



Computational Genomics Tutorial

Release '2020.1.8'

Sebastian Schmeier (<https://sschmeier.com>)

Jan 23, 2020

CONTENTS

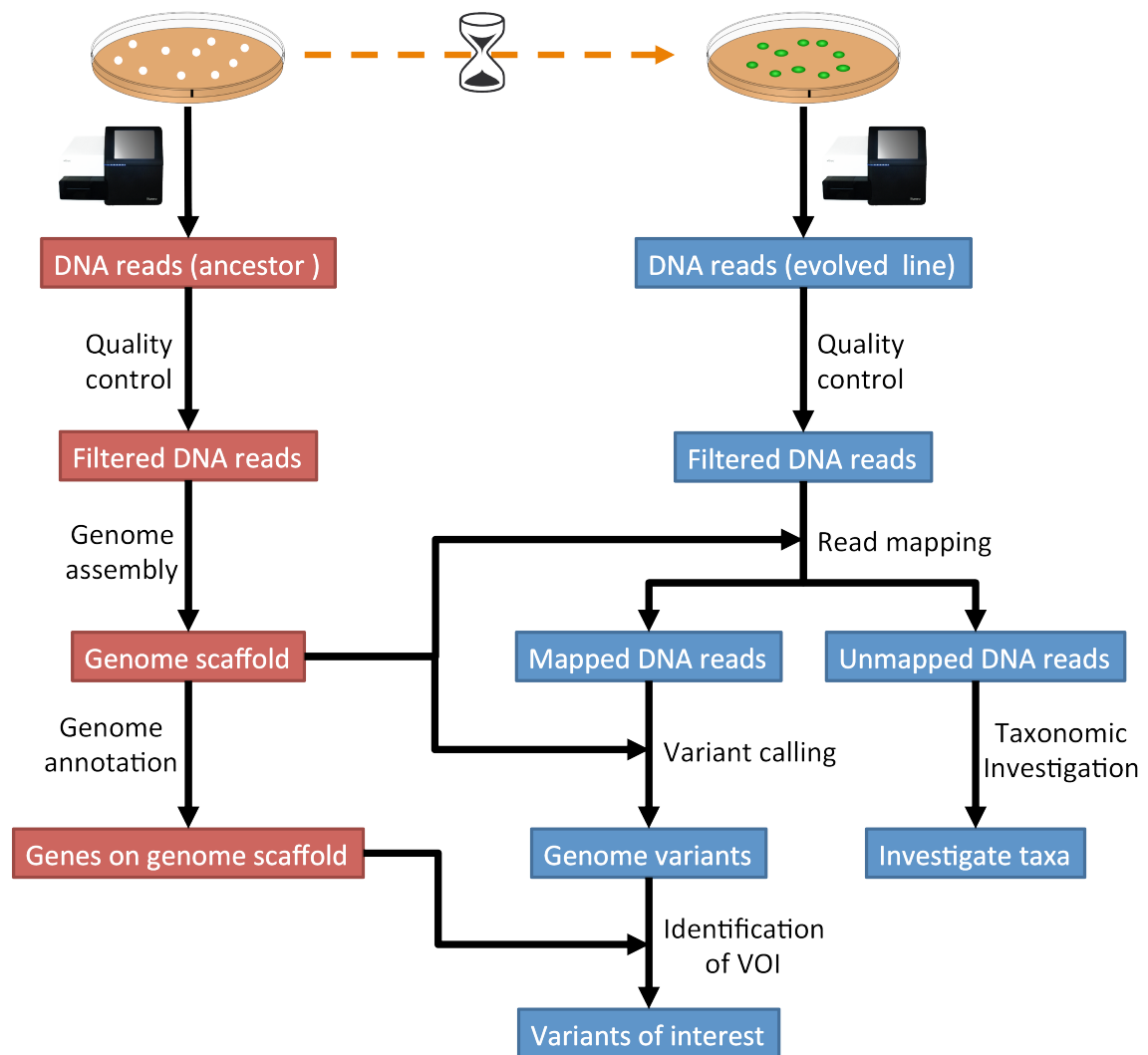
1	Introduction	3
1.1	The workflow	3
1.2	Learning outcomes	3
2	Tool installation	5
2.1	Install the conda package manager	5
2.2	Creating environments	6
2.3	Install software	6
2.4	General conda commands	7
3	Quality control	9
3.1	Preface	9
3.2	Overview	9
3.3	Learning outcomes	9
3.4	The data	9
3.5	The fastq file format	12
3.6	The QC process	12
3.7	PhiX genome	12
3.8	Adapter trimming	13
3.9	Quality assessment of sequencing reads	14
3.10	Run FastQC and MultiQC on the trimmed data	15
4	Genome assembly	19
4.1	Preface	19
4.2	Overview	19
4.3	Learning outcomes	19
4.4	Before we start	19
4.5	Creating a genome assembly	21
4.6	Assembly quality assessment	22
4.7	Compare the untrimmed data	22
4.8	Further reading	23
4.9	Web links	23
5	Read mapping	25
5.1	Preface	25
5.2	Overview	25
5.3	Learning outcomes	25
5.4	Before we start	25
5.5	Mapping sequence reads to a reference genome	27
5.6	BWA	27
5.7	The sam mapping file-format	29
5.8	Mapping post-processing	29
5.9	Mapping statistics	31
5.10	Sub-selecting reads	32

6	Taxonomic investigation	35
6.1	Preface	35
6.2	Overview	35
6.3	Before we start	35
6.4	Kraken2	37
6.5	Centrifuge	41
6.6	Visualisation (Krona)	44
7	Variant calling	47
7.1	Preface	47
7.2	Overview	47
7.3	Learning outcomes	47
7.4	Before we start	47
7.5	Installing necessary software	49
7.6	Preprocessing	49
7.7	Calling variants	49
7.8	Post-processing	49
8	Genome annotation	55
8.1	Preface	55
8.2	Overview	55
8.3	Learning outcomes	55
8.4	Before we start	55
8.5	Installing the software	57
8.6	Assessment of orthologue presence and absence	58
8.7	Annotation with Augustus	58
8.8	Annotation with Prokka	59
8.9	Interactive viewing	59
9	Orthology and Phylogeny	61
9.1	Preface	61
9.2	Learning outcomes	61
9.3	Before we start	61
9.4	Installing the software	62
9.5	Finding orthologues using BLAST	62
9.6	Performing an alignment	63
9.7	Building a phylogeny	63
9.8	Visualizing the phylogeny	64
10	Variants-of-interest	65
10.1	Preface	65
10.2	Overview	65
10.3	Learning outcomes	65
10.4	Before we start	65
10.5	General comments for identifying variants-of-interest	67
10.6	SnEff	67
11	Quick command reference	71
11.1	Shell commands	71
11.2	General conda commands	71
12	Coding solutions	73
12.1	QC	73
12.2	Assembly	74
12.3	Mapping	74
13	Downloads	77
13.1	Tools	77
13.2	Data	77

Warning: This is a pre-release version. This is work in progress with a new dataset. For a working version please see <https://genomics.sschmeier.com>

This is an introductory tutorial for learning computational genomics mostly on the Linux command-line. You will learn how to analyse next-generation sequencing (NGS) data. The data you will be using is real research data. The final aim is to identify genome variations in evolved lines of wild yeast that can explain the observed biological phenotypes. Currently [Sebastian¹](#) is teaching this material in the Massey University course [Genome Science²](#).

More information about other bioinformatics material and our research can be found on the webpages of the [Schmeier Group³](#) (<https://sschmeier.com>).



Note: A online version of this tutorial can be accessed at <https://genomics.sschmeier.com>.

¹ <https://sschmeier.com>

² https://www.massey.ac.nz/massey/learning/programme-course/course.cfm?course_code=203341

³ <https://sschmeier.com>

INTRODUCTION

This is an introductory tutorial for learning genomics mostly on the Linux command-line. Should you need to refresh your knowledge about either Linux or the command-line, have a look [here](http://linux.sschmeier.com/)⁴.

In this tutorial you will learn how to analyse next-generation sequencing (NGS) data. The data you will be using is actual research data. The experiment follows a similar strategy as in what is called an “experimental evolution” experiment [KAWECKI2012], [ZEYL2006]. The final aim is to identify the genome variations in evolved lines of *E. coli* that can explain the observed biological phenotype(s).

1.1 The workflow

The tutorial workflow is summarised in [Fig. 1.1](#).

1.2 Learning outcomes

During this tutorial you will learn to:

- Check the data quality of an NGS experiment
- Create a genome assembly of the ancestor based on NGS data
- Map NGS reads of evolved lines to the created ancestral reference genome
- Call genome variations/mutations in the evolved lines
- Annotate a newly derived reference genome
- Find variants of interest that may be responsible for the observed evolved phenotype(s)

⁴ <http://linux.sschmeier.com/>

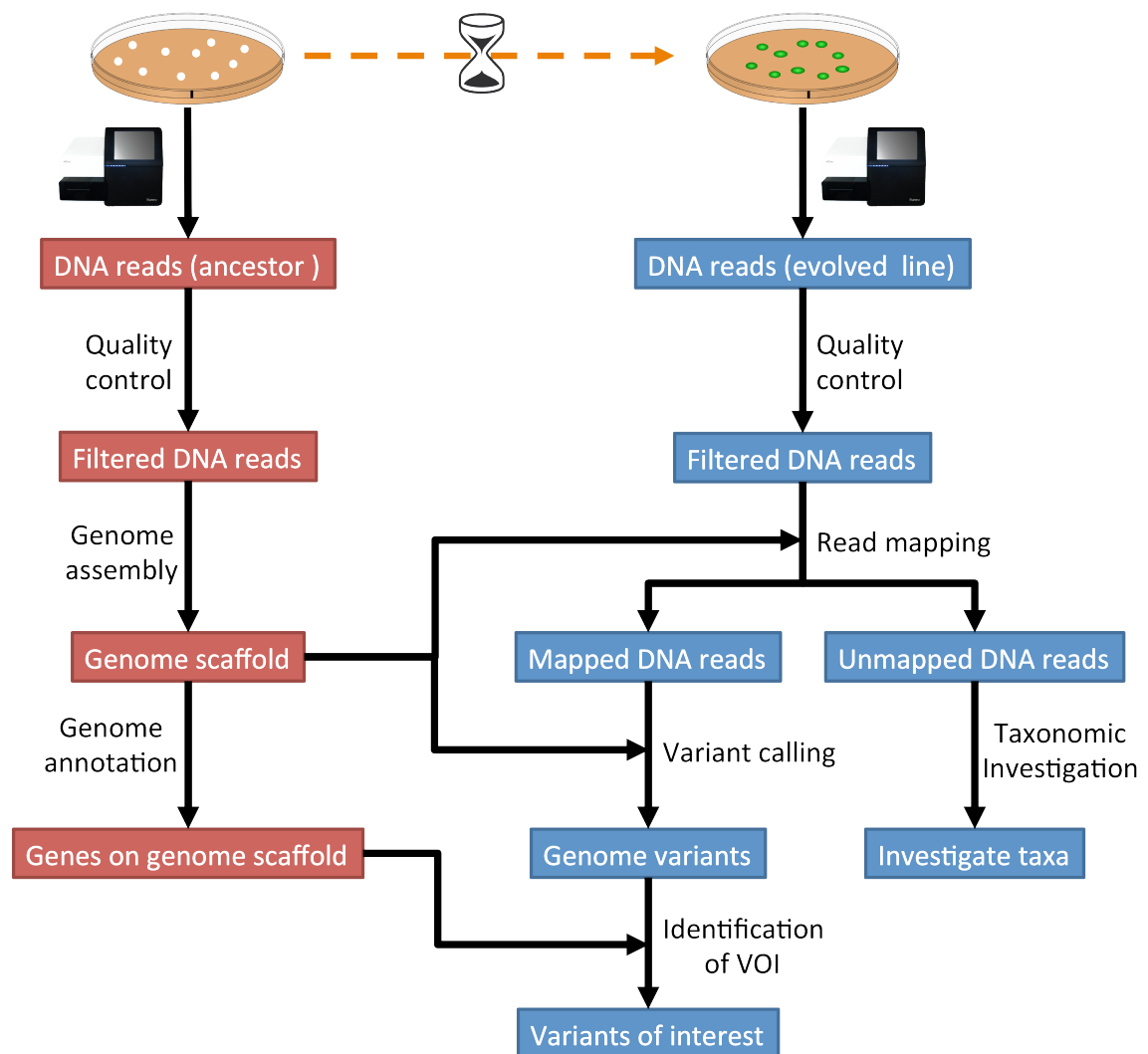


Fig. 1.1: The tutorial will follow this workflow.

TOOL INSTALLATION

2.1 Install the conda package manager

We will use the package/tool managing system [conda](http://conda.pydata.org/miniconda.html)⁷ to install some programs that we will use during the course. It is not installed by default, thus we need to install it first to be able to use it. Let us download [conda](http://conda.pydata.org/miniconda.html)⁸ first:

```
# download latest conda installer
$ wget https://repo.continuum.io/miniconda/Miniconda3-latest-Linux-x86_64.sh
```

Note: Should the conda installer download fail. Please find links to alternative locations on the [Downloads](#) (page 77) page.

Now lets install [conda](http://conda.pydata.org/miniconda.html)⁹:

```
# run the installer
$ bash Miniconda3-latest-Linux-x86_64.sh
```

After you accepted the license agreement [conda](http://conda.pydata.org/miniconda.html)¹⁰ will be installed. At the end of the installation you will encounter the following:

```
...
installation finished.
Do you wish the installer to initialize Miniconda3
by running conda init? [yes|no]
[no] >>>
```

Please type “yes” here. This will add some code to your `.bashrc` init file, which is important to work with [conda](http://conda.pydata.org/miniconda.html)¹¹ correctly.

Attention: Please close and reopen the terminal, to complete the installation.

After closing and re-opening the shell/terminal, we should be able to use the [conda](http://conda.pydata.org/miniconda.html)¹² command:

```
$ conda update --yes conda
```

⁷ <http://conda.pydata.org/miniconda.html>

⁸ <http://conda.pydata.org/miniconda.html>

⁹ <http://conda.pydata.org/miniconda.html>

¹⁰ <http://conda.pydata.org/miniconda.html>

¹¹ <http://conda.pydata.org/miniconda.html>

¹² <http://conda.pydata.org/miniconda.html>

2.1.1 Installing conda channels to make tools available

Different tools are packaged in what [conda](http://conda.pydata.org/miniconda.html)¹³ calls channels. We need to add some channels to make the bioinformatics and genomics tools available for installation:

```
# Install some conda channels
# A channel is where conda looks for packages
$ conda config --add channels defaults
$ conda config --add channels bioconda
$ conda config --add channels conda-forge
```

Attention: The order of adding channels is important. Make sure you use the shown order of commands.

2.2 Creating environments

We create a [conda](http://conda.pydata.org/miniconda.html)¹⁴ environment for some tools. This is useful to work **reproducible** as we can easily re-create the tool-set with the same version numbers later on.

```
$ conda create -n ngs python=3
# activate the environment
$ conda activate ngs
```

So what is happening when you type `conda activate ngs` in a shell. The `PATH` variable of your shell gets temporarily manipulated and set to:

```
$ echo $PATH
/home/guest/miniconda3/bin:/home/guest/miniconda3/condabin:...
$ conda activate ngs
$ echo $PATH
/home/guest/miniconda3/envs/ngs/bin:/home/guest/miniconda3/condabin: ...
```

Now it will look first in your environment's bin directory but afterwards in the general conda bin (`/home/guest/miniconda3/condabin`). So basically everything you install generally with conda (without being in an environment) is also available to you but gets overshadowed if a similar program is in `/home/guest/miniconda3/envs/ngs/bin` and you are in the `ngs` environment.

2.3 Install software

To install software into the activated environment, one uses the command `conda install`.

```
# install more tools into the environment
$ conda install package
```

Note: To tell if you are in the correct conda environment, look at the command-prompt. Do you see the name of the environment in round brackets at the very beginning of the prompt, e.g. `(ngs)`? If not, activate the `ngs` environment with `conda activate ngs` before installing the tools.

¹³ <http://conda.pydata.org/miniconda.html>

¹⁴ <http://conda.pydata.org/miniconda.html>

2.4 General conda commands

```
# to search for packages
$ conda search [package]

# To update all packages
$ conda update --all --yes

# List all packages installed
$ conda list [-n env]

# conda list environments
$ conda env list

# create new env
$ conda create -n [name] package [package] ...

# activate env
$ conda activate [name]

# deactivate env
$ conda deactivate
```


QUALITY CONTROL

3.1 Preface

There are many sources of errors that can influence the quality of your sequencing run [ROBASKY2014]. In this quality control section we will use our skill on the command-line interface to deal with the task of investigating the quality and cleaning sequencing data [KIRCHNER2014].

Note: You will encounter some **To-do** sections at times. Write the solutions and answers into a text-file.

3.2 Overview

The part of the workflow we will work on in this section can be viewed in [Fig. 3.1](#).

3.3 Learning outcomes

After studying this tutorial you should be able to:

1. Describe the steps involved in pre-processing/cleaning sequencing data.
2. Distinguish between a good and a bad sequencing run.
3. Compute, investigate and evaluate the quality of sequence data from a sequencing experiment.

3.4 The data

First, we are going to download the data we will analyse. Open a shell/terminal.

```
# create a directory you work in
$ mkdir analysis

# change into the directory
$ cd analysis

# download the data
$ wget http://compbio.massey.ac.nz/data/203341/data.tar.gz

# uncompress it
$ tar -xvzf data.tar.gz
```

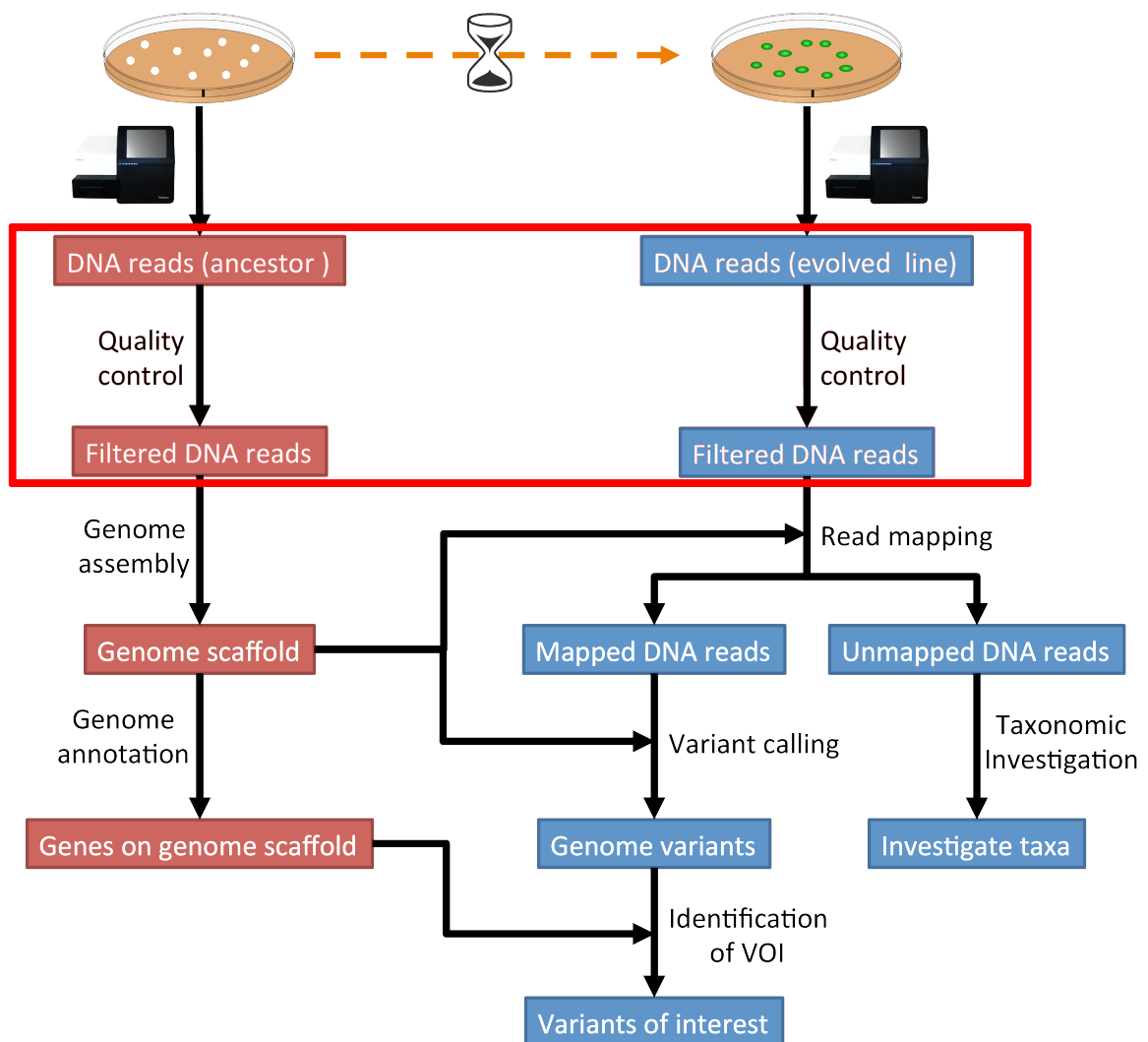


Fig. 3.1: The part of the workflow we will work on in this section marked in red.

Note: Should the download fail, download manually from [Downloads](#) (page 77).

The data is from a paired-end sequencing run data (see [Fig. 3.2](#)) from an [Illumina¹⁵](#) HiSeq [[GLENN2011](#)]. Thus, we have two files, one for each end of the read.

