

---

# Introduction to Linux for Phoenix Users

Steve Pederson  
Stephen Bent  
Dan Kortschak  
Ramona Rogers  
Robert Qiao  
Bowen Chen  
Exequiel Sepulveda



---

# Licensing

This work is licensed under a Creative Commons Attribution 3.0 Unported License and the below text is a summary of the main terms of the full Legal Code (the full licence) available at <http://creativecommons.org/licenses/by/3.0/legalcode>.

## You are free:

- to copy, distribute, display, and perform the work
- to make derivative works
- to make commercial use of the work

## Under the following conditions:

**Attribution** - You must give the original author credit.

## With the understanding that:

**Waiver** - Any of the above conditions can be waived if you get permission from the copyright holder.

**Public Domain** - Where the work or any of its elements is in the public domain under applicable law, that status is in no way affected by the license.

**Other Rights** - In no way are any of the following rights affected by the license:

- Your fair dealing or fair use rights, or other applicable copyright exceptions and limitations;
- The author's moral rights;
- Rights other persons may have either in the work itself or in how the work is used, such as publicity or privacy rights.

**Notice** - For any reuse or distribution, you must make clear to others the licence terms of this work.



---

# Contents

<b>Licensing</b>	<b>3</b>
<b>Contents</b>	<b>4</b>
<b>Workshop Information</b>	<b>7</b>
The Trainers . . . . .	8
Welcome . . . . .	9
Course Summary . . . . .	9
Using the Post-it Notes . . . . .	10
Providing Feedback . . . . .	10
Document Structure . . . . .	11
Computer Setup . . . . .	12
<b>Introduction to Linux Operating System</b>	<b>13</b>
Components of the Linux Operating System . . . . .	14
<b>Understanding file system on Phoenix</b>	<b>19</b>
Home directory . . . . .	20
Fast directory . . . . .	20
Mounting filesystems . . . . .	21
File and folder permissions . . . . .	21
Where are the applications on Phoenix? . . . . .	24
<b>Regular Expressions</b>	<b>25</b>
Introduction . . . . .	26
The command <code>grep</code> . . . . .	26
Pattern Searching . . . . .	27
<b>The Tools <code>sed</code> &amp; <code>awk</code></b>	<b>29</b>
<code>sed</code> : The Stream Editor . . . . .	30
Some Important Programming Concepts . . . . .	32
<code>awk</code> : A command and a language . . . . .	34
<b>Command Line Editors</b>	<b>37</b>
<code>vi</code> . . . . .	38
<code>emacs</code> . . . . .	39
<code>Nano</code> . . . . .	39
<b>Writing Scripts</b>	<b>45</b>
Shell Scripts . . . . .	46

---

Moving towards High Performance Computing . . . . .	51
<b>Using Modules</b>	<b>53</b>
Basic commands of modules . . . . .	54
Loading modules . . . . .	56
Using modules . . . . .	56
A complete exercise of modules . . . . .	57
<b>Writing Basic Slurm Scripts</b>	<b>59</b>
Basic template of a Slurm script . . . . .	60
Submitting Scripts . . . . .	61
Managing Scripts . . . . .	62
<b>Space for Personal Notes or Feedback</b>	<b>65</b>



---

## Workshop Information

## The Trainers

**Ms Ramona Rogers**

Computing Officer

IT Support

The University of Adelaide

South Australia

[ramona.rogers@adelaide.edu.au](mailto:ramona.rogers@adelaide.edu.au)

**Mr Robert Qiao**

Phoenix Support Team

The University of Adelaide

South Australia

[robert.qiao@adelaide.edu.au](mailto:robert.qiao@adelaide.edu.au)

**Mr Bowen Chen**

Phoenix Support Team

The University of Adelaide

South Australia

[bowen.chen@adelaide.edu.au](mailto:bowen.chen@adelaide.edu.au)

**Mr Exequiel Sepúlveda**

Phoenix Support Team

The University of Adelaide

South Australia

[exequiel.sepulvedaescobedo@adelaide.edu.au](mailto:exequiel.sepulvedaescobedo@adelaide.edu.au)



## Welcome

Thank you for your attendance & welcome to the Introduction to Linux for Phoenix Users Workshop. This is an offering by the University of Adelaide, Reserach Computing team which is a centrally funded initiative, with the aim of assisting & enabling researchers in their work.

This workshop is largaly based on the "Introduction to Linux" workshop delivered by The Bioinformatics Hub. The Hub has a web-page at <http://www.adelaide.edu.au/bioinformatics-hub/>, and **to be kept up to date on upcoming events and workshops, please join the internal Bioinformatics mailing list on <http://list.adelaide.edu.au/mailman/listinfo/bioinfo>.**

Phoenix has two main sources of information: <http://www.adelaide.edu.au/phoenix/> and or wiki page: [https://wiki.adelaide.edu.au/hpc/index.php/Main\\_Page](https://wiki.adelaide.edu.au/hpc/index.php/Main_Page).

Phoenix team has an active Slack team for discussing questions with the local community. Slack teams do require an invitation to join, so please email us the Hub on [hpcsupport@adelaide.edu.au](mailto:hpcsupport@adelaide.edu.au) to join the community. All are welcome.

Today's workshop has been put together based on previous material and courses prepared by Dr Stephen Bent (*University of Queensland*), with generous technical support & advice provided by Dr Nathan Watson-Haigh (*ACPFPG*) and Dr Dan Kortschak (*Adelaide University, Adelson Research Group*) and complemented by Phoenix team: Ramona Rogers, Robert Qiao, Bowen Chen and Exequiel Sepulveda.

We hope it will be useful in enabling you to continue and to advance your research.

## Course Summary

In today's workshop, the morning session will be spent introducing you to the basic tools and concepts required for data handling. In the afternoon session we'll develop these skills to a more advanced level, with progress in both sessions being made at your own pace. Some people may finish early today, but the majority of you probably won't.

You will use Phoenix all the time.

The majority of data handling and analysis required for research uses the *command line*, alternatively known as the terminal or the *bash shell*. This is a text-based interface in which commands must be typed, as opposed to the Graphical User Interfaces (aka GUIs) that most of us have become accustomed to. Being able to access your computer using these tools enables you to more fully utilise the power & capabilities of your machine, for both Linux & Mac operating systems, and to a lesser extent will even enable you to dig



deeper on a Windows system.

Whilst some of the tools we cover today may appear trivial, they are used on a daily basis by those working in the field. These basic tools are essential for writing what are known as *shell scripts*, which we will begin to cover in the afternoon session. These are essentially simple programs that utilise the inbuilt functions of the shell, and are used to automate processes such as de-multiplexing read libraries, or aligning reads to the genome. A knowledge of this simple type of programming and navigation is also essential for accessing the high-performance computing resources such as **phoenix**.

## Using the Post-it Notes

For today's session, you will be provided with 3 post-it notes of differing colours. Please use these to signal whether you need help or not by placing them on your monitors. We will interpret these as:

1. **Red** - Help! I can't make something work
2. **Yellow** - I'm working on something, but haven't made it yet
3. **Green** - I've finished the task I was working on

## Providing Feedback

While we endeavour to deliver a workshop with quality content and documentation in a venue conducive to an exciting, well run hands-on workshop with a bunch of knowledgeable and likable trainers, we know there are things we could do better.

Whilst we want to know what didn't quite hit the mark for you, what would be most helpful and least depressing, would be for you to provide ways to improve the workshop. i.e. constructive feedback. After all, if we knew something wasn't going to work, we wouldn't have done it or put it into the workshop in the first place!

Clearly, we also want to know what we did well! This gives us that "feel good" factor which will see us through those long days and nights in the lead up to such hands-on workshops!

## Document Structure

We have provided you with an electronic copy of the workshop's hands-on tutorial documents. We have done this for two reasons: 1) you will have something to take away with you at the end of the workshop, and 2) you can save time (mis)typing commands on the command line by using copy-and-paste.

*We advise you to use Acrobat Reader to view the PDF. This is because it properly supports some features we have implemented to ensure that copy-and-paste of commands works as expected. This includes the appropriate copy-and-paste of special characters like tilde and hyphens as well as skipping line numbers for easy copy-and-paste of whole code blocks.*



While you could fly through the hands-on sessions doing copy-and-paste, you will learn more if you use the time saved from not having to type all those commands, to understand what each command is doing!

