

**UNIVERSIDAD AUTÓNOMA METROPOLITANA**  
Unidad Iztapalapa

**Casa abierta al tiempo**

# 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

3

## 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 SR and Gracilibacteria Code
26. Pachysolen tannophilus Nuclear Code
27. Karyorelict Nuclear
28. Condylostoma Nuclear
29. Mesodinium Nuclear
30. Peritrich Nuclear
31. Blastocritidida 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	Ile	ACT	T	Thr	AAT	N	Asn	AGT	S	Ser
ATC	I	Ile	ACC	T	Thr	AAC	N	Asn	AGC	S	Ser
ATA	I	Ile	ACA	T	Thr	AAA	K	Lys	AGA	R	Arg
<b>ATG</b>	<b>M</b>	<b>Met i</b>	<b>ACG</b>	<b>T</b>	<b>Thr</b>	<b>AAG</b>	<b>K</b>	<b>Lys</b>	<b>AGG</b>	<b>R</b>	<b>Arg</b>
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

5

## Secuenciación (Sanger)

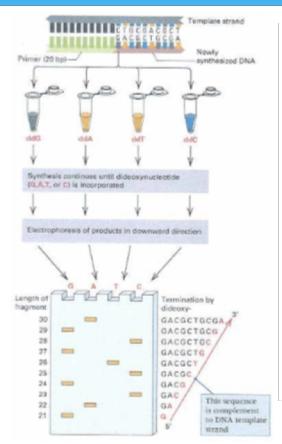


Figure 5.39  
Dideoxy method of DNA sequencing. Four DNA synthesis reactions are carried out in the presence of all normal nucleotides plus a small amount of one of the dideoxynucleotides (ddATP, ddTTP, ddCTP, ddGTP). Synthesis continues along a template strand until a dideoxynucleotide is incorporated. The products that result from termination at each dideoxynucleotide are indicated at the right. The fragments are separated by size by electrophoresis, and the sequence of the template strand can be read directly from the gel. In this example, the length of the primer needed

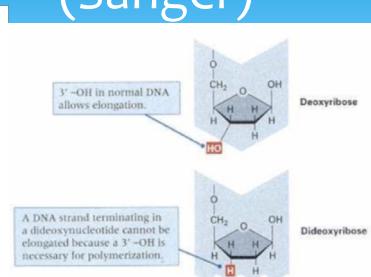


Figure 5.38  
Structures of normal deoxyribose and the dideoxyribose sugar used in DNA sequencing. The deoxyribose has a hydrogen atom (red) attached to the 3' carbon, in contrast with the hydroxyl group (red) at this position in deoxyribose. Because the 3' hydroxyl group is essential for the attachment of the next nucleotide in line in a growing DNA strand, the incorporation of a dideoxynucleotide immediately terminates synthesis.

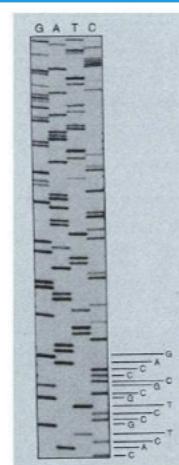
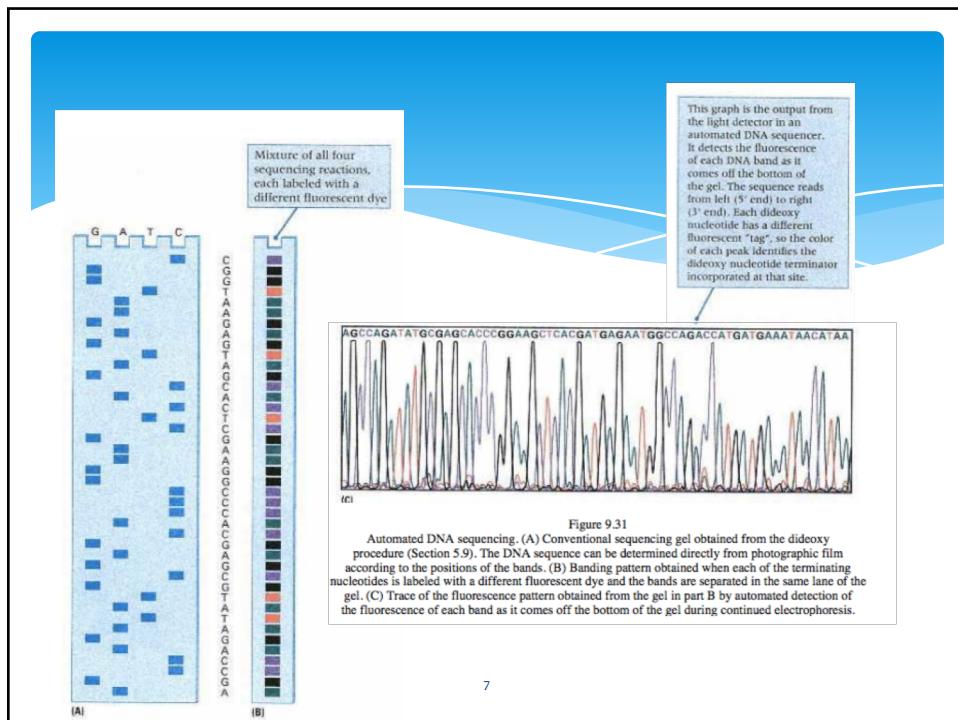
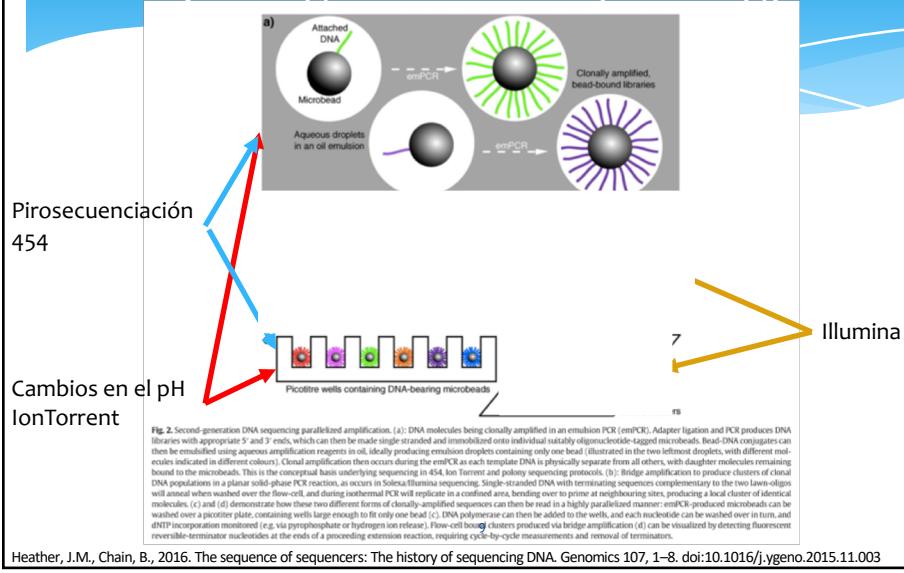


