



**Fundamentos y
herramientas
bioinformáticas
para análisis
genómicos**



Práctico ensamblaje y anotación genómica

Luisa Berná, PhD
lberna@pasteur.edu.uy

Unidad de Bioinformática - Laboratorio de Biología de Apicomplejos
Institut Pasteur de Montevideo

Laboratorio de Genómica Evolutiva - Facultad de Ciencias, UDELAR.

Objetivos

- ❖ Aprender a ensamblar genomas utilizando distintos tipos de secuencias y estrategias
- ❖ Afianzar conceptos sobre las distintas tecnologías de secuenciación
- ❖ Afianzar conceptos sobre algoritmos de ensamblaje
- ❖ Aprender a realizar anotación automática, y entender sus virtudes y limitaciones.

Cómo?

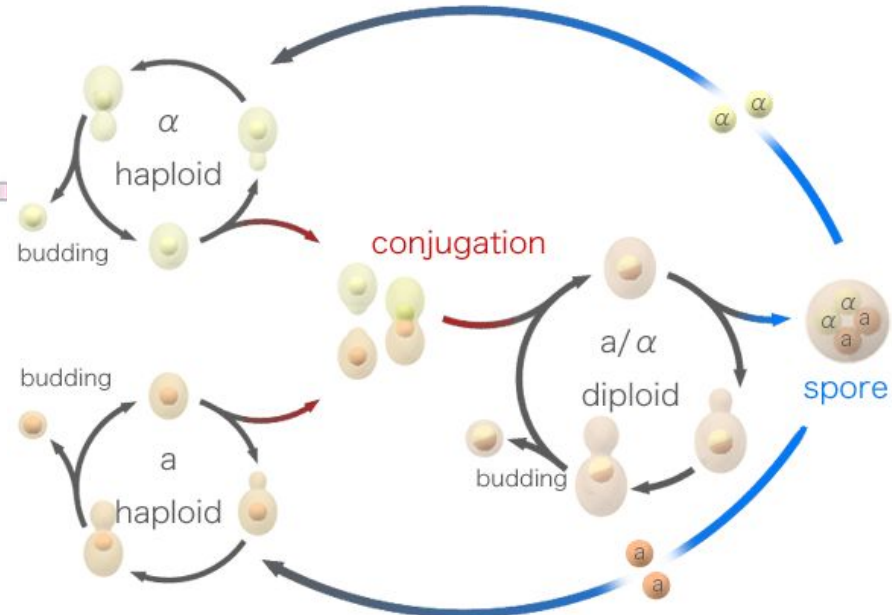
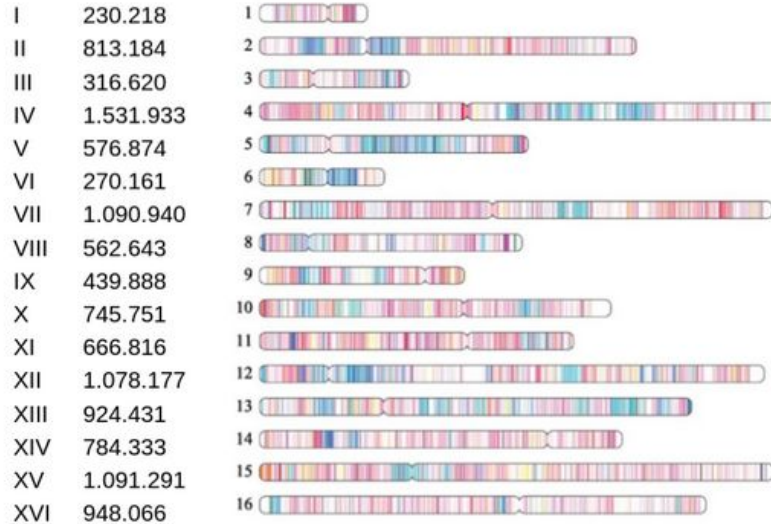
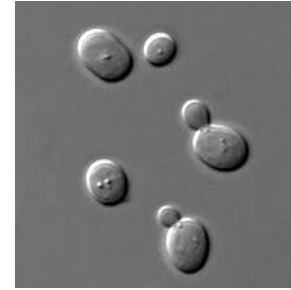
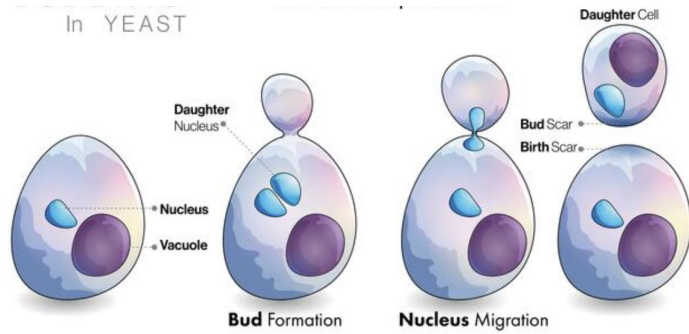
- ❖ Ensamblando el genoma de un organismo modelo (levadura) para el cual existen diferentes tipos de datos
 - Illumina (paired-end)
 - Nanopore
- ❖ Haciendo una predicción de genes
- ❖ A través de la práctica de distintas herramientas y la comparación y discusión de los resultados obtenidos

