

# INTRODUCTION TO SUPERCOMPUTING AT JSC

Andreas Smolenko

Benedikt Steinbusch

Alexandre Strube

Max Holicki

Ilya Zhukov

Jolanta Zjupa

21-24 November 2022

- Introduction
- Access
  - Getting a JSC account
  - Joining a compute time project
  - Login procedure
  - Checking System/Service Status
  - Further reading
- Unix shell basics
- Environment
  - Active project
  - File system points of interest
  - Further reading
- Software Modules
  - Further reading
- Custom software
  - Compiled languages
  - Scripting languages
- Transferring data
  - Download files from the web (supported only on login nodes!)
  - Transferring files and folders from/to cluster
  - Archiving files
- Budgeting
  - Job Accounting
  - Data Quotas
- Running jobs
  - Interactive mode
  - Batch mode
  - Affinity and multi-threading
  - Further reading
  - LLview - Detailed Job Reporting
- Using GPUs

- GPU Inspection During Execution
- GPU Affinity
- Network Architecture Study
- Further reading
- Useful Links
  - System Documentation
  - JSC Services
  - Job Reporting
  - Apply for Computing Time
  - Apply for a Data Project
  - JSC Course Programme
  - Supercomputing Support

## INTRODUCTION

The largest computers used for computational science have exhibited an **exponential increase in the rate of basic operations they can perform since at least the 1990s**. For more than a decade, this growth has been enabled, **not by increasing clock speeds of individual processing units**, but by assembling systems that consist of ever greater numbers of processing units. Scientific applications that are meant to run on these systems are expected to orchestrate many of these computational units to collaborate on solving a given computational problem. Building these kinds of applications is called parallel programming. Parallel programming will only be touched on briefly in this course, but Jülich Supercomputing Centre (JSC) offers several courses that teach various techniques related to the topic.

Scientists who want to run applications, be they custom made or third-party, on these systems are expected to know how to use these systems. Working through this guide will teach you how to

- access the systems available at JSC,
- navigate the file system,
- find pre-installed software,
- build your own software, and finally
- run software.

## ACCESS

This chapter will teach you how to log in to one of the systems at JSC.

## Getting a JSC account

A basic prerequisite to get access to the HPC system and other services at JSC is a JSC account. If you do not already have an account (they have the form <family name> <number>, e.g. steinbusch1), one can be created through [JSC's user portal JuDoor](#) (click the *Register* button).

## Joining a compute time project

To be allowed to log in to an HPC system, your JSC account needs to be a member of a computing time project that has an active budget on the system. This is the case if

- you have successfully [applied for test computing time](#) for a test project and are now the principal investigator (PI) of your own project, or
- you have successfully [applied for computing time](#) during one of our calls for project proposals and are now the principal investigator (PI) of your own project, or
- you have gained access to a project either by being invited by the PI or project administrator (PA) or by being granted access upon requesting to join a project through JuDoor.

We have created a computing time project for this course with a project id of training2230. To join the project, log in to [JuDoor](#) and click *Join a project* under the *Projects* heading. Enter the project id and, if you want to, a message to remind the PI/PA (one of the instructors) why you should be allowed to join the project. Afterwards the PI/PA will be automatically informed about your join request and can add you to the different systems available in the project. As soon as you are unlocked for the system, the system entry will be shown on your JuDoor main page. You have to accept our Usage Agreement before you can continue with the next step.

## Login procedure

Logging in to our systems is usually done through the [Secure Shell \(SSH\)](#) mechanism, although there are alternatives such as [UNICORE](#) and [JupyterLab](#). Our SSH configuration uses an authentication mechanism based on public and private keys rather than passwords. A pair of public and private keys has to be generated on your personal computer. The private key has to be protected by a passphrase. The public key is then registered for access to the system through JuDoor.

**NEVER SHARE YOUR PRIVATE KEY!!!**

Several software packages can be used for logging in through SSH. The procedure is documented below for some popular choices:

- **OpenSSH** which is a popular choice on GNU/Linux, macOS, and other Unix-like operating systems
- **PuTTY** which is a popular choice on Windows

## Generating a key pair with OpenSSH

OpenSSH is a set of command line tools, so open up a terminal. We suggest you start by creating a fresh pair of public and private key (a key pair). To generate a key pair enter the command below. The program asks for a passphrase. This passphrase is not used for authenticating to the remote system, but rather acts as an encryption key for the private part of the key pair stored on the local file system. In case the private key file is stolen by an attacker, they will not be able to use the key without knowing the passphrase, so make sure to use one that is **hard to guess**.

```
$ ssh-keygen -a 100 -t ed25519 -f ~/.ssh/id_ed25519
Generating public/private ed25519 key pair.
Enter passphrase (empty for no passphrase):
Enter same passphrase again:
Your identification has been saved in /Users/bsteinb/.ssh/id_ed25519.
Your public key has been saved in /Users/bsteinb/.ssh/id_ed25519.pub.
The key fingerprint is:
SHA256:tHin8v4j4cyVVe2BEWAinq/vlhFExupY+37s94216uA bsteinb@zam478
The key's randomart image is:
+--[ED25519 256]--+
|      .o+ o.o+. |
|      . +oo  ....|
|      o+      . ..|
|      =.o   .   .|
|      = S.oo      |
|      . +o+o      |
|      .=oo+.     .|
|      o*+oo.. oo|
|      .**+Eoo+o. |
+-----[SHA256]-----+
```

If the designated output file (`~/.ssh/id_ed25519`) already exists, the program asks to overwrite it. This is probably *not* what you want, since you might be using the key contained therein. Change the output name by using the arguments `-f ~/.ssh/id_ed25519_jsc` instead of `-f ~/.ssh/id_ed25519`. If you do so, keep in mind that your keys are in a non-default location for the remainder of the course.

Print the contents of the public key to the terminal by entering:

```
$ cat ~/.ssh/id_ed25519.pub
```

```
ssh-ed25519 AAAAC3NzaC1lZDI1NTE5AAAAIAVbTcHgGDpiLJ  
+Sn8DgeRzIRlzTESOMMGcr9UwZjLO user@systemname
```

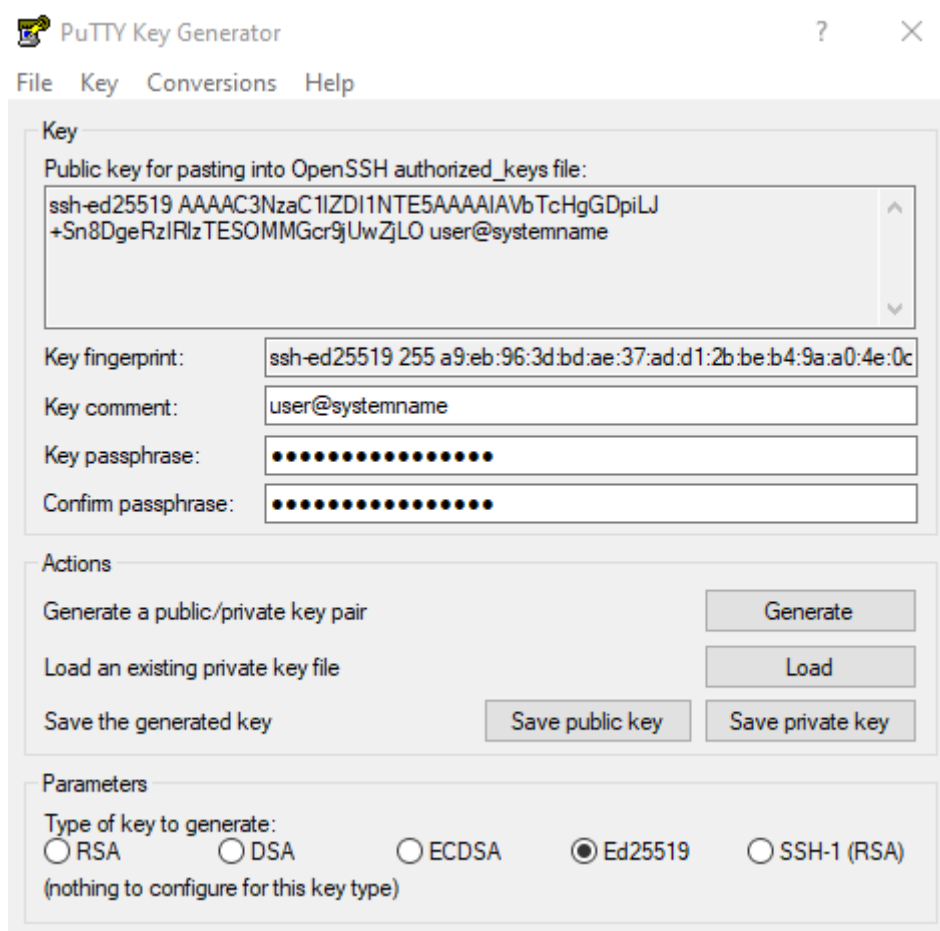
and copy it to the clipboard (do *not* copy the key above, that one is mine, yours will be different). Continue by **uploading the public key to JuDoor**.

## Generating a key pair with PuTTY

Open `puttygen.exe` to generate a key pair. Select *Ed25519* as the key type then click *Generate* and follow the instructions of the program. Once the key has been generated, enter a strong passphrase that cannot be guessed easily. This passphrase is used to encrypt the key while it is stored on disk so that it cannot be used if it is stolen.

Click *Save private key* to save the private key to a `.ppk` file.

Now copy the contents of the field *Public key for pasting into OpenSSH authorized\_keys file* to the clipboard.



Key generation with PuTTY

## Uploading the public key

Navigate to JuDoor and click on *Manage SSH-keys* next to the entry for the system you want to use under the *Systems* heading. Paste the public key into the form in the field labeled

*Your public key and options string*, but do not submit yet. As a further security measure, you have to specify the systems that your log in attempts will come from. This is done via an additional `from`-clause on your public key, that can contain single IP-addresses and address ranges as well as host names and even wildcard patterns based on either of these.

Specifying a `from`-clause is relatively easy if you have access to a system with a fixed IP-address or an IP-address that changes dynamically, but comes from a range of addresses that can be specified concisely. This is typically the case for systems which are connected to university or research centre networks (even via VPN when working from home). For example, systems connected to the network of Forschungszentrum Jülich will be assigned an IP-address from the range `134.94.0.0/16`, so a valid `from`-clause would be `from="134.94.0.0/16"`. Other institutions will use different address ranges, you should be able to find these out from your institutions network operations centre.

Sometimes, patterns based on host names will work better than those based on IP addresses. For example, Forschungszentrum Jülich assigns host names that end in either `fz-juelich.de` or `kfa-juelich.de`, so a valid `from`-clause could also be `from="*.fz-juelich.de,*.kfa-juelich.de"` (notice how multiple patterns can be combined with a comma in between). Once again, the host names assigned by other institutions will be different. To some extent, this scheme also works for home internet access. Internet providers typically assign IP addresses dynamically drawing from fragmented pools that are hard to specify completely in terms of address ranges, but they often assign host names which follow a pattern that can be found out. The command `nslookup <your IP>` will tell you the host name assigned to your system by the provider (find out your IP either from the JuDoor key upload form or by asking a search engine “what is my ip”). This host name might look something like `2909a2-ip.nrw.provider.net`. Chop name components off the beginning and replace them with `*` to come up with a pattern, e.g. `*.nrw.provider.net`.

