



# **INTRODUCTION TO BIOINFORMATICS FOR MALARIA MOLECULAR SURVEILLANCE**

Fighting Malaria Across Borders  
Pieter Moris



# Table of contents

<b>About</b>	<b>5</b>
Scope of the course . . . . .	5
<b>Acknowledgements</b>	<b>6</b>
<b>License</b>	<b>7</b>
<b>Introduction</b>	<b>8</b>
What is bioinformatics anyway? . . . . .	8
Computational thinking . . . . .	8
<b>I Introduction to the Unix shell</b>	<b>9</b>
<b>1 What is a CLI?</b>	<b>10</b>
1.1 Learning objectives . . . . .	10
Resources . . . . .	10
1.2 What is Unix? . . . . .	10
1.3 Why bother learning the Unix shell as a bioinformatician? . . . . .	11
1.4 Don't get discouraged . . . . .	12
1.5 A note on terminology . . . . .	12
<b>2 Setting up your own Unix shell</b>	<b>13</b>
2.1 Online Unix environment . . . . .	13
2.1.1 GitHub Codespaces . . . . .	13
2.1.2 Binder . . . . .	16
2.2 Local Unix environment . . . . .	16
2.2.1 Download the course files . . . . .	17
2.2.2 WSL installation . . . . .	18
2.2.2.1 Configuring WSL . . . . .	19
2.2.2.2 Accessing files across the Windows and WSL file systems . . . . .	19
2.2.2.3 Windows Terminal . . . . .	20
<b>3 Using the shell</b>	<b>21</b>
3.1 Interacting with the shell . . . . .	21
3.2 Command syntax . . . . .	21
3.3 Tips and hints . . . . .	22
<b>4 Navigating the Unix file system</b>	<b>23</b>
4.1 Layout of the Unix file system . . . . .	23
4.1.1 Home sweet home: ~ . . . . .	23
4.1.2 Where am I? . shortcuts . . . . .	25
4.2 Moving around the file system . . . . .	25
4.2.1 pwd: avoid getting lost . . . . .	25
4.2.2 cd: on the move . . . . .	25
4.2.3 ls: show me what you got . . . . .	26
4.3 Exercises . . . . .	27
4.4 Summary . . . . .	28



<b>5</b>	<b>Working with files and directories</b>	<b>29</b>
5.1	Examining files	29
5.1.1	cat: viewing short files	29
5.1.2	less: viewing large files	30
5.1.2.1	FASTA file format	30
5.1.3	head and tail: viewing the start or end of files	31
5.1.3.1	FASTQ file format	32
5.1.4	wc: counting lines	32
5.2	Editing files	33
5.3	Moving things around	34
5.3.1	cp: copying files and directories	34
5.3.2	Intermezzo: globbing and wildcards	34
5.3.3	mv Moving or renaming files and directories	35
5.4	Creation and destruction	36
5.4.1	Creating files	36
5.4.2	mkdir: creating directories	36
5.4.3	rm: removing things	36
5.5	Exercises	37
5.6	Summary	38
<b>6</b>	<b>More advanced commands</b>	<b>39</b>
6.1	Searching in files: grep	39
6.2	Tabular data: cut	41
6.2.1	SAM file format	42
6.3	File sizes: du	44
6.4	Compressed files: gzip	44
6.4.1	BAM file format	45
6.5	Downloading files: wget	46
6.6	Retrieving file names: basename	46
6.7	Sorting and removing duplicates: sort and uniq	47
6.8	Other commands	47
6.9	Exercises	48
6.10	Summary	48
<b>7</b>	<b>Streams, redirection and piping</b>	<b>50</b>
7.1	Streams	50
7.1.1	Redirecting output	50
7.1.1.1	VCF files	52
7.1.2	Redirecting errors	53
7.1.3	Input redirection	53
7.1.4	Overview of redirection operators	53
7.2	Piping	54
7.3	Further reading	56
7.4	Exercises	56
7.5	Summary	57
<b>8</b>	<b>Variables, loops and scripts</b>	<b>58</b>
8.1	Variables	58
8.2	Loops	58
8.3	Shell scripts	60
8.3.1	Creating a bash script	61
8.3.2	Running scripts	61
8.3.3	Making scripts executable	62
8.4	Overview of variant calling pipeline	63



8.5	A few final useful scripting concepts . . . . .	65
8.5.1	Command substitution . . . . .	65
8.5.2	Parameter expansion - manipulating strings . . . . .	66
8.6	Exercises . . . . .	66
8.7	Summary . . . . .	67
<b>II</b>	<b>R introduction</b>	<b>68</b>
<b>9</b>	<b>Other R resources</b>	<b>70</b>
9.1	Cheatsheets . . . . .	70
	<b>Appendices</b>	<b>71</b>
<b>A</b>	<b>Various Unix topics</b>	<b>71</b>
A.1	Tips and hints . . . . .	71
A.2	Overview of special syntax . . . . .	73
A.3	Understanding superusers, root and sudo . . . . .	74
A.4	What is \$PATH? . . . . .	74
A.5	Dealing with file permissions . . . . .	75
A.6	Working with remote machines via SSH . . . . .	75
A.7	Further reading . . . . .	76



# About

Welcome to the Bioinformatics Module of the Fighting Malaria Across Borders (FiMAB) international training programme, created by the Institute of Tropical Medicine (ITM) in Antwerp (Belgium) and supported by VLIR-UOS. Its primary goal is to support the implementation of targeted sequencing assays (in particular, AmpliSeq) to strengthen malaria molecular surveillance and help guide national control programmes. In conjunction with laboratory training, this bioinformatics course is intended to allow young academics around the globe to become familiar with molecular surveillance as a key activity to monitor transmission, sources of epidemics and the emergence and spread of drug resistance mutations in the *Plasmodium* parasite.

## Scope of the course

This course aims to provide an overview of key bioinformatics aspects related to performing molecular surveillance research in *Plasmodium*. It is divided into the following sections:

1. Introduction to the Unix shell (command line interface)
2. Introduction to R
3. Population genetics and molecular epidemiology for *Plasmodium* [to follow later]

Section 1-2 are online self-paced modules, whereas section 3 will include classroom lectures and practical sessions. Evaluation exercises will be conducted on [the ITM course page](#).



## Acknowledgements

Development of this course was supported by [VLIR-UOS](#) and the [Institute of Tropical Medicine](#).



## License

This work is licensed under [CC BY-NC-SA 4.0](#), which means that you are free to copy, redistribute or adapt any of its contents, provided that you **do** reference us (see [citation information below](#)) and the original license, **do** indicate if modifications were made, **do** distribute those contributions under the same license and **do not** use the material for commercial purposes.



# Introduction

## What is bioinformatics anyway?

Bioinformatics is a scientific discipline that develops and utilises computational methods to analyse large amounts of biological data, typically molecular sequence data such as DNA, RNA and amino acid sequences, as well as annotations describing those sequences. It spans the entire spectrum of collecting, storing and annotating data, to modelling, predicting and discovering new biological insights, from the level of individual molecules, to cells, organisms and populations. Bioinformatics is a highly interdisciplinary field that brings together many different types of researchers, including biologists, computer scientists, statisticians, clinicians and even chemists and physicists. A closely related term that you might see being used interchangeably with bioinformatics, is *computational biology*. However, there is no consensus on a clear distinction between the two<sup>1</sup>; like so many things in nature (e.g. the concept of a species), scientific disciplines are often hard to define and delineate.

<sup>1</sup> Source: <https://vcftools.sourceforge.net/VCF-poster.pdf>

A few examples of popular topics in bioinformatics are:

- Population genetics and molecular epidemiology (what we will be focusing on in this course)
- Analysis of gene/protein expression and regulation, e.g. in disease models (for RNA-seq spatial and single-cell approaches are become more prevalent)
- Structural bioinformatics, e.g. predicting the structure of proteins
- Network and systems biology: mapping and analysing the relationships between interacting biomolecules such as proteins, metabolites and their signal cascades, e.g. gene regulatory networks)
- Comparative genomics and phylogenetics: studying the ancestry of species, genes or entire genomes through time and space
- Genomic annotation and cataloguing genetic mutations and associated diseases in databases

## Computational thinking

We hope that this course can teach you a few computational thinking and problem solving skills that will help you along your bioinformatics journey. The learning curve in computational biology can be quite steep at times and the path is littered with arcane commands and obtuse syntax, but as you practice the concepts introduced in this course on your own, your command-line efficiency will improve and you will start to spot similarities across different types of environments and languages. Through this course, we hope to arm you with the necessary skills to make tasks like running custom analysis scripts or installing bioinformatics software seem a little less daunting.





## **Part I**

# **Introduction to the Unix shell**



# 1 What is a CLI?

This section of the course will introduce you to the general concept of the Unix command line interface (CLI) - as opposed to the graphical user interface (GUI) that you are familiar with - and Bash, one of the most ubiquitous Unix shells.

## 1.1 Learning objectives

- Knowledge of what a Unix shell and the CLI are and why/when they can be useful.
  - Setting up your own Unix environment.
  - Familiarity with basic `bash` commands for e.g., navigation, moving/copying and creating/deleting/modifying files and directories.
  - Introduction of a few more advanced commands and concepts like redirection, piping and loops.
  - First look at scripts and how they can be used in the context of DNA sequencing pipelines for variant calling.
- 

## Resources

This section of the course draws inspiration from the following resources:

- [Conor Meehan's UNIX shell tutorial](#) (CC BY-NC-SA 4.0)
  - [Mike Lee's Unix Crash Course](https://doi.org/10.21105/jose.00053) (<https://doi.org/10.21105/jose.00053>)
  - [Data Carpentry's Introduction to the Command Line for Genomics](https://doi.org/10.5281/zenodo.3260560) (<https://doi.org/10.5281/zenodo.3260560> CC-BY 4.0)
  - [Ronan Harrington's Bioinformatics Notebook](#) (MIT)
  - [A Primer for Computational Biology by Shawn T. O'Neil](#) (CC BY-NC-SA)
  - [Eric C. Anderson's Bioinformatics Handbook \(Chapter 4\)](#) (MIT)
  - [Course on UNIX and Genomic Data](#)
- 

## 1.2 What is Unix?

Unix is a family of operating systems, with one of their defining features being the *Unix shell*, which is both a **command line interface** and **scripting language**.

In simpler terms, shells look like what you see in the figure below and they are used to *talk* to computers using a CLI - i.e., through written text commands - instead of via a **graphical user interface** (GUI) where you primarily use a mouse cursor.



```
pmoris@3RRG9Y3: ~$ echo "This is the Bash shell!"
This is the Bash shell!
(base) pmoris@3RRG9Y3:~$ ls -l ./miniforge3/
total 116
-rw-r--r-- 1 pmoris pmoris 2473 Sep  4 01:12 LICENSE.txt
drwxr-xr-x 2 pmoris pmoris 4096 Nov 15 14:21 bin
drwxr-xr-x 2 pmoris pmoris 4096 Nov 15 14:21 compiler_compat
drwxr-xr-x 2 pmoris pmoris 4096 Nov 15 14:21 conda-meta
drwxr-xr-x 2 pmoris pmoris 4096 Nov 15 14:21 condabin
drwxr-xr-x 5 pmoris pmoris 4096 Nov 20 15:26 envs
drwxr-xr-x 5 pmoris pmoris 4096 Nov 15 14:21 etc
drwxr-xr-x 26 pmoris pmoris 4096 Nov 15 14:21 include
drwxr-xr-x 17 pmoris pmoris 12288 Nov 15 14:21 lib
drwxr-xr-x 3 pmoris pmoris 4096 Nov 15 14:21 libexec
drwxr-xr-x 3 pmoris pmoris 4096 Nov 15 14:21 man
drwxr-xr-x 258 pmoris pmoris 36864 Nov 20 15:26 pkgs
drwxr-xr-x 2 pmoris pmoris 4096 Nov 15 14:21 sbin
drwxr-xr-x 17 pmoris pmoris 4096 Nov 15 14:21 share
drwxr-xr-x 3 pmoris pmoris 4096 Nov 15 14:21 shell
drwxr-xr-x 3 pmoris pmoris 4096 Nov 15 14:21 ssl
drwxr-xr-x 3 pmoris pmoris 4096 Nov 15 14:21 x86_64-conda-linux-gnu
drwxr-xr-x 3 pmoris pmoris 4096 Nov 15 14:21 x86_64-conda_cos6-linux-gnu
(base) pmoris@3RRG9Y3:~$
```

Figure 1.1: Bash shell in WSL

There exist many different flavours of Unix, collectively termed “*Unix-like*”, but the ones you will most likely encounter yourself are Linux (which itself comes in many different varieties we call distributions, e.g. Debian, Ubuntu, Fedora, Arch, etc.) and MacOS. These operating systems come with a built-in Unix shell. While Windows also comes with a command line interface (Command Prompt and PowerShell), it is not a Unix shell and thus uses different syntax and commands. We’ll dig into how you can get your hands on a Unix shell on a Windows machine in a later section. The most ubiquitous Unix shell is **Bash**, which comes as the default on most Linux distributions.

### 1.3 Why bother learning the Unix shell as a bioinformatician?

Even if you are primarily a wet lab scientist, learning the basics of working with CLIs offers a number of advantages:

- **Automation:** CLIs and scripting excel at performing repetitive tasks, saving not only time, but also lowering the risk of mistakes. Have you ever tried manually renaming hundreds of files? Or adding an extra column to an Excel spreadsheet with millions of rows?
- **Reproducibility:** reproducibility is key in science and by using scripts (and other tools like git, package managers and workflow systems) you can ensure that your analyses can be repeated more readily. This is in stark contrast to the point-and-click nature of GUIs.
- **Built-in tools:** the Unix shell offers a plethora of tools for manipulating and inspecting large (text) files, which we often deal with in bioinformatics. E.g., DNA sequences are often stored as plain text files.
- **Availability of software:** many bioinformatics tools are exclusively built for Unix-like environments.
- **Access to remote servers:** Unix shells (usually bash) are the native language of most remote servers, High Performance Computing (HPC) clusters and cloud compute systems.
- **Programmatic access:** CLIs and scripts allows you to interact in various ways (e.g., via APIs) with data that is stored in large on-line databases, like those hosted by NCBI or EBI.



As a concrete example of what we will be using the shell for, consider the task of processing hundreds of *Plasmodium* DNA sequencing reads with the goal of determining the genetic variation in these samples (e.g., the presence of SNPs). Suppose we were to do this in a GUI program, where we would open each individual sample and subject it to a number of analyses steps. Even if each step were to only require a few seconds (in reality, minutes or even hours...), this would take quite a long time and be prone to errors (and quite boring!). With shell scripting, we can automate these repetitive steps and run the analysis without requiring human input at every step. Some of the key techniques we will use for this are:

- Navigating to directories and moving files around.
- Looping over a set of files, calling a piece of software on each of them.
- Extracting information from a particular location in a text file.
- Compressing and extracting files.
- Chaining commands: passing the output of one tool to another one. E.g., after aligning reads to a reference genome, the resulting output can be fed to the next step of the pipeline, the variant caller.
- Etc.

## 1.4 Don't get discouraged

Learning to use the shell, or learning programming languages and bioinformatics skills in general, can be daunting if you have had little experience with these types of tasks in the past. Don't worry though, just take your time and things will become easier over time as you gain more experience.

We do not expect you to be able to memorize every single command and all of its options. Instead, it is more important to be aware of the existence of commands to perform particular tasks, and to be able to independently retrieve information on how to use them when the need arises.

Finally, the appendix (Section A.1) of the course also contains a bunch of tips and tricks to keep in mind while learning your way around the shell.

## 1.5 A note on terminology

You will often see the terms `command line` (interface), `terminal`, `shell`, `bash`, `unix` (or `unix-like`) being thrown around more or less interchangeably (including in this course). Most of the time, it is not terribly important to know all the minute differences between them, but you can find an overview here if you are curious: <https://astrobiomike.github.io/unix/unix-intro>.



## 2 Setting up your own Unix shell

In order to get started, you need to get access to a Unix environment of your own. You can either work [locally](#) on your own computer or you can use the [online environment](#) that we have created. The latter comes with a bash shell and will immediately give you access to a bunch of files that we will use throughout the course and exercises.

There are also online playgrounds/simulations available to try out the Unix shell, for example <https://sandbox.bio/playground>. These are great to learn and we highly recommend checking out some of the tutorials on there, but the downside of course is that they are not true environments and you cannot interact with your own files. You can use this site to follow along while learning some of the basic unix commands, but you will need to switch to a different option for the exercises eventually.

### Shortcut to online environments

- [Codespace](#)
- [Binder](#)

### 2.1 Online Unix environment

We have provided two different options for getting access to an online unix environment: through [Binder](#) (free, but less powerful) or through [GitHub Codespaces](#) (free for 60 hours per month).

Both options will launch an environment containing all relevant training files, based on [this GitHub repository](#).

#### 2.1.1 GitHub Codespaces

To access Codespaces, you will first need to create a GitHub account via <https://github.com/signup>. Just follow the instructions and be sure to enable one of the [two-factor authentication options](#) (via a TOTP app like Authy, Google Authenticator or Microsoft Authenticator, or via text messages), otherwise you might not receive access to Codespaces.

