

Bioinformática aplicada a recursos genéticos

Dr. Miguel Ángel del Río Portilla
Profesor Investigador, Cátedra “Alfonso Villalobos”

Dra. Irene de los Ángeles Barriga Sosa
Profesor Investigador tiempo completo

Dra. Dra. Erika Magallón Gayón
Posdoctorado

Organizadores

- * Dra. Irene de los Ángeles Barriga Sosa
- * Dr. Miguel Ángel del Río Portilla
- * Departamento de Hidrobiología
- * UAM- Iztapalapa

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

Código Genético

- * El conjunto de correspondencia entre tripletes (tres nucleótidos) de DNA y los aminoácidos en proteínas(1)
- * Anteriormente se pensaba que era estático (2)

Códigos genéticos considerados en el NCBi (3)

1. The Standard Code
2. The Vertebrate Mitochondrial Code
3. The Yeast Mitochondrial Code
4. The Mold, Protozoan, and Coelenterate Mitochondrial Code and the Mycoplasma/Spiroplasma Code
5. The Invertebrate Mitochondrial Code
6. The Ciliate, Dasycladacean and Hexamita Nuclear Code
9. The Echinoderm and Flatworm Mitochondrial Code
10. The Euplotid Nuclear Code
11. The Bacterial, Archaeal and Plant Plastid Code
12. The Alternative Yeast Nuclear Code
13. The Ascidian Mitochondrial Code
14. The Alternative Flatworm Mitochondrial Code
16. Chlorophycean Mitochondrial Code
21. Trematode Mitochondrial Code
22. Scenedesmus obliquus Mitochondrial Code
23. Thraustochytrium Mitochondrial Code
24. Pterobranchia Mitochondrial Code
25. Candidate Division SR1 and Gracilibacteria Code
26. Pachysolen tannophilus Nuclear Code
27. Karyorelict Nuclear
28. Condylostoma Nuclear
29. Mesodinium Nuclear
30. Peritrich Nuclear
31. Blastocrithidia Nuclear

1. <https://www.ncbi.nlm.nih.gov/books/NBK21950/>
2. <https://www.ncbi.nlm.nih.gov/pubmed/1579111?dopt=Abstract>
3. <https://www.ncbi.nlm.nih.gov/Taxonomy/Utils/wprintgc.cgi>

Código estándar (1)

TTT	F	Phe	TCT	S	Ser	TAT	Y	Tyr	TGT	C	Cys
TTC	F	Phe	TCC	S	Ser	TAC	Y	Tyr	TGC	C	Cys
TTA	L	Leu	TCA	S	Ser	TAA	*	Ter	TGA	*	Ter
TTG	L	Leu i	TCG	S	Ser	TAG	*	Ter	TGG	W	Trp
CTT	L	Leu	CCT	P	Pro	CAT	H	His	CGT	R	Arg
CTC	L	Leu	CCC	P	Pro	CAC	H	His	CGC	R	Arg
CTA	L	Leu	CCA	P	Pro	CAA	Q	Gln	CGA	R	Arg
CTG	L	Leu i	CCG	P	Pro	CAG	Q	Gln	CGG	R	Arg
ATT	I	Ille	ACT	T	Thr	AAT	N	Asn	AGT	S	Ser
ATC	I	Ille	ACC	T	Thr	AAC	N	Asn	AGC	S	Ser
ATA	I	Ille	ACA	T	Thr	AAA	K	Lys	AGA	R	Arg
ATG	M	Met i	ACG	T	Thr	AAG	K	Lys	AGG	R	Arg
GTT	V	Val	GCT	A	Ala	GAT	D	Asp	GGT	G	Gly
GTC	V	Val	GCC	A	Ala	GAC	D	Asp	GGC	G	Gly
GTA	V	Val	GCA	A	Ala	GAA	E	Glu	GGA	G	Gly
GTG	V	Val	GCG	A	Ala	GAG	E	Glu	GGG	G	Gly

Secuenciación (Sanger)

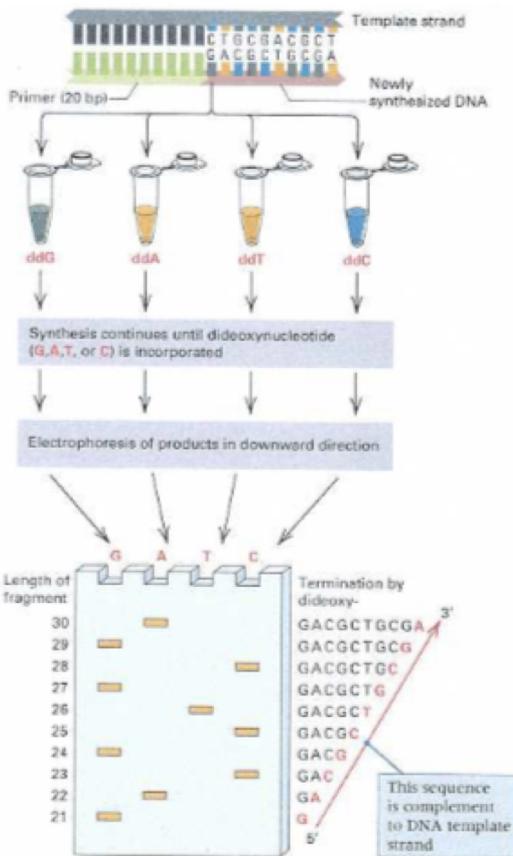


Figure 5.39

Dideoxy method of DNA sequencing. Four DNA synthesis reactions are carried out in the presence of all 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 at the right. The fragments are separated by size by electrophoresis, and the positions of the nucleotides are determined directly from the gel. In this example, the length of the primer needed

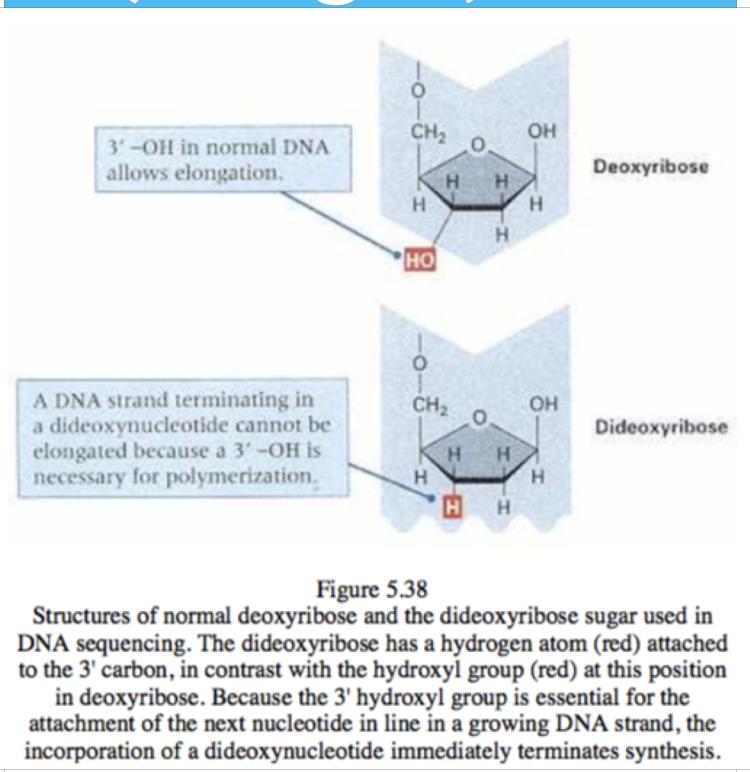


Figure 5.38

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.

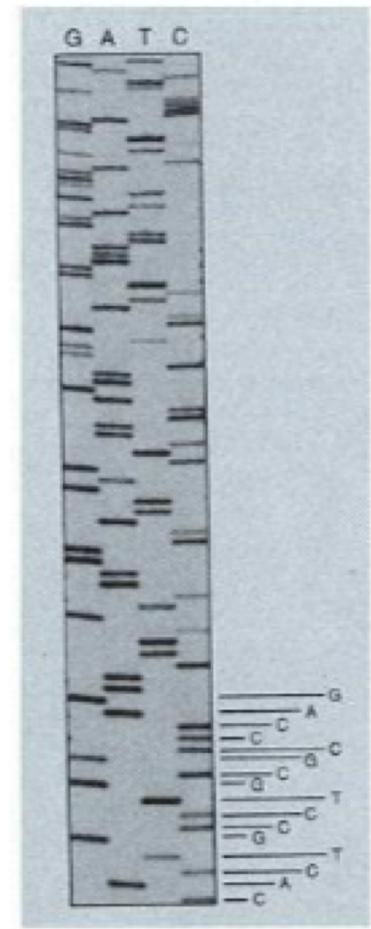


Figure 5.40

A section of a dideoxy sequencing gel. The sequence is read from the bottom to the top. Each horizontal row represents a single nucleotide position in the DNA strand synthesized from the template. The vertical columns result from termination by the dideoxy forms of G, A, T, or C. The sequence from the lower part of the gel is indicated.

