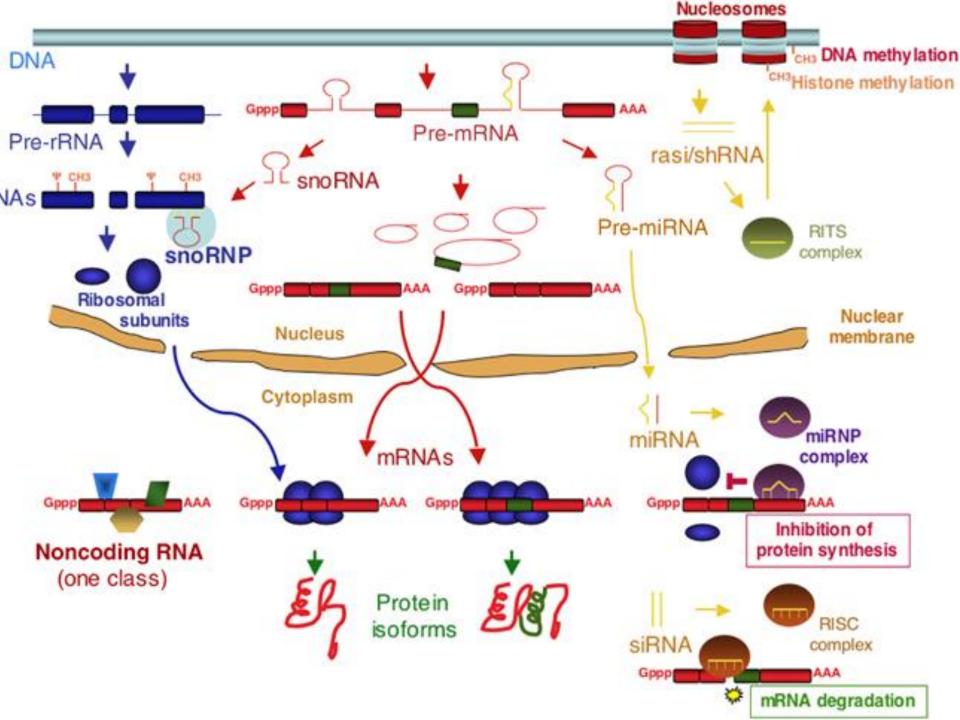
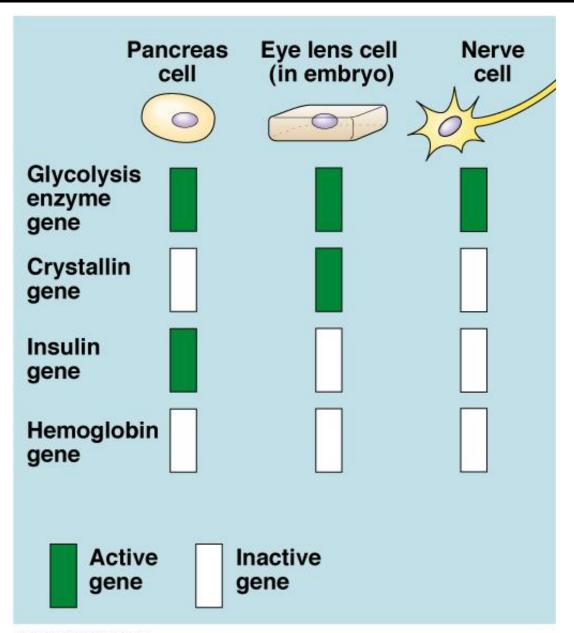
Estudio de Transcriptomas

Lo que el genoma tiene para decir



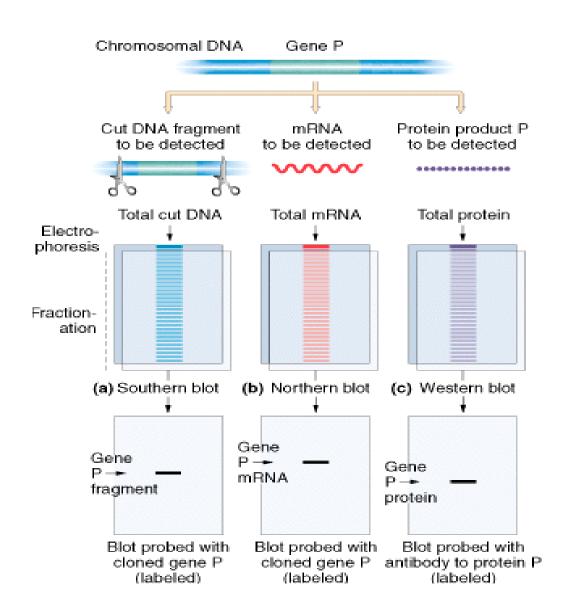
Expresión diferencial

- No todos los genes se expresan
 - existen distintos
 perfiles de ARNs y
 proteínas
- La expresión diferencial:
 - durante el desarrollo y diferenciación celular
 - en la homeostasis



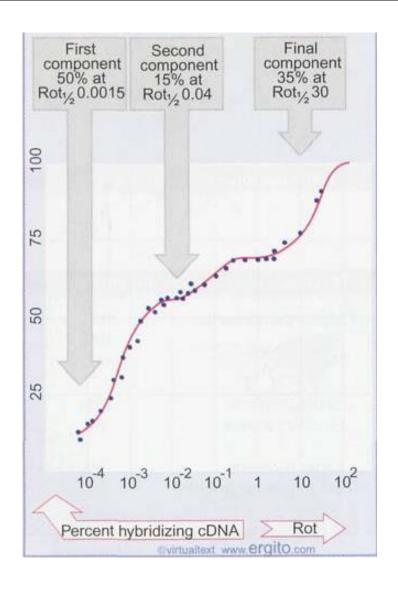
Estudios clásicos de expresión:

- Se focalizan en un número limitado de genes
 - Northern blot
 - RT PCR.

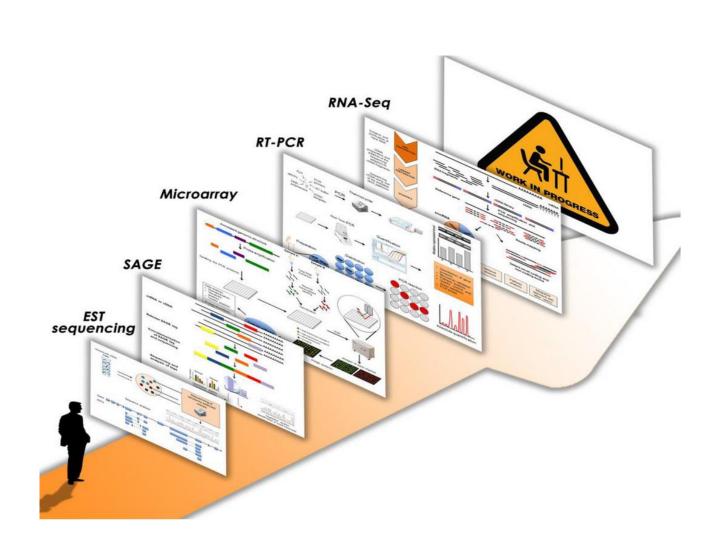


Curvas Rot

- Los ARN están presentes a diferentes abundancias
- Tipos de ARN (Okamuro & Goldberg)
 - Superabundantes
 - 15-90% de la masa de mRNA
 - <10 transcriptos diferentes
 - >5000 moléculas por célula
 - Abundantes
 - 50-75% de la masa de mRNA
 - ~200-1000 transcriptos diferentes (5% de la diversidad de mensajeros)
 - 500-2500 moléculas por célula
 - Raros
 - <25% de la masa
 - 95% de la diversidad
 - 1-10 por célula



Avances metodológicos

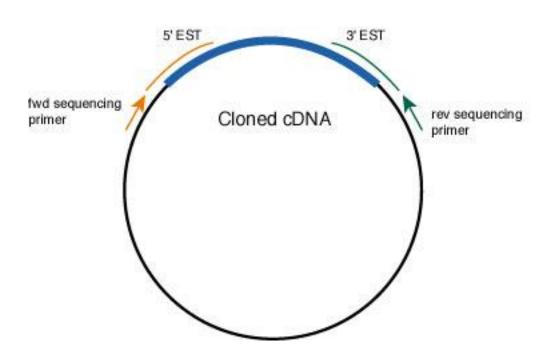


Expressed Sequence Tag (EST)