Add your `from`-clause in front of the public key you already pasted into the form. The result should be something like:

```
from="134.94.0.0/16" ssh-ed25519 AAAA [...]
```

Then click *Start upload of SSH keys*. It will take some time for the key you uploaded to JuDoor to be synched to the actual system. Eventually though, you will be able to log in. Once again, we have instructions for

- [OpenSSH](#)
- [PuTTY](#)

## Logging in with OpenSSH

To log in with OpenSSH, enter the following command:

```
$ ssh -i ~/.ssh/id_ed25519 <account name>@<system name>.fz-juelich.de
```

(Remember to change the location of the key `~/.ssh/id_ed25519` if you saved it to a non-default location.) For example, if I wanted to log in to JUWELS Cluster it would be:

```
$ ssh steinbusch1@juwels-cluster.fz-juelich.de
```

The following table lists the host names of login nodes for the different systems. Pick the one you want to use.

| <b>System</b>  | <b>Login node host name</b>  |
|----------------|------------------------------|
| JURECA DC      | jureca.fz-juelich.de         |
| JUWELS Cluster | juwels-cluster.fz-juelich.de |
| JUWELS Booster | juwels-booster.fz-juelich.de |
| JUSUF          | jusuf.fz-juelich.de          |

When connection for the first time, OpenSSH will prompt you to confirm the server key fingerprint:

```
The authenticity of host 'jusuf.fz-juelich.de (134.94.0.184)' can't be established.  
ECDSA key fingerprint is SHA256:tuswM7JtVcWNS5wRCVIfv1h4uRHReHIN77C4zTYaPjs.  
Are you sure you want to continue connecting (yes/no/[fingerprint])?
```

JSC publishes SSH fingerprints for its systems through JuDoor. You can find them on the page you used to upload your public key. Either compare the keys or, in newer versions of OpenSSH, you can paste the fingerprint from JuDoor into the prompt to confirm it.

Then you should see an informational message (the *message of the day*, **MOTD**) followed by a shell prompt similar to the following:

[illegible]

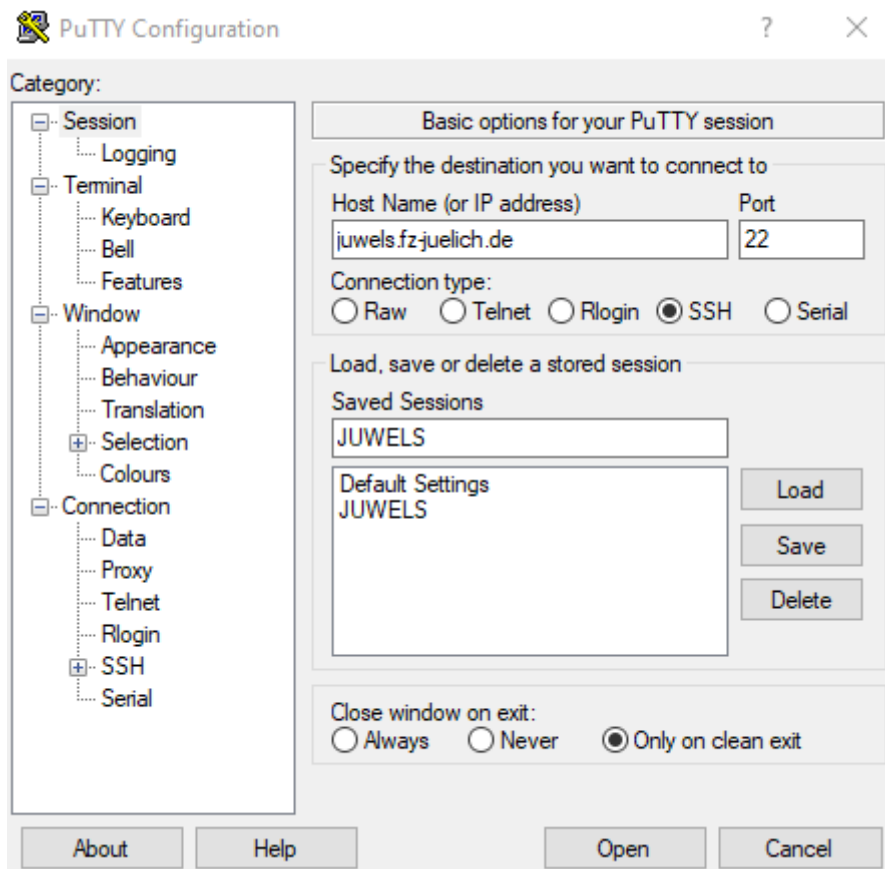
Known issues: <https://apps.fz-juelich.de/jsc/hps/juwels/known-issues.html>

```
*****
steinbusch1@jwlogin01:~ $
```

Once you have logged in successfully, you can continue with **Unix shell basics**.

## Logging in with PuTTY

Launch `putty.exe` to log in. Set the *Host name* for the system you want to connect to, e.g. `juwels-cluster.fz-juelich.de`.



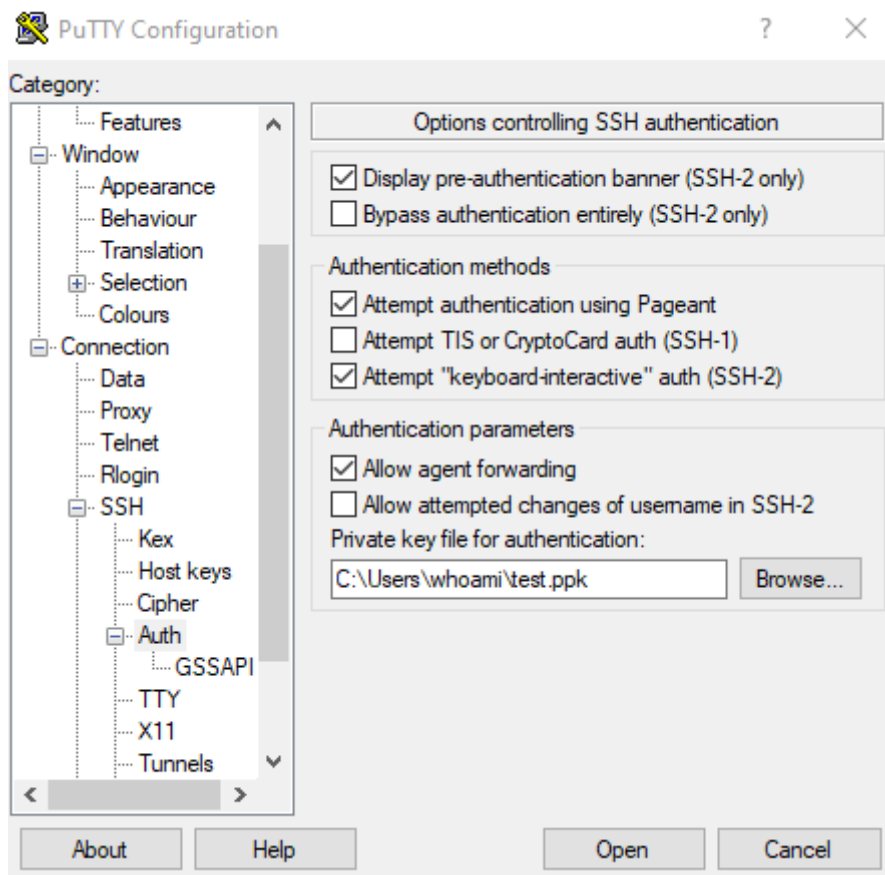
PuTTY session configuration

The following table lists the host names of login nodes for the different systems. Pick the one you want to use.

| System         | Login node host name         |
|----------------|------------------------------|
| JURECA DC      | jureca.fz-juelich.de         |
| JUWELS Cluster | juwels-cluster.fz-juelich.de |
| JUWELS Booster | juwels-booster.fz-juelich.de |
| JUSUF          | jusuf.fz-juelich.de          |

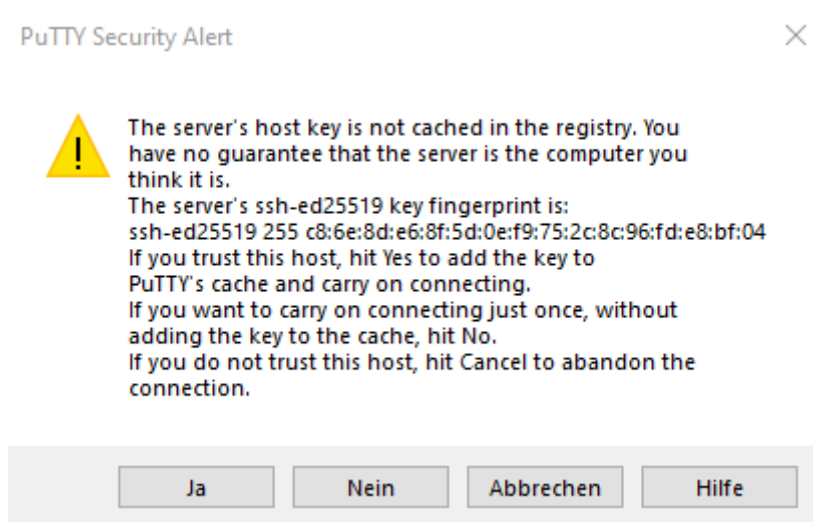
Navigate to *Connection > SSH > Auth* and under *Private key file for authentication*: select the key you just generated.





### PuTTY auth configuration

If you want to save this configuration, you can navigate back to the *Session* screen to give the session a name and save it. Now click *Open* to connect. When you connect for the first time, PuTTY will display a dialog like the following:



### PuTTY security alert

This is not an error, but a security feature. The server key fingerprint displayed in the dialog has to be verified by comparing it to the known good fingerprint. JSC publishes SSH fingerprints for its systems through JuDoor. You can find them on the page you used to upload your public key.

Once you have logged in successfully, you can continue with [Unix shell basics](#).

## JupyterLab

Alternatively, you can [use JupyterLab to log in](#). The authentication credentials are the same as for JuDoor. Once you have logged in, you need to create a JupyterLab instance by clicking *Add New JupyterLab*. On the next screen you must select which system you want to log in to, what project to use for accounting and what part of the system you want to log in to (more about this later), login nodes are the right choice for the moment. Startup of JupyterLab may take a while, but once it is done, you can launch a terminal running a shell on the system of your choice inside the browser. To do so, click *File > New > Terminal* and you should see a shell prompt similar to this:

```
[steinbusch1@jrl06 ~]$
```

## Checking System/Service Status

I cannot log in, the system is slow — the root cause of many problems can be found by checking the [JSC status webpage](#). Here you will find up-to-date status information on the services JSC provides, including upcoming planned maintenances. A trafficlight-colour system is used to indicate the state of a system or service, with green systems functioning as expected for most or all users. Yellow systems are degraded and this will impact many users. Red systems are strongly degraded which will impact most or all users. Finally, dark-red systems are unavailable. Do not fret that you see more systems than what you have access to, this is a general landing page for all our systems.

Below our cluster systems you can find information on our storage systems/tiers. If you cannot find files, a certain mount is unavailable or the system becomes unresponsive to filesystem commands like `ls` this is the place to check. Even further down is the status of the services JSC provides, like JuDoor or Jupyter-JSC. Finally there is the status of JSC-support. Check here if you cannot reach JSC support or they do not respond in a timely fashion, if there is no reported issue try contacting them again using a different mode of communication, *i.e.* telephone or email. At the bottom of the page is a description of the trafficlight icons.

You can get further information on the degradations of the systems and services by clicking on any of the system names, filesystem names or services. Give it a try. There you can also see older issues to help you diagnose problems that may have occurred a couple of days ago.

## Further reading

Our online documentation has more information on accessing the systems. It provides further examples of `from`-clauses, discusses configuration of SSH clients to set up shortcuts and gives hints for troubleshooting. If you want more details, you can find the documentation for our various systems here:

- [JUWELS documentation: Access](#)
- [JURECA documentation: Access](#)
- [JUSUF documentation: Access](#)

## UNIX SHELL BASICS

Whether you log in via OpenSSH or PuTTY or opening a Terminal in JupyterLab, you will be interacting with the system through a **Unix shell**. Unix shells are text based interfaces that prompt the user to input commands and display the result of executing those commands back to the user. The underlying concepts (the file system, executing programs, etc.) are probably familiar to you, but the text based interface can seem daunting at first. This section will teach you how to accomplish essential tasks on a Unix shell. If you are already familiar with this kind of interface, you may want to skip ahead to the [section describing the environment](#).