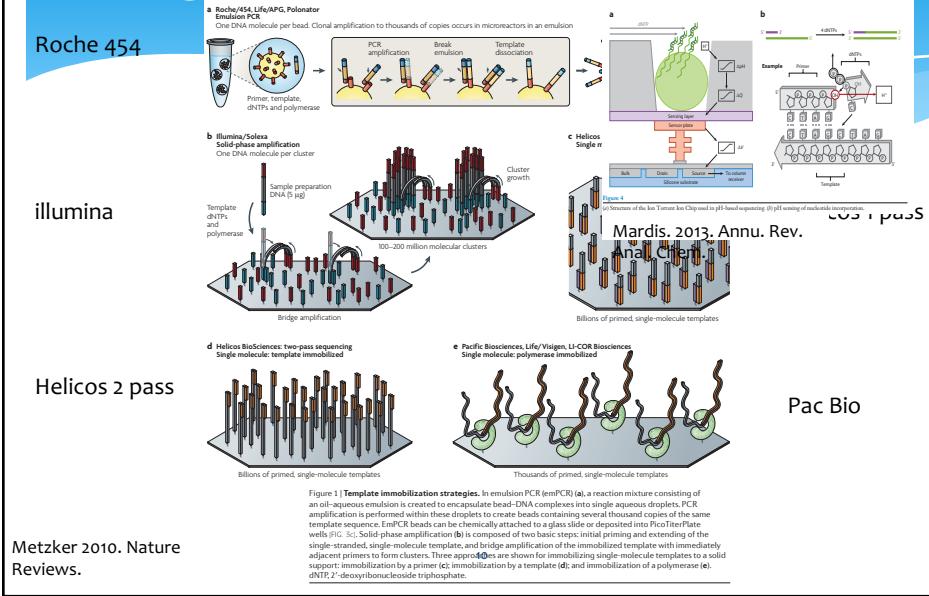
Figure 5.40  
A section of a dideoxy sequencing gel. The sequence is read from the bottom to the top. Each horizontal row represents a single nucleotide position in the DNA strand synthesized from the template. The vertical columns result from termination by the dideoxy forms of G, A, T, or C. The sequence from the lower part of the gel is indicated.



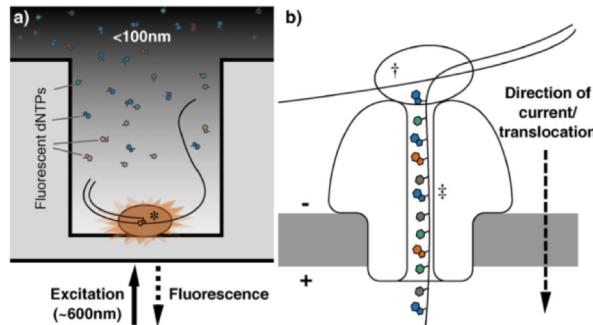
## Secuenciación de siguiente generación (Next generation sequencing)



