



Ensamblajes de Genomas Virales

Modulo_2 - Ejemplo Sars-CoV-2

Paola Rojas-Estevez

M.Sc Biología Computacional

Grupo Genómica de Microorganismos Emergentes Instituto Nacional de Salud de Colombia

Links a utilizar



https://github.com/BCVI/2do-Workshop-Genomica-Viral



https://colab.research.google.com/





Este repositorio se centra en la capacitación práctica en análisis de secuencia del genoma viral, principios científicos, análisis e interpretación de datos genómicos de secuenciación a gran escala.





2do Workshop: Bioinformática en Genómica Viral

Este repositorio se centra en la capacitación práctica en análisis de secuencia del genoma viral, principios científicos, análisis e interpretación de datos genómicos de secuenciación a gran escala.

Infraestructura computacional del curso

Este curso se desarrollará por medio de un servidor privado y usando Google Colab, un servicio de acceso libre.

Compromiso de tiempo - 2 semanas

2 Semanas: (1 semana de curso - 1 semanas de seguimiento) 4 horas por semana de tiempo personal y 1 hora de trabajo remoto con los instructores.

Fecha de inicio: El curso se desarrollará en modo remoto durante la semana del 18-21 de Febrero 2025.

Programa del curso

El programa cubrirá los siguientes temas generales:

- Introducción a Linux
- · Ensamble y Anotación de Genomas Virales con Referencia
- Identificación de Secuencias Virales (Metagenómica Viral y Profagos)
- · Clustering y Clasificación Taxonómica

Equipo de instructores

- · Paola Rojas-Estevez, Instituto Nacional de Salud, Bogotá Colombia.
- Laura Camelo-Valera, Universidad de McGill, Canadá, Montréal Canadá.
- Luis Alberto Chica Cárdenas, Washington University, St. Louis USA.
- Juan Manuel Hurtado Ramirez, Instituto de Biotecnología-UNAM, Cuernavaca México.
- Gamaliel López-Leal, Centro de Investigación en Dinámica Celular-UAEM, Cuernavaca México.



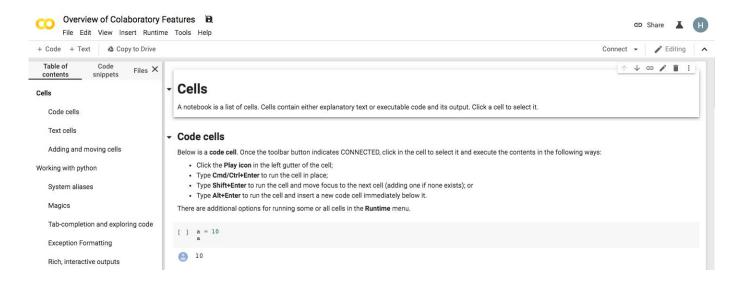
Repositorio GitHub

Carpeta Digital de Código

Google Colab

Entorno de Jupiter notebook personalizado por google

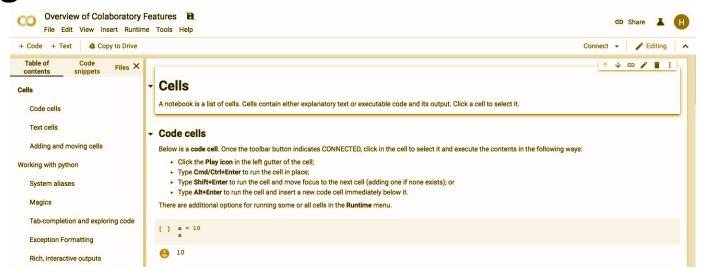




- Escribir Notas
- 2. Ejecutar líneas de código con Python

Google Colab





"Colaboratory" → permite programar y ejecutar Python/bash en tu navegador

- No requiere configuración
- Acceso a GPUs sin coste adicional
- Permite compartir contenido fácilmente