Afterwards, you can click [this link](#), optionally change the region to the one closest to you under *change options* (but leave the machine type on 2-core to remain eligible for 60 free hours!), and then press the *Create codespace* button.

Setting up the codespace can take a while, but eventually you will be greeted by a VSCode environment. The terminal (a bash shell) is accessible at the bottom (or by pressing the hamburger icon in the top left and selecting “new terminal”). This is where you will be able to run the Unix commands introduced in the next chapters.

You will only receive 60 hours of free usage of Codespaces per month. This means you should manually shutdown your codespace whenever you are done with it. Otherwise it will keep running for 30 more minutes (by default). Just closing your browser will not shut down the workspace. Instead, you need to manually shut it down from within the



## Create a new codespace

Codespace usage for this repository is paid for by pmoris

**Repository**  
To be cloned into your codespace: pmoris/FiMAB-bio...

**Branch**  
This branch will be checked out on creation: main

**Dev container configuration**  
Your codespace will use this configuration: Bioinformatics training environment

**Region**  
Your codespace will run in the selected region: Europe West

**Machine type**  
Resources for your codespace: 2-core

Create codespace

Figure 2.1: Creating a new codespace

## Setting up your codespace

```
✓ Image found.
🔧 Building container...
  Hide logs

#19 41.94 Unpacking libbinutils:amd64 (2.41.50.20231206-1) over (2.35.2-2)
...
#19 42.15 Preparing to unpack ../binutils-common_2.41.50.20231206-1_amd64
.deb ...
#19 42.21 Unpacking binutils-common:amd64 (2.41.50.20231206-1) over (2.35.
2-2) ...
#19 44.99 dpkg: libfreetype6-dev:amd64: dependency problems, but removing
anyway as you requested:
#19 44.99 libxft-dev:amd64 depends on libfreetype6-dev.
#19 44.99 libfontconfig-dev:amd64 depends on libfreetype6-dev (>= 2.8.1).
#19 44.99
(Reading database ... 26196 files and directories currently installed.)
#19 45.01 Removing libfreetype6-dev:amd64 (2.10.4+dfsg-1+deb11u1) ...
```

💡 Tip Customize your codespace using a devcontainer.json file. [Learn more](#)

Figure 2.2: Launching a new codespace

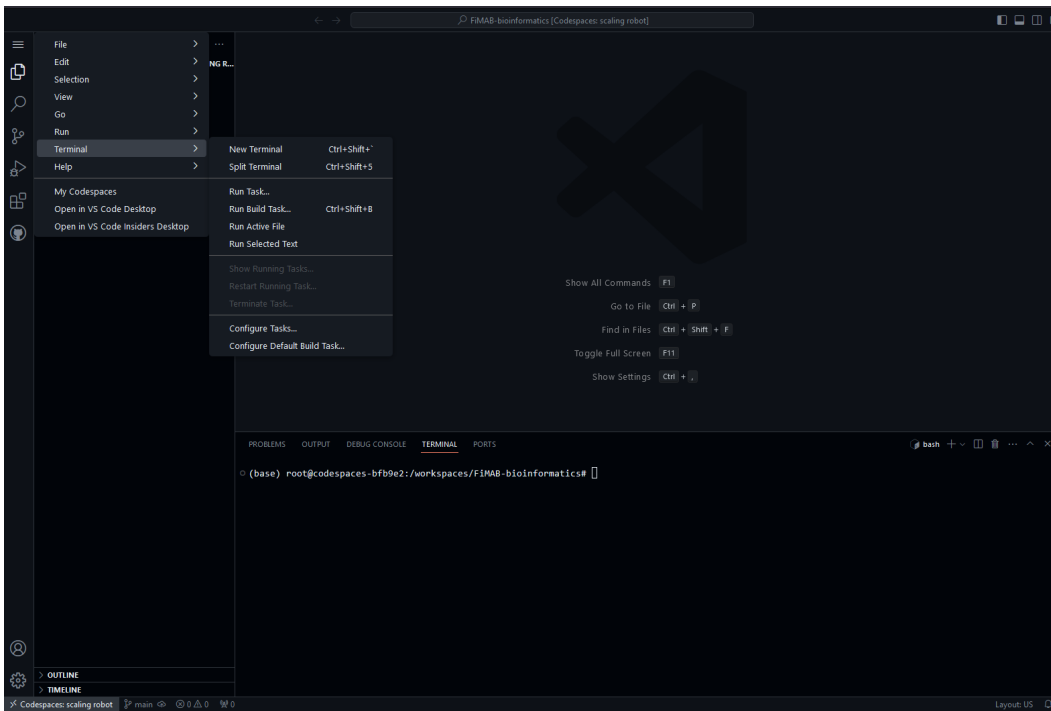


Figure 2.3: Terminal inside VSCode editor in GitHub Codespace

codespace (by clicking the >< button in the bottom left corner and selecting *stop current codespace*) or by browsing to <https://github.com/codespaces> and shutting it down from that page.

### **i** What are GitHub and git?

GitHub is a place to host code and software via a tool named *git*, which is a *version control system*. It allows you to keep track of the history of your code, easily revert changes and allows for collaborating with multiple people on the same project. We will not go into further detail on using version control, but for now just remember that it can play an important role in scientific reproducibility.

If you want to learn more about *git* already, you can have a look at the following resources:

- <https://happygitwithr.com/>
- <https://hwheeler01.github.io/CompBio/github/>
- <https://pmoris.github.io/git-workshop/> (self-promotion)

### **i** What is GitHub CodeSpaces

Similar to Binder, CodeSpaces are development environments that are hosted in the cloud. This is a paid service provided by GitHub/Microsoft, which offers 60 hours of free usage per individual per month. Instead of Jupyter notebooks, CodeSpaces use code editors, like VSCode and JetBrains IDEs, which come bundled with a *bash* terminal too.

You can find more info in the [GitHub CodeSpaces docs](#).



## 2.1.2 Binder

As an alternative to CodeSpaces, we have also created a Binder environment. It will operate similar to CodeSpaces and provide you with an online environment containing all of the required files as well as a Unix shell (bash).

### i What is Binder?

Binder is a service that allows people to share a customized compute environment based on a Git repository. It is mainly aimed at sharing Jupyter Notebooks (Python), but it also supports RStudio, Shiny and fortunately for us, a plain bash terminal too. You can find more info on the [Binder website](https://mybinder.org/).

To access the remote bash shell on Binder, browse to <https://mybinder.org/v2/gh/pmoris/FiMAB-bioinformatics/HEAD> and wait for the launcher to start. This process can take quite a while, so be patient.

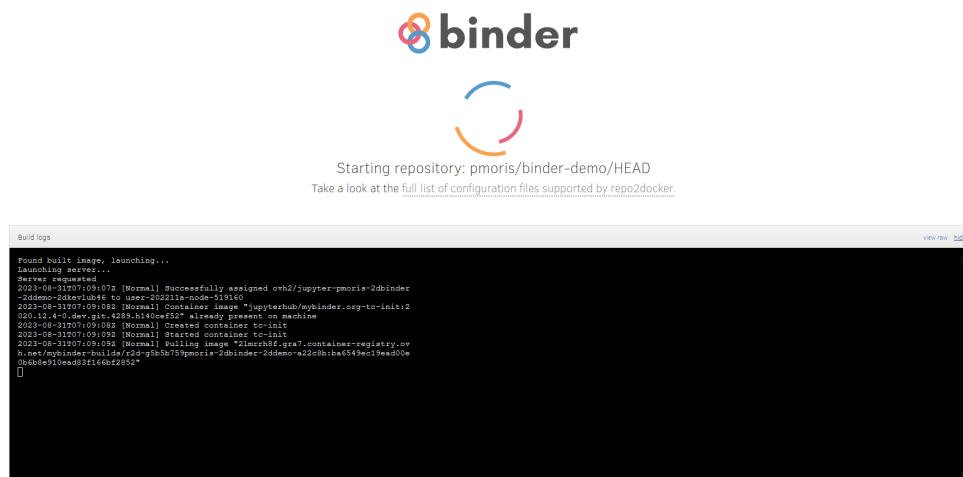


Figure 2.4: Launching the Binder environment

Eventually, you should be greeted by a screen (Jupyter Lab) with a number of launchers. Simply select the one labelled “Terminal” (a black square with a white \$) and you should be all set.

## 2.2 Local Unix environment

If you are using MacOS or Linux, then you will already have access to a Unix shell (either bash or zsh, which will mostly behave identical for our purposes). To access it, simply search for a program called `Terminal` (or search for anything resembling “command”, “prompt” or “shell”).

In case you are using a Windows machine, things are slightly more complex and different methods exist, each with their own pros and cons. You could use a fully-fledged virtual machine like [VirtualBox](#) to emulate a Linux machine within Windows. Or you could rely on the minimal bash emulator that comes bundled with [git for windows](#). However, nowadays we recommend that you use the [Windows Subsystem for Linux \(WSL\)](#), which was developed by Microsoft itself. In our opinion, it is one of the most polished methods to get access to a (nearly) full-featured Linux environment from within Windows, without the overhead of a full virtual machine or dual boot setup (dual boot means you install two different operating



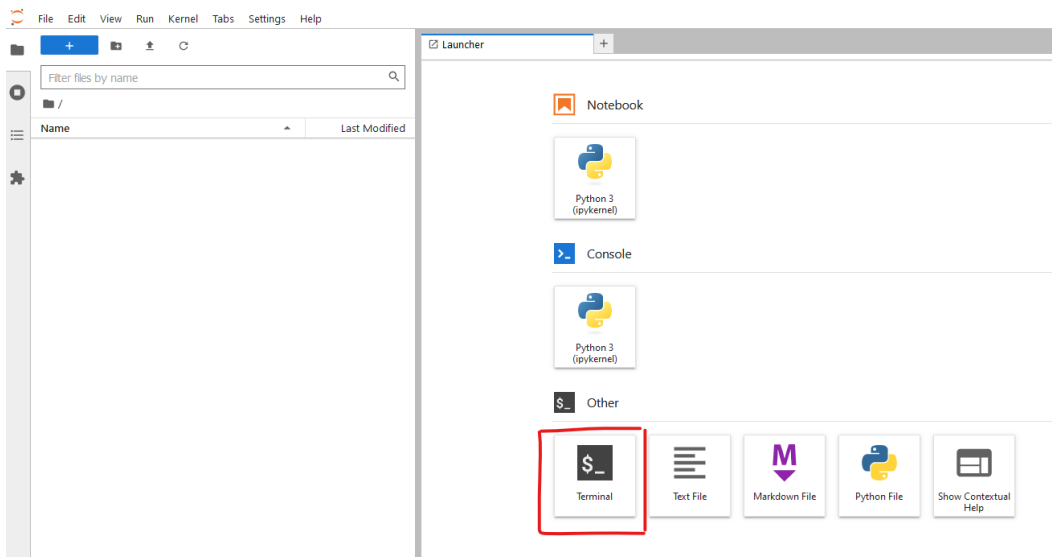


Figure 2.5: Starting a new Bash shell

systems on your machine, and you switch between them when booting). For instructions on how to set it up, you can refer to [this section](#).

### 2.2.1 Download the course files

Regardless of what type of local Unix environment you use, you will need to download the files that we will be using in our examples and exercises. You can do this directly on the command line or by manually downloading the files in the correct location.

1. Open your terminal and `cd` to a location where you want to place the training files.
2. Enter the command `git clone https://github.com/pmoris/FiMAB-bioinformatics.git`.
3. Afterwards, a new directory named `FiMAB-bioinformatics` will have been created.

It should look similar to this:

```
$ git clone https://github.com/pmoris/FiMAB-bioinformatics.git
Cloning into 'FiMAB-bioinformatics'...
remote: Enumerating objects: 7, done.
remote: Counting objects: 100% (7/7), done.
remote: Compressing objects: 100% (7/7), done.
remote: Total 7 (delta 0), reused 7 (delta 0), pack-reused 0
Receiving objects: 100% (7/7), done.
```

Alternatively,

1. Browse to <https://codeload.github.com/pmoris/FiMAB-bioinformatics/zip/refs/heads/main>
2. Save the `.zip` file in a directory accessible by your Unix environment. For Windows/WSL, the easiest option is to choose the Linux file system (e.g., `\\wsl.localhost\Ubuntu\home\pmoris`), which is accessible by clicking the Linux/WSL entry in your explorer.
3. Extract/unzip the file.

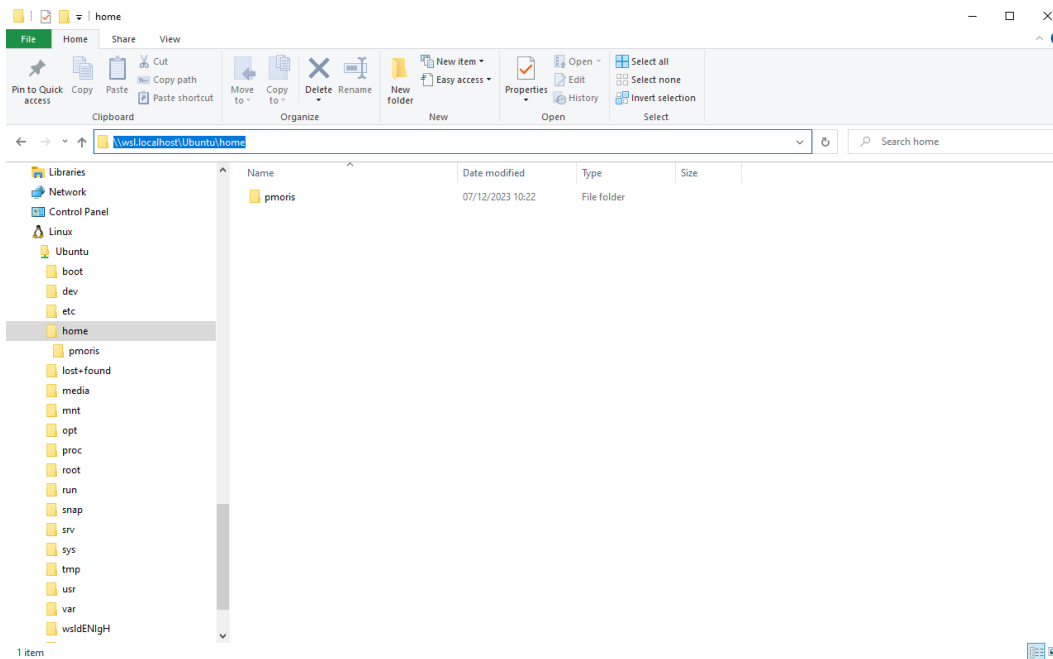


Figure 2.6: Linux file system inside Windows File Explorer

## 2.2.2 WSL installation

If you are using an updated version of Windows 10 (or 11), you should meet all the requirements and can simply follow the installation instructions listed here: <https://learn.microsoft.com/en-us/windows/wsl/install>. We recommend that you follow the instructions for WSL 2 (default), rather than the older WSL 1, and use the default Ubuntu 22.04 distribution (Linux comes in many different flavours, Ubuntu being one of the more popular ones).

Briefly:

- Open Windows PowerShell as administrator by right clicking your Windows Start Menu or searching for it in your list of applications.
- Type `wsl --install` and press enter.
- Afterwards, restart your PC.
- You can then launch WSL by searching for `wsl` or `Ubuntu` in your start menu.
- The first time you launch WSL, you will need to [configure it](#).

If you use software like RStudio or VSCode, you can tell these programs to use WSL as their built-in terminal from now on, instead of Command Prompt.

### WSL1 vs WSL2

WSL 2 is the newer version of WSL 1. For most tasks, WSL 2 tends to be much faster, hence why we (and Microsoft) recommend using it in favour of the previous version. However, WSL 2 is only faster when you interact with files that are stored directly on the WSL file system, rather than working directly on the Windows file system. More info on the distinction between these file systems can be found [further below](#) and in [Microsoft's WSL documentation](#).

You can switch between WSL1 and WSL2 on the fly by just calling `wsl --set-version <distro_name> 2` (or 1) in PowerShell, so feel free to experiment for yourself.

For a full overview of the differences, check out: <https://docs.microsoft.com/en-us/>



windows/wsl/compare-versions.

### 2.2.2.1 Configuring WSL

Microsoft also provides an excellent tutorial on setting up your WLS environment, which you can find [here](#).

---

After installing WSL and a Linux distribution, you will have access to it via its own built-in terminal emulator. It should be located in your Windows Start Menu with a name corresponding to the distribution that you installed, e.g. Ubuntu 20.04 LTS, or simply wsl.

The first time you run WSL, you will need to setup a Linux username and password. Note that while you are entering a password, nothing will appear on the screen, but this is intended (blind typing). The username will determine, among other things, the name of your home folder, whereas the password will grant you administrator rights (referred to as super users or admins in Linux land; the `sudo` command is used to invoke these rights, see Section A.3).

You will also need to upgrade the packages by running the following command: `sudo apt update && sudo apt upgrade`, followed by your password.

For more information, check the [docs](#).

### 2.2.2.2 Accessing files across the Windows and WSL file systems

#### **i** Note

Some of the information below might be a bit confusing at this point, but things should become more clear after working your way through the Unix section of this course.