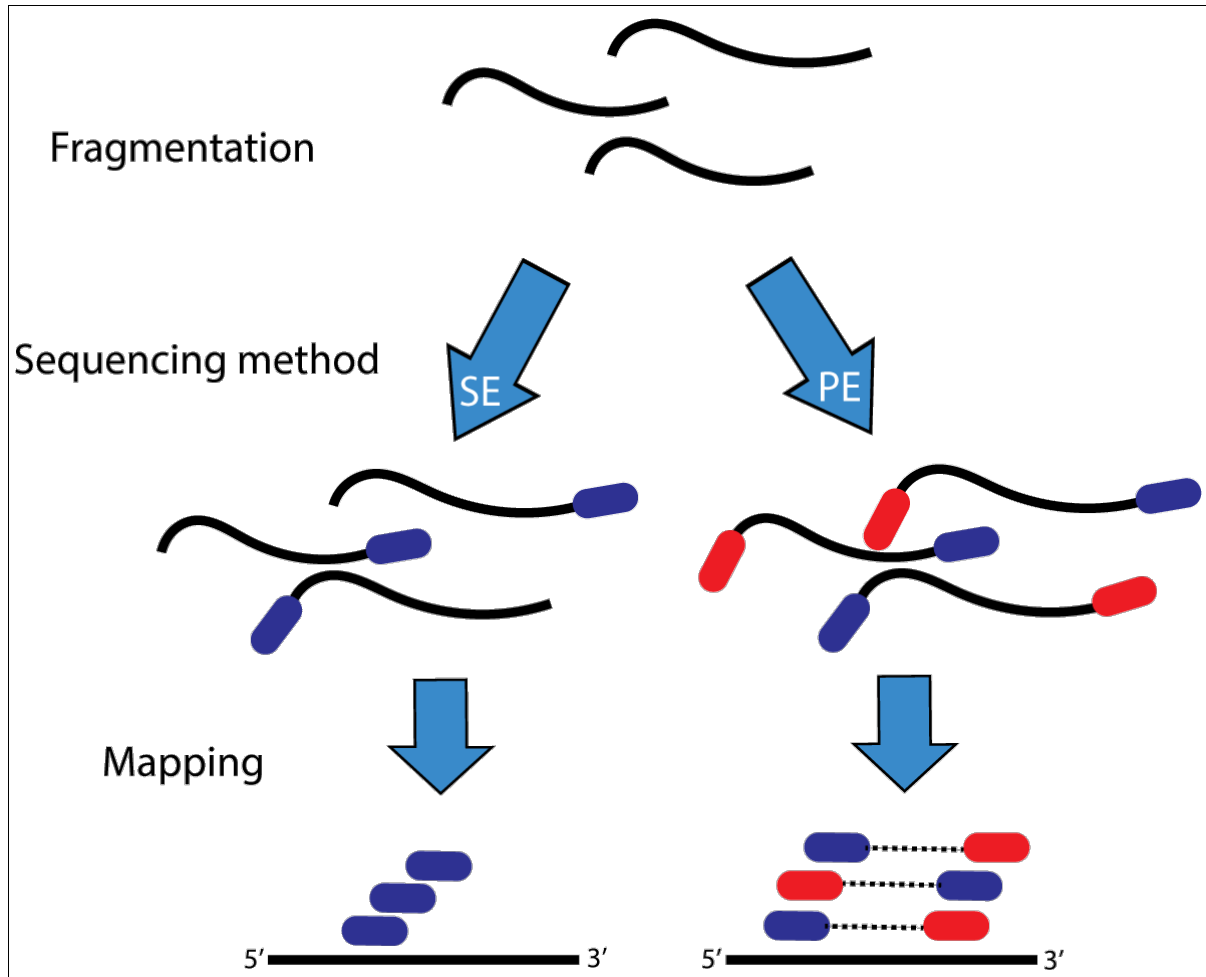


Fig. 3.2: Illustration of single-end (SE) versus paired-end (PE) sequencing.

If you need to refresh how [Illumina¹⁶](#) paired-end sequencing works have a look at the [Illumina technology webpage¹⁷](#) and this [video¹⁸](#).

Attention: The data we are using is “almost” raw data as it came from the machine. This data has been post-processed in two ways already. All sequences that were identified as belonging to the PhiX genome have been removed. This process requires some skills we will learn in later sections. [Illumina¹⁹](#) adapters have been removed as well already! The process is explained below and we are going to run through the process anyways.

¹⁵ <http://illumina.com>

¹⁶ <http://illumina.com>

¹⁷ http://www.illumina.com/technology/next-generation-sequencing/paired-end-sequencing_assay.html

¹⁸ <https://youtu.be/HMyCqWhwB8E>

¹⁹ <http://illumina.com>

3.4.1 Investigate the data

Make use of your newly developed skills on the command-line to investigate the files in data folder.

Todo:

1. Use the command-line to get some ideas about the file.
 2. What kind of files are we dealing with?
 3. How many sequence reads are in the file?
 4. Assume a genome size of ~4.6 MB. Calculate the coverage based on this formula: $C = LN / G$
-

- C: Coverage
- G: is the haploid genome length in bp
- L: is the read length in bp (e.g. 2x150 paired-end = 300)
- N: is the number of reads sequenced

3.5 The fastq file format

The data we receive from the sequencing is in fastq format. To remind us what this format entails, we can revisit the [fastq wikipedia-page](#)²⁰!

A useful tool to decode base qualities can be found [here](#)²¹.

Todo: Explain briefly what the quality value represents.

3.6 The QC process

There are a few steps one need to do when getting the raw sequencing data from the sequencing facility:

1. Remove PhiX sequences (we are not going to do this)
2. Adapter trimming
3. Quality trimming of reads
4. Quality assessment

3.7 PhiX genome

PhiX²² is a nontailed bacteriophage with a single-stranded DNA and a genome with 5386 nucleotides. PhiX is used as a quality and calibration control for [sequencing runs](#)²³. PhiX is often added at a low known concentration, spiked in the same lane along with the sample or used as a separate lane. As the concentration of the genome is known, one can calibrate the instruments. Thus, PhiX genomic sequences need to be removed before processing your data further as this constitutes a deliberate contamination [MUKHERJEE2015]. The steps involve mapping all reads to the “known” PhiX genome, and removing all of those sequence reads from the data.

²⁰ https://en.wikipedia.org/wiki/FASTQ_format

²¹ <http://broadinstitute.github.io/picard/explain-qualities.html>

²² https://en.wikipedia.org/wiki/Phi_X_174

²³ <http://www.illumina.com/products/by-type/sequencing-kits/cluster-gen-sequencing-reagents/phix-control-v3.html>

However, your sequencing provider might not have used PhiX, thus you need to read the protocol carefully, or just do this step in any case.

Attention: We are **not** going to do this step here, as the sequencing run we are using did not use PhiX. Please see the [Read mapping](#) (page 25) section on how to map reads against a reference genome.

3.8 Adapter trimming

The process of sequencing DNA via [Illumina](#)²⁴ technology requires the addition of some adapters to the sequences. These get sequenced as well and need to be removed as they are artificial and do not belong to the species we try to sequence. Generally speaking we have to deal with a trade-off between accuracy of adapter removal and speed of the process. Adapter trimming does take some time.

Also, we have generally two different approaches when trimming adapter:

1. We can use a tool that takes an adapter or list of adapters and removes these from each sequence read.
2. We can use a tool that predicts adapters and removes them from each sequence read.

For the first approach we need to know the adapter sequences that were used during the sequencing of our samples. Normally, you should ask your sequencing provider, who should be providing this information to you. [Illumina](#)²⁵ itself provides a [document](#)²⁶ that describes the adapters used for their different technologies. Also the [FastQC](#)²⁷ tool, we will be using later on, provides a [collection of contaminants and adapters](#)²⁸.

However, often (sadly) this information is not readily available, e.g. when dealing with public data. Thus, the second approach can be employed, that is, using a tool that predicts adapters.

Here, we are going to use the second approach with a tool called [fastp](#) to trim adapters **and** do quality trimming. [fastp](#) has a few characteristics which make it a great tool, most importantly: it is pretty fast, provides good information after the run, and can do quality trimming as well, thus saving us to use another tool to do this.

Quality trimming of our sequencing reads will remove bad quality called bases from our reads, which is especially important when dealing with variant identification.

```
# create env and install tools
$ conda create --yes -n qc fastp fastqc multiqc

# activate env
$ conda activate qc
```

Here, as an example we are trimming the sequence reads of the ancestor:

```
$ mkdir trimmed

$ fastp --detect_adapter_for_pe
      --overrepresentation_analysis
      --correction --cut_right --thread 2
      --html trimmed/anc.fastp.html --json trimmed/anc.fastp.json
      -i data/anc_R1.fastq.gz -I data/anc_R2.fastq.gz
      -o trimmed/anc_R1.fastq.gz -O trimmed/anc_R2.fastq.gz
```

²⁴ <http://illumina.com>

²⁵ <http://illumina.com>

²⁶ <https://support.illumina.com/downloads/illumina-customer-sequence-letter.html>

²⁷ <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

²⁸ https://github.com/csf-ngs/fastqc/blob/master/Contaminants/contaminant_list.txt

- `--detect_adapter_for_pe`: Specifies that we are dealing with paired-end data.
- `--overrepresentation_analysis`: Analyse the sequence collection for sequences that appear too often.
- `--correction`: Will try to correct bases based on an overlap analysis of read1 and read2.
- `--cut_right`: Will use quality trimming and scan the read from start to end in a window. If the quality in the window is below what is required, the window plus all sequence towards the end is discarded and the read is kept if its still long enough.
- `--thread`: Specify how many concurrent threads the process can use.
- `--html` and `--json`: We specify the location of some stat files.
- `-i data/anc_R1.fastq.gz -I data/anc_R2.fastq.gz`: Specifies the two input read files
- `-o trimmed/anc_R1.fastq.gz -O trimmed/anc_R2.fastq.gz`: Specifies the two desired output read files

Todo:

1. Run `fastp` also on the evolved samples.

Hint: Should you not get the commands together to trim the evolved samples, have a look at the coding solutions at [Code: fastp](#) (page 73). Should you be unable to run `fastp` at all to trim the data. You can download the trimmed dataset [here](#)²⁹. Unarchive and uncompress the files with `tar -xvzf trimmed.tar.gz`.

3.9 Quality assessment of sequencing reads

3.9.1 Installing FastQC

```
$ fastqc --help
```

FastQC - A high throughput sequence QC analysis tool

SYNOPSIS

```
fastqc seqfile1 seqfile2 .. seqfileN
```

```
fastqc [-o output dir] [--(no)extract] [-f fastq|bam|sam]
      [-c contaminant file] seqfile1 .. seqfileN
```

DESCRIPTION

FastQC reads a **set** of sequence files and produces from each one a quality control report consisting of a number of different modules, each one of which will **help** to identify a different potential **type** of problem in your data.

If no files to process are specified on the **command** line **then** the program will start as an interactive graphical application. If files are provided on the **command** line **then** the program will run with no user interaction required. In this mode it is suitable **for** inclusion into a standardised analysis pipeline.

²⁹ <http://compbio.massey.ac.nz/data/203341/trimmed.tar.gz>

3.9.2 FastQC manual

FastQC³⁰ is a very simple program to run that provides information about sequence read quality.

From the webpage:

“FastQC aims to provide a simple way to do some quality control checks on raw sequence data coming from high throughput sequencing pipelines. It provides a modular set of analyses which you can use to give a quick impression of whether your data has any problems of which you should be aware before doing any further analysis.”

The basic command looks like:

```
$ fastqc -o RESULT-DIR INPUT-FILE.fq(.gz) ...
```

- -o RESULT-DIR is the directory where the result files will be written
- INPUT-FILE.fq is the sequence file to analyze, can be more than one file.

Hint: The result will be a HTML page per input file that can be opened in a web-browser.

Hint: The authors of FastQC³¹ made some nice help pages explaining each of the plots and results you expect to see [here](#)³².

3.9.3 MultiQC

MultiQC³³ is an excellent tool to put FastQC³⁴ (and other tool) results of different samples into context. It compiles all FastQC³⁵ results and fastp stats into one nice web-page.

The use of MultiQC³⁶ is simple. Just provide the command with a directories where multiple results are stored and it will compile a nice report, e.g.:

```
$ multiqc DIRECTORY DIRECTORY ...
```

3.10 Run FastQC and MultiQC on the trimmed data

Todo:

1. Create a directory for the results -> trimmed-fastqc
2. Run FastQC on all **trimmed** files.
3. Visit the FastQC³⁷ website and read about sequencing QC reports for good and bad Illumina³⁸ sequencing runs.
4. Run MultiQC³⁹ on the trimmed-fastqc and trimmed directories

³⁰ <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

³¹ <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

³² <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/Help/3%20Analysis%20Modules/>

³³ <https://multiqc.info/>

³⁴ <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

³⁵ <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

³⁶ <https://multiqc.info/>

³⁷ <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

³⁸ <http://illumina.com>

³⁹ <https://multiqc.info/>

5. Compare your results to these examples (Fig. 3.3 to Fig. 3.5) of a particularly bad run (taken from the [FastQC⁴⁰](http://www.bioinformatics.babraham.ac.uk/projects/fastqc/) website) and write down your observations with regards to your data.
6. What elements in these example figures (Fig. 3.3 to Fig. 3.5) indicate that the example is from a bad run?

Hint: Should you not get it right, try the commands in *Code: FastQC* (page 73).

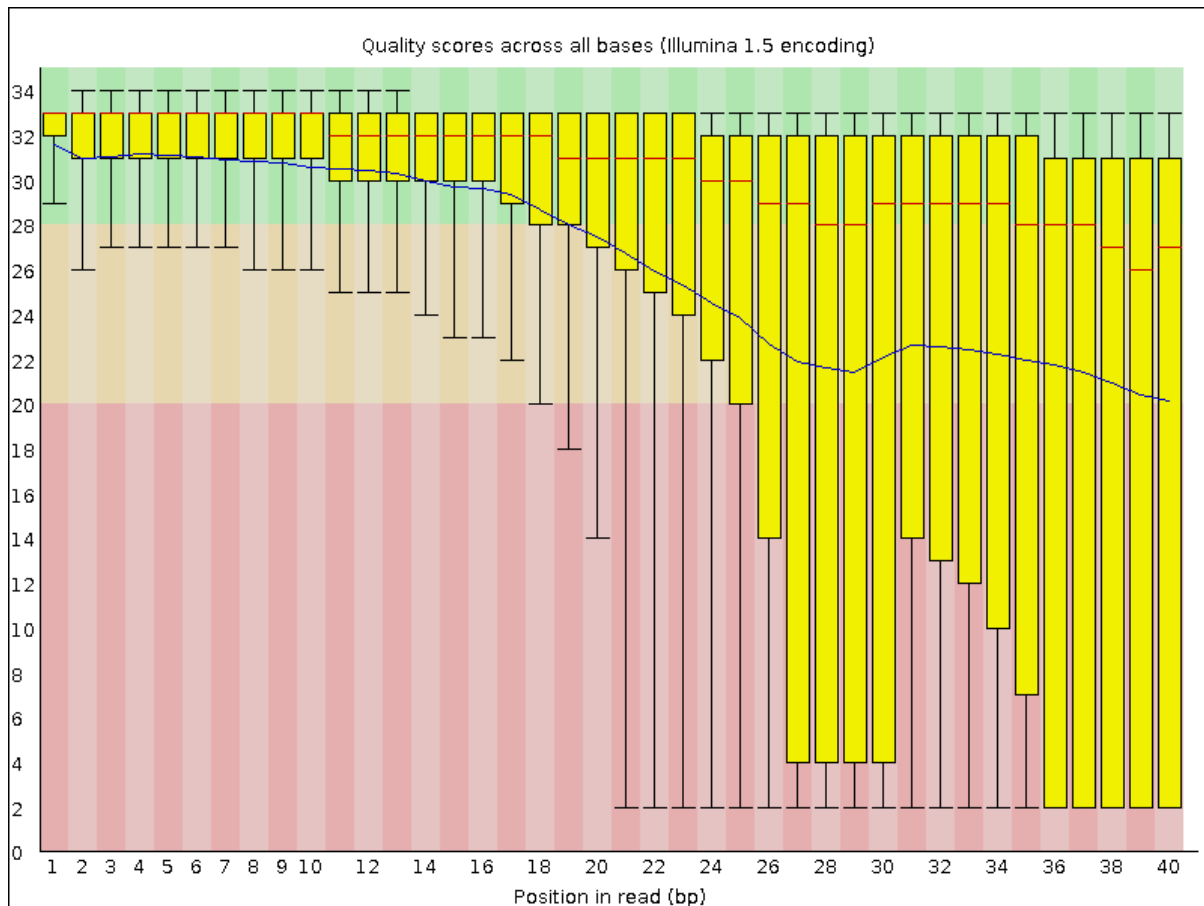


Fig. 3.3: Quality score across bases.

⁴⁰ <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

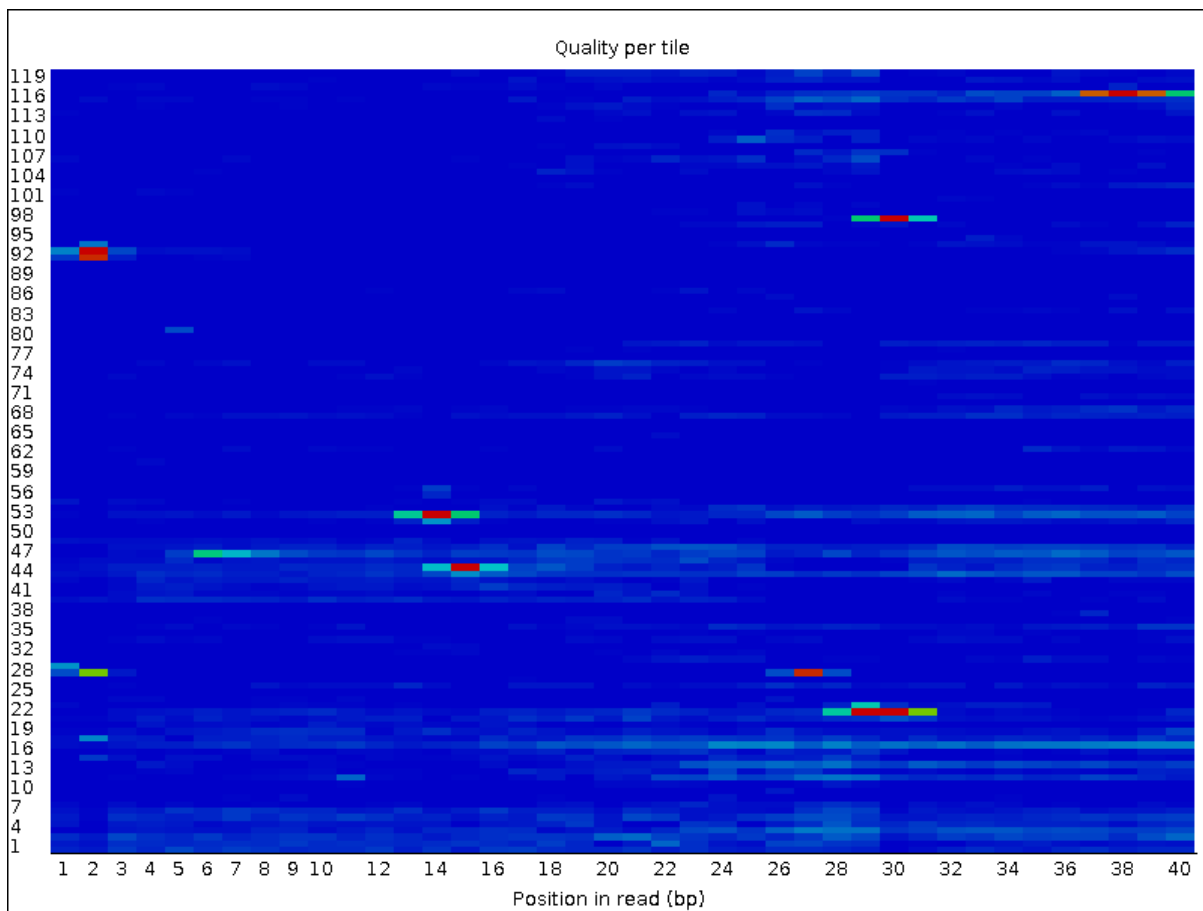


Fig. 3.4: Quality per tile.

