

# Single cell RNA-seq

Mónica Padilla

Alejandra Eugenia Medina Rivera, PhD.

Basado en:

Orchestrating single-cell analysis with Bioconductor (2020) Robert A. Amezquita... Stephanie C. Hicks.

Material de [Peter Hickey](#)

# Overview

- ¿Qué es el scRNA-seq?
- Captura de una sola célula
- Tecnologías de scRNA-seq y desafíos
- Procesamiento general de datos de scRNA-seq
  - Preprocesamiento de datos
  - Control de Calidad (QC)
  - Normalización
  - Imputación
  - Selección de atributos (feature selection)
  - Reducción de dimensiones
  - Integración de conjuntos de datos

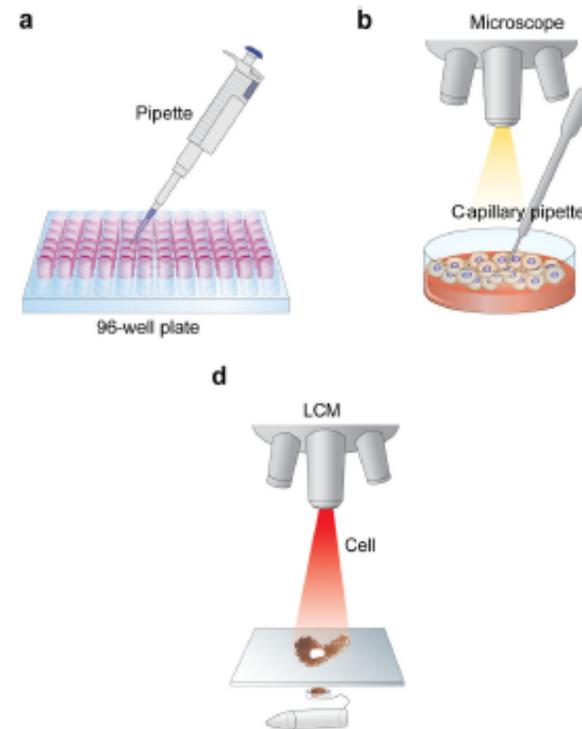
# ¿Qué es el scRNA-seq?

Tecnología que permite la disección de la expresión génica a la resolución de una sola célula.

# Captura de una sola célula

## Técnicas low-throughput

- **Dilución limitante**
  - Ineficiente
- **Capilar guiado por microscopio:**
  - Para muestras en suspensión o en tejido
  - Muestras de pocas o frágiles células
  - Consume tiempo, técnicamente retador
- **Microdissección por captura láser (LCM)**
  - Aislamiento de células específicas en tejido
  - Conserva la relación especial
  - Automático



Hwang, B., Lee, J. H., & Bang, D. (2018). Single-cell RNA sequencing technologies and bioinformatics pipelines. *Experimental & molecular medicine*, 50(8), 1-14. (editado)

# Captura de una sola célula

## Técnicas high-throughput

Ambas técnicas toman células de una suspensión, son automáticas y son *cost-effective*

- **Flow-activated cell sorting (FACS):**

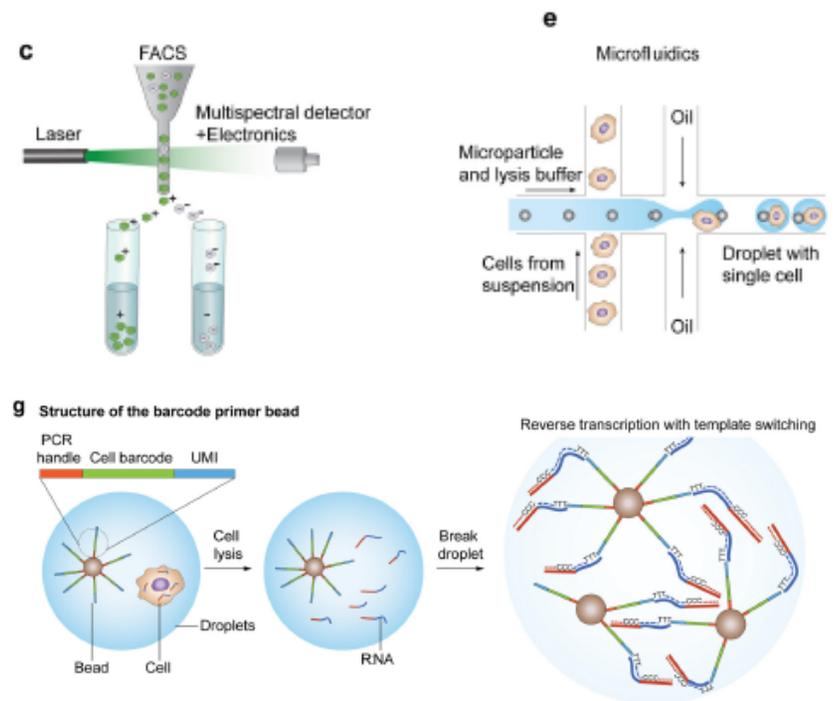
- Se requieren  $\geq 10,000$  células y anticuerpos (muy usado en inmunología).
- Puede dar lugar a pozos vacíos o con *doublets* (dos células).
- Sujeto a contaminación.

- **Microfluidic technology**

- Se requieren  $\geq 1000$  células
- Menor riesgo de contaminación
- Control de fluido y costo bajo de análisis

- **Microdroplet-based microfluidics**

- Uso de beads con *barcodes* únicos
- Encapsulamiento en *droplets*
- Etiqueta UMI



Hwang, B., Lee, J. H., & Bang, D. (2018). Single-cell RNA sequencing technologies and bioinformatics pipelines. *Experimental & molecular medicine*, 50(8), 1-14. (editado)

# Tecnologías y desafíos

## Pasos en común:

- Aislamiento de células
- Lisis celular
- Transcripción en reversa para la primera hebra del cDNA
  - Se estima que solo del 10 al 20% de los transcritos de retrotranscriben
  - Uso de una transcriptasa reversa modificada del virus de la leucemia murina
- Síntesis de la segunda hebra
  - Uso de *poly(A) tailing* o del mecanismo *template-switching*
- Amplificación de cDNA

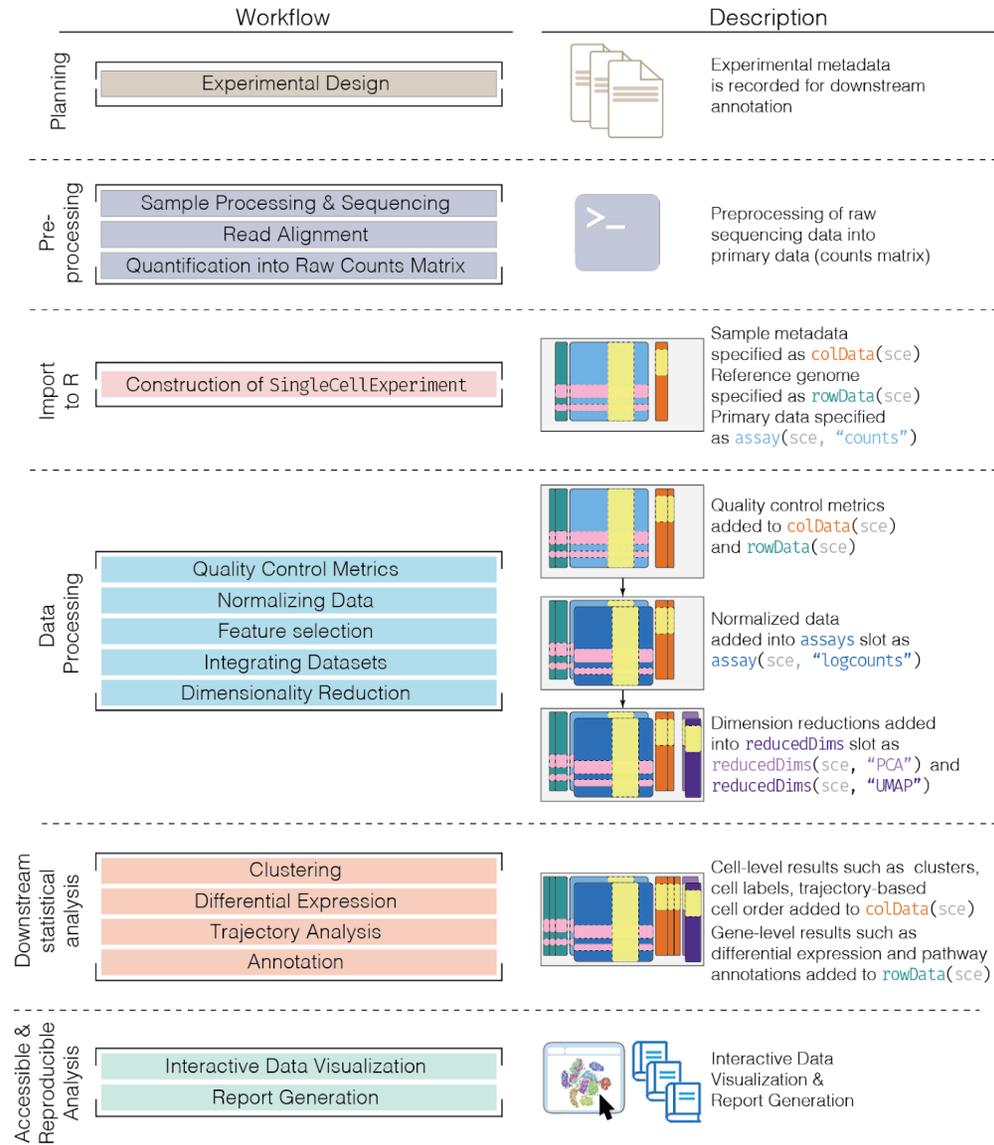
# Tecnologías y desafíos

- La captura de transcritos puede ser tanto del tamaño completo o de los extremos 3' o 5'
- Algunos métodos pueden capturar RNAs poliadenilados y no polineadenilados
- Uso de *spike-ins* y Unique molecular identifiers (identificador molecular único, UMI) para estimar la variación técnica entre células
- La cobertura de transcritos sesgada, la baja eficiencia de captura y la cobertura de secuenciación resultan en datos con más ruido que los de *bulk* RNA-seq.

**TABLE 1** | Summary of widely used scRNA-seq technologies.

Methods	Transcript coverage	UMI possibility	Strand specific	References
Tang method	Nearly full-length	No	No	Tang et al., 2009
Quartz-Seq	Full-length	No	No	Sasagawa et al., 2013
SUPeR-seq	Full-length	No	No	Fan X. et al., 2015
Smart-seq	Full-length	No	No	Ramskold et al., 2012
Smart-seq2	Full-length	No	No	Picelli et al., 2013
MATQ-seq	Full-length	Yes	Yes	Sheng et al., 2017
STRT-seq and STRT/C1	5'-only	Yes	Yes	Islam et al., 2011, 2012
CEL-seq	3'-only	Yes	Yes	Hashimshony et al., 2012
CEL-seq2	3'-only	Yes	Yes	Hashimshony et al., 2016
MARS-seq	3'-only	Yes	Yes	Jaitin et al., 2014
CytoSeq	3'-only	Yes	Yes	Fan H.C. et al., 2015
Drop-seq	3'-only	Yes	Yes	Macosko et al., 2015
InDrop	3'-only	Yes	Yes	Klein et al., 2015
Chromium	3'-only	Yes	Yes	Zheng et al., 2017
SPLIT-seq	3'-only	Yes	Yes	Rosenberg et al., 2018
sci-RNA-seq	3'-only	Yes	Yes	Cao et al., 2017
Seq-Well	3'-only	Yes	Yes	Gierahn et al., 2017
DroNC-seq	3'-only	Yes	Yes	Habib et al., 2017
Quartz-Seq2	3'-only	Yes	Yes	Sasagawa et al., 2018

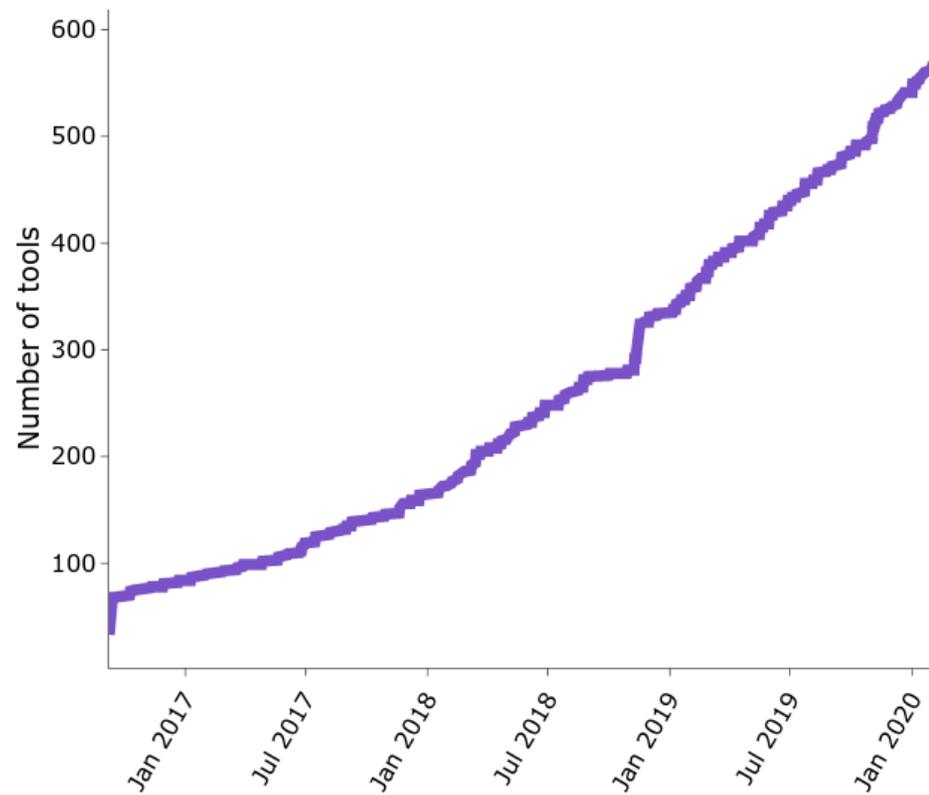
Chen, G., Ning, B., & Shi, T. (2019). Single-cell RNA-seq technologies and related computational data analysis. *Frontiers in genetics*, 10, 317.



# Procesamiento de datos

1. Preprocesamiento
2. Control de Calidad
3. Normalización
4. Imputación
5. Selección de atributos
6. Reducción de dimensiones
7. Integración de conjuntos de datos

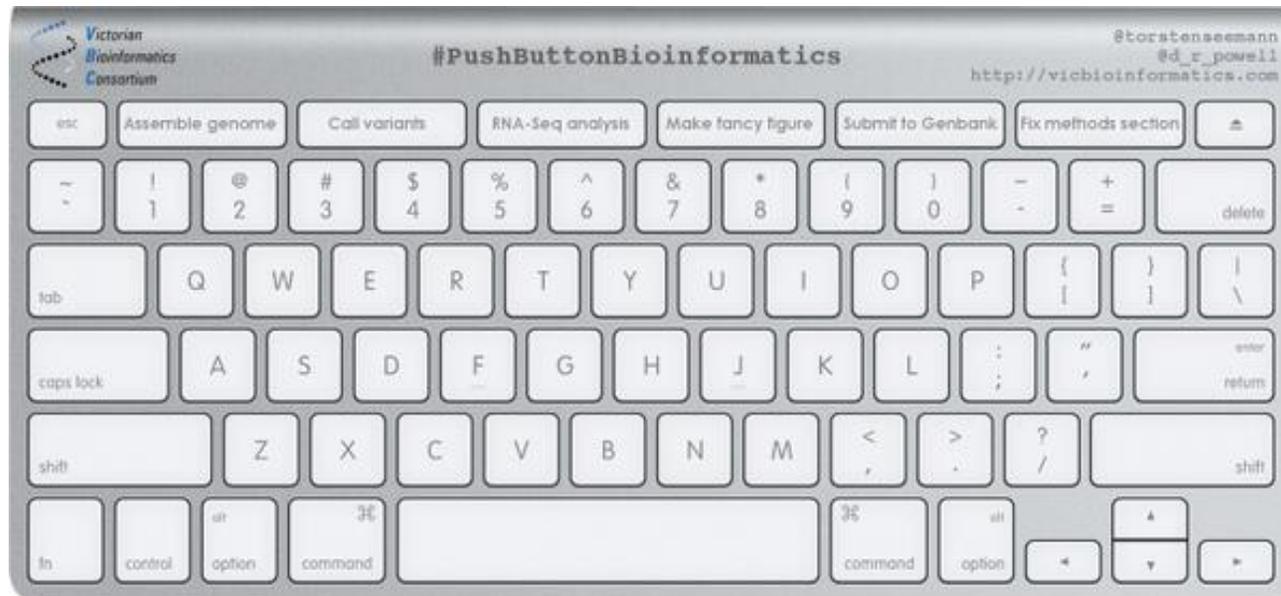
# Number of scRNA-seq tools over time



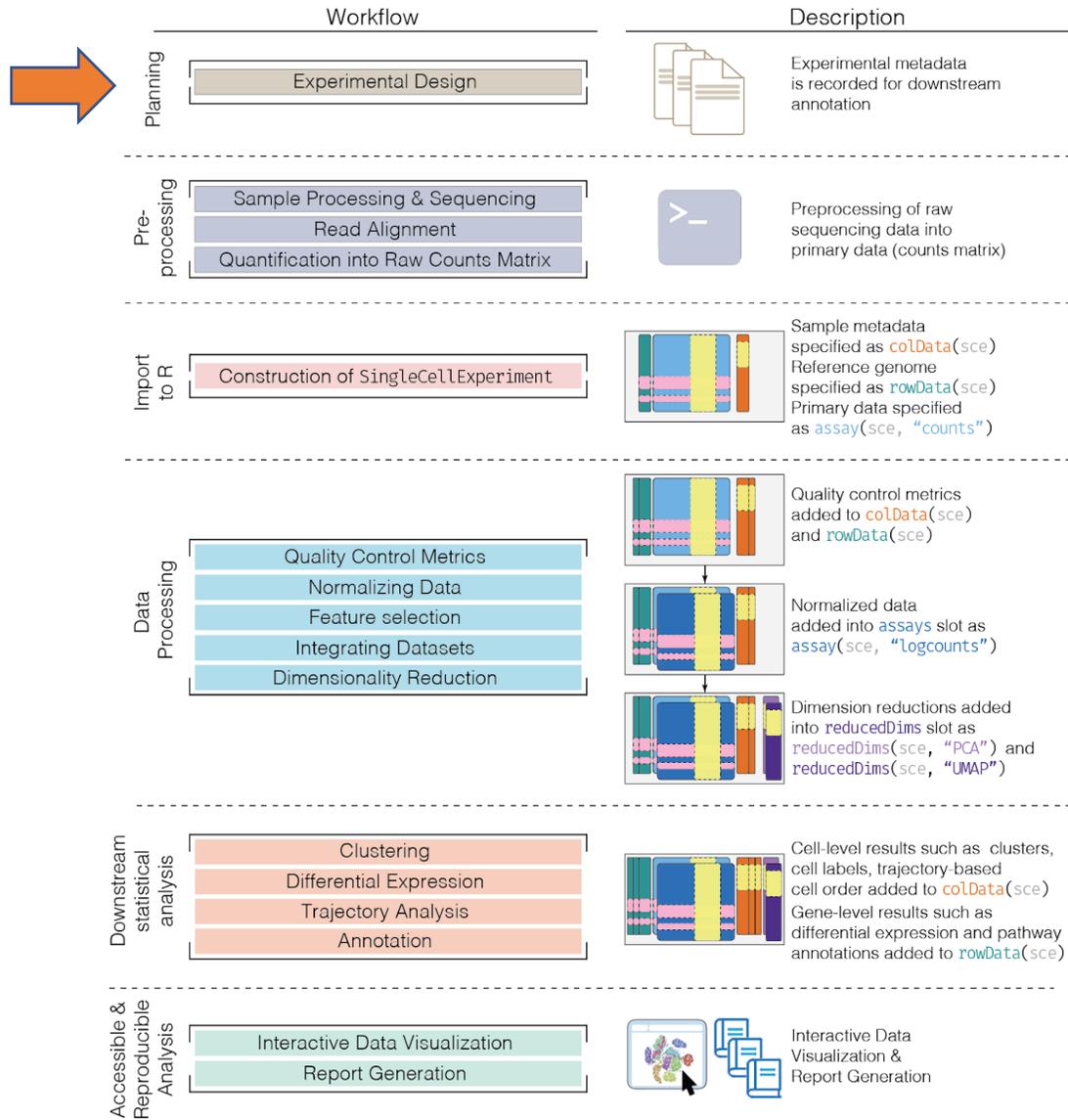
Modificado de [Peter Hickey](#)

<https://www.scrna-tools.org/analysis>

No hay recetas!!! Hay flujos de trabajo !



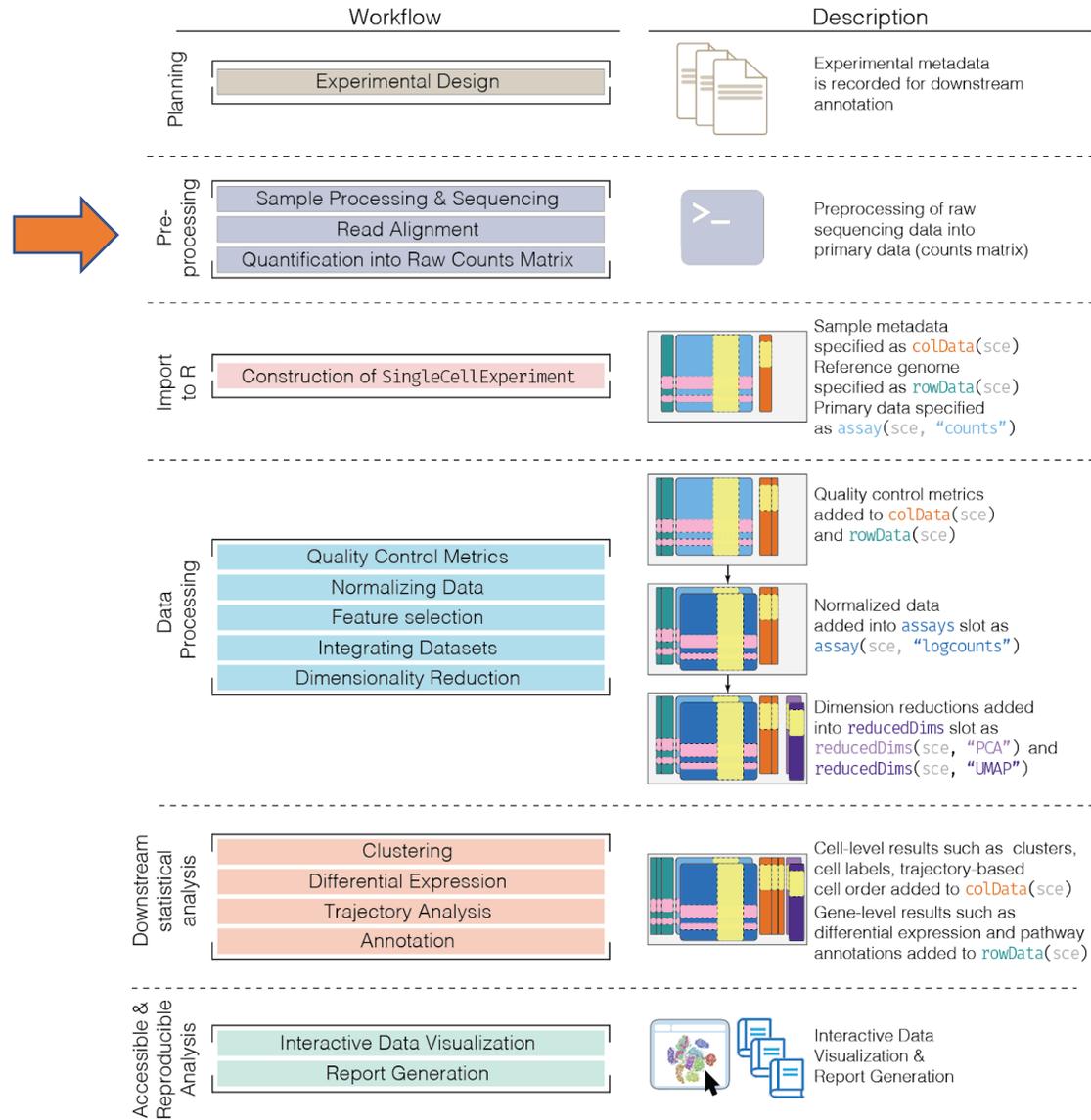
Modificado de [Peter Hickey](#)



Modificado de [Peter Hickey](#)

# Controles! Controles! Controles!

- Necesitamos controles positivos y negativos de las células y de las muestras.
- Necesitamos varias replicas biológicas.
- Los grupos experimentales no son batches. Los batches son técnicos.
- Las células individuales no pueden ser tratadas como replicas.
- Aquí los ahorros nos pueden salir muy caros después.



Modificado de [Peter Hickey](#)

# Procesamiento de datos: Preprocesamiento

- QC de secuenciación (FASTQC)  
<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>
- Alineamiento a transcriptoma de referencia
- Cuantificación por célula y por gen para la generación de la matriz de cuentas de expresión

# Procesamiento de datos: Preprocesamiento

- Herramientas para hacer la matriz de cuentas:
  - CellRanger para datos 10X
  - scPipe para CEL-Seq2 y otros tipos de datos
  - Kallisto y alevin para datos de droplets

# Cómo se ven los datos?

## Cuentas de Genes

	<b>Cell 1</b>	<b>Cell 2</b>	...	<b>Cell N</b>
<b>Gene 1</b>	0	1	...	0
<b>Gene 2</b>	1	3	...	0
...	...	...	...	...
<b>Gene M</b>	2	2		4

Modificado de [Peter Hickey](#)

# Cómo se ven los datos?

## Información de las células

	<b>Barcode</b>	<b>Donor</b>	...	<b>Treatment</b>
<b>Cell 1</b>	ACTGTA	D1	...	Drug
<b>Cell 2</b>	TGCATA	D1	...	Control
...	...	...	...	...
<b>Cell N</b>	CCTATA	D6		Drug

# Cómo se ven los datos?

## Información de los Genes

	<b>ID</b>	<b>Symbol</b>	...	<b>Chromosome</b>
<b>Gene 1</b>	ENSG00000155816	FMN2	...	1
<b>Gene 2</b>	ENSG00000229807	XIST	...	X
...	...	...	...	...
<b>Gene M</b>	ENSG00000139618	BRCA2		13

Modificado de [Peter Hickey](#)

# Cómo queremos que se vean los datos normalizados?

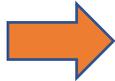
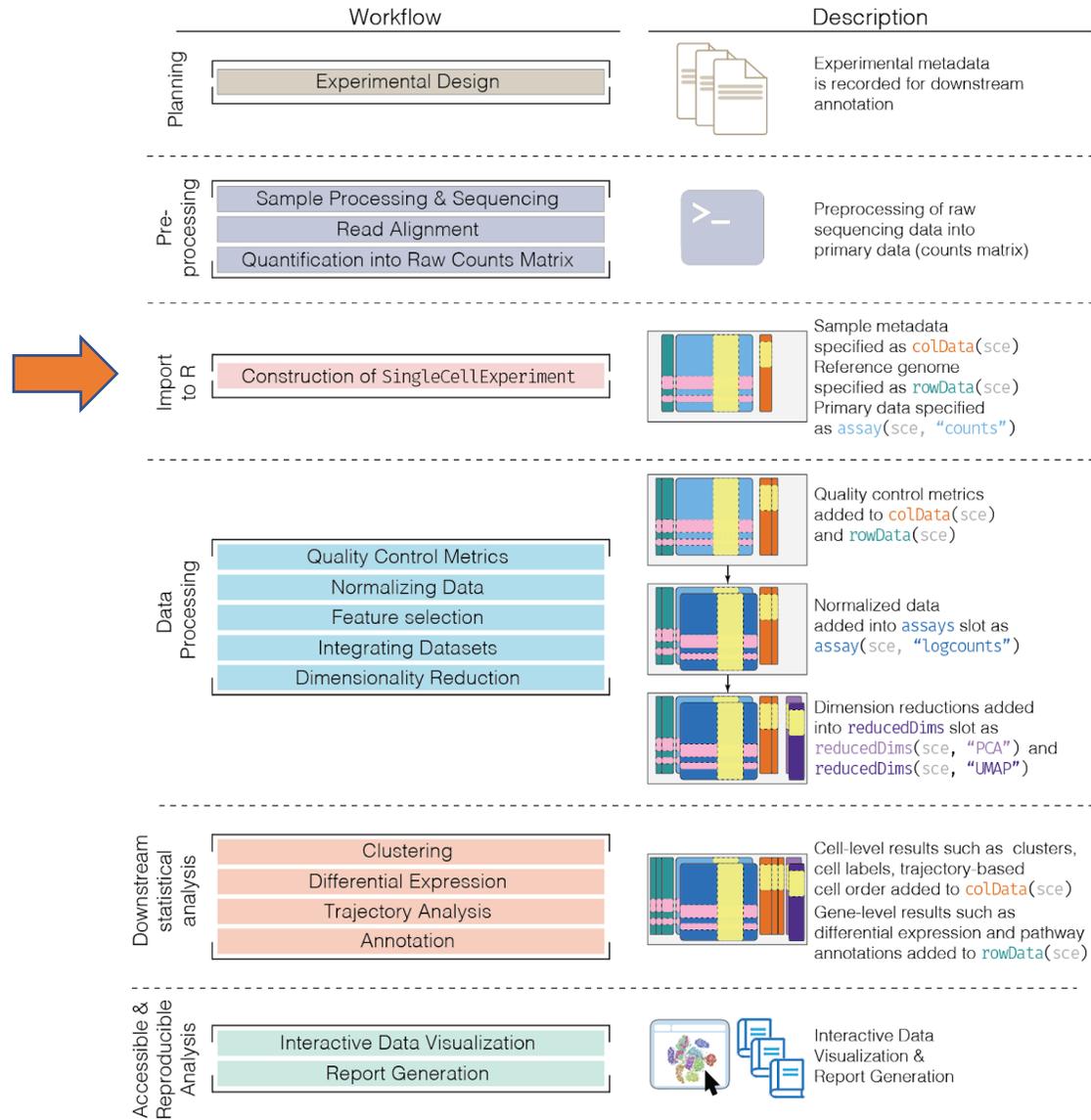
Normalización log de expression de genes

	<b>Cell 1</b>	<b>Cell 2</b>	...	<b>Cell N</b>
<b>Gene 1</b>	0	0.6	...	0
<b>Gene 2</b>	0.3	0.8	...	0
...	...	...	...	...
<b>Gene M</b>	0.35	0.67		2.1

# Cómo se verán nuestros datos cuando hagamos reducción de dimensiones?

	<b>PCA 1</b>	<b>PCA 2</b>	...	<b>PCA K</b>
<b>Cell 1</b>	0.93	1.28	...	0.03
<b>Cell 2</b>	0.32	1.22	...	0.09
...	...	...	...	...
<b>Cell N</b>	-0.66	1.00		0.15

	<b>t-SNE 1</b>	<b>t-SNE 2</b>
<b>Cell 1</b>	1.24	8.93
<b>Cell 2</b>	-0.33	7.85
...	...	...
<b>Cell N</b>	0.46	3.41



Modificado de [Peter Hickey](#)

# Hay que cargar nuestro experimento en R

Necesitamos objetos que nos permitan organizar toda la información

Información de genes

	<b>Cell 1</b>	<b>Cell 2</b>	...	<b>Cell N</b>
<b>Gene 1</b>	0	1	...	0
<b>Gene 2</b>	1	3	...	0
...	...	...	...	...
<b>Gene M</b>	2	2		4

Información de células

	<b>Barcode</b>	<b>Donor</b>	...	<b>Treatment</b>
<b>Cell 1</b>	ACTGTA	D1	...	Drug
<b>Cell 2</b>	TGCATA	D1	...	Control
...	...	...	...	...
<b>Cell N</b>	CCTATA	D6		Drug

Información de los genes

	<b>ID</b>	<b>Symbol</b>	...	<b>Chromosome</b>
<b>Gene 1</b>	ENSG00000155816	FMN2	...	1
<b>Gene 2</b>	ENSG00000229807	XIST	...	X
...	...	...	...	...
<b>Gene M</b>	ENSG00000139618	BRCA2		13

Modificado de [Peter Hickey](#)

# Hay que cargar nuestro experimento en R

Necesitamos objetos que nos permitan organizar toda la información

Información de genes

	<b>Cell 1</b>	<b>Cell 2</b>	...	<b>Cell N</b>
<b>Gene 1</b>	0	1	...	0
<b>Gene 2</b>	1	3	...	0
...	...	...	...	...
<b>Gene M</b>	2	2		4

Información de células

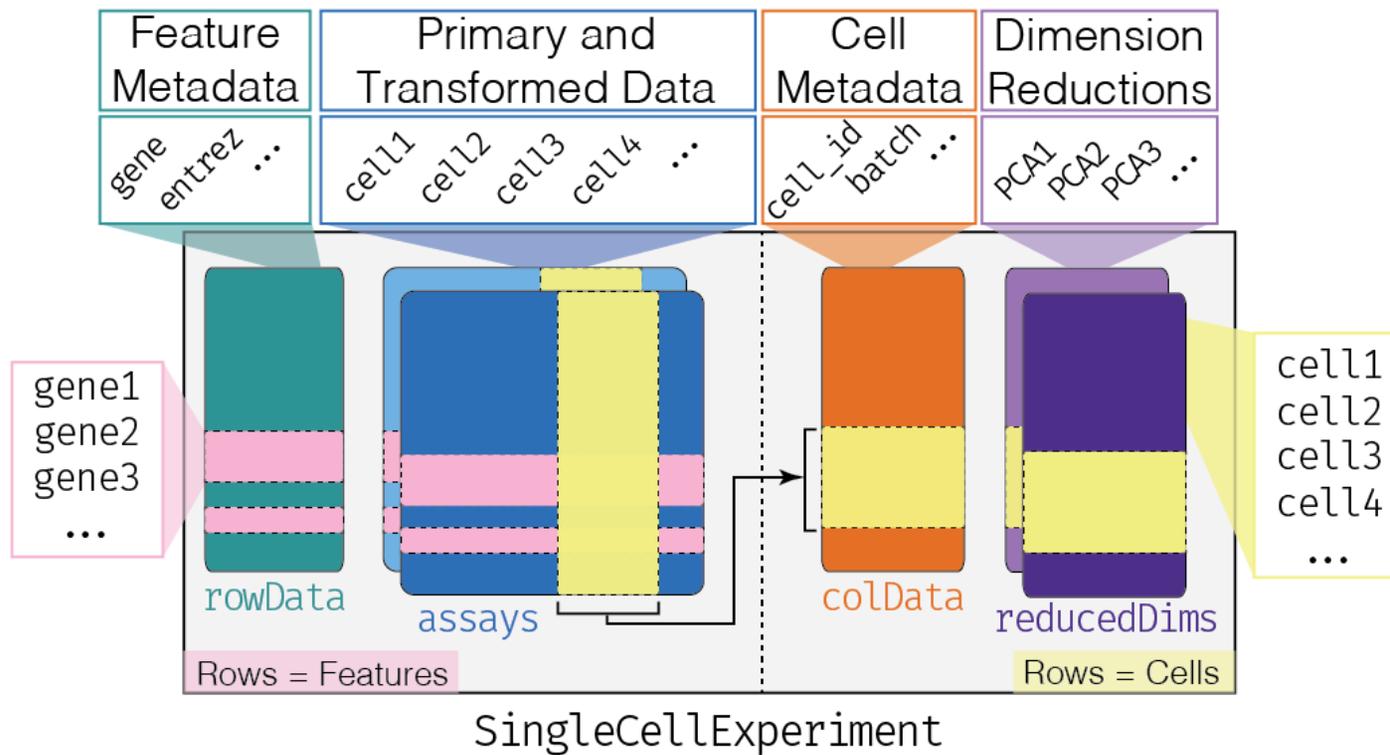
	<b>Barcode</b>	<b>Donor</b>	...	<b>Treatment</b>
<b>Cell 1</b>	ACTGTA	D1	...	Drug
<b>Cell 2</b>	TGCATA	D1	...	Control
...	...	...	...	...
<del><b>Cell N</b></del>	<del>CCTATA</del>	<del>D0</del>		<del>Drug</del>

Información de los genes

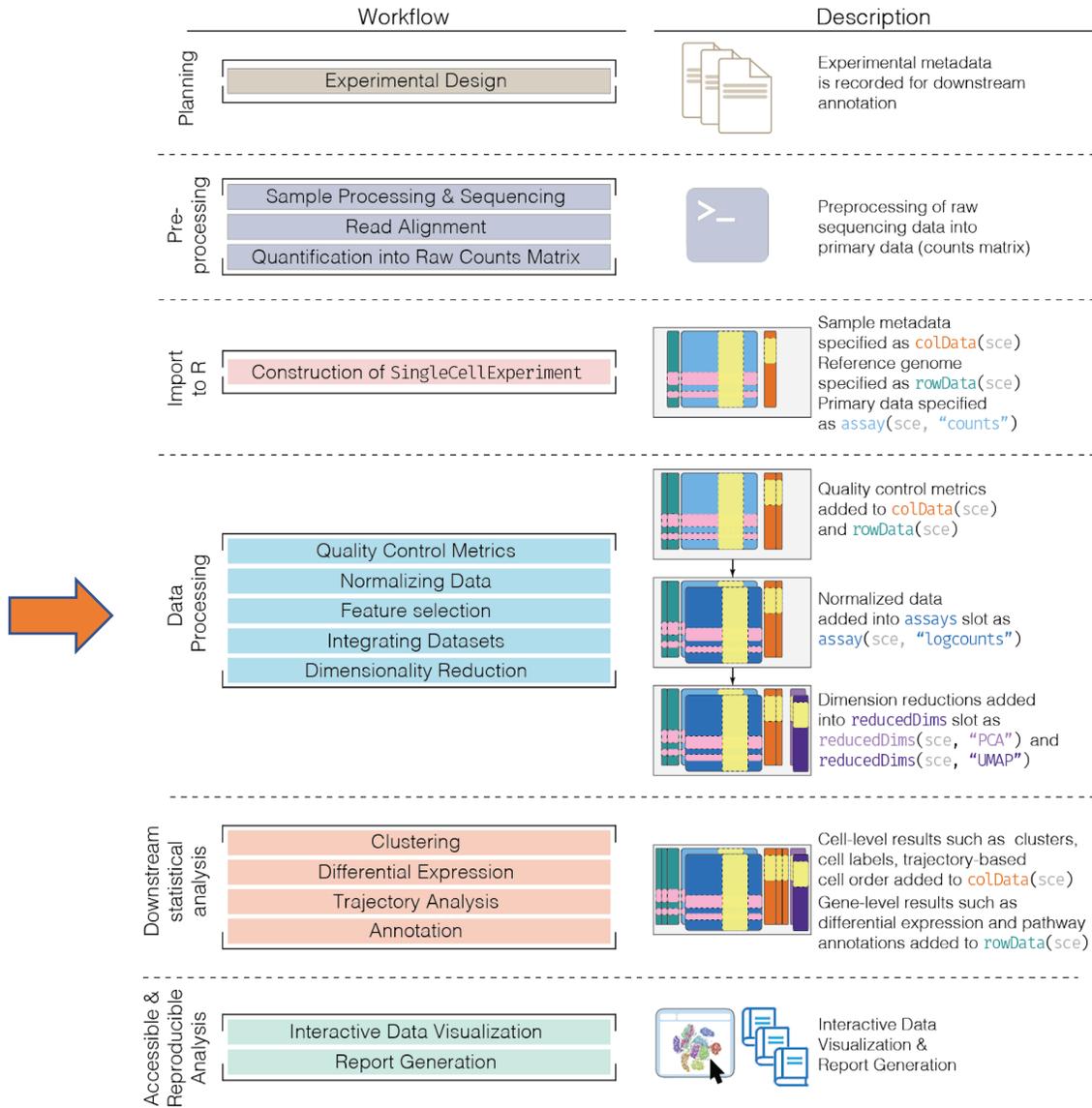
	<b>ID</b>	<b>Symbol</b>	...	<b>Chromosome</b>
<b>Gene 1</b>	ENSG00000155816	FMN2	...	1
<del><b>Gene 2</b></del>	<del>ENSG00000220007</del>	<del>XIST</del>	...	<del>X</del>
...	...	...	...	...
<b>Gene M</b>	ENSG00000139618	BRCA2		13

Modificado de [Peter Hickey](#)

# Hay que cargar nuestro experimento en R



Modificado de [Peter Hickey](#)



Modificado de [Peter Hickey](#)

# Procesamiento de datos: QC



Método de single-cell que esté de moda



# Procesamiento de datos: QC

Queremos eliminar

- células que sufrieron daño durante la disociación
- Células con poca cobertura (el RNA pudo haberse perdido en algún paso anterior)
- La ausencia de células y los *doublets*

# Procesamiento de datos: QC

Queremos eliminar

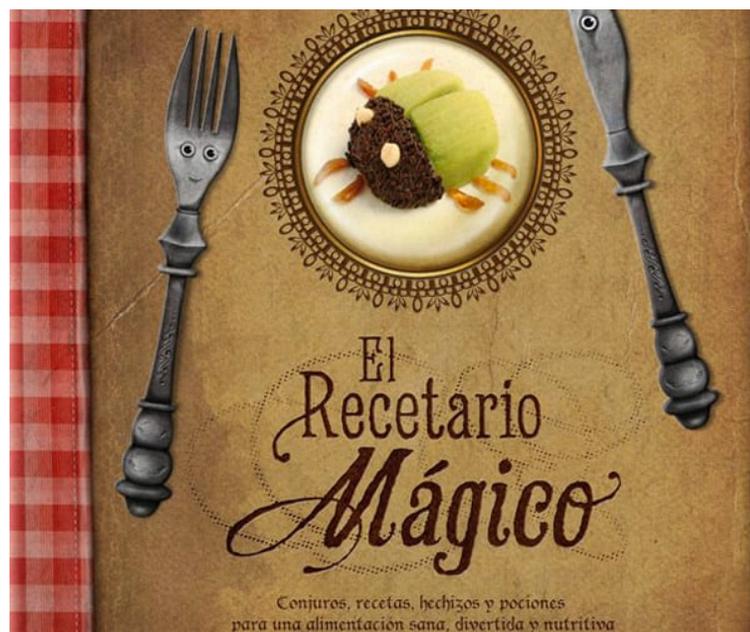
- células que sufrieron daño durante la disociación
- Células con poca cobertura (el RNA pudo haberse perdido en algún paso anterior)
- La ausencia de células y los *doublets*

Esto se suele manifestar como

- Células con un tamaño pequeño de librería (cuentas totales)
- Células con pocos genes expresados
- Una alta proporción de lecturas mitocondriales
- Células con alta proporción de *spike-ins*

No hay recetas!! Cada grupo de datos será distinto

Los umbrales de calidad de un experiment no necesariamente sirven en otro.



Modificado de [Peter Hickey](#)

# Procesamiento de datos: QC

Qué pasa si no hacemos bien nuestro QC?

- Vamos a ver grupos de "datos" que no serán nada biológico, pero podemos pensar que si.
- Veremos poblaciones heterogéneas que en verdad solo son ruido.
- Algunos genes se verán super expresados pero no significa nada.

# Procesamiento de datos: Normalización

Se observan diferencias de cobertura entre librerías. Esto parte de diferencias en la captura de cDNA o de la eficiencia de la amplificación. También hay un efecto por el número de distintos tipos celulares.

- Normalización por tamaño de librería
- Normalización por deconvolución
- Normalización refinada por *spike-ins/UMIs*, cuando están presentes

# Procesamiento de datos: Imputación

Necesario para vencer el reto de escasez de datos (los ceros).

**Por qué hay tantos ceros?**

	<b>Cell 1</b>	<b>Cell 2</b>	...	<b>Cell N</b>
<b>Gene 1</b>	0	0.6	...	0
<b>Gene 2</b>	0.3	0.8	...	0
...	...	...	...	...
<b>Gene M</b>	0.35	0.67		2.1

# Procesamiento de datos: Imputación

- Se han elaborado modelos que toman en cuenta los ceros en cuentas (***zero-inflated models***)
  - El nivel de *zero-inflated* depende del ensayo o el protocolo
  - Un buen modelo puede ser **ensayo dependiente**, lo que funcionó para un proyecto no necesariamente funcionen en otros.
- Los modelos de imputación han demostrado **generar falsos positivos** y **disminuir la reproducibilidad** de marcadores célula específicos
- SAVER es un modelo basado en bayesiana diseñados para datos con UMIs y MAGIC construye una gráfica basada en afinidad de markov.

# Procesamiento de datos: Selección de atributos

La meta es identificar genes con información biológica útil y quitar genes con ruido además de reducir el tamaño del *dataset*.

*Por qué?*

Por que queremos ver diferencias entre genes, células, si hay tipos celulares interesantes, etc.

# Procesamiento de datos: Selección de atributos

La meta es identificar genes con información biológica útil y quitar genes con ruido además de reducir el tamaño del *dataset*.

- Hay atributos (features) que darán información biológica (+)
- Otros atributos serán solo ruido (-)

# Procesamiento de datos: Selección de atributos

La meta es identificar genes con información biológica útil y quitar genes con ruido además de reducir el tamaño del *dataset*.

- Un enfoque sencillo es seleccionar los genes más variables acorde a su expresión.
  - Las transformaciones logarítmicas no logran una estabilización de la variancia
  - Uso de métrica ***deviance***, la cual cuantifica qué tan bien un gen dado se ajusta a un modelo nulo de expresión constante en todas las células. Se basa en UMIs.
- Podemos usar lo que sabemos de nuestro sistema y enfocarnos en genes relevantes.
- También existen métodos basados en *spike-ins* y en *dropouts*

# Procesamiento de datos: Reducción de dimensiones

Objetivo: identificar grupos de células parecidas.

