# Bioinformatics Training using *de novo* assembly as an exemplar

The bioinformatics training will comprise of the 3 phases.
1. Training on Galaxy to use command line software via a web interface. The objective is to
   a. become familiar with the bioinformatics steps for a particular process
   b. understand what software is used and how it is parameterized or configured
2. Reproduce the steps performed in Galaxy on the command line in order to become trained in how command line tools are used
3. To run batches of sequences through a pipeline that reproduces all the individual command line tools used from step 2.

# Genome Assembly

For this purpose we will initially start with the process of assembling contigs from raw reads. Contigs is a term that means contiguous DNA and refers to the consensus sequence that is formed when sequence reads (usually from fastq files) are 'stitched together' to form large regions from the genome. With short reads, repetitive sequences usually prevent complete closed genomes from being produced but instead the end result is usually smaller pieces of contiguous DNA that make up the most of the genome.

# Genome Assembly Tutorial

Galaxy is an open source, web-based platform for accessible, reproducible, and transparent computational biomedical research. It allows users without programming experience to easily specify parameters and run individual tools. It also captures run information so that any user can repeat and understand a complete computational analysis.

## Learning Objectives

By following this tutorial you will be able to:

- Be able to explain the principles of *de novo* assembly
- Login to a Galaxy server.
- Upload data to a Galaxy server from a file on your local computer
- Access and run the software tools applying them to data you have uploaded
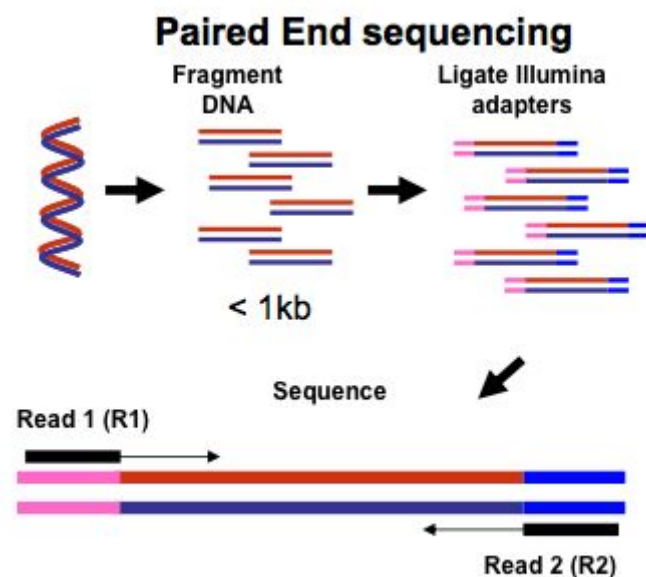- Assemble a genome from a pair of fastq files

## Stages in the Assembly Process
- Raw fastq QC assessment
- Fastq trimming
- Trimmed fastq QC assessment
- Assembly
- Assembly QC

## Introduction to Next Generation Sequencing

Next generation or high throughput sequencing involves massively parallel sequencing of small fragments of DNA that have been generated from an original nucleic acid source. This can be genomic DNA, PCR amplicons or cDNA generated from RNA. In this tutorial you will be working with sequence data that has been generated from whole genome DNA extracts. In this case the genomic DNA is fragmented and then Illumina adapter sequences ligated so that common sequencing primers can be used prime from the fragments and generate sequence. The fragmentation and ligation can either be done in one step (e.g with the enzymatic NextEra system that uses a process they call tagmentation) or in two steps such as with the TruSeq process where the DNA is sheared physically and the adapters ligated subsequently. The Illumina sequencing technology works best when the fragment size (also known as insert size) is 300 - 500bp. With paired end sequencing both ends of the fragment are sequenced but these reads (one in the forward and one in the reverse direction) may not meet. If the insert size is too small the reads can be, and for assembly, should be merged.

A nice video of this process can be found here.

