
Introduction to SDSC Systems, Launching and Managing Jobs

Mahidhar Tatineni (mahidhar@sdsc.edu)

SDSC Summer Institute

August 10, 2015

Getting Started

- **System Access – Logging in**
 - Linux/Mac – Use available ssh clients.
 - ssh clients for windows – Putty, Cygwin
 - <http://www.chiark.greenend.org.uk/~sgtatham/putty/>
 - Login hosts for the SDSC machines:
 - comet.sdsc.edu
 - gordon.sdsc.edu

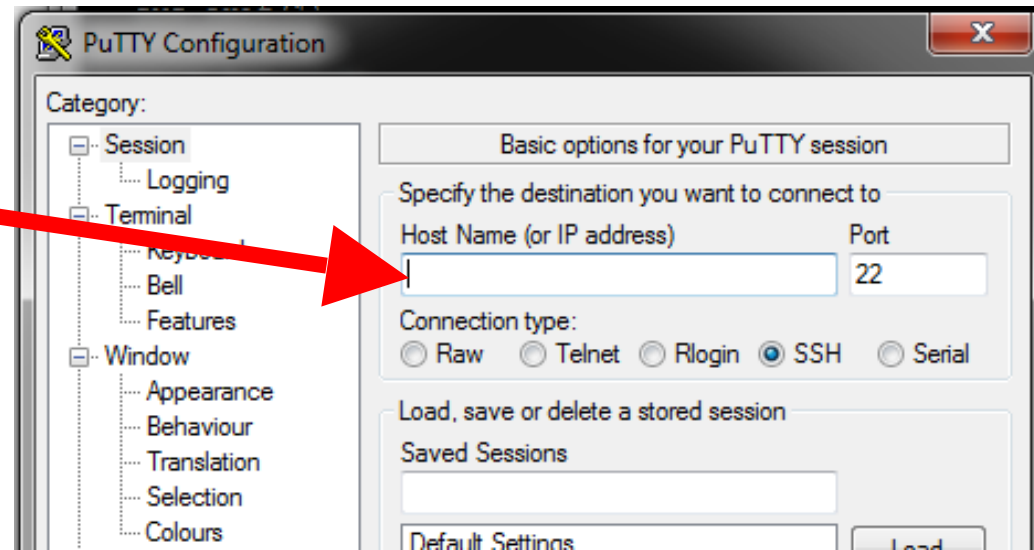
Logging into Comet

Mac/Linux:

```
ssh etrainXY@comet-ln1.sdsc.edu
```

Windows (PuTTY):

comet.sdsc.edu



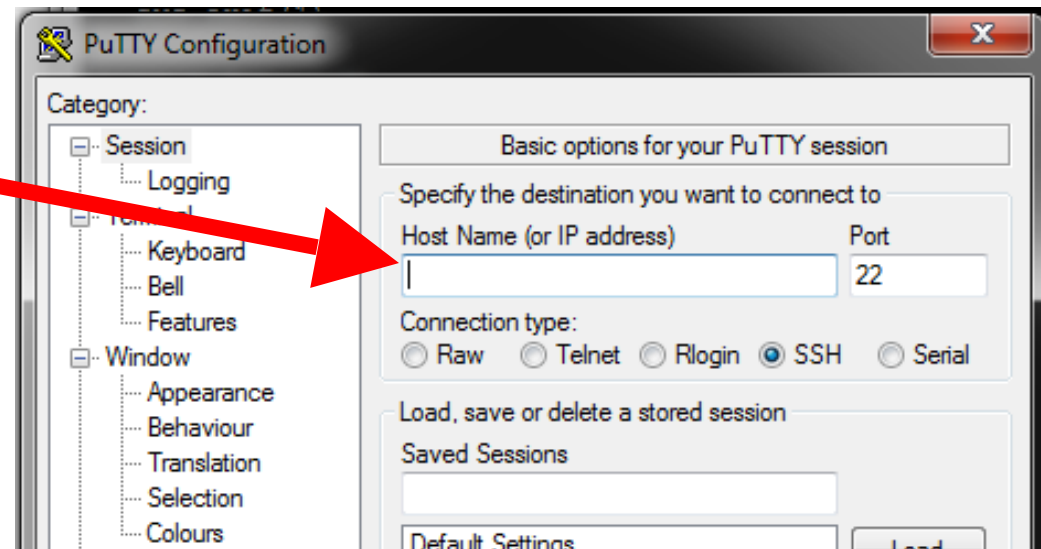
Logging into Gordon

Mac/Linux:

```
ssh etrainXY@gordon.sdsc.edu
```

Windows (PuTTY):

gordon.sdsc.edu



Training Accounts

- **Make sure we have the training intro exercises. Login to Comet and copy the files to your home directory:**

```
cp -r /share/apps/examples/SI2015/INTRO $HOME/
```

```
ls $HOME/SI2015/INTRO
```

```
COMET GORDON
```

Access Via Science Gateways (XSEDE)

- Community-developed set of tools, applications, and data that are integrated via a portal.
- Enables researchers of particular communities to use HPC resources through portals without the complication of getting familiar with the hardware and software details. Allows them to focus on the scientific goals.
- CIPRES gateway hosted by SDSC PIs enables large scale phylogenetic reconstructions using applications such as MrBayes, Raxml, and Garli. Enabled ~320 publications in 2013 and accounts for a significant fraction of the XSEDE users.
- NSG portal hosted by SDSC PIs enables HPC jobs for neuroscientists.

Data Transfer (scp, globus-url-copy)

- **scp** is o.k. to use for simple file transfers and small file sizes (<1GB). Example:

```
$ scp w.txt train40@gordon.sdsc.edu:/home/train40/w.txt  
100% 15KB 14.6KB/s 00:00
```
- **globus-url-copy** for large scale data transfers between XD resources (and local machines w/ a globus client).
 - Uses your XSEDE-wide username and password
 - Retrieves your certificate proxies from the central server
 - Highest performance between XSEDE sites, uses striping across multiple servers and multiple threads on each server.

7

Data Transfer – globus-url-copy

- **Step 1: Retrieve certificate proxies:**

```
$ module load globus
```

```
$ myproxy-logon -l xsedeusername
```

```
Enter MyProxy pass phrase:
```

```
A credential has been received for user xsedeusername in /tmp/  
x509up_u555555.
```

- **Step 2: Initiate globus-url-copy:**

```
$ globus-url-copy -vb -stripe -tcp-bs 16m -p 4 gsiftp://  
gridftp.ranger.tacc.teragrid.org:2811///scratch/00342/username/test.tar gsiftp://  
trestles-dm2.sdsc.xsede.org:2811///oasis/scratch/username/temp_project/test-  
gordon.tar
```

```
Source: gsiftp://gridftp.ranger.tacc.teragrid.org:2811///scratch/00342/username/
```

```
Dest: gsiftp://trestles-dm2.sdsc.xsede.org:2811///oasis/scratch/username/  
temp_project/
```

8

```
test.tar -> test-gordon.tar
```


Data Transfer – Globus

- **Works from Windows/Linux/Mac via globus online website:**
 - <https://www.globus.org>
- **Gordon and Trestles endpoints already exist. Authentication can be done using XSEDE-wide username and password for the NSF resources.**
- **Globus Connect application (available for Windows/Linux/Mac can turn your laptop/desktop into an endpoint.**