- Adams M. et al. Science 252:1651-1656 en 1991
- Lecturas de secuencias cortas, al azar y de un único pasaje por clon derivadas de bibliotecas de cDNA
- Alternativa mas barata al secuenciado genómico
- Propensas a errores y requieren un procesamiento bioinformático intensivo para dar información global del transcriptoma.
- Descubrimiento de genes nuevos, y determinación de estructura génica
- Se pueden usar para identificar el resto del gen en el genoma (anotación de genes nuevos)
- Pueden ser marcadas para localizarse en el cromosoma (marcadores físicos)(FISH).
- Verificación de las secuencias UTRs

Producción de ESTs

cDNAs parci	iales	
		AAAAAAAAAAA
		 ттттттт
		

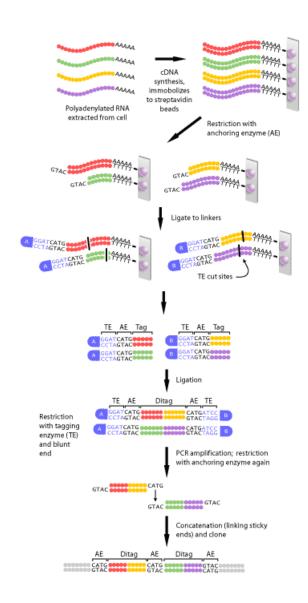


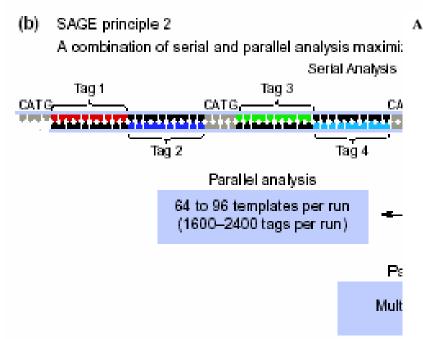
Problemas de los ESTs

- Las secuencias pueden contener errores tanto de cambio de base como de inserción/deleción
- Genes de baja expresión son difíciles de encontrar en la biblioteca de cDNA (normalización, secuenciado en masa)
- cDNAs difíciles de clonar
- No se representa el transcrito completo:
 - Solo un 11% de secuencias completas
 - 65% de regiones 3'
 - 25% de regiones 5'

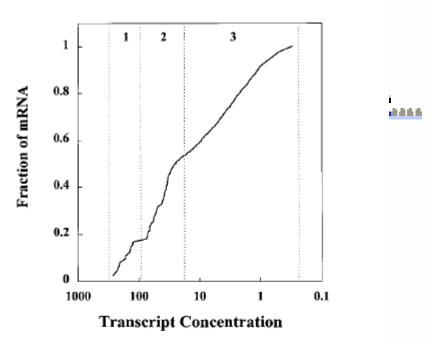
SAGE (Serial analysis of gene expression)

- 1995
- Secuenciado de pequeños tags
- Cada gen es identificado por un tag único:
 - Análisis cuantitativo de la expresión
- Requiere de un genoma secuenciado sobre el cual alinear los tags
- Búsqueda de nuevos genes









	Virtual Rot (SAGE)		Ret (Reassociation)	
Component	%mRNA	Copies/cell	%mRNA	Copies/cell
1	17	180	23	200
2	38	40	51	30
3	45	2.5	26	1.5

Figure 3. Virtual Rot

В

Precesamiento de los datos

CATGACCCACGAGCAGGGTACGATGATCATGGAAACCTATGCACCTTGGGTAGCACATG
TAG 1 TAG 2 TAG 3 TAG 4

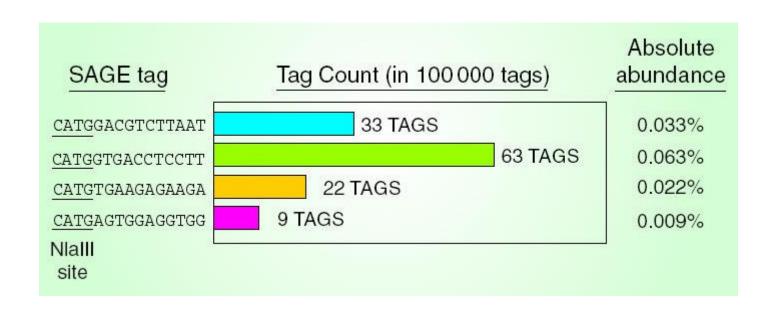
TAGGACGAGGGTGGACAATGCTCATG
TAG 5 TAG 6

Tag_Sequence	Count
ATCTGAGTTC	1075
GCGCAGACTT	125
TCCCCGTACA	112
TAGGACGAGG	92
GCGATGGCGG	91
TAGCCCAGAT	83
GCCTTGTTTA	80
GCGATATTGT	66
TACGTTTCCA	66
TCCCGTACAT	66
TCCCTATTAA	66
GGATCACAAT	55
AAGGTTCTGG	54
CAGAACCGCG	50
GGACCGCCCC	48

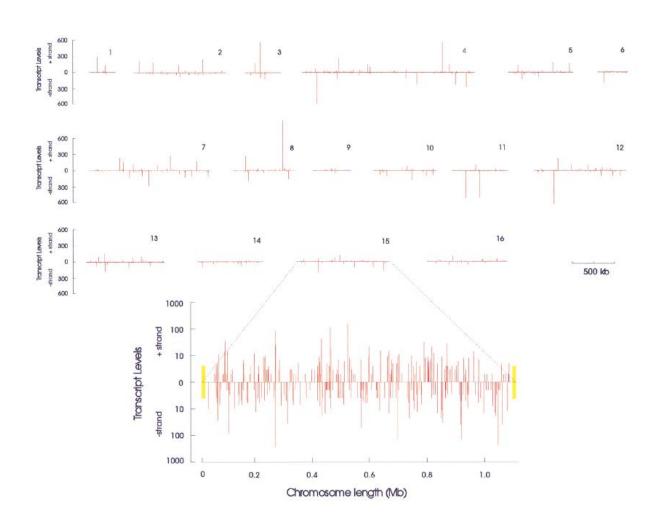
Tag_Sequence Count Gene Name					
ATATTGTCAA	5	translation elongation factor 1 gamma			
AAATCGGAAT	2	T-complex protein 1, z-subunit			
ACCGCCTTCG	1	no match			
GCCTTGTTTA	81	rpa1 mRNA fragment for r ribosomal protein			
GTTAACCATC	45	ubiquitin 52-AA extension protein			
CCGCCGTGGG	9	SF1 protein (SF1 gene)			
TTTTTGTTAA	99	NADH dehydrogenase 3 (ND3) gene			
GCAAAACCGG	63	rpL21			
GGAGCCCGCC	45	ribosomal protein L18a			
GCCCGCAACA	34	ribosomal protein S31			
GCCGAAGTTG	50	ribosomal protein S5 homolog (M(1)15D)			
TAACGACCGC	4	BcDNA.GM12270			

Análisis cuantitativo

 El número de veces que se secuencia un determinado tag puede ser usado como una medida de la cantidad relativa de mRNA original



Characterization of the Yeast Transcriptome



Problemas del SAGE

- Los tags son MUY pequeños (13 14 bp).
 - Si el gen no es conocido es casi imposible reconocer motivos funcionales
 - Especificidad: Múltiples genes pueden compartir el mismo tag. Esto ha sido mejorado aumentando el largo.
- Las enzimas de tipoIIS (generalmente la BsmFI) puede dar largos un poco diferentes lo cual genera problemas al interpretar los di-tags
- Algunos mRNAs pueden no tener el sitio de reconocimiento de la enzima
- Técnicamente complicado

Análisis a nivel masivo: Microarreglos de ADN

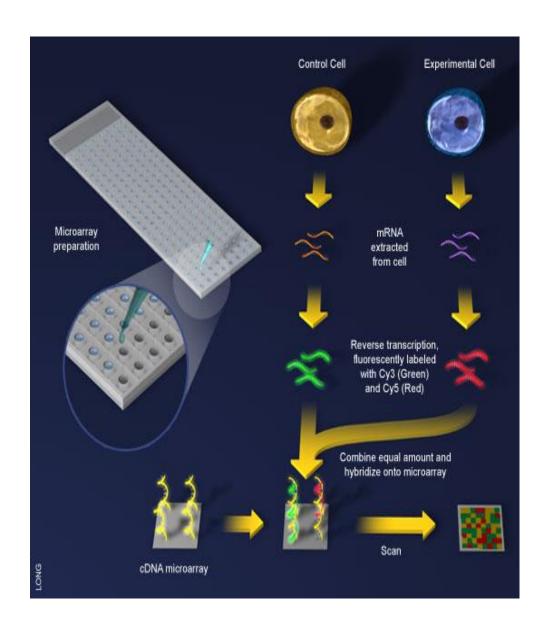
- Permiten el estudio simultaneo y comparativo del transcriptoma
- Siguen en mismo principio que el northern blot pero paralelizado masivamente
- Aunque la presencia de un mensajero no implica presencia de la proteína un cambio en la cantidad de mensajero esta mostrando una respuesta fisiológica:
 - Análisis comparativos de diferentes situaciones
- Permiten tener un pantallazo del transcriptoma a nivel global:
 - vías metabólicas completas
- Permiten también realizar inferencias funcionales (clustering)
- Aplicaciones clínicas:
 - Farmacogenómica
 - Marcadores pronósticos

Tipos de Arrays

De cDNA

- Se construyen imprimiendo un biblioteca de cDNA en una placa
- Son de hibridación comparativa ("de dos canales")
- Generalmente personalizados (diseño personal de la biblioteca a imprimir)
- Fragmentos largos de ADN (productos de PCR)
- Alrededor de 10.000 "features" por cm2
- De oligos (Chips de ADN)
 - Fragmentos cortos de ADN (Alrededor de 25pb)
 - Muy alta densidad (millones de features)
 - Son impresos "in situ":
 - Fotolitografía
 - Eliminan la necesidad de construir una biblioteca
 - De un solo canal (un chip por muestra)

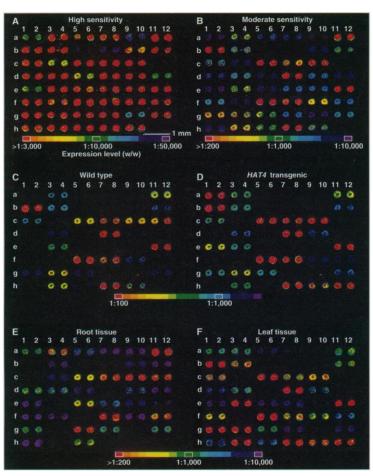
cDNA arrays



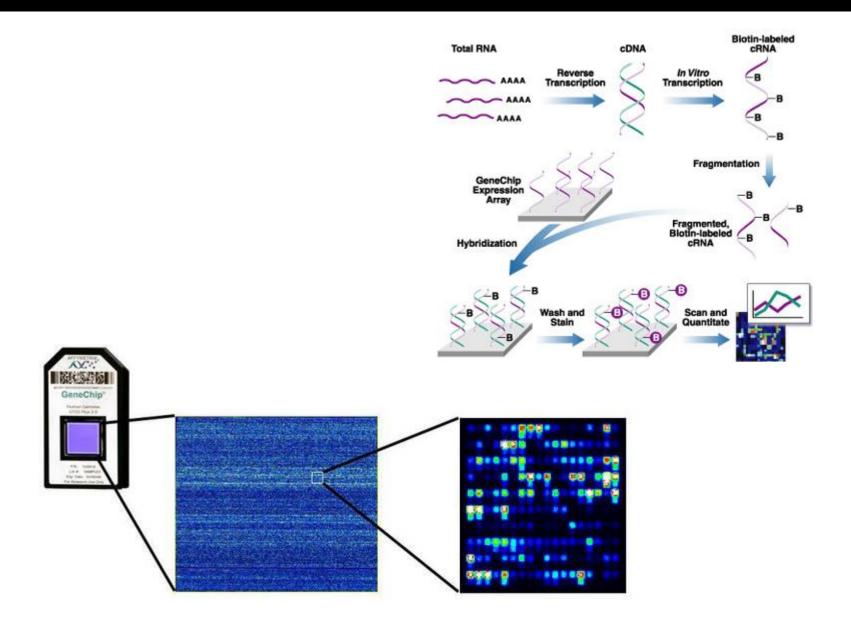
Quantitative Monitoring of Gene Expression Patterns with a Complementary DNA Microarray

Mark Schena,* Dari Shalon,*† Ronald W. Davis, Patrick O. Brown‡

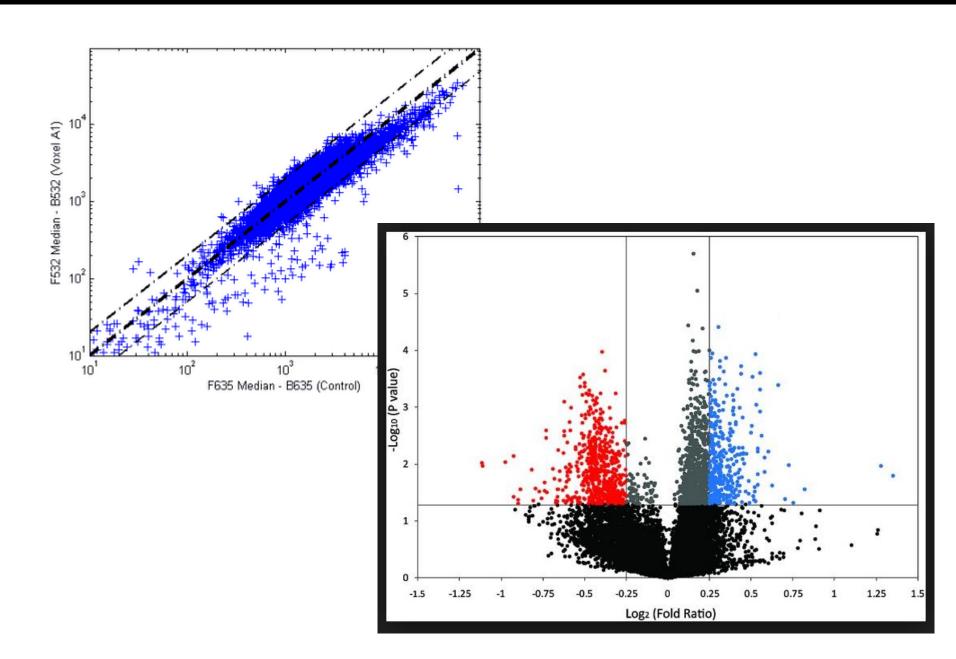
- 1995
- 48 sondas de cDNA por duplicado
- Hoy en día se producen arrays de decenas de miles de features



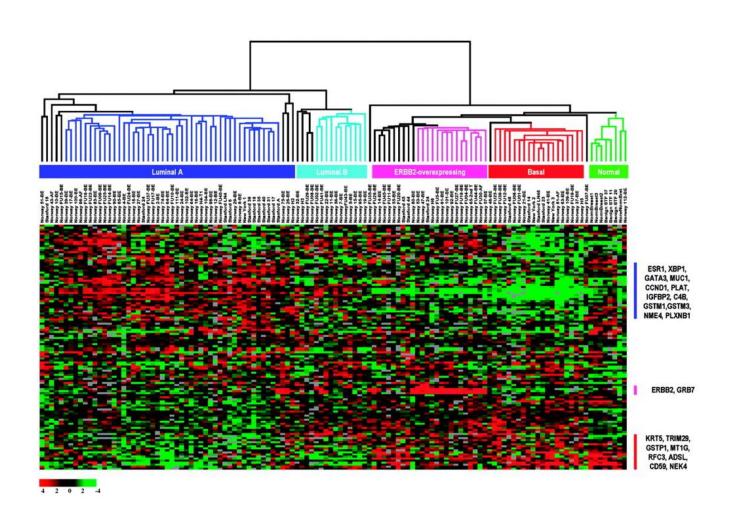
Affimetrix



Análisis de niveles de expresión

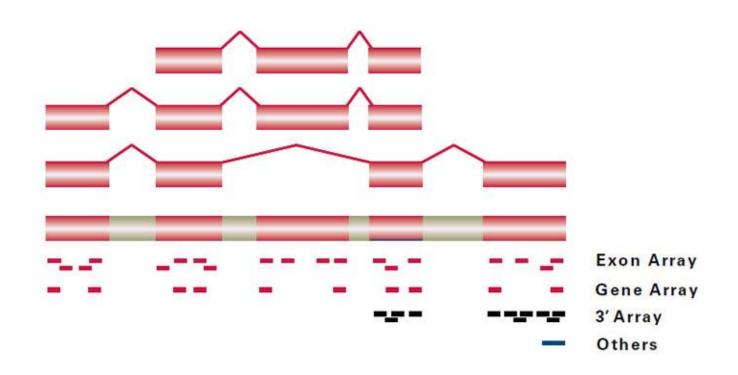


Análisis de niveles de expresión



Arrays de exones

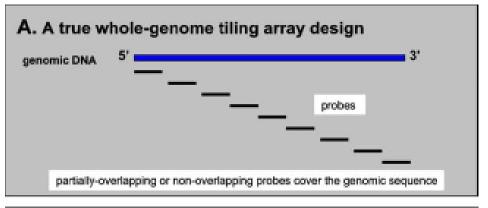
Permiten evaluar patrones de splicing

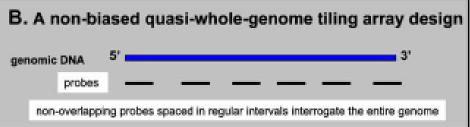


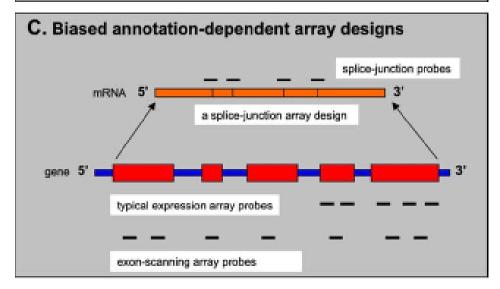
Veo solo lo que quiero ver...

- Desarrollo de los "tiling arrays" (Whole genome arrays, WGA)
 - La idea es abarcar el genoma entero
 - No es necesario conocimiento previo de que es lo que se transcribe
- Permitió avances significativos en el conocimiento de los transcriptomas
- Permitió abarcar un número mucho mayor de preguntas

Tiling arrays







Aplicaciones

Global Identification of Human Transcribed Sequences with Genome Tiling Arrays

```
Paul Bertone,<sup>1*</sup> Viktor Stolc,<sup>1,2*</sup> Thomas E. Royce,<sup>3</sup>
Joel S. Rozowsky,<sup>3</sup> Alexander E. Urban,<sup>1</sup> Xiaowei Zhu,<sup>1</sup>
John L. Rinn,<sup>3</sup> Waraporn Tongprasit,<sup>4</sup> Manoj Samanta,<sup>2</sup>
Sherman Weissman,<sup>5</sup> Mark Gerstein,<sup>3</sup>† Michael Snyder<sup>1,3</sup>†
```

- 13000 regiones transcritas previamente anotadas
- 8000 regiones transcritas sin anotación previa
 - 1500 serian exones previamente no reconocidos
 - 1500 en hebra anti sentido en los intrones
 - 5700 lejos de genes conocidos

Plataformas de secuenciación masiva

A. Sample Processing (Library preparation)



B. NGS Platform (Sequencing)



C. Quality Analysis



Data Analysis & Interpretation

TruSeq ¹

SureSelect, HaloPlex²

SeqCap EZ³

MIP 4

AmpliSeq, TargetSeq⁵

Hands-on-time: 6-48hours

SureSelect ²

SeqCap EZ⁶

Hands-on-time: 6-72hours

٠	Tru	seq,	Nex	tera 1	
100	0				

•	TruSeq	1
	TIME	

- SMRT⁷
- LT-Ion 5
- Ion-ExT Kit 5
- AmpliSeq⁵

Hands-on-time: 6-24hours

TruSea	1	

- SureSelect, HaloPlex²
- SeqCap EZ³
- MIP 4
- AmpliSeq, TargetSeq⁵

Hands-on-time: 6-48hours

Name		Run-time	Max length/ read	Output (Gb)/run
Roche	GS- FLX-454	10 h	400 bp	0.5-1
	MiSeq	5-55 h	2 × 300 bp	0.3-13
	HiSeq	10 h-11 d	2 × 150 bp	15-500
Illumina	NextSeq	11-30 h	2 × 150 bp	19-120
	HiSeqX	3 d	2 × 150 bp	1,800
	NovaSeq	48 h	2 × 150 bp	6,000
	PGM	3-7 h	400 bp	0.09-1.9
Ion Torrent	Proton	4-6 h	500 bp	12-88
200000000	S5	2.5-4h	400 bp	2-16
PacBio	RSII	2 h	3.000 bp	0.09
	Sequel	0.5-6h	20.000 bp	0.08-1.25
Oxford Nanopore	MinION	1 min-48 h	10.000 bp\$	44

CRITERIA

- Specificity: 95%-98%
- Depth: 200-1000×
- >20x for 95-98% of ROIs

CRITERIA

- Specificity: 75%-80%
- Depth: 100-200×
- >20x for 90-95% of ROIs

CRITERIA

- Specificity: 95%-98%
- Depth: 30-60×
- >20x for 90-95% of ROIs

CRITERIA

- Specificity: 95%-98%
- Depth: 100-200x *
- >20x for 90-95% of ROIs

CRITERIA

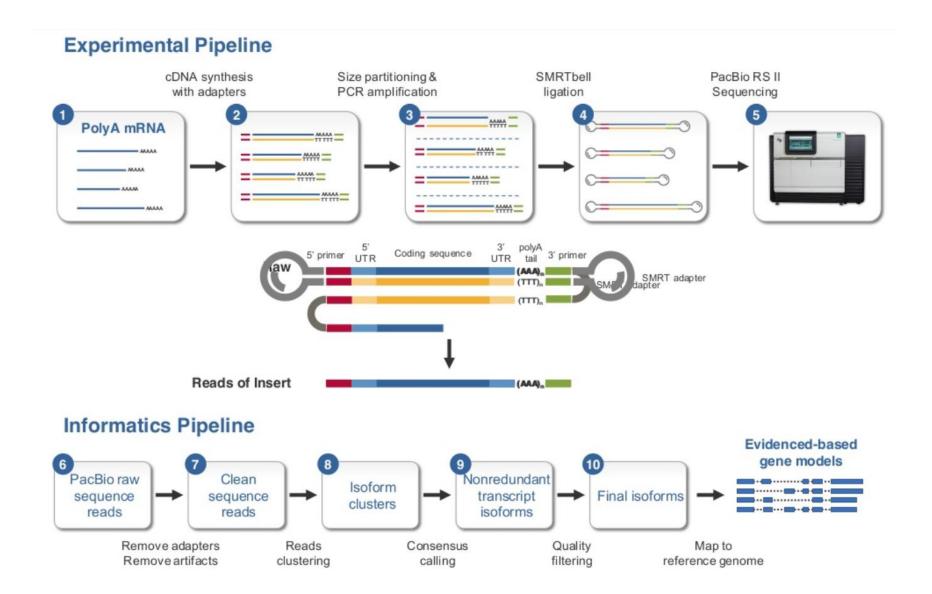
- Specificity: 95%-98%
 Depth: >2000×*
 >20x for 90-95% of ROIs

Sequencer Outputs are further processed

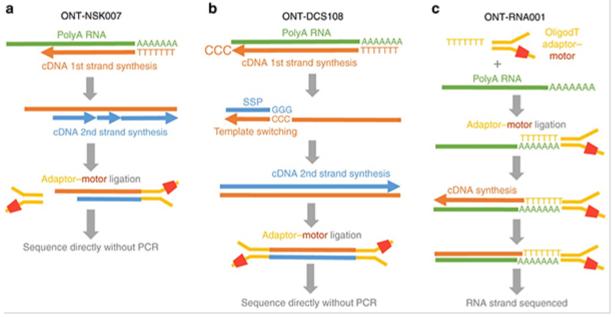
(Figure 2)

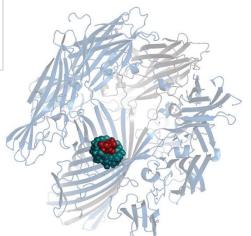
Gene-Panels

Secuenciadores de tercera generación: PacBio



Secuenciadores de tercera generación: Nanopore





Capacidad de secuenciar el transcrito entero

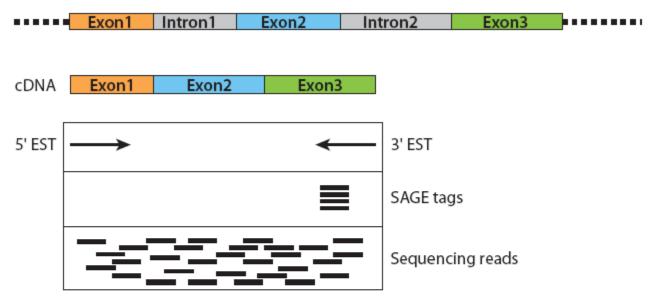
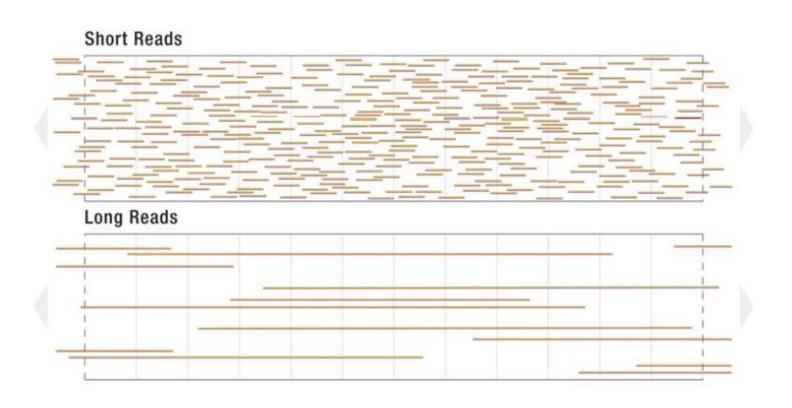


Figure 1

Gene model coverage by various sequencing-based methods for transcriptome analysis. Sanger-based expressed sequence tags (ESTs) are generated from the 3' or 5' end of the transcript, whereas SAGE tags represent short sequences at its 3' end. Randomly primed short reads generated by next-generation sequencers detect bases throughout the length of the transcript.

Ensamblaje

 Le elección de la plataforma va a estar sujeta a la disponibilidad de genoma completo



Anotación de regiones codificantes

- Se obtienen secuencias de todos los exones del gen
 - Descripción de nuevos exones
 - Descripción de nuevas variantes de splicing (pair ends sequencing)

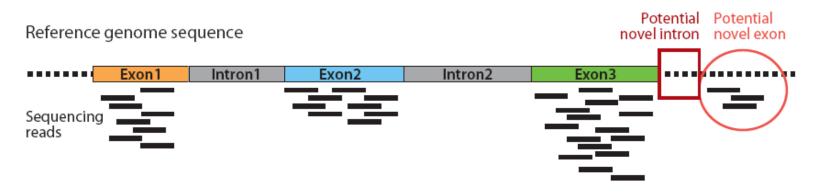
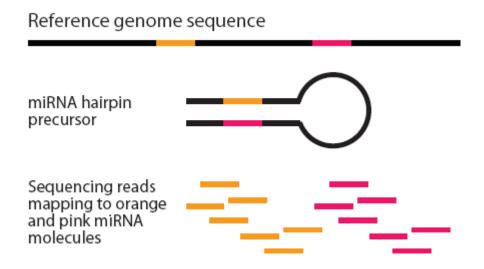


Figure 3

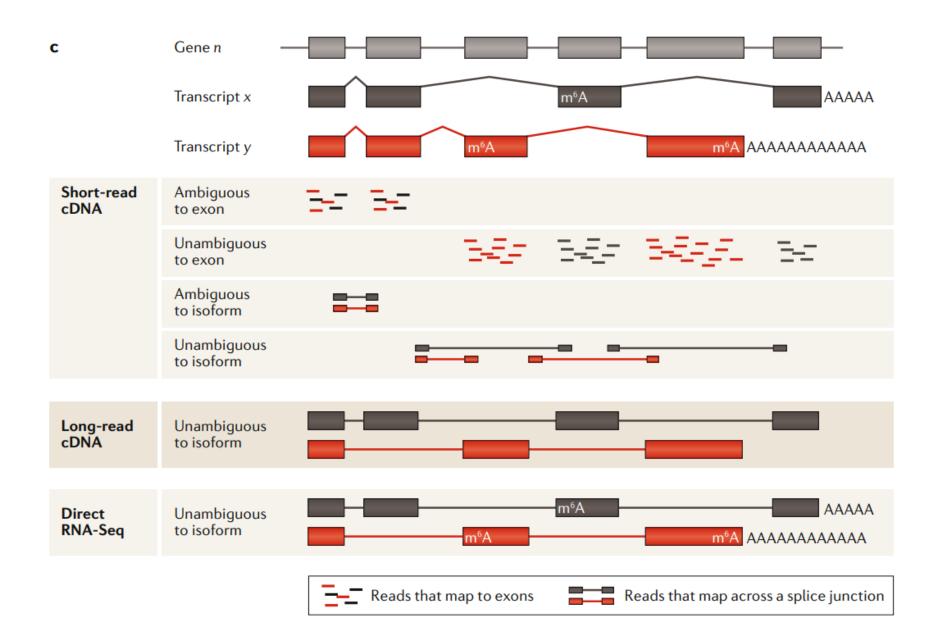
Protein-coding gene annotation using transcriptome sequencing data. This figure illustrates how novel exons and introns can be discovered by mapping transcriptome sequencing reads to an annotated reference genome sequence.

RNAs no codificantes

Secuenciado de regiones no anotadas no codificantes



Identificación de variantes

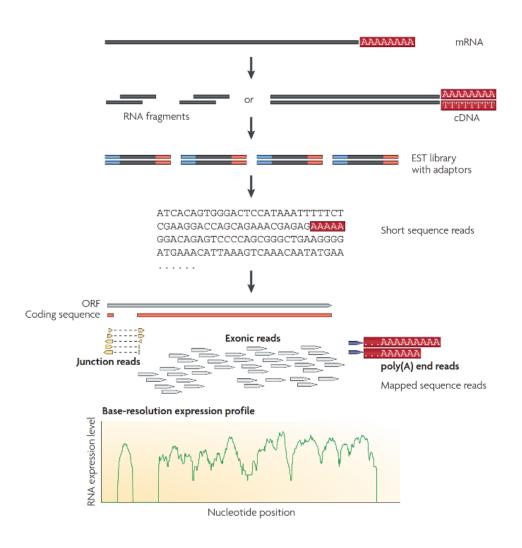


Segunda vs tercera generación

Sequencing technology	Platform	Advantages	Disadvantages	Key applications
Short-read cDNA	Illumina, Ion Torrent	 Technology features very high throughput: currently 100–1,000 times more reads per run than long-read platforms Biases and error profiles are well understood (homopolymers are still an issue for lon Torrent) A huge catalogue of compatible methods and computational workflows are available Analysis works with degraded RNA 	 Sample preparation includes reverse transcription, PCR and size selection adding biases to all methods Isoform detection and quantitation can be limited Transcript discovery methods require a de novo transcriptome alignment and/or assembly step 	Nearly all RNA-seq methods have been developed for short-read cDNA sequencing: DGE, WTA, small RNA, single-cell, spatialomics, nascent RNA, translatome, structural and RNA-protein interaction analysis, and more are all possible
Long-read cDNA	PacBio, ONT	 Long reads of 1–50 kb capture many full-length transcripts Computational methods for de novo transcriptome analysis are simplified 	 Technology features low-to-medium throughput: currently only 500,000 to 10 million reads per run Sample preparation includes reverse transcription, PCR and size selection (for some protocols), adding biases to many methods Degraded RNA analysis is not recommended 	Sequencing is particularly suited to isoform discovery, de novo transcriptome analysis, fusion transcript discovery, and MHC, HLA or other complex transcript analysis
Long-read RNA	ONT	 Long reads of 1–50 kb capture many full-length transcripts Computational methods for de novo transcriptome analysis are simplified Sample preparation does not require reverse transcription or PCR-reducing biases RNA base modifications can be detected Poly(A) tail lengths can be directly estimated from single-molecule sequencing 	 Technology features low throughput: currently only 500,000 to 1 million reads per run Sample preparation and sequencing biases are not well understood Degraded RNA analysis is not recommended 	 Sequencing is particularly suited to isoform discovery, de novo transcriptome analysis, fusion transcript discovery, and MHC, HLA or other complex transcript analysis Ribonucelotide modifications can be detected

RNA-seq

 El conteo de lecturas es cuantitativo lo que nos da idea de cantidad de expresión



