

# Genomas individuales: 454

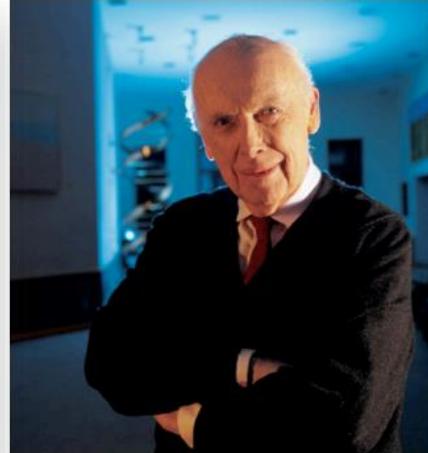
nature

Vol 452 | 17 April 2008 | doi:10.1038/nature06884

## LETTERS

### The complete genome of an individual by massively parallel DNA sequencing

David A. Wheeler<sup>1\*</sup>, Maithreyan Srinivasan<sup>2\*</sup>, Michael Egholm<sup>2\*</sup>, Yufeng Shen<sup>1\*</sup>, Lei Chen<sup>1</sup>, Amy McGuire<sup>3</sup>, Wen He<sup>2</sup>, Yi-Ju Chen<sup>2</sup>, Vinod Makhijani<sup>2</sup>, G. Thomas Roth<sup>2</sup>, Xavier Gomes<sup>2</sup>, Karrie Tartaro<sup>2†</sup>, Faheem Niazi<sup>2</sup>, Cynthia L. Turcotte<sup>2</sup>, Gerard P. Irvzyk<sup>2</sup>, James R. Lupski<sup>4,5,6</sup>, Craig Chinault<sup>4</sup>, Xing-zhi Song<sup>1</sup>, Yue Liu<sup>1</sup>, Ye Yuan<sup>1</sup>, Lynne Nazareth<sup>1</sup>, Xiang Qin<sup>1</sup>, Donna M. Muzny<sup>1</sup>, Marcel Margulies<sup>2</sup>, George M. Weinstock<sup>1,4</sup>, Richard A. Gibbs<sup>1,4</sup> & Jonathan M. Rothberg<sup>2†</sup>



Watson

- 6 meses luego de la publicacion del anterior
- Plataforma 454 FLX
- 24,5 billion bases en 93,2 millones de lecturas
- 7.4X de covertura
- Menos de 1 million de dolares
- Completado en menos de 5 meses

*Wheeler, D. A. et al. Nature (2008)*

# Genomas famosos!

Wednesday 06 April 2011



HOME NEWS SPORT FINANCE COMMENT CULTURE TRAVEL LIFESTYLE FASHION TECHNOLOGY  
UK | World | Politics | Obituaries | Royal Wedding | Earth | Science | Health News | Education |

Celebrity news

## Glenn Close has genome sequenced to publicise mental illness

American actress Glenn Close has joined a handful of celebrities to have their genome sequenced in the name of science.



Close is a founder of the nonprofit group BringChange2Mind, which raises awareness about mental illness. Photo: AP

Share:

Recommend

Tweet 0

Celebrity news  
News » World News » North America » USA »

IN NEWS



Naomi Campbell opens pop up shop

## Ozzy Osbourne genome sequenced

Genetic analysis of Black Sabbath star reveals he is more likely to experience hallucinations on marijuana and has increased risk of alcohol and cocaine addiction, researchers say

Sean Michaels

guardian.co.uk, Friday 5 November 2010 11.30 GMT  
[Article history](#)

Tweet 0

Share 814

b

Comments (29)

A [larger](#) | [smaller](#)

Music  
Ozzy Osbourne - Black Sabbath - Metal - Pop and rock

Culture

Science  
Genetics - Biology

More news

Related

8 May 2002  
Rocker at odds with MTV over second series

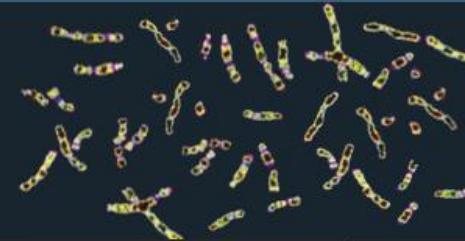
6 Oct 2010  
Ozzy Osbourne covers John Lennon song for charity



Man of visions ... Ozzy Osbourne with wife Sharon at Jon Stewart's Rally to Restore Sanity and/or Fear on Saturday. Photograph: Scott Gries/EMPICS Entertainment

# IGSR: The International Genome Sample Resource

Providing ongoing support for the 1000 Genomes Project data



Home   About   Data   Portal   Analysis   Contact   Browser   FAQ

Search 1000genomes



## IGSR and the 1000 Genomes Project



Populations: ● - African; ● - American; ● - East Asian; ● - European; ● - South Asian;

The International Genome Sample Resource (IGSR) was established to ensure the ongoing usability of data generated by the 1000 Genomes Project and to extend the data set. More information is available [about the IGSR](#).

## Links

[Announcements](#)

[IGSR Sample Collection Principles](#)

[1000 Genomes Project Publications](#)

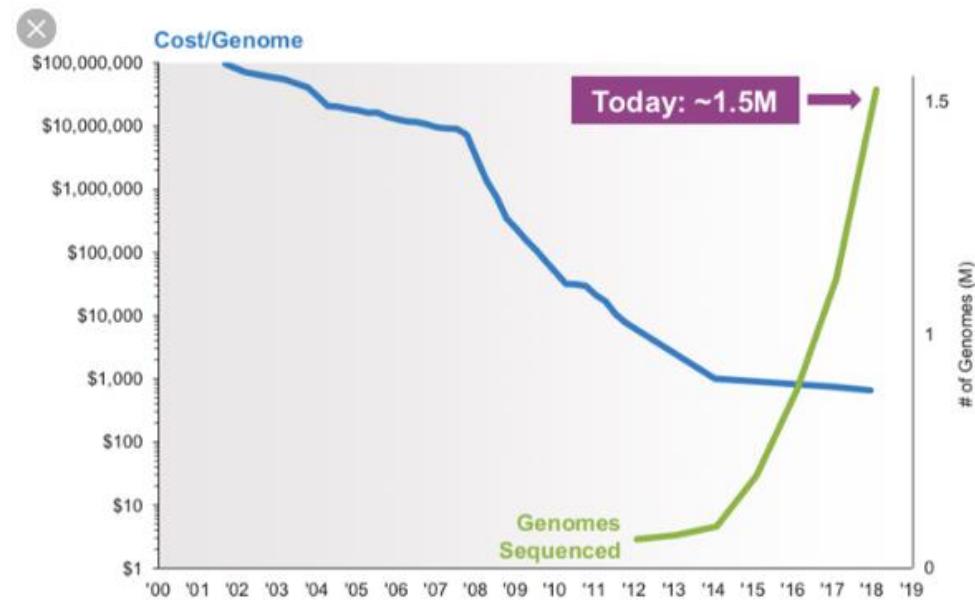
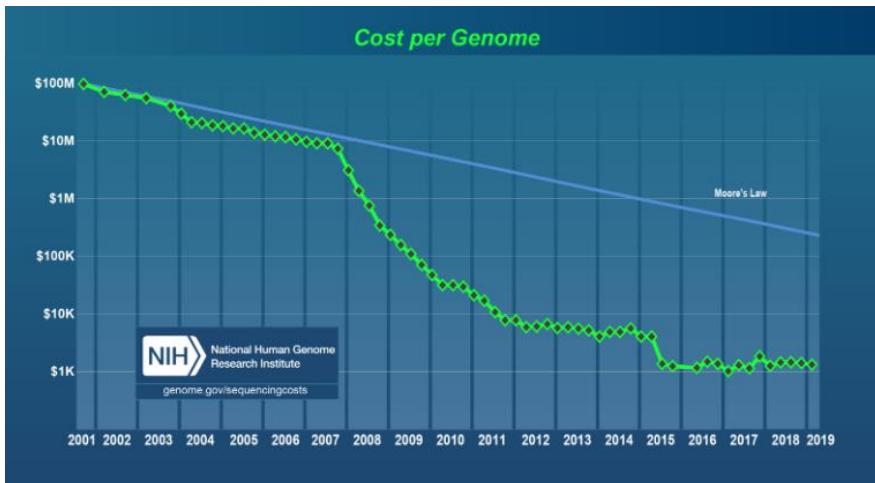
[File formats](#)

[Software tools](#)

[Download data](#)

[Twitter](#)

# Costos y número de genomas completos



- Hoy en día hay millones de genomas secuenciados

# “Personal Genomics”

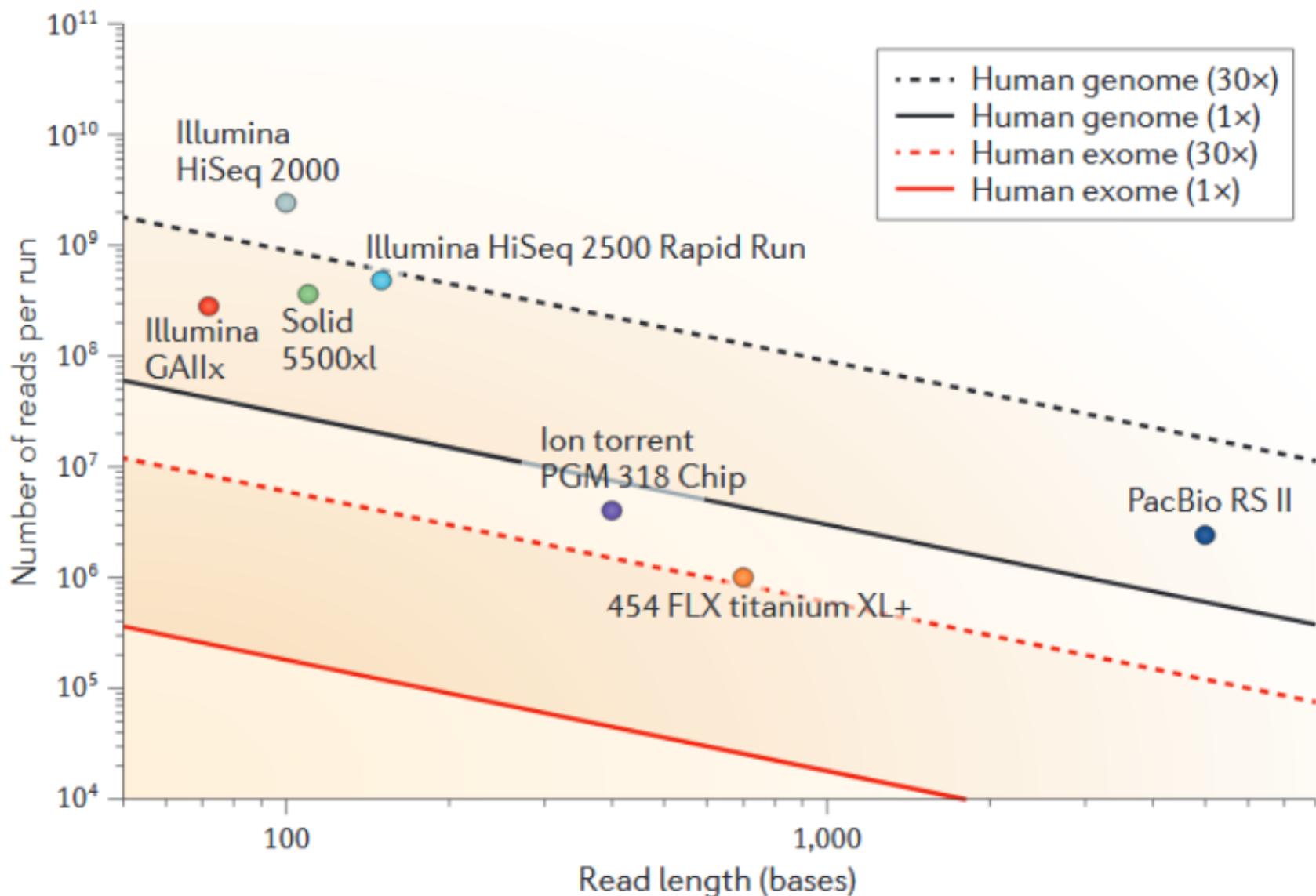
**23andMe competitor Veritas Genetics slashes price of whole genome sequencing 40% to \$600**

PUBLISHED MON, JUL 1 2019 • 9:30 AM EDT

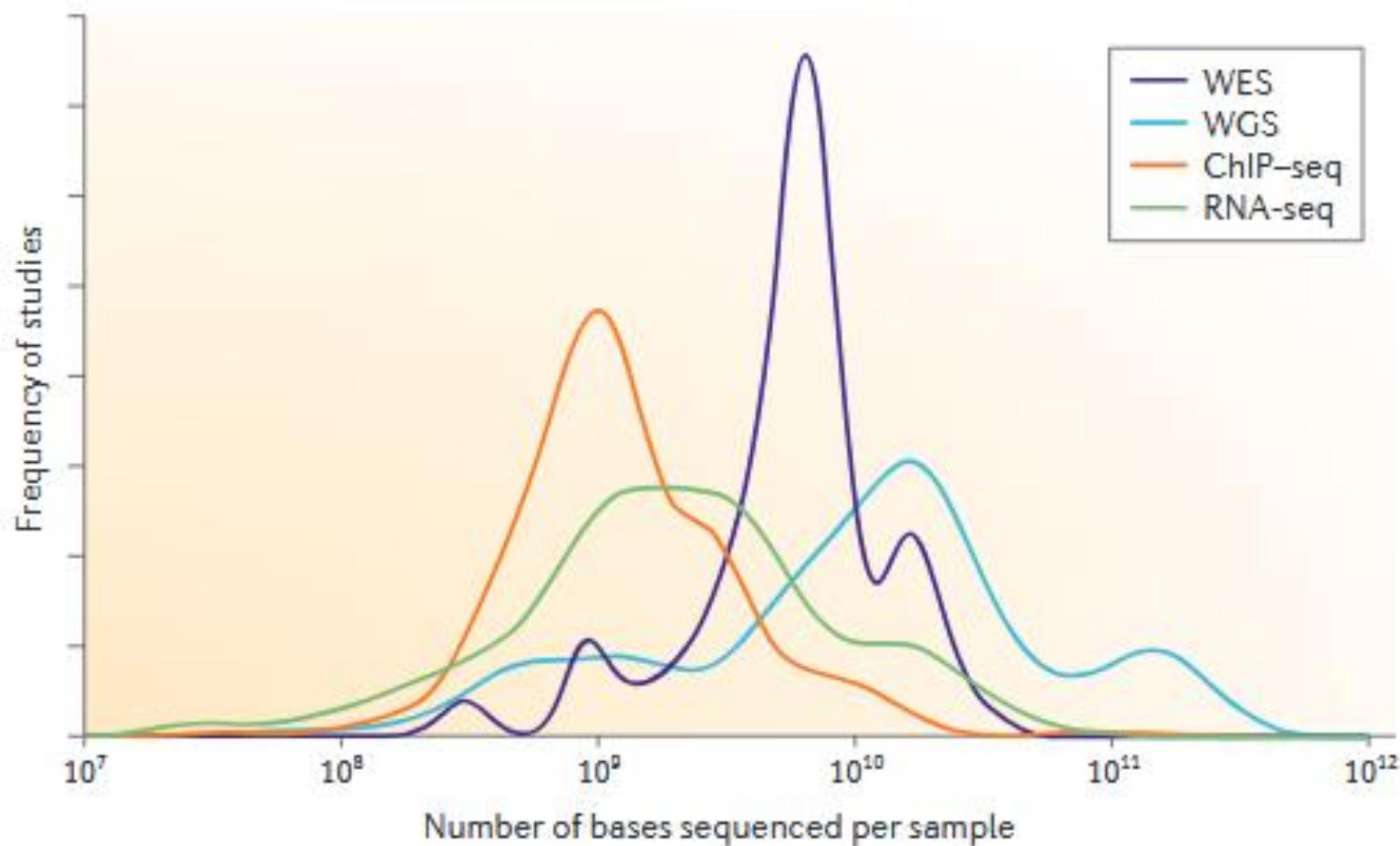
The image shows a screenshot of a website for Veritas Genetics. In the top right corner, the URL <https://www.veritasgenetics.com> is displayed. On the left side, the Veritas logo, consisting of a stylized 'V' icon followed by the word 'Veritas' and 'The Genome Company' below it, is visible. The main content area features a large image of a product box for 'my Genome by Veritas'. The box has a geometric pattern of red, yellow, and white. To the right of the box, promotional text reads: 'Get the most comprehensive genetic testing service there is.' and 'Now for \$599.' Below this text is a red button with the words 'ORDER NOW' in white. At the bottom of the page, there are five navigation links: 'I am a Physician', 'Disease Risk', 'Family', 'Health & Longevity', and 'I am Curious'.