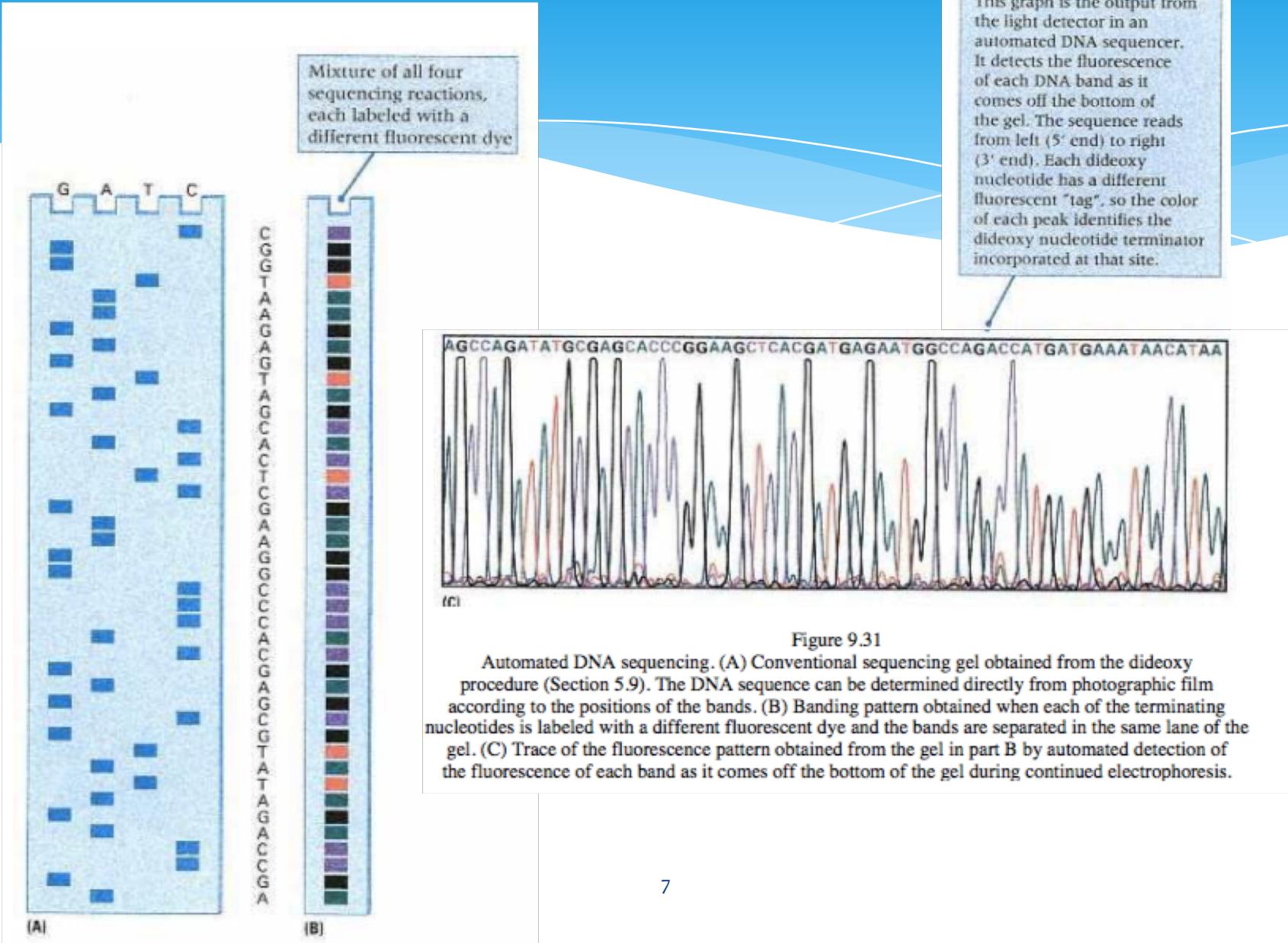


Figure 9.31
Automated DNA sequencing. (A) Conventional sequencing gel obtained from the dideoxy procedure (Section 5.9). The DNA sequence can be determined directly from photographic film according to the positions of the bands. (B) Banding pattern obtained when each of the terminating nucleotides is 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.

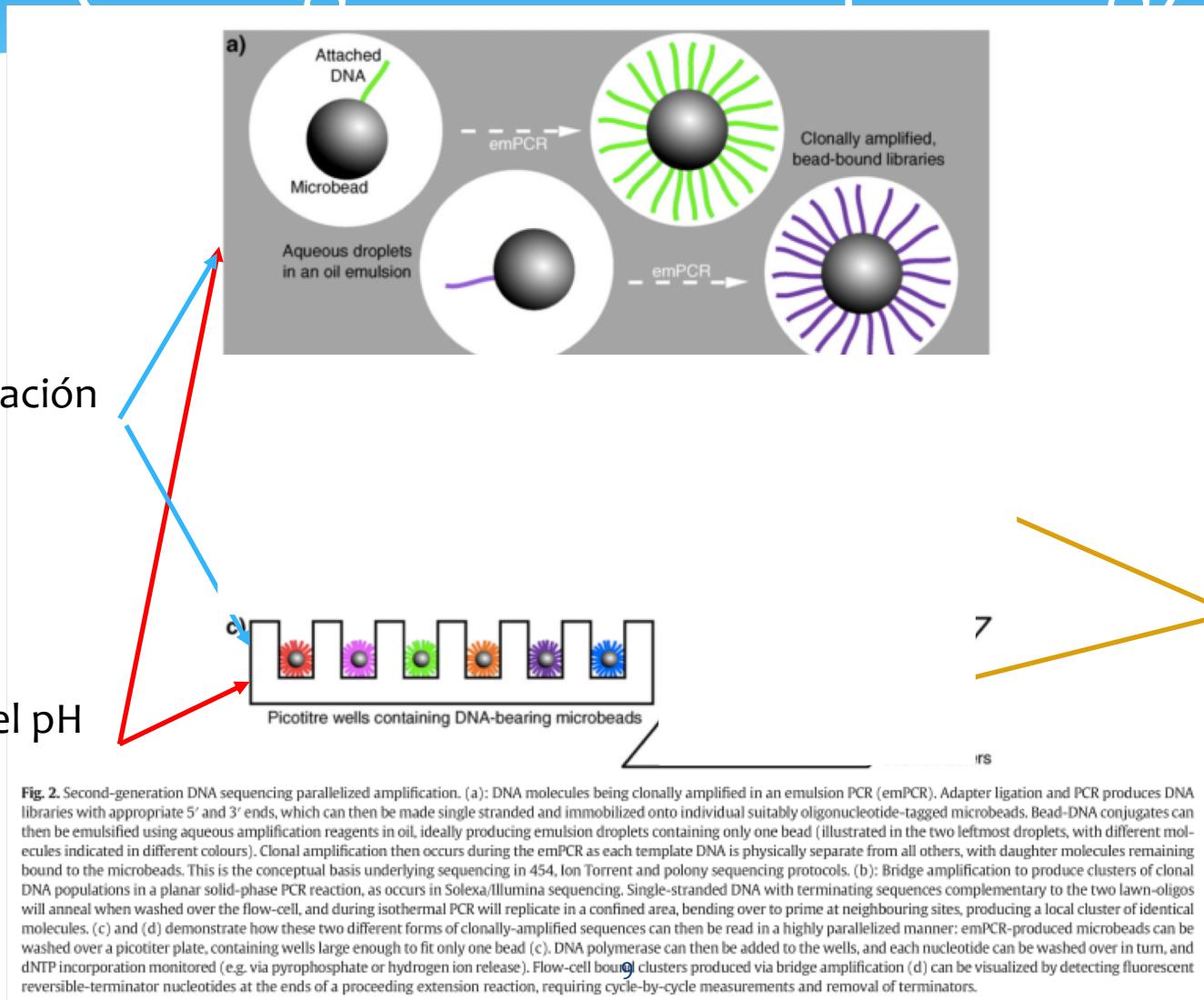
Siguiente generación Herramientas Genómicas

- * Proyecto del genoma Humano (1990-2003) (~3 Gb)
- * Secuenciación Sanger (600-800 pb)
- * Desarrollo de nuevas tecnologías
- * Next generation Sequencing
- * Secuencias de varios millones de fragmentos (lecturas) en una corrida
 - * Microsatélites
 - * SNP
 - * Genomas Mitocondriales



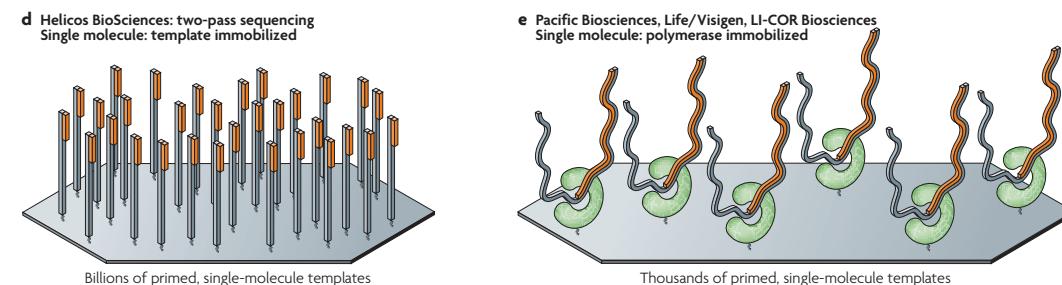
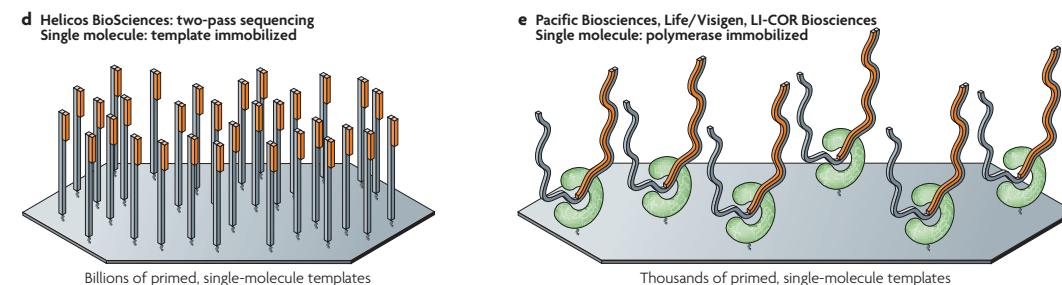
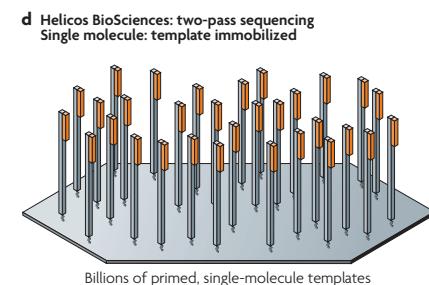
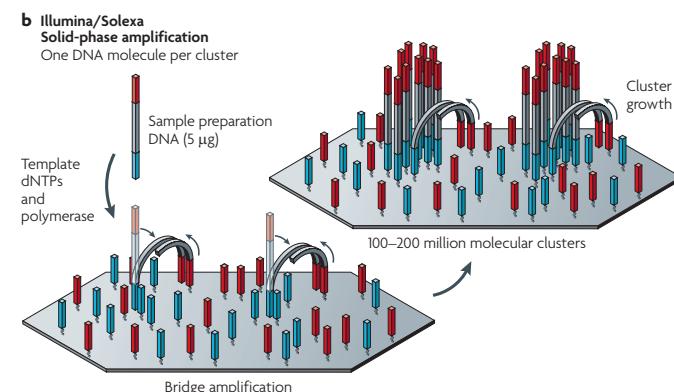
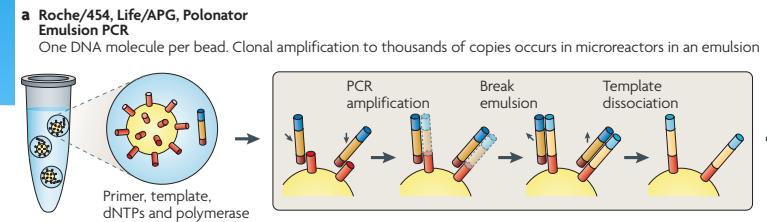
Indispensable para microsatélites y SNPs y análisis de RNA (RNAseq, miRNA...)

