



Programa de la asignatura:

Bioinformática

U1

Introducción a la bioinformática



DCSBA



BIOTECNOLOGÍA



Índice

Presentación de la Unidad	3
Propósitos de la unidad.....	4
Competencia específica	4
1.1. Genómica funcional.....	5
1.1.1. Dogma central de la Biología Molecular	4
1.1.2. Secuenciación genómica masiva.....	6
1.1.3. La era de las “Omicas”	9
1.1.4. Concepto de Bioinformática	10
1.2. Recursos disponibles en internet	11
1.2.1. Bases de datos: utilidad, concepto y clasificación	11
1.2.2. Secuencias de ADN o de proteínas como datos biológicos	13
1.2.3. La base de datos del NCBI	13
Actividades	28
Autorreflexiones	28
Cierre de la Unidad	29
Para saber más	30
Fuentes de consulta	30



Presentación de la Unidad

Las células son las unidades de construcción de todos los seres vivos. Dentro de cada célula existen dos macromoléculas fundamentales que contienen la información necesaria para definir las características biológicas de los organismos, llamadas por lo tanto, macromoléculas informacionales. Estas macromoléculas son los ácidos nucleicos y las proteínas, los cuales contienen información genética e información funcional, respectivamente. Dicha información está reunida en forma de secuencias de nucleótidos o de aminoácidos en un código determinado.

Para que esta información pueda reflejarse en un fenotipo, tiene que ser procesada en una serie de eventos específicos que tienen cierta dirección. El dogma central de la biología molecular establece la dirección de este flujo de información. El avance científico permitió descifrar la naturaleza de estos códigos. Estos códigos se reflejaron en secuencias de aminoácidos o de nucleótidos que tenían un cierto orden. Paralelamente, también se hizo posible dilucidar la estructura de las macromoléculas, con lo cual, su caracterización y manipulación fue posible, surgiendo la ingeniería genética o metodología del ADN recombinante. Entre las metodologías que sustentan a la ingeniería genética se destaca la técnica de secuenciación de ácidos nucleicos, que ha permitido conocer la secuencia de genomas completos en un periodo de tiempo relativamente corto. Debido al gran número de datos biológicos que la secuenciación genómica provee, se ha vuelto necesario el perfeccionamiento y uso de operaciones y sistemas basados en tecnologías de la información para almacenar, acceder y analizar este gran cúmulo de información.

La creación de bases de datos ha permitido el almacenamiento y acceso a los datos biológicos.

Así, de la necesidad de operar la gran cantidad de datos biológicos que estaban siendo generados por el avance científico y tecnológico de la biología molecular y la genética en las últimas décadas, surge la Bioinformática. En términos generales, la Bioinformática es un área interdisciplinaria de reciente creación que se encarga del análisis computacional de datos biológicos.



Propósitos de la unidad



- Reconocer el marco histórico que llevó al surgimiento y desarrollo de la Bioinformática.
- Describir la naturaleza de los datos generados y comprender la necesidad de crear bases de datos.
- Abordar el uso de una de las bases de datos de secuencias genéticas más comunes en la actualidad.

Competencia específica



Analizar la utilidad del uso de recursos informáticos actuales para almacenar, acceder y obtener información biológica del gran número de datos generados en los últimos años, a partir del estudio del contexto histórico.



1.1. Genómica funcional

Con el objeto de adentrarnos en el campo de la Bioinformática, será preciso, primeramente, puntualizar algunos conceptos básicos de biología molecular.

1.1.1. Dogma central de la Biología Molecular

El dogma central de la biología molecular constituye el paradigma que establece la dirección del flujo de la información genética para que ésta pueda ser expresada, permitiendo obtener un fenotipo a partir de un genotipo. Determina que en los sistemas celulares, el ADN tiene la capacidad de replicarse y de ser transcrito a ARN, el cual a su vez es traducido a proteínas. Fue establecido por Francis Crick (Ondarza, 1994; Nelson & Cox, 2005).

En asignaturas anteriores (Biología Molecular 1 y Genética Molecular Bacteriana) se describió en detalle la naturaleza de las macromoléculas y procesos involucrados en este flujo de información, por lo que no es el objetivo de esta Unidad describirlos. Sin embargo, es importante destacar, que a este esquema del dogma clásico, se han adicionado algunos procesos particulares que han demostrado que la unidireccionalidad del flujo de la información genética no es absoluta, como anteriormente se creía. Así, hoy en día se sabe que tanto el ADN como el ARN pueden funcionar como material genético, teniendo ambos la capacidad de hacer copias de sí mismos y que se puede obtener ADN a partir de ARN. Dichos descubrimientos se realizaron en virus, donde se demostró que su ARN puede replicarse mediante la acción de la enzima ARN polimerasa dependiente de ARN. Por otro lado, también en virus (específicamente en retrovirus), se determinó que por la acción de la enzima transcriptasa reversa, el ARN puede ser transcrito a ADN, proceso conocido como transcripción reversa o retro-transcripción. Para lograr el proceso de transcripción reversa, las enzimas reverso transcriptasas catalizan tres distintas reacciones: la síntesis de ADN dependiente de ARN; la degradación de ARN; y la síntesis de ADN dependiente de ADN. Estas enzimas, a diferencia de las ADN polimerasas, presentan una alta tasa de error (1 en cada 20,000 nucleótidos adicionados), lo cual genera más altas tasas de mutación y evolución, contribuyendo a la frecuente aparición de cepas nuevas en retrovirus (Nelson & Cox, 2005).

En la Figura 1 se muestra el esquema del dogma central de la biología molecular ampliada donde se compendian tanto los procesos generales que se llevan a cabo en todas las células, como aquéllos característicos de los virus. A pesar de que los fenómenos de replicación y transcripción reversa de ARN están restringidos a virus, su importancia está dada en términos evolutivos, ya que han contribuido a la generación de cambios en el genoma de las células, pues una vez que un virus ha infectado a un



hospedero, éste tiene la capacidad de integrar ADN retro-transcrito a partir de ARN en el genoma de la célula (Nelson & Cox, 2005; Krebs, et al., 2010).

El planteamiento del dogma central de la biología molecular no fue una tarea fácil, sino que sólo tras años de intensa investigación fue que pudo establecerse de manera correcta la dirección del flujo de la información biológica. Importantemente, este hecho sentó una de las piedras angulares que llevó al surgimiento de una nueva rama de la biología, que en un principio fue llamada biología computacional, para después constituir lo que actualmente conocemos como Bioinformática.

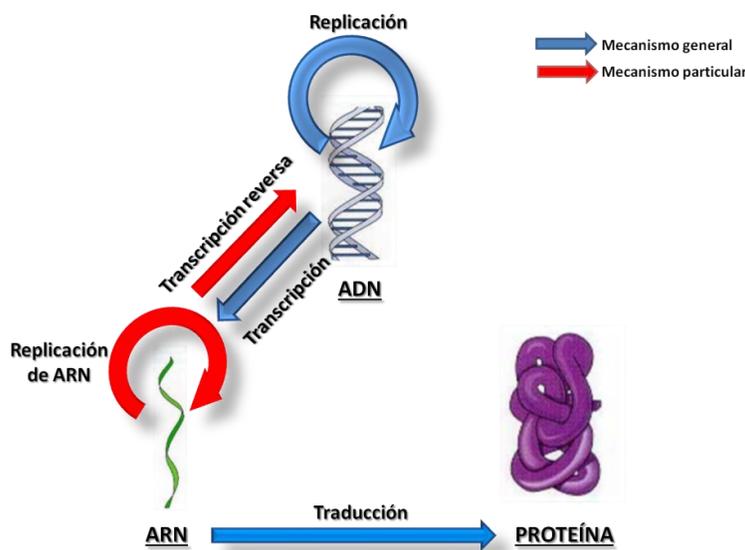


Figura 1. Esquema donde se representa el dogma central de la biología molecular mostrando el flujo de la información genética, señalando los procesos generales exhibidos por todas las células, como aquéllos procesos específicos llevados a cabo por virus.

