

# Maria del Saz Navarro



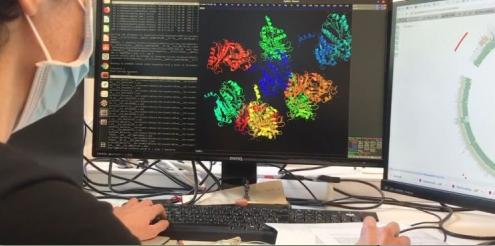
0000-0003-1551-5341



<https://github.com/mariadsn>



mariadsn89@gmail.com



## Bióloga

- **Técnico de laboratorio Microbiología**  
U. Sevilla
- **Tecnico Bioinformática**  
Centro Andaluz de Biología del Desarrollo (CABD)
- **Pre-PhD Bioinformática (DATAi group)**  
U. Pablo de Olavide

## Universidad Pablo de Olavide

Ctra. de Utrera, Sevilla

## Intelligent Data Analysis (DATAi)

### Members

Norberto Díaz Díaz

Carlos D. Barranco

Francisco A. Gómez Vela

Domingo S. Rodríguez Baena

Aurelio López Fernández

Juan José Díaz Montaña

### Ph.D student

José A. Lagares

María del Saz Navarro

Francisco Gonzalez Cabanes



UNIVERSIDAD

PABLO  
OLAVIDE



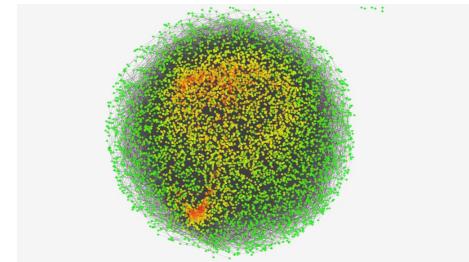
Universidad  
Internacional  
de Valencia