Like many operating systems, Unix provides an abstraction for storage media called a file system. Data of various types (text, images, executable code, etc.) is stored in files which can be organized in a tree-like hierarchy of directories that starts at a single root (the “root directory”). Objects in the file system (files or directories) are addressed using strings of characters called “paths” that list the directories one has to traverse to get to an object plus the objects name. The slash `/` serves as the separator between elements of a path and cannot itself appear in file or directory names. Some examples for paths are:

```
/etc  
/usr/bin/env  
/home/bsteinb/Documents
```

These paths are all “absolute paths”, meaning, they describe the location of an objects in relation to the root directory (which is represented by a single slash `/`):

- `etc` is a directory that is found inside of the root directory
- `env` is a file found in the directory `bin` which itself is found in the directory `usr` inside the root directory
- `Documents` is a directory in `bsteinb` which is a directory in `home` which is a directory in the root directory

Since absolute paths can become unwieldy in deep directory hierarchies, Unix also allows relative paths. To this end, every program (including the shell you are using) is executed in a “working directory” (which can be changed during the execution of the program). Relative path specifications are then interpreted in relation to this working directory. They are written without the initial slash /. Some examples for relative paths are:

```
Documents
bin/env
../etc/crontab
```

With a working directory of /usr, bin/env refers to /usr/bin/env, just like the absolute path above. The path component .. above has a special function. It refers to the parent (the containing directory) of a file system object and can appear in both relative and absolute paths. So ../etc/crontab refers to a file crontab in a directory etc that can be found in the parent directory of the current working directory.

/home/bsteinb/../janedoe/.bashrc can be simplified to  
/home/janedoe/.bashrc.

To find out the current working directory of the shell you are using, type:

```
$ pwd
```

The output should be something like:

```
/p/home/jusers/steinbusch1/juwels
```

which is the “home directory” associated with your account on the system. To list the contents of the working directory, execute the ls command:

```
$ ls
```

If you are working with a fresh user account, the output of this command might be empty, because there are no files (or only hidden files) in your home directory. To make ls display the hidden files as well, add the optional argument -a:

```
$ ls -a
```

The output should now be non-empty and contain files and directories with names that start with the period .. In Unix, whether a file system item is hidden or not is determined by the first character in its name being the period ..

ls -a is our first example of a more complex command invocation. It starts with the name of a command (so far, we have seen pwd and ls) followed by a list of arguments (here -a), all separated by spaces. ls can be used to list the contents of any directory, by specifying the path of the directory in the last position. To list the items in the /etc directory, type:

```
$ ls /etc
```

Most commands and the list of arguments they accept are documented in the Unix manual pages. They can be accessed through a command – `man` – that takes as its only argument the name of the manual page you want to read. For most commands there is a manual page with the same name as the command. To read the manual page for `ls`, type:

```
$ man ls
```

You can scroll through the manual page using the arrow keys. When you are done reading, close the manual by pressing `q` on the keyboard. To find manual pages for a specific topic, you can use the `apropos` command which searches the library of manual pages for a given keyword.

To change the working directory of your shell and all commands you invoke subsequently, use the `cd` command:

```
$ cd /
```

This will take you to the root directory. If you now execute `ls` without specifying a path, it should show you all items in the root directory, e.g.:

```
$ ls
arch  bin  dev  gpfs  lib   media  opt  proc  run  selinux  sys  usr
arch2 boot etc  home  lib64 mnt    p    root  sbin  srv      tmp  var
```

Invoking `cd` without an argument takes you back to your home directory:

```
$ cd
$ pwd
/p/home/jusers/steinbusch1/juwels
```

Alternatively, the path to your home directory is also available as the value of an “environment variable”. Environment variables map names (strings) to values (also strings) and can be seen as implicit input to commands while arguments on the command line are explicit inputs. The name of the environment variable that contains the path to your home directory is `HOME`. Its value can be inspected using the `printenv` command:

```
$ printenv HOME
/p/home/jusers/steinbusch1/juwels
```

The `printenv` command asks the environment for the value of the variable `HOME` (using the `getenv` function) and prints it to the terminal. In some situations it makes sense, to use the value of environment variables as explicit arguments to a command (e.g. if you want to `cd` to the value of `HOME`). This is supported by a shell mechanism called “variable expansion”: mention the name of a variable, prefixed by the dollar sign `$` in a command

line and the shell will substitute the value of the variable and pass that as an argument to the command:

```
$ cd $HOME
$ pwd
/p/home/jusers/steinbusch1/juwels
```

The `env` command can be used to inspect the environment. When invoked without any arguments it prints a list of all variables currently defined and their values.

```
$ env
[...]
HOME=/p/home/jusers/steinbusch1/juwels
[...]
```

`pwd`, `cd` and `ls` let you navigate the file system. The following commands can be used to make modifications to the file system. First is `mkdir` which allows you to make a directory:

```
mkdir <directory_path>
```

To create an empty file at a given location, use:

```
touch <file_path>
```

In your home directory, create two directories and a file:

```
$ mkdir dir1 dir2
$ touch dir1/file1
```

You can use `ls` to confirm that you have created two directories next to each other, one of which contains an empty file.

```
$ ls
dir1 dir2
$ ls dir1
file1
$ ls dir2
```

Files and directories can be moved, copied and deleted with the commands:

```
$ mv <source_path> <destination_path>
$ cp -r <source_path> <destination_path>
$ rm -r <path>
```

Make a copy of `dir1` and check that it also contains `file1`.

```
$ cp -r dir1 dir3
$ ls dir3
file1
```

Move the copy of the file into `dir2`.

```
$ mv dir3/file1 dir2
$ ls dir2
file1
$ ls dir3
```

Finally, remove all three directories.

```
$ rm -r dir1 dir2 dir3
$ ls
```

Lastly, we will mention one way of editing text files: the nano editor. To open a file in nano, type:

```
$ module load nano
$ nano <file_path>
```

(The `module` command will be explained in detail later on.)

To insert something into the file, just start typing. Save your changes by pressing *CTRL-O*. Exit the editor by pressing *CTRL-X*. The bottom part of the terminal will display more functions which can be reached using certain key bindings. Interaction with the editor, such as specifying a file name when saving, will also happen here.

## ENVIRONMENT

Now that you know about the basic Unix commands, this section will teach you about some of the peculiarities of the environment on the systems at JSC.

### Active project

The first point to talk about is the *active project*. You already know about accounts and computing time projects and by this point you should be a member of at least one project to have access to one of our systems. However, in general, a single user account can be a member of multiple computing time (“C”) projects (and also data projects (“D”)) at the same time. You can see the projects that you are currently a member of in your user profile on JuDoor, or, if you are logged in to one of the HPC systems, you can use the `jutil` command:

```
$ jutil user projects
```

| project    | unixgroup  | PI-uid   | project-type | budget-accounts |
|------------|------------|----------|--------------|-----------------|
| hello      | hello      | hellopi1 | D            | -               |
| chello     | chello     | hellopi1 | C            | hello           |
| training00 | training00 | coach2   | C            | training00      |

Certain system resources, like file system space and compute time, are associated with the projects that you are a member of. Performing actions that consume these resources, storing files or running a computation, have to be counted against the resource pool available to the project. This is done by storing files in certain locations or specifying a compute time budget when running computations. It is possible to explicitly specify a project, each time one of these actions is performed. For brevity's sake, one can also make one of the projects the “active project” and then all actions performed in the remainder of the session will implicitly be performed in the context of that project. This can also be done through the `jutil` command:

```
$ jutil env activate -p training2230 -A training2230
```

Now `training2230` is the active project. Any computational jobs will be accounted against its budget and the special file system locations associated with it can be reached through certain environment variables. More about that in the next section.

Hint: In case you are working on different compute budgets we recommend to set the budget explicitly as it is described later in the document to avoid using the “wrong” budget for a specific simulation job.

## File system points of interest

Every user account on the systems has a home directory (reachable through the `HOME` environment variable) where the user can store his personal files. However, there is a limit on the volume of data and also the number of files that can be stored in this directory. Furthermore, the file system performance in `HOME` is reduced. It is recommended to use `HOME` only for configuration files. More storage space is granted to computing time projects. At least two directories are created for each project:

- a `PROJECT` directory, that can store medium amounts of data, offers modest performance and is backed up regularly, and
- a `SCRATCH` directory, that offers high I/O bandwidth, should be used for input and output of computations (however, no back up is performed and files that have not been touched in 90 days get deleted automatically).

Data projects have access to other storage locations, e.g. the tape based `ARCHIVE` for long term storage of results.

The path of these directories is available as the value of environment variables of the form `<directory>_<project>`, e.g. `PROJECT_training2230` or `SCRATCH_training2230`. If you have activated a project in the previous section, you



will also have environment variables that are just PROJECT and SCRATCH that point to the respective directories of the active project.

Print the contents of PROJECT\_training2230 and PROJECT:

```
$ printenv PROJECT_training2230
/p/project/training2230
$ printenv PROJECT
/p/project/training2230
```

Change into that directory and see what is already there:

```
$ cd $PROJECT_training2230
$ ls
```

Inside the PROJECT directory, make a directory to contain the files that you work on. In order to avoid collisions, use your account name as the name of the directory (the USER environment variable contains your user name):

```
$ mkdir $USER
```

There is more information on [file system points of interest](#) in the documentation.

## Further reading

Our online documentation has more information on the system environment. It describes further file systems covering more specialised use cases and discusses transferring files to and from the systems via SSH and Git. If you want more details, you can find the documentation for our various systems here:

- [JUWELS documentation: Environment](#)
- [JURECA documentation: Environment](#)
- [JUSUF documentation: Environment](#)

## SOFTWARE MODULES

HPC centres will usually make some effort to provide software that is commonly used for scientific purposes. This includes compilers, parallel programming libraries like MPI, numerical libraries, and even complete simulation programs. These software packages form a hierarchy of dependencies (simulation programs use the numerical and parallel programming libraries and all of it is compiled with a certain compiler). Towards the bottom of this hierarchy, packages tend to be interchangeable (several compilers for C or Fortran, several libraries implement the MPI standard) and some of the higher up packages perform better for example when compiled with a certain compiler. Therefore,

it makes sense to offer a range of software packages that implement low level functions and then build a software landscape upon each combination of those low level packages. The two lowest levels in this hierarchy, compiler and MPI library together form a “toolchain”. To help keep the complexity of accessing these different collections of software in check, JSC uses a combination of **EasyBuild** and **Lmod** to build software and make it available as software modules. During a log in session, modules can be loaded and unloaded using the `module` command to use the software that is provided by them. When you log in, a set of default modules is loaded for you, e.g. on JUWELS:

```
$ module list
```

Currently Loaded Modules:

```
1) GCCcore/.9.3.0 (H)   3) binutils/.2.34 (H)
2) zlib/.1.2.11   (H)   4) StdEnv/2020
```

Where:

H: Hidden Module

To see what other modules can currently be loaded, type:

```
$ module avail
```

```
----- Core packages -----
Advisor/2020_update3
Autotools/20200321
Autotools/20200321 (D)

[...]

unzip/6.0
xpra/4.0.4-Python-3.8.5
zsh/5.8

----- Compilers -----
GCC/9.3.0 NVHPC/20.9-GCC-9.3.0 (g)
Intel/2020.2.254-GCC-9.3.0 NVHPC/20.11-GCC-9.3.0 (g,D)
NVHPC/20.7-GCC-9.3.0 (g) NVHPC/21.1-GCC-9.3.0 (g)

----- User-based install configuration -----
UserInstallations/easybuild

Where:
S: Module is Sticky, requires --force to unload or purge
g: built for GPU
L: Module is loaded
Aliases: Aliases exist: foo/1.2.3 (1.2) means that "module load foo/1.2" will load
foo/1.2.3
D: Default Module
```

Use "module spider" to find all possible modules and extensions.