1.1.2. Secuenciación genómica masiva

Una vez que se estableció que dentro de la célula existen macromoléculas informacionales que determinan su función y evolución, el siguiente paso fue descifrar el orden en que esta información se encontraba contenida. Así, se empezaron a desarrollar metodologías que permitieran conocer tanto la secuencia del ADN o ARN, conformada por nucleótidos, como la de las proteínas, conformada por aminoácidos.

Las técnicas de secuenciación de ADN fueron las que se desarrollaron más rápidamente, siendo la descrita por Sanger (1977) la más usada desde su creación, hace ya más de 30 años. Está basada en la síntesis química y se le conoce también como método de terminación de cadena. Debido a que en la asignatura de Biología Molecular II se



describió con detalle dicha técnica, no se ahondará más sobre este tema aquí. Sin embargo, en la sección “Para saber más” podrás encontrar un recurso que te ayudará a recordar el procedimiento de manera general.

A lo largo del tiempo, dicha técnica experimentó una importante transformación que impactó grandemente en el desarrollo de la biología actual. En sus inicios, la secuenciación Sanger permitió descifrar únicamente genes completos, pero más tarde, tras su optimización y automatización, llevó a la lectura de genomas enteros. Recordemos que el genoma es la colección completa de las secuencias de ADN que conforman a un organismo (Krebs, et al., 2010). Desde el punto de vista biológico, ya no sólo se quería conocer la secuencia de algunos cuantos genes aislados, si no que se pretendía conocer la secuencia de todos los genes que estaban presentes en los organismos. Así, el primer organismo vivo del cual se conoció por primera vez su secuencia completa en el año de 1995 fue el de la bacteria patógena causante de la gripe, *Haemophilus influenzae* (Fleischmann, et al., 1995). Este acontecimiento marcó el nacimiento de la era genómica, pues a partir de ese momento se empezaron a secuenciar rápidamente otros muchos, miles de genomas de distintos organismos y cuyo objetivo principal era conocer la secuencia de nucleótidos completa que los conformaban.

Así, las técnicas de secuenciación fueron optimizadas básicamente y principalmente para producir datos a gran escala, es decir, para poder secuenciar fragmentos de ADN muy grandes, pues debido a la complejidad de los organismos que se pretendía estudiar, era muy probable que estos estuvieran conformados por miles, hasta millones de nucleótidos. De esta forma, con la creación y puesta en marcha del proyecto de secuenciación del genoma humano, la demanda en la producción de secuencias de ADN aumentó considerablemente, por lo que las estrategias de secuenciación fueron optimizadas considerablemente. Pero el genoma humano no fue el único en el que la comunidad científica estaba interesada en secuenciar, sino que muchísimos otros organismos pretendían ser secuenciados, por lo que se necesitaban tener métodos que permitieran una secuenciación más barata y de alto rendimiento o eficiente, creándose entonces la secuenciación de segunda generación, también llamada secuenciación paralela masiva. Constituye una “tecnología que permite acumular información genética tanto de manera cualitativa como cuantitativa de cualquier tipo de ácido nucleico en una muestra dada, con un rendimiento enormemente alto” (Reis-Filho, 2009).

Usando distintas aproximaciones, esta tecnología ha sobrepasado la limitada capacidad de la secuenciación Sanger, permitiendo que millones de reacciones de secuenciación puedan llevarse a cabo de una vez y de forma paralela y no sólo 96, como en la secuenciación por terminación de cadena (Mardis, 2007). En la actualidad, existen en el mercado varias de estas tecnologías, las cuales son enlistadas en la Tabla 1, donde se anotan los datos más sobresalientes de cada tecnología que permiten compararla con la secuenciación Sanger tradicional, quedando de manifiesto la gran diferencia en el número



de secuencias obtenidas con las tecnologías de segunda generación, así como el tiempo requerido, que es mucho menor, obteniéndose con ello, un mayor número de secuencias en un tiempo relativamente corto y adicionalmente, a menor costo .

En los siguientes enlaces encontrarás videos explicativos sobre las metodologías de secuenciación de segunda generación comercialmente disponibles:

Secuenciación 454:

<http://www.my454.com/>

Solexa:

<http://www.youtube.com/watch?v=HMyCgWhwB8E#t=58>

SOLID:

http://www.appliedbiosystems.com/absite/us/en/home/applications-technologies/solid-next-generation-sequencing/next-generation-systems/solid-4-system.html?CID=FL-091411_solid4

Tabla 1. Tecnologías de segunda generación comercialmente disponibles.

Método	Longitud de la lectura (pares de bases)	Templados secuenciados	Producción de datos/día	Reacción de la secuencia	Referencia
ABI 3730xl	900 a 1,100	96	1 Mb/día	Método Sanger	www.appliedbiosystems.com
454 FLX Roche	400	1,000,000	400 Mb/corrida/7-8 h	Pirosecuenciación	www.rockefeller.edu/science.com
Illumina (Solexa)	36 a 175	40,000,000	>17 Gb/corrida/3-6 días	Terminador reverso	www.illumina.com
ABI SOLiD	50	85,000,000	10-15 Gb/corrida/6 días	Secuenciación de ligación	www.appliedbiosystems.com
Helicos Heliscope	30 a 35	800,000,000	21-28 Gb/corrida/8 días	Secuencia de molécula única por síntesis	www.helicosbio.com

Tabla modificada de Reis-Filho, 2009.



Actualmente, se cuenta con 7,581 proyectos de secuenciación terminados y 28,913 por terminar (<http://www.ncbi.nlm.nih.gov/>, 2013). Con esta revolución en las técnicas de secuenciación, uno puede apenas imaginar el gran número de datos que se han obtenido desde la creación de estas plataformas tecnológicas, lo cual puede reflejarse en el número de secuencias que se han depositado en diferentes fuentes de información en los últimos años. Como puede observarse en la gráfica de la Figura 2, el número de pares de bases y de secuencias que se depositaron en 10 años en centros de información especializados, aumentaron en número de manera considerable, duplicando el número de entradas aproximadamente cada 18 meses (Posada, 2009).