9

Data Transfer – Globus Online

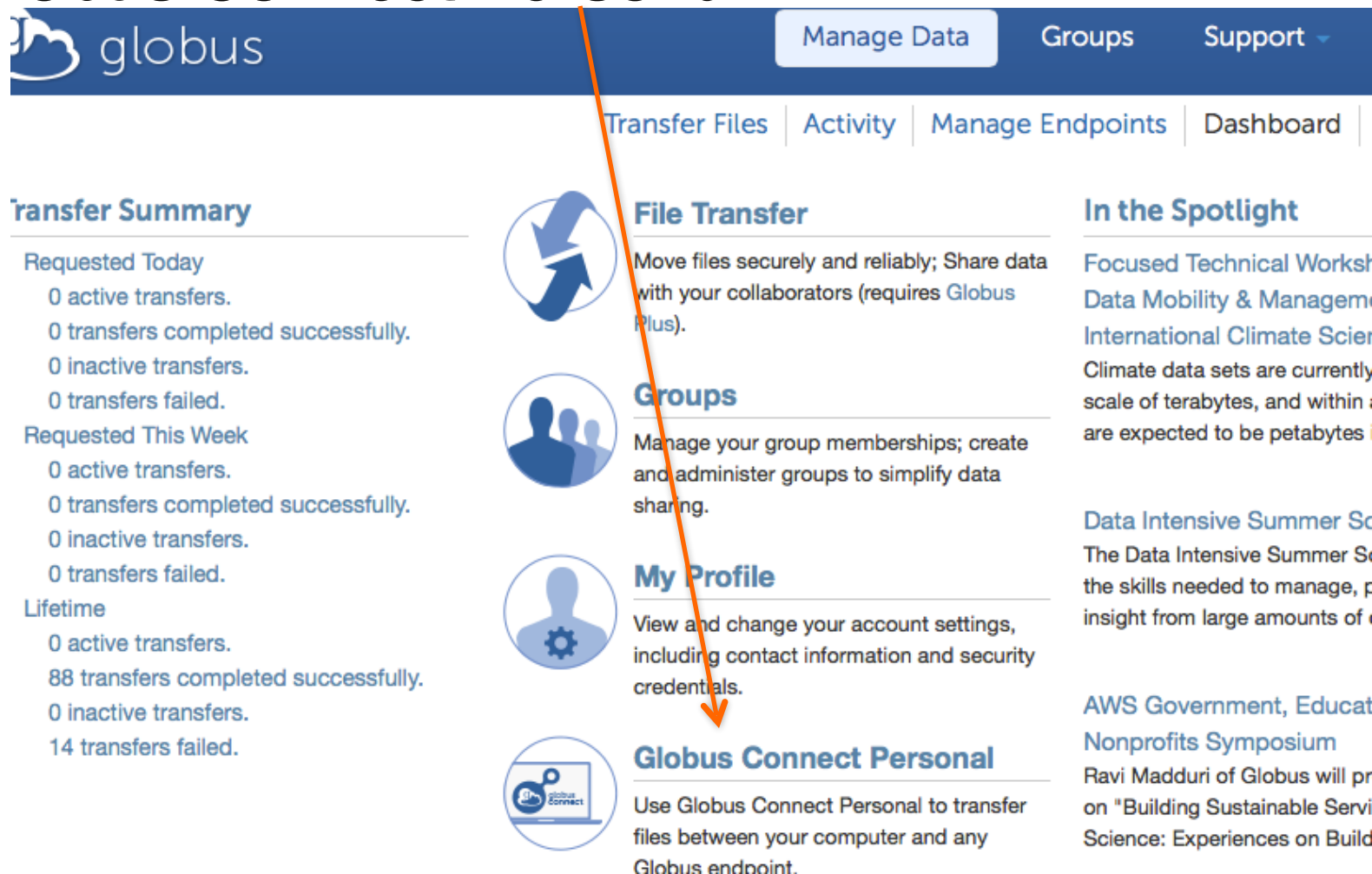
- **Step 1: Create a globus account (at www.globus.org)**

The screenshot shows the Globus Online website in a web browser. The browser's address bar displays "https://www.globus.org". The website has a dark blue header with the Globus logo (a white 'g' inside a cloud) and the word "globus" in white. Navigation links for "Products", "News", "About", and "Support" are visible, along with "Log In" and "Sign Up" buttons. The main content area features a large graphic with orange arrows forming a cycle around a central "BIG DATA" circle, with labels "share", "move", and "sync". Below this graphic, a large number "57,636,318,706 MB" is displayed, with "TRANSFERRED" in smaller text underneath. To the right, the text "Your research data where you need it." is shown. At the bottom, there are three links: "Researchers", "Resource Providers", and "How It Works".

10

Data Transfer – Globus

- Step 2: Set up local machine as endpoint using Globus Connect Personal.



The screenshot shows the Globus Connect Personal web interface. At the top is a dark blue header with the Globus logo on the left and navigation links 'Manage Data', 'Groups', and 'Support' on the right. Below the header is a secondary navigation bar with links 'Transfer Files', 'Activity', 'Manage Endpoints', and 'Dashboard'. The main content area is divided into three columns. The left column, titled 'Transfer Summary', contains statistics for 'Requested Today', 'Requested This Week', and 'Lifetime'. The middle column contains four sections: 'File Transfer' (with a double-headed arrow icon), 'Groups' (with a group of people icon), 'My Profile' (with a person and gear icon), and 'Globus Connect Personal' (with a laptop and plug icon). An orange arrow points from the 'Step 2' text in the list above to the 'Globus Connect Personal' section. The right column, titled 'In the Spotlight', contains three articles: 'Focused Technical Workshop on Data Mobility & Management', 'Data Intensive Summer School', and 'AWS Government, Education & Nonprofits Symposium'.

Transfer Summary

Requested Today

- 0 active transfers.
- 0 transfers completed successfully.
- 0 inactive transfers.
- 0 transfers failed.

Requested This Week

- 0 active transfers.
- 0 transfers completed successfully.
- 0 inactive transfers.
- 0 transfers failed.

Lifetime

- 0 active transfers.
- 88 transfers completed successfully.
- 0 inactive transfers.
- 14 transfers failed.

File Transfer

Move files securely and reliably; Share data with your collaborators (requires Globus Plus).

Groups

Manage your group memberships; create and administer groups to simplify data sharing.

My Profile

View and change your account settings, including contact information and security credentials.

Globus Connect Personal

Use Globus Connect Personal to transfer files between your computer and any Globus endpoint.

In the Spotlight

Focused Technical Workshop on Data Mobility & Management

International Climate Science

Climate data sets are currently on the scale of terabytes, and within the next few years are expected to be petabytes.


Data Intensive Summer School

The Data Intensive Summer School provides the skills needed to manage, process, and gain insight from large amounts of data.

AWS Government, Education & Nonprofits Symposium

Ravi Madduri of Globus will present on "Building Sustainable Service Science: Experiences on Building Sustainable Service Science".

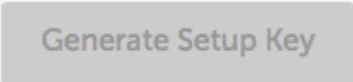
Data Transfer – Globus

 Create Globus Connect Personal Endpoint

Step 1 Get Your Globus Connect Personal Setup Key

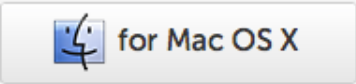
Please enter a unique name for your Globus Connect Personal endpoint to help you identify it.

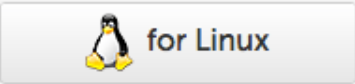
Endpoint Name: mahidhar#




Step 2 Download & Install Globus Connect Personal

Click one of the buttons below to download and install Globus Connect Personal for your operating system.

 for Mac OS X

 for Linux

 for Windows

Once downloaded, run the installer. When prompted, paste in the Setup Key to complete the installation.

Data Transfer – Globus

- Step 3: Pick Endpoints and Initiate Transfers!

Transfer Files | Activity | Manage Endpoints | Dashboard | F

Transfer Files Get Globus Connect
Turn your computer into

Endpoint xsede#gordon Go

Path /~/GNU/ Go

select all | none up one folder refresh list

gcc-4.6.1	Folder
gmp	Folder
gmp-4.3.2	Folder
mpc	Folder
mpc-0.9	Folder
mpfr	Folder
mpfr-3.0.1	Folder
gcc-4.6.1.tar.gz	89.1 MB
gcc-core-4.6.1.tar.gz	36.38 MB
gcc-fortran-4.6.1.tar	12.94 MB
gcc-g++-4.6.1.tar.gz	8.74 MB
gmp.tar	18.21 MB
mpc-0.9.tar.gz	552.69 kB
mpfr-3.0.1.tar.gz	1.41 MB

Endpoint mahidhar#laptop


Path /~/lambda/

select all | none up one folder refresh list

No files available to display

Data Transfer – Globus



mahidhar#laptop to xsede#gordon 

transfer completed 2 months ago



overview



event log

Task ID a79e3904-f26d-11e3-b567-12313940394d

Source mahidhar#laptop ?

Destination xsede#gordon

Status SUCCEEDED

User mahidhar

Requested 2014-06-12 01:10 pm

Deadline 2014-06-13 01:10 pm

Completed 2014-06-12 01:10 pm

Transfer Settings

- overwriting all files on destination
- verify file integrity after transfer
- transfer is not encrypted

Files	1
Directories	0
Bytes Transferred	297,821
Pending	0
Succeeded	1
Cancelled	0
Expired	0
Failed	0
Retrying	0
Skipped	0

[view debug data](#)

HPC Resources at SDSC

NSF : Comet, Gordon
Parallel Filesystems: Data Oasis
UCSD IDI Resource: TSCC



UC San Diego

SDSC
SAN DIEGO SUPERCOMPUTER CENTER



INDIANA UNIVERSITY



SDSC SAN DIEGO SUPERCOMPUTER CENTER

This work supported by the National Science Foundation, award ACI-1341698.