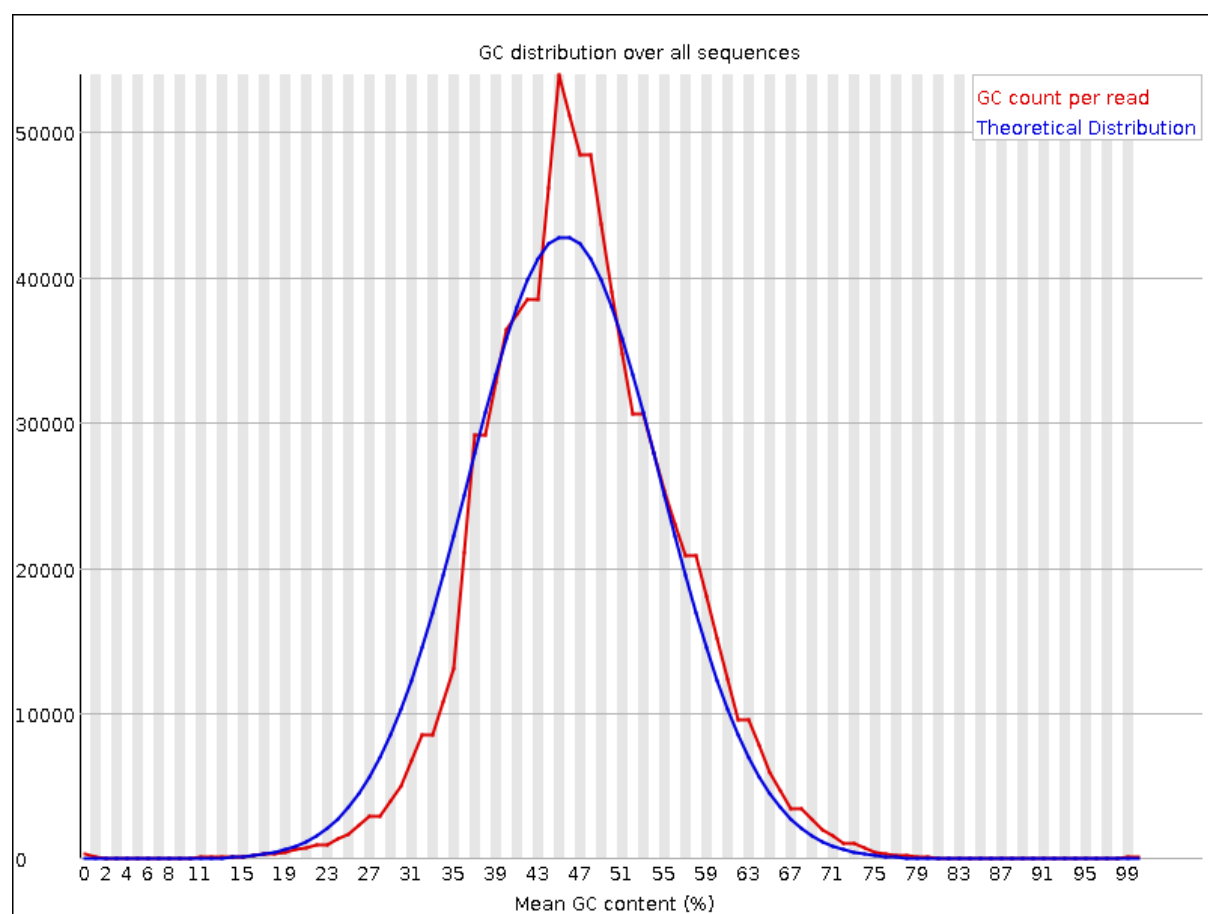


Fig. 3.5: GC distribution over all sequences.

GENOME ASSEMBLY

4.1 Preface

In this section we will use our skill on the command-line interface to create a genome assembly from sequencing data.

Note: You will encounter some **To-do** sections at times. Write the solutions and answers into a text-file.

4.2 Overview

The part of the workflow we will work on in this section can be viewed in [Fig. 4.1](#).

4.3 Learning outcomes

After studying this tutorial you should be able to:

1. Compute and interpret a whole genome assembly.
2. Judge the quality of a genome assembly.

4.4 Before we start

Lets see how our directory structure looks so far:

```
$ cd ~/analysis
$ ls -lF
```

```
data/
multiqc_data/
multiqc_report.html
trimmed/
trimmed-fastqc/
```

Attention: If you have not run the previous section [Quality control](#) (page 9), you can download the trimmed data needed for this section here: [Downloads](#) (page 77).

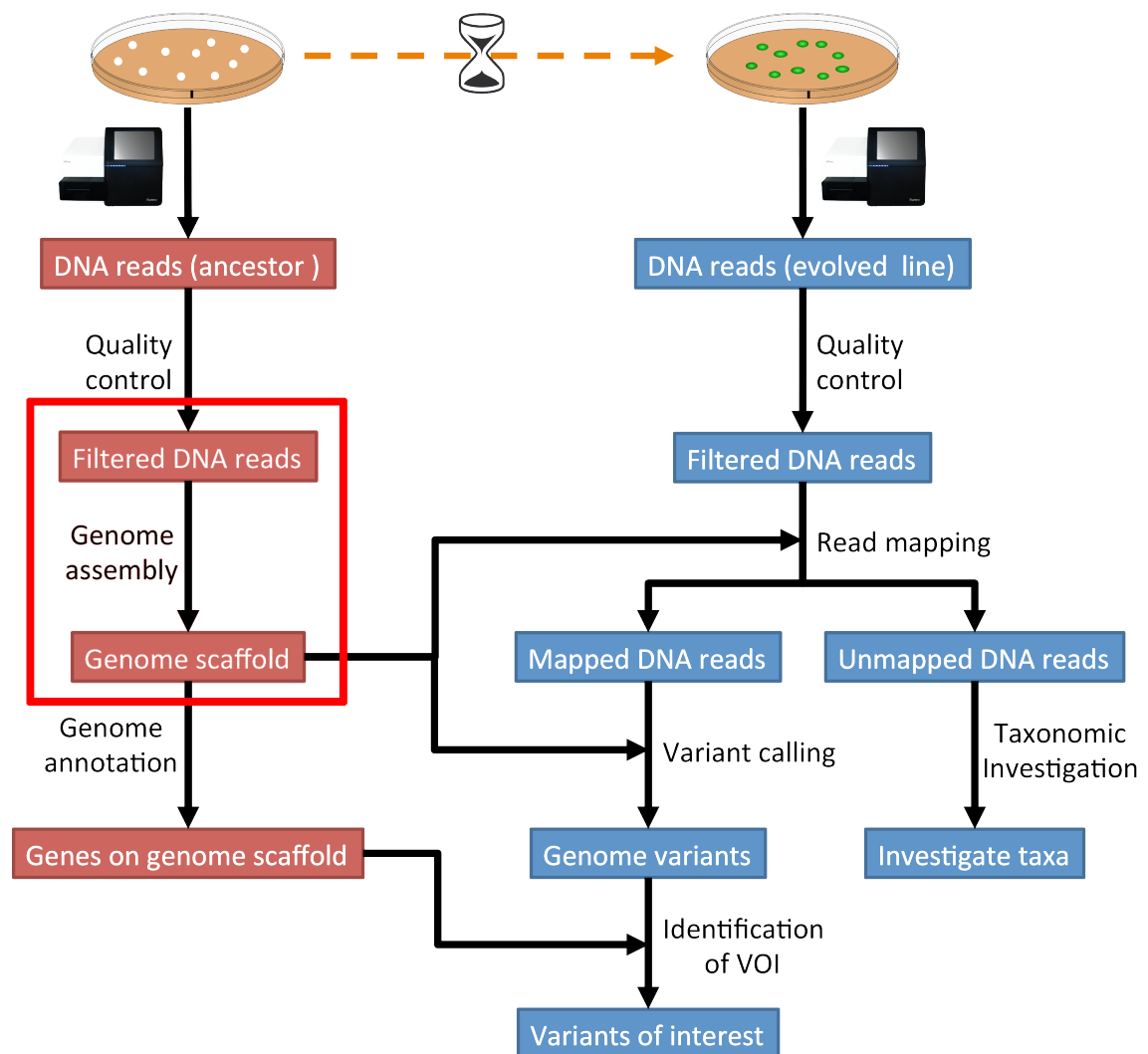


Fig. 4.1: The part of the workflow we will work on in this section marked in red.

4.5 Creating a genome assembly

We want to create a genome assembly for our ancestor. We are going to use the quality trimmed forward and backward DNA sequences and use a program called [SPAdes⁴⁵](#) to build a genome assembly.

Todo:

1. Discuss briefly why we are using the ancestral sequences to create a reference genome as opposed to the evolved line.
-

4.5.1 Installing the software

We are going to use a program called [SPAdes⁴⁶](#) fo assembling our genome. In a recent evaluation of assembly software, [SPAdes⁴⁷](#) was found to be a good choice for fungal genomes [ABBAS2014]. It is also simple to install and use.

```
$ conda create -n assembly spades quast
$ conda activate assembly
```

4.5.2 SPAdes usage

```
# change to your analysis root folder
$ cd ~/analysis

# first create a output directory for the assemblies
$ mkdir assembly

# to get a help for spades and an overview of the parameter type:
$ spades.py -h
```

Generally, paired-end data is submitted in the following way to [SPAdes⁴⁸](#):

```
$ spades.py -o result-directory -1 read1.fastq.gz -2 read2.fastq.gz
```

Todo:

1. Run [SPAdes⁴⁹](#) with default parameters on the ancestor's trimmed reads
 2. Read in the [SPAdes⁵⁰](#) manual about about assembling with 2x150bp reads
 3. Run [SPAdes⁵¹](#) a second time but use the options suggested at the [SPAdes⁵²](#) manual [section 3.4⁵³](#) for assembling 2x150bp paired-end reads. Use a different output directory `assembly/spades-150` for this run.
-

Hint: Should you not get it right, try the commands in [Code: SPAdes assembly \(trimmed data\)](#) (page 74).

⁴⁵ <http://bioinf.spbau.ru/spades>

⁴⁶ <http://bioinf.spbau.ru/spades>

⁴⁷ <http://bioinf.spbau.ru/spades>

⁴⁸ <http://bioinf.spbau.ru/spades>

⁴⁹ <http://bioinf.spbau.ru/spades>

⁵⁰ <http://bioinf.spbau.ru/spades>

⁵¹ <http://bioinf.spbau.ru/spades>

⁵² <http://bioinf.spbau.ru/spades>

⁵³ <http://cab.spbu.ru/files/release3.14.0/manual.html#sec3.4>

4.6 Assembly quality assessment

4.6.1 Assembly statistics

Quast⁵⁴ (Quality ASsessment Tool) [GUREVICH2013], evaluates genome assemblies by computing various metrics, including:

- N50: length for which the collection of all contigs of that length or longer covers at least 50% of assembly length
- NG50: where length of the reference genome is being covered
- NA50 and NGA50: where aligned blocks instead of contigs are taken
- miss-assemblies: miss-assembled and unaligned contigs or contigs bases
- genes and operons covered

It is easy with Quast⁵⁵ to compare these measures among several assemblies. The program can be used on their website⁵⁶.

```
$ conda install quast
```

Run Quast⁵⁷ with both assembly scaffolds.fasta files to compare the results.

```
$ quast -o assembly/quast assembly/spades-default/scaffolds.fasta assembly/spades-150/scaffolds.  
↪ fasta
```

Todo:

1. Compare the results of Quast⁵⁸ with regards to the two different assemblies.
2. Which one do you prefer and why?

4.7 Compare the untrimmed data

Todo:

1. To see if our trimming procedure has an influence on our assembly, run the same command you used on the trimmed data on the original untrimmed data.
2. Run Quast⁵⁹ on the assembly and compare the statistics to the one derived for the trimmed data set. Write down your observations.

Hint: Should you not get it right, try the commands in *Code: SPAdes assembly (original data)* (page 74).

⁵⁴ <http://quast.bioinf.spbau.ru/>

⁵⁵ <http://quast.bioinf.spbau.ru/>

⁵⁶ <http://quast.bioinf.spbau.ru/>

⁵⁷ <http://quast.bioinf.spbau.ru/>

⁵⁸ <http://quast.bioinf.spbau.ru/>

⁵⁹ <http://quast.bioinf.spbau.ru/>

4.8 Further reading

4.8.1 Background on Genome Assemblies

- How to apply de Bruijn graphs to genome assembly. [COMPEAU2011]
- Sequence assembly demystified. [NAGARAJAN2013]

4.8.2 Evaluation of Genome Assembly Software

- GAGE: A critical evaluation of genome assemblies and assembly algorithms. [SALZBERG2012]
- Assessment of de novo assemblers for draft genomes: a case study with fungal genomes. [ABBAS2014]

4.9 Web links

- Lectures for this topic: [Genome Assembly: An Introduction](#)⁶⁰
- [SPAdes](#)⁶¹
- [Quast](#)⁶²
- [Bandage](#)⁶³ (Bioinformatics Application for Navigating De novo Assembly Graphs Easily) is a program that visualizes a genome assembly as a graph [WICK2015].

⁶⁰ <https://dx.doi.org/10.6084/m9.figshare.2972323.v1>

⁶¹ <http://bioinf.spbau.ru/spades>

⁶² <http://quast.bioinf.spbau.ru/>

⁶³ <https://rrwick.github.io/Bandage/>

READ MAPPING

5.1 Preface

In this section we will use our skill on the command-line interface to map our reads from the evolved line to our ancestral reference genome.

Note: You will encounter some **To-do** sections at times. Write the solutions and answers into a text-file.

5.2 Overview

The part of the workflow we will work on in this section can be viewed in [Fig. 5.1](#).

5.3 Learning outcomes

After studying this section of the tutorial you should be able to:

1. Explain the process of sequence read mapping.
2. Use bioinformatics tools to map sequencing reads to a reference genome.
3. Filter mapped reads based on quality.

5.4 Before we start

Lets see how our directory structure looks so far:

```
$ cd ~/analysis
# create a mapping result directory
$ mkdir mappings
$ ls -lF
```

```
assembly/
data/
mappings/
multiqc_data/
multiqc_report.html
trimmed/
trimmed-fastqc/
```

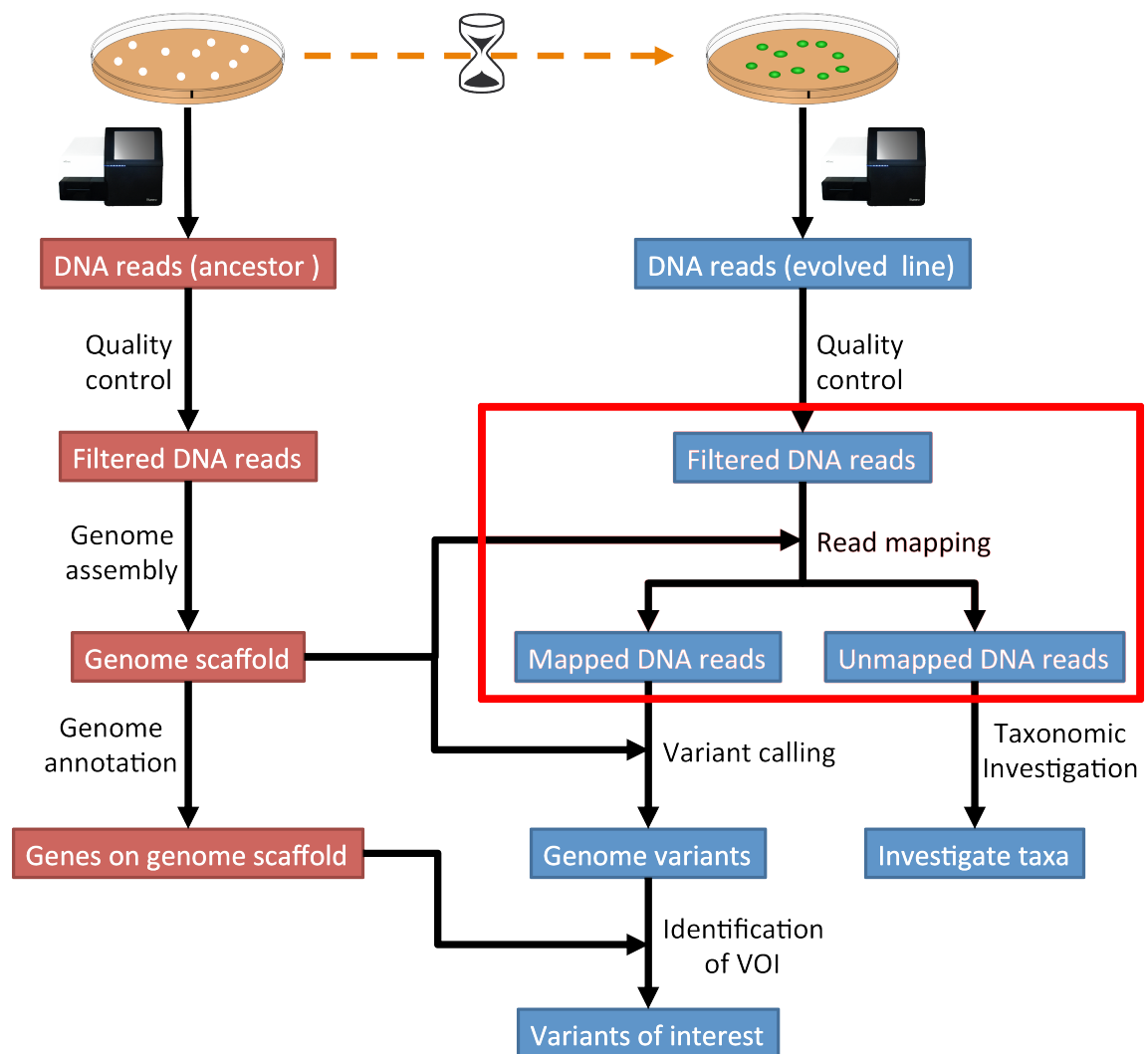


Fig. 5.1: The part of the workflow we will work on in this section marked in red.

Attention: If you have not run the previous sections on [Quality control](#) (page 9) and [Genome assembly](#) (page 19), you can download the trimmed data and the genome assembly needed for this section here: [Downloads](#) (page 77).

5.5 Mapping sequence reads to a reference genome

We want to map the sequencing reads to the ancestral reference genome. We are going to use the quality trimmed forward and backward DNA sequences of the evolved line and use a program called [BWA](#)⁷⁰ to map the reads.

Todo:

1. Discuss briefly why we are using the ancestral genome as a reference genome as opposed to a genome for the evolved line.
-

5.5.1 Downloading the reference genome assembly

Todo: In the assembly section at “[Genome assembly](#) (page 19)”, we created a genome assembly. However, we actually used sub-sampled data as otherwise the assemblies would have taken a long time to finish. To continue, please download the assembly created on the complete dataset ([Downloads](#) (page 77)). Unarchive and uncompress the files with `tar -xvzf assembly.tar.gz`.

5.5.2 Installing the software

We are going to use a program called [BWA](#)⁷¹ to map our reads to our genome.

It is simple to install and use.

```
$ conda create --yes -n mapping samtools bwa qualimap r-base
$ conda activate mapping
```

5.6 BWA

5.6.1 Overview

[BWA](#)⁷² is a short read aligner, that can take a reference genome and map single- or paired-end sequence data to it [LI2009]. It requires an indexing step in which one supplies the reference genome and [BWA](#)⁷³ will create an index that in the subsequent steps will be used for aligning the reads to the reference genome. While this step can take some time, the good thing is the index can be reused over and over. The general command structure of the [BWA](#)⁷⁴ tools we are going to use are shown below:

⁷⁰ <http://bio-bwa.sourceforge.net/>

⁷¹ <http://bio-bwa.sourceforge.net/>

⁷² <http://bio-bwa.sourceforge.net/>

⁷³ <http://bio-bwa.sourceforge.net/>

⁷⁴ <http://bio-bwa.sourceforge.net/>

```
# bwa index help
$ bwa index

# indexing
$ bwa index path/to/reference-genome.fa

# bwa mem help
$ bwa mem

# single-end mapping, general command structure, adjust to your case
$ bwa mem path/to/reference-genome.fa path/to/reads.fq.gz > path/to/aln-se.sam

# paired-end mapping, general command structure, adjust to your case
$ bwa mem path/to/reference-genome.fa path/to/read1.fq.gz path/to/read2.fq.gz > path/to/aln-pe.sam
```

5.6.2 Creating a reference index for mapping

Todo: Create an [BWA⁷⁵](#) index for our reference genome assembly. Attention! Remember which file you need to submit to [BWA⁷⁶](#).

Hint: Should you not get it right, try the commands in *Code: BWA indexing* (page 74).

Note: Should you be unable to run [BWA⁷⁷](#) indexing on the data, you can download the index from [Downloads](#) (page 77). Unarchive and uncompress the files with `tar -xvzf bwa-index.tar.gz`.

5.6.3 Mapping reads in a paired-end manner

Now that we have created our index, it is time to map the trimmed sequencing reads of our two evolved line to the reference genome.

Todo: Use the correct `bwa mem` command structure from above and map the reads of the two evolved line to the reference genome.

Hint: Should you not get it right, try the commands in *Code: BWA mapping* (page 74).

⁷⁵ <http://bio-bwa.sourceforge.net/>

⁷⁶ <http://bio-bwa.sourceforge.net/>

⁷⁷ <http://bio-bwa.sourceforge.net/>

Attention: The step of sam to bam-file conversion might take a few minutes to finish, depending on how big your mapping file is.

We will be using the [SAM flag](#)⁸⁴ information later below to extract specific alignments.

Hint: A very useful tools to explain flags can be found [here](#)⁸⁵.

Once we have bam-file, we can also delete the original sam-file as it requires too much space and we can always recreate it from the bam-file.

```
$ rm mappings/evol1.sam
```

5.8.2 Sorting

We are going to use [SAMtools](#)⁸⁶ again to sort the bam-file into **coordinate order**:

```
# convert to bam file and sort
$ samtools sort -O bam -o mappings/evol1.sorted.bam mappings/evol1.fixmate.bam

# Once it successfully finished, delete the fixmate file to save space
$ rm mappings/evol1.fixmate.bam
```

- -o: specifies the name of the output file.
- -O bam: specifies that the output will be bam-format

5.8.3 Remove duplicates

In this step we remove duplicate reads. The main purpose of removing duplicates is to mitigate the effects of PCR amplification bias introduced during library construction. **It should be noted that this step is not always recommended.** It depends on the research question. In SNP calling it is a good idea to remove duplicates, as the statistics used in the tools that call SNPs sub-sequently expect this (most tools anyways). However, for other research questions that use mapping, you might not want to remove duplicates, e.g. RNA-seq.

```
$ samtools markdup -r -S mappings/evol1.sorted.bam mappings/evol1.sorted.dedup.bam

# if it worked, delete the original file
$ rm mappings/evol1.sorted.bam
```

Todo: Figure out what “PCR amplification bias” means.

Note: Should you be unable to do the post-processing steps, you can download the mapped data from [Downloads](#) (page 77).

⁸⁴ <http://bio-bwa.sourceforge.net/bwa.shtml#4>

⁸⁵ <http://broadinstitute.github.io/picard/explain-flags.html>

⁸⁶ <http://samtools.sourceforge.net/>

5.9 Mapping statistics

5.9.1 Stats with SAMtools

Lets get an mapping overview:

```
$ samtools flagstat mappings/evol1.sorted.dedup.bam
```

Todo: Look at the mapping statistics and understand [their meaning](#)⁸⁷. Discuss your results. Explain why we may find mapped reads that have their mate mapped to a different chromosome/contig? Can they be used for something?

