# Shaheen III: Storage

Bilel Hadri, KAUST Supercomputing Lab

February 4th, 2025

جامعة الملك عبدالله
للعلوم والتقنية
King Abdullah University of
Science and Technology

# Shaheen Storage

"A supercomputer is a device for converting a CPU-bound problem into an I/O bound problem."
[Ken Batcher]

- **HPC systems consist of three main components:**
  - **Compute nodes**
  - **High-speed interconnect**
  - **I/O infrastructure**

- **Most optimization work on HPC applications is carried out on**
  - **Single node performance across many cores**
  - **Network performance ( communication)**
  - **I/O only when it becomes a real problem ( or when the sys-admin contact us...)**

2

# Shaheen Storage

- I/O is commonly used by scientific applications to achieve goals like:
    - storing numerical output for later analysis
    - loading initial conditions or datasets for processing
    - checkpointing to files that save the state of an application in case of system failure

- On Supercomputer, you will need parallel file system

3

# Shaheen Storage Why do we need parallel I/O?

- **Imagine a 24 hours simulation on 16 cores.**
  - o **1% of run time is serial I/O.**
  - o **You get the compute part of your code to scale to 6 nodes, ie 1152 cores.**
  - o **Up to 72x speedup in compute: I/O is 39% of run time ( 22'16" in computation and 14'24" in I/O).**

- **Parallel I/O is needed to**
  - o **Spend more time doing science**
  - o **Not waste resources**
  - o **Prevent affecting other users**
  - o **Easier to use and get the performance without major change in the workflow.**

# Shaheen III Hardware Specifications: Storage

| Characteristics<br>Tiered storage | Shaheen III /scratch storage<br>/scratch/username |
|---|---|
| Total Capacity (usable) | 32 PB |
| Capacity tier (HDD) | 25 PB |
| Capacity tier perf Read/Write | 330/260 GB/s |
| BW tier capacity | 6.8PB |
| BW Perf. tier Read/Write | 3750/2500 GB/s |
| IOPS tier capacity | 338 TB |
| IOPS tier IOPS (Read/Write) | 10+M IOPS |

| Characteristics | Shaheen III /projects<br>/project/kxxxxx |
|---|---|
| Total Capacity (usable) | 57 PB* |
| Capacity perf Read/Write | ~200/120 GB/s |

# Shaheen Storage

- **/scratch**
  - Read + Write on login and Compute nodes
  - Fast file system

- **/project**
  - Read only on Compute nodes
  - Read + Write on PPN and DTN
  - Keeping large file for the lifetime of the project

- **/home**
  - Only available on login nodes

- **https://docs.hpc.kaust.edu.sa/policy/shaheen3.html**

6

# Tiered Storage /scratch

Shaheen 3 uses tiered storage for Scratch (80TB quota per PI)

- Capacity Tier: default
    /scratch/username  or cd $SCRATCH
    10 TB quota per user

- Bandwidth Tier
    /scratch/username/bandwidth cd $SCRATCH_BW
    1 TB quota per user
    For files which do require high bandwidth, large files...

- IOPs Tier
    /scratch/username/iops or cd $SCRATCH_IOPS
    50 GB quota per user
    For small files with high IOPs pattern requirements

# Shaheen Data Transfer

- Data transfer between tiers: just **copy** the data with the `cp` command.

- Data transfer between scratch <-> project: SLURM job script:

```
user@login4:> cat dtn_copy_job.sh
#!/usr/bin/env bash
#SBATCH -N 1
#SBATCH -n 1
#SBATCH -t 20
#SBATCH -p dtn
#SBATCH --output=SLURM_LOGS/copy_out.txt
#SBATCH --error=SLURM_LOGS/copy_err.txt
#SBATCH --job-name=COPY_project_to_scratch
rsync <options> -v --progress /lustre2/project/k????? /scratch/user/
```

- For outside Shaheen, use scp command
`scp -r myusername@shaheen.hpc.kaust.edu.sa:/path/directory` .

- For large files, Globus via dtn6. Contact help@hpc.kaust.edu.sa for any guidance.

# Storage Best Practices

- /scratch meant for temporary use during the lifetime of the job.

- /project for persistent storage during the lifetime of the project

- Users must delete their unused data from scratch and move necessary files to project filesystem.