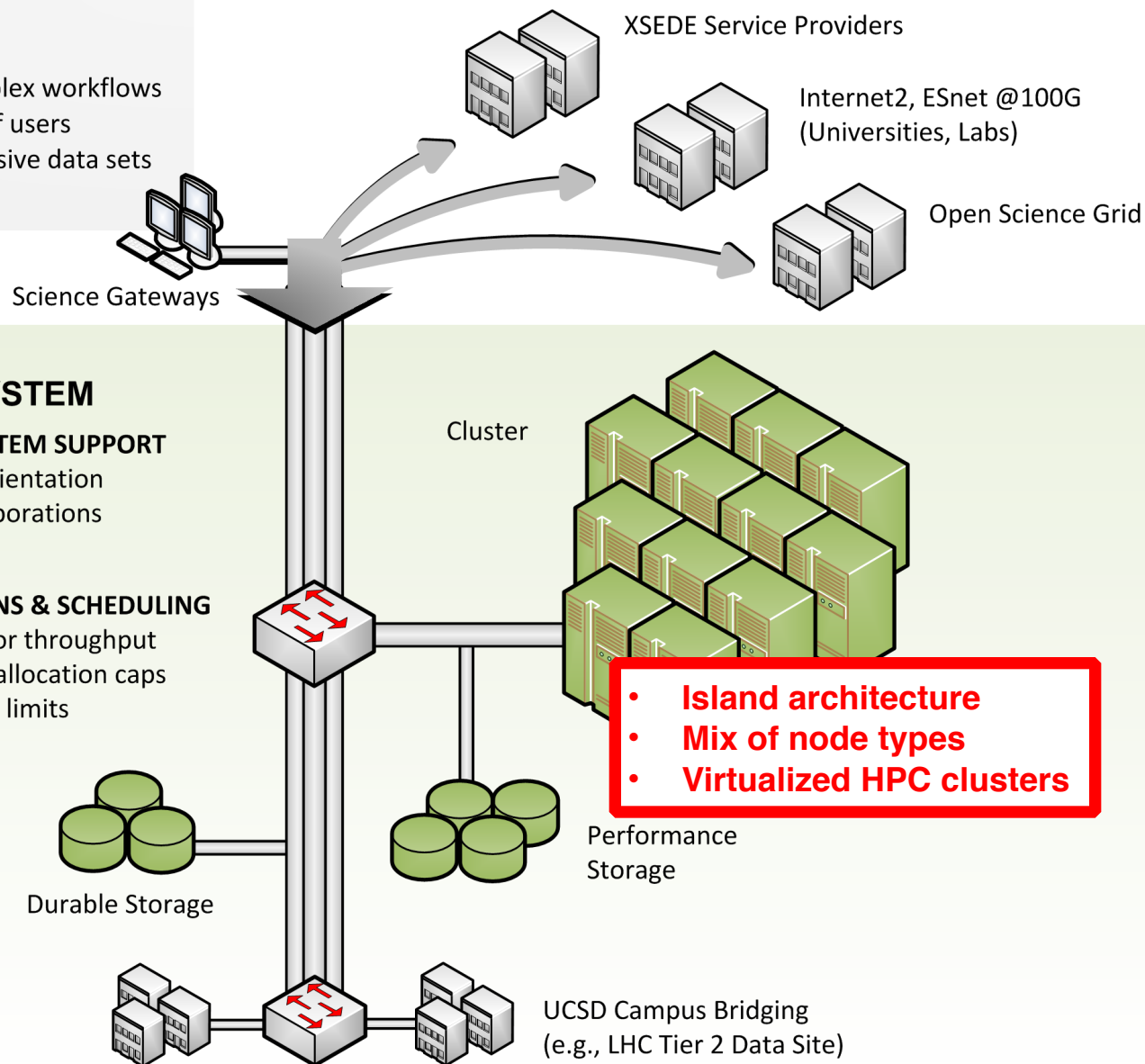
at the UNIVERSITY OF CALIFORNIA; SAN DIEGO



Comet Will Serve the 99%

CHALLENGES OUR PROPOSAL ADDRESSES

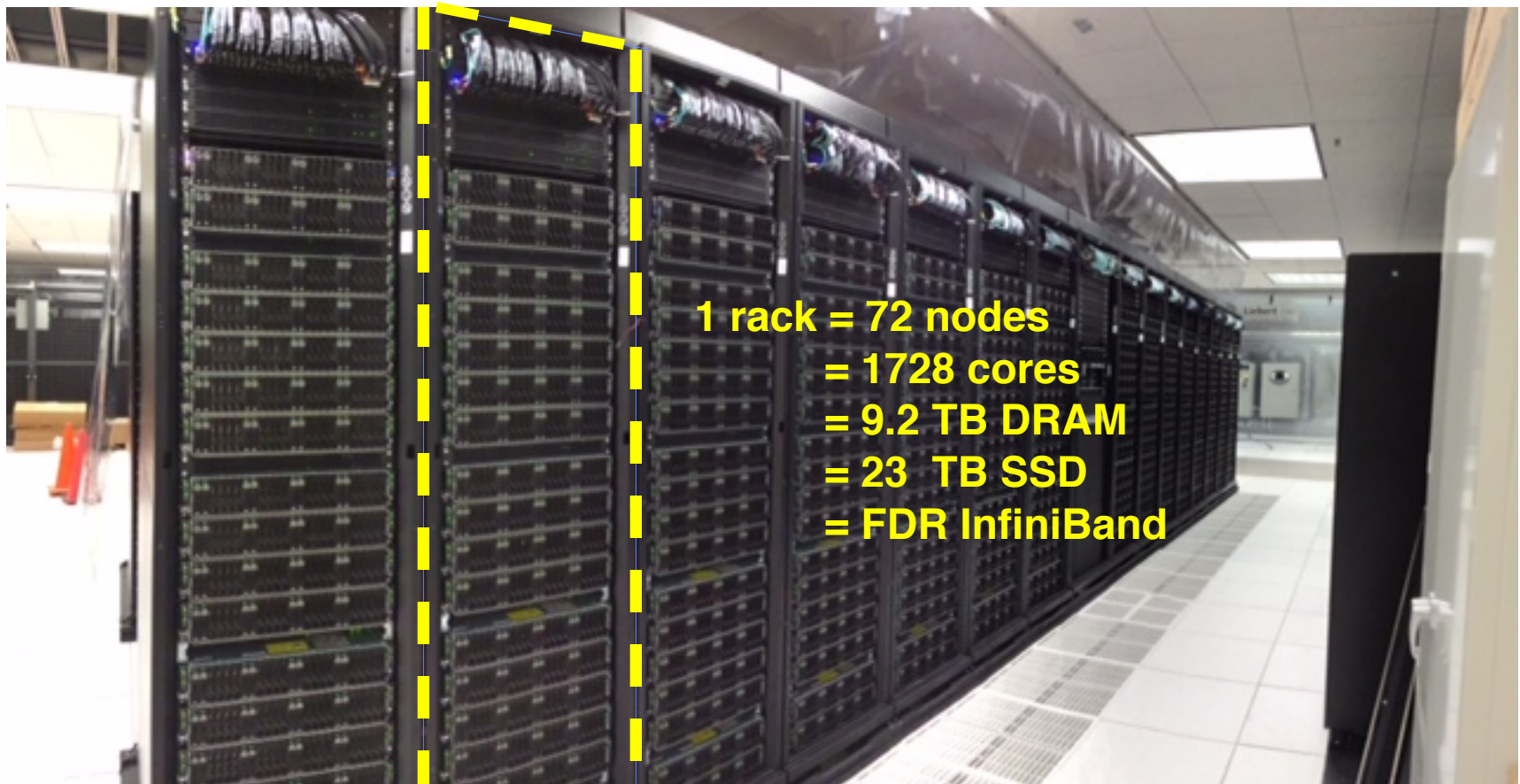
- ✓ Attract new users and communities
- ✓ Support diverse applications with complex workflows
- ✓ Ensure responsiveness for thousands of users
- ✓ Transfer, store, analyze, and share massive data sets
- ✓ Integrate with XSEDE