# “Personal Genomics”

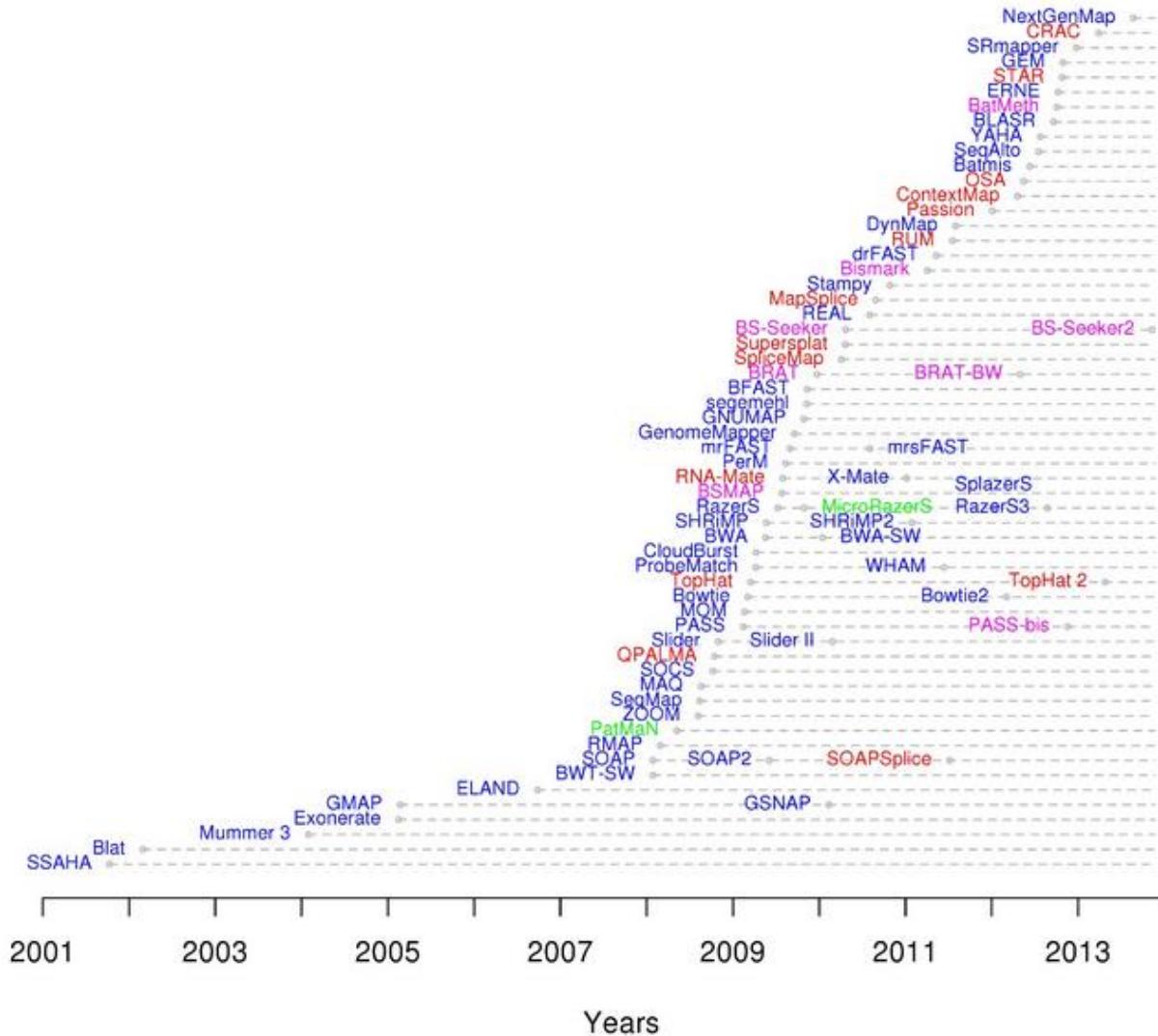
Genetic test company	blood sample	sent by post	USA <sup>[3]</sup>
ACGT, Inc	blood sample	sent by post	USA <sup>[4]</sup>
Admera Health LLC	blood sample	sent by post	USA <sup>[5]</sup>
Affiliated Genetics, Inc	blood sample	sent by post	USA <sup>[6]</sup>
Applied Biological Materials Inc	blood sample	sent by post	Canada <sup>[7]</sup>
CD Genomics	blood sample	sent by post	USA
CEN4GEN Institute for Genomics and Molecular Diagnostics	blood sample	sent by post	Canada <sup>[8]</sup>
ChunLab, Inc.	blood sample	sent by post	Korea, USA <sup>[9]</sup>
Core Life Sciences, Inc.	blood sample	sent by post	USA <sup>[10]</sup>
Complete Genomics	blood sample, frozen tissue, cell lines, or saliva <sup>[11]</sup>	send by post	USA
Clinical Microarray Core	blood sample	send by post	USA
DNA Link, Inc	blood sample	sent by post	Korea, USA <sup>[12]</sup>
Fulgent Diagnostics	blood sample	sent by post	USA <sup>[13]</sup>
Full Genomes Corporation	saliva sample	send by post	USA <sup>[14]</sup>
GENEWIZ, Inc.	blood sample	sent by post	USA <sup>[15]</sup>
Gene by Gene	blood sample (3 to 5cc)	send by post	USA
Genomix4Life	blood sample	sent by post	Italy <sup>[16]</sup>
Genomics Personalized Health	blood sample	sent by post	USA <sup>[17]</sup>
HIB Genomics Services Lab	blood sample	sent by post	USA <sup>[18]</sup>
Illumina	blood sample (4 to 8 ml vial) <sup>[19]</sup>	sent by post	USA <sup>[20]</sup>
Kinghorn Centre for Clinical Genomics	blood sample	sent by post	Australia <sup>[22]</sup>
Macrogen	blood sample	sent by post	Korea, Japan, USA, Netherlands <sup>[23]</sup>
Meports	saliva sample	sent by post	USA <sup>[24]</sup>
myGenomics, LLC	blood sample	sent by post	USA <sup>[25]</sup>
Novogene Corporation	blood sample	sent by post	USA, UK, Japan, China <sup>[26]</sup>
Quick Biology	blood sample	send by post	USA
Sequenom	blood sample	drawn on nearby sequenom location	USA
SBMRI Analytical Genomics Core Facility <sup>[28]</sup>	blood sample	sent by post	USA <sup>[29]</sup>
Sequentia Biotech SL	blood sample	sent by post	Spain <sup>[30]</sup>
SGI-DNA	blood sample	sent by post	USA <sup>[31]</sup>



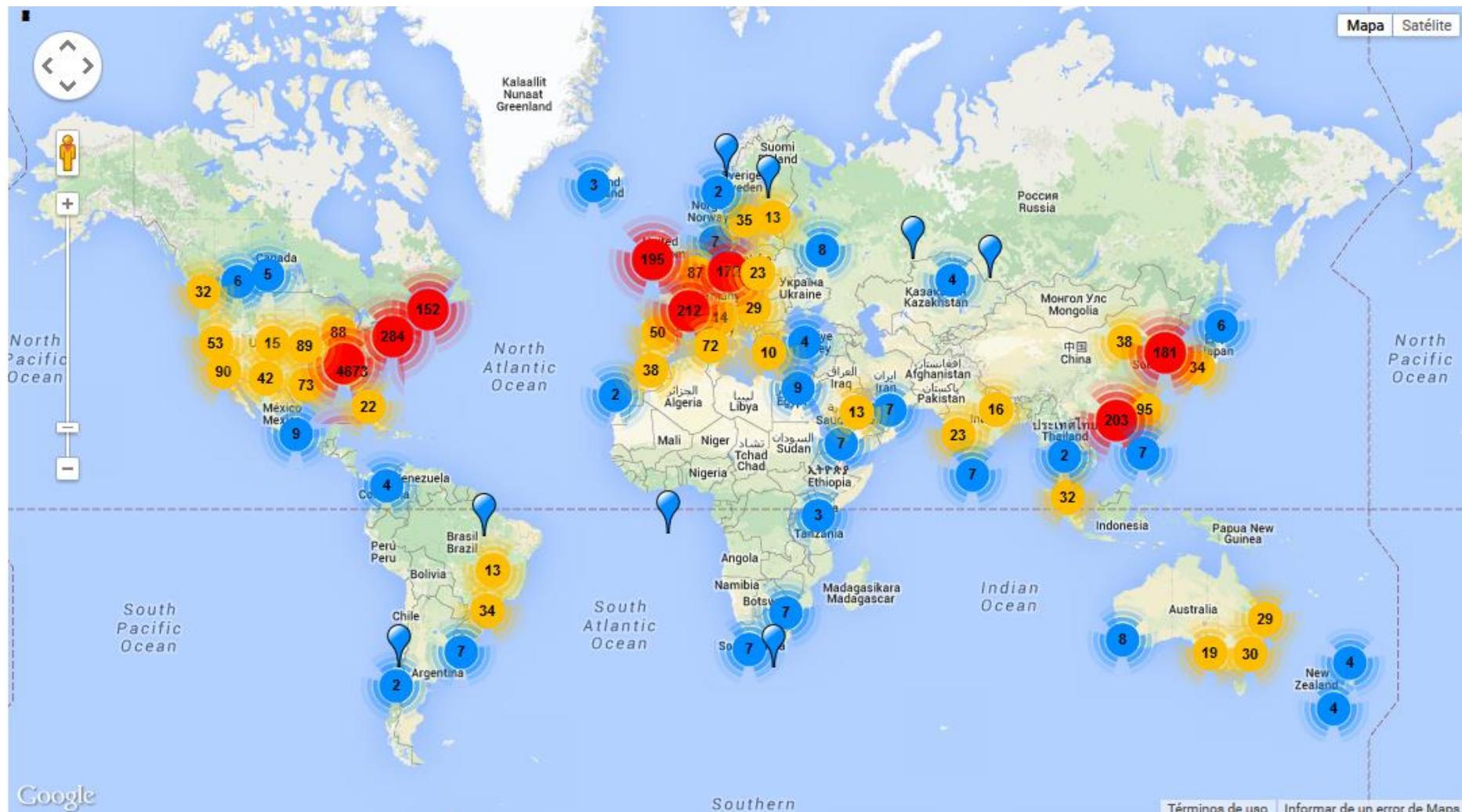
# Número de bases secuenciadas



# Software de análisis de secuencias cortas



# Mapa mundial de secuenciadores de segunda generación

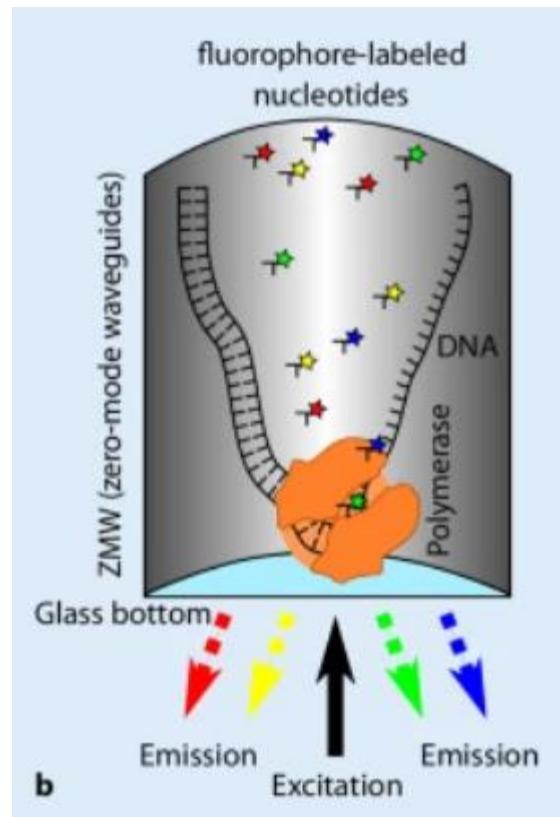


# Secuenciadores de tercera generación

- Secuenciadores de molécula única
  - No hay amplificación de la muestra
  - Mayor tamaño de las lecturas
  - Actualmente mayor costo por base

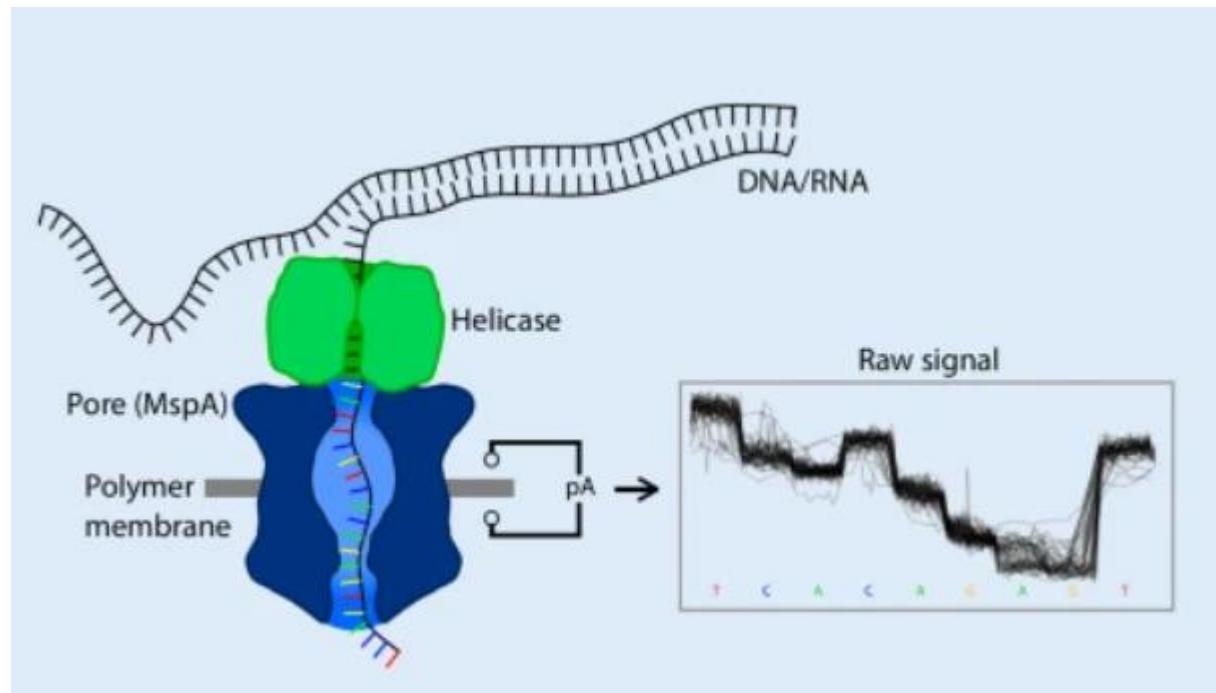
# PacBio

- single-molecule real-time (SMRT) sequencing
- Polimerasa unida a un pocillo



# Nanopore

- pasaje a través de un poro proteico modificado
- Se mide la variación de la corriente eléctrica a través del poro
- 450 bases por segundo



# The MinION™ device: a miniaturised sensing system

Nanopore sensing technology can be miniaturised into a portable device for electronic single-molecule sensing.

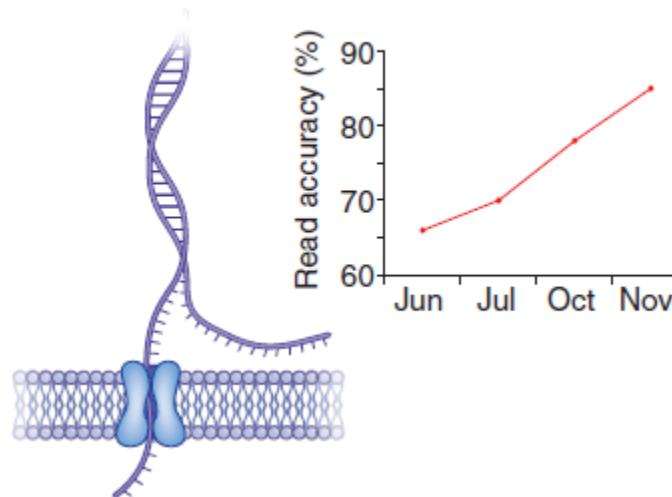
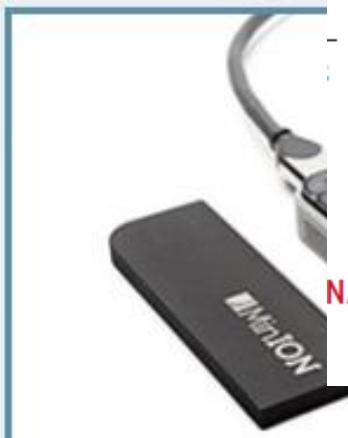
The MinION™ device is a small instrument that is compatible with consumable flow cells containing the proprietary sensor chip, Application-Specific Integrated Circuit (ASIC) and nanopores that are needed to perform a complete single-molecule sensing experiment. Plugging directly into a laptop or desktop computer through a USB port, it is a self-contained device to deliver real-time experimental data.

The MinION device is adapted for hundred participants in the MinION Access Program (MAP). It can handle long read lengths, real time data analysis and reporting. The MAP initially focuses on DNA sequencing.

The MinION is operated using a graphical user interface (GUI) and has the option to perform these analyses in the cloud.

Oxford Nanopore is focused on developing the MinION device to be compatible with complex samples such as blood/serum.

*The MinION device is a miniaturised sensing system, designed for single-molecule sensing and connected via a USB port of a laptop.*



**Figure 1 |** Accuracy of two-directional reads on the MinION nanopore sequencer has improved rapidly since the device was introduced to early access users (months shown for 2014).

NATURE METHODS | VOL.12 NO.4 | APRIL 2015 | 303

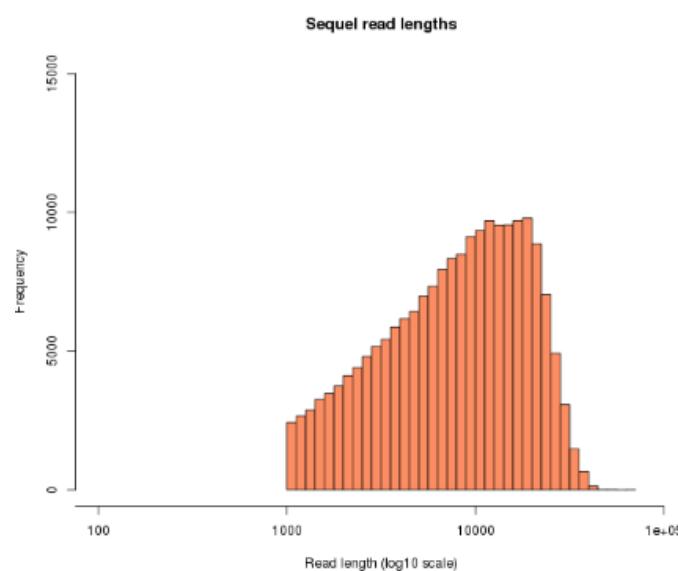
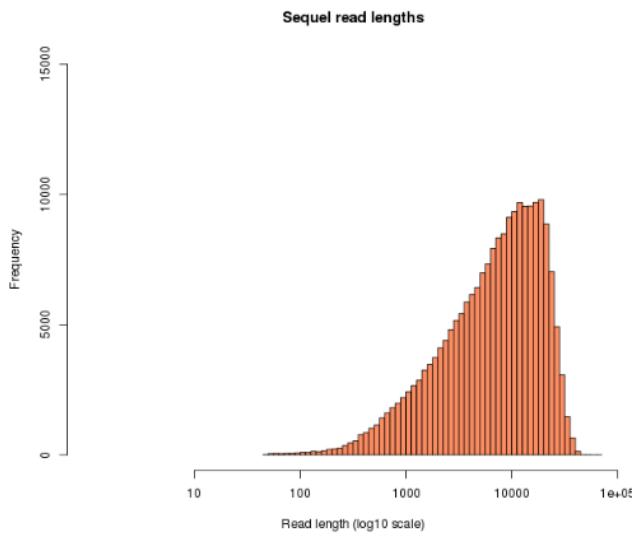
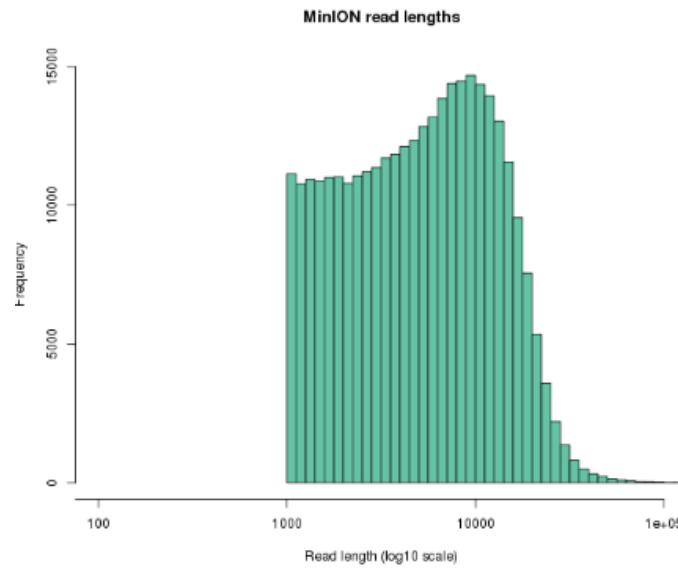
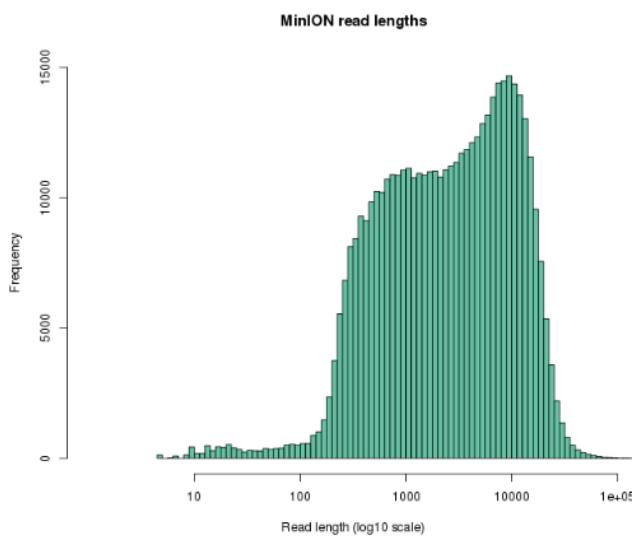
ng techniques. Currently, several groups are exploring how its features - including real-time data analysis and reporting - can be used to answer a range of biological questions.