Modelo de estudio

Para este práctico se utilizará la levadura (*Saccharomyces cerevisiae*). Este organismo tiene varias ventajas que lo hacen interesante para realizar un práctico de ensamblaje y anotación:

- ★ Ampliamente conocido e interesante (pan/vino/cerveza)
- ★ Eucariota (complejidad intermedia)
- ★ Tamaño genómico reducido (12Mb)
- ★ 16 Cromosomas
- ★ Genes con Intrones (~10%)
- ★ EXISTEN datos disponibles de varias plataformas (illumina/PacBio/Nanopore)

In YEAST



Complete genomic and transcriptional landscape analysis using third-generation sequencing: a case study of *Saccharomyces cerevisiae* CEN.PK113-7D

Piroon Jenjaroenpun^{1,†}, Thidathip Wongsurawat^{1,†}, Rui Pereira²,
Preecha Patumcharoenpol¹, David W. Ussery^{1,3}, Jens Nielsen^{2,4} and Intawat Nookaew^{1,2,3,*}

¹Department of Biomedical Informatics, College of Medicine, University of Arkansas for Medical Sciences, Little Rock, AR 72205, USA, ²Department of Biology and Biological Engineering, Chalmers University of Technology, Gothenburg SE-412 96, Sweden, ³Department of Physiology and Biophysics, College of Medicine, The University of Arkansas for Medical Sciences, Little Rock, AR 72205, USA and ⁴Novo Nordisk Foundation Center for Biosustainability, Technical University of Denmark, DK2800 Lyngby, Denmark

Received September 11, 2017; Revised January 03, 2018; Editorial Decision January 04, 2018; Accepted January 05, 2018

1 - Pre-procesamiento de secuencias

- Descargar datos // Obtener datos a partir de secuenciaciones
- Analizar la calidad de las secuencias
- Filtrar si es necesario
- Re-analizar la calidad para determinar si los filtros utilizados fueron adecuados

2 - Ensamblaje DE NOVO

Vamos a ensamblar *de novo* un genoma de *Saccharomyces cerevisiae* usando reads obtenidos a partir de tecnologías diferentes:

- Illumina
- Nanopore ONT

Vamos a utilizar distintos programas y comparar sus resultados

3- Evaluación de los ensamblajes

- + **infoseq**
- + **YASS**
- + **BLAST**
- + **Assemblytics**
- + **BUSCO**

4- Búsqueda de genes

- AUGUSTUS

Augustus [gene prediction]