## Next generation sequencing



## 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.

11

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

12

## Análisis de secuencias (Sanger)

### Procedimiento

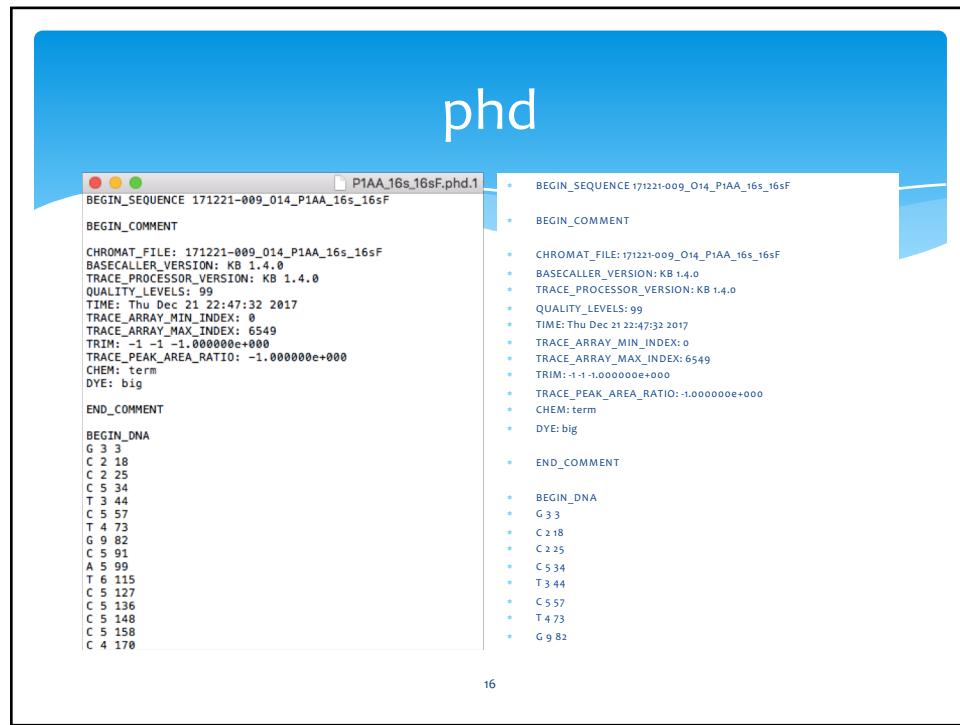
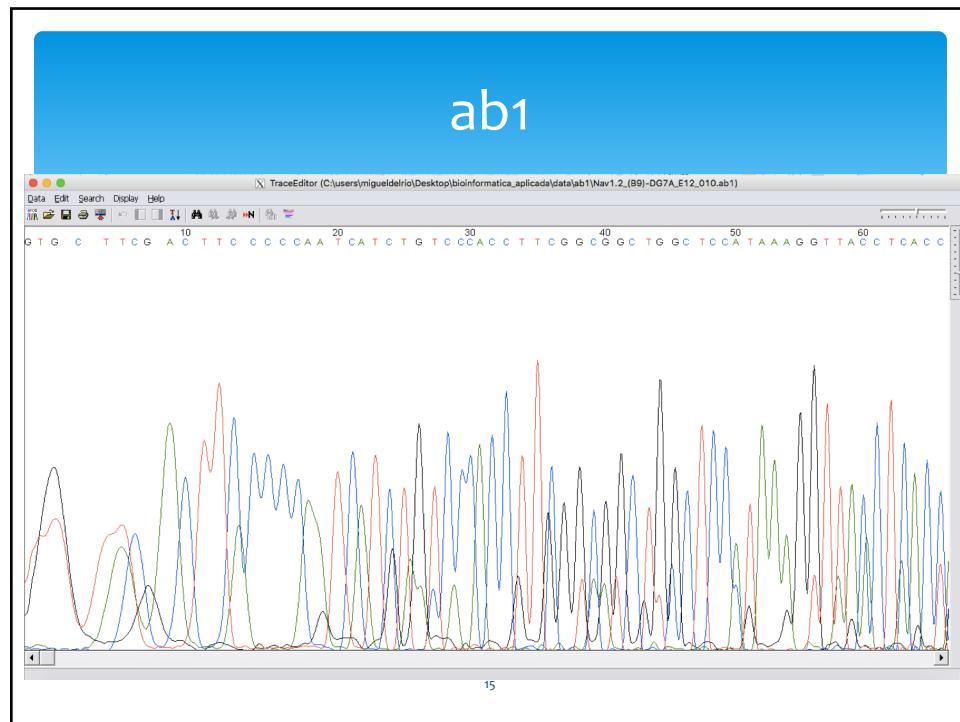
- \* Toma de muestra
  - \* Tejido: puede variar de una especie a otra (leta, branquia, músculo, etc.)
  - \* Conservación (etanol 95%, congelación)
- \* Extracción
  - \* Kit (Qiagen, Promega, etc.)
  - \* Precipitación con sales (LiCl)
- \* Amplificación
  - \* Dirigida al fragmento de interés Iniciadores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

13

## 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
■	16S_RO_1F-16AR_A07_001.seq	26/05/2017 11:20
■	16S_RO_1R-16BR_A08_002.ab1	26/05/2017 11:20
■	16S_RO_1R-16BR_A08_002.seq	26/05/2017 11:20
■	16S_RO_2F-16AR_B07_003.ab1	26/05/2017 11:20
■	16S_RO_2F-16AR_B07_003.seq	26/05/2017 11:20
■	P1AA_16s_16sF.ab1	21/12/2017 22:48
■	P1AA_16s_16sF.pdf	21/12/2017 22:49
■	P1AA_16s_16sF.phd.1	21/12/2017 22:48
■	P1AA_16s_16sF.txt	21/12/2017 22:49
■	P1AA_16s_16sR.ab1	21/12/2017 22:48
■	P1AA_16s_16sR.pdf	21/12/2017 22:48