Cloud-based data analysis allows for real-time data processing and reporting. The system is designed to be compatible with a range of sample types, including blood and serum.

The MinION device is designed to sense from complex samples such as blood/serum.



# Comparación de las plataformas



# Comparación de las plataformas

**Table 1** Comparison of long-read sequencing methods

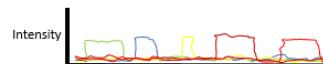
	Device	Device costs	Output max (avg) in Gb	Read length (avg/max)	Costs per Gb	Run time	Accuracy		Multiplexing capacity
							Long read	Consensus read	
ONT	Flongle	\$	2 (1)	5–35 kb/>2 Mb	\$\$\$	0.5–48 h	>Q10	>Q30	96
	MinION	\$	30 (15)		\$\$				
	GridION	\$\$	150 (75)		\$\$				
	PromethION	\$\$\$	15 Tb (4–6 Tb)		\$	0.5–72 h			
PacBio	RS II	\$\$\$\$	2 (1)	5–15 kb/>60 kb	\$\$\$\$	0.5–6 h	>Q10	–	384
	Sequel	\$\$\$	50 <sup>a</sup> /20 <sup>b</sup> (8–10)	5–30 kb/>200 kb	\$\$\$	0.5–20 h		>Q30	
	Sequel II	?	300 <sup>a</sup> /100 <sup>b</sup> (?)		?	0.5–30 h			

PacBio  
SMRT seq

DNA passes thru polymerase in an illuminated volume



Raw output is fluorescent signal of the nucleotide incorporation, specific to each nucleotide

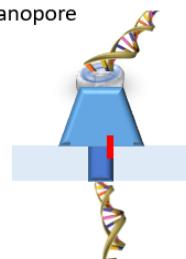


A,C,T,G have known pulse durations, which are used to infer methylated nucleotides



Oxford  
Nanopore

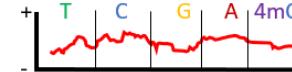
DNA passes thru nanopore



Raw output is electrical signal caused by nucleotide blocking ion flow in nanopore



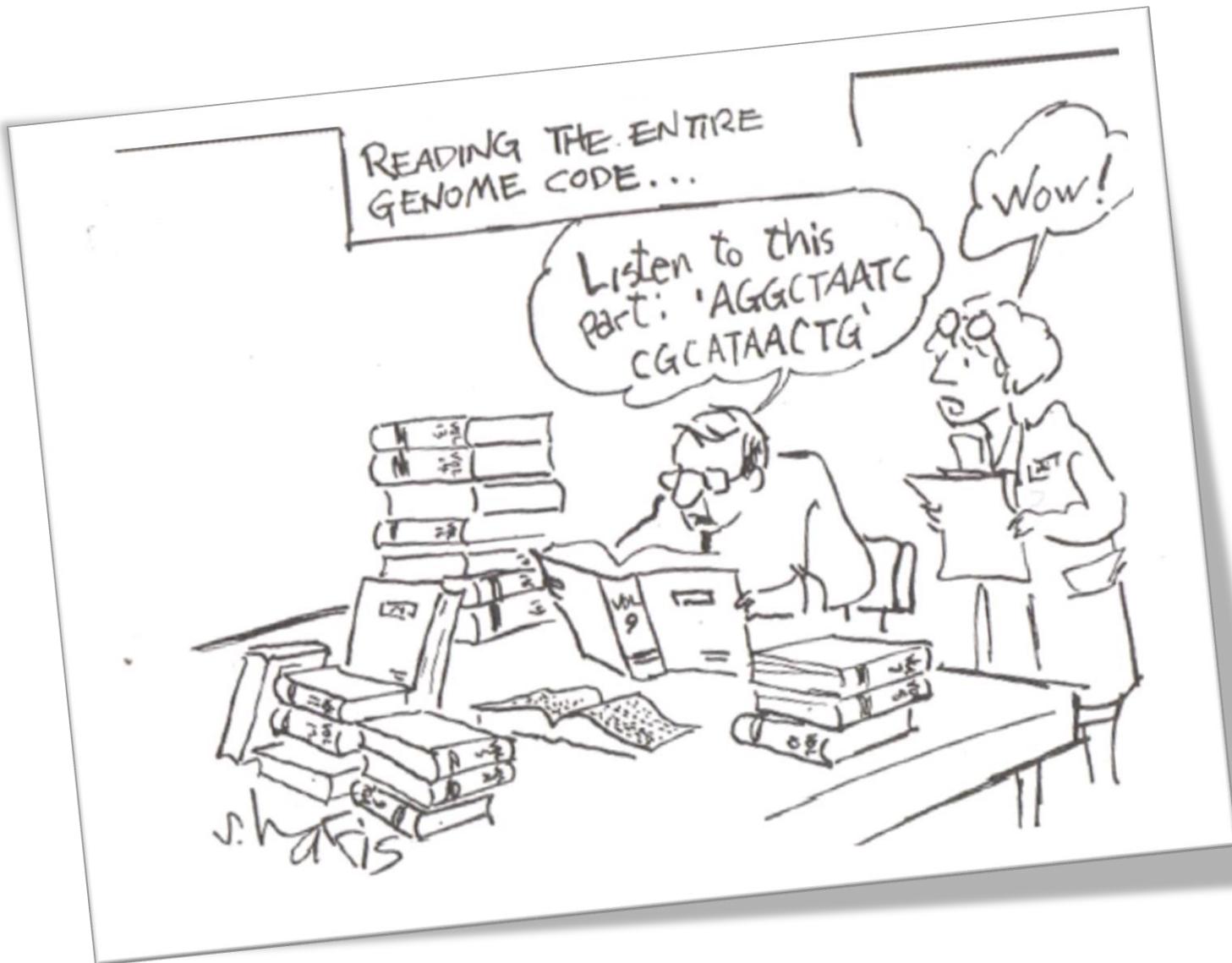
Each nucleotide has a specific electric “signature”



# Software de análisis de secuencias largas

Application	Tool		URL
Basecalling	Guppy	N	<a href="https://community.nanoporetech.com/downloads">https://community.nanoporetech.com/downloads</a>
	SMRT Analysis	P	<a href="https://www.pacb.com/support/software-downloads/">https://www.pacb.com/support/software-downloads/</a>
Alignment	BLASR	P	<a href="https://github.com/PacificBiosciences/blasr">https://github.com/PacificBiosciences/blasr</a>
	LAST	N/P	<a href="http://last.cbrc.jp/">http://last.cbrc.jp/</a>
	minimap2	N/P	<a href="https://github.com/lh3/minimap2">https://github.com/lh3/minimap2</a>
	NGMLR	N/P	<a href="https://github.com/philres/ngmlr">https://github.com/philres/ngmlr</a>
Assembly	Canu	N/P	<a href="https://github.com/marbl/canu">https://github.com/marbl/canu</a>
	FALCON	N/P	<a href="https://github.com/PacificBiosciences/FALCON">https://github.com/PacificBiosciences/FALCON</a>
	Flye	N/P	<a href="https://github.com/fenderglass/Flye">https://github.com/fenderglass/Flye</a>
	wtdbg2	N/P	<a href="https://github.com/ruanjue/wtdbg2">https://github.com/ruanjue/wtdbg2</a>
Error correction	HALC	P	<a href="https://github.com/lanl001/halc">https://github.com/lanl001/halc</a>
	Medeka	N	<a href="https://github.com/nanoporetech/medaka">https://github.com/nanoporetech/medaka</a>
	Nanopolish	N	<a href="https://github.com/jts/nanopolish">https://github.com/jts/nanopolish</a>
SV calling	Nplnv	N	<a href="https://github.com/haojingshao/nplnv">https://github.com/haojingshao/nplnv</a>
	Pbsv	P	<a href="https://github.com/PacificBiosciences/pbsv">https://github.com/PacificBiosciences/pbsv</a>
	Sniffles	N/P	<a href="https://github.com/fritzsedlazeck/Sniffles">https://github.com/fritzsedlazeck/Sniffles</a>
	SVIM	N/P	<a href="https://github.com/eldariont/svim">https://github.com/eldariont/svim</a>
SNP calling	DeepVariant	N/P	<a href="https://github.com/google/deepvariant">https://github.com/google/deepvariant</a>
	GATK	N/P	<a href="https://software.broadinstitute.org/gatk/">https://software.broadinstitute.org/gatk/</a>
	Medeka	N	<a href="https://github.com/nanoporetech/medaka">https://github.com/nanoporetech/medaka</a>
	Nanopolish	N	<a href="https://github.com/jts/nanopolish">https://github.com/jts/nanopolish</a>
DNA/RNA modifications	Nanopolish	N	<a href="https://github.com/jts/nanopolish">https://github.com/jts/nanopolish</a>
	SMRT Analysis	P	<a href="https://www.pacb.com/support/software-downloads/">https://www.pacb.com/support/software-downloads/</a>
	Tombo	N	<a href="https://github.com/nanoporetech/tombo">https://github.com/nanoporetech/tombo</a>
Tandem repeat analysis	NanoSatellite	N	<a href="https://github.com/arnederoeck/NanoSatellite">https://github.com/arnederoeck/NanoSatellite</a>
	nanoSTRique	N	–
	Pbsv	P	<a href="https://github.com/PacificBiosciences/pbsv">https://github.com/PacificBiosciences/pbsv</a>
	STRetch	P	<a href="https://github.com/Oshlack/STRetch">https://github.com/Oshlack/STRetch</a>

# Que hacemos con la secuencia??



# Formato Fastq

- Formato de salida de los secuenciadores de nueva generación
  - Incluye información de secuencia
  - Incluye información de calidad para cada base (confiabilidad)

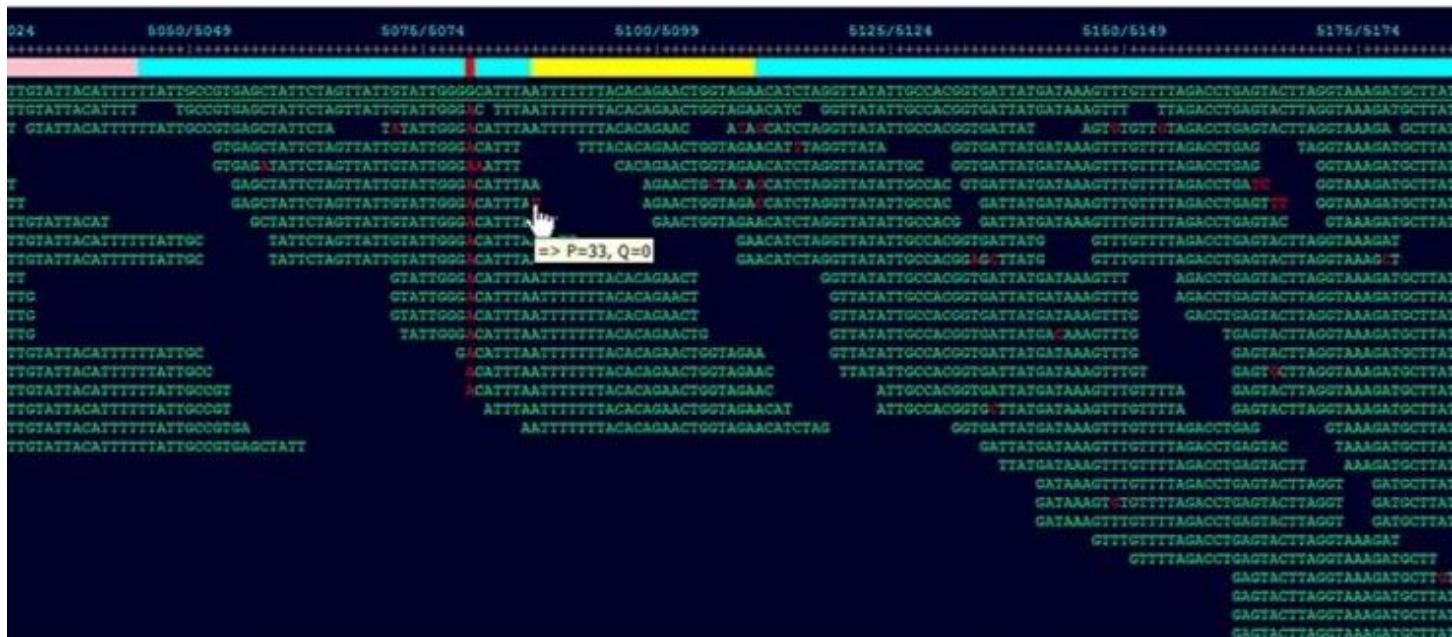
```
@SEQ_ID
GATTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTGTTCAACTCACAGTT
+
!'''*((((*++))%%%++)(%%%%).1***-+*'')**55CCF>>>>CCCCCCCC65
```

```
@SRR001666.1 071112_SLXA-EAS1_s_7:5:1:817:345 length=36
GGGTGATGGCCGCTGCCGATGGCGTCAAATCCCACC
+SRR001666.1 071112_SLXA-EAS1_s_7:5:1:817:345 length=36
IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII9IG9IC
```

Dec	Hx	Oct	Char		Dec	Hx	Oct	Html	Chr		Dec	Hx	Oct	Html	Chr		Dec	Hx	Oct	Html	Chr
0	0	000	<b>NUL</b>	(null)	32	20	040	&#32;	<b>Space</b>		64	40	100	&#64;	<b>Ø</b>		96	60	140	&#96;	`
1	1	001	<b>SOH</b>	(start of heading)	33	21	041	&#33;	<b>!</b>		65	41	101	&#65;	<b>A</b>		97	61	141	&#97;	<b>a</b>
2	2	002	<b>STX</b>	(start of text)	34	22	042	&#34;	<b>"</b>		66	42	102	&#66;	<b>B</b>		98	62	142	&#98;	<b>b</b>
3	3	003	<b>ETX</b>	(end of text)	35	23	043	&#35;	<b>#</b>		67	43	103	&#67;	<b>C</b>		99	63	143	&#99;	<b>c</b>
4	4	004	<b>EOT</b>	(end of transmission)	36	24	044	&#36;	<b>\$</b>		68	44	104	&#68;	<b>D</b>		100	64	144	&#100;	<b>d</b>
5	5	005	<b>ENQ</b>	(enquiry)	37	25	045	&#37;	<b>%</b>		69	45	105	&#69;	<b>E</b>		101	65	145	&#101;	<b>e</b>
6	6	006	<b>ACK</b>	(acknowledge)	38	26	046	&#38;	<b>&amp;</b>		70	46	106	&#70;	<b>F</b>		102	66	146	&#102;	<b>f</b>
7	7	007	<b>BEL</b>	(bell)	39	27	047	&#39;	<b>'</b>		71	47	107	&#71;	<b>G</b>		103	67	147	&#103;	<b>g</b>
8	8	010	<b>BS</b>	(backspace)	40	28	050	&#40;	<b>(</b>		72	48	110	&#72;	<b>H</b>		104	68	150	&#104;	<b>h</b>
9	9	011	<b>TAB</b>	(horizontal tab)	41	29	051	&#41;	<b>)</b>		73	49	111	&#73;	<b>I</b>		105	69	151	&#105;	<b>i</b>
10	A	012	<b>LF</b>	(NL line feed, new line)	42	2A	052	&#42;	<b>*</b>		74	4A	112	&#74;	<b>J</b>		106	6A	152	&#106;	<b>j</b>
11	B	013	<b>VT</b>	(vertical tab)	43	2B	053	&#43;	<b>+</b>		75	4B	113	&#75;	<b>K</b>		107	6B	153	&#107;	<b>k</b>
12	C	014	<b>FF</b>	(NP form feed, new page)	44	2C	054	&#44;	<b>,</b>		76	4C	114	&#76;	<b>L</b>		108	6C	154	&#108;	<b>l</b>
13	D	015	<b>CR</b>	(carriage return)	45	2D	055	&#45;	<b>-</b>		77	4D	115	&#77;	<b>M</b>		109	6D	155	&#109;	<b>m</b>
14	E	016	<b>SO</b>	(shift out)	46	2E	056	&#46;	<b>.</b>		78	4E	116	&#78;	<b>N</b>		110	6E	156	&#110;	<b>n</b>
15	F	017	<b>SI</b>	(shift in)	47	2F	057	&#47;	<b>/</b>		79	4F	117	&#79;	<b>O</b>		111	6F	157	&#111;	<b>o</b>
16	10	020	<b>DLE</b>	(data link escape)	48	30	060	&#48;	<b>0</b>		80	50	120	&#80;	<b>P</b>		112	70	160	&#112;	<b>p</b>
17	11	021	<b>DC1</b>	(device control 1)	49	31	061	&#49;	<b>1</b>		81	51	121	&#81;	<b>Q</b>		113	71	161	&#113;	<b>q</b>
18	12	022	<b>DC2</b>	(device control 2)	50	32	062	&#50;	<b>2</b>		82	52	122	&#82;	<b>R</b>		114	72	162	&#114;	<b>r</b>
19	13	023	<b>DC3</b>	(device control 3)	51	33	063	&#51;	<b>3</b>		83	53	123	&#83;	<b>S</b>		115	73	163	&#115;	<b>s</b>
20	14	024	<b>DC4</b>	(device control 4)	52	34	064	&#52;	<b>4</b>		84	54	124	&#84;	<b>T</b>		116	74	164	&#116;	<b>t</b>
21	15	025	<b>NAK</b>	(negative acknowledge)	53	35	065	&#53;	<b>5</b>		85	55	125	&#85;	<b>U</b>		117	75	165	&#117;	<b>u</b>
22	16	026	<b>SYN</b>	(synchronous idle)	54	36	066	&#54;	<b>6</b>		86	56	126	&#86;	<b>V</b>		118	76	166	&#118;	<b>v</b>
23	17	027	<b>ETB</b>	(end of trans. block)	55	37	067	&#55;	<b>7</b>		87	57	127	&#87;	<b>W</b>		119	77	167	&#119;	<b>w</b>
24	18	030	<b>CAN</b>	(cancel)	56	38	070	&#56;	<b>8</b>		88	58	130	&#88;	<b>X</b>		120	78	170	&#120;	<b>x</b>
25	19	031	<b>EM</b>	(end of medium)	57	39	071	&#57;	<b>9</b>		89	59	131	&#89;	<b>Y</b>		121	79	171	&#121;	<b>y</b>
26	1A	032	<b>SUB</b>	(substitute)	58	3A	072	&#58;	<b>:</b>		90	5A	132	&#90;	<b>Z</b>		122	7A	172	&#122;	<b>z</b>
27	1B	033	<b>ESC</b>	(escape)	59	3B	073	&#59;	<b>;</b>		91	5B	133	&#91;	<b>[</b>		123	7B	173	&#123;	<b>{</b>
28	1C	034	<b>FS</b>	(file separator)	60	3C	074	&#60;	<b>&lt;</b>		92	5C	134	&#92;	<b>\</b>		124	7C	174	&#124;	<b> </b>
29	1D	035	<b>GS</b>	(group separator)	61	3D	075	&#61;	<b>=</b>		93	5D	135	&#93;	<b>]</b>		125	7D	175	&#125;	<b>}</b>
30	1E	036	<b>RS</b>	(record separator)	62	3E	076	&#62;	<b>&gt;</b>		94	5E	136	&#94;	<b>^</b>		126	7E	176	&#126;	<b>~</b>
31	1F	037	<b>US</b>	(unit separator)	63	3F	077	&#63;	<b>?</b>		95	5F	137	&#95;	<b>_</b>		127	7F	177	&#127;	<b>DEL</b>

