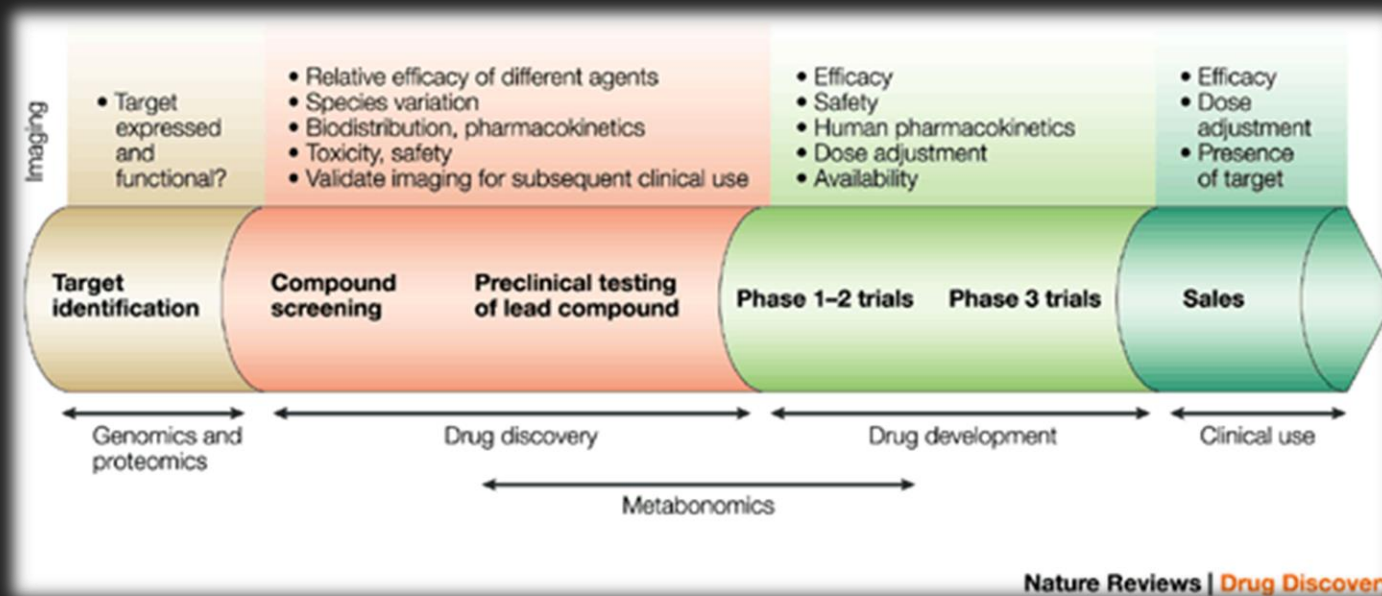


Combinar R y Python en Oncología

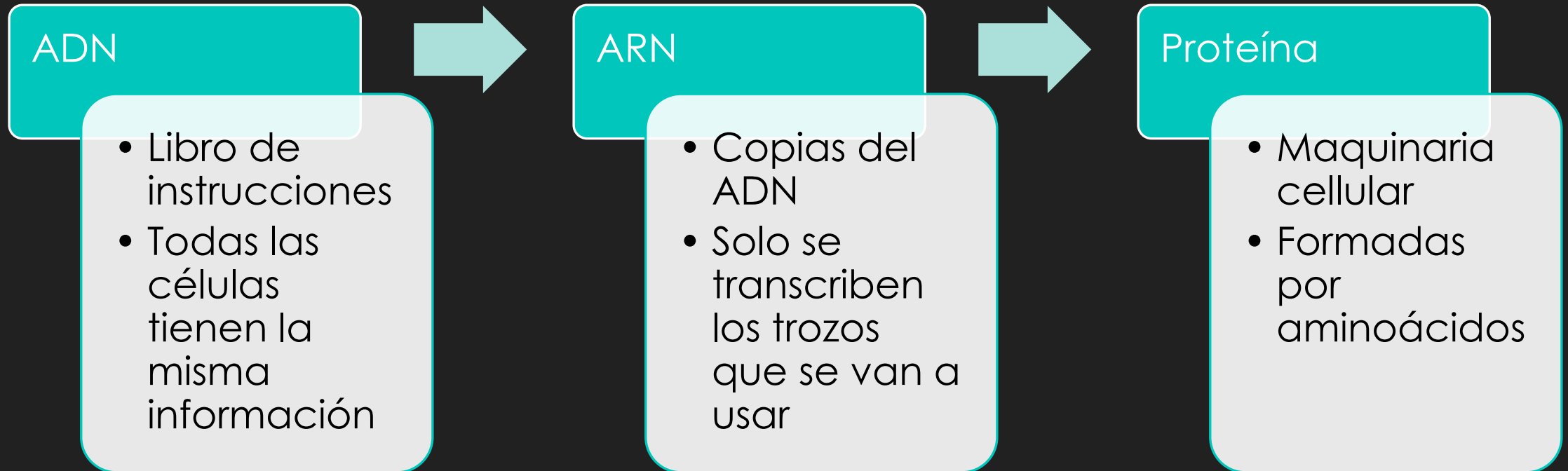
Amalia Martínez Segura

¿Qué hace un bioinformático