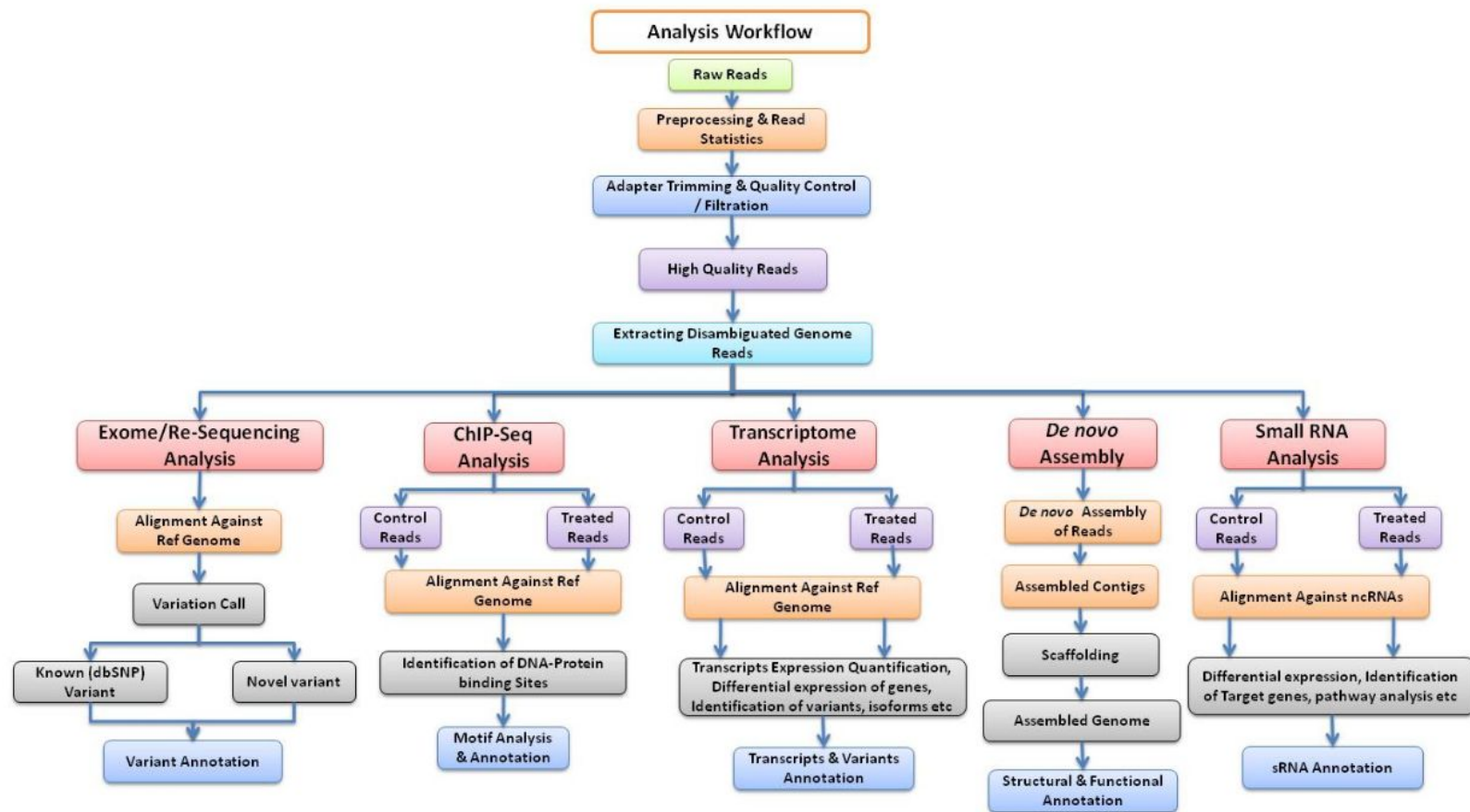
Bioinformatics Group of the Institute for Mathematics and Computer Science of the University of Greifswald

[web interface](#) [WebAUGUSTUS](#) [accuracy results](#) [download AUGUSTUS](#) [data sets](#) [predictions](#) [references](#)



AUGUSTUS is a program that predicts genes in eukaryotic genomic sequences. It can be [run on this web server](#), on a new web server for larger input files or be downloaded and run [locally](#). It is open source so you can compile it for your computing platform. You can now run AUGUSTUS on the German [MediGRID](#). This enables you to submit larger sequence files and allows to use protein homology information in the prediction. The MediGRID requires an instant easy registration by email for first-time users.

Formatos



Algunas definiciones

- **Single read sequencing** = DNA fragments read from one end.



- **Paired-end sequencing** = DNA fragments read from both ends. Reads are typically up to a few hundred bp apart, but may also overlap.



- **Mate pair sequencing** – the ends of long DNA fragments (typically several kb apart) are fused and then excised for sequencing.