Máster Universitario en **Big  
Data y Ciencia de Datos**

7 Diciembre 2023  
18:00h



# *Técnicas de Procesamiento de datos de Secuenciación de Alto Rendimiento*



**María del Saz**  
Pre-Ph.D



Minería de Datos

# Procesamiento de Datos **BIOLÓGICOS**

INTRO

CONCEPTOS

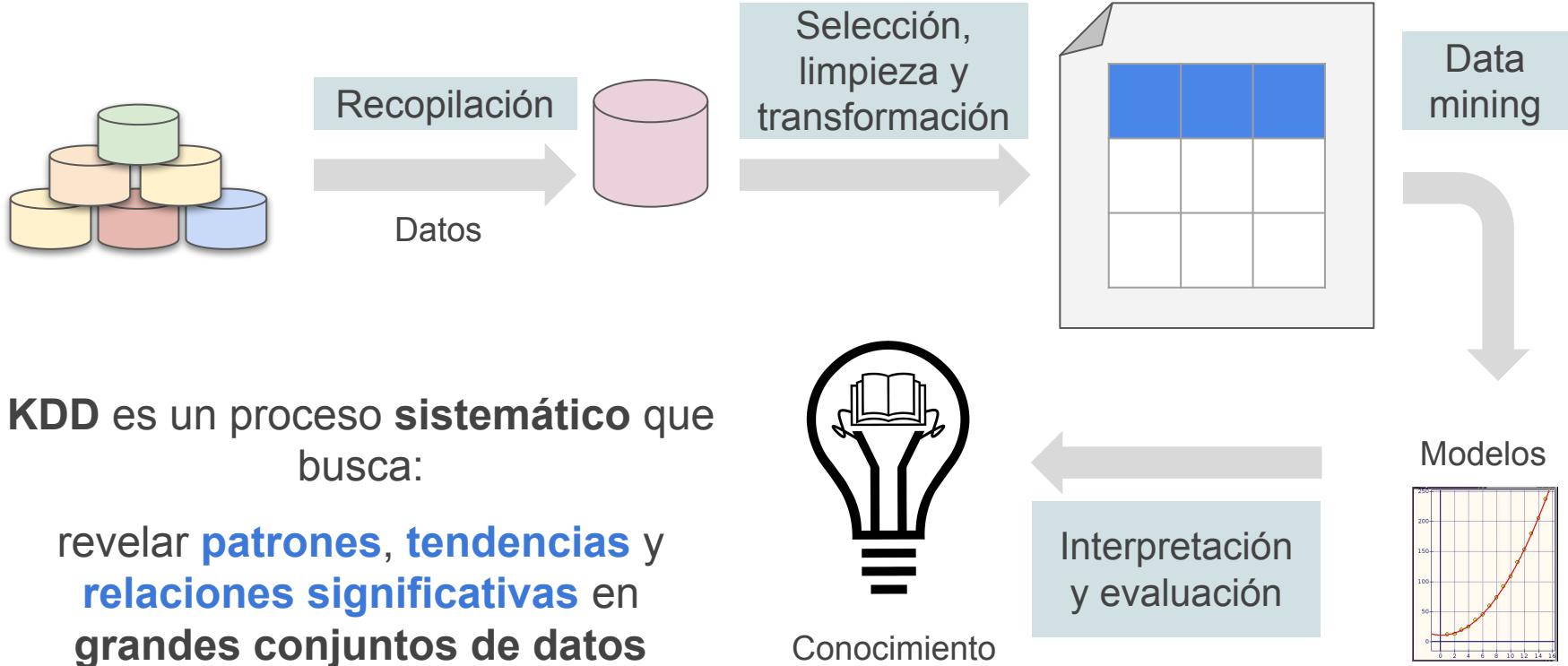
INVESTIGACIÓN

WORKSHOP

GDC

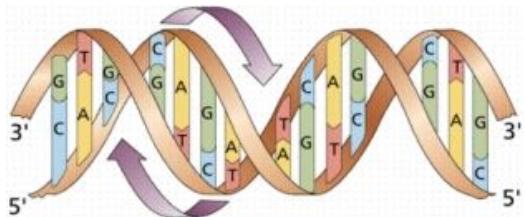
.fastq

# 1- Introducción / Minería de Datos



# Fundamentos Básicos en Biología

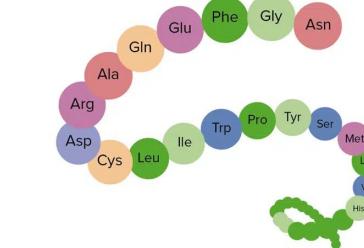
- 1- Dogma Central de la Biología
- 2- Secuenciación de Segunda Generación



```
>Gen 1
TCATATTGTTTACGTTGTCAAGCCTCAT
AGCCGAGTTGAACGTATACGCTCTGA
GTCAGACCTAAATCGTAGCTACACAAT
TCTGTGAATTTCTTGTGCGGTGAAAC
ACTTCCAATAAAAATCATATGGTGAGTACT
TTAAAAAAAATCTAGTCAAATAATGCTGAA
AAGAAATTGTGTGGGCAAAATTCAATGGG
CAAAACGCGATGCGGCTTTCTCAAAAT
GGCGGCCGGCTGCGTTTTCTCAAAAAT
GTGATGACGTATGCCTGTTTTTTTG
TTGCAATGAGGAATGGCTTAAAT...
```



```
>ANRm 1
TCATATTGTTTACGTTGTCAAGCCTCATAG
CCGGCAGTTGAACGTATACGCTCTGAGTC
AGACCTCGAAATCGTAGCTACACAATTCTG
TGAATTTCTTGTGCGGTGAAACACTTCC
AATAAAAATCAATATGCGTAATGTATCTCA
TCCATGTTGGTCAGGCTGGTGTCCAGATTGGA
AACGCCTGCTGGGAGCTACTGCTTGGAGC
ACGGCATCCAGCCGATGGCCAGATGCCGTC
TGACAAGACCCTGGCGAGGTGATGACTCG
TTCAACACCTCTCAGCGAGACTGGAGCTGG
CAAGCAC...
```



```
>Proteína 1
MRECISIHVGQAGVQIGNACWELYCLEHGIQP
DGQMPSDKTVGGGDDSFNTFFSETGAGKHVP
RAVFVDLEPTVVDEVRTGTYRQLFHPEQLITGK
EDAANNYARGHYTIGKEIVDLVLDRIKLAQDC
TGLQGFLIFHSFGGGTGSFTSLLMERLSVDY