en la industria farmacéutica?

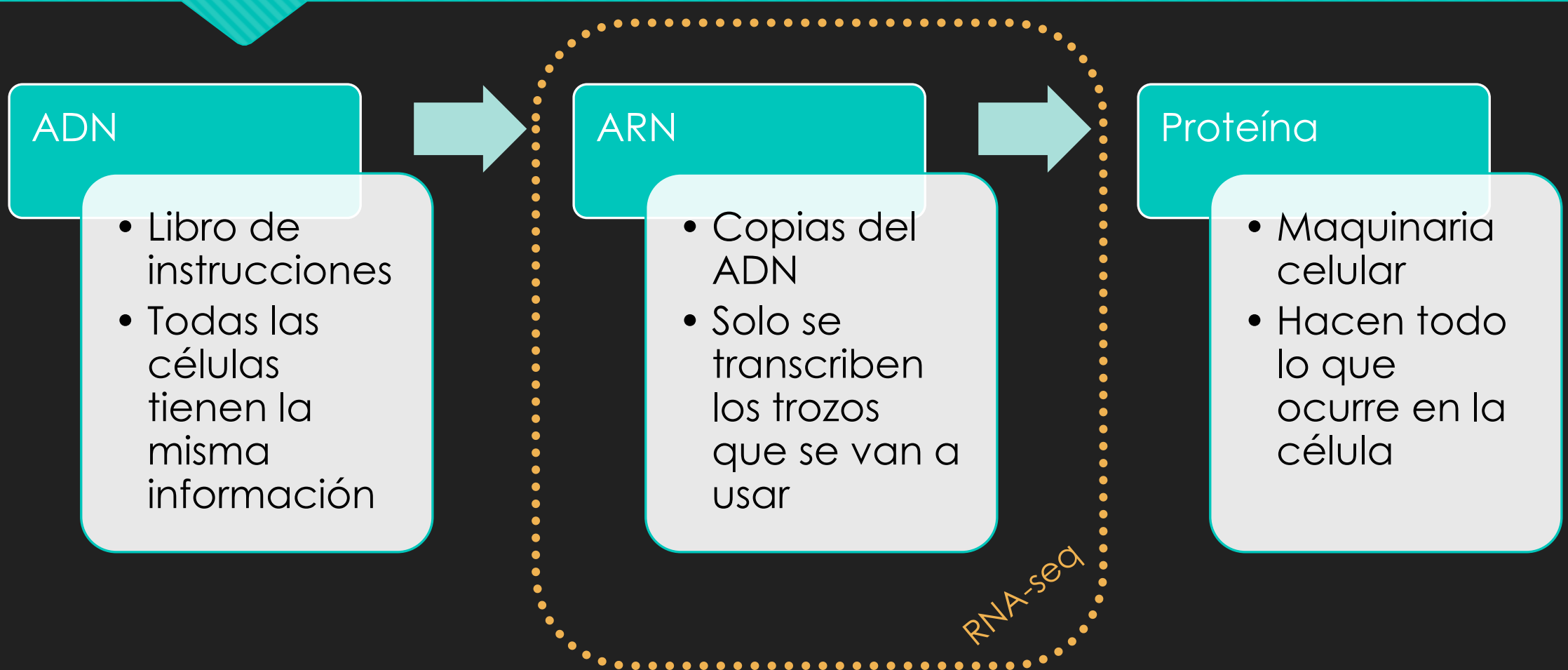


- Los bioinformáticos podemos participar en todos los estadios de desarrollo de un medicamento