- **Long reads** – Much longer reads! Can thus improve de novo assembly, mapping certainty, transcript isoform identification, and detection of structural variants

Formatos

- + FASTA
- + FASTQ
 - + Calidad PHRED

FASTA

Un formato de texto plano utilizado para representar secuencias de nucleótidos o proteínas. Los nucleótidos se representan con las letras A,C,G,T y los aminoácidos con el código de una letra:

- + **Encabezado**

descripción en una única línea que comienza con el símbolo '>'. No hay espacios entre el '>' y el comienzo de la descripción. Puede contener un simple identificador o una descripción más compleja y comentarios

- + **Cuerpo**

la propia secuencia representada por la sucesión de letras. Pueden ir todas en una línea o en bloques (en general de 80 caracteres)

FASTA

```
>seq1
```

```
GATCAACGCAAAGGACTAAGCACTGCTGCCAAAAGCCACCAGCCCCAGAGACAACAGAGG  
CTCCCAAATTTCTAGCCTCTGATCTCTGCCTCGGAACATTCTTGGGTCAAATAAATGTG  
CGATCGCTAGCTAAAACGTTTCGAT
```

```
>gi|5524211|gb|AAD44166.1| cytochrome b [Elephas maximus maximus]  
LCLYTHIGRNIYYGSYLYSETWNTGIMLLITMATAFMGYVLPWGQMSFWGATVITNLFS  
AIPYIGTNLV  
EWIWGGFSVDKATLNRFFAFHFILPFTMVALAGVHLTFLHETGSNNPLGLTSDSDKIPF  
HPYYTIKDFLG  
LLILILLLLLLLALLSPDMLGDPDNHMPADPLNTPLHIKPEWYFLFAYAILRSVPNKLGG  
VLALFLSIVIL  
GLMPFLHTSKHRSMMLRPLSQALFWTLTMDLLTLTWIGSQPVEYPYTIIGQMASILYFS  
IILAFPLIAGX  
IENY
```


multiFASTA

- + Un archivo puede contener más de una secuencia. Cada secuencia comienza con su encabezado y tiene su cuerpo
- + Un multifasta puede tener miles y miles de secuencias (por ejemplo uno con todas las proteínas de un genoma)

```
>SEQUENCE_1
MTEITAAMVKELRESTGAGMMDCKNALSETNGDFDKAVQLLREKGLGKAACKADRLAAEG
LVSVKVSDDFTIAAMRPSYLSYEDLDMTFVENEYKALVAELEKENEERRRLKDPNKPEHK
IPQFASRKQLSDAILKEAEEKIKEELKAQGKPEKIWDNIIPGKMNSFIADNSQLDSKLT
MGQFYVMDDKKTVEQVIAEKEKEFGGKIKIVEFICFEVGEGLKKTEDFAAEVAAQL
>SEQUENCE_2
SATVSEINSETDFVAKNDQFIALTKDTHAHIQSNLSQSVEELHSSTINGVKFEEYLSQI
ATIGENLVVRRFATLKAGANGVVNGYIHTNGRVGVVIAAACDSAEVASKSRDLLRQICMH
```

FASTQ Fasta Quality

- + Es el archivo que obtenemos de las distintas plataformas de secuenciación
- Posee información de secuencia (como el archivo FASTA:

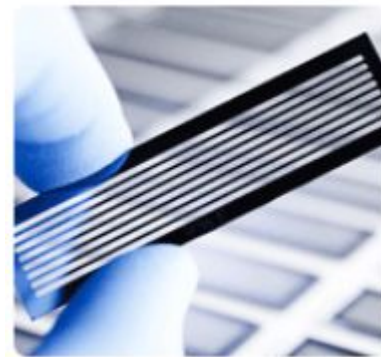
[illegible]

ezado '+' y

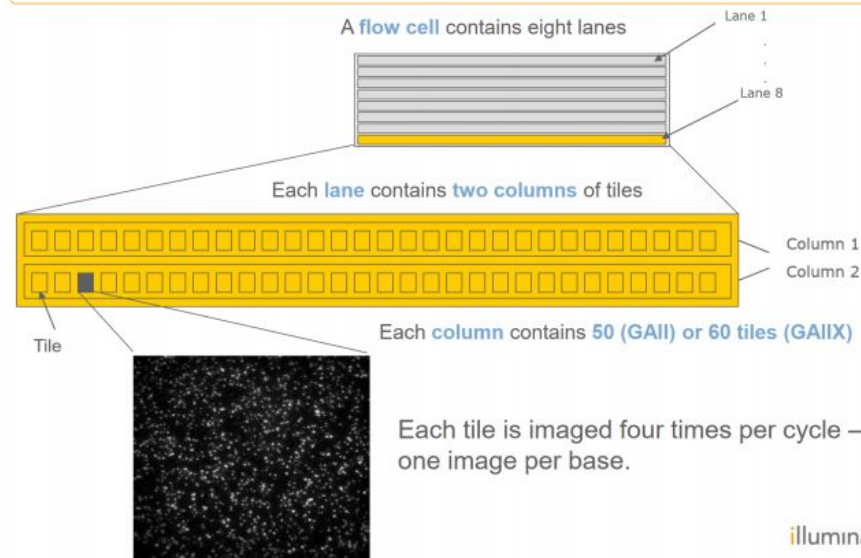
FASTQ encabezado

```
@EAS139:136:FC706VJ:2:2104:15343:197393 1:Y:18:ATCACG
```