Newer versions of WSL will automatically add a shortcut to the WSL file system in your Windows File Explorer (look for Tux, Linux' penguin mascot). The file path will look similar to `\\wsl$\\Ubuntu\\home\\<user name>\\Project`, indicating that Windows treats the WSL file system as a sort of network drive. You can also open a file location in Windows File Explorer from within a WSL terminal (e.g. after you browse to a particular directory `cd ~/my-project`) by simply using the command `explorer.exe .` (don't forget the dot!).

Vice versa, you can also access the Windows file system from within WSL because it is mounted under `/mnt/c`. So, you could for example do something like `cp /mnt/c/Users/<user name>/Downloads/file-downloaded-via-webbrowser ~/projects/filename`.

More information can be found in the [WSL documentation](#).



### 2.2.2.3 Windows Terminal

Even though WSL comes with its own terminal application, it is rather bare-bones and can make some operations like copying and pasting via CTRL+C/CTRL+V a bit tricky (you will need to use CTRL+SHIFT+C to copy and right mouse click to paste). Fortunately, Microsoft has also been working on a new terminal emulator that is much nicer to work with. Meet the [Windows Terminal](#).



## 3 Using the shell

### **i** Prior experience

You can skip this section and proceed directly to the exercises if you are already familiar with the basic syntax of Unix commands. We still recommend checking out the tips and hints in Section [A.1](#) and the special syntax overview in Section [A.2](#) though.

### 3.1 Interacting with the shell

When you launch your (Bash) shell, you will be greeted by what is called a shell prompt: a short snippet of text followed by a cursor, which indicates that the shell is waiting for input. The prompt can look different on different systems, but it often consists of your linux username followed by the name of your machine (like in the picture below) or sometimes just a single \$ symbol. When you see the prompt, you can enter commands interactively and execute them by pressing enter.



Figure 3.1: A bash shell prompt waiting for user input

Already note that you cannot use your mouse cursor to move around your terminal. You will need to use your arrow keys (or [shortcuts](#)) to move around while typing commands.

### 3.2 Command syntax

Unix commands generally follow the format:

`command [OPTIONS] argument`

where,

- `command` is the name of the (usually built-in) command that you want to execute.
- `[OPTIONS]` is a list of optional flags to modify the behaviour of the command. They are often preceded by a single (-) or double (--) dash.
- `argument` is a thing that your command can use. E.g., it can be a file name, a short piece of text (or *string*)

Try it yourself with the following command:

```
echo "Hello world!"
```



#### 💡 What did that do? (Click me to expand!)

`echo` is a command that simply prints a message to your screen (technically, to the *standard output stream* (stdout) of the terminal). `echo` is the command, telling the shell what we want to do. `"Hello world!"` is the target, in this case the message we want to print.

We place the message between quotes (`"`) because it contains spaces, and as you will see, spaces (and certain other special characters) can cause confusions. For now, just note that the message that gets printed, is whatever was written between the quotes, but not the quotes themselves.

```
$ echo "Hello world!"  
Hello world!  
$
```

### 3.3 Tips and hints

We have compiled a number of helpful tips in the appendix of this course (Section [A.1](#)), some of which will hopefully be helpful on your journey towards mastering the unix shell. For now, we recommend at the very least checking out the section on tab-completion and your command history. In fact, you can give it a try already. Just press the up arrow and see if you can recall your previous command! Next, try to type out `ec`, and press `<tab>`, to see auto-complete in action.



## 4 Navigating the Unix file system

### Prior experience

You can skip this section and proceed directly to the exercises if you are already familiar with the Unix directory structure and basic commands like `cd` and `ls`.

### Tip

We have provided a list of helpful tips and hints in the appendix: Section [A.1](#). Have a look already and refer back to it after you have worked your way through the next sections on navigation and basic commands.

Additionally, there is an overview of some of the most common symbols that are used by the Unix shell here: Section [A.2](#).

### 4.1 Layout of the Unix file system

All files and directories (or folders) in Unix are stored in a hierarchical tree-like structure, similar to what you might be used to on Windows or Mac (cf. File Explorer). The base or foundation of the directory layout in Unix is the *root* (`/`) (like the root of a tree). All other files and directories are built on top of this root location. When navigating the file system, it is also important to be aware of your current location. This is called the *working directory*.

The address of a particular file or directory is provided by its *filepath*: this is a sequence of location names separated by a forward slash (`/`), like `/home/user1`. Note that this differs from the convention in Windows, where backslashes (`\`) are used in file paths instead.

There are two types of file paths: *absolute* and *relative* paths.

- Absolute file path: this is the exact location of a file and is always built up from the root location. E.g., `/home/user1/projects/document.txt`.
- Relative file path: this is the relative address of a file compared to some other path. E.g., from the perspective of `/home/user1`, the file `document.txt` is located in `projects/document.txt`.

#### 4.1.1 Home sweet home: `~`

Another important location is the *home* directory. In general, every user has their own home directory, found in `/home/username`. A frequently used shortcut for this is the tilde symbol (`~`). Depending on the current user, this will refer to a particular directory under `/home/..`

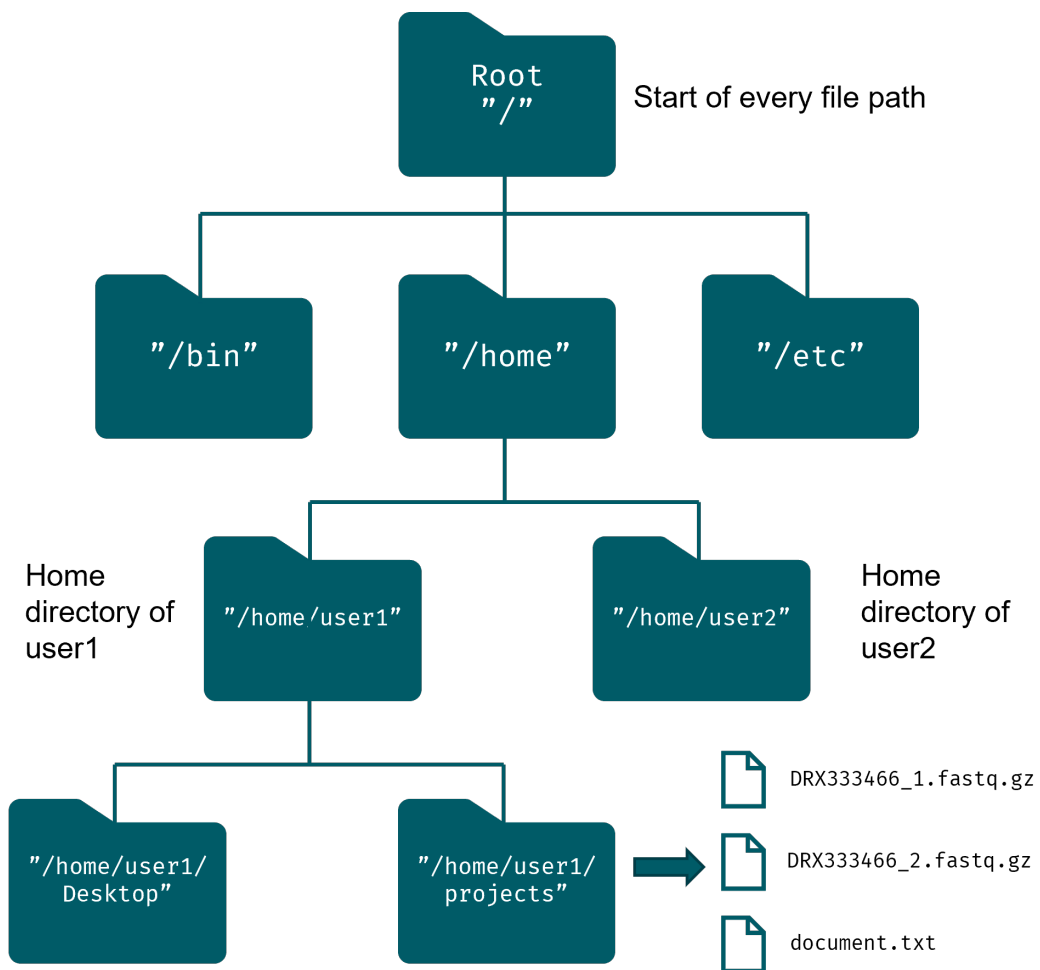


Figure 4.1: Overview of the Unix file system or directory layout





💡 How can user1 write the file path to document.txt using the ~ shortcut?

```
~/projects/document.txt
```

### 4.1.2 Where am I? . shortcuts

The dot (.) also has an important function in file paths:

- . represents the directory you are currently in, i.e. the working directory.
  - E.g., while inside the projects directory, any files inside can be accessed using either filename or ./filename.
- .. represents the parent directory of the working directory.
  - E.g., from /home/user1/Desktop, the relative path to file document.txt can be written as ../projects/document.txt.
  - These expressions can be nested; while inside the projects directory, ../../user2 can be used to access the user2 home directory.

## 4.2 Moving around the file system

In this section we will introduce a few essential commands that allow you to navigate the file system: pwd, cd and ls.

### 4.2.1 pwd: avoid getting lost

pwd stands for *print working directory* and it does exactly that: it allows you to figure out where you are in the file system. For example, in the figure above, user1 would generally find themselves in their home directory upon login:

```
$ pwd
/home/user1
```

### 4.2.2 cd: on the move

Next, there is the cd command. This is used to move between directories (the name derives from *change directory*). Simply follow the command name by a file path to navigate there: cd <filepath>. To move from user1's home directory to the projects directory:

```
cd projects
```

Note that you can use the special symbols we saw earlier as navigational shortcuts:

Command	Result
cd ~	Change to home directory (/home/username)
cd ..	Change to parent directory (e.g., go up 1 directory)
cd /	Change to the root location



Command	Result
---------	--------

### 4.2.3 `ls`: show me what you got

Finally, we have the `ls` command. Its name stands for *listing* and it will list the names of the files and directories in the current working directory. The basic structure of the command `ls [OPTIONS] <target>`, with `<target>` being an optional path to a directory.

To continue upon our previous example, from inside `/home/user1/projects` we would see:

```
$ ls
DRX333466_1.fastq.gz  DRX333466_2.fastq.gz  document.txt
```

Note that we did not specify a path, in which case `ls` will just list the contents of the current working directory. If we do specify a path, we will of course be shown the contents of that particular location:

```
$ ls /home
user1  user2
```

By default, the files and directories are listed in alphabetically order and depending on your terminal settings, files and directories might even be colour-coded differently.

`ls` also comes with a few handy optional flags to modify its behaviour:

Command	Result
<code>ls -l</code>	Show detailed list view
<code>ls -hl</code>	Show detailed list view and print file sizes in a human readable format
<code>ls -a</code>	List all files and directories, including <i>hidden</i> ones
<code>ls -lha</code>	Combine all options into one command
<code>ls --help</code>	Show more information on the <code>ls</code> command and its options
<code>ls</code>	

#### What are hidden files?

Earlier, we mentioned that `.` is used to refer to the current working directory, but it actually has a second function as well. Any file or directory name that starts with a dot (like `/home/user1/.ssh`) will be hidden and not displayed by default when using `ls`, hence the need for the `-a` flag.

Linux often hides system or configuration files to avoid cluttering up your (home) directory. We will not deal with hidden files directly in this course, but one of the situations where you might encounter them are when modifying your `.bashrc` file (e.g., when creating custom functions, aliases or tweaking your `PATH` Section A.4) or when managing SSH keys for remote server access Section A.6).

The `ls -l` command is particularly useful, because it shows all types of additional information.



### 💡 What do the different columns in the output of `ls -l` represent?

```
$ ls -l
total 83764
-rw-r--r-- 1 pmoris pmoris 14367565 Dec  7 09:39
↳ 3B207-2_S92_L001_R1_001.fastq.gz
-rw-r--r-- 1 pmoris pmoris 16622378 Dec  7 09:39
↳ 3B207-2_S92_L001_R2_001.fastq.gz
-rw-r--r-- 1 pmoris pmoris 13592342 Dec  7 09:39
↳ MRA1242_S28_L001_R1_001.fastq.gz
-rw-r--r-- 1 pmoris pmoris 15821981 Dec  7 09:39
↳ MRA1242_S28_L001_R2_001.fastq.gz
-rw-r--r-- 1 pmoris pmoris 12131772 Dec  7 09:39
↳ NK6_S57_L001_R1_001.fastq.gz
-rw-r--r-- 1 pmoris pmoris 13226198 Dec  7 09:39
↳ NK6_S57_L001_R2_001.fastq.gz
```

The first column represents the permissions of the files/folders. In a nutshell, these determine things like who can read or write (= modify, including deletion) particular files. There is a column for the owner, a group of users and everyone else. There is more info in the appendix (Section A.5). The next column showing a 1 for each entry, you can ignore for now (they represent hard links, a concept we will not dive into). The two names in the following columns are the *user* and the *group* owner of the file. Next is the size of the file in bytes. If we had used the `-h` flag, the size would have been shown in KB, MB or GB instead. Next we have the time of the last modification and finally the name of the file/directory.

## 4.3 Exercises

1. Navigate to your home directory and list all the files and folders there. Try typing the path with and without using the `~`. Rely on tab-completion to assist you and avoid typos (Section A.1).
2. Print the name of the current working directory to your screen.
3. List the contents of the `./training/data/fastq/` directory of the course files, without first moving there. Experiment with absolute and relative paths.
4. What is the most recent modification date of the file `Homo_sapiens.GRCh38.dna.chromosome.Y.truncated.fa` found in the `./training/unix-demo/` directory?
5. Try to search for the file `penguins.csv`: what is the absolute path to it on your machine?
6. Navigate to `./training/data/fastq/` to make it your working directory (double check using `pwd`!). What is the relative path to the `penguins.csv` file from here?
7. Suppose your working directory is still `./training/data/fastq/`. What will the result of `pwd` be after running each of the following commands in succession?
  - `cd ../`
  - `cd ../unix-demo/`
  - `cd files_to_loop_through/../../data/..`
  - `cd /`
  - `cd ~`



## 4.4 Summary

### Overview of concepts and commands

- *Absolute* versus *relative* file paths
- Root (/) and home directory (~)
- . represents the current working directory
- .. represents the parent directory
- pwd: print the path of the current working directory
- cd <path>: navigate to the given directory
- ls <path>: list files and directories in the given location
- Hidden files contain a . at the start of their name and are not visible by default



## 5 Working with files and directories

### Prior experience

You can skip this section and proceed directly to the exercises if you are already familiar with basic commands like `cp`, `mv`, `less` and `nano`.

### Tip

Remember that we have provided a list of helpful tips and hints in the appendix: Section A.1.

### 5.1 Examining files

#### 5.1.1 `cat`: viewing short files

The most basic command for viewing a file is the `cat <file>` command. It simply prints all of the contents of a file to the screen (= *standard output*).

```
$ cd training/unix-demo
$ cat short.txt
On the Origin of Species
```


```
BY MEANS OF NATURAL SELECTION,
```

```
OR THE PRESERVATION OF FAVOURED RACES IN THE STRUGGLE FOR LIFE.
```

```
By Charles Darwin, M.A., F.R.S.,
```

```
Author of "The Descent of Man," etc., etc.
```

```
Sixth London Edition, with all Additions and Corrections.
```

 Try using `cat` on the file named `long.txt` and see what happens (Click me to expand!)

The entire file (in this case, the entirety of the Origin of Species by Charles Darwin) is printed to the screen. This works, but is not very easy to navigate. Especially if you consider the fact that this text is still just tiny compared to some of the files that we deal with in bioinformatics; it is only ~0.03% of the size of the (rather short) human Y chromosome (~60 Mbp) that we will look at next.

While `cat` is very useful, it is clearly not suitable for large text files. Since long files are very prevalent - and not only in bioinformatics - we need an alternative. Enter the `less` command.



## 5.1.2 less: viewing large files

This tool is suitable for streaming very large files which would otherwise crash a normal text editor or program like Excel. `less` will open the contents of the file in a dedicated viewer, i.e. your terminal and prompt will be replaced by a unique interface for the `less` tool. You can exit this interface by pressing `q`.

Using `less`, we can have a look at the (truncated) version of the human Y chromosome (in FASTA format):

```
$ less Homo_sapiens.GRCh38.dna.chromosome.Y.truncated.fa
```

Figure 5.1: Opening a FASTA file in `less`

### Navigating inside `less`

- Use arrow keys to navigate. `space` and `b` can also be used to go forward and backwards, and `page up`/`page down` work as well.
- Press `g` to jump to the start of the file
- Press `G` (`shift + g`) to jump to the end of the file
- Type `/` followed by a string to search forward (`?[string]` for backwards search) and `n/N` for the previous/next match
- To exit, press `Q`
- Use the help command for more info: `less --help`

### 5.1.2.1 FASTA file format

#### DNA sequence file formats: FASTA

The FASTA file format (usually denoted by a `.fa` or `.fasta` file extension) is very common in bioinformatics. As you can see, the FASTA files contain a long stretch of nucleotides, which in our case represent the sequence of the human Y chromosome (or at least the first ~6,000,000 basepairs). The sequence itself is usually broken up over multiple lines. At the very top of the file there is a header or identifier, which always starts with the `>` symbol, followed by a short description. FASTA files can store



one or multiple sequences, each with their own header.

FASTA files are a type of text-based or *plain text files*, meaning that we can simply read them using a tool like `cat` or `less`. This seems obvious, but we will later encounter another file type, namely *binary files*, where this is not the case.