Secuenciación de siguiente generación (Next generation sequencing)



Next generation sequencing

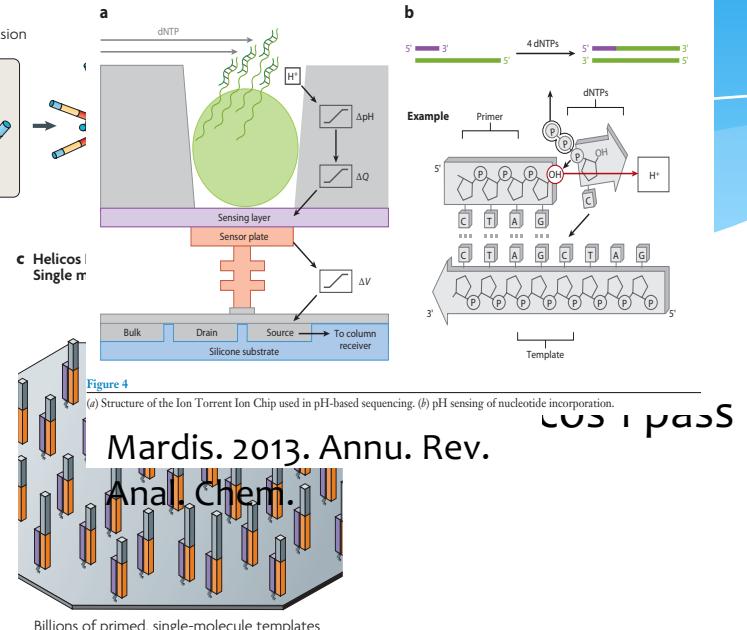
Roche 454



illumina

Helicos 2 pass

Metzker 2010. Nature Reviews.



Mardis. 2013. Annu. Rev.
Anal. Chem.

cos | pasS

Pac Bio

Figure 1 | Template immobilization strategies. In emulsion PCR (emPCR) (a), a reaction mixture consisting of an oil-aqueous emulsion is created to encapsulate bead-DNA complexes into single aqueous droplets. PCR amplification is performed within these droplets to create beads containing several thousand copies of the same template sequence. EmPCR beads can be chemically attached to a glass slide or deposited into PicoTiterPlate wells (FIG. 3c). Solid-phase amplification (b) is composed of two basic steps: initial priming and extending of the single-stranded, single-molecule template, and bridge amplification of the immobilized template with immediately adjacent primers to form clusters. Three approaches are shown for immobilizing single-molecule templates to a solid support: immobilization by a primer (c); immobilization by a template (d); and immobilization of a polymerase (e). dNTP, 2'-deoxyribonucleoside triphosphate.

Tercer generación de secuenciación

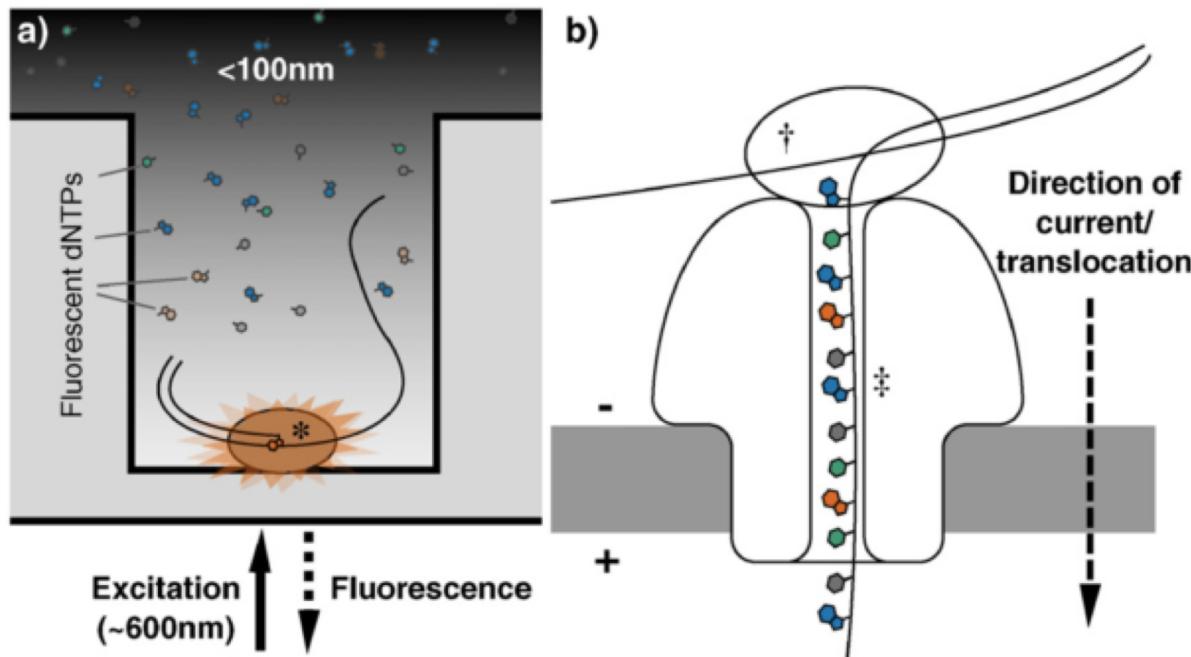


Fig. 3. Third-generation DNA sequencing nucleotide detection. (a): Nucleotide detection in a zero-mode waveguide (ZMW), as featured in PacBio sequencers. DNA polymerase molecules are attached to the bottom of each ZMW (*), and target DNA and fluorescent nucleotides are added. As the diameter is narrower 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.

Comparación

Análisis de secuencias (Sanger)

Procedimiento

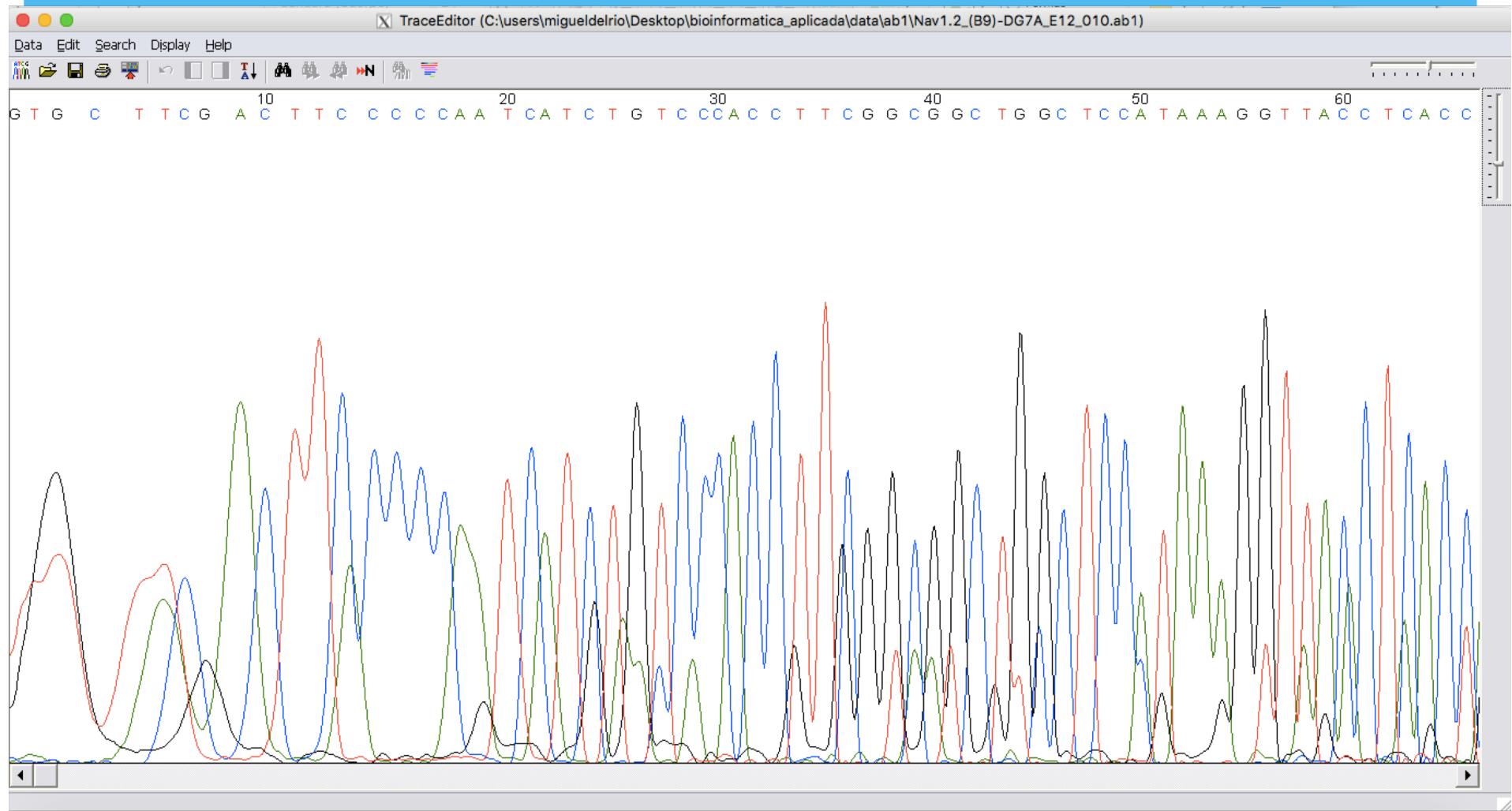
- * Toma de muestra
 - * Tejido: puede variar de una especie a otra (aleta, branquia, músculo, etc.)
 - * Conservación (etanol 95%, congelación)
- * Extracción
 - * Kit (Qiagen, Promega, etc.)
 - * Precipitación con sales (LiCl)
- * Amplificación
 - * Dirigida al fragmento de interés Iniciadores específicos
- * Corroboration de la amplificación
- * Limpieza del fragmento
 - * Kits
- * Envío de la muestra
 - * Langebio
 - * Macrogen
 - * Secuenciador propio

Archivo

- * AB1 formato AB
- * phd.1 formato Phred, (no todos los proveedores lo proporcionan)
- * pdf Electroferograma (cromatograma)
- * Seq, txt texto