# Resecuenciado de Genomas

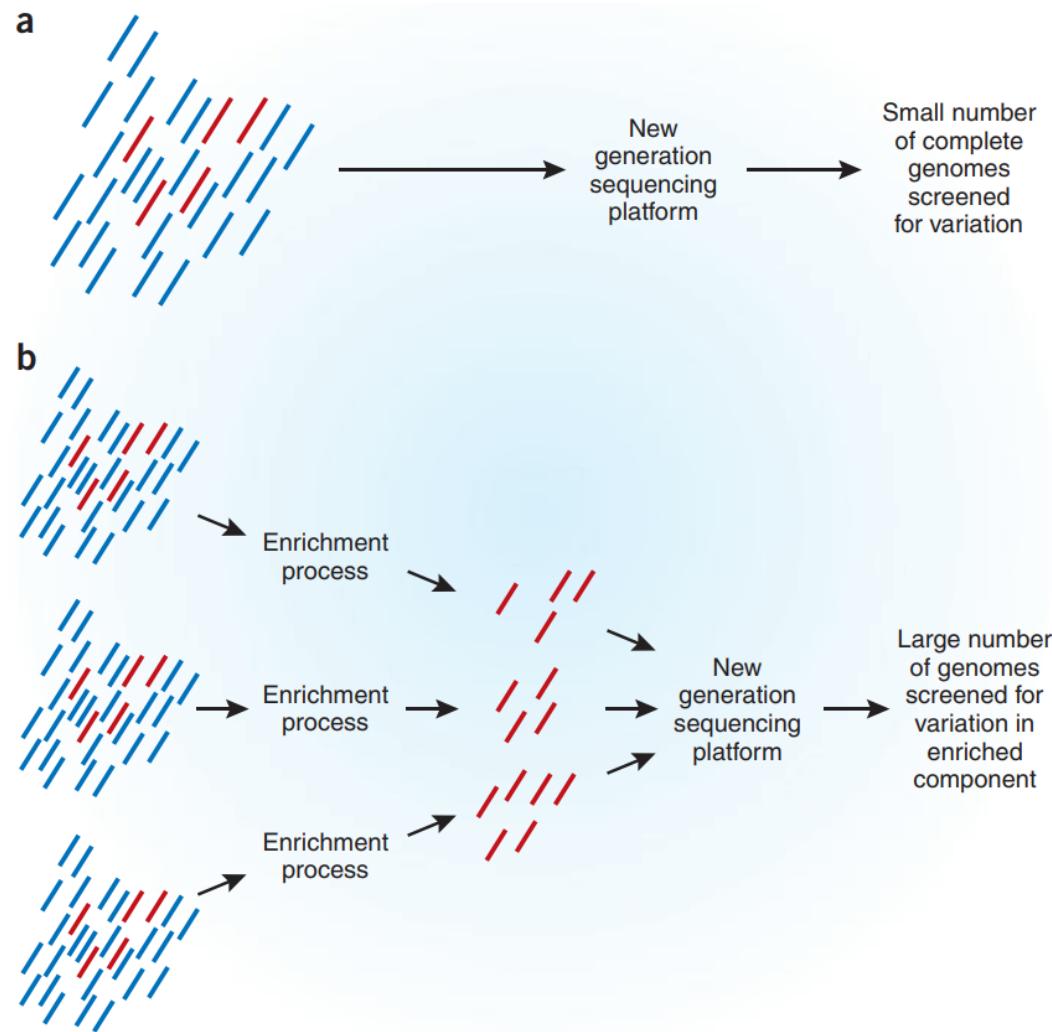
- Comparación a una referencia
  - Se obtienen las lecturas y se comparan a una referencia
    - Búsqueda de variantes, reordenamientos, etc



# Mapeo de lecturas



# Resecuenciado de Genomas



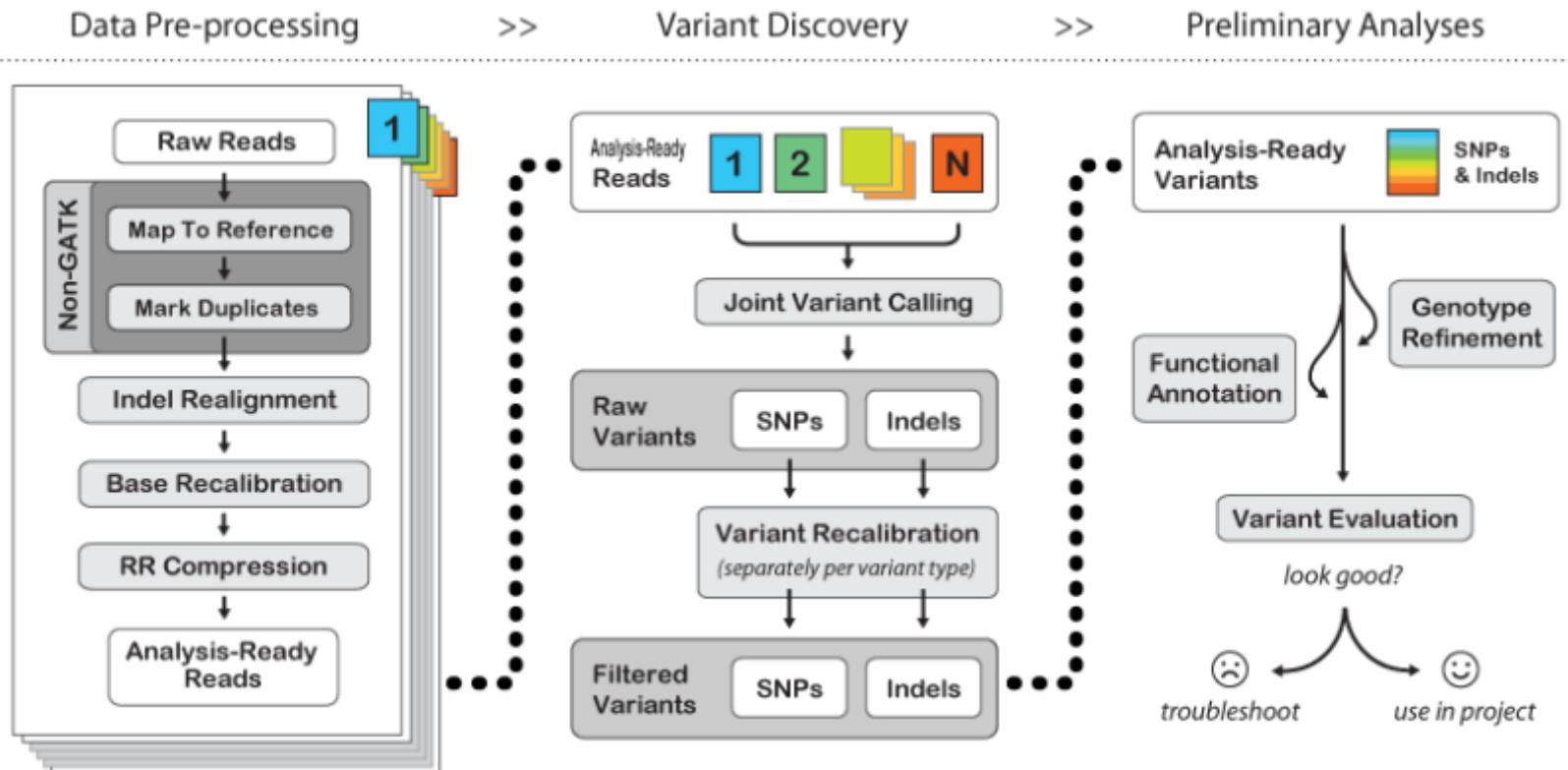
**Figure 1** Enrichment procedures allow redistribution of sequencing throughput from all of a small number of genomes (**a**) to a small component of a large number of genomes (**b**).

# Resecuenciado de Genomas

```
##fileformat=VCFv4.0
##fileDate=20090805
##source=myImputationProgramV3.1
##reference=1000GenomesPilot-NCBI36
##phasing=partial
##INFO=<ID=NS,Number=1,Type=Integer,Description="Number of Samples With Data">
##INFO=<ID=DP,Number=1,Type=Integer,Description="Total Depth">
##INFO=<ID=AF,Number=.,Type=Float,Description="Allele Frequency">
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele">
##INFO=<ID=DB,Number=0,Type=Flag,Description="dbSNP membership, build 129">
##INFO=<ID=H2,Number=0,Type=Flag,Description="HapMap2 membership">
##FILTER=<ID=q10,Description="Quality below 10">
##FILTER=<ID=s50,Description="Less than 50% of samples have data">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##FORMAT=<ID=HQ,Number=2,Type=Integer,Description="Haplotype Quality">
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT NA00001
NA00002 NA00003
20 14370 rs6054257 G A 29 PASS NS=3;DP=14;AF=0.5;DB;H2 GT:GQ:DP:HQ 0|0:4
8:1:51,51 1|0:48:8:51,51 1/1:43:5:,,.
20 17330 . T A 3 q10 NS=3;DP=11;AF=0.017 GT:GQ:DP:HQ 0|0:4
9:3:58,50 0|1:3:5:65,3 0/0:41:3
20 1110696 rs6040355 A G,T 67 PASS NS=2;DP=10;AF=0.333,0.667;AA=T;DB GT:GQ:DP:HQ 1|2:2
1:6:23,27 2|1:2:0:18,2 2/2:35:4
20 1230237 . T . 47 PASS NS=3;DP=13;AA=T GT:GQ:DP:HQ 0|0:5
4:7:56,60 0|0:48:4:51,51 0/0:61:2
20 1234567 microsat1 GTCT G,GTACT 50 PASS NS=3;DP=9;AA=G GT:GQ:DP 0/1:3
5:4 0/2:17:2 1/1:40:3
```

1. #CHROM
2. POS
3. ID
4. REF
5. ALT
6. QUAL
7. FILTER
8. INFO

# Resecuenciado de Genomas



# Resecuenciado de Genomas

## Effect prediction details

Detailed description of the effect predicted by SnpEff in the `Effect` and `Effect_Impact` sub-fields.

Notes:

- **Effect (Sequence Ontology)** Sequence ontology ([SO](#)) allows to standardize terminology used for assessing sequence changes and impact. This allows for a common language across all variant annotation programs and makes it easier to communicate using a uniform terminology. Starting from version 4.0 VCF output uses SO terms by default.
- **Effect (Classic)** These are the "classic" effect names used by SnpEff, these can be accessed using the `-classic` command line option.
- **Effect Impact** Effects are categorized by 'impact': {High, Moderate, Low, Modifier}. These are pre-defined categories to help users find more significant variants.

⚠ Impact categories must be used with care, they were created only to help and simplify the filtering process. Obviously, there is no way to predict whether a "high impact" or a "low impact" variant is the one producing a phenotype of interest.

Here is a list of effects and some brief explanations:

Effect Seq. Ontology	Effect Classic	Note & Example	Impact
coding_sequence_variant	CDS	The variant hits a CDS.	MODIFIER
chromosome	CHROMOSOME_LARGE_DELETION	A large part (over 1% or 1,000,000 bases) of the chromosome was deleted.	HIGH
duplication	CHROMOSOME_LARGE_DUPLICATION	Duplication of a large chromosome segment (over 1% or 1,000,000 bases).	HIGH
inversion	CHROMOSOME_LARGE_INVERSION	Inversion of a large chromosome segment (over 1% or 1,000,000 bases).	HIGH
coding_sequence_variant	CODON_CHANGE	One or many codons are changed e.g.: An MNP of size multiple of 3	LOW
inframe_insertion	CODON_INSERTION	One or many codons are inserted	MODERATE

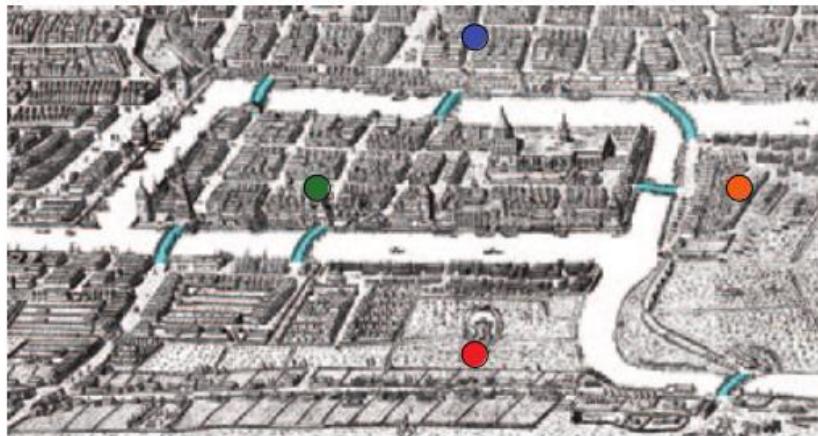
# Ensamblado de secuencias cortas

- Existen ensambladores específicos
  - VELVET
  - Abyss
  - SOAPdenovo
  - ...
- La mayoría se basan en encontrar el camino euleriano en un grafo de de Bruijn construido con k-mers solapantes obtenidos de las lecturas...

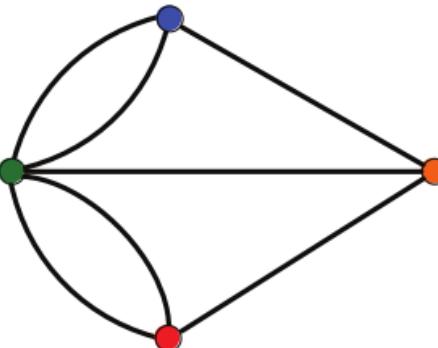
# Ensamblando secuencias cortas

- Problema de los puentes de Königsberg
  - Leonhard Euler (1736)

