

GenBank



Práctica 1 - Bioinformática

Julia Guerrero Viu Sergio Izquierdo Barranco Ubaldo Ramón Júlvez

- 1. Presentación GenBank
- 2. Consultas
- 3. Descarga datos

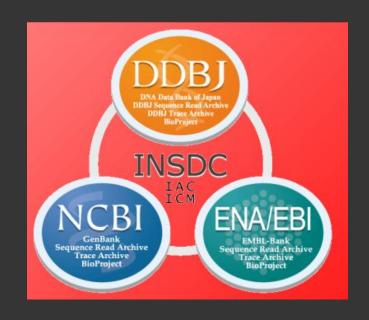
¿Qué es GenBank?

Colección totalmente pública de todas las secuencias de nucleótidos y su traducción a proteínas

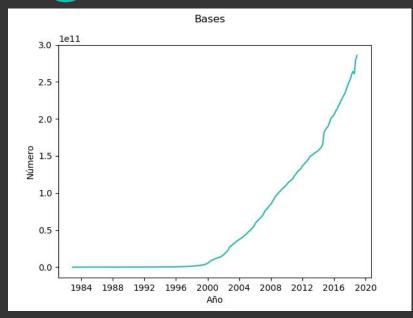
_

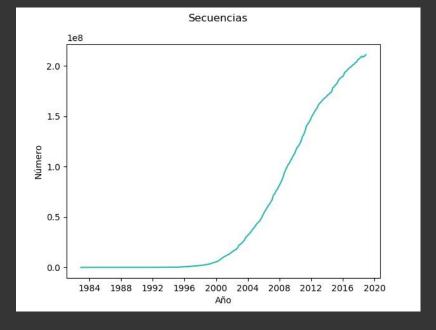
Historia GenBank

- 1979 Walter Goad in Los Álamos National Laboratory
- 1982 GenBank (NIH)
- 1992 NCBI
- INSDC (International Nucleotide Sequence Database Collaboration)



285.688.542.186 bases iy sigue creciendo!





_

Algunos datos

Cada **2 meses** sale una nueva versión de GenBank Cada **18 meses** se duplica el número de secuencias de GenBank

Recibe secuencias diariamente de más de

100.000 organismos internacionales distintos

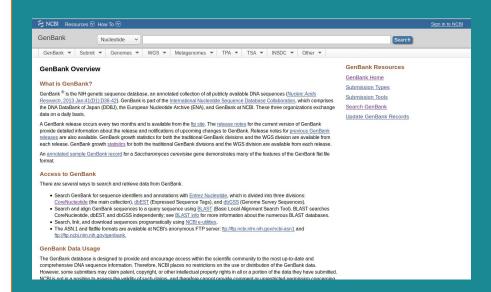
ACCEDER a GenBank

- Directamente BD (herramienta web)
- BLAST

 (herramienta de alineamiento de secuencias)
- Programas del NCBI
 - Servidor FTP
 (base de datos completa:

ftp://ftp.ncbi.nih.gov/genbank)

http://www.ncbi.nlm.nih.gov/genbank/

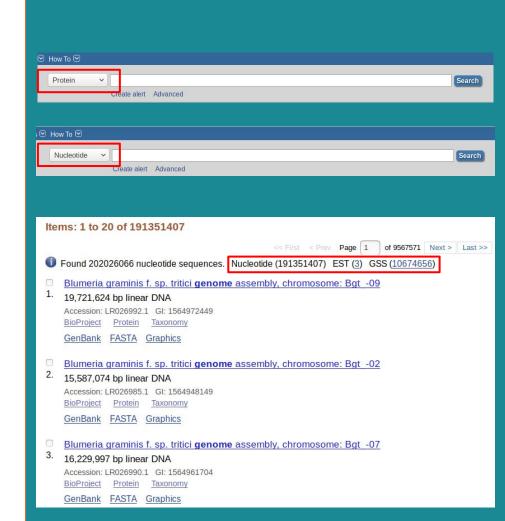


FUENTES (BBDD):