GKKSKLEFAIYPAQPVSTAVVEPYNSILTTHTTL
EHSDCAFMDNEAIYDICRRNLDIERPTYTNLN
RLIGQIVSSITASLRFDGALNVDLTEFQTNLVPY
PRIHFPLVTYAPVISAEKAYHEQLSVAEITNACF
EPANQMVVKCDPRHGKYMACCMLYRGDVVPK
DVNAAIATIKTKRTIQFVDWCPTGFKVGINYQPP
TVVPGGDLAKVQRAVCM
```

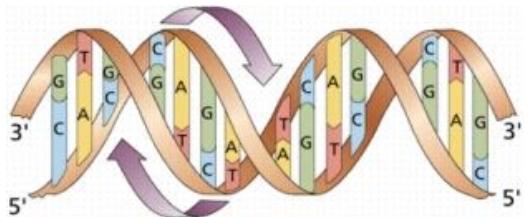
DNA

Transcripción

RNA

Traducción

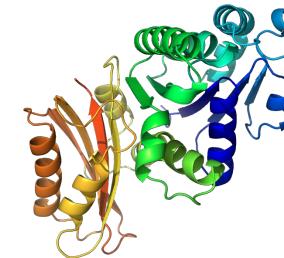
Protein



```
>Gen 1
TCATATTGTTTACGTTGTCAAGCCTCAT
AGCCGAGTTGAACGTATAACGCTCTGA
GTCAGACCTAAATCGTAGCTACACAAT
TCTGTGAATTTCTTGTGCGGTGAAAC
ACTTCCAATAAAAACTCATATGGTGAGTACT
TTAAAAAAAAAATCTAGTCAAATAATGCTGAA
AAGAAATTGTGTGGGCAAAATTCAATGGG
CAAAAACGCGATGCGGCTTTCTCAAAAT
GGCGGCCGGCTGCGTTTTCTCAAAAAT
GTGATGACGTATGCCTGTTTTTTTTTG
TTCGCAATGAGGAATGGCTTAAAT...
```



```
>ANRm 1
TCATATTGTTTACGTTGTCAAGCCTCATAG
CCGGCAGTTGAACGTATAACGCTCTGAGTC
AGACCTCGAAATCGTAGCTACACAATTCTG
TGAATTTCTTGTGCGGTGAAACACTTCC
AATAAAAAACTCAATATGCGTAATGTATCTCA
TCCATGTTGGTCAGGCTGGTGTCCAGATTGGA
AACGCCTGCTGGGAGCTACTGCTTGGAGC
ACGGCATCCAGCCGATGGCCAGATGCCGTC
TGACAAGACCCTGGCGAGGTGATGACTCG
TTCAACACCTCTCAGCGAGACTGGAGCTGG
CAAGCAC...
```



```
>Proteína 1
MRECISIHVGQAGVQIGNACWELYCLEHGIQP
DGQMPSDKTVGGGDDSFNTFFSETGAGKHVP
RAVFVDLEPTVVDEVRTGTYRQLFHPEQLITGK
EDAANNYARGHYTIGKEIVDLVLDRIKLAQDC
TGLQGFLIFHSFGGGTGSFTSLLMERLSVDY
GKKSKLEFAIYPAQPVSTAVVEPYNSILTTHTTL
EHSDCAFMDNEAIYDICRRNLDIERPTYTNLN
RLIGQIVSSITASLRFDGALNVDTFQTNLVPY
PRIHFPLVTYAPVISAEKAYHEQLSVAEITNACF
EPANQMVVKCDPRHGKYMACCMLYRGDVVPK
DVNAAIATIKTKRTIQFVDWCPTGFKGINYQPP
TVVPGGDLAKVQRAVCM
```