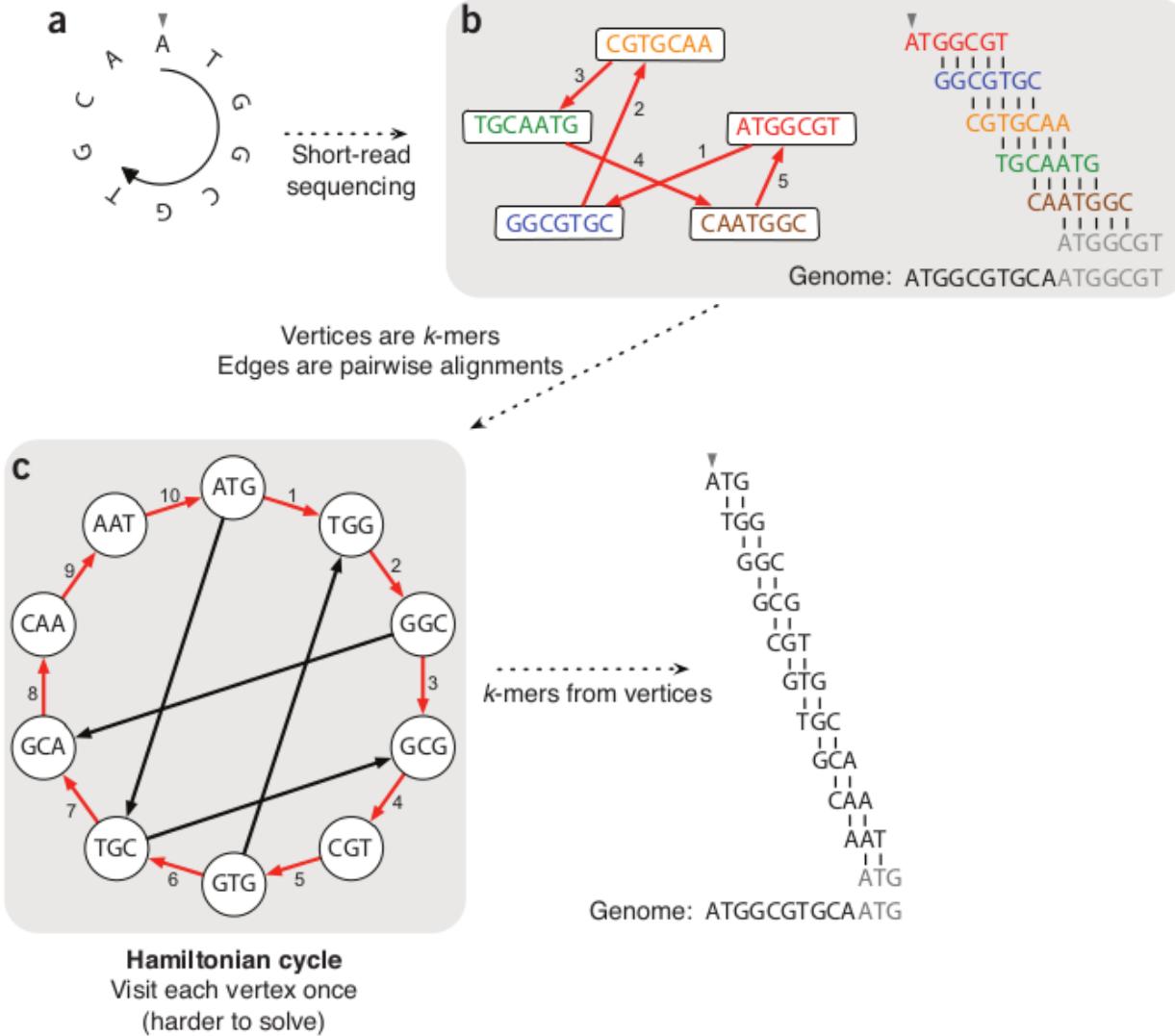
a



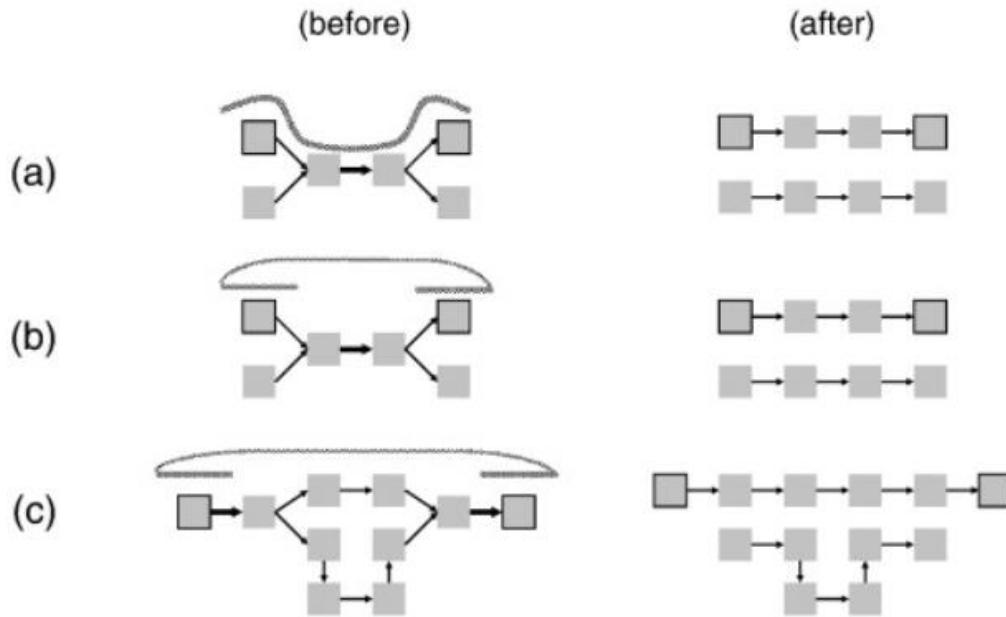
b



# Ensamblando secuencias cortas



# Resolución de nodos problemáticos

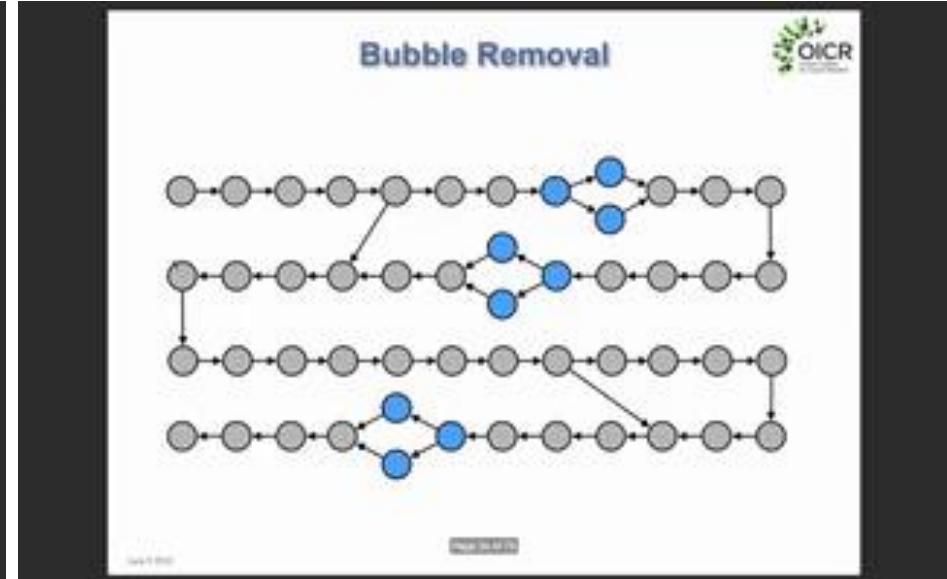
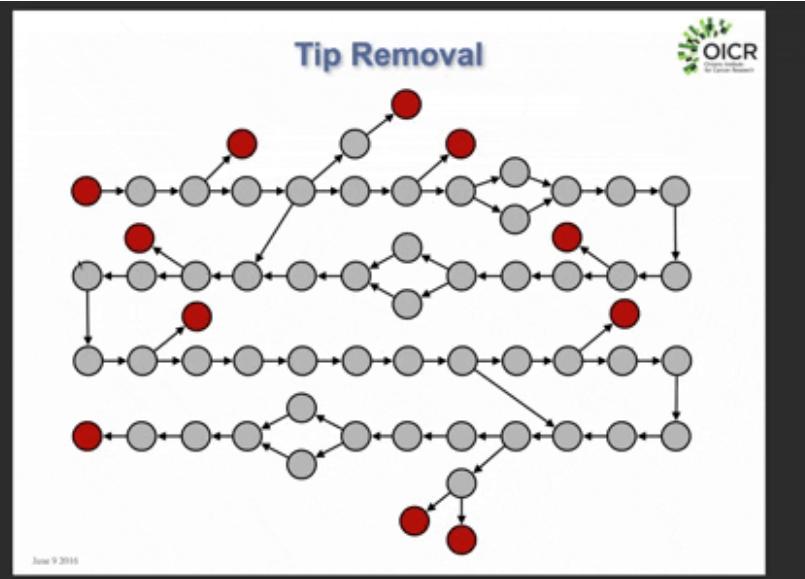


- Mapeo de las lecturas al grafo:
  - (a) Se mapea la lectura la cual atraviesa el nodo complejo
  - (b) Las lecturas pareadas se usan para resolver
  - (c) La distancia entre los pares se usan para resolver

# Simplificación del grafo

Remoción de puntas

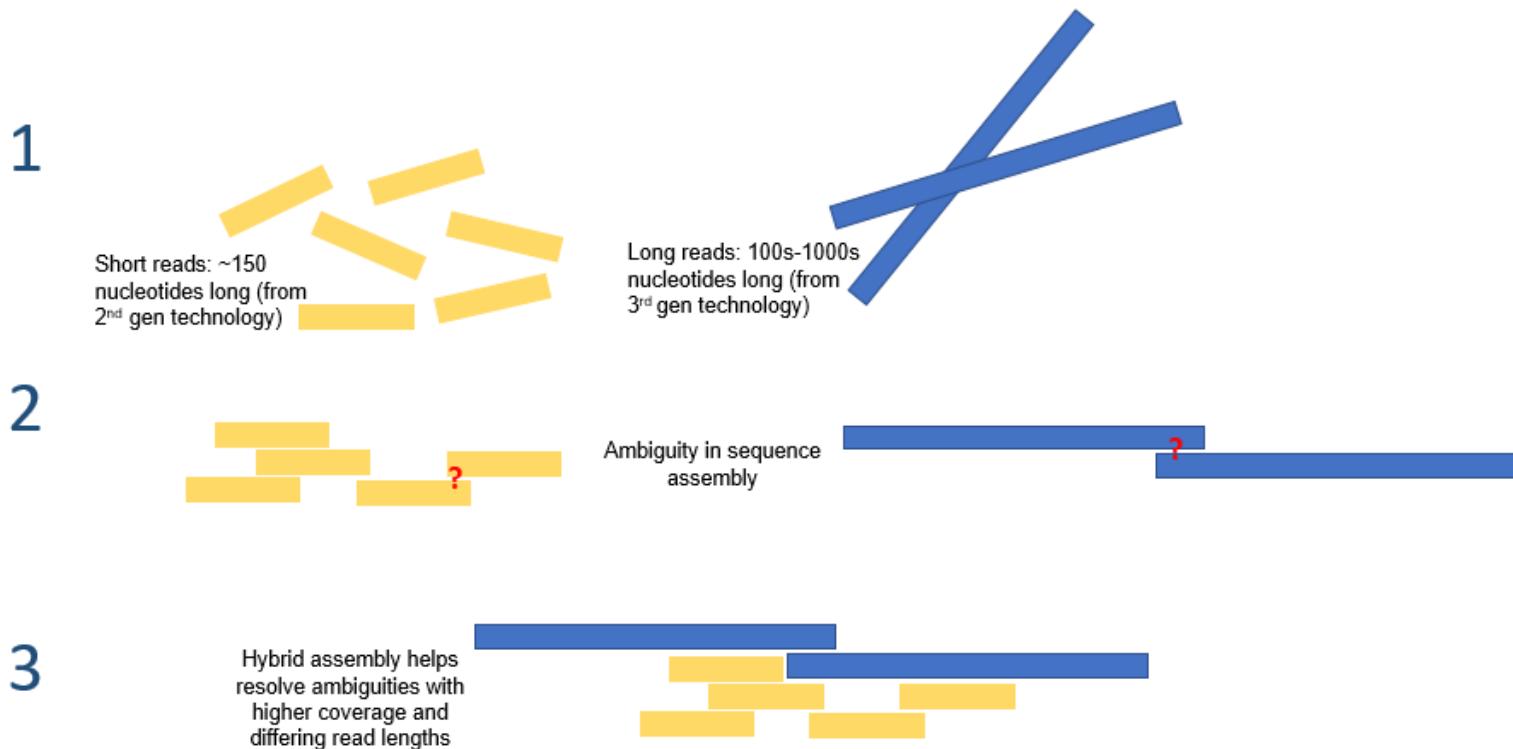
► Remoción de burbujas



- Mayormente basado en la profundidad de lecturas de los nodos

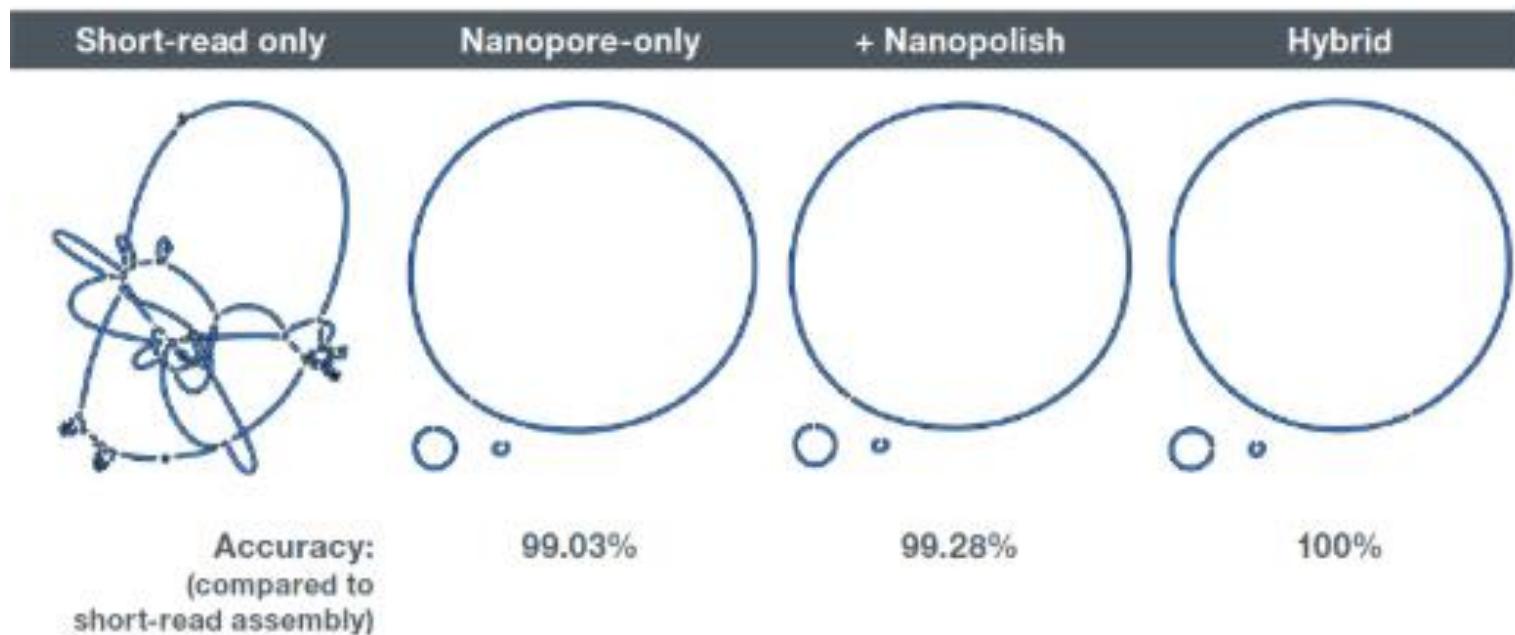
# Ensamblado con lecturas largas

- Único
- Hibrido



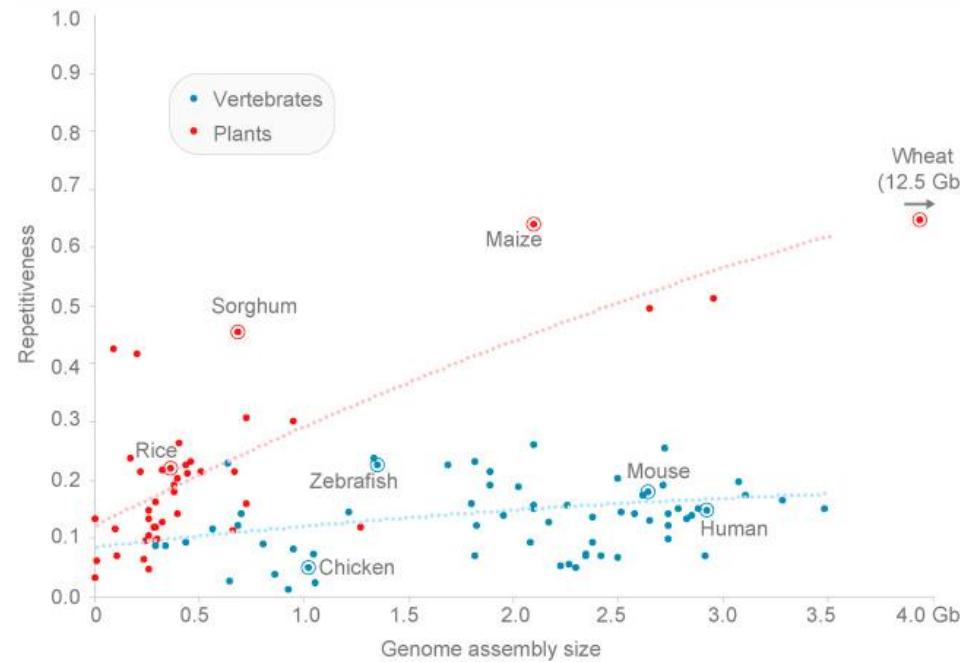
# Ensamblado con lecturas largas

- Único
- Hibrido



# Ensamblado con lecturas largas

- Gran impacto en genomas muy repetitivos



# El ensamblado es el comienzo del proceso

- Búsqueda de genes
  - Codificantes para proteínas
  - Pseudogenes
  - Genes para ARNs funcionales
    - tRNA
    - rRNA
    - snoRNA
    - snRNA
    - miRNA
  - Regiones reguladoras
  - Otras regiones funcionales

TCTGTGGATGTGGAATAC  
TCATATTCCATTCTGC  
AGGAGAAAAACGAGTGACT  
GACTTCATTGATGTACA  
GCTACTGACTATTGGAT  
TGCATGAGCGTGAAATCA  
GTGCGAATCAGACTCATG  
AACGCTGAAAAATAGTTG  
GGGATGGCACCGATTCCCG  
TCCAATCGCGATGGCGCA  
TATTGATGACGAAAACGG  
CGAAAAACGAATGTGCAAA  
CTATGCGTTATTAATGA  
TATGGAGTAGGACTGGCG  
TGAATCTGGTT.....

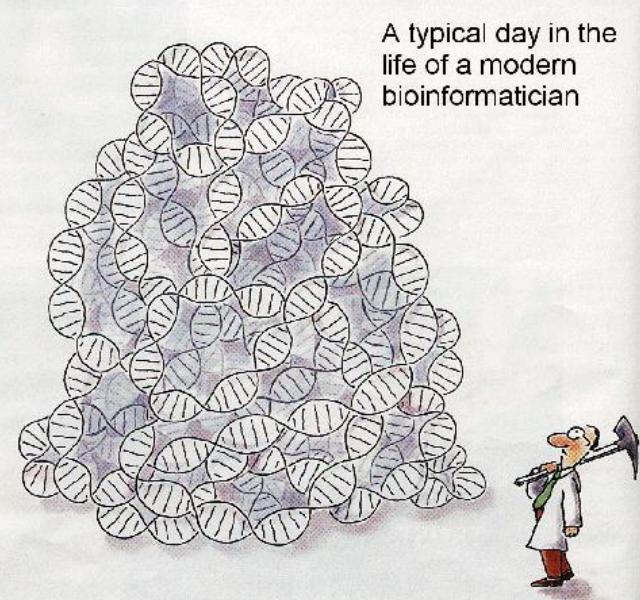
Que defino como gen?

Donde están?