- For overall quota on scratch:
```
lfs quota -uh $USER /scratch
```

- For quota on capacity tier on scratch:
```
lfs quota -uh $USER --pool capacity /scratch
```

- For quota on bandwidth tier on scratch:
```
lfs quota -uh $USER --pool bandwidth /scratch
```

- For quota on IOPS tier on scratch:
```
lfs quota -uh $USER --pool iops /scratch
```

9

# Storage Best Practices 2

- My usage of storage quota: command kuq

```
kuq
-----------------------------------------------------------------------------------------------
Filesystem quota limits for user hadrib
Tier             Filesystem     used    quota    limit    grace    files    quota    limit    grace
-----------------------------------------------------------------------------------------------
capacity         /scratch   119.1G      0k      10T       —      240782       0        0       —
bandwidth        /scratch   924.1G      0k       1T       —      240782       0        0       —
iops             /scratch   60.88M      0k      50G       —      240782       0        0       —
project          /project   7.778T      0k       0k       —      209596       0  1000000       —
-----------------------------------------------------------------------------------------------



lfs quota —uh $USER /scratch
Disk quotas for usr xxx (uid xxx):
     Filesystem     used    quota    limit    grace     files    quota    limit    grace
       /scratch    1.043T      0k      11T       —    240782        0  1024000        —


lfs quota —uh $USER /project
Disk quotas for usr hadrib (uid 129285):
     Filesystem     used    quota    limit    grace     files    quota    limit    grace
       /project   7.778T      0k       0k       —    209596        0  1000000        —
```

10

# Storage Best Practices 3

- Quota for /project, kpq project-id . Quota 80 TB per PI.

```
kpq k10005
----------------------------------
PI quota for : Principal Investigator
----------------------------------

Filesystem   used    quota    limit    grace    files    quota    limit    grace
/project     16k      0k       80T       -        4        0        0        -
```

- Don't be shy, ask for help. help@hpc.kaust.edu.sa

- Limit the number of files per directory

- Clean up after running a job. ( remove temporary, slurm output…)

- Important files shall be copied to your personal workstation.

# Who can use Shaheen III?

- Every Shaheen user must be an official member of at least one project, and every project must originate from an approved organization.

- PI needs to be a faculty/manager to endorse the project.

- Access to Shaheen is available to a select group of academic and industry partners.

- Following the terms and conditions

- Submitting all necessary documents

# Charges of core hours

- Unit of measure is core hours for the CPU nodes
  - Each node has 192 cores.
  - For exclusive usage, full node ( ie 192 cores) will be charged

- Check the usage of the project: sb k1xxxx

```
Title: KSL computational scientists
PI : Saber Feki
-----------------------------------
Project k01    expiry: 2030-12-31
-----------------------------------
Allocations:
2024-02-08                100000000
-----------------------------------
Total allocations:        100000000
Core hours used:            1300532
Remaining balance:         98699468
-----------------------------------
```

# Queue

- FIFO policy with backfilling as long as
  - core hours available
  - Maximum of number of jobs/nodes not reached

- 5 QoS
  - workq  ( entire system up to 4608 nodes)
    - Limit per user 2048 nodes max, 500 job running, 1000 max job queued
  - shared ( 16 nodes in total)
  - debug ( 4 nodes in total)
    - 1 node per user
  - 72hours: up 128 nodes partition shared by all users
    - 16 nodes maximum per user, 3 job running, maximum 32 job queued
  - ppn nodes: post processing nodes: 15 nodes

# Storage quotas for each tiers,

- /scratch
  - Total aggregated PI 80 TB default
  - Per user 1M files, 10 TB in /scratch/username and 1 TB in bandwidth and 50 GB in IOPS

- /project
  - Total aggregated PI 80 TB default
  - Per user 1M files

- Check:
  - Check the quota with kpq, kuq  on project and overall
  - `lfs quota -uh $USER /scratch`
  - `Disk quotas for user username (uid 123456):`

| Filesystem | used | quota | limit | grace | files | quota | limit | grace |
|---|---|---|---|---|---|---|---|---|
| /scratch | 402.7G | 0k | 11T | - | 123339 | 0 | 1024000 | - |