 16S_RO_1F-16AR_A07_001.ab1	26/05/2017 11:20	236 KB
 16S_RO_1F-16AR_A07_001.seq	26/05/2017 11:20	611 bytes
 16S_RO_1R-16BR_A08_002.ab1	26/05/2017 11:20	236 KB
 16S_RO_1R-16BR_A08_002.seq	26/05/2017 11:20	613 bytes
 16S_RO_2F-16AR_B07_003.ab1	26/05/2017 11:20	236 KB
 16S_RO_2F-16AR_B07_003.seq	26/05/2017 11:20	615 bytes
 P1AA_16s_16sF.ab1	21/12/2017 22:48	337 KB
 P1AA_16s_16sF.pdf	21/12/2017 22:49	43 KB
 P1AA_16s_16sF.phd.1	21/12/2017 22:48	6 KB
 P1AA_16s_16sF.txt	21/12/2017 22:49	602 bytes
 P1AA_16s_16sR.ab1	21/12/2017 22:48	355 KB
 P1AA_16s_16sR.pdf	21/12/2017 22:48	41 KB

ab1



phd

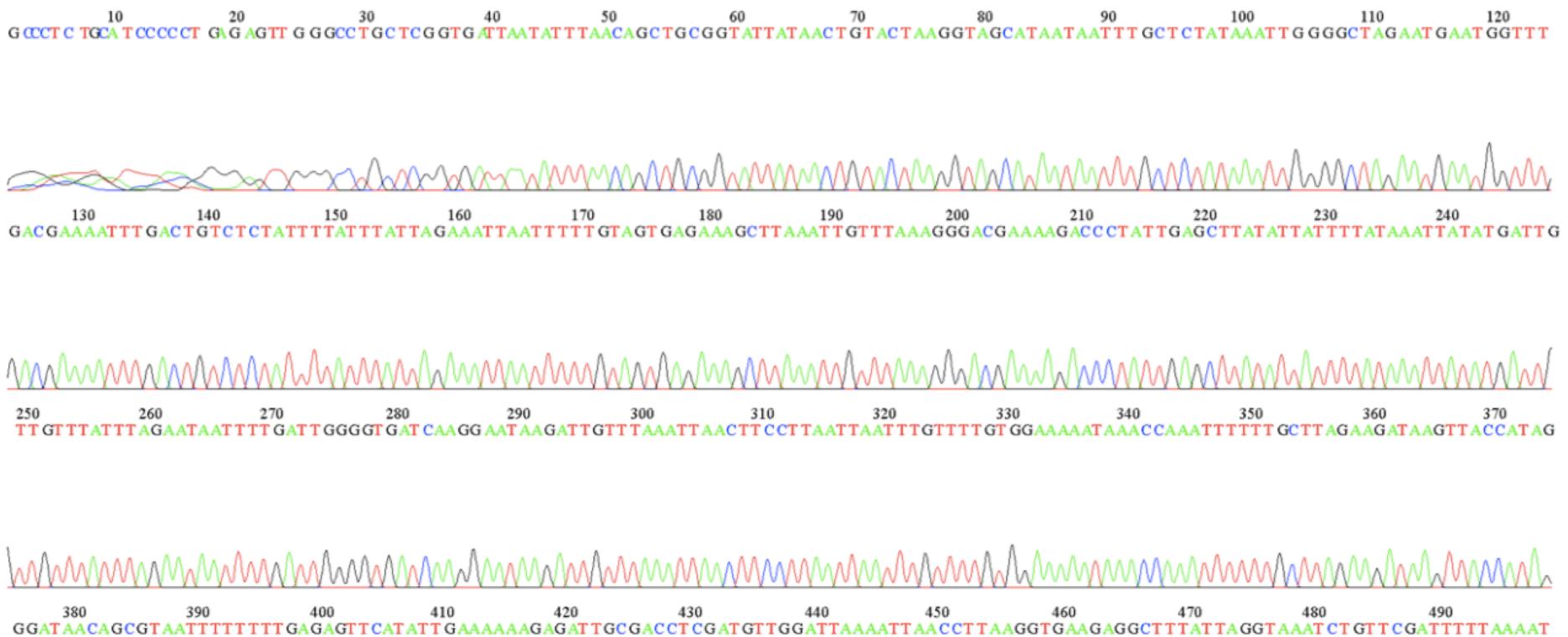
```
BEGIN_SEQUENCE 171221-009_014_P1AA_16s_16sF.phd.1  
  
BEGIN_COMMENT  
  
CHROMAT_FILE: 171221-009_014_P1AA_16s_16sF  
BASECALLER_VERSION: KB 1.4.0  
TRACE_PROCESSOR_VERSION: KB 1.4.0  
QUALITY_LEVELS: 99  
TIME: Thu Dec 21 22:47:32 2017  
TRACE_ARRAY_MIN_INDEX: 0  
TRACE_ARRAY_MAX_INDEX: 6549  
TRIM: -1 -1 -1.000000e+000  
TRACE_PEAK_AREA_RATIO: -1.000000e+000  
CHEM: term  
DYE: big  
  
END_COMMENT  
  
BEGIN_DNA  
G 3 3  
C 2 18  
C 2 25  
C 5 34  
T 3 44  
C 5 57  
T 4 73  
G 9 82  
C 5 91  
A 5 99  
T 6 115  
C 5 127  
C 5 136  
C 5 148  
C 5 158  
C 4 170
```

```
* BEGIN_SEQUENCE 171221-009_014_P1AA_16s_16sF  
  
* BEGIN_COMMENT  
  
* CHROMAT_FILE: 171221-009_014_P1AA_16s_16sF  
* BASECALLER_VERSION: KB 1.4.0  
* TRACE_PROCESSOR_VERSION: KB 1.4.0  
* QUALITY_LEVELS: 99  
* TIME: Thu Dec 21 22:47:32 2017  
* TRACE_ARRAY_MIN_INDEX: 0  
* TRACE_ARRAY_MAX_INDEX: 6549  
* TRIM: -1 -1 -1.000000e+000  
* TRACE_PEAK_AREA_RATIO: -1.000000e+000  
* CHEM: term  
* DYE: big  
  
* END_COMMENT  
  
* BEGIN_DNA  
* G 3 3  
* C 2 18  
* C 2 25  
* C 5 34  
* T 3 44  
* C 5 57  
* T 4 73  
* G 9 82
```

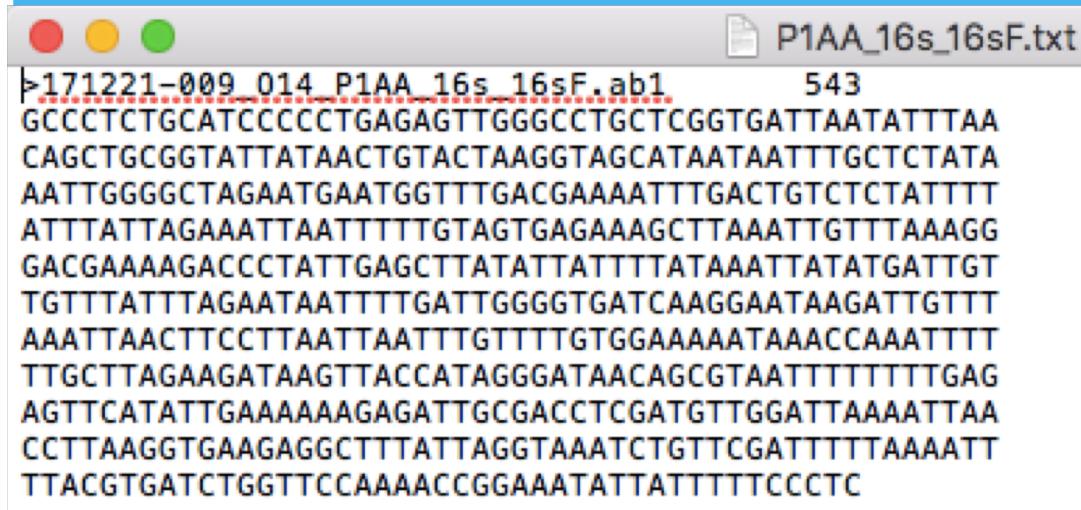
pdf electroferogramma

File: PIAA_16s_16sF.ab1 Run Ended: 2017/12/21 22:25:5 Signal G:1500 A:2400 C:1621 T:3672

Sample: PIAA_16s_16sF Lane: 49 Base spacing: 16.79731 543 bases in 18649 scans Page 1 of 2



Seq, txt



Archivos

Formato:

fasta

>171221-009_O14_P1AA_16s_16sF.ab1 543
GCCCTCTGCATCCCCCTGAGAGTTGGCCTGCTCGGTGATTAAATATTAA
CAGCTCGGGTATTATAACTGTACTAAGGTAGCATAATAATTGCTCTATA
AATTGGGGCTAGAACATGAATGGTTGACGAAAATTGACTGTCTATT
ATTATTAGAAATTAAATTGGTAGTGAGAAAGCTAAATTGTTAAAGG
GACGAAAAGACCCATTGAGCTTATTATTAAATTATGATTGT
TGTTTATTAGAATAATTGATTGGGTGATCAAGGAATAAGATTGTT
AAATTAACTCCCTTAATTAAATTGTTTGTGAAAAATAACCAAATT
TTGCTTAGAAGATAAGTACCATAGGGATAACAGCGTAATTGTTTGTGAG
AGTCATATTGAAAAAAAGAGATTGCGACCTCGATGTTGGATTAAAATTAA
CCTTAAGGTGAAGAGGCTTATTAGGTAAATCTGTTGATTGTTAAAATT
TTACGTGATCTGGTCCAAAACCGGAAATTATTGTTCCCTC

