

# Servidores genómicos

Enrique Blanco García

PID\_00165220



Universitat Oberta  
de Catalunya

[www.uoc.edu](http://www.uoc.edu)



# Índice

<b>Introducción.....</b>	<b>5</b>
<b>Objetivos.....</b>	<b>7</b>
<b>1. Introducción.....</b>	<b>9</b>
<b>2. Filosofía de la navegación genómica.....</b>	<b>13</b>
<b>3. El navegador genómico UCSC.....</b>	<b>19</b>
3.1. Estructura .....	19
3.2. Pistas .....	25
3.3. Anotaciones propias .....	36
3.4. ENCODE .....	42
3.5. Instalación local .....	44
<b>4. El navegador genómico ENSEMBL.....</b>	<b>50</b>
<b>Resumen.....</b>	<b>55</b>
<b>Actividades.....</b>	<b>57</b>
<b>Ejercicios de autoevaluación.....</b>	<b>57</b>
<b>Solucionario.....</b>	<b>59</b>
<b>Bibliografía.....</b>	<b>60</b>



## Introducción

La secuenciación de los genomas de múltiples especies constituye un hito sin precedentes en la historia del progreso científico. Gracias al trabajo conjunto de miles de investigadores de todo el mundo, podemos explorar desde nuestro propio ordenador estas secuencias biológicas para descifrar las claves del código genético. El acceso universal a toda esta información está transformando radicalmente la investigación en biología molecular y biomedicina. En consecuencia, en el camino para explicar numerosos enigmas hasta ahora sin resolver, el enfoque científico clásico está derivando en nuevas aproximaciones más pragmáticas. Dicho de otro modo, en lugar de orientar el foco de investigación hacia un gen en particular, ahora podemos trabajar a gran escala, analizando simultáneamente todos los genes involucrados en un proceso biológico concreto.

El progreso exponencial experimentado por la ciencia en la última década se sustenta, en gran parte, en la implementación de eficientes herramientas de transmisión de información sobre la plataforma de la Red. Hoy día, para un investigador resulta extremadamente sencillo explorar desde su propio ordenador personal la totalidad de las anotaciones biológicas realizadas sobre una región concreta del genoma por otros miembros de la comunidad científica. Los navegadores genómicos son una poderosa herramienta para inferir computacionalmente nuevo conocimiento a partir de los datos aportados por otros medios más tradicionales, como los resultados obtenidos en un entorno experimental. Los enormes avances conseguidos en los estudios de regulación genómica más recientes, de hecho, no pueden comprenderse sin la contribución capital de esta clase de aplicaciones.

Mediante una serie de convenios razonablemente establecidos, cualquier bioinformático puede obtener en pocos minutos el cartografiado completo del genoma para realizar nuevas contribuciones. Con estas regulaciones generalmente aceptadas por la comunidad de investigadores, el conocimiento existente sobre el genoma de una especie (la secuencia de nucleótidos y el mapa de anotaciones sobre ésta) resulta continuamente actualizado, gracias a la incesante actividad científica que intenta caracterizar con mayor precisión cada escenario biológico que ocurre dentro del entorno celular.

A lo largo de este módulo mostraremos al estudiante los mecanismos esenciales para explorar con garantías los genomas utilizando estos portales web, incidiendo especialmente en la correcta interpretación del increíble volumen de información que contienen. Numerosos servidores genómicos surgieron desde la publicación de los primeros genomas. En las próximas líneas tomaremos como referencia el navegador de la Universidad de Santa Cruz de California (UCSC) en Estados Unidos. Múltiples proyectos internacionales surgidos con

posterioridad a la consecución de la secuencia del genoma humano depositan sus datos en este repositorio, cuyo manejo resulta enormemente sencillo para cualquier investigador. No obstante, los principales fundamentos de manejo que enseñaremos para este servidor pueden aplicarse en cualquier otra aplicación similar. Por consiguiente, para complementar estas informaciones, mencionaremos otros recursos genómicos de notable relevancia, como el portal ENSEMBL, diseñado en el Instituto Europeo de Bioinformática (EBI) en Hinxton (Reino Unido).

## Objetivos

Con el programa de contenidos establecido en este módulo, una vez finalizada la etapa de aprendizaje, el estudiante debe lograr la consecución de los siguientes objetivos relacionados con el manejo de los principales servidores genómicos existentes:

1. Conocer el elenco de bancos de datos genómicos existentes.
2. Analizar cómo se almacena computacionalmente un genoma.
3. Manejar las opciones básicas de los navegadores genómicos.
4. Aprender a interpretar los datos suministrados por éstos.
5. Enriquecer las anotaciones existentes con nuevas informaciones.
6. Acceder a las anotaciones del proyecto internacional ENCODE.



## 1. Introducción

El genoma de una célula es un repositorio de secuencias de ADN empaquetado en forma de cromosomas. Este material hereditario codifica los genes que una vez descifrados resultan útiles para la síntesis de proteínas. Junto con los genes, cohabitan en el genoma otros elementos que regulan la activación de los mismos, y proporcionan, además, una cierta estructuración a la cromatina. Para modelar este complejo escenario biológico dentro de un entorno informático, la secuencia de nucleótidos de cada cromosoma de un genoma debe transformarse, en primer lugar, en un fichero de texto. Junto con las secuencias, es necesario cartografiar cada cromosoma aportando un segundo tipo de datos denominados anotaciones. Las anotaciones son necesarias para indicar la ubicación exacta de aquellos elementos cifrados en una secuencia concreta. El catálogo de elementos biológicos identificables dentro de un genoma está constituido primordialmente por:

- Genes.
- Sitios de unión de factores de transcripción.
- Inicios de transcripción.
- Marcas de modificación de histonas.
- Regiones repetitivas.
- Polimorfismos.

En función de los nuevos datos aportados por investigaciones más recientes, constantemente ampliamos y mejoramos nuestro conocimiento sobre cualquier función biológica o componente celular. En consecuencia, parece natural pensar que tanto la secuencia como las anotaciones de cualquier genoma necesitan renovarse. Para favorecer la reproducibilidad, por regla general se suele distribuir la última versión de cada genoma (secuencia y anotaciones), junto con el repositorio de anteriores distribuciones. Cada nueva versión de un genoma establece una codificación propia de referencia, mientras que las coordenadas de un elemento funcional pueden variar entre versiones cuando la secuencia de base resulta modificada, posiblemente debido a mejoras en la secuenciación de ese organismo. Como veremos más adelante, existen ficheros de conversión de coordenadas entre versiones (en inglés, *liftover*) para comparar correctamente distintas anotaciones. Es fundamental mencionar cuál estamos empleando en nuestros análisis bioinformáticos (la tabla 1 muestra las distribuciones genómicas más recientes de varios organismos modelo).

### Ved también

Una descripción detallada de los genomas y los elementos funcionales más comunes codificados en su interior puede encontrarse en la asignatura *Fundamentos de biología molecular*.

Por otro lado, en la asignatura *Fundamentos de informática en entornos bioinformáticos* aparecen varios ejemplos de anotación computacional de regiones genómicas.

### Ved también

Un listado exhaustivo de los genomas secuenciados en la actualidad puede encontrarse en la asignatura *Fundamentos de biología molecular*.

Una distribución concreta de un genoma viene definida por:

- La secuencia de nucleótidos de cada cromosoma determinada por una versión concreta del ensamblado producido por el consorcio correspondiente.
- La serie de anotaciones pertenecientes a distintas características biológicas cartografiadas dentro de la secuencia de los cromosomas por diferentes entidades de investigación.

Tabla 1. Distribuciones actualizadas de los genomas más utilizados.

Genoma	Ensamblado	Anotaciones	Fecha
<i>Homo sapiens</i>	GRCh37	hg19	Febrero 2009
<i>Pan troglodytes</i>	CGSC 2.1	panTro2	Marzo 2006
<i>Bos taurus</i>	Baylor 4.0	bosTau4	Octubre 2007
<i>Mus musculus</i>	NCBI37	mm9	Julio 2007
<i>Rattus norvegicus</i>	Baylor 3.4	rn4	Noviembre 2004
<i>Gallus gallus</i>	WUGSC 2.1	galGal3	Mayo 2006
<i>Danio rerio</i>	Zv8	danRer6	Diciembre 2008
<i>Fugu rubripes</i>	JGI 4.0	fr2	Octubre 2004
<i>Drosophila melanogaster</i>	BDGP R5	dm3	Abril 2006
<i>Anopheles gambiae</i>	IAGC MOZ2	anoGam1	Febrero 2003
<i>Caenorhabditis elegans</i>	WS190	ce6	Mayo 2008
<i>Saccharomyces cerevisiae</i>	SGD	sacCer2	Junio 2008

Cuando analiza la secuencia de una región cromosómica para explorar su contenido o anotar nuevos elementos, el bioinformático accede a estas informaciones mediante un programa especial denominado servidor genómico (en inglés, *genome browser*). Generalmente, tanto la secuencia como las anotaciones de las múltiples versiones de un genoma se distribuyen de forma pública a través de la Red. Con estos programas, todo el conocimiento aportado por distintos grupos de investigación sobre el genoma está integrado en una sola herramienta, facilitando enormemente el intercambio de información. Los servidores genómicos proporcionan los datos que poseen sobre cada genoma a través de potentes interfaces gráficos que favorecen su legibilidad. Estas fotografías del genoma, no obstante, son generadas a partir de simples ficheros de texto que contienen las coordenadas de las anotaciones. Estos archivos suelen distribuirse también de forma separada para realizar tratamientos bioinformáticos de forma local.

#### Lectura complementaria

M. S. Cline; W. J. Kent (2009). "Understanding genome browsing". *Nature Biotechnology* (núm. 27 págs. 153-155).

Un **servidor genómico** es una aplicación informática que proporciona herramientas para navegar eficientemente a lo largo de la secuencia de los genomas, permitiendo el acceso a las anotaciones cartografiadas en su interior.

Dada su enorme versatilidad, existen varios puntos de entrada en un servidor genómico para acceder a la información de un elemento en particular. En la tabla 2 mostramos algunos ejemplos a la hora de localizar la misma región genómica que contiene un gen de interés para nosotros. Según el método escogido y el nivel de detalle de nuestra búsqueda, es probable que el servidor genómico identifique más de un posible resultado. En ese caso, la exploración visual de las alternativas debería ser suficiente para elegir correctamente la región deseada. Una vez realizada satisfactoriamente la búsqueda, el navegador genómico nos proporcionará una fotografía de las anotaciones disponibles en forma de pistas o carriles (en inglés, *tracks*). Con este mecanismo de visualización, el bioinformático puede comparar fácilmente diferentes elementos anotados sobre una misma región genómica.

#### Distintos puntos de entrada

Dependiendo de los datos que conozcamos sobre el objeto de la búsqueda, debemos escoger el modo más adecuado de acceder a él. Generalmente, es posible acceder desde varios lugares distintos a la misma información.

Tabla 2. Puntos de entrada en un servidor genómico.

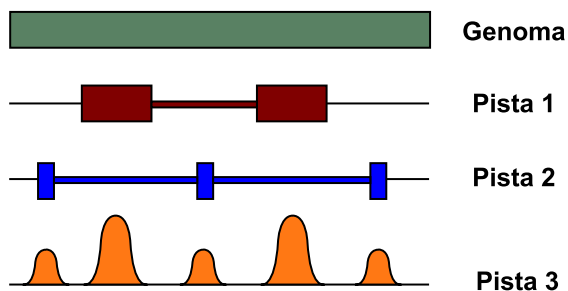
Clase	Formato	Ejemplo
Región	Coordenadas	chr2:80529003-80531487
Gen	Abreviatura	<i>LRRTM1</i>
Gen	Nombre completo	<i>leucine-rich repeat transmembrane neuronal</i>
Gen	Código RefSeq	NM_178839
Proteína	Identificador	LRRTM1
Tránsito	Identificador	BC045113
Secuencia	FASTA	ATGGATTTCCTGCTGCTCGGTC...

Una **pista de información** contiene la colección de elementos biológicos de una determinada clase anotados en un servidor genómico para una región específica del genoma. Cada anotación está definida generalmente por las coordenadas exactas necesarias para posicionar dicho elemento en el interior de la secuencia, y por un valor numérico asociado al grado de representatividad de éste.

El servidor genómico proyecta sobre cada pista de su visor genómico la ubicación de los elementos biológicos conocidos en ambas orientaciones de la molécula de ADN simultáneamente. Las anotaciones en hebras distintas pueden presentarse gráficamente integradas dentro de la misma pista o estar separadas en cuadrantes distintos de la imagen, según el navegador. La propia secuencia del genoma es una pista cuyo contenido corresponde a la base observada

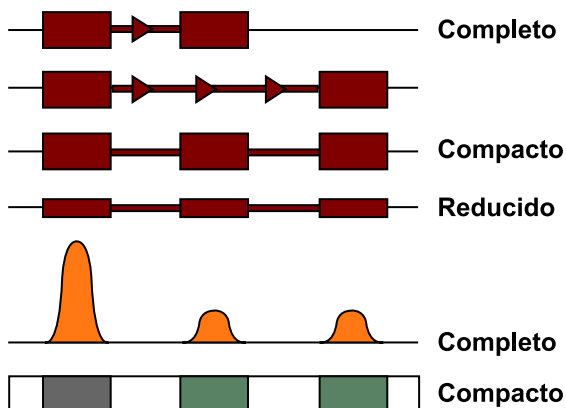
en cada nucleótido de la región visualizada (únicamente cuando el tamaño de ésta sea inferior a la resolución de la imagen, el navegador muestra explícitamente la secuencia de bases). Para gestionar de forma eficiente la proyección simultánea de varias pistas, el navegador superpone la información en diferentes niveles (ver figura 1). El usuario puede elegir el elenco de pistas que desea visualizar en cada momento, así como en qué orden desea disponerlas en la ventana gráfica.

Figura 1. Representación gráfica de anotaciones en forma de pistas.



En tiempo real, el servidor compone una fotografía del genoma en función de nuestras necesidades a partir de los ficheros de anotaciones almacenados en su propia base de datos. El usuario puede configurar esta superposición gráfica de pistas, escogiendo, entre las disponibles, aquellas que más le interesan. Por regla general, las anotaciones están agrupadas en bloques conceptualmente relacionados. Una vez establecido el contenido del nuevo mapa de anotaciones, el navegador realiza una actualización de la imagen. También es posible establecer el nivel de detalle gráfico de cada pista para conseguir representaciones gráficas que faciliten una comparación óptima (desplegar o compactar pistas, ver figura 2). Existen aplicaciones auxiliares en los navegadores genómicos para realizar estas comparaciones de forma cuantitativa, teniendo en cuenta la correlación entre la ubicación de los elementos de diferentes anotaciones.

Figura 2. Mostrando diversos niveles de información con pistas.



#### Leyenda figura 1

Superposición de varios niveles de información en una misma región genómica (genes, sitios de unión a factores de transcripción y marcas de metilación de histonas, respectivamente).

#### Ved también

Los ficheros de anotaciones pueden administrarse con gestores de bases de datos, como se explica en la asignatura *Fundamentos de informática en entornos bioinformáticos*.

#### Bloques de opciones

Los bloques de opciones disponibles pueden cambiar en función de la distribución del genoma habilitado.

#### Leyenda figura 2

Mostramos dos clases de pistas. En el primer bloque visualizamos un gen con dos isoformas, agrupando esta información de tres maneras posibles. En el segundo bloque resumimos el perfil de una determinada modificación de histonas asociada al gen anterior con un gráfico de densidad más compacto.

## 2. Filosofía de la navegación genómica

En la actual era post-genómica, los investigadores han desplazado el foco de atención desde las secuencias individuales hacia el análisis a gran escala del genoma. Con el objetivo de tener una visión más amplia de las anotaciones agrupadas en torno a una región del genoma, aparecen los servidores genómicos, que integran distintas fuentes de información dispersas por la Red. Estos sistemas redefinen, por tanto, la comunicación con sus complejos repositorios de vastos volúmenes de información.