Problema: Nuestros datos tienen demasiados atributos como para hacerlo de forma eficiente (muchas células, muchos genes, más nuestro diseño experimental).

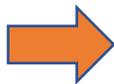
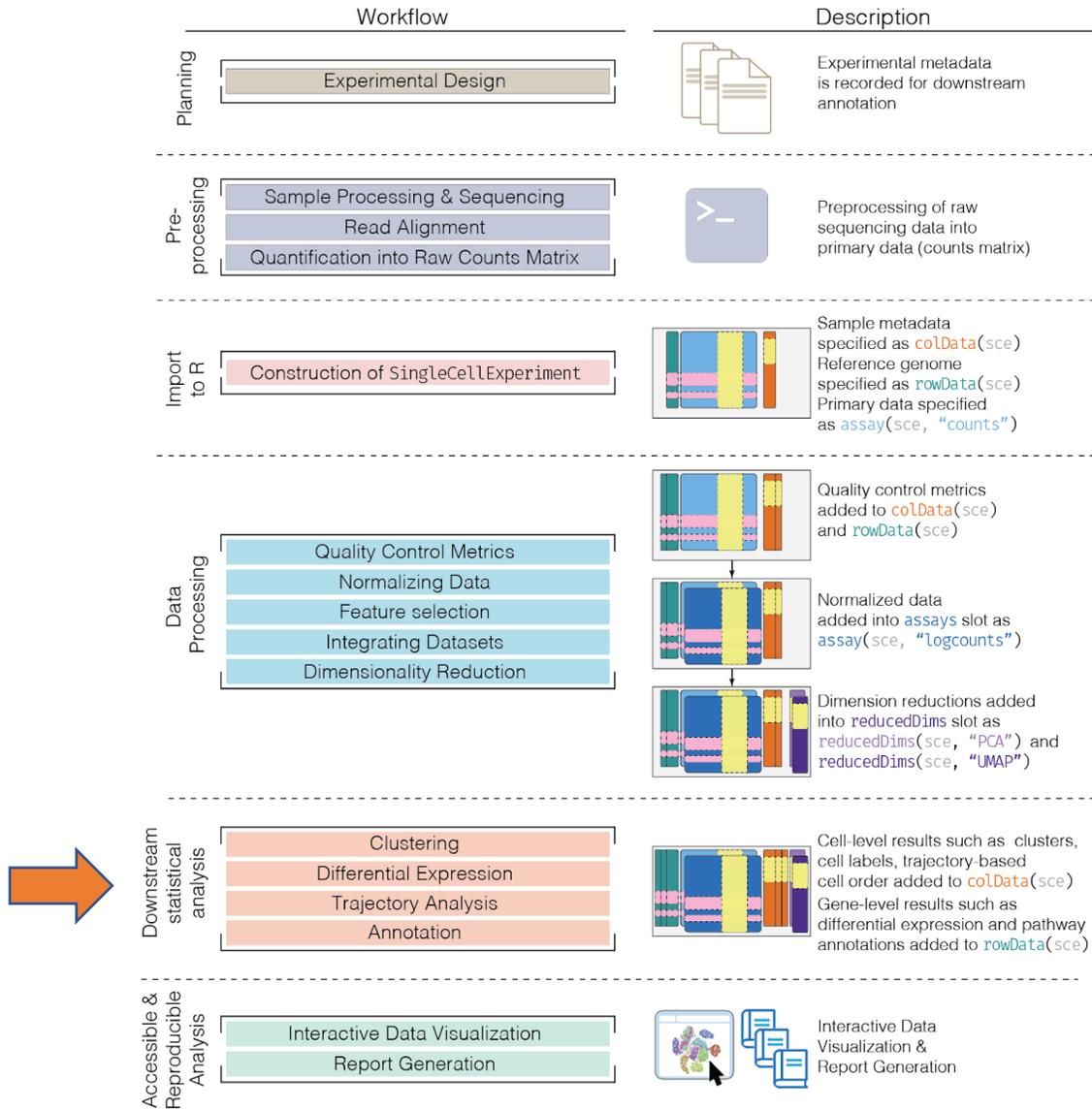
# Procesamiento de datos: Reducción de dimensiones

- Genes distintos pueden estar correlacionados si el mismo proceso biológico los afecta.
- La meta es crear representaciones de los datos en dimensiones pequeñas pero que preservan una estructura significativa.
- Análisis de Componentes Principales (PCA)
- Para propósitos de visualización también existen:
  - T-distributed stochastic neighbor embedding(t-SNE)
  - Uniform Manifold Approximation and Projection (UMAP)

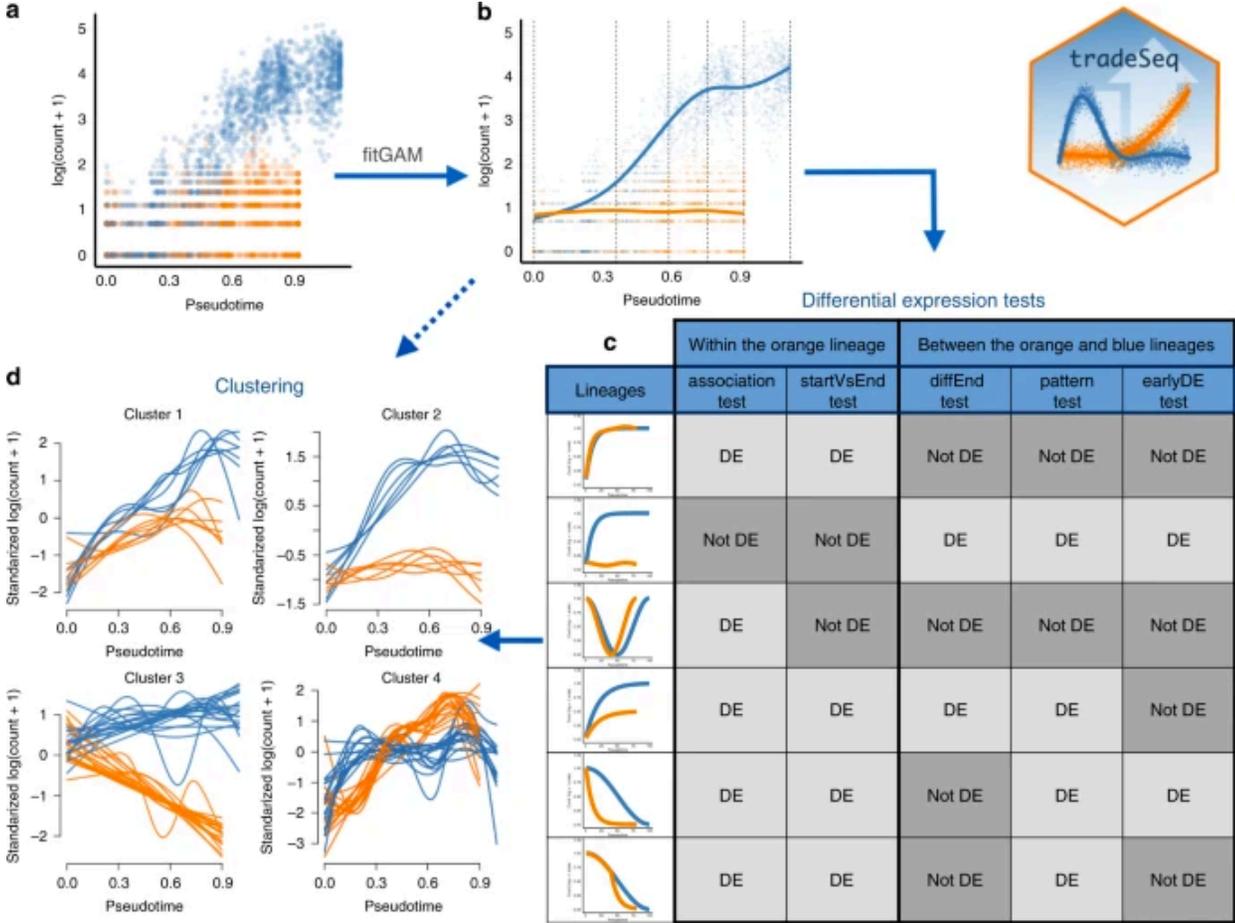
# Procesamiento de datos: Integración de conjuntos de datos

En proyectos grandes de scRNA-seq es necesario generar los datos a través de múltiples *batches*.