FASTA

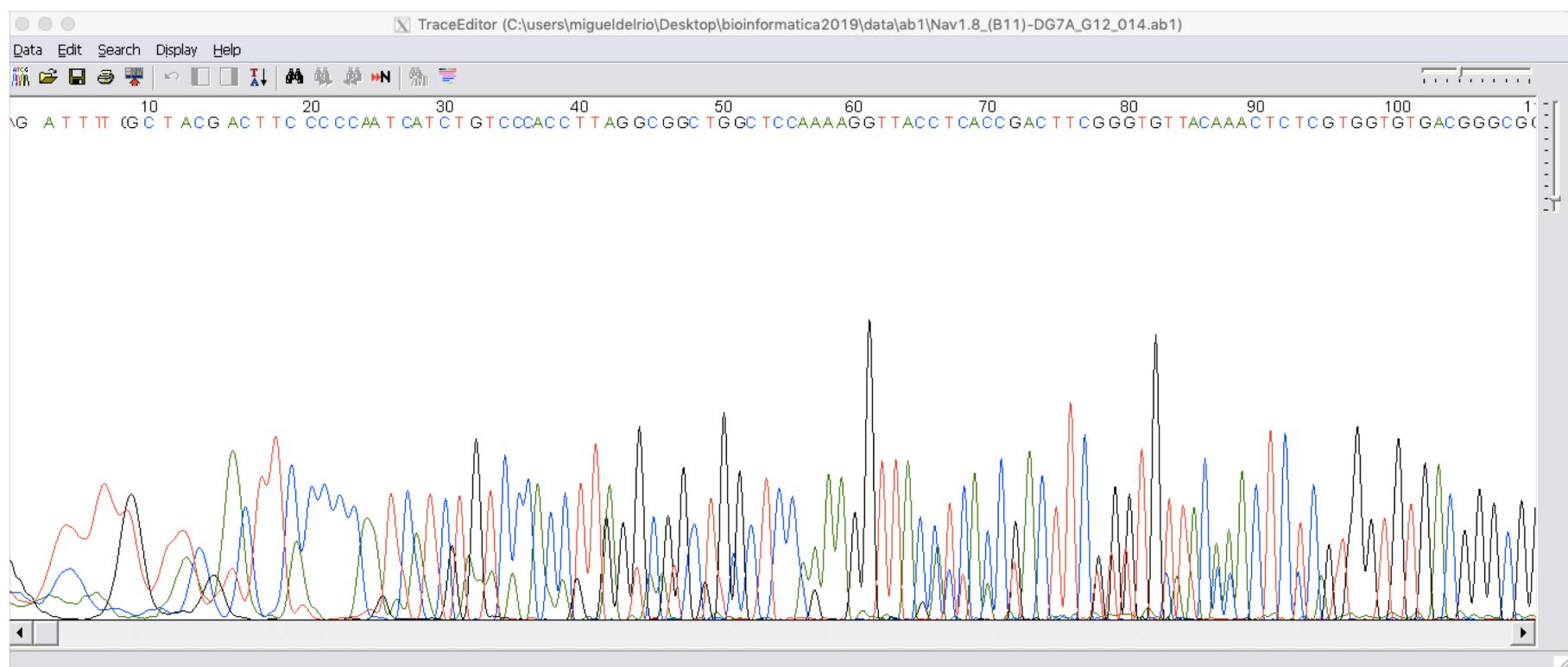
- * Archivo de texto
 - * Secuencias de DNA
 - * Secuencias de amino ácidos (péptidos)
 - * Línea de descripción, que contiene identificadores
 - * Secuencia (una o varias líneas)
- > Secuencia de amino ácidos
- MTEITAAMVKELRESTGAGMMDCKNALSETNGDFDKAVQLLREKGLGK
AAKKADRLAAEGLVSVKVSDDFTIAAMRPSYLS
- > Secuencia de DNA
- ATGACTGAAATTACTGCTGCAATGGTAAAAGAACTCCCGCGAAAGTACAG
GCCCGGGGATGATGGATTGTAAAAATGCTTGAGTGAAACTAATGGAG
ATTTGATAAAGCAGTACAACTTAACAGAGAAAAAGGTTAGGTAAGGC
TGCTAAAAAGCAGATAGACTTGCTGCAGAAGGTTGGTAAGTGTAAA
AGTA

Secuencia de DNA (y RNA)

Código de ácido nucléico	Significado
A	<u>Adenosina</u>
C	<u>Citosina</u>
G	<u>Guanina</u>
T	<u>Timidina</u>
U	<u>Uracilo</u>
R	G A (<u>puRina</u>)
Y	T C (<u>pirimidina/pYrimidine</u>)
K	G T (<u>cetona/Ketone</u>)
M	A C (<u>grupo aMino</u>)
S	G C (interacción fuerte/ Strong interaction)
W	A T (interacción débil/ Weak interaction)
B	G T C (no A) (B viene tras la A)
D	G A T (no C) (D viene tras la C)
H	A C T (no G) (H viene tras la G)
V	G C A (no T, no U) (V viene tras la U)
N	A G C T (cuálquiera/aNy)
X	máscara
-	²⁰ hueco (gap) de longitud indeterminada

ab1

- * Archivo con la información de salida (electroferograma) de la secuenciación Sanger del secuenciador de Applied Biosystems



CSV

- * Archivo de texto, valores separados por comas
(comma separated values)

especie,org,long_total

especie1,1,100

especie1,2,89

especie1,3,114

especie1,4,88

especie1,5,124

especie1,12,101

fasta (.fa, .fas, .fasta, .faa)

GenBank (nucleotido)

* >NM_001273529.2 Drosophila melanogaster vasa (vas), transcript variant B, mRNA

```
CACTAGATTTCGGTACTTTAACAGATCCTTCGGT
TTTGCCTTGCAGAAGTGATCTGAACCTATCA
AAAGTTGTAAGGTAAATACATAAAAGTAAAAAGAATT
ATTGCTCTGAAAGGCAGGCCAAATTAAAAAA
AAAATATCAATATGTCTGACGACTGGATGATGAGCC
CATTGTTGATACTCGCGCGCCCGCGTGGAGA
TTCGACCCATGATGAGCACACGCCAAGACCTTCAG
CGCGAAGCTGAAGGCGATGGTGTGGAGGGAGC
GGTGGTGAAGGCGCGCTACCAAGGAGGAAATCGA
GATGTGTCGGAAGGATCGCGGAGGCAGAGGAG
GAGGAGCTGGAGGTATCGAGGAGGAAATCGAGATG
GAGGGGGCTTCCACGGTGGACGTGGAGGGAGA
AAGGGACTTCCCGGGTGGAGAAGGCGGCTCCGG
TGGACAAGGCGGCTCCCGGGTGGACAAGGCGG
TCCCGGGGGACAAGGCGGCTCCGTGGAGAA
GGCGGCTTCCCGGGTGGCTGTACGAAAACGAGG
```

Swiss-Prot

* >tr|M9PBB5|M9PBB5_DROME Vasa, isoform B
OS=Drosophila melanogaster GN=vas PE=3 SV=1
MSDDWDDEPIVDTRGARGGDWSDDDETAKSFSGEAE
GDGVGGSGEGGGYQGGNRDVFGR
IGGGRGGGAGGYRGGNRDGGGFHGGRREGERDFRGG
EGGFRGGQGGSRGQGGSRGQGG
FRGEGGFRGRLYENEDGDERGRLDREERGGERRG
RLDREERGGERGERGDGGFARRRR
NEDDINNNNNIVEDVERKREFYIPPEPSNDAIEIFSSGIA
SGIHFSKYNNIPVKVTGSDV
PQPIQHFTSADLRDIIDNVNKSGYKIPTPIQKCSIPVISS
GRDLMACAQTSQGKTAFL
LPILSKLLEDPHELELGRPQVIVSPTRELAIQIFNEARK
FAFESYLKIGIVYGCTSFRH
QNECITRGCHVVIATPGRLDFDRTFITFEDTRFVVLD
EADRMLDMGFSEDMRRIIMTHV

fastaq

- * Archivos de texto con información de la secuenciación masiva
 - 1. Encabezado “@”,
 - 2. secuencia,
 - 3. Separador “+”,
 - 4. calidad de cada nucleótido

```
@M00313:49:00000000-A41BL:1:1101:15868:1583 1:N:0:356
GNATGAAATGTTCAAAATGTCTTGGCATGCTGAAGCAGAGTTCTTCATCAGAACT
AACGGGCCAAATCCTGAAAAACAACCCCCACACCATAACCCCCCCCACCACCACCA
AACTTTAAACTTGATACAGTGCAATCAGACAAATACTGTTCTCCTGGCAACTGCCAAA
CCCAGACTCATCCAGACAGAGAAGTGTGACTGGTCACTCCAGAGAATATGTGTAC
+
B#>>AABFFFDFGGGGGGGGHHHAEGHFGGFHHGHHHHHHHHEGHFHHGHG
HHFEEGEEEEFGHFHHFGGHFC?FFF1EE?EF?1CGHGB3?>E/EEEGHGHGFHAG
HHGGHHHHHHFFDGGHHHD<FDGFHCHHGFF1?=FFFFFGECH.<FFFCCH
H.<EGHEGC;GGHECFHGGCG0;00;CE;BFFGGFFFGGGFGBF;CFB0
```

gb

* Archivo con la información del GenBank

LOCUS AB191108 503 bp DNA linear INV 23-MAR-2005
DEFINITION Argonauta argo gene for 16S rRNA, partial sequence.
ACCESSION AB191108
VERSION AB191108.1
KEYWORDS .
SOURCE Argonauta argo
ORGANISM Argonauta argo
Eukaryota; Metazoa; Lophotrochozoa; Mollusca; Cephalopoda;
Coleoidea; Neocoledoidea; Octopodiformes; Octopoda; Incirrata;
Argonautidae; Argonauta.
REFERENCE 1
AUTHORS Takumiya,M., Kobayashi,M., Tsuneki,K. and Furuya,H.
TITLE Phylogenetic Relationships among coleoid cephalopods in Japanese
waters
JOURNAL Unpublished
REFERENCE 2 (bases 1 to 503)
AUTHORS Takumiya,M., Kobayashi,M., Tsuneki,K. and Furuya,H.
TITLE Direct Submission
JOURNAL Submitted (24-SEP-2004) Hidetaka Furuya, Osaka University,
Department of Biology, Graduate School of Science; 1-1,
Machikaneyama, Toyonaka, Osaka 560-0043, Japan
(E-mail:hfuruya@bio.sci.osaka-u.ac.jp, Tel:81-6-6850-6775,
Fax:81-6-6850-5817) 25

out

- * Archivo de texto, usualmente se utiliza como indicador de salida (blast.out)
- * No tiene un formato particular
- * Adquiere el formato del programa que lo genera

tab

- * Archivo separado por tabuladores. Similar a csv, pero en vez de comas se usa el tabulador

Phred score

Es una medida de calidad en la identificación de los nucleótidos generados por algún método de secuenciación.

El Phred score o Q se define como una propiedad que está relacionada logarítmicamente con las probabilidades de error en la identificación de las bases (P).

$$Q = -10 \log_{10} P$$

Phred Quality Score	Probability of Incorrect Base Call	Base Call Accuracy	
10	1 in 10	90%	
20	1 in 100	99%	Sanger
30	1 in 1,000	99.9%	
40	1 in 10,000	99.99%	
50	1 in 100,000	99.999%	

NCBI

NCBI Resources How To

All Databases Search

Welcome to NCBI

The National Center for Biotechnology Information advances science and health by providing access to biomedical and genomic information.

[About the NCBI](#) | [Mission](#) | [Organization](#) | [NCBI News & Blog](#)

Submit
Deposit data or manuscripts into NCBI databases

Download
Transfer NCBI data to your computer

Learn
Find help documents, attend a class or watch a tutorial

Develop
Use NCBI APIs and code libraries to build applications

Analyze
Identify an NCBI tool for your data analysis task

Research
Explore NCBI research and collaborative projects

Popular Resources

- PubMed
- Bookshelf
- PubMed Central
- PubMed Health
- BLAST
- Nucleotide
- Genome
- SNP
- Gene
- Protein
- PubChem

NCBI News & Blog

ClinVar Unveils New, More Intuitive Variation Display
18 Dec 2017

ClinVar, NCBI's database of clinically relevant genetic variations with

NIH Data Hackathon on campus – January 22-24, 2018
18 Dec 2017

From January 22-24, 2018, the NCBI will host a data science hackathon on

January 10 NCBI Minute: QuickBLASTP — a program for rapidly finding high-scoring protein matches in large databases

NCBI Home

Resource List (A-Z)

- All Resources
- Chemicals & Bioassays
- Data & Software
- DNA & RNA
- Domains & Structures
- Genes & Expression
- Genetics & Medicine
- Genomes & Maps
- Homology
- Literature
- Proteins
- Sequence Analysis
- Taxonomy
- Training & Tutorials
- Variation

Uniprot

The mission of UniProt is to provide the scientific community with a comprehensive, high-quality and freely accessible resource of protein sequence and functional information.