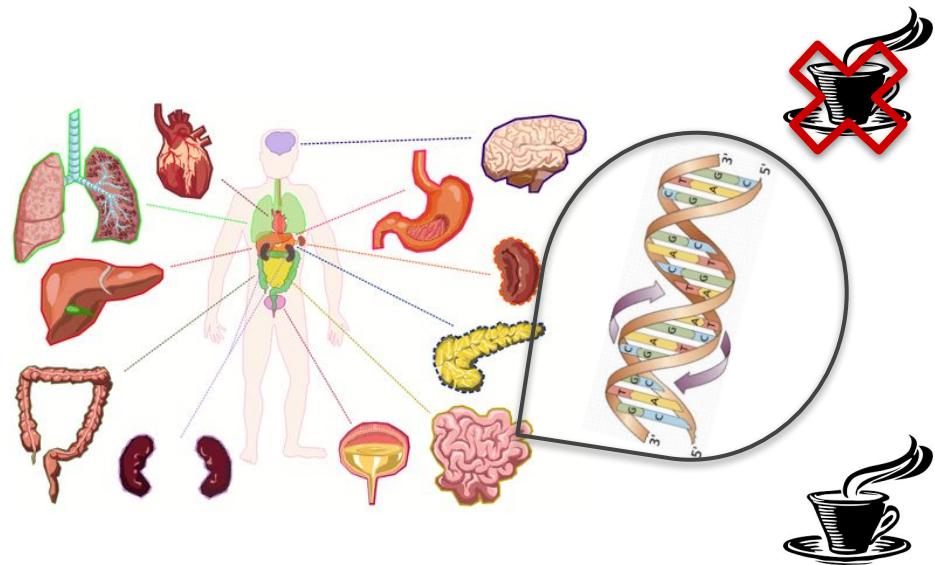
DNA

Transcripción

RNA

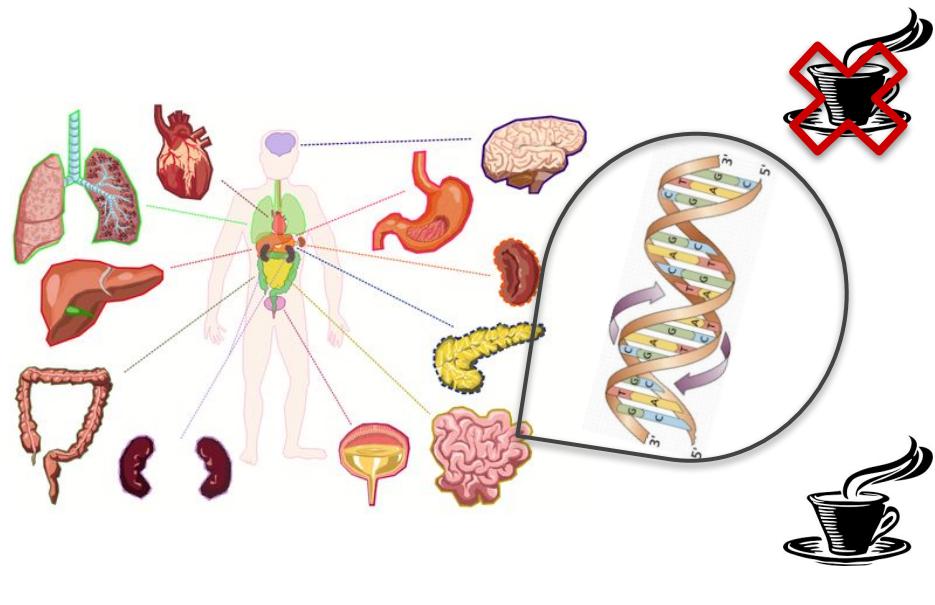
Traducción

Protein



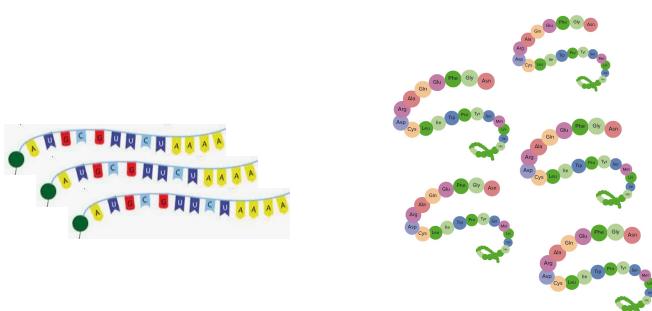
DNA → RNA → Protein

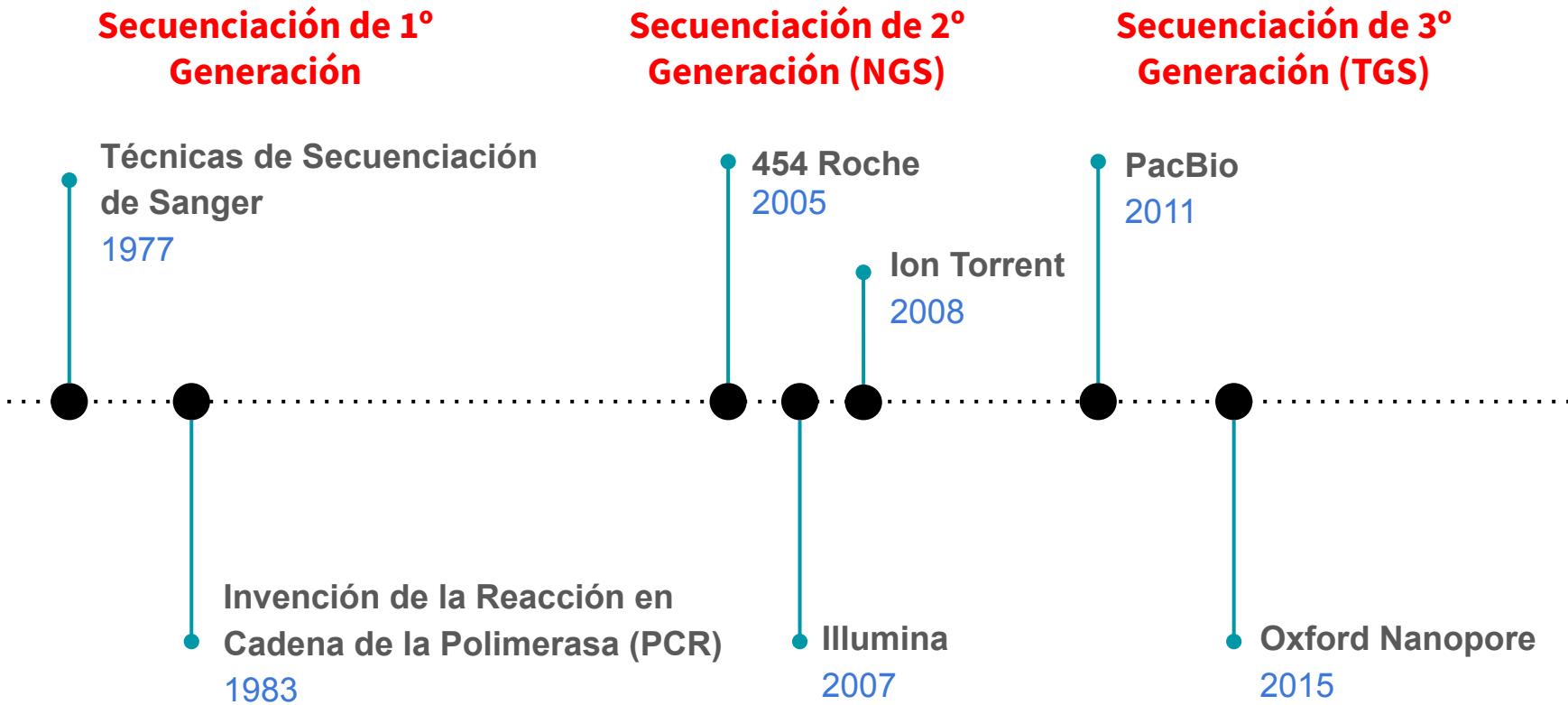
Transcripción                      Traducción



DNA → RNA → Protein

Transcripción                      Traducción

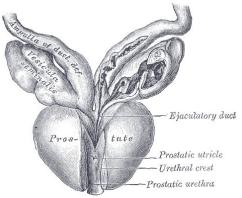




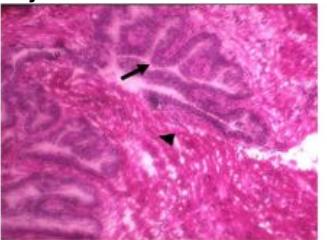
# Investigación Biológica

- 1- Diseño Experimental
- 2- Fases del Estudio

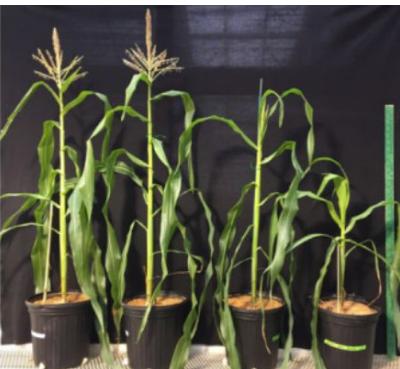
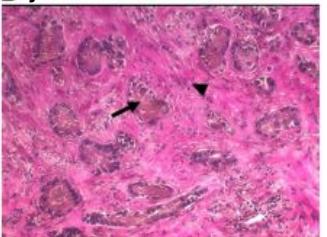
## 3- Investigación / Diseño Experimental