Actualmente navegamos por heterogéneos paisajes genómicos codificados en el interior de los cromosomas, concentrando nuestro interés únicamente sobre determinadas áreas. Desde el genoma completo hasta la secuencia individual de una proteína, el usuario de estos servicios puede recuperar todas las anotaciones con diversos niveles de detalle. Lógicamente, esta amalgama de recursos también es accesible de forma aislada, menguando sin embargo la efectividad del análisis bioinformático fuera de estos entornos integrados.

Como muestra la figura 3, desde la representación gráfica servida por el navegador genómico podemos visitar las anotaciones particulares suministradas por otro tipo de recursos más específicos. Por ejemplo, cuando analizamos los genes anotados en una región genómica en concreto, podemos acceder desde la pista individual a una nueva pantalla que contiene información recopilada por múltiples fuentes de información. Muchos de estos repositorios, que denominaremos **bases de datos primarias**, constituyeron en su momento el germen de los actuales servidores genómicos.

Cuando no era posible llevar a cabo análisis globales de un genoma completo, la unidad de búsqueda de información era precisamente la secuencia individual (cuya anotación era suministrada por diferentes miembros de la comunidad científica). Estas primeras colecciones de secuencias, a pesar de su enorme valor, contenían frecuentes errores que podían conducir a importantes excesos de redundancia. Estos catálogos reciben actualmente un tratamiento de validación mucho más cuidadoso, basado en la verificación manual llevada a cabo por expertos entrenados en este tipo de tareas. A continuación ofrecemos una selección de los recursos primarios más populares. A lo largo de estos materiales analizaremos detalladamente el contenido de estos repositorios de información genómica, cuyo acceso está integrado dentro de cualquiera de los servidores genómicos existentes.

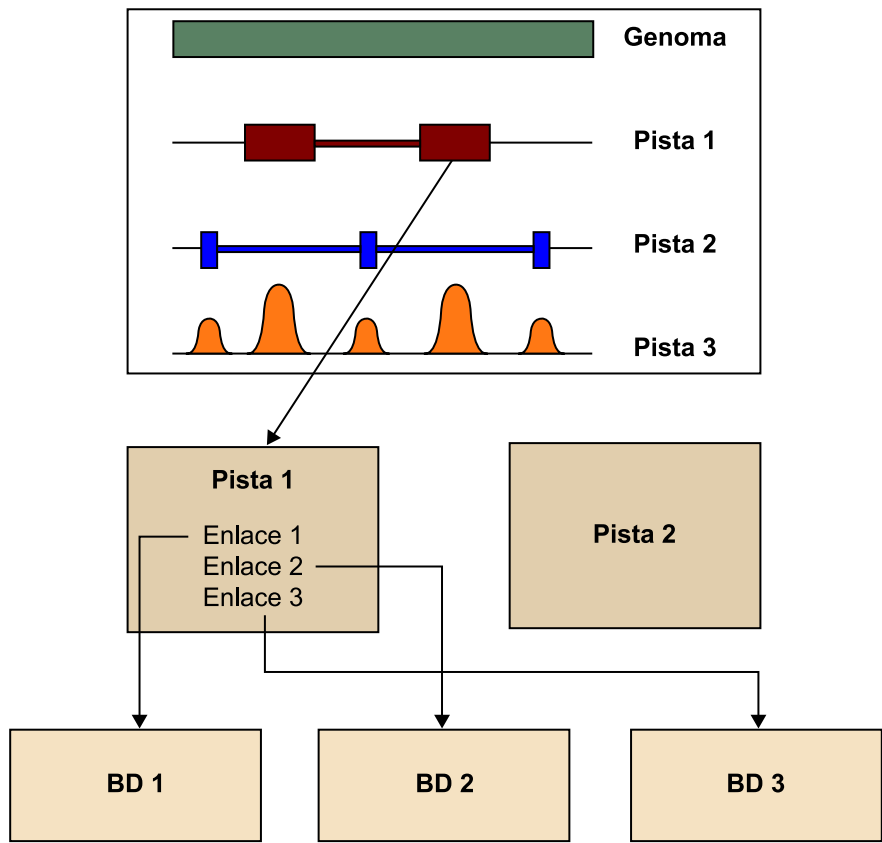
### Lectura complementaria

C. A. Ouzounis; A. Valencia (2003). "Early bioinformatics: the birth of a discipline - a personal view". *Bioinformatics* (núm. 19, págs. 2176-2190).

### Nucleic Acids Research

Esta prestigiosa revista científica publica semestralmente ediciones especiales sobre nuevos servidores y bases de datos en el ámbito genómico y proteómico.

Figura 3. Paradigma de la comunicación con navegadores genómicos.



**Leyenda figura 3**

Desde su propio ordenador, el investigador viaja a través de varios niveles de información, que se sirve mediante distintas bases de datos integradas en el navegador genómico.

Tabla 3. Listado de recursos genómicos esenciales

Recursos primarios	Dirección
GenBank	<a href="http://www.ncbi.nlm.nih.gov/genbank">http://www.ncbi.nlm.nih.gov/genbank</a>
RefSeq	<a href="http://www.ncbi.nlm.nih.gov/refseq">http://www.ncbi.nlm.nih.gov/refseq</a>
Unigene	<a href="http://www.ncbi.nlm.nih.gov/unigene">http://www.ncbi.nlm.nih.gov/unigene</a>
Homologene	<a href="http://www.ncbi.nlm.nih.gov/homologene">http://www.ncbi.nlm.nih.gov/homologene</a>
Gene Ontology	<a href="http://www.geneontology.org">http://www.geneontology.org</a>
Kyoto Encyclopedia of Genes and Genomes	<a href="http://www.genome.jp/kegg">http://www.genome.jp/kegg</a>
Transfac	<a href="http://www.gene-regulation.com">http://www.gene-regulation.com</a>
Jaspar	<a href="http://jaspar.genereg.net">http://jaspar.genereg.net</a>
ORegAnno	<a href="http://www.oreganno.org">http://www.oreganno.org</a>
Gene Expression Omnibus	<a href="http://www.ncbi.nlm.nih.gov/geo">http://www.ncbi.nlm.nih.gov/geo</a>
dbSNP	<a href="http://www.ncbi.nlm.nih.gov/projects/SNP">http://www.ncbi.nlm.nih.gov/projects/SNP</a>
OMIM	<a href="http://www.ncbi.nlm.nih.gov/omim">http://www.ncbi.nlm.nih.gov/omim</a>
Pubmed	<a href="http://www.ncbi.nlm.nih.gov/pubmed">http://www.ncbi.nlm.nih.gov/pubmed</a>

Navegadores genómicos	Dirección
UCSC	<a href="http://genome.ucsc.edu">http://genome.ucsc.edu</a>
ENSEMBL	<a href="http://www.ensembl.org">http://www.ensembl.org</a>
NCBI	<a href="http://www.ncbi.nlm.nih.gov/sites/genome">http://www.ncbi.nlm.nih.gov/sites/genome</a>
Vertebrate Genome Annotation	<a href="http://vega.sanger.ac.uk">http://vega.sanger.ac.uk</a>
Mouse Genome Informatics Database	<a href="http://www.informatics.jax.org">http://www.informatics.jax.org</a>
Rat Genome Database	<a href="http://rgd.mcw.edu">http://rgd.mcw.edu</a>
FlyBase	<a href="http://flybase.org">http://flybase.org</a>
WormBase	<a href="http://www.wormbase.org">http://www.wormbase.org</a>
Saccharomyces Genome Database	<a href="http://www.yeastgenome.org">http://www.yeastgenome.org</a>
Arabidopsis Information Resource	<a href="http://www.arabidopsis.org">http://www.arabidopsis.org</a>
Rice Genome Annotation Project	<a href="http://rice.plantbiology.msu.edu">http://rice.plantbiology.msu.edu</a>
International Wheat Genome Sequencing Consortium	<a href="http://www.wheatgenome.org">http://www.wheatgenome.org</a>
E.coli database	<a href="http://www.genome.wisc.edu">http://www.genome.wisc.edu</a>
Gbrowse	<a href="http://gmod.org/wiki/Gbrowse">http://gmod.org/wiki/Gbrowse</a>

Navegadores genómicos avanzados	Dirección
VISTA	<a href="http://pipeline.lbl.gov">http://pipeline.lbl.gov</a>
rVISTA	<a href="http://rvista.dcode.org">http://rvista.dcode.org</a>
enhancerVISTA	<a href="http://enhancer.lbl.gov">http://enhancer.lbl.gov</a>
modENCODE	<a href="http://www.modencode.org">http://www.modencode.org</a>
International Cancer Genome Consortium	<a href="http://www.icgc.org">http://www.icgc.org</a>
International HapMap Project	<a href="http://www.hapmap.org">http://www.hapmap.org</a>
Human Epigenome Consortium	<a href="http://www.epigenome.org">http://www.epigenome.org</a>

El analista bioinformático, como vemos, dispone de un amplio abanico de servidores genómicos. Los **servidores genéricos** integran la secuencia y las anotaciones de múltiples genomas, presentando todos los datos dentro de un marco común uniforme que facilita su manipulación. El navegador genómico de UCSC o la plataforma ENSEMBL son las aplicaciones genéricas más populares. Todas estas herramientas suelen cruzar sus pistas de anotaciones, facilitando de este modo la integración y visibilidad de diferentes fuentes de información desde cualquier lugar. El código de la mayoría de estas plataformas se distribuye gratuitamente junto con todos sus bancos de datos. Profundizando en esta filosofía, la herramienta Gbrowse ofrece todas las funciones de un navegador genómico genérico adaptable a cualquier genoma, cediendo al usuario la responsabilidad de introducir las anotaciones que deben respetar ciertas pautas sobre el formato de los ficheros.

En determinadas especies se han implementado **navegadores específicos** suministrados por consorcios de secuenciación para acceder exclusivamente a sus anotaciones. Estos recursos están optimizados exclusivamente para trabajar con ese genoma, convirtiéndose en un repositorio de referencia para la comunidad de estudio de esa especie. Los ejemplos más conocidos de estos programas están dedicados a distintos organismos modelos como FlyBase para la mosca de la fruta, los navegadores de ratón y rata, o los portales para la secuenciación de numerosas especies vegetales. Los servidores genéricos, no obstante, comparten con éstos las mismas distribuciones de los genomas (ver figura 4).

### Lecturas complementarias

W. J. Kent (2002). "The human genome browser at UCSC". *Genome Research* (núm. 12, págs. 996-1006).

P. Flicek y otros (2010). "Ensembl's 10th year". *Nucleic Acids Research* (núm. 38, págs. D557-D562).

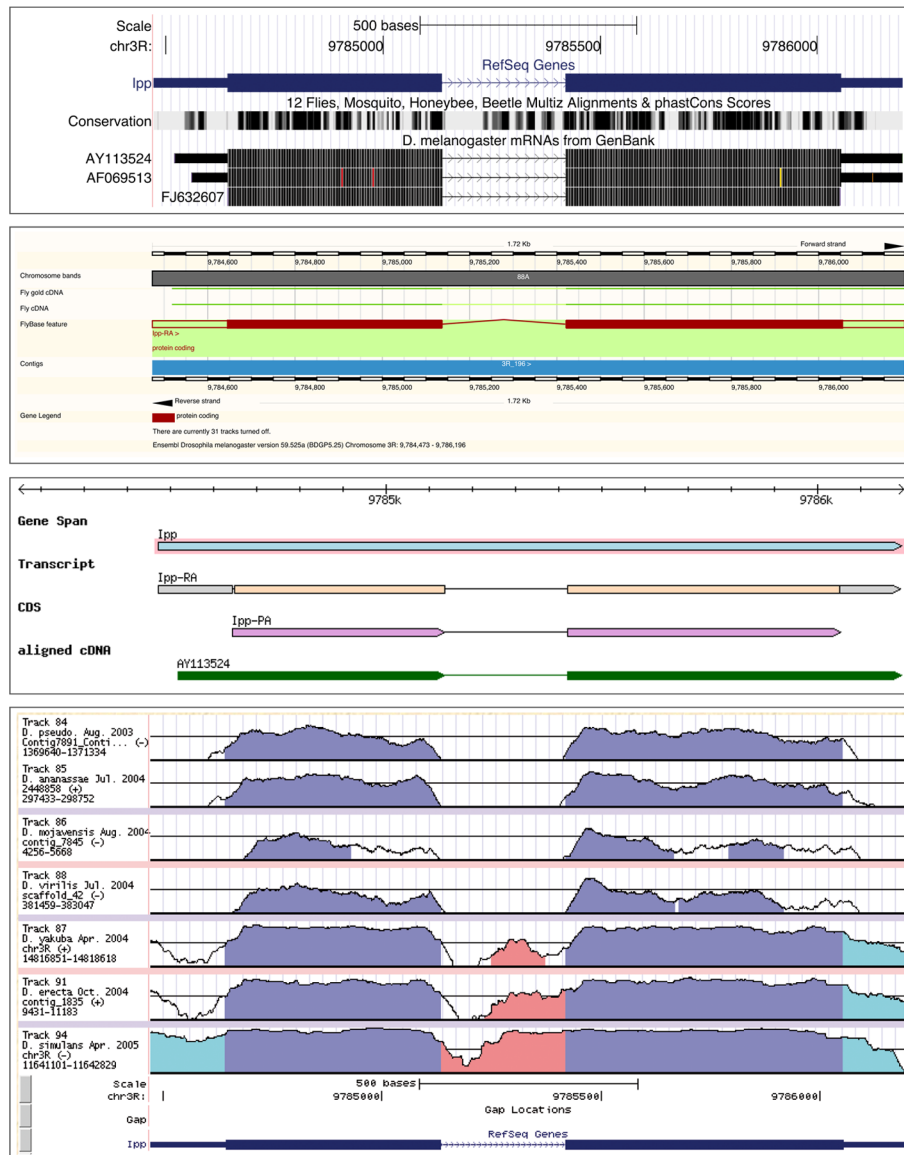
D. A. Benson y otros (2008). "GenBank". *Nucleic Acids Research* (núm. 36, págs. D25-30).

L. D. Stein y otros (2002). "The generic genome browser: a building block for a model organism system database". *Genome Research* (núm. 12, págs. 1599-1610).

The FlyBase Consortium (2009). "FlyBase: enhancing Drosophila Gene Ontology annotations". *Nucleic Acids Research* (núm. 37, págs. D555-D559).

Mouse Genome Database Group (2008). "The Mouse Genome Database (MGD): mouse biology and model systems". *Nucleic Acids Research* (núm. 36, págs. D724-728).

Figura 4. Visualización de una región genómica.

**Leyenda figura 4**

El gen *lpp* visualizado a través de distintos navegadores genómicos. Desde el margen superior de la imagen hacia el inferior: UCSC, ENSEMBL, FlyBase y VISTA.

La comparación entre las anotaciones de varios genomas en regiones ortólogas puede proporcionar una información muy valiosa sobre la ubicación de ciertos elementos funcionales conservados evolutivamente. Por ejemplo, la identificación de algunas secuencias preservadas a lo largo de millones de años permite refinar la anotación de genes y elementos asociados a su regulación. VISTA es uno de los portales pioneros en realizar este tipo de comparaciones, facilitando enormemente la búsqueda a los investigadores.

En otro orden de cosas, el análisis en profundidad del genoma humano resulta potencialmente interesante dentro del campo de la biomedicina y la salud, debido a que la mayoría de enfermedades documentadas poseen un importante componente genético. No es casual, por tanto, que sean necesarios también navegadores genómicos que proporcionen un mapeado completo de las

**Ved también**

El concepto de genómica comparada se trata en profundidad en la asignatura *Fundamentos de biología molecular*.

anotaciones referentes a estos desórdenes hereditarios. Ejemplos de estos portales biomédicos son el proyecto de secuenciación del cáncer o la detección del conjunto de polimorfismos existentes en el ser humano.

Dada su enorme popularidad, fraguada básicamente por su simplicidad de uso, a continuación analizaremos el navegador genómico de UCSC como ejemplo de servidor de propósito general. Profundizar en su estudio nos servirá también para investigar el contenido de muchas de las bases de anotaciones primarias. El portal genómico de UCSC distribuye, además de un vasto conjunto de genomas, las anotaciones generadas en el marco de distintos proyectos internacionales de secuenciación, aumentando por tanto su relevancia como objeto de estudio.