The commands to enter at a terminal look something like this:

```
1 tophat --solexa-quals -g 2 --library-type fr-unstranded -j \  
  annotation/Danio_rerio.Zv9.66.spliceSites -o tophat/ZV9_2cells \  
  genome/ZV9 data/2cells_1.fastq data/2cells_2.fastq
```

The following styled code is not to be entered at a terminal, it is simply to show you the syntax of the command. You must use your own judgement to substitute in the correct arguments, options, filenames etc

```
| tophat [options]* <index_base> <reads_1> <reads_2>
```

The following icons are used in the margin, throughout the documentation to help you navigate around the document more easily:



Important



For reference



Follow these steps



Questions to answer



Warning - STOP and read



Bonus exercise for fast learners



Advanced exercise for super-fast learners

## Computer Setup



We will all be working on **Phoenix** directly, which is the University of Adelaide's High Performance Computing (HPC) system. The software client **Bitvise** or **putty** which you will have already installed, enables us to access these machines in a familiar Desktop style, even though the majority of our time will be spent within the terminal.



In case you need to install Bitvise, download the installer from here: [Bitvise SSH Client](#)

You will need to open a SSH session to **phoenix**. First, we need to create a session with the basic parameters

1. Hostname phoenix.adelaide.edu.au
2. Username your student or staff ID

Now we have created the session, you will be asked for your password.



Now that you are connected, you will notice we are now in the head node of **Phoenix**. Welcome to **Phoenix**.

---

# Introduction to Linux Operating System

Primary Author(s):

Ramona Rogers, Research Services, University of Adelaide

[ramona.rogers@adelaide.edu.au](mailto:ramona.rogers@adelaide.edu.au)

Contributor(s):

Exequiel Sepulveda, Research Services, University of Adelaide

[exequiel.sepulvedaescobedo@adelaide.edu.au](mailto:exequiel.sepulvedaescobedo@adelaide.edu.au)

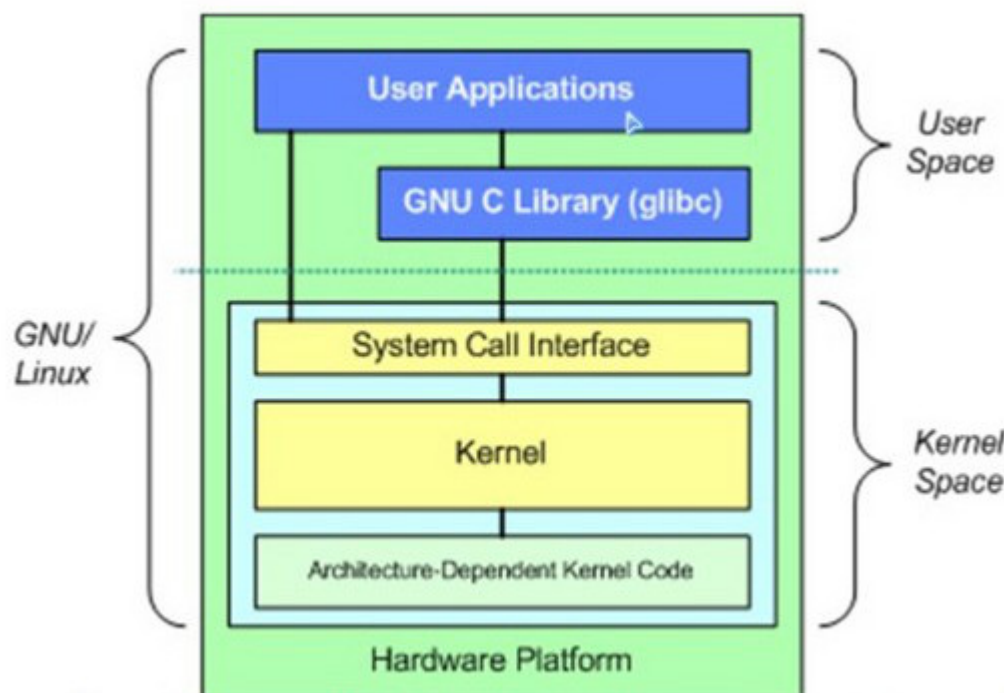
Linux has become the most used operating system for servers, and for High Performance Computing (HPC) in particular. That is the motivation to bring you this workshop. The main goal is to show what a Linux OS is and how to use it, focus on Phoenix, the HPC facilities for supercomputing at The University of Adelaide.

## Components of the Linux Operating System

An operating system is a dedicated computer program whose primary purpose is to execute other programs, to manage the hardware and system resources for those programs. The operating system sits on top of the hardware and it is on what the general programs run. You probably are familiarised with Windows operating system, but there are many others, among the most popular are: Linux, macOS and Android.

Phoenix uses Linux as operating system.

The Linux operating system can be depicted as a series of layers as hardware, device-drivers, kernel, shells and special applications.



## Hardware

In our case the hardware are the parts of the computer like:

- Central processing unit or CPU. The CPU executes programs, which are located in the memory. Modern CPU has many cores enabling a first layer of parallelism.
- Random access memory, in short RAM or memory. The memory allows fast access to the instructions and data that the CPU needs and because the RAM is limited it will store only information for immediate use. Permanent data is kept in auxiliary storage on devices of the input/output subsystem.
- Input/output subsystem contains devices called peripherals that provide input and output for the system as well as the user, and short-term and long-term storage for the processes and data. These devices include among others:

- Hard disk(s)
- Various added cards (ie video, sound, USB-ports)
- Monitor(s)
- Keyboard and mouse
- Network

On Phoenix the most relevant parts and devices are: Hard disks, RAM, Cores (CPU) and network.

## Device drivers

Device drivers are special software used to communicate with peripheral devices. They are called kernel modules and are usually located in one of the subdirectories of kernel. When the device is requested the module is loaded otherwise is not.

The Kernel it is a small program, which is the main part of the operating system and is responsible for controlling system resources. It is compiled from source files, object files and configurable parameters. The kernel it is loaded into memory at very early stage when the system boots and initialised. Once the initialization is completed and various system daemons (processes performing system functions) have been created and made executable, the kernel waits for requests:

- user programs requesting services from the kernel through system calls.
- hardware devices getting kernel response through interrupts.
- when programs are running, the operating system creates a process to handle it.

## Shell

Shells are commandline language interpreters that execute commands from the standard input device (keyboard) or from file(s). The shell is the way to interface or communicate with the operating system.

The most popular shell today is the Bourne Again Shell (BASH, originates from Born shell used in the old unix systems, called sh). There are other types like: csh, tcsh and ksh among others. The most used shell on Phoenix is bash.

## **Programs and Applications**

Finally, the last layer is the programs and applications, which are used by users and the shell through the kernel commands to generate output or results. A program is an executable file, and a process is an instance of the program in execution.

## **Linux accounts**

In Linux, there are three types of accounts and all are using the programs available in the operating system based on their access and rights to the system:

- The root account is also called superuser, which has complete and unrestricted access to the system. Which means can run any commands but can also do big problems. On Phoenix, only administrators have root privileges.
- The system accounts are those accounts used for the operation of the system-specific components for example mail, database and secure shell daemon accounts. These accounts are usually needed for specific functions in the system and modifying them can affect the integrity or stability of the system.
- User accounts provide interactive access to the system for users and groups of users. Generally, users are assigned to this category of accounts and normally have limited access to critical system files and directories to protect the system, nevertheless, users have full privileges to their folders and files.

## **Graphical User Interface**

Windows OS users are very familiarised with GUI and it is for sure, the most critical difference with Linux, and Phoenix in particular.

Despite most Linux distributions has a GUI, Phoenix, as most of HPC facilities, does not support GUI. Therefore, all interaction within Phoenix must be done using a terminal. There are many terminal programs for Windows OS, such as Putty and Bitvise. All Unix based OS (macOS and Linux) already have a terminal program to access any Linux system.



## Filesystems

Linux OS data is organised in a folder structure, where the root folder is ”/”. Under this folder many subfolders are organised to support the system. Any folder can be a physical storage unit (hard disk) or a network storage. Examples of common subfolders of a Linux OS:

- /usr for common files and executables shared among users
- /etc system configurations
- /shared for shared resources among system and users
- /home the parent folder of users local folders

Phoenix has specific folders:

- /apps where applications and modules are installed
- /fast where parallel storage for user are located
- /data similar to fast but shared by users

Any file or folder has three levels of access. The owner, the group and others. On Phoenix, by default, any file or folder created by a user will have full permissions to its owner and no permissions for users in the same group or to any other users. If you want to change the access to any file or folder, you need to grant explicitly.



