in point-to-point calls with a very low transfer rate. This suggests load imbalance is causing synchronization overhead; use an MPI profiler to investigate.

## I/O

A breakdown of the 0.0% I/O time:

| | | |
|---|---|---|
| Time in reads | 0.0% | |
| Time in writes | 0.0% | |
| Effective process read rate | 0.00 bytes/s | |
| Effective process write rate | 0.00 bytes/s | |

No time is spent in I/O operations. There's nothing to optimize here!

## Threads

A breakdown of how multiple threads were used:

| | | |
|---|---|---|
| Computation | 0.0% | |
| Synchronization | 0.0% | |
| Physical core utilization | 3.1% | |
| System load | 3.1% | |

No measurable time is spent in multithreaded code.

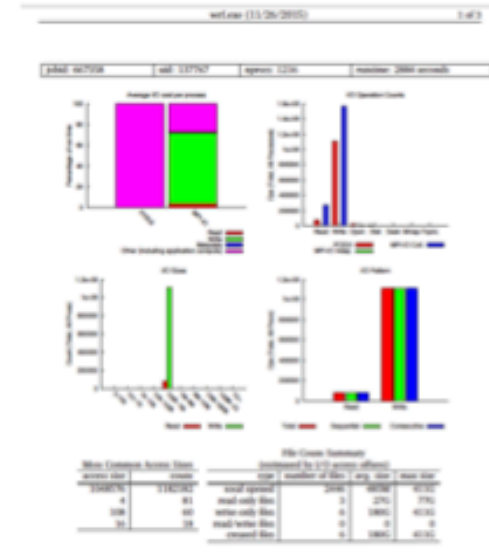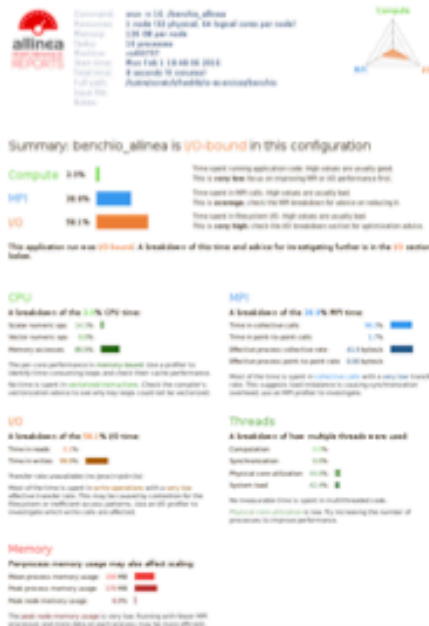Physical core utilization is low. Try increasing the number of processes to improve performance.

## Memory

Per-process memory usage may also affect scaling:

| | | |
|---|---|---|
| Mean process memory usage | 31.0 MiB | |
| Peak process memory usage | 31.2 MiB | |
| Peak node memory usage | 1.0% | |

The peak node memory usage is very low. Running with fewer MPI processes and more data on each process may be more efficient.

## Energy

A breakdown of how the 17.0 Wh was used:

| | | |
|---|---|---|
| CPU | 69.6% | |
| System | 30.4% | |
| Mean node power | 128 W | |
| Peak node power | 151 W | |

Significant time is spent waiting for memory accesses. Reducing the CPU clock frequency could reduce overall energy usage.

# How is my IO ?

- Use profiling and charactirization tools
  - Allinea report,
  - Craypat profiling
  - Darshan
  - Contact CS team at KSL

Table 1:

| Time% | | Time | Imb. | Imb. | | Calls | Functionnction |
| | | | Time | Time% | | | Sourceurce |
| | | | | | | | Linene |

| 100.0% | 13,461.594081 | -- | -- | 666,344.0 | Total |
|--------|---------------|------|------|-----------|-------|
| 32.1% | 4,326.121649 | -- | -- | 3,072.0 | mpi_barrier_(sync) |
| 24.4% | 3,284.591116 | -- | -- | 48,630.0 | MPI_FILE_WRITE_ALL |
| 14.0% | 1,884.152065 | -- | -- | 71,930.0 | h5dwrite_c_ |
| 3 | | | | | line.334 |
| 12.7% | 1,704.005636 | -- | -- | 88,516.0 | nc4_put_vara_tc |
| 3 | | | | | line.1431 |
| 9.9% | 1,338.717666 | -- | -- | 49,539.0 | write_var |
| 3 | | | | | line.2262 |
| 3.0% | 397.666538 | -- | -- | 128.0 | mpi_init_(sync) |

======================== Additional details ========================

# Debugging

# Valgrind4hpc

- Valgrind4hpc debugging tool helps in the detection of memory leaks and errors in parallel applications.