For the sorted bam-file we can get read depth for at all positions of the reference genome, e.g. how many reads are overlapping the genomic position.

```
$ samtools depth mappings/evol1.sorted.dedup.bam | gzip > mappings/evol1.depth.txt.gz
```

Todo: Extract the depth values for contig 20 and load the data into R, calculate some statistics of our scaffold.

```
$ zcat mappings/evol1.depth.txt.gz | egrep '^NODE_20_' | gzip > mappings/NODE_20.depth.txt.gz
```

Now we quickly use some [R](#)⁸⁸ to make a coverage plot for contig NODE20. Open a [R](#)⁸⁹ shell by typing R on the command-line of the shell.

```
x <- read.table('mappings/NODE_20.depth.txt.gz', sep='\t', header=FALSE, strip.white=TRUE)

# Look at the beginning of x
head(x)

# calculate average depth
mean(x[,3])
# std dev
sqrt(var(x[,3]))

# mark areas that have a coverage below 20 in red
plot(x[,2], x[,3], col = ifelse(x[,3] < 20, 'red', 'black'), pch=19, xlab='postion', ylab='coverage')

# to save a plot
png('mappings/covNODE20.png', width = 1200, height = 500)
plot(x[,2], x[,3], col = ifelse(x[,3] < 20, 'red', 'black'), pch=19, xlab='postion', ylab='coverage')
dev.off()
```

The result plot will be looking similar to the one in [Fig. 5.2](#)

Todo: Look at the created plot. Explain why it makes sense that you find relatively bad coverage at the beginning and the end of the contig.

⁸⁷ <https://www.biostars.org/p/12475/>

⁸⁸ <https://www.r-project.org/>

⁸⁹ <https://www.r-project.org/>

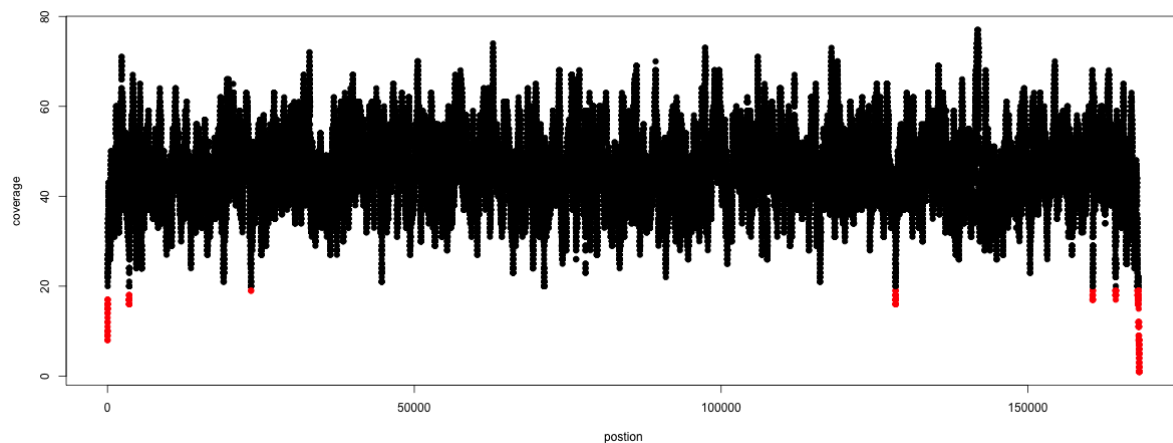


Fig. 5.2: A example coverage plot for a contig with highlighted in red regions with a coverage below 20 reads.

5.9.2 Stats with QualiMap

For a more in depth analysis of the mappings, one can use [QualiMap](#)⁹⁰ [OKO2015].

[QualiMap](#)⁹¹ examines sequencing alignment data in SAM/BAM files according to the features of the mapped reads and provides an overall view of the data that helps to detect biases in the sequencing and/or mapping of the data and eases decision-making for further analysis.

Run [QualiMap](#)⁹² with:

```
$ qualimap bamqc -bam mappings/ev01.sorted.dedup.bam
# Once finished open result page with
$ firefox mappings/ev01.sorted.dedup_stats/qualimapReport.html
```

This will create a report in the mapping folder. See this [webpage](#)⁹³ to get help on the sections in the report.

Todo: Investigate the mapping of the evolved sample. Write down your observations.

5.10 Sub-selecting reads

It is important to remember that the mapping commands we used above, without additional parameters to sub-select specific alignments (e.g. for [Bowtie2](#)⁹⁴ there are options like `--no-mixed`, which suppresses unpaired alignments for paired reads or `--no-discordant`, which suppresses discordant alignments for paired reads, etc.), are going to output all reads, including unmapped reads, multi-mapping reads, unpaired reads, discordant read pairs, etc. in one file. We can sub-select from the output reads we want to analyse further using [SAMtools](#)⁹⁵.

⁹⁰ <http://qualimap.bioinfo.cipf.es/>

⁹¹ <http://qualimap.bioinfo.cipf.es/>

⁹² <http://qualimap.bioinfo.cipf.es/>

⁹³ http://qualimap.bioinfo.cipf.es/doc_html/analysis.html#output

⁹⁴ <http://bowtie-bio.sourceforge.net/bowtie2/index.shtml>

⁹⁵ <http://samtools.sourceforge.net/>

Todo: Explain what concordant and discordant read pairs are? Look at the [Bowtie2⁹⁶](#) manual.

5.10.1 Concordant reads

We can select read-pair that have been mapped in a correct manner (same chromosome/contig, correct orientation to each other, distance between reads is not stupid).

Attention: We show the command here, but we are not going to use it.

```
$ samtools view -h -b -f 3 mappings/evol1.sorted.dedup.bam > mappings/evol1.sorted.dedup.concordant.
↪bam
```

- -b: Output will be bam-format
- -f 3: Only extract correctly paired reads. -f extracts alignments with the specified [SAM flag⁹⁷](#) set.

Todo: Our final aim is to identify variants. For a particular class of variants, it is not the best idea to only focus on concordant reads. Why is that?

5.10.2 Quality-based sub-selection

In this section we want to sub-select reads based on the quality of the mapping. It seems a reasonable idea to only keep good mapping reads. As the SAM-format contains at column 5 the *MAPQ* value, which we established earlier is the “MAPping Quality” in Phred-scaled, this seems easily achieved. The formula to calculate the *MAPQ* value is: $MAPQ = -10 * \log_{10}(p)$, where p is the probability that the read is mapped wrongly. However, there is a problem! **While the MAPQ information would be very helpful indeed, the way that various tools implement this value differs.** A good overview can be found [here⁹⁸](#). Bottom-line is that we need to be aware that different tools use this value in different ways and the it is good to know the information that is encoded in the value. Once you dig deeper into the mechanics of the *MAPQ* implementation it becomes clear that this is not an easy topic. If you want to know more about the *MAPQ* topic, please follow the link above.

For the sake of going forward, we will sub-select reads with at least medium quality as defined by [Bowtie2⁹⁹](#):

```
$ samtools view -h -b -q 20 mappings/evol1.sorted.dedup.bam > mappings/evol1.sorted.dedup.q20.bam
```

- -h: Include the sam header
- -q 20: Only extract reads with mapping quality ≥ 20

Hint: I will repeat here a recommendation given at the source [link¹⁰⁰](#) above, as it is a good one: If you unsure what *MAPQ* scoring scheme is being used in your own data then you can plot out the *MAPQ* distribution in a BAM file using programs like the mentioned [QualiMap¹⁰¹](#) or similar programs. This will at least show you the range and frequency with which different *MAPQ* values appear and may help identify a suitable threshold you may want to use.

⁹⁶ <http://bowtie-bio.sourceforge.net/bowtie2/index.shtml>

⁹⁷ <http://bio-bwa.sourceforge.net/bwa.shtml#4>

⁹⁸ <https://sequencing.qcfail.com/articles/mapq-values-are-really-useful-but-their-implementation-is-a-mess/>

⁹⁹ <http://bowtie-bio.sourceforge.net/bowtie2/index.shtml>

¹⁰⁰ <https://sequencing.qcfail.com/articles/mapq-values-are-really-useful-but-their-implementation-is-a-mess/>

¹⁰¹ <http://qualimap.bioinfo.cipf.es/>

Todo: Please repeat the whole process for the second evolved strain => mapping and post-processing.

Note: Should you be unable to process the second evolved strain look at the coding solutions here:
Code: Mapping post-processing (page 74)

5.10.3 Unmapped reads

We could decide to use [Kraken2](#)¹⁰² like in section *Taxonomic investigation* (page 35) to classify all unmapped sequence reads and identify the species they are coming from and test for contamination.

Lets see how we can get the unmapped portion of the reads from the bam-file:

```
$ samtools view -b -f 4 mappings/evol1.sorted.dedup.bam > mappings/evol1.sorted.unmapped.bam
# we are deleting the original to save space,
# however, in reality you might want to save it to investigate later
$ rm mappings/evol1.sorted.dedup.bam

# count the unmapped reads
$ samtools view -c mappings/evol1.sorted.unmapped.bam
```

- -b: indicates that the output is BAM.
- -f INT: only include reads with this [SAM flag](#)¹⁰³ set. You can also use the command `samtools flags` to get an overview of the flags.
- -c: count the reads

Lets extract the fastq sequence of the unmapped reads for read1 and read2.

```
$ samtools fastq -1 mappings/evol1.sorted.unmapped.R1.fastq.gz -2 mappings/evol1.sorted.unmapped.R2.
↪ fastq.gz mappings/evol1.sorted.unmapped.bam
# delete not needed files
$ rm mappings/evol1.sorted.unmapped.bam
```

¹⁰² <https://www.ccb.jhu.edu/software/kraken2/>

¹⁰³ <http://bio-bwa.sourceforge.net/bwa.shtml#4>

TAXONOMIC INVESTIGATION

6.1 Preface

We want to investigate if there are sequences of other species in our collection of sequenced DNA pieces. We hope that most of them are from our species that we try to study, i.e. the DNA that we have extracted and amplified. This might be a way of quality control, e.g. have the samples been contaminated? Lets investigate if we find sequences from other species in our sequence set.

We will use the tool [Kraken2](#)¹⁰⁶ to assign taxonomic classifications to our sequence reads. Let us see if we can id some sequences from other species.

Note: You will encounter some **To-do** sections at times. Write the solutions and answers into a text-file.

6.2 Overview

The part of the workflow we will work on in this section can be viewed in [Fig. 6.1](#).

6.3 Before we start

Lets see how our directory structure looks so far:

```
$ cd ~/analysis
$ ls -1F
```

```
assembly/
data/
mappings/
multiqc_data
trimmed/
trimmed-fastqc/
```

Attention: If you have not run the previous section [Read mapping](#) (page 25), you can download the unmapped sequencing data needed for this section here: [Downloads](#) (page 77).

¹⁰⁶ <https://www.ccb.jhu.edu/software/kraken2/>

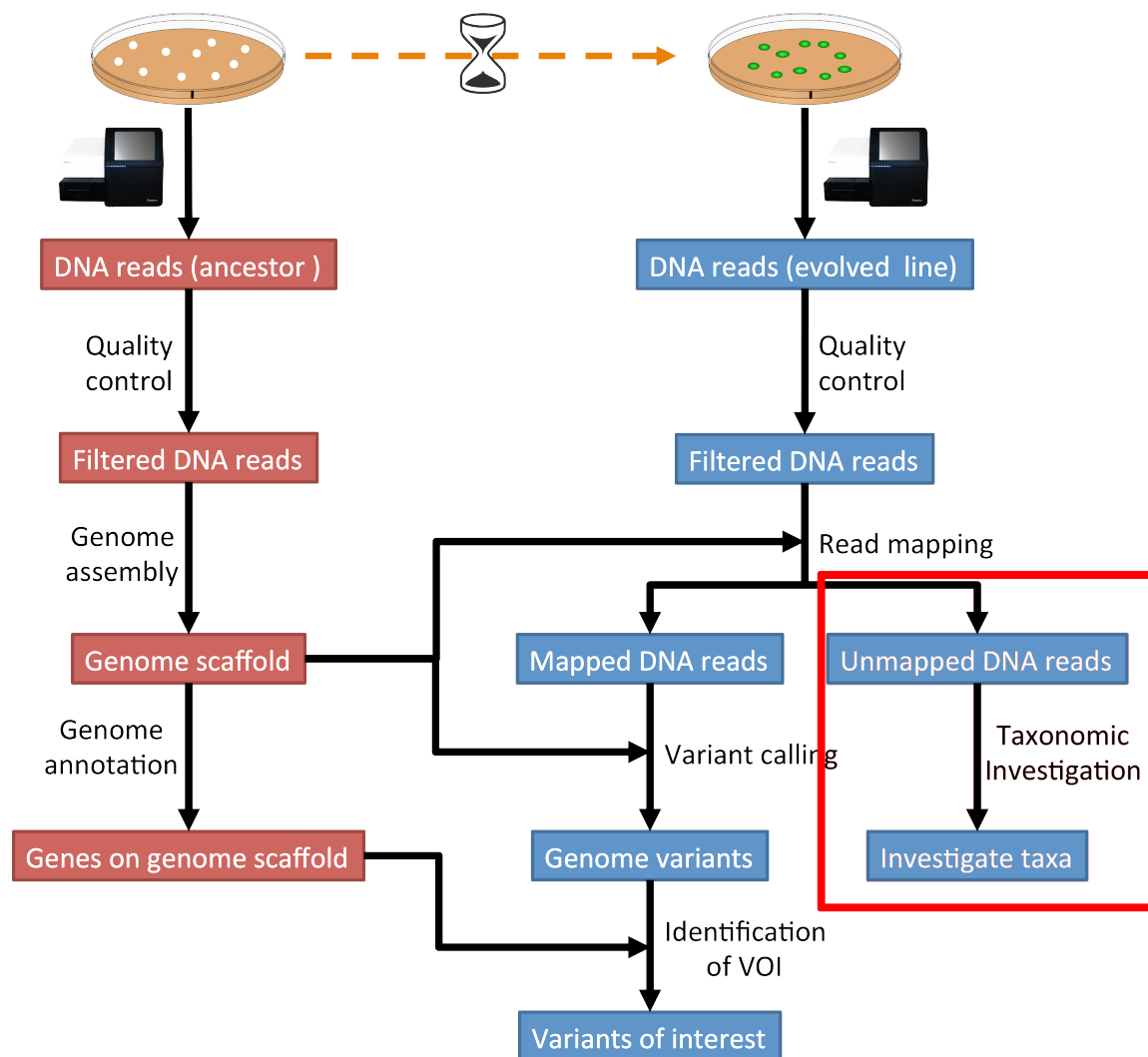


Fig. 6.1: The part of the workflow we will work on in this section marked in red.

6.4 Kraken2

We will be using a tool called [Kraken2](#)¹⁰⁷ [WOOD2014]. This tool uses k-mers to assign a taxonomic labels in form of [NCBI Taxonomy](#)¹⁰⁸ to the sequence (if possible). The taxonomic label is assigned based on similar k-mer content of the sequence in question to the k-mer content of reference genome sequence. The result is a classification of the sequence in question to the most likely taxonomic label. If the k-mer content is not similar to any genomic sequence in the database used, it will not assign any taxonomic label.

6.4.1 Installation

Use conda in the same fashion as before to install [Kraken2](#)¹⁰⁹. However, we are going to install kraken into its own environment:

```
$ conda create --yes -n kraken kraken2 bracken
$ conda activate kraken
```

Now we create a directory where we are going to do the analysis and we will change into that directory too.

```
# make sure you are in your analysis root folder
$ cd ~/analysis

# create dir
$ mkdir kraken
$ cd kraken
```

Now we need to create or download a [Kraken2](#)¹¹⁰ database that can be used to assign the taxonomic labels to sequences. We opt for downloading the pre-build “minikraken2” database from the [Kraken2](#)¹¹¹ website:

```
$ curl -O ftp://ftp.ccb.jhu.edu/pub/data/kraken2_dbs/minikraken2_v2_8GB_201904_UPDATE.tgz

# alternatively we can use wget
$ wget ftp://ftp.ccb.jhu.edu/pub/data/kraken2_dbs/minikraken2_v2_8GB_201904_UPDATE.tgz

# once the download is finished, we need to extract the archive content:
$ tar -xvzf minikraken2_v2_8GB_201904_UPDATE.tgz
```

Attention: Should the download fail. Please find links to alternative locations on the [Downloads](#) (page 77) page.

Note: The “minikraken2” database was created from bacteria, viral and archaea sequences. What are the implications for us when we are trying to classify our sequences?

¹⁰⁷ <https://www.ccb.jhu.edu/software/kraken2/>

¹⁰⁸ <https://www.ncbi.nlm.nih.gov/taxonomy>

¹⁰⁹ <https://www.ccb.jhu.edu/software/kraken2/>

¹¹⁰ <https://www.ccb.jhu.edu/software/kraken2/>

¹¹¹ <https://www.ccb.jhu.edu/software/kraken2/>

6.4.2 Usage

Now that we have installed [Kraken2¹¹²](#) and downloaded and extracted the minikraken2 database, we can attempt to investigate the sequences we got back from the sequencing provider for other species as the one it should contain. We call the [Kraken2¹¹³](#) tool and specify the database and fasta-file with the sequences it should use. The general command structure looks like this:

```
$ kraken2 --use-names --threads 4 --db PATH_TO_DB_DIR --report example.report.txt example.fa >
↳ example.kraken
```

However, we may have fastq-files, so we need to use `--fastq-input` which tells [Kraken2¹¹⁴](#) that it is dealing with fastq-formatted files. The `--gzip-compressed` flag specifies that input-files are compressed. In addition, we are dealing with paired-end data, which we can tell [Kraken2¹¹⁵](#) with the switch `--paired`. Here, we are investigating one of the unmapped paired-end read files of the evolved line.

```
$ kraken2 --use-names --threads 4 --db minikraken2_v2_8GB_201904_UPDATE --fastq-input --report
↳ evol1 --gzip-compressed --paired ../mappings/evol1.sorted.unmapped.R1.fastq.gz ../mappings/evol1.
↳ sorted.unmapped.R2.fastq.gz > evol1.kraken
```

This classification may take a while, depending on how many sequences we are going to classify. The resulting content of the file `evol1.kraken` looks similar to the following example:

```
C      7001326F:121:CBVVLANXX:1:1105:2240:12640      816      251      816:9 171549:5 816:5
↳ 171549:3 2:2 816:5 171549:4 816:34 171549:8 816:4 171549:2 816:10 A:35 816:10 171549:2 816:4
↳ 171549:8 816:34 171549:4 816:5 2:2 171549:3 816:5 171549:5 816:9
C      7001326F:121:CBVVLANXX:1:1105:3487:12536      1339337 202      1339337:67 A:35 1339337:66
U      7001326F:121:CBVVLANXX:1:1105:5188:12504      0      251      0:91 A:35 0:91
U      7001326F:121:CBVVLANXX:1:1105:11030:12689      0      251      0:91 A:35 0:91
U      7001326F:121:CBVVLANXX:1:1105:7157:12806      0      206      0:69 A:35 0:68
```

Each sequence classified by [Kraken2¹¹⁶](#) results in a single line of output. Output lines contain five tab-delimited fields; from left to right, they are:

1. C/U: one letter code indicating that the sequence was either classified or unclassified.
2. The sequence ID, obtained from the FASTA/FASTQ header.
3. The taxonomy ID [Kraken2¹¹⁷](#) used to label the sequence; this is **0** if the sequence is unclassified and otherwise should be the [NCBI Taxonomy¹¹⁸](#) identifier.
4. The length of the sequence in bp.
5. A space-delimited list indicating the lowest common ancestor (in the taxonomic tree) mapping of each k-mer in the sequence. For example, `562:13 561:4 A:31 0:1 562:3` would indicate that:
 - the first 13 k-mers mapped to taxonomy ID #562
 - the next 4 k-mers mapped to taxonomy ID #561
 - the next 31 k-mers contained an ambiguous nucleotide
 - the next k-mer was not in the database
 - the last 3 k-mers mapped to taxonomy ID #562

¹¹² <https://www.ccb.jhu.edu/software/kraken2/>

¹¹³ <https://www.ccb.jhu.edu/software/kraken2/>

¹¹⁴ <https://www.ccb.jhu.edu/software/kraken2/>

¹¹⁵ <https://www.ccb.jhu.edu/software/kraken2/>

¹¹⁶ <https://www.ccb.jhu.edu/software/kraken2/>

¹¹⁷ <https://www.ccb.jhu.edu/software/kraken2/>

¹¹⁸ <https://www.ncbi.nlm.nih.gov/taxonomy>

Note: The [Kraken2](#)¹¹⁹ manual can be accessed [here](#)¹²⁰.

6.4.3 Investigate taxa

We can use the webpage [NCBI TaxIdentifier](#)¹²¹ to quickly get the names to the taxonomy identifier. However, this is impractical as we are dealing potentially with many sequences. [Kraken2](#)¹²² has some scripts that help us understand our results better.

Because we used the [Kraken2](#)¹²³ switch `--report FILE`, we have got also a sample-wide report of all taxa found. This is much better to get an overview what was found.

The first few lines of an example report are shown below.