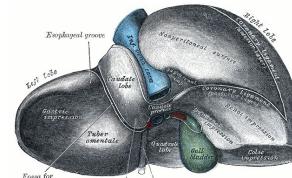
Tejido Prostático Normal



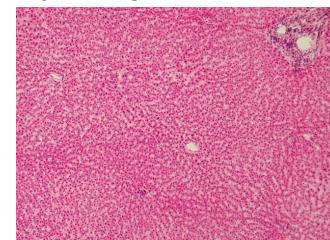
Tejido Prostático Tumoral



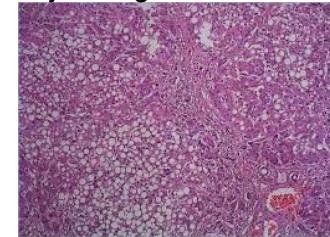
Ambiente      Sequía



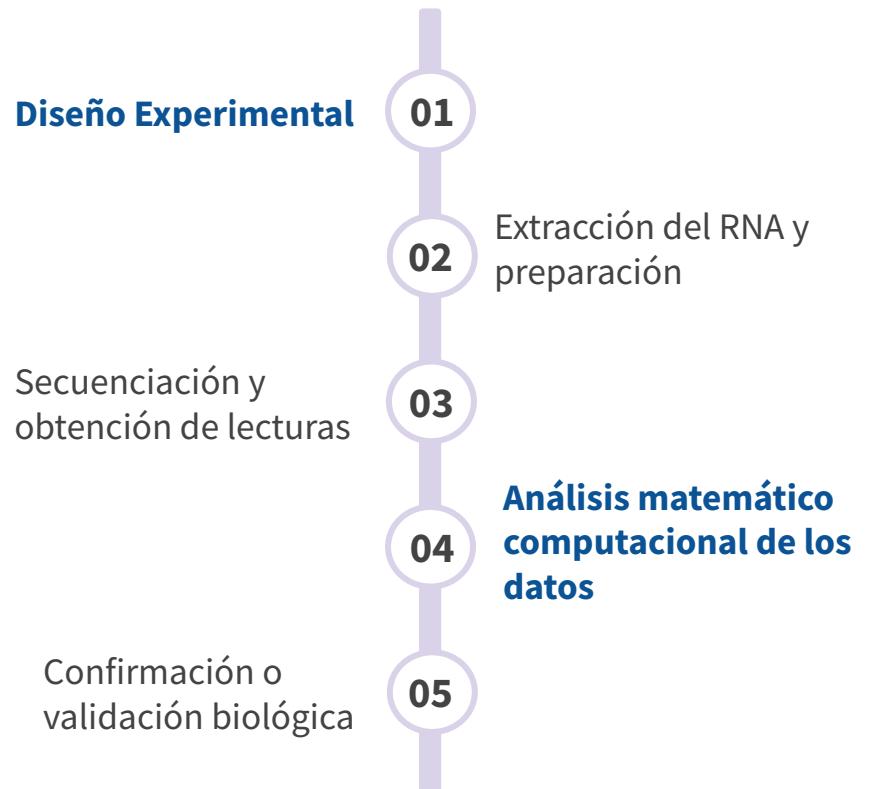
Tejido Hígado Normal



Tejido Hígado Cirrosis



Estudio transcriptómico  
basado en **RNA-Seq**



## Estudio transcriptómico basado en **RNA-Seq**

### Diseño Experimental

01

Secuenciación y obtención de lecturas

02

Extracción del RNA y preparación

03

Análisis matemático computacional de los datos

04

Confirmación o validación biológica

05

02



03



.fastq

Identifier @HWI-EAS209\_0006\_FC706VJ:5:58:5894:21141#ATCACG/1  
Sequence TTAATTGGTAAATAATCTCTAAAGCTTAGANTTTACCTNNNNNNNNNTAGTTCTTGAGA  
+ sign & identifier +HWI-EAS209\_0006\_FC706VJ:5:58:5894:21141#ATCACG/1  
Quality scores efcfffffcfeefffcfffffdff`feed`^]\_Ba\_`^[\_YBBBBBBBBBBRTT\`]] dddd`  
Base T phred Quality ] = 29

## Estudio transcriptómico basado en **RNA-Seq**

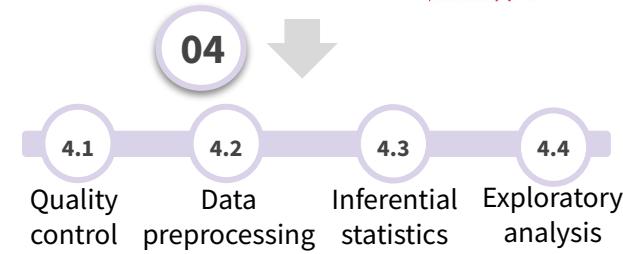
### Diseño Experimental

- 01
- 02 Extracción del RNA y preparación
- 03 Secuenciación y obtención de lecturas
- 04 Análisis matemático computacional de los datos
- 05 Confirmación o validación biológica



.fastq

Identifier @HWI-EAS209\_0006\_FC706VJ:5:58:5894:21141#ATCACG/1  
Sequence TTAATTGGTAAATAATCTCTAAATGCTTAGATNTTACCTNNNNNNNTAGTTCTTGAGA  
+ sign & identifier +HWI-EAS209\_0006\_FC706VJ:5:58:5894:21141#ATCACG/1  
Quality scores efcfffffcfeffffccfffffdff`feed)`\_Ba\_`^\_[YBBBBBBBBBBRTTV]`{}` dddd`  
Base T phred Quality ] = 29



# Estudio transcriptómico

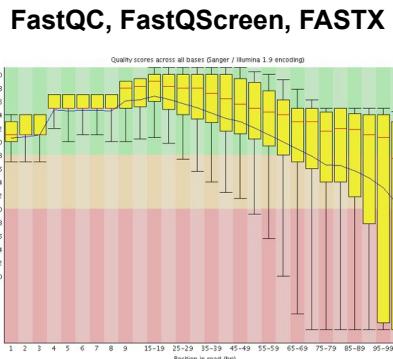
- 1- .fastq
- 2- Objetivo



[https://github.com/mariadsn/GeneExpressionAnalysis\\_Seminar](https://github.com/mariadsn/GeneExpressionAnalysis_Seminar)



## Quality control

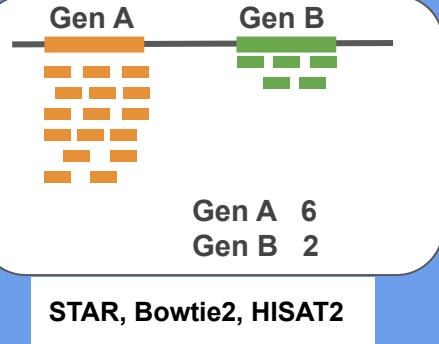


4.1

## Inferential statistics

4.2

## Data preprocessing

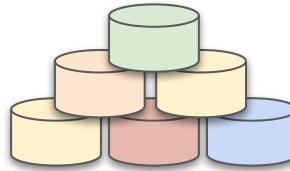


4.3

## Exploratory analysis

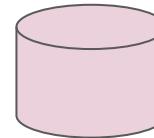
4.4

## 4- Workshop / Objetivo



Recopilación

Datos



## Gene Expression Omnibus (GEO)

Noviembre 2023

Organism	Platforms	Samples
<i>Homo sapiens</i>	6,257	3,585,720
<i>Mus musculus</i>	2,746	2,035,017
<i>Drosophila melanogaster</i>	402	115,391
<i>Arabidopsis thaliana</i>	416	94,502
<i>Saccharomyces cerevisiae</i>	639	90,187