#### Lectura complementaria

C. Mayor y otros (2000). "VISTA: visualizing global DNA sequence alignments of arbitrary length". *Bioinformatics* (núm. 16, págs. 1046-1047).

### 3. El navegador genómico UCSC

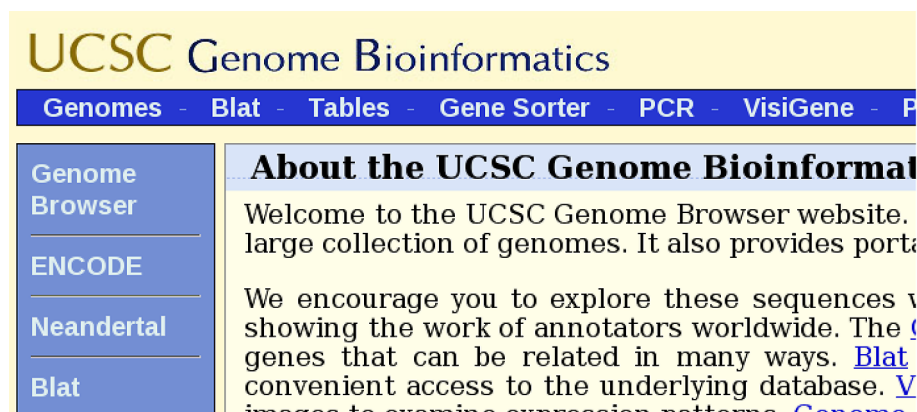
#### 3.1. Estructura

El navegador de la Universidad de Santa Cruz de California (UCSC, de ahora en adelante) fue diseñado con el objeto de construir una herramienta que proporcionara al usuario un acceso sencillo a las anotaciones existentes sobre el genoma humano, que estaba siendo secuenciado en aquel momento. Desde su nacimiento, este portal ha incorporado numerosas mejoras a su implementación inicial, convirtiéndose en un recurso de referencia, utilizado, por ejemplo, para distribuir los resultados de la primera fase del proyecto ENCODE. Todas las anotaciones sobre un genoma base de referencia se suministran en forma de pistas de datos, gráficamente representadas sobre una línea horizontal en paralelo con la propia secuencia de cada cromosoma.

#### Lectura complementaria

W. J. Kent y otros (2002). "The human genome browser at UCSC". *Genome Research* (núm. 12, págs. 996-1006).

Figura 5. Pantalla inicial de UCSC.



Para iniciar la exploración debemos seleccionar primero el organismo y la versión del genoma que deseamos explorar. A continuación, es necesario introducir información suficiente para que el navegador genómico pueda identificar el objeto de nuestra búsqueda (por ejemplo, genes, regiones o cromosomas). Vamos a centrarnos en un gen humano denominado *LRRTM1* (abreviatura de *leucine-rich repeat transmembrane neuronal*).

Figura 6. Parámetros de búsqueda en UCSC.

clade	genome	assembly	
Mammal	Human	Feb. 2009 (GRCh37/hg19)	LRRTM1

En función del tipo de identificador especificado, los resultados de la búsqueda serán más o menos concretos, siendo preciso incluir en algunas ocasiones información adicional para caracterizar mejor aquello que deseamos localizar.

En la tabla 4 mostramos distintos ejemplos con diferentes tipos de valores aceptados por el motor de búsqueda del navegador UCSC para acabar encontrando el mismo gen en todos los casos.

Tabla 4. Objetos de búsqueda de UCSC.

Nombre	Ejemplo
Cromosoma	chr2
Región	chr2:80529002-80531487
Región	chr2:80529002+2485
Gen	<i>LRRTM1</i>
Tránsito	NM_178839
Proteína	NP_849161

Una vez acotados los parámetros de búsqueda, el sistema procesa estos datos y proporciona un listado de todas las pistas que posean algún nexo en común con la información suministrada. En este caso, encontramos siete pistas relacionadas con el gen *LRRTM1* en la distribución hg19 del genoma humano (a noviembre del 2010). Para mostrar cómo interpretar esta información, procederemos a enumerar los distintos resultados obtenidos, analizando la ficha desde la parte superior de ésta (ver figura 7):

- Pista UCSC Genes: Dos isoformas de este gen incluidas dentro de la anotación del genoma humano servida por el propio navegador UCSC.
- Pista RefSeq Genes: Una única anotación producida por el consorcio RefSeq.
- Pista Non-human RefSeq Genes: Cuatro formas ortólogas de este gen identificadas en las anotaciones servidas por RefSeq para otras especies.
- Pista ENCODE Gencode: Siete isoformas analizadas por el consorcio ENCODE.
- Pista Human mRNA: Dos transcritos humanos obtenidos experimentalmente.
- Pista Non-Human mRNA: Dos transcritos experimentales de otras especies.
- Pista Vega: Siete formas alternativas reconocidas por el consorcio VEGA.

#### RefSeq

El proyecto RefSeq produce un conjunto depurado de los transcritos pertenecientes a todos los genes conocidos, evitando la redundancia y las ambigüedades existentes dentro del resto de bases de datos.

Figura 7. Listado de pistas para el gen *LRRTM1* (resumen).

## UCSC Genes

LRRTM1 (uc002sok.1) at chr2:80529003-80531487  
 - leucine rich repeat transmembrane neuronal 1  
 LRRTM1 (uc002soj.3) at chr2:80515481-80531487  
 - leucine rich repeat transmembrane neuronal 1

## RefSeq Genes

LRRTM1 at chr2:80529003-80531487  
 - (NM\_178839) leucine-rich repeat transmembrane neuronal

## Non-Human RefSeq Genes

LRRTM1 at chr2:80515605-80531552  
 - (NM\_001133111) leucine-rich repeat transmembrane neuronal

...

## ENCODE Gencode Manual Gene Annotations (May 2010)

LRRTM1 at chr2:80515483-80531518

...

## Human Aligned mRNA Search Results

BC045113 - Homo sapiens leucine rich repeat transmembrane neuronal 1, mRNA (cDNA clone MGC:47886 IMAGE:5200499), complete cds.

...

## Non-Human Aligned mRNA Search Results

BC027803 - Mus musculus, Lrrtml protein

...

## Vega Protein Coding Annotations

OTHUMT00000252294 at chr2:80515483-80531518

...

Como en este ejemplo, suele suceder con cierta frecuencia que para el mismo gen obtenemos anotaciones divergentes producidas por diferentes sistemas de anotación. Para acceder a la representación gráfica del genoma humano que contiene nuestro gen, vamos a seleccionar el enlace asociado a ese gen en la pista de RefSeq (RefSeq Genes en la figura 7), cuyas anotaciones constituyen un estándar de referencia reconocido actualmente por toda la comunidad científica debido a su precisión. Una vez el servidor reconoce nuestra petición, abre la pantalla de navegación principal (ver figura 8) para, automáticamente, enviarnos a la región (80,529,003-80,531,487) del cromosoma 2 que contiene nuestro gen (según el repositorio de RefSeq).

El visor genómico está centrado por defecto en la región que contiene exactamente todos los exones de nuestro gen. Observamos los dos exones del gen *LRRTM1* en forma de cajas rectangulares junto con otras pistas básicas de información (la pista de conservación entre vertebrados y los ARN mensajeros humanos disponibles). Para indicar el sentido de traducción de cada gen, el servidor introduce múltiples flechas a lo largo de los intrones. En este caso, la

**Primera conexión**

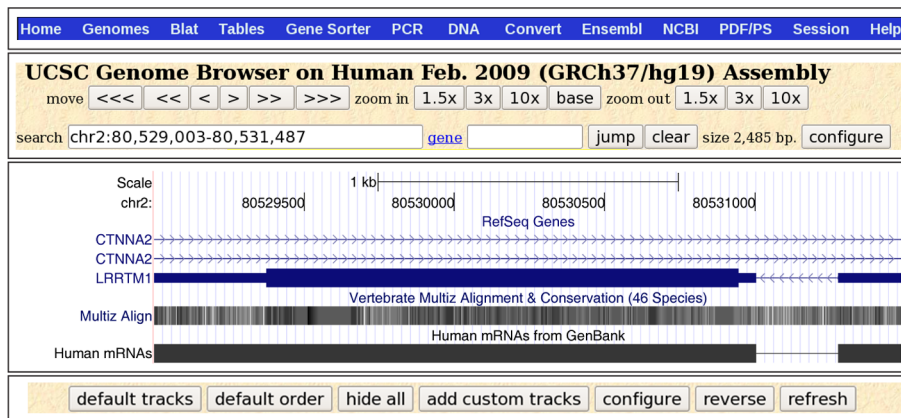
La primera conexión a UCSC suele mostrar un conjunto predefinido de pistas.

traducción de este gen debe realizarse de derecha a izquierda (hebra negativa de la molécula de ADN). En ambos extremos del gen, los exones presentan un trazo más fino para distinguir la fracción codificante (CDS, *coding sequence* en inglés) de la fracción no traducible del transcrito (UTR, *UnTRanslated*). Curiosamente, nuestro gen *LRRTM1* está ubicado en el interior del intrón de otro gen denominado *CTNNA2*, codificado en la hebra contraria.

### Ved también

La estructura de los genes se analiza con más detalle en la asignatura *Fundamentos de biología molecular*.

Figura 8. Pantalla inicial de navegación de UCSC.



Los navegadores suelen proporcionar un conjunto sencillo de botones de movimiento (figura 8). UCSC permite realizar movimientos utilizando:

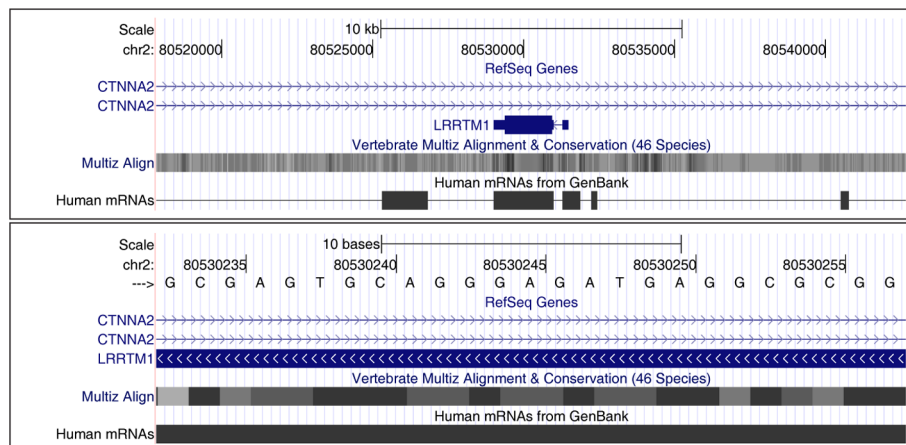
- Ampliaciones de posición con diferente resolución para acercar o alejar la escena.
- Desplazamientos horizontales de diferente tamaño a izquierda o derecha de la localización actual.
- Saltos hacia una región concreta mediante la caja de coordenadas.
- Búsquedas rápidas dirigidas por el identificador de cualquier elemento catalogado en este servidor.

Una vez accedemos a una pista concreta desde la pantalla inicial, el navegador nos muestra la anotación centrada ocupando toda la ventana gráfica. Para modificar la panorámica podemos alterar la perspectiva usando tres modos de resolución (1.5x, 3x y 10x). Por ejemplo, ocurre en ocasiones que es necesario utilizar *zoom out* para obtener una visión global del paisaje genómico que rodea a nuestro gen de interés. En otras circunstancias puede ser relevante acceder a fragmentos locales de las anotaciones empleando *zoom in*.

### Botón Base

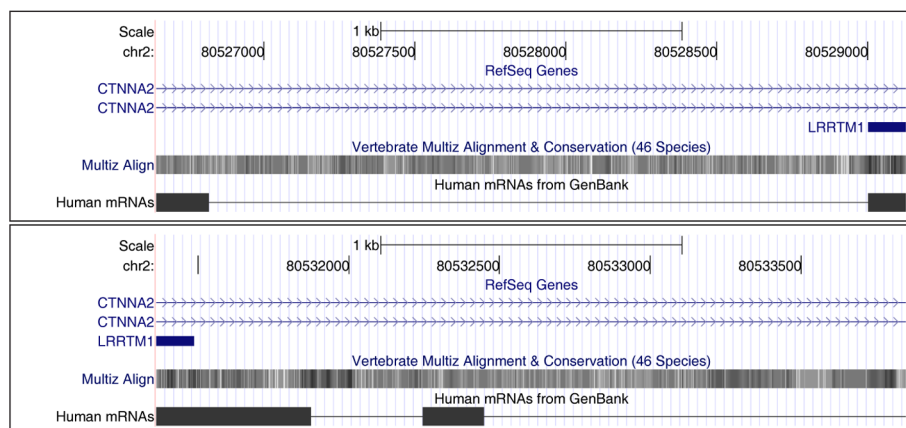
El botón Base intensifica al máximo el acercamiento para visualizar directamente la secuencia de nucleótidos.

Figura 9. Jugando con el enfoque de UCSC.



Los botones de desplazamiento (</<</<< y >/>>/>>) generan una nueva fotografía de las regiones del genoma inmediatamente anteriores o posteriores a esta ventana. Cada movimiento refresca automáticamente la imagen, actualizando el contenido de las pistas que estemos explorando. Si combinamos la herramienta de ampliación para fijar el ancho de la ventana gráfica (número de bases) con el desplazamiento de ésta (solapamiento entre la ventana actual y la siguiente), dispondremos de un mecanismo muy efectivo para visitar cada sección del cromosoma.

Figura 10. Desplazamiento por el genoma con UCSC.

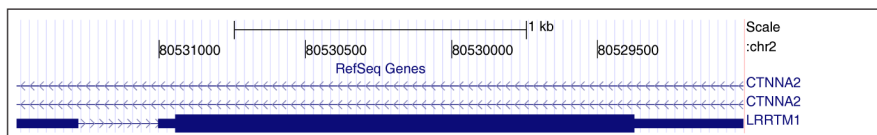


El botón Reverse, ubicado en la parte inferior (ver figura 8), permite obtener la anotación gráfica de la región actual seleccionando como referencia de coordenadas la hebra alternativa del cromosoma. Como podemos ver en la figura 11, con esta opción podemos visualizar nuestro gen *LRRTM1* mostrando los exones de acuerdo con su propia orientación (obteniendo la ordenación natural de izquierda a derecha con el margen izquierdo de la imagen indicando el punto de inicio de la transcripción del gen).

#### Botón Configure

El botón Configure permite el acceso a una pantalla donde podemos modificar los parámetros gráficos del visor genómico.

Figura 11. Cambiando la hebra de ADN con UCSC.

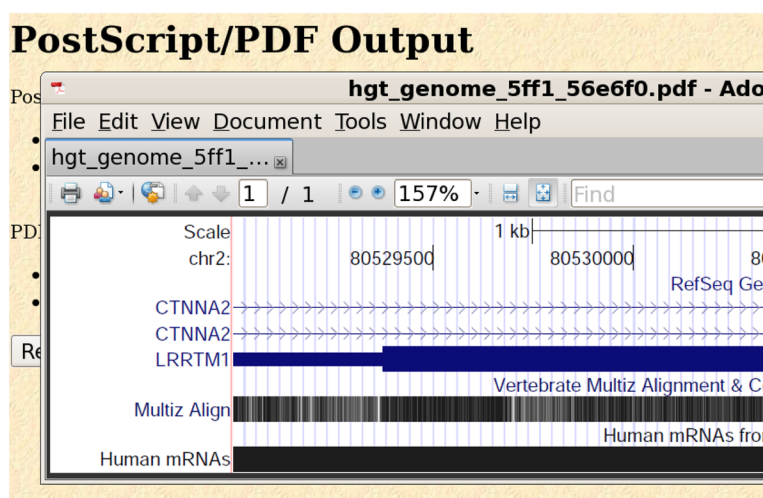


El menú azul superior de la pantalla principal (figura 8) contiene enlaces hacia varias aplicaciones complementarias. El botón denominado PDF/PS permite realizar una fotografía de la región actualmente visitada, incluyendo todas las pistas activadas en ese instante. Esta opción es verdaderamente útil para generar imágenes de alta calidad que podremos modificar para incluir en nuestras publicaciones científicas (figura 12).