83.56	514312	514312	U	0	unclassified
16.44	101180	0	R	1	root
16.44	101180	0	R1	131567	cellular organisms
16.44	101180	2775	D	2	Bacteria
13.99	86114	1	D1	1783270	FCB group
13.99	86112	0	D2	68336	Bacteroidetes/Chlorobi group
13.99	86103	8	P	976	Bacteroidetes
13.94	85798	2	C	200643	Bacteroidia
13.94	85789	19	O	171549	Bacteroidales
13.87	85392	0	F	815	Bacteroidaceae

The output of `kraken-report` is tab-delimited, with one line per taxon. The fields of the output, from left-to-right, are as follows:

1. **Percentage** of reads covered by the clade rooted at this taxon
2. **Number of reads** covered by the clade rooted at this taxon
3. **Number of reads** assigned directly to this taxon
4. A rank code, indicating (U)nclassified, (D)omain, (K)ingdom, (P)hylum, (C)lass, (O)rder, (F)amily, (G)enus, or (S)pecies. All other ranks are simply “-“.
5. [NCBI Taxonomy](#)¹²⁴ ID
6. The indented scientific name

Note: If you want to compare the taxa content of different samples to another, one can create a report whose structure is always the same for all samples, disregarding which taxa are found (obviously the percentages and numbers will be different).

We can create such a report using the option `--report-zero-counts` which will print out all taxa (instead of only those found). We then sort the taxa according to taxa-ids (column 5), e.g. `sort -n -k5`.

The report is not ordered according to taxa ids and contains all taxa in the database, even if they have not been found in our sample and are thus zero. The columns are the same as in the former report, however, we have more rows and they are now differently sorted, according to the [NCBI Taxonomy](#)¹²⁵ id.

¹¹⁹ <https://www.ccb.jhu.edu/software/kraken2/>

¹²⁰ <https://www.ccb.jhu.edu/software/kraken2/index.shtml?t=manual>

¹²¹ https://www.ncbi.nlm.nih.gov/Taxonomy/TaxIdentifier/tax_identifier.cgi

¹²² <https://www.ccb.jhu.edu/software/kraken2/>

¹²³ <https://www.ccb.jhu.edu/software/kraken2/>

¹²⁴ <https://www.ncbi.nlm.nih.gov/taxonomy>

¹²⁵ <https://www.ncbi.nlm.nih.gov/taxonomy>

6.4.4 Bracken

Bracken¹²⁶ stands for Bayesian Re-estimation of Abundance with Kraken, and is a statistical method that computes the abundance of species in DNA sequences from a metagenomics sample [LU2017]. **Bracken**¹²⁷ uses the taxonomy labels assigned by **Kraken2**¹²⁸ (see above) to estimate the number of reads originating from each species present in a sample. **Bracken**¹²⁹ classifies reads to the best matching location in the taxonomic tree, but does not estimate abundances of species. Combined with the Kraken classifier, **Bracken**¹³⁰ will produce more accurate species- and genus-level abundance estimates than **Kraken2**¹³¹ alone.

The use of **Bracken**¹³² subsequent to **Kraken2**¹³³ is optional but might improve on the **Kraken2**¹³⁴ results.

Installation

We installed **Bracken**¹³⁵ already together with **Kraken2**¹³⁶ above, so it should be ready to be used. We also downloaded the **Bracken**¹³⁷ files together with the minikraken2 database above, so we are good to go.

Usage

Now, we can use **Bracken**¹³⁸ on the **Kraken2**¹³⁹ results to improve them.

The general structure of the **Bracken**¹⁴⁰ command look like this:

```
$ bracken -d PATH_TO_DB_DIR -i kraken2.report -o bracken.species.txt -l S
```

- -l S: denotes the level we want to look at. S stands for species but other levels are available.
- -d PATH_TO_DB_DIR: specifies the path to the **Kraken2**¹⁴¹ database that should be used.

Let us apply **Bracken**¹⁴² to the example above:

```
$ bracken -d minikraken2_v2_8GB_201904_UPDATE -i evol1.kraken -l S -o evol1.bracken
```

The species-focused result-table looks similar to this:

name	taxonomy_id	taxonomy_lvl	kraken_assigned_reads	added_reads	new_est_reads	
↩ fraction_total_reads						
Streptococcus sp. oral	taxon 431	712633	S	2	0	2 0.00001
Neorhizobium sp. NCHU2750	1825976	S	0	0	0	0.00000
Pseudomonas sp. MT-1	150396	S	0	0	0	0.00000
Ahniella affigens	2021234	S	1	0	1	0.00000
Sinorhizobium sp. CCBAU 05631	794846	S	0	0	0	0.00000
Cohnella sp. 18JY8-7	2480923	S	1	0	1	0.00000

(continues on next page)

¹²⁶ <https://ccb.jhu.edu/software/bracken/index.shtml>

¹²⁷ <https://ccb.jhu.edu/software/bracken/index.shtml>

¹²⁸ <https://www.ccb.jhu.edu/software/kraken2/>

¹²⁹ <https://ccb.jhu.edu/software/bracken/index.shtml>

¹³⁰ <https://ccb.jhu.edu/software/bracken/index.shtml>

¹³¹ <https://www.ccb.jhu.edu/software/kraken2/>

¹³² <https://ccb.jhu.edu/software/bracken/index.shtml>

¹³³ <https://www.ccb.jhu.edu/software/kraken2/>

¹³⁴ <https://www.ccb.jhu.edu/software/kraken2/>

¹³⁵ <https://ccb.jhu.edu/software/bracken/index.shtml>

¹³⁶ <https://www.ccb.jhu.edu/software/kraken2/>

¹³⁷ <https://ccb.jhu.edu/software/bracken/index.shtml>

¹³⁸ <https://ccb.jhu.edu/software/bracken/index.shtml>

¹³⁹ <https://www.ccb.jhu.edu/software/kraken2/>

¹⁴⁰ <https://ccb.jhu.edu/software/bracken/index.shtml>

¹⁴¹ <https://www.ccb.jhu.edu/software/kraken2/>

¹⁴² <https://ccb.jhu.edu/software/bracken/index.shtml>

(continued from previous page)

Bacillus velezensis	492670	S	4	4	8	0.00002
Actinoplanes missouriensis	1866	S	2	8	10	0.00002

The important column is the `new_est_reads`, which gives the newly estimated reads.

6.5 Centrifuge

We can also use another tool by the same group called [Centrifuge](#)¹⁴³ [KIM2017]. This tool uses a novel indexing scheme based on the Burrows-Wheeler transform (BWT) and the Ferragina-Manzini (FM) index, optimized specifically for the metagenomic classification problem to assign a taxonomic labels in form of [NCBI Taxonomy](#)¹⁴⁴ to the sequence (if possible). The result is a classification of the sequence in question to the most likely taxonomic label. If the search sequence is not similar to any genomic sequence in the database used, it will not assign any taxonomic label.

Note: I would normally use [Kraken2](#)¹⁴⁵ and only prefer [Centrifuge](#)¹⁴⁶ if memory and/or speed are an issue .

6.5.1 Installation

Use conda in the same fashion as before to install [Centrifuge](#)¹⁴⁷:

```
$ conda create --yes -n centrifuge centrifuge
$ conda activate centrifuge
```

Now we create a directory where we are going to do the analysis and we will change into that directory too.

```
# make sure you are in your analysis root folder
$ cd ~/analysis

# create dir
$ mkdir centrifuge
$ cd centrifuge
```

Now we need to create or download a [Centrifuge](#)¹⁴⁸ database that can be used to assign the taxonomic labels to sequences. We opt for downloading the pre-build database from the [Centrifuge](#)¹⁴⁹ website:

```
$ curl -O ftp://ftp.ccb.jhu.edu/pub/infphilo/centrifuge/data/p_compressed+h+v.tar.gz

# # alternatively we can use wget
$ wget ftp://ftp.ccb.jhu.edu/pub/infphilo/centrifuge/data/p_compressed+h+v.tar.gz

# once the download is finished, we need to extract the archive content
# It will extract a few files from the archive and may take a moment to finish.
$ tar -xvzf p_compressed+h+v.tar.gz
```

¹⁴³ <http://www.ccb.jhu.edu/software/centrifuge/index.shtml>

¹⁴⁴ <https://www.ncbi.nlm.nih.gov/taxonomy>

¹⁴⁵ <https://www.ccb.jhu.edu/software/kraken2/>

¹⁴⁶ <http://www.ccb.jhu.edu/software/centrifuge/index.shtml>

¹⁴⁷ <http://www.ccb.jhu.edu/software/centrifuge/index.shtml>

¹⁴⁸ <http://www.ccb.jhu.edu/software/centrifuge/index.shtml>

¹⁴⁹ <http://www.ccb.jhu.edu/software/centrifuge/index.shtml>

Attention: Should the download fail. Please find links to alternative locations on the [Downloads](#) (page 77) page.

Note: The database we will be using was created from bacteria and archaea sequences only. What are the implications for us when we are trying to classify our sequences?

6.5.2 Usage

Now that we have installed [Centrifuge](#)¹⁵⁰ and downloaded and extracted the pre-build database, we can attempt to investigate the sequences we got back from the sequencing provider for other species as the one it should contain. We call the [Centrifuge](#)¹⁵¹ tool and specify the database and fastq-files with the sequences it should use. The general command structure looks like this:

```
$ centrifuge -x p_compressed+h+v -1 example.1.fq -2 example.2.fq -U single.fq --report-file report.
↳txt -S results.txt
```

Here, we are investigating paired-end read files of the evolved line.

```
$ centrifuge -x p_compressed+h+v -1 ../mappings/evol1.sorted.unmapped.R1.fastq -2 ../mappings/
↳evol1.sorted.unmapped.R2.fastq --report-file evol1-report.txt -S evol1-results.txt
```

This classification may take a moment, depending on how many sequences we are going to classify. The resulting content of the file `evol1-results.txt` looks similar to the following example:

readID	seqID	taxID	score	2ndBestScore	hitLength	queryLength	numMatches
M02810:197:000000000-AV55U:1:1101:15316:8461					cid 1747	1747 1892	0
↳103	135	1					
M02810:197:000000000-AV55U:1:1101:15563:3249					cid 161879	161879 18496	0
↳151	151	1					
M02810:197:000000000-AV55U:1:1101:19743:5166					cid 564 564	10404 10404	117
↳151	2						
M02810:197:000000000-AV55U:1:1101:19743:5166					cid 562 562	10404 10404	117
↳151	2						

Each sequence classified by [Centrifuge](#)¹⁵² results in a single line of output. Output lines contain eight tab-delimited fields; from left to right, they are according to the [Centrifuge](#)¹⁵³ website:

1. The read ID from a raw sequencing read.
2. The sequence ID of the genomic sequence, where the read is classified.
3. The taxonomic ID of the genomic sequence in the second column.
4. The score for the classification, which is the weighted sum of hits.
5. The score for the next best classification.
6. A pair of two numbers: (1) an approximate number of base pairs of the read that match the genomic sequence and (2) the length of a read or the combined length of mate pairs.
7. A pair of two numbers: (1) an approximate number of base pairs of the read that match the genomic sequence and (2) the length of a read or the combined length of mate pairs.
8. The number of classifications for this read, indicating how many assignments were made.

¹⁵⁰ <http://www.ccb.jhu.edu/software/centrifuge/index.shtml>

¹⁵¹ <http://www.ccb.jhu.edu/software/centrifuge/index.shtml>

¹⁵² <http://www.ccb.jhu.edu/software/centrifuge/index.shtml>

¹⁵³ <http://www.ccb.jhu.edu/software/centrifuge/index.shtml>

6.5.3 Investigate taxa

Centrifuge report

The command above creates a [Centrifuge](http://www.ccb.jhu.edu/software/centrifuge/index.shtml)¹⁵⁴ report automatically for us. It contains an overview of the identified taxa and their abundances in your supplied sequences (normalised to genomic length):

name	taxID	taxRank	genomeSize	numReads	numUniqueReads	abundance
<i>Pseudomonas aeruginosa</i>	287	species	22457305	1	0	0.0
<i>Pseudomonas fluorescens</i>	294	species	14826544	1	1	0.0
<i>Pseudomonas putida</i>	303	species	6888188	1	0.0	
<i>Ralstonia pickettii</i>	329	species	6378979	3	2	0.0
<i>Pseudomonas pseudoalcaligenes</i>	330	species	4691662	1	1	0.0171143

Each line contains seven tab-delimited fields; from left to right, they are according to the [Centrifuge](http://www.ccb.jhu.edu/software/centrifuge/index.shtml)¹⁵⁵ website:

1. The name of a genome, or the name corresponding to a taxonomic ID (the second column) at a rank higher than the strain.
2. The taxonomic ID.
3. The taxonomic rank.
4. The length of the genome sequence.
5. The number of reads classified to this genomic sequence including multi-classified reads.
6. The number of reads uniquely classified to this genomic sequence.
7. The proportion of this genome normalized by its genomic length.

Kraken-like report

If we would like to generate a report as generated with the former tool [Kraken2](https://www.ccb.jhu.edu/software/kraken2/)¹⁵⁶, we can do it like this:

```
$ centrifuge-kreport -x p_compressed+h+v evolved-6-R1-results.txt > evolved-6-R1-kreport.txt
```

0.00	0	0	U	0	unclassified
78.74	163	0	-	1	root
78.74	163	0	-	131567	cellular organisms
78.74	163	0	D	2	Bacteria
54.67	113	0	P	1224	Proteobacteria
36.60	75	0	C	1236	Gammaproteobacteria
31.18	64	0	O	91347	Enterobacterales
30.96	64	0	F	543	Enterobacteriaceae
23.89	49	0	G	561	Escherichia
23.37	48	48	S	562	Escherichia coli
0.40	0	0	S	564	Escherichia fergusonii
0.12	0	0	S	208962	Escherichia albertii
3.26	6	0	G	570	Klebsiella
3.14	6	6	S	573	Klebsiella pneumoniae
0.12	0	0	S	548	[Enterobacter] aerogenes
2.92	6	0	G	620	Shigella
1.13	2	2	S	623	Shigella flexneri
0.82	1	1	S	624	Shigella sonnei
0.50	1	1	S	1813821	Shigella sp. PAMC 28760
0.38	0	0	S	621	Shigella boydii

¹⁵⁴ <http://www.ccb.jhu.edu/software/centrifuge/index.shtml>

¹⁵⁵ <http://www.ccb.jhu.edu/software/centrifuge/index.shtml>

¹⁵⁶ <https://www.ccb.jhu.edu/software/kraken2/>

This gives a similar (not the same) report as the [Kraken2¹⁵⁷](#) tool. The report is tab-delimited, with one line per taxon. The fields of the output, from left-to-right, are as follows:

1. Percentage of reads covered by the clade rooted at this taxon
2. Number of reads covered by the clade rooted at this taxon
3. Number of reads assigned directly to this taxon
4. A rank code, indicating (U)nclassified, (D)omain, (K)ingdom, (P)hylum, (C)lass, (O)rder, (F)amily, (G)enus, or (S)pecies. All other ranks are simply “-”.
5. NCBI Taxonomy ID
6. The indented scientific name

6.6 Visualisation (Krona)

We use the [Krona¹⁵⁸](#) tools to create a nice interactive visualisation of the taxa content of our sample [ONDOV2011]. [Fig. 6.2](#) shows an example (albeit an artificial one) snapshot of the visualisation [Krona¹⁵⁹](#) provides. [Fig. 6.2](#) is a snapshot of the [interactive web-page](#) similar to the one we try to create.

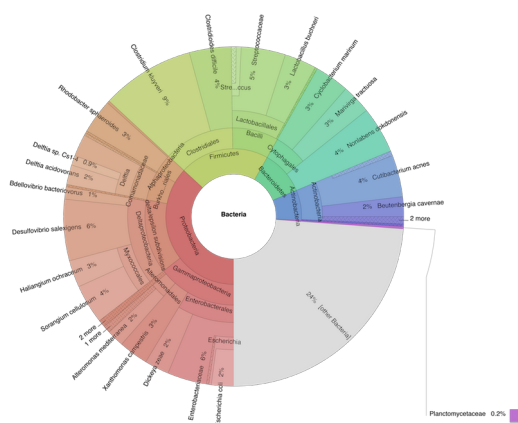


Fig. 6.2: Example of an Krona output webpage.

6.6.1 Installation

Install [Krona¹⁶⁰](#) with:

```
$ conda create --yes -n krona krona
$ conda activate krona
```

First some house-keeping to make the [Krona¹⁶¹](#) installation work. Do not worry to much about what is happening here.

```
# we delete a symbolic link that is not correct
$ rm -rf ~/miniconda3/envs/ngs/opt/krona/taxonomy

# we create a directory in our home where the krona database will live
```

(continues on next page)

¹⁵⁷ <https://www.ccb.jhu.edu/software/kraken2/>

¹⁵⁸ <https://github.com/marbl/Krona/wiki>

¹⁵⁹ <https://github.com/marbl/Krona/wiki>

¹⁶⁰ <https://github.com/marbl/Krona/wiki>

¹⁶¹ <https://github.com/marbl/Krona/wiki>

(continued from previous page)

```
$ mkdir -p ~/krona/taxonomy

# now we make a symbolic link to that directory
$ ln -s ~/krona/taxonomy ~/miniconda3/envs/ngs/opt/krona/taxonomy
```

6.6.2 Build the taxonomy

We need to build a taxonomy database for [Krona](#)¹⁶². However, if this fails we will skip this step and just download a pre-build one. Lets first try to build one.

```
$ ktUpdateTaxonomy.sh ~/krona/taxonomy
```

Now, if this fails, we download a pre-build taxonomy database for krona:

```
# Download pre-build database
$ curl -O http://compbio.massey.ac.nz/data/taxonomy.tab.gz

# we unzip the file
$ gzip -d taxonomy.tab.gz

# we move the unzipped file to our taxonomy directory we specified in the step before.
$ mv taxonomy.tab ~/krona/taxonomy
```

Attention: Should this also fail we can download a pre-build database on the [Downloads](#) (page 77) page via a browser.

6.6.3 Visualise

Now, we use the tool `ktImportTaxonomy` from the [Krona](#)¹⁶³ tools to create the html web-page. We first need build a two column file (`read_id<tab>tax_id`) as input to the `ktImportTaxonomy` tool. We will do this by cutting the columns out of either the [Kraken2](#)¹⁶⁴ or [Centrifuge](#)¹⁶⁵ results:

```
# Kraken2
$ cd kraken
$ cat evl1.kraken | cut -f 2,3 > evl1.kraken.krona
$ ktImportTaxonomy evl1.kraken.krona
$ firefox taxonomy.krona.html

# Centrifuge
$ cd centrifuge
$ cat evl1-results.txt | cut -f 1,3 > evl1-results.krona
$ ktImportTaxonomy evl1-results.krona
$ firefox taxonomy.krona.html
```

What happens here is that we extract the second and third column from the [Kraken2](#)¹⁶⁶ results. Afterwards, we input these to the [Krona](#)¹⁶⁷ script, and open the resulting web-page in a browser. Done!

¹⁶² <https://github.com/marbl/Krona/wiki>

¹⁶³ <https://github.com/marbl/Krona/wiki>

¹⁶⁴ <https://www.ccb.jhu.edu/software/kraken2/>

¹⁶⁵ <http://www.ccb.jhu.edu/software/centrifuge/index.shtml>

¹⁶⁶ <https://www.ccb.jhu.edu/software/kraken2/>

¹⁶⁷ <https://github.com/marbl/Krona/wiki>

VARIANT CALLING

7.1 Preface

In this section we will use our genome assembly based on the ancestor and call genetic variants in the evolved line [NIELSEN2011].

7.2 Overview

The part of the workflow we will work on in this section can be viewed in [Fig. 7.1](#).

7.3 Learning outcomes

After studying this tutorial section you should be able to:

#. Use tools to call variants based on a reference genome. #, Be able to describe what influences the calling of variants.

7.4 Before we start

Lets see how our directory structure looks so far:

```
$ cd ~/analysis
$ ls -lF
```

```
assembly/
data/
kraken/
mappings/
multiqc_data/
trimmed/
trimmed-fastqc/
```

Attention: If you have not run the previous section on [Read mapping](#) (page 25), you can download the mapped data needed for this section here: [Downloads](#) (page 77).

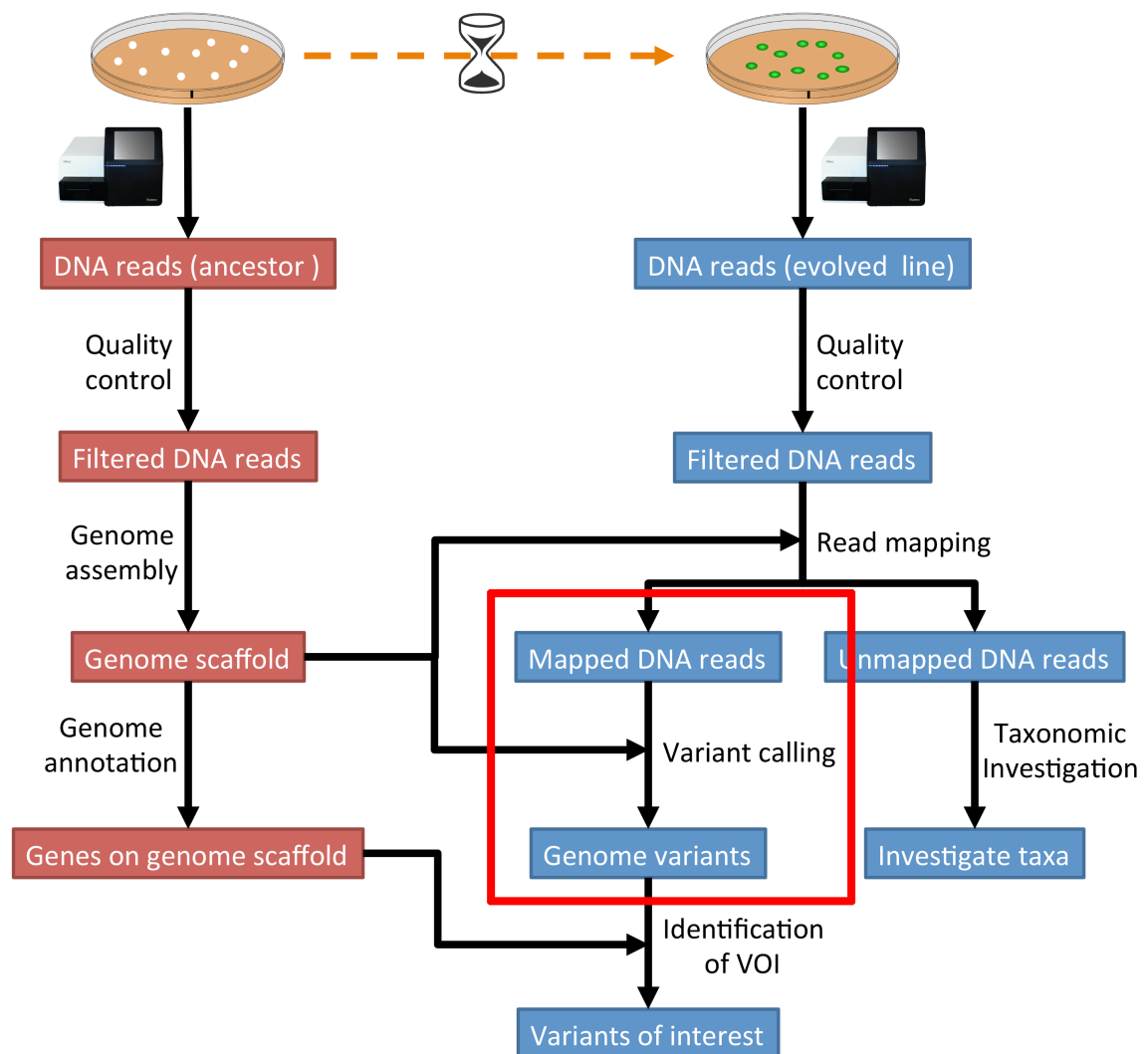


Fig. 7.1: The part of the workflow we will work on in this section marked in red.