16

# Applications available

- **Over 500 version of libraries, tool and applications available**

- **Module avail ; module avail –S nameofsoftware**

```
-------------------------------------------- /opt/cray/pe/perftools/23.09.0/modulefiles ---------------------------------------------
perftools          perftools-lite     perftools-lite-events perftools-lite-gpu     perftools-lite-hbm   perftools-lite-loops perftools-preload

----------------------------------------------------------- /opt/cray/pe/modulefiles -----------------------------------------------------------
PrgEnv-aocc/8.4.0(default)   cpe-cuda/23.05              cray-hdf5-parallel/1.12.2.3        cray-mpixlate/1.0.2(default)   cray-python/3.10.10(default)   gcc/10.3.0
PrgEnv-aocc/8.5.0            cpe-cuda/23.09(default)     cray-hdf5-parallel/1.12.2.7(default) cray-mpixlate/1.0.3          cray-python/3.11.5            gcc/11.2.0
PrgEnv-cray/8.4.0(default)   cpe-cuda/23.12             cray-hdf5-parallel/1.12.2.9        cray-mrnet/5.1.0              cray-stat/4.12.0             gcc/12.2.0(default)
PrgEnv-cray/8.5.0            cray-R/4.2.1.2(default)     cray-libpals/1.2.12(default)       cray-mrnet/5.1.1(default)     cray-stat/4.12.1(default)    gcc-native/12.3(default)
PrgEnv-gnu/8.4.0(default)    cray-R/4.3.1               cray-libsci/23.05.1.4             cray-mrnet/5.1.2             cray-stat/4.12.2            gdb4hpc/4.15.0
PrgEnv-gnu/8.5.0            cray-ccdb/5.0.0            cray-libsci/23.09.1.1(default)     cray-netcdf/4.9.0.3          cray-ucx/1.14.0(default)     gdb4hpc/4.15.1(default)
PrgEnv-intel/8.4.0(default)  cray-ccdb/5.0.1(default)    cray-libsci/23.12.5              cray-netcdf/4.9.0.7(default)   cray-ucx/2.7.0-1            gdb4hpc/4.16.0
PrgEnv-intel/8.5.0          cray-ccdb/5.0.2            cray-libsci_acc/23.12.0(default)   cray-netcdf/4.9.0.9          cray-ucx/default(default)    iobuf/2.0.10(default)
PrgEnv-nvhpc/8.4.0(default)  cray-cti/2.18.0           cray-mpich/8.1.26                cray-netcdf-hdf5parallel/4.9.0.3 craype/2.7.21              papi/7.0.0.2
PrgEnv-nvhpc/8.5.0          cray-cti/2.18.1(default)    cray-mpich/8.1.27(default)        cray-netcdf-hdf5parallel/4.9.0.7(default) craype/2.7.23       papi/7.0.1.1(default)
PrgEnv-nvidia/8.4.0(default) cray-cti/2.18.2           cray-mpich/8.1.28               cray-netcdf-hdf5parallel/4.9.0.9 craype/2.7.30(default)    papi/7.0.1.2
atp/3.15.0                 cray-dsmml/0.2.2(default)   cray-mpich-abi/8.1.26            cray-openshmemx/11.6.0       craype-dl-plugin-ftr/22.06.1.2(default) perftools-base/23.05.0
atp/3.15.1(default)        cray-dyninst/12.2.0        cray-mpich-abi/8.1.27(default)    cray-openshmemx/11.6.1(default) craype-dl-plugin-py3/21.02.1.3  perftools-base/23.09.0(default)
atp/3.15.2                 cray-dyninst/12.3.0(default) cray-mpich-abi/8.1.28           cray-openshmemx/11.7.0       craype-dl-plugin-py3/21.04.1    perftools-base/23.12.0
cce/16.0.0                 cray-dyninst/12.3.1        cray-mpich-abi-pre-intel-5.0/8.1.26  cray-pals/1.2.12(default)    craype-dl-plugin-py3/22.06.1.2  sanitizers4hpc/1.1.0
cce/16.0.1(default)        cray-fftw/3.3.10.4         cray-mpich-abi-pre-intel-5.0/8.1.27(default) cray-parallel-netcdf/1.12.3.3 craype-dl-plugin-py3/22.08.1  sanitizers4hpc/1.1.1(default)
cce/17.0.0                 cray-fftw/3.3.10.5(default)  cray-mpich-abi-pre-intel-5.0/8.1.28  cray-parallel-netcdf/1.12.3.7(default) craype-dl-plugin-py3/22.09.1  sanitizers4hpc/1.1.2
cpe/23.05                  cray-fftw/3.3.10.6         cray-mpich-ucx/8.1.26            cray-parallel-netcdf/1.12.3.9 craype-dl-plugin-py3/22.12.1(default) valgrind4hpc/2.13.0
cpe/23.09(default)         cray-hdf5/1.12.2.3         cray-mpich-ucx/8.1.27(default)    cray-pmi/6.1.11             craypkg-gen/1.3.28           valgrind4hpc/2.13.1(default)
cpe/23.12                  cray-hdf5/1.12.2.7(default)  cray-mpich-ucx/8.1.28           cray-pmi/6.1.12(default)      craypkg-gen/1.3.30(default)    valgrind4hpc/2.13.2
                          cray-hdf5/1.12.2.9         cray-mpixlate/1.0.1.11           cray-pmi/6.1.13             craypkg-gen/1.3.31