Vamos a intentar encontrar  
patrones

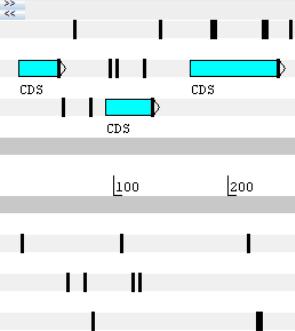
```
Nothing selected
>>
<<
V E Y S Y F H F C R R K R V T D F I * C T A T D Y L D C M S V K S V R I R L M N A E K + L G M A R F P S N R D G A Y * * R K R R K R M C K L C V Y # * Y G V G
W N T H I S I S A G E N E * L T S F D V Q L L T I W I A * A * N Q C E S D S * T L K N S W G W H D S R P I A M A H I D D E N G E N E C A N Y A F I N D M E + D
G I L I F P F L Q E K T S D * L H L M Y S Y * L F G L H E R E I S A N Q T H E R * K I V G D G T I P V Q S R U R I L M T K T A K T N V Q T H R L L M I W S R T
TGGAAATACTCATATTCCATTTCTGCAGGAGAAAACGAGUCAGTGACTTCATTTGATGTACAGCTACTGACTATTTGGATTGCATGAGCGTAAATCAGCTCATGAACCGCTGAAAAAATAGTTGGGATGCCACGATCCCGTCAATCCGATGCCATATTGATGCCAAACGCCAAATGCTCAAATCTGCCTTAAATGATATGGTAGGACT
[20] [40] [60] [80] [100] [120] [140] [160] [180] [200] [220] [240]
ACCTTATGAGTATAAAGGTAAGACGCTCTCTTTGCTCACTGACTGAACTAACATGTCATGACTGATAAAACCTAACGCTACTCGCACCTTACGGCTTAGCTGACTCTTGCGACTTTTTCAACCCCTACCGTGCTAACGGCAGGTTAGCGCTACCGCTATAACTACTGCTTTGCGCTTTGCTTACGTTGATACGCCAATAATTACTATACCTCATCCGAG
H F V * I E M E A P S F S H S V E N S T C S S V I Q I A H A H F * H S D S E H V S F F L Q P H C S E R G I A I A C I S S S F P S F S H A F + A N I L S I S Y S C
P I S M N G N R C S F V L S Q S * K I Y L + Q S N P N C S R S I L A F * V * S R Q F I T P S P V I G T W D R H R M N I V F V A F V F T C V I R K N I I H L L V
S Y E Y K W K Q L L F R T V S K M Q H V A V S + K S Q M L T F D T R I L S M F A S F Y N P I A R N G D L R S P A Y Q H R F R R F R I H L S H T # # H Y P T P S
<<
```

A typical day in the  
life of a modern  
bioinformatician



Entry:  GENETEST.SEQ  GENETEST.TAB

Nothing selected



V E Y S Y F H F C R R K R V T D F I \* C T A T D Y L D C M S V K S V R I R L M N A E K + L G M A R F P S N R D G A Y \* \* R K R R K R N C K L C V Y # \* Y G V G I  
W N T H I S I S A G E N E \* D L T S F D V Q L L T I W I A \* A \* N Q C E S D S \* T L K N S W G W H D S R P I A M A H I D D E N G E N E C A N Y A F I N D M E + D  
G I L I F P F L Q E K T S D \* L H L M Y S Y \* L F G L H E R E I S A N O T H E R \* D K I V G D G T I P V Q S R U R I L M T K T A K T N V Q T M R L L M I W S R T  
T G G A A T A C T C A T A T T C C A T T T C T G C A G G A G A A A C G A G U G A C T G A C T T C A T T T G A T G T G C A T G A G G C T G A A A T C A G T G C C A A T C A G C T G A A A A T A G T T G G G A T G G C A C G A T T C C C G T C C A A T C G C A T G G C G A T A T T G A T G A C G A A A A C G G C G A A A C G A A T G T G C A A A C T A T G C G T T T A T T A A T G A T A T G G A G T A G G A C T  
A C C T T A T G A T G A T A A A G G T A A A G A C G T C T C T T T G C T C A C T G A C T G A A G T A A A C T A C A T G T G A T G A C T G A C T G A T G A C T G A T G A C T G A C T T T G C T C A C G C T T A C C G C T A A G G C C A G G T T A C C G C T A A C T A C T G C T T T T G C C G T T T T G C T T A C C G T T T G A T A C G C A A A T A A T T A C T A T A C C T C A T C C G A  
H F V \* I E M E A P S F S H S V E N S T C S S V I Q I A H A H F \* H S D S E H V S F F L Q P H C S E R G I A I A C I S S S F P S F S H A F + A N I L S I S Y S C  
P I S H N G N R C S F V L S Q S \* K I Y L + Q S N P N C S R S I L A F \* V \* S R Q F I T P S P V I G T W D R H R M N I V F V A F V F T C V I R K N I I H L L V  
S Y E Y K W K Q L L F R T V S K M Q H V A V S + K S Q M L T F D T R I L S M F A S F Y N P I A R N G D L R S P A Y Q H R F R R F R I H L S H T # # H Y P T P S

20 40 60 80 100 120 140 160 180 200 220 240

CDS 17 52  
CDS 93 134  
CDS 167 244

# Anotación del genoma secuenciado:

- Los genes para ARNs funcionales son difíciles de encontrar.
  - Pequeños
  - Sin ORFs
- Las regiones codificantes tienen características identificables
  - En procariotas se identifican de manera relativamente sencilla dada la alta densidad de genes
  - En eucariotas es mas complejo debido a que los exones son pequeños comparados con los intrones y regiones intergenicas

# Lowe Lab

## tRNAscan-SE Search Server



Search for tRNA genes in genomic sequence

## tRNAscan-SE 1.21

The principles underlying the tRNAscan-SE program are described in:

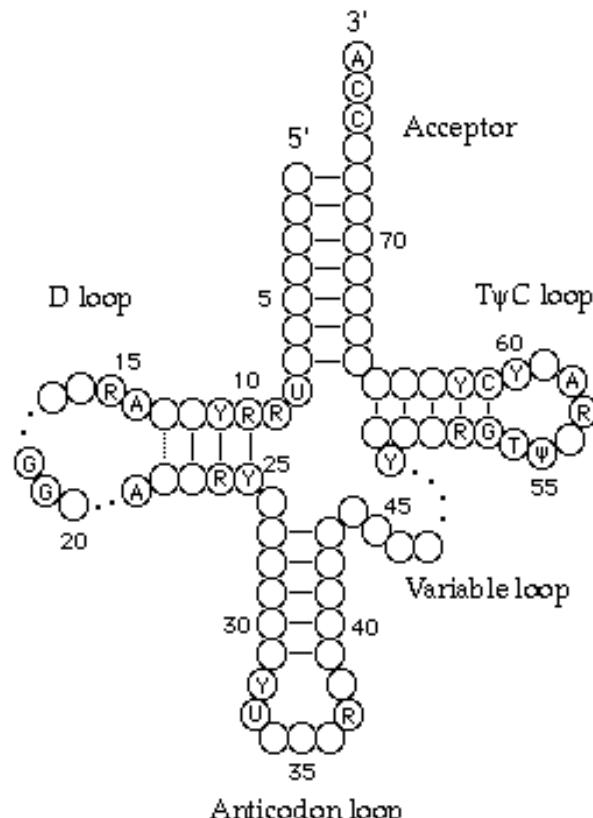
Lowe, T.M. and Eddy, S.R. (1997)

tRNAscan-SE: a program for improved detection of transfer RNA genes in genomic sequence.

Nucleic Acids Res, 25, 955-964.

Instructions for using the tRNAscan-SE server and interpreting the output can be found in the [tRNAscan-SE README file](#).

If you would like to run tRNAscan-SE locally, you can get the UNIX [source code](#) (gzip'd tar file).



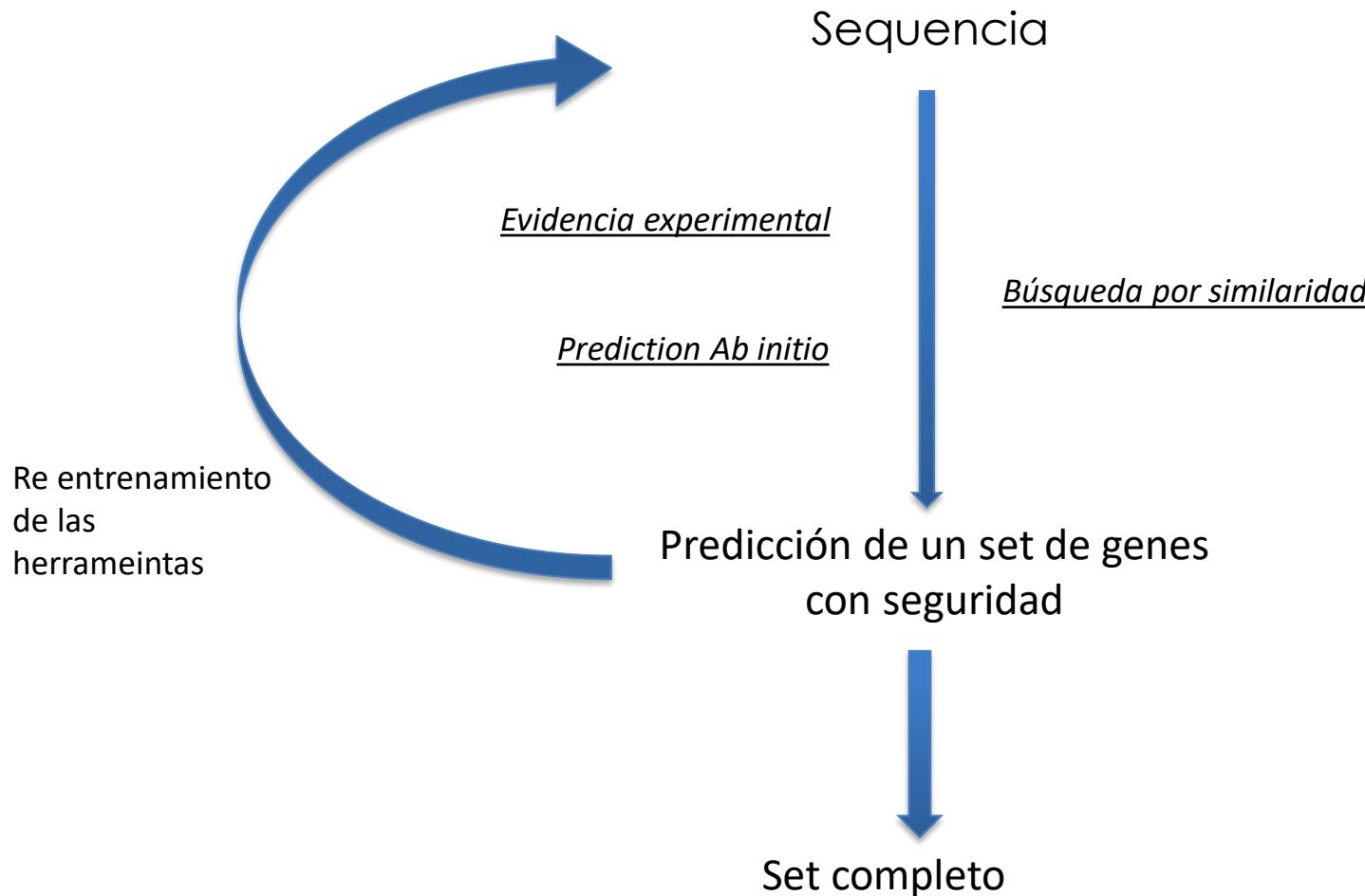
# Anotación del genoma secuenciado:

- Búsqueda de ORFs
  - Marcos de lectura abiertos
    - Comienzan con un codón de inicio ATG
    - Finalizan en un codón de terminación
  - Son las **posibles** regiones codificantes

# Diferentes aproximaciones

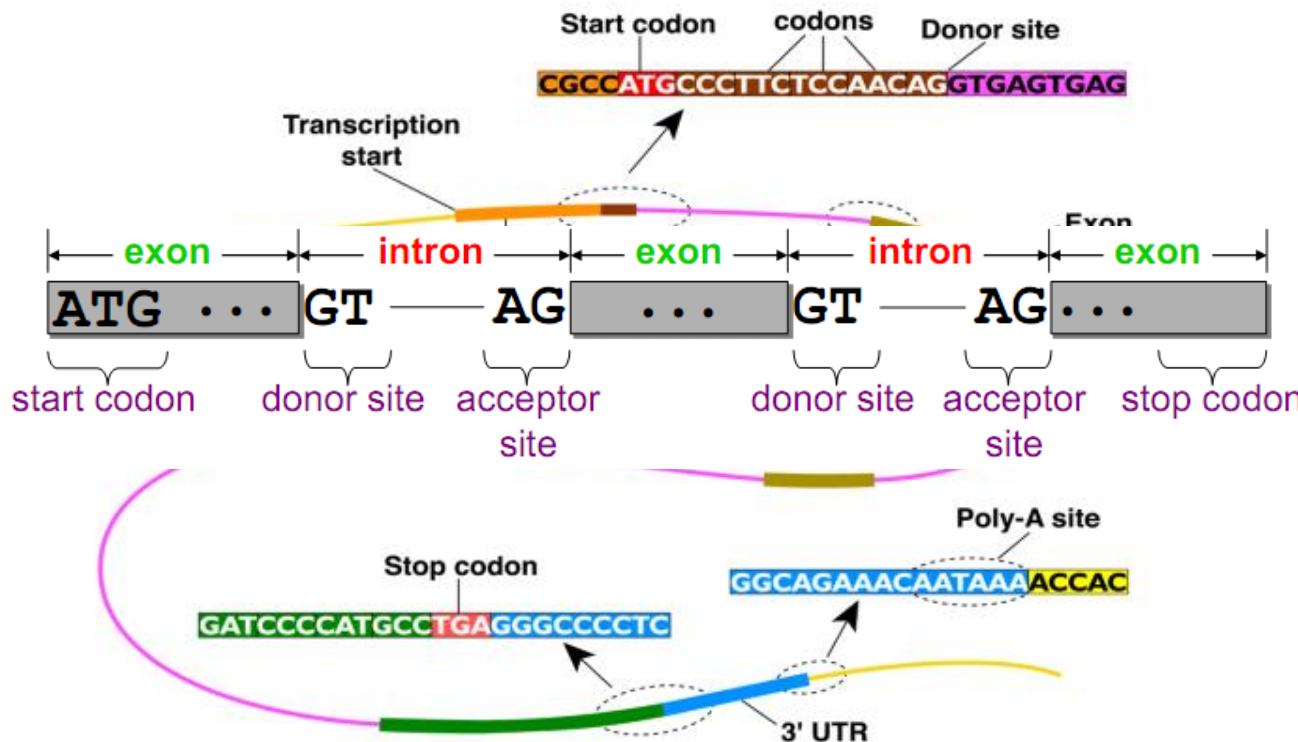
- Métodos predictivos
  - Buscan directamente en la secuencia de acuerdo a reglas definidas
- Métodos comparativos
  - Comparan con otras secuencias y extraen información (conservación)
- Métodos experimentales
  - Utilizan información obtenida de forma experimental

# Procedimiento típico



# Métodos predictivos

- Búsqueda de señales: motivos cortos delimitan las secuencias codificantes



# El contenido en GC varía entre los distintos genomas y dentro de un genoma determinado

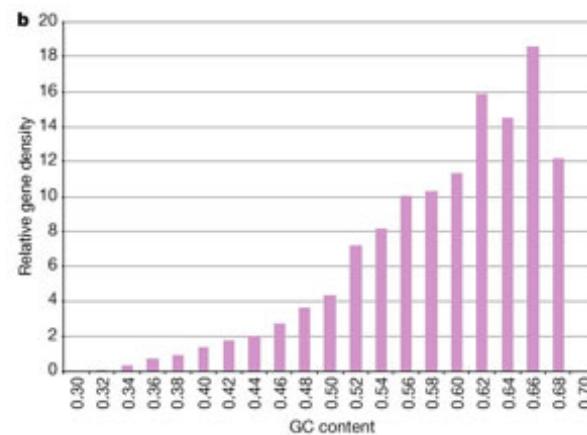
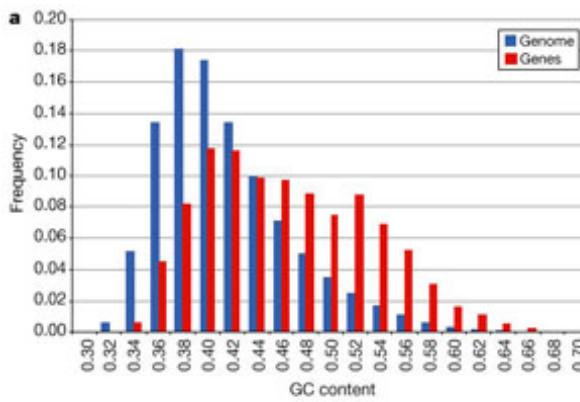
*Escherichia coli*

