

# Introducción a la Bioinformática

## Práctica 1: BLAST y Recuperación de Secuencias

### Recuperación de Secuencias

La recuperación de secuencias, es decir la búsqueda y obtención de secuencias de interés en bases de datos, es una de las tareas más comunes en bioinformática. A primera vista puede parecer una tarea sencilla, pero llegar a hacerlo de una manera realmente efectiva requiere de cierto conocimiento y destreza.

Esta práctica cubrirá con cierta extensión esta labor, y al final de ella serémos capaces de extraer la información precisa de las bases de datos más comunes, de una manera eficiente.

### **NCBI**

Una de las bases de datos más conocidas presentes en el NCBI es el GenBank. Esta base de datos consta de 59,750,386,305 bases en 54,584,635 entradas en las divisiones más comunes de GenBank (EST/UniGene, STS, GSS HTGS) y 63,183,065,091 bases en 12,465,546 entradas en la división WGS (Febrero 2006).

A continuación veremos una de las formas más sencillas de acceder a la información presente en GenBank y el NCBI en general.

Acceda al sitio web del NCBI ubicado en la siguiente dirección:

<http://www.ncbi.nlm.nih.gov/>

The screenshot shows the NCBI homepage with the search bar set to "All Databases" and the "Go" button visible. The left sidebar contains navigation links for "SITE MAP", "About NCBI", "GenBank", "Literature databases", "Molecular databases", and "Genomic biology". The main content area features several informational boxes: "What does NCBI do?", "Whole Genome Association", "100 Gigabases", and "PubMed Central". A "Hot Spots" list is on the right.

Realizaremos una búsqueda de HIV-1. Asegurese de que ha definido una búsqueda en todas las bases de datos en el menú desplegable ubicado en la esquina superior izquierda y digite el término “**HIV-1**”. A continuación presione el botón “**go**”.

This screenshot shows the same NCBI homepage, but with the search dropdown menu open. The "All Databases" option is highlighted with a red box. The search text "HIV-1" is entered in the search box, and the "Go" button is visible. The rest of the page content remains the same as in the previous screenshot.

Obviamente es posible escoger cualquiera de las posibilidades ofrecidas en el menú. Son de destacar Pubmed, Protein y Nucleotide, con las cuales buscamos directamente en la base de datos de bibliografía, DNA o proteínas respectivamente.

Unos segundos después seremos llevados a la página web del sistema **ENTREZ** del NCBI, desde donde tendremos una perspectiva general de la información relacionada con nuestra secuencia presente en el NCBI.

The screenshot shows the NCBI Entrez search engine interface. At the top, there is a navigation bar with links for HOME, SEARCH, SITE MAP, PubMed, All Databases, Human Genome, GenBank, Map Viewer, and BLAST. Below this is a search bar with the text "HIV-1" and buttons for GO, CLEAR, and Help. The main content area displays search results across various databases, organized into two columns. Each result includes a count, an icon, and a brief description. A legend at the bottom indicates that gray background colors for result counts indicate one or more terms not found.

Count	Database Name	Description
50871	PubMed	biomedical literature citations and abstracts
13510	PubMed Central	free, full text journal articles
45	Site Search	NCBI web and FTP sites
490	Books	online books
126	OMIM	online Mendelian Inheritance in Man
none	OMIA	Online Mendelian Inheritance in Animals
167470	Nucleotide	sequence database (GenBank)
168532	Protein	sequence database
5	Genome	whole genome sequences
680	Structure	three-dimensional macromolecular structures
1	Taxonomy	organisms in GenBank
7742	SNP	single nucleotide polymorphism
1098	Gene	gene-centered information
469	HomoloGene	eukaryotic homology groups
7	PubChem Compound	unique small molecule chemical structures
833	PubChem Substance	deposited chemical substance records
2	Genome Project	genome project information
145	UniGene	gene-oriented clusters of transcript sequences
6	CDD	conserved protein domain database
2624	3D Domains	domains from Entrez Structure
130	UniSTS	markers and mapping data
1813	PopSet	population study data sets
5561	GEO Profiles	expression and molecular abundance profiles
13	GEO DataSets	experimental sets of GEO data
1	Cancer Chromosomes	cytogenetic databases
126	PubChem BioAssay	bioactivity screens of chemical substances
none	GENSAT	gene expression atlas of mouse central nervous system
251	Probe	sequence-specific reagents
none	Journals	detailed information about the journals indexed in PubMed and other Entrez databases
237	NLM Catalog	catalog of books, journals, and audiovisuals in the NLM collections
66	MeSH	detailed information about NLM's controlled vocabulary

■ - Result counts displayed in gray indicate one or more terms not found

De esta manera es posible saber qué información existe para nuestro término de búsqueda en todo el sitio web del NCBI (ej, 168532 entradas de proteínas, 1 entrada en la sección de taxonomía y 167470 entradas de nucleótidos).

También es posible acceder directamente al sitio web de ENTREZ a través de la siguiente dirección: <http://www.ncbi.nlm.nih.gov/gquery/gquery.fcgi>

ENTREZ es, de manera sencilla, el sistema que mantiene unida toda la información presente en el NCBI, algo así como el "GOOGLE" del NCBI, y es quien realiza la búsqueda de nuestro término a través de todas las bases de datos presentes en el NCBI.

A continuación presione el hipervínculo de la sección Genome. Espere unos segundos. Se encontrará con una página de resultados similar a la imagen a su izquierda.

Por el momento dejaremos esta búsqueda ahí y la retomaremos más adelante.

Regrese a la página principal del NCBI: <http://www.ncbi.nlm.nih.gov/> Realice de nuevo la búsqueda por **HIV-1**, pero esta vez asegúrese de escoger la sección “**genome**” y no “**all databases**”. Espere unos segundos y analice la página de resultados que obtiene.