----------------------------------------------------------- /opt/cray/modulefiles -----------------------------------------------------------
chapel/1.30.0                                          cudatoolkit/23.3_11.8                                xpmem/2.6.2-2.5_2.27__gd067c3f.shasta(default)
cray-lustre-client-ofed/2.15.0.7_rc2_cray_25_ga33b7d9-2.5_5.24__ga33b7d9745.shasta(default) cudatoolkit/23.3_12.0(default)
cudatoolkit/23.3_11.0                                 libfabric/1.15.2.0(default)

----------------------------------------------------------- /opt/modulefiles -----------------------------------------------------------
amduprof/4.2.850           aocc/4.2.0(default)       intel/19.0.5.281         intel/2023.1.0(default)    intel-classic/2023.1.0(default) intel-oneapi/2023.1.0(default)

------------------------------------------------ /opt/cray/pe/craype-targets/default/modulefiles ------------------------------------------------
craype-accel-amd-gfx90a craype-accel-host      craype-arm-grace       craype-hugepages2M     craype-hugepages16M   craype-hugepages128M   craype-network-none   craype-x86-genoa      craype-x86-rome       craype-x86-trento
craype-accel-amd-gfx908 craype-accel-nvidia70  craype-hugepages1G     craype-hugepages32M    craype-hugepages256M  craype-network-ofi     craype-x86-milan      craype-x86-spr
craype-accel-amd-gfx940 craype-accel-nvidia80  craype-hugepages2G     craype-hugepages8M     craype-hugepages64M   craype-hugepages512M  craype-network-ucx     craype-x86-milan-x    craype-x86-spr-hbm

----------------------------------------------------------- /sw/ex109genoa/modulefiles -----------------------------------------------------------
abinit/9.6.2          berkeleygw/2.1         converge/3.0.27_udf    egsnrc/2020            gpaw/24.1.0           molpro/2012.1p16       packmol/20.14.3        smina/20220122        vaspkit/1.5.1
abinit/9.10.3         berkeleygw/3.1.0       converge/3.1.4         egsnrc/2023            grace/5.1.25          moltemplate/2.20.20    paraview/5.11.2        sod/0.47              vasputil/6.1
adf/2019.301         bio/blast-2.15.0       converge/3.1.4_udf     eigen/3.3.7           gromacs/2023.1        mopac/22.1.0           paraview/5.12.0-egl    sod/0.52              vesta/3.4.3
airss/0.9.4          bio/bwa/0.7.17         converge/3.1.5         elk/6.3.2             gsl/2.6-cce16.0.1     motif/2.3.8           paraview/5.12.0-mesa(default) spack/0.20.0     virtualflow/15.7
alamode/1.3.0        bio/edta/2.2.0         converge/3.1.5_udf     elk/9.2.12            gsl/2.6-gcc12.2.0     mpifileutils/0.11.1   paraview/5.12.1       sumo/2.3.8            visit/3.3.3
alamode/1.4.2        bio/gatk/4.1.6         converge/3.1.6         elpa/2023.05.001      gulp/6.0             mrcc/2020-02-22_mpi    pytorch/2.2.1         sumo/2.3.8_saber      vmd/1.9.3
almabte/1.3.2        bio/java/8u401         converge/3.1.6_udf     espresso/6.4.1        gulp/6.2             mrcc/2020-02-22_omp    perturbo/2.2.0        tbmodels/1.4.3        vmd/1.9.4
amber/14            bio/java/21.0.2        converge/3.1.7         espresso/6.4.1_environ  horovod/0.28.1-tf212-torch1131 mrcc/2023-08-28_mpi phono3py/1.18.2     tdep/20231024         wannier90/1.2
amber/14_mpi        bio/plink/1.9b         converge/3.1.7_udf     espresso/6.8          horovod/0.28.1-torch221(default) mrcc/2023-08-28_omp phonopy/2.21.0      tecplot/2023r2       wannier90/2.1.0
amber/14_mpi_plumed  bio/plink/2.0a         converge/3.1.8         espresso/6.8_elpa      ifermi/0.3.3          multiwfn/3.6          polyrate/17-C_rp      tensorflow/2.12(default) wannier90/3.1.0
amber/18            bio/samtools/1.8       converge/3.1.8_udf     espresso/6.8_libxc     java/19.0.1          multiwfn/3.8dev       polyrate/17-C_vrc     thirdorder/1.1.1      wannierberri/0.15.0
amber/18_mpi        bio/samtools/1.18      converge/3.1.11        espresso/7.2          jdftx/1.7.0          music/4.0            psi4/1.8.0            timemory/3.0(default) wanniertools/2.5.1
amber/23_mpi        bio/tabix/0.2.6        converge/3.1.11_udf    espresso/7.2_elpa      jmol/14.31.44        mysw/0.0.1           psi4/1.9             totalview/2023.3.10   wanniertools/2.7.0