51%

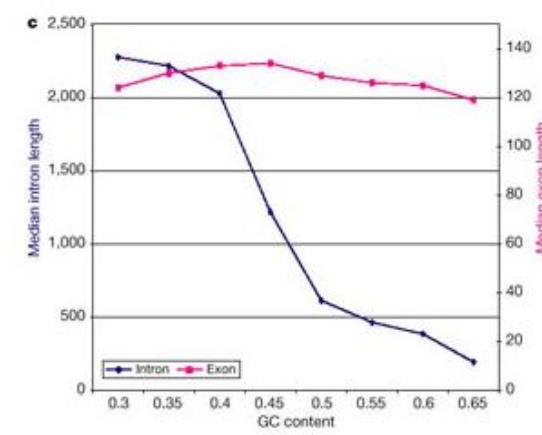
*Campylobacter jejuni*

31%

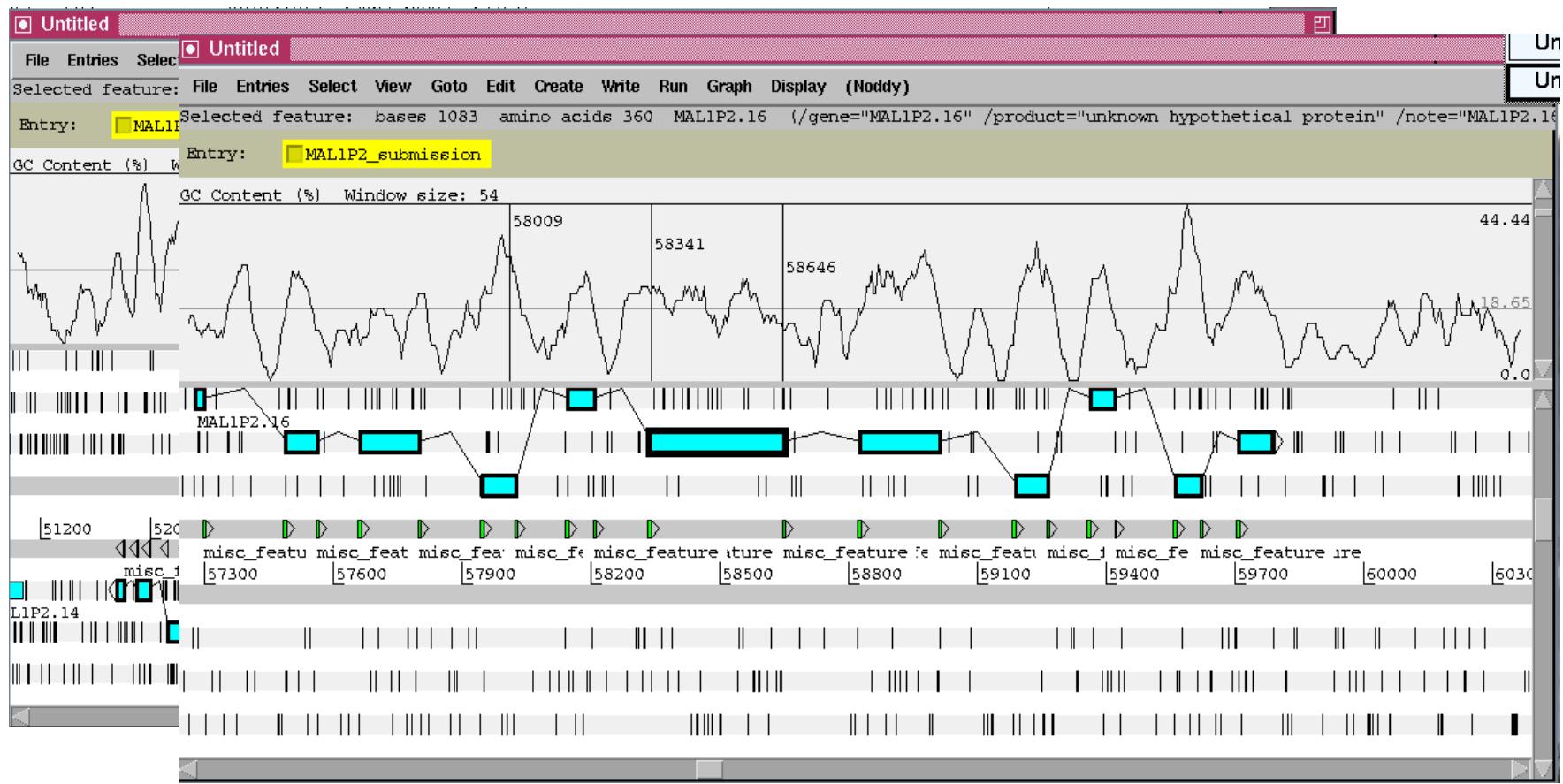
- La variación esta mas restringida en las regiones codificantes
- Estas regiones tienen un contenido GC mas alto (fenómeno mas claro en genomas ricos en AT)



		Second nucleotide			
		U	C	A	G
First nucleotide	U	UUU Phe UUC UUA Leu UUG	UCU UCC UCA Ser UCG	UAU Tyr UAC UAA STOP UAG STOP	UGU Cys UGC UGA STOP UGG Trp
	C	CUU CUC Leu CUA CUG	CCU CCC CCA Pro CCG	CAU His CAC CAA Gln	CGU CGC CGA Arg CGG
A	A	AUU Ile AUC AUA Met AUG	ACU ACC ACA Thr ACG	AAU Asn AAC AAA Lys AAG	AGU Ser AGC AGA Arg AGG
	G	GUU GUC Val GUA GUG	GCU GCC GCA Ala GCG	GAU Asp GAC GAA Glu GAG	GGU GGC GGA Gly GGG
		Third nucleotide			
		U	C	A	G



# Visualización en Artemis



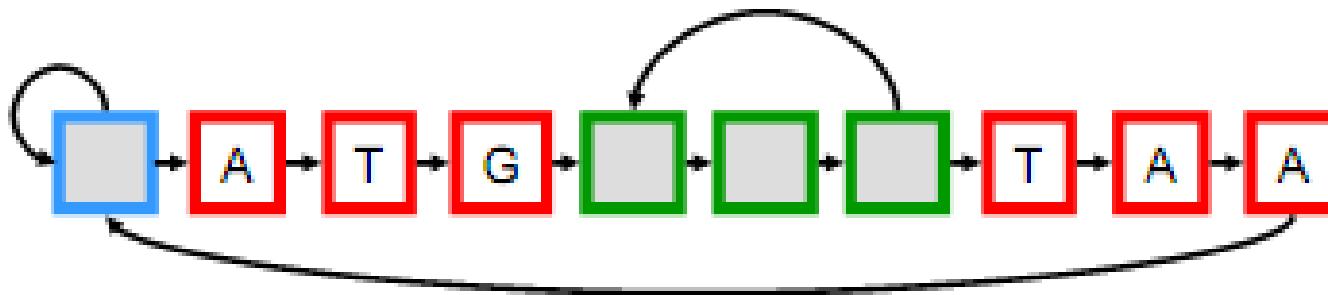
# Métodos predictivos

- Los primeros predictores hacían uso de reglas como las que comentamos antes.
  - Problemas:
    - Muchas reglas!
    - Muchas excepciones!
    - Los predictores se hacen demasiados complejos
- El análisis estadístico de las secuencias mejora las predicciones
  - Modelos ocultos de Markov

# Modelos estadísticos (HMM)

- Los predictores más recientes toman en cuenta modelos estadísticos para realizar la anotación
- Los observables son cada una de las bases
- Cada observable puede estar en diferentes estados
  - Exón
  - Intrón
  - Intergénica
- Cada uno de los estados tiene probabilidades de emisión particulares para cada una de las bases
- Los parámetros se “aprenden” de genes conocidos

# Modelos estadísticos (HMM)



AAAGC ATG CAT TTA ACG AGA GCA CAA GGG CTC TAA TGCCG

The sequence of states is an annotation of the generated string – each nucleotide is generated in intergenic, start/stop, coding state

# Modelos estadísticos (HMM)

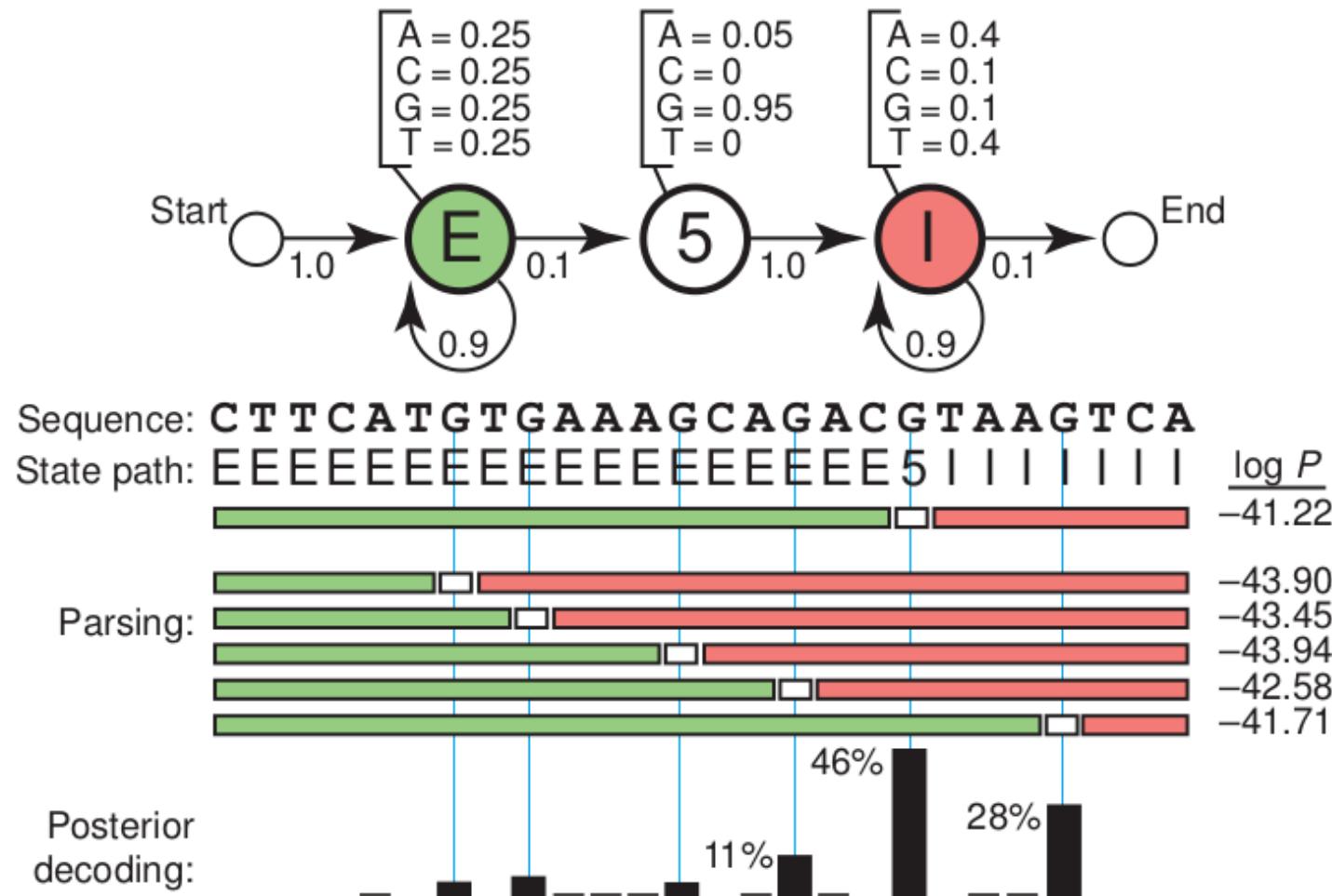


Figure 1 A toy HMM for 5' splice site recognition. See text for explanation.

# Modelos estadísticos (HMM)

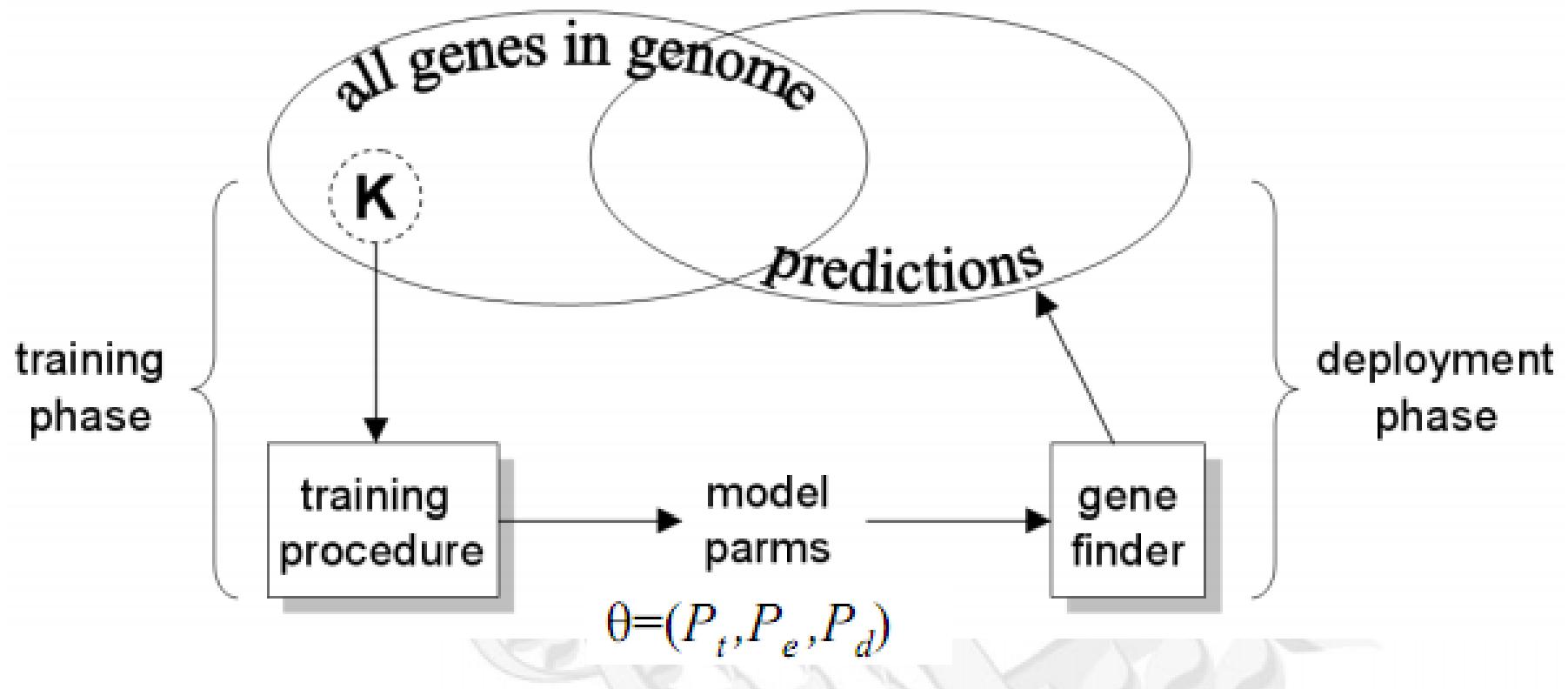
CGATATTCTGATTCTACGCGCGTATACTAGCTTATCTGATC  
011111112222222111111222211111112222111110

		to state		
		0	1	2
from state	0	0 (0%)	1 (100%)	0 (0%)
	1	1 (4%)	21 (84%)	3 (12%)
	2	0 (0%)	3 (20%)	12 (80%)

		symbol			
		A	C	G	T
in state	1	6 (24%)	7 (28%)	5 (20%)	7 (28%)
	2	3 (20%)	3 (20%)	2 (13%)	7 (47%)

- Los parámetros se “aprenden” de genes conocidos

# Modelos estadísticos (HMM)



# Predicción basada en homología

- Las secuencias codificantes son mas conservadas que las no codificantes debido a la presión selectiva
- Alineamiento de secuencias conocidas de genomas relacionados.
  - Podemos traducir nuestro genoma en los 6 marcos de lectura y alinear los péptidos obtenidos a proteínas conocidas (Blastx)
  - Nos puede dar idea de la función de la secuencia detectada

# Predictión basada en homología

**VISTA** Tools for Comparative Genomics

About Us Cite Us Contact Us

VISTA Home Custom Alignment Browser Enhancer DB Downloads Publications Training Help

VISTA is a comprehensive suite of programs and databases for comparative analysis of genomic sequences. There are two ways of using VISTA - you can submit your own sequences and alignments for analysis (VISTA servers) or examine pre-computed whole-genome alignments of different species.

**Submit Your Sequences**

mVISTA



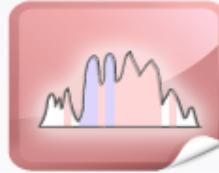
**mVISTA**  
Align and compare your sequences from multiple species

**rVISTA**  
Locate regulatory sequences in your data using comparative sequence analysis and transcription factor binding site search.

**GenomeVISTA**  
Compare your sequences against whole-genome assemblies.

**Precomputed Alignments**

VISTA Browser



**VISTA-Point**  
Access complete data and visual presentation of pairwise and multiple alignments of whole genome assemblies.

**VISTA Browser**  
Examine pre-computed pairwise and multiple alignments of whole genome assemblies.

**Whole Genome rVISTA**  
Identify transcription factor binding sites that are

New tool from VISTA family!

 **VISTA Region Viewer (RVViewer)** is an interactive on-line tool for comparing and prioritizing genomic intervals.

**Updates**

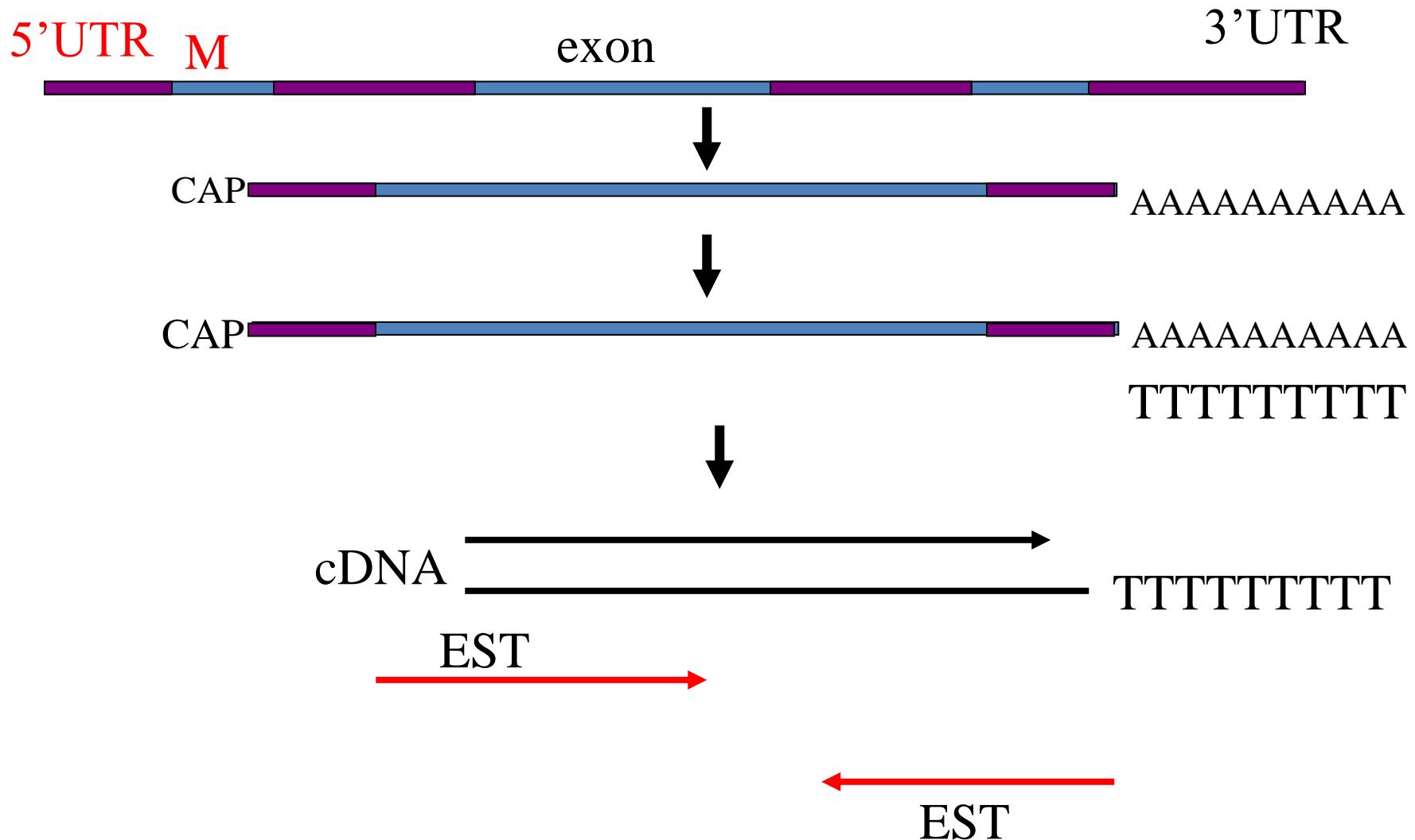
**September 2011**  
Updated the Cassava, Foxtail millet, Wine grape, and Maize assemblies, and added new plants: Columbine, Sweet orange, Eucalyptus, and Peach.