### Don't Panic!!!

More details will come in the next chapters about the practical use of the Filesystem.



---

# Understanding file system on Phoenix

Primary Author(s):

Bowen Chen, Technology Service, University of Adelaide  
[bowen.chen@adelaide.edu.au](mailto:bowen.chen@adelaide.edu.au)

Contributor(s):

Bowen Chen, Technology Service, University of Adelaide  
[bowen.chen@adelaide.edu.au](mailto:bowen.chen@adelaide.edu.au)

## Home directory

On Linux, each user has a personal directory, i.e. home directory, to store his own files and data, as well as directories. Your home directory on Phoenix is located at `/home/<userid>`. Here, `/<userid>` means `aXXXXXXXX`, your id.

When you login to Phoenix, your user shell is located in your home directory. The absolute pathname of your home directory, i.e. `/home/<userid>`, is also stored in the environmental variable, `$HOME`. Tilde symbol, `~`, also represent your home directory. You can use following commands to change into your home directory.



```
1 cd ~
2 cd $HOME
3 cd /home/<userid>
```



The `/home` file system is backed up, and is solely intended for the files that define your user environment and irrecoverable data such as source code. The default quota of your home directory is 10 GB.



Don't use `$HOME` for launching jobs or active job data. The `/home` file system hardware is not designed to support the intensive file access generated by the many hundreds of jobs that run on the Phoenix compute cores, you should use the `/fast` file system for job input/output instead.

## Fast directory

The `/fast` directory is designed to improve the data handling performance. It is supported by high performance Lustre filesystem and is intended for active job data and immediate data storage.



The `/fast` directory can handle the intensive file access created by the compute jobs that run on Phoenix. It is designed to launch jobs and store personal/research data. The default quota of your fast directory is 4 TB.

Your personal fast directory is located at `/fast/users/<userid>` The absolute pathname of your fast directory is also stored in the environmental variable, `$FASTDIR`. A symbolic link to that directory can be found in your home directory, i.e. `~/fastdir`. Thus, you can use following commands to change into your fast directory:



```
1 cd ~/fastdir
2 cd $FASTDIR
3 cd /fast/users/<userid>
```

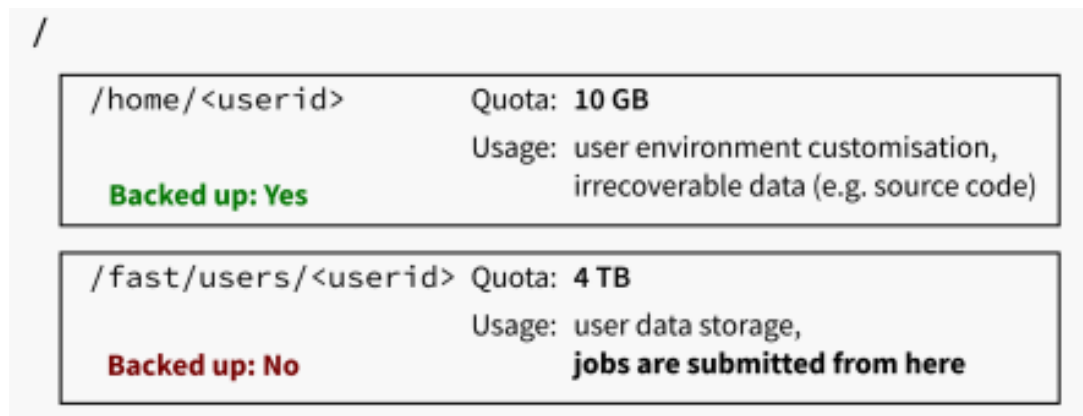


Figure 1: Home and fast directory on Phoenix

## Mounting filesystems



"A file system or filesystem is used to control how data is stored and retrieved. Without a file system, information placed in a storage medium would be one large body of data with no way to tell where one piece of information stops and the next begins. By separating the data into pieces and giving each piece a name, the information is easily isolated and identified." (Wikipedia) Mounting a filesystem simply means making the particular filesystem accessible at a certain point in the Linux directory tree.



You can use following commands to observe the mounting point and other information about /home and /fast on Phoenix

```
1 mount | egrep home
2 mount | egrep fast
```

## File and folder permissions

Since Linux is a multi-user OS that is based on concepts of file ownership and permissions to provide security, at the file system level. Using the command, `chmod`, can help us to change permissions if you are the owner of the files or directories.

Before to use it, you need to remember several symbols and letters with the `chmod`.

### Identities

- u - the user who owns the file, i.e. owner.
- g - the group to which the user belongs, i.e. group
- o - others (not the owner nor the owner's group)
- a - everyone or all (u, g, and o)



## Permissions

- r - read access
- w - write access
- x - execute access

## Actions

- + - adds the permission
- - - removes the permission
- = - makes it the only permission



The first thing we need to do is to create the file, `foo.txt`. Use the following commands to create the file and show its default permission.

```
1 touch foo.txt
2 ls -l foo.txt
```

When using `ls -l` command, you can see some thing like this:

```
-rw-r--r-- 1 a1695820 CATS_STUDENTS 5 Sep 26 14:34 foo.txt
```

Let's focus on `-rw-r--r--`. To explain it, parentheses will be added.

In `-(rw-)(r-)(r-)`, the first pair of parentheses means the users's permission; the second means the group's permissions; the last is for all the others.

Then using the `echo` command to write some words.

```
1 echo "On my way to supercomputing" > foo.txt
```

Using the following command to change permissions and show current permission of it. By typing `u-r`, you are going to remove read permission for the user, i.e. owner, and group from the file `foo.txt`.

```
1 chmod ug-r foo.txt
2 ls -l foo.txt
```

Now try to read the file with `cat` command

```
1 cat foo.txt
```



When you execute commands on Phoenix, you may also meet error message like: ... Permission denied. Considering the reason caused it and how should we fix it based on what we introduced.



You should know what went wrong. Now let's fix the issue.

```
1 ls -l foo.txt
2 cat foo.txt
3 chmod u+r foo.txt
4 ls -l foo.txt
5 cat foo.txt
```

Another way to change permissions uses numeric representation. Each permission setting can be represented by a numerical value:

- r = 4
- w = 2
- x = 1
- - = 0

For example, if foo.txt has following permissions settings:  
- (rw-) (rw-) (r-)

Let's compute the numeric representation. The numeric representation for the user is six(4+2+0), for the group is six(4+2+0), and for others is four(4+0+0). Thus, the permissions setting is 664.

Below is a list of common settings, numerical and meanings:

Setting	Numerical	Meaning
-rw-----	600	Only owner can read and write
-rw-r--r--	644	Only owner can read and write; group and others can read only
-rwx-----	700	Only owner can read, write and execute
-rwxr-xr-x	755	Owner can read, write and execute; group and others can read and execute
-rwx--x--x	711	Owner can read, write and execute; group and others can execute only
-rw-rw-rw	666	All can read and write. (Be careful with this setting)
-rwxrwxrwx	777	All can read, write and execute. (All can modify the file. Be careful.)

Belowing is a list of common settings for directories:

Setting	Numerical	Meaning
-rw-----	600	Only owner can read and write in the directory
-rwxr-xr-x	755	Owner can read, write and execute in the directory; group and others can read and execute

## Where are the applications on Phoenix?

The applications on Phoenix are located in the /apps folder. The directory stores applications, as well as related module files. We will talk about modules on Phoenix latter.



Let's observe the /apps directory. Firstly, list the directory's permission and what are in it.

```
1 ls -ld /apps
2 ls -l /apps
```

You should see the two directories, modules and software. Calculating the numerical permission of the two. Second, list what are in the directory, module

```
1 ls -l /apps/modules
```

Finally, list what are in the directory, software.

```
1 ls -l /apps/software
```

If too many things are listed, using pipe with less

```
1 ls -l /apps/modules | less
```

When using less, press f key for going forward and press b key for going backward. If you want to quit, press q.



