

Galaxy for virologist training

Exercise 3: Illumina Assembly 101

Title	Galaxy
Training dataset:	PRJEB43037 - In August 2020, an outbreak of West Nile Virus affected 71 people with meningoencephalitis in Andalusia and 6 more cases in Extremadura (south-west of Spain), causing a total of eight deaths. The virus belonged to the lineage 1 and was relatively similar to previous outbreaks occurred in the Mediterranean region. Here, we present a detailed analysis of the outbreak, including an extensive phylogenetic study. This is one of the outbreak samples.

Questions:

- What is assembly?
- How can I evaluate my assembly?

Objectives:

- Understand assembly concept
- Learn how to interpret assembly quality control metrics

Estimated time: 40 min

1. Description

Sometimes, we don't have a reference genome to map against, or we want to reconstruct a genome without any bias caused by a reference. In such cases, we need to do a *de novo assembly*. This type of analysis tries to reconstruct the original genome without any template, using only the reads. Some considerations:


- When we assemble, the longer the reads are and the longer the size of the library fragments the easier it gets for the assembler. That's why PacBio or Nanopore are recommended for assembly. Think of it like a puzzle, the bigger the pieces, the easier it is to form the image.
- It's almost impossible to reconstruct the entire genome of a large-genome microorganism with only one sequencing, although it can be done for smaller ones, like viruses.
- Assembly is not recommended for amplicon based libraries due to the depth of coverage unevenness and the amplicons intrinsic bias.

2. Upload data to galaxy

Training dataset

- Experiment info: PRJEB43037, WGS, Illumina MiSeq, paired-end
- Fastq R1: [ERR5310322_1](ftp://ftp.sra.ebi.ac.uk/vol1/fastq/ERR531/002/ERR5310322/ERR5310322_1.fastq.gz) - url :
`ftp://ftp.sra.ebi.ac.uk/vol1/fastq/ERR531/002/ERR5310322/ERR5310322_1.fastq.gz`
- Fastq R2: [ERR5310322_2](ftp://ftp.sra.ebi.ac.uk/vol1/fastq/ERR531/002/ERR5310322/ERR5310322_2.fastq.gz) url :
`ftp://ftp.sra.ebi.ac.uk/vol1/fastq/ERR531/002/ERR5310322/ERR5310322_2.fastq.gz`
- Reference genome NC_009942.1: [fasta](#) -- [gff](#)

Create new history

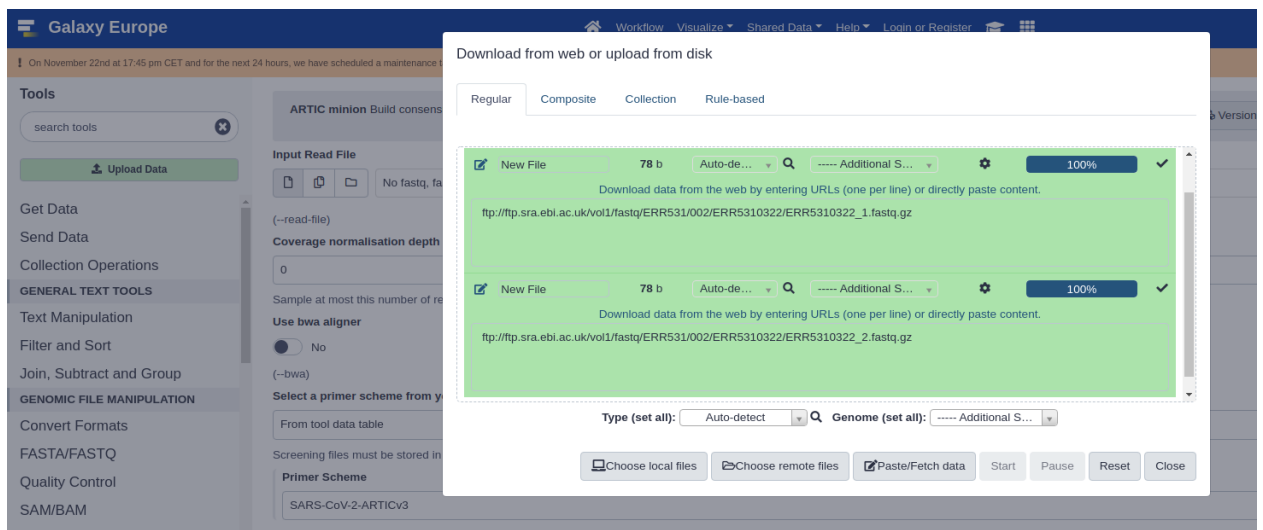
- Click the  icon at the top of the history panel and create a new history with the name `illumina assembly 101 tutorial` as explained [here](#)