UniProtKB
UniProt Knowledgebase

Swiss-Prot (556,388)
Manually annotated and reviewed.

TrEMBL
(102,248,261)
Automatically annotated and not reviewed.

UniRef
Sequence clusters

UniParc
Sequence archive

Proteomes

Supporting data

Literature citations

Cross-ref. databases

Taxonomy

Diseases

Subcellular locations

Keywords

News

BLOG Twitter Facebook RSS

Forthcoming changes
Planned changes for UniProt

UniProt release 2017_12
Swiss-Prot in the sky with psilocybin: the biosynthesis pathway of a psychedelic drug unveiled

UniProt release 2017_11
Sex determination in insects: 50 ways to achieve sex-specific splicing

News archive

Getting started

Text search
Our basic text search allows you to search all the

UniProt data

Download latest release
Get the UniProt data

Protein spotlight

When The Mind Bends

<http://www.uniprot.org>

ENA

EMBL-EBI

ENASearch Examples: BN000065, histone Advanced Sequence

Services Research Training About us

Home Search & Browse Submit & Update Software About ENA Support

European Nucleotide Archive

The European Nucleotide Archive (ENA) provides a comprehensive record of the world's nucleotide sequencing information, covering raw sequencing data, sequence assembly information and functional annotation. [More about ENA](#)

Access to ENA data is provided through the browser, through search tools, large scale file download and through the API.

Text Search

Examples: BN000065, histone

Search Advanced search

Sequence Search

Enter or paste a nucleotide sequence or accession number

Popular

- Submit and update
- Sequence submissions
- Genome assembly submissions
- Submitting environmental sequences
- Citing ENA data
- Rest URLs for data retrieval
- Rest URLs to search ENA

Latest ENA news

21 Dec 2017: [ENA services over the holiday period](#)

Between Friday 22nd December and Tuesday 2nd January ENA services such as submissions and retrieval...

21 Dec 2017: [ENA release 134 expected early January](#)

The last release of assembled and annotated sequences for 2017 (134) has been particularly...

<https://www.ebi.ac.uk/ena>

DDBJ, DNA Database of Japan

 **DDBJ** DNA Data Bank of Japan 

Japanese

Google カスタム検索

About DDBJ How to Use Report/Statistics FAQ Contact Us

RSS DDBJ Twitter Mail Magazine

DDBJ Service

Data Submission Search / Analysis Super Computer ftp.ddbj.nig.ac.jp

Hot Topics News Archive

News Release PR Maintenance Operation All

2018.01.04 Release of genome data of California poppy (*Eschscholzia californica* subsp. *californica*)
2017.12.25 Release of genome data of red seabream (*Pagrus major*)
2017.12.22 D-way submission services will be unavailable soon (Dec. 22 19:00) (Only BioProject/BioSample submission services are resumed on Dec. 25)
2017.12.20 DDBJ Rel. 111.0, DAD Rel. 81.0 Completed
2017.12.18 Release of TSA data of brown planthopper (*Nilaparvata lugens*)

 
International Nucleotide Sequence Database Collaboration

 **National Institute of Genetics**
Research Organization of Information and Systems

 大学共同利用機関法人 情報・システム研究機構
Research Organization of Information and Systems

 **JBI portal**
Japan alliance for Bioscience Information

<http://www.ddbj.nig.ac.jp>

NCBI

NCBI Resources How To

All Databases Search

Welcome to NCBI

The National Center for Biotechnology Information advances science and health by providing access to biomedical and genomic information.

[About the NCBI](#) | [Mission](#) | [Organization](#) | [NCBI News & Blog](#)

Submit Deposit data or manuscripts into NCBI databases 

Download Transfer NCBI data to your computer 

Learn Find help documents, attend a class or watch a tutorial 

Develop Use NCBI APIs and code libraries to build applications 

Analyze Identify an NCBI tool for your data analysis task 

Research Explore NCBI research and collaborative projects 

Popular Resources

PubMed
Bookshelf
PubMed Central
PubMed Health
BLAST
Nucleotide
Genome
SNP
Gene
Protein
PubChem

NCBI News & Blog

ClinVar Unveils New, More Intuitive Variation Display 18 Dec 2017
ClinVar, NCBI's database of clinically relevant genetic variations with 
NIH Data Hackathon on campus – January 22-24, 2018 18 Dec 2017
From January 22-24, 2018, the NCBI will host a data science hackathon on 
January 10 NCBI Minute: QuickBLASTP — a program for rapidly finding high-scoring protein matches in large databases 

Bases de datos

Nombre	Dirección	Datos
NCBI	https://www.ncbi.nlm.nih.gov	
GenBank	https://www.ncbi.nlm.nih.gov/genbank/	Base de datos
Nucleotide	https://www.ncbi.nlm.nih.gov/nucleotide/	nucleótidos
PubMed	https://www.ncbi.nlm.nih.gov/pubmed	Referencias bibliográficas
NR	https://www.ncbi.nlm.nih.gov/protein/	Secuencia de proteínas no redundante
Swiss-Prot, UniProtKB	http://www.uniprot.org	Universal Protein Resource
KEGG (Kyoto Encyclopedia of Genes and Genomes)	http://www.genome.jp/kegg/	Rutas metabólicas

Identificadores en diferentes bases de datos



Database name	Identifier syntax
GenBank	accession locus (gb accession locus)
EMBL Data Library	emb accession locus
DDBJ, DNA Database of Japan	dbj accession locus
NBRF PIR	pir entry
Protein Research Foundation	prf name
SWISS-PROT	sp accession entry name
Brookhaven Protein Data Bank	pdb entry chain
Patents	pat country number
GenInfo Backbone Id	bbs number
General database identi era	gnl database identi er
NCBI Reference Sequence	ref accession locus
Local Sequence identi er	lcl identifier



McEntyre, Jo, Ostell. 2012. The NCBI Handbook³⁵

https://www.ncbi.nlm.nih.gov/books/NBK21101/pdf/Bookshelf_NBK21101.pdf

Datos en servidores (nube)

- * Servicios de almacenamiento
- * Servicios de renta de equipo
- * Servicio de procesamiento



Recomendaciones al usar servidores

1. Los datos nunca están seguros en la red
2. Anotar el servidor, la base de datos y la versión del programa utilizado
3. Anotar los números de identificación de las secuencias
4. Anotar los parámetros del programa
5. Guarde los resultados de internet inmediatamente
6. Utilice los valores de E
7. Asegúrese de que puede confiar en los alineamientos
8. Utilice diferentes programas para verificar los resultados extremos
9. No utilice métodos que no se han publicado
10. Las bases de datos no son como los vinos (actualícela)
11. No confíe ciegamente en los programas gratuitos
12. Utilice las herramientas en el momento adecuado (Bitting the bullet at the right time)

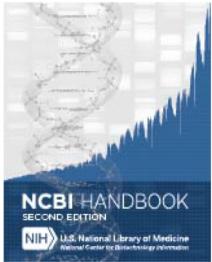
Manuales

NCBI Resources How To

mdeorio My NCBI Sign Out

Bookshelf Books Search

Browse Titles Advanced Help

 The NCBI Handbook, 2nd edition

Bethesda (MD): National Center for Biotechnology Information (US); 2013-.
[Copyright and Permissions](#)

Search this book

< Prev Next > [f](#) [t](#) [g+](#)

Views

- PubReader
- Print View
- Cite this Page
- PDF version of this title (14M)
- Disable Glossary Links

Related information

- NLM Catalog

Recent Activity

Turn Off Clear

- The NCBI Handbook
- Drosophila melanogaster vasa (vas), transcript variant B, mRNA Nucleotide
- Vasa, isoform B Drosophila melanogaster AND (animals[filter]) (12) Nucleotide
- Haliotis diversicolor diversicolor 16S ribosomal RNA gene, partial sequence; mit... Nucleotide
- 16s haliotis AND (animals[filter]) (284) Nucleotide

<https://www.ncbi.nlm.nih.gov/books/NBK143764/>

Basic Local Alignment Search Tool (BLAST)

The screenshot shows the NCBI BLAST search interface. At the top, there's a navigation bar with links like 'Correos-IAS', 'Apple', 'iCloud', 'Facebook', 'Google Maps', 'Noticias', 'Populares', 'PersonalCicese', 'NCBI', 'Cicese', 'Save to Mendeley', 'How to ThinkLCSc', 'News', 'Galaxy', 'UABC', and 'European Nucleotide Archive < EMBL-EBI'. Below the bar, the URL 'blast.ncbi.nlm.nih.gov' is visible. The main header includes the NIH logo, 'U.S. National Library of Medicine', 'NCBI National Center for Biotechnology Information', and user links for 'mdelrio', 'My NCBI', and 'Sign Out'. A 'BLAST®' logo is on the left, and a navigation menu on the right offers 'Home', 'Recent Results', 'Saved Strategies', and 'Help'. A 'NEWS' sidebar highlights 'IgBLAST 1.8.0 released' with a link to 'More BLAST news...'. The main content area features three large buttons: 'Nucleotide BLAST' (nucleotide → nucleotide), 'blastx' (translated nucleotide → protein), and 'tblastn' (protein → translated nucleotide). To the right is a 'Protein BLAST' button (protein → protein). Below these are sections for 'BLAST Genomes' with a search bar and dropdowns for 'Human', 'Mouse', 'Rat', and 'Microbes', and a 'Search' button.

Github

The screenshot shows the GitHub homepage. At the top, there is a navigation bar with links for "Pull requests", "Issues", "Marketplace", and "Explore". Below the navigation bar, a prominent green banner reads "Learn Git and GitHub without any code!". It includes a brief description: "Using the Hello World guide, you'll create a repository, start a branch, write comments, and open a pull request." Two buttons are present: a green "Read the guide" button and a white "Start a project" button. The main content area displays a timeline of activity from the user "sr320". The activity items are:

- sr320 created a repository [sr320/talk-bh-2017](#) 28 days ago
- sr320 forked [sr320/semester-biology](#) from [datacarpentry/semester-biology](#) on 22 Nov 2017
- sr320 created a repository [sr320/course-fish507-2018](#) on 14 Nov 2017
- sr320 created a repository [RobertsLab/project-crab](#) on 8 Nov 2017

To the right of the activity feed, there is a sidebar with a "Game Off 2017 winners" section featuring a blue icon with an 'A' and a link to "View 37 new broadcasts". Another sidebar lists "Repositories you contribute to" with two entries: "sr320/paper-pano-go" (2 stars) and "ContinuumIO/anaconda-issues" (191 stars). At the bottom right, there is a "Your repositories" section with a count of 32, a "New repository" button, and a search bar with the placeholder "Find a repository...".