Use "module keyword key1 key2 ..." to search for all possible modules matching any of the "keys".

The available modules are grouped into categories:

- Core packages, which are independent of the choice of toolchain
- Compilers, which are the first ingredient of a toolchain
- Architectures, that can be used to load software for different processor architectures, this category does not exist on all systems

Go ahead and load a compiler:

```
$ module load GCC
```

If you now run `module avail` again, you will notice two additional software categories:

```
$ module avail
```

```
----- MPI runtimes available for GNU compilers -----  
[...]
```

```
----- Packages compiled with GNU compilers -----  
[...]
```

These contain modules that depend on (or were built with) the GCC module that you just loaded. Loading one of the available MPI modules will complete your choice of a toolchain and make more software available:

```
$ module load OpenMPI  
$ module avail
```

```
----- OpenMPI settings -----  
mpi-settings/CUDA-low-latency    mpi-settings/CUDA (L,D)
```

```
----- Packages compiled with OpenMPI and GCC compilers -----  
[...]
```

If you are looking for a particular piece of software that you know the name of, rather than rummaging through all the toolchains, you can use the `module spider` subcommand, as the output of `module avail` suggests:

```
$ module spider LAMMPS
```

```
-----  
LAMMPS:  
-----
```

Description:

LAMMPS is a classical molecular dynamics code, and an acronym for Large-scale Atomic/Molecular Massively Parallel Simulator. LAMMPS has potentials for solid-state materials (metals, semiconductors) and

soft matter (biomolecules, polymers) and coarse-grained or mesoscopic systems. It can be used to model atoms or, more generically, as a parallel particle simulator at the atomic, meso, or continuum scale. LAMMPS runs on single processors or in parallel using message-passing techniques and a spatial-decomposition of the simulation domain. The code is designed to be easy to modify or extend with new functionality.

#### Versions:

LAMMPS/24Dec2020-CUDA

LAMMPS/24Dec2020

LAMMPS/29Oct2020-CUDA

LAMMPS/29Oct2020

---

For detailed information about a specific "LAMMPS" package (including how to load the modules) use the module's full name.

Note that names that have a trailing (E) are extensions provided by other modules.  
For example:

```
$ module spider LAMMPS/29Oct2020
```

---

#### Loading the LAMMPS module with OpenMPI loaded fails:

```
$ module load LAMMPS
```

```
Lmod has detected the following error:  These module(s) or  
extension(s) exist but cannot be loaded as requested: "LAMMPS"
```

```
Try: "module spider LAMMPS" to see how to load the module(s).
```

module spider with a specific module version provides details on how the module can be loaded:

```
$ module spider LAMMPS/24Dec2020
```

---

```
LAMMPS: LAMMPS/24Dec2020
```

---

#### Description:

LAMMPS is a classical molecular dynamics code, and an acronym for Large-scale Atomic/Molecular Massively Parallel Simulator. LAMMPS has potentials for solid-state materials (metals, semiconductors) and soft matter (biomolecules, polymers) and coarse-grained or mesoscopic systems. It can be used to model atoms or, more generically, as a parallel particle simulator at the atomic, meso, or continuum scale. LAMMPS runs on single processors or in parallel using message-passing techniques and a spatial-decomposition of the simulation domain. The code is designed to be easy to modify or extend with new functionality.

You will need to load all module(s) on any one of the lines below before the

"LAMMPS/24Dec2020" module is available to load.

```
GCC/9.3.0 ParaStationMPI/5.4.7-1
Intel/2020.2.254-GCC-9.3.0 ParaStationMPI/5.4.7-1
```

Help:

Description

=====

LAMMPS is a classical molecular dynamics code, and an acronym for Large-scale Atomic/Molecular Massively Parallel Simulator. LAMMPS has potentials for solid-state materials (metals, semiconductors) and soft matter (biomolecules, polymers) and coarse-grained or mesoscopic systems. It can be used to model atoms or, more generically, as a parallel particle simulator at the atomic, meso, or continuum scale. LAMMPS runs on single processors or in parallel using message-passing techniques and a spatial-decomposition of the simulation domain. The code is designed to be easy to modify or extend with new functionality.

More information

=====

- Homepage: <https://lammps.sandia.gov/>
- Site contact: [a.kreuzer@fz-juelich.de](mailto:a.kreuzer@fz-juelich.de)

The problem is that LAMMPS is only available in toolchains which include ParaStationMPI. It is not necessary to reload the entire toolchain, it is enough to reload the MPI runtime:

```
$ module load ParaStationMPI
$ module load LAMMPS
```

Specific modules can be unloaded again using the `module unload` command. To unload (almost) all modules and start with a fresh environment, use `module purge`.

The `module` command is part of the Lmod software package. It comes with its own help document which you can access by running `module help` and a [user guide is available online](#).

The JUWELS system is special in terms that it consist of multiple system modules (as opposed to software modules) based on different compute technologies. The software we provide on JUWELS is also split into different hierarchies, one per system module. As JUWELS uses different login nodes for the different system modules (Cluster and Booster), the correct software collection is loaded automatically based on which login node you use.

## Further reading

Our online documentation has more information on software modules. It lists the basic tool chains (compiler + communication library + math library) available on our systems and discusses using older software stages. If you want more details, you can find the documentation for our various systems here:

- [JUWELS documentation: Software Modules](#)
- [JURECA documentation: Software Modules](#)
- [JUSUF documentation: Software Modules](#)

## CUSTOM SOFTWARE

For some, the software that is made available via the module system is enough to do their daily work. Others will want to bring their own software to the systems. This chapter will teach you how to run software distributed as source code for both compiled programming languages and scripting languages.

### Compiled languages

For the three most common compiled languages in scientific computing, C, C++, and Fortran, the basic workflow is very similar. Open the file `hellompi.c` in the nano editor (or a different editor of your choice). (nano is available as a module, if you want to use it, type `module load nano`.)

```
$ nano hellompi.c
```

Paste the following listing into the file, save and close the editor.

```
#include <stdio.h>
#include <mpi.h>

int main(int argc, char* argv[]) {
    MPI_Init(&argc, &argv);

    int r, s;
    MPI_Comm_rank(MPI_COMM_WORLD, &r);
    MPI_Comm_size(MPI_COMM_WORLD, &s);
    printf("hello from process %d of %d\n", r, s);

    MPI_Finalize();
}
```

Once you have a compiler and an MPI library loaded (e.g. `module load GCC OpenMPI`), the file can be compiled as follows:

```
$ mpicc -std=c11 -o hellompi hellompi.c
```

We will explain how to run the program in a later chapter.

A lot of software is not compiled and installed by invoking the compiler directly, but by using a build system. GNU make is installed from the operating system package sources and GNU autotools as well as CMake are available as modules. More exotic build systems are also available, as are compilers for other languages like Go or Rust.

## Scripting languages

Scripting languages have become more popular in scientific computing recently. Modules are available for Python and Julia.

### Python

The Python interpreter can be loaded as a module as well as the mpi4py package that allows you to use MPI from your Python programs.

```
$ module load Python mpi4py
```

Edit a file `hellompi.py`:

```
($ module load nano)
$ nano hellompi.py
```

And paste the following content into it, then save and exit the editor.

```
from mpi4py import MPI

r = MPI.COMM_WORLD.rank
s = MPI.COMM_WORLD.size

print(f"hello from process {r} of {s}")
```

We will explain how to run the program in a later chapter.

More Python packages are available as modules. For scientific computing, the SciPy-Stack collection is especially interesting.

## TRANSFERRING DATA

With increasing supercomputer performance the data produced through simulation increases. Data management needs to be considered for every compute time project. In some workflows it could be necessary to get access to storage technology with improved IO bandwidth. At JSC there are **several storage technologies** serving different needs. Access is granted for some of the storage technologies through **application for a data**

**project** which can be submitted at any time. For large data transfers to or from the supercomputer infrastructure at JSC the system **JUDAC** with the address `judac.fz-juelich.de` delivers maximal bandwidth performance.

## Download files from the web (supported only on login nodes!)

`wget` is a simple file downloader that allows downloading files using HTTP, HTTPS, and FTP protocols. `wget` supports a number of options allowing to download multiple files, resume downloads, limit the bandwidth, recursive downloads, download in the background, etc.

Here is the typical syntax

```
$ wget <url link to the file>
```

## Transferring files and folders from/to cluster

### scp

`scp` allows to copy files over a secure, encrypted network connection. As `scp` command uses SSH to transfer data, it requires a password for authentication.

Copy file to the cluster

```
[from your laptop] $ scp [options] /path/to/file/filename <account name>@<system name>.fz-juelich.de:/path/where/to/copy
```

Download file from the cluster

```
[from your laptop] $ scp [options] <account name>@<system name>.fz-juelich.de:/path/to/the/file /path/where/to/save
```

To recursively copy a directory, use the `-r` (recursive) option.

### rsync

If you already experienced with `scp`, you can test `rsync`. The `rsync` utility provides many advanced features for file transfer.

The syntax is similar to `scp`. Here is an example of file transfer to the cluster with commonly used options

```
[from your laptop] $ rsync -avzP /path/to/file/filename <account name>@<system name>.fz-juelich.de:/path/where/to/copy
```



## Download file from the cluster

```
[from your laptop] $ rsync -avzP <account name>@<system name>.fz-juelich.de:/path/to/the/file /path/where/to/save
```

Where \* -a (archive) preserves the date and times, and permissions of the files; \* -v (verbose) option gives verbose output to help monitor the transfer; \* -z (compression) option compresses the file during transit to reduce size and transfer time; \* -P (partial/progress) option preserves partially transferred files in case of an interruption and also displays the progress of the transfer.

## SSHFS

SSHFS allows you to mount a remote filesystem using SFTP.

To mount a remote filesystem you can do the following

- Make sure that **SSHFS** is installed on your local machine, e.g.

```
[from your laptop] $ which sshfs
```

Output will show where SSHFS is installed. If the result is empty, you need to install it (or tell the shell which directories to search for executable files).

- Create a directory which will be your mounting point

```
[from your laptop] $ mkdir <mountpoint>
```

- To mount remote directory

```
[from your laptop] $ sshfs <account name>@<system name>.fz-juelich.de:/path/to/directory /path/to/mounting/point
```

- To unmount the filesystem

```
[from your laptop] $ fusermount -u <mountpoint>
```

On BSD and macOS, to unmount the filesystem

```
[from your laptop] $ umount <mountpoint>
```

## Alternatives

- On Windows you can use various clients, e.g. **WinSCP**, **FileZilla**, PuTTY, etc.
- **UFTP (UNICORE FTP)** is a file transfer tool similar to Unix' FTP. Its main features include high-performance file transfers from client to server (and vice versa), list directories, make/remove files or directories, sync files and data sharing. In

addition, users can easily share their data even with users who do not have Unix-level access to the data.

- **GridFTP** is an extension of **FTP** used within large science projects. It includes features like parallelized FTP streams, fault tolerancy, download of portions of data and authentication and encryption for file transfers.

## Archiving files

One of the biggest problems we often encounter when transferring data between remote HPC systems is the transfer of large numbers of files. There is an overhead involved in transferring each individual file, and when transferring a large number of files, this overhead in combination slows down the data transfer dramatically.

This issue can be solved by archiving multiple files into a smaller number of larger files before transferring the data. It is also possible to combine archiving with compression to reduce the amount of data we need to transfer, thereby speeding up the transfer. This can be done for example with tar utility.

Here is an example of archiving all data from a specific directory

```
$ tar -cvf <achive name>.tar /path/to/data/to/be/archived
```

Extract data from the archive

```
$ tar -xvf <achive name>.tar
```

Where \* -c (create) create new archive; \* -v (verbose) option gives verbose output to help monitor the archiving process; \* -f (file) filename of the archive; \* -x (extract) extract files from an archive.

