# Introduction to the Piz Daint environment

Directive Based GPU Programming
*Vasileios Karakasis, CSCS*
May 14–15, 2018

# Overview

- Accessing CSCS
- Compiling my code
- Running my code
- Editing my code
- Transferring files from/to CSCS
- Repository of the course

# Piz Daint

Computing nodes

Piz Daint is a hybrid cluster of Cray XC40/XC50 nodes

- Hybrid nodes (XC50)
  - 5320 total
  - Intel® Xeon® E5-2690 v3 @ 2.60GHz (12 cores, 64GB RAM, Haswell)
  - NVIDIA® Tesla® P100 16GB (Pascal)

- Multicore nodes
  - 1813 nodes
  - Two Intel® Xeon® E5-2695 v4 @ 2.10GHz (2 x 18 cores, 64/128 GB RAM, Broadwell)

- Login nodes
  - 5 total
  - Intel® Xeon® CPU E5-2650 v3 @ 2.30GHz (10 cores, 256 GB RAM, Haswell)

- Aries routing and communications ASIC, and Dragonfly network topology

# Piz Daint

Filesystems

- `/scratch`: High performance Lustre filesystem accessible from the computing nodes
  - Environment variable $SCRATCH points to it
  - Total capacity: 6.2 PB
  - Must be used for heavy I/O
- `/users`: GPFS filesystem for the users' homes
- `/project`, `/store`: Long-term storage for computational projects

More on `https://user.cscs.ch/storage/file_systems/`

# Accessing Piz Daint

- Accessible through SSH
- Piz Daint is not directly accessible from the outside world:
  - ela → daint10x → nidxxxxx

Two-steps process:

1. Login to the frontend, forwarding X11 (will be needed the second day)
2. Move to the login nodes of Piz Daint

```
# Login to the frontend first
ssh -Y courseXX@ela.cscs.ch
ssh -Y daint
```

# Programming Environments

Cray Linux Programming Environment

- 4 compilers available: CCE, GNU, INTEL, PGI
- 4 predefined Programming Environments:
  - PrgEnv-cray (default), PrgEnv-gnu, PrgEnv-intel, PrgEnv-pgi
  - echo $PE_ENV to get the current PrgEnv
- 3 wrappers available: ftn (Fortran), cc (C), CC (C++)
  - Required for compiling MPI programs
  - They set appropriate optimisation flags for the target architecture (CPU or GPU)
  - They provide a sort of portability across the programming environments

# Managing programming environments

Daint uses *Environment Modules* (TMod) for managing the programming environments and the software packages:

- Dynamic modification of a user's environment via *modulefiles*.
- All programming environments and software on Daint is available through modules.
- The compiler wrappers will detect the loaded programming environment and automatically set the correct flags and libraries.

# Managing programming environments

Listing modules

- `module list`

# Managing programming environments

Switching programming environments

- Switch to the PGI programming environment
- `module switch`

# Managing programming environments

Switching back to the Cray programming environment

# Managing programming environments

Loading and unloading modules

- `module load [MODULE_NAME]`
- `module unload [MODULE_NAME]`

# Managing programming environments

Checking available modules

- The `daint-gpu` makes available the CSCS software stack built for the hybrid nodes of the system

# Managing programming environments

Checking available modules

- Check available versions of a software

# Managing programming environments

Get information about a module

    – Environment variables set, paths etc.

# Managing programming environments

Get help for a module

# Running on Piz Daint

The job scheduler

Piz Daint uses native SLURM for running jobs on the compute nodes. There are two ways of submitting a job:

1. Interactively from the login nodes using the `srun` command.
2. By submitting a job script using the `sbatch` command.

# Running on Piz Daint

Using the `srun` command

Necessary and useful options:

- `-C gpu`: requests allocation on the hybrid (GPU) nodes (required)
- `--reservation=openacc`: the reservation for our course to avoid waiting times
  - Reservation is valid for the two days of the course from 07:00-19:00.
- `-N 2`: number of compute nodes (default is 1)
- `-n 2`: number of MPI tasks (default is 1)
- `-t 5`: maximum duration of the job (default is 30min)
  - Allows to get an allocation quicker
  - Job will be killed if time limit is reached
  - Maximum time slot for a job is 24h

More on `https://user.cscs.ch/getting_started/running_jobs/`

# Running on Piz Daint

Using the `srun` command



```
course00@daint104:~> srun -Cgpu -t1 -N2 hostname
srun: job 7502915 queued and waiting for resources
srun: job 7502915 has been allocated resources
nid04296
nid04297
course00@daint104:~> srun -Cgpu -t1 -n2 hostname
srun: job 7503014 queued and waiting for resources
srun: job 7503014 has been allocated resources
nid01979
nid01979
course00@daint104:~>
```

# Running on Piz Daint

Using the `sbatch` command



```
● ● ●                                  2. ssh
course00@daint104:~> cat job.sh
#!/bin/bash
#SBATCH -J 'my_first_job'
#SBATCH -C gpu
#SBATCH -N 2
#SBATCH -t 1
#SBATCH -o myjob.out
#SBATCH -e myjob.err

echo "My job id is $SLURM_JOB_ID"
hostname
course00@daint104:~> sbatch job.sh
Submitted batch job 7503300
course00@daint104:~> squeue -j 7503300
   JOBID    USER ACCOUNT        NAME ST REASON    START_TIME         TIME  TIME_LEFT NODES CPUS
 7503300 course00   crs03  my_first_job CG None    20:44:18           0:06       0:54     2    48
course00@daint104:~> squeue -u $USER
   JOBID    USER ACCOUNT        NAME ST REASON    START_TIME         TIME  TIME_LEFT NODES CPUS
course00@daint104:~> ▉
```

CSCS

**ETH**zürich

# Running on Piz Daint

Using the sbatch command – Examinining the output

# Running on Piz Daint

Other useful commands

- `squeue [OPTIONS]`: Check the status of the system job queue
    - Useful options: `-u [USERNAME]`, `-j [JOBID]`
- `scancel [JOBID]`: Cancel a job
- `scontrol`: Detailed information about partitions, reservations, computing nodes etc.

# Running on Piz Daint

Other useful commands

# Editing files

- `vim` or `gvim` (X version)
- `emacs -nw` or just `emacs` (X version)
- `gedit` (X only)

# **Moving data to/from CSCS**

- scp: Remote copy over SSH
  - Getting a file: scp course00@ela.cscs.ch:remotefile localfile
  - Getting a directory: scp -r course00@ela.cscs.ch:remotedir localdir
  - Sending a file: scp localfile course00@ela.cscs.ch:remotefile
  - Sending a directory: scp localdir course00@ela.cscs.ch:remotedir
- rsync: Synchronize files remotely over SSH
  - rsync -avz course00@ela.cscs.ch:remotedir/ localdir/
  - rsync -avz localdir/ course00@ela.cscs.ch:remotedir/
  - Pay attention to the slashes! *rsync behaves differently with or without slashes.*