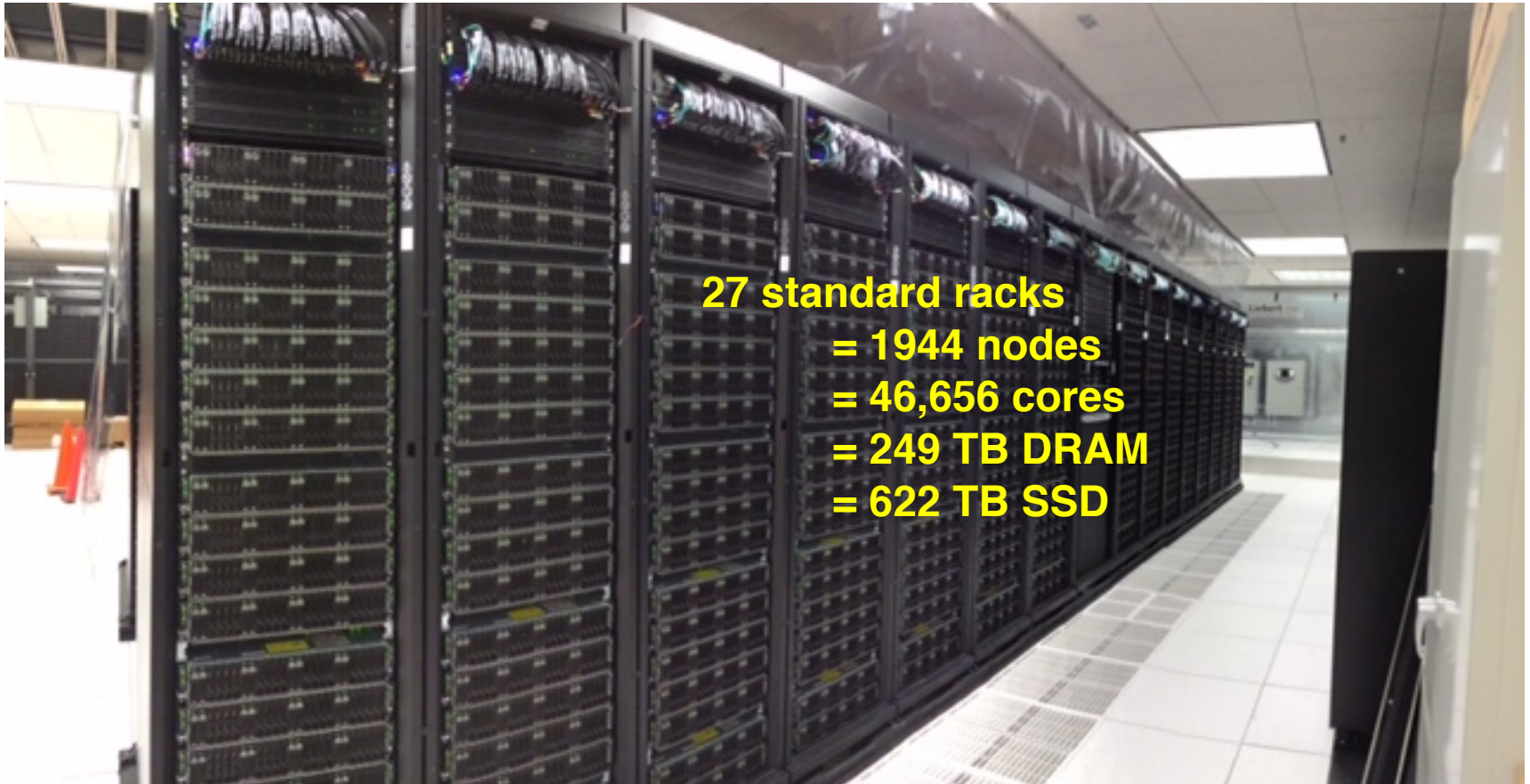
Comet: System Characteristics

- **Total peak flops ~2.1 PF**
- **Dell primary integrator**
 - Intel Haswell processors w/ AVX2
 - Mellanox FDR InfiniBand
- **1,944 standard compute nodes (46,656 cores)**
 - Dual CPUs, each 12-core, 2.5 GHz
 - 128 GB DDR4 2133 MHz DRAM
 - 2*160GB GB SSDs (local disk)
- **36 GPU nodes**
 - Same as standard nodes *plus*
 - Two NVIDIA K80 cards, each with dual Kepler3 GPUs
- **4 large-memory nodes (June 2015)**
 - 1.5 TB DDR4 1866 MHz DRAM
 - Four Haswell processors/node
- **Hybrid fat-tree topology**
 - FDR (56 Gbps) InfiniBand
 - Rack-level (72 nodes, 1,728 cores) full bisection bandwidth
 - 4:1 oversubscription cross-rack
- **Performance Storage (Aeon)**
 - 7.6 PB, 200 GB/s; Lustre
 - Scratch & Persistent Storage segments
- **Durable Storage (Aeon)**
 - 6 PB, 100 GB/s; Lustre
 - Automatic backups of critical data
- **Home directory storage**
- **Gateway hosting nodes**
- **Virtual image repository**
- **100 Gbps external connectivity to Internet2 & ESNet**

~67 TF supercomputer in a rack



And 27 single-rack supercomputers



Gordon – A Data Intensive Supercomputer

- Designed to accelerate access to massive amounts of data in areas of genomics, earth science, engineering, medicine, and others
- Emphasizes memory and IO over FLOPS.
- Appro integrated 1,024 node Sandy Bridge cluster
- 300 TB of high performance Intel flash
- Large memory supernodes via vSMP Foundation from ScaleMP
- 3D torus interconnect from Mellanox
- In production operation since February 2012
- Funded by the NSF and available through the NSF Extreme Science and Engineering Discovery Environment program (XSEDE)

SDSC



ScaleMP™

XSEDE
Extreme Science and Engineering
Discovery Environment

SDSC

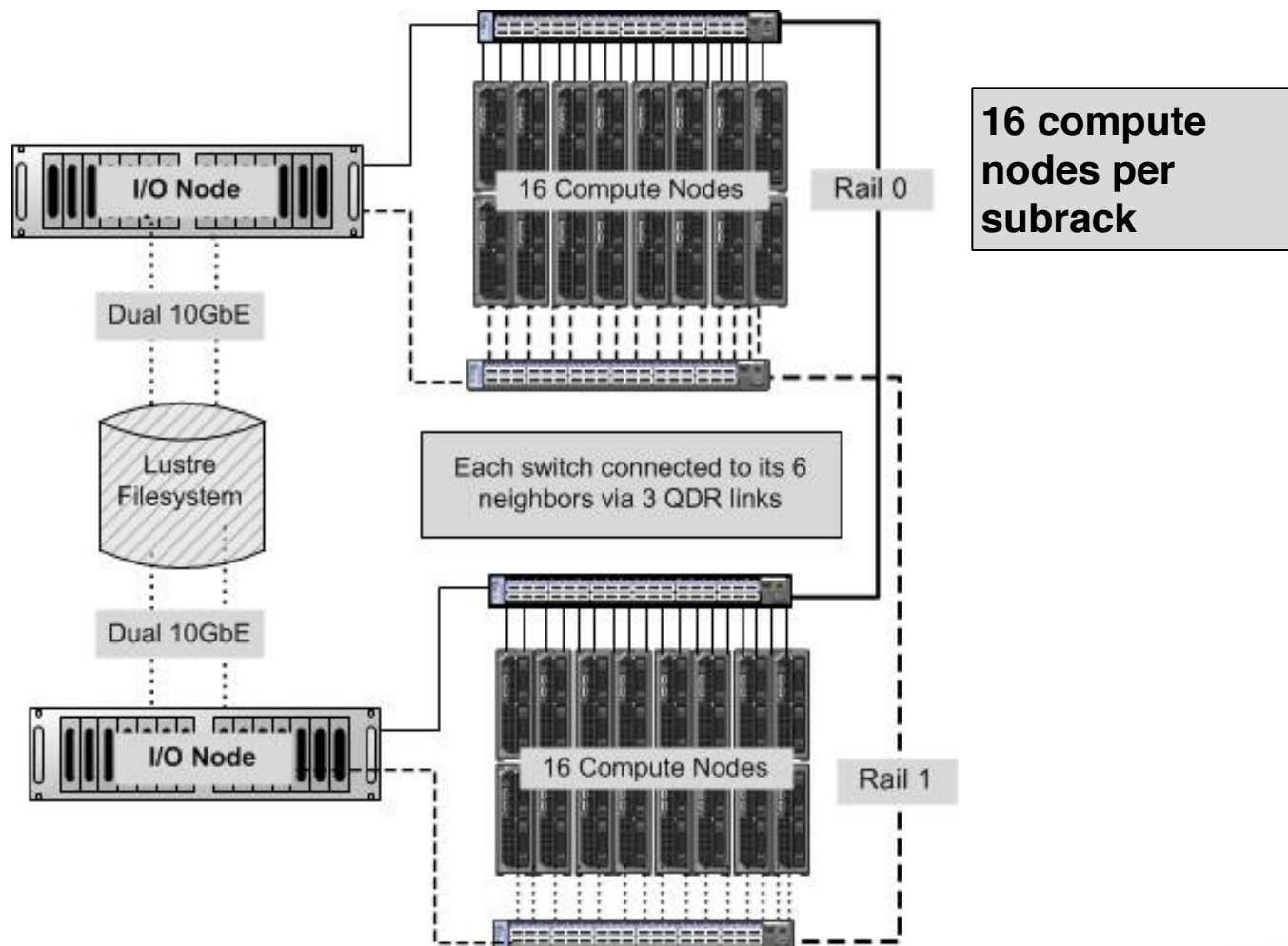
SAN DIEGO SUPERCOMPUTER CENTER *at the* UNIVERSITY OF CALIFORNIA, SAN DIEGO