7.5 Installing necessary software

Tools we are going to use in this section and how to install them if you have not done it yet.

```
# activate the env
$ conda create --yes -n var samtools bamtools freebayes bedtools vcflib rtg-tools bcftools_
↪matplotlib
$ conda activate var
```

7.6 Preprocessing

We first need to make an index of our reference genome as this is required by the SNP caller. Given a scaffold/contig file in fasta-format, e.g. scaffolds.fasta which is located in the directory assembly/, use [SAMtools](#)¹⁷² to do this:

```
$ samtools faidx assembly/scaffolds.fasta
```

Furthermore we need to pre-process our mapping files a bit further and create a bam-index file (.bai) for the bam-file we want to work with:

```
$ bamtools index -in mappings/evol1.sorted.dedup.q20.bam
```

Let's also create a new directory for the variants:

```
$ mkdir variants
```

7.7 Calling variants

7.7.1 Freebayes

We can call variants with a tool called [freebayes](#)¹⁷³. Given a reference genome scaffold file in fasta-format, e.g. scaffolds.fasta and the index in .fai format and a mapping file (.bam file) and a mapping index (.bai file), we can call variants with [freebayes](#)¹⁷⁴ like so:

```
# Now we call variants and pipe the results into a new file
$ freebayes -p 1 -f assembly/scaffolds.fasta mappings/evol1.sorted.dedup.q20.bam > variants/evol1.
↪freebayes.vcf
```

- -p 1: specifies the ploidy level. *E.Coli* are haploid.

7.8 Post-processing

7.8.1 Understanding the output files (.vcf)

Let's look at a vcf-file:

```
# first 10 lines, which are part of the header
$ cat variants/evol1.freebayes.vcf | head
```

¹⁷² <http://samtools.sourceforge.net/>

¹⁷³ <https://github.com/ekg/freebayes>

¹⁷⁴ <https://github.com/ekg/freebayes>

```
##fileformat=VCFv4.2
##fileDate=20200122
##source=freeBayes v1.3.1-dirty
##reference=assembly/scaffolds.fasta
##contig=<ID=NODE_1_length_348724_cov_30.410613,length=348724>
##contig=<ID=NODE_2_length_327290_cov_30.828326,length=327290>
##contig=<ID=NODE_3_length_312063_cov_30.523209,length=312063>
##contig=<ID=NODE_4_length_202800_cov_31.500777,length=202800>
##contig=<ID=NODE_5_length_164027_cov_28.935175,length=164027>
##contig=<ID=NODE_6_length_144088_cov_29.907986,length=144088>
```

Lets look at the variants:

```
# remove header lines and look at top 4 entires
$ cat variants/evoll.freebayes.vcf | grep -v '##' | head -4
```

```
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT unknown
NODE_1_length_348724_cov_30.410613 375 . A C 0 . AB=0;ABP=0;
AC=0;AF=0;AN=1;AO=3;CIGAR=1X;DP=21;DPB=21;DPRA=0;EPP=3.73412;EPPR=3.49285;GTI=0;LEN=1;MEANALT=1;
MQM=44;MQMR=40.3333;NS=1;NUMALT=1;ODDS=63.5226;PAIRED=1;PAIREDR=1;PAO=0;PQA=0;PQR=0;PRO=0;QA=53;
QR=414;RO=18;RPL=2;RPP=3.73412;RPPR=7.35324;RPR=1;RUN=1;SAF=3;SAP=9.52472;SAR=0;SRF=14;SRP=15.074;
SRR=4;TYPE=snp GT:DP:AD:RO:QR:AO:QA:GL 0:21:18,3:18:414:3:53:0,-29.6927
NODE_1_length_348724_cov_30.410613 393 . T A 0 . AB=0;ABP=0;
AC=0;AF=0;AN=1;AO=2;CIGAR=1X;DP=24;DPB=24;DPRA=0;EPP=7.35324;EPPR=6.56362;GTI=0;LEN=1;MEANALT=1;
MQM=36;MQMR=42.9545;NS=1;NUMALT=1;ODDS=127.074;PAIRED=1;PAIREDR=1;PAO=0;PQA=0;PQR=0;PRO=0;QA=21;
QR=717;RO=22;RPL=2;RPP=7.35324;RPPR=3.0103;RPR=0;RUN=1;SAF=2;SAP=7.35324;SAR=0;SRF=17;SRP=17.2236;
SRR=5;TYPE=snp GT:DP:AD:RO:QR:AO:QA:GL 0:24:22,2:22:717:2:21:0,-57.4754
NODE_1_length_348724_cov_30.410613 612 . A C 2.32041e-15 .
AB=0;ABP=0;AC=0;AF=0;AN=1;AO=3;CIGAR=1X;DP=48;DPB=48;DPRA=0;EPP=9.52472;EPPR=11.1654;GTI=0;LEN=1;
MEANALT=1;MQM=60;MQMR=60;NS=1;NUMALT=1;ODDS=296.374;PAIRED=1;PAIREDR=0.977778;PAO=0;PQA=0;PQR=0;
PRO=0;QA=53;QR=1495;RO=45;RPL=0;RPP=9.52472;RPPR=3.44459;RPR=3;RUN=1;SAF=3;SAP=9.52472;SAR=0;
SRF=19;SRP=5.37479;SRR=26;TYPE=snp GT:DP:AD:RO:QR:AO:QA:GL 0:48:45,3:45:1495:3:53:0,-129.869
```

The fields in a vcf-file are described in he table (Table 7.1) below:

Table 7.1: The vcf-file format fields.

Col	Field	Description
1	CHROM	Chromosome name
2	POS	1-based position. For an indel, this is the position preceding the indel.
3	ID	Variant identifier. Usually the dbSNP rsID.
4	REF	Reference sequence at POS involved in the variant. For a SNP, it is a single base.
5	ALT	Comma delimited list of alternative seugence(s).
6	QUAL	Phred-scaled probability of all samples being homozygous reference.
7	FILTER	Semicolon delimited list of filters that the variant fails to pass.
8	INFO	Semicolon delimited list of variant information.
9	FORMAT	Colon delimited list of the format of individual genotypes in the following fields.
10+	Sample(s)	Individual genotype information defined by FORMAT.

7.8.2 Statistics

Now we can use it to do some statistics and filter our variant calls.

First, to prepare our vcf-file for querying we need to index it with tabix:

```
# compress file
$ bgzip variants/ev011.freebayes.vcf
# index
$ tabix -p vcf variants/ev011.freebayes.vcf.gz
```

- -p vcf: input format

We can get some quick stats with `rtg vcfstats`:

```
$ rtg vcfstats variants/ev011.freebayes.vcf.gz
```

Example output from `rtg vcfstats`:

```
Location                : variants/ev011.freebayes.vcf.gz
Failed Filters           : 0
Passed Filters           : 35233
SNPs                     : 55
MNPs                     : 6
Insertions               : 3
Deletions                : 5
Indels                   : 0
Same as reference        : 35164
SNP Transitions/Transversions: 0.83 (25/30)
Total Haploid            : 69
Haploid SNPs             : 55
Haploid MNPs             : 6
Haploid Insertions       : 3
Haploid Deletions        : 5
Haploid Indels           : 0
Insertion/Deletion ratio : 0.60 (3/5)
Indel/SNP+MNP ratio      : 0.13 (8/61)
```

However, we can also run `BCFtools`¹⁷⁵ to extract more detailed statistics about our variant calls:

```
$ bcftools stats -F assembly/scaffolds.fasta -s - variants/ev011.freebayes.vcf.gz > variants/ev011.
↪ freebayes.vcf.gz.stats
```

- -s -: list of samples for sample stats, "-" to include all samples
- -F FILE: faidx indexed reference sequence file to determine INDEL context

Now we take the stats and make some plots (e.g. [Fig. 7.2](#)) which are particular of interest if having multiple samples, as one can easily compare them. However, we are only working with one here:

```
$ mkdir variants/plots
$ plot-vcfstats -p variants/plots/ variants/ev011.freebayes.vcf.gz.stats
```

- -p: The output files prefix, add a slash at the end to create a new directory.

¹⁷⁵ <http://www.htslib.org/doc/bcftools.html>

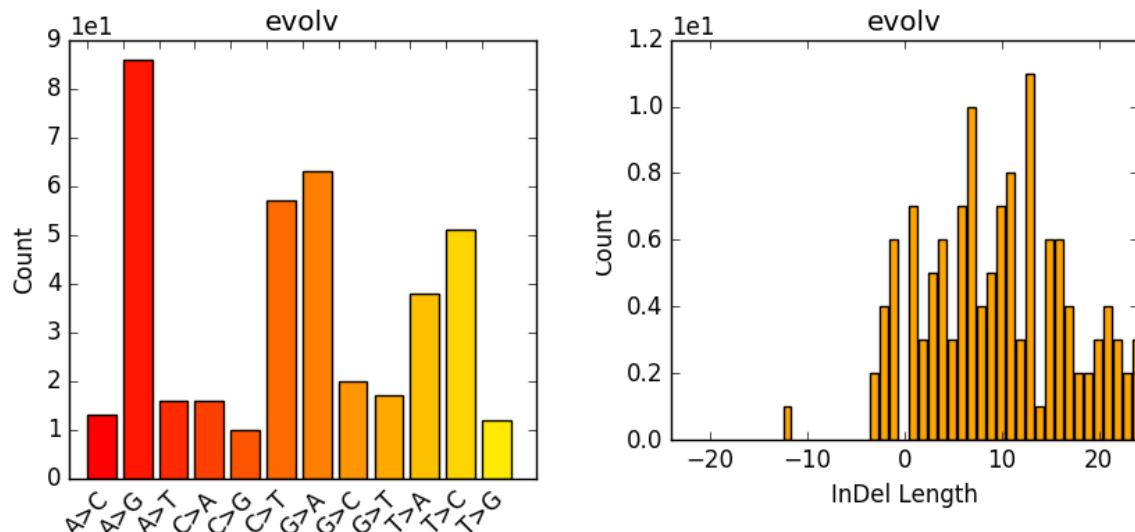


Fig. 7.2: Example of plot-vcfstats output.

7.8.3 Variant filtration

Variant filtration is a big topic in itself [OLSEN2015]. There is no consensus yet and research on how to best filter variants is ongoing.

We will do some simple filtration procedures here. For one, we can filter out low quality reads.

Here, we only include variants that have quality > 30.

```
# use rtg vcffilter
$ rtg vcffilter -q 30 -i variants/evol1.freebayes.vcf.gz -o variants/evol1.freebayes.q30.vcf.gz
```

- -i FILE: input file
- -o FILE: output file
- -q FLOAT: minimal allowed quality in output.

or use `vcflib`¹⁷⁶:

```
# or use vcflib
$ zcat variants/evol1.freebayes.vcf.gz | vcffilter -f "QUAL >= 30" | gzip > variants/evol1.
freebayes.q30.vcf.gz
```

- -f "QUAL >= 30": we only include variants that have been called with quality >= 30.

Quick stats for the filtered variants:

```
# look at stats for filtered
$ rtg vcfstats variants/evol1.freebayes.q30.vcf.gz
```

`freebayes`¹⁷⁷ adds some extra information to the vcf-files it creates. This allows for some more detailed filtering. This strategy will NOT work on calls done with e.g. `SAMtools`¹⁷⁸/`bcftools` mpileup called variants. Here we filter, based on some recommendation from the developer of `freebayes`¹⁷⁹:

```
$ zcat variants/evol1.freebayes.vcf.gz | vcffilter -f "QUAL > 1 & QUAL / AO > 10 & SAF > 0 & SAR > 0 & RPR > 1 & RPL > 1" | bgzip > variants/evol1.freebayes.filtered.vcf.gz
```

¹⁷⁶ <https://github.com/vcflib/vcflib#vcflib>

¹⁷⁷ <https://github.com/ekg/freebayes>

¹⁷⁸ <http://samtools.sourceforge.net/>

¹⁷⁹ <https://github.com/ekg/freebayes>

- $QUAL > 1$: removes really bad sites
- $QUAL / AO > 10$: additional contribution of each obs should be 10 log units ($\sim Q10$ per read)
- $SAF > 0$ & $SAR > 0$: reads on both strands
- $RPR > 1$ & $RPL > 1$: at least two reads “balanced” to each side of the site

```
$ tabix -p vcf variants/evol1.freebayes.filtered.vcf.gz
```

This strategy used here will do for our purposes. However, several more elaborate filtering strategies have been explored, e.g. [here](#)¹⁸⁰.

Todo: Look at the statistics. One ratio that is mentioned in the statistics is transition transversion ratio (ts/tv). Explain what this ratio is and why the observed ratio makes sense.

Todo: Call and filter variants for the second evolved strain, similarly to what we described here for the first strain. Should you be unable to do it, check the code section: *Code: Variant calling* (page 75).

¹⁸⁰ <https://github.com/ekg/freebayes#observation-filters-and-qualities>

GENOME ANNOTATION

8.1 Preface

In this section you will predict genes and assess your assembly using [Augustus](http://augustus.gobics.de)¹⁸³ and [BUSCO](http://busco.ezlab.org)¹⁸⁴, as well as [Prokka](https://github.com/tseemann/prokka)¹⁸⁵.

Note: You will encounter some **To-do** sections at times. Write the solutions and answers into a text-file.

8.2 Overview

The part of the workflow we will work on in this section can be viewed in [Fig. 8.1](#).

8.3 Learning outcomes

After studying this section of the tutorial you should be able to:

1. Explain how annotation completeness is assessed using orthologues
2. Use bioinformatics tools to perform gene prediction
3. Use genome-viewing software to graphically explore genome annotations and NGS data overlays

8.4 Before we start

Lets see how our directory structure looks so far:

```
$ cd ~/analysis
$ ls -1F
```

```
assembly/
data/
kraken/
mappings/
multiqc_data/
trimmed/
trimmed-fastqc/
variants/
```

¹⁸³ <http://augustus.gobics.de>

¹⁸⁴ <http://busco.ezlab.org>

¹⁸⁵ <https://github.com/tseemann/prokka>

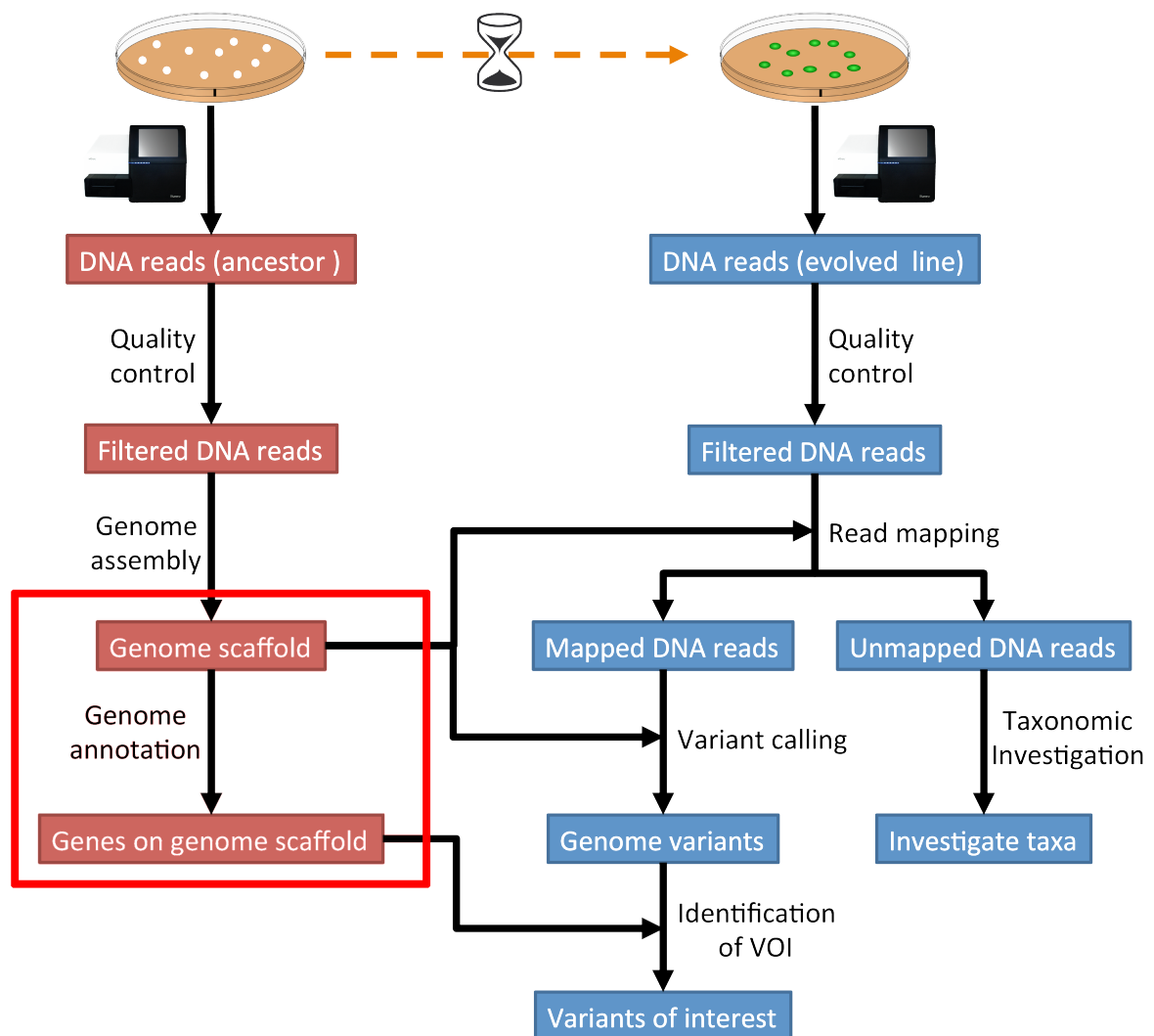


Fig. 8.1: The part of the workflow we will work on in this section marked in red.

8.5 Installing the software

```
# activate the env
$ conda create --yes -n anno busco
```

This will install both the [Augustus](#)¹⁸⁶ [STANKE2005] and the [BUSCO](#)¹⁸⁷ [SIMAO2015] software, which we will use (separately) for gene prediction and assessment of assembly completeness, respectively.

Make a directory for the annotation results:

```
$ mkdir annotation
$ cd annotation
```

We need to get the database that [BUSCO](#)¹⁸⁸ will use to assess orthologue presence absence in our genome annotation. [BUSCO](#)¹⁸⁹ provides a command to list all available datasets and download datasets.

```
$ busco --list-datasets
```

```
INFO:   Downloading information on latest versions of BUSCO data...
```

```
#####
```

```
Datasets available to be used with BUSCOv4 as of 2019/11/27:
```

```
bacteria_odb10
- acidobacteria_odb10
- actinobacteria_phylum_odb10
  - actinobacteria_class_odb10
    - corynebacteriales_odb10
    - micrococcales_odb10
    - propionibacteriales_odb10
    - streptomycetales_odb10
    - streptosporangiales_odb10
  - coriobacteriia_odb10
    - coriobacteriales_odb10
...
```

[BUSCO](#)¹⁹⁰ will download the dataset when starting an analysis.

We also need to place the configuration file for this program in a location in which we have “write” privileges. Do this recursively for the entire config directory, placing it into your current annotation directory:

```
$ cp -r ~/miniconda3/envs/anno/config/ .
```

¹⁸⁶ <http://augustus.gobics.de>

¹⁸⁷ <http://busco.ezlab.org>

¹⁸⁸ <http://busco.ezlab.org>

¹⁸⁹ <http://busco.ezlab.org>

¹⁹⁰ <http://busco.ezlab.org>

8.6 Assessment of orthologue presence and absence

BUSCO¹⁹¹ will assess orthologue presence absence using `blastn`¹⁹², a rapid method of finding close matches in large databases (we will discuss this in lecture). It uses `blastn`¹⁹³ to make sure that it does not miss any part of any possible coding sequences. To run the program, we give it

- A fasta format input file
- A name for the output files
- The name of the lineage database against which we are assessing orthologue presence absence (that we downloaded above)
- An indication of the type of annotation we are doing (genomic, as opposed to transcriptomic or previously annotated protein files).
- The config file to use

```
$ busco -i ../assembly/scaffolds.fasta -o my_anno -l bacteria_odb10 -m geno --config config/config.ini
```

Navigate into the output directory you created. There are many directories and files in there containing information on the orthologues that were found, but here we are only really interested in one: the summary statistics. This is located in the `short_summary*.txt` file. Look at this file. It will note the total number of orthologues found, the number expected, and the number missing. This gives an indication of your genome completeness.