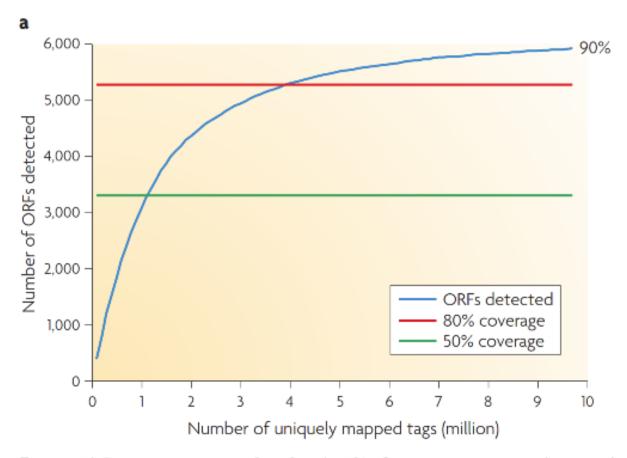
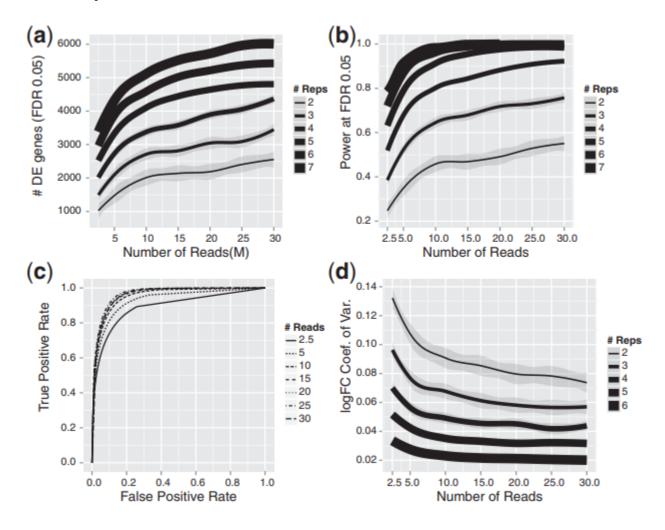
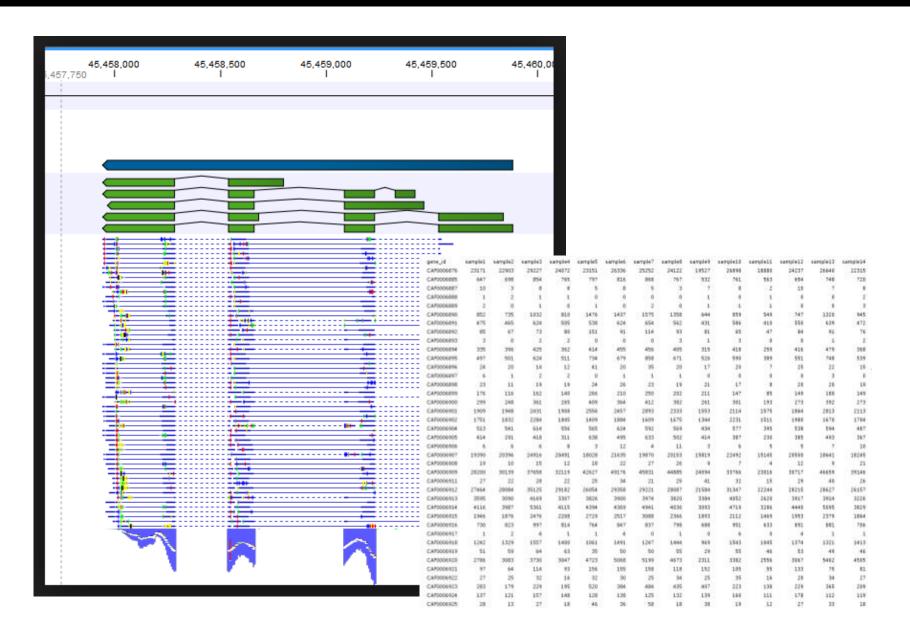
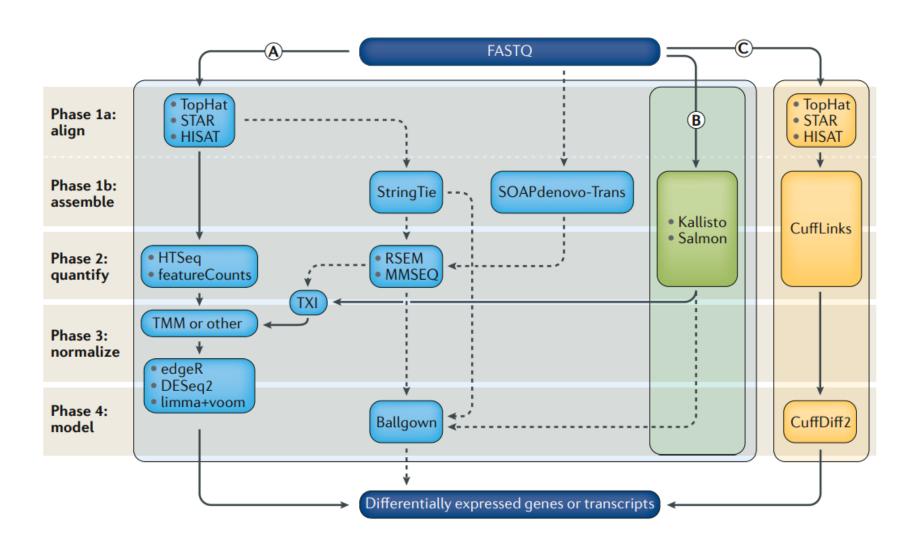


Figure 5 | Coverage versus depth. a | 80% of yeast genes were detected at 4 million uniquely mapped RNA-Seq reads, and coverage reaches a plateau afterwards despite the increasing sequencing depth. Expressed genes are defined as having at least four independent reads from a 50-bp window at the 3' end. Data is taken from REF. 18.

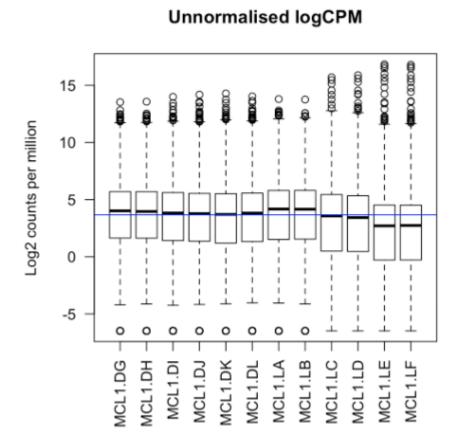
Replicas vs profundidad



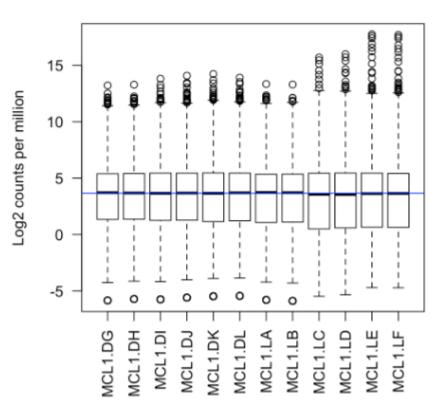




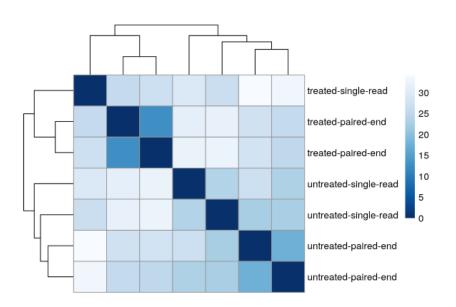
- Normalización:
 - Corrige efecto de profundidad diferente entre muestras

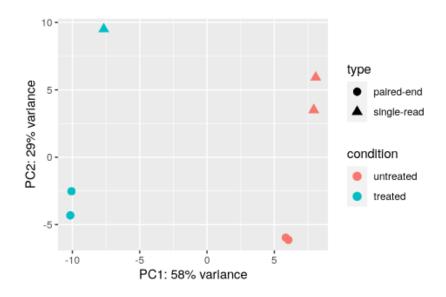


Voom transformed logCPM

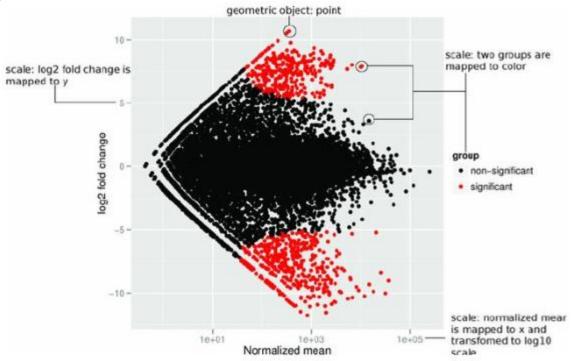


- Análisis exploratorios:
 - Correlación entre muestras

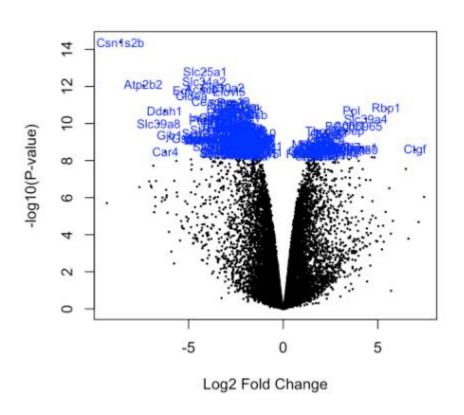




- Expresión diferencial
 - MA plot



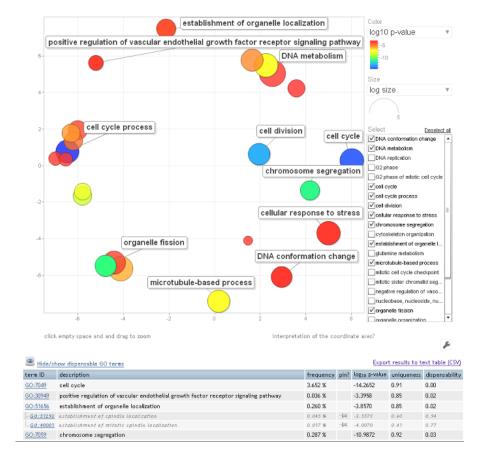
- Expresión diferencial
 - Volcano plot

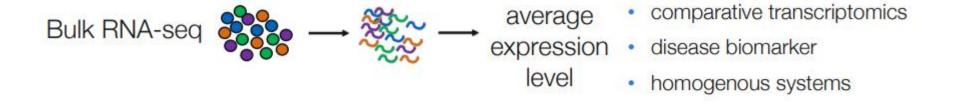


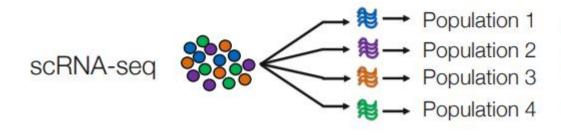
(0)	А	В	С	D	E	F	G
1	transcript_name	baseMean	log2FoldChange	IfcSE	stat	pvalue	padj
2	30484730	3,144.32943360273	5.04026145518364	0.196928866418254	25.5943252345683	1.7644314301418e-144	4.54023495604088e-140
3	30561786	1,413.35213135671	3.41626574251214	0.13451370327947	25.3971577558488	2.71095008913392e-142	3.4879083846797e-138
4	30504829	880.771473715433	4.32396956861188	0.173416033458984	24.9340818283367	3.17818540613069e-137	2.72603556235183e-133
5	30544189	8,258.20272934674	3.24090885472581	0.130554314699064	24.8242186571643	4.91076326741851e-136	3.15909400993032e-132
6	30483127	1,295.09139872873	6.55681270831358	0.273206432704986	23.9994814301966	2.81566948590279e-127	1.44905614422501e-123
7	30551621	1,511.46596367277	3.3786619006878	0.142163737793915	23.7659895070121	7.50972165530951e-125	3.22066929390707e-121
8	30475802	1,898.44230304467	6.95150465936128	0.296772084097406	23.4237148029047	2.45035765738168e-121	9.00751474853505e-118
9	30480296	1,210.07334890955	3.48353842881072	0.154359693570149	22.5676687238797	9.00712662522495e-113	2.89714227900361e-109
10	30543987	710.020123515779	5.68496687104772	0.252361926081687	22.5270386833534	2.2553967057206e-112	6.44842978128917e-109
11	30517326	3,611.6823066722	3.51111066113919	0.156500582481462	22.4351283903694	1.78792563628964e-111	4.60069024730051e-108
12	30535740	559.981220962047	4.09061739312591	0.182682690256112	22.3919266099654	4.71750598855994e-111	1.1035533099784e-107
13	30544292	1,472.41459635489	4.80156003343999	0.216812606426031	22.1461293814487	1.13676216585688e-108	2.43759700431909e-105
14	30499835	457.358467856158	6.26052947903444	0.2955910524082	21.1796988712259	1.46968872158091e-99	2.90907924490154e-96
15	30548836	1,891.57726951951	5.50297046375875	0.263949523079027	20.8485713463901	1.57041671532408e-96	2.88642592276565e-93
16	30477139	1,224.48273150217	3.84512112082617	0.184655160791858	20.8232529453122	2.66467630015296e-96	4.57116337036906e-93
17	30518455	417.579205665653	4.7316312598701	0.229901300395872	20.5811417844206	4.05028224251837e-94	6.51386641653017e-91
18	30554782	755.049051445421	3.60242774699136	0.180563057495959	19.9510785702772	1.46689640228573e-88	2.22036342491861e-85
19	30487155	1,808.40178314046	4.28163601075905	0.215972721453325	19.8248926158222	1.81568244319052e-87	2.59561892378769e-84
20	30499588	608.66999182642	7.19339714094185	0.363520473282165	19.78814859035	3.76624508057103e-87	5.10068517964494e-84

0	A	В	С	D	E	F	G
1	transcript_name	baseMean	log2FoldChange	IfcSE	stat	pvalue	padj
2	30484730	3,144.32943360273	5.04026145518364	0.196928866418254	25.5943252345683	1.7644314301418e-144	4.54023495604088e-140
3	30561786	1,413.35213135671	3.41626574251214	0.13451370327947	25.3971577558488	2.71095008913392e-142	3.4879083846797e-138
4	30504829	880.771473715433	4.32396956861188	0.173416033458984	24.9340818283367	3.17818540613069e-137	2.72603556235183e-133
5	30544189	8,258.20272934674	3.24090885472581	0.130554314699064	24.8242186571643	4.91076326741851e-136	3.15909400993032e-132
6	30483127	1,295.09139872873	6.55681270831358	0.273206432704986	23.9994814301966	2.81566948590279e-127	1.44905614422501e-123
7	30551621	1,511.46596367277	3.3786619006878	0.142163737793915	23.7659895070121	7.50972165530951e-125	3.22066929390707e-121
8	30475802	1,898.44230304467	6.95150465936128	0.296772084097406	23.4237148029047	2.45035765738168e-121	9.00751474853505e-118
9	30480296	1,210.07334890955	3.48353842881072	0.154359693570149	22.5676687238797	9.00712662522495e-113	2.89714227900361e-109
10	30543987	710.020123515779	5.68496687104772	0.252361926081687	22.5270386833534	2.2553967057206e-112	6.44842978128917e-109
11	30517326	3,611.6823066722	3.51111066113919	0.156500582481462	22.4351283903694	1.78792563628964e-111	4.60069024730051e-108
12	30535740	559.981220962047	4.09061739312591	0.182682690256112	22.3919266099654	4.71750598855994e-111	1.1035533099784e-107
13	30544292	1,472.41459635489	4.80156003343999	0.216812606426031	22.1461293814487	1.13676216585688e-108	2.43759700431909e-105
14	30499835	457.358467856158	6.26052947903444	0.2955910524082	21.1796988712259	1.46968872158091e-99	2.90907924490154e-96
15	30548836	1,891.57726951951	5.50297046375875	0.263949523079027	20.8485713463901	1.57041671532408e-96	2.88642592276565e-93
16	30477139	1,224.48273150217	3.84512112082617	0.184655160791858	20.8232529453122	2.66467630015296e-96	4.57116337036906e-93
17	30518455	417.579205665653	4.7316312598701	0.229901300395872	20.5811417844206	4.05028224251837e-94	6.51386641653017e-91
18	30554782	755.049051445421	3.60242774699136	0.180563057495959	19.9510785702772	1.46689640228573e-88	2.22036342491861e-85
19	30487155	1,808.40178314046	4.28163601075905	0.215972721453325	19.8248926158222	1.81568244319052e-87	2.59561892378769e-84
20	30499588	608.66999182642	7.19339714094185	0.363520473282165	19.78814859035	3.76624508057103e-87	5.10068517964494e-84