Github

- * Abrir una cuenta
- * Buscar mdelrio1
- * En usuarios, seguirme

The screenshot shows the GitHub user profile for the account `mdelrio1`. The top navigation bar includes links for Pull requests, Issues, Marketplace, and Explore. The main content area displays a summary of repository statistics: 1 Repository, 8 Code, 13 Commits, 1 Issue, 1 Wiki, and 1 User. Below this, a single repository result is shown: `mdelrio1/mdelrio1github`, which was last updated on July 18, 2017.

Repositories	1
Code	8
Commits	
Issues	13
Wikis	1
Users	1

[Advanced search](#) [Cheat sheet](#)

1 repository result

[mdelrio1/mdelrio1github](#)

Updated on 18 Jul 2017

mdelrio1 / bioinformatica_aplicada

Code Issues Pull requests Projects Wiki Insights Settings

archivos del curso bioinformatica aplicada UAM-I

Add topics

4 commits 1 branch 0 releases 1 contributor

Branch: master New pull request Create new file Upload files Find file Clone or download

mdelrio1 actualizar Latest commit 0953a0e 6 minutes ago

.gitattributes Initial commit 24 days ago

00temario.ipynb temario 24 days ago

GTGCTTCGACTTCCCCAATCATCTGTCACCTTCGGCGCTGGCTCCATAAGGTTACCTCACCGACTTCGGGTGTTACAAACTCTCGTGGTGACGGGCGGTGTACAAGGCCGGAACGTATTCAACCGGGCATGCTGATCCGAATTAC
AACCGATTCCAGCTCACCATTCAACTGCAAACCTGAACTGAAAACAGATTCTGGAATTGGCTTAACCTCCGGTTCCCTGCCCTTCTCTGTCATTGTACACCTGTAACCCAGGTATAAGGGCATGATTTGACCTC
ATCCCCATTCCCTCCAGTAA

Copiar y pegar la secuencia en el blastn

Code Issues Pull requests Projects Wiki Insights Settings

Branch: master bioinformatica_aplicada / ejercicio1.txt Find file Copy path

mdelrio ejer1 f5f93b7 10 minutes ago

0 contributors

Executable File | 2 lines (1 sloc) | 335 Bytes

Raw Blame History

1 GTGCTTCGACTTCCCCAATCATCTGTCACCTTCGGCGCTGGCTCCATAAGGTTACCTCACCGACTTCGGGTGTTACAAACTCTCGTGGTGACGGGCGGTGTACAAGGCCGGAACGT

Blastn

NIH U.S. National Library of Medicine NCBI National Center for Biotechnology Information mdelrio My NCBI Sign Out

BLAST® » blastn suite Standard Nucleotide BLAST

blast blastp blastx tblastn tblastx Enter Query Sequence Enter accession number(s), gi(s), or FASTA sequence(s) CACCGACTTCGGGTGTTACAACTCTGGTGTGACGGGCGGTGTACAAGGCCAAGTATTACCGCGCATGCTGATCCGAATTACAACCGATTCCACGCTTCACGCATTGCAAACCTGCAATCCGAACTGAAAACAGATTGTGGAATTGGCTTAACCTCCCCGGTTCCCTTGTCTGCCATTGTACCACTGTGTACCCCAGGTATAAGGGGCATGATGATGTCATCCCCATCCCTCCAGTAA

Or, upload file Seleccionar archivo ning n archivo seleccionado Job Title Enter a descriptive title for your BLAST search Align two or more sequences

Choose Search Set Database Human genomic + transcript Mouse genomic + transcript Nucleotide collection (nr/nt) Organism Optional Enter organism name or id--completions will be suggested Exclude Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown. Models (XM/XP) Uncultured/environmental sample sequences Sequences from type material Entrez Query Optional Enter an Entrez query to limit search YouTube Create

Genomic plus Transcript Human genomic plus transcript (Human G+T) Mouse genomic plus transcript (Mouse G+T)

Other Databases Nucleotide collection (nr/nt) 16S ribosomal RNA sequences (Bacteria and Archaea) Reference RNA sequences (refseq_rna) RefSeq Representative genomes (refseqRepresentative_genomes) RefSeq Genome Database (refseq_genomes) Whole-genome shotgun contigs (wgs) Expressed sequence tags (est) Sequence Read Archive (SRA) Transcriptome Shotgun Assembly (TSA) High throughput genomic sequences (HTGS) Patent sequences (pat) Protein Data Bank (pdb) Reference genomic sequences (refseq_genomic) Human RefSeqGene sequences (RefSeq_Gene) Genomic survey sequences (gss) Sequence tagged sites (dbsts)

Ejercicio 1

- * Someta la secuencia al GenBank

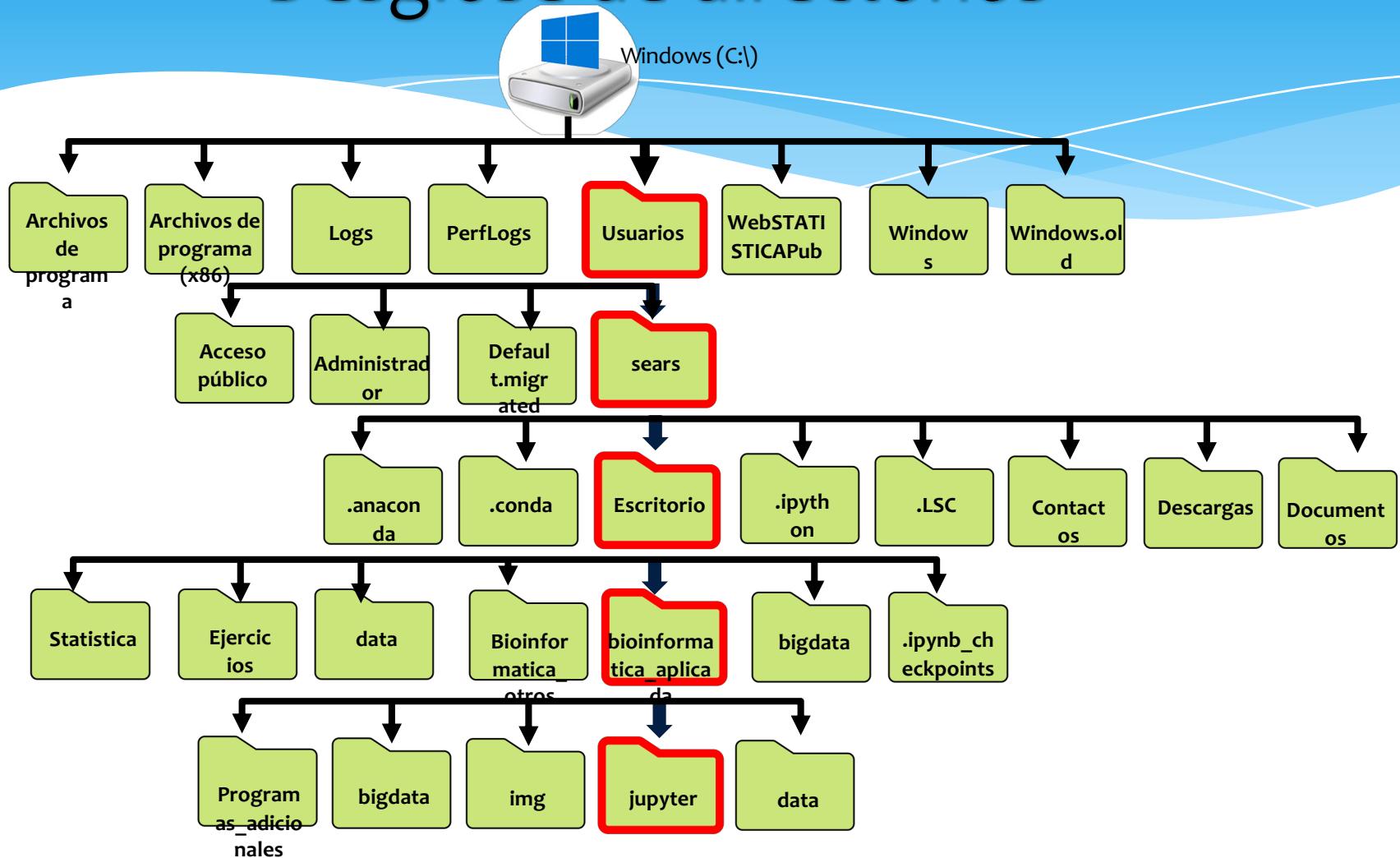
CTTCCCCAATCATCTGTCCCACCTCGCCGGCTGGCTCCATAAA
GGTTACCTCACCGACTTCGGGTGTTACAAACTCTCGTGGTGTGAC
GGGCGGTGTGTACAAGGCCGGAACGTATTCACCGCGGCATGC
TGATCCGCAATTACAACCGATTCCAGCTTCACGCATTCAAGTTGC
AAAATGCAATCCGAACTGAAAACAGATTGTGGAATTGGCTAA
CCTCCCGGTTCCCTGCCCTTGTTCATTGTACACACGTGTG
TACCCCAGGTATAAGGGGCATGATGATTGACGTCATC

- * Utilizando las bases de datos
 - * Nt
 - * 16S microbial
- * Con sus compañeros de mesa discuta las diferencias entre la información proporcionada en cada búsqueda

Cuestionario

- * Defina los siguientes términos:
 - * Alineamiento
 - * Marcador/Calificación máxima (max score)
 - * Marcador total (total score)
 - * Cobertura de búsqueda
 - * Valor de E (E value)
 - * Identidad (ident)
 - * Acceso
- * ¿En qué formatos de archivo se puede descargar la información?