FASTA files are commonly used in genomics to store the reference genome of the organism we are studying. A reference genome can be used as a template to which we map (or align) new DNA sequence reads we have generated.

```
>Pf3D7_01_v3 | organism=Plasmodium_falciparum_3D7 |
version=2020-09-01 | length=640851 | S0=chromosome
TGAACCCCTAAAACCTAAACCCTAAACCCTAAACCCTGAACCCCTAAACCCTGAACCCCTAAA
CCCTAAACCCTGAACCCCTAAACCCTAAACCCTGAACCCCTAAACCCTGAAACCTAAACCCT
GAACCCCTAAACCCTGAACCCCTGAACCCCTAAACCCTAAACCCTAAACCCTAAACCCTGAACCC
CTAAACCCTGAACCCCTGAACCCCTAAACCCTGAACCCCTAAACCCTAAACCCTGAACCCCTAA
ACCCTGAACCCCTAAACCCTAAACCCTGAACCCCTGAACCCCTAAACCCTAAACCCTAAACCCT
TAAACCCTAAACCCTGAACCCCTAAACCCTAAACCCTAAACCCTAAACCCTGAACCCCTTACT
TTTCATTTCTTCTTCTTATCTTCTTACTTTTCATTCTTTACTCTTACTTACTTAGTCTTA
CTTACTTACTCTTACTTACTTACTCTTATCTTCTTACTTTTCATTCTTAGTCTTACTTA
...
```

### 5.1.3 head and tail: viewing the start or end of files

Sometimes we are not interested in viewing the entire file, but just the first few or last lines. The commands `head` and `tail` were created for exactly this use case. The basic usage is simply `head <filepath>`, but there again are a few optional flags that can alter the default behaviour.

Command	Result
<code>head file</code>	Print the first 10 lines of a file
<code>tail file</code>	Print the last 10 lines of a file
<code>head -n # file</code>	Display the first # lines of a file
<code>tail -n # file</code>	Display the final # lines of a file
<code>tail -n +2 file</code>	Display all lines except for the first one (i.e., perform tail, but start at line 2)

Let us inspect the first lines of one of the (uncompressed) FASTQ files in the `unix-demo` directory:

```
$ head PF0512_S47_L001_R1_001.fastq
@M05795:43:000000000-CFLMP:1:1101:16134:1717 1:N:0:47
TTGGTCAAGATCTTTTACATTCCATGCACACAAAAGAATTCTTCTACTTGTCTGATCCTTTTTCATTATATTTATTATCTTTTTTTATTTT
+
-8B<CGGFGGGDGGGGGEFGGGAFFFCFGGDC888C,CF@FFG,6FC,C,,,,<<,<<CCE,,<C,,,<6B,;B,<,<A@EEC+,:@,,S
@M05795:43:000000000-CFLMP:1:1101:20605:1731 1:N:0:47
GAAAAAGGAAGAGAATTGAACTTTTGGCAGCAAACTCAAACATTATAAGTGAAATTAAGATGCCCAAGTCTGTGCTCAATCTCATTTTTGT
+
-@CCCGGGGGGGGCGFG@<EFGGGGAFFFEggc8FEF9,,,CE,C,,,<,,,<C,,,<,,6;,,,;C,<,6;C,<,,,,:@,5,8+,
@M05795:43:000000000-CFLMP:1:1101:9135:1768 1:N:0:47
CGTTAAATCTTGCTCCTCATCACTACTAACCTTTTGTTCATTCTCATCACAAATATTATCCTTATCTTCATTATCTACTTCATCTACATTA
```



### 5.1.3.1 FASTQ file format

#### 💡 DNA sequence file formats: FASTQ

Aside from FASTA files, another typical DNA sequence format is FASTQ (extension `.fastq` or `.fq`), which is used to store the raw output of high-throughput sequencing (like AmpliSeq) in the form of short read fragments. Like FASTA, it is text-based format, but instead of just identifiers and sequences, it also contains quality scores associated with each nucleotide. Each read is described by four lines of text. A single read might look like this:

```
@SEQ_ID
ACTACTAGGATTGAGGACGTCCTCCCAACAGGGAGTTGGTTGGGCGCCCGTGCCGTCATGTCCGATCGCTATCTACGTCTAGTACTAGAGA
+
"H85<EI4A533D;E1A56C@@GHI=BFGIIH6;F=3: :HGF8C;9/>;EI?E4I(F?FID<CBAFFD69E:BB>+##<58H:/<>IE;881&
```

Line	Description
1	identifier: always starts with '@' and contains information about the read (e.g., instrument, lane, multiplex tag, coordinates, etc.)
2	The sequence of nucleotides making up the read
3	Always begins with a '+' and sometimes repeats the identifier
4	Contains a string of ASCII characters that represent the quality score for each base (i.e., it has the exact same length as line 2)

FASTQ files often come in pairs, which are usually named the same with a slightly different suffix (e.g., `sample_1_R1.fastq` and `sample_1_R2.fastq`). These pairs are reads of the same fragment in the opposite direction. In a nutshell, paired-end sequencing is used because the additional information provided by reads being paired can help with mapping repetitive regions of the genome.

💡 Try inspecting the contents of one of the `.fastq.gz` files in the `training/data/fastq` directory (Click me to expand!)

The `less` command most likely behaves as we expect it to, but if we were to try `cat`, `head` or `tail`, you would see a lot of gibberish being printed to your screen. The reason is that these FASTQ files are compressed using `gzip` (which is why the file extension ends in `.gz`). Because of this, they are no longer plain text files, but compressed binary versions. We will learn more about compressed files and how to deal with them in a later section of this course. In a nutshell though, compressed files either need to be unpacked or they require a tool that was designed to handle them (e.g., `zcat`, `zgrep`, `zless`. In most linux distributions, `zless` is called automatically when you try to `less` a file with the `.gz` extension, which is why the text seemed normal.

### 5.1.4 `wc`: counting lines

Sometimes we're not interested in the specific contents of a file, but only in how long it is in terms of text (not file size). For this we can use the `wc` command: it can count the number of lines, words and characters in a text file. By default, it prints all of this information, but by providing the `-l` flag, you can tell the command to only return the number of lines. Taking





the example of our FASTQ file again, we see:

```
# number of lines, words and characters
$ wc PF0512_S47_L001_R1_001.fastq
582940    728675 69598721 PF0512_S47_L001_R1_001.fastq

# number of lines only
$ wc -l PF0512_S47_L001_R1_001.fastq
582940 PF0512_S47_L001_R1_001.fastq
```

💡 How many reads are there in this FASTQ file? (Click me to expand!)

Each read in a FASTQ file consists of four lines (see Section 5.1.3.1). Therefore, we can simply divide the output of `wc -l` by four to figure out the number of reads. In this case:

$$\frac{582,940}{4} = 145,735$$

reads.

## 5.2 Editing files

You can edit files directly on the command line, i.e. without opening them in a text editor like Notepad(++) or VSCode, by using the `nano` command. This can come in quite handy in a variety of situations, like fixing small errors in your code before running it or to editing configuration files. Similar to `less`, `nano` will open a special editor interface where you can edit text files.

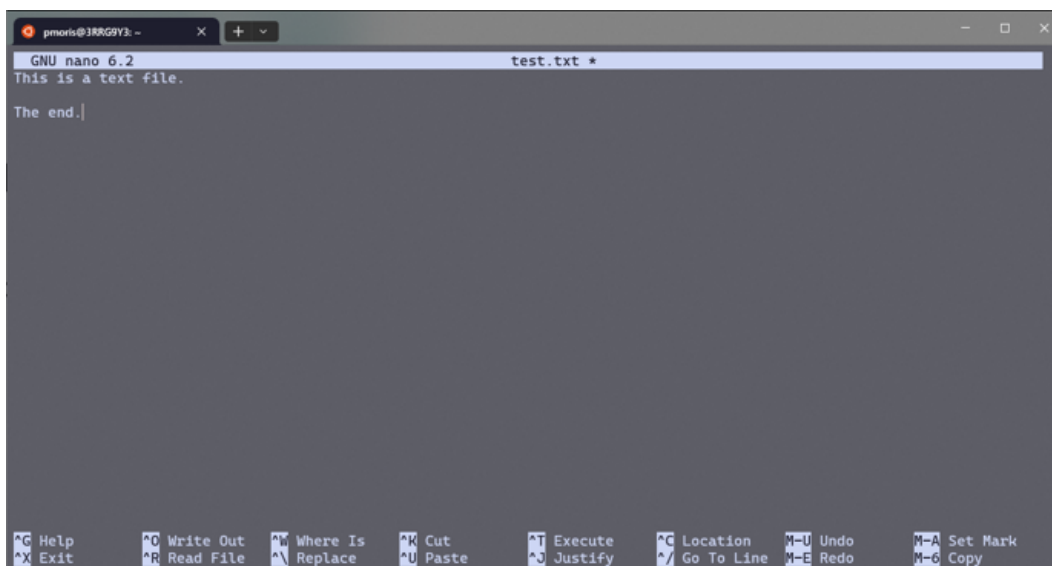


Figure 5.2: The nano text editor

💡 Navigating inside nano

- Your mouse pointer won't work. Use arrow keys to move instead.
- To save, press `ctrl+o`, followed by return/enter.
- To exit, press `ctrl+x`, followed by return/enter.



There exist many other editors, one of the most beloved, yet notorious ones, being vim. It is quite a bit more powerful, but also more complex. Even [closing vim](#) has become somewhat of a meme because it can be difficult to figure out (it's <escape> followed by :q and 'enter/return').

## 5.3 Moving things around

Now that we have spent some time on inspecting files, let us move on to moving them around.

### 5.3.1 cp: copying files and directories

cp stands for copy and it does exactly what it says on the tin. It can copy files, as well as directories to a new location. For files, the syntax is as follows:

```
cp path/to/source_file path/to/destination
```

Where source is the original file that you want to copy and destination is the new path where you want to place the copy. If the destination is a directory, the file will be placed inside of it with the same name as the original file. If the destination does not exist yet, it will be used as the new name for the copy.

When we want to move around directories instead of files, we need to add the -r flag (short for --recursive).

```
cp -r path/to/source_directory path/to/destination_file
```

You can even copy multiple files at the same time!

```
$ cp file_1 file_2 file_3 /path/to/destination  
  
$ ls /path/to/destination  
file_1 file_2 file_3
```

### 5.3.2 Intermezzo: globbing and wildcards

Now seems like a good time to introduce the concept of the globbing and wildcards. Globbing allows you to perform operations on multiple files. By providing specific patterns, the shell will be able to expand them into a list of matching file names. The patterns are built using wildcards, one of the most common ones being the asterisk \*.

How does this work? Well, \* can represent any number of other characters. For example, the string \*.txt can match all file names ending with .txt in your directory. Let's look at a concrete example, using the ls command we saw earlier:

```
$ ls  
Homo_sapiens.GRCh38.dna.chromosome.Y.truncated.fa  
↪ PF0512_S47_L001_R1_001.fastq files_to_copy files_to_delete  
↪ files_to_move long.txt short.txt penguins.csv
```



```
$ ls *.txt
long.txt  short.txt
```

As you can see, we can make `ls` list only those files that match a particular pattern, instead of showing all the files in the directory. What happens behind the scenes is that `*.txt` is expanded to `long.txt short.txt`. This means that the command that the shell eventually sees is actually `ls long.txt short.txt`.

Similarly, we can combine wildcards with the new `cp` command.

```
$ cp *.txt ..

$ ls ..
```

💡 What do you think this previous command will do? (Click me to expand!)

`*.txt` will be expanded to a list of all `.txt` files in the current working directory. The `cp` command will then try to copy each of those files to the destination, which is `..` in this case. As we saw before, `..` represents the parent directory of the current directory (see Section A.2).

This means that the command is equivalent to `cp long.txt short.txt /absolute/path/to/parent_directory` and will move all the `.txt` files in the current directory to its parent directory.

Another type of wildcard is `[...]`. This is used to supply a list of possible character matches. For example, the glob pattern `[bcr]at` would match `bat`, `cat` and `rat`.

There are a number of other wildcards, but even `*` alone will prove to be very useful. If you'd like to find out more, have a look [at this resource](#). Also note that globbing looks similar to *regular expressions*, but while related, these two concepts behave slightly differently. We will not dive into regular expressions here though, but we will mention them again when we talk about the search tool `grep`.

To summarise, globbing is an extremely powerful tool that will allow you to more easily target multiple files. We will rely on the power of globbing a lot going forward.

### 5.3.3 `mv` Moving or renaming files and directories

The `mv` (move) command behaves very similar to the `cp` command, the main difference being that the former allows you to move rather than copy files and directories. Also note that `mv` is used to *rename* files as well.

```
# move around/rename a particular file
mv <source_file> <destination_file>

# move a directory
mv <path/to/source_directory> <path/to/destination_directory>
```



## 5.4 Creation and destruction

We will end this section by teaching you how to create and delete files or directories.

### 5.4.1 Creating files

There are several ways of creating new files in Unix, but one of them is the `nano` command that we already introduced earlier. If you provide a file name that does not yet exist, `nano` will create the file for you.

```
$ ls

$ nano new_file.txt
# inside nano, use ctrl+x to save the file and then close the
↵ editor via ctrl+x

$ ls
new_file.txt
```

Another option is to use the `touch /path/to/file` command. This will just create a new empty file at the specified location.

### 5.4.2 `mkdir`: creating directories

`mkdir` stands for *make directory* and it does just that:

```
$ mkdir new_dir

$ ls
new_dir
```

One useful optional flag is `-p/--parents`: this allows you to create multiple nested (parent) directories in one go. For example, if we're inside an empty directory, we could call:

```
mkdir -p my/new/multi/level/directory
```

And all the intermediate directories would be automatically created.

### 5.4.3 `rm`: removing things

#### ⚠ Watch out...

Be careful while learning your way around the command-line. The Unix shell will do *exactly* what you tell it to, often without hesitation or asking for confirmation. This means that you might accidentally move, overwrite or delete files without intending to do so. For example, when creating, copying or moving files, they can overwrite existing ones if you give them the same name. Similarly, when a file is deleted, it will be removed completely, without first passing by a recycle bin.

**No matter how much experience you have, it is a good idea to remain cautious**



### when performing these types of operations.

For the purposes of learning, if you are using your own device instead of a cloud environment, we recommend that you work in a dedicated playground directory or even create a new user profile to be extra safe. And like always, backups of your important files are invaluable regardless of what you are doing.

The `rm` command (*remove*) is used to delete files and directories. Be warned though, once deleted, things are really gone. There is no recycle bin or trash folder where you can restore deleted items!

```
# for files:
rm <file path>

# for directories
rm -r <directory path>
```

For files, this works as expected, but for directories you need to provide the `-r` flag (or `--recursive`). This tells Unix to remove the directory recursively, i.e. all of its contents need to be removed as well. If you don't use this option, you will see the following warning:

```
rm directory
rm: cannot remove 'test/': Is a directory
```

### **i** Protected files

Sometimes, files will be protected and you will get another warning message when you try to remove them. If you are really sure that you want to delete them, you can type `y` and press enter. Alternatively, you can cancel the operation (by entering `n` or by pressing `ctrl+c`) and try again, but this time providing the `-f/--force` option.

```
# create a new empty file
$ touch protected-file
# change its permissions so that it is protected against
↪ writing and deleting (see appendix for more info on
↪ file permissions)
$ chmod a-w protected-file
# try to remove it
$ rm protected-file
rm: remove write-protected regular empty file
↪ 'protected-file'? n
# use the --force flag
$ rm -f protected-file
```

## 5.5 Exercises

1. Create a new directory named "my\_dir" inside the `./training/unix-demo` directory. Next, without using `cd` first, create another directory named `my_sub_dir`



inside of it. Finally, again without using `cd`, create a final directory named `my_sub_sub_dir` inside of that one.

2. Read the last 20 lines of the FASTA file in the `./training/unix-demo` directory.
3. Create a new text file named `lines` inside `my_sub_dir` using `nano`. Store the number of lines of the file `long.txt` inside. Then read it using `cat` and `less`.
4. Navigate to the `files_to_copy` directory and copy its contents to the `my_sub_dir` directory. What is the relative path of the destination to use?
5. Move the file under `files_to_move` to its parent directory.
6. Remove all the files under `files_to_delete` using a glob pattern.
7. Rename the directory `files_to_delete` to `empty_dir`.
8. List the contents of the `./training/unix-demo` directory.

## 5.6 Summary

### Overview of concepts and commands

- FASTA file format is used to store DNA sequences
- FASTQ file format is used to store sequence reads and their quality scores
- Compressed files (e.g., `.gz`) are smaller in file size, but are no longer plain text files and require special tools.
- The permission of files can be set to prevent users from reading, writing or (re)moving them.

Command	Result
<code>cat &lt;path/to/file&gt;</code>	print the content of files
<code>less &lt;path/to/file&gt;</code>	read the contents of (large) files in a special viewer
<code>head/tail &lt;path/to/file&gt;</code>	view the first or last lines of a file
<code>wc &lt;path/to/file&gt;</code>	display the line/word/character count of a file
<code>nano &lt;path/to/file&gt;</code>	open a file (or create a new file) in the nano text editor
<code>cp [-r] &lt;source&gt; &lt;destination&gt;</code>	copy a file/directory to a new location
<code>mv [-r] &lt;source&gt; &lt;destination&gt;</code>	move a file/directory to a new location (or rename it)
<code>rm [-r] &lt;path/to/file_or_directory&gt;</code>	permanently remove a file/directory
<code>mkdir &lt;path/to/directory&gt;</code>	create a new directory



## 6 More advanced commands

### Prior experience

You can skip this section and proceed directly to the exercises if you are already familiar with commands like `grep`, `sort` and `cut`.