---

# Regular Expressions

Primary Author(s):

Stephen Bent, Robinson Research Institute, University of Adelaide

Steve Pederson, Bioinformatics Hub, University of Adelaide

[stephen.pederson@adelaide.edu.au](mailto:stephen.pederson@adelaide.edu.au)

Contributor(s):

Dan Kortschak, Adelson Research Group, University of Adelaide

[dan.kortschak@adelaide.edu.au](mailto:dan.kortschak@adelaide.edu.au)

Bowen Chen, Technology Service, University of Adelaide [bowen.chen@adelaide.edu.au](mailto:bowen.chen@adelaide.edu.au)

## Introduction

*Regular expressions* are a powerful & flexible way of searching for text strings amongst a large document or file. Most of us are familiar with searching for a word within a file using software such as MS Word or Excel, but regular expressions allow us to search for these with more power and flexibility, particularly in Linux files and directories management. Instead of searching strictly for a word or text string, we can search using less strict matching criteria. For example, we could search for a string that is either `slurm-3725` or `slurm-3825` by using the patterns `slurm-3[78]25` or `slurm-3(7|8)25`. These two patterns will search for an `slurm-3`, followed by either a 7 or 8, then followed strictly by a 25. Similarly a match to `slurm-37725` can be found by using the patterns `slurm-3[78][78]25` or `slurm-3[78]{2}25`.

Whilst the bash shell has a great capacity for searching a file to matches to regular expressions, this is where languages like *perl* and *python* offer a great degree more power, such as more complex data structure, object-oriented programming(OOP) and more libraries.

## The command `grep`

The built-in command which searches using regular expressions in the terminal is `grep`. This function searches a file or input on a line-by-line basis, making patterns split across lines more difficult to find, which is one place that a programming language like Python or Perl would become preferable. The `man grep` page contains more detail on regular expressions under the `REGULAR EXPRESSIONS` header (scroll down a few pages). As can be seen in the `man` page, the command follows the form

```
grep [OPTIONS] 'pattern' filename
```

The option `-E` is preferable as it stands for Extended, which we can think of as “Easier”. As well as the series of conventional numbers and characters that we are familiar with, we can match to characters with special meaning, as we saw above where enclosing the two letters in brackets gave the option of matching either. The `-E` option opens up the full set of wild-card characters, and can also be called simply by using `egrep` instead of `grep -E`. This is the default version that many of us use.

Special Character	Meaning
\w	match any letter or digit, i.e. a word character
\s	match any white space character, includes spaces, tabs & end-of-line marks
\d	match any digit from 0 to 9
.	matches any single character
+	matches one or more of the preceding character (or pattern)
*	matches zero or more of the preceding character (or pattern)
?	matches zero or one of the preceding character (or pattern)
{x} or {x,y}	matches x or between x and y instances of the preceding character
^	matches the beginning of a line (when not inside square brackets)
\$	matches the end of a line
()	contents of the parentheses treated as a single pattern
[]	matches any one of the characters inside the brackets
[^]	matches anything other than any of the characters in the brackets
	either the string before or the string after the "pipe" (use parentheses)
\	don't treat the following character in the way you normally would. This is why the first three entries in this table started with a backslash, as this gives them their "special" properties, whereas placing a backslash before a '.' symbol will enable it to function as an actual dot/full-stop.

## Pattern Searching

In this section we'll learn the basics of using the **egrep** command & what forms the output can take. Using the command to download the file, `regular_express.txt`

```
cp /apps/examples/training_linux/regular_express.txt .
```

This is simply a text file with several words on every line.



Make sure the `txt` file in your working directory or else `egrep` won't be able to find the file.

Now let's try a few searches to get a feel for the basic syntax of the command. Using the previous table of special characters, try to describe what you're searching for on your notes **BEFORE** you enter the command. Do the results correspond with what you expected to see?

```
1 | egrep -n 't[ae]st' regular_express.txt
```

```
1 | egrep -n 'oo' regular_express.txt
```

```
1 | egrep -n '[^g]oo' regular_express.txt
```

```
1 | egrep -n '[^a-z]oo' regular_express.txt
```

```
1 | egrep -n '[0-9]' regular_express.txt
```

```
1 | egrep -n '[^[:lower:]]oo' regular_express.txt
```

```
1 | egrep -n '^[a-z]' regular_express.txt
```



In the above, we were changing the pattern to extract different results from the files. Now we'll try a few different options to change the output, whilst leaving the pattern unchanged. If you're unsure about some of the options, don't forget to consult the `man` page.

```
1 | egrep -n 'the' regular_express.txt
```

```
1 | egrep -vn 'the' regular_express.txt
```

```
1 | egrep -in 'the' regular_express.txt
```

---

# The Tools sed & awk

Primary Author(s):

Stephen Bent, Robinson Research Institute, University of Adelaide

Steve Pederson, Bioinformatics Hub, University of Adelaide

[stephen.pederson@adelaide.edu.au](mailto:stephen.pederson@adelaide.edu.au)

Contributor(s):

Dan Kortschak, Adelson Research Group, University of Adelaide

[dan.kortschak@adelaide.edu.au](mailto:dan.kortschak@adelaide.edu.au)

## sed: The Stream Editor



One very useful command in the terminal is **sed**, which is short for *stream editor*. Instead of the **man** page for **sed** the **info sed** page is larger but a little easier to digest. This is a very powerful command which can be a little overwhelming at first. If using this for your own scripts & you can't figure something out, remember "Google is your friend" & sites like [www.stackoverflow.com](http://www.stackoverflow.com) are full of people wrestling with similar problems to you. You can be certain you're not the first person to be stumped by a problem & these are great places to start looking for help. Even advanced programmers use Google & Stack Overflow to find solutions.

For today, there are two key **sed** functionalities that we want to introduce.

1. Using **sed** to alter the contents of a file/input;
2. Using **sed** to print regions of a file

### Altering a file or other input

**sed** uses *regular expressions* that we have come across under the **grep** section, and we can use these to replace strings or characters within a text string. The command works in the form

```
sed SCRIPT INPUT
```

and the script section is where all the action happens. Input can be given to **sed** as either a file, or just as a text stream via **stdin** using the *pipe* symbol that we have already introduced.



In the following example the script section begins with an 's' to indicate that we are going to make a substitution. The beginning of the first pattern (i.e. the *regexp* we are searching for) is denoted with the backslash, with the identical delimiter indicating the replacement pattern, and this is in turn completed with the same delimiter. Try this simple example from the link <http://www.grymoire.com/Unix/Sed.html> which is a very detailed & helpful resource about the usage **sed**. Here we are sending the input to the command via the pipe, so no 'INPUT' section is required:

```
1 | echo Sunday | sed 's/day/night/'
```

Here you are passing **sed** the string Sunday, and **sed** takes day and turns it into night. **sed** will only replace the first instance of the string on any line, so try:

```
1 | echo Sundayday | sed 's/day/night/'
```

Note that it only replaced the first instance of day and left the second. However, you can make it 'global', where it switches every instance by using the 'g' option at the end of the pattern like this:

```
1 | echo Sundayday | sed 's/day/night/g'
```

You can also 'capture' parts of the pattern in parentheses and access that in the second part of the regular expression (what you are switching to) using \1, \2, etc., to denote the number of the captured string, in the order they were captured. If you want to match 'ATGNNNTGA', where N is any base, and just output these three bases you could try the following:

```
1 | echo 'ATGCCAGTA' | sed -r 's/ATG(.{3})GTA/\1/g'
```

Clearly, we have just given this command a sequence so we know exactly what to expect. However, hopefully this demonstrates the concept of extracting a subset of the sequence.

Or if we needed to replace those three bases with an expanded repeat of them, you could do the following where we capture the undefined string between ATG & GTA, and expand it:

```
1 | echo 'ATGCCAGTA' | sed -r 's/ATG(.{3})GTA/ATG\1\1\1GTA/g'
```

The \1 take the contents of the first parenthesis and uses it in the substitution, even though you don't know what the bases are. Note that the '-r' option was set for these operations, which turns on extended regular expression capabilities. This can be a powerful tool & multiple parentheses can also be used:

```
1 | echo 'ATGCCAGTA' | sed -r 's/(ATG)(.{3})(GTA)/\3\2\2\1/g'
```

In this last command we switched the order of the first & last triplet, and expanded the middle unknown string twice. Note how quickly this starts to look confusing though! Taking care to be clear when writing these types of procedures can be an important idea when you have to go back & re-read your code a year or two later. (Yes this will happen a lot!!!)