Harmonized Cancer Datasets  
Genomic Data Commons Data Portal

Get Started by Exploring:

[Projects](#) [Exploration](#) [Analysis](#) [Repository](#)

e.g. BRAF, Breast, TCGA-BLCA, TCGA-A5-A0G2

Data Portal Summary Data Release 38.0 - August 31, 2023

PROJECTS

79

FILES

986.114

PRIMARY SITES

69

GENES

22,534

CASES

44,451

MUTATIONS

2,930,136



Expression Atlas

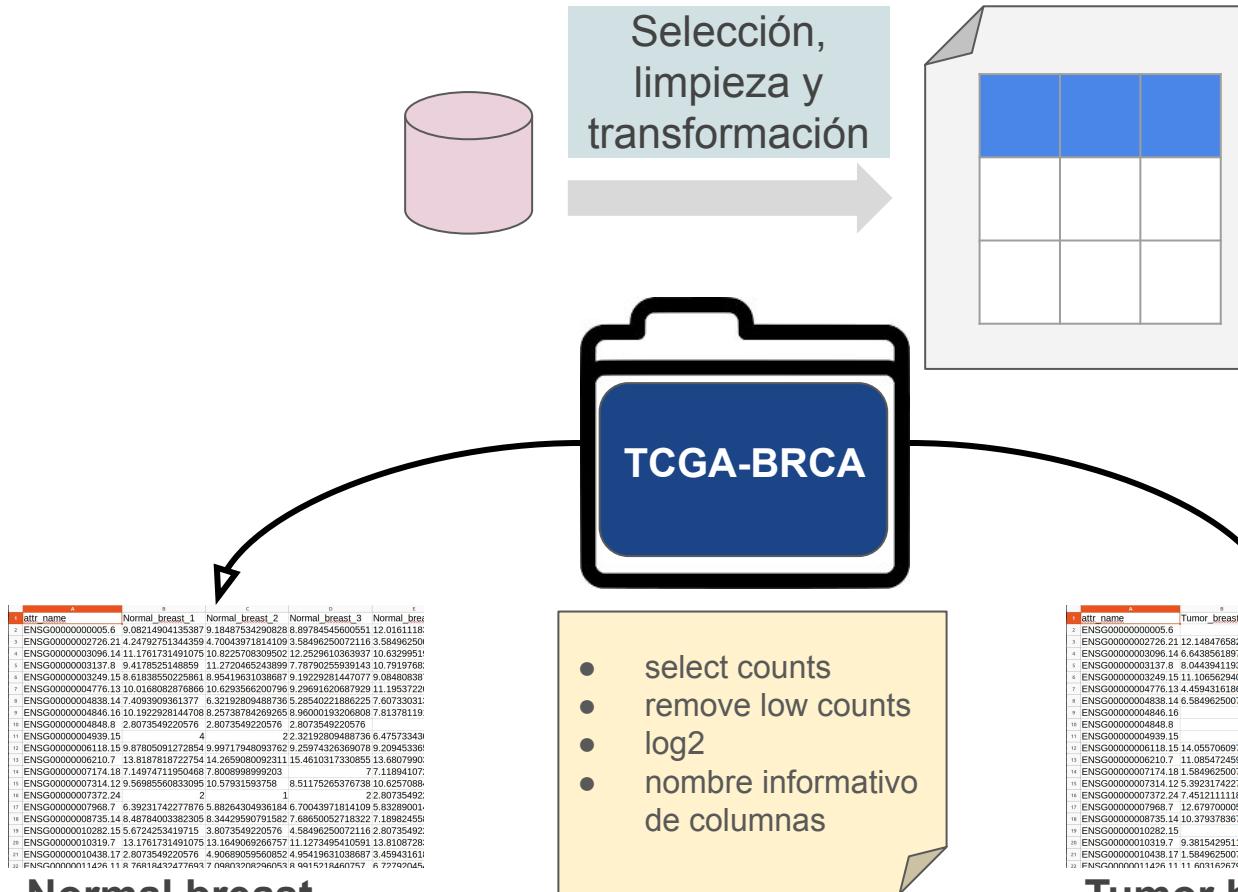




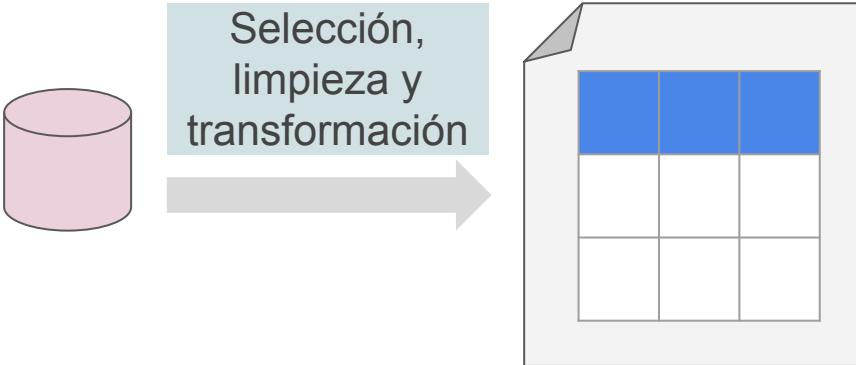
## Objetivos del Workshop

- Anotación funcional de genes
- Obtención de biomarcadores

## 4- Workshop / Objetivo



## 4- Workshop / Objetivo

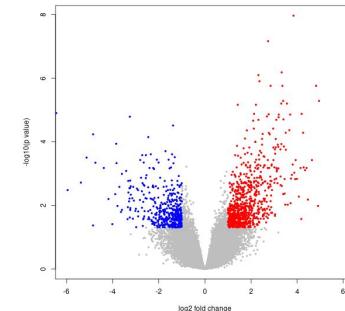


## Normal breast

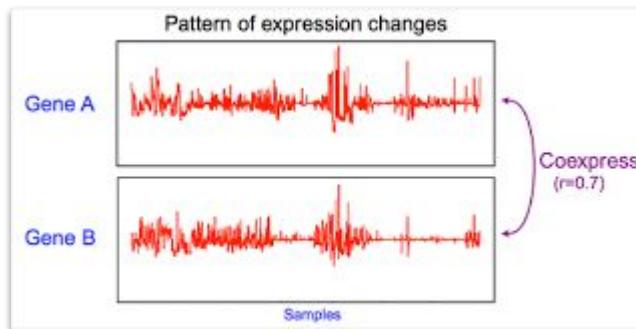