Notebooks Colab





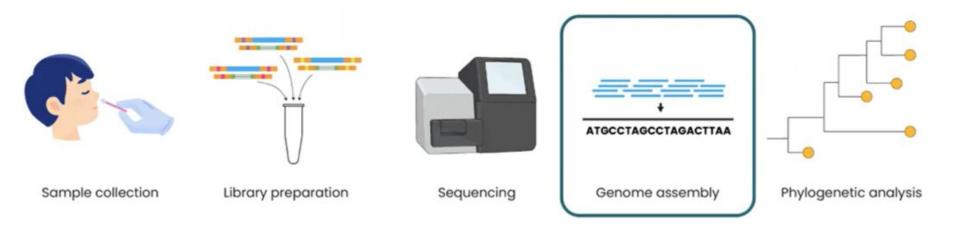




Ensamblajes de Genomas Virales

Modulo 2 - Calidad / Ensamblaje

Workflow epidemiología genómica



Secuenciación Genómica









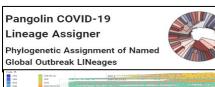
Vigilancia Genómica

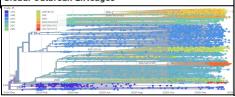
DENV, CHIKV, ZIKV, Rabies, Measles, POXVIRUS, H1N1, VIH

Metagenómica de virus emergentes y virodiversidad



Nanopore



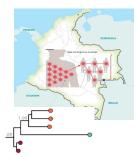






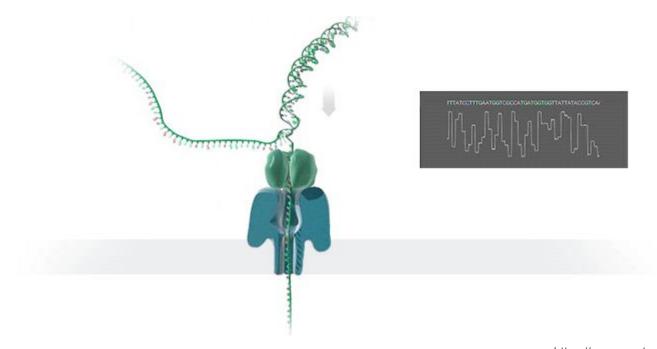


SARS-CoV-2



Secuenciación Nanopore

Archivos fast5/pods5: datos de señal eléctrica



Paso a Paso Análisis

1. Llamado de bases guppy_basecaller

2. Control de Calidad artic guppyplex

3. Ensamblaje artic minion



4. Estimación del Linaje Pangoline / Nextclade



5. Revisión de Mutaciones NEXTSTRAIN



Control de Calidad

Formato fastq

Todos los secuenciadores producen datos en un formato llamado **fastq**. Todas las secuencias con un fastq están representadas por 4 líneas:

La calidad de las secuencias se representa como un carácter del código ASCII. Consulte <u>aquí</u> para obtener una explicación. Los valores numéricos corresponden a los valores de calidad phred

Cargar los notebooks



Abrir cuade	erno			
Ejemplos	>	Escribe una URL de GitHub o busca por organización o usuario	ositorios privados	Copiar URL del curso
Recientes	>	https://github.com/BCVI/Taller-Virtual-Bioinformatica-Genomica-Viral.git	Q	
Google Drive	>	Repositorio: Rama:		
GitHub	>	Ruta		
Subir	>	Notebook/Modulo_2/Module_2_Part1_QC-Nanopore.ipynb		Notebook por módulo
		Notebook/Modulo_2/Modulo_2_Part2_Reconstruccion_genoma_Nanopore_SARS_Co	Q C	Notebook por modulo
-				



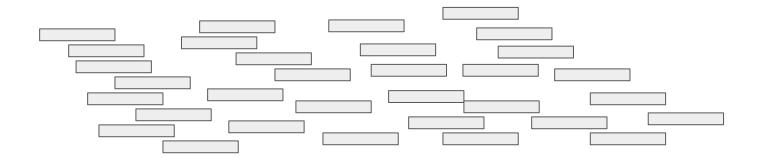


Ensamblajes de Genomas Virales

Módulo 2 - Práctica 2 [Ensamblaje]