The use of backslashes to delineate each section of the script used by sed is the most common convention, but we are not restricted to it. We could have used any 'wild-card' type character to follow the 's', such as '\*' ':' or '%' although this must be consistent across all sections of the script, and should only be used if the backslash itself is part of the search or replacement string.

## Displaying a region from a file

The command `sed` can also be used to replicate some of the functionality of the `head` & `grep` commands, but with a little more power at your fingertips. By default `sed` will print the entire input stream it receives, but setting the option `-n` will turn this off. Try this by adding an `n` immediately after the `-r` in one of the above lines & you will notice you receive no output. This is useful if we wish to restrict our output to a subset of lines within a file, and by including a `p` at the end of the script section, only the section matching the results of the script will be printed.



Make sure you are in the correct directory & we can look through the `regular_express.txt` file again.

```
1 | sed -n '1,10p' regular_express.txt
```

This will print the first 10 lines, like the `head` command will by default. However, we could now print any range of lines we choose. Try this by changing the script to something interesting like `'15,21p'`.

We could also restrict the range to specific lines by using the `sed` increment operator `'~'`.

```
1 | sed -n '1~5p' regular_express.txt
```

This will print every 5th line, beginning at the first.



We can also make `sed` operate like `grep` by making it only print the lines which match a pattern.

```
1 | sed -rn '/g[o]*d/ p' regular_express.txt
```

Note however, that the line numbers are not present in this output.

## Some Important Programming Concepts

Before moving on to `awk`, we need to quickly recap two of the most widely used techniques in programming:

1. The `for` loop
2. Logical tests using an `if` statement



## For Loops



A **for** loop is what we use to cycle through an input one item at a time. As a simple example, we could print each number from a set of numbers.

```
1 | for i in 1 2 3; do (echo -e $i); done
```



In the above code, the fragment before the semi-colon asked the program to cycle through the values 1, 2 & 3, letting the variable ‘i’ take each value in order of appearance. First  $i = 1$ , then  $i = 2$  & finally  $i = 3$ . If you’re wondering why we chose ‘i’, it just seemed like a sensible choice for an integer. We simply needed to choose a name for a *variable* which we would pass the values to.

After assigning each value to *i*, was the instruction on what to do for each value. Note that the value of the variable ‘i’ was *prefaced by the dollar sign (\$)*. *This is how the bash shell knows it is a variable, not the letter ‘i’*. The command **done** then finished the **do** command. All commands like **do**, **if** or **case** have completing statements, which respectively are **done**, **fi** & **esac**.

Another important concept which was glossed over in the previous paragraph is that of a *variable*. These are essentially just ‘*placeholders*’ which have a value that can change (hence the name). In the above loop, the same operation was performed on the variable *i*, but the value changed from 1 to 2 to 3. Variables in shell scripts can hold numbers or text strings and don’t have to be formally defined as in some other languages.



An alternate approach could be to make a breathtaking claim about some files. Here we’ll use the variable called ‘f’, which seems sensible for a filename.

```
1 | cd /apps/examples/training_linux/scripting/  
2 | for f in $(ls); do (echo -e "I can see the file $f"); done  
3 | cd -
```

Note, that we’ve also assigned the output from the command **ls** to this variable, by using the **\$( )** syntax. The use of double quotes for the **echo** command also allows us to refer to the values held by *f*. Single quotes at this point would only return the characters ‘\$f’.

## If Statements

If statements are those which only have a binary ‘yes’ or ‘no’ response, or more correctly a TRUE/FALSE response. For example, we could specify things like:

- `if (i > 1) then do something, or`
- `if (fileName==bob.txt) then do something else`



Notice that in the second if statement, there was a double equals sign. This is the programmers way of saying *compare* the first argument with the second argument. A single equals sign is generally interpreted by a program as *assign* the value of the first argument to be the second argument. This use of ‘double operators’ is very common, notably you will see `&&` to represent the command ‘and’, and `||` to represent ‘or.’ A final useful trick to be aware of is the use of an exclamation mark to reverse a command. A good example of this is the use of the command ‘`!=`’ as the representation of *not equal to* in a logical test.

## awk: A command and a language

Moving on to `awk`, this is a very useful tool which can be used either as a command, as well as functioning as it’s own language. We’ll just use it as a command today, and it is extremely useful for dealing with tab- or comma-separated files, such as we often see in biological data.



The basic structure of an `awk` command is:

```
awk '/<pattern>/' file
```

`awk` will then search the file and output any line containing the regular expression pattern (kind of like `grep`). With `awk`, you can also do:

```
awk '\{<code>\}' file
```

where you can put a program, or set of instructions in the curly braces. In the code, you can specify values from different columns of the file by using the numbers `$1`, `$2`, etc., (or you can use `$0` for the whole line). Values can also be returned in the output by using the command `print` followed by the field number.



We have that `regular_express.txt` file and we’ve already looked at it a little, so let’s pull out some particular features! Make your terminal as wide as the screen, then change into the appropriate directory & enter

```
1 awk '{if ($3=="is") print $0}' regular_express.txt
```

Here we’ve specified that the third field must be `is`.

We could make it a little more complex and just look for genes in a given region. **In the following line, the symbol ‘\’ has been placed here to indicate it is a single line, extending beyond the width of the page. Do not enter this character!**

```
1 awk '{if ( ($3=="gene") && ($4 > 10000) && ($4 < 20000) ) print $0}' \
    regular_express.txt
```



In the above code, `$3=="gene"` asks for the entry in the third field to be “gene.” The next two fragments request for the values in the fourth field (i.e. `$4`) to be between 10000 & 20000. Notice that each these three commands were enclosed in a pair of brackets within an outer pair of brackets. This gave a command of the form:

```
( (Condition1) && (Condition2) && (Condition3) )
```

After this came the fragment `print $0` which asked `awk` to print the entire line if the 3 conditions are true. You’ve just written (& hopefully understood) a computer program!



Another example (what does this do?):

```
1 awk '{if (($5 - $4 > 1000) && ($3 == "gene")) print $0}' \
    regular_express.txt
```

If you don’t want to output all of the columns, you can specify which ones to output. While we’re at it, let’s save the output as a file:

```
1 awk '{if (($5 - $4 > 1000) && ($3 == "gene")) print $1, $2, $4, $5, \
    $9}' regular_express.txt > awkout.txt
```



The command `awk` has a pretty serious set of in-built commands which can be used in the code sections as above. Although it looks a little overwhelming, there is a detailed page <http://www.grymoire.com/Unix/Awk.html> which gives a rundown on the full capabilities of the language. One command that we may find helpful is `length`, which counts the number of characters in a line.



---

# Command Line Editors

Primary Author(s):

Robert Qiao, Research Services Phoenix Team, University of Adelaide  
[robert.qiao@adelaide.edu.au](mailto:robert.qiao@adelaide.edu.au)

Contributor(s):

Dan Kortschak, Adelson Research Group, University of Adelaide  
[dan.kortschak@adelaide.edu.au](mailto:dan.kortschak@adelaide.edu.au)

In addition to the streaming editor named `sed` as you have studied earlier, which operates on a file in a non-interactive mode according to a set of instructions that you specify on the command line or in a script. It can be used for "canned" editing tasks that must be applied in the same way to many files. More commonly, we use interactive editors and there are three powerful interactive text editors that are commonly used on Unix computers: `vi`, `pico`, and `emacs`.

Although `vi`, `pico`, `emacs` are extremely powerful, to get a firm grasp takes some efforts and most of times we just need to open, amend and save a file without remembering the keyboard shortcuts. Luckily, there is such editor called `nano` and that is what we will use for the class today. The `nano` editor has its own set of keyboard shortcuts of course and in this guide I aim to help you to understand the meaning of all those special keystrokes you can use to make your life easier when using `nano`.

But before we dive into the `nano`, let's generally go over some of the streamline editors like `vi`, `emacs`.

## vi



`Vi` (pronounced "vee-eye") is the standard editor for Unix systems. It is universally available on Unix systems. `Vi` is a screen editor. It treats your computer screen as a window into the file. You move the window around to view different parts of the file. You move the cursor to the location on the screen where you want to make a change; or optionally, you specify some kind of global change. `Vi` updates your screen to reflect changes that you make in the file. Actually, it works on a copy of the file in memory, and only updates the file on disk when you tell it to, such as when you end the editing session.