Todo: Is it necessarily true that your assembly is incomplete if it is missing some orthologues? Why or why not?

8.7 Annotation with Augustus¹⁹⁴

We will use `Augustus`¹⁹⁵ to perform gene prediction. This program implements a hidden markov model (HMM) to infer where genes lie in the assembly you have made. To run the program you need to give it:

- Information as to whether you would like the genes called on both strands (or just the forward or reverse strands)
- A “model” organism on which it can base it’s HMM parameters on (in this case we will use E.coli)
- The location of the assembly file
- A name for the output file, which will be a .gff (general feature format) file.
- We will also tell it to display a progress bar as it moves through the genome assembly.

```
$ augustus --progress=true --strand=both --species=E_coli_K12 --AUGUSTUS_CONFIG_PATH=config ../assembly/scaffolds.fasta > augustus.gff
```

Note: Should the process of producing your annotation fail, you can download a annotation manually from [Downloads](#) (page 77). Remember to unzip the file.

¹⁹¹ <http://busco.ezlab.org>

¹⁹² https://blast.ncbi.nlm.nih.gov/Blast.cgi?PAGE_TYPE=BlastSearch

¹⁹³ https://blast.ncbi.nlm.nih.gov/Blast.cgi?PAGE_TYPE=BlastSearch

¹⁹⁴ <http://augustus.gobics.de>

¹⁹⁵ <http://augustus.gobics.de>

8.8 Annotation with Prokka¹⁹⁶

Install Prokka¹⁹⁷:

```
$ conda create --yes -n prokka prokka
$ conda activate prokka
```

Run Prokka¹⁹⁸:

```
$ prokka --kingdom Bacteria --genus Escherichia --species coli --outdir annotation assembly/
↳ scaffolds.fasta
```

Your results will be in the annotation directory with the prefix PROKKA.

8.9 Interactive viewing

We will use the software IGV¹⁹⁹ to view the assembly, the gene predictions you have made, and the variants that you have called, all in one window.

8.9.1 IGV²⁰⁰

```
$ conda activate anno
$ conda install --yes igv
```

To run IGV type:

```
$ igv
```

This will open up a new window. Navigate to that window and open up your genome assembly:

- **Genomes -> Load Genome from File**
- Load your assembly (scaffolds.fasta), not your gff file.

Load the tracks:

- **File -> Load from File**
- Load your unzipped vcf file from last week
- Load your unzipped gff file from this week.

At this point you should be able to zoom in and out to see regions in which there are SNPs or other types of variants. You can also see the predicted genes. If you zoom in far enough, you can see the sequence (DNA and protein).

If you have time and interest, you can right click on the sequence and copy it. Open a new browser window and go to the blastn homepage. There, you can blast your gene of interest (GOI) and see if blast can assign a function to it.

The end goal of this lab will be for you to select a variant that you feel is interesting (e.g. due to the gene it falls near or within), and hypothesize as to why that mutation might have increased in frequency in these evolving populations.

¹⁹⁶ <https://github.com/tseemann/prokka>

¹⁹⁷ <https://github.com/tseemann/prokka>

¹⁹⁸ <https://github.com/tseemann/prokka>

¹⁹⁹ <http://software.broadinstitute.org/software/igv/>

²⁰⁰ <http://software.broadinstitute.org/software/igv/>

ORTHOLOGY AND PHYLOGENY

9.1 Preface

In this section you will use some software to find orthologue genes and do phylogenetic reconstructions.

9.2 Learning outcomes

After studying this tutorial you should be able to:

1. Use bioinformatics software to find orthologues in the NCBI database.
2. Use bioinformatics software to perform sequence alignment.
3. Use bioinformatics software to perform phylogenetic reconstructions.

9.3 Before we start

Lets see how our directory structure looks so far:

```
$ cd ~/analysis
$ ls -1F
```

```
annotation/
assembly/
data/
kraken/
mappings/
multiqc_data/
trimmed/
trimmed-fastqc/
variants/
```

Make a directory for the phylogeny results (in your analysis directory):

```
$ mkdir phylogeny
```

Download the fasta file of the *S. cerevisiae* TEF2 gene to the phylogeny folder:

```
$ cd phylogeny
$ curl -O http://compbio.massey.ac.nz/data/203341/s_cerev_tef2.fas
```

Note: Should the download fail, download manually from [Downloads](#) (page 77).

9.4 Installing the software

```
# activate the env
conda activate ngs

conda install blast
```

This will install a [BLAST²⁰³](#) executable that you can use to remotely query the NCBI database.

```
conda install muscle
```

This will install [MUSCLE²⁰⁴](#), alignment program that you can use to align nucleotide or protein sequences.

We will also install [RAxML-NG²⁰⁵](#), a phylogenetic tree inference tool, which uses maximum-likelihood (ML) optimality criterion. However, there is no conda repository for it yet. Thus, we need to download it manually.

```
wget
https://github.com/amkozlov/raxml-ng/releases/download/0.5.1/raxml-ng_v0.5.1b_linux_x86_64.zip

unzip raxml-ng_v0.5.1b_linux_x86_64.zip

rm raxml-ng_v0.5.1b_linux_x86_64.zip
```

9.5 Finding orthologues using BLAST

We will first make a [BLAST²⁰⁶](#) database of our current assembly so that we can find the orthologous sequence of the *S. cerevisiae* gene. To do this, we run the command `makeblastdb`:

```
# create blast db
makeblastdb in ../assembly/spades_final/scaffolds.fasta dbtype nucl
```

To run [BLAST²⁰⁷](#), we give it:

- `-db`: The name of the database that we are BLASTing
- `-query`: A fasta format input file
- A name for the output files
- Some notes about the format we want

First, we blast without any formatting:

```
blastn db ../assembly/spades_final/scaffolds.fasta query s_cerev_tef2.fas > blast.out
```

This should output a file with a set of [BLAST²⁰⁸](#) hits similar to what you might see on the [BLAST²⁰⁹](#) web site.

Read through the output (e.g. using `nano`) to see what the results of your [BLAST²¹⁰](#) run was.

Next we will format the output a little so that it is easier to deal with.

²⁰³ <https://blast.ncbi.nlm.nih.gov/Blast.cgi>

²⁰⁴ <http://www.ebi.ac.uk/Tools/msa/muscle/>

²⁰⁵ <https://github.com/amkozlov/raxml-ng>

²⁰⁶ <https://blast.ncbi.nlm.nih.gov/Blast.cgi>

²⁰⁷ <https://blast.ncbi.nlm.nih.gov/Blast.cgi>

²⁰⁸ <https://blast.ncbi.nlm.nih.gov/Blast.cgi>

²⁰⁹ <https://blast.ncbi.nlm.nih.gov/Blast.cgi>

²¹⁰ <https://blast.ncbi.nlm.nih.gov/Blast.cgi>

```
blastn db ../assembly/spades_final/scaffolds.fasta query s_cerev_tef2.fas evaluate 1e-100 outfmt "6_
↩length sseq" > blast_formatted.out
```

This will yield a file that has only the sequences of the subject, so that we can later add those to other fasta files. However, the formatting is not perfect. To adjust the format such that it is fasta format, open the file in an editor (e.g. nano) and edit the first line so that it has a name for your sequence. You should know the general format of a fasta-file (e.g. the first line start with a ">").

Hint: To edit in vi editor, you will need to press the escape key and "a" or "e". To save in vi, you will need to press the escape key and "w" (write). To quit vi, you will need to press the escape key and "q" (quit).

Next, you have to replace the dashes (signifying indels in the [BLAST²¹¹](#) result). This can easily be done in vi: Press the escape key, followed by: `:%s/\-//g`

Now we will [BLAST²¹²](#) a remote database to get a list of hits that are already in the NCBI database.

Note: It turns out you may not be able to access this database from within BioLinux. In such a case, download the file named `blast.fas` and place it into your `~/analysis/phylogeny/` directory.

```
curl -O http://compbio.massey.ac.nz/data/203341/blast_u.fas
```

Append the fasta file of your yeast sequence to this file, using whatever set of commands you wish/know.

Note: Should the download fail, download manually from [Downloads](#) (page 77).

9.6 Performing an alignment

We will use [MUSCLE²¹³](#) to perform our alignment on all the sequences in the [BLAST²¹⁴](#) fasta file. This syntax is very simple (change the filenames accordingly):

```
muscle in infile.fas out your_alignment.aln
```

9.7 Building a phylogeny

We will use [RAxML-NG²¹⁵](#) to build our phylogeny. This uses a maximum likelihood method to infer parameters of evolution and the topology of the tree. Again, the syntax of the command is fairly simple, except you must make sure that you are using the directory in which [RAxML-NG²¹⁶](#) sits.

The arguments are:

- `-s`: an alignment file
- `-m`: a model of evolution. In this case we will use a general time reversible model with gamma distributed rates (GTR+GAMMA)
- `-n`: outfile-name

²¹¹ <https://blast.ncbi.nlm.nih.gov/Blast.cgi>

²¹² <https://blast.ncbi.nlm.nih.gov/Blast.cgi>

²¹³ <http://www.ebi.ac.uk/Tools/msa/muscle/>

²¹⁴ <https://blast.ncbi.nlm.nih.gov/Blast.cgi>

²¹⁵ <https://github.com/amkozlov/raxml-ng>

²¹⁶ <https://github.com/amkozlov/raxml-ng>

- -p: specify a random number seed for the parsimony inferences

```
raxmlHPC -s your_alignment.aln -m GTRGAMMA n yeast_tree p 12345
```

9.8 Visualizing the phylogeny

We will use the online software [Interactive Tree of Life \(iTOL\)](#)²¹⁷ to visualize the tree. Navigate to this homepage. Open the file containing your tree (*bestTree.out), copy the contents, and paste into the web page (in the Tree text box).

You should then be able to zoom in and out to see where your yeast taxa is. To find out the closest relative, you will have to use the [NCBI taxa page](#)²¹⁸.

Todo: Are you certain that the yeast are related in the way that the phylogeny suggests? Why might the topology of this phylogeny not truly reflect the evolutionary history of these yeast species?

²¹⁷ <http://itol.embl.de/upload.cgi>

²¹⁸ https://www.ncbi.nlm.nih.gov/Taxonomy/TaxIdentifier/tax_identifier.cgi

VARIANTS-OF-INTEREST

10.1 Preface

In this section we will use our genome annotation of our reference and our genome variants in the evolved line to find variants that are interesting in terms of the observed biology.

Note: You will encounter some **To-do** sections at times. Write the solutions and answers into a text-file.

10.2 Overview

The part of the workflow we will work on in this section can be viewed in [Fig. 10.1](#).

10.3 Learning outcomes

After studying this section of the tutorial you should be able to:

1. Identify variants of interests.
2. Understand how the variants might affect the observed biology in the evolved line.

10.4 Before we start

Lets see how our directory structure looks so far:

```
cd ~/analysis  
ls -lF
```

```
annotation/  
assembly/  
data/  
kraken/  
mappings/  
phylogeny/  
SolexaQA/  
SolexaQA++  
trimmed/  
trimmed-fastqc/  
trimmed-solexaqa/  
variants/
```

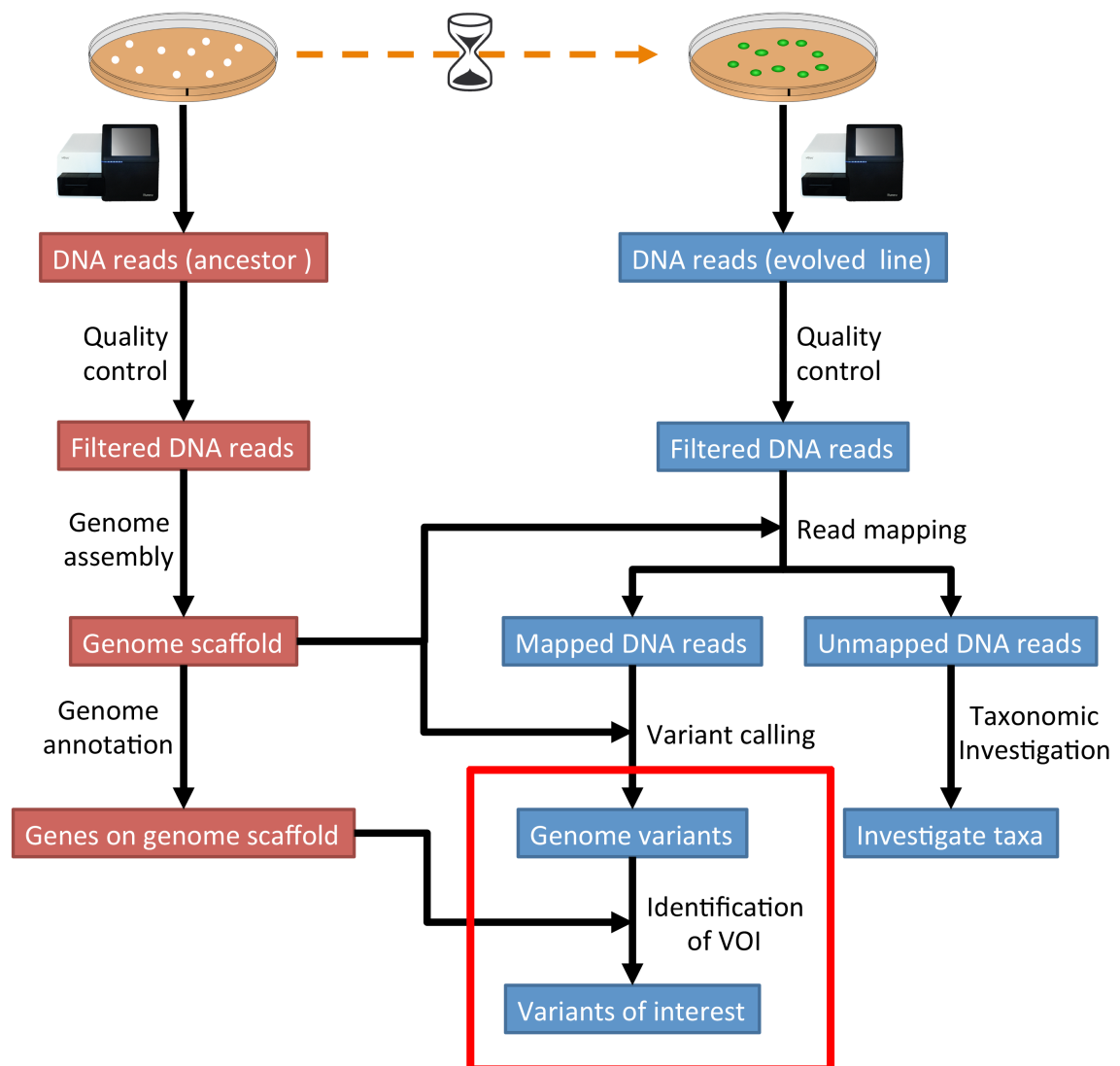


Fig. 10.1: The part of the workflow we will work on in this section marked in red.

10.5 General comments for identifying variants-of-interest

Things to consider when looking for variants-of-interest:

- The quality score of the variant call.
 - Do we call the variant with a higher than normal score?
- The mapping quality score.
 - How confident are we that the reads were mapped at the position correctly?
- The location of the SNP.
 - SNPs in larger contigs are probably more interesting than in tiny contigs.
 - Does the SNP overlap a coding region in the genome annotation?
- The type of SNP.
 - substitutions vs. indels

10.6 SnpEff

We will be using [SnpEff²¹⁹](#) to annotate our identified variants. The tool will tell us on to which genes we should focus further analyses.

10.6.1 Installing software

Tools we are going to use in this section and how to install them if you not have done it yet.

```
# activate the env
conda activate ngs

# Install these tools into the conda environment
# if not already installed
conda install snpeff
conda install genomtools-genomtools
```

Make a directory for the results (in your analysis directory) and change into the directory:

```
mkdir voi

# change into the directory
cd voi
```

10.6.2 Prepare SnpEff database

We need to create our own config-file for [SnpEff²²⁰](#). Where is the snpEff.config:

```
find ~ -name snpEff.config
/home/manager/miniconda3/envs/ngs/share/snpeff-4.3.1m-0/snpEff.config
```

This will give you the path to the snpEff.config. It might be looking a bit different then the one shown here, depending on the version of [SnpEff²²¹](#) that is installed.

Make a local copy of the snpEff.config and then edit it with an editor of your choice:

²¹⁹ <http://snpeff.sourceforge.net/index.html>

²²⁰ <http://snpeff.sourceforge.net/index.html>

²²¹ <http://snpeff.sourceforge.net/index.html>

```
cp /home/manager/miniconda3/envs/ngs/share/snpEff-4.3.1m-0/snpEff.config .
nano snpEff.config
```

Make sure the data directory path in the snpEff.config looks like this:

```
data.dir = ./data/
```

There is a section with databases, which starts like this:

```
#-----
# Databases & Genomes
#
# One entry per genome version.
#
# For genome version 'ZZZ' the entries look like
#   ZZZ.genome           : Real name for ZZZ (e.g. 'Human')
#   ZZZ.reference        : [Optional] Comma separated list of URL to site/s Where information
#   ↪ for building ZZZ database was extracted.
#   ZZZ.chrName.codonTable : [Optional] Define codon table used for chromosome 'chrName' (Default:
#   ↪ 'codon.Standard')
#
#-----
```

Add the following two lines in the database section underneath these header lines:

```
# my yeast genome
yeastanc.genome : WildYeastAnc
```

Now, we need to create a local data folder called ./data/yeastanc.

```
# create folders
mkdir -p ./data/yeastanc
```

Copy our genome assembly to the newly created data folder. The name needs to be sequences.fa or yeastanc.fa:

```
cp ../assembly/spades_final/scaffolds.fasta ./data/yeastanc/sequences.fa
gzip ./data/yeastanc/sequences.fa
```

Copy our genome annotation to the data folder. The name needs to be genes.gff (or genes.gtf for gtf-files).

```
cp ../annotation/your_new_fungus.gff ./data/yeastanc/genes.gff
gzip ./data/yeastanc/genes.gff
```

Now we can build a new SnpEff²²² database:

```
snpEff build -c snpEff.config -gff3 -v yeastanc > snpEff.stdout 2> snpEff.stderr
```

Note: Should this fail, due to gff-format of the annotation, we can try to convert the gff to gtf:

```
# using genomertools
gt gff3_to_gtf ../annotation/your_new_fungus.gff -o ./data/yeastanc/genes.gtf
gzip ./data/yeastanc/genes.gtf
```

Now, we can use the gtf annotation to build the database:

²²² <http://snpeff.sourceforge.net/index.html>

```
snpeff build -c snpeff.config -gtf22 -v yeastanc > snpeff.stdout 2> snpeff.stderr
```

10.6.3 SNP annotation

Now we can use our new [Snpeff²²³](#) database to annotate some variants, e.g.:

```
snpeff -c snpeff.config yeastanc ../variants/evolved-6.freebayes.filtered.vcf.gz > evolved-6.freebayes.filtered.anno.vcf
```

[Snpeff²²⁴](#) adds ANN fields to the vcf-file entries that explain the effect of the variant.

Note: If you are unable to do the annotation, you can download an annotated vcf-file from [Downloads](#) (page 77).

10.6.4 Example

Lets look at one entry from the original vcf-file and the annotated one. We are only interested in the 8th column, which contains information regarding the variant. [Snpeff²²⁵](#) will add fields here :

```
# evolved-6.freebayes.filtered.vcf (the original), column 8
AB=0.5;ABP=3.0103;AC=1;AF=0.5;AN=2;AO=56;CIGAR=1X;DP=112;DPB=112;DPRA=0;EPP=3.16541;EPPR=3.16541;
↳GTI=0;LEN=1;MEANALT=1;MQM=42;MQMR=42;NS=1;NUMALT=1;ODDS=331.872;PAIRED=1;PAIREDR=1;PAO=0;PQA=0;
↳PQR=0;PRO=0;QA=2128;QR=2154;RO=56;RPL=35;RPP=10.6105;RPPR=3.63072;RPR=21;RUN=1;SAF=30;SAP=3.63072;
↳SAR=26;SRF=31;SRP=4.40625;SRR=25;TYPE=snp

# evolved-6.freebayes.filtered.anno.vcf, column 8
AB=0.5;ABP=3.0103;AC=1;AF=0.5;AN=2;AO=56;CIGAR=1X;DP=112;DPB=112;DPRA=0;EPP=3.16541;EPPR=3.16541;
↳GTI=0;LEN=1;MEANALT=1;MQM=42;MQMR=42;NS=1;NUMALT=1;ODDS=331.872;PAIRED=1;PAIREDR=1;PAO=0;PQA=0;
↳PQR=0;PRO=0;QA=2128;QR=2154;RO=56;RPL=35;RPP=10.6105;RPPR=3.63072;RPR=21;RUN=1;SAF=30;SAP=3.63072;
↳SAR=26;SRF=31;SRP=4.40625;SRR=25;TYPE=snp;ANN=T|missense_variant|MODERATE|CDS_NODE_40_length_1292_
↳cov_29.5267_1_1292|GENE_CDS_NODE_40_length_1292_cov_29.5267_1_1292|transcript|TRANSCRIPT_CDS_NODE_
↳40_length_1292_cov_29.5267_1_1292|protein_coding|1/1|c.664T>A|p.Ser222Thr|664/1292|664/1292|222/
↳429||WARNING_TRANSCRIPT_INCOMPLETE,T|intragenic_variant|MODIFIER|GENE_NODE_40_length_1292_cov_29.
↳5267_1_1292|GENE_NODE_40_length_1292_cov_29.5267_1_1292|gene_variant|GENE_NODE_40_length_1292_cov_
↳29.5267_1_1292|||n.629A>T|||||
```

When expecting the second entry, we find that [Snpeff²²⁶](#) added annotation information starting with ANN=T|missense_variant|. . . . If we look a bit more closely we find that the variant results in a amino acid change from a threonine to a serine (c.664T>A|p.Ser222Thr). The codon for serine is TCN and for threonine is ACN, so the variant in the first nucleotide of the codon made the amino acid change.