### Tip

Remember that we have provided a list of helpful tips and hints in the appendix: Section A.1.

### 6.1 Searching in files: `grep`

Being able to search through (long) text files is incredibly useful in a wide range of scenarios. For this, we make use of the `grep` command. It is a very powerful and complex command, with many different options to tweak its behaviour, but even just the basic version can already be a lifesaver. The basic syntax is

```
grep "PATTERN" <path/to/target_file>
```

For example, we can look for the word “*evolution*” in the Origin of Species:

```
$ grep "evolution" long.txt
evolutionists that mammals are descended from a marsupial form;
↪ and if
At the present day almost all naturalists admit evolution under
↪ some
evolutionists; but there is no need, as it seems to me, to
↪ invoke any
Everyone who believes in slow and gradual evolution, will of
↪ course
of gradual evolution, through the preservation of a large number
↪ of
a strong disbeliever in evolution, but he appears to have been
↪ so much
historian will recognise as having produced a revolution in
↪ natural
the fact would be fatal to the theory of evolution through
↪ natural
the revolution in our palæontological knowledge effected by the
opposed to the admission of such prodigious geographical
↪ revolutions
has thus been arrived at; and the belief in the revolution of
↪ the earth
```







```
grep -c "Adelie" penguins.csv
152
```

In some situations we want to search through multiple files simultaneously. This is where the `-r/--recursive` flag comes in. It allows us to target a directory and search through all of its contents (including subdirectories). Let us try searching for the same DNA sequence as before, but this time targeting all the files in the `unix-demo` directory:

```
$ grep -r "AACCGGGGT" .
./PF0512_S47_L001_R1_001.fastq:AGCCATACCAAGACCACAATTCTGAAGAGGAAACAAAACAAAAAATAATTAAAA
./Homo_sapiens.GRCh38.dna.chromosome.Y.truncated.fa:AATAAACCGGGGTGATACCACCACTTCCAGGTCCCA
./Homo_sapiens.GRCh38.dna.chromosome.Y.truncated.fa:CTGGAGTCAGGACGTGAGCCGACTTGCTTAAAAATAATC
./Homo_sapiens.GRCh38.dna.chromosome.Y.truncated.fa:ACAGCTAACCGGGGTTTTAGTATATGTGCCACATCTCTGT
./Homo_sapiens.GRCh38.dna.chromosome.Y.truncated.fa:TATGATCGTGCCACTGCACTTCAACCGGGGTGACAAAGCC
```

💡 How would you search through `.txt` files only? (Click me to expand!)

Instead of using the `-r` flag, we can also rely on globbing again (see Section 5.3.2). To search for the string “needle” in all `.txt` files in a particular folder, we can do the following:

```
grep "needle" path/to/directory_with_txt_files/*.txt
```

We already mentioned *regular expressions* in the previous section: they allow you to search for particular patterns that can match more than one exact string of text. This is [tremendously useful](#), but we will not dive deeply into how they work during this course. If you are interested, you can check out an [excellent tutorial here](#).

One special pattern that we will introduce is the start or end of line anchor. You can search for patterns that start at the beginning of a line by prefixing the pattern with a `^` symbol. Similarly, you can search for patterns that occur at the end of a line by using a `$` symbol. For example, if we search for `grep ">@" read.fastq` in a FASTQ file, we will only retrieve lines that start with the `@` symbol, ignoring any `@` symbols that are used elsewhere on a line.

⚠️ The above command could still return matches that are not read headers. Why? (Click me to expand!)

As we saw in Section 5.1.3.1, every read in a FASTQ file is represented by four lines, the first of which is the header and which always starts with an `@` symbol. So far so good, but the fourth line, which contains the quality score of each base in the sequence, can also start with an `@` symbol, because `@` is just another score value that could happen to occur for the first base in the sequence.

We will see more elaborate use cases for `grep` when we introduce the Unix concepts of piping and redirection.

## 6.2 Tabular data: cut

We already encountered tabular data (the penguin dataset in `.csv` format) when talking about `grep`. Tabular data files like `.csv` are a very common format, and not just in bioin-



formatics.

#### 💡 Tabular data and .csv files

Tabular data files are usually plain text files, where each row corresponds to a record (e.g., an individual penguin), and each column represents a particular field (e.g., species, flipper length, body mass, etc.). The columns can be separated by different *field delimiters* or *separators*. In .csv files, these are usually commas (*comma separated values*), but they can also be TABS (.tsv) or semicolons (;).

A particularly useful Unix tool for manipulating tabular data files, is `cut`. It allows us to extract particular columns from these files. The syntax is as follows:

```
cut [OPTIONS] target_file
```

Option	Effect
<code>-d ",/--delimiter ";"</code>	Change the default delimiter (TAB) to another character like ,
<code>-f 1</code>	Select the first column
<code>-f 2,3</code>	Select the second and third column
<code>-f 1-3,6</code>	Select columns one through three and columns six
<code>--complement -f 1</code>	Select all columns <i>except</i> for the first one
<code>-r</code>	Recursive search through all files in a folder

### 6.2.1 SAM file format

Aside from .csv and .tsv files, there are many bioinformatics file formats that also follow a tabular lay-out. One of these is the [SAM file format](#).

#### 💡 SAM: Sequence Alignment/Map Format

The SAM file format is a tab-delimited text file that stores information about the alignment of sequence reads to a reference genome.

When considering the steps that are taken during the variant calling analysis of sequence reads, SAM files result when the reads inside a FASTQ file are processed by an alignment tool like `bwa` and `minimap`. These tools take each individual read corresponding to a particular sample and try to map them to their most likely position in a reference genome.

SAM files consists of an optional header section followed by an alignment section.

The header lines always start with an @ symbol and contain information about e.g., the reference genome that was used and the sorting of the alignments. This section is not yet tab-delimited.

The alignment section contains a tab-delimited line for each sequence read that was aligned to the reference. There are 11 mandatory fields, containing information on the position of the read in the reference genome (i.e., where it was mapped), the sequence itself, its quality score (= the quality of each base pair of the sequence during sequence calling, cf. FASTQ format), its mapping score (= how well the read aligned to the reference genome), etc.

Col	Field	Type	Brief description
1	QNAME	String	Query template NAME



2	FLAG	Int	bitwise FLAG
3	RNAME	String	References sequence NAME
4	POS	Int	1- based leftmost mapping POSition
5	MAPQ	Int	MAPping Quality
6	CIGAR	String	CIGAR string
7	RNEXT	String	Ref. name of the mate/next read
8	PNEXT	Int	Position of the mate/next read
9	TLEN	Int	observed Template LENgth
10	SEQ	String	segment SEQUENCE
11	QUAL	String	ASCII of Phred-scaled base QUALity+33

We will not dive into the details of each of these fields for now, but if you are interested you can check out the [official documentation](#) or this [useful write-up](#).

A typical SAM file will look something like this:

```
@SQ      SN:ref  LN:45
@SQ      SN:ref2 LN:40
r001     163     ref      7      30      8M4I4M1D3M      =
37       39      TTAGATAAAGAGGATACTG      *
XX:B:S,12561,2,20,112
r002     0       ref      9      30      1S2I6M1P1I1P1I4M2I
*        0       0        AAAAGATAAGGGATAAA      *
r003     0       ref      9      30      5H6M      *      0
0        AGCTAA  *
r004     0       ref      16     30      6M14N1I5M      *
0        0       ATAGCTCTCAGC      *
r003     16     ref      29     30      6H5M      *      0
0        TAGGC   *
r001     83     ref      37     30      9M      =      7
-39     CAGCGCCAT      *
x1       0       ref2     1      30      20M      *      0
0        aggtttttataaaacaaataa      ??????????????????????
x2       0       ref2     2      30      21M      *      0
0        ggtttttataaaacaaataatt      ??????????????????????
x3       0       ref2     6      30      9M4I13M *      0
0        ttataaaacAAATaattaagtctaca
????????????????????????????????
x4       0       ref2     10     30      25M      *      0
0        CaaaTaattaagtctacagagcaac
????????????????????????????????
x5       0       ref2     12     30      24M      *      0
0        aaTaattaagtctacagagcaact      ??????????????????????
x6       0       ref2     14     30      23M      *      0
0        Taattaagtctacagagcaacta ??????????????????????
```

This example has two header lines denoting the reference genome contigs and their lengths, followed by 12 aligned reads.

Now inspect the SAM file in `./training/unix-demo/PF0302_S20.sort.sam` and try to identify the different sections that make up the file.



## 6.3 File sizes: du

We already saw that the `ls -lh` can be used to figure out the file size of files in a particular directory. `du` is another tool to do this, but it operates on individual files or directories directly. Like `ls`, it also provides the `-h/--human-readable` option to return file sizes in KB/MB/GB, so it is generally recommended to always use this option. When used on a file, it will simply return its size, but when used on a directory, it will output information for all files, as well as the total file size of the entire directory (the final line of the output).

```
# targetting an individual file
$ du -h Homo_sapiens.GRCh38.dna.chromosome.Y.truncated.fa
5.9M    Homo_sapiens.GRCh38.dna.chromosome.Y.truncated.fa

# targetting a directory
$ du -h training/data/
284M    training/data/fastq
62M     training/data/reference
346M    training/data/
```

## 6.4 Compressed files: gzip

We already introduced the concept of file compression when talking about the FASTQ files in the `training/data/fastq` directory. As a reminder, compressed files are binary files (as opposed to human-readable plain text files) that are used to reduce the file size for more efficient storage. Many of the files that we use in bioinformatics tend to be compressed. Some of the tools we use, will not work on compressed files (e.g., try to `cat` a compressed file and see what happens), so we either need to 1) use specialized tools that expect compressed files as their input, or 2) decompress or extract the files first.

For gzipped (`.gz`) files specifically, we can do this via the `gzip` and `gunzip` commands. The former allows us to create a gzip-compressed version of the file, whereas the latter will extract one back to a plain text file. The basic syntax is `gzip/gunzip <path/to/file>`, but a very useful option is the `-k/--keep` flag. Without it, compressing a file would replace the uncompressed file with the new compressed one (vice versa: extracting would replace the compressed version with the extracted one), but when using the flag both files will be retained.

💡 Try compressing the FASTQ file in the `unix-demo` directory. By how much does its file size change? (Click me to expand!)

```
$ du -h PF0512_S47_L001_R1_001.fastq
67M    PF0512_S47_L001_R1_001.fastq
$ gzip --keep PF0512_S47_L001_R1_001.fastq
$ du -h PF0512_S47_L001_R1_001.fastq.gz
17M    PF0512_S47_L001_R1_001.fastq.gz
```

After compressing the FASTQ file with `gzip`, it shrunk to less than a third of its original size.

Do note that there exist other types of file compression besides `gzip`, like `.zip/.7zip`. In unix we also often make use of `tar` (which technically is not a compression tool, but a file



archiver). File compression and tar can even be combined, leading to files with suffixes like `.tar.gz`. This allows us to compress entire directories, instead of only individual files.

To extract these so called *tarballs*, we need to use the `tar` command:

```
# extract .tar.gz archive
$ tar -xzf tar_archive.tar.gz
tar_archive/
tar_archive/3
tar_archive/2
tar_archive/1

$ ls tar_archive
1 2 3

# create .tar.gz archive
$ tar -czvf new_archive.tar.gz <path/to/target_directory>
```

This command is notorious for how arcane its option flags are, but you can either try to remember it using a mnemonic (*“eXtract/Compress Ze Vucking Files”*, pronounced like a B-movie vampire) or the meaning of the individual flags (z tells tar that we are using gzip compression, -v stands for --verbose to make the command show more information and output, -c/x switches between compression and extraction mode, -f is the last option and points to the tar file). And of course, the correct syntax is only a [google/tldr](#) search or `tar --help` call away.

### 6.4.1 BAM file format

A final example of a binary file that you might encounter in the bioinformatic analysis of AmpliSeq sequencing is the BAM file format:

#### BAM: Binary Alignment Map

A BAM file is nothing more than a compressed, binary representation of a SAM file. It is used to reduce the file size of SAM files and also improve the speed of specific operations like sorting or retrieving information from the file.

Many bioinformatics tools can handle BAM file natively. However, similar to gzipped files, BAM files are binary and we can no longer preview them using tools like `cat` and `less`.

One of the tools that is often used to convert between the two types of alignment file formats is [samtools](#), which is maintained by the same group of people who manage the format specification for SAM/BAM files.

Aside from BAM, you might also encounter the [CRAM file format](#). This is a more recent and even more highly optimized compressed alternative for sequence alignments, but it is not as commonly used (yet).

`samtools` is a commonly used program to manage SAM/BAM files, which we will introduce later when discussing the structure of variant calling pipelines.



## 6.5 Downloading files: wget

`wget` is a command that allows you to download files from a particular web address or URL and place them in your working directory. While there are several optional flags, in its most basic form the syntax is simply: `wget URL`. This command is not only useful when automating certain tasks, but also crucial if you ever find yourself in a Unix environment that does not have a GUI at all (e.g., compute clusters or cloud servers).

💡 Try downloading the *Plasmodium falciparum* 3D7 reference genome in FASTA format from [PlasmoDB](#) and store it in `./training/data/results`. (Click me to expand!)

- At the top of the page, click on *Data -> Download data files*.
- Search for *falciparum 3D7* and then narrow down your search by selecting the most recent release and the FASTA file format. Alternatively, you can click on the [Download Archive](#) link in the top and navigate the file directory to the current release.
- The file name of the 3D7 reference genome in FASTA format is `PlasmoDB-66_Pfalciparum3D7_Genome.fasta`.
- Right click the file and copy its URL to your clipboard: [https://plasmodb.org/common/downloads/release-66/Pfalciparum3D7/fasta/data/PlasmoDB-66\\_Pfalciparum3D7\\_Genome.fasta](https://plasmodb.org/common/downloads/release-66/Pfalciparum3D7/fasta/data/PlasmoDB-66_Pfalciparum3D7_Genome.fasta).
- Create and navigate to the output directory (`mkdir -p ./training/data/results` and `cd ./training/data/results`)
- Download the file here using the command: `wget https://plasmodb.org/common/downloads/release-66/Pfalciparum3D7/fasta/data/PlasmoDB-66_Pfalciparum3D7_Genome.fasta`

Two optional flags that you might find useful are: 1) `-o` allows you to rename the download file, and 2) `-P <path/to/directory>` saves the file in a directory of your choice instead of the current working directory. Of course, these are just small convenient timesavers, since you can always `cd` to a particular location and use `mv` to rename the file afterwards.

Lastly, an alternative to `wget` that you might encounter at some point is `curl`. On the whole, it acts quite similar to `wget` for the most part.

## 6.6 Retrieving file names: basename

`basename` is a rather simple command: if you give it a long file path, it will return the final section (i.e., the file name).

```
# starting in the `training` directory
$ pwd
/home/pmoris/itg/FiMAB-bioinformatics/training

# get the file name for the reference genome we just downloaded
$ basename
data/reference/PlasmoDB-65_Pfalciparum3D7_Genome.fasta
PlasmoDB-65_Pfalciparum3D7_Genome.fasta
```

We can also use this command to remove a particular suffix from a filename:



```
$ basename PF0512_S47_L001_R1_001.fastq .fastq
PF0512_S47_L001_R1_001
```

At this point in time, it might not seem particularly useful to be able to extract the file name of a file, but when we introduce the concept of for loops and bash scripting, it will become more clear why this can be so useful.

## 6.7 Sorting and removing duplicates: sort and uniq

The final two commands that we will introduce are yet again tools to manipulate plain text files. The first is `sort`, which does exactly that you expect it to. It can sort all the rows in a text file. Its syntax is:

```
sort [OPTIONS] <./path/to/file>`
```

There are optional flags that allow you to choose the type of order to use (e.g., numerical `-n/--numeric-sort` instead of alphabetical), reverse the order of the sort (`-r/--reverse`) or ignore capitals (`-f/--ignore-case`).

The second command, `uniq`, is used to remove duplicate lines in a file. It also offers the option to count the frequency of each unique line.

```
# given the following file
$ cat file_with_duplicates.txt
a
a
b
b
b
b
b
c

$ uniq -c file_with_duplicates.txt
2 a
4 b
1 c
```

These two commands are often used in conjunction, because `uniq` on its own is not capable of filtering out identical lines that are not adjacent. So to truly remove all duplicate lines in a file, we would first need to sort it: `sort unsorted_file.txt | uniq`. Also note that the `sort` command also offers a flag `-u`, which will make it return only unique lines after sorting.

In the next section, we will introduce a method of combining these commands in a more convenient way than running them one by one and without needing to create any intermediary files.

## 6.8 Other commands

Of course, there exist many more Unix commands than the ones we introduced here. We will end this section by briefly mentioning two that you might run into at some point `awk`



and `sed`. Both of them allow you to search and replace patterns in text files, and with `awk` you can even perform more complex operations including calculations. We will not dive into them here, but keep them in the back of your mind for the future.