## Tumor breast



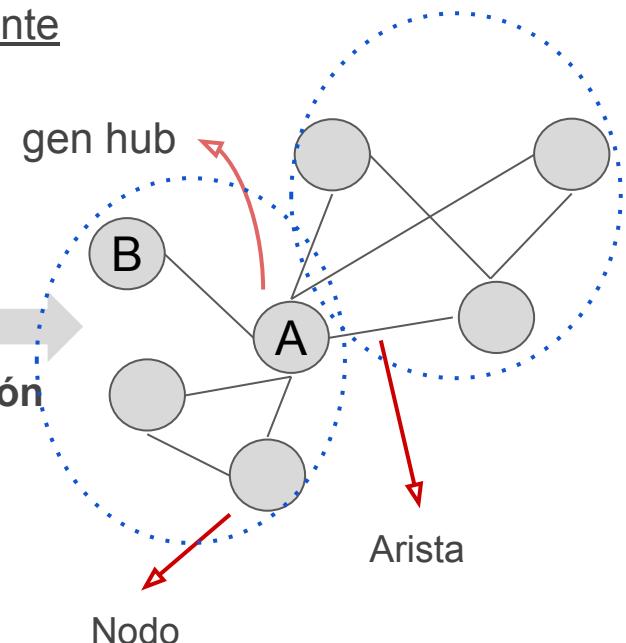
## Expresión diferencial de genes



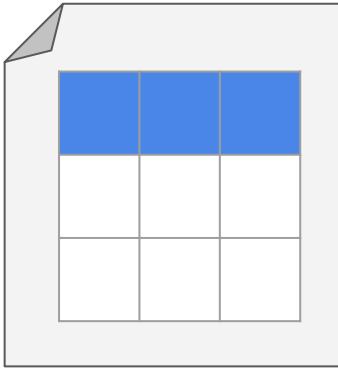
Genes con patrones de expresión similares (co-expresados) están asociados funcionalmente



Red de Co-expresión génica



## 4- Workshop / Objetivo



# Data mining



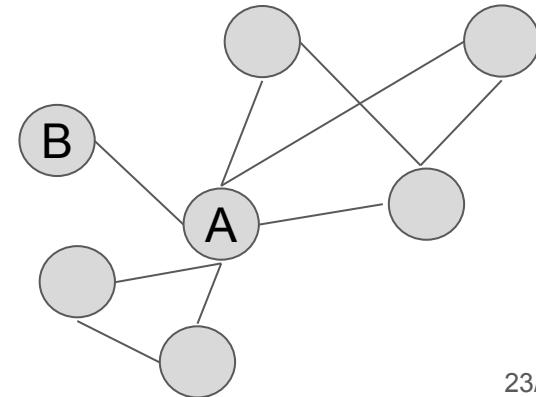
## Normal breast

## Tumor breast

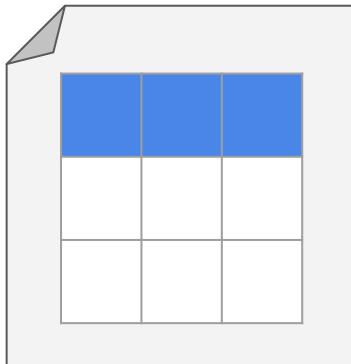
## Tools

- EnGNet
  - WGCNA
  - NetMiner
  - DiffCoEx
  - RNASeqNet

## Redes de coexpresión



## 4- Workshop / Objetivo



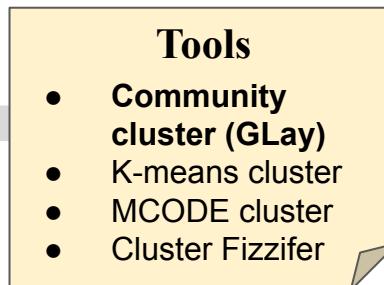
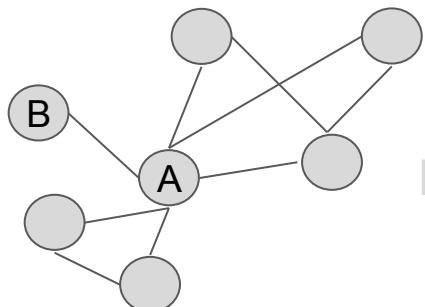
Data  
mining