106 New whole-genome plant alignments are added to VISTA Browser.

Clade	Genome	Release	Position
Vertebrate	Human	Mar. 2006	chr11:5,200,001-5,250,000

VISTA-Point

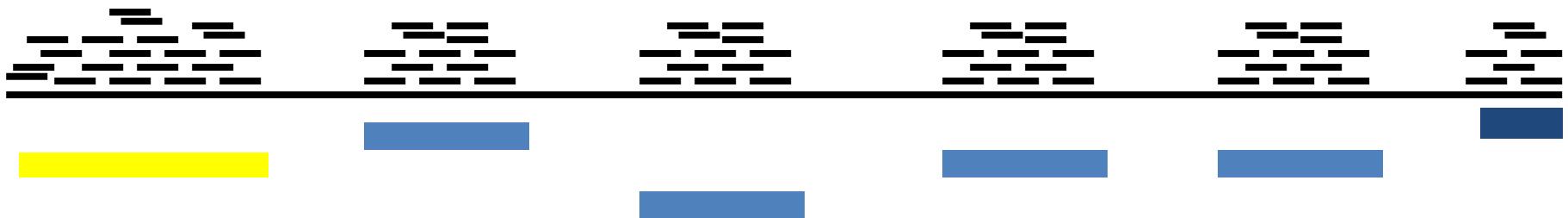
# Evidencia experimental ESTs



# Evidencia experimental RNAseq

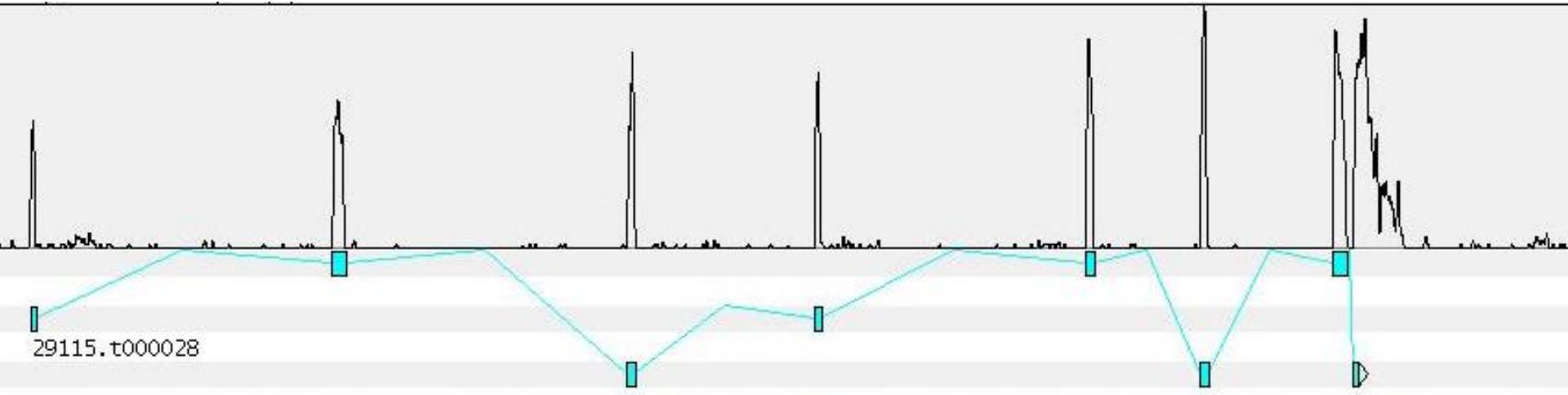


**Mapeo a la secuencia genómica**

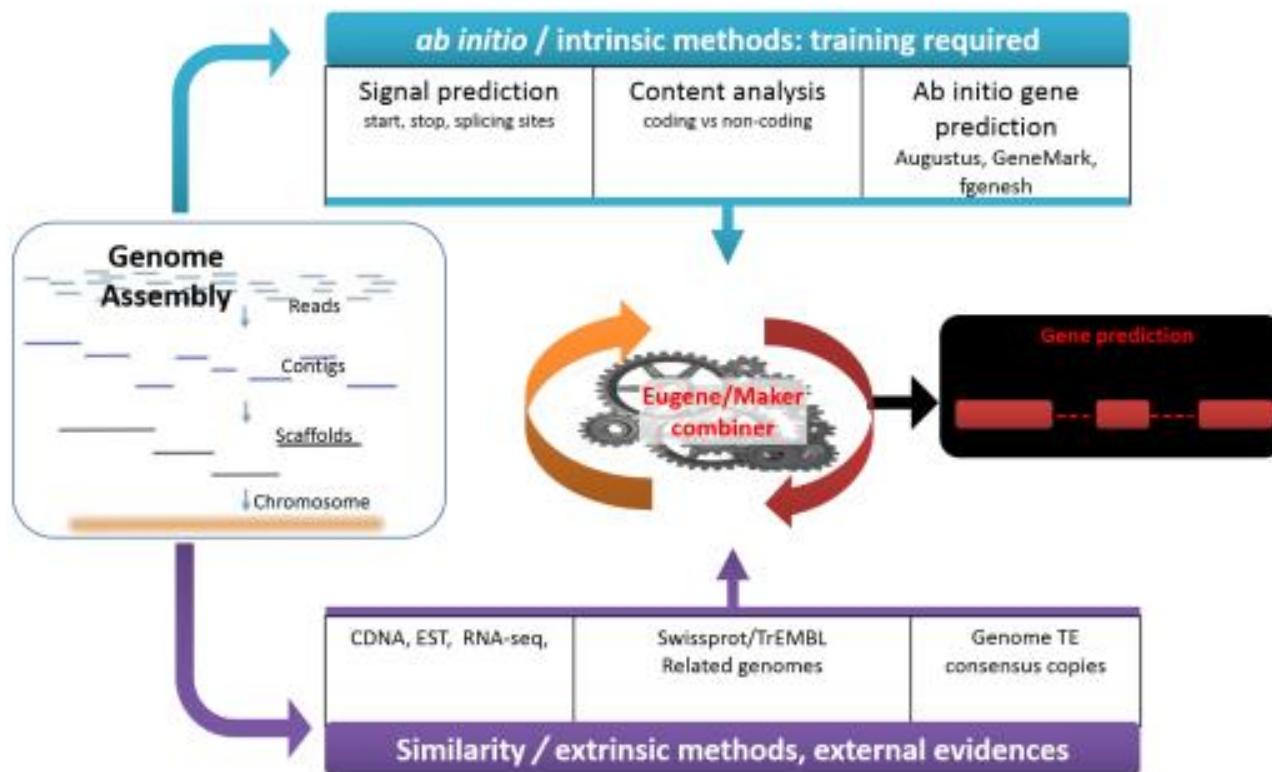


Genomic DNA sequence - assembly

# Evidencia experimental RNAseq



# Anotación: Combinación de evidencia y curado



# Anotación: Combinación de evidencia y curado

Gene prediction  
(SNAP)



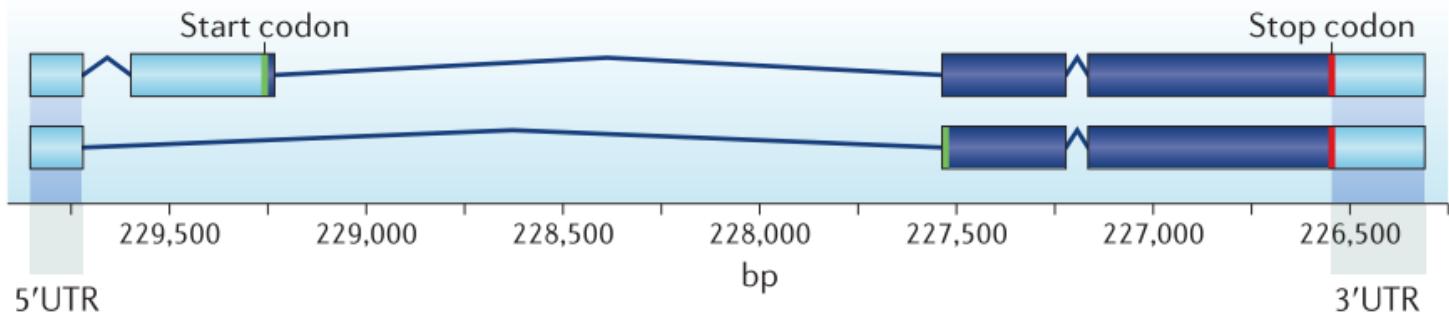
mRNA or EST evidence  
(Exonerate)



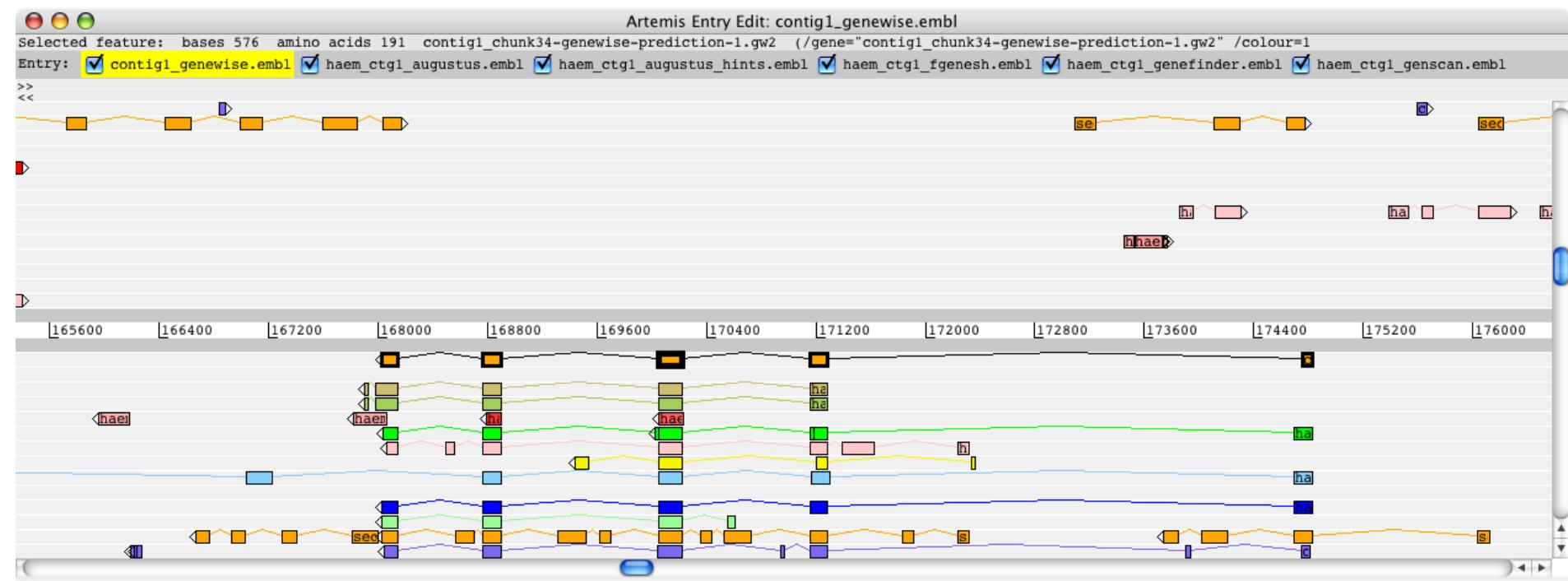
Protein evidence  
(BLASTX)



Gene annotation resulting  
from synthesizing all  
available evidence  
(two alternative splice forms)



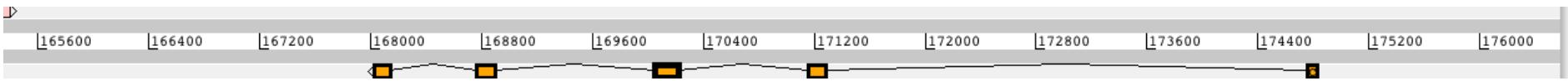
# Anotación: Combinación de evidencia y curado



pale brown  
brown-green  
pink/red blocks  
bright green  
pale pink  
yellow  
pale blue  
red  
dark blue  
jade green  
orange  
purple

hit to *H. contortus* EST cluster in Nembase found using PASA  
hit to *H. contortus* individual ESTs in NCBI database found using PASA  
hits to Uniprot  
twinscan prediction (**homology based**)  
snap prediction (ab initio)  
hmmgene prediction (ab initio)  
genscan prediction (ab initio)  
genefinder (ab initio)  
fgenesh prediction (ab initio)  
augustus hints prediction (**homology based**)  
augustus prediction (ab initio)  
genewise prediction (**homology based**)

# Que sigue???

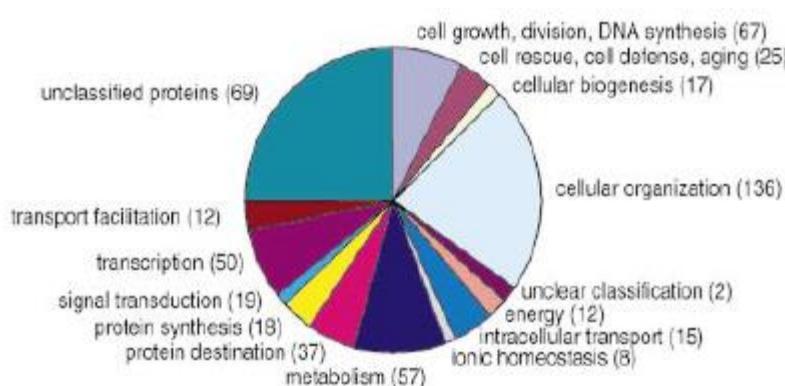


Para que sirve?

**Anotación funcional**

# Anotación funcional

- Búsqueda de similitud de secuencia
  - Inferir función por posible homología
    - A nivel global
      - Blast
    - Dominios y motivos
      - InterPro
        - » PROSITE, Pfam (HMM), PRINTS, ProDom, SMART, TIGRFAMs, PIRSF, SUPERFAMILY, GENE3D y PANTHER
- Señales particulares
  - Localización subcelular (SignalP)
  - Dominios transmembrana
- Asignación de categorías en “gene ontologies” (GOs)



# Anotación funcional

- Búsqueda de similitud de estructura
  - La estructura 3D de las proteínas es mas conservada que su secuencia
  - Se buscan las estructuras mas compatibles con la secuencia que se esta analizando
    - “*Protein threading*”

The screenshot shows the RaptorX web interface. At the top, there is a logo consisting of a stylized protein structure made of colored spheres and the text "RaptorX". To the right of the logo are links for "New Job", "Documentation", "My Jobs", "About", and "Xu Group". Below the header, it says "Jobs for user morten.kallberg@...". A navigation bar at the top indicates "Page 1 of 6. previous | next".

**Cfr Again**

Sequence	Status	Predictions
<a href="#">cfr</a>	Completed	Secondary Structure, Structure, Disorpred,

**Cfr**

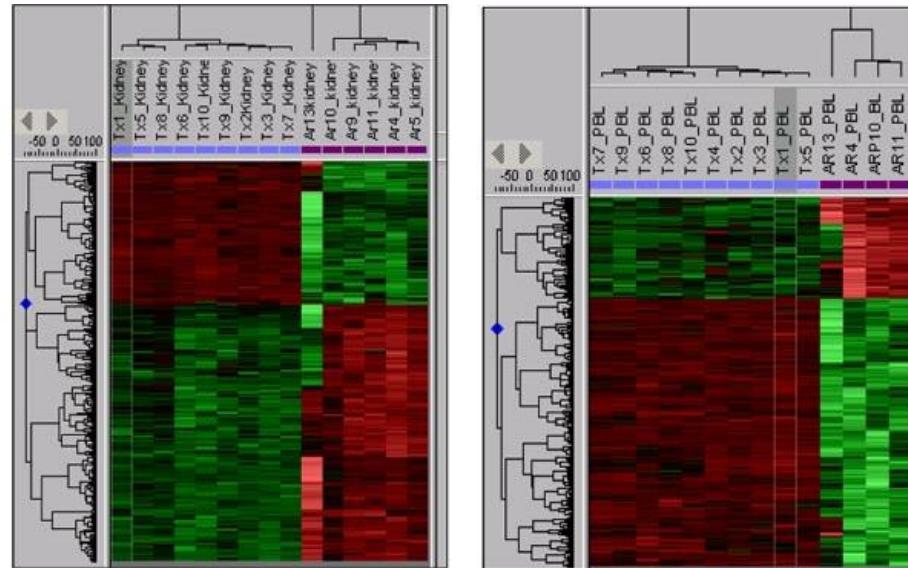
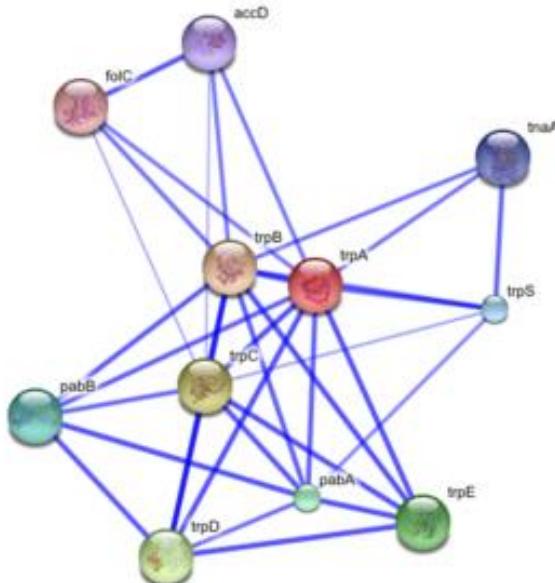
Sequence	Status	Predictions
<a href="#">sp P13569 CFTR_HUMAN</a>	Completed	Secondary Structure, Structure, Disorpred,

# Anotación funcional

- Se asignan nombres de acuerdo a la evidencia:
  - “Putative”: alta similitud a genes conocidos o dominios funcionales
  - “Expressed”: con evidencia de EST (o RNA-seq) sin evidencia de función
  - “Hypothetical Protein Conserved”: Predicha por software de anotación y con secuencia similar en bases de datos
  - “Hypothetical Protein”: Solo predicha por software de anotación

# Anotación funcional

- *Guilt by association:*
  - Se anota de acuerdo a asociaciones entre genes
    - Interacciones proteína - proteína
    - Perfiles de expresión
  - Es una anotación general



# Anotación funcional

- *Guilt by association:*



# Anotación funcional

