



UNIVERSIDAD AUTÓNOMA METROPOLITANA
Casa abierta al tiempo

UAM
Casa abierta al tiempo

Secuenciación masiva y anotación de genes.

Dr. Miguel Ángel del Río Portilla
Cátedra Divisional
"Alejandro Villalobos"

Secuenciación masiva y anotación de genes

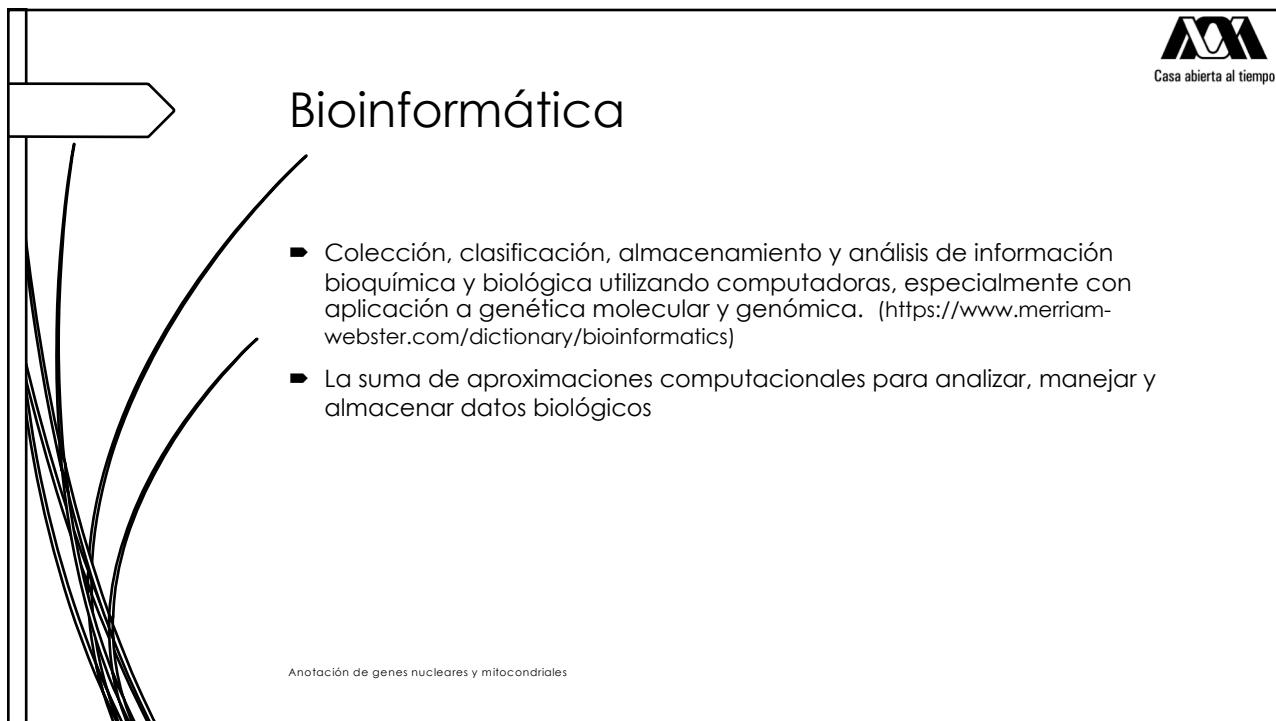
► ORGANIZADORES:

- Dra. Irene de los Ángeles Barriga Sosa
- Dr. Miguel Angel del Río Portilla

► PROFESORES

- Dr. Miguel Angel del Río Portilla.
mdelrio@cicese.mx,
mdelrio@xanum.uam.mx
- M. en C. Eduardo Zúñiga León
eduardoz@xanum.uam.mx
- Dra. Dra. Erika Magallón Gayón

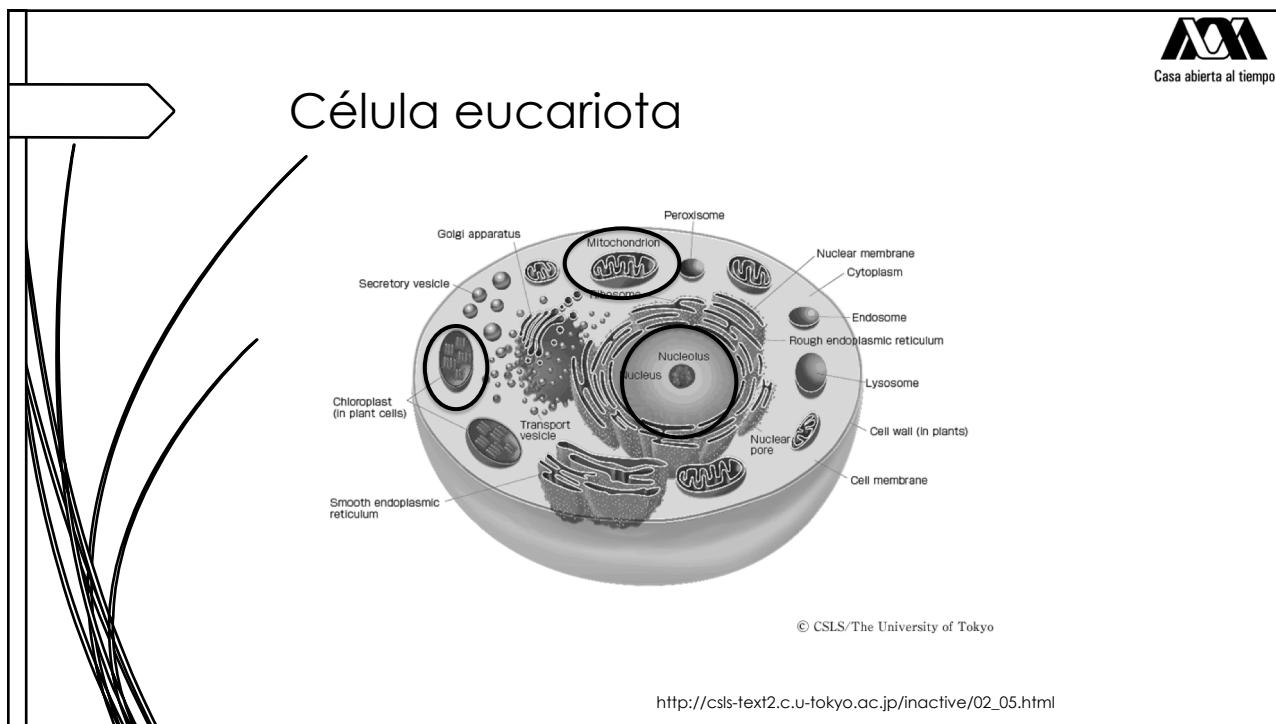
Anotación de genes nucleares y mitocondriales



Bioinformática

- Colección, clasificación, almacenamiento y análisis de información bioquímica y biológica utilizando computadoras, especialmente con aplicación a genética molecular y genómica. (<https://www.merriam-webster.com/dictionary/bioinformatics>)
- La suma de aproximaciones computacionales para analizar, manejar y almacenar datos biológicos

Anotación de genes nucleares y mitocondriales



Célula eucariota

© CSLS/The University of Tokyo

http://csls-text2.c.u-tokyo.ac.jp/inactive/02_05.html

Genoma

- Gen
 - Cromosoma
 - Genoma
- RNA (mRNA transcripto)
 - Transcriptoma
- Proteína
 - Proteoma
- Interacción de genes
 - Epigenética
- Tejido y condición específico

TRANSCRIPTOME
NHGRI FACT SHEETS
genome.gov

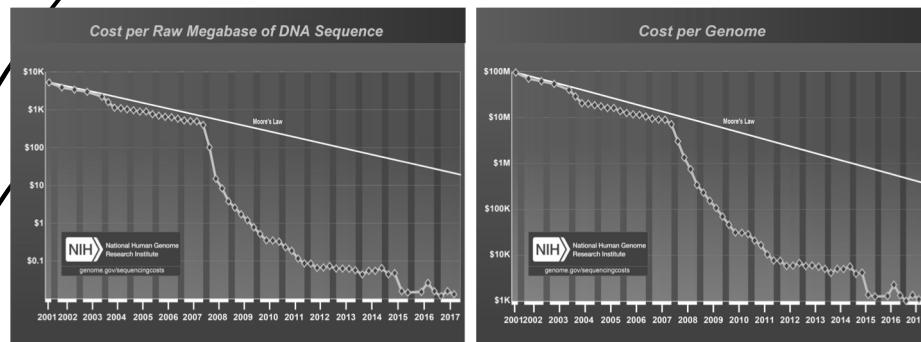
The diagram shows a cell with a nucleus and cytoplasm. In the nucleus, DNA is transcribed by RNA polymerase into an antisense strand and a sense strand, forming an RNA Transcript. This transcript moves to the cytoplasm where it is translated by ribosomes into proteins.

Herramientas Genómicas

- Proyecto del genoma Humano (1990-2003)
- Secuenciación Sanger
- Desarrollo de nuevas tecnologías
- Next generation Sequencing
- Secuencias de varios millones de fragmentos (lecturas) en una corrida
 - Microsatélites
 - SNP
 - Genomas Mitochondriales

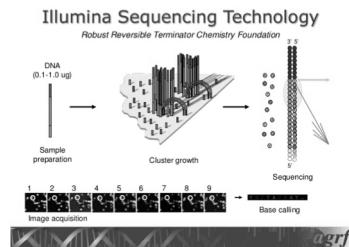
Indispensable para microsatélites y SNPs y análisis de RNA (RNAseq, miRNA...)

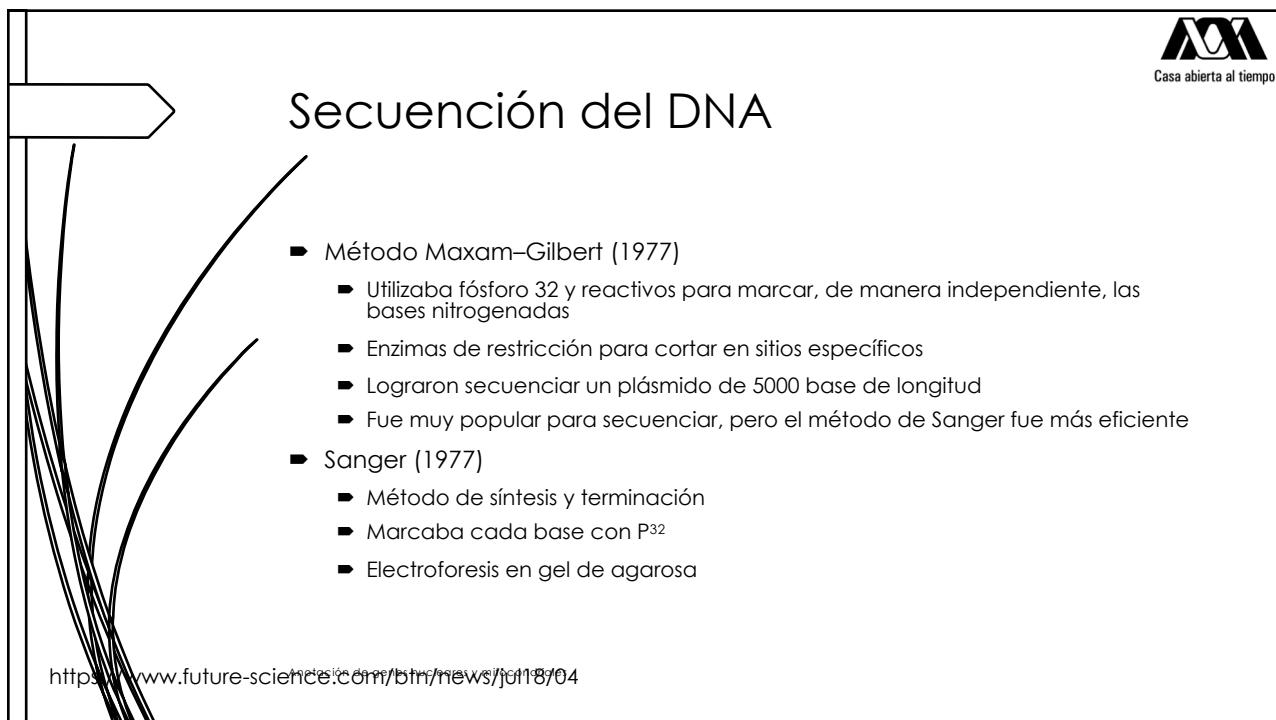
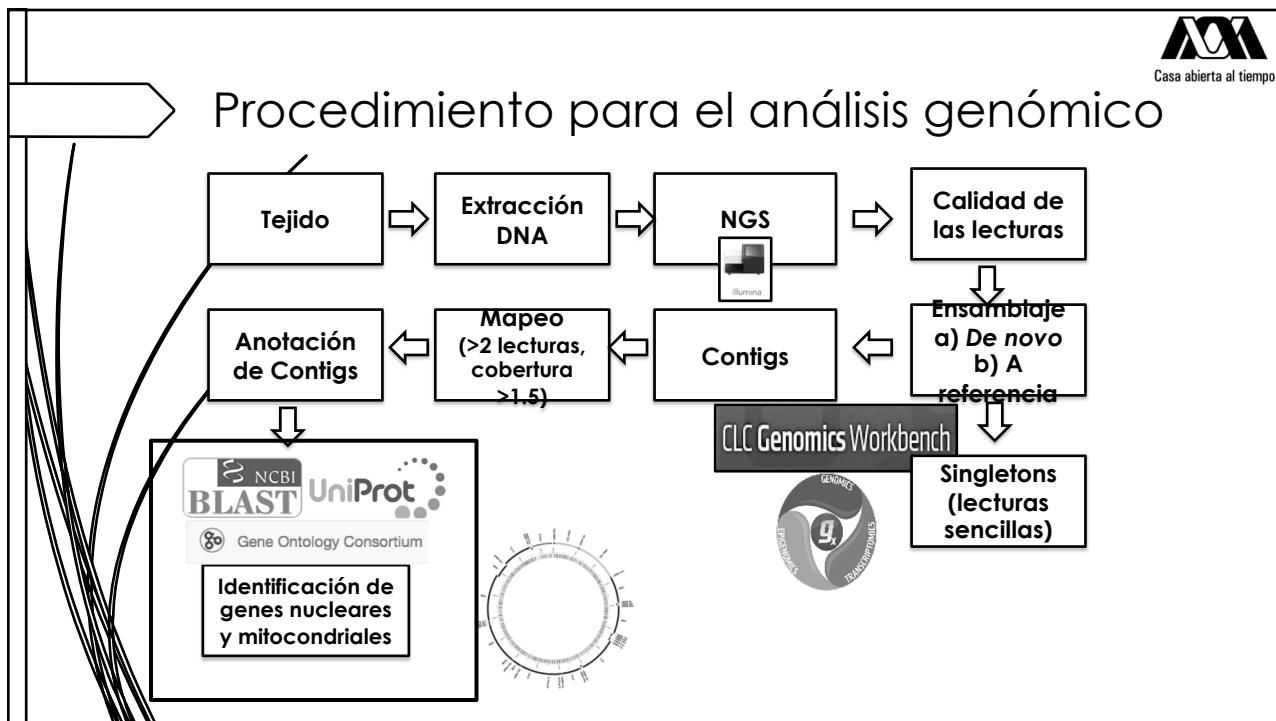
Costos de la secuenciación masiva



NGS análisis genómico y transcriptómico

- ▶ Obtención de genoma
 - ▶ Nuclear (identificación del número de genes)
 - ▶ Mitochondrial (estructura)
- ▶ Microsatélites
 - ▶ Genética de poblaciones
 - ▶ Identificación de progenitores
 - ▶ Localización de caracteres cuantitativos de interés (QTLs)
- ▶ SNP
 - ▶ Identificación de progenitores
 - ▶ Localización de caracteres cuantitativos de interés (QTLs)
- ▶ Secuenciación de mRNA (transcriptoma: tejido/condición)
 - ▶ Genes de interés
 - ▶ Respuesta a enfermedades





Secuenciación Sanger. Secuenciación por terminación

A Deoxyribonucleotide

A Di-Deoxyribonucleotide

Anotación de genes nucleares y mitocondriales

Secuenciación

Deoxyribose

3'-OH in normal DNA allows elongation.

primer C A T G G A C T A A T G G A dd
G T A C C T G A T T A C C T A

primer C A T G G A C T A A dd
G T A C C T G A T T A C C T A

primer C A T G G A C T A dd
G T A C C T G A T T A C C T A

primer C A T G G A dd
G T A C C T G A T T A C C T A

Length of fragment
30
29
28
27
26
25
24
23
22
21

Dideoxy n carries the chain of one of the bases along the template. The bases are indicated by the dideoxynucleotides, and the positions of the nucleotides are determined directly from the gel. In this example, the length of the primer needed

Hartl, D. L., Jones, E. W. (1998). *Genetics: Principles and analysis*. Jones and Barlett Publishers.
Anotación de genes nucleares y mitocondriales

Secuenciación

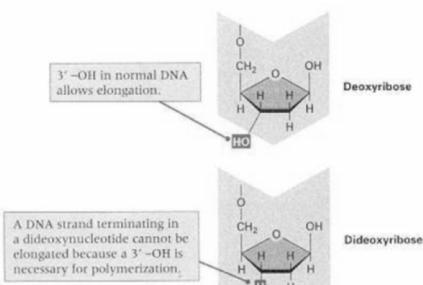
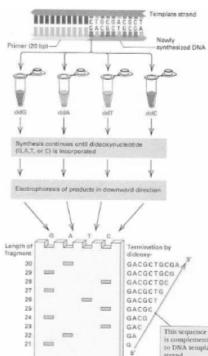


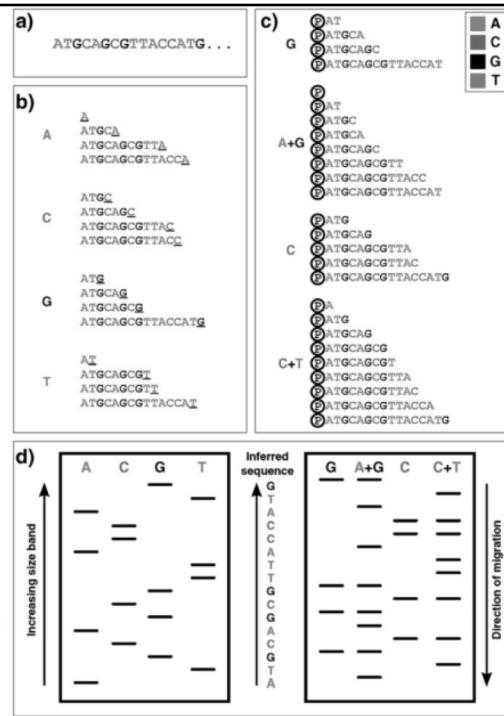
Figure 5.38

Structures of normal deoxyribose and the dideoxyribose sugar used in DNA sequencing. The dideoxyribose has a hydrogen atom (red) attached to the 3' carbon, in contrast with the hydroxyl group (red) at this position in deoxyribose. Because the 3' hydroxyl group is essential for the attachment of the next nucleotide in line in a growing DNA strand, the incorporation of a dideoxynucleotide immediately terminates synthesis.

Hartl, D. L., & Jones, E. W. (1998). *Genetics: Principles and analysis*. Jones and Barlett Publishers.
Anotación de genes nucleares y mitocondriales



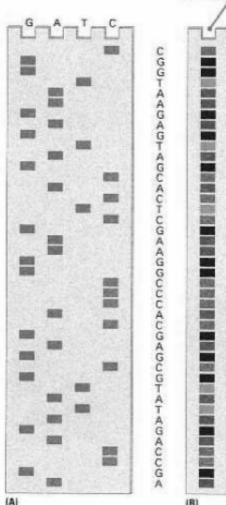
Dideoxy method of DNA sequencing. DN terminus reactions are carried out in the presence of d. normal nucleotides plus a small amount of one of the dideoxynucleotides containing G, A, T, C. Synthesis continues along the template strand until a dideoxynucleotide is incorporated. The products that result from termination at each dideoxynucleotide are indicated on the right. The fragments are separated by size by electrophoresis and visualized by staining with ethidium bromide and viewed directly from the gel. In this example, the length of the primer needed



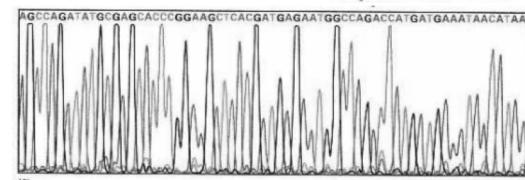
Electroforesis capilar



Mixture of all four sequencing reactions, each labeled with a different fluorescent dye

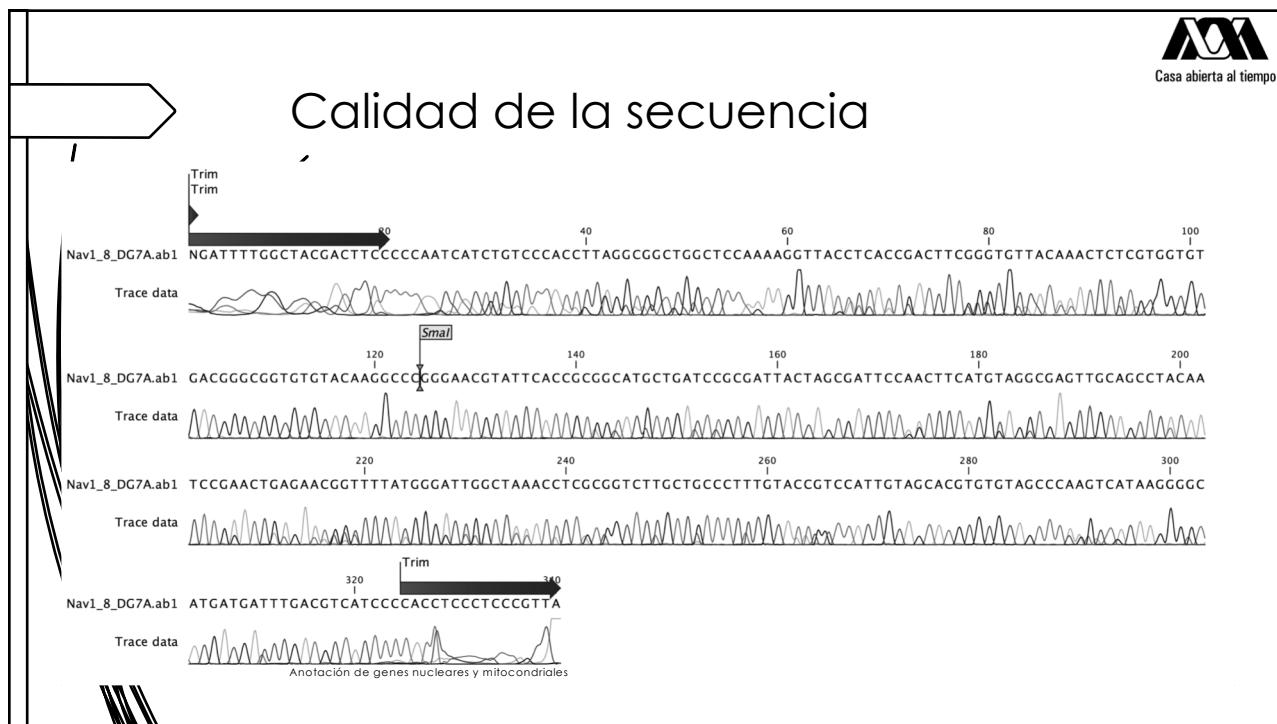


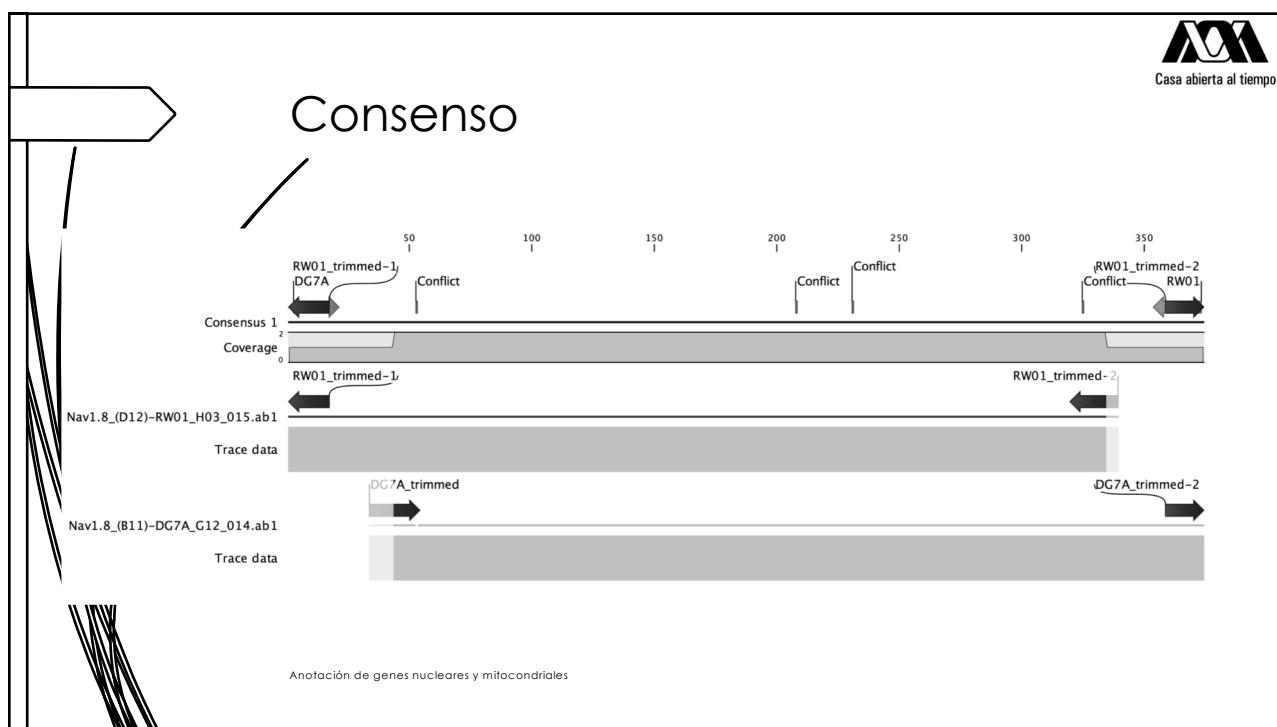
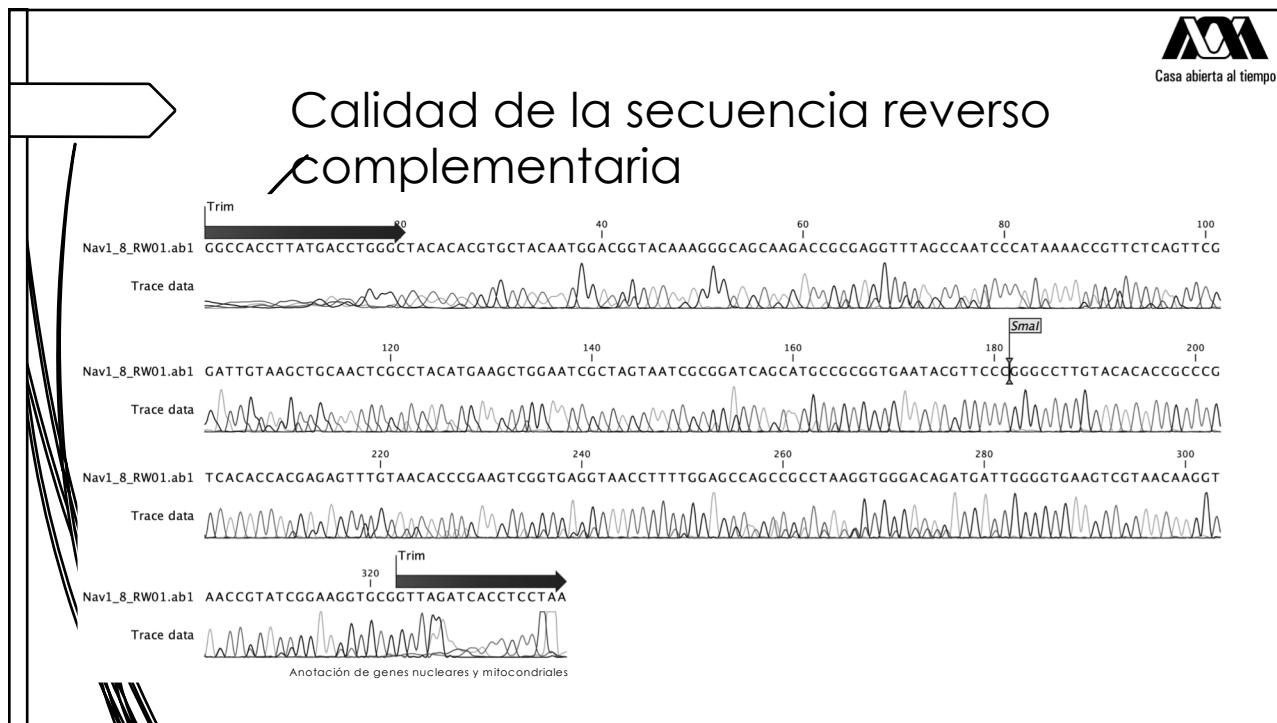
This graph is the output from the light detector in an automated DNA sequencer. It detects the fluorescence of each DNA band as it comes off the bottom of the gel. The sequence reads from left (5' end) to right (3' end). Each dideoxy nucleotide has a different fluorescent "tag", so the color of each peak identifies the dideoxynucleotide terminator incorporated at that site.

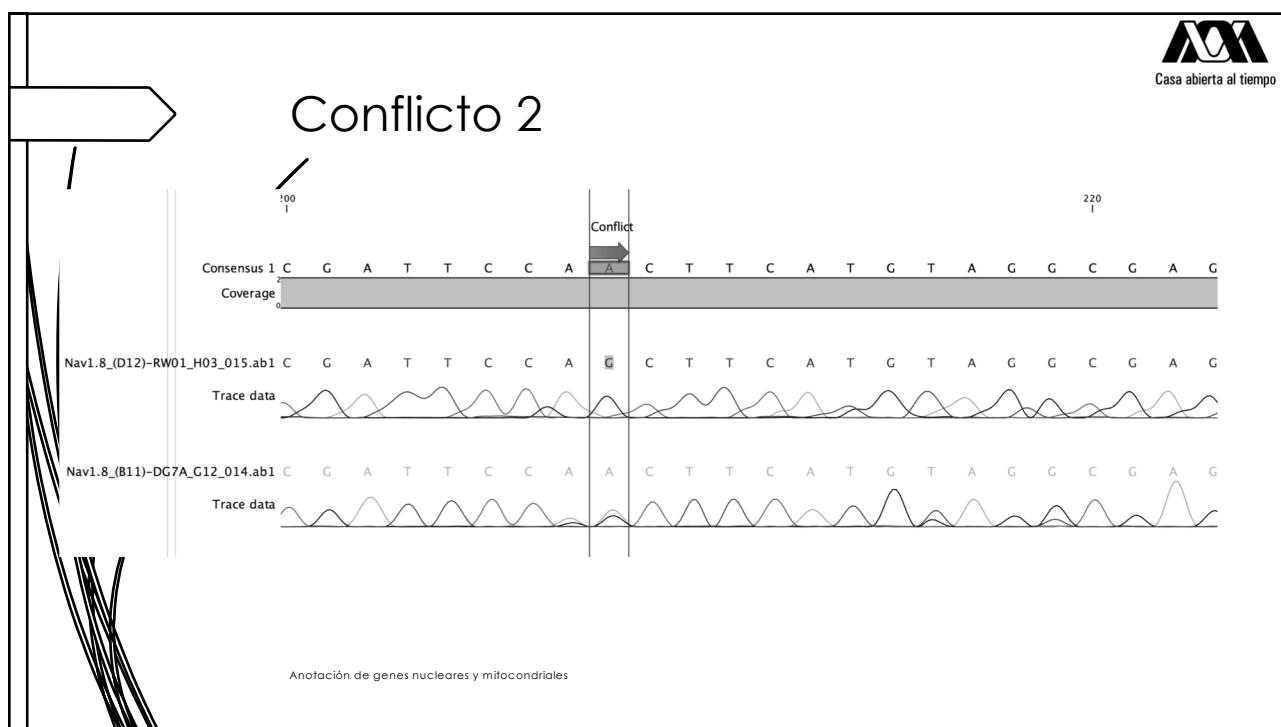
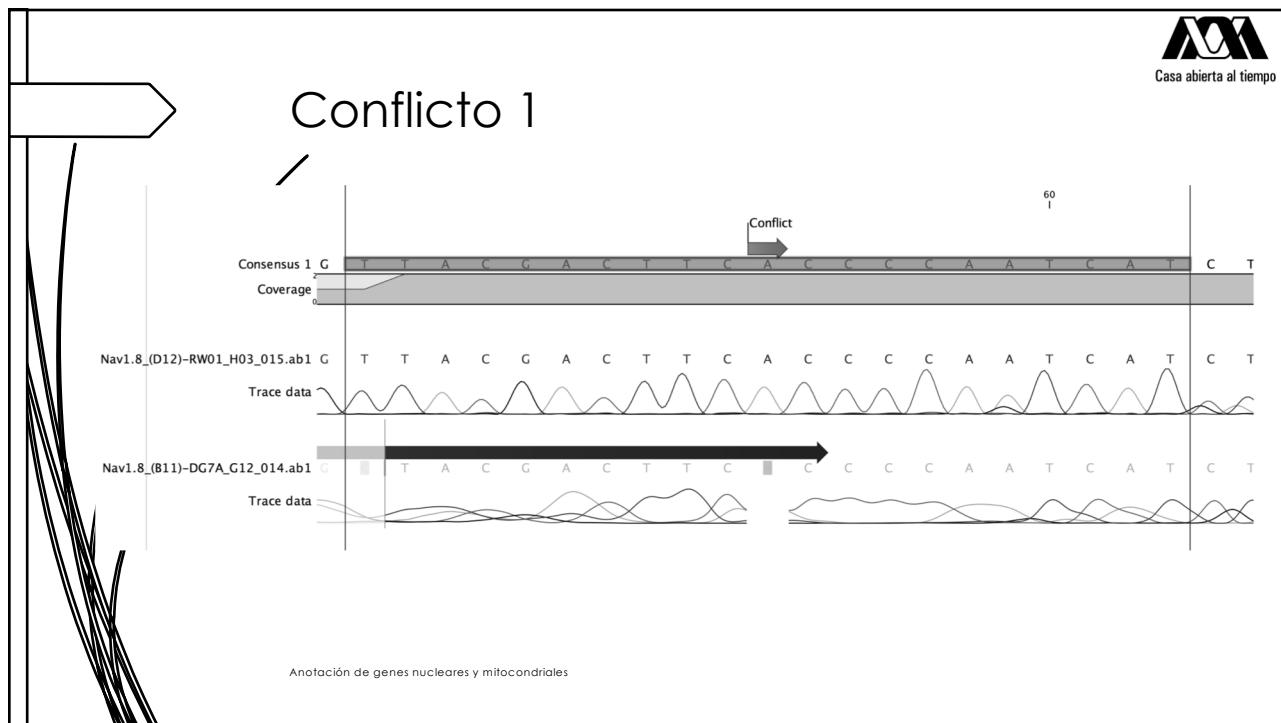


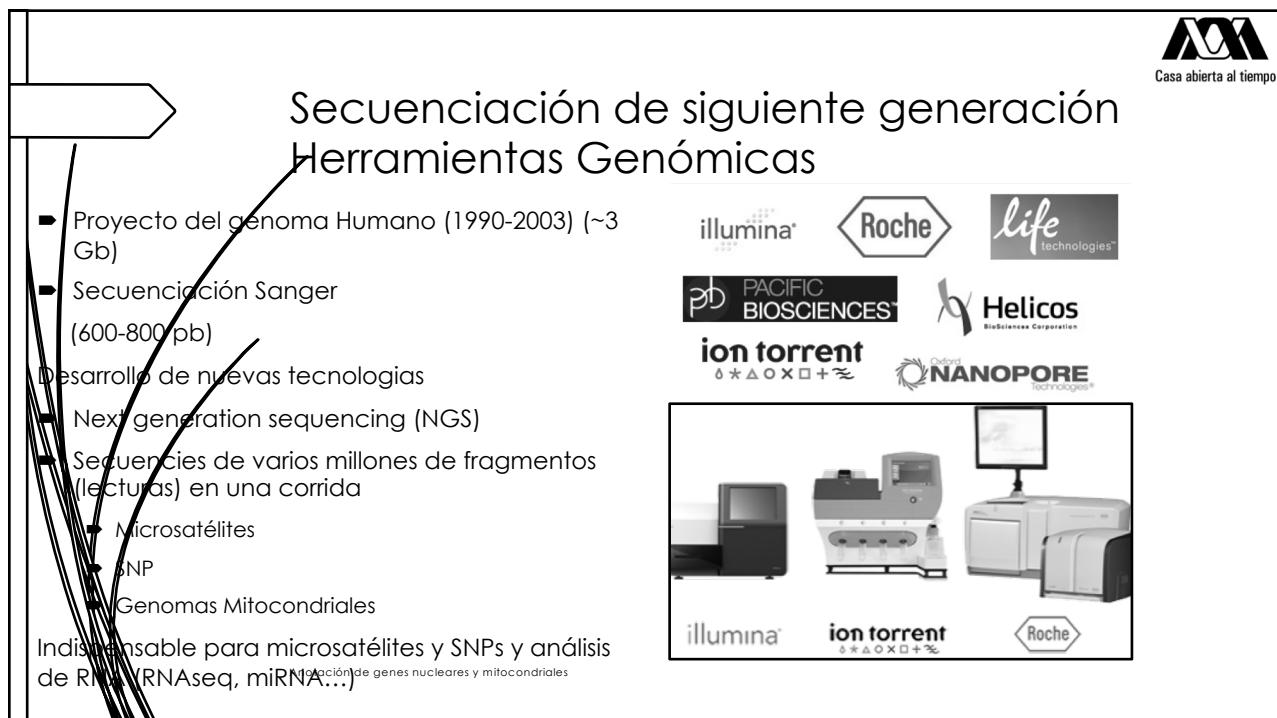
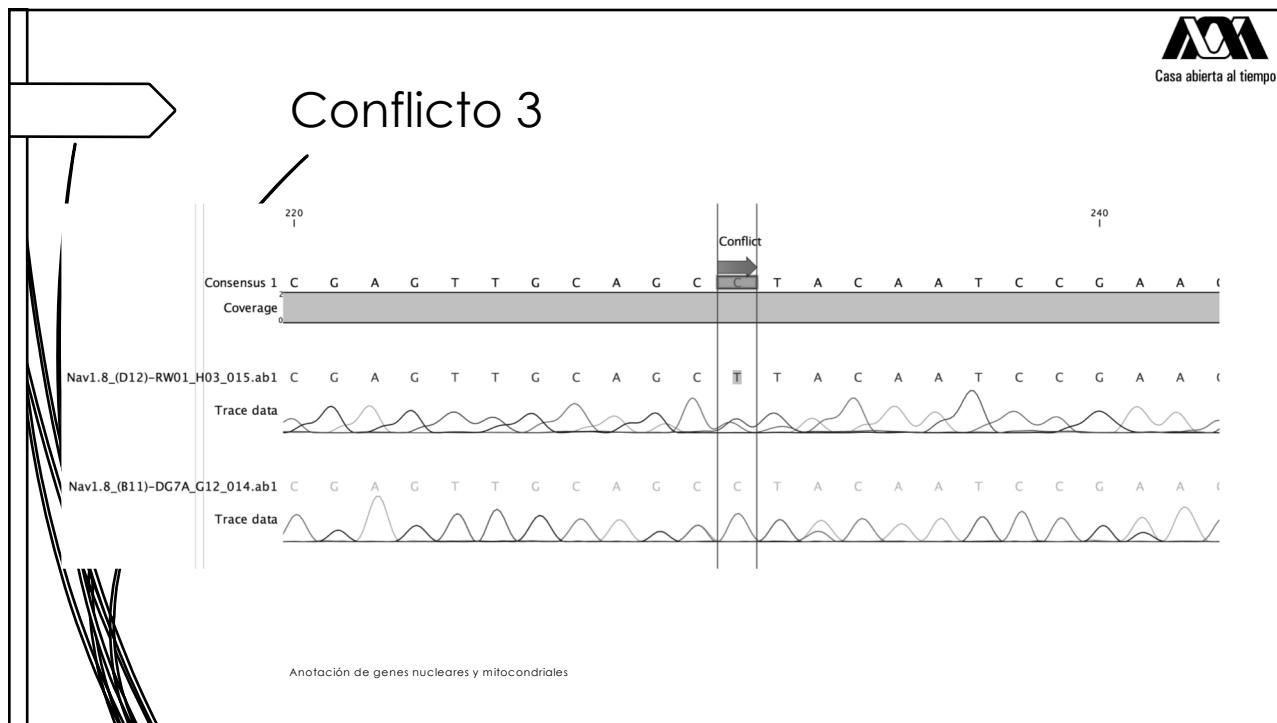
Automated DNA sequencing. (A) Conventional sequencing gel obtained from the dideoxy procedure (Section 5.9). The DNA sequence can be determined directly from photographic film according to the positions of the bands. (B) Banding pattern obtained when each of the terminating nucleotides is labeled with a different fluorescent dye and the bands are separated in the same lane of the gel. (C) Trace of the fluorescence pattern obtained from the gel in part B by automated detection of the fluorescence of each band as it comes off the bottom of the gel during continued electrophoresis.

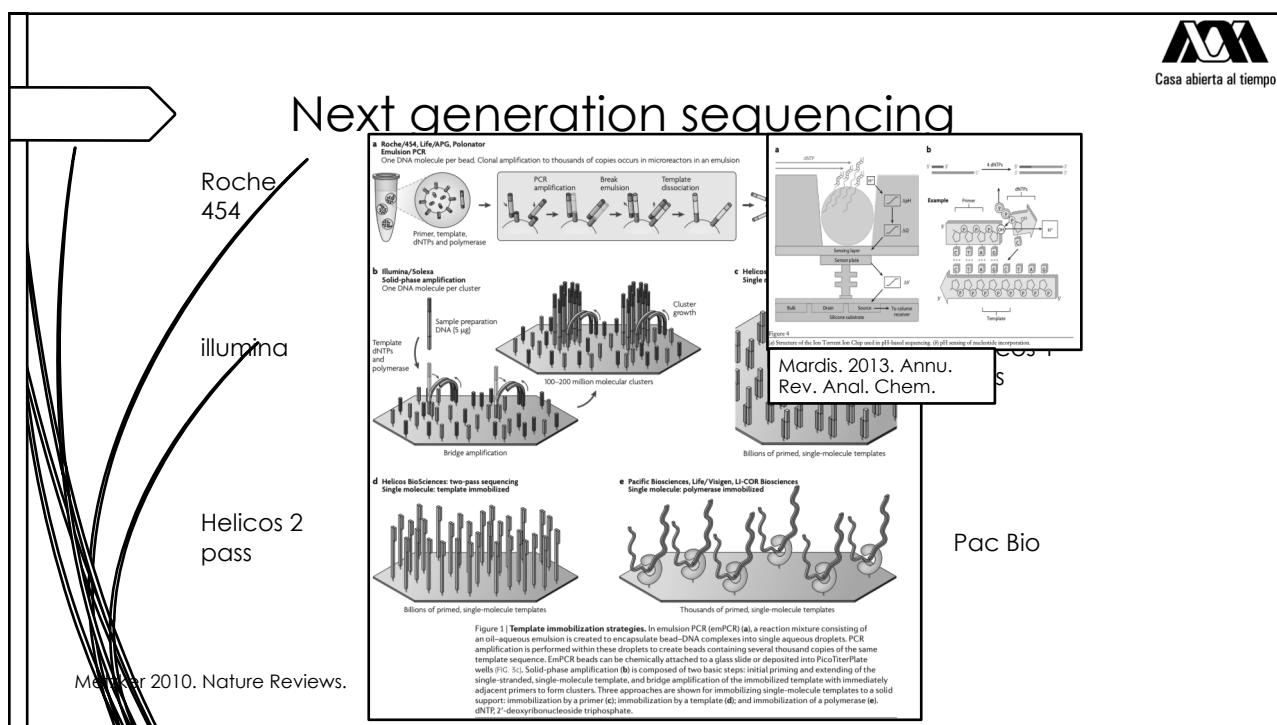
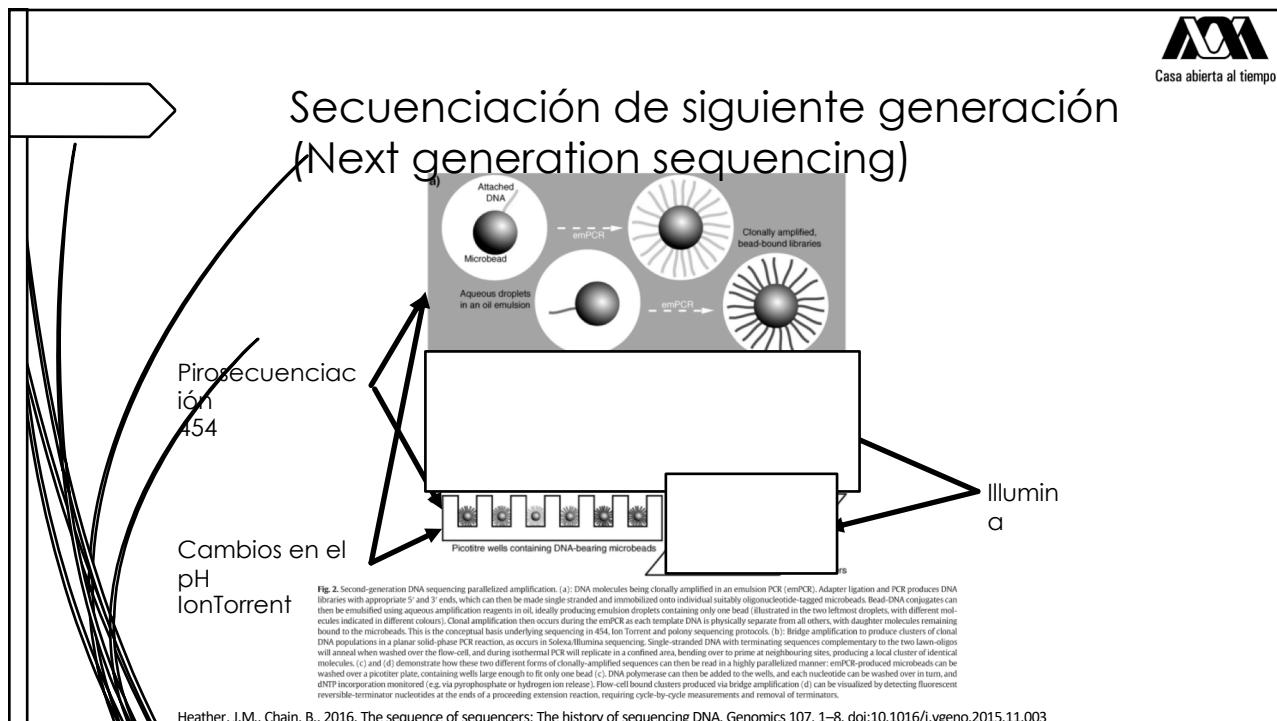
s nucleares y mitocondriales











Tercer generación de secuenciación

Fig. 3. Third-generation DNA sequencing nucleotide detection. (a): Nucleotide detection in a zero-mode waveguide (ZMW), as featured in PacBio sequencers. DNA polymerase molecules are attached to the bottom of each ZMW (*), and target DNA and fluorescent nucleotides are added. As the diameter is narrower than the excitation light's wavelength, illumination rapidly decays travelling up the ZMW: nucleotides being incorporated during polymerisation at the base of the ZMW provide real-time bursts of fluorescent signal, without undue interference from other labelled dNTPs in solution. (b): Nanopore DNA sequencing as employed in ONT's MinION sequencer. Double stranded DNA gets denatured by a processive enzyme (†) which ratchets one of the strands through a biological nanopore (‡) embedded in a synthetic membrane, across which a voltage is applied. As the ssDNA passes through the nanopore the different bases prevent ionic flow in a distinctive manner, allowing the sequence of the molecule to be inferred by monitoring the current at each channel.

Heather, J.M., Chain, B., 2016. The sequence of sequencers: The history of sequencing DNA. *Genomics* 107, 1–8. doi:10.1016/j.ygeno.2015.11.003

Comparación

Platform	Read length (bp)	Throughput	Reads	Runtime	Error profile	Instrument cost (US\$)	Cost per Gb (US\$, approx.)
<i>Sequencing by ligation</i>							
SOLiD 5500 Wildfire	50 (SE)	80 Gb	~700M*	6 d*	≤0.1%, AT bias [‡]	NA [§]	\$130 [†]
	75 (SE)	120 Gb					
	50 (SE)*	160 Gb*					
SOLiD 5500xl	50 (SE)	160 Gb	~1.4 B*	10 d*	≤0.1%, AT bias [‡]	\$251,000 [†]	\$70 [†]
	75 (SE)	240 Gb					
	50 (SE)*	320 Gb*					
BGISEQ-500 FCS ¹⁵⁵	50–100 (SE/PE)*	8–40 Gb*	NA	24 h*	≤0.1%, AT bias [‡]	\$250 (REF. 155)	NA
BGISEQ-500 FCL ¹⁵⁵	50–100 (SE/PE)*	40–200 Gb*	NA	24 h*	≤0.1%, AT bias [‡]	\$250,000 (REF. 155)	NA

Anotación de genes nucleares y mitocondriales

Comparación

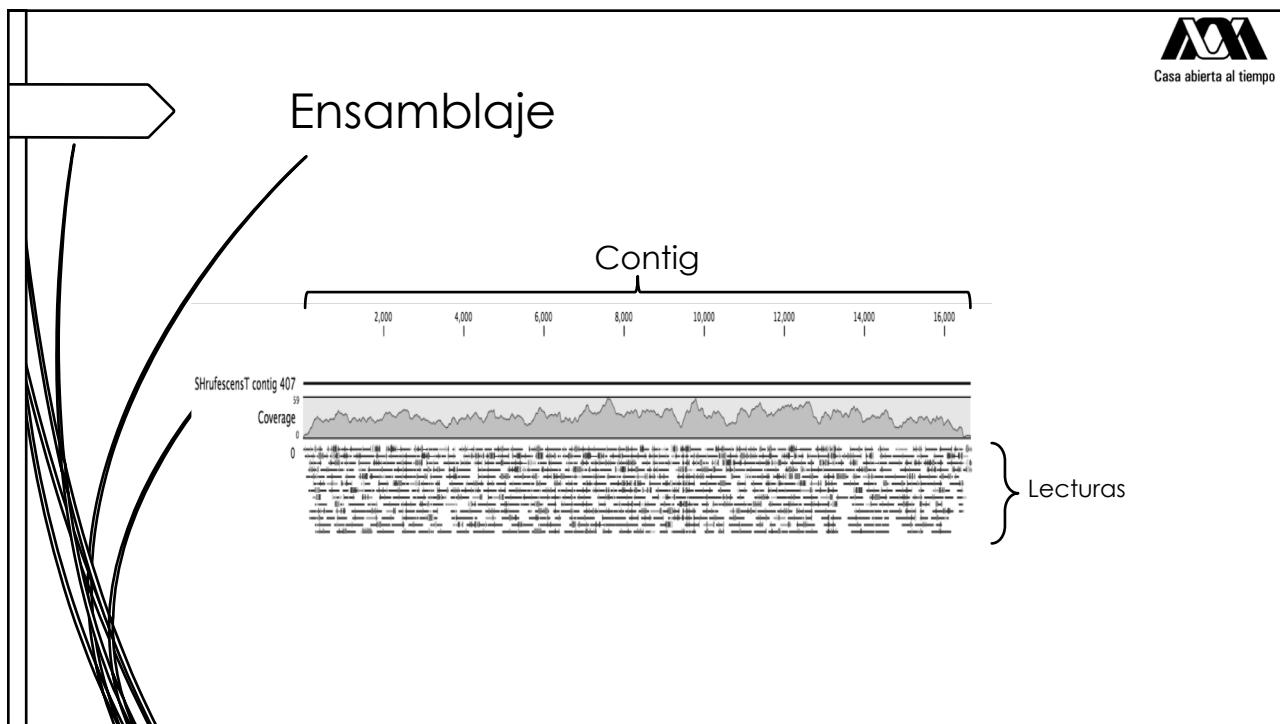
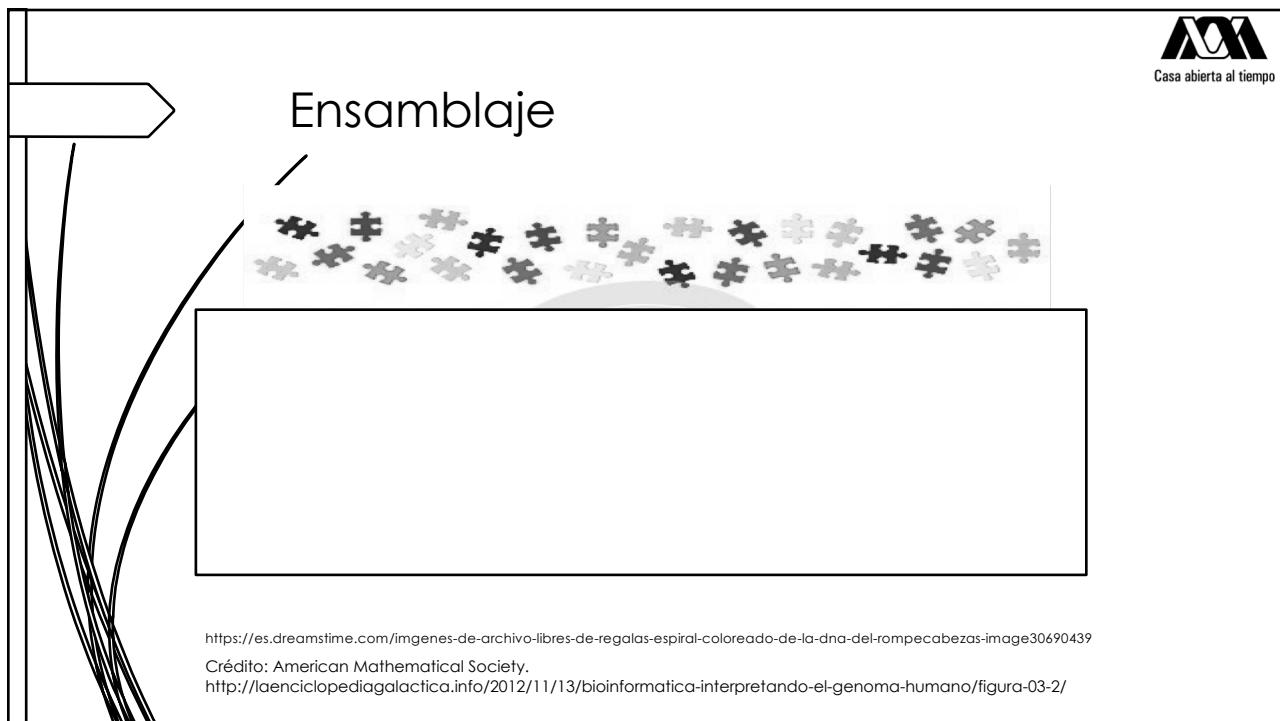
Table 1 | Summary of NGS platforms

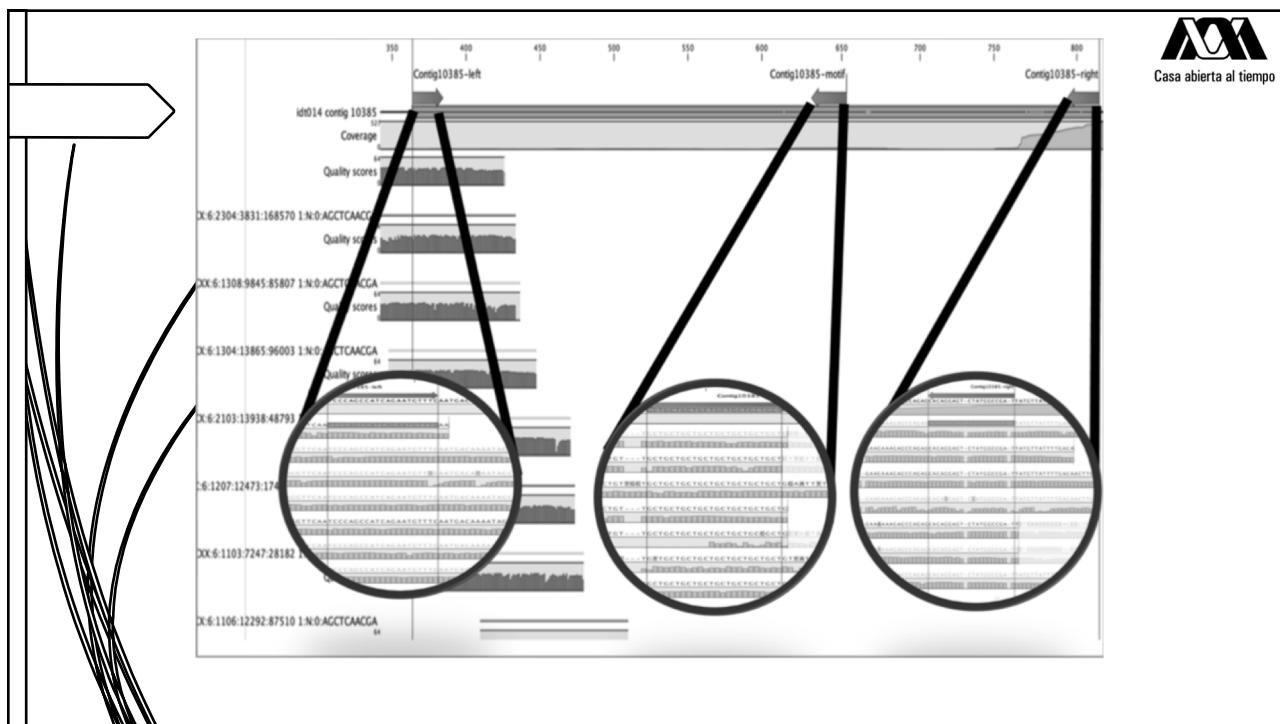
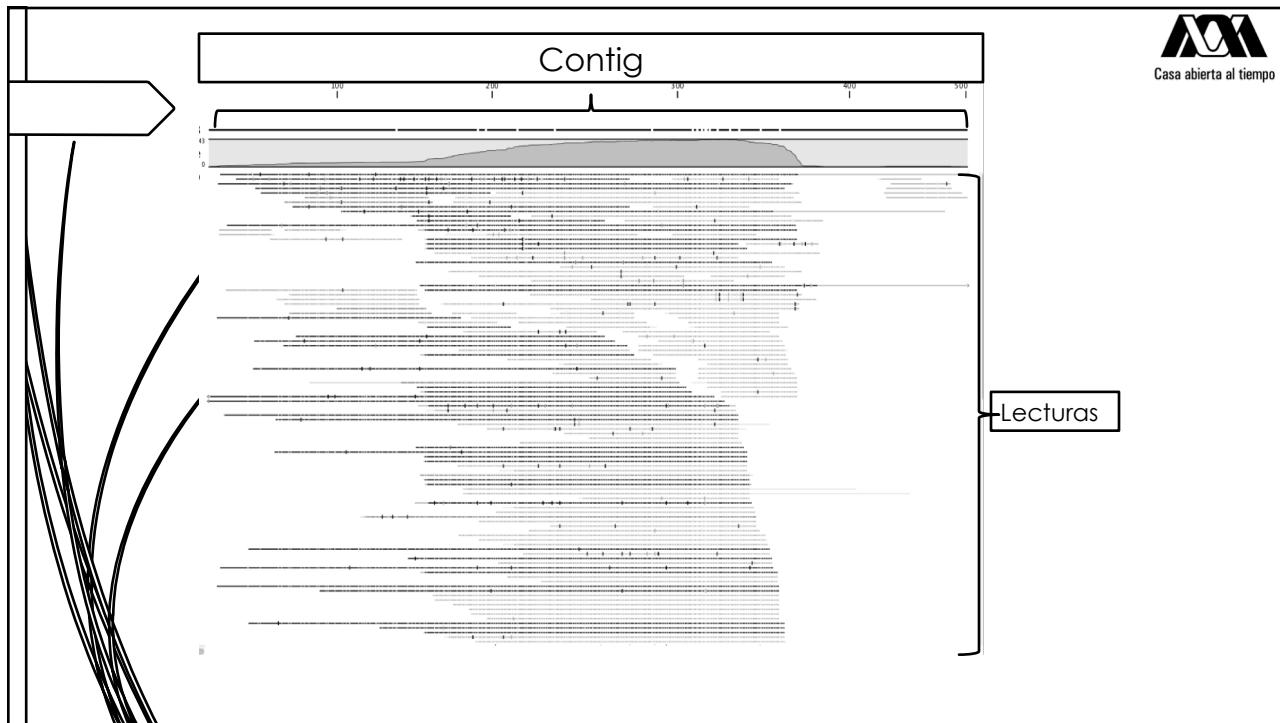
Platform	Read length (bp)	Throughput	Reads	Runtime	Error profile	Instrument cost (US\$)	Cost per Gb (US\$, approx.)
Sequencing by synthesis: CRT							
Illumina MiniSeq Mid output	150 (SE)*	2.1–2.4 Gb*	14–16M*	17 h*	<1%, substitution [†]	\$50,000 (REF. 118)	\$200–300 (REF. 118)
Illumina MiniSeq High output	75 (SE)	1.6–1.8 Gb	22–25M (SE)*	7 h	<1%, substitution [†]	\$50,000 (REF. 118)	\$200–300 (REF. 118)
	75 (PE)	3.3–3.7 Gb	44–50M (PE)*	13 h			
	150 (PE)*	6.6–7.5 Gb*		24 h*			
Illumina MiSeq v2	36 (SE)	540–610 Mb	12–15 M (SE)	4 h	0.1%, substitution [†]	\$99,000 [‡]	~\$1,000
	25 (PE)	750–850 Mb	24–30M (PE)*	5.5 h			\$996
	150 (PE)	4.5–5.1 Gb		24 h			\$212
	250 (PE)*	7.5–8.5 Gb*		39 h*			\$142 [‡]
Illumina MiSeq v3	75 (PE)	3.3–3.8 Gb	44–50M (PE)*	21–56 h*	0.1%, substitution [†]	\$99,000 [‡]	\$250
	300 (PE)*	13.2–15 Gb*					\$110 [‡]
Illumina NextSeq 500/550 Mid output	75 (PE)	16–20 Gb	Up to 260 M (PE)*	15 h	<1%, substitution [†]	\$250 [‡]	\$42
	150 (PE)*	32–40 Gb*		26 h*			\$40 [‡]

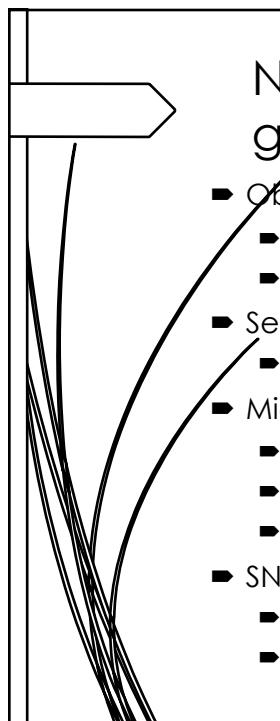
Comparación

Table 1 | Summary of NGS platforms

Platform	Read length (bp)	Throughput	Reads	Runtime	Error profile	Instrument cost (US\$)	Cost per Gb (US\$, approx.)
Sequencing by synthesis: CRT							
Illumina HiSeq2500 v2 Rapid run	36 (SE)	9–11Gb	300 M (SE)*	7 h	0.1%, substitution [‡]	\$690 [‡]	\$230
Illumina HiSeq2500 v3	50 (PE)	25–30 Gb	600 M (PE)*	16 h			\$90
	100 (PE)	50–60 Gb		27 h			\$52
	150 (PE)	75–90 Gb		40 h			\$45
	250 (PE)*	125–150 Gb*		60 h*			\$40 [‡]
	36 (SE)	47–52 Gb	1.5 B (SE)	2 d	0.1%, substitution [‡]	\$690 [‡]	\$180
Illumina HiSeq2500 v4	50 (PE)	135–150 Gb	3 B (PE)*	5.5 d			\$78
	100 (PE)*	270–300 Gb		11 d*			\$45 [‡]
	36 (SE)	64–72 Gb	2 B (SE)	29 h	0.1%, substitution [‡]	\$690 [‡]	\$150
Illumina HiSeq3000/4000	50 (PE)	180–200 Gb	4 B (PE)*	2.5 d			\$58
	50 (SE)	105–125 Gb	2.5 B (SE)*	1–3.5 d*	0.1%, substitution [‡]	\$740/\$900 (REF. 156)	\$50
	75 (PE)	325–375 Gb					\$31
	150 (PE)*	650–750 Gb*					\$22 (REF. 157)







NGS análisis genómico y transcriptómico

■ Obtención de genoma

- Mitocondrial
- Nuclear

■ Secuenciación de mRNA (transcriptoma: tejido/condición)

- Genes de interés

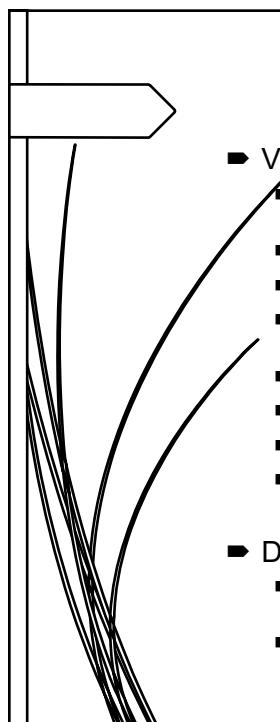
■ Microsatélites

- Genética de poblaciones
- Identificación de progenitores
- Localización de caracteres cuantitativos de interés (QTLs)

■ SNP

- Identificación de progenitores
- Localización de caracteres cuantitativos de interés (QTLs)

■ Casa abierta al tiempo



NGS

■ Ventajas

- Obtención de una gran cantidad de secuencias en una sola corrida a diferencia del método de Sanger
- Análisis de muchas secuencias a la vez
- Análisis cuantitativo del transcriptoma (expresión de genes)
- Búsqueda de marcadores moleculares más fácilmente (microsatélites, Polimorfismo de un solo nucleótido [SNP])
- Biomas (metagenómica)
- Costo va reduciéndose cada vez más.
- Uso en especies no modelo
- Uso en muestras ambientales

■ Desventajas

- Requiere programas específicos para el análisis de las secuencias (bioinformática)
- Capacidad de cómputo y almacenamiento

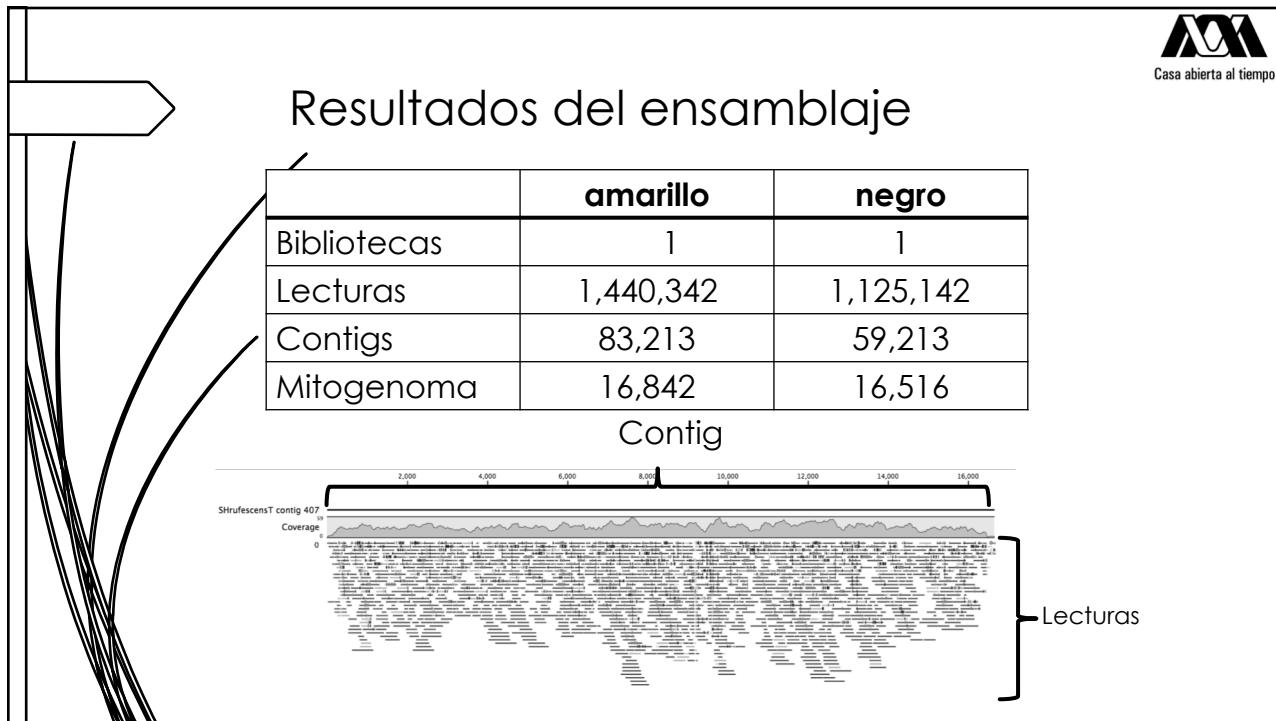
■ Casa abierta al tiempo



ACM
Casa abierta al tiempo



ACM
Casa abierta al tiempo



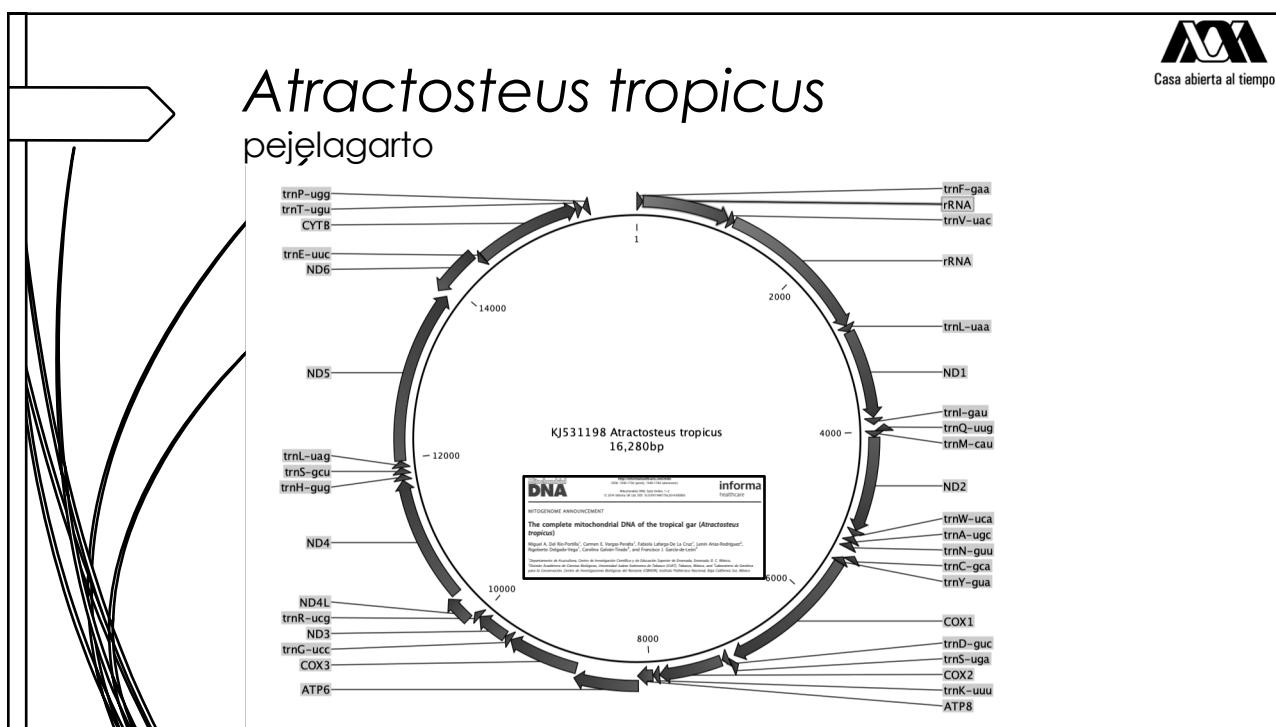
Genes mitocondriales

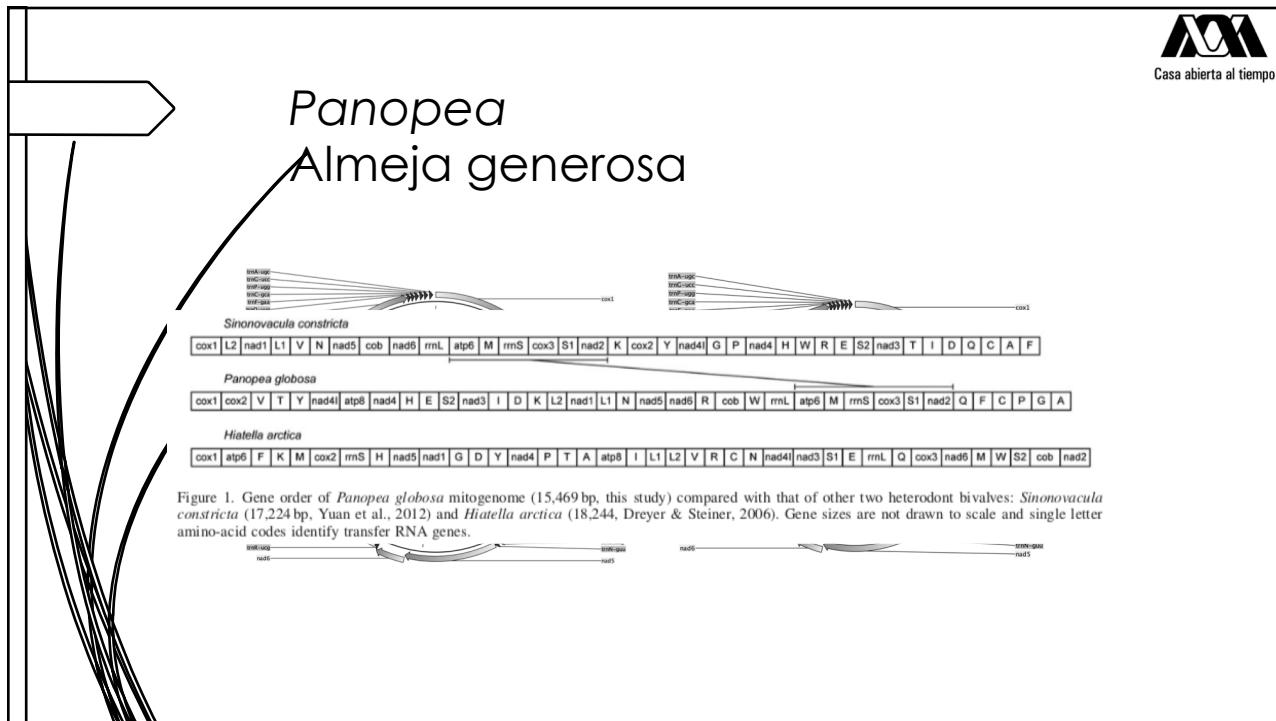
Table 1. The 13 protein, two ribosomal RNA and 22 tRNA genes typically found in animal mitochondrial genomes

Protein encoded	Designation for animal mtDNA	Synonym
Cytochrome oxidase subunit I, II, III	<i>COI, COII, COIII</i>	<i>cox1, cox2, cox3</i>
Cytochrome b apoenzyme	<i>Cytb</i>	<i>cob</i>
NADH dehydrogenase subunits 1-6, 4L	<i>ND1-6, 4L</i>	<i>nad1-6, 4L</i>
ATP synthase subunits 6, 8	<i>A6, A8 or ATP6, ATP8</i>	<i>atp6, atp8</i>
Large ribosomal subunit RNA	<i>IrRNA</i>	<i>rnl</i>
Small ribosomal subunit RNA	<i>srRNA</i>	<i>rns</i>
18 Transfer RNAs each specifying a single amino acid	Corresponding one-letter amino acid code	<i>trnX</i>
Two transfer RNAs specifying leucine	Differentiated by codon recognized, <i>L(CUN)</i> and <i>L(UUR</i>)	Differentiated by subscript
Two transfer RNAs specifying serine	Differentiated by codon recognized, <i>S(AGN)</i> and <i>S(UCN)</i>	Differentiated by subscript

For historical reasons these have designations commonly used only for animal systems. To aid comparison with other works, column three lists the synonymous gene labels used for non-animal systems.

Boore, J.L., 1999. Animal mitochondrial genomes. *Nucleic Acids Research*, 27(8), pp.1767-1780.





Ejemplo

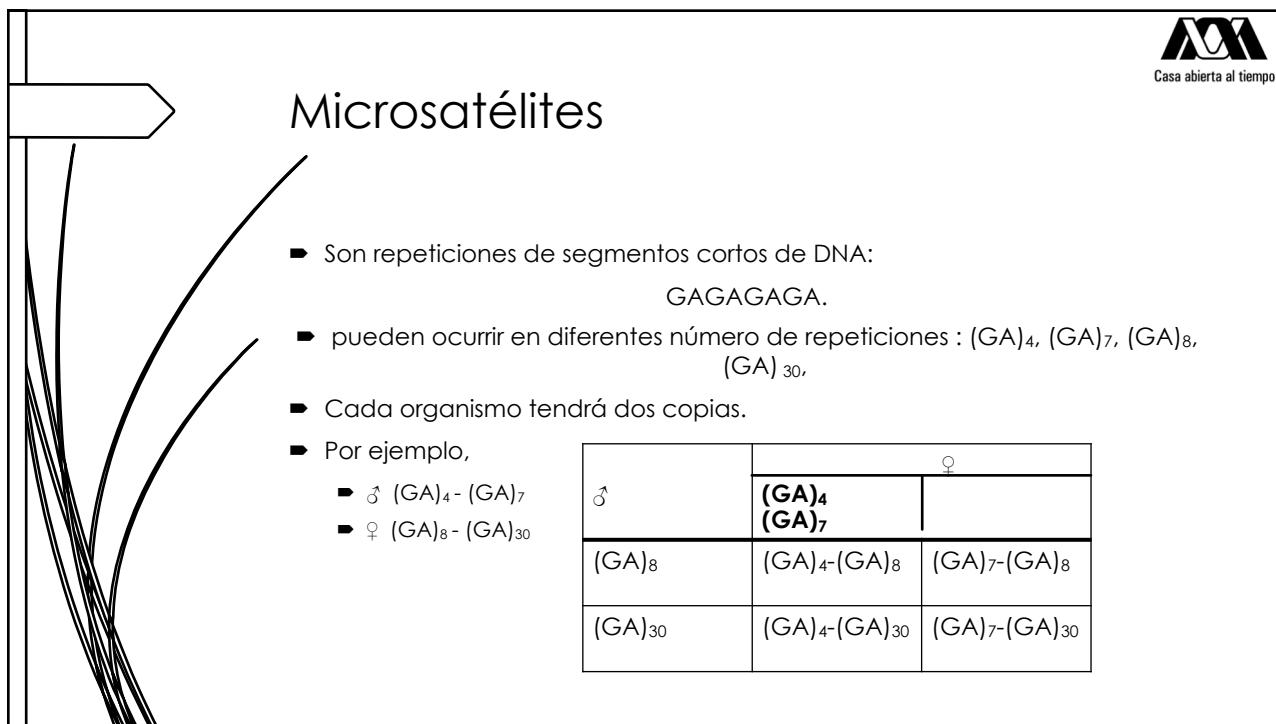
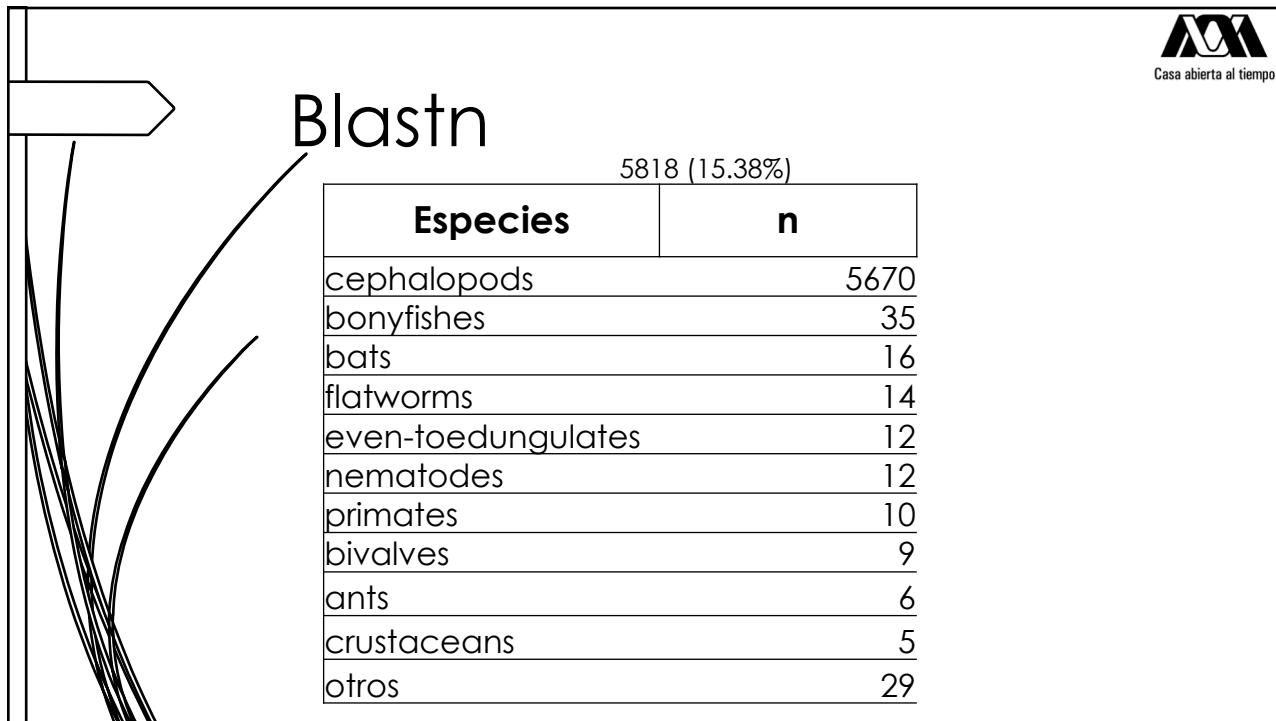
Análisis bioinformático de datos de secuenciación masiva para el diseño de marcadores microsatélites en dos especies de *Octopus* (Cephalopoda: Octopodidae)

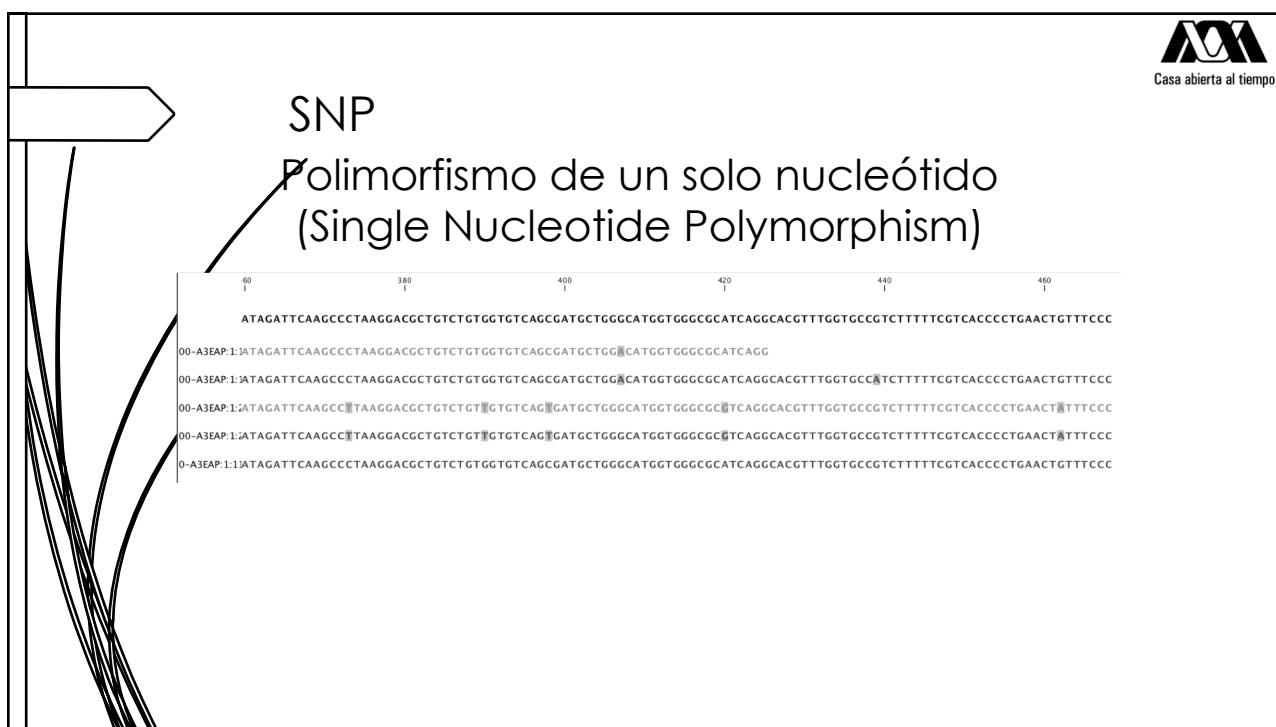
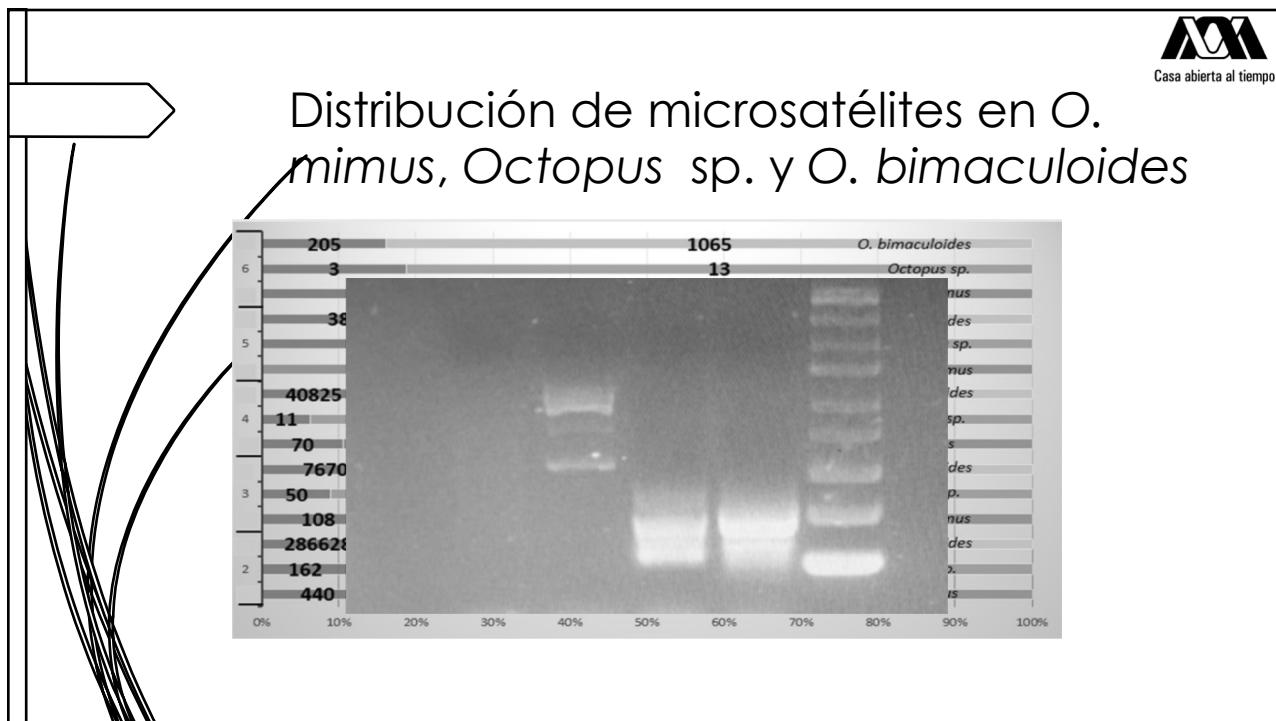
Cecilia Daniela Gutiérrez Hernández.
Asesores,
Dra. Irene de los Ángeles Barriga Sosa.
Dr. Miguel A. del Río Portilla

	O. mimus	Octopus sp.	O. bimaculoides
Lecturas (Reads)	12,714,080	13,393,466	-
N75 (bp)	269	263	-
N50 (bp)	304	298	-
N25 (bp)	353	348	-
Mínimo (bp)	200	200	-
Máximo (bp)	15,078	10,449	-
Promedio (bp)	312	306	-
Contigs (n)	37,784	21,609	-
Total (bp)	11,778,014 (0.49%)	6,613,744 (0.28%)	2,383,190,000
pb entre microsatélite	3044	3353	1034

O. mimus

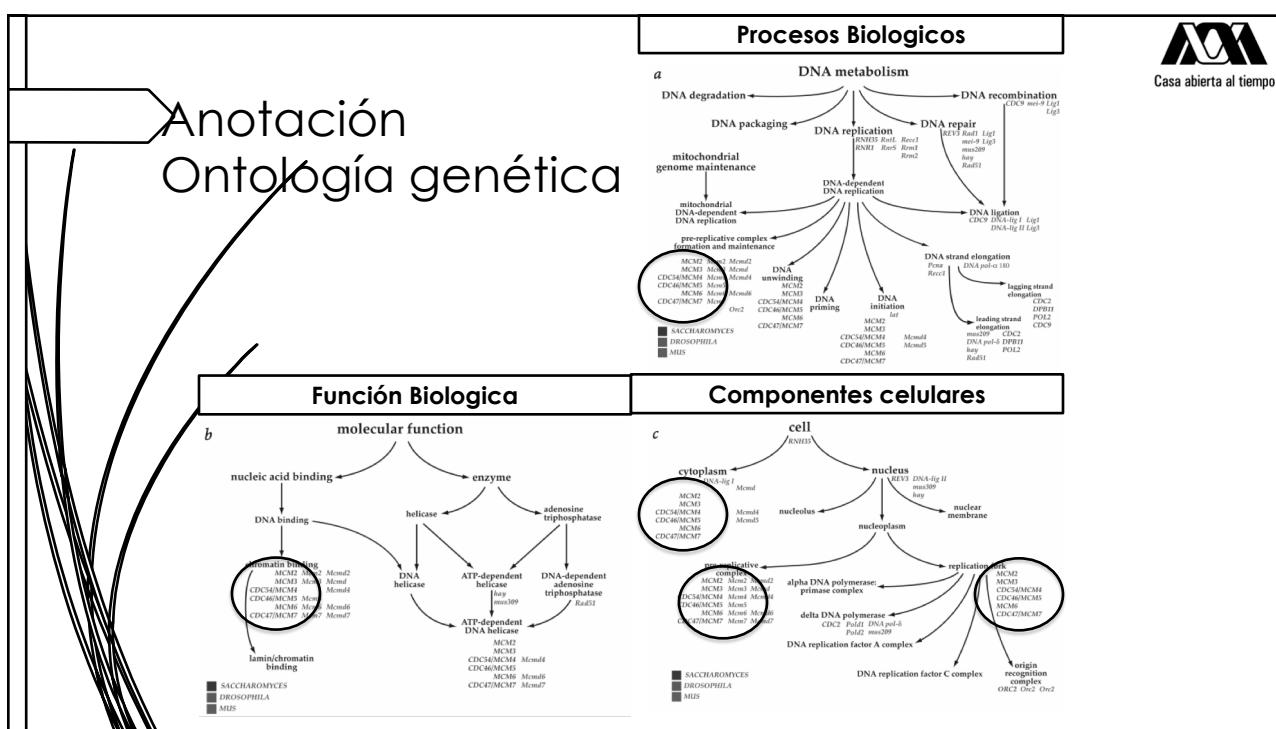
	Largo	GC	CpG
n	37821.000	37821.000	37821.000
Promedio	311.061	35.043	0.573
Desviación estandar	146.922	5.639	0.411
Mínimo	122.000	12.397	0.000
25%	259.000	31.222	0.234
50%	291.000	34.672	0.539
75%	331.000	38.484	0.846
Máximo	15078.000		







ACCA
Casa abierta al tiempo



Anotación de los contigs: Blast, Uniprot

Identificar la localización o función de una secuencia de DNA

blast2go

FastAnnotator

NCBI BLAST UniProt

Gene Ontology Consortium

Panopea transcriptome GO annotation

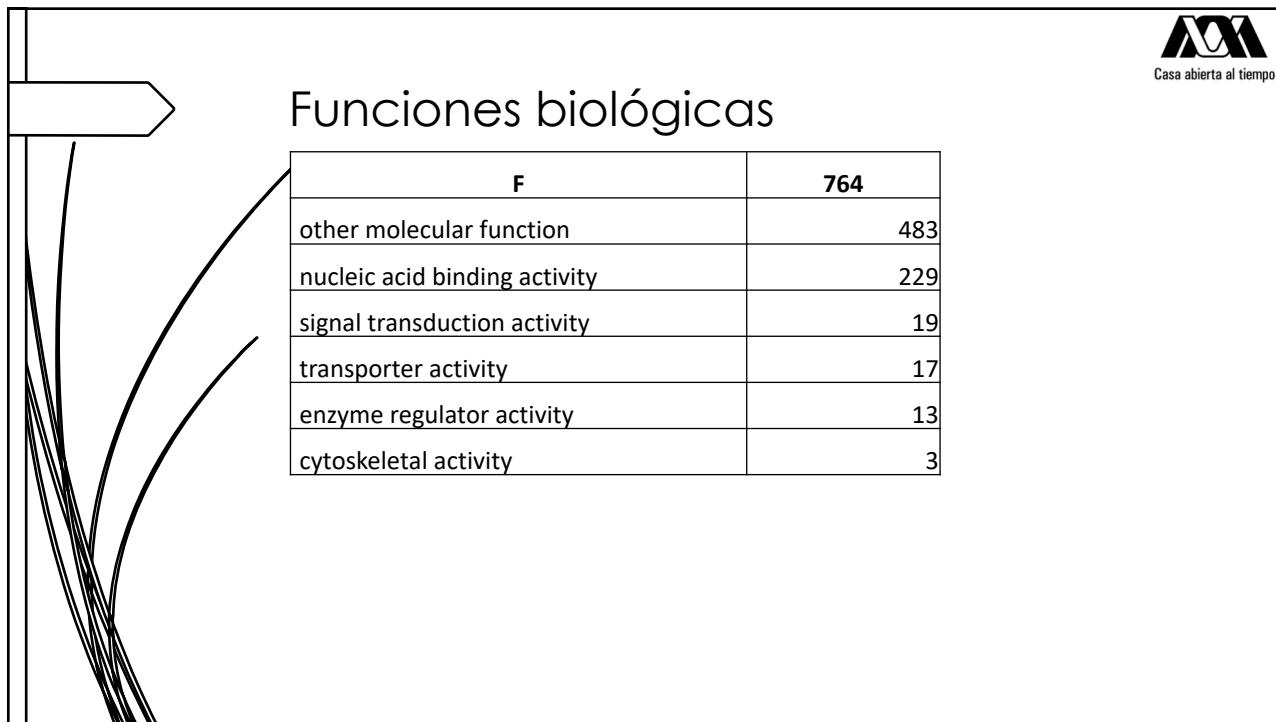
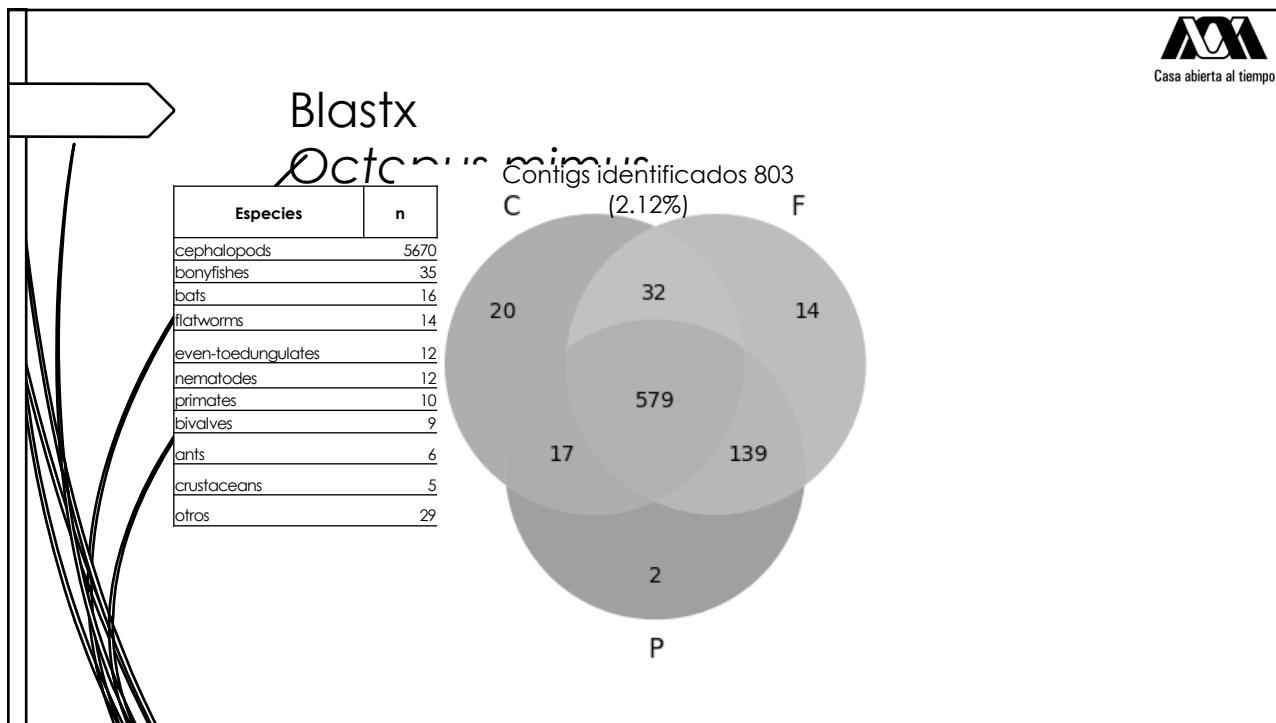
Protein-protein Interaction Plot

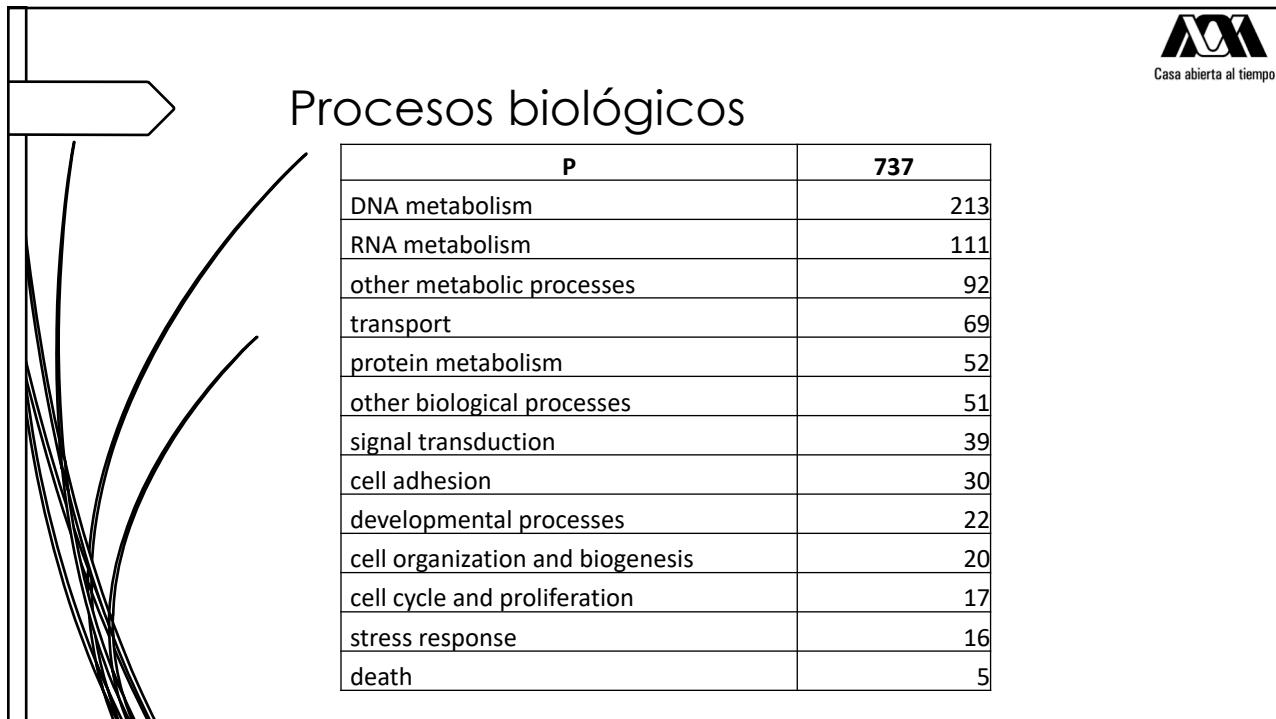
Clasificación de los genes
Ontología génica

Procesos Biológico

Función Molecular

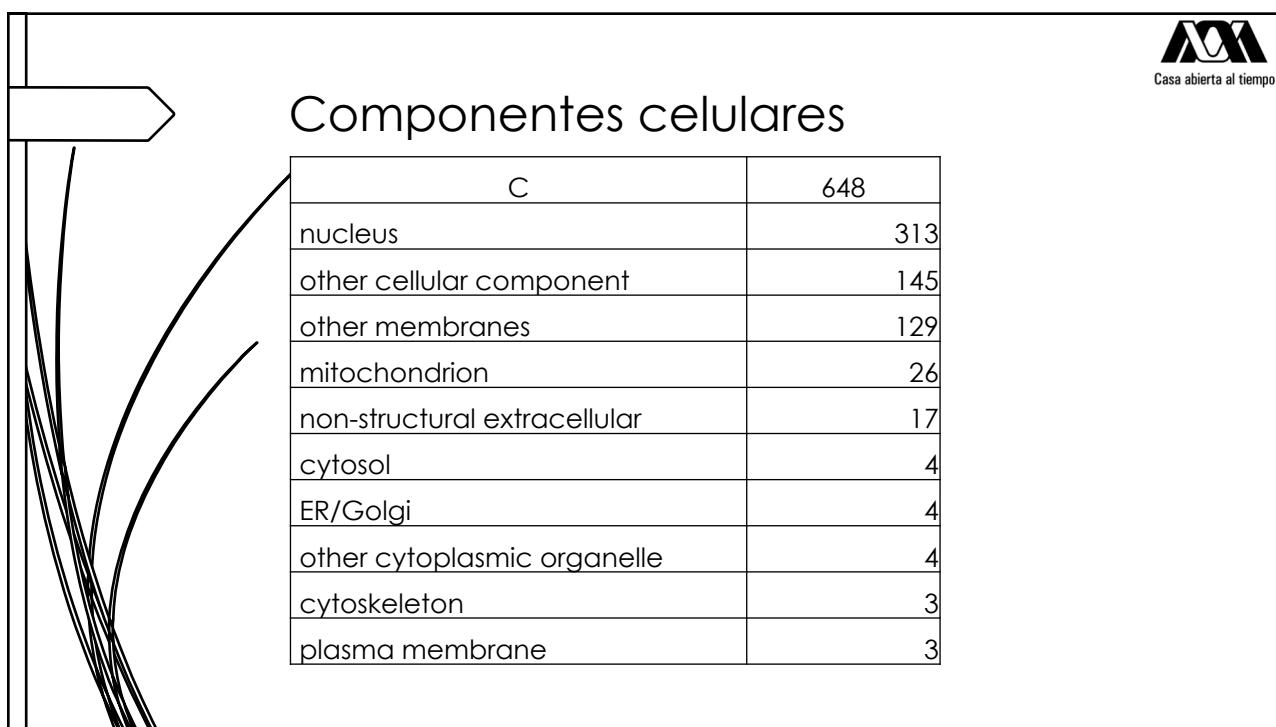
Componentes celulares





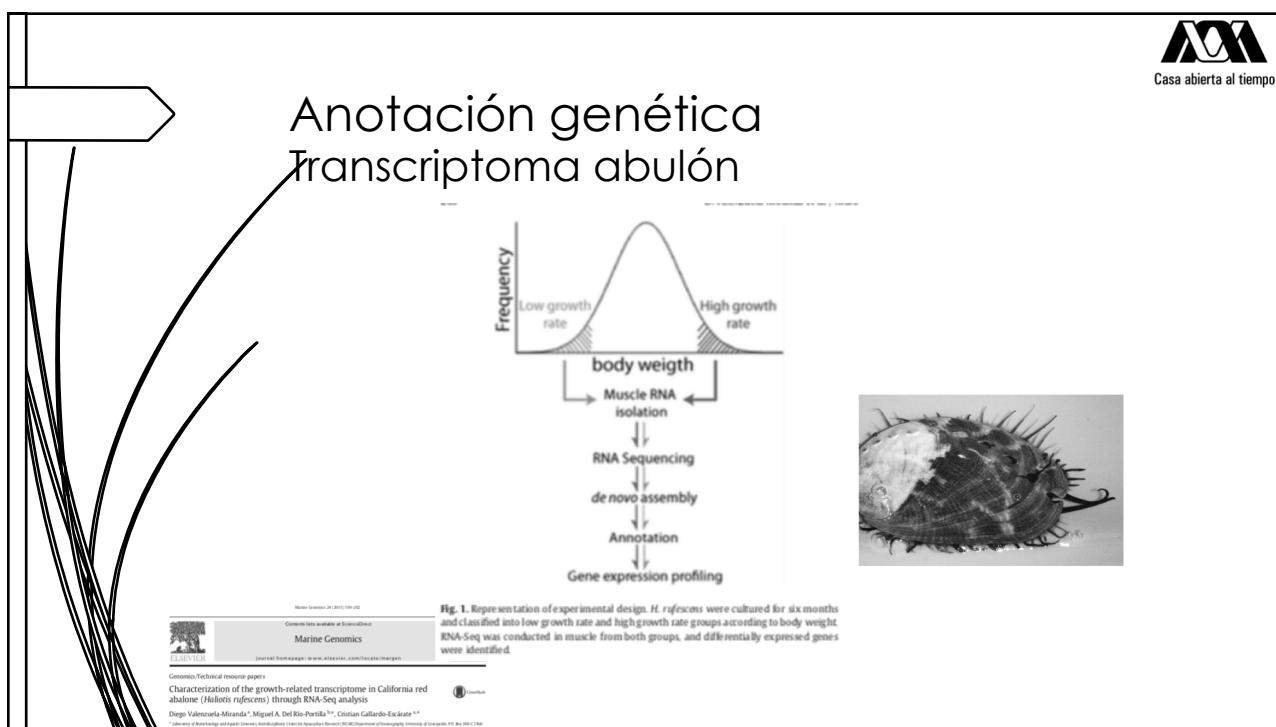
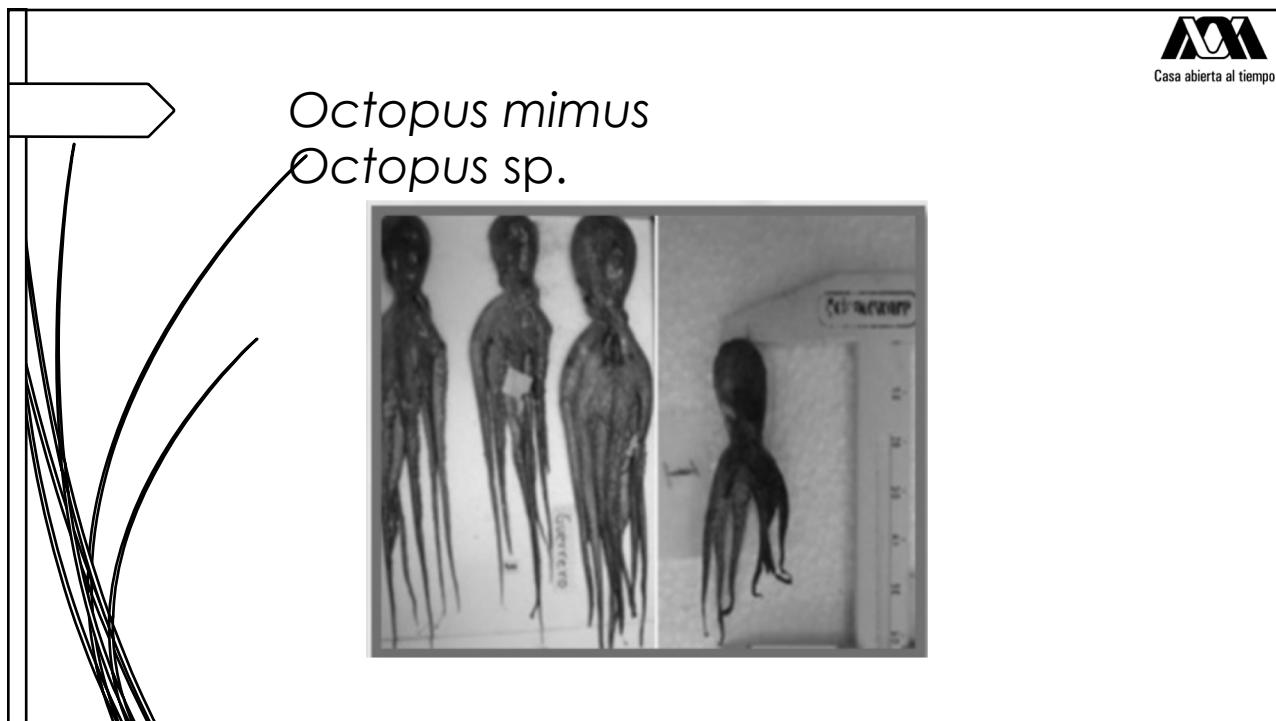
Procesos biológicos

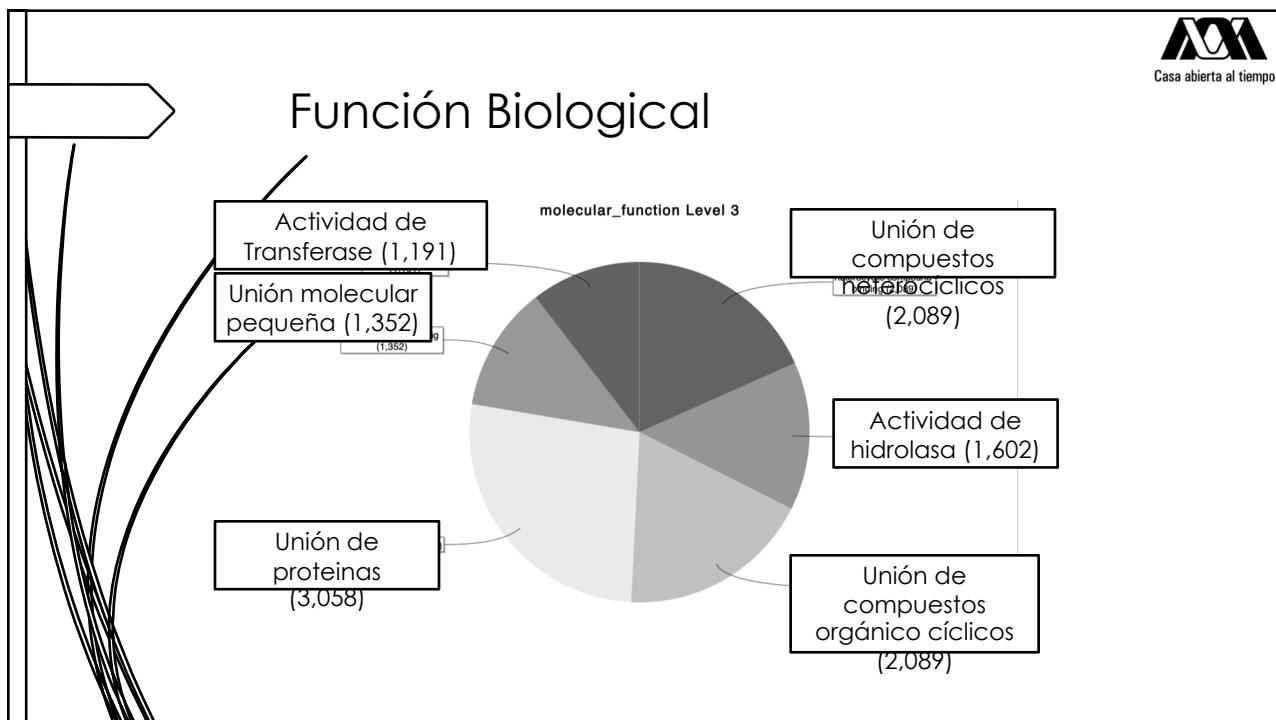
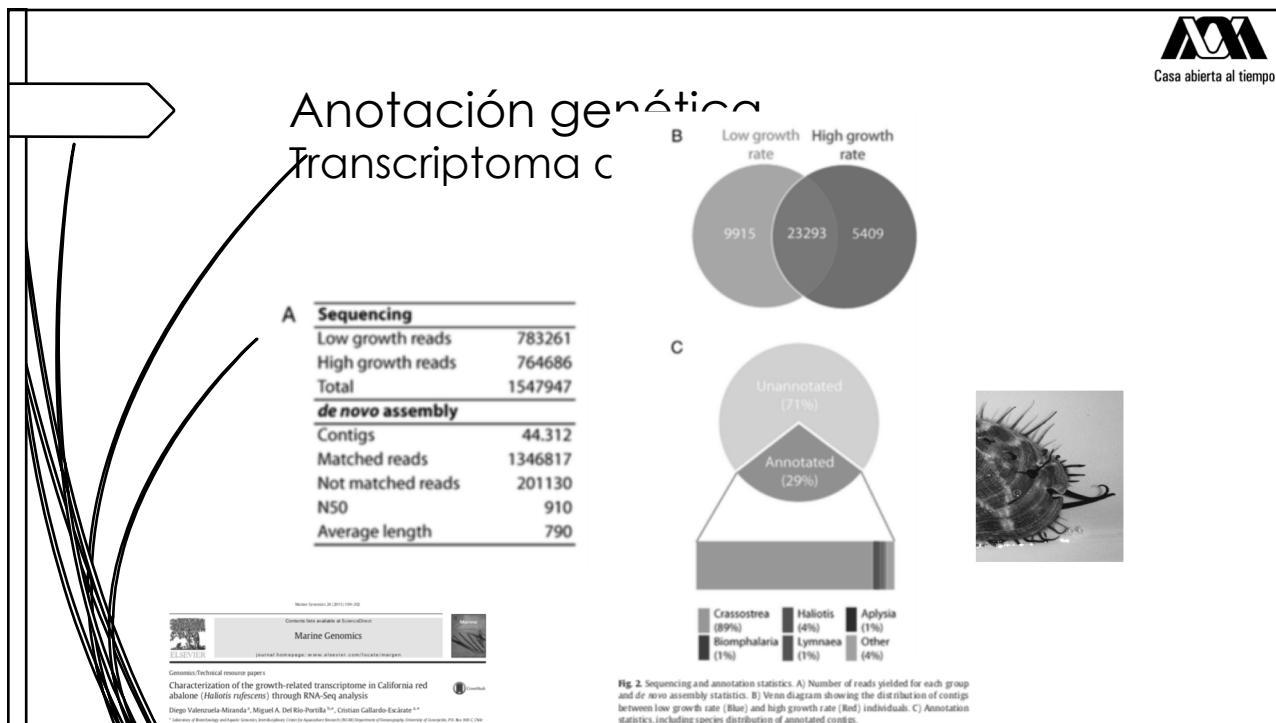
P	737
DNA metabolism	213
RNA metabolism	111
other metabolic processes	92
transport	69
protein metabolism	52
other biological processes	51
signal transduction	39
cell adhesion	30
developmental processes	22
cell organization and biogenesis	20
cell cycle and proliferation	17
stress response	16
death	5

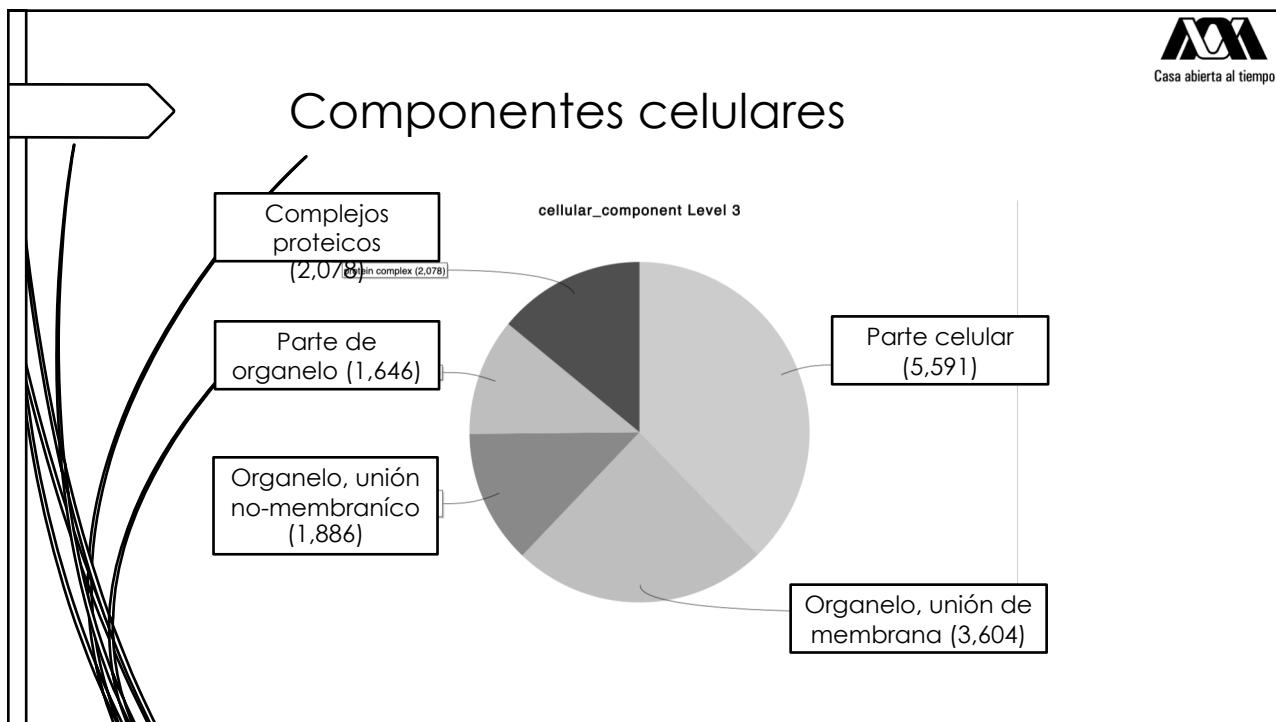
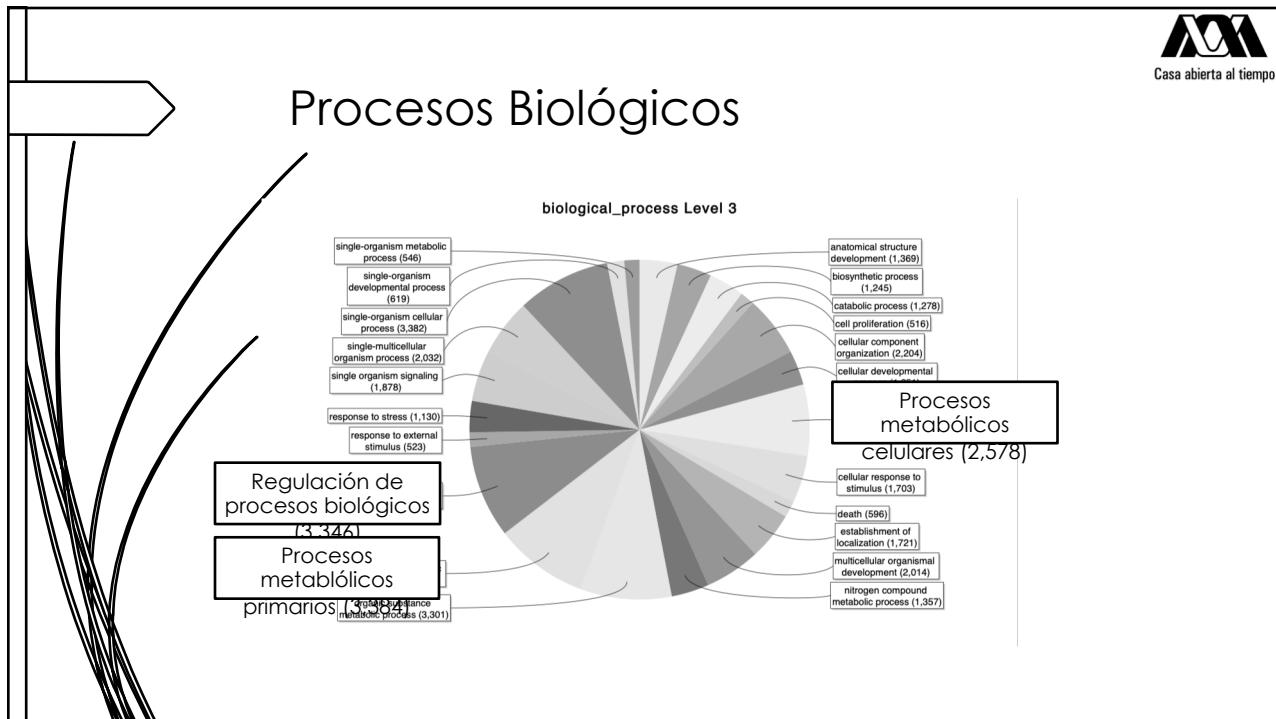


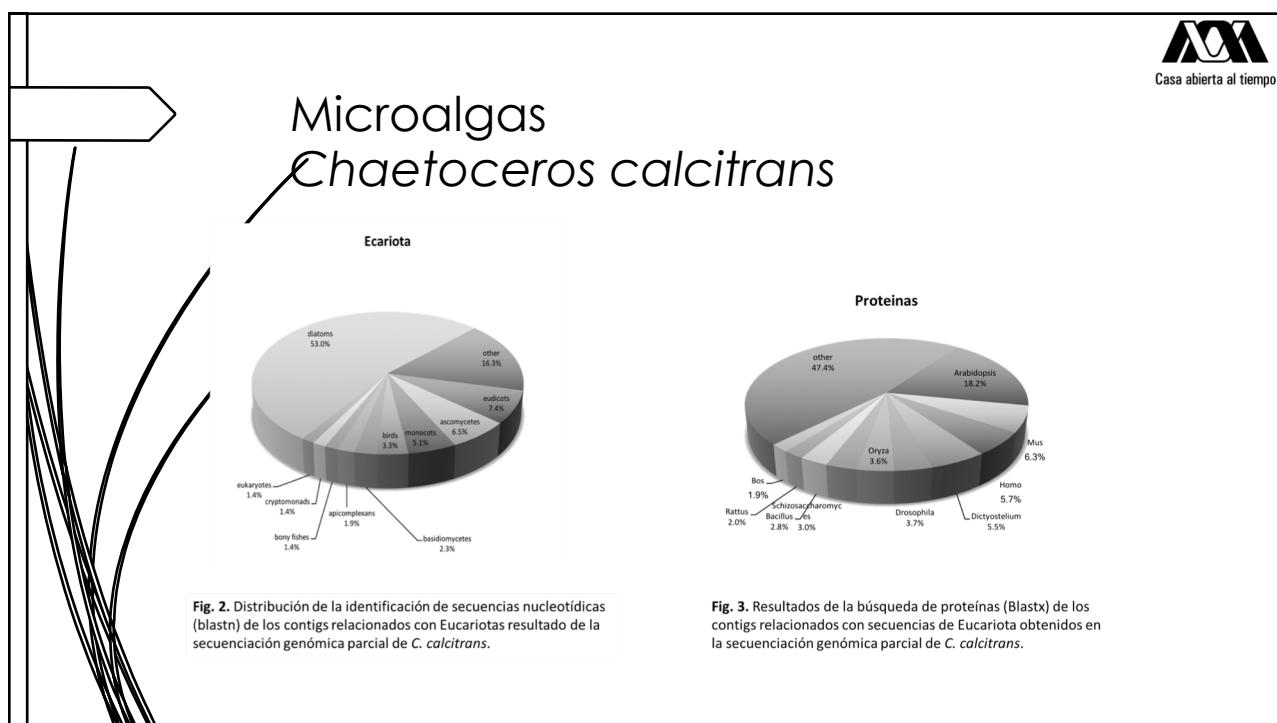
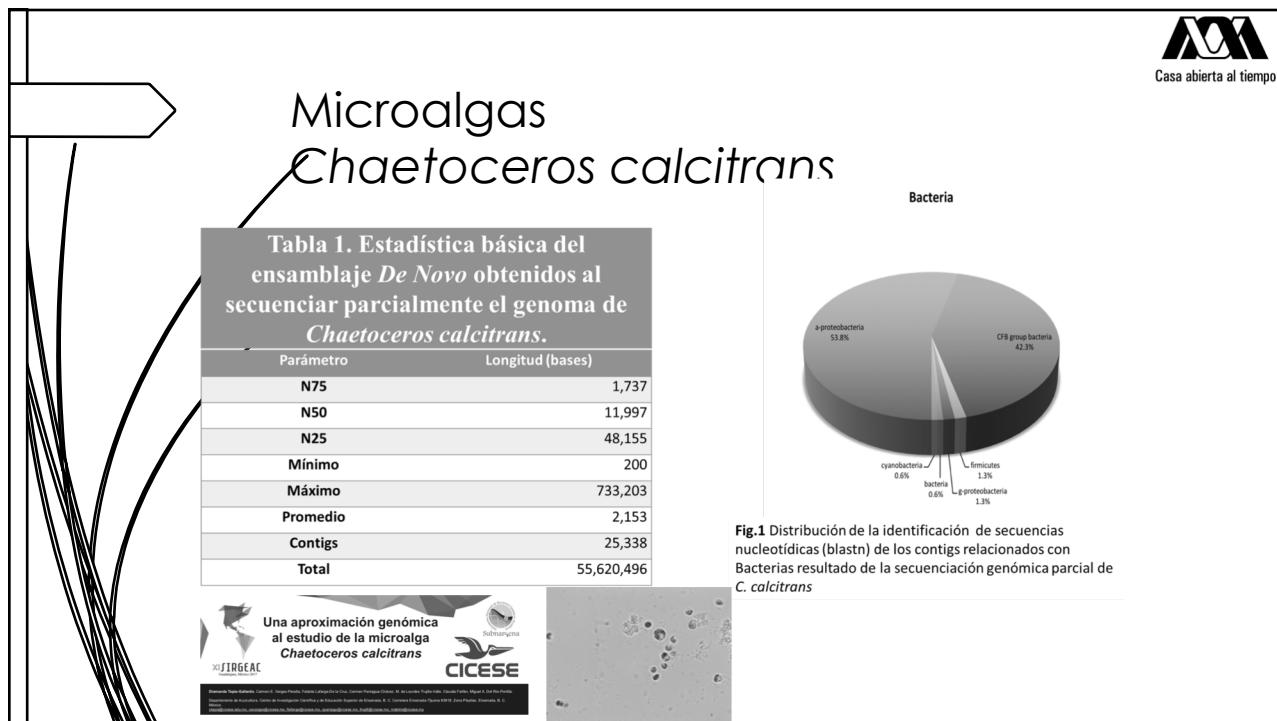
Componentes celulares

C	648
nucleus	313
other cellular component	145
other membranes	129
mitochondrion	26
non-structural extracellular	17
cytosol	4
ER/Golgi	4
other cytoplasmic organelle	4
cytoskeleton	3
plasma membrane	3











Bash

- Bash (Bourne-again shell) es un programa que sirve para interpretar las instrucciones en una línea de comandos.
- Se originó de palabras (born-again significa "nacido de nuevo") sobre la Bourne shell (sh), que fue uno de los primeros intérpretes importantes de Unix.

Anotación de genes nucleares y mitocondriales

<https://es.wikipedia.org/wiki/Bash>

Lenguajes de programación

Depends on who you ask, of course

Will Python dethrone Java this year? Programming language rankings say yes

Rating (%)

IT

August 2, 2018 Gabriela Motroc

<https://jaxenter.com/will-python-dethrone-java-147499.html>

reddit Twitter LinkedIn Facebook Google+

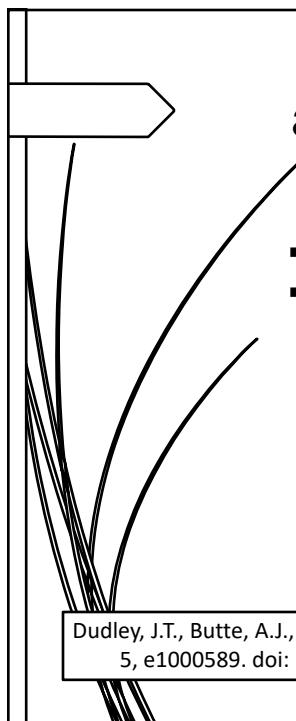
Annotación de genes nucleares y mitocondriales

<https://tiobe.com/tiobe-index>

Bibliotecas

- BioPERL
- BioPython
- BioRuby
- BioJava
- Bio-Conductor (R)

Annotación de genes nucleares y mitocondriales



¿Por qué es importante documentar?

- Uso de bitácoras electrónicas
- “In fact, many journals will now require you to publish your source code along with a manuscript.” (Dudley and Butte, 2009)

Dudley, J.T., Butte, A.J., 2009. A quick guide for developing effective bioinformatics programming skills. PLoS Comput. Biol. 5, e1000589. doi: 10.1371/journal.pcbi.1000589

ACM
Casa abierta al tiempo



Programas para NGS



Milicchio, F., Rose, R., Bian, J., Min, J., Prosperi, M., 2016. Visual programming for next-generation sequencing data analytics. BioData Min. 9, 16. doi:10.1186/s13040-016-0095-3

ACM
Casa abierta al tiempo

The diagram features a central vertical line with several diagonal lines branching off to the left, creating a complex, organic shape that looks like a brain or a network of connections.

Casa abierta al tiempo

- The Importance of Building Your Technology Toolbox
- · The Benefits and Opportunities of Open Source Communities
- · The Importance of UNIX Skills
- · Keeping Projects Documented and Manageable
- · Preserving Your Source Code
- · Embracing Parallel Computing Paradigms
- · Structuring Data for Speed and Scalability
- · Understand the Capabilities of Hardware
- · Embracing Standards and Interoperability
- · Value Your Time

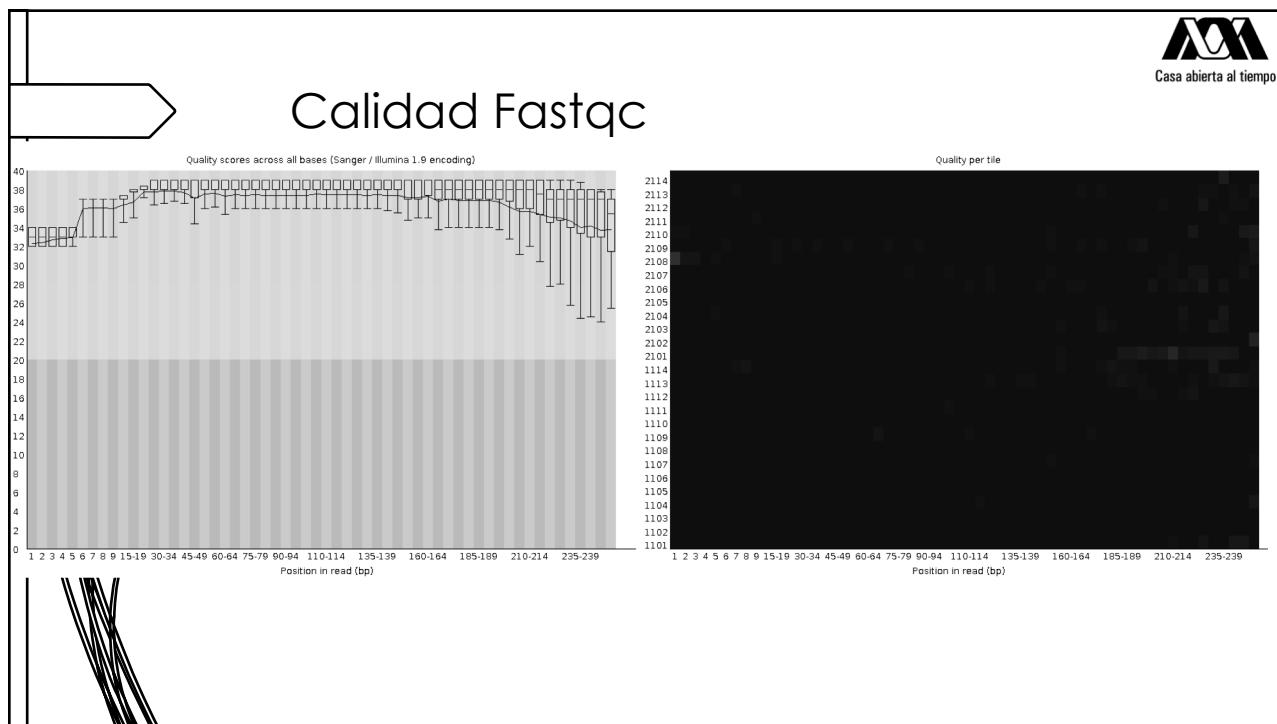
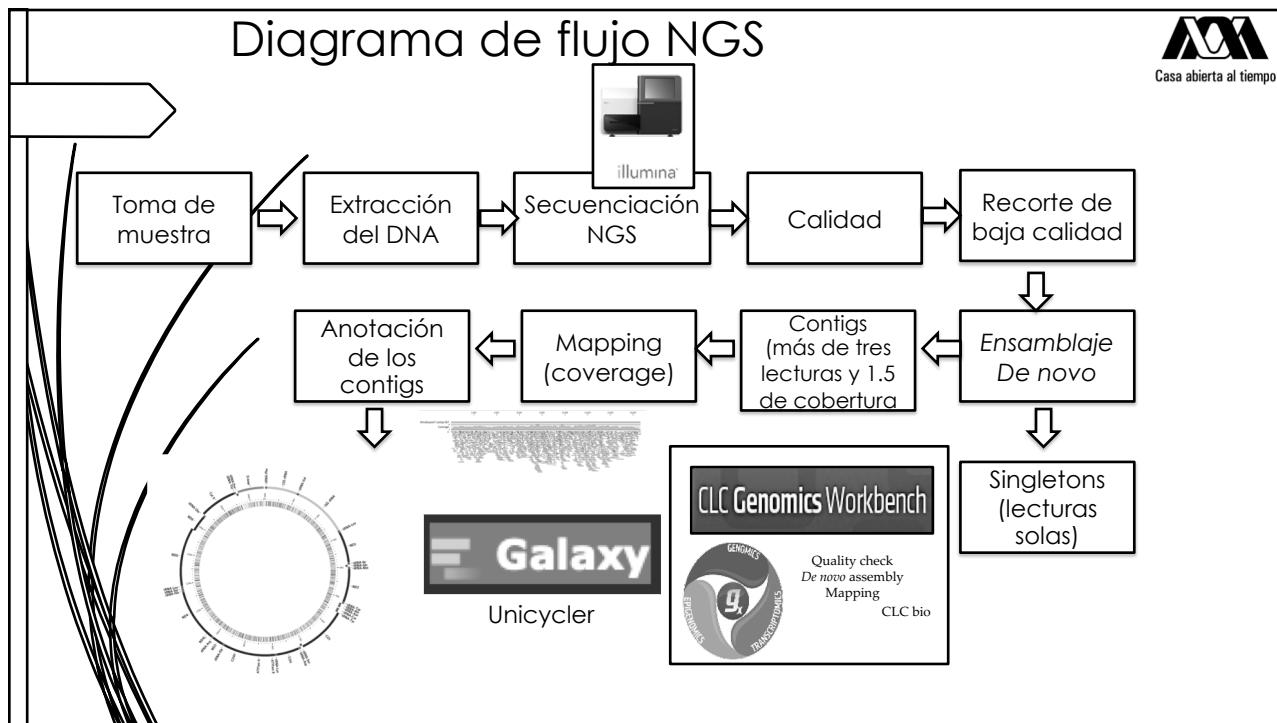
Dudley, J.T., Butte, A.J., 2009. A quick guide for developing effective bioinformatics programming skills. PLoS Comput. Biol. 5, e1000589. doi: 10.1371/journal.pcbi.1000589

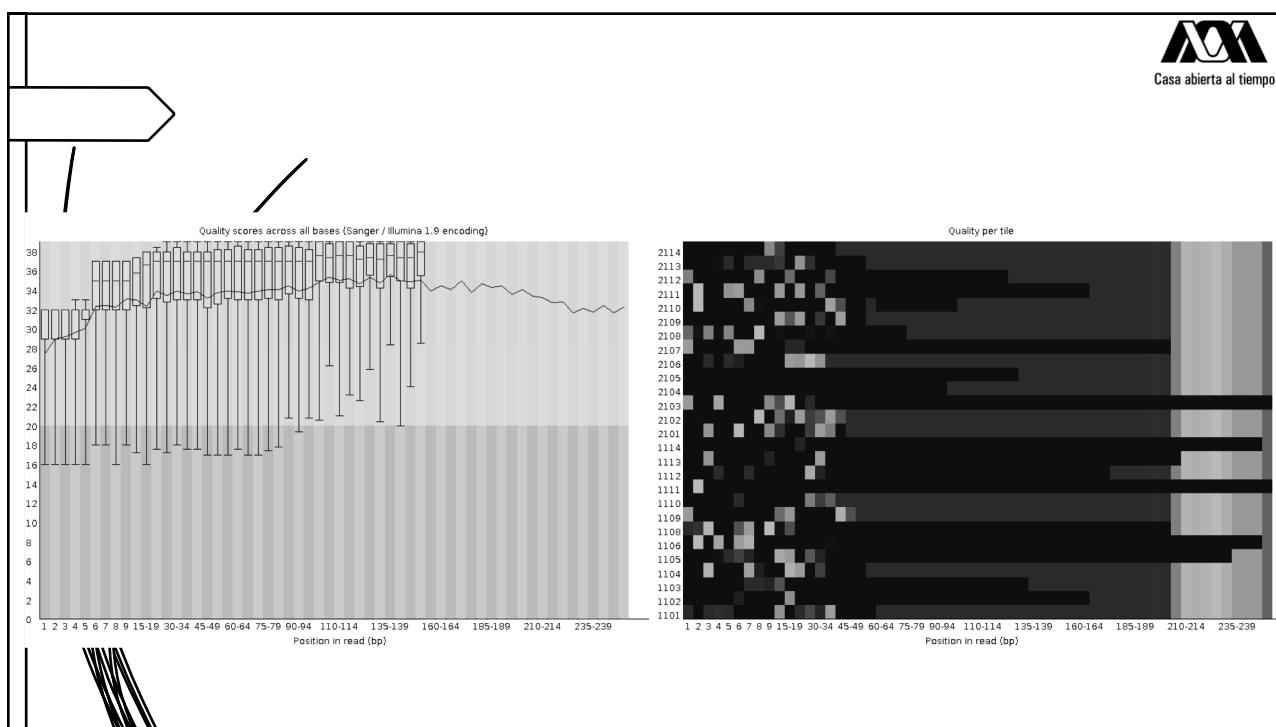
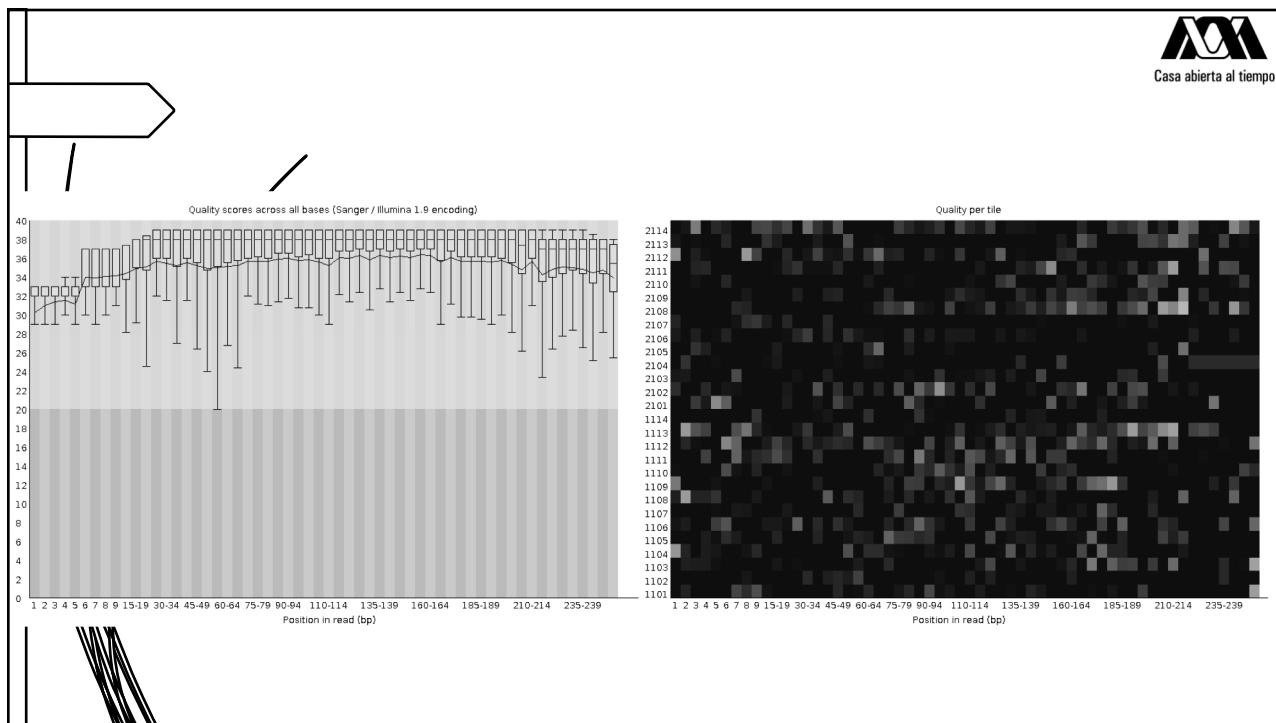
The diagram features a central vertical line with several diagonal lines branching off to the left, creating a complex, organic shape that looks like a brain or a network of connections.

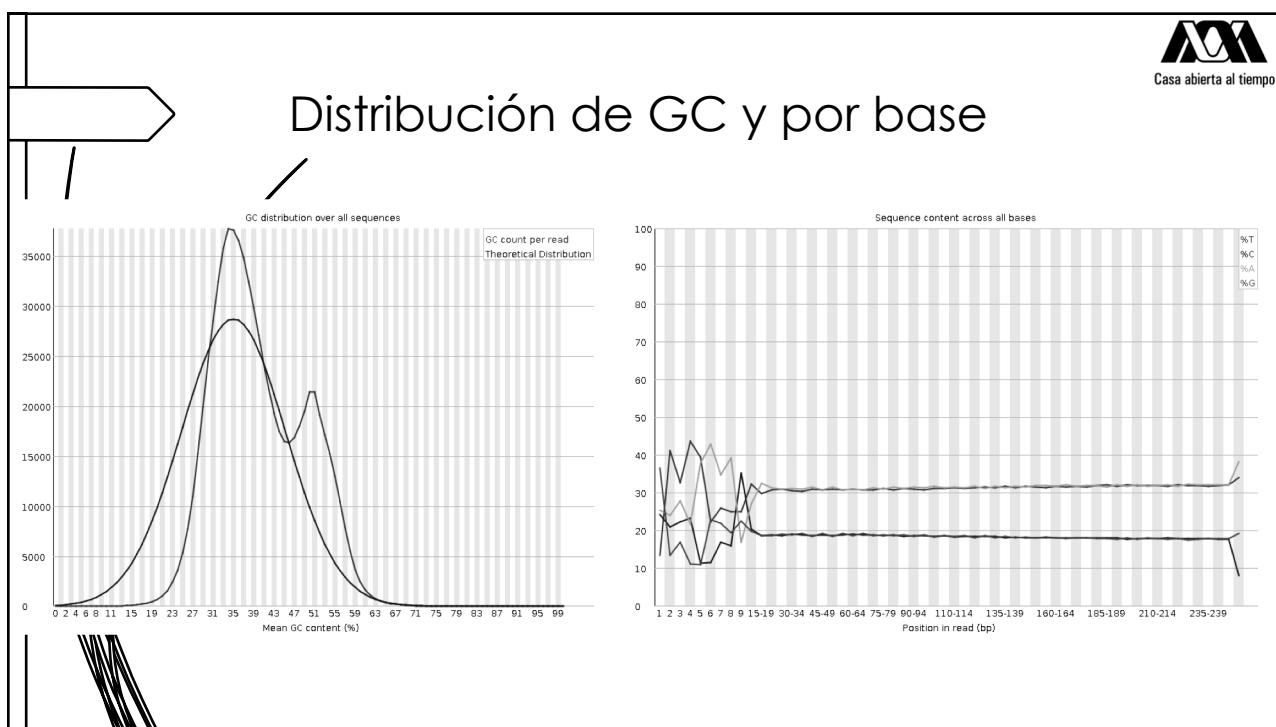
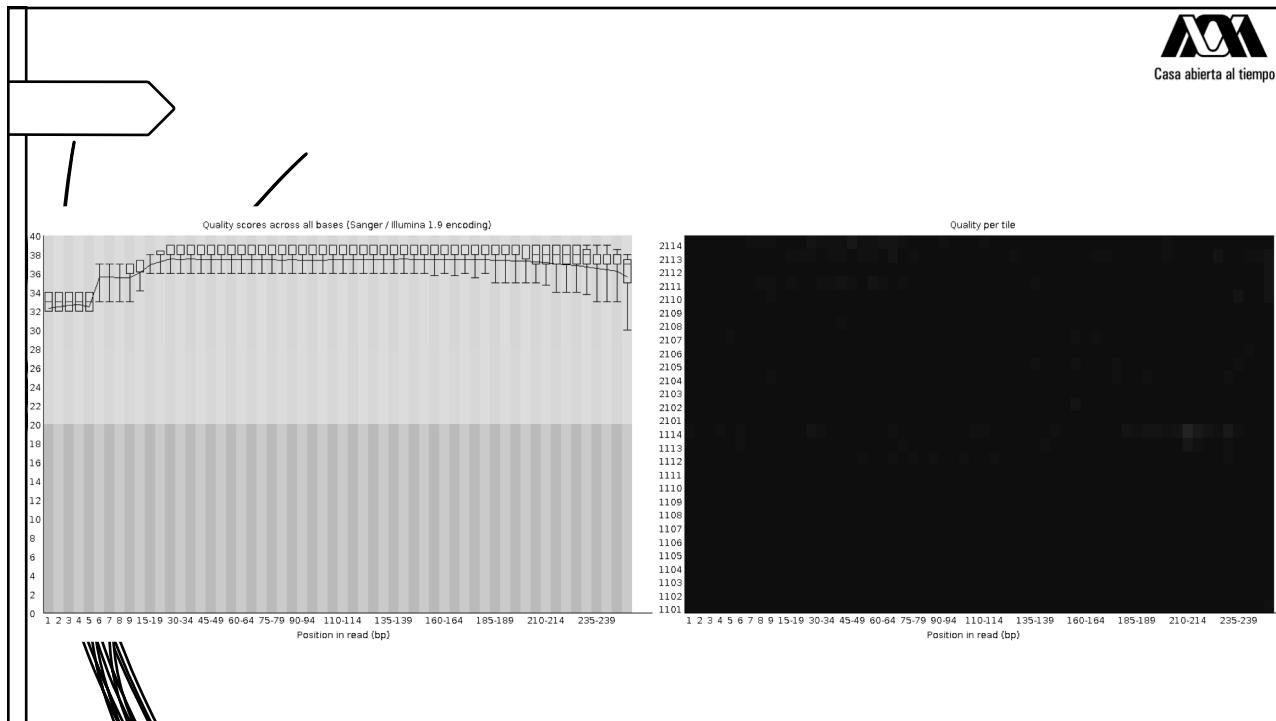
Casa abierta al tiempo

- Desarrollo de programas
- · “Consequently it is no surprise that many successful bioinformatics apps are written by biologists who lack formal computer science training, · as they un · doubtedly put scientific utility ahead of architectural elegance and completeness.”

Dudley, J.T., Butte, A.J., 2009. A quick guide for developing effective bioinformatics programming skills. PLoS Comput. Biol. 5, e1000589. doi: 10.1371/journal.pcbi.1000589

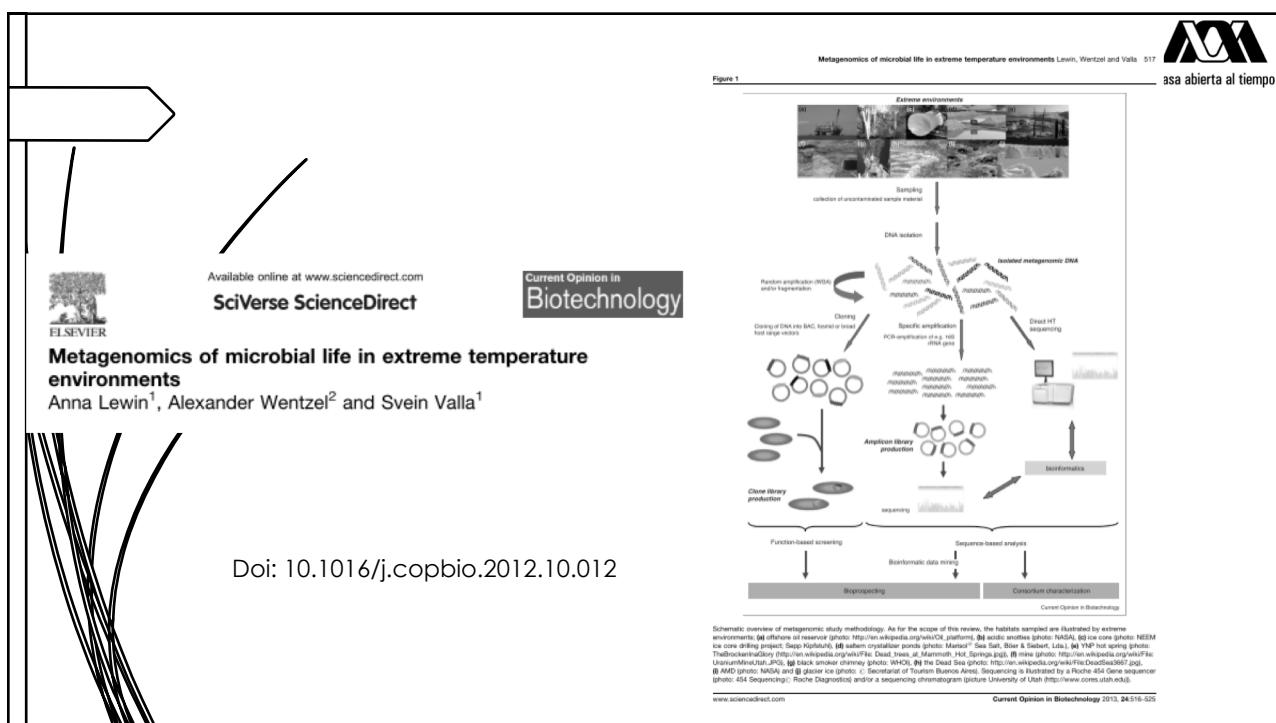


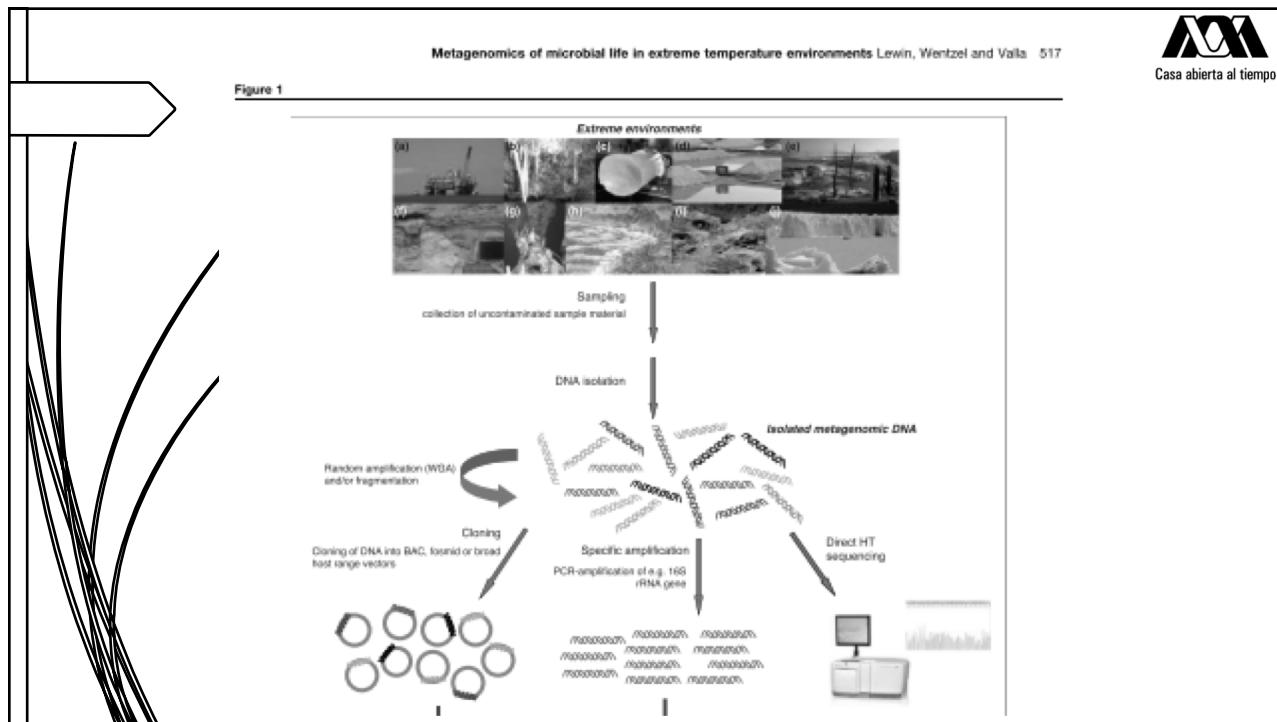
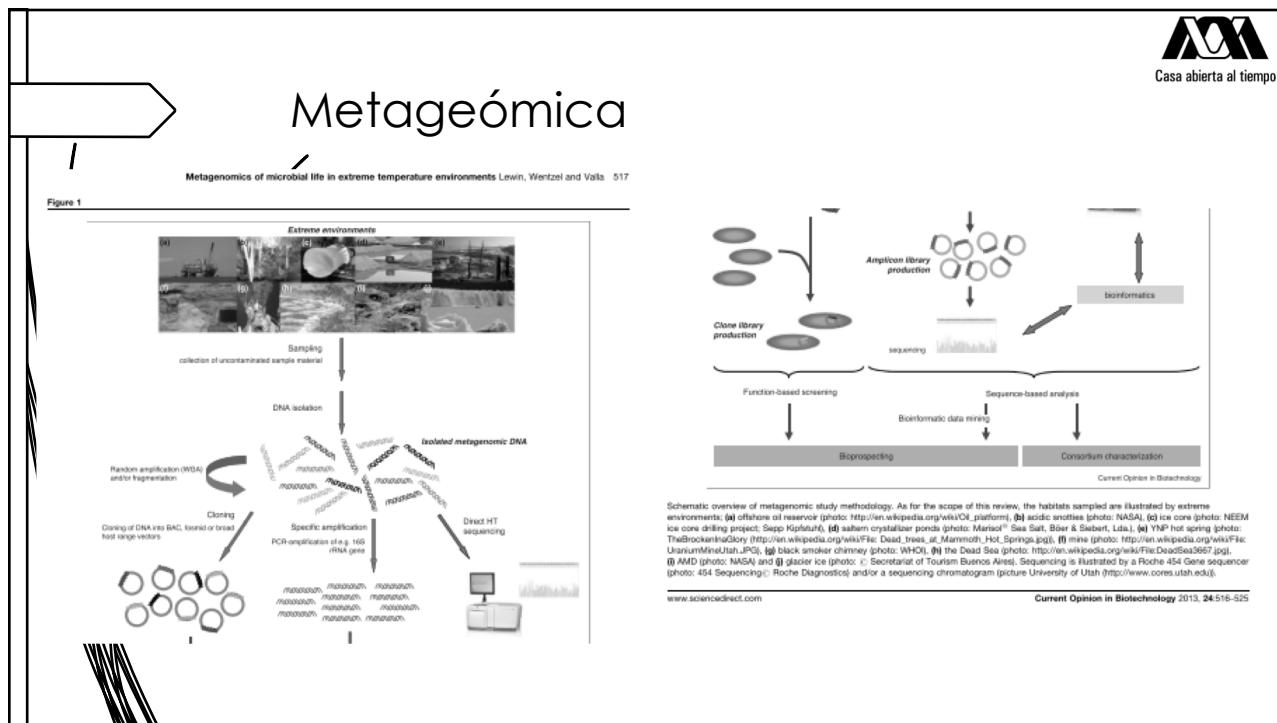


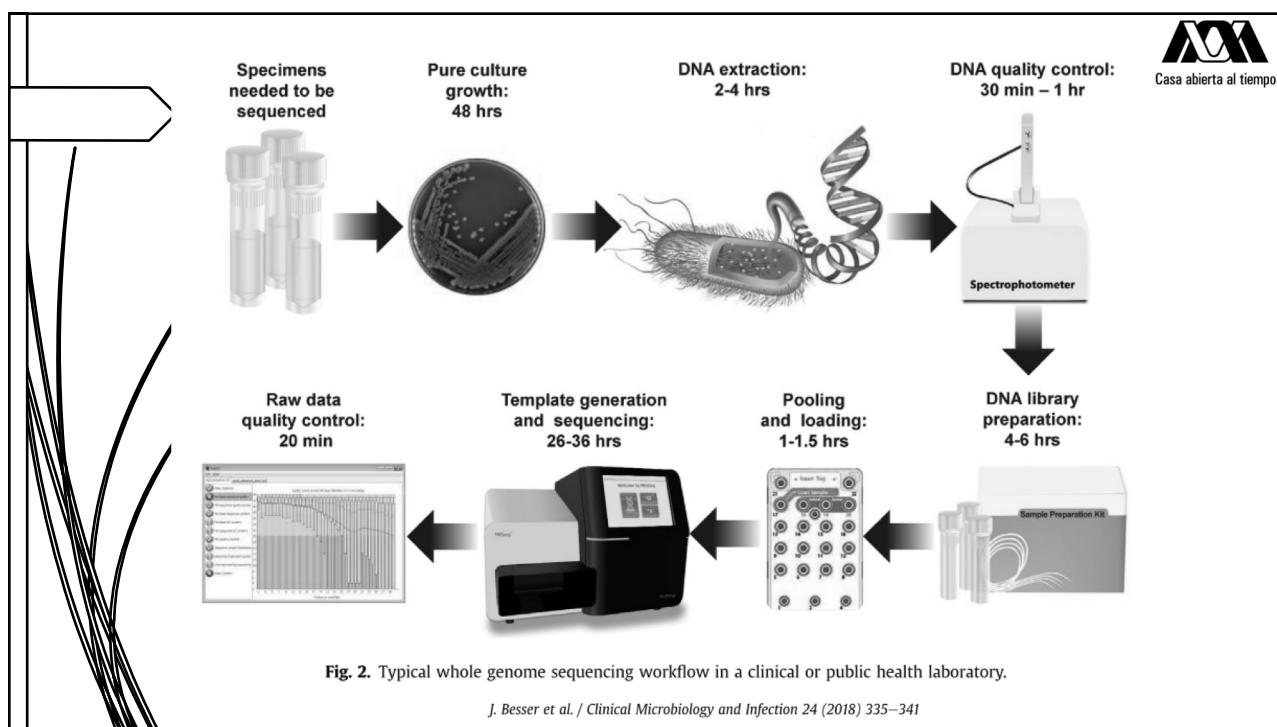
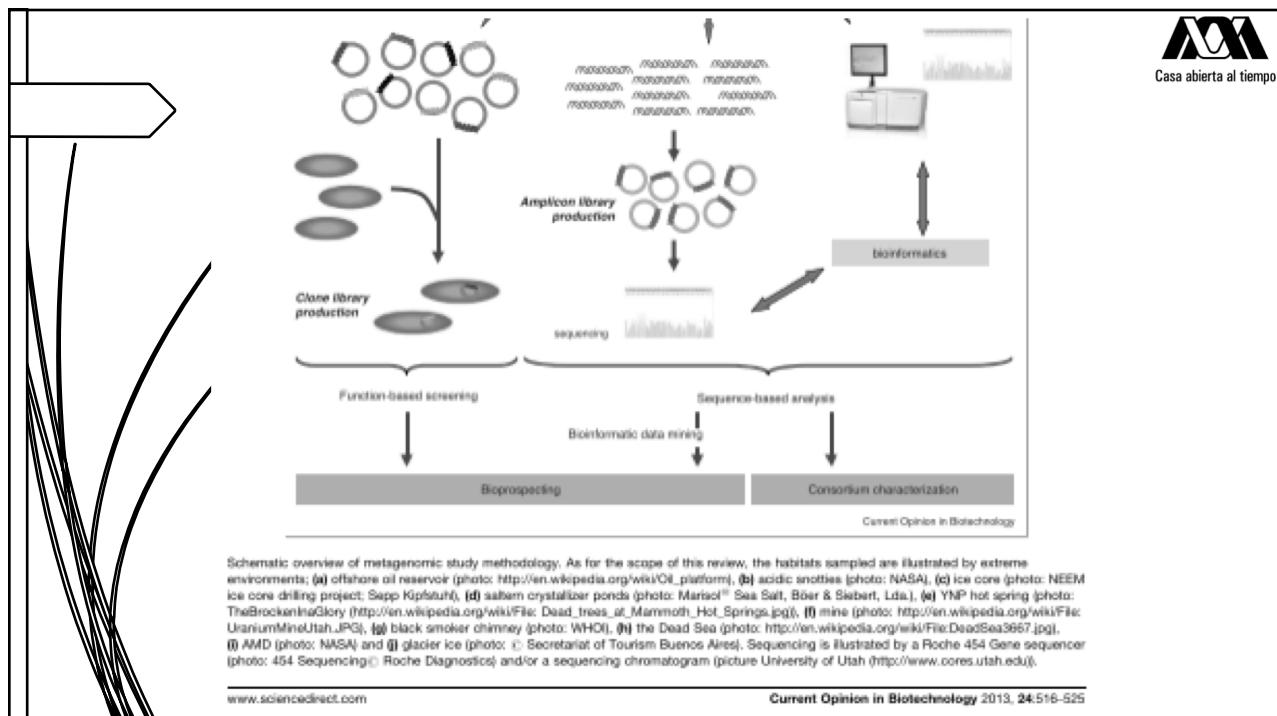




ACCA
Casa abierta al tiempo







**Table 2**

Characteristics, strengths and weaknesses of commonly used sequencing platforms

Platform \ Instrument	Throughput range (Gb) ^a	Read length (bp)	Strength	Weakness
<i>Sanger sequencing</i>				
ABI 3500/3730	0.0003	Up to 1 kb	Read accuracy and length	Cost and throughput
<i>Illumina</i>				
MiniSeq	1.7–7.5	1×75 to ×150	Low initial investment	Run and read length
MiSeq	0.3–15	1×36 to 2×300	Read length, scalability	Run length
NextSeq	10–120	1×75 to 2×150	Throughput	Run and read length
HiSeq (2500)	10–1000	×50 to ×250	Read accuracy, throughput,	High initial investment, run
NovaSeq 5000/6000	2000–6000	2×50 to ×150	Read accuracy, throughput	High initial investment, run
<i>IonTorrent</i>				
PGM	0.08–2	Up to 400	Read length, speed	Throughput, homopolymers ^c
S5	0.6–15	Up to 400	Read length, speed,	Homopolymers ^c
Proton	10–15	Up to 200	Speed, throughput	Homopolymers ^c
<i>Pacific BioSciences</i>				
PacBio RSII	0.5–10 ^b	Up to 60 kb	Read length, speed (Average 10 kb, N50 20 kb)	High error rate and initial
Sequel	5–10 ^b	Up to 60 kb	Read length, speed (Average 10 kb, N50 20 kb)	High error rate
<i>Oxford Nanopore</i>				
MinION	0.1–1	Up to 100 kb	Read length, portability	High error rate, run length,

^a The throughput ranges are determined by available kits and run modes on a per run basis. As an example of a 15-GB throughput, thirty-five 5-MB genomes can be sequenced to a minimum coverage of 40× on the Illumina MiSeq using the v3 600 cycle chemistry.

^b Per one single-molecule real-time cell.

^c Results in increased error rate (increased proportion of reads containing errors among all reads) which in turn results in false-positive variant calling.

Doi: 10.1016/j.cmi.2017.10.013