Modelos

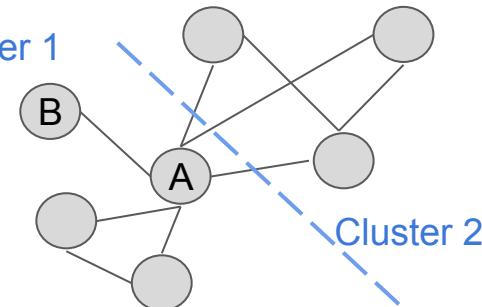


Redes de coexpresión

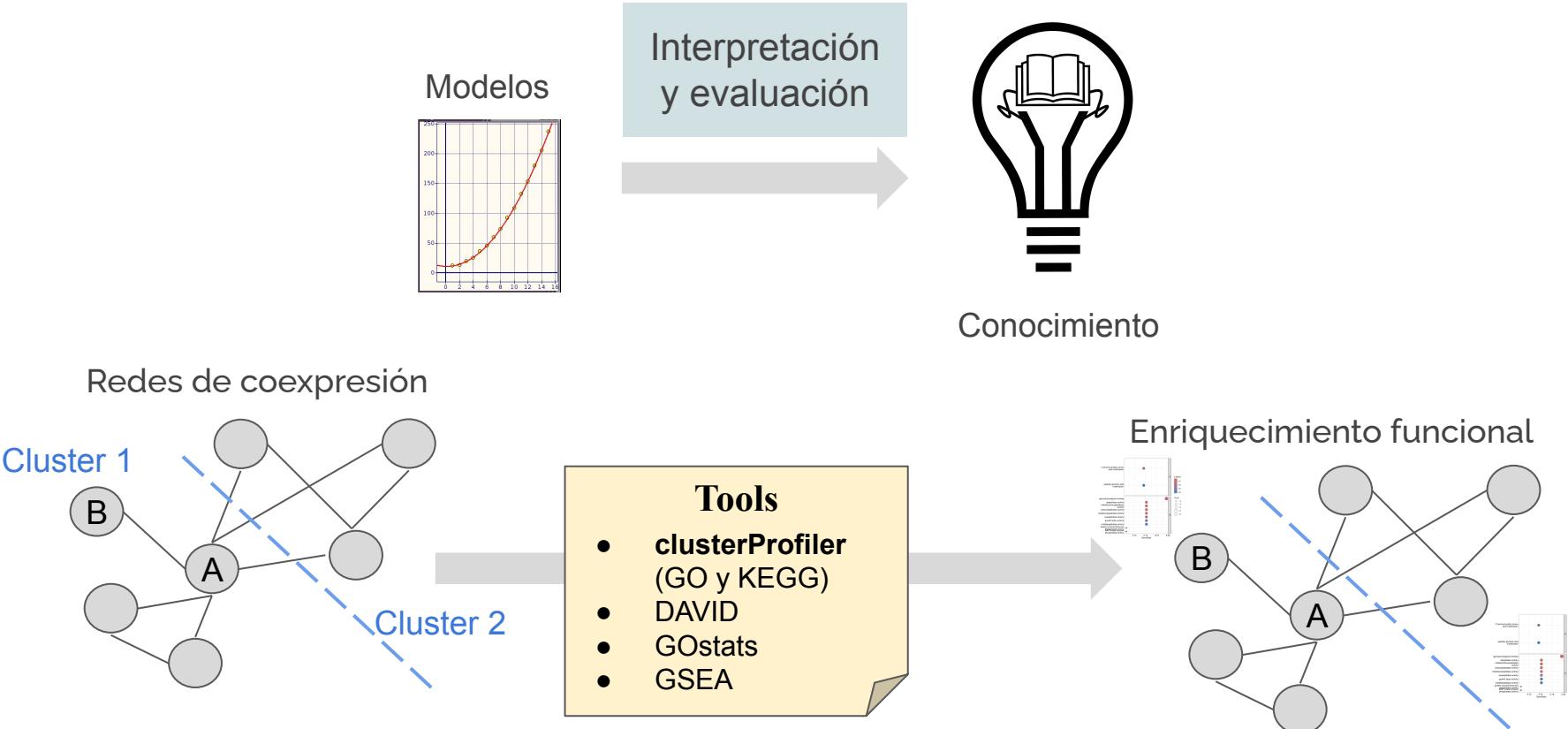


Clusterización de la red

Cluster 1

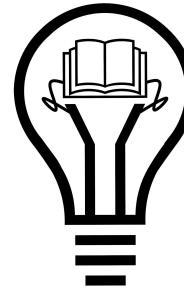


## 4- Workshop / Objetivo



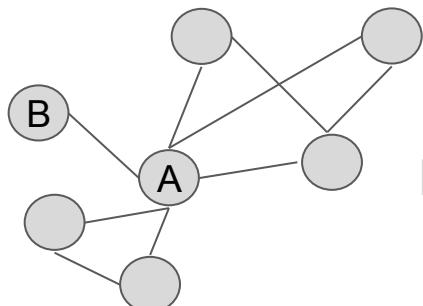
## 4- Workshop / Objetivo

Modelos

Interpretación  
y evaluación

Conocimiento

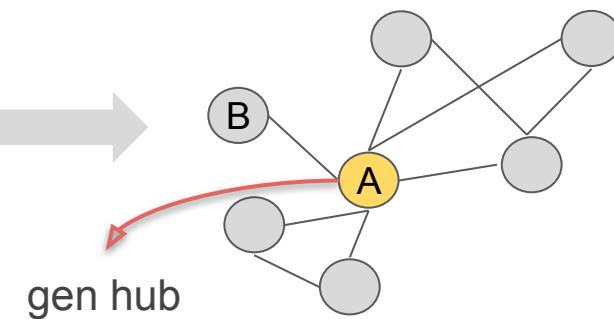
Redes de coexpresión



## Tools

- Degree centrality
- Betweenness centrality

Enriquecimiento funcional





[https://github.com/mariadsn/GeneExpressionAnalysis\\_Seminar](https://github.com/mariadsn/GeneExpressionAnalysis_Seminar)



## Quality control

4.1

## Data preprocessing

4.2

## Inferential statistics

4.3

## Exploratory analysis

4.4