- **Compile and link with  -g option** , then allocate and follow the steps shown bellow.

```
salloc -N 1 module

unload darshan xalt

module load valgrind4hpc

export CTI_WLM_IMPL=slurm export CTI_LAUNCHER_NAME=srun

valgrind4hpc -n2 --launcher-args="--hint=nomultithread --ntasks=2" --valgrind-args="--track-
origins=yes --leak-check=full" ./my_exe
```

- Here is a clean output. Otherwise, follow the instructions to detect the memory leaks:

```
RANKS: <0,1>

HEAP SUMMARY: in use at exit: 0 bytes in 0 blocks

All heap blocks were freed -- no leaks are possible ERROR SUMMARY: 0 errors from 0 contexts
(suppressed 19)
```

- To run your program and debug it across multiple nodes, allocate the desired number of nodes and then update accordingly the parameters in the launcher-args similar to the option for the srun/sbatch script.

- More information is available in the man pages of valgrind and valgrind4hpc.

# Different tools available

- Several tools for C/C++/Fortran debugging tools:
  - gdb4hpc
  - valgrind4hpc
  - sanitize4hpc

# gdb4hpc

- gdb4hpc (Cray Line Mode Parallel Debugger) is a GDB-based parallel debugger, developed by Cray.

- It can debug with CCE, PGI, GNU and Intel Fortran, C and C++ compilers.

- **gdb4hpc** also includes comparative debugger technology that enables programmers to compare data structures between two executing applications. Cray, however, recommends accessing the comparative debugger technology through the new Cray Comparative Debugger (CCDB) with graphical user interface (GUI) that enhances the parallel debugging capabilities of **gdb4hpc**.

- More info in man pages

```
module load gdb4hpc
```

**Note: need to unload xalt**

# Debugging with ARM/DDT

```
> ftn -g -O0 -qopenmp -o exe_mpi_debug test_hib.f90
> salloc -N 1
> module load arm-forge/19.1.4
> ddt exe_mpi_debug
```

# Debugging with ARM/DDT

**Use the GUI to navigate within the code and check the variable**

# Debugging with Totalview

- Compile with –g option as usual
- Allocate the node

```
>module load totalview
>tv8 &
```
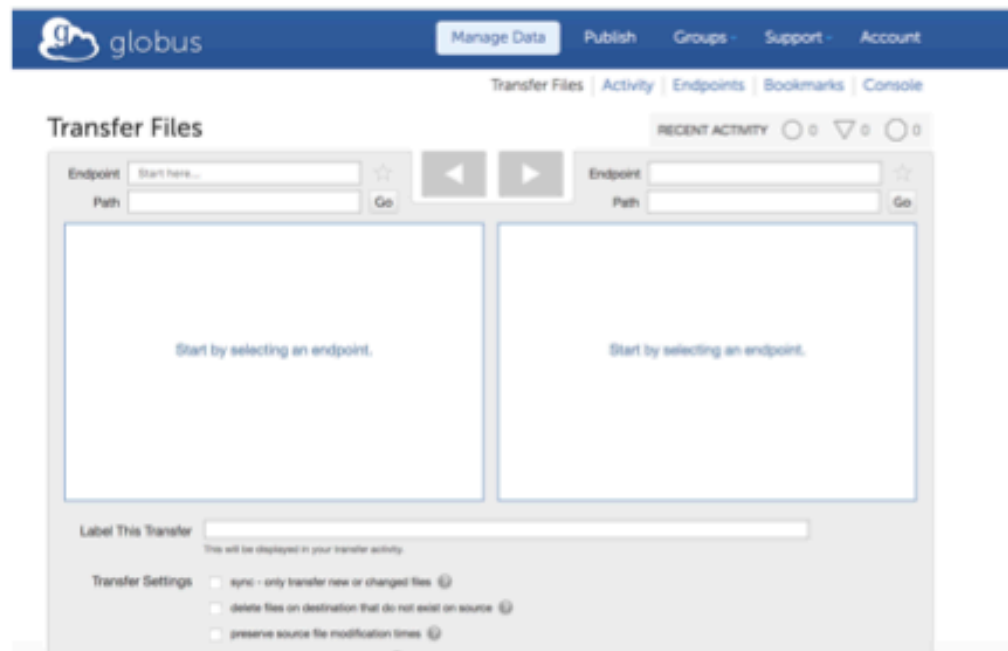
# Transferring Files