**Seq, txt**

P1AA\_16s\_16sF.txt

```
>171221-009_014_P1AA_16s_16sF.ab1      543
GCCCTCTGCATCCCCCTGAGAGTTGGCCTGCTCGGTGATTAATATTAA
CAGCTCGGTATTATAACTGTACTAAGGTAGCATATAATTGGCTTATA
AATTGGGCTAGAATGAATGGTTGACGAAAATTGACTGTCTCTATT
ATTATTAGAAAATTAAATTGGTGTAGTGAGAAAGCTTAAATTGGTTAAAGG
GACGAAAAGACCTTATTGAGCTTATTATTATAAATTATATGATTGT
TGTTTATTAGAATAATTGGTGTGGAAAGGAAATAAACCAAATT
AAATTAACTCCCTAAATTAAATTGGTGTGGAAAGGAAATAAACCAAATT
TTGCTTAGAGATAAGTTACCATAGGGATAACAGCGTAATTGGAG
AGTTCATATTGAAAAAGAGATTGCGACCTCGATGTTGATTAATTAA
CCTTAAGGTGAAGAGGCTTTATTAGTAAATCTGTCGATTTTAAATT
TTACGTGATCTGGTCCAAAACCGGAAATATTATTTTCCCTC
```

Archivos

Formato: fasta

```
>171221-009_014_P1AA_16s_16sF.ab1      543
GCCCTCTGCATCCCCCTGAGAGTTGGCCTGCTCGGTGATTAATATTAA
CAGCTCGGTATTATAACTGTACTAAGGTAGCATATAATTGGCTTATA
AATTGGGCTAGAATGAATGGTTGACGAAAATTGACTGTCTCTATT
ATTATTAGAAAATTAAATTGGTGTAGTGAGAAAGCTTAAATTGGTTAAAGG
GACGAAAAGACCTTATTGAGCTTATTATTATAAATTATATGATTGT
TGTTTATTAGAATAATTGGTGTGGAAAGGAAATAAACCAAATT
AAATTAACTCCCTAAATTAAATTGGTGTGGAAAGGAAATAAACCAAATT
TTGCTTAGAGATAAGTTACCATAGGGATAACAGCGTAATTGGAG
AGTTCATATTGAAAAAGAGATTGCGACCTCGATGTTGATTAATTAA
CCTTAAGGTGAAGAGGCTTTATTAGTAAATCTGTCGATTTTAAATT
TTACGTGATCTGGTCCAAAACCGGAAATATTATTTTCCCTC
```

18

# 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
- ```
MTEITAAMVKELRESTGAGMMCKNALSETNGDFDKAVQLLREKGLGK
AAKKADRLAEGLVSVKSDDFTIAAMRPSYLS
```
- > Secuencia de DNA
- ```
ATGACTGAAATTACTGCTGCAATGGAAAAGAACTCCCGCAAAGTACAG
GCGCGGGGATGATGGATTGAAAAATGCTTGAGTGAACATAATGGAG
ATTTGATAAAGCAGTACAACCTTTAAGAGAAAAAGGTTTAGGTAGGC
TGCTAAAAAACAGATAGACTTGCTGCAGAAGGTTGGTAAGTGTAAA
AGTA
```

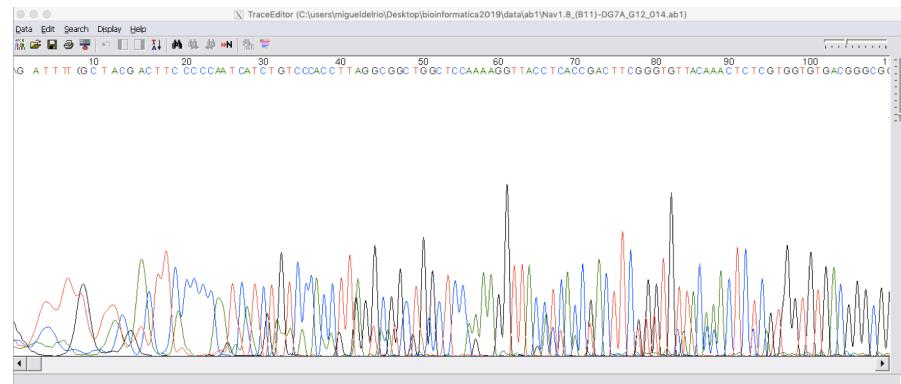
19

# Secuencia de DNA (y RNA)

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

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

## GenBank (nucleotido)

Swiss-Prot

* >NM_001273529.2 Drosophila melanogaster vasa (vas), transcript variant B, mRNA	CACTAGATTTGGTACTTTAACAGATCCTTCGGT TTGGCTTGGCGAAGGTGATCTGAACATTATCA AAGTTTGTAAGGTAATACATAAAAGTAAAAAGAATTAA TTTTGCTCTTGAAGGCAGGCCAAATTAAAAAAA AAAATATCAATATGTCGACCGATCGGATGAGCC CATTGTTGATACTCGCCGGCCGGCTGGAGA TTGGACCGATGATGAGGACACGGCCAAGAGCTTCAG CGGCCGAAGCTGAAGGCGATGGTGTGGAGGGAGC GGTGTGAAGGGCGGCCAACAGGAGGAATCGA GATGTTGGAAAGGATCGGGGAGGGACAGGAG GAGGAGCTGGAGGTTACGGAGGAATATCGGAGT GAGGGGGCTTCACGGTGGACGTGCGAGGGAGA AAGGGACTTCGGCGTGGAGAACGGCGCTCCCG TGGACAAGGGCGCTCCGGCTGGACAAGGCC TCCCGGGGGGAGACAAGGGCTTCTGGTGGAGAA GGCGGCTTCGGCGTGGCTGACGAAAACAGGG	* >tr M9PBB5 M9PBB5_DROME Vasa, isoform B OS=Drosophila melanogaster GN=vas PE=3 SV=1 MSDDWDPEIVDTRGARGDWSDDEDTAKSFSGEAE GDGVGGGGGEGGGYQQGGNNDVFGRR IGGGRGGGAGGYRGGNRDGGFHGRRERGERDFRGG EGGFRGGQGGSRGGQGGSRGGQGG FRGEGGGFRGRLYENEDGDRRGLRDLREERGGERRG RLDREERGERGERDGGFARRRR NEDINNNNNIVEDVERKREFYIPPEPSNDAEIFSSGIA SGIHFSKYNINPVKTGSDV PQPIQHFTSADLIRDIDNVNKSGYKIPPTIQKCSIPV GRDLMAACQTSGSKTAFL LPLISKLLELDPEHEELGRPQVVIVSPRELAIQIFNEAR FAFESYLKIGIVYGGTSFRH QNECITRGCHVIATPGRLLDFVDRTFITFEDTRFVVL EADRMLDMGFSEDMRRIMTHV
--	--	--

23

# fastaq

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

24

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; Neocoleoidea; 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

27

### 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%
30	1 in 1,000	99.9%
40	1 in 10,000	99.99%
50	1 in 100,000	99.999%

Sanger

The screenshot shows the NCBI homepage with a blue header. The main content area features a "Welcome to NCBI" section with links to "Submit", "Download", "Learn", "Develop", "Analyze", and "Research". To the left is a sidebar with a "NCBI Home" menu containing links like "Resource List (A-Z)", "All Resources", "Chemicals & Bioassays", etc. On the right, there's a "Popular Resources" section with links to "PubMed", "Bookshelf", "PubMed Central", "PubMed Health", "BLAST", "Nucleotide", "Genome", "SNP", "Gene", "Protein", and "PubChem". A "NCBI News & Blog" section is also present.

29

The screenshot shows the UniProt homepage with a blue header. The main content area features a "UniProtKB" search bar and a message about transitioning to HTTPS. Below this, there are sections for "UniProtKB", "UniRef", "UniParc", "Proteomes", "Supporting data", and "News". The "Supporting data" section includes links for "Literature citations", "Cross-ref. databases", "Taxonomy", "Diseases", and "Keywords". The "News" section has a "Forthcoming changes" link. At the bottom, there are links for "Getting started", "YouTube", "UniProt data", "Text search", "Download latest release", and "Protein spotlight". The URL "http://www.uniprot.org" is displayed at the bottom.

30

The screenshot shows the ENA homepage with a blue header containing the word "ENA". Below the header is the EMBL-EBI logo and the ENA logo with the text "European Nucleotide Archive". A search bar at the top right contains the placeholder "Examples: BN000065, histone" and a "Search" button. The main content area includes a "Text Search" section with a search input field and a "Sequence Search" section with a text area for pasting sequences. To the right, there's a "Popular" sidebar with links to "Submit and update", "Sequence submissions", "Genome assembly submissions", "Submitting environmental sequences", "Citing ENA data", "Rest URLs for data retrieval", and "Rest URLs to search ENA". Another sidebar on the right displays "Latest ENA news" with entries for December 21, 2017, and January 21, 2018. At the bottom of the page is the URL <https://www.ebi.ac.uk/ena>.

The screenshot shows the DDBJ homepage with a blue header containing the text "DDBJ, DNA Database of Japan". Below the header is the DDBJ logo with the text "DNA Data Bank of Japan" and "30". A navigation bar with tabs for "About DDBJ", "How to Use", "Report/Statistics", "FAQ", and "Contact Us" is visible. On the left, there are links for "RSS", "DDBJ Twitter", and "Mail Magazine". Below these are logos for "DDBJ", "INSDC", "NCBI", and "EMBL-EBI", with the text "International Nucleotide Sequence Database Collaboration". On the right, there are icons for "Data Submission", "Search / Analysis", "Super Computer", and "ftp.ddbj.nig.ac.jp". A "Hot Topics" section lists recent news items, and a "News Archive" section provides links for "News", "Release", "PR", "Maintenance", "Operation", and "All". At the bottom of the page is the URL <http://www.ddbj.nig.ac.jp>.

The screenshot shows the main NCBI homepage. At the top, there's a navigation bar with links for 'NCBI Resources' and 'How To'. Below the header, the NCBI logo is displayed with the text 'National Center for Biotechnology Information'. A search bar is present with the placeholder 'All Databases'. On the left, a sidebar lists categories such as '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', and 'Variation'. The main content area features a 'Welcome to NCBI' section with a brief description and links to 'About the NCBI', 'Mission', 'Organization', and 'NCBI News & Blog'. It also includes sections for '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), and 'Research' (Explore NCBI research and collaborative projects). On the right, there's a 'Popular Resources' sidebar with links to PubMed, Bookshelf, PubMed Central, PubMed Health, BLAST, Nucleotide, Genome, SNP, Gene, Protein, and PubChem. Below that is a 'NCBI News & Blog' section with recent articles like 'ClinVar Unveils New, More Intuitive Variation Display' and 'NIH Data Hackathon on campus – January 22-24, 2018'. At the bottom of the page, there's a footer with the number '33'.

## Bases de datos

Nombre	Dirección	Datos
NCBI	<a href="https://www.ncbi.nlm.nih.gov">https://www.ncbi.nlm.nih.gov</a>	
GenBank	<a href="https://www.ncbi.nlm.nih.gov/genbank/">https://www.ncbi.nlm.nih.gov/genbank/</a>	Base de datos
Nucleotide	<a href="https://www.ncbi.nlm.nih.gov/nucleotide/">https://www.ncbi.nlm.nih.gov/nucleotide/</a>	nucleótidos
PubMed	<a href="https://www.ncbi.nlm.nih.gov/pubmed">https://www.ncbi.nlm.nih.gov/pubmed</a>	Referencias bibliográficas
NR	<a href="https://www.ncbi.nlm.nih.gov/protein/">https://www.ncbi.nlm.nih.gov/protein/</a>	Secuencia de proteínas no redundante
Swiss-Prot, UniProtKB	<a href="http://www.uniprot.org">http://www.uniprot.org</a>	Universal Protein Resource
KEGG (Kyoto Encyclopedia of Genes and Genomes)	<a href="http://www.genome.jp/kegg/">http://www.genome.jp/kegg/</a>	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 Handbook.  
[https://www.ncbi.nlm.nih.gov/books/NBK21101/pdf/Bookshelf\\_NBK21101.pdf](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

**Amazon EC2**

 Google Cloud Platform



37

## Recomendaciones al usar servidores

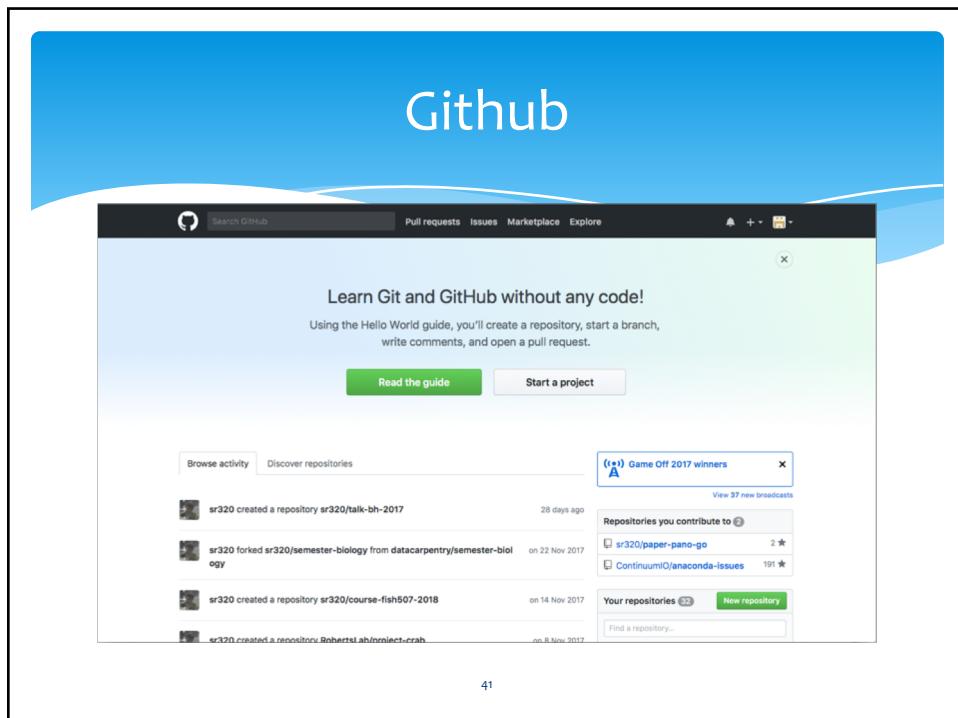
1. Los datos nunca están seguros en la red
2. Anotar el servidor, la base de datos y la versión del programa utilizado
3. Anotar los números de identificación de las secuencias
4. Anotar los parámetros del programa
5. Guarde los resultados de internet inmediatamente
6. Utilice los valores de E
7. Asegúrese 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)

Claverie y Notredame. 2007. Bioinformatics for dummies. 2nd Ed. Wile Publishing, Inc.

38

The screenshot shows a blue header with the word "Manuales". Below it is a white content area. At the top left is the NCBI logo and a search bar. The main content is titled "The NCBI Handbook, 2nd edition". It features a thumbnail image of the book cover, which is a dark blue cover with the title and some abstract diagrams. Below the thumbnail is a brief description: "Bethesda (MD): National Center for Biotechnology Information (US); 2013-. Copyright and Permissions". To the right of the description is a sidebar with links: "Views", "PubReader", "Print View", "Cite this Page", "PDF version of this title (14M)", and "Disable Glossary Links". Further down the sidebar are sections for "Related Information" (listing "NLM Catalog") and "Recent Activity" (listing several recent searches like "The NCBI Handbook", "Drosophila melanogaster vasa (vas)", etc.). At the bottom of the content area is a URL: <https://www.ncbi.nlm.nih.gov/books/NBK143764/>.

The screenshot shows a blue header with the text "Basic Local Alignment Search Tool (BLAST)". Below it is a white content area. At the top is the NCBI header with the NIH logo, a search bar, and navigation links. The main content area is titled "Basic Local Alignment Search Tool". It has a sub-section "Web BLAST" featuring three buttons: "Nucleotide BLAST" (nucleotide > nucleotide), "blastx" (translated nucleotide > protein), and "tblastn" (protein > translated nucleotide). Another sub-section "Protein BLAST" (protein > protein) is also shown. At the bottom is a section "BLAST Genomes" with a search bar and buttons for "Human", "Mouse", "Rat", and "Microbes". A news box on the right says "igBLAST 1.8.0 released" and "A new version of igBLAST is now available. Wed, 15 Nov 2017 16:00:00 EST". There is also a link "More BLAST news...".



41



42

archivos del curso bioinformática aplicada UAM-I

4 commits 1 branch 0 releases 1 contributor

**.gitattributes** Initial commit 24 days ago

**0temario.pynb** temario 24 days ago

Copiar y pegar la secuencia en el blastn

bioinformatica\_aplicada / ejercicio1.txt

0 contributors

Executable File 2 lines (1 sloc) | 335 Bytes

Raw Blame History

```
GTGCTTCGACTTCCCCAACATCTGTGCCACCTCGGCCTGGCTCATAAAGGTTACCTCAGGACTCTCGTGTTGACGGCCGGTGTACAAGGCCGGAACGT
```

43

U.S. National Library of Medicine > NCBI National Center for Biotechnology Information

BLAST® > blastn suite Standard Nucleotide BLAST

Enter accession number(s), gie(s), or FASTA sequence(s)

Job Title

Align two or more sequences

Choose Search Set Database: Human genomic + transcript

Organism Optional

Exclude Optional

Limit to Optional

Entrez Query Optional

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 (refseq\_representative\_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

```
CTTCCCCAATCATCTGTCCCACCTCGCGGCTGGCTCCATAAA
GGTTACCTCACCGACTCAGGGTGTACAAACTCTCGTGGTGTGAC
GGGCGGTGTGTACAAGGCCGGAACGTATTCACCGCGGCATGC
TGATCCGAATTACAACCGATTCCAGCTTCACGCATTCAAGTTGC
AAACTGCAATCCGAACTGAAAACAGATTGTGGAATTGGCTAA
CCTCCCGGTTTCCCTGCCCTTGTCTGTCATTGTACACGTGTG
TACCCCAGGTATAAGGGCATGATGATTGACGTCATC
```

\* Utilizando las bases de datos

- \* Nt

- \* 16S microbial

\* Con sus compañeros de mesa discuta las diferencias entre la información proporcionada en cada búsqueda