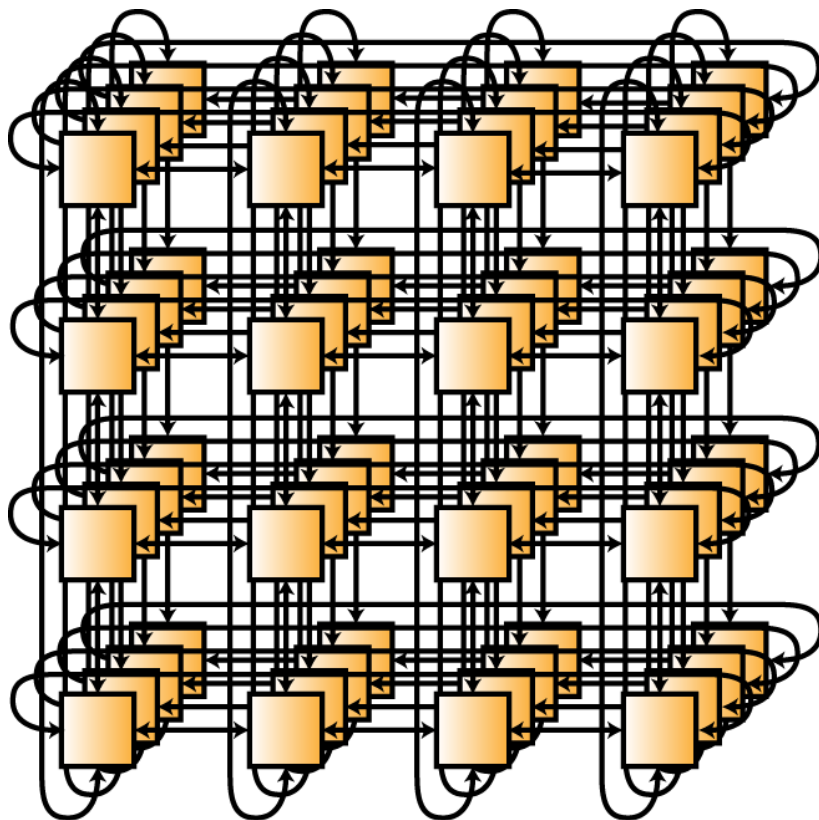
Gordon System Specification

INTEL SANDY BRIDGE COMPUTE NODE	
Sockets	2
Cores	16
Clock speed	2.6
DIMM slots per socket	4
DRAM capacity	64 GB
INTEL FLASH I/O NODE	
NAND flash SSD drives	16
SSD capacity per drive/Capacity per node/total	300 GB / 4.8 TB / 300 TB
Flash bandwidth per drive (read/write)	270 MB/s / 210 MB/s
Flash bandwidth per node (write/read)	4.3 /3.3 GB/s
SMP SUPER-NODE	
Compute nodes	32
I/O nodes	2
Addressable DRAM	2 TB
Addressable memory including flash	12TB
GORDON	
Compute Nodes	1,024
Total compute cores	16,384
Peak performance	341TF
Aggregate memory	64 TB
INFINIBAND INTERCONNECT	
Aggregate torus BW	9.2 TB/s
Type	Dual-Rail QDR InfiniBand
Link Bandwidth	8 GB/s (bidirectional)
Latency (min-max)	1.25 μ s – 2.5 μ s
DISK I/O SUBSYSTEM	
Total storage	/oasis/scratch (1.6 PB), /oasis/projects/nsf(1.5PB)
I/O bandwidth	100 GB/s
File system	Lustre

Subrack Level Architecture



3D Torus of Switches



- Linearly expandable
- Simple wiring pattern
- Short Cables- Fiber Optic cables generally not required
- Lower Cost :40% as many switches, 25% to 50% fewer cables
- Works well for localized communication
- Fault Tolerant within the mesh with 2QoS Alternate Routing
- Fault Tolerant with Dual-Rails for all routing algorithms

3rd dimension wrap-around not shown for clarity

Software – Applications, Libraries



- Users can manage environment via modules.
- Applications packaged into “Rocks Rolls” that can built and deployed on any of the SDSC systems. Benefits wider community deploying software on their Rocks clusters.
- Efficient system administration pooling software install/testing efforts from different projects/machines – Comet benefits from work done for Trestles, Gordon, and Triton Shared Computing Cluster (TSCC).
- Users benefit from a familiar applications environment across SDSC systems => can easily transition from Trestles to Comet.
- Rolls available for all installed applications on Trestles, Gordon. Updated recently for newer compiler, application versions. Listed on next slide.

List of Compilers, Applications, Libraries

Category	List of Software/Libraries
Compilers	<i>Intel, PGI, GNU, MVAPICH2, OPENMPI</i>
Bioinformatics	<i>BamTools, BEAGLE, BEAST, BEAST 2, bedtools, Bismark, BLAST, BLAT, Bowtie, Bowtie 2, BWA, Cufflinks, DPPDiv, Edena, FastQC, FastTree, FASTX-Toolkit, FSA, GARLI, GATK, GMAP-GSNAP, IDBA-UD, MAFFT, MrBayes, PhyloBayes, Picard, PLINK, QIIME, RAxML, SAMtools, SOAPdenovo2, SOAPsnp, SPAdes, TopHat, Trimmomatic, Trinity, Velvet</i>
Chemistry	<i>CPMD, CP2K, GAMESS, Gaussian, MOPAC, NWChem, Q-Chem, VASP</i>
Molecular Dynamics	<i>AMBER, Gromacs, LAMMPS, NAMD</i>
Engineering	<i>ABAQUS</i>
Data Analysis/ Analytics	<i>Hadoop 1, Hadoop 2 (with YARN), Spark, R, Weka, KNIME</i>
Visualization	<i>VisIt, IDL</i>
Numerical libraries	<i>ATLAS, FFTW, GSL, LAPACK, MKL, ParMETIS, PETSc, ScaLAPACK, SPRNG, Sundials, SuperLU, Trilinos</i>
Debugging/Profiling	<i>DDT, PAPI, TAU, mpiP, Valgrind</i>

SDSC HPC Resources: Running Jobs

Running Batch Jobs

- Comet uses SLURM for scheduling. The bulk of our examples today will be on Comet.
- Gordon and TSCC use the TORQUE/PBS resource manager for running jobs. On Gordon we have the Catalina scheduler to control the workload.
- Both schedulers allows the user to submit one or more jobs for execution, using parameters specified in a job script.
- Earlier we copied the examples into our home directory:
/home/\$USER/SI2015/INTRO

Comet – Compiling/Running Jobs

- **Copy and change to COMET workshop directory:**
`cd /home/$USER/SI2015/INTRO/COMET`
- **Verify modules loaded:**
`module list`
Currently Loaded Modulefiles:
1) gnutools/2.69 2) intel/2013_sp1.2.144 3) mvapich2_ib/2.1
- **Compile the MPI hello world code:**
`cd /home/$USER/SI2015/INTRO/COMET/MPI`
`mpif90 -o hello_mpi hello_mpi.f90`
- **Verify executable has been created:**
`ls -lt hello_mpi`
`-rwxr-xr-x 1 mahidhar sdsc 721912 Mar 25 14:53 hello_mpi`

Comet: Hello World on compute nodes

The submit script is hellompi-slurm.sb:

```
#!/bin/bash
#SBATCH --job-name="hellompi"
#SBATCH --output="hellompi.%j.%N.out"
#SBATCH --partition=compute
#SBATCH --nodes=2
#SBATCH --ntasks-per-node=24
#SBATCH --export=ALL
#SBATCH -t 01:30:00
```

#This job runs with 2 nodes, 24 cores per node for a total of 48 cores.
#ibrun in verbose mode will give binding detail

```
ibrun -v ./hello_mpi
```

Comet: Hello World on compute nodes

IBRUN: Command is ../hello_mpi

IBRUN: Command is /share/apps/examples/MPI/hello_mpi

...