# Transfer files: Use Globus

- scp/ftp for small size ( in order of KB)

- For larger file, use Globus, especially for moving data in & out of Shaheen http://www.globus.org/ ( Free)
  - Reliable & easy-to-use web-based service:
  - Email notification of success or failure

- Globus extensive documentation  https://docs.globus.org
  - Web based interaction with service
  - REST/API for scripted interactions with service
  - Globus Connect Server & Personal for setting up additional remote  endpoints such your personal laptop/ workstation

- Globus on Shaheen. Look for Shaheen End point Point
  - Within Campus: choose dm2.hpc.kaust.edu.sa
  - Outside Campus: choose dm1.hpc.kaust.edu.sa

# Globus

**Core Labs and Research Infrastructure**

- Connect to globus.org
- Sign in or create an account
- Use Shaheen  dm2 when inside KAUST and dm1 when connected externally.

# Transferring with Globus

# File Status notification ( email and web-interface)

Core Labs and
Research Infrastructure

## Activity

☰ Task List

✔ NERSC Cori to Shaheen cdl2
transfer completed a month ago

| ⓘ Overview | ☰ Event Log |
|---|---|

| | | | | |
|---|---|---|---|---|
| Task ID | ec24a144-2a81-11e8-b7fa-0ac6873fc732 | | Files | 1 |
| Owner | Bilel Hadri (hadri@globusid.org) | | Directories | 0 |
| Source | NERSC Cori ⓘ | | Bytes Transferred | 103.57 MB |
| | owner: nersc@globusid.org | | Effective Speed | 6.07 MB/s |
| Destination | Shaheen cdl2 ⓘ | | Pending | 0 |
| | owner: shaheen@globusid.org | | Succeeded | 2 |
| Condition | SUCCEEDED | | Cancelled | 0 |
| Requested | 2018-03-18 10:56 am | | Expired | 0 |
| Completed | 2018-03-18 10:57 am | | Failed | 0 |
| Transfer Settings | • verify file integrity after transfer | | Retrying | 0 |
| | • transfer is not encrypted | | Skipped | 0 |
| | • overwriting all files on destination | | | view debug data |

# Tips & Summary

# Best Practices for Performance (1)

- Check the system details thoroughly
  - Never assume ! ( Login nodes different than compute)

- Choose a compiler and MPI to build your application
  - All are not same !  Rely on the latest versions

- Start with some basic compiler flags and try additional flags one at a time
  - Optimization is incremental !  Benchmarking and testing is a must

- Use the built-in/optimized  libraries and tools to save time and improve performance
  - Libraries Tools are your friends !
  - By doing the different steps of optimizations:
  - You can achieve huge speedup ( o(10x ) and more ) by using Optimized Mathematics libraries (Cray, MKL)
  - Optimizing the cache and memory

# Best Practices for Performance (2)

- Don't Reinvent the wheel ! Several tools are available for debugging and performance

- Test your application at every level to arrive at an optimized code
  - Check correctness !

- Customize your runtime environment to achieve desired goals
  - Play with the number of threads, memory and core affinity

- Profile and adjust optimization and runtime environments accordingly
  - Start with small and short runs

- **READ the manual and/or attend the tutorials/workshops !**

- **Visit https://www.hpc.kaust.edu.sa/training**

# Best Practices

- Use adequately your allocation
  - Check your core hours, sb kxxxx ,sb_user kxxxx
  - Check your quota usage kuq, kpq
  - Prepare in advance the project proposal


- Shaheen is a shared resource
  - Be kind to your neighbor users
  - Don't run on login.


- Follow the terms and conditions
  - Don't share your account with others.

# KSL CS Team

- Need help: send a ticket [help@hpc.kaust.edu.sa](mailto:help@hpc.kaust.edu.sa)
  - Help us to help you :D
    - Provide details.
    - Which HPC system?
    - What is the problem? When did it happen? What modules were loaded? How did you try to fix or work around it? Send the error and job script.


- **Acknowledge KAUST Supercomputing Lab and HPC resources used in your papers.**

**bilel.hadri@kaust.edu.sa**

@mnoukhiya    @KAUST_HPC

Thank You !    ! شكراً