Dogma fundamental de la biología



Dogma fundamental de la biología



Datos de secuenciación masiva

Secuenciación



Datos sin procesar

Header	Sequence	Quality
@HWI-ST227:389:C4WA2ACXX:7:1204:2272:59979	GGAGGAAGGTCTCTCGTCCCTCTTTCATATAAGGGAATGGCTGAAT	+
FFFFHHHHHHJJJJJJJJJJJJJIGIGIGIGIJJJJJJJJJJJJJJ	@HWI-ST227:389:C4WA2ACXX:7:1205:15214:42893	
GAGGATCCCAGGAGGAGAAGGTCTCGCTCCTTTCATCTAAGGGA	+	
12BAFB?A:3<AE1@<FF;1*(EG*)?0?DBD>9BF9B*?#####	@HWI-ST227:389:C4WA2ACXX:8:2208:2467:44624	
AAAGAGGAGAGAGGACCATCCTCCCTGGGATCCTCAGAAGTCTACT	+	
BDDA:DB?2AA@FC>F?EEGC<FED>GFD;?GBB?<?F99*/9?9?		

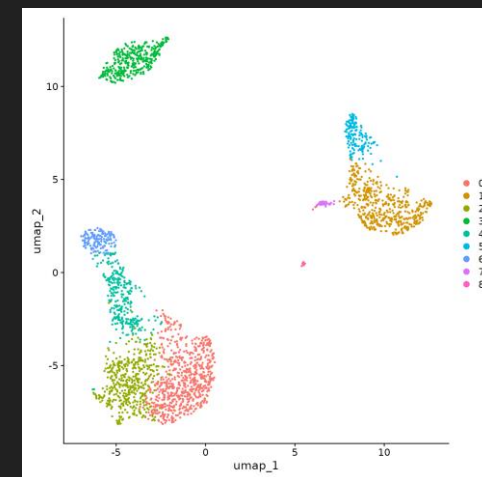
Matriz numérica

Células

100

Metadata

Clustering y visualización



Pasos básicos del análisis

1. QC de los datos del secuenciador
2. Transformar los datos del secuenciador en una matriz numérica
3. QC de las células
4. Reducción de dimensionalidad
5. Análisis posteriores

Pasos básicos del análisis

1. QC de los datos del secuenciador
2. Transformar los datos del secuenciador en una matriz numérica
3. QC de las células
4. Reducción de dimensionalidad y visualización
5. Análisis posteriores

Command line

R o Python



Orquestrando workflows con Snakemake

- “Tool to create **reproducible and scalable** data analyses. Workflows are described via a human readable, **Python** based language. ”

```
rule select_by_country:  
    input: "data/worldcitiespop.csv"  
    output: "by-country/{country}.csv"  
    shell: "xsv search -s Country '{wildcards.country}' {input} > {output}"
```

```
Snakemake -c 1 by-country/Spain.csv  
wildcard = 'Spain'
```


Orquestrando workflows con Snakemake

- “Tool to create **reproducible and scalable** data analyses. Workflows are described via a human readable, **Python** based language. ”

```
rule select_by_country:  
    input: "data/worldcitiespop.csv"  
    output: "by-country/{country}.csv"  
    shell: "xsv search -s Country '{wildcards.country}' {input} > {output}"
```

```
Snakemake -c 1 by-country/Spain.csv by-country/UnitedKingdom.csv  
wildcard = ['Spain', 'UnitedKingdom']
```