At the end of the sequencing process when paired end sequencing is used each sample sequenced will produce a pair of fastq files corresponding to one file with all the read 1s from the fragments and in the second file all the read 2s. The order of the reads in these files is the same so that the R1 and R2 for each fragment are synchronised. However since the fragments are generated randomly in the case of library preparation from whole genomes they are not in any order in relation to the genome. This can only be achieved by giving each read context through *de novo* assembly or aligning (usually known as mapping) to a close reference genome.
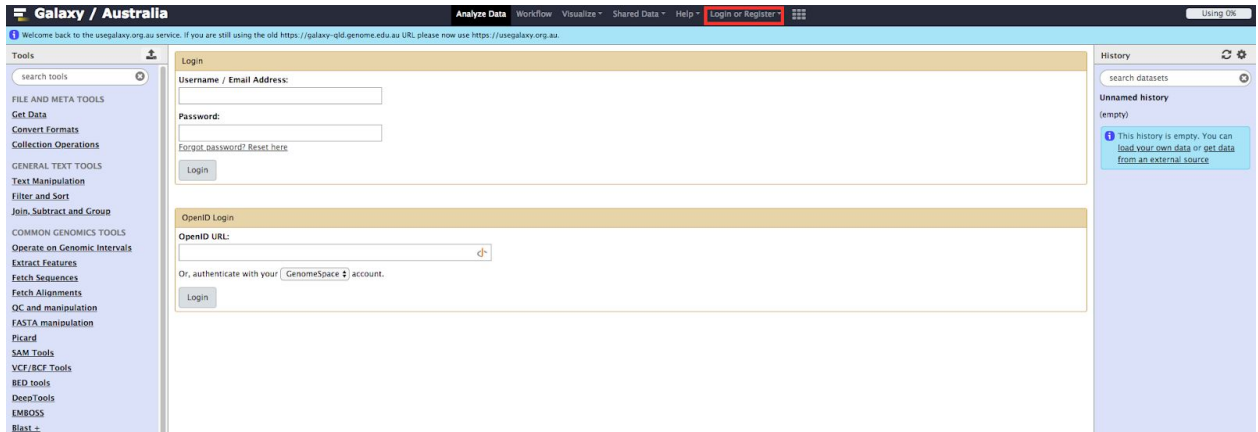
### Fastq Format

A FASTQ file normally has four lines per sequence fragment.

Line 1 begins with a '@' character and is followed by a sequence identifier and an optional description (like a FASTA title line).
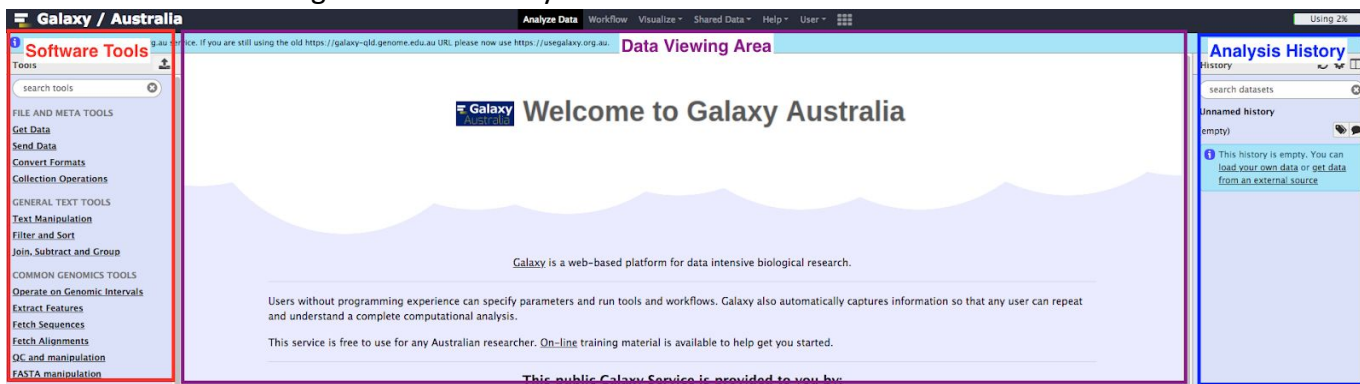Line 2 contains the sequence as letters that represent the nucleotides.
Line 3 begins with a '+' character and is optionally followed by the same sequence identifier (and any description) again.
Line 4 encodes the quality values for the sequence in Line 2, and contains the same number of symbols as letters in the sequence.

So for example  FASTQ file containing a single sequence might look like this:

```
@SEQ_ID
GATTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT
+
!''*(((((***+))%%%++)(%%%%).1***-+*''))**55CCF>>>>>>CCCCCCC65
```

The quality metrics are encoded as shown below (taken from https://en.wikipedia.org/wiki/FASTQ_format#Encoding). Most sequence data is now in Phred+33 format where quality ranges from 0 to 40 and is encoded as ! to J.

```
  SSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSS...................................................
  .........................XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX..............................
  ...............................IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII.........................
  ..................................JJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJ.......................
  LLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLL....................................................
  !"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz{|}~
  |                         |    |         |                              |                          |
  33                        59   64        73                             104                        126
  0........................26...31.......40
                           -5....0........9.............................40
                                0.......9.............................40
                                    3.....9...........................41
  0.2.....................26...31........41


S - Sanger        Phred+33,  raw reads typically (0, 40)
X - Solexa        Solexa+64, raw reads typically (-5, 40)
I - Illumina 1.3+ Phred+64,  raw reads typically (0, 40)
J - Illumina 1.5+ Phred+64,  raw reads typically (3, 41)
    with 0=unused, 1=unused, 2=Read Segment Quality Control Indicator (bold)
    (Note: See discussion above).
L - Illumina 1.8+ Phred+33,  raw reads typically (0, 41)
```