#### Captura de imágenes

Esta opción, en combinación con la inclusión de nuestras propias pistas de datos, resulta especialmente útil para la composición de imágenes sobre anotaciones genómicas.

Figura 12. Fotografiando una región del genoma.



El marco de trabajo establecido por el navegador genómico delimita una serie de acciones que nos permiten modificar dinámicamente la representación gráfica de los elementos anotados en una región genómica. No debemos olvidar, sin embargo, que todas estas anotaciones en forma de pistas contienen implícitamente información sobre la localización exacta de dichos elementos. Para obtener directamente la secuencia de nucleótidos fotografiados en la actual ventana, podemos utilizar el botón denominado DNA, con el que accedemos a la secuencia real sobre la cual estamos superponiendo las anotaciones, junto con las regiones flanqueantes a ambos lados de ésta (ver figura 13).

Figura 13. Extracción de la secuencia de una región genómica.

**Get DNA in Window**

**Get DNA for**

Position

Note: if you would prefer to get DNA for features of a particular track or table, try the [Table Browser](#) using the output format sequence.

**Sequence Retrieval Region Options:**

Add  extra bases upstream (5') and  extra downstream (3')

Note: if a feature is close to the beginning or end of a chromosome and upstream/downstream bases are added, they may be truncated in order to avoid extending past the edge of the chromosome.

**Sequence Formatting Options:**

☒ All upper case.  
☐ All lower case.  
☐ Mask repeats: ☒ to lower case ☐ to N  
☐ Reverse complement (get '-' strand sequence)

Note: The "Mask repeats" option applies only to "get DNA", not to "extended case/color options".

3.2. Pistas

Las pistas contienen anotaciones suministradas por distintos consorcios y grupos de investigación. Para mostrar únicamente aquella información relevante en cada caso, el usuario puede configurar dinámicamente el inventario de pistas que desea estudiar, seleccionando el modo más conveniente de visualización para trabajar con cada una durante el análisis. La mayoría de los servidores genómicos permiten modificar la configuración actual de las pistas mediante un conjunto de opciones de visibilidad agrupadas en distintos bloques conceptuales (por ejemplo, genes, regulación, conservación, etc.).

En UCSC podemos modificar simultáneamente el comportamiento de varias pistas, siendo necesario presionar posteriormente el botón Refresh (refrescar) para reflejar el efecto final de la nueva configuración sobre la región mostrada en el visor. Para cada distribución de un genoma existe un conjunto determinado de pistas con las anotaciones registradas sobre su secuencia. En consecuencia, el navegador ofrece un conjunto de opciones ligeramente distinto entre diferentes genomas. No obstante, los bloques principales de opciones se suelen conservar dentro de los servicios ofrecidos por el mismo servidor. El listado de bloques proporcionados por UCSC para la versión hg19 del genoma humano puede encontrarse en la tabla 5. Utilizaremos como referencia esta distribución a lo largo de las explicaciones incluidas en esta sección.

Modificar el orden relativo entre pistas

El usuario puede modificar el orden relativo entre las pistas mostradas en el visor genómico de UCSC arrastrándolas verticalmente hacia su nueva posición.

Tabla 5. Bloques de pistas de UCSC.

Nombre	Descripción
Mapping and sequencing tracks	Ensamblado y secuenciación
Phenotype and disease associations	Fenotipos y enfermedades

Nombre	Descripción
<i>Genes and gene prediction tracks</i>	Genes y predicción génica
<i>Mrna and est tracks</i>	Tránscritos experimentales
<i>Expression</i>	Expresión génica
<i>Regulation</i>	Regulación génica
<i>Comparative genomics</i>	Genómica comparada
<i>Variation and repeats</i>	Polimorfismos y regiones repetitivas

En función del tipo de elemento biológico, el espacio físico en la pantalla o el tipo de análisis comparativo, en ocasiones debemos modificar la manera en que debe plasmarse gráficamente en nuestro navegador genómico la información sobre las anotaciones. Disponemos de cinco modos de visualización para mostrar una determinada pista sobre la región genómica actual (enumerados de menor a mayor cantidad de información):

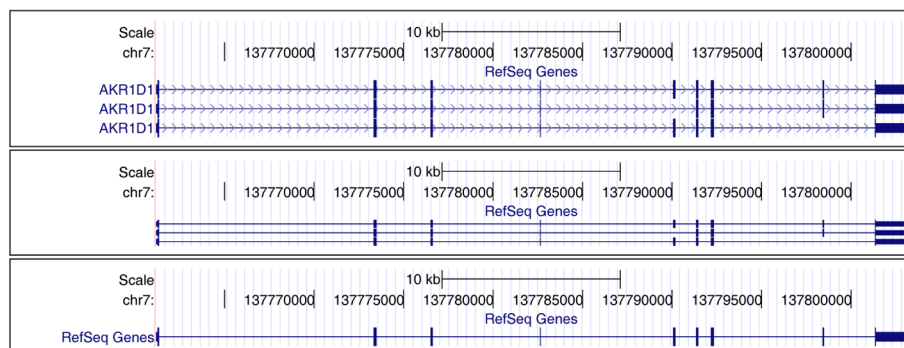
- *Hide* (oculto): no muestra las anotaciones (suele ser la opción por defecto).
- *Dense* (compacto): muestra las anotaciones comprimidas en una sola pista.
- *Squish* (reducido): muestra las anotaciones en varias pistas reducidas.
- *Pack* (agrupado): muestra las anotaciones en varias pistas agrupadas.
- *Full* (completo): muestra toda la información disponible para esa pista.

Cada una de estas opciones permite activar gráficamente un conjunto específico de informaciones sobre la pista de trabajo. La visualización de las distintas formas alternativas de un gen es un ejemplo típico que denota la importancia de realizar una selección apropiada de estos datos. En la figura 14 se puede verificar cómo podemos plegar y desplegar a nuestro gusto la representación gráfica de las tres isoformas anotadas para el gen humano *AKR1D1*, según el modo de visualización seleccionado para la pista RefSeq Genes.

#### Ved también

Ofrecemos más información sobre las formas alternativas de los genes en la asignatura *Fundamentos de biología molecular*.

Figura 14. Configurando la visualización de una pista.



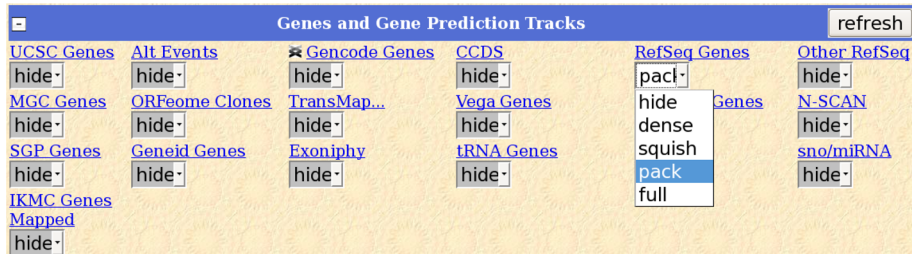
#### Leyenda figura 14

Primero desplegamos la información básica utilizando *pack*, después compactamos las tres isoformas cambiando a *squish* y finalmente plegamos las tres anotaciones en una sola pista con *dense*.

Para empezar, vamos a proceder a estudiar el bloque de opciones Genes and Gene Prediction Tracks (pistas de genes y predicción génica). Este bloque, mostrado en la figura 15, nos permite acceder a todos los datos de que disponen

diferentes repositorios sobre la región que estamos visualizando en un momento dado. En particular, podemos encontrar la pista de anotaciones RefSeq Genes producida por el consorcio RefSeq.

Figura 15. Bloque de pistas de anotación de genes de UCSC.



UCSC ofrece información específica sobre cada pista, accesible mediante el enlace coloreado en azul que muestra el nombre de cada pista. Por ejemplo, cuando seleccionamos RefSeq Genes en el bloque de opciones de información génica aparece una nueva pantalla que informa sobre el origen de esas anotaciones:

Figura 16. Información asociada a la pista RefSeq Genes.

RefSeq Genes Track Settings

[View table schema](#)

Data last updated: 2010-10-13

Description

The RefSeq Genes track shows known human protein-coding and non protein-coding genes taken from the NCBI RNA reference sequences collection (RefSeq). The data underlying this track are updated daily.

Methods

RefSeq RNAs were aligned against the human genome using blat; those with an alignment of less than 15% were discarded. When a single RNA aligned in multiple places, the alignment having the highest base identity was identified. Only alignments having a base identity level within 0.1% of the best and at least 96% base identity with the genomic sequence were kept.

Credits

This track was produced at UCSC from RNA sequence data generated by scientists worldwide and curated by the NCBI RefSeq project.

References

Kent WJ. BLAT - the BLAST-like alignment tool. *Genome Res.* 2002 Apr;12(4):656-64.

Pruitt KD, Tatusova T, Maglott DR. NCBI Reference Sequence (RefSeq): a curated non-redundant sequence database of genomes, transcripts and proteins. *Nucleic Acids Res.* 2005 Jan 1;33(Database issue):D501-4.

Para analizar los datos concretos asociados a las anotaciones gráficas mostradas en el visor, el usuario debe seleccionar con el ratón el elemento de la pista que desea explorar. Por ejemplo, podemos acceder a una fuente adicional de información asociada a cada anotación RefSeq cuando seleccionamos alguno de los exones del gen *LRRTM1* (figura 8). Podemos ver un resumen del contenido de esta entrada en la figura 17. Descendiendo desde la parte superior de la ficha, en primer lugar nos encontramos con el identificador asignado por RefSeq a este gen (NM\_178839). A continuación, el servidor nos proporciona los enlaces que ha recopilado para este gen en distintas bases de datos primarias (por ejemplo, Entrez, Pubmed, OMIM, GeneCards, etc.).

Figura 17. Anotaciones para la pista RefSeq Genes.

RefSeq Gene *LRRTM1*  
 RefSeq: [NM\\_178839.4](#) Status: Validated  
 Description: Homo sapiens leucine rich repeat  
                   transmembrane neuronal 1 (*LRRTM1*), mRNA.  
 CCDS: [CCDS1966.1](#)  
 CDS: [3' complete](#)  
 OMIM: [610867](#)  
 Entrez Gene: [347730](#)  
 PubMed on Gene: [LRRTM1](#)  
 PubMed on Product: [leucine-rich repeat transmembrane neuronal](#)  
 GeneCards: [LRRTM1](#)  
 AceView: [LRRTM1](#)  
 Stanford SOURCE: [NM\\_178839](#)

Summary of *LRRTM1*  
 Position: [chr2:80529003-80531487](#)  
 Band: 2p12  
 Genomic Size: 2485  
 Strand: -  
 Gene Symbol: *LRRTM1*  
 CDS Start: complete  
 CDS End: complete  
 Links to sequence:  
   \* [Predicted Protein](#)  
   \* [mRNA Sequence may be different from the genomic sequence](#)  
   \* [Genomic Sequence from assembly](#)  
   \* [CDS FASTA alignment from multiple alignment](#)

Data last updated: 2010-10-11

#### Visualización dinámica

Podemos ejecutar acciones diferentes cuando pinchamos en distintas partes de la imagen dinámica, en función de su ubicación.

#### Nota

Recomendamos al estudiante el análisis práctico detallado de todos los enlaces de esta entrada.

Antes de continuar analizando el contenido de la ficha de anotaciones que ofrece RefSeq para nuestro gen, proponemos al estudiante profundizar el recorrido en alguno de estos repositorios. El enlace hacia Entrez Gene (347730), por ejemplo, nos permite acceder a otra ficha de información recopilada por este recurso también gestionado, como RefSeq, por el instituto National Center for Biotechnology Information (NCBI, centro nacional para la información biotecnológica de Estados Unidos). Apreciamos en esta nueva ficha que cada anotación recopilada sobre el gen *LRRTM1* posee un enlace que nos permite llegar a la fuente original de esos datos. En la figura 18 mostramos únicamente el resumen y la bibliografía referente a este gen, omitiendo el resto de infor-

#### Lectura complementaria

D. L. Wheeler y otros (2004). "Database resources of the NCBI: update". *Nucleic Acids Research* (núm. 32, págs. D35-D40).

maciones. El usuario puede estudiar la estructura exónica del gen o acceder al navegador genómico suministrado por el propio NCBI para navegar por la región colindante.

Figura 18. Información en Entrez sobre el gen *LRRTM1*.

**LRRTM1 leucine rich repeat transmembrane neuronal 1 [ *Homo sapiens* ]**  
Gene ID: 347730, updated on 19-Sep-2010

**Summary**

**Official Symbol** LRRTM1 provided by [HGNC](#)

**Official Full Name** leucine rich repeat transmembrane neuronal 1 provided by [HGNC](#)

**Primary source** [HGNC:19408](#)

**Locus tag** UNQ675/PRO1309

**See related** [Ensembl:ENSG00000162951](#); [HPRD:14323](#); [MIM:610867](#)

**Gene type** protein coding

**RefSeq status** VALIDATED

**Organism** [Homo sapiens](#)

**Lineage** Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Haplorrhini; Catarrhini; Hominidae; Homo

**Also known as** FLJ32082; LRRTM1

**Bibliography**

**Related articles in PubMed**

1. [Polymorphisms in leucine-rich repeat genes are associated with autism spectrum disorder susceptibility in populations of European ancestry.](#) Sousa I, et al. Mol Autism, 2010 Mar 25. PMID 20678249.
2. [Supporting evidence for LRRTM1 imprinting effects in schizophrenia.](#) Ludwig KU, et al. Mol Psychiatry, 2009 Aug. PMID 19626025.
3. [Understanding the genetics of behavioural and psychiatric traits will only be achieved through a realistic assessment of their complexity.](#) Francks C. Laterality, 2009 Jan. PMID 19125367.
4. [Where and what is the right shift factor or cerebral dominance gene? A critique of Francks et al. \(2007\).](#) Crow TJ, et al. Laterality, 2009 Jan. PMID 19125366.
5. [Editorial commentary: is LRRTM1 the gene for handedness?](#) McManus C, et al. Laterality, 2009 Jan. PMID 19125365.

[See all \(13\) citations in PubMed](#)

Especialmente interesante por su importancia en estudios de análisis masivo de datos genómicos resulta la anotación de las funciones biológicas que desempeña este gen en la célula que podemos ver en la tabla 6:

Tabla 6. Funciones biológicas del gen *LRRTM1*.

Función	Evidencia
<i>Protein binding</i>	computacional
<i>Axon</i>	bibliográfica
<i>NOT cell surface</i>	bibliográfica
<i>Endoplasmic reticulum</i>	bibliográfica
<i>Endoplasmic reticulum membrane</i>	computacional
<i>Growth cone</i>	bibliográfica
<i>Integral to membrane</i>	computacional
<i>Membrane</i>	computacional

Desde aquí podemos acceder a la secuencia de este transcrito depositada en GenBank. Este repositorio pionero ofrece información sobre millones de secuencias genómicas y proteómicas. Cualquier ficha de GenBank está estructu-

### Gene Ontology

Gene Ontology es un diccionario de términos que permite describir cientos de eventos biológicos distintos. Estas definiciones se emplean para caracterizar los genes en base a evidencias experimentales.

### Manipular las secuencias con GenBank

GenBank proporciona también herramientas para manipular las secuencias (por ejemplo, extraer subsecuencias dentro de un cierto rango de coordenadas).

rada siguiendo un formato notablemente rígido. En cada línea, el tipo de característica anotada se indica primero, mientras que su valor concreto se escribe a continuación. Estas fichas incluyen generalmente información descriptiva y anotaciones sobre una determinada región para acabar con la secuencia completa de ésta. Podemos apreciar en la figura 19 las coordenadas de los dos exones del ARN mensajero del gen *LRRTM1* (1..211 y 212..2212), así como su fracción codificante (271..1839).