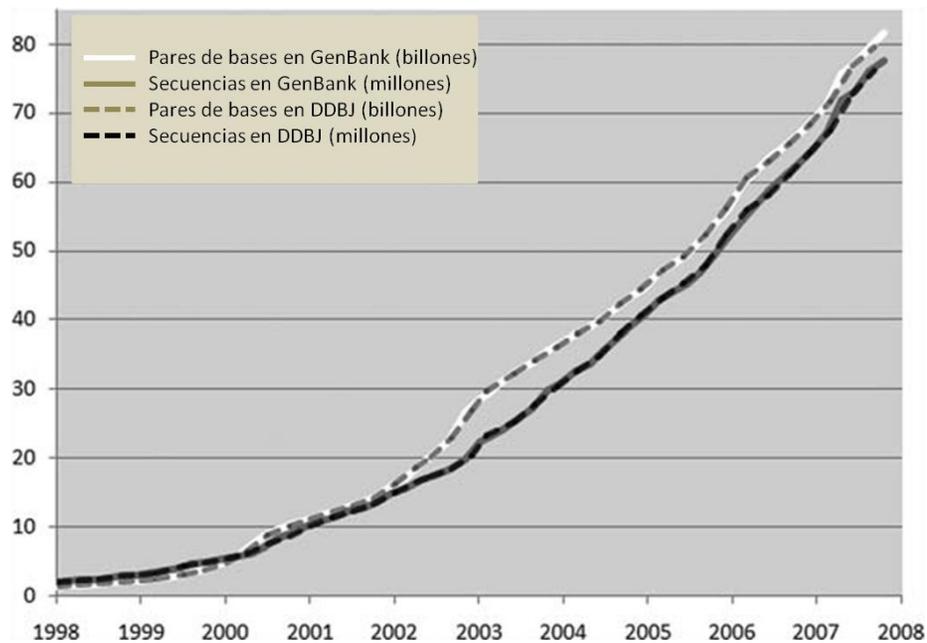


Figura 2. Crecimiento en el número de secuencias que se introdujeron en dos de las bases de datos más populares, el Gen Bank y DDBJ a lo largo de 10 años (modificada de Posada, 2009).

1.1.3. La era de las “Omicas”

El objetivo inicial de la era genómica fue determinar la secuencia genética completa de los organismos para hacer un análisis de su organización estructural, es decir, cómo estaban constituidos los genes, si poseían regiones regulatorias, si presentaban exones, dónde se encontraban localizados y por supuesto para qué proteínas codificaban, así como el estudio de secuencias no codificantes. Sin embargo, con el estudio individual de estas secuencias genéticas no es posible determinar si se expresarán o no y en qué condiciones, cuál sería su papel en la célula, cuál sería su mecanismo de acción o cómo ayudarían a conformar a la célula, etc. Por lo anterior, se volvió necesario el desarrollo de metodologías que permitieran analizar y explotar toda esta información generada, surgiendo entonces la **genómica funcional o era post-genómica**, que pretende



determinar la función de los genes y descubrir su comportamiento en los sistemas biológicos, procurando estudiarlos no como entes individuales, sino siendo parte de un todo en un momento determinado.

La genómica funcional está relacionada íntimamente con otras “-ómicas”. Dicho sufijo, aplicado a las ciencias de la vida, se refiere a la medición de la totalidad de las moléculas biológicas e información en un nivel biológico particular (Schneider & Orchard, 2011). Su principal objetivo es tener una visión global de cuándo y cómo ocurren diversos procesos en los sistemas biológicos. Las omicas han incidido en una amplia variedad de campos y entre las principales encontramos a la transcriptómica, la proteómica, la metabolómica, la interactómica. Cada una de estas se vale de metodologías particulares. La transcriptómica se refiere al estudio del transcriptoma, es decir, a los ARN mensajeros o transcritos presentes en células, tejidos u organismos. Por lo tanto, estudia la expresión de los genes. Por otro lado, la proteómica, también analiza la expresión de genes, pero a nivel de proteína. La metabolómica se refiere al estudio del metaboloma, siendo este último el conjunto de metabolitos (moléculas de bajo peso molecular) presentes en un tiempo determinado. El interactoma estudia la interacción que puede ocurrir entre las proteínas en una condición específica (Kandpal, et al., 2009). Mientras que el genoma de un organismo no varía en distintas condiciones ambientales, el transcriptoma, proteoma, metaboloma o interactoma variarán dependiendo de las distintas condiciones de crecimiento, el estado fisiológico, el desarrollo, etc. Así, estas metodologías permiten el análisis de sistemas biológicos completos.

Es claro entonces, que tras el desarrollo en las últimas décadas de las metodologías mencionadas, se han generado un inmenso número de datos. Sin embargo, su generación sólo cobraría sentido si se cuenta con recursos que permitan analizarlos y entenderlos. Por lo que, paralelamente al desarrollo científico y tecnológico, se deben mejorar también las metodologías analíticas y desarrollar herramientas que permitan manejar estos datos.

1.1.4. Concepto de Bioinformática

La bioinformática es una disciplina que se originó con el propósito utilitario de introducir orden en el conjunto masivo de datos que se estaban generando gracias a las nuevas tecnologías de la biología molecular, entre ellas la secuenciación a gran escala (Polanski & Kimmel, 2007). Se basa en usar aproximaciones computacionales de la tecnología de la información como software, simulación gráfica, algoritmos (fórmulas matemáticas), para responder problemas biológicos. Lo anterior implica tomar ventaja de un gran y complejo número de datos de forma rigurosa, con el firme propósito de que las conclusiones biológicas concebidas sean válidas. Su potencial incluso está ligado a guiar de manera eficiente, el diseño experimental en el laboratorio (Baxevanis, Ouellette, 2001).



1.2. Recursos disponibles en internet

Con el auge de las técnicas de secuenciación, un gran número de secuencias biológicas fueron generadas en las últimas décadas, resultando imposible mantener toda esta gran cantidad de información disponible en forma impresa, por lo que se volvió necesario empezar a digitalizarla, con el propósito de poder acceder a ella de una manera mucho más rápida. Además, para cumplir el mismo fin de accesibilidad, también era necesario mantenerla ordenada. Una de las formas de conservar a la información en este formato es mediante la creación de **bases de datos**. Una base de datos es, en términos simples, una **colección de datos dispuestos de manera ordenada**.

1.2.1. Bases de datos: utilidad, concepto y clasificación

El primer esfuerzo que involucró la recopilación de secuencias biológicas fue realizado a mediados de los 60's en EUA por Margaret Dayhohh y sus colaboradores, quienes reunieron todas las secuencias y estructuras tridimensionales de proteínas conocidas en un atlas de secuencias y estructuras de proteínas que fueron clasificadas de acuerdo con su similitud en secuencia (Posada, 2009). Conforme la información biológica correspondiente a secuencias se siguió acumulando con el tiempo, finalmente en 1982, el Laboratorio Europeo de Biología Molecular o EMBL (por sus siglas en inglés, ***E**uropean **M**olecular **B**iology **L**aboratory*) generó una base de datos unificada de secuencias biológicas. Posteriormente, otros dos centros se unieron a este esfuerzo, el Centro Nacional de Información en Biotecnología o NCBI (por sus siglas en inglés, ***N**ational **C**enter for **B**iotec**n**ology **I**nformation*) y el Banco de Datos de ADN de Japón o DDBJ (por sus siglas en inglés, ***D**NA **D**ata **B**ank of **J**apan*). El Instituto de Bioinformática Europeo tiene su sede en Hinxton, Inglaterra, mientras que el Centro Nacional para la Información en Biotecnología se encuentra en Bethesda, Maryland, EUA. El Banco de Datos de ADN de Japón es mantenido por el Instituto Nacional de Genética y el Centro para la Información Biológica o NIG/CIB (por sus siglas en inglés, ***N**ational **I**nstitute of **G**enetics, **C**enter for **I**nformation **B**iology*), los cuales se localizan en Mishima, Japón (Lemey et al., 2009). Es preciso decir que en la actualidad, cada una de estas bases de datos se nutre principalmente de los datos introducidos por científicos de manera individual. Cabe mencionar que desde hace aproximadamente 18 años, estos 3 grupos han formado un consorcio colaborativo denominado Colaboración Internacional de Bases de Datos de Secuencias de Nucleótidos o INSDC (por sus siglas en inglés, ***I**nternational **N**ucleotide **S**equence **D**atabase **C**ollaboration*), que tiene el propósito de compartir toda la información depositada en sus bases de datos respectiva y establecer estándares que permitan garantizar una uniformidad en los datos, de hecho, cada servidor actualiza recíprocamente su información diariamente.