EAS139	the unique instrument name
136	the run id
FC706VJ	the flowcell id
2	flowcell lane
2104	tile number within the flowcell lane
15343	'x'-coordinate of the cluster within the tile
197393	'y'-coordinate of the cluster within the tile
1	the member of a pair, 1 or 2 (<i>paired-end or mate-pair reads only</i>)



Flow Cell Images



FASTQ Calidad

PHRED SCORE es una medida de calidad **Q**

P es la probabilidad de que un nucleótido haya sido asignado incorrecto:

$$Q = -10 \log_{10} P \quad P = 10^{\frac{-Q}{10}}$$

Phred quality scores are logarithmically linked to error probabilities

Phred Quality Score	Probability of incorrect base call	Base call accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1000	99.9%
40	1 in 10,000	99.99%
50	1 in 100,000	99.999%
60	1 in 1,000,000	99.9999%

FASTQ Calidad - código ASCII

[illegible]

FASTQ Calidad - código ASCII + 33

Caracteres ASCII imprimibles			
32	espacio	64	@
33	!	65	A
34	"	66	B
35	#	67	C
36	\$	68	D
37	%	69	E
38	&	70	F
39	'	71	G
40	(72	H
41)	73	I
42	*	74	J
43	+	75	K
44	,	76	L
45	-	77	M
46	.	78	N
47	/	79	O
48	0	80	P
49	1	81	Q
50	2	82	R
51	3	83	S
52	4	84	T
53	5	85	U
54	6	86	V
55	7	87	W
56	8	88	X
57	9	89	Y
58	:	90	Z
59	;	91	[
60	<	92	\
61	=	93]
62	>	94	^
63	?	95	_

@SSR34543212 *Saccharomyces cerevisiae*

ATTCGCCAGGTCTAG

+
33 34 2 19 29 ← Phred score Q

@SSR34543212 *Saccharomyces cerevisiae* CEN.PK113-7D

ATTCGCCAGGTCTAG

+

BC#4>RTWX!![KU>

Base1: Phred value=33, sumamos 33 = 66 (ascii 66 = B)

Base2: Phred value=34, sumamos 33 = 67 (ascii 67 = C)

Base3: Phred value=2, sumamos 33 = 35 (ascii 35 = #)

Base4: Phred value=19, sumamos 33 = 52 (ascii 52 = 4)

Base15: Phred value= 29, sumamos 33 = 62 (ascii 62 = >)

Qué pasa con los reads largos?

Mayor proporciones de errores

Errores aleatorios

Mejoras en los algoritmos que realizan la asignación de bases

Repositorios públicos - banco de datos

- **SRA (Sequence Reads Archive)**

Un repositorio público con datos de secuenciaciones de distintas plataformas

- **GEO (Gene Expression Omnibus)**

Un repositorio público con datos de expresión génica (Microarrays, NGS)

SRA

SRA ▾

[Advanced](#)

Search

[Help](#)

COVID-19 is an emerging, rapidly evolving situation.

Get the latest public health information from CDC: <https://www.coronavirus.gov>.

Get the latest research from NIH: <https://www.nih.gov/coronavirus>.

Find NCBI SARS-CoV-2 literature, sequence, and clinical content: <https://www.ncbi.nlm.nih.gov/sars-cov-2/>.



SRA - Now available on the cloud

Sequence Read Archive (SRA) data, available through multiple cloud providers and NCBI servers, is the largest publicly available repository of high throughput sequencing data. The archive accepts data from all branches of life as well as metagenomic and environmental surveys. SRA stores raw sequencing data and alignment information to enhance reproducibility and facilitate new discoveries through data analysis.