- Entrez Protein

- Entrez Nucleotide
 - Nucleotide (core)
 - EST and GSS
 - -> se van a fusionar en 2019

https://ncbiinsights.ncbi.nlm.nih.gov/2018/07/30/upcoming-changes-est-gss-databases/



Divisiones de GenBank

Submissions phases

Organismal Divisions

Database	Division	BLAST	Example
ВСТ	Bacterial sequences	nr, month	
PRI	Primate sequences	nr, month	Human Phase 3
ROD	Rodent sequences	nr, month	
MAM	Other mammalian sequences	nr, month	
VRT	Other vertebrate sequences	nr, month	
INV	Invertebrate sequences	nr, month	Drosophila, C. elegans Phase 3
PLN	Plant and Fungal sequences	nr, month	Arabidopsis Phase 3
VRL	Viral sequences	nr, month	
PHG	Phage sequences	nr, month	
RNA	Structural RNA sequences	nr, month	
SYN	Synthetic and chimeric sequences	nr, month	
UNA	Unannotated sequences	nr, month	

Functional Divisions

Database	Division	BLAST	Example
EST	Expressed Sequence Tags	dbest, month	
STS	Sequence Tagged Sites	dbsts, month	
GSS	Genome Survey Sequences	dbgss, month	
HTG	High Throughput Genomic sequences	htgs, month	All Organisms: Phase 0, 1, and 2

- Phase 0 sequences are single-few pass reads of a single clone (not contigs usually).
- Phase 1 sequences are unfinished, unordered, and contain gaps.
- Phase 2 sequences are unfinished, ordered, and can contain one or more gaps.
- Phase 3 sequences are high quality finished sequences that do not contain gaps.

Consultas en GenBank



- Consultas de frase
 - 1. Proteínas
 - 2. Genes
 - 3. Autores

kinase

Smith JR

Felis catus

"duchenne muscular dystrophy"

- Operadores Binarios
 - OR, AND, NOT
- Filtros
 - Especie
 - Tipo de molécula (DNA, RNA)
 - Fecha
 - etc.

contactin AND neurofascin

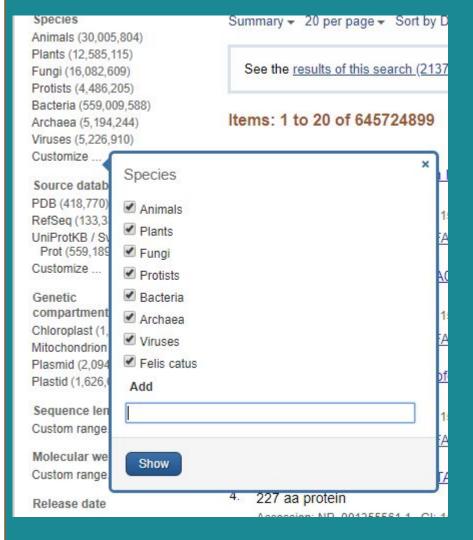
"Felis catus"[Organism]

biomol_genomic[PROP]

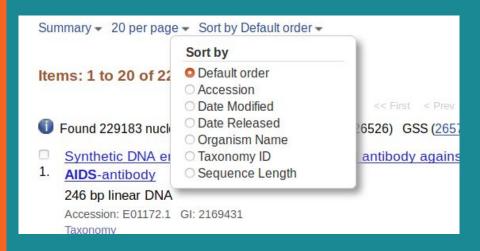
all[filter]

- Más filtros

Species Animals Archaea See LOC113060289 (CO Bacteria contactin reference seque Fungi Plants Protists Items: 1 to 20 of 343 Viruses Customize ... Molecule types Columba livia breed Da genomic DNA/RNA (337) 4,078,870 bp linear DN mRNA (6) Accession: AKCR02000055 Customize ... BioProject BioSample Source databases GenBank FASTA Grai INSDC (GenBank) (49) RefSeg (292) Crassostrea gigas stra Customize ... whole genome shotgur Sequence Type 450,376 bp linear DNA Nucleotide (343) Accession: NW 011937980 Assembly BioProject Sequence length GenBank FASTA Grai Custom range... Release date Amazona aestiva isola Custom range... 3,954,725 bp linear DN Accession: LMAW0100044 Revision date Assembly BioProject Custom range... GenBank FASTA Gran Clear all Tupaia chinensis unpla Show additional filters 8,920,270 bp linear DN - Filtros personalizados