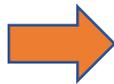
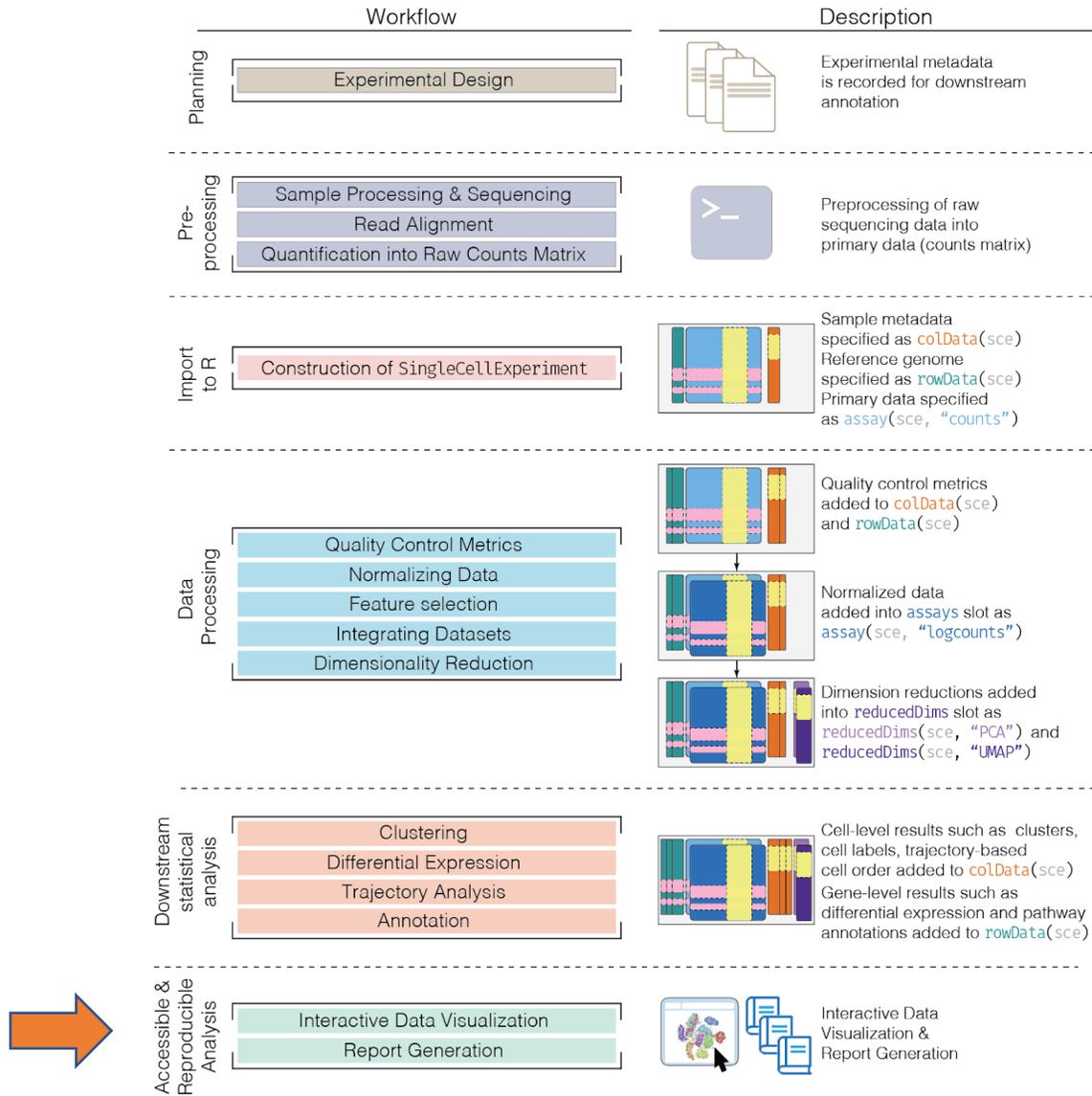
- **La presencia de batches se puede detectar en un PCA** de valores de expresión logarítmicos
- Se han desarrollado métodos *bespoke* para la corrección de *batches* en datos de scRNA-seq que no requieren información a priori sobre la composición de la población
- Otros métodos que tratan de corregir los efectos de múltiples *batches* son **MNN** (*mutual nearest neighbor*) y **kBET** (*k-nearest neighbor batch effect test*)



# Análisis de trayectorias



<https://www.nature.com/articles/s41467-020-14766-3>

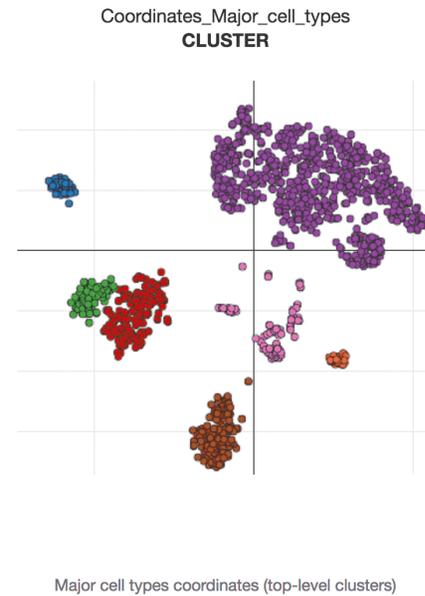
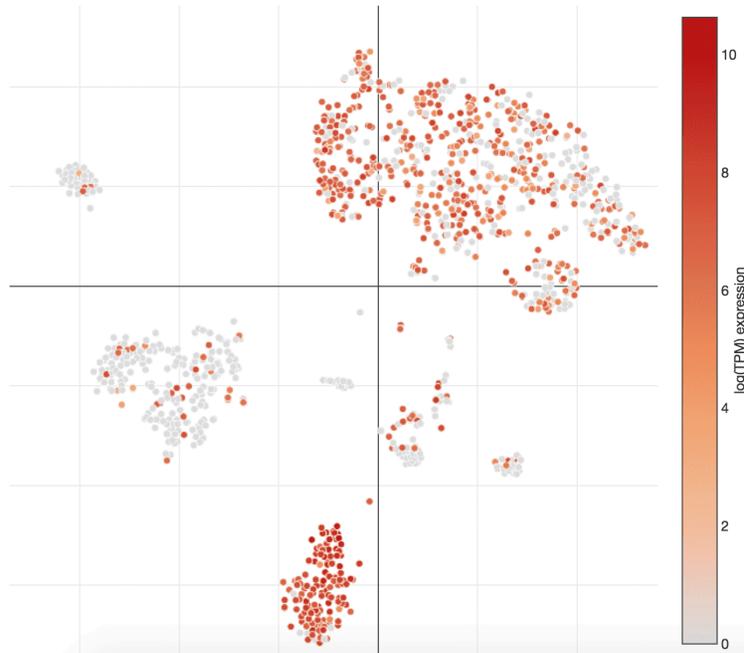


# Gene Expression for *Gad1*

Summary Explore Download Analysis Settings

Gad1 Distribution Scatter

View Options



### Load cluster

Coordinates\_Major\_cell\_types

### Select annotation

CLUSTER

### Subsampling threshold

All Cells

### Distribution

#### Plot Type

Violin Plot

#### Data Points

All

Toggle Annotations

> Scatter

> Heatmap