...

IBRUN: MPI binding policy: **compact/core for 1 threads per rank (12 cores per socket)**

IBRUN: Adding MV2_CPU_BINDING_LEVEL=core to the environment

IBRUN: Adding MV2_ENABLE_AFFINITY=1 to the environment

IBRUN: Adding MV2_DEFAULT_TIME_OUT=23 to the environment

IBRUN: Adding **MV2_CPU_BINDING_POLICY=bunch** to the environment

...

...

IBRUN: Added 8 new environment variables to the execution environment

IBRUN: Command string is [**mpirun_rsh -np 48 -hostfile /tmp/rssSvaauJA -export /share/apps/examples/MPI/hello_mpi**]

node 18 : Hello world

node 13 : Hello world

node 2 : Hello world

node 10 : Hello world

Compiling OpenMP Example

- **Change to the examples directory:**

```
cd /home/$USER/SI2015/INTRO/COMET/OPENMP
```

- **Compile using `-openmp` flag:**

```
ifort -o hello_openmp -openmp hello_openmp.f90
```

- **Verify executable was created:**

```
[mahidhar@comet-08-11 OPENMP]$ ls -lt hello_openmp  
-rwxr-xr-x 1 mahidhar sdsc 750648 Mar 25 15:00 hello_openmp
```


OpenMP job script

```
#!/bin/bash
#SBATCH --job-name="hell_openmp"
#SBATCH --output="hello_openmp.%j.%N.out"
#SBATCH --partition=compute
#SBATCH --nodes=1
#SBATCH --ntasks-per-node=24
#SBATCH --export=ALL
#SBATCH -t 01:30:00

#SET the number of openmp threads
export OMP_NUM_THREADS=24

#Run the job using mpirun_rsh
./hello_openmp
```

Output from OpenMP Job

\$ more hello_openmp.out

```
HELLO FROM THREAD NUMBER = 7
HELLO FROM THREAD NUMBER = 6
HELLO FROM THREAD NUMBER = 9
HELLO FROM THREAD NUMBER = 8
HELLO FROM THREAD NUMBER = 5
HELLO FROM THREAD NUMBER = 4
HELLO FROM THREAD NUMBER = 0
HELLO FROM THREAD NUMBER = 12
HELLO FROM THREAD NUMBER = 14
HELLO FROM THREAD NUMBER = 3
HELLO FROM THREAD NUMBER = 13
HELLO FROM THREAD NUMBER = 10
HELLO FROM THREAD NUMBER = 11
HELLO FROM THREAD NUMBER = 2
HELLO FROM THREAD NUMBER = 1
HELLO FROM THREAD NUMBER = 15
```

Running Hybrid (MPI + OpenMP) Jobs

- Several HPC codes use a hybrid MPI, OpenMP approach.
- **“ibrun”** wrapper developed to handle such hybrid use cases. Automatically senses the MPI build (mvapich2, openmpi) and binds tasks correctly.
- **“ibrun –help”** gives detailed usage info.
- **hello_hybrid.c** is a sample code, and **hello_hybrid.cmd** shows “ibrun” usage.

hello_hybrid.cmd

```
#!/bin/bash
#SBATCH --job-name="hellohybrid"
#SBATCH --output="hellohybrid.%j.%N.out"
#SBATCH --partition=compute
#SBATCH --nodes=2
#SBATCH --ntasks-per-node=24
#SBATCH --export=ALL
#SBATCH -t 01:30:00
```

#This job runs with 2 nodes, 24 cores per node for a total of 48 cores.

We use 8 MPI tasks and 6 OpenMP threads per MPI task

export OMP_NUM_THREADS=6

ibrun --npernode 4 ./hello_hybrid

Hybrid Code Output

[user@comet-08-11 HYBRID]\$ **more hellohybrid.435461.comet-17-41.out**

Hello from thread 0 out of 6 from process 2 out of 8 on comet-17-41.local

Hello from thread 3 out of 6 from process 2 out of 8 on comet-17-41.local

Hello from thread 4 out of 6 from process 2 out of 8 on comet-17-41.local

Hello from thread 5 out of 6 from process 2 out of 8 on comet-17-41.local

Hello from thread 0 out of 6 from process 3 out of 8 on comet-17-41.local

Hello from thread 2 out of 6 from process 2 out of 8 on comet-17-41.local

Hello from thread 1 out of 6 from process 3 out of 8 on comet-17-41.local

Hello from thread 2 out of 6 from process 3 out of 8 on comet-17-41.local

...

...

Hello from thread 4 out of 6 from process 7 out of 8 on comet-17-42.local

Hello from thread 2 out of 6 from process 7 out of 8 on comet-17-42.local

Hello from thread 3 out of 6 from process 7 out of 8 on comet-17-42.local

Hello from thread 5 out of 6 from process 7 out of 8 on comet-17-42.local

Hello from thread 1 out of 6 from process 6 out of 8 on comet-17-42.local

Using SSD Scratch

```
#!/bin/bash
```

```
#SBATCH --job-name="localscratch"
```

```
#SBATCH --output="localscratch.%j.%N.out"
```

```
#SBATCH --partition=compute
```

```
#SBATCH --nodes=1
```

```
#SBATCH --ntasks-per-node=24
```

```
#SBATCH --export=ALL
```

```
#SBATCH -t 01:30:00
```

```
export WKDIR=`pwd`
```

```
#Change to local scratch (SSD) and run IOR benchmark
```

```
cd /scratch/$USER/$SLURM_JOBID
```

```
#Run IO benchmark
```

```
ibrun -np 4 $WKDIR/IOR.exe -F -t 1m -b 4g -v -v > IOR_native_scratch.log
```

```
cp IOR_native_scratch.log $WKDIR/
```

Using SSD Scratch

- **Snapshot on the node during the run:**

```
$ pwd
```

```
/scratch/mahidhar/435463
```

```
$ ls -lt
```

```
total 22548292
```

```
-rw-r--r-- 1 mahidhar hpss 5429526528 May 15 23:48 testFile.00000001
```

```
-rw-r--r-- 1 mahidhar hpss 6330253312 May 15 23:48 testFile.00000003
```

```
-rw-r--r-- 1 mahidhar hpss 5532286976 May 15 23:48 testFile.00000000
```

```
-rw-r--r-- 1 mahidhar hpss 5794430976 May 15 23:48 testFile.00000002
```

```
-rw-r--r-- 1 mahidhar hpss      1101 May 15 23:48 IOR_native_scratch.log
```

- **Performance from single node (in log file copied back):**

- Max Write: 250.52 MiB/sec (262.69 MB/sec)

- Max Read: 181.92 MiB/sec (190.76 MB/sec)

39

Comet GPU Nodes

2 NVIDIA K-80 Cards (4 GPUs total) per node.

Hands On Example using OpenACC

GPU Node – OpenACC Example

- **Setup PGI compiler:**

- module purge
 - module load pgi

- **Compiling:**

- pgcc -acc -Minfo=accel laplace2d.c

- laplace:

- 58, Generating copy(A[:, :])

- Generating create(Anew[:, :])

- 63, Generating Tesla code

- 64, Loop is parallelizable

- 66, Loop is parallelizable

- Accelerator kernel generated

- 64, #pragma acc loop gang /* blockIdx.y */

- 66, #pragma acc loop gang, vector(128) /* blockIdx.x threadIdx.x */

- 70, Max reduction generated for error

- 74, Generating Tesla code

- 75, Loop is parallelizable

- 77, Loop is parallelizable

- Accelerator kernel generated

- 75, #pragma acc loop gang /* blockIdx.y */

- 77, #pragma acc loop gang, vector(128) /* blockIdx.x threadIdx.x */

GPU Nodes – OpenACC script

```
#!/bin/bash  
#SBATCH --job-name="hpl"  
#SBATCH --output="hpl.%j.%N.out"  
#SBATCH --partition=gpu  
#SBATCH --nodes=1  
#SBATCH --ntasks-per-node=24  
#SBATCH --export=ALL  
#SBATCH --gres=gpu:1  
#SBATCH -t 01:00:00
```

./laplace_openacc



SAN DIEGO SUPERCOMPUTER CENTER

at the UNIVERSITY OF CALIFORNIA; SAN DIEGO



GPU Nodes – OpenACC Example

main()

Jacobi relaxation Calculation: 4096 x 4096 mesh

0, 0.250000

100, 0.002397

200, 0.001204

300, 0.000804

400, 0.000603

500, 0.000483

600, 0.000403

700, 0.000345

800, 0.000302

900, 0.000269

total: 15.914405 s

Gordon : Filesystems

- **Lustre filesystems – Good for scalable large block I/O**
 - Accessible from both native and vSMP nodes.
 - /oasis/scratch/gordon – 1.6 PB, peak measured performance ~50GB/s on reads and writes.
 - /oasis/projects – 1.5PB
- **SSD filesystems**
 - /scratch local to each native compute node – 300 GB each.
 - /scratch on vSMP node – 4.8TB of SSD based filesystem.
 - Few “bigflash” nodes with 4.8TB of SSD space.
- **NFS filesystems (/home)**

Gordon: Compiling/Running Jobs

- **Job Queue basics:**
 - Gordon uses the TORQUE/PBS Resource Manager with the Catalina scheduler to define and manage job queues.
 - Native/Regular compute (Non-vSMP) nodes accessible via “normal” queue.
 - vSMP node accessible via “vsmp” queue.
- **Workshop examples illustrate use of Gordon for two use cases.**
 - **hello_native.cmd** – script for running hello world example on native nodes (using MPI).
 - vSMP example to access the virtual large memory node.