To create a compressed archive using tar we add the -z option and add the .gz extension to the file to indicate it is compressed

```
$ tar -czvf <achive name>.tar.gz /path/to/data/to/be/archived
```

Please note that data compression and decompressing can take longer than transferring the un-compressed data.

The extract compressed files from the archive you can use the same way as for uncompressed data as tar recognizes it is compressed and decompresses and extracts at the same time.

```
$ tar -xvf <achive name>.tar.gz
```

# BUDGETING

There is a large amount of users involved in using supercomputing resources. In the application phase it is made sure that they are in need of this amount of computing resource. Would we let everyone use all the resources, monopolizing of the compute time and/or data storage capabilities would be quickly happening. With budgeting this fate is prevented. There are budgets on compute time, the amount of data and the number of files stored. This ensures that every user can use a portion of the supercomputing facilities at JSC.

## Job Accounting

Each computing time project has been granted a certain amount of compute time (core hours) on an HPC system. This budget is split monthly over the runtime of a project so that a regular project that runs for 12 months has 1/12 of the total amount of the granted core-h available each month. To allow further flexibility we have established a “3-month-window”: Core hours that have not been used in the previous month can be used in the current month and will be lost in the next month if they are not used in the current month. Whereby in the current month you can also use the quota of the next month but with a decreased priority of the submitted jobs. The priority will be further decreased if you have used up even the quota of the next month.

Users are charged for complete nodes they occupy, regardless of the number of CPUs used since the requested compute nodes for your application are not shared among users. The compute time used for one job will be accounted by the following formula:  $\text{\#nodes} * \text{\#AvailableCoresPerNode} * \text{walltime}$ .

Jobs that run on nodes equipped with GPUs are charged in the same way. Independent of the usage of the GPUs the available cores on the host CPU node are taken into account.

Detailed information of each job can be found in KontView which is accessible via the button ‘show extended statistics’ for each project in [Judoor](#).

Alternatively, you can execute the following command on the login nodes to query your CPU quota usage: `jutil user cpuquota`. Further information can be found in the “Accounting” chapter of the corresponding [System Documentation](#).

## Data Quotas

There are limitations on the amount of data and the number of files on each [file system](#). The usage of the data within a project is visualized in [JuDoor](#). Following the links within

JuDoor to KontView, more detailed statistic on the data usage are visualized.

Applications for a data project, giving access to other data storage facilities than PROJECT or SCRATCH of the compute time projects, can be submitted according to the information [here](#). Applications for data projects are processed in a rolling manner. Therefore, you can apply at any point should you see the need during a compute time project for such a data project.

## RUNNING JOBS

Up to now, you have been working on the log in nodes of the system. These nodes are set aside for working interactively on tasks that are needed to prepare your computations, such as compiling your applications, moving input data into place, and writing configuration files for your programs. Since the number of log in nodes for each system is small and they are shared between all users, we ask you to keep the resource consumption on these systems as low as possible. Building software should be restricted to using only a few processes in parallel, simulations and post-processing jobs should be run on the compute nodes. Use the `who` command to see who else is logged in to the log in node you are currently using:

```
$ who
steinbusch1 pts/71      2021-03-11 09:51 (pool-148-54.vpn.kfa-juelich.de)
[...]
$ who | wc -l
59
```

Unlike the log in nodes, users are not given free access to the compute nodes at any time. Instead they form a pool of resources managed by the resource manager software. Due to our collaboration with the company [Partec](#) we use “psslurm” which is based on [Slurm](#) and optimized for our systems to manage these resources. To run a computation on the compute nodes, you have to specify to the resource manager what amount of resources you need and for which duration. Once the resources have become available, you will be allowed to execute programs on them. Two modes of operation are possible:

- interactive mode where programs can be run on the allocated resources from a shell, possibly repeatedly, and
- batch mode where a shell script describing the commands to run as part of a computation is handed off to the resource manager for asynchronous execution.

## Interactive mode

## One-shot

The `srun` command is used to execute commands on a set of allocated resources. If no resources are currently allocated, `srun` can infer from its command line arguments what resources are needed, request them from the resource manager and defer the execution of the associated commands until the resources are available. After the associated commands have been run, the resources are relinquished and running further commands will have to ask for resources again. This one-shot mode can be useful when you want to interactively run a few quick jobs with varying sets of resources allocated for them. Run the `hostname` command to see how `srun` will run commands on different nodes than the log in nodes. On JURECA and JUSUF, use this command (Important: do not forget to replace `YYYYMMDD`, where `YYYY` and `MM` and `DD` are the current year and month and day in the Gregorian calendar, e.g. 20221122):

```
$ hostname
jrlogin09.jureca
$ srun -A training2230 --reservation hands-on-YYYYMMDD hostname
srun: job 3472578 queued and waiting for resources
srun: job 3472578 has been allocated resources
jrc0454
```

For the JUWELS Cluster and JUWELS Booster, there are a few differences: The name of the reservation on JUWELS Cluster is `hands-on-cluster-YYYYMMDD` and `hands-on-booster-YYYYMMDD` on JUWELS Booster. To submit to JUWELS Cluster, you want to be logged in to the Cluster login nodes:

```
$ hostname
jwlogin02.juwels
$ srun -A training2230 --reservation hands-on-cluster-YYYYMMDD hostname
srun: job 9792359 queued and waiting for resources
srun: job 9792359 has been allocated resources
jwc06n213.juwels
```

To submit to JUWELS Booster, you want to be logged in to the Booster login nodes and you have to specify the number of GPUs you want to use

```
$ hostname
jwlogin24.juwels
$ srun -A training2230 --reservation hands-on-booster-YYYYMMDD --gres gpu:4 hostname
srun: job 4575092 queued and waiting for resources
srun: job 4575092 has been allocated resources
jwb0053.juwels
```

Please keep these differences in mind if you are using JUWELS Booster, they will not be repeated in further examples.

Invocations of the `srun` command have the following syntax:

```
$ srun <srun options...> <program> <program options...>
```

Above we have seen four `srun` options:

- `-A` (short for `--account`) to charge the resources consumed by the computation to the budget allotted to this course (if you have used `jutil env activate -A training2230` earlier on, you do not need this)
- `--reservation` to use nodes which have been set aside for this course. For this course we have active reservations for the following systems: JURECA, JUWELS Cluster, JUWELS Booster and JUSUF. For JURECA and JUSUF use the following reservation: `hands-on-YYYYMMDD`. To work on JUWELS Cluster or Booster modules, you have to use `hands-on-cluster-YYYYMMDD` or `hands-on-booster-YYYYMMDD` respectively. Do not forget to replace `YYYYMMDD`, where `YYYY` and `MM` and `DD` are the current year and month and day in the Gregorian calendar, e.g. `20221122`.
- `--partition` specifies which set of compute nodes to request resources from. We typically group nodes of the same hardware type into a partition.
- `--gres` specifies additional resources, other than compute nodes, in this case the presence of four GPUs in the compute nodes.

For the `<program>` we used `hostname` with no arguments of its own.

To run more parallel instances of a program, increase the number of Slurm *tasks* using the `-n` option to `srun`:

```
$ srun --label -A training2230 --reservation hands-on-cluster-YYYYMMDD -n 10 hostname
srun: job 3472812 queued and waiting for resources
srun: job 3472812 has been allocated resources
8: jwc00n002.juwels
9: jwc00n002.juwels
0: jwc00n002.juwels
1: jwc00n002.juwels
6: jwc00n002.juwels
3: jwc00n002.juwels
5: jwc00n002.juwels
2: jwc00n002.juwels
7: jwc00n002.juwels
4: jwc00n002.juwels
```

If you do not tell Slurm that your commands are multi-threaded (`hostname` is not), it will assume each task only needs a single CPU core and pack as many as possible into a node. Note also the `--label` option to `srun` which prefixes every line of output by a number that identifies the task that generated the output.

Running more tasks than will fit on a single node will allocate two nodes and split the tasks between nodes:

```
$ srun --label -A training2230 --reservation hands-on-cluster-YYYYMMDD -n 100 hostname
srun: job 3473040 queued and waiting for resources
srun: job 3473040 has been allocated resources
 0: jwc00n007.juwels
[...]
50: jwc00n008.juwels
[...]
```

Allocations always contain entire nodes exclusively. So your jobs should request a number of tasks that is divisible by the number of tasks which can fit on a node to avoid losing parts of your budget.

You can now also use `srun` to run the `hellompi` program introduced in the previous section on deploying custom software:

```
$ srun -A training2230 --reservation hands-on-cluster-YYYYMMDD -n 5 ./hellompi
srun: job 3471349 queued and waiting for resources
srun: job 3471349 has been allocated resources
hello from process 4 of 5
hello from process 0 of 5
hello from process 3 of 5
hello from process 1 of 5
hello from process 2 of 5
```

## Interlude: Partitions

The systems at JSC typically provide more than one pool of resources, called *partitions*. The resources in the different partitions might have different hardware characteristics or cater to different use cases.

The previous examples were run on the default partition of the system you are using, batch on JUWELS Cluster and JUSUF Cluster, booster on JUWELS Booster and dc-cpu on JURECA. You can find out what partitions the different systems have in the documentation for [JURECA](#), [JUWELS](#), and [JUSUF](#).

Of particular interest are the development partitions on each system (look for `devel` in their name). These consist of a small number of nodes which are set aside to prioritise small and short jobs which are typically run as part of development work on your application rather than production use of the system.

Try running the previous two examples using `hostname` on the development partition of your system by specifying it through `srun`'s `-p` option. Remove the `--reservation` option, because the reservation does not include nodes from the development partition.

We will have a look at other partitions later.

## Interactive allocation

If, instead of requesting resources anew everytime you want to run a command on the compute nodes, you want to hold on to a specific set of resources and quickly dispatch a series of commands to run on them, you can use the `salloc` command in combination with `srun`. To do so, you specify the amount of resources you will need for your computations when calling `salloc`. `salloc` will request these resources from the resource manager and block until they are available. Then it will launch a new shell for you from which you can call `srun`, possibly multiple times, to dispatch commands onto the allocated resources.

In the previous section you took a task-centric approach to requesting resources by using the `-n` command line argument to `srun` to specify a number of tasks you want to run. This approach also works with `salloc` – in fact the way you specify resources is mostly the same between all different modes Slurm supports. However, since the number of CPU cores is always rounded up to the next multiple of the number of CPU cores in a single node, it might make sense to take a hardware centric approach to requesting resources. Using the `-N` command line argument, you can request a number of nodes from the resource manager (remember to specify `--gres gpu:4` for JUWELS Booster):

```
$ salloc -A training2230 --reservation hands-on-cluster-YYYYMMDD -N 1
salloc: Pending job allocation 3475519
salloc: job 3475519 queued and waiting for resources
salloc: job 3475519 has been allocated resources
salloc: Granted job allocation 3475519
salloc: Waiting for resource configuration
salloc: Nodes jwc00n014 are ready for job
$
```

At the new shell prompt, you can use `srun` to run commands without having to specify resources again:

```
$ srun hostname
jwc00n014.juwels
```

By default, Slurm assumes that your program is single-threaded, but still only launches one task per allocated node. This can be changed by specifying the CPUs per task with the `-c` argument.

```
$ srun -c 1 hostname
jwc00n014.juwels
[...]
jwc00n014.juwels
```



If you want to run several commands on a node without having to go through `srun` each time, you can use `srun` to launch a shell on the node:

```
$ srun --pty --cpu-bind=none /bin/bash
$ hostname
jwc00n014.juwels
$ exit
```

When using `srun` in one-shot mode, your account is charged for the time it takes to run the associated command. With `salloc` your account is charged for the duration of time you spend in the shell launched by `salloc` (and commands launched by that shell). Once you are done with the allocated resources, do not forget to exit from the shell:

```
$ exit
salloc: Relinquishing job allocation 3475519
salloc: Job allocation 3475519 has been revoked.
$ printenv SLURM_JOB_ID
$
```

If the `printenv SLURM_JOB_ID` prints a number, then you are still inside the allocation.