Figura 19. Entrada de la secuencia NM\_178839 en GenBank.

```

LOCUS       NM_178839 2217 bp mRNA linear PRI 19-SEP-2010
DEFINITION  Homo sapiens leucine rich repeat
              transmembrane neuronal 1 (LRRTM1),
              mRNA.
ACCESSION   NM_178839
VERSION     NM_178839.4  GI:86990455
SOURCE      Homo sapiens (human)
  ORGANISM  Homo sapiens
            Eukaryota; Metazoa; Chordata; Craniata; Vertebrata;
            Euteleostomi; Mammalia; Eutheria; Euarchontoglires;
            Primates; Haplorrhini;
            Catarrhini; Hominidae; Homo.
REFERENCE   2 (bases 1 to 2217)
  AUTHORS   Ludwig,K.U., Mattheisen,M., Muhleisen,T.W., Roeske,D.,
            Schmal,C., Breuer,R., Schulte-Korne,G., Muller-Myhsok,B.,
            Nothen,M.M., Hoffmann,P., Rietschel,M. and Cichon,S.
  TITLE     Supporting evidence for LRRTM1 imprinting effects
            in schizophrenia
  JOURNAL   Mol. Psychiatry 14 (8), 743-745 (2009)
  PUBMED    19626025
  REMARK    GeneRIF: Results strengthen the evidence for an association
            of imprinted alleles of LRRTM1 with schizophrenia. Weaker
            supportive evidence was also obtained for a possible
            association of LRRTM1 with human brain asymmetry.
FEATURES             Location/Qualifiers
     source           1..2217
                     /organism=Homo sapiens
                     /mol_type=mRNA
                     /db_xref=taxon:9606
                     /chromosome=2
                     /map=2p12
     gene             1..2217
                     /gene=LRRTM1
                     /gene_synonym=FLJ32082
                     /note=leucine rich repeat transmembrane neuronal 1
     exon             1..211
                     /gene=LRRTM1
                     /gene_synonym=FLJ32082
                     /inference=alignment:Splice
                     /number=1
     exon             212..2212
                     /gene=LRRTM1
                     /gene_synonym=FLJ32082
                     /inference=alignment:Splice
                     /number=2
     CDS              271..1839
                     /gene=LRRTM1
                     /translation=MDFLLLGLCLYWLLRRPSGVVLCLLGACFQMLPAAPS
                     GCPQLCRCEGRLLYCEALNLTEAPHNLSGLLGLSLRYNSLSELRAQGFTG
                     ...
                     QHGTGTFEPATVALPGGEHAENAVQIHKVVTGTMALIFSFLIVVLVLYVSW
                     KCFPASLRQLRQCFVTQRRKQKQKQTMHQMAAMSAQEYYVDYKPNHIEGA
                     LVIINEYGSCCTCHQQPARECEV
ORIGIN
  1 ttgcgggtcc tagaagtcgc ctccccgcct tgccggccgc...
```

Dentro de la ficha de anotaciones de RefSeq que proporciona el navegador UCSC (ver figura 17), observamos también una interfaz denominada Links to sequence (enlaces a la secuencia) que permite extraer directamente la secuencia de nucleótidos y aminoácidos de cada gen. Mediante el enlace Predicted protein (proteína predicha) podemos ver el producto de la traducción de este gen, en formato FASTA (figura 20).

Figura 20. Secuencia de la proteína *LRRTM1*.

```
>NP_849161 length = 522
MDFLLGLCLYWLRLRRPSGVVLCILGACFQMLPAAPSGCPQLCRCEGRLLYCEALNLTEAPHNLSGLLG
LSLRYNLSSELRAQFTGLMQLTWLYLDHNHICSVQGDFAFKLRRVKELTLSSNQITQLPNTTFRMPN
LRVDLSYNKLQALAPDLFHGLRKLTTLHMRANAIQFVPVRIQDCRSLKFLDIGYNQLKSLARNSFAG
LFKLTELHLEHNDLVKVNFAHFPRILSLHSLCLRRNKVAIVVSSLDWVWNLEKMDLSGNEIEYMEPHVF
ETVPHLQSLQLDSNRLTYIEPRILNSWKSLSITLAGNLWDCGRNVICALASWLNNFQGRYDGNLQCASP
EYAQGEDVLDVAVAFHLCEDGAEPSTGHLLSAVTNRSDLGPPASSATTLADGGEGQHDGTTFEPATVALP
GGEHAENAVQIHKVVGTMALIFSLIVLVLYVSWKCFPASLRQLRQCFVTQRRKQKQKQTMHQMAM
SAQEYVVDYKPNHIEGALVIINEYGSCTCHQQPARECEV
```

#### FASTA

Una secuencia en formato FASTA está compuesta por un encabezamiento, cuyo primer carácter es el símbolo ">", que informa de su origen biológico, seguida de la propia secuencia, agrupada en líneas con el mismo número de caracteres.

Como comentábamos anteriormente, podemos extraer diferentes partes de la secuencia de este gen, haciendo uso en este caso de la opción Genomic sequence from assembly (secuencia genómica del ensamblado). El navegador genómico nos proporciona ahora el acceso a una nueva pantalla, donde podemos elegir los componentes de la estructura génica que más nos interesen (figura 21). De este modo, es posible obtener un fragmento de la secuencia promotora del gen (Upstream), de la región no traducible inicial o final (5'UTR y 3'UTR), de la región codificante (CDS), o de la secuencia del único intrón de este gen (Introns). Cuando activamos varias opciones simultáneamente podemos combinar las secuencias resultantes que constituyen el propio gen. Cada elemento génico puede separarse en diferentes secuencias en formato FASTA o integrarse en una única secuencia de salida para su posterior análisis.

Figura 21. Extraer diferentes regiones de un gen con UCSC.

### Genomic Sequence Near Gene

#### Get Genomic Sequence Near Gene

Note: if you would prefer to get DNA for more than one feature of this track at a time, try the [Table Browser](#) using the output format sequence.

**Sequence Retrieval Region Options:**

☐ Promoter/Upstream by  bases

☒ 5' UTR Exons

☒ CDS Exons

☒ 3' UTR Exons

☐ Introns

☐ Downstream by  bases

☒ One FASTA record per gene.

☐ One FASTA record per region (exon, intron, etc.) with  extra bases upstream (5') and  extra downstream (3')

☐ Split UTR and CDS parts of an exon into separate FASTA records

Note: if a feature is close to the beginning or end of a chromosome and upstream/downstream bases are added, they may be truncated in order to avoid extending past the edge of the chromosome.

**Sequence Formatting Options:**

☒ Exons in upper case, everything else in lower case.

☐ CDS in upper case, UTR in lower case.

☐ All upper case.

☐ All lower case.

☐ Mask repeats: ☒ to lower case ☐ to N

En la figura 22 mostramos la secuencia completa del gen *LRRTM1*, presentada en el sentido de lectura original para facilitar la comparación con otros genes anotados en la hebra positiva. Cada exón de este gen está incluido en un fichero FASTA distinto. Utilizamos letras minúsculas para denotar la región no traducible y letras mayúsculas para indicar el fragmento codificante, que empieza con el codón ATG y finaliza con el codón TGA. El navegador UCSC genera suficiente información en el encabezamiento de cada secuencia para una fácil identificación. Podemos reconocer aquí la distribución del genoma (hg19), la pista (refGene equivale a RefSeq Genes), el código RefSeq (NM\_178839), el exón (0 o 1) y su ubicación genómica.

#### Datos en cualquier orientación

UCSC realiza automáticamente la operación de complementar las bases e invertir su secuencia para mostrar los datos en cualquier orientación.

Figura 22. Secuencia del transcrito del gen *LRRTM1*.

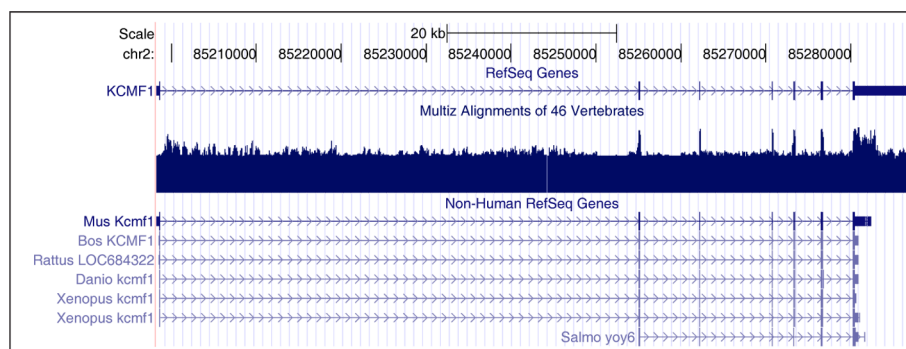
```
>hg19_refGene_NM_178839_0 range = chr2:80531277-80531487
ttgcgggtcctagaaagtcgcctccccgccttgccggccgccttgagccccgagccgagcag
caaagtgagacattgtgagcctgccagatccgcggccgcggaccggggctgcctcggaaca
cagaggggtcttctctcgccctgcataataattagcctgcacacaaagggagcagctgaatgga
ggtgtcactctctggaagag
>hg19_refGene_NM_178839_1 range = chr2:80529003-80531003
atttctgaccgagcgcttccaatggacattctccagtctctctggaagattctcgctaATGG
ATTTCTGCTGCTCGGTCTCTGTCTATACTGGCTGCTGAGGAGGCCCTCGGGGTGGTCTTGT
GTCTGCTGGGGGCTGCTTTTCTAGATGCTGCCCCGCCGCCAGCGGTGCCGCAGCTGTGCC
GGTGCAGAGGGCGGCTGCTGTACTGCGAGGCGCTCAACCTCACCGAGGCGCCCCACAACCTGT
CCGGCCTGCTGGGCTGTCTCTGCGCTACAACAGCCTCTCGGAGCTGCGCGCCGCCAGTTCA
CGGGGTTAATGCAGCTCACGTGGCTCTATCTGGATCACAATCACATCTGCTCCGTGCAGGGGG
ACGCCTTTTCAGAACTGCGCCGAGTTAAGGAACCTCACGTGAGTTCCAACCAGATCACCCAAC
TGCCCAACACCACCTTCCGGCCCATGCCCAACCTGCGCAGCGTGGACCTCTCGTACAACAAGC
TGCAGGCGCTGCGCCCCGACCTCTTCCACGGGCTGCGGAAGCTCACACGCTGCATATGCGGG
CCAACGCCATCCAGTTTGTGCCCGTGCATCTTCCAGGACTGCCGCAGCCTCAAGTTTCTCG
ACATCGGATACAATCAGCTCAAGAGTCTGGCGCGCAACTCTTTCGCCGGCTTGTTTAAGCTCA
CCGAGCTGCACCTCGAGCACAACGACTTGGTCAAGGTGAACTTCGCCCCTTCCCGCGCCTCA
TCTCCCTGCACTCGCTCTGCCTGCGGAGGAACAAGGTGGCCATTGTGGTCACTCGCTGGACT
GGGTTTGAACCTGGAGAAAATGGACTTGTGCGGCAACGAGATCGAGTACATGGAGCCCCATG
TGTTTCGAGACCGTGCCGCACCTGCAGTCCCTGCAGCTGGACTCCAACCGCCTCACCTACATCG
AGCCCCGGATCCTCAACTCTTGAAGTCCCTGACAAGCATACCCTGGCCGGGAACCTGTGGG
ATTGCGGGCGCAACGTGTGTGCCCTAGCCTCGTGGCTCAACAACCTTCAGGGGCGCTACGATG
GCAACTTGCACTGCGCCAGCCGAGTACGCACAGGGCGAGGACGTCCTGGACGCCGTGTACG
CCTTCCACCTGTGCGAGGATGGGGCCGAGCCACCAGCGGCCACCTGCTCTCGGCCGTACCA
ACCGCAGTGATCTGGGGCCCCCTGCCAGCTCGGCCACCACGCTCGCGGACGGCGGGAGGGGC
AGCACGACGGCACATTCGAGCCTGCCACCGTGGCTCTTCCAGGCGGCGAGCACGCCGAGAAGC
CCGTGCAGATCCACAAGGTGGTCACGGGACCATGGCCCTCATCTTCTCCTTCTCATCGTGG
TCCTGGTGCTCTACGTGTCTGGAAGTGTTCCTCCAGCCAGCCTCAGGCAGCTCAGACAGTGCT
TTGTACGCAGCGCAGGAAGCAAAAGCAGAAACAGACCATGCATCAGATGGCTGCCATGTCTG
CCCAGGAATACTACGTTGATTACAAACCGAACCACATTGAGGGAGCCCTGGTGATCATCAACG
AGTATGGCTCGTGATCTGCCACCAGCAGCCCGGAGGGAATGCGAGGTGTGAttgtccaggt
ggctctcaacccatgcgctaccaataacgctgggcagccgggaagggccggcgggcaccagg
ctgggggtctccttctgtctgtctgatgtctccttgactgaaactttaaggggatctctccc
agagacttgacatttttagctttattgtgtcttaaaaaacaaaagcgaattaaaacacacaaaa
aaccaccacccacacaccttcaggacagctctatcttaatttcataatgagaactccttctccc
tttgaagatctgtccatattcaggaatctgagagtgtaaaaaagggtaccaatcattgattttt
ttttttttgtgaaactaaaatgttttaaaataaaatagcatttacagt
```

Junto con los transcritos humanos revisados manualmente, el proyecto RefSeq produce anotaciones de referencia para varias especies de vertebrados. Esta información resulta extremadamente útil para identificar regiones funcionales conservadas a lo largo de la evolución. De este modo, mediante las oportunas comparaciones entre las distintas anotaciones de RefSeq, el propio navegador UCSC habilita una pista adicional dentro del bloque de anotaciones génicas denominada Other RefSeq para reconocer las formas ortólogas de cada gen en otros genomas (figura 23). Con esta información, el usuario debe cambiar el genoma de referencia dentro del navegador y posteriormente utilizar el nuevo código suministrado por RefSeq para recuperar la secuencia homóloga de dicho ARN mensajero.

#### Ved también

Los alineamientos de secuencias se explican detalladamente en el módulo "Comparación de secuencias".

Figura 23. Formas ortólogas del gen humano *KCMF1* en RefSeq.



#### Leyenda figura 23

Pueden apreciarse nítidamente los picos de conservación sobre las regiones exónicas del gen en todas las especies.