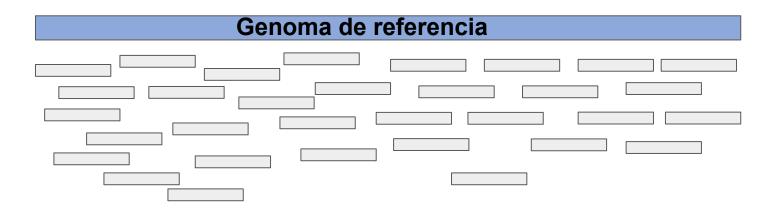
Métodos para ensamblar un Genoma

Ensamblaje de novo, Se infiere la relación de las lecturas por sobrelapamiento En otras palabras, Comparamos lecturas entre sí.



Métodos para ensamblar un Genoma

¿Qué pasa si en cambio comparamos las lecturas con el genoma de referencia?



Software

BWA (Burrows-Wheeler Aligner)



y Uso: Lecturas cortas de Illumina

Bowtie2

https://github.com/BenLangmead/bowtie2

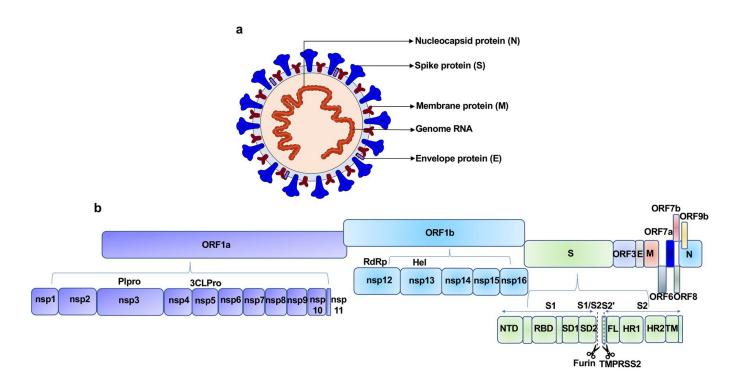
P Uso: Alineamiento rápido de lecturas cortas

Minimap2

https://github.com/lh3/minimap2

📌 Uso: Lecturas largas de Nanopore y PacBio

29,903 bases - Genoma de SARS-CoV-2



La importancia de las mutaciones



Guppy - Basecalling

guppy_basecaller -i_mnt/fast5 -s mnt/output_dir -c dna_r9.4.1_450bps_hac.cfg -m template_r9.4.1_450bps_hac.jsn -barcode_kits "EXP-NBD196" --require_barcodes_both_ends -compress_fastq -trim_barcodes_x "cuda:0" --gpu_runners_per_device [num] --num_callers [num] --chunks_per_runner [num] --chunk_size_[num]

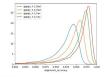
Archivos de entrada (fast5) y salida (fastq)

Opciones del programa

Modelo para redes neuronales (fast o hac)

Barcode demultiplexing

Opciones de optimización



tps://mirror.oxfordnanoportal.com/software/analysis/ont-guppy_6.0.1_linux64.tar.gz

Estimación del linaje

PANGOLIN:/ https://pangolin.cog-uk.io https://cov-lineages.org/resources/pangolin/installation.html

Importante mantener actualizado

5. Revisión de mutaciones



NEXTSTRAIN: https://clades.nextstrain.org/

Disponible en web y línea de comandos

8. Someter a base de datos

- Mínimo 50% de cobertura (mínimo GISAID)
- Se relacionan secuencias ingresadas con metadatos, info pacientes, información geográfica, institución que provee secuencia.
- https://www.gisaid.org/
- Códigos ISO departamentos: https://en.wikipedia.org/wiki



2. Control de calidad

- Se corre por cada barcode o muestra
- Remoción de lecturas quiméricas y filtro por tamaño <u>artic guppyplex --min-length 400 --max-length 700 --directory output_directory/barcode03</u> <u>--prefix run_name</u>
 - Directorio donde se guardan las secuencias filtradas
 - Prefijo con el que se guardan los archivos de salida



https://artic.network/ncov-2019/ncov2019-bioinformatics-sop.html

6. Métricas

- Profundidad promedio por posición
 samtools depth aln_\${Barcode}-sorted.bam
 - Cobertura

samtools mpileup aln_\${Barcode}-sorted.bam

O

Usar columna "Missing (N)" del análisis en Nextclade 1-N/29903

Archivo BAM ordenado producto del alineamiento

3. Ensamblaje (Minimap2 y Nanopolish)

artic minion --normalise 200 --threads 4 --scheme-directory

Se corre por cada barcode o muestra

~/artic-ncov2019/primer_schemes --read-file run_name_\${barcode}.fastq --fast5 directory
path to fast5 --sequencing-summary path to sequencing summary.txt nCoV-2019/V3

samplename

Directorio donde están los datos crudos fast5

Prefijo con el que se guardamos los archivos de salida en el paso de control de calidad

Sequencing summary de Guppy. Único para todos los barcodes.

Nombre de la muestra (Usualmente el barcode que se analiza)



7. Reportar secuencias

COLUMNAS	DESCRIPCIÓN	EJEMPLO		
Institución que secuencia	Nombre de la institución que realizó el experimento de secuenciación	Instituto Nacional de Salud		
Código lab secuenciación	Asignado durante secuenciación	INS-COV-VG-005		
Tipo de muestreo	Descriptor estudio o tipo muestreo (Rutinario, Probabilístico, Hospitalizados, Fallecidos, Vacunados, Investigación propia,)	Probabilístico		
Lab_origen_de_la_muestra	Nombre del laboratorio de donde proviene la muestra	SURA, IPS Saludpass, Laboratorio clínico		
Código de origen de la muestra	Código que asignó el laboratorio de origen (Si está disponible)	C-346465; 2021081563113		
Documento de identidad	Número de identificación del paciente (Si está disponible)	1105111111		
Procedencia	Departamento de procedencia en código ISO (Ver tabla abajo)	DC; ARA; SAP		
Fecha de muestra	Fecha de toma de muestra (Si está disponible)	12/02/2022		
Frameshifts	Cambios en el marco de lectura (reportado por nextclade o su anotación)	ORF3a:257-276		
Deleciones	Deleciones en el genoma (reportado por nextclade o su anotación)	26158-26162		
Inserciones	Inserciones en el genoma (reportado por nextclade o su anotación)	21991:ACT		
		N:T205I,N:D377Y,ORF1a:T1055A,ORF1a:		
		T1538I,ORF1a:T3255I,ORF1a:Q3729R,OR		
Sustituciones aminoácidos	Todas las sustituciones de aminoácidos en el genoma (reportado por nextclade o su anotación)	F1b:P314L,ORF1b:A1131V,ORF1b:P1342S		
Sustituciones aminoácidos		ORF3a:Q57H,ORF3a:N257X,ORF8:T11K,O		
		RF8:P38S,ORF8:S67F,S:T95I,S:R346K,S:E4		
		RAK S-NS01V S-D614G S-D6R1H S-D050N		







Cargar los notebooks



Abrir cuade	erno			
Ejemplos	>	Escribe una URL de GitHub o busca por organización o usuario	ositorios privados	Copiar URL del curso
Recientes	>	https://github.com/BCVI/Taller-Virtual-Bioinformatica-Genomica-Viral.git	Q	
Google Drive	>	Repositorio: Rama:		
GitHub	>	Ruta		
Subir	>	Notebook/Modulo_2/Module_2_Part1_QC-Nanopore.ipynb		Notebook por módulo
		Notebook/Modulo_2/Modulo_2_Part2_Reconstruccion_genoma_Nanopore_SARS_Co	Q C	Notebook por modulo
-				

Gracias

Contacto:

<u>crojas@ins.gov.co</u> <u>paolarojasestevez@gmail.com</u>