ams/2023.103        blitz/1.0.2            converge/3.2.1         espresso/7.2_environ   koopmans/1.0.1       mytools/1.01          py4vasp/0.7.4         towhee/8.2.3          wham/2.0.11
amset/0.4.18        boltztrap/1.2.5        converge/3.2.1_udf     espresso/7.2_libxc     ktf/0.1.5            namd/2.14_gcc         pyprocar/6.1.7        turbomole/7.1         wien2k/21.1_elpa
amset/0.4.18_parallel boltztrap2/24.1.1    converge/3.2.4         espresso/7.3          kwant/1.4.3          namd/2.14_smp_gcc     pyrx/0.9.9            turbomole/7.1_mpi     wien2k/21.1_scalapack
amset/0.4.20        boost/1.85             converge/4.0.0         excimontec/1.0.0      lammps/6Aug2023       nb06/6.0_g09          python/3.10.13(default) turbomole/7.1_smp   wien2k/23.2_elpa
ansys/v22R1-fluids   castep/21.11          converge/4.0.0_udf     exciting/neon21       lev00/4.01           nb06/6.0_g16          pytorch/1.13.1(default) uspex/10.5         wien2k/23.2_scalapack
ansys/v23R1-fluids   chemshell/21.0.2      converge/4.0.2         exciting/nitrogen14    libxc/4.2.3          nwchem/6.8.1          pytorch/2.2.1         vampire/6.0           wrf/4.5.2_cray
ansys/v23R1-structures chemshell/23.0.1    cp2k/2023.2            fermisurfer/2.2.1      libxc/4.3.4          nwchem/7.2.2          qchem/5.4_mpi         vasp/5.4.4            wrf/4.5.2_intel
ansys/v23R2-electronics chimera/1.16       critic2/1.1dev         fermisurfer/2.4.0      libxc/5.2.3          octopus/11.3          qchem/5.4_smp         vasp/5.4.4.pl2         xcrysden/1.5.60
ansys/v23R2-fluids   cif2cell/1.2.10       critic2/1.1stable      ffmpeg/6.1.1(default)  libxc/5.1.7          octopus/13.0          qchem/6.1             vasp/5.4.4.pl2_dftd4   xcrysden/1.6.2
aocl/4.2.0(default)  cmake/3.28.3          crystal14/1.0.3        fhiaims/210716_3       libxc/5.2.3          oommf/2.0alpha3       qchem/6.2             vasp/5.4.4.pl2_libbeef  xtb/6.4.1
arm-forge/23.1.2     code/2024.1.4         cuby4/4               fhiaims/221103        libxc/6.2.2          oommf/2.1alpha0       quantumkite/1.0       vasp/5.4.4.pl2_nbo     xtb/6.5.0
ase/3.19.0          columbus/7.2          deepspeed/0.14.0       fhiaims/231212_1       likwid/5.3.0         openbabel/3.1.1       quantumkite/1.1       vasp/5.4.4.pl2_occmat   xtb/6.6.7.0
ase/3.22.1          columbus/7.2.2         dftbplus/21.2          flex/2.6.4            lobster/5.0.0        openfoam/2212         qvasp/2.23            vasp/5.4.4.pl2_optaxis  xthi/0.1
atk/2019.03sp1      comsol/6.2            dftd4/2.5.0            fourphonon/1.1         materstudio/2023      openmolcas/23.10      raspa2/2.0.3          vasp/5.4.4.pl2_transopt yambo/5.0.4
atompaw/4.2.0.3      comsol/6.2_tmp        fourphonon/20211001    gamess/30Sept2022R2    milo/1.0.3           openmx/3.9.9          seissol/1.1.4_gnu     vasp/5.4.4.pl2_vaspsol  yambo/5.0.4_slepc
autodockvina/1.2.3   converge/3.0.13       dlpoly/4.09            gamess/30Sept2023R2    mkl/2024.0.0         orca/5.0.4           shengbte/1.5.0        vasp/5.4.4.pl2_vtst178  yambo/5.2.1
bader/1.04          converge/3.0.13_udf    dlpoly/4.09_plumed     gaussian09/d.01       mohid/19.10          orca/6.0.0           siesta/4.1.5          vasp/5.4.4.pl2_z2pack   yambo/5.2.1_slepc
bader/1.05          converge/3.0.25        dlpoly/5.1.0          gaussian16/c.02       mohid/23.10          ovito/2.9.0           siesta/psml204        vasp/6.4.2            z2pack/2.2.0
bands4vasp/0.4      converge/3.0.25_udf    dlpoly/5.1.0_plumed    dssp/2.3.0           molden/6.9           ovito/3.10.1          simmate/0.13.2        vasp/6.4.2_optaxis     zendnn/4.1
bazel/6.1.0          converge/3.0.26        eddp/0.2              go/1.21.3            molden/7.2.1          p4vasp/0.3.30         simmate/0.16.1        vasp/6.4.2_scpc
bazel/6.5.0(default) converge/3.0.26_udf    edmftf/Apr2021        gollum2/2.0           moleculargsm/20240115 paccham/20200322      simmate/0.17.0        vasp/6.4.2_vtst198
bazel/7.1.0          converge/3.0.27        edmftf/Jan2019        gpaw/22.1.0           molgw/3.2            packmol/18.169        singularity/4.0.1     vaspkit/1.4.1
```