## 6.9 Exercises

1. Visit PlasmoDB again and find the download URL for the *Plasmodium vivax* P01 reference genome sequence in FASTA format. Download it via the command line and store it in `./training/data/reference`.
2. Report the file size of this reference genome in MBs.
3. Find out how many lines of text the file contains.
4. Search through the file for the `>` character, which is used to denote every chromosome/contig. Use a single command to count them.
5. Compress the FASTQ file `PF0512_S47_L001_R1_001.fastq` in the `unix-demo` directory using `gzip`.
6. Navigate to the directory `./training/data/fastq` and in a single command, extract the forward (`PF0097_S43_L001_R1_001.fastq.gz`) and reverse (`PF0097_S43_L001_R2_001.fastq.gz`) of the `PF0097_S43` sample, without removing the compressed files. Hint: use globbing!
7. Search both FASTQ files for the read fragment with identifier `@M05795:43:000000000-CFLMP:1:1101:21518:5742:N:0:43` using a single command.
8. Compare the file sizes of the two compressed and uncompressed FASTQ files.
9. Extract the columns containing the *island* and *flipper length* of each penguin from the `./training/unix-demo/penguins.csv` file.
10. Count how often the sequence `CATCATCATCATCAT` occurs in the FASTA file `./training/unix-demo/Homo_sapiens.GRCh38.dna.chromosome.Y.truncated.fa`.
11. Which command can be used to extract the name of the SAM file `./training/unix-demo/PF0302_S20.sort.sam` without the `.sam` suffix?
12. Which reference genome was used to create the SAM file?

## 6.10 Summary



### Overview of concepts and commands

- Tabular data: `.csv`, `.tsv`
- The SAM/BAM file formats store sequence reads aligned to a reference genome.

Command	Result
<code>grep &lt;pattern&gt;</code>	Search through a (very large) file for the supplied pattern
<code>du &lt;-h&gt;</code>	Check how much space a file or directory
<code>&lt;path/to/file_or_directory&gt;</code>	occupies





<code>cut -f [--delimiter ","] &lt;path/to/file&gt;</code>	Extract columns from tabular data using the specified delimiter
<code>gzip / gunzip (-keep) &lt;path/to/file&gt;</code>	Compress or extract a gzip compressed file (.gz)
<code>wget &lt;url&gt;</code>	Downloads a file from the URL to the current directory
<code>basename &lt;path/to/file&gt;</code>	Returns the name of the file without the path prefix
<code>basename &lt;path/to/file&gt;</code>	Returns the name of the file and remove the provided suffix



## 7 Streams, redirection and piping

### Tip

Remember that we have provided a list of helpful tips and hints in the appendix: Section A.1.

Up until now, most of the commands that we have used, printed their output directly to the terminal screen. But what if we want to save that output to a file? Or similarly, if we want to run program X on the output of program Y? This is where redirection and piping come into play. But first, we will have to briefly introduce the concept of streams.

### 7.1 Streams

The general flow of Unix commands is that we supply a specific input on the terminal, which is supplied to a command, and any output is printed back to the terminal screen. In other words, processes have three different *data streams* connected to them.



Figure 7.1: Unix input and output streams

The output for most commands, like that of `echo`, `cat` and `ls`, is called the *standard output* or *stdout*, and it is printed to our terminal screen by default. However, there exists another output stream in Unix, namely the *standard error* or *stderr*. This stream will contain error or warning messages produced by commands, and it is also printed to the terminal screen by default. There also exists an input stream, called *standard input* or *stdin*, which provides the data that is fed into a program.

Redirection and piping allows us to make these data streams go to or come from another file or process, instead of the terminal. Connecting these streams in different combination allows us to perform all kinds of useful operations.

#### 7.1.1 Redirecting output

One of the most common uses of redirection is redirecting the output to a file. For this, we make use of the greater than operator `>`:



```
$ ls > redirected_output.txt
some      files      in      a      directory

$ cat redirected_output.txt
some      files      in      a      directory
```

In the example above, the output of `ls` was not printed to the screen, but redirected to a file named `ls_output.txt`. Note that if the file does not exist, it will be created for us. However, if the file already exists, it will be overwritten (i.e., its contents will be removed entirely and replaced by our new output).

A related operator is `>>`. It will behave similar, with the difference being that `>>` will instead append its output to existing files, rather than overwriting them.

```
$ ls > redirected_output.txt
some      files      in      a      directory

$ echo "A second line!" >> redirected_output.txt

$ cat redirected_output.txt
some      files      in      a      directory
A second line!
```

Technically, whenever we use redirection, we are targeting a specific stream. Stdout is the default stream, so in the previous examples, `>` and `>>` were actually shorthand for `1>` and `1>>`.

As a more concrete example, we can use the redirection operator to store the results of a `grep` search.

```
$ grep "contig" ampliseq-variants.vcf > vcf-contigs.txt
$ cat vcf-contigs.txt
##contig=<ID=Pf3D7_01_v3,length=640851,assembly=PlasmoDB-44_Pfal Ciparum3D7_Genome.fasta>
##contig=<ID=Pf3D7_02_v3,length=947102,assembly=PlasmoDB-44_Pfal Ciparum3D7_Genome.fasta>
##contig=<ID=Pf3D7_03_v3,length=1067971,assembly=PlasmoDB-44_Pfal Ciparum3D7_Genome.fasta>
##contig=<ID=Pf3D7_04_v3,length=1200490,assembly=PlasmoDB-44_Pfal Ciparum3D7_Genome.fasta>
##contig=<ID=Pf3D7_05_v3,length=1343557,assembly=PlasmoDB-44_Pfal Ciparum3D7_Genome.fasta>
##contig=<ID=Pf3D7_06_v3,length=1418242,assembly=PlasmoDB-44_Pfal Ciparum3D7_Genome.fasta>
##contig=<ID=Pf3D7_07_v3,length=1445207,assembly=PlasmoDB-44_Pfal Ciparum3D7_Genome.fasta>
##contig=<ID=Pf3D7_08_v3,length=1472805,assembly=PlasmoDB-44_Pfal Ciparum3D7_Genome.fasta>
##contig=<ID=Pf3D7_09_v3,length=1541735,assembly=PlasmoDB-44_Pfal Ciparum3D7_Genome.fasta>
##contig=<ID=Pf3D7_10_v3,length=1687656,assembly=PlasmoDB-44_Pfal Ciparum3D7_Genome.fasta>
##contig=<ID=Pf3D7_11_v3,length=2038340,assembly=PlasmoDB-44_Pfal Ciparum3D7_Genome.fasta>
##contig=<ID=Pf3D7_12_v3,length=2271494,assembly=PlasmoDB-44_Pfal Ciparum3D7_Genome.fasta>
##contig=<ID=Pf3D7_13_v3,length=2925236,assembly=PlasmoDB-44_Pfal Ciparum3D7_Genome.fasta>
##contig=<ID=Pf3D7_14_v3,length=3291936,assembly=PlasmoDB-44_Pfal Ciparum3D7_Genome.fasta>
##contig=<ID=Pf3D7_API_v3,length=34250,assembly=PlasmoDB-44_Pfal Ciparum3D7_Genome.fasta>
##contig=<ID=Pf_M76611,length=5967,assembly=PlasmoDB-44_Pfal Ciparum3D7_Genome.fasta>
```

In the example above, we searched for lines containing the word `contig` in a `.vcf` file and stored the output in a file called `vcf-contigs.txt`, instead of just printing the output to the screen. But what is a VCF file anyway?



### 7.1.1.1 VCF files

#### Variant Call Format (VCF)

VCF is the de facto file format for storing gene sequence variation data. It is a plain text file with tab-delimited columns preceded by header lines with additional meta-data (starting with ##), similar to the structure of a BAM file.

An example of a VCF file is provided below<sup>1</sup>:

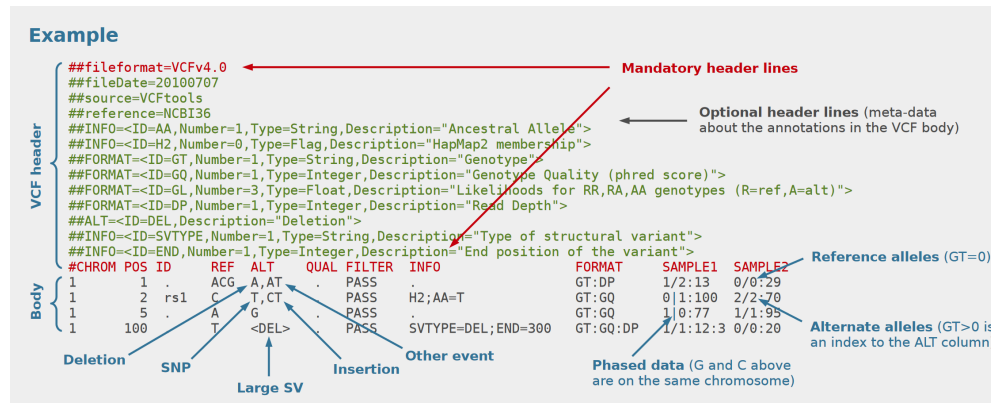


Figure 7.2: Description of Variant Call Format (VCF)

For now, the most important aspect to remember is that each of the lines in the body of the file store information on the presence of an indel at a particular position in the genome. There are eight mandatory fields for each variant, but the format is flexible and additional fields can be used to store extra information:

Field	Name	Description
1	CHROM	The name of the sequence (typically a chromosome) against which a variant is compared.
2	POS	The reference position.
3	ID	The identifier of the variant.
4	REF	The reference base occurring at this position in the reference sequence.
5	ALT	The list of alternative alleles found in your samples at this position.
6	QUAL	A quality score associated with the inference of the given alleles.
7	FILTER	Indicates which filters the variant has passed. Used as a quality control.
8	INFO	Additional info about the variant can be stored here as key-value pairs.

For a more in-depth view, we refer to the following excellent resources:

- <https://gatk.broadinstitute.org/hc/en-us/articles/360035531692-VCF-Variant-Call-Format>
- [https://davetang.github.io/learning\\_vcf\\_file/](https://davetang.github.io/learning_vcf_file/)

Similar to SAM/BAM files, there is also a binary version of VCF, namely BCF (Binary variant Call Format). Moreover, like samtools was developed to handle SAM/BAM files, there is dedicated program named bcftools that can be used to work with VCF/BCF files.



### 7.1.2 Redirecting errors

Without redirection, most commands print their error and warning messages to the terminal screen. If we use `>` (or `1>`) to redirect the output stream, the `stderr` will still be printed to the screen (and not be stored in the file that we are redirecting to). To redirect and store the error messages, we need to specify `stderr` as the stream via `2<`.

```
# cat on two files normally prints the output of both,
# but in this case fake_file does not exist
$ cat real_file fake_file
foobar
cat: fake_file: No such file or directory

# when we redirect stdout to a file
# stderr is still printed to the screen
$ cat real_file fake_file > redirected_output.txt
cat: fake_file: No such file or directory

# when we redirect stderr to a file
# stdout is still printed to the screen
$ cat real_file fake_file 2> redirected_errors.txt
foobar
```

If we want to **redirect both stdout and stderr**, we have to use a slightly more complex command:

```
# redirect both stdout and stderr to a file
$ cat real_file fake_file > redirection.txt 2>&1

$ ls redirection.txt
foobar
cat: fake_file: No such file or directory
```

### 7.1.3 Input redirection

The input stream (`stdin`) can also be redirected. Most commands like `cat` can open and process a file, but some commands cannot operate directly on a file. Instead, they need to be supplied with data directly. This is where input redirection (`<`) comes in. We have not yet encountered any commands that need to work on however, so the example below would work equally well without input redirection.

```
$ cat < input_file.txt
lines in
input_file.txt
```

### 7.1.4 Overview of redirection operators

---

<sup>1</sup>Source: <https://vcftools.sourceforge.net/VCF-poster.pdf>



Redirection operator	Result
<code>command &gt; file</code>	write stdout to file, overwriting if file exists
<code>command &gt;&gt; file</code>	write stdout to file, appending if file exists
<code>command 2&gt; file</code>	write stderr to file, overwriting if file exists
<code>command &gt; file 2&gt;&amp;1</code>	write both stdout and stderr to file, overwriting if file exists
<code>command &lt; file</code>	read input from file and pass it to command

## 7.2 Piping

Piping allows us to redirect the output of one command, to the input of another command. It is the more common way of redirecting the input stream. Pipes can chain multiple commands one after another so that a complex series of steps can be run in one go, without any intermediary output files. In its simplest form, piping looks like this:

```
command_1 | command_2
```

Where the first command produces some kind of output that can be used by the second one. For example, we could pipe the output of `ls` to `grep` to search through a list of directories and files:

```
$ ls
PF0080_S44_L001_R1_001.fastq.gz  PF0157_S55_L001_R2_001.fastq.gz
↪ PF0329_S56_L001_R1_001.fastq.gz
PF0080_S44_L001_R2_001.fastq.gz  PF0275_S68_L001_R1_001.fastq.gz
↪ PF0329_S56_L001_R2_001.fastq.gz
PF0097_S43_L001_R1_001.fastq.gz  PF0275_S68_L001_R2_001.fastq.gz
↪ PF0512_S47_L001_R1_001.fastq.gz
PF0097_S43_L001_R2_001.fastq.gz  PF0302_S20_L001_R1_001.fastq.gz
↪ PF0512_S47_L001_R2_001.fastq.gz
PF0157_S55_L001_R1_001.fastq.gz  PF0302_S20_L001_R2_001.fastq.gz

# search through the list of files in the
# current directory for a particular sample name
$ ls | grep "PF0275"
PF0275_S68_L001_R1_001.fastq.gz
PF0275_S68_L001_R2_001.fastq.gz
```

Note that in this case the syntax of `grep "pattern" <file>` changes slightly: we only supply the pattern, and the target file is now replaced by the stdout stream of `ls`. Piping makes it a lot more convenient to manipulate and chain commands in this manner. Doing the same thing without making use of piping takes a lot more work:

```
# write the output of ls to a file
$ ls > ls_output.txt

$ grep "PF0275" ls_output.txt

$ rm ls_output.txt
```



Clearly, the above approach is not very convenient. Especially if you consider the fact that you can chain as many pipes and redirections as you want: `command < input.txt | command | command > output.txt`. Let's take a look at a few more examples:

💡 Try counting the number of `grep` matches using a pipe instead of the `-c` flag. (Click me to expand!)

```
# count the number of matches in grep search results
$ grep "pattern" file | wc -l
```

`grep` will return a single line for each match that it finds. These lines are passed to the stdin of `wc -l`, which will count the number of lines.

💡 How can we remove non-consecutive duplicate lines from a file? (Click me to expand!)

```
# consider a file with non-consecutive duplicate lines
$ cat file.txt
a
a
b
c
b

# using uniq will only remove the consecutive duplicate
↪ line
$ uniq file.txt
a
b
c
b

# if we first sort and then run uniq, we get the desired
↪ output
$ sort file | uniq
a
b
c
```

`sort` will sort the file alphabetically, which results in all duplicates on consecutive lines. If we then pipe this output into `uniq`, all duplicates will be removed.

Here is another example of how to use pipes, this time applied to the SAM files that we saw earlier. SAM files sometimes contain a few lines of additional information - called the header and starting with an `@` symbol - before the start of the tab-delimited alignments (1 read per line). If we were to use the `cut` command to extract a particular column from such a file, the first few lines matching the header would cause problems, because these lines do not correspond to the tabular structure of the rest of the file.



💡 How can we use a pipe to only apply the `cut -f10` command to only the lines after the header? (Click me to expand!)

```
# count the number of lines in the header
$ grep -c "@" alignment-with-header.sam
2

# use tail -n +3 to print the output of the file starting
# the 3rd line, then pipe this output into the cut command
$ tail -n +3 alignment-with-header.sam | cut -f10
TTAGATAAAGAGGATACTG
AAAAGATAAGGGATAAA
AGCTAA
ATAGCTCTCAGC
TAGGC
CAGCGCCAT
aggttttataaaacaaataa
ggttttataaaacaaataatt
ttataaaacAAATaattaagtctaca
CaaaTaattaagtctacagagcaac
aaTaattaagtctacagagcaact
Taattaagtctacagagcaacta
```

We use the `tail` command to first extract the parts of the file that we are interested in, and then feed this output into the `cut` command to select a particular column (in this case, the 10th column corresponds to the sequence).

As a final tip on the usages of piping, consider that you can pipe the output of any command to `| less`. This is extremely convenient whenever the output of a particular command is too long and does not fit on your terminal screen. Of course, in some situations you are probably better off storing the output in a file using a `stdout >` redirection.

## 7.3 Further reading

Redirection operations and pipes can be combined in many more complex ways than what we saw here. For example, in case we want to redirect output to both a file and the terminal, we can make use of the `tee` command, [as described here](#). It is even possible to create more complex nested processes, where you feed the output of multiple different commands into a single command: `diff <(ls old) <(ls new)`; this is called *process substitution*.

You do not need to concern yourself with learning these more advanced concepts for the time being, but just keep in mind that whatever you want to do, the Unix shell likely offers a way of doing it.

---

## 7.4 Exercises

1. Search for the DNA sequence “*aacct*” in the truncated human Y chromosome FASTA file and store the output in a file called `aacct-hits.txt`.