Entonces, las bases de datos permiten manejar información de diversos tipos. La mayoría de ellas está conformada por un archivo de información, una organización lógica de esa información y herramientas para acceder a ella.

Dependiendo del tipo de datos que opera, el tipo de información implicada puede pertenecer a uno o varios de tres distintos niveles:

- En el primer nivel, la información únicamente está almacenada.
- En el segundo nivel, hay un acceso a esa información, ya que las bases de datos tienen cierta estructuración que lo permite.
- En el tercer nivel se consigue el análisis de los datos, en este nivel, que depende de los dos anteriores, se obtiene nueva información mediante la comparación de los datos disponibles.

Es importante mencionar que aquellas bases de datos que permiten un rápido acceso a la información poseen una estructura, donde la unidad estructural más pequeña es un registro, el cual contiene información relacionada, en un formato determinado. A dicho registro generalmente se le asigna un código que lo caracteriza, con el fin de poder localizarlo fácilmente en el futuro (Baxevanis, Ouellette, 2001).

Las bases de datos en biología molecular incluyen aquellas que contienen secuencias de ácidos nucleicos, de aminoácidos, modelos de estructuras, datos de expresión, entre otros, constituyendo una herramienta Bioinformática indispensable para mantener y recuperar información biológica. Pueden ser clasificadas de la siguiente forma (Lesk, 2002):

1. **Bases de datos de registros de información biológica:** Contienen una gran cantidad de información adicional, como por ejemplo de qué organismo proviene la secuencia, el método de secuenciación utilizado, etc. Ejemplos de estas incluyen bases de datos con los siguientes contenidos:
 - Secuencias de ácidos nucleicos y proteínas
 - Estructuras de ácidos nucleicos y proteínas
 - Patrones de expresión de proteínas
2. **Bases de datos derivadas:** Contienen información colectada a partir de las bases de datos de registros de información biológica y de los análisis de su contenido: Ejemplos de estas incluyen bases de datos con los siguientes contenidos:
 - Motivos de secuencia (motivos señal característicos de una familia de proteínas)
 - Mutaciones y variaciones en secuencias de ADN y proteínas.



1.2.2. Secuencias de ADN o de proteínas como datos biológicos

El desarrollo de técnicas para secuenciar ADN, como la terminación en cadena de Sanger mencionada en el tema anterior, introdujo un nuevo tipo de dato: las secuencias de ADN. Una vez definido este dato biológico, la técnica se utilizó para adquirir miles de secuencias de organismos muy diversos. La información que se origina directamente del laboratorio se conoce como información primaria. Así, los datos primarios son tomados en un laboratorio sobre una molécula allí existente y no son una recopilación de datos de otros autores. Una vez que una secuencia es obtenida en el laboratorio y compartida en una base de datos, esta debe de ser curada.

Las secuencias de ADN se refinan para asegurar que son de buena calidad y además se debe adicionar información relacionada con ella, como por ejemplo la posición de genes y proteínas que codifica. A este proceso se le conoce como **curación**. El proceso de curación está pensado para dar valor agregado a una secuencia. Si se elimina cualquier otro dato y sólo se deja la secuencia o por el contrario, se incluye en un registro toda la información posible, pero sin que esta haya sido verificada, producirá que la secuencia se quede sin interpretación o que los datos no tengan credibilidad (Baxevanis, Ouellette, 2001).

Una vez que se crearon y adquirieron los datos deseados, estos pueden pasar por distintos procesos, siendo uno de los más importantes, su uso. La adquisición de datos tiene un costo que se justifica por el uso que se le dará. Cuando se mencionan secuencias en los artículos científicos, lo más común es que se pretenda probar alguna hipótesis mediante dicha secuencia. Un ejemplo sería establecer la hipótesis de que nuestra secuencia de estudio actúa como una región promotora de la transcripción, por lo que este será el uso que le daremos en nuestro contexto particular. Sin embargo, el uso inicial de un dato puede resultar distinto al que se le dé subsecuentemente. Usemos el ejemplo de los marcadores genéticos de la mosca de la fruta (*Drosophila*) que Morgan encontró y usó para explicar la herencia genética. Esos mismos marcadores se utilizaron mucho tiempo después para secuenciar el genoma de esta mosca. Así, se puede hacer uso de información ya depositada, en contextos distintos al original.

1.2.3. La base de datos del NCBI

Con el fin de que te familiarices con el uso de las bases de datos más importantes en biología, en este tema nos enfocaremos en la aplicación práctica de la búsqueda de secuencias en una de las principales bases de datos, el NCBI.

Como ya se mencionó anteriormente, en los servidores del NCBI se encuentra almacenada la mayor cantidad de información primaria en biología molecular. Fue



fundada hace aproximadamente 30 años con la intención de almacenar información biológica, hacer investigación en biología computacional, desarrollar herramientas de análisis de información biológica.

Las bases de datos principales que son encontradas en el NCBI son las siguientes:

- Bases de datos de literatura científica
- Bases de datos moleculares
- Bases de datos de genomas

Y los principales tipos de datos biológicos que se pueden analizar son: secuencias de nucleótidos, secuencias de proteínas, estructuras en tres dimensiones, genes, expresión de genes y taxonomía.

Cabe mencionar que esta base de datos se encuentra soportada por los Institutos Nacionales de Salud de EUA, de los que depende a su vez, la Librería Nacional de Medicina. Así, al ser una extensión de una librería médica, su enfoque principal está dirigido a la salud.

Búsqueda de una secuencia de nucleótidos

Ahora es el momento de visitar por primera vez la página principal del NCBI.

Primeramente, ingresa a la página principal del sitio: <http://www.ncbi.nlm.nih.gov/>. Debido a que esta URL será consultada en la Unidad siguiente, te recomendamos añadirla a tus favoritos del navegador de Internet.

Desmenuzando la URL tenemos que la extensión nlm viene de ***N**ational **L**ibrary of **M**edicine*. La extensión nih viene de ***N**ational **I**nstitutes of **H**ealth*. La terminación gov es asignada sólo a entidades gubernamentales, lo cual significa que esta base de datos se sustenta con recursos públicos.

Así, la interface general de la página principal del NCBI es la siguiente:



Nuestro primer ejercicio consiste en acceder a la secuencia de nucleótidos de una proteasa con potencial biotecnológico, la papaína. Dicha enzima ha sido utilizada para aclarar la cerveza.

Esta opción de búsqueda se nutre de la base de datos del GenBank, la base de datos genómica más conocida, que es mantenida por el NCBI y contiene todas las secuencias de ácidos nucleicos y aminoácidos. Constituye una base de datos con una gran cantidad de información, que sería imposible abordar en su totalidad en este curso. Por lo tanto, en esta Unidad sólo te guiaremos para que te introduzcas en una de las aplicaciones que te permitirán obtener secuencias con las que trabajarás en la siguiente Unidad.