- Ordenación de resultados
 - Fecha
 - Organismo
 - Longitud secuencia
 - Taxonomía



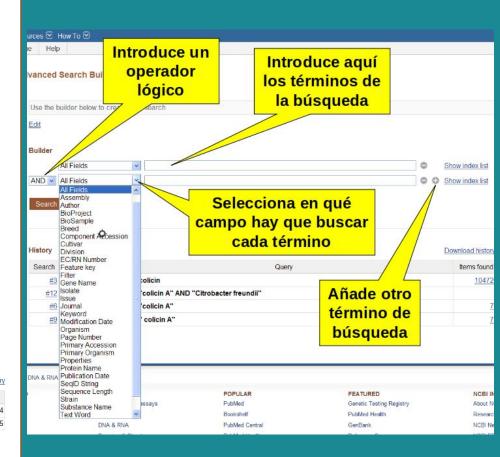
- Búsqueda **AVANZADA**



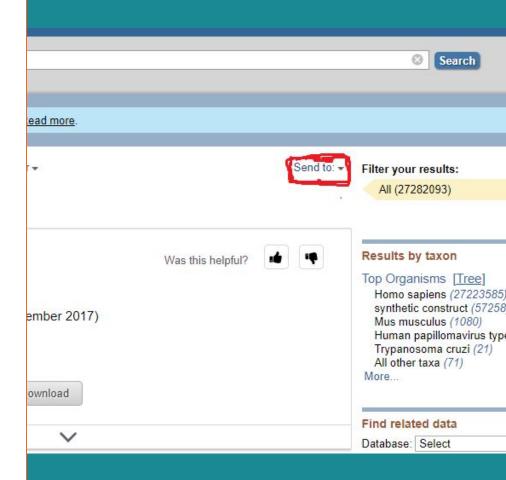
- Campos
- Operadores lógicos
- Historial

History Download history O				
Search	Add to builder	Query	Items found	Time
<u>#4</u>	Add	Search genome	191351407	11:25:04
<u>#1</u>	Add	Search aids	226526	11:24:45

Consulta y modifica consultas anteriores



- Descarga de datos



Formato GenBank:

- Cabecera
- Referencias
- Features
- Secuencia

```
LOCUS
            AP018036
                                                       linear BCT 10-JAN-2019
                                     327 bp
           Mycobacterium tuberculosis DNA, complete genome, strain: HN-506.
DEFINITION
            AP018036 REGION: 1646..1972
ACCESSION
VERSION
            AP018036.1
DBI TNK
            BioProject: PRJDB5607
            BioSample: SAMD00076258
KEYWORDS
SOURCE
            Mycobacterium tuberculosis
           Mycobacterium tuberculosis
 ORGANISM
            Bacteria; Actinobacteria; Corynebacteriales; Mycobacteriaceae;
            Mycobacterium; Mycobacterium tuberculosis complex.
REFERENCE
            Wada, T., Hijikata, M., Maeda, S., Hang, N.T.L., Thuong, P.H.,
  AUTHORS
            Hoang, N.P., Hung, N.V. and Keicho, N.
            Complete Genome Sequences of Three Representative Mycobacterium
 TITLE
            tuberculosis Beijing Family Strains Belonging to Distinct Genotype
            Clusters in Hanoi, Vietnam, during 2007 to 2009
            Genome Announc 5 (27), e00510-17 (2017)
  JOURNAL
   PUBMED
            28684565
            DOI:10.1128/genomeA.00510-17
  REMARK
            Publication Status: Online-Only
REFERENCE
            2 (bases 1 to 327)
            Wada, T., Hijikata, M., Maeda, S., Hang, N.T., Thuong, P.H., Hoang, N.P.,
  AUTHORS
            Hung.N.V. and Keicho.N.
 TITLE
            Direct Submission
            Submitted (06-APR-2017) Contact: Takayuki Wada Institute of Tropical
            Medicine, Department of International Health; Sakamoto 1-12-4,
            Nagasaki, Nagasaki 852-8523, Japan
            This project was organized by Department of Pathophysiology and
COMMENT
            Host Defense, The Research Institute of Tuberculosis, Japan
            Anti-Tuberculosis Association, Tokyo, Japan, collaborating with
            Hanoi Lung Hospital, Hanoi, Vietnam.
            This project was supported by J-GRID, granted by Japan Agency for
            Medical Research and Development (AMED) and Ministry of Education,
            Culture, Sports, Science and Technology (MEXT).
            Annotated using prokka 1.11 from http://www.vicbioinformatics.com.
            Annotated at D-FAST https://dfast.nig.ac.jp
```