## Using Galaxy to assemble a Genome

### Logging onto Galaxy, navigating the system and uploaded data

1. Go to the Galaxy Australia website at https://usegalaxy.org.au/
   If you haven't registered, do so via the Register link and then Login

2. There are three main regions in the Galaxy window



Software tools are selected by searching by name on the 'Software Tools' region on the left side of the window. Clicking on them brings up the options for running the software in the central 'Data viewing area' and these software tools can be applied to data seen in the 'Analysis History' region on the right.

3. When you first login you will have no data to analyse, so let's enter some data. For the purposes of this tutorial we will use a pair of fastq.gz (gzipped fastq) files from small *Salmonella typhi*.
   These can be downloaded from the following 2 links:
   https://drive.google.com/file/d/14sGYl1hnkJwCKQmSzwKgh_t_Cda0pST_
   https://drive.google.com/file/d/14pv42xRDKyDZNtR_VNz_XFV9fjyoBxMl
   Download these 2 files to a directory on your computer

4. On the left hand side click on Get Data->Upload File or on the upload file icon

5. Drag the 2 fastq files you downloaded onto the pop up window that appears. In the Type dropdown for each file make sure you select 'fastqsanger.gz' and then click on the blue 'Start' button



Depending on your internet speed it may now take a few minute to upload the data. The files will then appear in the History on the right hand section of the screen in a small box that is first grey whilst it is waiting to be processed, then yellow as it is processing and finally green when the file has been registered in the galaxy system. In the screenshot below you can see 1 data file that has finished uploading and being processed and another that has uploaded and is now being processed so that it can be accessed on Galaxy.

**Raw fastq QC assessment**

1. Now that you have data in the galaxy server you can run software on that data. Behind the scenes it is using CLI (command line interface) UNIX software but galaxy provides a web browser interface to those software. With the 2 fastq files we have uploaded we will first assess the quality using the

fastqc tool. Type fastqc into the search tools box on the left



2. Click on the FastQC link. This will bring up a screen in the central data viewing area where you can select parameters for the software. In this case you just want to run fastqc with its standard parameters but first make sure you select the _1 fastq file (Read 1):



Then click the blue Execute button.

On the right hand side you will see the two output files for the fastqc software start to be processed as fastqc is run using Galaxy

When fastqc completes the boxes will turn green. Click on the eye icon  in the box that says 'FastQC …. Web page'

3. In the central data view you will now see the output that has been created when running fastqc



If you click on 'Per base sequence quality' you will see that the quality of this fastq file is very high. Anything above 25 is usually considered good quality.

However if you click on the adapter content link you will see that there is some contamination of the Nextera Transposase sequence suggesting that the sequencing inserts may be small in some of the sequence data. In this case the sequencing has read all the way through from one side of the insert to the other into the adapter sequences on the other side of the sequence insert. Since this is not sequence from the originating DNA this should be removed.

Before going further run fastqc on the second fastq file. Have a look at the report from this and note if there are any significant differences from the read 1 data. At the end of this your history should look something like this:

## Fastq trimming

1.  Because of the adapter contamination we will use a program called Trimmomatic to remove any low quality data and the adapters. Search for this by typing 'Trimmomatic' into the search box on the left. Click on the link that says 'Trimmomatic flexible read trimming tool for Illumina NGS'

    In the central area the software should now be configured with a few parameters
    a.  First select Paired end
    b.  Then ensure that you have both the read 1 and read 2 files selected (_1 and _2)
    c.  Then click yes to perform the ILLUMINACLIP step and
    d.  choose the Nextera (paired-ended) adapter sequences
    e.  The other parameters should be left as 2, 30, 10 and 8 respectively

By specifying this parameter we are telling Trimmomatic to remove Illumina adapters

2. Next you should add some more parameters for trimming. These will be
   a. SLIDINGWINDOW (remove areas where the average quality is less than 20 across 4 base)
   b. LEADING (cut off low quality bases at the 5' end below 25)
   c. TRAILING (cut off low quality bases at the 3' end below 25)
   d. MINLEN (remove read pairs if either is less than 30 bases after trimming)