## Batch mode

If the system is relatively quiet and you are asking for a small amount of resources (or working on the `devel` partitions), `salloc` or one-shot `srun` should allow you to work with the system more or less interactively. Large production jobs on the other hand might have to wait an uncomfortably long time for resources and so running them interactively is not really convenient. Imagine you `salloc` a large number of nodes and while you wait you decide to go have lunch. If the allocation comes through while you are away you will still be charged for the resources even if they idle.

Also, if the systems were only used interactively, resource utilization would drop off in the late hours of the evening and ramp up in the mornings.

To enable better resource utilization and allow users to schedule jobs asynchronously, Slurm offers a batch mode through the `sbatch` command. It too requests resources from the resource manager, but unlike `salloc` which presents you with an interactive shell prompt from which you can call `srun`, `sbatch` runs commands from a shell script (the “job script”) without needing user intervention. The resources can be specified as command line arguments to `sbatch`, same as with `salloc` and `srun`, but can also be described in the job script. Open a new shell script in the editor:

```
($ module load nano)
$ nano testjob.sh
```

And enter the following script:

```
#!/bin/bash
#SBATCH --account=training2230
#SBATCH --reservation=hands-on-cluster-YYYYMMDD
#SBATCH --nodes=2
#SBATCH --cpus-per-task=1
#SBATCH --output=mpi-out.%j
#SBATCH --error=mpi-err.%j
#SBATCH --time=00:05:00
```

```
module load GCC ParaStationMPI
```

```
srun ./hellompi
```

Remember to specify `gpu:4 gres` for JUWELS Booster.

Then save the script and submit it for execution with:

```
$ sbatch testjob.sh
Submitted batch job 3476793
```

After the first line (the shebang line) the script contains specially formatted comments that act like arguments to `sbatch`. These arguments are written in their long form. Previously, you used the short form (e.g. `-N` is the same as `--nodes`). After the block of comments come regular shell commands. Inside the job script, we use the `module` command to make the software modules needed by the job programs available (here the compiler with its runtime libraries and an MPI library). The tasks are once again created using the `srun` command which works the same as before.

The job created by `sbatch` has to wait in a queue until the necessary resources become available. Use the `squeue` command to inspect the queue:

```
$ squeue -u $USER
```

| JOBID   | PARTITION | NAME     | USER     | ST | TIME | NODES | NODELIST(Reason) |
|---------|-----------|----------|----------|----|------|-------|------------------|
| 3476793 | batch     | testjob. | steinbus | PD | 0:00 | 2     | (Priority)       |

You might have to wait for a while, but eventually your job will be run. While your job is pending in the queue or already running you can execute another command to retrieve further information about your job:

```
$ scontrol show job <JOBID>
```

Once it is running, you will find two files next to the job script, `mpi-err.XXXXXXX` and `mpi-out.XXXXXXX` where `X` are decimal digits. These contain what was written to the standard error and output streams by your job. Should you need access to the hardware

during job execution check out the `sgoto --help` command to log into a compute node during job execution.

## Affinity and multi-threading

Computers today are typically equipped with multi-core CPUs which can work on multiple streams of instructions at the same time. The operating system is in charge of deciding which program gets to use which CPU core at a given point in time. Usually, it will let those programs which need access to resources run wherever resources are available, meaning one and the same program can end up using different CPU cores at different points in time. On a desktop machine this is not a problem. In fact it is a good thing, since we typically run far more programs than we have CPU cores available.

In an HPC setting things are different in that the workloads are adapted to use a number of processes or threads which matches the number of CPU cores (normally, you will have  $n\_processes \times n\_threads = n\_nodes \times n\_CPU\_cores\_per\_node$ ). If there is exactly one process or thread per CPU core, it would be wasteful to shuffle them around between different CPU cores. In order to avoid this shuffling, the resource manager assigns to the processes that it spawns an *affinity mask*. An affinity mask is a set of numbers identifying the CPU cores a process is allowed to use. By default, Slurm assumes that the processes you create are single threaded and gives each process access to a single CPU core. Allocate a node for playing around with this mechanism:

```
$ salloc -A training2230 --reservation hands-on-cluster-YYYYMMDD -N 1
salloc: Pending job allocation 3499694
salloc: job 3499694 queued and waiting for resources
salloc: job 3499694 has been allocated resources
salloc: Granted job allocation 3499694
salloc: Waiting for resource configuration
salloc: Nodes jwc00n001 are ready for job
```

Use the `numactl` command to inspect the affinity masks created by Slurm:

```
$ srun --label numactl --show
0: policy: default
0: preferred node: current
0: physcpubind: 0
0: cpubind: 0
0: nodebind: 0
0: membind: 0 1
```

The identifiers of accessible CPU cores are listed in `physcpubind`. Here, the single process that is created has access to a single CPU core, 0. Now, confirm that different processes will get access to different CPU cores:

```
$ srun --label -n 3 numactl --show
2: policy: default
2: preferred node: current
2: physcpubind: 1
2: cpubind: 0
2: nodebind: 0
2: membind: 0 1
1: policy: default
1: preferred node: current
1: physcpubind: 24
1: cpubind: 1
1: nodebind: 1
1: membind: 0 1
0: policy: default
0: preferred node: current
0: physcpubind: 0
0: cpubind: 0
0: nodebind: 0
0: membind: 0 1
```

The three processes get access to CPU cores 0, 1, and 24 respectively. If your processes are not single-threaded, you will have to give them access to more CPU cores (otherwise all threads will run on the same CPU core). This can be done using Slurm's `--cpus-per-task` parameter, or `-c`:

```
$ srun --label -c 24 numactl --show
1: policy: default
1: preferred node: current
1: physcpubind: 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47
1: cpubind: 1
1: nodebind: 1
1: membind: 0 1
0: policy: default
0: preferred node: current
0: physcpubind: 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23
0: cpubind: 0
0: nodebind: 0
0: membind: 0 1
2: policy: default
2: preferred node: current
2: physcpubind: 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71
2: cpubind: 0
2: nodebind: 0
2: membind: 0 1
3: policy: default
3: preferred node: current
3: physcpubind: 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95
3: cpubind: 1
3: nodebind: 1
3: membind: 0 1
```

Note how once you specify the number of CPU cores per task, Slurm switches its behavior from creating one process per node to filling the node with as many processes as possible. Each process gets access to 24 different CPU cores.

Copy the following small program into a file `hellohybrid.c`:

```
#include <stdio.h>
#include <mpi.h>
#include <omp.h>

int main(int argc, char* argv[]) {
    MPI_Init(&argc, &argv);

    int r, s;
    MPI_Comm_rank(MPI_COMM_WORLD, &r);
    MPI_Comm_size(MPI_COMM_WORLD, &s);
    #pragma omp parallel
    if (!omp_get_thread_num())
        printf(
            "hello from process %d of %d, using %d threads\n",
            r, s, omp_get_num_threads()
        );

    MPI_Finalize();
}
```

And compile it with:

```
$ mpicc -fopenmp -o hellohybrid hellohybrid.c
```

Now run the program:

```
$ srun ./hellohybrid
hello from process 0 of 1, using 1 threads
```

Again, using default settings, Slurm creates a single process and restricts it to a single CPU core. The **OpenMP** run time library supports shared-memory multiprocessing and allows to query the number of CPU cores accessible to the process. It creates just as many threads (here only one). If you specify a number of CPU cores per process this changes:

```
$ srun -c 24 ./hellohybrid
hello from process 2 of 4, using 24 threads
hello from process 0 of 4, using 24 threads
hello from process 3 of 4, using 24 threads
hello from process 1 of 4, using 24 threads
```

Once more, Slurm fills the node with four processes having appropriate affinity masks. The OpenMP run time figures out that each process is allowed to use 24 CPU cores and creates a team of threads to fill those CPU cores.

**IMPORTANT:** Do not forget to exit your salloc session at this point.

## JSC Affinity Tools

Since we are using psslurm we have implemented a few options different than the default in Slurm. For this reason we are offering two tools that can help you to understand the process affinity on our systems:

1. The command line executable: `psslurmgetbind`
2. An online [pinning tool](#)

Further information can be found in the “Processor Affinity” chapter of the corresponding [System Documentation](#).

## Further reading

Our online documentation has more information on working with the resource manager. It has detailed lists with the hardware available in various partitions as well as job limits. Also, it discusses advanced topics like multiple job steps, dependency chains and heterogeneous jobs. If you want more details, you can find the documentation for our various systems here:

- [JUWELS documentation: Batch system](#)
- [JURECA documentation: Batch system](#)
- [JUSUF documentation: Batch system](#)

You can also have a look at [the official Slurm documentation](#).

## LLview - Detailed Job Reporting

LLview is an excellent tool that provides an overview of currently running and finished jobs, including detailed job reports plus obscure error messages that are hard to find for users. Your jobs crash and you do not know why? This is the first place to check. There is a website for each of our large systems. You can find the link for every system at the lower left corner of the [documentation webpage of LLview](#). To begin check out your system of interest, ideally one you have run jobs on.

## Currently Active Jobs

When opening LLview by default it will first show you the list of your currently active jobs, either pending or running. If you are the Principal Investigator (PI) or Project

Administrator (PA) for a project and you are in the project view you will see all active jobs from your project. Project mentors also have access to this view.

You filter the list based on the filters below any of the column headings. Clicking on any column heading will cause the jobs to be sorted in ascending or descending order, an arrow will appear next to the column title indicating either ascending (upwards-pointing arrow) or descending (downwards-pointing arrow). In the case of a sort conflict the submission time of a job is used to resolve the conflict, with jobs that were submitted more recently appearing above jobs that started earlier, if sorted ascending. By default jobs are sorted according to ascending job start submit time with your most recently submitted jobs at the top.

Values in red in the list indicate that something may be wrong. For example the average load on a node may be high. Note that this is reported in fractions of the utilization of a single core. A 1.0 therefore means that a single core was fully utilized. This value should ideally be as close as possible to the number of cores a node has.

Important to note is also the state of the job, which you can find on the right. The job can either be pending, if it has been submitted, but is not yet running, or running. Sometimes it can also be completed (CMPL) or error if the job has finished successfully or errored out respectively. These jobs will then be shortly removed from this list.

Clicking on any of the jobs will cause the graphs to be populated at the bottom of the page if the job ran for more than a couple of minutes, depends on how often LLview is able to query the system state. The “Load on Node” gives you an idea of the evolution of the load placed on a node over time. The other two graphs show the evolution of I/O bandwidth and number of I/O operations per second for various storage tiers, you can change the storage tier by clicking on one of them in the bottom right (HOME, PROJECT, SCRATCH, FASTDATA).

The two right-most columns can contain small pictograms, one of a chart and the other is the Adobe PDF icon. If they are available it means that job reports are available for the jobs. The chart pictogram takes you to a webpage-based interactive report. The PDF icon downloads a non-interactive PDF report. These contain detailed information on the job as well as error codes. They with a textual header listing important information regarding the job. Then, if the job ran long enough graphs of various metrics, like CPU and GPU usage, are presented. Finally a list of nodes and error messages is included at the end. The error messages are especially important as these can be Slurm diagnostic error messages that might be difficult for users to find. For more information see the [detailed-report documentation](#).

## Jobs Ended Today

This is the same view as the currently active jobs view but for jobs that have finished in the last 24 hours.

## Jobs < 3 weeks

This is the same view as the currently active jobs view but for jobs that have finished in the last three weeks.