Announcement

[NIH Request for Information \(RFI\) on SRA data format changes and plans.](#)

Getting Started

[How to Submit](#)[How to search and download](#)[How to use SRA in the cloud](#)[Submit to SRA](#)

Tools and Software

[Download SRA Toolkit](#)[SRA Toolkit Documentation](#)[SRA-BLAST](#)[SRA Run Browser](#)[SRA Run Selector](#)

Related Resources

[Submission Portal](#)[Trace Archive](#)[dbGaP Home](#)[BioProject](#)[BioSample](#)

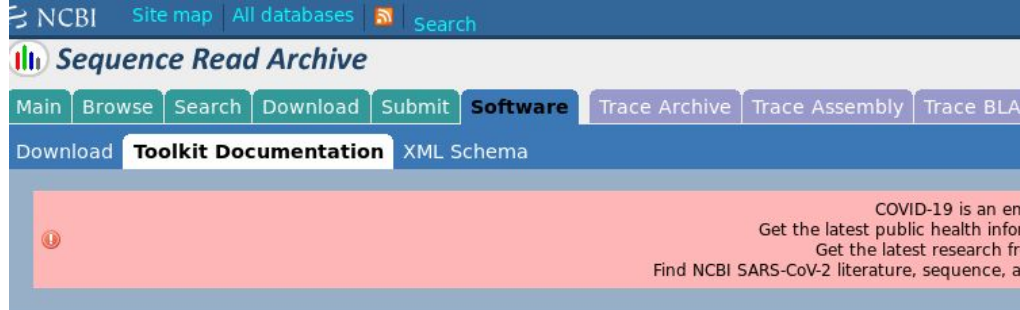
fastq-dump

Convierte datos SRA en
formato FASTQ

fastq-dump [opciones]
<IDENTIFICADOR>

Ej:

fastq-dump --splite-files
SRR6074044.sra



SRA Toolkit Documentation

[SRA Toolkit Installation and Configuration Guide](#)

[Protected Data Usage Guide](#)

Frequently Used Tools:

[fastq-dump](#): Convert SRA data into fastq format

[prefetch](#): Allows command-line downloading of SRA, dbGaP, and ADSP data

[sam-dump](#): Convert SRA data to sam format

[sra-pileup](#): Generate pileup statistics on aligned SRA data

[vdb-config](#): Display and modify VDB configuration information

[vdb-decrypt](#): Decrypt non-SRA dbGaP data ("phenotype data")

Additional Tools:

[abi-dump](#): Convert SRA data into ABI format (csfasta / qual)

[illumina-dump](#): Convert SRA data into Illumina native formats (qseq, etc.)

[sff-dump](#): Convert SRA data to sff format

[sra-stat](#): Generate statistics about SRA data (quality distribution, etc.)

[vdb-dump](#): Output the native VDB format of SRA data.

[vdb-encrypt](#): Encrypt non-SRA dbGaP data ("phenotype data")

[vdb-validate](#): Validate the integrity of downloaded SRA data

fastq-dump

Convierte datos SRA en
formato FASTQ

fastq-dump [opciones]
<IDENTIFICADOR>

Ej:

fastq-dump --splite-files
SRR6074044.sra

The read was processed for clean up before it was size selected with the BluePippin system with a cut-off value of 9000 bp. We used one SMRTcell™ to sequence the DNA library on the PacBio Sequel instrument using the Sequel 2.0 polymerase and 600 min of movie time. The high quality PacBio reads are deposited in an **SRA** database under BioProject:PRJNA398797, SRP116559.

SRA [Create alert](#) [Advanced](#) [Help](#)

COVID-19 is an emerging, rapidly evolving situation.
Get the latest public health information from CDC: <https://www.coronavirus.gov>.
Get the latest research from NIH: <https://www.nih.gov/coronavirus>.
Find NCBI SARS-CoV-2 literature, sequence, and clinical content: <https://www.ncbi.nlm.nih.gov/sars-cov-2/>.