**Trimmomatic Operation**



These extra settings are added using the **➕ Insert Trimmomatic Operation** button. Once you have entered all of these click the blue 'Execute button'

11. When trimmomatic has finished you will see 4 files have been produced it is the R1_paired and R2_paired files that we will use for assembly. These are the fastq files that have been trimmed for quality and adapter sequences. They are called _paired since these are the paired reads that have been kept. For downstream processes we do not want any unpaired reads.
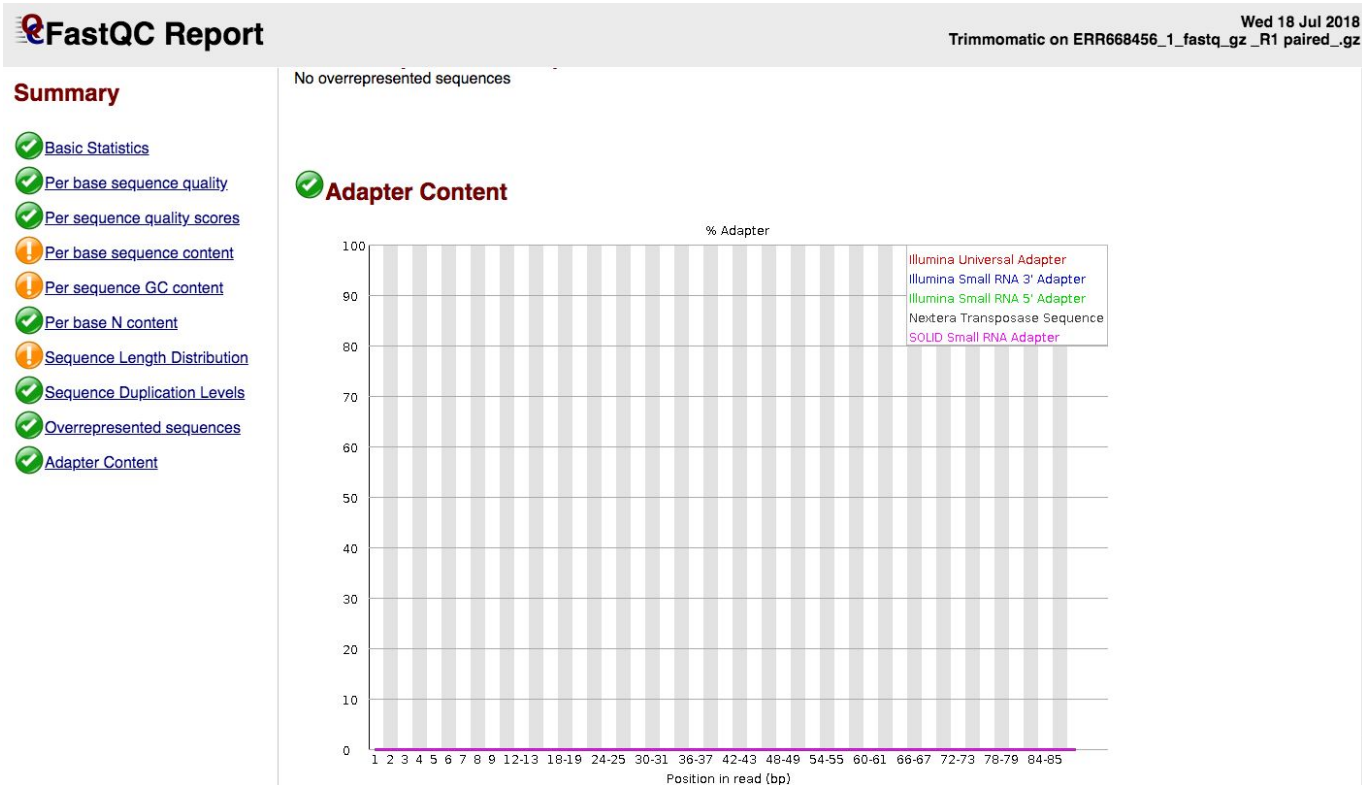


Before assembling them however, you should first check on the quality.

**Trimmed fastq QC assessment**

1. Run fastqc again but this time select first the trimmed R1 paired file and then the trimmed R2 paired file.



Make sure you do this for both files.