A quick protein [BLAST²²⁷](#) of the CDS sequence where the variant was found (extracted from the genes.gff.gz) shows that the closest hit is a translation elongation factor from a species called [Candida dubliniensis²²⁸](#) another fungi.

²²³ <http://snpeff.sourceforge.net/index.html>

²²⁴ <http://snpeff.sourceforge.net/index.html>

²²⁵ <http://snpeff.sourceforge.net/index.html>

²²⁶ <http://snpeff.sourceforge.net/index.html>

²²⁷ <https://blast.ncbi.nlm.nih.gov/Blast.cgi>

²²⁸ https://en.wikipedia.org/wiki/Candida_dubliniensis

²²⁹ <https://blast.ncbi.nlm.nih.gov/Blast.cgi>

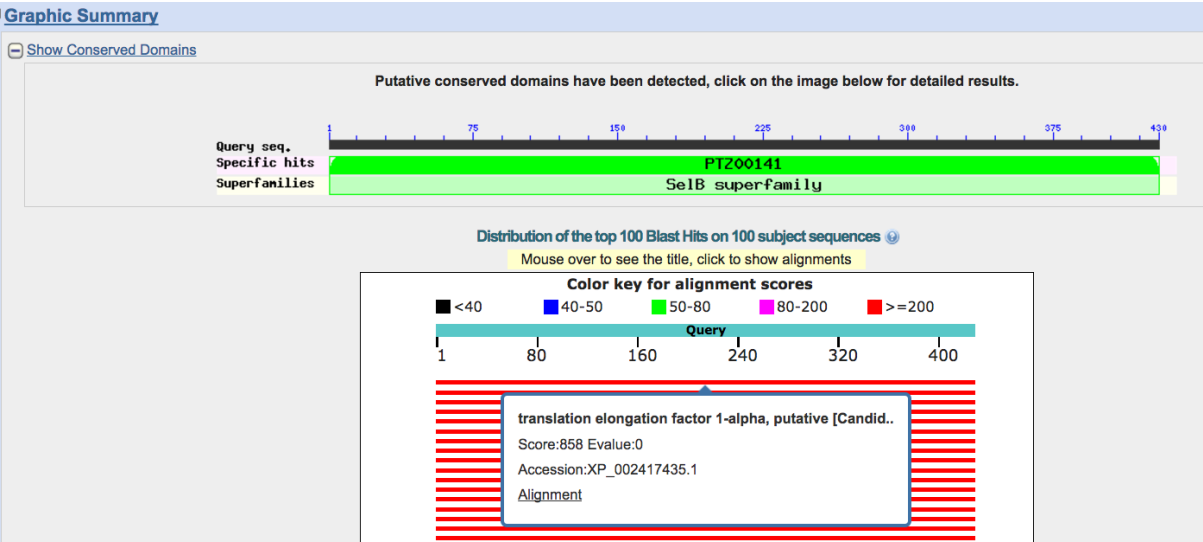


Fig. 10.2: Results of a BLAST²²⁹ search of the CDS.

QUICK COMMAND REFERENCE

11.1 Shell commands

```
# Where in the directory tree am I?
pwd

# List the documents and sub-directories in the current directory
ls

# a bit nicer listing with more information
ls -laF

# Change into your home directory
cd ~

# Change back into the last directory
cd -

# Change one directory up in the tree
cd ..

# Change explicitly into a directory "temp"
cd temp

# Quickly show content of a file "temp.txt"
# exist the view with "q", navigate line up and down with "k" and "j"
less temp.txt

# Show the beginning of a file "temp.txt"
head temp.txt

# Show the end of a file "temp.txt"
tail temp.txt
```

11.2 General conda commands

```
# To update all packages
conda update --all --yes

# List all packages installed
conda list [-n env]

# conda list environments
conda env list
```

(continues on next page)

(continued from previous page)

```
# create new env
conda create -n [name] package [package] ...

# activate env
conda activate [name]

# deactivate env
conda deactivate
```

CODING SOLUTIONS

12.1 QC

12.1.1 Code: fastp

```
# run fastp like this on the ancestor:
fastp --detect_adapter_for_pe --overrepresentation_analysis --correction --cut_right --html_
↳ trimmed/anc.fastp.html --json trimmed/anc.fastp.json --thread 2 -i data/anc_R1.fastq.gz -I data/
↳ anc_R2.fastq.gz -o trimmed/anc_R1.fastq.gz -O trimmed/anc_R2.fastq.gz

# run the evolved samples through fastp
fastp --detect_adapter_for_pe --overrepresentation_analysis --correction --cut_right --html trimmed/
↳ evol1.fastp.html --json trimmed/evol1.fastp.json --thread 2 -i data/evol1_R1.fastq.gz -I data/
↳ evol1_R2.fastq.gz -o trimmed/evol1_R1.fastq.gz -O trimmed/evol1_R2.fastq.gz

fastp --detect_adapter_for_pe --overrepresentation_analysis --correction --cut_right --html trimmed/
↳ evol2.fastp.html --json trimmed/evol2.fastp.json --thread 2 -i data/evol2_R1.fastq.gz -I data/
↳ evol2_R2.fastq.gz -o trimmed/evol2_R1.fastq.gz -O trimmed/evol2_R2.fastq.gz
```

12.1.2 Code: FastQC

Create directory:

```
mkdir trimmed-fastqc
```

Run FastQC:

```
fastqc -o trimmed-fastqc trimmed/*.fastq.gz
```

Run MultiQC

```
multiqc trimmed-fastqc trimmed
```

Open |multiqc| report html webpage:

```
firefox multiqc_report.html
```

12.2 Assembly

12.2.1 Code: SPAdes assembly (trimmed data)

```
spades.py -o assembly/spades-150/ --careful -1 trimmed/anc_R1.fastq.gz -2 trimmed/anc_R2.fastq.gz
```

12.2.2 Code: SPAdes assembly (original data)

```
spades.py -o assembly/spades-original/ --careful -1 data/anc_R1.fastq.gz -2 data/anc_R2.fastq.gz
```

12.2.3 Code: Quast

```
quast -o assembl/quast assembly/spades-150/scaffolds.fasta assembly/spades-original/scaffolds.fasta
```

12.3 Mapping

12.3.1 Code: BWA indexing

Index the genome assembly:

```
bwa index assembly/scaffolds.fasta
```

12.3.2 Code: BWA mapping

Run bwa mem:

```
# trimmed data
bwa mem assembly/scaffolds.fasta trimmed/evol1_R1.fastq.gz trimmed/evol1_R2.fastq.gz > mappings/
↪evol1.sam
bwa mem assembly/scaffolds.fasta trimmed/evol2_R1.fastq.gz trimmed/evol2_R2.fastq.gz > mappings/
↪evol2.sam
```

12.3.3 Code: Mapping post-processing

```
#
# Evol 1
#

# fixmate and compress to bam
samtools sort -n -O sam mappings/evol1.sam | samtools fixmate -m -O bam - mappings/evol1.fixmate.bam
rm mappings/evol1.sam
# sort
samtools sort -O bam -o mappings/evol1.sorted.bam mappings/evol1.fixmate.bam
rm mappings/evol1.fixmate.bam
# mark duplicates
samtools markdup -r -S mappings/evol1.sorted.bam mappings/evol1.sorted.dedup.bam
rm mappings/evol1.sorted.bam
# extract q20 mappers
samtools view -h -b -q 20 mappings/evol1.sorted.dedup.bam > mappings/evol1.sorted.dedup.q20.bam
# extract unmapped
```

(continues on next page)

(continued from previous page)

```

samtools view -b -f 4 mappings/evol1.sorted.dedup.bam > mappings/evol1.sorted.unmapped.bam
rm mappings/evol1.sorted.dedup.bam
# covert to fastq
samtools fastq -1 mappings/evol1.sorted.unmapped.R1.fastq.gz -2 mappings/evol1.sorted.unmapped.R2.
↪ fastq.gz mappings/evol1.sorted.unmapped.bam
# delete not needed files
rm mappings/evol1.sorted.unmapped.bam

#
# Evol 2
#

samtools sort -n -O sam mappings/evol2.sam | samtools fixmate -m -O bam - mappings/evol2.fixmate.bam
rm mappings/evol2.sam
samtools sort -O bam -o mappings/evol2.sorted.bam mappings/evol2.fixmate.bam
rm mappings/evol2.fixmate.bam
samtools markdup -r -S mappings/evol2.sorted.bam mappings/evol2.sorted.dedup.bam
rm mappings/evol2.sorted.bam
samtools view -h -b -q 20 mappings/evol2.sorted.dedup.bam > mappings/evol2.sorted.dedup.q20.bam
rm mappings/evol2.sorted.dedup.bam

```

12.3.4 Code: Variant calling

```

# index genome
samtools faidx assembly/scaffolds.fasta
mkdir variants

#
# Evol 1
#

# index mappings
bamtools index -in mappings/evol1.sorted.dedup.q20.bam

# calling variants
freebayes -p 1 -f assembly/scaffolds.fasta mappings/evol1.sorted.dedup.q20.bam > variants/evol1.
↪ freebayes.vcf
# compress
bgzip variants/evol1.freebayes.vcf
# index
$ tabix -p vcf variants/evol1.freebayes.vcf.gz

# filtering
zcat variants/evol1.freebayes.vcf.gz | vcfilter -f "QUAL > 1 & QUAL / AO > 10 & SAF > 0 & SAR > 0 &
↪ RPR > 1 & RPL > 1" | bgzip > variants/evol1.freebayes.filtered.vcf.gz
tabix -p vcf variants/evol1.freebayes.filtered.vcf.gz

#
# Evol 2
#

# index mappings
bamtools index -in mappings/evol2.sorted.dedup.q20.bam

# calling variants
freebayes -p 1 -f assembly/scaffolds.fasta mappings/evol2.sorted.dedup.q20.bam > variants/evol2.
↪ freebayes.vcf
# compress
bgzip variants/evol2.freebayes.vcf

```

(continues on next page)

(continued from previous page)

```
# index
$ tabix -p vcf variants/evol2.freebayes.vcf.gz

# filtering
zcat variants/evol2.freebayes.vcf.gz | vcffilter -f "QUAL > 1 & QUAL / AO > 10 & SAF > 0 & SAR > 0 &
↪ RPR > 1 & RPL > 1" | bgzip > variants/evol2.freebayes.filtered.vcf.gz
tabix -p vcf variants/evol2.freebayes.filtered.vcf.gz
```

DOWNLOADS

13.1 Tools

- Miniconda installer [[EXTERNAL²³⁰](#)]
- Minikraken database [[EXTERNAL²³¹](#)]
- Centrifuge database [[EXTERNAL²³²](#)]
- Krona taxonomy database [[DROPBOX²³³](#)]
- RAxML-NG [[EXTERNAL²³⁴](#) | [DROPBOX²³⁵](#)]

13.2 Data

- *Quality control* (page 9): Raw data-set [[DROPBOX²³⁶](#)]
- *Quality control* (page 9): Trimmed data-set [[DROPBOX²³⁷](#)]
- *Genome assembly* (page 19): Assembled data-set [[DROPBOX²³⁸](#)]
- *Read mapping* (page 25): Mapping index (bwa) [[DROPBOX²³⁹](#)]
- *Read mapping* (page 25): Mapped data [[DROPBOX²⁴⁰](#)]
- *Variant calling* (page 47): Called/filtered variants [[DROPBOX²⁴¹](#)]
- *Genome annotation* (page 55): GFF annotation file [[DROPBOX²⁴²](#)]
- *Orthology and Phylogeny* (page 61): *S. cerevisiae* TEF2 gene file [[`DROPBOX <>`](#)]
- *Orthology and Phylogeny* (page 61): BLAST file [[`DROPBOX <>`](#)]
- *Variants-of-interest* (page 65): SnpEff annotated vcf-file [[`DROPBOX <>`](#)]

²³⁰ https://repo.continuum.io/miniconda/Miniconda3-latest-Linux-x86_64.sh

²³¹ ftp://ftp.ccb.jhu.edu/pub/data/kraken2_dbs/minikraken2_v2_8GB_201904_UPDATE.tgz

²³² ftp://ftp.ccb.jhu.edu/pub/infphilo/centrifuge/data/p_compressed+h+v.tar.gz

²³³ <https://www.dropbox.com/s/cwf1qc5zyq65yvn/taxonomy.tab.gz?dl=0>

²³⁴ https://github.com/amkozlov/raxml-ng/releases/download/0.3.0/raxml-ng_v0.3.0b_linux_x86_64.zip

²³⁵ https://www.dropbox.com/s/iliws53ri5z4y69/raxml-ng_v0.3.0b_linux_x86_64.zip?dl=0

²³⁶ <https://www.dropbox.com/s/3vu1mct230ewhw1/data.tar.gz?dl=0>

²³⁷ <https://www.dropbox.com/s/y3xsggn0glb6ter/trimmed.tar.gz?dl=0>

²³⁸ <https://www.dropbox.com/s/h906x9maw879t5s/assembly.tar.gz?dl=0>

²³⁹ https://www.dropbox.com/s/ii3vbdj9yn916k4/mapping_idx.tar.gz?dl=0

²⁴⁰ <https://www.dropbox.com/s/8bporden0o230oo/mappings.tar.gz?dl=0>

²⁴¹ <https://www.dropbox.com/s/lraipofsvkl1md/variants.tar.gz?dl=0>

²⁴² <https://www.dropbox.com/s/16p9tb22lsvqxbg/annotation.tar.gz?dl=0>

LIST OF FIGURES

1.1	The tutorial will follow this workflow.	4
3.1	The part of the workflow we will work on in this section marked in red.	10
3.2	Illustration of single-end (SE) versus paired-end (PE) sequencing.	11
3.3	Quality score across bases.	16
3.4	Quality per tile.	17
3.5	GC distribution over all sequences.	18
4.1	The part of the workflow we will work on in this section marked in red.	20
5.1	The part of the workflow we will work on in this section marked in red.	26
5.2	A example coverage plot for a contig with highlighted in red regions with a coverage below 20 reads.	32
6.1	The part of the workflow we will work on in this section marked in red.	36
6.2	Example of an Krona output webpage.	44
7.1	The part of the workflow we will work on in this section marked in red.	48
7.2	Example of plot-vcfstats output.	52
8.1	The part of the workflow we will work on in this section marked in red.	56
10.1	The part of the workflow we will work on in this section marked in red.	66
10.2	Results of a BLAST search of the CDS.	70

LIST OF TABLES

5.1	The sam-file format fields.	29
7.1	The vcf-file format fields.	50

BIBLIOGRAPHY

- [KAWECKI2012] Kawecki TJ et al. Experimental evolution. *Trends in Ecology and Evolution* (2012) 27:10⁵
- [ZEYL2006] Zeyl C. Experimental evolution with yeast. *FEMS Yeast Res*, 2006, 685–691⁶
- [GLENN2011] Glenn T. Field guide to next-generation DNA sequencers. *Molecular Ecology Resources* (2011) 11, 759–769 doi: 10.1111/j.1755-0998.2011.03024.x⁴¹
- [KIRCHNER2014] Kirchner et al. Addressing challenges in the production and analysis of Illumina sequencing data. *BMC Genomics* (2011) 12:382⁴²
- [MUKHERJEE2015] Mukherjee S, Huntemann M, Ivanova N, Kyrpides NC and Pati A. Large-scale contamination of microbial isolate genomes by Illumina PhiX control. *Standards in Genomic Sciences*, 2015, 10:18. DOI: 10.1186/1944-3277-10-18⁴³
- [ROBASKY2014] Robasky et al. The role of replicates for error mitigation in next-generation sequencing. *Nature Reviews Genetics* (2014) 15, 56-62⁴⁴
- [ABBAS2014] Abbas MM, Malluhi QM, Balakrishnan P. Assessment of de novo assemblers for draft genomes: a case study with fungal genomes. *BMC Genomics*. 2014;15 Suppl 9:S10.⁶⁴ doi: 10.1186/1471-2164-15-S9-S10. Epub 2014 Dec 8.
- [COMPEAU2011] Compeau PE, Pevzner PA, Tesler G. How to apply de Bruijn graphs to genome assembly. *Nat Biotechnol*. 2011 Nov 8;29(11):987-91⁶⁵
- [GUREVICH2013] Gurevich A, Saveliev V, Vyahhi N and Tesler G. QUAST: quality assessment tool for genome assemblies. *Bioinformatics* 2013, 29(8), 1072-1075⁶⁶
- [NAGARAJAN2013] Nagarajan N, Pop M. Sequence assembly demystified. *Nat Rev Genet*. 2013 Mar;14(3):157-67⁶⁷
- [SALZBERG2012] Salzberg SL, Phillippy AM, Zimin A, Puiu D, Magoc T, Koren S, Treangen TJ, Schatz MC, Delcher AL, Roberts M, Marçais G, Pop M, Yorke JA. GAGE: A critical evaluation of genome assemblies and assembly algorithms. *Genome Res*. 2012 Mar;22(3):557-67⁶⁸
- [WICK2015] Wick RR, Schultz MB, Zobel J and Holt KE. Bandage: interactive visualization of de novo genome assemblies. *Bioinformatics* 2015, 10.1093/bioinformatics/btv383⁶⁹

⁵ <http://dx.doi.org/10.1016/j.tree.2012.06.001>

⁶ <http://doi.org/10.1111/j.1567-1364.2006.00061.x>

⁴¹ <http://doi.org/10.1111/j.1755-0998.2011.03024.x>

⁴² <http://doi.org/10.1186/1471-2164-12-382>

⁴³ <http://doi.org/10.1186/1944-3277-10-18>

⁴⁴ <http://doi.org/10.1038/nrg3655>

⁶⁴ <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4290589/>

⁶⁵ <http://dx.doi.org/10.1038/nbt.2023>

⁶⁶ <http://bioinformatics.oxfordjournals.org/content/29/8/1072>

⁶⁷ <http://dx.doi.org/10.1038/nrg3367>

⁶⁸ <http://genome.cshlp.org/content/22/3/557.full?sid=59ea80f7-b408-4a38-9888-3737bc670876>

⁶⁹ <http://bioinformatics.oxfordjournals.org/content/early/2015/07/11/bioinformatics.btv383.long>

- [LI2009] Li H, Durbin R. (2009). Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*. 25 (14): 1754–1760.¹⁰⁴
- [OKO2015] Okonechnikov K, Conesa A, García-Alcalde F. Qualimap 2: advanced multi-sample quality control for high-throughput sequencing data. *Bioinformatics* (2015), 32, 2:292–294.¹⁰⁵
- [KIM2017] Kim D, Song L, Breitwieser FP, Salzberg SL. Centrifuge: rapid and sensitive classification of metagenomic sequences. *Genome Res*. 2016 Dec;26(12):1721-1729¹⁶⁸
- [LU2017] Lu J, Breitwieser FP, Thielen P, Salzberg SL. Bracken: estimating species abundance in metagenomics data. *PeerJ Computer Science*, 2017, 3:e104, doi:10.7717/peerj-cs.104¹⁶⁹
- [ONDOV2011] Ondov BD, Bergman NH, and Phillippy AM. Interactive metagenomic visualization in a Web browser. *BMC Bioinformatics*, 2011, 12(1):385.¹⁷⁰
- [WOOD2014] Wood DE and Steven L Salzberg SL. Kraken: ultrafast metagenomic sequence classification using exact alignments. *Genome Biology*, 2014, 15:R46. DOI: 10.1186/gb-2014-15-3-r46¹⁷¹.
- [NIELSEN2011] Nielsen R, Paul JS, Albrechtsen A, Song YS. Genotype and SNP calling from next-generation sequencing data. *Nat Rev Genetics*, 2011, 12:433-451¹⁸¹
- [OLSEN2015] Olsen ND et al. Best practices for evaluating single nucleotide variant calling methods for microbial genomics. *Front. Genet.*, 2015, 6:235.¹⁸²
- [SIMAO2015] Simao FA, Waterhouse RM, Ioannidis P, Kriventseva EV and Zdobnov EM. BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*, 2015, Oct 1;31(19):3210-2²⁰¹
- [STANKE2005] Stanke M and Morgenstern B. AUGUSTUS: a web server for gene prediction in eukaryotes that allows user-defined constraints. *Nucleic Acids Res*, 2005, 33(Web Server issue): W465–W467.²⁰²

¹⁰⁴ <https://doi.org/10.1093%2Fbioinformatics%2Fbtp324>

¹⁰⁵ <https://doi.org/10.1093/bioinformatics/btv566>

¹⁶⁸ <https://www.ncbi.nlm.nih.gov/pubmed/27852649>

¹⁶⁹ <https://peerj.com/articles/cs-104/>

¹⁷⁰ <http://www.ncbi.nlm.nih.gov/pubmed/21961884>

¹⁷¹ <http://doi.org/10.1186/gb-2014-15-3-r46>

¹⁸¹ <http://doi.org/10.1038/nrg2986>

¹⁸² <https://doi.org/10.3389/fgene.2015.00235>

²⁰¹ <http://doi.org/10.1093/bioinformatics/btv351>

²⁰² <https://dx.doi.org/10.1093/nar/gki458>