Gordon – Compiling/Running Jobs

- **Change to workshop directory:**

```
cd /home/$USER/SI2015/INTRO/GORDON
```

- **Verify modules loaded:**

```
module list
```

Currently Loaded Modulefiles:

1) intel/2013.1.117 2) mvapich2_ib/1.9 3) gnubase/1.0

- **Compile the MPI hello world code:**

```
mpif90 -o hello_world hello_mpi.f90
```

- **Verify executable has been created:**

```
ls -lt hello_world
```

```
-rwxr-xr-x 1 mahidhar hpss 735429 May 15 21:22 hello_world
```

Gordon: Hello World on native (non-vSMP) nodes

The submit script (located in the workshop directory) is hello_native.cmd

```
#!/bin/bash
#PBS -q normal
#PBS -N hello_native
#PBS -l nodes=4:ppn=1:native
#PBS -l walltime=0:10:00
#PBS -o hello_native.out
#PBS -e hello_native.err
#PBS -V
##PBS -M youremail@xyz.edu
##PBS -m abe
#PBS -A gue998
cd $PBS_O_WORKDIR
mpirun_rsh -hostfile $PBS_NODEFILE -np 4 ./hello_world
```

Gordon: Output from Hello World

- **Submit job using “qsub hello_native.cmd”**

```
$ qsub hello_native.cmd
```

```
845444.gordon-fe2.local
```

- **Output:**

```
$ more hello_native.out
```

```
node      2 : Hello world
```

```
node      1 : Hello world
```

```
node      3 : Hello world
```

```
node      0 : Hello world
```

```
Nodes:    gcn-15-58 gcn-15-62 gcn-15-63 gcn-15-68
```


Sample VSMP queue script

```
#!/bin/bash
#PBS -q vsmp
#PBS -N hello_vsmp
#PBS -l nodes=1:ppn=16:vsmp
#PBS -l walltime=0:10:00
#PBS -o hello_vsmp.out
#PBS -e hello_vsmp.err
#PBS -V
##PBS -M youremail@xyz.edu
##PBS -m abe
#PBS -A gue998
cd $PBS_O_WORKDIR
export LD_PRELOAD=/opt/ScaleMP/libvsmpclib/0.1/
lib64/libvsmpclib.so
export PATH="/opt/ScaleMP/numabind/bin:$PATH"
export KMP_AFFINITY=compact,verbose,0,`numabind --
offset 8`
export OMP_NUM_THREADS=8
./hello_vsmp
```

Bundling Jobs

- HPC systems typically have 16 or more cores per node. On several systems these nodes are provided on a “non-shared” basis.
- Often users have several jobs to run but may not be able to use all the cores on a node for each job => bundling jobs can help throughput.
- SDSC User Services developed bundler scripts to enable such use cases:
 - <https://github.com/sdsc/sdsc-user/tree/master/bundler>
- Can bundle serial and OpenMP jobs.

Job Bundler Example Script

```
#!/bin/bash
```

```
...
```

```
#PBS -N bundler.gordon
```

```
#PBS -q normal
```

```
#PBS -l nodes=1:ppn=8:native
```

```
#PBS -l walltime=1:00:00
```

```
#PBS -v Catalina_maxhops=None,QOS=0
```

```
TASKS=tasks           # the name of your tasks list
```

```
cd $PBS_O_WORKDIR
```

```
module load python/2.7.3   # necessary for bundler.py on Gordon
```

```
mpirun_rsh -export \
```

```
-np $PBS_NP \
```

```
-hostfile $PBS_NODEFILE \
```

```
/home/diag/opt/mpi4py/mvapich2/intel/1.3.1/lib/python/mpi4py/bin/python-mpi \
```

```
/home/diag/opt/sdsc-user/bundler/bundler.py $TASKS
```

Job Bundler Example - Tasks

\$ more tasks

/bin/date > task1.out

/bin/hostname > task2.out

cat /proc/cpuinfo > cpuinfo.dat

df | grep oasis > Lustre.fs.info

df | grep scratch > task5.out

ls /scratch/\$USER > task6.out

ps -ef | grep \$USER > task7.out

/usr/bin/w > task8.out

Output from bundler

Got 8 task slots

MPI process 6 on gcn-3-58.sdsc.edu running task [ps -ef | grep \$USER > task7.out]

MPI process 2 on gcn-3-58.sdsc.edu running task [cat /proc/cpuinfo > cpuinfo.dat]

MPI process 3 on gcn-3-58.sdsc.edu running task [df | grep oasis > Lustre.fs.info]

MPI process 0 on gcn-3-58.sdsc.edu running task [/bin/date > task1.out]

MPI process 7 on gcn-3-58.sdsc.edu running task [/usr/bin/w > task8.out]

MPI process 4 on gcn-3-58.sdsc.edu running task [df | grep scratch > task5.out]

MPI process 5 on gcn-3-58.sdsc.edu running task [ls /scratch/\$USER > task6.out]

MPI process 1 on gcn-3-58.sdsc.edu running task [/bin/hostname > task2.out]

Add Data Analysis to Existing Compute Infrastructure



**Physical
Compute**

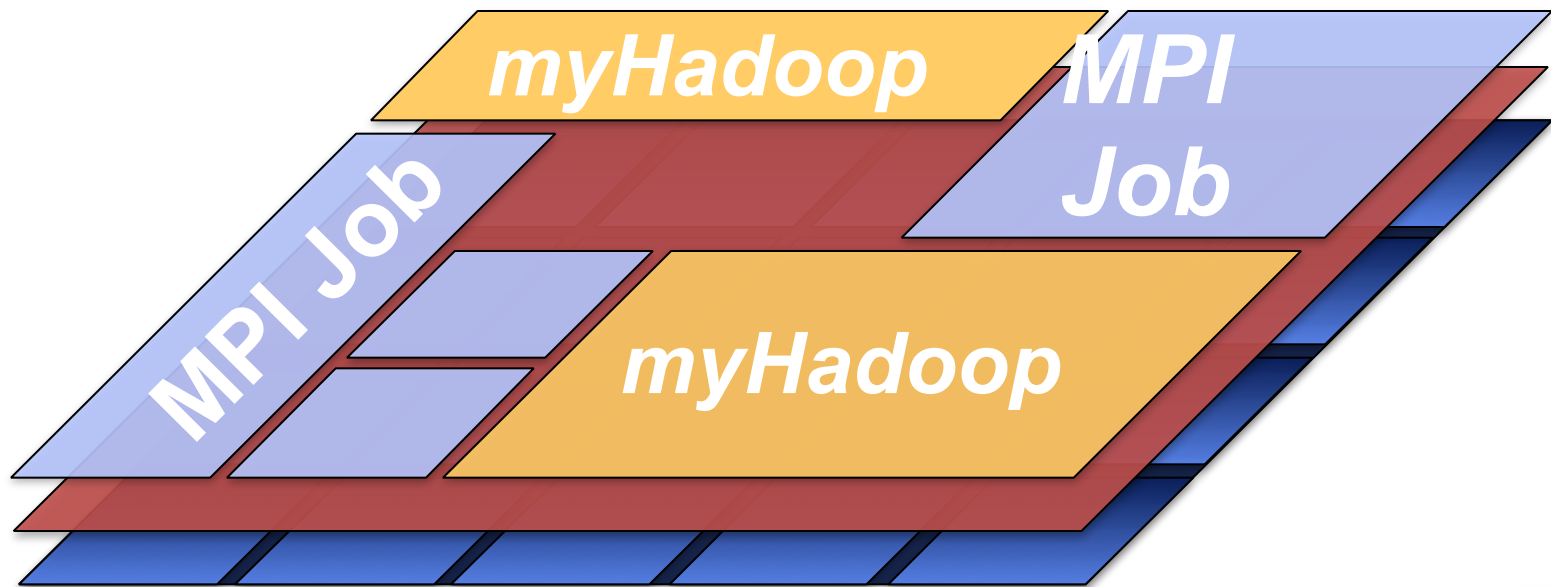
Add Data Analysis to Existing Compute Infrastructure