2. Count the number of chromosomes in the *P. falciparum* reference genome fasta file.
  3. Store the chromosome identifiers of the *P. falciparum* reference genome fasta file in a file.
  4. Store the last 40 lines of PF0512\_S47\_L001\_R1\_001.fastq in a file named PF0512\_S47\_L001\_R1\_001.subset.fastq.
  5. How many penguin records are there for each island in penguins.csv? Hint: Try to do it in one go, without grep, by combining multiple pipes (cut, sort and uniq).
  6. How can you count the number of unique commands in your command history?
  7. Extract the header information from ampliseq-variants.vcf, sort it alphabetically, and store it in a file named vcf-header.txt.
- 

## 7.5 Summary



### Overview of concepts and commands

- VCF files are used to store genetic variant information
- Data streams: stdin, stdout and stderr
- Redirecting output to a file to replace (>) or append (>>)
- Redirecting errors to a file (2>)
- Piping the output of one command to the input of another command via |



## 8 Variables, loops and scripts

### Tip

Remember that we have provided a list of helpful tips and hints in the appendix: Section A.1.

### 8.1 Variables

*Variables* are placeholder names to refer to specific values. You can use them as shortcuts to refer back to a specific value or file path. Moreover, they are easy to set and re-use in bash scripts, which we'll introduce later. To set a variable, we assign a value to a name using the equals sign `=`. Afterwards, we can always recall the value via the variable's name, prefixed with a dollar sign `$`. While not strictly necessary, it is good practice to also enclose the name of the variable in `{}`, because it makes creating new variable names a lot easier during scripting.

```
$ my_value="Plasmodium falciparum"

$ echo ${my_value}
Plasmodium falciparum
```

As before, we use spaces around the value that we are assigned to the variable. This allows us to use spaces and other special characters inside our value. The variable in our example was a piece of text (a string), but we can also store things like integers (`x=101`) or booleans (`true/false`).

The main reason we introduced the concept of variables is because they play an important role in (for) loops, so let us move on to that topic now.

### 8.2 Loops

Loops provide a powerful method of repeating a set of operations multiple times. They are integral to automation and being able to process large numbers of samples in bioinformatics pipelines, but also very convenient for performing other tasks like renaming a bunch of files. The idea is a bit similar to the concept of globbing, but loops offer a lot more flexibility and control over the process.

The most common type of loops are probably *for loops*:

```
$ for nucleotide in A C T G \
> do echo ${nucleotide} \
> done
A
```



```
C  
T  
G
```

There are a number of different things going on here:

- The for loop consists out of three different sections:
  1. `for nucleotide in A C T G`: this tells bash that we want to start a for loop. It also defines the range of values that our loop will iterate over, in this case the characters A, C, T, and G. Finally, it creates a new variable called `nucleotide`. During every pass or round of the loop, its value will change to one of the values defined in the loop's range.
  2. `do echo $i`: this is the body of the loop. It always starts with `do` and is then followed by one or more commands. In the body, you can make use of the loop variable `$nucleotide`.
  3. `done`: this notifies bash that the body and loop definition end here.
- This is the first time that we see a multi-line bash command, where we split across new lines using a backslash symbol (`\`). We could have just as well written this statement on a single line (`for i in a c t g; do echo $i; done`), using colons (`;`) to mark the end of each section of the loop.

Instead of looping over a set of words, we can also loop over a range of values:

```
$ for i in {1..3} \  
> do echo ${i} >> loop.txt \  
> done  
  
$ cat for_loop.txt  
1  
2  
3
```

Also note that we used a different name for our loop variable this time around; you can use any name you like, but `i` is a very common placeholder. To make your scripts easier to read, it is good to stick with a reasonable short, but informative name.

💡 As a reminder, what would have happened if we had used `>` instead of `>>`? (Click me to expand!)

The for loop `for i in {1..3}; do echo ${i} > loop.txt; done` is basically equivalent to running the following three commands one after another:

```
echo 1 > loop.txt  
echo 2 > loop.txt  
echo 3 > loop.txt
```

As we saw in the previous sections, the redirection `>` will always overwrite the contents of its destination file. So in this case, the file would only contain the final number of the loop, namely 3. Loops always run in the order defined in their range.

Another common for loop pattern is the following one, which is used to loop over a set of files. It combines the for loop syntax with glob patterns (Section 5.3.2):



```
$ ls ./directory
sample_1_R1.fastq  sample_1_R2.fastq  text_file.txt

$ for fq in ./directory/*.fastq; do wc -l ${fq}; du -h ${fq};
  ↪ done
582940 ./directory/sample_1_R1.fastq
67M    ./directory/sample_1_R1.fastq
462334 ./directory/sample_1_R2.fastq
54M    ./directory/sample_1_R2.fastq
```

There are a few important things to note here:

1. The glob pattern `./directory/*.fastq` will be expanded by the shell to a list of all files ending with `.fastq`. Consequently, the for loop will only iterate of the FASTQ files and the `.txt` file is ignored.
2. Inside the execution statement of the loop, the current file is referred to via the variable `${fq}`.
3. We used a semicolon `(;)` to write the loop statement on a single line.
4. Unlike the previous example, the body of this for loop contains multiple commands: first the filename is printed to the screen using `echo`, then the number of lines in the file is printed to the screen.

---

Loops, combined with scripting, are incredibly useful when performing more advanced operations, like performing the bioinformatics analysis of DNA sequencing reads. For example, to process the DNA reads generated by an AmpliSeq assay and identify the genetic variants (variant calling analysis), the following steps would be performed by looping over the FASTQ file corresponding to each sample:

```
for every FASTQ file:
  1. Perform a quality control step
  2. Map the reads to the reference genome
  3. Call the variants in the alignment
```

---

There actually exists another type of loop, namely the *while loop*. These behave similar, but instead of going through a list of items or a range of numbers, the loop will continue for as long as a certain condition is met. You can find more information [here](#) in case you are interested.

## 8.3 Shell scripts

All the topics that we have covered so far, were performed interactively on the command line. However, we can also write scripts that contain a series of commands, loops and variables, which can be executed in one go. That way, you can queue up a bunch of long-running processes and don't need to stick around to start up each next step in the process. Moreover, it allows you to reuse the same set of operations in the future. Scripts can even be written in such a way that they can be called using different options, similar to how we can provide different optional arguments to bash commands.



### 8.3.1 Creating a bash script

Shell scripts are nothing more than **executable** text files written in a specific scripting language, in our case **bash**. We can write bash scripts in any type of text editor (like Notepad or VS Code), but we can also do it directly on the command line, by making use of an editor like **nano** or **vim**.

The only requirements for bash scripts is that the first line contains a shebang directive, like `#!/usr/bin/env bash`. When the script is executed, this line tells your machine to run the script using the **bash** interpreter (i.e., that it is a **bash** script and needs to be treated accordingly). By convention, shell scripts are saved with the `.sh` extension.

### 8.3.2 Running scripts

A very simple script might look a bit like this:

```
#!/usr/bin/env bash

echo "My first script!"
```

Inside a script, you can use any valid **bash** statement that would work on the command line. This includes all the commands we introduced up until now, structures like for loops, output redirection, pipes, etc.

To run a bash script, you can simply execute the **bash** command and point it to a script. If we save the two lines above to `script.sh` file and then executing it by running `bash script.sh`, the single command inside of it will be executed and printed to the screen:

```
$ bash script.sh
My first script!
```

Below is a slightly more complex example:

```
#!/usr/bin/env bash

echo "This is an example script."

echo "The script was executed from the directory: "
pwd

# this is a comment

echo "Running a for loop"

for i in {1..5}
do echo i
done

echo "This is a grep command"

# long commands can be split over multiple lines
# this grep command counts the number of times "tttataaaaaaac"
```



```
# occurs in the current directory and all of its subdirectories
# while ignoring case
grep -i \
    -r \
    "tttataaaaaaac" \
    .
```

#### Tip

Note that the indentation that we used is not strictly necessary for loops to work, but it does help with legibility and it is common practice to do this, especially in scripts. You can use hashes (“#”) to comment out a line. Use this to describe what your code is doing. Your future self will be grateful! You can also use backslash (“\”) to split a long command over multiple lines, making it easier to read your script. E.g., for listing `-options` on consecutive lines.

If we save and run the script above, we get the following output.

```
This is an example script.
The script was executed from the directory:
/home/pmoris/itg/FiMAB-bioinformatics/training/unix-demo
This is a for loop
i
i
i
i
i
This is a grep command
./PF0512_S47_L001_R1_001.fastq:CTAACTACAATGAAGACAAAAATATTATGTATATGTACCCAAATGAACCAAATTATAAGGA
./PF0512_S47_L001_R1_001.fastq:CTAACTACAATGAAGACAAAAATATTATGTATATGTACCCAAATGAACCAAATTATAAGGA
./Homo_sapiens.GRCh38.dna.chromosome.Y.truncated.fa:TGAGATTGGATTTTTAAACATTAATATGGCGTGTTACATT
./my_script.sh:# this grep command counts the number of times
↪ "tttataaaaaaac"
./my_script.sh: "tttataaaaaaac" \
Exiting script
```

#### Can you spot what is wrong with this script? (Click me to expand!)

The output of the for loop section of the script is five lines of the character “i”, which is not what we wanted. The problem is that inside of the body of the for loop, we used the character “i”, instead of referencing the for loop variable using `${i}`.

Lastly, note that not all scripts are bash scripts, you can also create R, Python or other types of scripts.

### 8.3.3 Making scripts executable

In the previous examples we ran scripts by invoking them using the `bash` command. However, we can also make the file *executable* by using the `chmod` command. Once a file has been marked as executable for a given user, it can be run by simply typing the file name, without any other command in front.



```
$ ls -l my_script.sh
-rw-r--r-- 1 pmoris pmoris 488 Jan  4 14:35 my_script.sh

$ chmod +x my_script.sh

$ ls -l my_script.sh
-rwxr-xr-x 1 pmoris pmoris 488 Jan  4 14:35 my_script.sh

$ my_script.sh
<script output>
```

You can read more about file permissions in Section A.5.

#### **i** Different methods of running script

There are actually a few different ways to execute a shell script:

1. Prefixing the interpreter: `bash path/to/script.sh`. For this method, we explicitly execute the `bash` command and point it to the location of a script.
2. Directly running a script in the current working directory: `./script.sh`. This option has two requirements: first, the file needs to be made executable (as we saw above). Furthermore, as a safety precaution, we need to prefix the name of the file with `./`, to let the system know that we are trying to run a script that resides in the current working directory, as opposed to one that is globally accessible.
3. Directly running a globally accessible script: `script.sh`. This option also requires the file to be executable, but on top of that it will only work for files that are found in the list of directories making up your `PATH` (see Section A.4). These are pre-defined (or custom) directories like `/usr/bin` that usually store all the common commands that we have used until now, as well as any other software that you might install yourself.

## 8.4 Overview of variant calling pipeline

The chart below provides a high-level overview of the steps involved in a basic variant calling pipeline. Throughout the course, we have already introduced some of the file formats that are used (green). We will explore these steps and the associated tools, in more detail during the in-person courses later on.

For now though, we have included a few example scripts that run through the above steps, to showcase how the building blocks that we have learned so far can be used to orchestrate a more complex analysis.

#### **i** Note

Inspect the scripts stored in `./training/scripts` and try to make sense of the general steps that they describe. You can gloss over the specifics of what the specialized commands like `fastqc` or `bwa mem` do; instead, try to focus on the general structure and syntax of the scripts, the use of patterns like for loops, directory navigation, how arguments are provided to commands, etc.

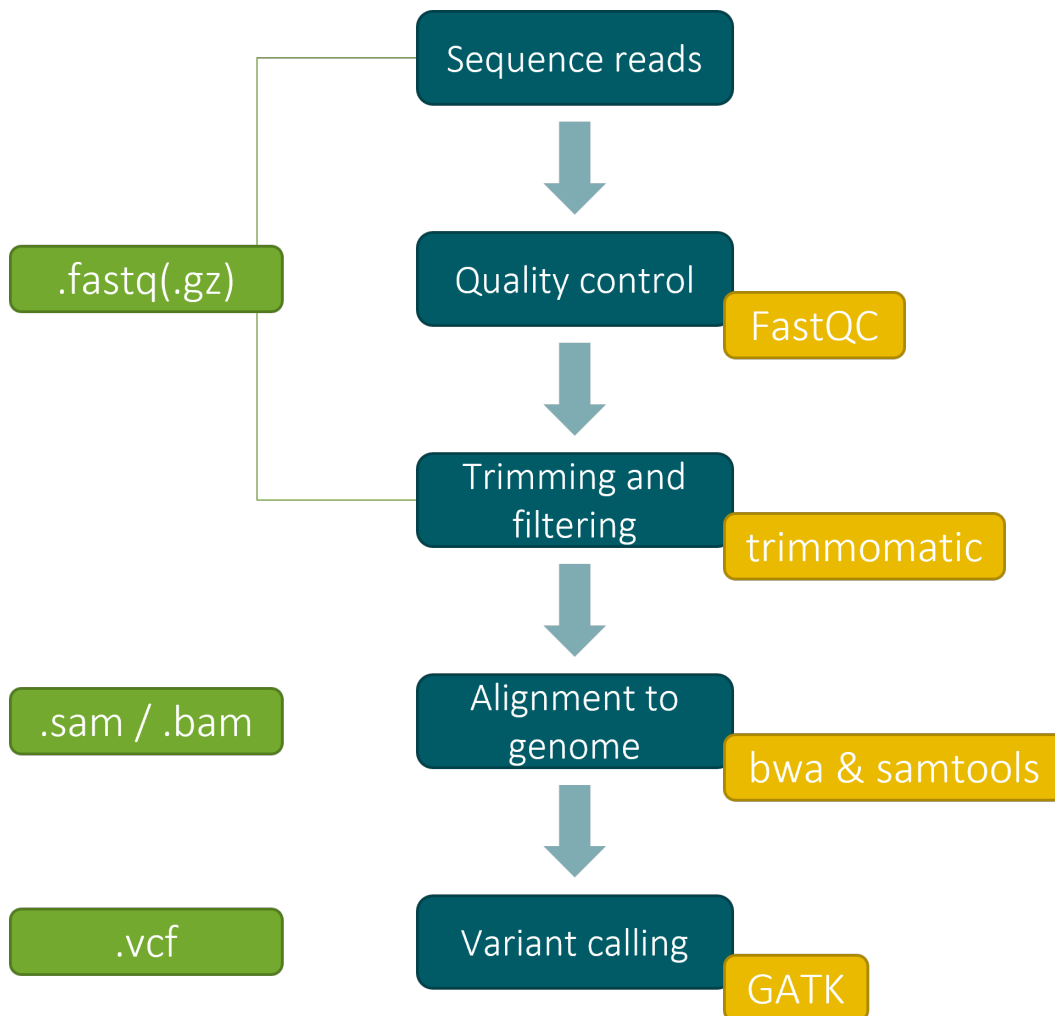


Figure 8.1: Variant calling pipeline





💡 What does the following line do? `sample_name=$(basename ${read_1} _R1_001.fastq.gz)` Hint: `$(command)` provided a method to run a command inside another statement, so try figuring out what the command between the brackets does first.

This line is present in most of the example scripts inside of a for loop that iterates over a set of FASTQ files, each time processing a pair of R1/R2 files.

```
# first navigate to the directory containing fastq files
cd ./training/data/fastq/

# loop over pairs of fastq files
for read_1 in *_R1_001.fastq.gz
do
    sample_name=$(basename ${read_1} _R1_001.fastq.gz)
    ...
```

During each iteration of the loop `${read_1}` will correspond to the file path of a specific FASTQ file.

When we call `basename` on it with the extra argument `_R1_001.fastq.gz`, we will receive back the name of the file with that suffix removed. E.g.:

```
$ basename PF0080_S44_L001_R1_001.fastq.gz _R1_001.fastq.gz
PF0080_S44_L001
```

The last step we do is running this command inside a command substitution: `$(command)`. When using a command substitution, the output of the command inside the brackets will just be passed along to the command line. In our case, we try to assign the value `$(basename ${read_1} _1.fastq.gz)` to a new variable named `sample_name`. The value will then contain the output of its command substitution, namely `PF0080_S44_L001`.

After running the above statement, we now are able to more easily construct the name of the first and second read pair:

```
$ echo ${read_1} ${sample_name}_R2_001.fastq.gz
PF0080_S44_L001_R1_001.fastq.gz PF0080_S44_L001_R2_001.fastq.gz
```

The first read we already had, but to create the second one we concatenate the sample name with the new suffix `_R2_001.fastq.gz`.

Paired sequence data is usually named in such a way that it allows to access pairs of files in this way.

## 8.5 A few final useful scripting concepts

### 8.5.1 Command substitution

In the callout box above we already introduced the concept of using *command substitutions* and the `basename` command to manipulate filenames and paths during a for loop, in order to be able to iterate over pairs of FASTQ files with similar, but predictable file names (e.g., `read_R1_001.fastq` and `read_R2_001.fastq`). To reiterate, the general concept is that we only loop over R1 of each pair, use `basename` to extract the filename without



the `R1_001.fastq.gz` suffix *and* without the preceding file path prefix, and then use this to construct the known file path to R2.

```
for read_1 in *_R1_001.fastq.gz
do
    sample_name=$(basename ${read_1} _R1_001.fastq.gz)

    command ${read_1} "/file/path/to/${sample_name}_R2_001.fast.qz"
...

```

A simpler example would be letting `echo` output the result of another command, like `date`:

```
echo "Today is $(date)."
```

You can also use it to store the output of a command in a variable, like we did for `basename`:

```
count=$(grep -c "#" variants.vcf)
```

## 8.5.2 Parameter expansion - manipulating strings

Bash contains all sorts of tricks to manipulate strings of text. These can be used in a similar way to `basename`, in order to strip of file extensions, parts of names or file paths and many other more complex tasks. For more details, we refer to <https://mywiki.woledge.org/BashFAQ/073>.

## 8.6 Exercises

1. Create a for loop over the files in `./training/unix-demo/files_to_loop_through` that prints the first line of each file to the screen.
2. Create a bash script that does the same thing.
3. On a single-line, run the script, but sort its output in reverse order (check `sort --help` to check how) and store the new output in a file called `loop_sort.txt`.
4. Create a bash script with a for loop that prints the name of each read file in the `./training/data/fastq` directory.
5. Modify the previous bash script so that it also 1) creates a single new directory named `fastq_meta` and 2) creates a new file in that directory, one for each FASTQ file, which contains two lines: the first with the number of lines in the FASTQ file and the second with its file size.
6. Create a bash script with a for loop that prints the sample name for each pair of reads in the `./training/data/fastq` directory (i.e., half as many names as in the previous exercise).
7. Create a bash script that:
  - Counts the number of header lines in `./training/unix-demo/ampliseq-variants.vcf` and store this number as a variable.
  - Extracts the contents of the VCF file after the header lines (i.e., the tabular section) and store it in a separate file.
  - Creates a for loop to extract the first eight columns and store them each in a separate file named `vcf_column_#.txt` (where `#` is the column number).



## 8.7 Summary

### Overview of concepts and commands

- Variables can be assigned via `name=value` and referenced via `${name}`
- For loops are used to iterate over a list of items or files
- Scripts can be used to combine multiple commands into a single set of instructions that can be re-used.
- Manipulating file names during a loop can be done using `basename` or parameter expansion.
- Command substitutions (`$(command)`) can be used to insert the output of a particular command (e.g., `basename`) into the middle of other commands or during the construction of new strings of text.



## **Part II**

# **R introduction**



As your first introduction to R, we will make use of [Data Carpentry's Intro to R and RStudio for Genomics course](#). This self-paced workshop walks through the basics of R - and its most popular IDE, RStudio - in the context of genomics. The focus lies on the basic syntax of R, wrangling tabular data using the [tidyverse](#) - where we will encounter the VCF file format once again (Section 7.1.1.1) - and visualization using [ggplot2](#).

The Data Carpentry workshop will sometimes refer to running RStudio in a cloud environment with pre-installed packages and files, but when going through these materials on your own, we recommend installing R and RStudio on your own machine. You can do so by following the instructions listed here: <https://rstudio-education.github.io/hopr/starting.html>. You will also need to download the following two files, which are used throughout the workshop: [combined\\_tidy\\_vcf.csv](#) and [Ecoli\\_metadata.xlsx](#).

If you do run into trouble installing R and RStudio on your own computer, you could instead:

- make use of the free version of Posit Cloud (formerly known as RStudio Server): <https://posit.cloud/plans/free>. However, note that you only receive a limited number of hours of use as a free user.
- Alternatively, you can again make use of GitHub codespace, but you will need to select a different variant this time: [RStudio Server](#).
- As a final option, you could make use of a binder RStudio instance, such as [this one](#).



## 9 Other R resources

Here is a collection of other useful resources for learning R.

- [Introduction to Data Science by Rafael A Irizarry](#): a clear overview of the most important aspects of data analysis in R, covering basic syntax, data manipulation using the tidyverse, visualizations and a few productivity tools.
- [The Epidemiologist R Handbook](#): an extensive R reference manual for various epidemiological analyses.
  - [Applied Epi Tutorials](#): they also offer free self-paced R tutorials accessible online (requires registration). If you prefer learning in an interactive manner, then these are highly recommended. Also see [swirl - interactive tutorials in R](#).
- [Software Carpentry's R for Reproducible Scientific Analysis workshop](#): this workshop goes over some of the other aspects of R that were not touched upon in the genomics workshop, e.g., data structures, loops, functions, etc.
- [Jenny Bryan's Stat 545](#): an accessible guide to data wrangling, exploration and analysis in R. It also covers a few important aspects for newcomers related to [managing projects and treating your code as absolute rather than R's history](#), as well as several other adjacent topics
- [R for Data Science by Hadley Wickham](#): this book is a bit more advanced, but it is an excellent resource nonetheless. Hadley Wickham is one of the driving behind tidyverse and the author of the amazing [Tidy Data paper](#).
- [The Pirate's Guide to R](#): slightly whimsical and quite thorough.
- [R Graphics Cookbook](#): tons of examples for almost any kind of graphical representation you can imagine.
- [ggplot2: Elegant Graphics for Data Analysis](#): a complete guide to the grammar of ggplot2.
- [PeterMac Data Science: Intro to R Tidyverse Course](#): another introduction to the tidyverse, this time using the analysis of RNA-seq data as an example.

### 9.1 Cheatsheets

Everyone loves cheatsheets!

- [ggplot2](#)
- [dplyr data transformations](#)
- [data tidying using tidyr and tibble](#)



## A Various Unix topics

### A.1 Tips and hints

#### Naming conventions and cases

Never (never!) use spaces in your file or directory names. This will only lead to pain... Instead, use hyphens (-) or underscores (\_) to separate words. E.g., `my_first_script` and `3B207-2_S92_L001_R1_001.fastq.gz`. Additionally, unlike in Windows, in Unix everything is case-sensitive. Thus, `/home/documents != /home/Documents`. Be mindful of this when naming or pointing to files/directories.

#### Autocompletion and command history

Avoid unnecessary typing and just make things easy for yourself! While typing commands in the shell, you can almost always use the tab key for **auto-completion**. This will automatically type out paths, file names or known commands. If there exist multiple matches, a single press of tab will not appear to do anything, but if you press the button twice, a list of possible options will appear on your screen. This is incredibly useful, not only for speeding things up, but also for avoiding typos when dealing with long or complex file names. An equally useful tool is your **command history**. While on the shell prompt, pressing the up arrow (↑) will bring up your most recent previous command. Pressing it again will cycle through the entire history, in reverse order. You can also search through your history by pressing `ctrl+r` allows you to search through your command history. Just start typing and you will see the search try to narrow down on the command that you are looking for. Once you find it, just press enter to run it directly or tab to copy it to your prompt (in case you still want to change it). The search form will look like this: `(reverse-i-search)`world': echo "Hello world!"`

#### Copying and pasting

Copying and pasting might work slightly different to what you are used to, depending on the terminal application that you are using. If `ctrl+c` and `ctrl+v` do not appear to work, you can try `ctrl+shift+c` and `ctrl+shift+v` instead. Often times, the mouse middle or right click can also be used for pasting. For the native WSL terminal specifically, you can refer to this site for more info: <https://devblogs.microsoft.com/commandline/copy-and-paste-arrives-for-linuxwsl-consoles/>

#### Don't panic when you lose control of your shell!

If a command seems to hang or get stuck, your terminal becomes unresponsive, or if you tried to print a very large text file to your screen, you can use `CTRL+C` to interrupt almost any operation and regain control.



Similarly, CTRL+D is an often used shortcut for exiting/logging out (e.g., when dealing with remote servers or nested shells).

In some cases, like when using an interactive terminal program such as the text editors `nano` and `vim` or a text viewer like `less`, you will only be able to exit them using that particular program's shortcut keys (CTRL+X, : followed by q and enter, and Q, for these applications respectively). For more info on terminal programs, check out [Fantastic terminal programs and how to quit them](#).

### ! Watch out...

Be careful while learning your way around the command-line. The Unix shell will do *exactly* what you tell it to, often without hesitation or asking for confirmation. This means that you might accidentally move, overwrite or delete files without intending to do so. For example, when creating, copying or moving files, they can overwrite existing ones if you give them the same name. Similarly, when a file is deleted, it will be removed completely, without first passing by a recycle bin.

**No matter how much experience you have, it is a good idea to remain cautious when performing these types of operations.**

For the purposes of learning, if you are using your own device instead of a cloud environment, we recommend that you work in a dedicated playground directory or even create a new user profile to be extra safe. And like always, backups of your important files are invaluable regardless of what you are doing.

### 💡 Google, -h/--help and comments are your friends.

At the beginning things will be awkward, so don't worry about having to search for the same information multiple times. That is all part of the learning process. Moreover, being able to retrieve information when you are in need of a particular command is more useful than memorizing everything.

It can still be a good idea though to keep a list of commands that you often use, but have a difficult time committing to memory.

Many commands will also display a short help text when called with the `-h/--help` flag. For some tools, you will need to call them without any arguments to display the help. Lastly, for some tools, you can use the `man <tool_name>` to open up an even more in-depth manual.

### 💡 Make your scripts easier to read by using comments and breaking up long lines

Remember that you can always write comments inside of your scripts by starting a line with `#`. That way you can add a short explainer or extra info to the different sections of a script. Code that seems clear while you are writing it, has the unfortunate tendency of becoming much more confusing when you refer back to it at a later time.





```
#!/usr/bin/env bash

#####
# Script to map fastq files to reference genome with bwa #
#####

# make sure to run this script from within the directory
↪ where it is stored!

# move to the directory containing the fastq files
cd ../data/fastq

# create a directory to store the results and store the
↪ path as a variable
output_dir="../results/bwa"
mkdir -p ${output_dir}

...
```

Additionally, you can break up long commands using a `\` to make them easier to read. You can do this both in scripts or on the command line. E.g.,

```
bwa mem \
  ../reference/PlasmoDB-65_Pfalciparum3D7_Genome.fasta \
  ${read_1} \
  ${sample_name}_R2_001.fastq.gz
```

**i** <https://explainshell.com/> & <https://tldr.inbrowser.app>

The first website is tremendously useful for figuring out what a command and all of its options mean. Whereas the second shows you a quick summary of the most command usages of a particular command.

Use both of these to your advantage! But do not forget that most commands also have a built-in help page that can be accessed using the `--help` flag (in some cases just typing the command without any arguments also shows some help information).

## A.2 Overview of special syntax

The table below gives you an overview of some of the special characters that we will encounter. You do not need to memorize them, but you can always refer back to this section if you see a symbol later on and are not quite sure what its purpose is.

Symbol	Name	Uses
/	Forward slash	File path separator or root location/file path
\	Back slash	Split long command to a new line and escape special characters (+ file path separator in Windows)
~	Tilde	Shortcut for home directory in file paths



Symbol	Name	Uses
	Pipe or vertical bar	Chains the output of one command to the input of another one (piping)
#	Hash	Part of the shebang at the top of scripts #! and used for comments in shell scripts
\$	Dollar sign	Used to access variables in bash
*	Asterisk or wildcard	Globbing operator
>	Greater than symbol	Redirect output of a command (>> redirect and append instead of overwriting)
<	Less than symbol	Redirect input to a command
.	Dot	In the context of a path, it represents the current working directory
..	Double dot	In the context of a path, it represents the parent directory of the working directory

### A.3 Understanding superusers, root and sudo

Superusers are special users on Unix systems with additional privileges (cf. Windows administrator). By convention, the default name for a superuser on Linux is `root` (don't confuse this with the root of the filesystem `/`). Superusers have access to all files and directories on the file system, including any critical components that make the system work and any files owned by other users. Moreover, several tasks like software installation and modifying system configuration require administrative privileges. It would be a security risk to run as a superuser constantly, but how can regular users install software then? Or what about modifying the permissions of a file that you (accidentally) removed your own access to? This is where the `sudo` command comes into play.

`sudo` stands for "superuser do" and it allows regular users, who have been added to the `sudo` user group, to run individual commands as if they were the root user. A common use-case is using `sudo` to install new software via `apt` (on Debian-like systems) or `dnf` (on Fedora/CentOS): `sudo apt install ncbi-blast+`.

With great power comes great responsibility, so exercise caution and think twice before running a command with `sudo`. Only do so when you are sure you know what the end result will be.

### A.4 What is \$PATH?

The `$PATH` is a way of letting your computer know where specific tools or other special locations are stored on your file system. Unless you tell it explicitly, it won't know where to find any new software you install. Fortunately, most methods of installing software automatically take care of this for you, but every now and then you will need to manually add things to your `$PATH`. If you don't, you will be greeted by messages like `Command 'python' not found, did you mean:.`

The `$PATH` is nothing more than a list of locations on your computer. Everything that is found in those locations, will become available to use directly on the CLI without having to type out its full location. Even the basic Unix commands, like `ls` and `cd` are only known to your shell because they are in a location that is indexed by your path.



In the following example, we will demonstrate how you can add a custom directory with scripts to your \$PATH, making them callable from anywhere.

```
# show the contents of PATH
echo $PATH
/home/pmoris/miniforge3/bin:/home/pmoris/miniforge3/condabin:/usr/local/sbin:/usr/local/bin:/usr/local/games:/usr/games
↪ Files (x86)/Common
↪ Files/Oracle/Java/javapath:/mnt/c/windows/system32:/mnt/c/windows:/mnt/c/windows/System32
↪ Files/Git/cmd:/mnt/c/Program
↪ Files/dotnet:/mnt/c/Users/pmoris/AppData/Local/Programs/Quarto/bin:/mnt/c/Users/pmoris/AppData/Local/Programs/Common-IO
↪ VS Code/bin:/snap/bin

# temporarily add directory of scripts to PATH
$ ls ~/itg/FiMAB-bioinformatics/training/scripts
call_variants.sh  download_reference.sh  map.sh  remove_dups.sh
↪ trim.sh
$ export
↪ PATH="$PATH:~/itg/FiMAB-bioinformatics/training/scripts"

# now these scripts can be invoked directly without having to
↪ type out their full location
# e.g., map.sh works just as well as
↪ ~/itg/FiMAB-bioinformatics/training/scripts/map.sh
```

To make these changes permanent, you'd have to add that export statement to your .bashrc file (stored in your home directory). This file is run every time you launch a new shell, so that will allow the \$PATH to be modified every time during startup.

You can find more information on modifying the PATH [here](#).

Lastly, be careful when modifying your PATH. If you mess it up, it can cause all kinds of havoc.

## A.5 Dealing with file permissions

You can find an excellent explanation on file permissions [here](#).

## A.6 Working with remote machines via SSH

In some cases, you will need to work on a Linux machine that is physically located somewhere else, i.e. a remote server. Access to these is usually managed via a command-line tool called SSH (or a stand-alone GUI tool like Putty in Windows). The syntax of the ssh command is as follows:

```
ssh username@domain
```

Where username is a name given to you by the admin of the system and domain is the address of the server (can be a URL or an IP address).

The connection is secured via *SSH keys*: a pair of files used for authentication stored in ~/.ssh.

1. Public key: e.g., id\_rsa.pub or id\_ed25519.pub, located on the remote server.



2. Private file: e.g., `id_rsa` or `id_ed25519.pub`, located on your own machine. \*Never share this file with anyone else!\*\*

Instructions to generate new SSH keys can be found <https://docs.github.com/en/authentication/connecting-to-github-with-ssh/generating-a-new-ssh-key-and-adding-it-to-the-ssh-agent>.

A common problem when connecting is that the file permissions of your keys or credential files are messed up. This can happen if you generate them in Windows and later move them to a Linux file system. To fix this, check <https://superuser.com/a/215506>.

More info on remote servers can be found [here](#).

Lastly, keep in mind that when working on remote servers, it is essential to use `screen` or `tmux` for long-running jobs. Otherwise, they will be interrupted when you disconnect or even when the network briefly fails.

## A.7 Further reading

There are still many facets of Unix that we have not covered here yet, but we hope that you can pick these up on your own if we were to ever need them. For now though, the basics that we covered here, will hopefully already go a long way in helping you use Unix and read scripts written by others.

A selection of more advanced topics to explore at a later time could be:

- Installing software on Linux machines, [from source or via apt \(Ubuntu\)](#), `dnf` (Fedora) or `pacman` (Arch), or through tools like [conda/Miniforge](#) and [bioconda](#).
- **The `find` command**: useful for finding files in your file system.
- **`if/else` conditions**: allow you to execute parts of scripts if certain conditions are fulfilled.
- Command substitution and process substitution: run commands in a subshell to allow more complex ways of redirecting and piping outputs/inputs.
- `sed/awk` are commands that let you do things like search and replace, calculations or other manipulations based on all the lines of a file.
- the `join` command: combining columns of (multiple) tabular data files in particular ways.
- Environment variables: variables that are always available, like `PATH`.
- [Using `screen` or `tmux`](#) to spawn persistent background terminal sessions. This allows you to run commands on a remote server and shutdown your own machine, without the remote process being interrupted.
- Learn about the difference in line endings on Unix (`\n`) and Windows (`\r\n`), how to switch between them (`dos2unix`) and how to set a preference in your editors like RStudio.
- Learn about `sudo`, allowing you to perform actions as an administrator.
- Regular expressions, to power up your `grep` searches and `sed/awk` commands.