Upload data


- Import and rename the read files `ERR5310322_1` and `ERR5310322_2`
 1. Click in upload data.
 2. Click in paste/fetch data
 3. Copy url for fastq R1 (select and Ctrl+C) and paste (Ctrl+V).
 4. Click in Start.
 5. Wait until the job finishes (green in history)
 6. Do the same for fastq R2.

Table of Contents

- [1. Description](#)
- [2. Upload data to galaxy](#)
 - [Training dataset](#)
 - [Create new history](#)
 - [Upload data](#)
 - [Assemble reads with Spades](#)
 - [Assembly quality control with Quast](#)



- Rename R1 and R2 files.

1. Click in the  in the history for ERR5310322_1.fastq.gz
2. Change the name to ERR5310322_1
3. Do the same for R2.

Name

Info

Annotation

Add an annotation or notes to a dataset; annotations are available when a history is viewed.




Database/Build

----- Additional Species Are Below -----

7: GCF_000875385.1_Viral
Proj30293_genomic.fna.gz

1 sequences

format: fasta.gz, database: ?

display with IGV local

```
>NC_009942.1 West Nile virus lineage 1, con
AGTAGTTGGCTGTGTGAGCTGACAACTAGTAGTTTGTG
TAGCAGGAAATCTCATATCTAAGAAACCAAGAGGCCGCG
CCGCGCTGTGTCTTATTTGACTGAAAGAGGCTATTTGAC
GCTCTTGGGTTTCTCAGGTTTCAGAGCAATTGCTCGAGCC
```

6: ERR5310322_2

5: ERR5310322_1






- Import the reference genome and GFF file.

```
https://ftp.ncbi.nlm.nih.gov/genomes/all/GCA/009/858/895/GCA_009858
https://ftp.ncbi.nlm.nih.gov/genomes/all/GCA/009/858/895/GCA_009858
```

Descargar de la red o cargar desde disco

Regular Composite Collection Rule-based

You added 1 file(s) to the queue. Add more files or click 'Start' to proceed.

Name	Size	Type	Genome	Settings	Status
 New File	267 b	Auto-de...	unspecified (?)		0%
Download data from the web by entering URLs (one per line) or directly paste content. <pre> oi.nlm.nih.gov/genomes/all/GCF/000/875/385/GCF_000875385.1_ViralProj30293/GCF_000875385.1_ViralProj30293_genomic.fna.gz oi.nlm.nih.gov/genomes/all/GCF/000/875/385/GCF_000875385.1_ViralProj30293/GCF_000875385.1_ViralProj30293_genomic.gff.gz </pre>					
Type (set all): Auto-detect Genome (set all): unspecified (?)					
<div>  Elegir archivos locales  Choose remote files  Paste/Fetch data Start Pause Reset Close </div>					

- Rename the reference genome and gff file.

- Name

GCF_000875385.1_ViralProj30293_genomic.fna.gz

Info

Annotation

Add an annotation or notes to a dataset; annotations are available when a history is viewed.

Database/Build

----- Additional Species Are Below -----
- 7: GCF_000875385.1_ViralProj30293_genomic.fna.gz

1 sequences

format: **fasta.gz**, database: ?

display with IGV local

```

>NC_009942.1 west Nile virus lineage 1, co
AGTAGTTCGGCTGTGTGACGTGACAACTTAGTAGTGTGT
TAGACACAGAGATCTCGATATCTTAGAAACAGAGAGCCGCG
CCGCGGTGTGTCTCTTATTTGAACTGAAGAGGCGCTATTTAA
GCTCTCTTGGCGTCTTCAAGTTCACAGCAATTGCTCCGACAC

```

6: ERR5310322_2

5: ERR5310322_1

1. Search **Spades** in the search tool box and select *mvirSPAdes de novo assembler for transcriptomes, metatranscriptomes and metaviromes*
2. Single-end or paired-end short-reads > Paired-end: individual datasets
3. **FASTQ RNA-seq file(s): forward reads:** ERR53103221; **_FASTQ RNA-seq file(s): reverse reads:** ERR5310322_2
4. Select optional output file(s) > Scaffolds stats
5. Click execute and wait.