Main advantages of `vi` are including:

- `Vi` is universally available on Unix systems. It has been around so long in a stable form that it is essentially bug free. Many clones have been written for other kinds of computers.
- `Vi` has many powerful commands that utilize just the alphanumeric keys – it does not require special function keys.
- `Vi` is a small program that does not require a lot of system memory or CPU time. It works very fast, even on large files.
- While `vi` is not programmable, it has a simple way to let other Unix programs, such as the `sort` utility, work on selected portions of your file. This adds the functionality of all those programs to the editor.

- Vi is completely terminal device independent. It will work with any kind of terminal. A system file describes the capabilities and control sequences of each kind of terminal for vi. All the program needs to know is what type of terminal you have. When you log in, if pangea cannot figure out what kind of terminal you have, it will prompt you to specify a terminal type. The most common type is the vt100, which most modern terminals and PC communications software emulate.

The chief disadvantage of vi is that it is touchy. That is, every single key you touch on the keyboard seems to do something, often something mysterious. There is a rich set of single character commands to learn.

## emacs



Emacs is a widely used editor available on many types of systems. The GNU version is installed on pangea. This editor is very large and fully programmable with a built-in object oriented language. As a result, it has become more of a shell than an editor. Using "macros", it is possible to do many non-editing functions from within emacs, including compiling and debugging programs and browsing the web.

Emacs is not described in these notes because it is too complicated and bloated for simple editing. Entire books are available to teach emacs.

If you are already familiar with emacs, you can use it on pangea and most other Unix systems on campus. Simply type the shell command

```
1 | emacs
```

## Nano

Now, let's have a more detailed look of Nano.



To setup this part of the course, please login to your Phoenix account and download the example file to your /fastdir directory now

```
1 | cp -r /apps/examples/training_linux/ ~/fastdir/
```

please check to make sure you command blow contains valid 4 files. If you see error messages, please indicate to tutors, you should get this step fixed before carry on further

```
1 | ls -al ~/fastdir/training_linux
```

## How To Get Nano

The nano editor is available by default in all the most popular Linux distributions and you can run it with one simple command:

```
1 nano
```

The above command will simply open a new file. You can type into the window, save the file and exit.

## How To Open A New File And Give It A Name Using Nano

Whilst simply running nano is ok you might want to give your document a name before starting. To do this simply give the filename after the nano command.

```
1 nano myfile.txt
```

You can, of course, supply a complete path to create a file anywhere on your Linux system (as long as you have the permissions to do so).

```
1 nano ~/fastdir/training_linux/myfile.txt
```

## How To Open An Existing File Using Nano

You can use the same command as the one above to open an existing file.

Simply run nano with the path to the file you wish to open.

To be able to edit the file you must have permissions to edit the file otherwise, it will open as a readonly file (assuming you have read permissions).

```
1 cd ~/fastdir/training_linux
2 nano seqData.fastq
```

You can, of course, supply a complete path to open a file anywhere on your Linux system (as long as you have the permissions to do so).

```
1 nano ~/fastdir/training_linux/seqData.fastq
```



## How To Save A File Using Nano

You can add text to the nano editor simply by typing the contents directly into the editor. Saving the file, however, requires the use of a keyboard shortcut.

To save a file in nano press `ctrl` and `O` at the same time.

If your file already has a name you just need to press `enter` to confirm the name otherwise you will need to enter the filename that you wish to save the file as.

## How To Save A File In DOS Format Using Nano

To save the file in DOS format press `ctrl` and `O` to bring up the filename box. Now press `alt` and `d` for DOS format.

## How To Save A File In MAC Format Using Nano

To save the file in MAC format press `ctrl` and `O` to bring up the filename box. Now press `alt` and `m` for MAC format.

## How To Backup A File Before Saving It In Nano

If you want to save the changes to a file that you are editing but you want to keep a backup of the original press `ctrl` and `O` to bring up the save window and then press `alt` and `B`.

The word `[backup]` will appear in the filename box.

After you have finished editing your file you will want to leave the nano editor.

To exit nano simply press `ctrl` and `x` at the same time.

If the file hasn't been saved you will be prompted to do so. If you select "Y" then you will be prompted to enter a file name.

## Displaying Cursor Position Using Nano

If you wish to know how far down a document you are within nano you can press the `ctrl` and `c` keys at the same time.

The output is shown in the following format:

```
1 | line 5/11 (54%), col 10/100 (10%), char 100/200 (50%)
```

## How To Read A File Using Nano

If you opened nano without specifying a filename you can open a file by pressing `ctrl` and `r` at the same time.

You are now able to specify a filename to read into the editor. If you already have text loaded into the window the file you read in will append itself to the bottom of your current text.

If you want to open the new file in a new buffer press `alt` and `f`.

## How To Search And Replace Using Nano

To start a search within nano press `ctrl` and `\`.

To turn off replace press `ctrl` and `r`. You can turn on replace again by repeating the keystroke.

To search for text enter the text you wish to search for and press return.

To search backwards through the file press `ctrl` and `r` to bring up the search window. Press `alt` and `b`.

To force case sensitivity bring up the search window again and then press `alt` and `c`. You can turn it off again by repeating the keystroke.

Nano wouldn't be a Linux text editor if it didn't provide a way to search using regular expressions. To turn regular expressions on bring up the search window again and then press `alt` and `r`.

You can now use regular expressions for searching for text.

## Check Your Spelling Within Nano

Again nano is a text editor and not a word processor so I'm not sure why spelling is a key feature of it but you can indeed check your spellings using the `ctrl` and `t` keyboard shortcut.

## Some Goodies



There are a number of switches you can specify when using nano. The best ones are covered below. You can find the rest by reading the nano manual.

- nano -B (backs up the file prior to editing it)
- nano -E (converts tabs to spaces when editing)
- nano -c (constantly show the cursor position stats)

- nano -i (automatically indents new lines to the same position as the previous line)
- nano -k (toggle cut so that it cuts from cursor position instead of the whole line)
- nano -m (provides mouse support to the editor)
- nano -v (opens file as readonly)



---

# Writing Scripts

## Primary Author(s):

Stephen Bent, Robinson Research Institute, University of Adelaide

Steve Pederson, Bioinformatics Hub, University of Adelaide

[stephen.pederson@adelaide.edu.au](mailto:stephen.pederson@adelaide.edu.au)

Robert Qiao, Research Services Phoenix Team, University of Adelaide

[robert.qiao@adelaide.edu.au](mailto:robert.qiao@adelaide.edu.au)

## Contributor(s):

Dan Kortschak, Adelson Research Group, University of Adelaide

[dan.kortschak@adelaide.edu.au](mailto:dan.kortschak@adelaide.edu.au)

Jimmy Breen, Robinson Research Institute & Bioinformatics Hub, University of Adelaide

[jimmy.breen@adelaide.edu.au](mailto:jimmy.breen@adelaide.edu.au)

We often need to perform repetitive tasks, or need to perform complex series of procedures. Rather than typing instruction receptively and steering the screen waiting for each process to finish before providing next instruction, writing the set of instructions into a script and interpreter/compiler does the waiting and steering for us is a very powerful way of liberating our time. They are also an excellent way of ensuring the commands you have used in your research are retained for future reference. Keeping copies of all electronic processes to ensure reproducibility is a very important component of any research. Writing scripts requires an understanding of several key concepts which form the foundation of much computer programming, so let's walk our way through a few of them.

## Shell Scripts

Now that we've been through just some of the concepts & tools we can use when writing scripts, it's time to tackle one of our own where we can bring it all together.



Every bash script begins with what is known as a *shebang*, which we would commonly recognise as a hash sign followed by an exclamation mark, i.e `#!`. This is immediately followed by `/bin/bash`, which tells the interpreter to run the command `bash` in the directory `/bin`. (This is actually where the program `bash` lives on a Linux system.) This opening sequence is vital & tells the computer how to respond to all of the following commands. As a string this looks like:

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



The hash symbol generally functions as a comment character in scripts. Sometimes we can include lines in a script to remind ourselves what we're trying to do, and we can preface these with the hash to ensure the interpreter doesn't try to run them. Its presence as a comment here, followed by the exclamation mark, is specifically looked for by the interpreter but beyond this specific occurrence, comment lines are generally ignored by scripts & programs.

## Some Example Scripts

Let's now look at some simple scripts. These are really just examples of some useful things you can do & may not really be the best scripts from a technical perspective. Hopefully they give you some pointers so you can get going



To setup this part of the course, please login to your Phoenix account and download the example file to your /fastdir directory now

```
1 | cp -r /apps/examples/training_linux/ ~/fastdir/
```

please check to make sure you command blow contains valid 4 files. If you see error messages, please indicate to tutors, you should get this step fixed before carry on further

```
1 | ls -al ~/fastdir/training_linux
```

## A Simple Example to Start



**Don't try to enter these commands directly in the terminal!!!** They are designed to be placed in a script which we will do after we've inspected the contents of the script. First, just have a look through the script & make sure you understand what the script is doing.

Also remember that any long lines of code may be automatically broken into new lines on your page by the '\ ' character. We don't want to enter this character when we create our script. Note that the line numbers on the left of the code don't change when this happens, e.g. line 9.

Before we go any further, have a look at the following script.

```
#!/bin/bash
#
# First we'll declare some variables with some text strings
ME='Put your name here'
MESSAGE='This is your first script'

# Now well place these variables into a command to get some output
echo -e "Hello ${ME}\n${MESSAGE}\nWell Done!"
```



Firstly, you may notice some lines that begin with the # character. These are *comments* which have no impact on the execution of the script, but are written so you can understand what you were thinking when you wrote it. If you look at your code 6 months from now, there is a very strong chance that you won't recall exactly what you were thinking, so these comments can be a good place just to explain something to the future version of

yourself. There is a school of thought which says that you write code primarily for humans to read, not for the computer to understand.



In the above script, there are two variables. Although we have initially set them to be one value, they are still variables. What are their names?



First we'll create an empty file which will become our script. We'll give it the suffix `.sh` as that is the common convention for bash scripts.

```
1 | cd ~/firstname
2 | touch wellDone.sh
```



Now using the text editor *nano*, enter the above code into this file *setting your actual name as the ME variable*, and save it by using `Ctrl+o`, which is indicated as `␣` in the nano screen.

```
1 | nano wellDone.sh
```



Once you're finished, you can exit the *nano* editor by hitting `Ctrl+x`.

Another coding style which can be helpful is the enclosing of each variable name in curly braces every time the value is called. Whilst not being strictly required, this can make it easy for you to follow in the future when you're looking back. Variables have also been names using strictly upper-case letters. This is another optional coding style, but can also make things clear for you as you look back through your work. Most command line tools use strictly lower-case names, so this is another reason the upper-case variable names can be helpful.



Unfortunately, this script cannot be executed yet but we can easily enable execution of the code inside the script. If you recall the flags from earlier which denoted the read/write/execute permissions of a file, all we need to do is set the execute permission for this file. First we'll look at the files in the folder using `ls -l` and note these triplets should be `rw-` for the user & the group you belong to. To make this script executable, enter the following in your terminal.

```
1 | cd ~/firstname
2 | chmod +x wellDone.sh
3 | ls -l
```





Notice that the third flag in the triplet has now become an **x**. This indicates that we can now execute the file in the terminal. As a security measure, Linux doesn't allow you to execute a script from within the same directory so to execute it enter the following:

```
1 | ./wellDone.sh
```

### Making a Small Change



Now let's change the variable **ME** in the script to read as

```
1 | ME=$1
```

and save this as **wellDone2.sh**. (You may like to create this first using **cp**) You'll now need to set the execute permission again.

```
1 | chmod +x wellDone2.sh
```



This time we have set the script to *receive input from stdin* (i.e. the terminal), and we will need to supply a value, which will then be placed in the variable **ME**. Choose whichever random name you want and enter the following

```
1 | ./wellDone2.sh Boris
```



As you can imagine, this style of scripting can be useful for iterating over multiple objects. A trivial example, which builds on a now familiar concept would be to try the following.

```
1 | for n in Boris Fred; do (./wellDone2.sh $n); done
```

## A more complicated script

Here's a more complicated script with some more formal procedures. This is a script which will extract only the ProbeFeature features from the .gff file we have been working with, and export them to a separate file. Look through each line carefully & write down your understanding of what each line is asking the program to do.

```
#!/bin/bash
# Declare some helpful variables
FILEDIR=~/.fastdir/training_linux/scripting
FILENAME=NC_015214.gff
OUTFILE=NC_015214_CDS.txt
# Make sure the directory exists
if [ -d ${FILEDIR} ]; then
    echo Changing to ${FILEDIR}
    cd ${FILEDIR}
else
    echo Cannot find directory ${FILEDIR}
    exit 1
fi

# If the file exists, extract the important ProbeFeature data
if [ -a ${FILENAME} ]; then
    echo Extracting ProbeFeature data from ${FILEDIR}/${FILENAME}
    echo "SeqID Source Start Stop Strand Tags" > ${OUTFILE}
    awk '{if (($3=="ProbeFeature")) print $1, $2, $4, $5, $7, $9}' \
        ${FILENAME} >> \
        ${OUTFILE}
else
    echo Cannot find ${FILENAME}
    exit 1
fi
```

Notice that this time we didn't require a file to be given to the script. We defined it within the script, as we did for the output file.



The directory & file checking stages were of the form `if [...]`. This is a curious command that checks for the presence of something. The options `-d` & `-a` specify a directory or file respectively.



Will the above script generate a tab, comma or space delimited text file?



Open the `gedit` text editor & save the blank file in your directory as `extract_CDS.sh`. Now write this above script into the editor, but *taking care to use the directory where you have the .gff file stored in the appropriate place*. Once you have written the script, save it & close it. Now make it executable and run it.

## Moving towards High Performance Computing

### High Performance Computing



In current genomics era, where we regularly work with large datasets, the amount of resources available on desktop computers are often insufficient to enable your script finish in a reasonable time. For these datasets, it maybe useful to work on a high performance computing system, which will enable your data or command to be run simultaneously on > 8 threads, i.e. in parallel. It is possible to gain access to large computing resources through the University's Phoenix HPC <https://www.phoenix.adelaide.edu.au>, enabling analysis of large datasets efficiently.



Having many users on one machine at one time also means that there needs to be a system which determines who runs what and when. Phoenix uses a scheduler "SLURM", whereby users submit jobs to a queue, and then executed when the appropriate resources on the machine become available.

A typical slurm job script contains extra parameters such as:

- `-n`: The number of threads to allow
- `--time`: The maximum time it is allowed to takes to completion
- `--mem`: The amount of memory to allocate to the job

```
1 #!/bin/bash
2 #SBATCH -p batch
```

```
3 #SBATCH -N 1
4 #SBATCH -n 8
5 #SBATCH --time=20:00:00
6 #SBATCH --mem=20GB
7
8 # Execution code going below
9 <execution code>
```

We don't need to write this script. It is included here as a simple example of a real world script as used by Phoenix users. This script may look a little intimidating at first, but slowly work your way through each line & try to understand what each line is specifying. More detailed info will be included in the last chapter of this course.

---

# Using Modules

Primary Author(s):

Exequiel Sepulveda, Research Services, University of Adelaide  
[exequiel.sepulvedaescobedo@adelaide.edu.au](mailto:exequiel.sepulvedaescobedo@adelaide.edu.au)

Contributor(s):

Robert Qiao, Research Services, University of Adelaide [robert.qiao@adelaide.edu.au](mailto:robert.qiao@adelaide.edu.au)

There are many applications that need to change or define configurations to work properly. For example, they need to add the folder where binaries are located to the PATH environmental variable. If an application has many binaries located in /apps/MYAPP/bin, setting the PATH variable can be done by:

```
1 export PATH=$PATH:/apps/MYAPP/bin
```

Doing manually these modifications may be very tedious. The situation worsen if there are many applications that need to do the same with many potential conflicts and side effects. To keep environmental variables under control, Phoenix has an program that manages applications in a safe way. This program is called "modules"

## Basic commands of modules

Phoenix has installed hundreds of applications for users to use. By default, all applications are located in /apps/software and each one has a module configuration located in /apps/modules/all. Modules program has an unique executable named module.



Try to get the help from module program:

```
module --help
```

The relevant options of module command are:

Option	Description of function	Useful options
<code>avail [name]</code>	Display the list of modules	-r for using regular expressions, -d for listing only the default version of each module
<code>spider [name]</code>	Explore and list modules	-r for using regular expressions, -d for listing only the default version of each module
<code>load name</code>	Load the specific module name	
<code>unload name</code>	Unload the specific module name	
<code>list</code>	List all loaded modules	
<code>show name</code>	Show the information of the module name	
<code>purge</code>	Unload all loaded modules to have a fresh starts	



The `avail/spider` option should be always used with a text to search for, because listing all modules installed may be very slow and useless. The correct way to use `avail/spider` is using a text after, for example, to search Perl modules:

```
module avail Perl
module spider Perl
```



Spider and Avail are similar but the output format is different. As an exercise, try to find all matlab modules using spider and avail options.

```
module spider matlab
module avail matlab
```



Repeat the same exercise to search for Python modules installed on Phoenix.

In general, most of modules follow this pattern: `MODULE/VERSION[-TOOLCHAIN]`. For example, for the module `Python/3.6.1-foss-2016b`: `MODULE=Python`, `VERSION=3.6.1` and `TOOLCHAIN=foss-2016b`. Which toolchain should be selected is critical to understand.



Toolchain is a set of compiler, utilities and libraries (all they are modules!). The general rule is to use the newest official toolchain, which is `foss-2016b`. On Phoenix there are two main toolchains, `foss` and `intel`. The `foss` toolchain includes GNU compilers, whereas `intel` includes Intel compilers. You should use modules with the same toolchain. If you need to

use modules with a different toolchain, a best practice is to use the purge option before changing toolchain.

## Loading modules

Once you have found the module you want to load, the next step is loading it.

For example, search for R modules and load the newest one:



```
module spider R
```

```
    Versions
```

```
R/3.2.1-foss-2015b
```

```
R/3.3.0-foss-2016uofa
```

```
R/3.4.0-foss-2016b
```

```
module load R/3.4.0-foss-2016b
```

To see what happens after the module is loaded, use the option list:

```
module list
```

Why are there many loaded modules? It is because R/3.4.0-foss-2016b module has many prerequisites that are automatically loaded as well.



In most the cases, module configurations include automatically all modules needed. In few cases, because a precise control is necessary, this should be manually done.

## Using modules

Once you have loaded modules, the current session will be correctly configured to use those modules. For example, you could load any Python module and check if the version is the correct one:

```
module load Python
```

```
python --version
```



Be careful with modules containing programs that are part of the operating system. For example, C/C++ compiler, Python, Perl, among others. If you don't load the right module, you may end using the wrong version.





What is the Python version included in the operating system?

## A complete exercise of modules

In this section you are asked to complete a full exercise to reach a master level of modules. Run the traditional "Hello World!" program in four different languages and versions.

Python 2.X:

```
print "Hello World, from Python 2.x"
```

Python 3.X:

```
print("Hello World, from Python 3.x")
```

R:

```
print("Hello World, from R")
```

Matlab:

```
display("Hello World, from matlab");
```

And these are the examples of how to run programs on those languages:

R:

```
echo 'COMMAND' | R --no-save
```

Matlab:

```
matlab -noFigureWindows -nodisplay -r 'COMMAND; exit;'
```

Python:

```
python -c COMMAND
```



The exercise is to find the modules to run each program and execute them. Replace COMMAND accordingly to each language.



---

# Writing Basic Slurm Scripts

Primary Author(s):

Exequiel Sepulveda, Research Services, University of Adelaide  
[exequiel.sepulvedaescobedo@adelaide.edu.au](mailto:exequiel.sepulvedaescobedo@adelaide.edu.au)

Contributor(s):

Robert Quiao, Research Services, University of Adelaide  
[robert.quiao@adelaide.edu.au](mailto:robert.quiao@adelaide.edu.au)

Phoenix is a shared HPC facility, therefore, all its resources are not for exclusive use of a particular user. The Phoenix architecture includes a head node and many computing nodes. The head node is where you are placed, after connect to Phoenix and is not designed for real computing.

The computing nodes are only available to our scheduler Slurm, which receives many scripts from all users and decides when and where to run those scripts according to the resources available and account priorities. When Slurm receives a new job (or script) to run, he evaluates when is the sooner available time to run the job and assigns one or many computing nodes. The job will be ran by computing nodes, but not by the head node.



Slurm accounting, resources and priorities are complex enough. We have a workshop dedicated to this.

Therefore, to use the real power of Phoenix you need to write a slurm script. The script template is very simple: it is a normal script as you have written many, but with several Slurm meta-parameters at the beginning of the script.

## Basic template of a Slurm script

A useful template for any slurm job is like this:

```
1  #!/bin/bash
2  #SBATCH -p batch
3  #SBATCH -N 1
4  #SBATCH -n 1
5  #SBATCH --time=00:05:00
6  #SBATCH --mem=1GB
7
8  #If you want to get feedback by email
9  #SBATCH --mail-type=ALL
10 #SBATCH --mail-user=firstname.lastname@adelaide.edu.au
11
12 #LOAD HERE ALL MODULES YOU NEED
13 #module load MODULE1
14 #module load MODULE2
15 #module load MODULE3
16
17 #EXECUTE HERE WHATEVER YOU WANT
18 #hostname
```

Let's have a look to each Slurm meta-parameter:

`#SBATCH -p batch` indicates the queue where the job will be submitted to.

Phoenix has several queues, for example batch, test and bigmem. The batch queue is the default one for almost all kind of jobs. The queue test is useful to test jobs since the waiting times are reduced but also resources are limited just for testing.



`#SBATCH -N 1` indicates the number of computing nodes needed to run the job. Usually should be set to 1. `#SBATCH -n 1` indicates the number of cores needed to run the job. Usually should be set to 1.



Most of programs are serial and do not have parallel capabilities. If you do not know the capability of a program, assume the program is serial and, therefore, nodes and cores should be set to 1. If the program can use multicores, nodes should be set to 1 and cores to the desired number of cores, but less or equals to 32, which is the current maximum. If the program can use Message Passing Interface (MPI) and multicores, nodes should be set to equals or greater than 1, and cores to the desired number of cores, but less or equals to nodes\*32.



`#SBATCH --time=00:05:00` indicates the maximum walltime of the job. In this example the maximum time will be 5 minutes.

`#SBATCH --mem=1GB` indicates the maximum amount of memory. In this example the maximum memory will be 1 Gigabytes.



Both time and mem meta-parameters are hard limits parameters. This means, if the job spends more time or consumes more memory, the job will be cancelled. There is a trade-off. You would use conservative values, but if you overestimate them, the scheduler will take longer to executed it.

Getting notifications by e-mail is optional, but it is very recommended to know the progress of a job.



`#SBATCH --mail-type=ALL` indicates that Slurm will notify any state change, such as Cancellation, Success and Failure. `#SBATCH --mail-user=firstname.lastname@adelaide.edu.au` indicates the email address to send those notifications.

## Submitting Scripts

Once you have a script ready to submit, you need to "submit" the script to the scheduler by the use of the command `sbatch`. Let's assume there is a script `job_script.sh`. To submit it to the schedule, just use the following command:

```
$ sbatch jobs_script.sh
Submitted batch job 3853482
```

In return you get a job number or *JOBID*, 3853482 is the example above, which is very important to record. Also, a new file with name `slurm-JOBID.out` is generated with the output of the script.



Create and submit a script that shows the hostname. Hint: the command for asking the hostname is `hostname`.

## Managing Scripts

Finally, after a job submission, you would need to manage its execution and get the information of completion. There are four basic commands to manage jobs:

Command	Description of function	Options
<code>squeue</code>	Display the list of queued jobs	<code>-u USER</code>
<code>scancel</code>	Cancel the execution of JOBID	JOBID
<code>scontrol show job</code>	Show useful information of a queued or running job	JOBID
<code>rcstat</code>	Show useful information of a job, specially for a completed job	JOBID

Now, from the last exercise in the Modules chapter (the execution of Hello World example in four languages), the final exercise:



Adapt that script as a Slurm script and submit it to the scheduler.

You can monitor its execution.



The command `rcstat` is useful to get the job statistics back, such as the real time spent, cpu and memory used, among others. You should use this information to adjust those values for the next submissions in order to request the right resources as closer as possible to the real ones.





---

Space for Personal Notes or Feedback

[illegible]

[illegible]

[illegible]

[illegible]