17

# Acknowledging

**Terms and Conditions regarding Research Publications**

Whenever the results of research conducted on the HPC systems at KAUST are published, or the research involved personnel from KAUST Supercomputing Laboratory (KSL), Principal Investigators (PIs) are required to acknowledge the usage of the HPC systems at KAUST and/or the involvement of KSL personnel in their research in their publications. For example, the following statement could be used: **"For computer time, this research used Shaheen III managed by the Supercomputing Core Laboratory at King Abdullah University of Science & Technology (KAUST) in Thuwal, Saudi Arabia**.

# Best Practices / Getting help

- When submitting a ticket to help@hpc.kaust.edu.sa requesting help, you will likely get faster resolution by supporting a few best practices:
  - Where possible, provide helpful details that can help speed the process. For example: Project ID, relevant directories, job scripts, jobIDs, modules at compile/runtime, login name, etc.

- One issue per ticket. Do not add unrelated questions on existing tickets

- Do not open multiple tickets on the same unresolved topic.

- Let us know if your issue is fixed or you solved it(and let us know what worked!)

- MOST IMPORTANT: do not hesitate to contact us  help@hpc.kaust.edu.sa;

# Best Practices

- Check the documentations: https://docs.hpc.kaust.edu.sa/policy/shaheen3.html

- Read the announcements/newsletter sent by emails and available in the website:
  - hpc.kaust.edu.sa/

**Thanks !**

# Agenda


**Shaheen III Survey**

22