Desglose de directorios

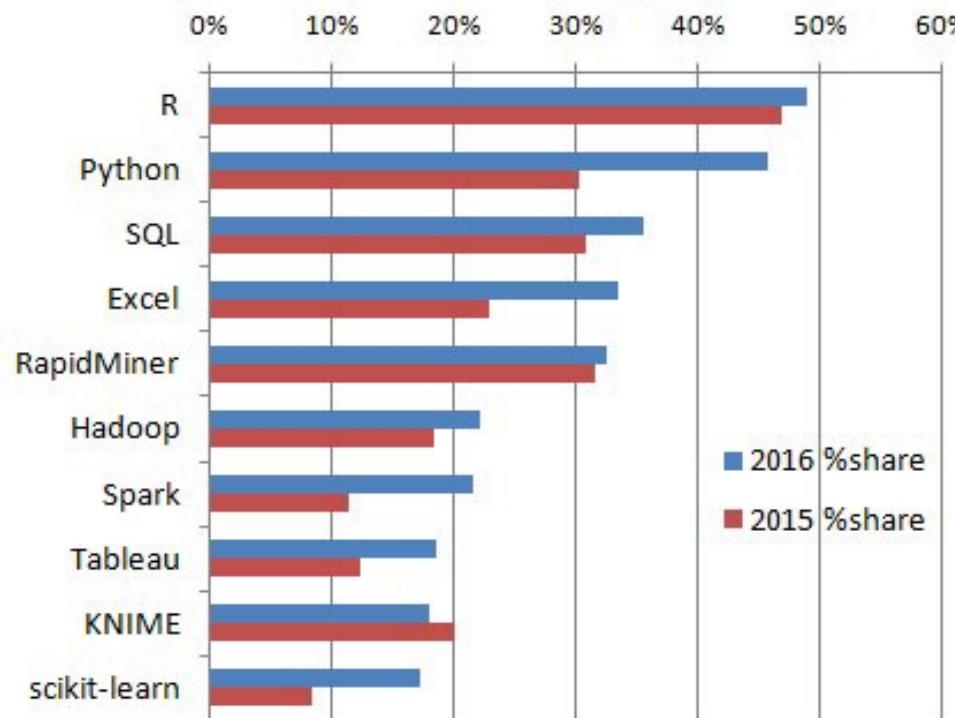


Lenguajes

- * Python
- * R
- * Perl
- * Ruby
- * Java
- * Linux
- * C++

Python

KDnuggets Analytics/Data Science 2016 Software Poll, top 10 tools



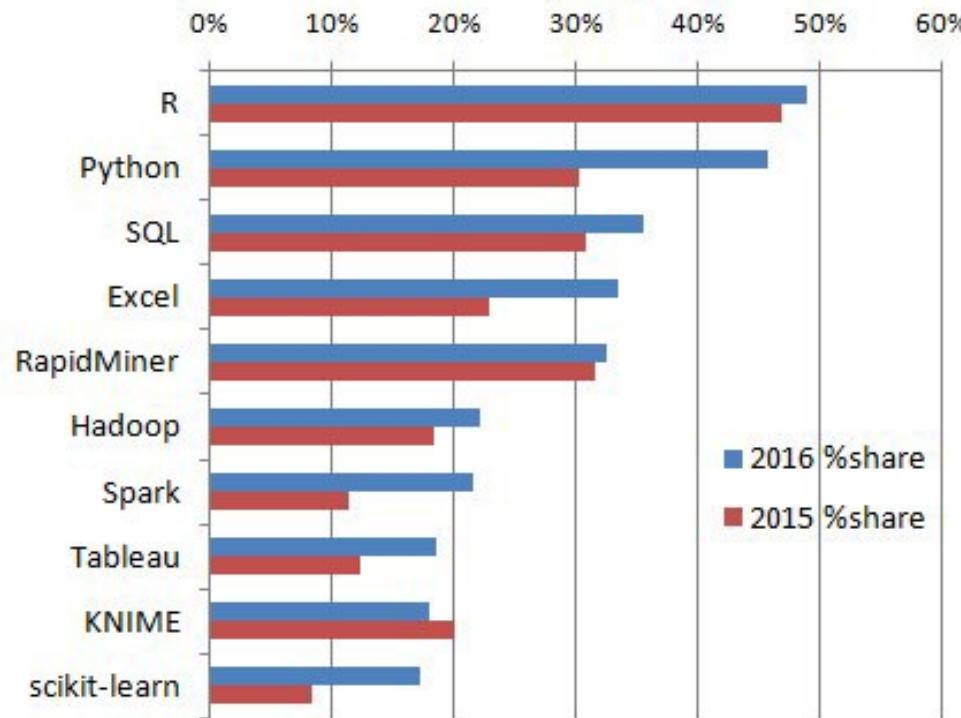
Python es el lenguaje más popular utilizado por “científicos de datos”
Python está ganando popularidad en apreciación general y es el lenguaje más práctico para construir productos
Python es una herramienta poderosa para procesamiento de datos a mediana escala.
Python también tiene la ventaja de una comunidad rica y que ofrece una gran cantidad de herramientas y kits
Python no es el lenguaje con mayor

- * Python
 - * python.org
 - * <http://www.python-course.eu/index.php>
- * Google
 - * <https://developers.google.com/edu/python/>

- * Python tiene una comunidad sana, activa y de gran soporte
- * Python tiene corporaciones que dan gran soporte
- * Python tiene “Big Data”
- * Python tiene un gran número y sorprendente de biblioteca (libraries).
- * Python es confiable y eficiente
- * Python es accesible

R

KDnuggets Analytics/Data Science 2016 Software Poll, top 10 tools



R ha estado desde 1997 como una alternativa gratuita de programas estadísticos caros, como Matlab o SAS.

En los últimos años, R se ha vuelto el niño dorado de la ciencias de datos (Data Science)

Lo usan Google, Facebook, Bank of America, y the New York Times. Sus utilerías comerciales se han dispersado enormemente.

Otros lenguajes

- * Perl
- * Ruby

Perl vs. Python

Use Python

- * Cuando necesite usar el código más de una vez
- * Siempre que haya una posibilidad remota que alguien más (i.e. colegas) utilicen su código
- * Siempre que necesite usar funciones u objetos. La mayoría de los códigos de Perl no necesariamente incluyen las funciones, porque son más difíciles de escribir que en Python. Por ejemplo en Perl, necesita aprender cómo pasar de referencias a variables, etc. y ello conduce a códigos que son más difíciles de entender (en Perl).
- * Si es el primer lenguaje de programación que está aprendiendo. Se sugiere que inicie con Python,
- * Python es mucho más limpio que Perl
- * Python está diseñado para respetar buenas prácticas que cualquier programador debería conocer (ver The Zen of Python, <https://www.python.org/dev/peps/pep-0020/>)
- * Python tiene soporte para estructura tabular de datos (Data Frames), siempre que tenga trabajos con tablas o análisis de datos use Python (Pandas) o R
- * Python tiene soporte para “machine learning algorithms”.

Python

- * “Es uno de los lenguajes de programación más populares hoy en día a pesar de ser un idioma relativamente viejo. Fue creado a finales de los ochentas por Guido van Rossum --que dentro de la comunidad de Python es conocido como el Benevolente Dictador Vitalicio-- y su nombre está inspirado en el grupo de comedia británico Monthly Python.”

<https://hipertextual.com/2011/02/zen-python>

Python vs Ruby vs Perl

- * Perl for one-liners in a pipe. Python for everything else
- * Python vs Ruby vs Perl
- * Right folks! Time for another religious war. Or, a Halo style three way battle.
- * On the left, we have Perl, the grand-daddy of the scripting world. But don't let his age fool you, he still has a few tricks up his sleeve.
- * On the right we have Ruby, the young kid on the block. He may be the youngest, but he sure knows how to throw his weight around.
- * In the middle we have Python. Entering middle age, his pot belly is beginning to show, but get too close, and he'll knock you out in one.
- * Come ladies and gentlemen, place bets. Who will win?

El zen de Python:

Hermoso es mejor que feo.
Explícito es mejor que implícito.
Simple es mejor que complejo.
Complejo es mejor que complicado.
Sencillo es mejor que anidado.
Escaso es mejor que denso.
La legibilidad cuenta.
Los casos especiales no son lo suficientemente especiales para romper las reglas.
Lo práctico le gana a la pureza.
Los errores no debe pasar en silencio.
A menos que sean silenciados.
En cara a la ambigüedad, rechazar la tentación de adivinar.
Debe haber una - y preferiblemente sólo una - manera obvia de hacerlo.
Aunque esa manera puede no ser obvia en un primer momento a menos que seas holandés.
Ahora es mejor que nunca.
Aunque "nunca" es a menudo mejor que "ahora mismo".
Si la aplicación es difícil de explicar, es una mala idea.
Si la aplicación es fácil de explicar, puede ser una buena idea.
Los espacios de nombres son una gran idea ¡hay que hacer más de eso!

Hermoso es mejor que feo.
Explícito es mejor que implícito.
Simple es mejor que complejo.
Complejo es mejor que complicado.
Plano es mejor que anidado.
Escaso es mejor que denso.
Cuenta la legibilidad.
Los casos especiales no son lo suficientemente especial como para romper las reglas.
Aunque sentido práctico supera pureza.
Los errores nunca debe pasar en silencio.
A menos que explícitamente silenciados.
Ante la ambigüedad, rechaza la tentación de adivinar.
Debería haber una - y preferiblemente sólo una - manera obvia de hacerlo.
Aunque esa manera puede no ser obvia al principio a menos que seas holandés.
Ahora es mejor que nunca.
Aunque nunca es a menudo mejor que la * justo * ahora.
Si la implementación es difícil de explicar, es una mala idea.
Si la implementación es fácil de explicar, puede ser una buena idea.
Namespaces son una gran idea de fanfarria - Vamos a hacer más de esos!

Jupyter

Jupyter Documentation 4.1

« Try Jupyter Optional: Installing Kernels »

Search



Jupyter Notebook Quickstart

- Try Jupyter
- Installing Jupyter Notebook
- Optional:* Installing Kernels
- Running the Notebook
- Migrating from IPython Notebook

Architecture Guides

Narratives and Use Cases

IPython

Installation, Configuration, and Usage

Community Guides

Contributor Guides

Release Notes

Reference

Installing Jupyter Notebook

Contents

- Prerequisite: Python
- Installing Jupyter using Anaconda and conda
- *Alternative for experienced Python users:* Installing Jupyter with pip

This information explains how to install the Jupyter Notebook and the IPython kernel.

Prerequisite: Python

While Jupyter runs code in many programming languages, **Python** is a requirement (Python 3.3 or greater, or Python 2.7) for installing the Jupyter Notebook.

We recommend using the [Anaconda](#) distribution to install Python and Jupyter. We'll go through its installation in the next section.

Installing Jupyter using Anaconda and conda

For new users, we **highly recommend** installing Anaconda. Anaconda conveniently installs Python, the Jupyter Notebook, and other commonly used packages for scientific computing and data science.

Use the following installation steps:

1. Download [Anaconda](#). We recommend downloading Anaconda's latest Python 3 version (currently Python 3.5).
2. Install the version of Anaconda which you downloaded, following the instructions on the download page.
3. Congratulations, you have installed Jupyter Notebook. To run the notebook:

```
jupyter notebook
```

See [Running the Notebook](#) for more details.

<https://jupyter.readthedocs.io/en/latest/install.html#install>

Anaconda

 **ANACONDA**
Powered by Continuum Analytics

Anaconda Cloud Documentation Blog Contact  [DOWNLOAD](#)

What Is Anaconda? Products Support & Solutions Community About Resources

DOWNLOAD ANACONDA DISTRIBUTION

Version: 4.4.0 | Release Date: May 31, 2017

Download for:   

High-Performance Distribution Easily install 1,000+ data science packages	Package Management Manage packages, dependencies and environments with conda	Portal to Data Science Uncover insights in your data and create interactive visualizations
---	--	--

Download for Your Preferred Platform

 Windows |  macOS |  Linux
59

<https://www.continuum.io/downloads>