 Análisis de sobrerrepresentación de términos de ontología

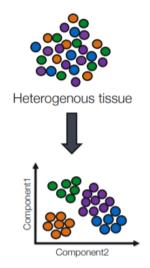




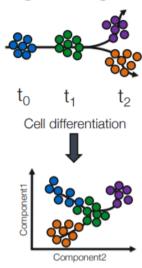


- define heterogeneity
- identify rare cell population
- cell population dynamics

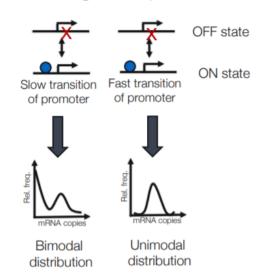




Lineage tracing study

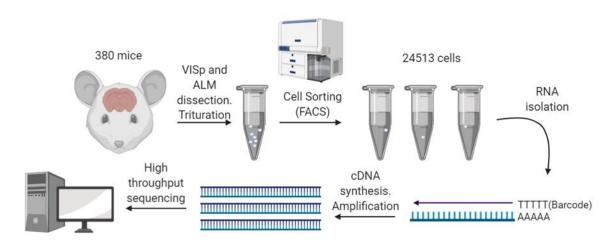


Stochastic gene expression

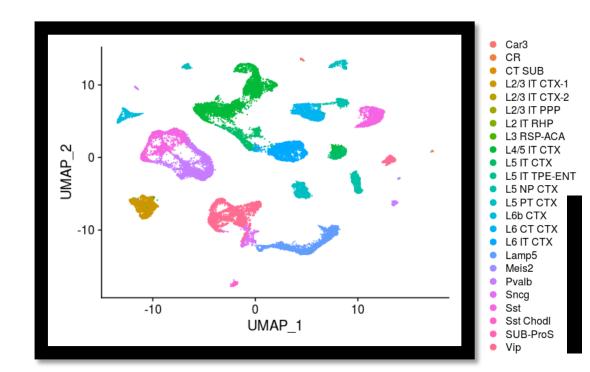


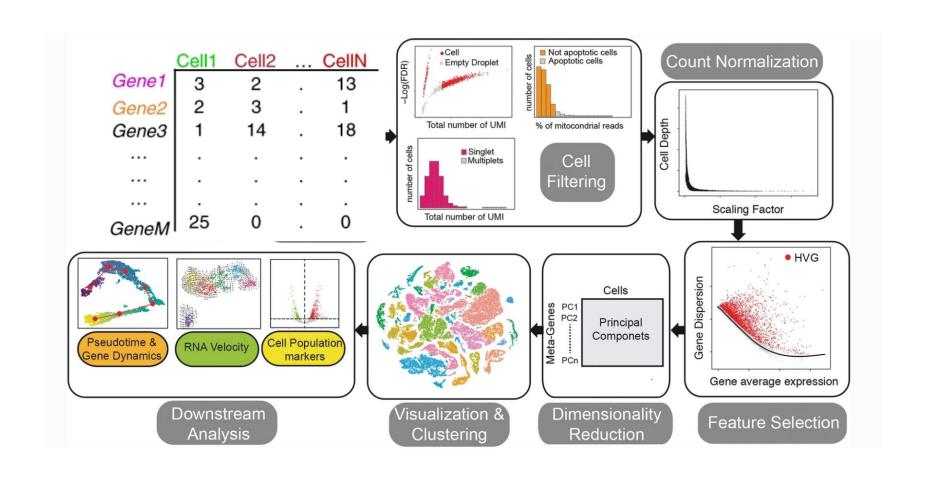
Liu S and Trapnell C. Single-cell transcriptome sequencing: recent advances and remaining challenges, F1000 Research 2016 (doi: 10.12688/f1000research.7223.1)

Junker and van Oudenaarden; Every Cell Is Special: Genome-wide Studies Add a New Dimension to Single-Cell Biology, Cell 2014 (doi: 10.1016/j.cell.2014.02.010)

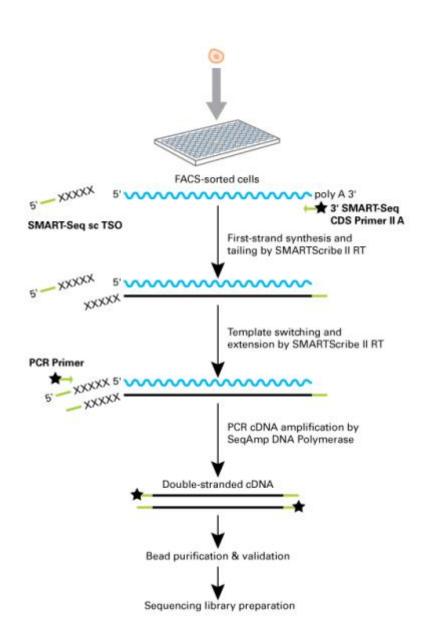


1.8 M mapped reads/cell 10000 genes detected/cell



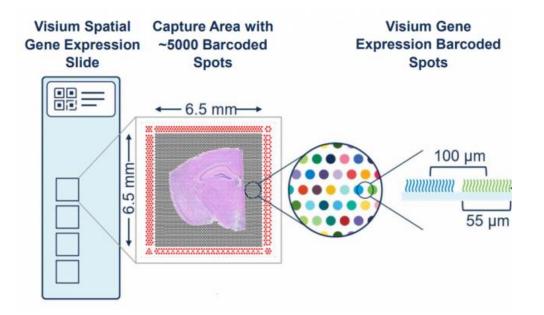


Smart-seq

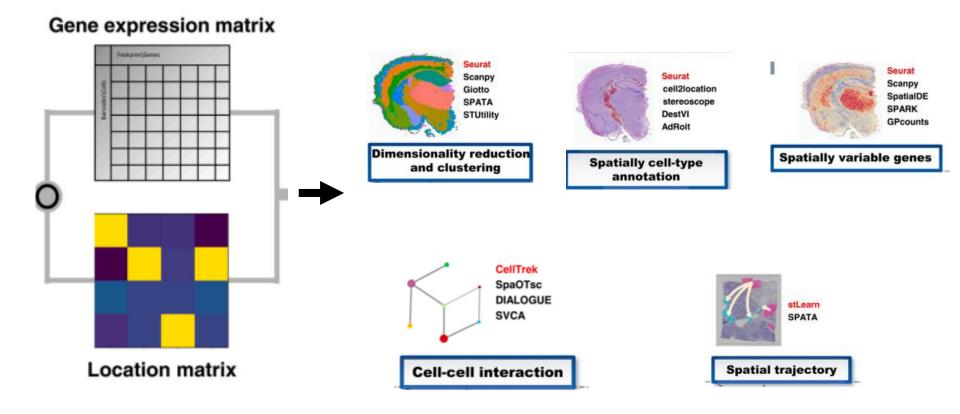


Transcriptómica espacial

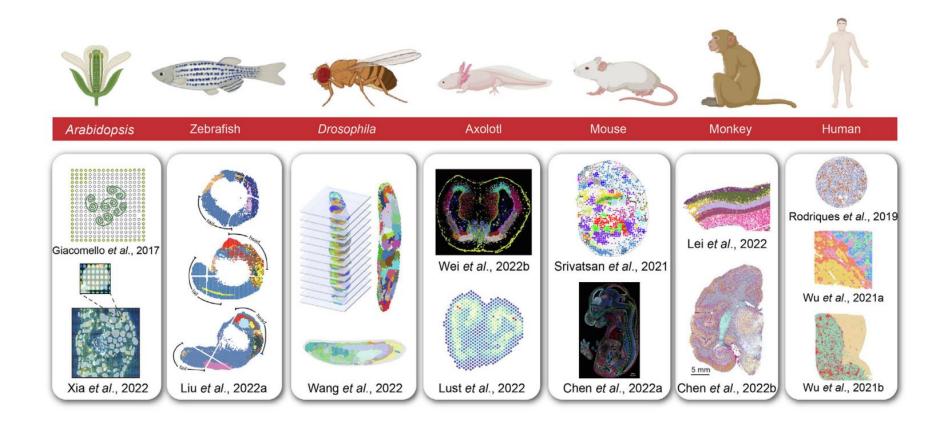
- Combina RNA-seq de pequeña región (célula única) con coordenadas
- Visium (10x)
 - Método más utilizado de transcriptómica especial
 - Lisis y transferencia de ARNs a sitios de captura
 - La resolución viene dada por el tamaño del sitio de captura.



Transcriptómica espacial



Transcriptómica espacial



Preguntas????