45

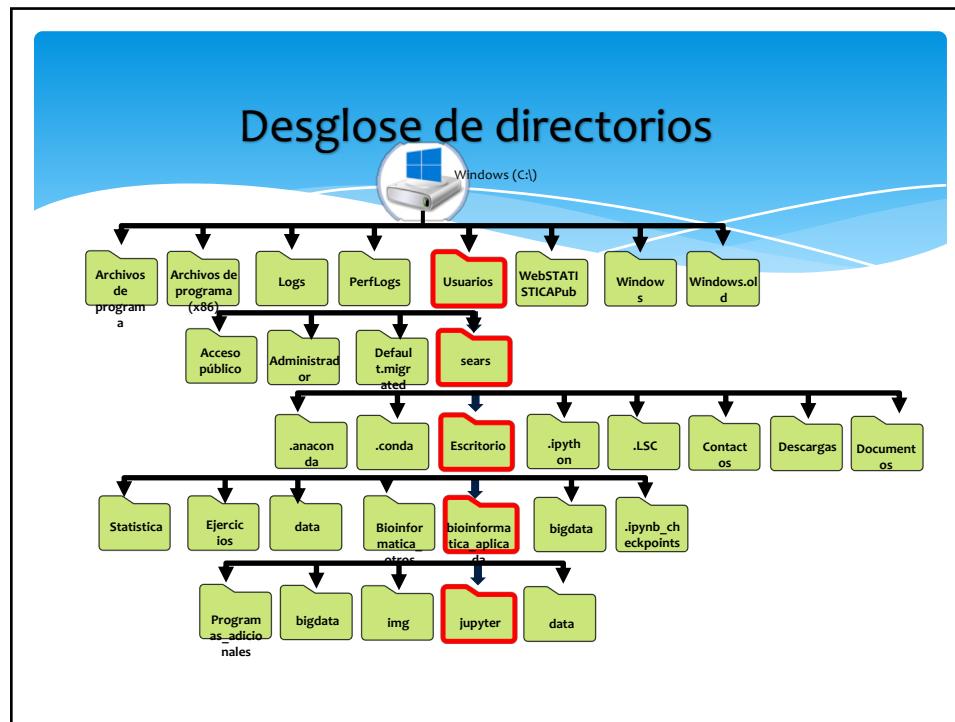
## Cuestionario

\* 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

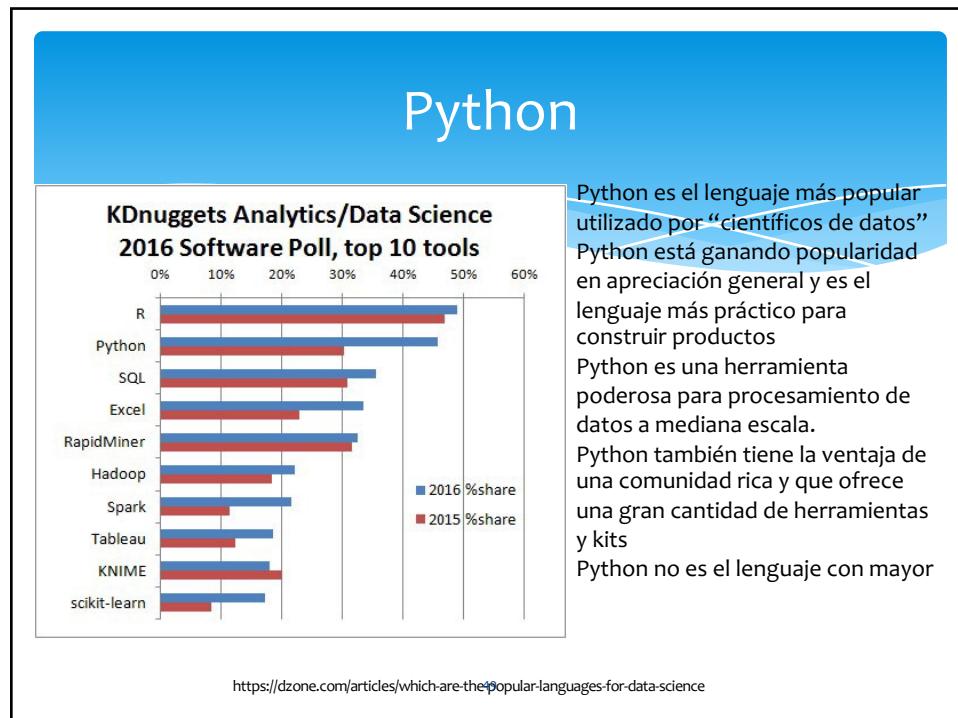
\* ¿En qué formatos de archivo se puede descargar la información?

46



## Lenguajes

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

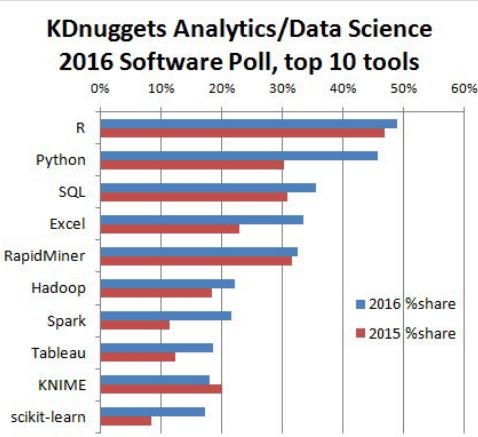


- \* Python
  - \* [python.org](http://python.org)
  - \* <http://www.python-course.eu/index.php>
  - \* Google
  - \* <https://developers.google.com/edu/python/>
- 50

- \* 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

<http://www.kdnuggets.com/2017/07/6-reasons-python-suddenly-super-popular.html>

## R



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 utilidades comerciales se han dispersado enormemente.

<https://dzone.com/articles/which-are-the-popular-languages-for-data-science>

## Otros lenguajes

- \* Perl
- \* Ruby

53

## 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”.

54

Giovanni M Dall'Olio <https://www.biostars.org/p/13972/>

## 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 Monty Python.”

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

## Python vs Ruby vs Perl

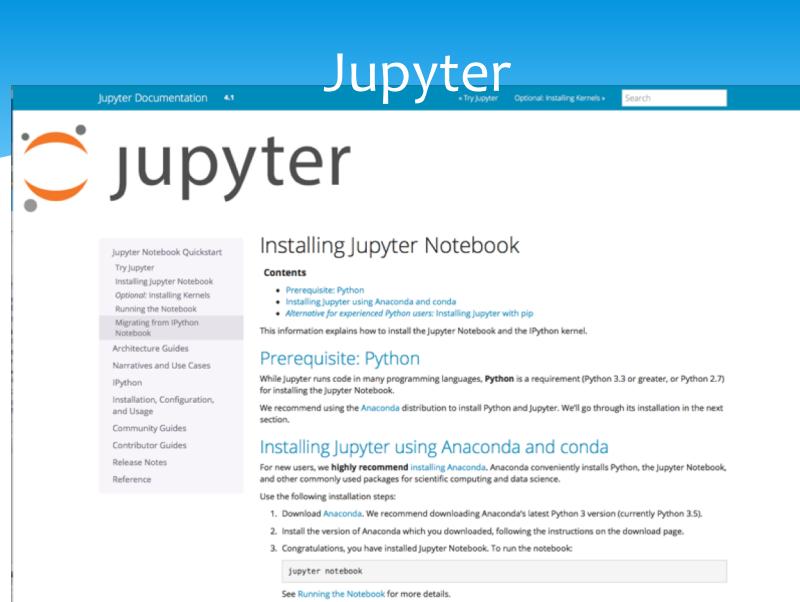
- \* Perl for one-liners in a pipe. Python for everything else
- \* Python vs Ruby vs Perl
- \* Right folks! Time for another religious war. Or, a Halo style three way battle.
- \* On the left, we have Perl, the grand-daddy of the scripting world. But don't let his age fool you, he still has a few tricks 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 especiales 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.  
 Deberá 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!

57



The screenshot shows the Jupyter Documentation website with the title 'Jupyter' at the top. On the left, there's a sidebar with links like 'Jupyter Notebook Quickstart', 'Try Jupyter', 'Installing Jupyter Notebook', etc. The main content area has a heading 'Installing Jupyter Notebook'. It includes a 'Contents' section with a list of links, a 'Prerequisite: Python' section with a note about Python version requirements, and a 'Installing Jupyter using Anaconda and conda' section with instructions for new users. At the bottom, there's a link to 'Running the Notebook'.

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

58



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