Summary ▾ 20 per page ▾ [Send to:](#) [Filters: Manage Filters](#)

[Saccharomyces cerevisiae strain:CEN.PK113-7D Genome sequencing and assembly - BioProject](#)

We employed and combined the long reads sequencing technologies (MinION and PacBio) and short reads (Illumina). The de novo assembly of the reads derived from the three technologies was performed, resulted in the completion of the yeast genome.
Genome sequencing and assembly project
Accession: PRJNA398797

[Send results to Blast](#)

Search results
Items: 10

- ☐ [Direct RNA-seq, poly-A, ethanol, replicate 4](#)
1 OXFORD_NANOPORE (MinION) run: 858,903 spots, 876.5M bases, 751.5Mb downloads
Accession: SRX3449623
- ☐ [Direct RNA-seq, poly-A, ethanol, replicate 3](#)
2 OXFORD_NANOPORE (MinION) run: 455,282 spots, 437.8M bases, 375.9Mb downloads
Accession: SRX3449622
- ☐ [Direct RNA-seq, poly-A, ethanol, replicate 2](#)
3 OXFORD_NANOPORE (MinION) run: 911,329 spots, 938.8M bases, 803.3Mb downloads
Accession: SRX3449621

Search in related databases

Database	Access		all
	public	controlled	
BioSample			
BioProject	1		1
dbGaP			
GEO Datasets			

Find related data

Database:

Search details

PRJNA398797[All Fields]

[See more...](#)

Recent activity

Obtención de datos de SRA

- Usando PREFETCH

prefetch SRR8922830

(genera un archivo SRR8922830.sra)

- Usando fastq-dump

fastq-dump -Z -X 10 SRR8922830

(imprime en pantalla (opción -Z) las primeras 10 secuencias (-X 10))

- Usando WGET

wget <https://sra-downloadb.st-v.a.ncbi.nlm.nih.gov/sos1/sra-pub-run-12/SRR6352892/SRR6352892.1>

- Usando FTP

ftp://ftp-trace.ncbi.nih.gov/sra/sra-instant/reads/ByRun/sra/SRR/SRR607/SRR6074044/SRR6074044.sra

Usando WGET

Sequence Read Archive

Main Browse Search Download Submit Software Trace Archive Trace Assembly Trace BLAST

Studies Samples Analyses **Run Browser** Run Selector Provisional SRA



COVID-19 is an emerging, rapidly evolving situation.
Get the latest public health information from CDC: <https://www.coronavirus.gov>.
Get the latest research from NIH: <https://www.nih.gov/coronavirus>.
Find NCBI SARS-CoV-2 literature, sequence, and clinical content: <https://www.ncbi.nlm.nih.gov/sars-cov-2/>.

Direct RNA-seq, poly-A, ethanol, replicate 4 (SRR6352892)

Metadata Analysis Reads **Data access**

SRA archive data

SRA archive data is normalized by the SRA load process and used by the [SRA Toolkit](#) to read and produce formats like FASTQ, SAM, etc. The default toolkit configuration enables it to find and retrieve SRA runs by accession.

Public SRA files are now available from GCP and AWS cloud platforms as well as from NCBI. Access to most data in the cloud requires a user account with the cloud service provider. The user's account will incur costs for cloud compute or to copy data outside of the specified cloud service region.

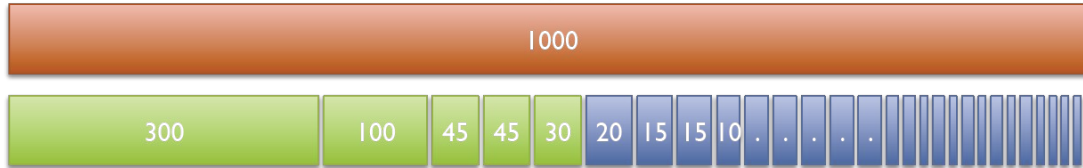
Type	Size	Location	Name	Free Egress	Access Type
run	769,543 Kb	NCBI	https://sra-downloadb.be-md.ncbi.nlm.nih.gov/sos2/sra-pub-run-11/SRR6352892/SRR6352892.1	worldwide	anonymous
		NCBI	https://sra-downloadb.st-va.ncbi.nlm.nih.gov/sos1/sra-pub-run-12/SRR6352892/SRR6352892.1	worldwide	anonymous
		AWS	s3://sra-pub-run-5/SRR6352892/SRR6352892.1	s3.us-east-1	aws identity
		GCP	gs://sra-pub-run-5/SRR6352892/SRR6352892.1	gs.US	gcp identity

N50 size

Def: 50% of the genome is in contigs larger than N50

Example: 1 Mbp genome

50%



N50 size = 30 kbp

$(300k + 100k + 45k + 45k + 30k = 520k \geq 500kbp)$

Note:

N50 values are only meaningful to compare when base genome size is the same in all cases