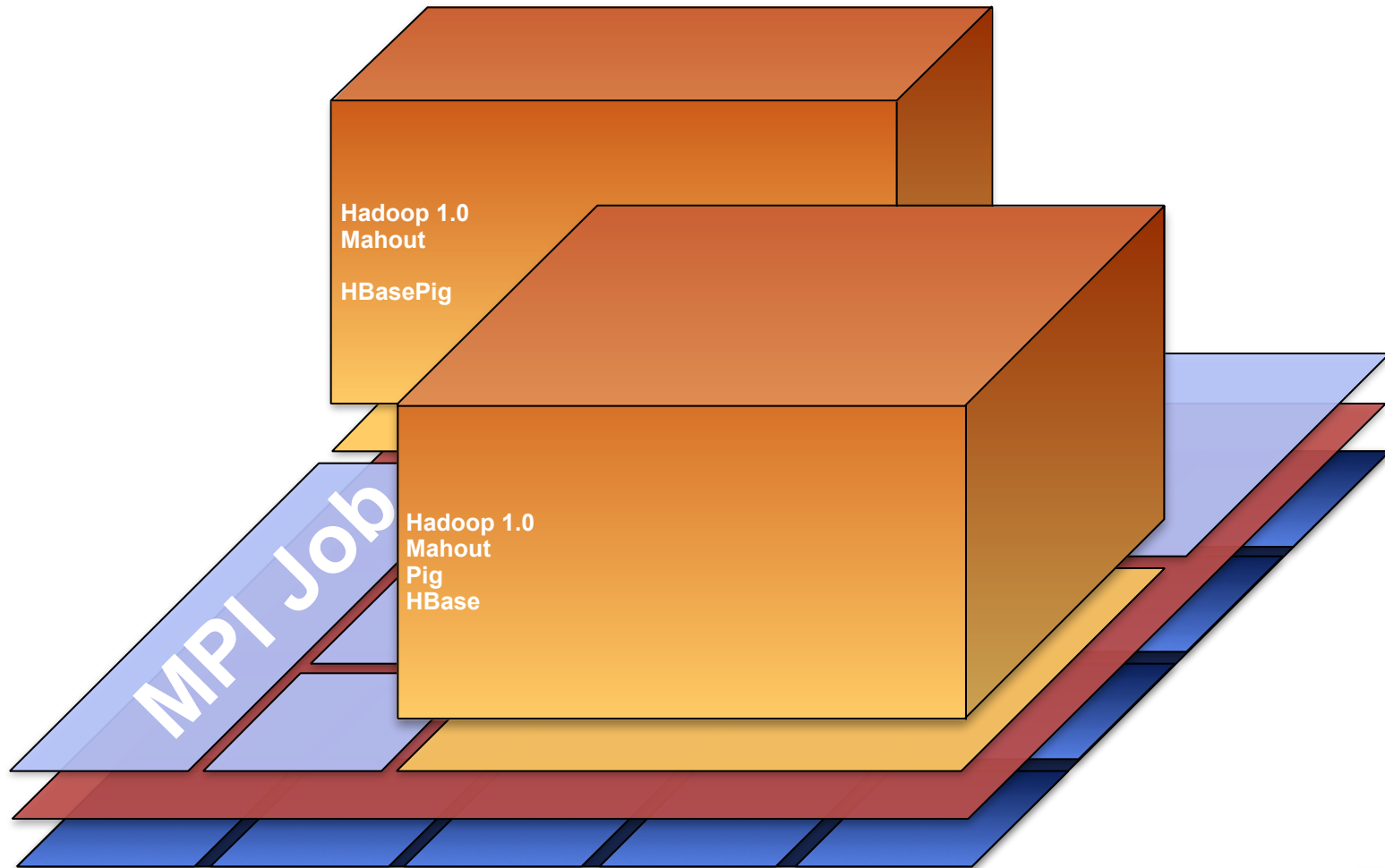
Resource Manager

(Torque, SLURM, SGE)

Add Data Analysis to Existing Compute Infrastructure



Add Data Analysis to Existing Compute Infrastructure



Anagram Example

- **Source:**

https://code.google.com/p/hadoop-map-reduce-examples/wiki/Anagram_Example

- **Uses Map-Reduce approach to process a file with a list of words, and identify all the anagrams in the file**
- **Code is written in Java. Example has already been compiled and the resulting jar file is in the example directory.**

Anagram Example – Submit Script

```
#!/bin/bash
#SBATCH --job-name="Anagram"
#SBATCH --output="Anagram.%j.%N.out"
#SBATCH --partition=compute
#SBATCH --nodes=4
#SBATCH --ntasks-per-node=24
#SBATCH --export=ALL
#SBATCH -t 01:00:00
export WRKDIR=`pwd`
myhadoop-configure.sh
start-all.sh
hadoop dfs -mkdir input
hadoop dfs -copyFromLocal $WRKDIR/SINGLE.TXT input/
hadoop jar $WRKDIR/AnagramJob.jar input/SINGLE.TXT output
hadoop dfs -copyToLocal output/part* $PBS_O_WORKDIR
stop-all.sh
myhadoop-cleanup.sh
```

Anagram Example – Sample Output

cat part-00000

...

aabcdelmnu	manducable,ambulanced,
aabcdeorrsst	broadcasters,rebroadcasts,
aabcdeorrst	rebroadcast,broadcaster,
aabcdkrsw	drawbacks,backwards,
aabcdkrw	drawback,backward,
aabceeehlnsst	teachableness,cheatableness,
aabceeeelnrsstu	uncreatableness,untraceableness,
aabceeeelrrt	recreatable,retraceable,
aabceehlt	cheatable,teachable,
aabceellr	lacerable,clearable,
aabceelnrtu	uncreatable,untraceable,
aabceelorrstv	vertebrosacral,sacrovertebral,

...

...

Spark on Comet

```
#!/bin/bash
#SBATCH --job-name="graphx-demo"
#SBATCH --output="graphx-demo.%j.%N.out"
#SBATCH --partition=compute
#SBATCH --nodes=4
#SBATCH --ntasks-per-node=24
#SBATCH --export=ALL
#SBATCH -t 00:30:00
```

```
### Environment setup for Hadoop and Spark
module load spark
export PATH=/opt/hadoop/2.6.0/sbin:$PATH
export HADOOP_CONF_DIR=$HOME/mycluster.conf
export WORKDIR=`pwd`
```

```
myhadoop-configure.sh
```

```
### Start HDFS. Starting YARN isn't necessary since Spark will be running in
### standalone mode on our cluster.
start-dfs.sh
```

```
### Load in the necessary Spark environment variables
source $HADOOP_CONF_DIR/spark/spark-env.sh
```

```
### Start the Spark masters and workers. Do NOT use the start-all.sh provided
### by Spark, as they do not correctly honor $SPARK_CONF_DIR
myspark start
```

```
### Copy the data into HDFS
hdfs dfs -mkdir -p /user/$USER
hdfs dfs -put $WORKDIR/facebook_combined.txt /user/$USER/
```

```
run-example org.apache.spark.examples.graphx.LiveJournalPageRank facebook_combined.txt --numEPart=8
```

```
### Shut down Spark and HDFS
myspark stop
stop-dfs.sh
```

```
### Clean up
myhadoop-cleanup.sh
```

Summary, Q/A

- **Access options** – ssh clients, XSEDE User Portal, Science Gateways.
- **Data Transfer options** – scp, globus-url-copy (gridftp), globus online, and XSEDE User Portal File Manager.
- **SLURM queues on Comet:** compute, shared, gpu, gpu-shared, and debug.
- **PBS queues on Gordon:** normal (native, non-vSMP) and vsmp.
- **Use SSD local scratch where possible.** Excellent for codes like Gaussian, Abaqus.
- **Hadoop, Spark clusters** can be launched within SLURM (Comet) and PBS/Torque (Gordon) framework using myHadoop set up.