Herramientas

- spades**
- Cargar Datos
- Show Sections

SPAdes genome assembler for genomes of regular and single-cell projects

metaviralSPAdes extract and assembly viral genomes from metagenomic data

metaplasmidSPAdes extract and assembly plasmids from metagenomic data

plasmidSPAdes extract and assembly plasmids from WGS data

metaSPAdes metagenome assembler

rnaviralSPAdes de novo assembler for transcriptomes, metatranscriptomes and metaviromes

biosyntheticSPAdes biosynthetic gene cluster assembly

coronaSPAdes SARS-CoV-2 de novo genome assembler

rnaSPAdes de novo transcriptome assembler

Shovill Faster SPAdes assembly of Illumina reads

The screenshot shows the Galaxy web interface for the tool **rnaviralSPAdes**. The top navigation bar includes "Herramientas" (Tools) and "History". The main panel displays the tool's description: "rnaviralSPAdes de novo assembler for transcriptomes, metatranscriptomes and metaviromes (Galaxy Version 3.15.4+galaxy2)". Below this, the "Operation mode" is set to "Assembly and error correction". A note states: "To run read error correction, reads should be in FASTQ format." The "Single-end or paired-end short-reads" section has a dropdown menu set to "Paired-end: individual datasets". Another note says: "It assumes that all samples belong to the same library. If you want to use samples from two different libraries, include the second library as additional set of short-reads." Under the heading "FASTQ RNA-seq file(s): forward reads", there are three input fields. The first field contains the sample ID "NC_009942.1". The second field contains "ERR5310322_2". The third field contains "ERR5310322_1". To the right of these fields are icons for uploading files and saving the job. Below this section, the "FASTQ RNA-seq file(s): reverse reads" section is visible, with similar input fields. At the bottom, the "Type of paired-reads" dropdown is set to "Default (--pe)". On the right sidebar, the "History" panel shows a list of previous jobs, including "illumina assembly 101 tutorial" and several NCBI accession numbers.

The screenshot shows the Galaxy Europe web interface. At the top, there's a navigation bar with logos and user information. The main area displays the 'Set Phred quality offset' section for the SPAdes tool. A red arrow points to the 'Execute' button. On the right, a 'History' panel lists previous runs.

Herramientas

spades x

Cargar Datos

Show Sections

SPAdes genome assembler for genomes of regular and single-cell projects

metaviralSPAdes extract and assembly viral genomes from metagenomic data

metaplasmidSPAdes extract and assembly plasmids from metagenomic data

plasmidSPAdes extract and assembly plasmids from WGS data

rnaSPAdes metagenome assembler

rnairalSPAdes de novo assembler for transcriptomes, metatranscriptomes and metaviromes

biosyntheticSPAdes biosynthetic gene cluster assembly

coronaSPAdes SARS-CoV-2 de novo genome assembler

7 Execute

Set Phred quality offset

Auto

Phred quality offset in the input reads. Default: auto-detect (--phred-offset)

Select optional output file(s)

Select/Unselect all

- Assembly graph
- Assembly graph with scaffolds
- Contigs
- Scaffolds

6

Contigs paths

Corrected reads

Contigs stats

Log

Scaffolds paths

Scaffolds stats

What it does

SPAdes - St. Petersburg genome assembler - is an assembly toolkit containing various assembly pipelines.

rnairalSPAdes is a pipeline specially designed for de novo assembler tailored for RNA viral datasets (transcriptome, metatranscriptome and metavirome).

Input

SPAdes takes as input paired-end reads, mate-pairs and single (unpaired) reads in FASTA and FASTQ. For IonTorrent data SPAdes also supports unpaired reads in unmapped BAM format (like the one produced by Torrent Server). However, in order to run read error correction, reads should be in FASTQ or BAM format. Sanger, Oxford Nanopore and PacBio CLR reads can be

History

buscar conjuntos de datos x

illumina assembly 101 tutorial

31.4 MB

4 : NC_009942.1.GFF #gff

3 : NC_009942.1 #reference

2 : ERR5310322_2 #reverse

1 : ERR5310322_1 #forward

Warning

:coffee::fork_and_knife::clock330: **Assembly takes time!** There is no such thing as Assembly in real time. It can take anywhere between 90 minutes and two hours.

Questions:

Click the :eye: icon in the history: Spades Contigs stats.

► How many contigs has been assembled?

Click the :eye: icon in the history: Spades scaffolds.

Assembly quality control with Quast

1. Search Quast in the search tool box.
2. rnaviralSpades Scaffolds
3. Use a reference genome: Yes. Select the NC_009942.1 fasta file previously loaded.
4. Genomic feature positions in the reference genome > NC_009942. gff file previously loaded.

Galaxy Europe Flujo de Trabajo Visualizar Datos Compartidos Ayuda Usuario

Herramientas ☆

quast 1

Cargar Datos

Show Sections

Quast Genome assembly Quality 2

rnaQUAST A Quality Assessment Tool for De Novo Transcriptome Assemblies

FLUJOS DE TRABAJO

Todos los flujos de trabajo

Quast Genome assembly Quality (Galaxy Version 5.2.0+galaxy0) ☆

Use customized names for the input files?

No, use dataset names

They will be used in reports, plots and logs

Contigs/scaffolds file

12: rnaviralSPAdes on data 2 and data 1: Scaffolds

11: rnaviralSPAdes on data 2 and data 1: Contigs

3: NC_009942.1 (as fasta)

3

Reads options

Disabled

Currently, the supported read types are Illumina unpaired, paired-end and mate-pair reads, PacBio SMRT, and Oxford Nanopore long reads.

Type of assembly

Genome

Use a reference genome?

Yes

Many metrics can't be evaluated without a reference. If this is omitted, QUAST will only report the metrics that can be evaluated without a reference.

Reference genome

12: rnaviralSPAdes on data 2 and data 1: Scaffolds

11: rnaviralSPAdes on data 2 and data 1: Contigs

3: NC_009942.1 (as fasta)

5

Galaxy Europe Flujo de Trabajo Visualizar Datos Compartidos Ayuda Usuario

Herramientas ☆

quast x

Cargar Datos

Show Sections

Quast Genome assembly Quality

rnaQUAST A Quality Assessment Tool for De Novo Transcriptome Assemblies

FLUJOS DE TRABAJO

Todos los flujos de trabajo

Genomic feature positions in the reference genome

4: NC_009942.1 GFF (as bed)

Gene coordinates for the reference genome (--features)

Operon positions in the reference genome

Nothing selected

Operon coordinates for the reference genome (--operons)

Compute k-mer-based quality metrics?

No

It is recommended for large genomes. This may significantly increase memory and time consumption on large genomes (--k-mer-stats)

Generage Circos plot

No

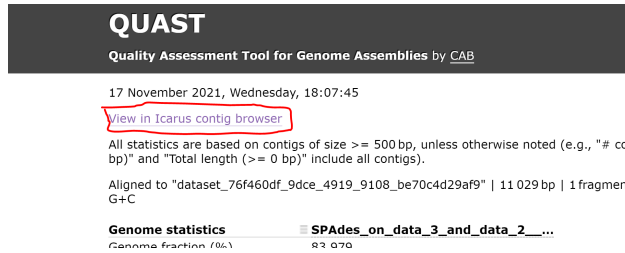
Plot Circos version of Icarus contig alignment viewer (--circos)

7 Execute

1. Click the :eye: icon Quast HTML report.

- ▶ How much of or reference genome have we reconstructed?
- ▶ How many contigs do we have greater than 1000 pb?
- ▶ How long is the largest contig in the assembly?
- ▶ Which is the N50?

2. Open the Icarus viewer in the quast report.



- ▶ How did the contig align against our reference genome?

This training history is available at: <https://usegalaxy.eu/u/s.varona/h/illumina-assembly-101-tutorial>