## Live

Here you can see a live view of the system and how jobs are distributed across the supercomputer on a rack level on the left. At the bottom on the left is a scheduling prediction which shows when which jobs are expected to be scheduled. Large jobs also have their names displayed.

On the right you can see the color-coded queue of all jobs, including currently running and finalizing jobs. Clicking on a column title sorts the list either ascending (upwards-pointing arrow) or descending (downwards-pointing arrow).

If this tab has a drop down then you can see different queues, for example for the JUWELS cluster you will be able to see the batch and GPU queues. Note that only the main queues are shown. Queues like the devel queue for development are not shown. To see these queues you will have to login to the systems and use `squeue`.

## Queue Tab

In the queue tab you can see the queue for all partitions you can submit to. Unlike (Live) the results are not filtered according to the partition and the queue is displayed in the typical LLview fashion.

## Further reading

[Here](#) you can find the documentation for LLview.

## USING GPUS

All systems at JSC have nodes which are accelerated by General Purpose Graphics Processing Units (GPGPUs or just GPUs). Since the GPUs are all made by NVIDIA, using them is accomplished through their **CUDA SDK**. CUDA is available as a module:

```
$ module load CUDA
```



This example is executed on the JUWELS booster. To demonstrate how to compile and run a program that uses GPUs, we will use one of the examples included in CUDA and load additionally the compiler NVHPC and the MPI implementation ParaStationMPI. The samples directory of the CUDA installation has a plethora of exemplaric codes you can play and learn with:

```
$ module load NVHPC ParaStationMPI
$ cp -r $EBROOTCUDA/samples $PROJECT_training2230/$USER
$ cd $PROJECT_training2230/$USER/samples/0_Simple/simpleMPI
$ make
/p/software/juwelsbooster/stages/2022/software/psmpi/5.5.0-1-NVHPC-22.1/bin/mpicxx -
I../common/inc -o simpleMPI_mpi.o -c simpleMPI.cpp
/p/software/juwelsbooster/stages/2022/software/CUDA/11.5/bin/nvcc -ccbin g++ -
I../common/inc -m64 --threads 0 --std=c++11 -gencode arch=compute_35,code=sm_35 -
gencode arch=compute_37,code=sm_37 -gencode arch=compute_50,code=sm_50 -gencode
arch=compute_52,code=sm_52 -gencode arch=compute_60,code=sm_60 -gencode
arch=compute_61,code=sm_61 -gencode arch=compute_70,code=sm_70 -gencode
arch=compute_75,code=sm_75 -gencode arch=compute_80,code=sm_80 -gencode
arch=compute_86,code=sm_86 -gencode arch=compute_86,code=compute_86 -o simpleMPI.o -c
simpleMPI.cu
nvcc warning : The 'compute_35', 'compute_37', 'compute_50', 'sm_35', 'sm_37' and
'sm_50' architectures are deprecated, and may be removed in a future release (Use -Wno-
deprecated-gpu-targets to suppress warning).
/p/software/juwelsbooster/stages/2022/software/psmpi/5.5.0-1-NVHPC-22.1/bin/mpicxx -o
simpleMPI simpleMPI_mpi.o simpleMPI.o -
L/p/software/juwelsbooster/stages/2022/software/CUDA/11.5/lib64 -lcudart
mkdir -p ../bin/x86_64/linux/release
cp simpleMPI ../bin/x86_64/linux/release
```

This sample shows how to compile a combination of C++, MPI and CUDA code. There should now be an executable called `simpleMPI` inside the `simpleMPI` directory. To run the program, use `srun` like before:

```
$ srun -A training2230 -p <gpu partition> --gres gpu:4 -N 1 -n 4 ./simpleMPI
[...]
Running on 4 nodes
Average of square roots is: 0.667305
PASSED
```

**Note:** In this output *nodes* are meaning *MPI tasks*. The developers seem to have assumed implicitly that only one GPU with one MPI task is located on one node when executing this software.

You have to specify a partition that contains nodes equipped with GPUs, `-p develgpus` for JUWELS and JUSUF, `-p dc-gpu-devel` for JURECA, or `-p develbooster` for JUWELS Booster, and you have to specify how many GPUs you want those nodes to have, `--gres gpu:4` (or `--gres gpu:1` on JUSUF).

## GPU Inspection During Execution

First of all, LLview is already an elaborated tool to monitor your jobs and extracts most of the data being relevant for many use cases of monitoring with small effort from the user side. Nevertheless, logging into the compute nodes during job execution is easy and comfortable and in some cases needed.

In the following bash session on the JUWELS booster a job is initiated through `srun` in the background of this login node session (& at the end of the command). Just hit enter after this line to retrieve the normal command line. The job is waiting 600 seconds or 10 minutes. After logging into the compute node, through `sgoto`, we show the usage of the GPUs with `nvidia-smi`, which can be exchanged with anything you would like to do on the compute node during job execution. Afterwards we log out from the compute node, put the executed `srun` command from the background to the foreground with `fg` and cancel this execution by hitting `CTRL-C` a couple of times until the normal command line is available.

```
$ srun -N 1 -n 1 -t 00:10:00 -A training2230 -p develbooster --gres=gpu:4 sleep 600 &
[1] 25114
```

```
srun: job 5535332 queued and waiting for resources
```

```
srun: job 5535332 has been allocated resources
```

```
$ sgoto 5535332 0
```

```
$ nvidia-smi
```

Thu May 12 08:49:34 2022

|                                                                             |        |              |            |                                 |  |                        |  |         |  |
|-----------------------------------------------------------------------------|--------|--------------|------------|---------------------------------|--|------------------------|--|---------|--|
| NVIDIA-SMI 510.47.03      Driver Version: 510.47.03      CUDA Version: 11.6 |        |              |            |                                 |  |                        |  |         |  |
| GPU Name                   Persistence-M                                    |        |              |            | Bus-Id                   Disp.A |  | Volatile Uncorr. ECC   |  |         |  |
| Fan    Temp    Perf    Pwr:Usage/Cap                                        |        | Memory-Usage |            |                                 |  | GPU-Util    Compute M. |  | MIG M.  |  |
| =====                                                                       |        |              |            |                                 |  |                        |  |         |  |
| 0                                                                           | NVIDIA | A100-SXM...  | On         | 000000000:03:00.0 Off           |  | 0                      |  |         |  |
| N/A                                                                         | 44C    | P0           | 55W / 400W | 0MiB / 40960MiB                 |  | 0%                     |  | Default |  |
|                                                                             |        |              |            |                                 |  | Disabled               |  |         |  |
| -----                                                                       |        |              |            |                                 |  |                        |  |         |  |
| 1                                                                           | NVIDIA | A100-SXM...  | On         | 000000000:44:00.0 Off           |  | 0                      |  |         |  |
| N/A                                                                         | 44C    | P0           | 54W / 400W | 0MiB / 40960MiB                 |  | 0%                     |  | Default |  |
|                                                                             |        |              |            |                                 |  | Disabled               |  |         |  |
| -----                                                                       |        |              |            |                                 |  |                        |  |         |  |
| 2                                                                           | NVIDIA | A100-SXM...  | On         | 000000000:84:00.0 Off           |  | 0                      |  |         |  |
| N/A                                                                         | 45C    | P0           | 58W / 400W | 0MiB / 40960MiB                 |  | 0%                     |  | Default |  |
|                                                                             |        |              |            |                                 |  | Disabled               |  |         |  |
| -----                                                                       |        |              |            |                                 |  |                        |  |         |  |
| 3                                                                           | NVIDIA | A100-SXM...  | On         | 000000000:C4:00.0 Off           |  | 0                      |  |         |  |
| N/A                                                                         | 44C    | P0           | 58W / 400W | 0MiB / 40960MiB                 |  | 0%                     |  | Default |  |
|                                                                             |        |              |            |                                 |  | Disabled               |  |         |  |
| -----                                                                       |        |              |            |                                 |  |                        |  |         |  |

```

+-----+
| Processes: |
| GPU  GI  CI          PID   Type   Process name                      GPU Memory |
|      ID  ID                                   Usage                      |
+=====+
| No running processes found |
+-----+

$ exit
logout
$ fg
srun -N 1 -n 1 -t 00:10:00 -A training2230 -p develbooster --gres=gpu:4 sleep 500
^Csrun: sending Ctrl-C to StepId=5535332.0
srun: forcing job termination
srun: Job step aborted: Waiting up to 6 seconds for job step to finish.

```

`sgoto` takes the job id as first argument and the node number within the job as second argument where the counting starts with 0. `nvidia-smi` prints some useful information about available GPUs on a node, like temperature, memory usage, currently running processes and power consumption.

## GPU Affinity

On systems with more than one GPU per node, a choice presents itself of which GPU should be visible to which application task. This is controlled through the environment variable `CUDA_VISIBLE_DEVICES`, which can be set to a comma separated list of integers identifying devices to be visible to a task. You can manually define this variable before running your tasks with `srun` if the pinning is going to be the same for every task.

Let us investigate further on this with a practical example. First, we prepare a device query example.

```

$ cd $PROJECT_training2230/$USER/samples/1_Uutilities/deviceQueryDrv
make
/p/software/juwelsbooster/stages/2022/software/CUDA/11.5/bin/nvcc -ccbin g++ -
I../common/inc -m64 --threads 0 --std=c++11 -gencode
arch=compute_35,code=compute_35 -o deviceQueryDrv.o -c deviceQueryDrv.cpp
nvcc warning : The 'compute_35', 'compute_37', 'compute_50', 'sm_35', 'sm_37' and
'sm_50' architectures are deprecated, and may be removed in a future release (Use -Wno-
deprecated-gpu-targets to suppress warning).
/p/software/juwelsbooster/stages/2022/software/CUDA/11.5/bin/nvcc -ccbin g++ -m64
-gencode arch=compute_35,code=compute_35 -o deviceQueryDrv deviceQueryDrv.o -
L/p/software/juwelsbooster/stages/2022/software/CUDA/11.5/lib64/stubs -lcuda
nvcc warning : The 'compute_35', 'compute_37', 'compute_50', 'sm_35', 'sm_37' and
'sm_50' architectures are deprecated, and may be removed in a future release (Use -Wno-
deprecated-gpu-targets to suppress warning).
mkdir -p ../bin/x86_64/linux/release
cp deviceQueryDrv ../bin/x86_64/linux/release

```

This will create the executable `deviceQueryDrv`. During the execution of `deviceQueryDrv` all visible CUDA devices are queried. The following sbatch script `gpuAffinityTest.sbatch` written for the JUWELS Booster executes the assisting bash script `gpuAffinityTest.bash` which in turn executes `deviceQueryDrv`.

```
#!/bin/bash
#SBATCH --ntasks=<number of MPI tasks>
#SBATCH --nodes=1
#SBATCH --time=00:01:00
#SBATCH --partition=develbooster
#SBATCH --gres=gpu:4
#SBATCH -A training2230
```

```
module load CUDA NVHPC ParaStationMPI
```

```
srun bash gpuAffinityTest.bash
```

The in parallel executed helper script `gpuAffinityTest.bash` will be needed to print the environment variable `CUDA_VISIBLE_DEVICES` for every MPI task initiated.

```
#!/bin/bash

#export CUDA_VISIBLE_DEVICES=<comma-separated list of visible gpus>
echo "MPI task" $SLURM_PROCID "with CUDA_VISIBLE_DEVICES =" $CUDA_VISIBLE_DEVICES

./deviceQueryDrv
```

The environment variable `SLURM_PROCID` contains the current MPI task ID. The definition of the environment variable `CUDA_VISIBLE_DEVICES` is not yet manually done. By uncommenting the commented line within `gpuAffinityTest.bash`, `CUDA_VISIBLE_DEVICES` can be defined manually for every task.

Execute this example for `ntasks=1` and study the output file.

```
MPI task 0 with CUDA_VISIBLE_DEVICES = 0,1,2,3
./deviceQueryDrv Starting...
```

```
CUDA Device Query (Driver API) statically linked version
Detected 4 CUDA Capable device(s)
```

```
Device 0: "NVIDIA A100-SXM4-40GB"
[...]
Device PCI Domain ID / Bus ID / location ID:    0 / 3 / 0
[...]
Device 1: "NVIDIA A100-SXM4-40GB"
[...]
Device PCI Domain ID / Bus ID / location ID:    0 / 68 / 0
[...]
Device 2: "NVIDIA A100-SXM4-40GB"
```

```
[...]
Device PCI Domain ID / Bus ID / location ID:    0 / 132 / 0
[...]
Device 3: "NVIDIA A100-SXM4-40GB"
[...]
Device PCI Domain ID / Bus ID / location ID:    0 / 196 / 0
[...]
> Peer-to-Peer (P2P) access from NVIDIA A100-SXM4-40GB (GPU0) -> NVIDIA A100-SXM4-40GB
(GPU1) : Yes
> Peer-to-Peer (P2P) access from NVIDIA A100-SXM4-40GB (GPU0) -> NVIDIA A100-SXM4-40GB
(GPU2) : Yes
[...]
Result = PASS
```

The value for `CUDA_VISIBLE_DEVICES` at the beginning and the different Bus IDs, representing the 4 GPUs, are of importance. For this single MPI task all 4 GPUs are visible. At the end of the file you can also see the successful interconnectivity tests of the GPUs.

If the environment variable `CUDA_VISIBLE_DEVICES` is not defined by you, `srun` will provide a default:

- for jobs with a single task (`-n 1`) all devices will be visible  
`CUDA_VISIBLE_DEVICES=0,1,2,3`
- for all other jobs, only a single device will be visible per task, with the same device being visible to multiple tasks if there are more tasks than GPUs

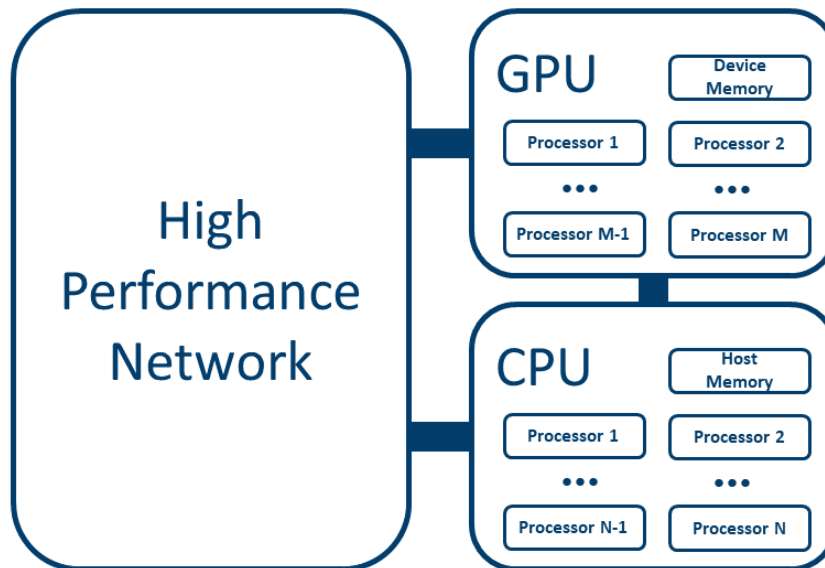
By playing a little with the number of tasks within the scripts stated above you can study the behaviour of the GPU pinning and confirm the default. If this default is not suited to your needs you can uncomment the line

```
export CUDA_VISIBLE_DEVICES=<comma-separated list of visible gpus>
```

and define `CUDA_VISIBLE_DEVICES` as you wish.

## Network Architecture Study

The JUWELS Booster delivers a network infrastructure accelerating direct data exchange between the GPUs. These GPUs have internal hardware to store data and are directly connected to the high-performance network.



Having this in mind, the traditional data exchange between two GPUs, with an intermediate hop of data on host memory, will lead to reduced performance. **CUDA-awareness** of an MPI implementation is a vital part to increase data exchange performance between GPUs. At the supercomputing infrastructure from JSC there are implementations for CUDA-aware MPI, like ParaStationMPI and OpenMPI, preinstalled. CUDA-awareness enables to pass a pointer to data on the GPU directly to an MPI-directive.

The following example `mpiBroadcasting.cpp` performs three different measurements for data exchange by use of `MPI_Bcast`. `MPI_Bcast` broadcasts data from one MPI process to other MPI processes. In the source code below, at first, data between host memories is exchanged. Secondly, data between GPUs is exchanged by hopping intermediately onto the host memory. At last, data between GPUs is exchanged by use of the direct network connection between the GPUs.

```
#include <stdio.h>
#include <string.h>
#include <mpi.h>
#include <time.h>
#include "cuda_runtime.h"

int main(int argc, char *argv[])
{
    clock_t start, end; // Time stamps
    double cpu_time_used;
    int myrank;
    int N=1000000000; // # elements to broadcast per repetition
```

```

//# broadcasting several repetitions since maximal size of
// elements to send is restricted through MPI
int Nbroadcast=20;

double *x = new double[N]; // Allocate space on host memory
double *d_x; // Array on device
cudaMalloc(&d_x, N*sizeof(double)); // Allocate space on device

for (int i=0;i<N;i++) x[i] = 1.0f; // prefilling data into allocated memory
// send data into device
cudaMemcpy(d_x, x, N*sizeof(double), cudaMemcpyHostToDevice);

MPI_Init(&argc, &argv);
MPI_Comm_rank(MPI_COMM_WORLD, &myrank);

// Initial unmeasured broadcasting due to
// setup offsets in initializing connections
for(int i=0;i<Nbroadcast;i++) {
    MPI_Bcast(x, N, MPI_DOUBLE, 0, MPI_COMM_WORLD);
    MPI_Barrier(MPI_COMM_WORLD);
    MPI_Bcast(d_x, N, MPI_DOUBLE, 0, MPI_COMM_WORLD);
    MPI_Barrier(MPI_COMM_WORLD);
}

// host to host memory measurement
start = clock(); // set the start time
for(int i=0;i<Nbroadcast;i++) {
    MPI_Bcast(x, N, MPI_DOUBLE, 0, MPI_COMM_WORLD);
    MPI_Barrier(MPI_COMM_WORLD);
}
end = clock(); // set the end time
// compute cpu time elapsed during the broadcasting
cpu_time_used = ((double) (end - start)) / CLOCKS_PER_SEC;
if (myrank == 0) printf("Broadcasting to all host memories \
    took %f seconds. \n", cpu_time_used);

// device to device with intermediate copy to/from host
start = clock();
for(int i=0;i<Nbroadcast;i++) {
    cudaMemcpy(x, d_x, N*sizeof(double), cudaMemcpyDeviceToHost);
    MPI_Bcast(x, N, MPI_DOUBLE, 0, MPI_COMM_WORLD);
    cudaMemcpy(d_x, x, N*sizeof(double), cudaMemcpyHostToDevice);
    MPI_Barrier(MPI_COMM_WORLD);
}
end = clock();
cpu_time_used = ((double) (end - start)) / CLOCKS_PER_SEC;
if (myrank == 0) printf("Broadcasting to all GPUs took %f seconds \
    with intermediate copy to host memory. \n", cpu_time_used);

// device to device through direct network connection of the GPUs
start = clock();
for(int i=0;i<Nbroadcast;i++) {

```

```

    MPI_Bcast(d_x, N, MPI_DOUBLE, 0, MPI_COMM_WORLD);
    MPI_Barrier(MPI_COMM_WORLD);
}
end = clock();
cpu_time_used = ((double) (end - start)) / CLOCKS_PER_SEC;
if (myrank == 0) printf("Broadcasting to all GPUs took %f \
seconds. \n", cpu_time_used);

// Release allocated memory space on host and device
cudaFree(d_x);
delete x;

MPI_Finalize();
return 0;
}

```

The initial broadcasts are needed to let the network establish connections between the MPI tasks. Some implementations of MPI are setting up network connections between MPI tasks only at first data exchange. This is an offset which is not planned to be measured here. `MPI_Barrier` directs all MPI tasks to wait until all data was broadcasted. As a result there are three times measured and printed. This example is executed on 2 nodes with 4 tasks on every node, where each task occupies one GPU.

*Important:* Note that we switch the compiler at this stage, when you compare to the previous instructions of this chapter. It is worth it mentioning that you should use the same modules for compilation which you are planning to use for the execution.

```

$ module load NVHPC CUDA OpenMPI
$ mpicxx -O0 -I$CUDA_HOME/include -L$CUDA_HOME/lib64 -lcudart -lcuda mpiBroadcasting.cpp
$ srun -N 2 -n 8 -t 01:00:00 -A training2230 -p booster --gres=gpu:4 ./a.out
Broadcasting to all host memories took 4.526835 seconds.
Broadcasting to all GPUs took 7.481972 seconds with intermediate copy to host memory.
Broadcasting to all GPUs took 2.625439 seconds.

```

The parameter `-O0` deactivates any optimizations performed by the compiler, which is needed since a powerful compiler could know at compile time that the same data is initialized for all tasks and then sent around. This could lead to a deletion of the MPI directives at compile time leading to extremely small but erroneous time measurements. The data exchange directly from one GPU to another GPU is the fastest. Furthermore, the CPUs on the JUWELS Booster nodes have a relatively small compute performance, to avoid too much overhead and unnecessary power consumption. These nodes are designed such that as much workload and data exchange as possible should be performed by the GPUs.

You can study the source code and play around with this setup. This will give you valuable insights on how to develop your own software for execution on the JUWELS Booster.



## Further reading

Our online documentation has more information on software modules. It lists the basic tool chains (compiler + communication library + math library) available on our systems and discusses using older software stages. If you want more details, you can find the documentation for our various systems here:

- [JUWELS documentation: GPU Computing](#)
- [JURECA documentation: GPU Computing](#)
- [JUSUF documentation: GPU Computing](#)

The [CUDA SDK](#) documentation gives you detailed information about how to develop CUDA code. There are also excellent articles in the web for learning CUDA like [An Even Easier Introduction to CUDA](#) or [An Introduction to CUDA-Aware MPI](#).

The JSC regularly offers [CUDA courses for HPC](#) being an ideal starting point to get into the topic.

## USEFUL LINKS

In this chapter, you can find a useful collection of links to get more information about several topics. The slides of workshops held at JSC (also from this introductory workshop) you can find [here](#). For the future, the regularly updated master of this working document you can find [here](#). Keep in mind that some template words are left within the document on purpose. You will need to adjust them for your testcase.

## System Documentation

JSC offers documentation for the production systems:

- [JUWELS](#)
- [JURECA](#)
- [JUSUF](#)

## JSC Services

- [JSC Service Status](#)
- [JuDoor](#)
- [Jupyter Lab](#)
- [HDF Cloud](#)

## Job Reporting

The Job Reporting service gives you access to PDF reports which contain certain performance metrics that the system automatically collects about your jobs. It also includes an overview over the system utilization and queue. You can access the Job Reporting service for the different systems here:

- [JUWELS](#)
- [JURECA](#)
- [JUSUF](#)

## Apply for Computing Time

The JSC web site describes [how to apply for computing time](#).

## Apply for a Data Project

The JSC web site describes [how to apply for a data project](#).

## JSC Course Programme

JSC offers many courses throughout the year covering topics such as parallel programming, machine learning, and visualization. Please have a look at the [course programme on the JSC web site](#).

## Supercomputing Support

Our high-level support team supports the users in case of problems on our systems, e.g. porting of the application, parallelisation and performance issues as well as usage of the HPC system. So if you are having a question, you cannot sort out by yourself, by working through this document or by having a look into the documentation, just drop a mail to [sc@fz-juelich.de](mailto:sc@fz-juelich.de).