Usando la pestaña desplegable del lado superior izquierdo donde se leen las palabras “All Databases”, despliega la pestaña y selecciona la opción de “Nucleotide”:



The screenshot shows the NCBI website interface. At the top, there is a search bar with the word "Search" highlighted in a blue box. Below the search bar, there are several navigation menus and content sections. The "All Databases" dropdown menu is open, showing a list of databases including dbGaP, dbVar, Epigenomics, EST, Gene, Genome, GEO DataSets, GEO Profiles, GSS, HomoloGene, MedGen, MeSH, NCBI Web Site, NLM Catalog, Nucleotide, OMA, OMM, PMC, PopSet, and Probe. The "Nucleotide" option is highlighted in blue. The main content area features a section titled "Genotypes and Phenotypes" with a sub-section "Data from Genome Wide Association studies that link genes and diseases. See study variables, protocols, and analysis." Below this, there is a "Popular Resources" section with links to PubMed, Bookshelf, PubMed Central, PubMed Health, BLAST, Nucleotide, Genome, SNP, Gene, Protein, and PubChem. The "NCBI Announcements" section includes news items such as "RefSeq release 63 now available" and "Taxonomy database now shows type material, sequences from type specimens and strains now labeled in Entrez".

Una vez realizado lo anterior, anota el nombre del gen que quieres buscar, en este caso, papain (por papaína) y presiona "Search" o Enter. Recuerda que el idioma de las palabras de búsqueda debe ser el inglés.



Acto seguido, en la ventana se desplegarán todas aquellas secuencias donde se encontró la palabra *papain*. En este caso, el número de secuencias encontrado fue de 12,857. Este dato puede ser checado al inicio de los resultados, donde dice “Results”: Selecciona con el cursor la primera de las secuencias mostradas.



Search results for 'papain' showing 14,001 sequences. The first result is highlighted:

1. [Carica papaya papain mRNA, complete cds](#)
1,292 bp linear mRNA
Accession: M15203.1 GI: 167390
GenBank FASTA Graphics

Una vez que seleccionaste la secuencia, este es el tipo de registro de resultados que se muestra y que es común para todas las secuencias:

Record details for **Carica papaya papain mRNA, complete cds** (GenBank: M15203.1).

LOCUS CPAPAP 1292 bp mRNA linear PLN 19-JUN-1996

DEFINITION Carica papaya papain mRNA, complete cds.

ACCESSION M15203

VERSION M15203.1 GI:167390

KEYWORDS cysteine proteinase; endopeptidase; papain.

SOURCE Carica papaya (papaya)

ORGANISM Carica papaya

REFERENCE Cohen, L.W., Coghlan, V.M. and Dihel, L.C. Cloning and sequencing of papain-encoding cDNA Gene 48 (2-3), 219-227 (1986)

COMMENT Original source text: C.papaya (Coimbatore 2) cDNA to mRNA, clone A7. Draft entry and computer-readable sequence of [1] kindly provided by L.W.Cohen (07-APR-1987). There is a double stop codon, with the second located at position 1083 to 1085. A possible regulatory sequence of '(at)agaa' is located at position 1111 to 1132. A possible polyadenylation signal sequence is located at position 1155 to 1164.

FEATURES

source
1..1292
/organism="Carica papaya"
/mol_type="mRNA"
/db_xref="taxon:3649"
c1..1292
/product="papain mRNA"
45..1082
/EC number="3.4.22.2"



```

Firefox
http://www.ncbi.nlm.nih.gov/nuccore/M152031
Carica papaya papain mRNA, complete...
45..1082
/EC_number="3.4.22.2"
/Note="papain precursor"
/codon_start=1
/protein_id="BAB02650.1"
/db_xref="GI:167391"
/translation="MAMIFSIKLLFVAICLFVYVGLSFGDFSIQVYQNDLSTERL
IQLFESWMLKHNKIKYKNIDEKIRFPIFKDNLKVIDETNKKNSYMLGLNVFADMSND
EFKRYTIGIIRKIRTNLNEYSEQLLDCDRSYCNGGYPWSALQLVAQYGHYRNT
YFVEGVQRVCRSREKGPAAKTDVQRVQVYNEGALLYSIANQPVSVVLEAKGDFQL
TRGGIFVGFQCHKVDHAAVAAGVGPNTILIKNSWGTGNGENYIRIKRGTGNSYVGC
LITSSFFVFNW"
mat_peptide 444..1079
/product="papain"
/EC_number="3.4.22.2"
ORIGIN 1 bp downstream of BamHI site.
1 atccattccc acttaagaag taaaaagata tagctagtgt cacaatggct atgatacctt
61 caatttccaa gtgcttttt gtgcgaatg gctcttttg tcatatgggt ttgctatttg
121 gtcagatttc taatgtgaga taacttccaa atgacttggc atccactgaa agactattc
181 agctatttga atcgtggatg tgaagcaca ataatattta taagatattt gatgaaaaa
241 totacagatt tgaattttt aaagataatc taaatatat tgatgagaca aataaagaaa
301 ataacagtta ttgctttgga taaatgtgt ttgctgatat gagcaatgat gaattcaaa
361 aaaagtatac tggcttattt gctggaatt atacaacaac cgaactatca taagaagaag
421 tgcctaatga tggtagtga aatatccogg agtatgoga ttggagacaa aaaggagctg
481 tcaactcctg aaaaaatcag ggtctttgtg gtagctgttg gcaattctca cgtctgttaa
541 ctacagaggg aaactcaag atcagaaatg ggaactcaaa tgaactctca gactcaaac
601 tgccttgaag cgaacagcgt agctcaggtt gaaatgagc taactccttg agtgcactc
661 aattagtggc toaatatggt attcaactca gaaactacta cccatattag gaaatgcaac
721 gttattgtcg ctcaaggagg aaaggtcctt atgcagcaaa aacgatggg gttgcacaag
781 tgcacaacata taatgaaggg gctctcttat atcaattgc aaaccaact gtgagcgttg
841 tccctgaagc tgcagaaa gatttccaat tatataagq ggaatattt ttggggccat
901 gcggaaacaa agtagatcat gcaatgcaag caatgggta tggaaacaa tacatactca
961 taagaattc atgggttaca ggtcgggttg aaatggata tatagaattc aaagaagca
1021 ctgaaactc ctatggaga tggcatttt atacaagtc atccatctat gttaaactc
1081 gatgagatca cggctttcat aaatccctt atatatatat atatatatag aactgtata
1141 ctactcgtgt gttgaataa taatgagag gattaataa ttgttaaac ctatatata
1201 cagttgtgtg tgaacaact ttgaatcgt tttaataa ttgttaaa ttgtgtttt
1261 gattgaataa acttttcat atactttat gc
//
You are here: NCBI > DNA & RNA > Nucleotide Database Write to the Help Desk
GETTING STARTED RESOURCES POPULAR FEATURED NCBI INFORMATION
NCBI Education Chemicals & Bioassays PubMed Genetic Testing Registry About NCBI
NCBI Help Manual Data & Software Bookshelf PubMed Health Research at NCBI

```

Cada registro de secuencias de ácidos nucleicos es desplegado en GenBank con una serie de datos acompañantes que a continuación se describen (Lemey, et al., 2009):

En la primera línea se incluye el LOCUS, el número de bases que contiene la secuencia, la molécula que se secuenció (ADN, ARN, ARNm, etc.), la división a la que pertenece, que en este caso se refiere al grupo de organismos específico, PLN: secuencias de algas, fúngicas o vegetales y finalmente, la fecha de la última modificación.

LOCUS: Anteriormente era una clave que involucraba en las primeras tres letras el organismo del que provenía la secuencia y otras tres o cuatro letras para designar el nombre del gen o la probable función. En nuestro ejemplo, el organismo es *Carica papaya* (CPA), y el gen es la *pap*ína (PAP), obteniendo el código CPAPAP. Sin embargo, fue complicado mantener un mismo LOCUS porque había casos en los que se descubría que en realidad ese gen codificaba o funcionaba para algo distinto a lo que originalmente se había pensado, con lo cual habría que cambiarle el nombre. Además, esta clave podía variar entre las tres bases de datos, es decir, EMBL, DDBJ y Gen Bank. Desde el 2006 entró en desuso en EMBL.

DEFINITION: Indica el nombre científico del organismo al que corresponde la secuencia así como una breve descripción del gen o fragmento de la secuencia.



ACCESSION: Muestra el número que se le asignó en la base de datos al momento de ingresarse por el usuario. Permanece constante a lo largo de nuevas liberaciones de secuencia y es el mismo en las tres bases de datos.

VERSION: Se construye a partir del número asignado a **ACCESSION** más un número que se adiciona cada vez que la secuencia se modifica. Con este número, uno puede acceder a versiones anteriores de la secuencia.

GI: Abreviación de GenInfo. Número asignado únicamente en la base de datos del GenBank y que permite identificar solamente una secuencia particular. Es un identificador único.

Además de estos tipos de identificadores, cada secuencia codificante (CDS por CoDing Sequence) tiene un identificador de proteína (ptotein_id), el cual es compartido por las tres bases de datos.

KEYWORDS: Línea opcional que puede utilizarse para listar palabras adicionales que relacionen la secuencia con posibles búsquedas sistemáticas.

SOURCE: Indica el origen biológico de la secuencia.

ORGANISM: Muestra la clasificación biológica del organismo del cual proviene la secuencia lo más detallada posible, desde lo más general (dominio) hasta lo más particular (especie, si es el caso).

REFERENCE: Comentarios literarios que incluyen:

AUTHORS: La (s) persona (s) que obtuvo y está publicando las secuencias.

TITLE: El nombre del trabajo de investigación publicado en donde se hace referencia a esta secuencia.

JOURNAL: La fuente (Revista, Libro, etc.) en donde aparece publicada la secuencia.

PUBMED: Numero de referencia asignado en esta base de datos. Esta es una interface al servicio de citas bibliográficas del Medline.

REMARK: Aquí puede indicarse el estatus de la publicación, esto es, enviado, aceptado, en prensa, etc.

REFERENCE: Comentarios adicionales que involucran a la secuencia:

AUTHORS: Nombre de quien obtuvo la secuencia (no necesariamente es el mismo al autor que publica).

TITLE: Indica si el autor capturo directamente dicha secuencia.

JOURNAL: Indica la fecha en que se sometió el artículo a publicación, así como la dirección postal del autor principal.



FEATURES: Indica el número de bases que comprenden esta secuencia; un resumen de la clasificación del organismo fuente y de qué tipo de gen (o fragmento de este) se trata.

ORIGIN: Muestra la secuencia completa en el lenguaje de bases nitrogenadas, esto es, Timina (T), Citosina (C), Guanina (G), Adenina (A) y Uracilo (U).

Finalmente, es el momento de obtener tu secuencia y guardarla en un formato que te permitirá manipularla posteriormente. En la ventana del registro de los resultados, selecciona la opción FASTA.

The screenshot shows the NCBI GenBank entry for Carica papaya papain mRNA, complete cds (M15203.1). The page is displayed in a Firefox browser window. The main content area shows the entry details, including the accession number M15203.1, the title 'Carica papaya papain mRNA, complete cds', and the FASTA format selected in the 'Display Settings' dropdown. A green arrow points to the 'FASTA' option. The 'FEATURES' section shows the source and location of the sequence. The right sidebar contains various analysis tools and related information.

El formato FASTA es un formato estándar con el que trabajan muchos programas que analizan secuencias de ADN. Inicia con el signo de mayor que (>), seguido de una breve descripción de la secuencia de longitud variable. El siguiente renglón corresponde a la secuencia de nucleótidos.



Introduce el término papain y presiona “Search” o Enter.

Selecciona el primer resultado. ¿Cuántos resultados obtuviste? El registro de resultados es muy parecido al obtenido anteriormente cuando buscabas tu secuencia de nucleótidos. En este caso, se mostrará la secuencia en aminoácidos.



papain precursor [Carica papaya] - Protein - NCBI

http://www.ncbi.nlm.nih.gov/protein/AA802650.1

GenBank: AAB02650.1
FASTA Graphics

Go to

LOCUS AAB02650 345 aa linear PLN 19-JUN-1996
DEFINITION papain precursor [Carica papaya].
ACCESSION AAB02650
VERSION AAB02650.1 GI:167391
DBSOURCE locus CPAPAP accession M15203.1
KEYWORDS .
SOURCE Carica papaya (papaya)
ORGANISM Carica papaya
Eukaryota; Viridiplantae; Streptophyta; Embryophyta; Tracheophyta; Spermatophyta; Magnoliophyta; eudicotyledons; core eudicotyledons; rosids; malvales; Brassicales; Caricaceae; Carica.
1 (residues 1 to 345)
REFERENCE
AUTHORS Cohen,L.W., Coghlan,V.M. and Dihel,L.C.
TITLE Cloning and sequencing of papain-encoding cDNA
JOURNAL Gene 48 (2-3), 219-227 (1986)
PUBMED 2881845
COMMENT Draft entry and computer-readable sequence of [1] kindly provided by L.W.Cohen (07-APR-1987).
There is a double stop codon, with the second located at position 1083 to 1085. A possible regulatory sequence of '(at)9agaa' is located at position 1111 to 1132. A possible polyadenylation signal sequence is located at position 1155 to 1164.

FEATURES
source Location/Qualifiers
1..345 /organism="Carica papaya" /db_xref="taxon:3643"
1..345 /EC_number="3.4.22.2" /name="papain precursor"

Protein

Customize view

Analyze this sequence
Run BLAST
Identify Conserved Domains
Highlight Sequence Features
Find in this Sequence

Protein 3D Structure
Binding Of Chloromethyl Ketone Substrate Analogues To Crystalline Papain
PDB: 5PAD
Source: Carica papaya, synthetic construct
Method: X-Ray Diffraction
Resolution: 2.8 Å
See all 24 structures...

Identical proteins for AAB02650.1
Sequence 286 from patent US 8211421 [AFO13052]
Sequence 7 from patent US 6759044 [AAV22125]
Sequence 7 from patent US 6254869 [AAE75015]
See all...

papain precursor [Carica papaya] - Protein - NCBI

http://www.ncbi.nlm.nih.gov/protein/AA802650.1

Region 48..103
/region_name="Inhibitor_I29"
/note="Cathepsin propeptide inhibitor domain (I29); smar00848"
/db_xref="CDD:197916"
mat_peptide 134..345
/product="papain"
/EC_number="3.4.22.2"
Region 135..342
/region_name="Peptidase_C1A"
/note="Peptidase C1A subfamily (MEROPS database nomenclature); composed of cysteine peptidases (CPs) similar to papain, including the mammalian CPs (cathepsins B, C, F, H, L, K, O, S, V, X and W). Papain is an endopeptidase with specific substrate preferences; ccd2248"
Site /db_xref="CDD:30292"
order(152,158,292,308)
/site_type="active"
Site /db_xref="CDD:30292"
order(200..201,266,290,293,338)
/site_type="other"
/note="S2 subsite"
/db_xref="CDD:30292"
CDS 1..345
/coded_by="M15203.1:45..1082"

ORIGIN
1 mamapskkl lfvasclfvv mslsfqdfsi vqyqndlts terllqlfes vmkhnkiyk
61 nidekiyrfe ifkdnklyid etnkknnsyv lqmvfadms ndefkekyg siagnyctte
121 layeevlndg dvnipeyvdw rkgqevtpvk nqsgogscwa fsavvtiegi ikirtqnlne
181 yseqllded rrsyqenggy pwalqlvqa yqihyntyp yegvqvrcs rekqpyaakt
241 dgvrqvqpm egallysian qpavvlea qkdfqlryg ifvqpcqkv dhavaavgyg
301 pnyiliknsw gtwgengyi rikrqtgnsy gwqglytsf ypvkn
//

LinkOut to external resources
MODBASE, Database of Comparative Protein Structure M [MODBASE, Database of Comparat...]
Evolutionary Trace of Functional Site [Evolutionary Trace of Functio...]

Related information
BLINK
Related Sequences
Identical Proteins
BioAssay by Target (Identical Proteins, List)
BioAssay by Target (Identical Proteins, Summary)
CDD Search Results
Conserved Domains (Concise)
Conserved Domains (Full)
Domain Relatives
Encoding mRNA
Full text in PMC
Nucleotide
PubMed
Related Structures (List)
Related Structures (Summary)
Structure
Taxonomy

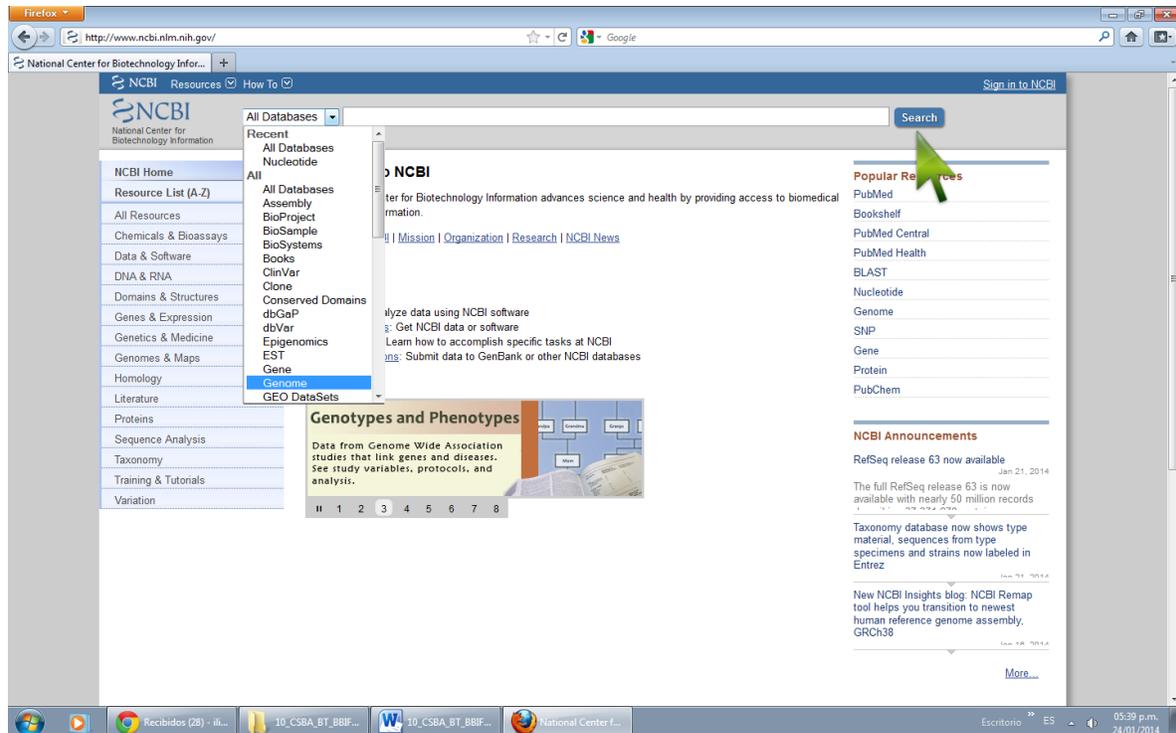
Guarda tu secuencia de aminoácidos de la misma forma que lo hiciste para la de nucleótidos, ¿verdad que es una tarea fácil? ¡Enhorabuena por la secuencia de aminoácidos que acabas de obtener!



Búsqueda de genomas secuenciados

Finalmente, con el objeto de conocer aquéllos organismos para los cuales su genoma ha sido completamente secuenciado, utilizaremos otra de las herramientas del NCBI.

Nuevamente, en la página principal, en la pestaña “All Databases”, selecciona la opción Genome.



Una vez en esta ventana, en el apartado de “External Resources”, selecciona la opción de GOLD-Genomes Online Database.



The screenshot shows the NCBI Genome homepage. At the top, there is a search bar with the text 'Genome' and a 'Search' button. Below the search bar, there is a navigation menu with links for 'Limits' and 'Advanced'. The main content area is titled 'Genome' and includes a sub-header: 'This resource organizes information on genomes including sequences, maps, chromosomes, assemblies, and annotations.' Below this, there are several sections of links:

- Using Genome:** Help, Browse by Organism, Download / FTP, Submit a genome.
- Custom resources:** Human Genome, Microbes, Organelles, Plants, Viruses.
- Other Resources:** Assembly, BioProject, BioSample, Map Viewer, Protein Clusters.
- Genome Tools:** BLAST the Human Genome, Genomic groups BLAST, NCBI remap, Genome Decoration Page.
- Genome Annotation and Analysis:** Eukaryotic Genome Annotation, Prokaryotic Genome Annotation, PASC (Pairwise Sequence Comparison), TaxPlot (3-way Genome Comparison).
- External Resources:** GOLD - Genomes Online Database, Eukaryotic Genome Browser, Genomes at Sanger, Large-scale Genome Sequencing (NHGRI).

Una vez desplegada la información de esta página, en el apartado “Isolate Genomes” presiona “Complete Projects” para poder acceder a todos los proyectos de secuenciación genómica completamente terminados o “Incomplete Projects” para conocer cuáles están en curso.



Genomes Online Database

Welcome to the Genomes OnLine Database

GOLD Genomes Online Database, is a World Wide Web resource for comprehensive access to information regarding genome and metagenome sequencing projects, and their associated metadata, around the world.