Un workflow básico en snakemake

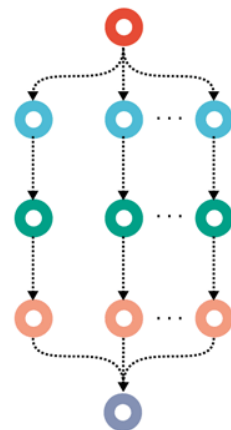
a

```
1 ● configfile: "config.yaml"
2
3 ● rule all:
4   input:
5     expand(
6       "results/plots/{country}.hist.pdf",
7       country=config["countries"]
8     )
9
10 ● rule download_data:
11   output:
12     "data/worldcitiespop.csv"
13   log:
14     "logs/download.log"
15   conda:
16     "envs/curl.yaml"
17   shell:
18     "curl -L https://burntsushi.net/stuff/worldcitiespop.csv > {output} 2> {log}"
19
20 ● rule select_by_country:
21   input:
22     "data/worldcitiespop.csv"
23   output:
24     "results/by-country/{country}.csv"
25   log:
26     "logs/select-by-country/{country}.log"
27   conda:
28     "envs/xsv.yaml"
29   shell:
30     "xsv search -s Country '{wildcards.country}' "
31     "{input} > {output} 2> {log}"
32
```

Legend

- domain knowledge
- technical knowledge
- Snakemake knowledge
- trivial

b



c

```
import sys
sys.stderr = open(snakemake.log[0], "w")

import matplotlib.pyplot as plt
import pandas as pd

cities = pd.read_csv(snakemake.input[0])

plt.hist(cities["Population"], bins=50)

plt.savefig(snakemake.output[0])
```

```
33 ● rule plot_histogram:
34   input:
35     "results/by-country/{country}.csv"
36   output:
37     "results/plots/{country}.hist.svg"
38   container:
39     "docker://faizanbashir/python-datascience:3.6"
40   log:
41     "logs/plot-hist/{country}.log"
42   script:
43     "scripts/plot-hist.py"
44
45 ● rule convert_to_pdf:
46   input:
47     "{prefix}.svg"
48   output:
49     "{prefix}.pdf"
50   log:
51     "logs/convert-to-pdf/{prefix}.log"
52   wrapper:
53     "0.47.0/utils/cairosvg"
```

Anotación de tipos celulares

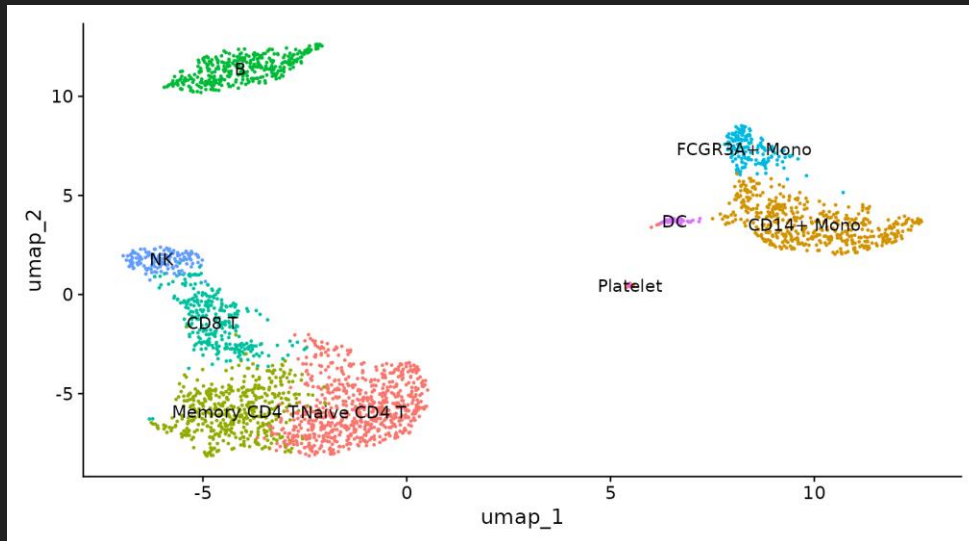
○ Anotación manual

- Usando marcadores conocidos en la literatura
- Sesgada

○ Anotación automática



- Menos sesgos
- Puede que no haya referencia para los tipos celulares en tu muestra



Habla con tu
biólogo de
confianza!



Otras aplicaciones

- ¿Hay algún cambio en los tipos celulares presents en la muestra?
 - Uso de herramientas de anotación automática basadas en datasets de referencia anotados manualmente
 - Pseudotiempo
- ¿Hay cambios en la expresión génica el tratamiento vs el control?
 - Análisis de expresión diferencial
- Buscar correlatos con variables de importancia clínica – por ejemplo, si el numero de linfocitos T en una muestra correlaciona con la supervivencia de un paciente de cáncer

¿Preguntas?