Formato GenBank:

- Cabecera
- Referencias
- Features
- Secuencia

```
FEATURES
                     Location/Oualifiers
     source
                     1..327
                     /organism="Mycobacterium tuberculosis"
                     /mol type="genomic DNA"
                     /strain="HN-506"
                     /host="Homo sapiens"
                     /db xref="taxon:1773"
                     /country="Viet Nam:Hanoi"
                     /collection date="2009-03-24"
                     1..327
     gene
                     /locus tag="HN506 00002"
                     1..327
     CDS
                     /locus tag="HN506 00002"
                     /inference="ab initio prediction:Prodigal:2.6"
                     /inference="similar to AA sequence:UniProtKB:P59801"
                     /codon start=1
                     /transl table=11
                     /product="insertion element IS6110 uncharacterized 12.0
                     kDa protein"
                     /protein id="BAX47168.1"
                     /translation="MSGGSSRRYPPELRERAVRMVAEIRGQHDSEWAAISEVARLLGV
                     GCAETVRKWVROAOVDAGARPGTTTEESAELKRLRRDNAELRRANAILKTASAFFAAE
                     LDRPAR"
ORIGIN
        1 atgtcaggtg gttcatcgag gaggtacccg ccggagctgc gtgagcgggc ggtgcggatg
       61 gtcgcagaga tccgcggtca gcacgattcg gagtgggcag cgatcagtga ggtcgcccgt
      121 ctacttggtg ttggctgcgc ggagacggtg cgtaagtggg tgcgccaggc gcaggtcgat
      181 gccggcgcac ggcccgggac cacgaccgaa gaatccgctg agctgaagcg cttgcggcgg
      241 gacaacqccg aattqcqaaq qqcqaacqcq attttaaaqa ccqcqtcqqc tttcttcqcq
      301 gccgagctcg accggccagc acgctaa
```

Formato Fasta:

- Cabecera
- Secuencia

>AP018036.1 Mycobacterium tuberculosis DNA, complete genome, strain: HN-506 TTGACCGATGACCCCGGTTCAGGCTTCACCACAGTGTGGAACGCGGTCGTCTCCGAACTTAACGGCGACC CTAAGGTTGACGACGGACCCAGCAGTGATGCTAATCTCAGCGCTCCGCTGACCCCTCAGCAAAGGGCTTG GCTCAATCTCGTCCAGCCATTGACCATCGTCGAGGGGTTTGCTCTGTTATCCGTGCCGAGCAGCTTTGTC CAAAACGAAATCGAGCGCCATCTGCGGGCCCCGATTACCGACGCTCTCAGCCGCCGACTCGGACATCAGA TCCAACTCGGGGTCCGCATCGCTCCGCCGGCGACCGACGACGACGACACTACCGTGCCGCCTTCCGA AAATCCTGCTACCACATCGCCAGACACCACAACCGACAACGACGAGATTGATGACAGCGCTGCGGCACGG GGCGATAACCAGCACAGTTGGCCAAGTTACTTCACCGAGCGCCCGCACAATACCGATTCCGCTACCGCTG CGCCGCCGCCTTGGCGATCGCAGAAGCACCCGCCCGCGCTTACAACCCCCTGTTCATCTGGGGCGAGTCC GGTCTCGGCAAGACACCCTGCTACACGCGGCAGGCAACTATGCCCAACGGTTGTTCCCGGGAATGCGGG TCAAATATGTCTCCACCGAGGAATTCACCAACGACTTCATTAACTCGCTCCGCGATGACCGCAAGGTCGC ATTCAAACGCAGCTACCGCGACGTAGACGTGCTGTTGGTCGACGACATCCAATTCATTGAAGGCAAAGAG GGTATTCAAGAGGAGTTCTTCCACACCTTCAACACCTTGCACAATGCCAACAAGCAAATCGTCATCTCAT CTGACCGCCCACCCAAGCAGCTCGCCACCCTCGAGGACCGGCTGAGAACCCGCTTTGAGTGGGGGCTGAT CACTGACGTACAACCACCCGAGCTGGAGACCCGCATCGCCATCTTGCGCAAGAAAGCACAGATGGAACGG CTCGCGGTCCCCGACGATGTCCTCGAACTCATCGCCAGCAGTATCGAACGCAATATCCGTGAACTCGAGG GCGCGCTGATCCGGGTCACCGCGTTCGCCTCATTGAACAAAACACCAATCGACAAAGCGCTGGCCGAGAT TGTGCTTCGCGATCTGATCGCCGACGCCAACACCATGCAAATCAGCGCGGCGACGATCATGGCTGCCACC GCCGAATACTTCGACACTACCGTCGAAGAGCTTCGCGGGCCCGGCAAGACCCGAGCACTGGCCCAGTCAC GACAGATTGCGATGTACCTGTGTCGTGAGCTCACCGATCTTTCGTTGCCCAAAATCGGCCAAGCGTTCGG CCGTGATCACACACCGTCATGTACGCCCAACGCAAGATCCTGTCCGAGATGGCCGAGCGCCGTGAGGTC TTTGATCACGTCAAAGAACTCACCACTCGCATCCGTCAGCGCTCCAAGCGCTAGCACGGCGTGTTCTTCC GACAACGTTCTTAAAAAAACTTCTCTCTCCCAGGTCACACCAGTCACAGAGATTTGAACCGCCCCGGCAT GTCCGGAGACTCCAGTTCTTGGAAAGGATGGGGTCATGTCAGGTGGTTCATCGAGGAGGTACCCGCCGGA GCTGCGTGAGCGGCGGTGCGGATGGTCGCAGAGATCCGCGGTCAGCACGATTCGGAGTGGGCAGCGATC AGTGAGGTCGCCCGTCTACTTGGTGTTGGCTGCGCGGAGACGGTGCGTAAGTGGGTGCGCCAGGCGCAGG

Limitar el número de resultados a descargar.

Items: 1 to 20 of 27282093

Selected: 3

<< First < Prev Page 1

- Filters activated: Homo sapiens. Clear all
- Found 28475230 nucleotide sequences. Nucleotide (27282093) GSS (119313)
- Homo sapiens ATP synthase F1 subunit alpha (ATP5F1A), RefSeqGen
- 1. 35,986 bp linear DNA

Accession: NG_041769.2 GI: 1573294352

GenBank FASTA Graphics

- Homo sapiens aspartoacylase (ASPA), RefSeqGene on chromosome 1
- 36,760 bp linear DNA

Accession: NG_008399.2 GI: 1573293807

GenBank FASTA Graphics

- Homo sapiens aldehyde dehydrogenase 2 family member (ALDH2), Re
- 3. 12

57,646 bp linear DNA

Accession: NG_012250.2 GI: 1573289586

GenBank FASTA Graphics

- Homo sapiens adducin 2 (ADD2), RefSegGene on chromosome 2
- 118,460 bp linear DNA

Accession: NG 029481.2 GI: 1573279893

GenBank FASTA Graphics