Metagenomes

- Classification
 - Studies: 439
 - Samples: 4133

Isolate Genomes

- Complete Projects: 12725
- Incomplete Projects: 26227
- Targeted Projects: 1008

Genome Distribution

- Project Type
- Sequencing Status
- Phylogenetix

1. Register

2. Annotate

3. Publish

©2012 The Regents of the University of California
Disclaimer | Credits

U.S. DEPARTMENT OF ENERGY Office of Science

Si seleccionamos los que ya están completos, obtendremos la siguiente información:

Complete Genome Projects: 12725

Archaeal: 317 Bacterial: 12096 Eukaryal: 312

Finished: 2676 Permanent Draft: 9849

GOLD ID	ORGANISM	DOMAIN	INFORMATION	SIZE	CHROM #	PLASM #	GC %	DATA	SEQUENCING CENTER	GENOME DATABASE	PUBLICATION	COMPLETION DATE
Gi0047944	Mycobacterium tuberculosis PRO5	B	ACTINOBACTERIA Taxonomy Entrez	4320 Kb 4301 orfs			65%	AOMG000000000	Universiti Teknologi MARA			
Gi18672	Rhizobium sp. JGI 001013-F22	B	PROTEOBACTERIA-ALPHA Taxonomy Entrez	3194 Kb 3364 orfs			60%	AUSB000000000	DOE Joint Genome Institute		Unpublished 2014-01-15	2014-01-15
Gi18667	Bacillus sp. JGI 001011-F15	B	FIRMICUTES Taxonomy Entrez	2048 Kb 2242 orfs			33%	ATZW000000000	DOE Joint Genome Institute		Unpublished 2014-01-15	2014-01-15
Gi0048415	Escherichia coli E1777	B	PROTEOBACTERIA-GAMMA Taxonomy Entrez	5404 Kb 5431 orfs			50%	ASYO000000000	Institute of Microbiology, Chinese Academy of Sciences		Unpublished 2014-01-14	2014-01-14
Gi09736	Brucella melitensis bv. 1 M28-12	B	PROTEOBACTERIA-ALPHA Taxonomy Entrez	3286 Kb 3216 orfs			57%	AFFA000000000	Microbial Genome Center, Beijing CASPMI - Institute of Microbiology, Chinese Academy of Sciences		J_Bacteriol. 193 (141): 3674-5 2014-01-14	2014-01-14
Gi09745	Brucella melitensis bv. 1 MS	B	PROTEOBACTERIA-ALPHA Taxonomy Entrez	3283 Kb 3215 orfs			57%	AFBZ000000000	Microbial Genome Research Center, CAS Key Laboratory of Pathogenic Microbiology and Immunology, Institute of Microbiology, Chinese Academy of Sciences CASPMI - Institute of Microbiology, Chinese Academy of Sciences		J_Bacteriol. 193 (141): 3674-5 2014-01-14	2014-01-14

La catalogación de los organismos con genomas secuenciados se encuentra de acuerdo con su taxonomía, es decir, Archea, Eubacteria o Eukarya. Se muestra la lista de los organismos cuyo genoma se ha completado, incluyendo algunas características



relevantes como el tamaño del genoma, el número de cromosomas, el lugar de secuenciación, entre otros.

Video explicativo para usar otra de las bases de datos del NCBI, el PubMed, que contiene información bibliográfica.

<http://www.youtube.com/watch?v=2FDjQ6vuARg>

Actividades

La elaboración de las actividades estará guiada por tu figura académica, mismo que te indicará, a través de la Planificación de actividades, la dinámica que tú y tus compañeros (as) llevarán a cabo, así como los envíos que tendrán que realizar.

Para el envío de tus trabajos usarás la siguiente nomenclatura: BIIN_U1_A1_XXYZ, donde BIIN corresponde a las siglas de la asignatura, U1 es la etapa de conocimiento, A1 es el número de actividad, el cual debes sustituir considerando la actividad que se realices, XX son las primeras letras de tu nombre, Y la primera letra de tu apellido paterno y Z la primera letra de tu apellido materno.

Autorreflexiones

Para la parte de **autorreflexiones** debes responder las *Preguntas de Autorreflexión* indicadas por tu figura académica y enviar tu archivo. Cabe recordar que esta actividad tiene una ponderación del 10% de tu evaluación.

Para el envío de tu autorreflexión utiliza la siguiente nomenclatura:

BIIN_U1_ATR_XXYZ, donde BIIN corresponde a las siglas de la asignatura, U1 es la unidad de conocimiento, XX son las primeras letras de tu nombre, y la primera letra de tu apellido paterno y Z la primera letra de tu apellido materno



Cierre de la Unidad

Durante la presente Unidad te adentraste en el contexto histórico mediante el cual la Bioinformática fue desarrollándose.

Primeramente, se incluyeron procesos adicionales al ya conocido dogma central de la biología molecular, con el fin de establecer el flujo de la información genética en su sentido más amplio.

Posteriormente se hizo una breve descripción de las tecnologías de secuenciación, desde la descrita por Sanger hace ya más de 30 años, hasta las más novedosas técnicas de la actualidad. La secuenciación Sanger constituyó el caballito de batalla para los esfuerzos genómicos iniciales y hoy en día, gracias al avance de la tecnología, es posible acceder a técnicas de secuenciación de segunda generación, las cuales permiten secuenciar una enorme cantidad de fragmentos de ADN de manera eficiente en un tiempo relativamente corto y a menores costos.

Finalmente, utilizaste la base de datos del NCBI, con el fin de aprender a utilizarla y entender el tipo de información que puedes obtener. Así, buscaste y accediste a una secuencia de nucleótidos y a otra de aminoácidos, además accediste a la lista de aquéllos organismos cuyos genomas han sido completamente secuenciados, así como de aquéllos que están en proceso de secuenciación.

Toda esta información te será de gran utilidad para las siguientes Unidades, en las cuales seguirás buscando secuencias para ahora sí poder analizarlas, ¡felicidades por terminar esta primera Unidad!



Para saber más



Video explicativo sobre la técnica de secuenciación de Sanger:

<http://www.dnalc.org/files/swfs/animationlib/sangerseq.exe>

Fuentes de consulta



- Baxevanis, A. E. & Ouellette, B. F. F. (Eds.). (2001). *Bioinformatics: A Practical Guide to the Analysis of Genes and Proteins. Second Edition*. New York: John Wiley & Sons, Inc.
- Fleischmann, R. D., Adams, M. D., White, O., et al. (1995). *Whole-genome random sequencing and assembly of Haemophilus influenzae Rd. Science* 269(5223):496-512.
- Kandpal, R. P., Saviola, B., Felton, J. (2009). *The era of 'omics unlimited. Biotechniques* 46(5):351-355.



- Krebs, J. E., Goldstein, E. S., Kilpatrick, S. T. (2010). *Lewin's essential genes*. Second Edition. Massachusetts: Jones and Bartlett Publishers.
- Lemey, P., Salemi, M., Vandamme A. M. (Eds.). (2009). *The Phylogenetic Handbook. A practical Approach to Phylogenetic Analysis and Hypothesis testing*. Second Edition. Cambridge.
- Lesk, A. M. (2002). *Introduction to Bioinformatics*. New York: Oxford University Press.
- Mardis, E. R. (2007). *The impact of next-generation sequencing technology on genetics*. *Trends in Genetics* 24(3):133-141.
- Nelson, D. L. & Cox, M. M. (2005). *Lehninger. Principles of Biochemistry*. Fourth Edition. New York: W.H. Freeman and Company.
- Ondarza, R. N. (1994). *Biología molecular. Antes y después de la doble hélice*. México: Siglo veintiuno editores.
- Polanski, A. & Kimmel, M. (2007). *Bioinformatics*. New York: Springer.
- Posada, D. (Ed.). (2009). *Bioinformatics for DNA sequence analysis*. New York: Humana Press.
- Reis-Filho, J. S. (2009). *Next-generation sequencing. Breast Cancer Research* 11(Suppl 3):S12 (doi: 10.1186/bcr2431).
- Sanger, F., Nicklen, S., Coulson, A. R. (1977). *DNA sequencing with chain-terminating inhibitors*. *Proceeding of the Natural Academy of Science U S A* 74(12):5463-5467.
- Schneider, M. V. & Orchard, S. Chapter 1. *Omics Technologies, Data and Bioinformatics Principles*, en Mayer, B. (Ed.). (2011). *Bioinformatics for Omics Data: Methods and Protocols*. Springer Science.