2. Look at the output and you will see that the adapters have been removed



**Assembly**

1. Finally we come to assembly. Search for SPAdes in the tool search box.
   Enter the following options
   a. Single-cell: No
   b. Run only assembly: Yes
   c. Careful correction: No
   d. Automatically choose k-mer values: No
   e. K-mers to use: 21,33,43,53,63,75
   f. Coverage Cutoff: Off

g. Library type: Paired-end/ Single Reads
h. select the R1 paired and R2 paired files from the previous Trimmomatic step

Click on the blue Execute button to run SPAdes



Please note that because Galaxy is a shared resource it may sometime take a while for the SPAdes assembly to start (boxes turn yellow). Patience may be required :)
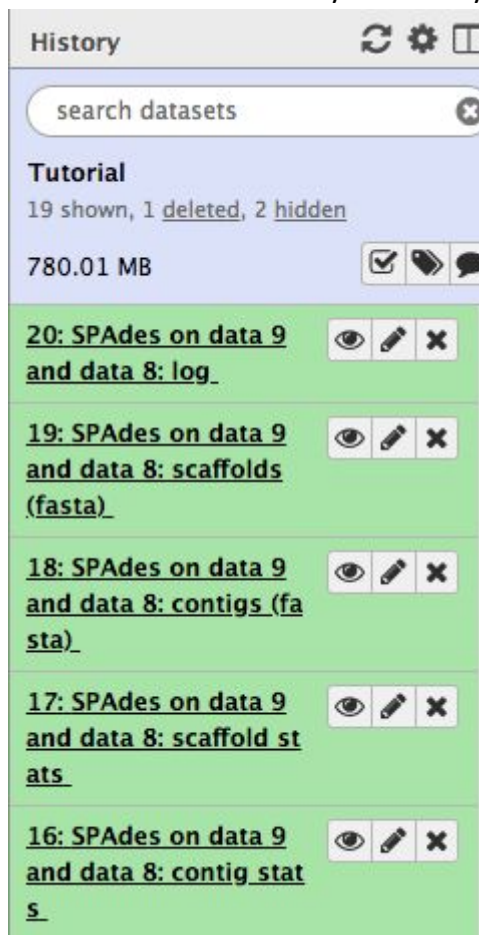
**While you are waiting**

You may want to read
● the original SPAdes publication: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3342519
● the SPAdes manual: http://spades.bioinf.spbau.ru/release3.11.1/manual.html

- And if you are really keen this primer about bacterial WGS
  http://cmr.asm.org/content/29/4/881.full.pdf

2. When SPAdes is finished your history should look something like this



The outputs you will see are

**scaffolds.fasta** contains resulting scaffolds in fasta format (contigs joined by paired end read information but where the gaps are padded with Ns, see this article for some more details https://genome.jgi.doe.gov/help/scaffolds.jsf )

**contigs.fasta** contains resulting contigs in fasta format

**scaffolds stats** length and coverage information about each resulting scaffold

**contigs stats** length and coverage information about each resulting contig

**Assembly QC assessment**

3. Now you will assess the quality of the assembly using Quast. Type Quast into the tool search bar. Select the SPAdes **contigs** file and specify a reference genome size of 4900000 (approximate size of *S.typhi*). Click on execute to run the Quast quality assessment. The report.tsv output will look like this

| Assembly | SPAdes_on_data_9_and_data_8__contigs__fasta_ |
|---|---|
| Assembly | SPAdes_on_data_9_and_data_8__contigs__fasta_ |
| # contigs (>= 0 bp) | 171 |
| # contigs (>= 1000 bp) | 78 |
| Total length (>= 0 bp) | 4911403 |
| Total length (>= 1000 bp) | 4886054 |
| # contigs | 90 |
| Largest contig | 488881 |
| Total length | 4894750 |
| Estimated reference length | 4900000 |
| GC (%) | 52.02 |
| N50 | 144430 |
| NG50 | 144430 |
| N75 | 66210 |
| NG75 | 66210 |
| L50 | 12 |
| LG50 | 12 |
| L75 | 24 |
| LG75 | 24 |
| # N's per 100 kbp | 0.00 |

The critical figures to look at are
   a. # (number) of contigs: smaller the better
   b. Total length: Should be approximately the size of the genome expected for the species
   c. N50/NG50: Larger the better (see this article for an explanation)

**Conclusions**

Well done you have completed a genome assembly and simple QC. The next step will be to dive into the command line and do the same from there.

**Version Control Table**

| Title | Genome Assembly Tutorial | | | |
|---|---|---|---|---|
| **Description** | A document describing how to perform *de novo* genome assembly for a single bacterial genome starting from a pair of Fastq files generated with short read sequencing using the Galaxy web platform | | | |
| **Created By** | Anthony Underwood | | | |
| **Date Created** | 27th July 2018 | | | |
| **Maintained By** | Anthony Underwood | | | |
| **Version Number** | **Modified By** | **Modifications Made** | **Date Modified** | **Status** |
| 1.0 | Anthony Underwood | First version | 27th July 2018 | First Live Version |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |