



Curso Teórico-Práctico

EPIDEMIOLOGÍA GENÓMICA

UNA HERRAMIENTA PARA
FORTALECER LA VIGILANCIA DE
AGENTES INFECCIOSOS



Día 3: *Análisis Metagenómicos y Plataformas para Divulgación*
Bogotá / Septiembre 25 del 2024

Análisis de metagenómica (*mNGS*)

Paola Rojas - Estevez
M.Sc Biología Computacional
Grupo Genómica de Microorganismos Emergentes
Instituto Nacional de Salud de Colombia

Flujo de trabajo de mNGS

Diseño del estudio



- Muestras
- Controles
- Almacenamiento
- Condiciones de laboratorio

Extracción de ácidos nucleicos (1-2 hrs)



- Tipos
- Controles
- Control de calidad

Preparación de librerías (1-2 días)



- Remoción hospedero
- Fragmentación
- Generación de dsDNA
- End-repair
- Ligación de adaptadores
- Barcoding
- Limpieza
- Pooling
- Control de calidad

Secuenciación (0.5-3 días)



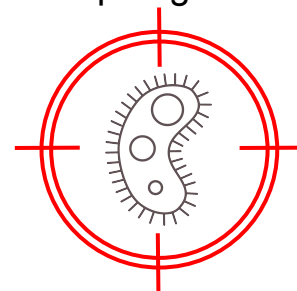
- Plataforma
- Longitud de reads
- Concentración de carga

Análisis de datos (0.5-4 hrs)



- Control de calidad de datos
- Remoción hospedero
- Alineamiento
- Blast
- Puntuación de hits

Identificación de patógenos



Flujo de trabajo de mNGS

Diseño del estudio



Preparación de librerías



Secuenciación

Extracción de ácidos nucleicos

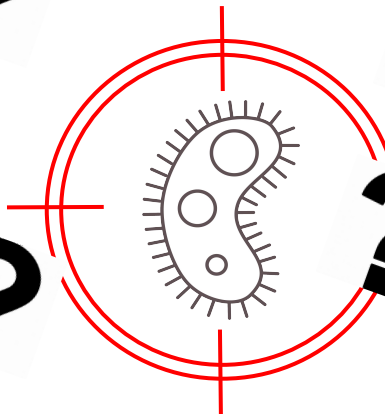


Análisis de datos

(0.5-4 hrs)



Identificación de patógenos



- Control de calidad de datos
- Remoción hospedero
- Alineamiento
- Blast
- Puntuación de hits



- Software gratuito
- Cloud-based pipeline
- Número libre de muestras
- Número libre de trabajos
- Pipeline de fondo
- Instructivo detallado
- Asistencia técnica y servicios:
help@czid.org

Módulos

Módulo	Datos de entrada	Soporte de tecnología de secuenciación
Metagenómica (mNGS)	Datos metagenómicos imparciales <i>(sin datos de amplicón/16 s/18 s)</i>	Illumina y Nanopore
Resistencia a los antimicrobianos	Secuenciación del genoma completo y mNGS	Illumina
SARS-CoV-2	Secuenciación de amplicones	Illumina y Nanopore
Genoma de consenso viral	Secuenciación de amplicones y WGS	Illumina

¿Cómo abordar el análisis de datos de secuencias metagenómicas de próxima generación (mNGS) a través de CZ ID?

Descripción general del análisis de datos de mNGS



Cargar datos

Realizar QC

- ☐ Verificar calidad de la muestra
- ☐ Evaluar controles de entrada
- ☐ Considerar contaminación usando modelos de fondo

Explorar Resultados y Analizar Datos

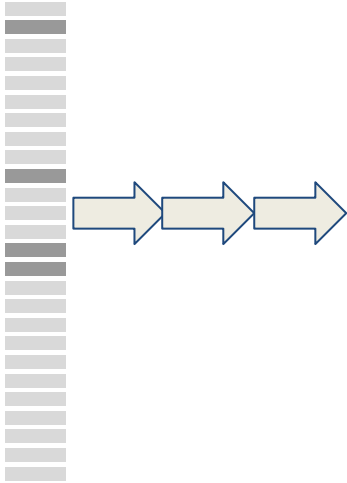
Identificar Composición Microbiana

- ☐ Explorar Reporte de Muestra interactivo
- ☐ Usar BLAST para confirmar coincidencias
- ☐ Generar genomas de consenso vira

Comparar Múltiples Muestras

- ☐ Crear heatmaps interactivos para visualizar taxa abundantes en varias muestras
- ☐ Construir árboles filogenéticos y colorear ramas por metadata

Raw reads



Medical Detectives

Patient 010 (CSF) ▾

Metagenomic Antimicrobial Resistance (Deprecated) Consensus Genome

Taxon name Name Type: Scientific ▾ Background: MedicalDetectives_WaterCtrls ▾ Categories ▾ Threshold filters: 1 ▾ Read Specificity: Specific Only ▾ Annotation ▾

NT rPM >= 10 ✕

112 rows passing the above filters, out of 3100 total rows [CLEAR FILTERS](#)

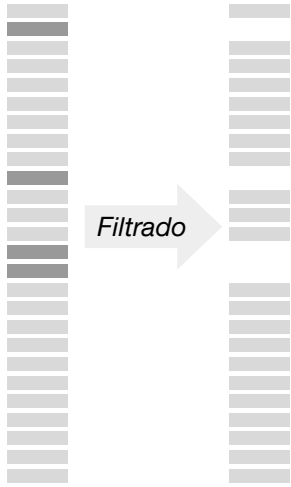
	Taxon	Score ▾	Z ...	rPM	r	con...	con...	%...	L	E v...	NT NR
▾	Alphavirus (1 viral species: ● 1)	72,976,252	100.0	3,650.9	19,344	1	19,139	100.0	11,671.6	10 ⁻³⁰⁵	
	Chikungunya virus Known Pathogen	72,976,252	100.0	3,646.7	19,344	1	19,139	99.9	2,450.2	10 ⁻³⁰⁵	
>	Tetraparvovirus (1 viral species:)	46,706,236	100.0	2,366.2	12,537	2	12,102	99.1	4,677.3	10 ⁻²⁹⁶	
	Providencia (1 bacterial species: ● 1)	14,332,531	100.0	1,433.3	7,594	1	7,594	99.8	1,354.0	10 ⁻³⁰⁸	
	Lactobacillus (1 bacterial species:)	2,502,625	100.0	28.9	153	7	79	98.6	209.0	10 ⁻¹⁰⁸	
>	Rothia (1 bacterial species:)	2,364,849	100.0	238.4	1,263	6	1,127	96.5	2,086.8	10 ⁻²⁷²	

Objetivo Pipeline

- **Asignar lecturas** filtradas de calidad a **taxones** mediante mapeo o alineación con bases de datos de referencia
- **Ensamblar** contigs y **mapear** lecturas a contigs **para aumentar el soporte para taxones** identificados (alineación basada en ensamblaje)

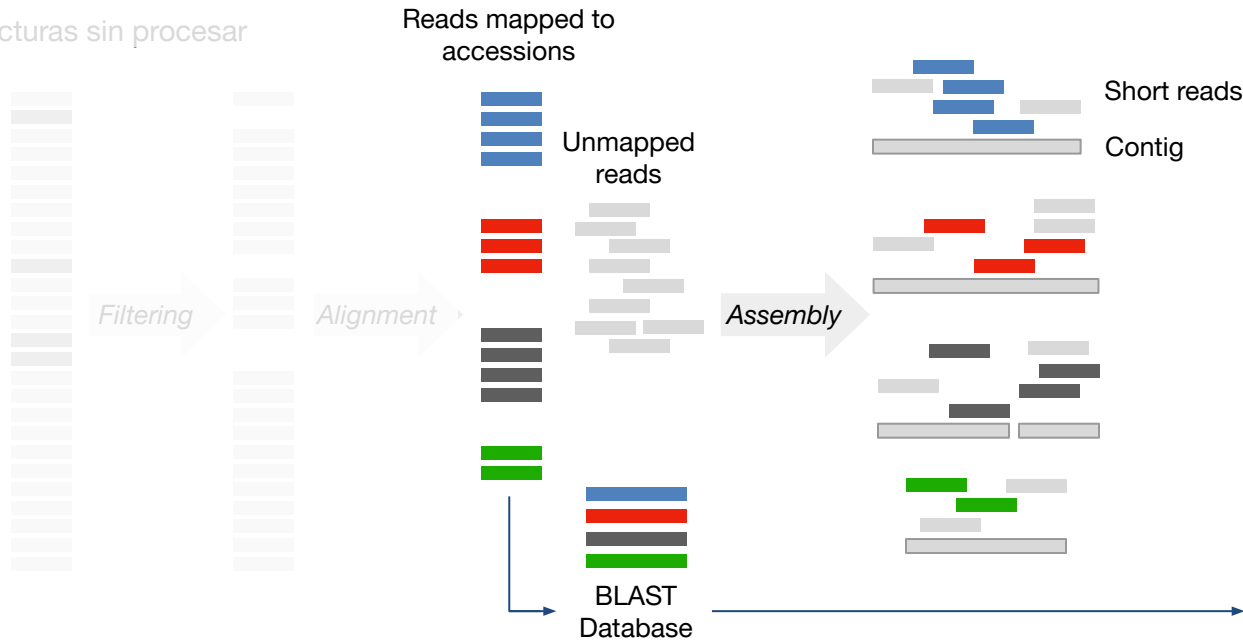
Lecturas sin procesar





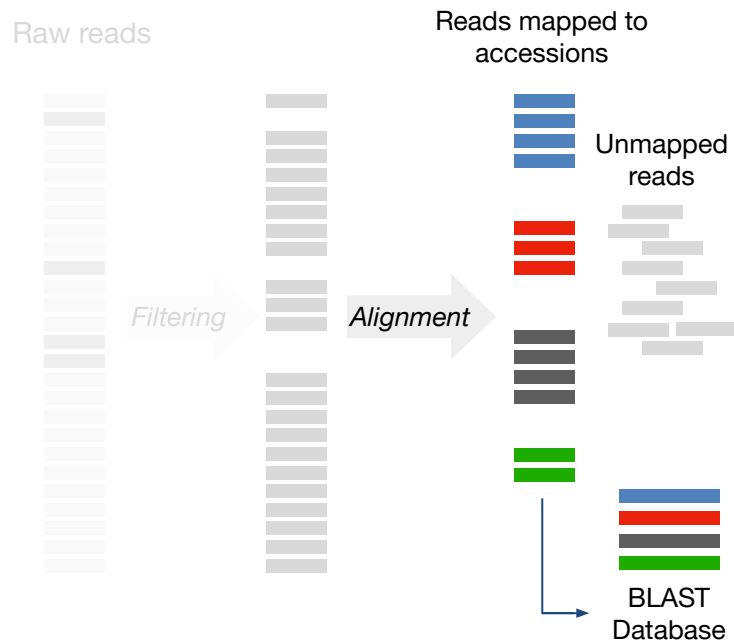
Paso 1: Control de calidad y filtrado del host

- Filtrar lecturas de **baja calidad/ complejidad y adaptadores** (fastp)
- **Eliminación** de lecturas del **host** (Bowtie seguido de HISAT2)
- **Eliminación** de lecturas **humanas** (Bowtie seguido de HISAT2)
- Contraer **lecturas duplicadas** (CZID-dedup)
- **Submuestreo** a 2 millones de lecturas o 1 millón de pares de lecturas



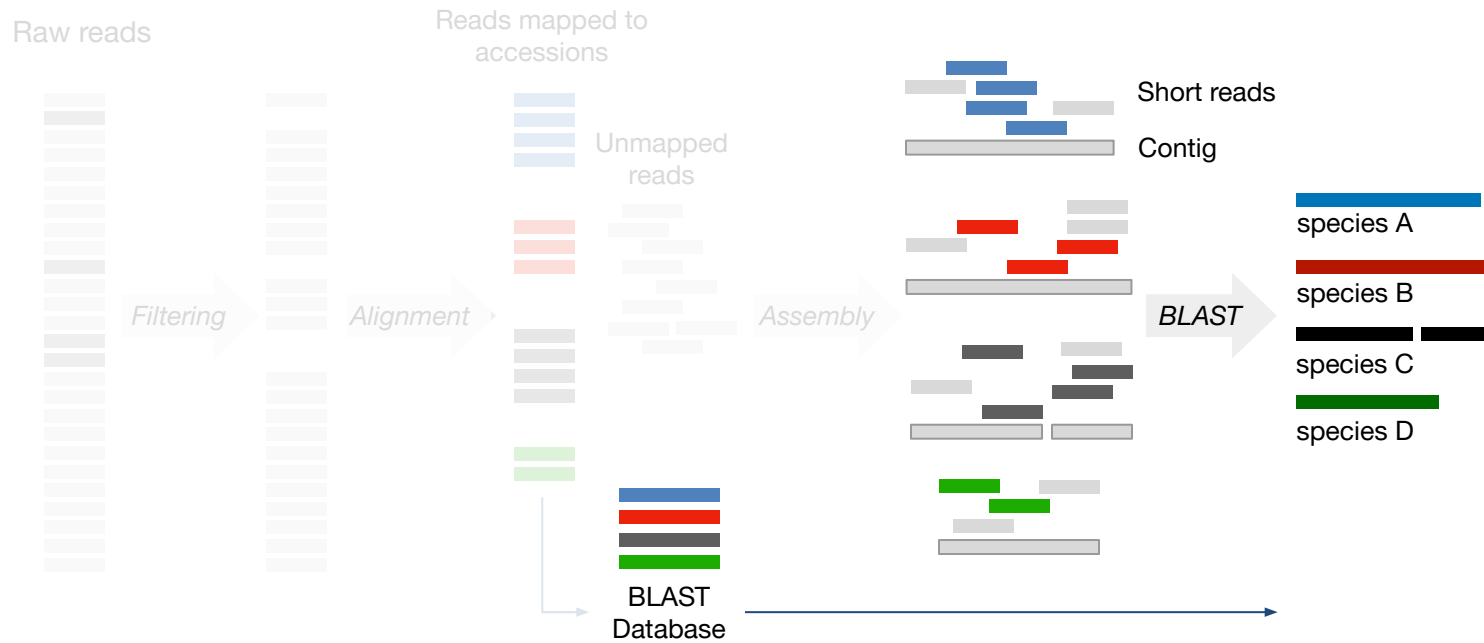
Paso 3: ensamblar contigs y asignar lecturas a contigs

- Ensamble las lecturas en contigs mediante **ensamblaje de novo** (SPAdes).
- Alinee las lecturas con los contigs ensamblados (Bowtie2).



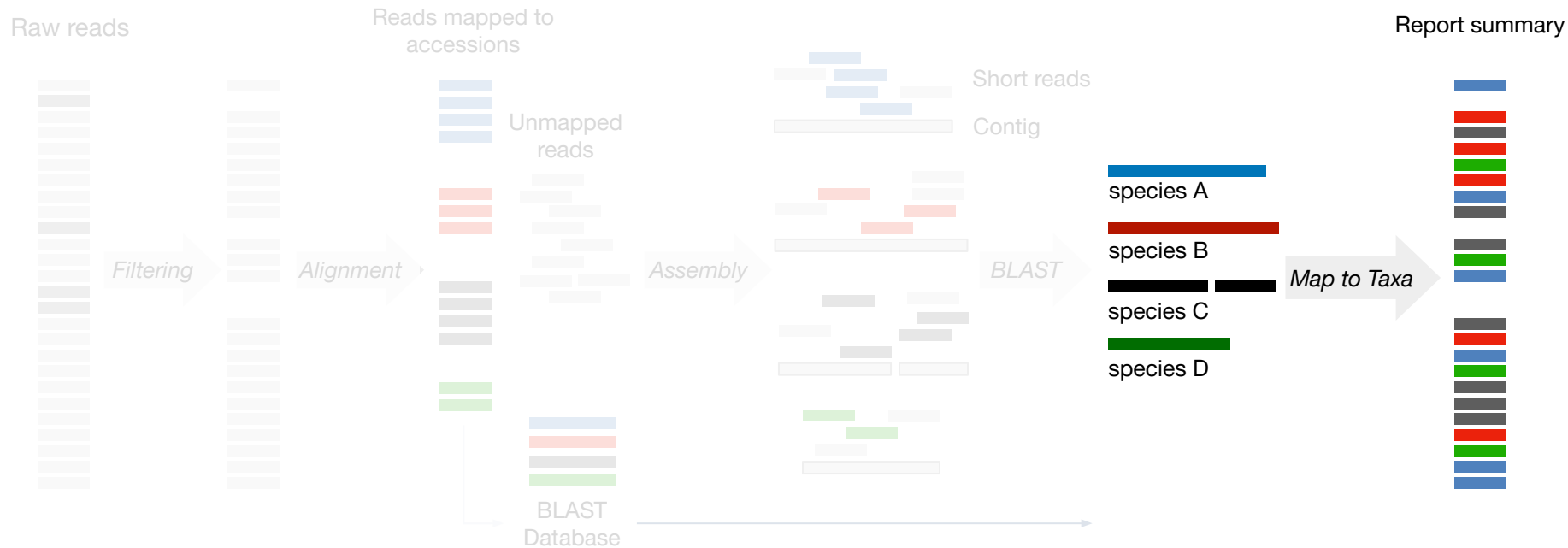
Paso 2: Alineamiento de lecturas con las bases de datos del NCBI

- Mapeo de las lecturas en bd de nucleótidos (**NT**) y proteínas (**NR**) para obtener una asignación preliminar para cada lectura (**minimap2** y **DIAMOND**).
- Genere una base de datos BLAST personalizada con todos los taxones identificados con asignaciones preliminares.



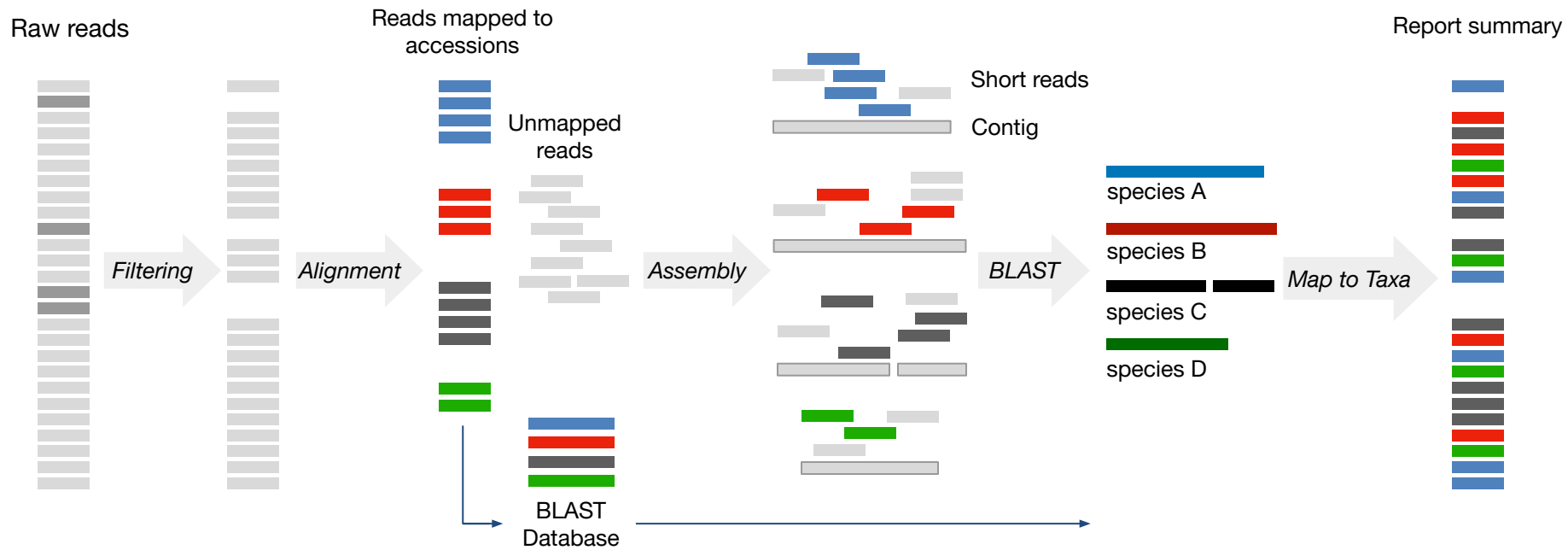
Paso 4: Alinear los contigs con la base de datos personalizada

- Se alinea los contigs en la bd generada con las asignaciones preliminares (BLAST)



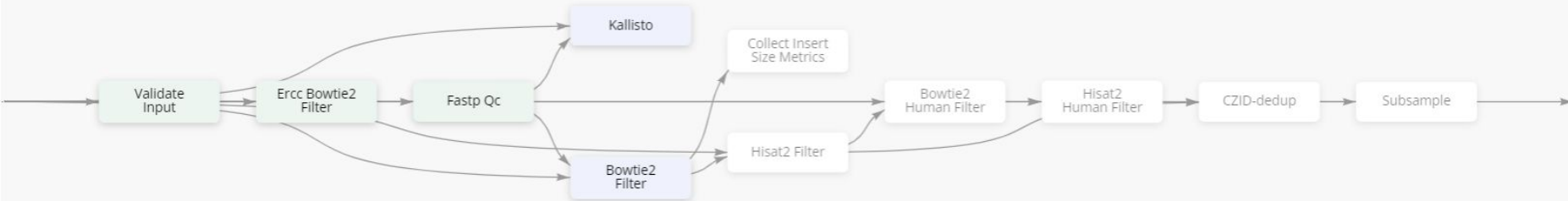
Paso 5: Mapeo a taxones

- **Asigna a cada lectura una accesión final** en función de:
 1. Accesión de contig, o
 2. Accesión lectura no se ensambló en un contig
- **Calcular estadísticas** (p. ej., lecturas por millón para cada taxón)

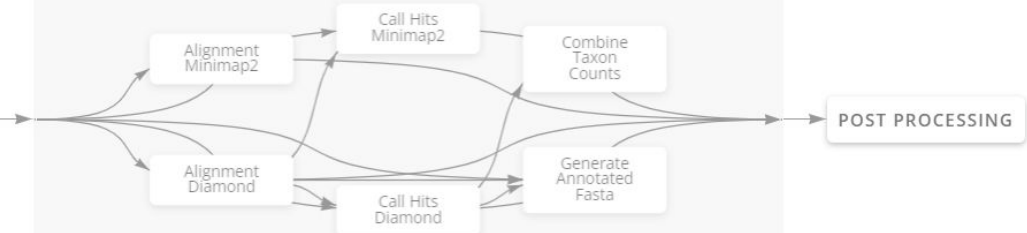


Pipeline

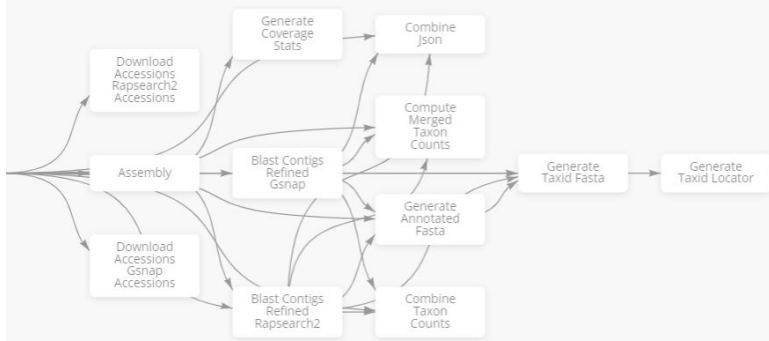
HOST FILTERING



ALIGNMENT



POST PROCESSING



Carga de Datos

Para cargar nuevas muestras en CZ ID, [inicie sesión](#) en la aplicación con su correo electrónico y contraseña.



Haga clic [aquí](#) para obtener instrucciones paso a paso sobre cómo cargar datos en CZ ID para el análisis mNGS.





Informe de la Muestra: Métricas y análisis

La muestra del informe presenta métricas para ayudar a identificar taxones de interés.

NT = DB nucleótidos
NR = DB proteínas

<i>Si yo veo:</i>	Bacteria	Virus	Eukaryote
Solo NT	Región Conservada (16s)	?	Región Conservada (18s)
Solo NR	No es un hit Real	Podría ser un virus divergente	No es un hit Real
Ambos	Hit real	Hit real	Hit real

¿Búsqueda más estricta de secuencias similares?

¿Qué pasa con las secuencias no codificantes?

La muestra del informe presenta métricas para ayudar a identificar taxones de interés.

NT = nucleotide database
NR = protein database

<i>Si yo veo:</i>	Bacteria	Virus	Eukaryote
Solo NT	Podría estar afectando la región 16s	Potencialmente afectando la región no codificante	Región Conservada (18s)
Solo NR	No es un hit Real	Podría ser un virus divergente	No es un hit Real
Ambos	Hit real	Hit real	Hit real

Identificar patógenos conocidos, debe centrarse en los valores de las coincidencias en la base de datos NT. Por el contrario, si sospecha de un patógeno nuevo, debe centrarse en los valores de la base de datos NR.

Añadir métricas complementarias a cada taxón

rPM	Por cada millón de lecturas secuenciadas nos indica el número de lecturas que se alinean con el taxón	Alto rPM = Más confianza en el resultado
r	Número de reads alineado por taxón	Alto r = Más confianza en el resultado
Contig	Número de contigs ensamblados	contig largos = Más confianza
Contig r	Número de Lecturas que se alinean a contigs ensamblados	Alto r = Más confianza
L	Longitud promedio de alineación de contigs y lecturas	Largo= Más confianza, unless it is divergent
E value	Significativo: el valor medio esperado de alineamiento	Bajo = Más confianza en el resultado

Los filtros son útiles para eliminar resultados falsos, pero no son perfectos

Usar los siguientes filtros es un buen punto de partida:

- **NT L \geq 50 [eliminar taxones con longitudes de alineamiento cortas]**
 - Solo quiero ver taxones con al menos 50 pb que se alineen con la base de datos de nucleótidos en NCBI.
- **NT rPM \geq 10 [eliminar taxones con baja abundancia]**
 - Solo quiero ver taxones con al menos 10 lecturas por millón que se alineen con la base de datos de nucleótidos. Las lecturas por millón son un valor normalizado que tiene en cuenta el total de lecturas.
- **NR rPM \geq 1 [eliminar taxones que solo tienen coincidencias con regiones no codificantes, por ejemplo, 16S/18S]**
 - Solo quiero ver taxones con al menos 1 lectura que se alinee con la base de datos de proteínas en NCBI.



Análisis de mNGS con una y múltiples muestras



mNGS para estudiar patógenos transmitidos por vectores

Actividad: Análisis de una muestra



Culex erythrothorax

- ❑ Un solo individuo
- ❑ Secuenciación: Illumina
- ❑ Este mosquito se alimentó de sangre

**¡Veamos qué taxones
podemos encontrar!**



Tiempo: ~ 30 minutos para la actividad y 10 minutos para la discusión

https://github.com/instituto-nacional-de-salud/Curso-internacional-de-epidemiologia-genomica-2024/blob/main/Sesiones%20Pr%C3%A1cticas/D%C3%ADa%203/mNGS/Actividad_1_%20Pat%C3%B3genos%20transmitidos%20por%20vectores.docx

Paso a Paso

cz ID

My DataPublicUpload

0

Projects680Samples26356

Taxon Filter

SAMPLE

CMS_001_RNA_A_S1

Project: Mosquito proj PHA4GE 2023

TAXON

680 projects

DESCRIPTION

No description

1. Acceder: [CMS 001 RNA A S1](#)

Threshold filters: 3

NT rPM

<=

10

X

NR rPM

>=

1

X

NR L (alignment length in ...)

>=

50

X

+ ADD THRESHOLD

Cancel

Apply

- 1. ¿Cuántas filas de taxones están presentes ahora?
- 2. ¿Cuál es el taxón más abundante según las lecturas por millón (rPM)?

¿Qué pasa con los virus divergentes?

Las mutaciones se acumulan más rápido a nivel de nucleótidos.

NR rPM ≥ 20 [Solo quiero ver taxones con 20 o más lecturas alineadas con la base de datos de proteínas en NCBI]

NR % ID ≤ 80 [Solo quiero ver taxones con menos del 80% de identidad con taxones en la base de datos de proteínas en NCBI]

1. ¿Cuántos géneros hay? ¿Cuáles son?

A que no sabías que ...

1. Este conjunto de datos proviene del artículo "[Single mosquito metatranscriptomics identifies vectors, emerging pathogens and reservoirs in one assay](#)".

LOCUS QRW41732 399 aa linear ENV 19-MAY-2021
DEFINITION MAG: nucleoprotein [Miglotas virus].
ACCESSION QRW41732
VERSION QRW41732.1
DBLINK BioProject: [PRJNA605178](#)
BioSample: [SAMN14051441](#)
Sequence Read Archive: [SRR11035377](#)
DBSOURCE accession [MW434677.1](#)
KEYWORDS ENV; Metagenome Assembled Genome; MAG.
SOURCE Miglotas virus (mosquito metagenome)
ORGANISM [Miglotas virus](#)
Viruses; Riboviria; Orthornavirae; Negarnaviricota;
Polyploviricotina; Ellioviricetes; Bunyavirales; Phasmaviridae;
Orthophasmavirus; Orthophasmavirus miglotasense.

nuevo virus que no ha sido incluido en la base de datos de NCBI ...

Pariente más cercano:
Porcentaje de identidad (38%) con el culex orthophasmavirus.

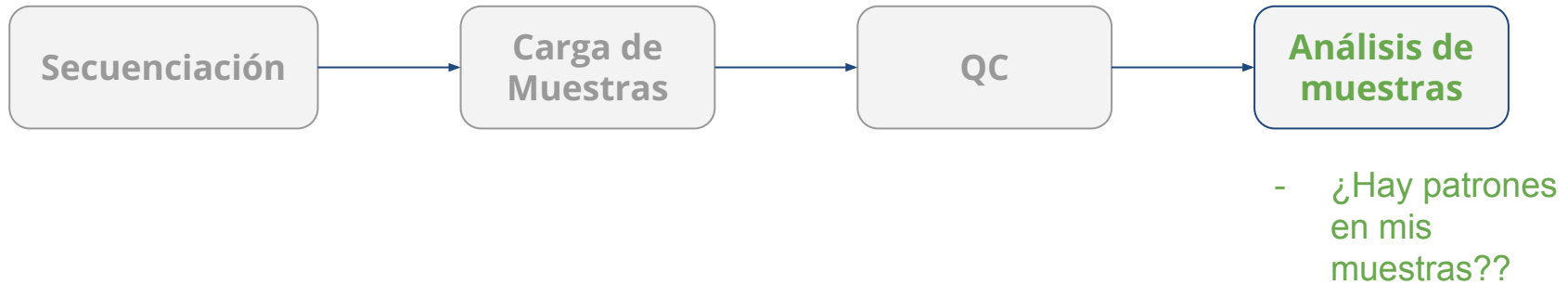
	Description	Max Score	Total Score	Query Cover	E value	Per. Ident	Accession
<input type="checkbox"/>	nucleoprotein [Miglotas virus]	826	826	57%	0.0	100.00%	QRW41732.1
<input type="checkbox"/>	nucleoprotein [Miglotas virus]	825	825	57%	0.0	99.75%	QRW41744.1
<input type="checkbox"/>	nucleoprotein [Miglotas virus]	824	824	57%	0.0	99.75%	YP_010840373.1
<input type="checkbox"/>	hypothetical protein QK753_s3gp1 [Miglotas virus]	269	269	18%	6e-83	100.00%	YP_010840372.1
<input type="checkbox"/>	hypothetical protein [Miglotas virus]	267	267	18%	3e-82	99.23%	QRW41755.1
<input type="checkbox"/>	hypothetical protein [Miglotas virus]	267	267	18%	4e-82	99.23%	QRW41737.1
<input type="checkbox"/>	nucleocapsid [Aedes japonicus bunyavirus 2]	261	261	50%	2e-76	41.29%	WCJ14327.1
<input type="checkbox"/>	nucleocapsid protein [Culex phasma-like virus]	255	255	53%	1e-73	38.74%	QHA33851.1

- ☐ Métricas para el análisis
- ☐ Identificación de hits
- ☐ Utilizar filtros

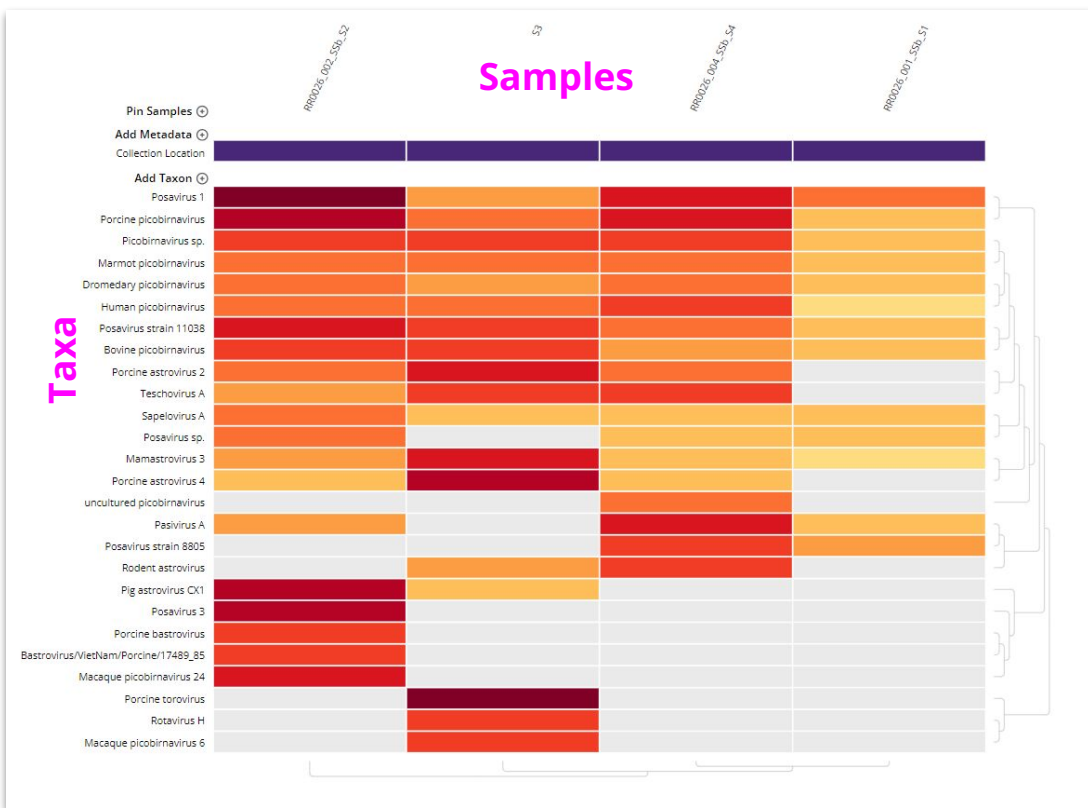
Análisis de mNGS múltiples muestras – Heatmap



Análisis de cohortes



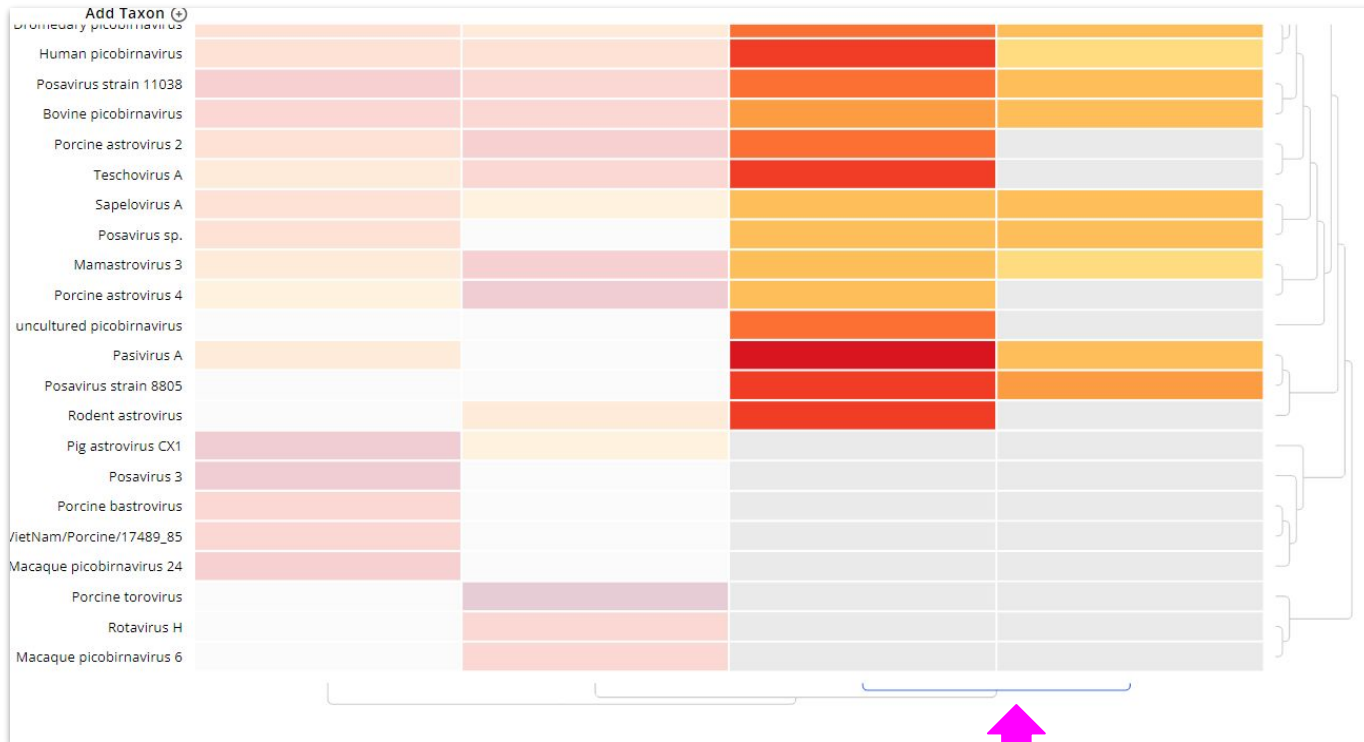
Interpretando un heatmap



Less
abundant

More
abundant

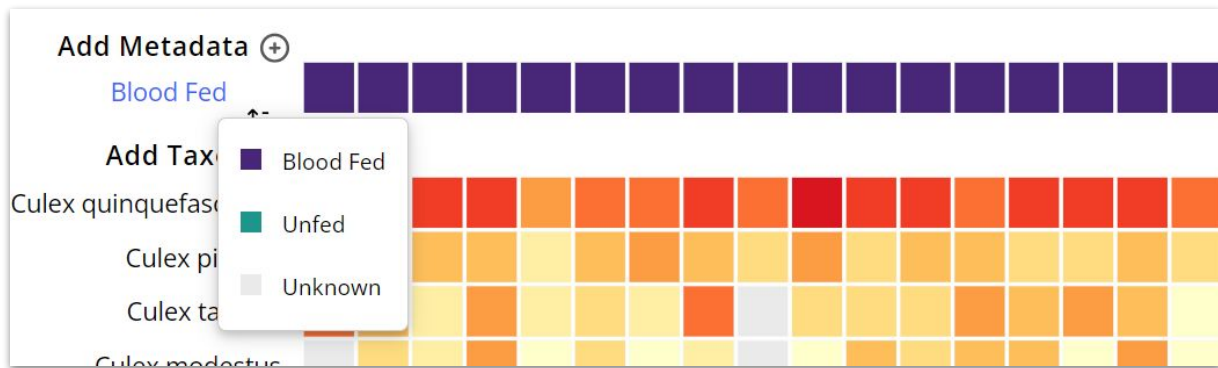
Las muestras se organizan automáticamente mediante agrupamiento jerárquico.



Agrupamientos de taxones en el dendrograma: los taxones que tienen más probabilidades de aparecer juntos en abundancias similares se agrupan

Muestras de grupos de dendrogramas: las muestras con taxones similares en abundancias similares se agrupan

Sort samples by metadata



Es posible que desee organizar sus muestras por diferentes campos de metadatos en lugar de agruparlas jerárquicamente.

Puede agregar metadatos al mapa de calor y ordenarlos. En este ejemplo, compararemos los taxones presentes en mosquitos alimentados con sangre y mosquitos no alimentados.

Aplicar filtros

Al igual que en el informe de una muestra, deberá aplicar filtros para eliminar los resultados falsos.

También puede comparar los controles de agua con las muestras para ver los contaminantes antes de agregar un filtro de puntuación **Z-score filter**.

The interface displays several configuration options for data filtering:

- Taxon Level: Species
- Categories
- Sort Taxa: Cluster
- Sort Samples: Cluster
- Metric: NT rPM
- Background: NID human CSF HC
- Threshold filters
- Read Specificity: Spec...
- Scale: Log
- Taxa per Sample: 10
- Color scale: 0 to 34k

A dialog box for adding a threshold filter is open, showing:

- NT rPM
- Operator: >=
- Value: 10
- + ADD THRESHOLD
- Buttons: Cancel, Apply

Below the dialog, a list of sample identifiers is visible, including:

- 1_S4
- 1_S4
- 1_S14
- 1_S6
- 1_S16
- 1_S18
- 1_S20
- 1_S21
- 1_S21
- 1_S3
- 1_S4
- 1_S6
- 1_S9
- 1_S10
- b_S124_L004
- b_S125_L004
- b_S126_L004
- b_S128_L004
- b_S129_L004
- b_S130_L004
- b_S131_L004
- b_S132_L004
- b_S133_L004
- b_S134_L004
- b_S135_L004
- b_S136_L004
- b_S137_L004
- b_S138_L004
- b_S185_L004
- b_S9_L004

¿Qué pregunta puedo responder con este diseño experimental??

- A. ¿Qué patógenos transmitidos por vectores están circulando en el ecosistema?
- B. ¿Hay virus nuevos presentes y cuáles son?
- C. Todas las anteriores
- D. Ninguna de las anteriores

Actividad: Análisis de Múltiples Muestras



Culex erythrothorax

- ❑ 33 individuos
- ❑ Secuenciación: Illumina
- ❑ 17 alimentados con sangre
- ❑ 2 controles de agua

¡Hagamos un análisis de múltiples muestras utilizando un Heatmap!



Tiempo: ~ 30 minutos para la actividad y 10 minutos para la discusión

[ENLACE GITHUB DONDE ESTÁ EL ARCHIVO.](#)

Publicaciones Relevantes

- ❑ [Single mosquito metatranscriptomics identifies vectors, emerging pathogens and reservoirs in one assay](#)
- ❑ [IDseq—An open source cloud-based pipeline and analysis service for metagenomic pathogen detection and monitoring](#)
- ❑ [Discovering disease-causing pathogens in resource-scarce Southeast Asia using a global metagenomic pathogen monitoring system](#)
- ❑ [Case Report: Cambodian National Malaria Surveillance Program Detection of Plasmodium knowlesi](#)
- ❑ [The Perfect Storm of 2019: An immunological and phylodynamic analysis of Cambodia's unprecedented dengue outbreak](#)
- ❑ [Metagenomic Pathogen Sequencing in Resource-Scarce Settings: Lessons Learned and the Road Ahead](#)
- ❑ [Pathogen genomics in public health](#)
- ❑ [IDseq—An open source cloud-based pipeline and analysis service for metagenomic pathogen detection and monitoring](#)
- ❑ [Pathogen genomics in public health](#)
- ❑ [Unbiased metagenomic sequencing for pediatric meningitis in Bangladesh reveals neuroinvasive Chikungunya virus outbreak and other unrealized pathogens](#)

Gracias

Contacto:

crojas@ins.gov.co

paolarojasestevez@gmail.com

Algunas de las diapositivas usadas en esta presentación hacen parte del material disponible en Help Center de CZID

<https://chanzuckerberg.zendesk.com/hc/en-us>

The coverage visualization is valuable in calling hits

Alphapapillomavirus 7 Coverage

[View read-level visualization >](#)

MF288662.1 - Human papillomavirus type 18 isolate 1359486_N-P, c... ▾

10 viewable accessions (46 total) ⓘ



Reference NCBI Entry [MF288662.1 - Human p...](#)

Aligned Contigs 1

Coverage Depth 936.8x

Max Alignment Length 7857

Reference Length 7857

Aligned Loose Reads 0

Coverage Breadth 100.0%

Avg. Mismatched % 0.0%



Reference Accession

Contigs (1)

The coverage visualization is valuable in calling hits

View the accession

Alphapapillomavirus 7 Coverage

MF288662.1 - Human papillomavirus type 18 isolate 1359486_N-P, c... ▾

10 viewable accessions (46 total) ⓘ

[View read-level visualization >](#)



Reference NCBI Entry [MF288662.1 - Human p...](#)

Aligned Contigs 1

Coverage Depth 936.8x

Max Alignment Length 7857

Reference Length 7857

Aligned Loose Reads 0

Coverage Breadth 100.0%

Avg. Mismatched % 0.0%



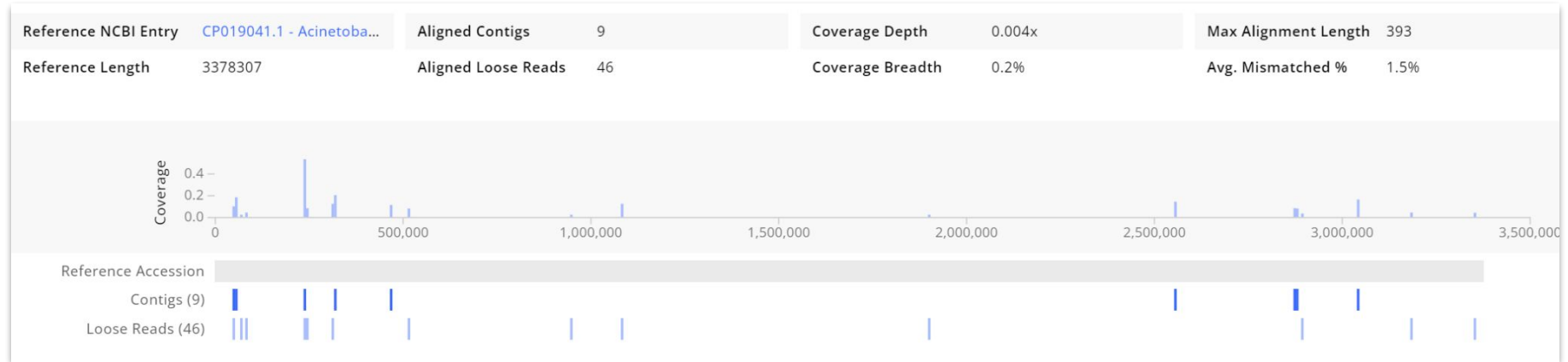
Reference Accession

Contigs (1)

Pay attention to
accession length

Take into consideration the depth and breadth of the coverage

- This is great coverage of bacterium
- Typically see less coverage, you will look for multiple peaks across the genome



This is an example of poor coverage. I would not be comfortable calling this hit real

Providencia rettgeri Coverage

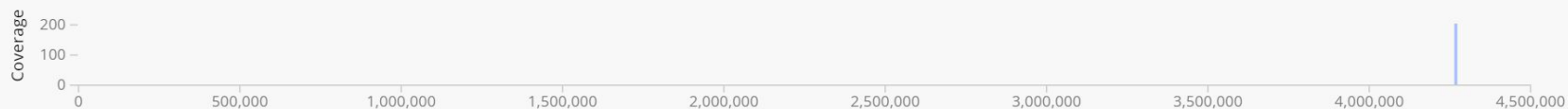
[View read-level visualization >](#)

CP029736.1 - *Providencia rettgeri* strain AR_0082 chromosome, co... ▼

1 viewable accessions



Reference NCBI Entry	CP029736.1 - Providenc...	Aligned Contigs	1	Coverage Depth	0.4x	Max Alignment Length	1354
Reference Length	4309166	Aligned Loose Reads	0	Coverage Breadth	0.0%	Avg. Mismatched %	0.2%

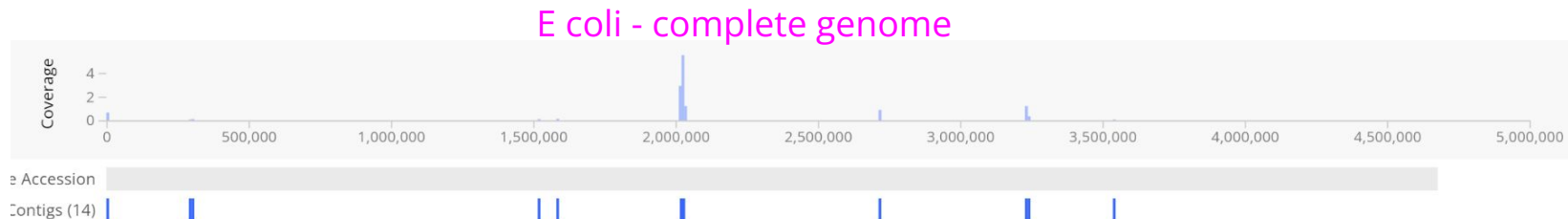


Reference Accession

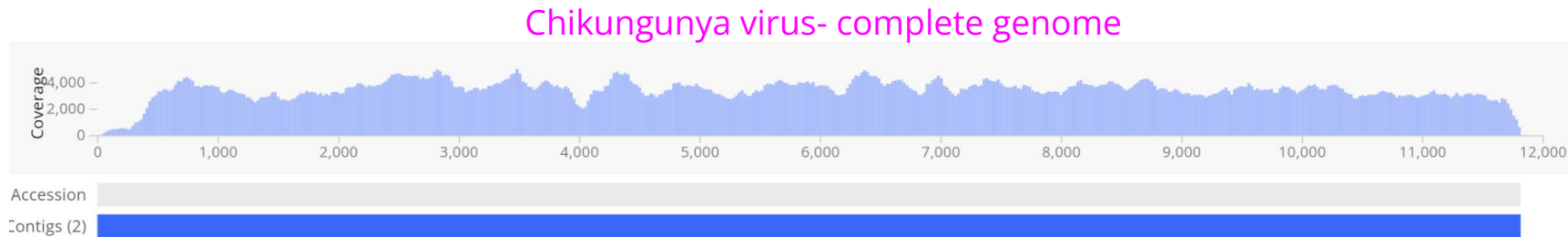
Contigs (1)

Which would you not consider a true hit?

A



B



C