Seguramente ya se ha percatado de que la página de resultados es idéntica a la que se obtuvo mediante el vínculo “**genome**” de la primera búsqueda, hecha en el sistema ENTREZ.

Realice nuevamente la búsqueda en el sistema ENTREZ y explore las diversas entradas que muestra la página de resultados (ej, Protein, UniGene, OMIM, Pubmed etc.). Corrobore dichos resultados con los que arrojan las búsquedas con las opciones en el menú desplegable del sitio web del NCBI.

## Accediendo a las secuencias

Ya que ha experimentado con las diferentes bases de datos que ofrece el NCBI y la manera más común de realizar búsquedas en ellas, es momento de conocer la manera en que podemos acceder a los datos que queremos obtener con nuestra búsqueda.

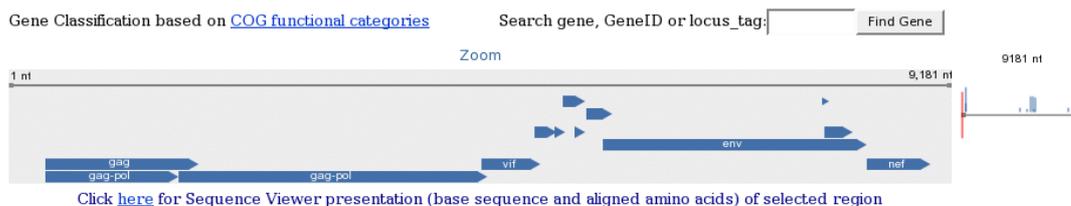
Realice nuevamente la búsqueda de HIV-1 en la sección genome. Encontrará 5 entradas acompañadas de una breve descripción. Siga el hipervínculo para la entrada con identificador: **NC\_001802**

En este momento debe encontrarse en una página web cuyo contenido es similar a la siguiente gráfica:

[Genome](#) > [Viruses](#) > *Human immunodeficiency virus 1, complete genome*

**Lineage:** [Viruses](#) ; [Retro-transcribing viruses](#) ; [Retroviridae](#) ; [Orthoretrovirinae](#) ; [Lentivirus](#) ; [Primate lentivirus group](#) ; [Human immunodeficiency virus 1](#)

Genome Info:	Features:	BLAST homologs:	Links:	Review Info:
Refseq: <a href="#">NC_001802</a>	Genes: <a href="#">9</a>	COG	<a href="#">Genome Project</a>	Publications: <a href="#">[1]</a>
GenBank: <a href="#">AF033819</a>	Protein coding: <a href="#">9</a>	3D Structure	<a href="#">Refseq FTP</a>	Refseq Status: <b>Reviewed</b>
Length: <b>9,181 nt</b>	Structural RNAs: <b>None</b>	TaxMap	GenBank FTP	Seq.Status: <b>Completed</b>
GC Content: <b>42%</b>	Pseudo genes: <b>None</b>	TaxPlot	<a href="#">BLAST</a>	Sequencing center: <b>NLM, NIH, USA, Bethesda</b>
% Coding: <b>93%</b>	Others: <b>7</b>	GenePlot	TraceAssembly	Completed: <b>1998/01/22</b>
Topology: <b>linear</b>	Contigs: <a href="#">1</a>	<a href="#">gMap</a>	CDD	<a href="#">Organism Group</a>
Molecule: <b>ssRNA</b>			Other genomes for species: <a href="#">843</a>	



El cuadro que observa resume la información relacionada con el genoma que hemos buscado. Gracias a éste sabemos que cuenta con 9 genes, que codifican 9 proteínas y que su longitud es de 9181 nucleótidos. Entre otras cosas.

Este tipo de resumen es necesario cuando tratamos de acceder a este tipo de información, es decir si lo que buscamos es simplemente una proteína o secuencia de ADN, por lo general no seremos llevados a un cuadro de resumen como éste sino directamente a la entrada de dicha secuencia.

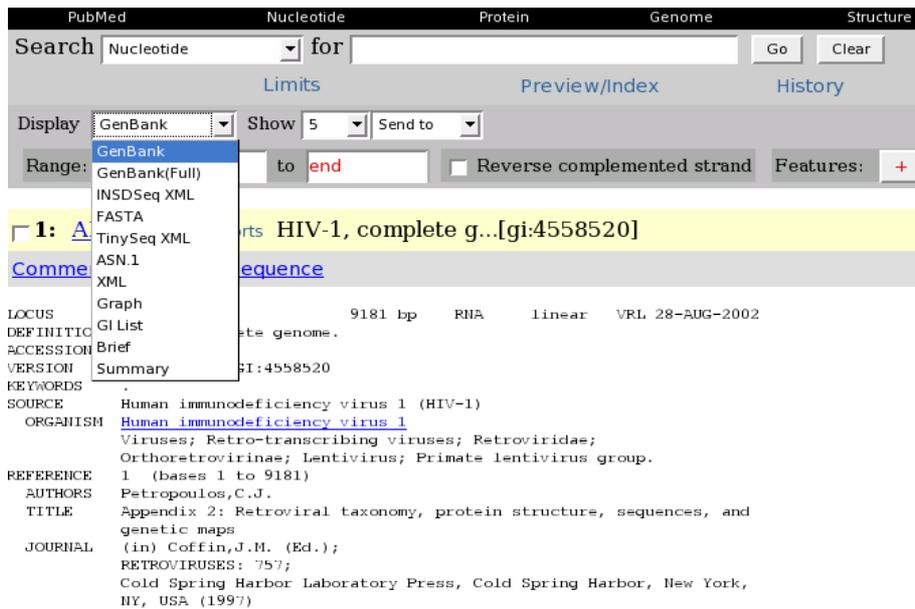
## Estudiando la entrada para HIV-1

Además de nuestro interés por conocer algunas características del genoma que consultamos, nos resulta interesante obtener también su secuencia completa, para poder acceder a dicha información tenemos que consultar la entrada en genBank para dicho genoma.

Esto se hace siguiendo el hipervínculo al genBank: **AF033819**.

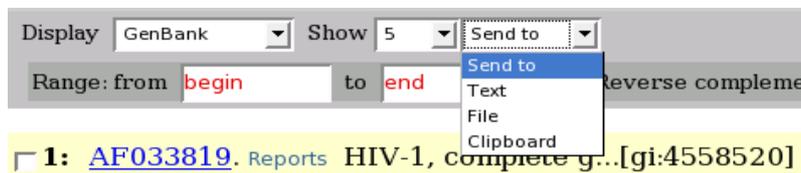
Es importante anotar la diferencia que existe entre el NCBI en general y el genBank, como podemos verlo el genBank es una de las bases de datos del NCBI, si siguiéramos los enlaces a proteínas seríamos llevados a la sección de proteínas del NCBI, cuya información proviene de las bases de datos de UNIPROT.





Explore cada uno de los formatos del menú desplegable. ¿Qué diferencias y semejanzas encuentra en cada uno de ellos? Preste especial atención al formato FASTA. ¿Por qué razón cree que este es el formato más usado en bioinformática?

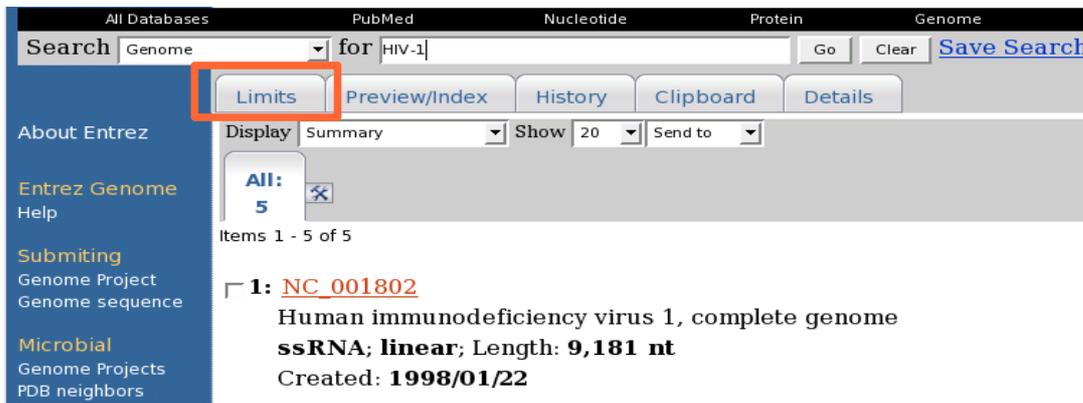
Justo al lado derecho del menú desplegable para los diferentes formatos existe otro menú, el cual nos permite elegir el lugar al que queremos enviar nuestros resultados.



Explore cada una de las opciones allí mostradas. Describa las diferencias que encuentra.

# Sequence Retrieval System

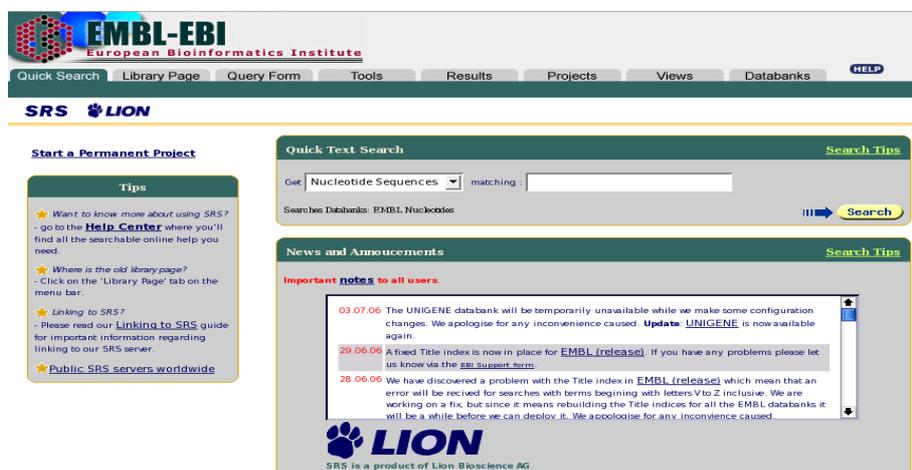
A pesar de todas las virtudes del sistema ENTREZ del NCBI, algunas veces la búsqueda de información allí puede tornarse tediosa y casi imposible, esto se debe principalmente a las pocas posibilidades que este sistema ofrece para filtrar los resultados, la cual se hace por medio de la opción “**limits**”, cuyas opciones en realidad son un poco “limitadas”.



Existe sin embargo una alternativa excelente para la búsqueda de secuencias biológicas, que nos permite controlar casi todos los aspectos de nuestra búsqueda, esta alternativa es el **Sistema de Recuperación de Secuencias (SRS)**. Este sistema fue desarrollado teniendo en mente precisamente esta labor de recuperar secuencias biológicas de una manera efectiva, de allí su diseño y sus capacidades.

En este taller trabajaremos con el SRS ofrecido por el **Instituto Europeo de Bioinformática (EBI)**, cabe anotar que existen muchos servidores SRS alrededor del mundo que ofrecen sus servicios de manera gratuita (ej., el servidor srs del **CBIB**: <http://srs.ibun.unal.edu.co:8080/srs81/>).

Visite la siguiente URL: <http://srs.ebi.ac.uk>



Una manera sencilla de consultar el SRS es mediante la casilla Quick Text Search. En dicha casilla es posible realizar búsquedas en diversas bases de datos disponibles en el menú desplegable.

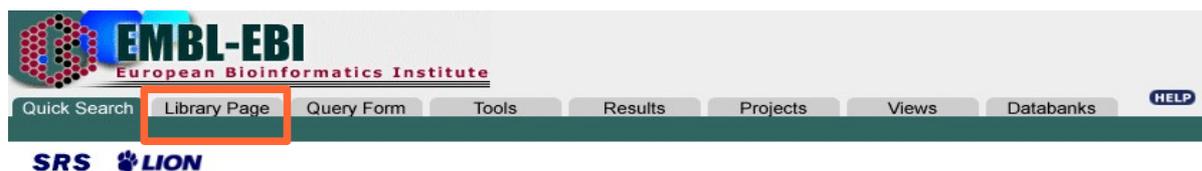


Por ejemplo seleccionando la opción “**Nucleotide Sequences**” realizaremos nuestra búsqueda en la base de datos de DNA EMBL (homóloga al genBank y al DDBJ).

Realice la búsqueda por HIV-1 con diferentes opciones del menú desplegable.

Hasta este momento el SRS parece ser bastante menos completo el sitio web del NCBI, pero ahora empezaremos a comprobar donde radica todo su potencial.

Seleccione la pestaña Library Page ubicada en la parte superior de su pantalla



A continuación será llevado a la sección del SRS donde se describen cada una de las bases de datos que componen el sistema. Como puede ver el SRS comprende muchas bases de datos a la vez y esa es una de sus principales virtudes, por esta razón al SRS se le conoce algunas veces como una “base de datos de bases de datos”, pues a través de este sistema podemos consultar múltiples bases de datos al mismo tiempo, de acuerdo a nuestras necesidades particulares.

Como puede darse cuenta el SRS es similar al sistema ENTREZ del NCBI, en el sentido en que nos permite consultar muchas bases de datos al mismo tiempo, pero esta vez no restringidos únicamente a aquellas con las que cuenta el NCBI sino a virtualmente cualquier base de datos.

El número de bases de datos con las que cuenta el SRS depende de cada implementación, es decir el administrador del SRS determina qué bases de datos quiere o no incluir en su sistema.

- ☐ **Literature, Bibliography and Reference Databases**
  - [all](#)  [TAXONOMY](#)  [GENETICCODE](#)  [OMIM](#)  [MEDLINE](#)
  - [Patent Abstracts](#)  [Karyn's Genomes](#)
  - Literature, Bibliography and Reference Databases - subsections**
  - [MEDLINE \(Updates\)](#)  [MEDLINE \(Main Release 2006\)](#)  [MED2PUB](#)
- ☐ **Gene Dictionaries and Ontologies**
- ☐ **Nucleotide sequence databases**
  - [all](#)  [EMBL](#)  [PATENT\\_DNA](#)  [IMGTLIGM-DB](#)  [IMGTHLA](#)
  - [IPD-KIR](#)  [EMBL \(Contig\)](#)  [Genome Reviews](#)  [EMBL \(Contigs expanded\)](#)
  - [RefSeq Genome DB](#)  [LiveLists](#)  [EMBL ID/Accession Mapping](#)
  - Nucleotide sequence databases - subsections**
  - [all](#)  [EMBL \(Updates\)](#)  [EMBL \(Release\)](#)  [EMBL \(Whole Genome Shotgun\)](#)
  - [EMBL \(Whole Genome Shotgun release\)](#)  [EMBL \(Whole Genome Shotgun release\)](#)
  - [EMBL \(Contig release\)](#)  [EMBL \(Contig release\)](#)
  - [EMBL \(Contigs expanded updates\)](#)  [EMBL \(Annotated Contigs updates\)](#)
  - [EMBL \(Annotated Contigs updates\)](#)  [RefSeq Genomes](#)
- ☐ **Nucleotide related databases**
- ☐ **UniProt Universal Protein Resource**
  - [all](#)  [UniProtKB](#)  [UniProtKB/Swiss-Prot](#)  [UniProtKB/TrEMBL](#)  [UniRef100](#)  [UniRef90](#)
  - [UniRef50](#)  [UniParc](#)
- ☐ **Other protein sequence databases**
- ☐ **Protein function, structure and interaction databases**
- ☐ **Enzymes, reactions and metabolic pathway databases**
- ☐ **Mutation and SNP databases**
- ☐ **Biological Resources Catalogues (CABRI)**
- ☐ **Mapping databases**
- ☐ **Other databases**
- ☐ **User owned databases**
- ☐ **Application result databases**
- ☐ **EMBOSS result databases**
- ☐ **Eurofir Food data**
- ☐ **EMBLCDS Grouped By**

The latest release of the EMBL nucleotide sequence database. Releases are on a quarterly basis. Updates to this release can be found in EMBL(Updates).  
*To obtain comprehensive information on this databank, click the link*

Póse el cursor del mouse por alguna de las entradas, después de unos segundos una casilla de texto explicativo aparecerá. ¿Qué tipo de información proveen las bases de datos EMBL (Contig Updates), UniprotKB/Swissprot?

Al seguir el enlace a cualquiera de estas bases de datos obtendremos mayor información acerca de esta, como el número de entradas presentes, fecha de actualización etc. Sin embargo, por ahora nuestro interés es el de seleccionar algunas bases de datos para realizar nuestras búsquedas.

Seleccione las casillas pertenecientes a las bases de datos de “**UniprotKB/Swissprot**” y “**UniprotKB/TrEMBL**”. Cerciorese de que estas sean las únicas bases de datos seleccionadas.

A la izquierda de su pantalla encontrará la casilla “Search Options” la cual nos permitirá seleccionar el nivel de profundidad de nuestra búsqueda, Por ser esta la primera vez que trabajamos con este sistema seleccionaremos la forma estándar de búsqueda.

Presione el botón “**Standard query Form**” de la casilla “**Search Options**”.



Esta acción le llevará al formulario estándar de búsqueda en el SRS.

The screenshot shows the SRS search interface with four numbered annotations:

- 1**: Points to the search input fields under the "Your search terms" section, where multiple search terms can be entered.
- 2**: Points to the "Search Options" section, which includes a dropdown for "with: & (AND)", a checked "Use wildcards" checkbox, and a "Get results of type:" dropdown set to "Entry".
- 3**: Points to the "Result Display Options" section, which includes a "View results using:" dropdown set to "UniprotView", a "Create a view" option, and a "Show 30 results per page" dropdown.
- 4**: Points to the "Create a view" section, which includes a "Choose 1 or more fields:" list (ID, EntryName, AccessionNumber, Creation Date, Seq Mod Date, Annot Mod Date, Description), a "Display As:" dropdown set to "Table", and a "Sequence Format:" dropdown set to "SWISS".

El cual consta de 4 partes fundamentales.

1. **Campos de búsqueda**, donde podemos entrar nuestros términos de búsqueda de acuerdo a cualquiera de las opciones presentes en los respectivos menús desplegables.
2. **Opciones de búsqueda**, donde podemos definir, entre otras cosas, el tipo de conector lógico (booleano) a utilizar para los términos definidos en 1.
3. **Opciones para mostrar los resultados**, donde podemos definir el número de resultados que queremos por página, así como el formato de salida, ya sea alguno de los definidos en el menú desplegable o mediante la creación de una vista personalizada (opción “**create view**”).
4. **Crear vista**, esta opción trabaja en conjunto con la opción 3, y acá podemos definir el tipo de campos que queremos ver en nuestra página de resultados.

Para nuestro ejemplo, tenemos interés en seleccionar todas las proteínas de superficie conocidas de *Plasmodium falciparum* con actividad inmunogénica, relacionadas con el merozoito.

Defina estos criterios en la sección “**campos de búsqueda**” de acuerdo a la siguiente imagen:

Fields you can search	Your search terms
In a single field, you can separate multiple values by &,  , !	
	
<b>i</b> Organism Name	plasmodium falciparum
<b>i</b> Keywords	merozoite
<b>i</b> Description	surface antigen
<b>i</b> AllText	

A continuación presione el botón “**search**” ubicado en la parte superior de esta sección y espere unos segundos.

Seguramente en este momento ya tenga una visión más exacta de las posibilidades que ofrece el SRS y sus principales diferencias con el sistema ENTREZ. Primero, pudimos definir exactamente no solamente la base de datos que queríamos consultar, sino las secciones específicas de esta. Además de esto pudimos también definir exactamente los términos de búsqueda en secciones específicas de las entradas, lo cual nos da un completo control sobre los resultados que queremos obtener.

Cree usted que existe alguna manera de realizar esta misma consulta en el sistema ENTREZ?

Juegue un poco con las diferentes opciones de formatos que ofrece el SRS en la sección **3**, del formulario de búsqueda. Intente también creando su propio formato de salida con la opción **create view** y la sección **4**.

Encuentre todas las proteínas nucleares hipotéticas de SACCHAROMYCES CEREVISIAE, y muestre la información en formato fasta.

# BLAST: Basic local Alignment Search Tool

BLAST es un **algoritmo** para comparación (alineamiento) de secuencias. Más exactamente se encuentra clasificado dentro de los algoritmos para alineamiento local.

Existen varias “implementaciones” de este algoritmo, una de las más conocidas es la realizada por el NCBI, el **NCBI-BLAST**.

Otra implementación muy conocida de este algoritmo es la realizada por la Universidad de Washington el **WU-BLAST** (<http://blast.wustl.edu/>).

Es importante notar que el NCBI-BLAST **no** se utiliza cuando estamos realizando búsquedas de secuencias por palabras clave o términos de búsqueda. Este se utiliza cuando estamos buscando secuencias similares a la nuestra en las diferentes bases de datos del NCBI. Entonces, existe una clara diferencia entre consultar el NCBI para **obtener una o varias secuencias** requeridas y consultarlo para **buscar coincidencias** de nuestra secuencia con otras.

El NCBI-BLAST es el programa/algoritmo que usa por defecto el NCBI para realizar búsquedas de secuencias en sus bases de datos. En esta sección trabajaremos con dicha implementación.

Visite el sitio web del NCBI-BLAST: <http://www.ncbi.nlm.nih.gov/BLAST/>

The screenshot shows the NCBI BLAST website interface. At the top, it says "NCBI + BLAST" and "Latest news: 7 May 2006 : BLAST 2.2.14 released". The main content area is divided into several sections:

- About:** The Basic Local Alignment Search Tool (BLAST) finds regions of local similarity between sequences. The program compares nucleotide or protein sequences to sequence databases and calculates the statistical significance of matches. BLAST can be used to infer functional and evolutionary relationships between sequences as well as help identify members of gene families.
- Nucleotide:**
  - Quickly search for highly similar sequences (megablast)
  - Quickly search for divergent sequences (discontiguous megablast)
  - Nucleotide-nucleotide BLAST (blastn)
  - Search for short, nearly exact matches
  - Search trace archives with megablast or discontiguous megablast
- Protein:**
  - Protein-protein BLAST (blastp)
  - Position-specific iterated and pattern-hit initiated BLAST (PSI- and PHI-BLAST)
  - Search for short, nearly exact matches
  - Search the conserved domain database (rpsblast)
  - Protein homology by domain architecture (cdart)
- Translated:**
  - Translated query vs. protein database (blastx)
  - Protein query vs. translated database (tblastn)
  - Translated query vs. translated database (tblastx)
- Genomes:**
  - Human, mouse, rat, chimp, cow, pig, dog, sheep, cat
  - Chicken, puffer fish, zebrafish
  - Fly, honey bee, other insects
  - Microbes, environmental samples
  - Plants, nematodes
  - Fungi, protozoa, other eukaryotes
- Special:**
  - Search for gene expression data (GEO BLAST)
  - Align two sequences (bl2seq)
  - Screen for vector contamination (VecScreen)
- Meta:**
  - Retrieve results

On the left side, there is a navigation menu with the following categories:

- About:** Getting started, News, FAQs
- More info:** NAR 2004, NCBI Handbook, The Statistics of Sequence Similarity Scores
- Software:** Downloads, Developer info
- Other resources:** References, NCBI Contributors, Mailing list, Contact us

realizar una búsqueda BLAST puede resumirse en 3 sencillos pasos:

1. Seleccionar el tipo de programa BLAST a usar (blastp, blastn, blastx, tblastx,tblastn).
2. Introducir nuestra secuencia pregunta (**query sequence**, en términos BLAST).
3. Seleccionar la base de datos en la que queremos buscar.

Opcionalmente podemos controlar la salida de los resultados, modificando algunas de las opciones de salida.

## Selección del programa BLAST

Revise cuidadosamente las tres primeras secciones de los programas BLAST (Nucleotide, Protein y Translated). Describa la funcionalidad de cada uno de ellos, de acuerdo a las descripciones de los mismos<sup>1</sup>.

Siga el enlace a “**blastp**” y digite el identificador: **NP\_057849** en la casilla de texto **Search**.

The screenshot shows the top part of the BLAST search interface. At the top, there is a text input field containing the sequence identifier 'NP\_057849'. Below this field is a blue button labeled 'Search', which is highlighted with a red rectangular box. Underneath the search box, there are several options: a link 'Set subsequence' followed by 'From:' and 'To:' input fields; a 'Choose database' dropdown menu currently set to 'nr'; a 'Do CD-Search' checkbox which is checked; and a 'Now:' section with three buttons: 'BLAST!' (highlighted in blue), 'Reset query', and 'Reset all'.

The screenshot shows the 'Options for advanced blasting' section. It includes several settings: a 'Limit by' dropdown menu set to 'entrez query' with a note 'or select from: All organisms'; a 'Compositional adjustments' dropdown menu set to 'No adjustment'; a 'Choose filter' section with three checkboxes: 'Low complexity' (checked), 'Mask for lookup table only' (unchecked), and 'Mask lower case' (unchecked); and an 'Expect' value input field set to '10'.

<sup>1</sup> El siguiente enlace puede ayudarle a comprender más dichas descripciones:  
<http://www.ncbi.nlm.nih.gov/blast/producttable.shtml#pstab>

Existen diversas formas de ingresar nuestra secuencia para realizar búsquedas BLAST. Una de ellas es digitando un identificador conocido por el NCBI. Otra manera de hacerlo es ingresando directamente la secuencia en esta misma casilla, ya sea “cruda” o en formato FASTA.

El menú desplegable “**choose database**”, permite seleccionar alguna de las bases de datos permitidas para nuestro tipo de búsqueda. Seleccione **Swissprot**.

La opción “**Do-CD-Search**” que viene por defecto seleccionada le dice a BLAST que también realice una búsqueda de **Dominios Comunes** para dicha proteína. Utilizaremos las opciones (**Options**) por defecto.

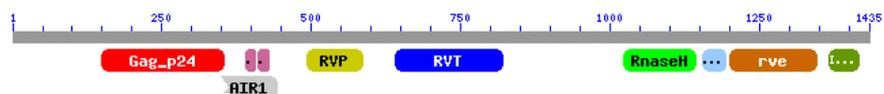
Presione el Botón **Blast** y espere unos segundos.

Esta acción le llevará a una página intermedia entre la página de resultados y el formulario de consulta, es una página de “formateo” de resultados, en la cual además es posible ver los dominios comunes encontrados para nuestra secuencia. No nos preocuparemos mucho por el formato de salida, y dejaremos las opciones por defecto, así que...

Your request has been successfully submitted and put into the Blast Queue.

Query = gj|28872819 (1435 letters)

Putative conserved domains have been detected, click on the image below for detailed results.



The request ID is 1152137026-7229-22113003689.BLASTQ4

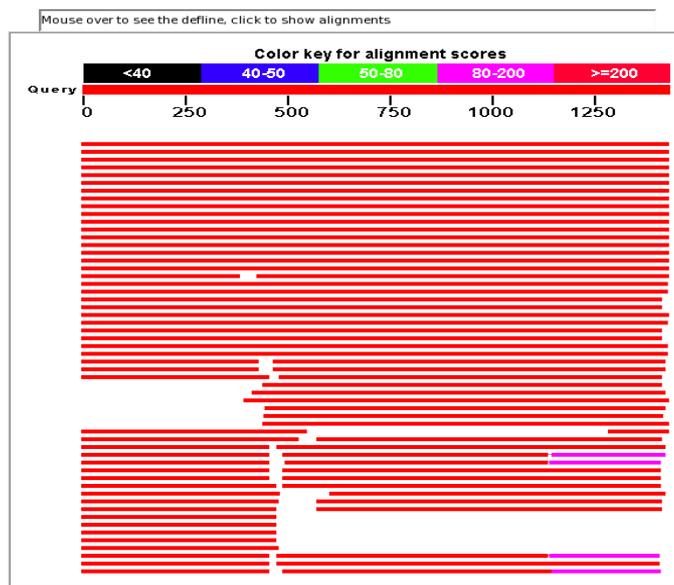
[Format!](#) or [Reset all](#)

...presione el botón **Format** y espere unos segundos (algunas veces es necesario esperar un poco más).

En este momento nuestra secuencia está siendo comparada contra cada una de las entradas en las bases de datos que escogimos (3,782,570 secuencias). En unos momentos aparecerá una página de resultados con las secuencias que BLAST ha encontrado que son muy similares (o idénticas) a nuestra secuencia.

La primera parte de la página de resultados muestra la siguiente gráfica:

**Distribution of 114 Blast Hits on the Query Sequence**



Esta nos permite ver la distribución de los alineamientos. Cada banda debajo del mapa representa una secuencia de la base de datos que resultó ser muy similar a la secuencia de búsqueda. De esta manera podemos ver la extensión de los alineamientos.

Resulta evidente que para la secuencia que hemos elegido BLAST encontró varias secuencias idénticas a la nuestra, por esta razón predomina el color rojo (notar el código de colores en el mapa).

La segunda sección corresponde a las descripciones de los alineamientos.

<a href="#">gi 120836 sp P18800 GAG_HV1ND</a>	Gag polyprotein (Pr55Gag) [Cont...	<a href="#">786</a>	0.0
<a href="#">gi 120829 sp P12495 GAG_HV1Z2</a>	Gag polyprotein (Pr55Gag) [Cont...	<a href="#">783</a>	0.0
<a href="#">gi 120838 sp P04594 GAG_HV1MA</a>	Gag polyprotein (Pr55Gag) [Cont...	<a href="#">764</a>	0.0
<a href="#">gi 120843 sp P24736 GAG_HV1U4</a>	Gag polyprotein (Pr55Gag) [Cont...	<a href="#">749</a>	0.0
<a href="#">gi 120844 sp P05889 GAG_HV1W2</a>	Gag polyprotein (Pr55Gag) [Cont...	<a href="#">727</a>	0.0
<a href="#">gi 120885 sp P17282 GAG_SIVCZ</a>	Gag polyprotein [Contains: Core...	<a href="#">689</a>	0.0
<a href="#">gi 120890 sp P19504 GAG_SIVSP</a>	Gag polyprotein [Contains: Core pr	<a href="#">535</a>	5e-151
<a href="#">gi 120853 sp P04590 GAG_HV2RO</a>	Gag polyprotein (Pr55Gag) [Cont...	<a href="#">529</a>	3e-149
<a href="#">gi 120854 sp P20874 GAG_HV2ST</a>	Gag polyprotein (Pr55Gag) [Cont...	<a href="#">527</a>	1e-148
<a href="#">gi 120889 sp P12496 GAG_SIVS4</a>	Gag polyprotein [Contains: Core pr	<a href="#">526</a>	2e-148
<a href="#">gi 120848 sp P24106 GAG_HV2CA</a>	Gag polyprotein (Pr55Gag) [Cont...	<a href="#">524</a>	8e-148
<a href="#">gi 120850 sp P18041 GAG_HV2G1</a>	Gag polyprotein (Pr55Gag) [Cont...	<a href="#">521</a>	5e-147
<a href="#">gi 2495241 sp Q74119 GAG_HV2KR</a>	Gag polyprotein (Pr55Gag) [Con...	<a href="#">521</a>	7e-147
<a href="#">gi 120849 sp P15832 GAG_HV2D2</a>	Gag polyprotein (Pr55Gag) [Cont...	<a href="#">518</a>	4e-146
<a href="#">gi 120852 sp P05891 GAG_HV2NZ</a>	Gag polyprotein (Pr55Gag) [Cont...	<a href="#">517</a>	1e-145
<a href="#">gi 120847 sp P18095 GAG_HV2BE</a>	Gag polyprotein (Pr55Gag) [Cont...	<a href="#">516</a>	2e-145
<a href="#">gi 399525 sp P31634 GAG_SIVMS</a>	Gag polyprotein [Contains: Core pr	<a href="#">516</a>	2e-145
<a href="#">gi 120846 sp P17756 GAG_HV2D1</a>	Gag polyprotein (Pr55Gag) [Cont...	<a href="#">516</a>	2e-145
<a href="#">gi 120851 sp P12450 GAG_HV2SB</a>	Gag polyprotein (Pr55Gag) [Cont...	<a href="#">514</a>	8e-145
<a href="#">gi 120887 sp P05894 GAG_SIVMI</a>	Gag polyprotein [Contains: Core pr	<a href="#">510</a>	2e-143
<a href="#">gi 130597 sp P16088 POL_FIVPE</a>	Pol polyprotein [Contains: Prot...	<a href="#">507</a>	1e-142
<a href="#">gi 120884 sp P05892 GAG_SIVVT</a>	Gag polyprotein [Contains: Core...	<a href="#">507</a>	1e-142

Esta es una lista de las secuencias encontradas (ordenadas de acuerdo a su valor E), en cuatro columnas: identificador, breve descripción de la secuencia, bit score y valor E.

La tercera sección corresponde a los alineamientos hechos por BLAST de nuestra secuencia pregunta (**query**) y la secuencia encontrada (**subject**).

Estos alineamientos son mostrados en un formato convencional de alineamiento pareado, mostrando a demás el porcentaje de identidad, el valor E del alineamiento, el número de gaps y el Score.

```
> gi|400822|sp|P31822|POL\_F1V2 Pol polyprotein [Contains: Protease (Retropeps
H (RT); Deoxyuridine 5'-triphosphate
nucleotidohydrolase (dUTPase); Integrase (IN)]
Length=1124

Score = 494 bits (1273), Expect = 7e-139
Identities = 286/673 (42%), Positives = 404/673 (60%), Gaps = 39/673 (5%)

Query 492 TLWQRPLVTIKIGGQKLEALLDGTGADDTVLEEMSLP-----GRWKFPMIGGIGGFIVKVRQ 546
          TL +RP + I + G + LLDGTGAD T+L K MIG +GG +
Sbjct 46 TLERRPEIQIFVNGHPITKFLDGTGADITILNRKDFQIGNSIENGKQNMIG-VGGGKRGTH 104

Query 547 YDQLLIEICGHKAIGTVLVGPTPV-----NIIGRNLLTQIGCTLNFP--ISPIETVP 596
          Y + +EI + G V ++GR+ + + L I V
Sbjct 105 YINVHLEIRDENYRNQCIFGNVVCVLEDNSLIQPLLRDNIKIFNIRLVMQAISEKIPIVK 164

Query 597 VKLKPMDGPKVKQWPLTEEKIKALVEICTEMEKEGKISKIGPENPYNTPVFAIKKKDST 656
          V++K GP+VKQWPL+ EKI+AL +I +E EGK+ + P NP+NTPVFAIKKK S
Sbjct 165 VRMKDPTQGFQVQWPLSNEKIEALTDIVERLESEGKVKRADPHNFWNTPVFAIKKK-SG 223

Query 657 KWRKLVDFRELNRKTQDFWEVQLGIPHPAGLKKKKSVTVLVDVGDAYFVPLDEDFRKYTA 716
          KWR L+DFR LNK T EVQLG+PHPAGL+ KK VTVLD+GDAYF++PLD D+ YTA
Sbjct 224 KWRMLIDFRVLNKLTDKGAEVQLGLPHPAGLQMKKQVTVLDIGDAYFTIPLDPPYAPYTA 283

Query 717 FTIPSINNETPGIRYQYNVLPQGWGSPALFQSSMTIKILEFFRKQNPEDIVIYQYMDLIV 776
          FT+P NH PG RY + LPQGW SP I+QS++ IL+PF KQN ++ IYQYMD+Y+
Sbjct 284 FTLPRKNNAGPGRRYVWC SLPQGWVLSPLIYQSTLNNIIQPFIRKQNSLIDYQYMDDIYI 343

Query 777 GSDLEIGQHRTKIEELRQHLRLRWGLTTPDKKHQKEPPFLWNGYELHPDKWTVQPIVLPEK 836
          GS+L +H+ K+EELR+ LL WG TP+ K Q+EPP+ WNGYELHP W++Q L
Sbjct 344 GSNLHKKEHKQKVEELRKLRLRWGFETPEDKLQEEPPYKNGYELHPLIWSIQKQLEIP 403

Query 837 DSWTVNDIQKLVGKLNWASQIYPGIVRQLCKLLRGTALTEVIPLT====leleNREI 896
          + T+H++QKL GK+HWASQ P + +++L ++RG + L + T EA+ E+ + +E
Sbjct 404 ERPTLNEIQKLAGKLNWASQTIPDLSIKELTNMIRGDQKLDSTREWTVEAKREVQKAKEA 463

Query 897 LKEPVHGVYDPSKDLIAEIQKQGGQWNTYQIQE--PFFKILKTKGYARMRGAAHTNDVKQL 955
          ++ YDPP++ L A++ G Q YQ+YQ+ P L GK R + N
Sbjct 464 IETQAQLNYYDPNRLYAKLSLVGPHQIC YQVYQKNPEHILWYGKINRQKKAENTCDIA 523

Query 956 TEAVQKITTESIVIVGKTPKFKLP IQE+twetwte YWQATWI-----PEWEFVNT 1006
          A KI ESI+ GK P +++P +E W++ I PE EF++
Sbjct 524 LRACYKIREESIIRIGKEPVYEIPASREA-----WESNLIRSPYLKAPPPEVEFIIHA 575
```

Al seguir cada uno de los links de las secuencias alineadas seremos llevados a la página de resultados para dicha secuencia.

Seleccione las casillas de selección de las cuatro primeras secuencias mostradas y a continuación presione el botón “Get selected sequences” que se encuentre al inicio de los alineamientos. ¿Qué obtiene al realizar esta acción?

Explore las otras posibilidades allí ofrecidas. Especialmente la opción “tree view”. ¿Qué utilidad cree usted que puede tener esta opción?

Realice nuevamente esta búsqueda, pero esta vez modifique el formato en la página intermedia, seleccionando la opción “**pairwise with identities**” en el menú desplegable en “**alignment view**”. ¿Qué diferencias y similitudes encuentra en el formato de salida con respecto al parámetro por defecto: **Pairwise**?

**Format**

Show  [Graphical Overview](#)  [Linkout](#)  [Sequence Retrieval](#)  [NCBI-g](#) Alignment  in  [format](#)

[CDS feature](#)

[Masking Character](#)  [Masking Color](#)

Number of: [Descriptions](#)  [Alignments](#)  [Graphic overview](#)

[Alignment view](#)

[Format for PSI-BLAST](#)

[Limit results by entrez query](#)

[Expect value range:](#)

Guía elaborada por Andrés M. Pinzón V., del **Centro de Bioinformática** del Instituto de Biotecnología en la Universidad Nacional de Colombia y está distribuida bajo licencia:



Bogotá Colombia - Julio de 2006.

**Cualquier sugerencia o inquietud dirigirla a:**  
[ampinzonv@unal.edu.co](mailto:ampinzonv@unal.edu.co) ó [andrespinzon@gmail.com](mailto:andrespinzon@gmail.com)