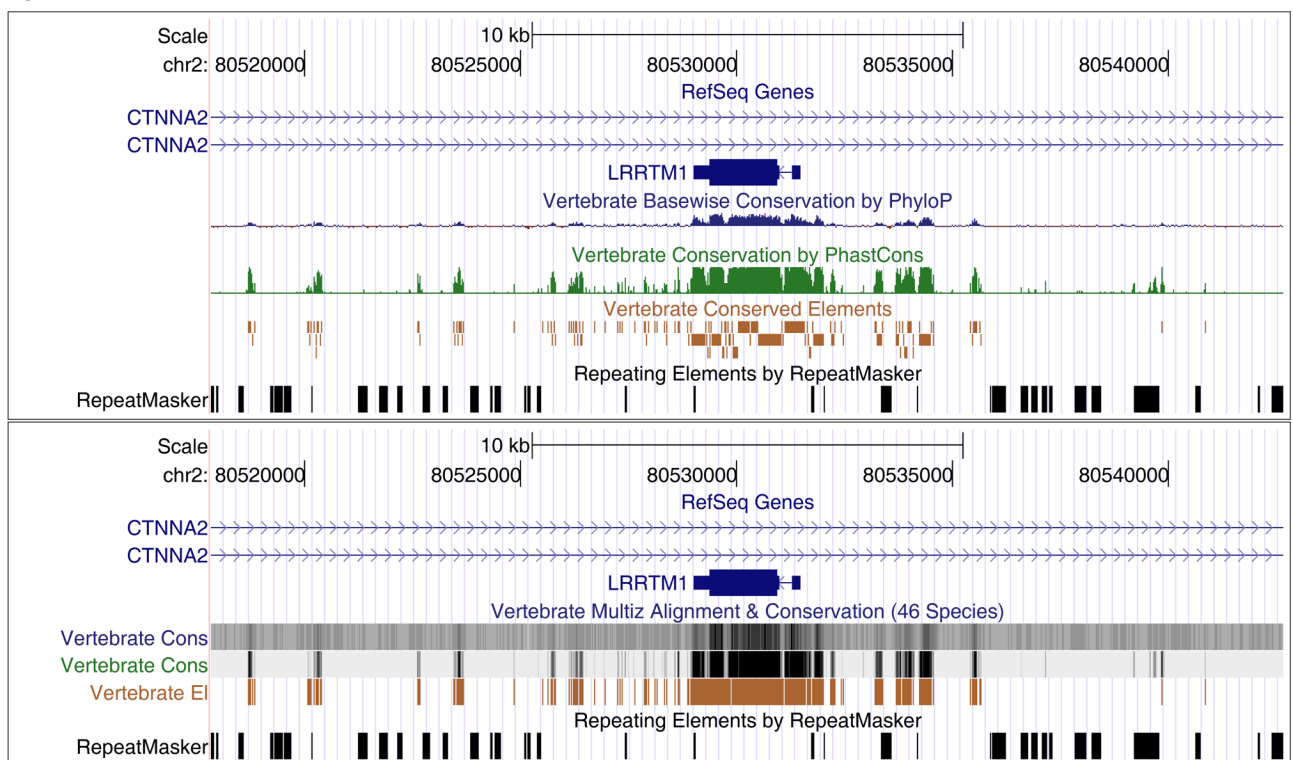
Los bloques de pistas de genómica comparada (en inglés, *comparative genomics*) y de predicción génica son las secciones de anotaciones suministradas por el servidor genómico UCSC más frecuentemente utilizadas por los analistas bioinformáticos. Como podemos ver en la figura 24, disponemos de múltiples comparaciones entre la secuencia que estamos utilizando de referencia (humano, en este caso la distribución NCBI36/hg18) y los genomas ensamblados de otros vertebrados. La comparación entre organismos cuya evolución difiere a partir de determinados puntos del árbol filogenético permite descubrir aquellos elementos funcionales conservados entre distintas secuencias. Para proporcionar estos datos con celeridad, todos los análisis están precalculados de antemano en el interior del servidor. Con estas pistas, por ejemplo, podemos establecer la existencia de regiones codificantes a través de los alineamientos de las secuencias homólogas (ver figura 25).

Figura 24. Bloque de genómica comparada.

Comparative Genomics refresh

Conservation	28-Way Cons	28-Way Base Cons	28-Way Most Cons	17-Way Cons	17-Way Most Cons
<a href="#">hide</a>	<a href="#">hide</a>	<a href="#">hide</a>	<a href="#">hide</a>	<a href="#">hide</a>	<a href="#">hide</a>
<a href="#">Cons Indels</a>	<a href="#">Tetraodon Ecores</a>	<a href="#">Chimp Chain/Net</a>	<a href="#">Orangutan Chain/Net</a>	<a href="#">Rhesus Chain/Net</a>	<a href="#">Marmoset Chain/Net</a>
<a href="#">MmCf</a>	<a href="#">hide</a>	<a href="#">hide</a>	<a href="#">hide</a>	<a href="#">hide</a>	<a href="#">hide</a>
<a href="#">Guinea pig Chain/Net</a>	<a href="#">Rat Chain/Net</a>	<a href="#">Mouse Chain/Net</a>	<a href="#">Dog Chain/Net</a>	<a href="#">Cat Chain/Net</a>	<a href="#">Horse Chain/Net</a>
<a href="#">hide</a>	<a href="#">hide</a>	<a href="#">hide</a>	<a href="#">hide</a>	<a href="#">hide</a>	<a href="#">hide</a>
<a href="#">Cow Chain/Net</a>	<a href="#">Opossum Chain/Net</a>	<a href="#">Platypus Chain/Net</a>	<a href="#">Lizard Chain/Net</a>	<a href="#">Zebra finch Chain/Net</a>	<a href="#">Chicken Chain/Net</a>
<a href="#">hide</a>	<a href="#">hide</a>	<a href="#">hide</a>	<a href="#">hide</a>	<a href="#">hide</a>	<a href="#">hide</a>
<a href="#">X. tropicalis Chain/Net</a>	<a href="#">Zebrafish Chain/Net</a>	<a href="#">Medaka Chain/Net</a>	<a href="#">Stickleback Chain/Net</a>	<a href="#">Fugu Chain/Net</a>	<a href="#">Tetraodon Chain/Net</a>
<a href="#">hide</a>	<a href="#">hide</a>	<a href="#">hide</a>	<a href="#">hide</a>	<a href="#">hide</a>	<a href="#">hide</a>
<a href="#">Lamprey Chain/Net</a>	<a href="#">Lancelet Chain/Net</a>	<a href="#">S. purpuratus Chain/Net</a>			
<a href="#">hide</a>	<a href="#">hide</a>	<a href="#">hide</a>			

Figura 25. Conservación de secuencias funcionales entre vertebrados.



Con el navegador UCSC se pueden configurar numerosos parámetros para incluir diferentes comparaciones. Para ello, el usuario debe pinchar sobre el enlace coloreado en azul que indica el nombre de cada pista de este bloque. La pista Conservation contiene alineamientos entre la mayoría de los grupos de especies más representativos. En su interior tenemos, por ejemplo, la opción de incorporar o retirar genomas en el visor genómico, modificar el color del gráfico de picos en función de las pautas de lectura o activar diferentes modos de alineamiento para la identificación de regiones conservadas:

Figura 26. Configuración de la pista Conservation.

**Multiz Alignments Configuration**

**Species selection:**

**Primate**

☒ chimp ☒ gorilla ☒ orangutan ☒ rhesus ☒ baboon  
☒ marmoset ☒ tarsier ☒ mouse lemur ☒ bushbaby

**Placental Mammal**

☒ tree shrew ☒ mouse ☒ rat ☒ kangaroo rat ☒ guinea pig  
☒ squirrel ☒ rabbit ☒ pika ☒ alpaca ☒ dolphin  
☒ cow ☒ horse ☒ cat ☒ dog ☒ microbat  
☒ megabat ☒ hedgehog ☒ shrew ☒ elephant ☒ rock hyrax  
☒ tenrec ☒ armadillo ☒ sloth

**Vertebrate**

☒ wallaby ☒ opossum ☒ platypus ☒ chicken ☒ zebra finch  
☒ lizard ☒ x. tropicalis ☒ tetraodon ☒ fugu ☒ stickleback  
☒ medaka ☒ zebrafish ☒ lamprey

**Multiple alignment base-level:**  
☐ Display bases identical to reference as dots  
☒ Display chains between alignments

**Codon Translation:**  
 Default species to establish reading frame:   
☐ No codon translation  
☒ Use default species reading frames for translation  
☐ Use reading frames for species if available, otherwise no translation  
☐ Use reading frames for species if available, otherwise use default species

**Phylogenetic Tree:**

- Primates
  - Human
  - Chimp
  - Gorilla
  - Orangutan
  - Rhesus
  - Baboon
  - Marmoset
  - Tarsier
  - Mouse lemur
  - Bushbaby
  - Tree shrew
- Placental mammals
  - Mouse
  - Rat
  - Kangaroo rat
  - Guinea Pig
  - Squirrel
  - Rabbit
  - Pika
  - Alpaca
  - Dolphin
  - Cow
  - Horse
  - Cat
  - Dog
  - Microbat
  - Megabat
  - Hedgehog
  - Shrew
  - Elephant
  - Rock hyrax
  - Tenrec
  - Armadillo
  - Sloth
  - Wallaby
  - Opossum
  - Platypus
  - Chicken
  - Zebra finch
  - Lizard
  - X.tropicalis
  - Tetraodon
  - Fugu
  - Stickleback
  - Medaka
  - Zebrafish
  - Lamprey
- Vertebrates
  - Lamprey

La herramienta BLAT (accesible desde la barra de menú superior de color azul) complementa perfectamente el bloque de pistas de conservación. BLAT (*Blast-like alignment tool*) es un programa de alineamiento de secuencias que identifica las regiones más similares a aquella propuesta por el usuario en un determinado genoma. A diferencia de otros programas de alineamiento que veremos más adelante, BLAT está especializado en la detección de secuencias prácticamente idénticas (hecho que aumenta su velocidad de respuesta de forma notable). Podemos usar esta aplicación para identificar, por ejemplo, cuál es la ubicación de cualquier secuencia codificada dentro del mismo genoma.

### Lectura complementaria

W. J. Kent (2002). "BLAT: the BLAST-like alignment tool". *Genome Research* (núm. 12, págs. 656-664).


Figura 27. Funcionamiento de BLAT.

**BLAT Search Genome**

Genome:  Assembly:  Query type:

☒ CDS  
 ATGGATTTCCTGCTGCTCGGTCTCTGTCTACTAGGCTGCTGAGGAGGCC

Scale 20 bases  
 chr2: | 80530900 | 80530905 | 80530910 | 80530915 | 80530920 | 80530925 | 80530930 | 80530935 | 80530940 |  
 ---> G G C C T C C T C A G C A G C C A G T A T A G A C A G A G A C C G A G C A G C A G G A A T C C A T  
 Your Sequence from Blat Search

Blat Sequence 

RefSeq Genes

CTNNA2  
 CTNNA2  
 LRRTM1

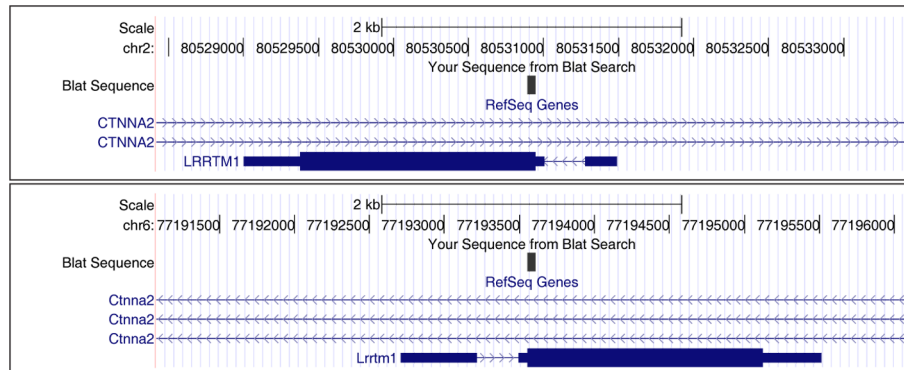
### Leyenda figura 27

Estamos identificando el inicio de la región codificante del gen *LRRTM1*.

BLAT proporciona una clasificación de posibles ubicaciones de nuestra secuencia ordenadas por similitud. Aquella secuencia que seleccionamos como resultado es posteriormente integrada como una nueva pista dentro del navegador genómico. El empleo de BLAT resulta especialmente interesante para loca-

lizar la posición de la región más similar a nuestra secuencia en otro genoma. Con este procedimiento podemos detectar, para el gen humano *LRRTM1*, su correspondiente ortólogo anotado dentro del genoma de ratón:

Figura 28. Utilizando BLAT para identificar genes ortólogos.



#### Leyenda figura 28

Presentamos el genoma humano en la parte superior de la imagen y el genoma de ratón en la parte inferior.

### 3.3. Anotaciones propias

Los servidores genómicos proporcionan un enorme volumen de información. El estudio exhaustivo de la cartografía existente para los distintos elementos funcionales codificados en el genoma resulta fundamental para abordar cualquier problema biológico. Cualquier grupo de investigación, sin embargo, genera habitualmente sus propios resultados durante el proceso de elaboración de hipótesis en función de los experimentos llevados a cabo. En consecuencia, es fundamental que el bioinformático conozca los entresijos del manejo de estos navegadores genómicos para enriquecer las informaciones originales con los resultados experimentales obtenidos en el laboratorio.

La introducción de anotaciones propias dentro de un navegador genómico resulta esencial para realizar comparaciones entre el conocimiento existente y los nuevos resultados experimentales, favoreciendo además la generación de representaciones gráficas de alta calidad.

Una vez el usuario introduce sus propias pistas, éstas son automáticamente integradas dentro de la imagen proporcionada por el servidor. Este hecho abre todo un abanico de opciones de trabajo para realizar el análisis comparativo de las pistas. Podemos importar múltiples pistas propias dentro del navegador, utilizando el servidor para estudiar el comportamiento de nuestros resultados. De este modo, todas las herramientas disponibles para llevar a cabo comparaciones entre distintas anotaciones, tanto a nivel cualitativo como cuantitativo, pueden aprovecharse directamente para contrastar el conocimiento de referencia con nuevos datos experimentales observados para cualquier característica genómica.

Dentro del entorno de trabajo del navegador UCSC, las pistas introducidas en el sistema por el usuario reciben el nombre de *Custom tracks* (en inglés, *pistas adaptadas*). Desde la pantalla principal del visor genómico, el usuario debe presionar el botón Add custom tracks para incorporar una nueva pista de datos. El navegador UCSC asocia un código interno a la conexión realizada desde nuestro ordenador, gestionando internamente las pistas estándar que estamos visualizando y el conjunto de pistas propias que hemos cargado en el sistema. De este modo, estos resultados únicamente pueden ser visualizados desde nuestro ordenador, no siendo accesibles para el resto de la comunidad de usuarios.

Una vez la operación de carga de nuestra pista finaliza satisfactoriamente, podemos visualizar todos sus elementos a lo largo del genoma junto con el resto de anotaciones convencionales. Esta integración completa de los nuevos datos resulta muy interesante para futuras operaciones. Por regla general, el navegador mantiene en memoria nuestras pistas durante uno o dos días únicamente, para no almacenar un volumen excesivo de datos. Dado que estos datos son procesados como cualquier pista, podemos hacer uso de las opciones habituales. Resultan particularmente útiles para la mayoría de usuarios la generación de figuras en alta resolución y la comparación con otras bases de datos primarias.

El usuario puede crear sus propias pistas en numerosos formatos dentro de UCSC. En estos materiales vamos a focalizar nuestro interés en los formatos más populares para la integración de anotaciones genómicas propias. Previamente a la carga en el sistema, el navegador UCSC procesa sintácticamente nuestras pistas para verificar su corrección. En caso de detectarse un error de formato, el navegador nos informará apropiadamente para subsanarlo con facilidad.

Cada pista debe poseer la siguiente estructura dividida en tres componentes.

- Configuración del navegador: listado de pistas visibles.
- Configuración de la pista adaptada: parámetros de visualización.
- Anotaciones de la pista adaptada: valores identificados.

En primer lugar, el usuario debe configurar el comportamiento inicial del navegador empleando la palabra clave *browser*. Para ello, es necesario introducir las coordenadas de la región genómica e indicar el conjunto de pistas convencionales que deseamos activar para comparar con nuestras anotaciones (utilizando distintos modos de visualización).

#### Nota

El estudiante puede encontrar abundante información sobre la carga de nuevas pistas dentro de las páginas de ayuda del navegador UCSC.

#### Sesiones protegidas

El servidor UCSC ofrece también un sistema de sesiones de usuario protegidas bajo contraseña, facilitando la comunicación de datos entre diversos colaboradores.

#### Ficheros comprimidos

UCSC permite la carga directa de ficheros comprimidos para facilitar una rápida transferencia de los datos a través de la Red.

#### Ved también

La asignatura *Fundamentos de informática en entornos bioinformáticos* instruye sobre la manipulación de ficheros con LINUX.

Figura 29. Encabezamiento de pistas propias (1).

```
browser position cromosoma:inicio-final
browser [hide/dense/pack/full] pista 1
browser [hide/dense/pack/full] pista 2
...
browser [hide/dense/pack/full] pista n
```

A continuación, dentro del mismo fichero, introduciendo la palabra clave *track* debemos definir los parámetros de visualización de la nueva pista: asignar un nombre y una descripción, definir un modo de visibilidad o seleccionar un color.

Figura 30. Encabezamiento de pistas propias (2).

```
track name=nombre description=descripcion visibility=[0-4] color=R,G,B
```

Finalmente, debemos aportar las coordenadas de nuestras anotaciones genómicas para integrarlas en la visualización final. Cada línea corresponde exclusivamente a una anotación efectuada por el usuario, introduciéndose el carácter tabulador `\t` para separar los distintos atributos. Podemos emplear varios formatos para codificar esta información. Para elementos genómicos delimitados por pares de coordenadas es necesario utilizar el formato BED o el estándar GFF. Para mostrar el comportamiento de una determinada característica biológica a lo largo del genoma debemos especificarlo en el formato WIG.

#### Versión del genoma

Recomendamos precaución con la versión del genoma. Nuestras propias anotaciones deben obtenerse con esa misma secuencia de referencia.

El formato BED (del inglés, *browser extensible format*, formato extensible del navegador) es útil para plasmar en pantalla aquellas anotaciones constituidas por elementos genómicos definidos por un rango de coordenadas. Dentro de este grupo de características biológicas están clasificados los transcritos, los exones o los sitios de unión a factores de transcripción. En su formato más simple, podemos codificar anotaciones proporcionando su localización en el genoma. Para añadir atributos adicionales (por ejemplo, el nombre, el color o un valor asociado a su fiabilidad) debemos extender la línea actual con esos valores. Es posible introducir un rango de coordenadas adicional para definir distintos grosores en la anotación gráfica (por ejemplo, regiones codificantes y no traducibles de los genes). En la figura 31, se puede observar la sintaxis completa de este formato:

#### Nota

El usuario puede guardar más de una pista dentro de un mismo fichero. Es posible también codificar los datos en diferentes formatos.

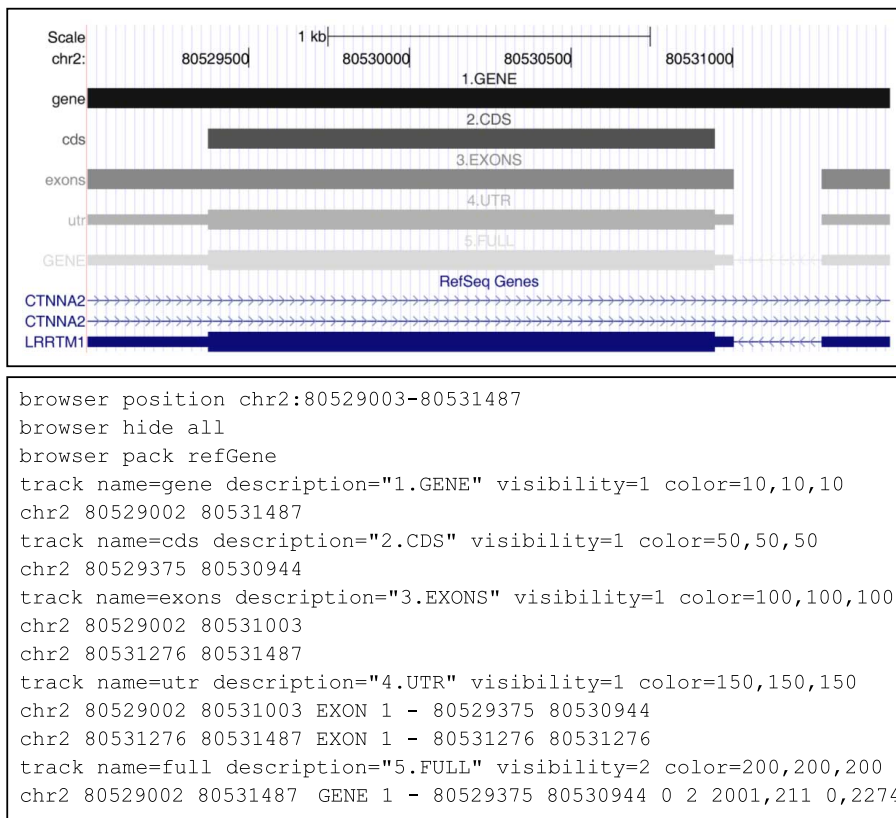
Figura 31. Definición básica y extensión del formato BED.

```
cromosoma inicio final
cromosoma inicio final nombre score hebra inicio2 final2
```

Hemos ilustrado con un ejemplo real el comportamiento de las anotaciones codificadas en formato BED (ver figura 32). Tomando como referencia la información que conocemos sobre la ubicación del gen *LRRTM1*, analizado a lo largo de este capítulo, hemos generado nuestras propias pistas para resaltar distintas regiones en su interior. A partir de las coordenadas para el propio transcrito, sus exones y la región codificante, podemos generar tres pistas di-

ferentes (pistas 1, 2 y 3 en la figura 32), resaltando un área diferente del gen en cada caso. Para denotar con un trazo de distinto grosor la frontera entre la región codificante y la región no traducible, debemos extender el formato BED para añadir una segunda pareja de coordenadas (inicio del CDS, pista 4 figura 32). Finalmente, para emular la propia representación de la pista de RefSeq (con los intrones marcados con una línea recta entre exones y la orientación del gen indicada mediante flechas), enriquecemos el formato añadiendo más información a continuación sobre el número de exones y la longitud de éstos (pista 5, figura 32).

Figura 32. Anotaciones en distintas versiones del formato BED.



El formato WIG (del inglés *wiggle*, agitado) es útil para integrar en el navegador genómico información sobre elementos cuya ubicación más probable viene definida por una distribución continua de valores que varía en intensidad a lo largo de distintas regiones del genoma. Los resultados de experimentos de secuenciación masiva referentes a modificaciones de histonas y reconocimiento de factores de transcripción son claros ejemplos de esta familia de anotaciones. Para codificar estos elementos, el bioinformático debe asignar un valor numérico a cada posición de una determinada región genómica. Esta cantidad representa la intensidad de la señal biológica anotada en esta nueva pista. Cuando la distancia física entre dos anotaciones consecutivas de nuestra pista es variable se debe utilizar la palabra clave *variableStep*. En cambio, cuando los propios datos están ubicados a intervalos regulares de distancia podemos

### Lecturas complementarias

S. L. Berger (2007). "The complex language of chromatin regulation during transcription". *Nature* (núm. 447, págs. 407-412).

A. Barski y otros (2007). "High-resolution profiling of histone methylations in the human genome". *Cell* (núm. 129, págs. 823-837).

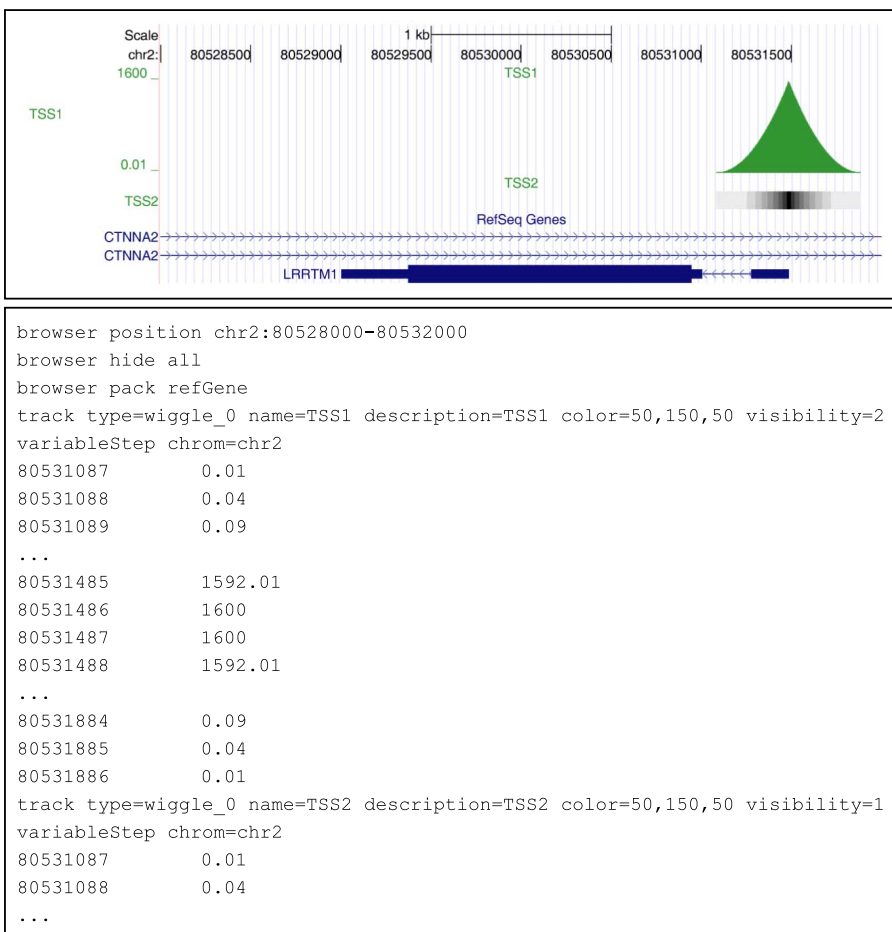
introducir la instrucción *fixedStep*. Cuando un cierto rango de posiciones contiguas presenta el mismo valor, es posible resumir varias anotaciones en una única línea mediante la palabra clave *span*.

Figura 33. Definiciones del formato WIG.

```
variableStep chrom=cromosoma
coordenada    valor
coordenada    valor
coordenada    valor
...
variableStep chrom=cromosoma span=largo
coordenada    valor
...
fixedStep chrom=cromosoma start=inicio span=largo step=paso
coordenada    valor
...
```

En la figura 34 procedemos a cargar una pista en formato WIG para anotar una hipotética función de probabilidad asociada a la localización del posible inicio de transcripción de nuestro gen *LRRTM1* (mostramos las mismas anotaciones utilizando dos modos distintos de visualización).

Figura 34. Anotaciones en formato WIG.



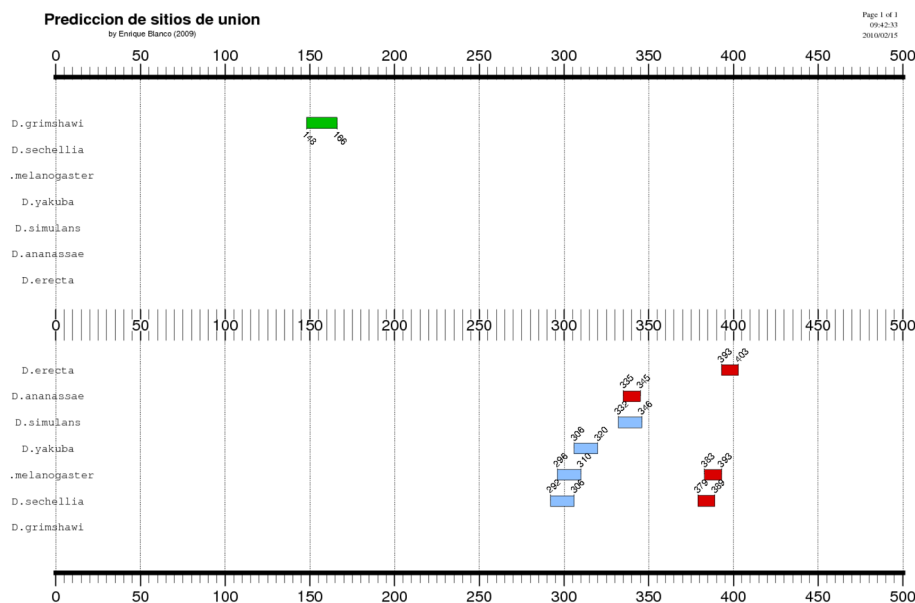
El estándar GFF (del inglés *general feature format*, formato de características generales) fue concebido para permitir la transmisión eficiente de anotaciones génicas entre diferentes aplicaciones bioinformáticas. Permite generar representaciones gráficas similares al formato BED, dedicando cada línea a un elemento genómico distinto. Posee la particularidad de que podemos asociar un identificador propio individualmente a cada anotación. Si utilizamos un mismo nombre para un conjunto de anotaciones relacionadas, el servidor interpretará que éstas pertenecen al mismo grupo (por ejemplo, los exones de un gen). Es posible asignar una puntuación a cada registro, lo que se refleja gráficamente en cajas de diferentes grosores o con una gama de colores de distintas intensidades. El formato GFF presenta la siguiente sintaxis básica:

Figura 35. Definición del formato GFF.

```
# comentario
secuencia origen caracteristica inicio final valor hebra pauta grupo
secuencia origen caracteristica inicio final valor hebra pauta grupo
secuencia origen caracteristica inicio final valor hebra pauta grupo
...
```

Este formato está muy extendido entre los desarrolladores de aplicaciones bioinformáticas. Existen, por ejemplo, numerosos programas que procesan anotaciones en formato GFF para construir representaciones gráficas de calidad (como los mapas de anotaciones producidas en publicaciones que anuncian la secuenciación de los genomas más populares). Por ejemplo, gracias al programa GFF2PS, las predicciones pueden volcarse en PostScript, fácilmente convertible a otros formatos más populares. El eje central delimita si las predicciones se encuentran en la hebra positiva o negativa de la molécula de ADN:

Figura 36. Representación gráfica de predicciones con GFF2PS.



### Página web

Se puede encontrar más información sobre el estándar GFF en la página web del Wellcome Trust Sanger Institute.

### Lectura complementaria

J. F. Abril; R. Guigó (2000). "gff2ps: visualizing genomic annotations". *Bioinformatics* (núm. 8, págs. 743-744).

### Leyenda figura 36

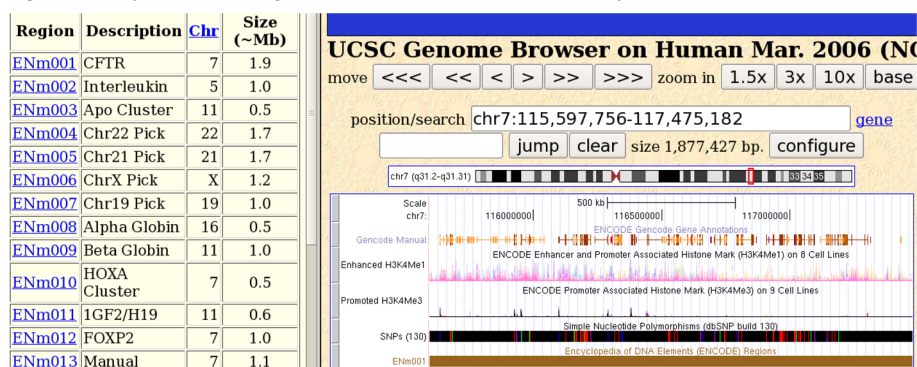
Para personalizar la representación gráfica, el programa GFF2PS permite introducir un fichero de configuración de colores. Ver el manual de usuario del programa para más detalles.

### 3.4. ENCODE

La obtención de la secuencia del genoma humano proporcionó a la comunidad científica mundial el acceso inmediato a una enorme cantidad de información inédita hasta ese momento. Este hito histórico reveló también la deficiente calidad de las anotaciones existentes. Un consorcio internacional de investigadores fundó entonces el proyecto ENCODE (del inglés *Encyclopedia of DNA Elements*, enciclopedia de elementos del ADN), con la misión de producir para varias clases de tejidos una anotación mucho más precisa de todos los componentes codificados en el interior del genoma humano (por ejemplo, genes, transcritos, microARN, elementos regulatorios, etc.). Para lograr este objetivo, el mismo consorcio centró una parte sustancial de sus esfuerzos en el desarrollo de nuevas técnicas experimentales y computacionales que permitieron progresar espectacularmente en múltiples áreas de conocimiento. El navegador genómico UCSC fue escogido como herramienta para servir estos datos de forma pública. Las pistas de datos ENCODE se depositaban en este servidor inmediatamente después de llevar a cabo su anotación. Estos datos, no obstante, poseían un periodo de cuarentena (nueve meses desde su liberación) que impedía utilizarlas antes en otra publicación, para proteger a sus autores originales.

En una primera fase del proyecto, entre los años 2003 y 2007, el consorcio ENCODE dirigió sus esfuerzos sobre el uno por ciento del genoma humano (la fase piloto del proyecto). Para acceder a las anotaciones de esta primera fase, se debe seleccionar el enlace ENCODE pilot project en la página principal de UCSC. Una vez allí, se puede escoger cualquiera de las 44 regiones habilitadas en esta fase del proyecto y abrir el visor genómico convencional (ver figura 37).

Figura 37. Explorando las regiones humanas anotadas en la fase piloto de ENCODE.



Para acceder a través de UCSC a la anotación completa del genoma humano que está llevando a cabo este consorcio público, debemos seleccionar el enlace ENCODE desde la página web inicial de este portal. Este proyecto está contribuyendo sustancialmente a incrementar nuestro conocimiento respecto a los patrones de expresión génica y la regulación de la transcripción en diferentes

#### Lecturas complementarias

**The ENCODE Project Consortium (2007).** "Identification and analysis of functional elements in 1% of the human genome by the ENCODE pilot project". *Nature* (núm. 447, págs. 799-816).

**K. R. Rosenbloom y otros (2010).** "ENCODE whole-genome data in the UCSC Genome Browser". *Nucleic Acids Research* (núm. 38, págs. D620-D625).

#### Lectura complementaria

**J. Harrow y otros (2006).** "GENCODE: producing a reference annotation for ENCODE". *Genome Biology* (supl. 1 núm. S4, págs. 1-9).

tejidos. En particular, el subproyecto GENCODE produjo un catálogo de genes de referencia mediante la revisión manual de anotaciones automáticas generadas por varios grupos de investigación.

Para gestionar todo este volumen de información, los diseñadores de UCSC han creado recientemente un nuevo formato de pista denominado superpista (*supertrack*, en inglés) que aglutina varias pistas de datos. La secuenciación masiva en varias líneas celulares de aproximadamente una docena de modificaciones postranscripcionales de histonas puede visualizarse de forma óptima empleando este tipo de pistas. Cada marca posee su propia distribución a lo largo del genoma, presentando un patrón de picos asociable a distintos elementos genómicos. Si, como es esperable, crece el número de modificaciones estudiadas con este procedimiento, parece evidente que la disposición horizontal clásica de pistas presentadas en el visor genómico de UCSC podría resultar impracticable. Para optimizar el espacio útil en pantalla y facilitar la interpretación de la información, UCSC proporciona estos datos integrados en una superpista, fusionando los patrones observados en todos los tejidos para una única marca (ver figura 38):

Figura 38. Configuración de la superpista H3K4Me1 en ENCODE.

**Enhanced H3K4Me1 Track Settings**

**ENCODE Enhancer and Promoter Associated Histone Mark (H3K4Me1) on 8 Cell Lines (ENCODE Regulation)**

Display mode:   [Reset to defaults](#)

Overlay method:

Type of graph:

Track height:  pixels (range: 11 to 100)

Vertical viewing range: min:  max:  (range: 0 to 10000)

Data view scaling:  Always include zero:

Transform function: Transform data points by:

Windowing function:  Smoothing window:  pixels

Draw y indicator lines: at y = 0.0:  at y =

[Graph configuration help](#)

All subtracks:

List subtracks: ☒ only selected/visible ☐ all (8 of 8 selected) [Restricted Until](#)

Subtrack	Description	Schema	Restricted Until
<input checked="" type="checkbox"/> Gm12878	Enhancer and Promoter Associated Histone Mark (H3K4Me1) on Gm12878 cells from ENCODE	<a href="#">schema</a>	2009-10-05
<input checked="" type="checkbox"/> H1ES	Enhancer and Promoter Associated Histone Mark (H3K4Me1) on H1 ES cells from ENCODE	<a href="#">schema</a>	2010-06-30
<input checked="" type="checkbox"/> HMEC	Enhancer and Promoter Associated Histone Mark (H3K4Me1) on HMEC cells from ENCODE	<a href="#">schema</a>	2010-06-28
<input checked="" type="checkbox"/> HSMM	Enhancer and Promoter Associated Histone Mark (H3K4Me1) on HSMM cells from ENCODE	<a href="#">schema</a>	2010-09-16
<input checked="" type="checkbox"/> HUVEC	Enhancer and Promoter Associated Histone Mark (H3K4Me1) on HUVEC cells from ENCODE	<a href="#">schema</a>	2009-10-05
<input checked="" type="checkbox"/> K562	Enhancer and Promoter Associated Histone Mark (H3K4Me1) on K562 cells from ENCODE	<a href="#">schema</a>	2009-10-05
<input checked="" type="checkbox"/> NHEK	Enhancer and Promoter Associated Histone Mark (H3K4Me1) on NHEK cells from ENCODE	<a href="#">schema</a>	2009-10-06
<input checked="" type="checkbox"/> NHLF	Enhancer and Promoter Associated Histone Mark (H3K4Me1) on NHLF cells from ENCODE	<a href="#">schema</a>	2010-06-28

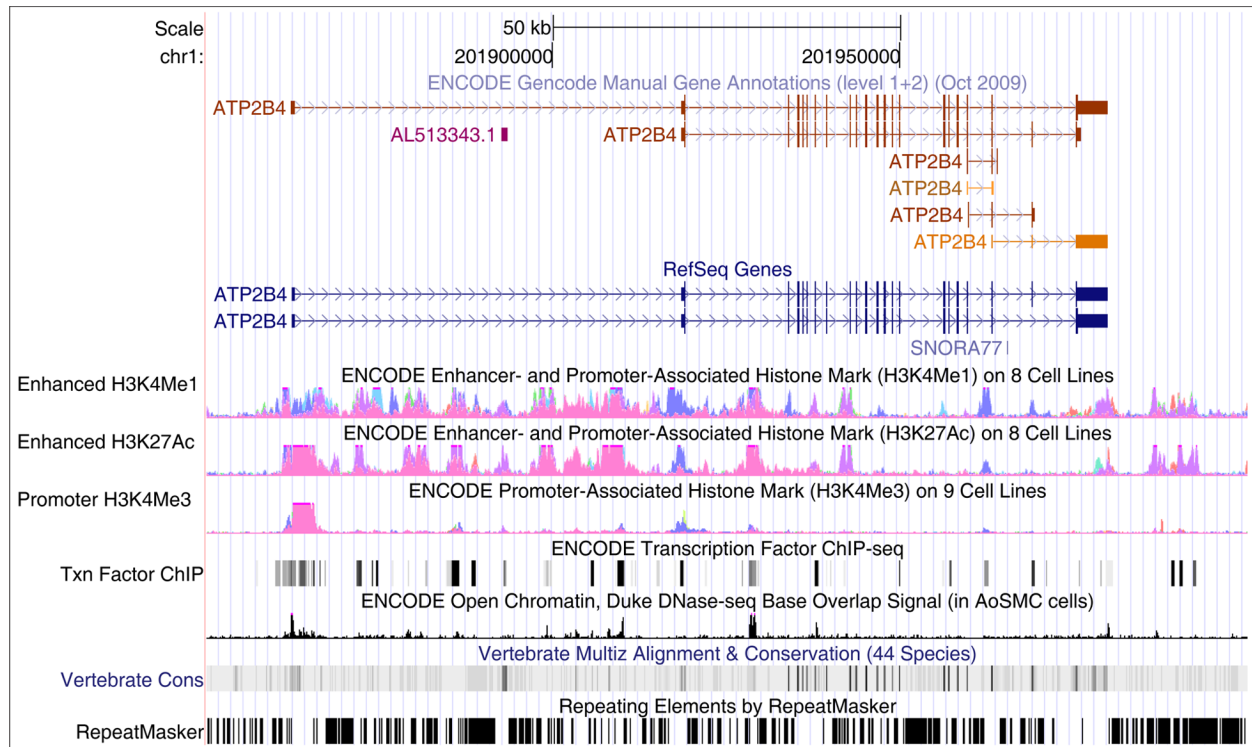
8 of 8 selected

#### Marcas de histonas

La modificación H3K4Me3 marca los inicios de transcripción, mientras que H3K4Me1 marca regiones reguladoras y H3K27Me3, las regiones de represión génica.

Estas anotaciones más sofisticadas se integran en una superpista denominada ENCODE Integrated Regulation Track (pista de regulación de ENCODE). En esta superpista podemos manipular simultáneamente datos de secuenciación sobre histonas, factores de transcripción y transcritos, constituyendo una fuente de información inagotable para las investigaciones genómicas durante las próximas décadas. En la figura 39 puede observarse la precisión de estas anotaciones para detectar un amplio conjunto de formas alternativas de un mismo gen previamente desconocidas.

Figura 39. Región del genoma humano anotada por el consorcio ENCODE.



La marca ENCODE, como signo de excelencia científica, está en periodo de expansión hacia otros genomas, con la misión de proporcionar conjuntos similares de anotaciones para otros organismos. El proyecto modENCODE, en particular, representa una revolución similar dentro de la investigación de dos organismos modelo: la mosca de la fruta *Drosophila melanogaster* y el gusano *Caenorhabditis elegans*. Las anotaciones de este nuevo consorcio internacional están ya a disposición de ambas comunidades científicas gracias al portal dedicado a este proyecto (ver figura 41).

### 3.5. Instalación local

UCSC sirve anotaciones para más de 50 genomas de distintas especies. Para ejecutar protocolos de análisis bioinformático a gran escala que involucran miles de genes, resulta más conveniente obtener una copia local de la distribución de los genomas estudiados. El enlace Download genomes (descargar genomas) en la página principal permite recuperar los ficheros de texto interpretados por la interfaz gráfica de UCSC (secuencias y anotaciones):

#### Lecturas complementarias

S. E. Celniker y otros (2009). Unlocking the secrets of the genome. *Nature* (núm. 459, págs. 927-930).

The modENCODE Consortium (2010). "Identification of functional elements and regulatory circuits by *Drosophila* modENCODE". *Science* (núm. 330, págs. 1787-1797).

#### Ved también

Repasar estos conceptos en los materiales de la asignatura *Fundamentos de informática en entornos bioinformáticos*.



- El conjunto completo de datos de un cromosoma en particular (*data set by chromosome*).
- Las anotaciones de cada pista en ficheros individuales (*annotation database*).
- Los ficheros de conversión entre distintas versiones (*liftOver files*).
- Las comparaciones con otros genomas (*pairwise alignments/multiple alignments*).

Cada enlace nos abre el acceso a un nuevo directorio para escoger los ficheros deseados.

Figura 42. Descarga de anotaciones disponibles para el genoma humano.

Human Genome Feb. 2009 (hg19, GRCh37)

- \* [Full data set](#)
- \* [Data set by chromosome](#)
- \* [Annotation database](#)
- \* [GC percent data](#)
- \* [Protein database for hg19](#)
- \* [SNP131-masked FASTA files](#)
- \* [LiftOver files](#)
- \* [Pairwise Alignments](#)
- \* [Multiple Alignments](#)

El conjunto de secuencias de nucleótidos de los cromosomas en cada genoma está empaquetado en un único archivo comprimido. En la figura 43 podemos ver el contenido una vez descomprimido de la distribución hg19 del genoma humano, centrando nuestro interés en acceder a la secuencia del cromosoma 17. Para ello, hemos procedido a descargar y desempaquetar el fichero `chromFa.tar.gz`, visualizando después el fichero `chr17.fa` presentado en formato FASTA.

#### **Secuencias enmascaradas**

También están disponibles las secuencias enmascaradas de cada cromosoma, filtrando las regiones de baja complejidad.

Figura 43. Descarga de la secuencia de los cromosomas humanos.

```

chromFa.tar.gz          20-Mar-2009 09:21 905M
chromFaMasked.tar.gz    20-Mar-2009 09:30 477M
est.fa.gz               14-Oct-2010 01:07 1.4G
refMrna.fa.gz           14-Oct-2010 01:07 36M
upstream1000.fa.gz      14-Oct-2010 01:07 7.2M
xenoMrna.fa.gz          14-Oct-2010 01:08 910M

chromFa/
chr10.fa  chr14.fa  chr18.fa  chr21.fa  chr4.fa  chr8.fa
chr11.fa  chr15.fa  chr19.fa  chr22.fa  chr5.fa  chr9.fa
chr12.fa  chr16.fa  chr1.fa  chr2.fa  chr6.fa  chrX.fa
chr13.fa  chr17.fa  chr20.fa  chr3.fa  chr7.fa  chrY.fa

>chr17
AAGCTTCTCACCCCTGTTCTCCTGCATAGATAATTGCATGACAATTGCCTTGT
CCCTGCTGAATGTGCTCTGGGGTCTCTGGGGTCTCACCCACGACCAACTC
CCTGGGCCTGGCACCAGGGAGCTTAACAAACATCTGTCCAGCGAATACCT
GCATCCCTAGAAGTGAAGCCACCGCCCAAAGACACGCCCATGTCCAGCTT
AACCTGCATCCCTAGAAGTGAAGGCACCGCCCAAAGACACGCCCATGTCC
AGCTTATTCTGCCCAGTTCTCTCCAGAAAGGCTGCATGGTTGACACACA
GTGcctgcgacaaagctgaatgctatcatcttaaaaaactccttgctggttt
gagaggcagaaaatgatatactcatagttgctttactttgcatattttAAA
ATTGTGACTTTTCATGGCATAAATAATACTGGTTTATTACAGAAGCACTAG
...

```

La carpeta Annotation database contiene todas las pistas utilizadas por la interfaz gráfica del navegador UCSC (mostramos una pequeña parte de estas en la figura 44). En particular, el fichero `refGene.txt.gz`, que contiene las anotaciones de cada gen humano según el consorcio RefSeq, es procesado por UCSC para generar la pista gráfica RefSeq Genes (ejemplo incluido en la misma imagen):

#### El fichero `refGene.txt`

Cada línea de éste archivo contiene la anotación suministrada por RefSeq para un transcrito concreto. Las formas alternativas del mismo gen están codificadas en líneas distintas.

Figura 44. La base de datos de anotaciones de UCSC.

```

all_bacends.sql          24-May-2009 11:40 2.1K
all_bacends.txt.gz       24-May-2009 11:40 93M
all_est.sql              12-Sep-2010 16:33 2.3K
all_est.txt.gz           12-Sep-2010 16:34 374M
...
refGene.sql              10-Oct-2010 20:40 2.0K
refGene.txt.gz           10-Oct-2010 20:40 3.7M
refLink.sql              10-Oct-2010 20:43 1.8K
refLink.txt.gz           10-Oct-2010 20:43 8.3M
...
xenoRefGene.sql          10-Oct-2010 21:05 2.0K
xenoRefGene.txt.gz       10-Oct-2010 21:05 13M
xenoRefSeqAli.sql        10-Oct-2010 21:14 2.3K
xenoRefSeqAli.txt.gz     10-Oct-2010 21:14 13M

1199 NM_178839 chr2 - 80529002 80531487 80529375 80530944 2
80529002,80531276, 80531003,80531487, 0 LRRTM1 cml1 cml1 0,-1,

```

Esta es la descripción de los atributos pertenecientes a cada registro de este fichero (tomando como ejemplo la línea de la anotación correspondiente al gen *LRRTM1*):

Tabla 7. Atributos del fichero *refGene.txt*.

Campo	Ejemplo	Descripción
1	1199	Identificador interno
2	NM_178839	Código asignado por RefSeq
3	chr2	Cromosoma
4	-	Hebra
5	80529002	Inicio del transcrito
6	80531487	Fin del transcrito
7	80529375	Inicio de la región codificante
8	80530944	Fin de de la región codificante
9	2	Número total de exones
10	80529002,80531276,	Coordenadas iniciales de los exones
11	80531003,80531487,	Coordenadas finales de los exones
12	0	Puntuación
13	<i>LRRTM1</i>	Nombre común abreviado del gen
14	cmpl	Anotación del inicio del CDS (completa)
15	cmpl	Anotación del fin del CDS (completa)
16	0,-1,	Pauta de lectura de los exones

Recomendamos precaución para interpretar las anotaciones de genes codificados en la hebra negativa de la molécula de ADN. En estos casos, las coordenadas iniciales y finales de las anotaciones deben intercambiarse para tener en cuenta este concepto (el campo número cinco contendría el final del gen y el campo número seis el inicio). Una conversión similar debería realizarse con las coordenadas de cada exón individualmente.

El terminal de la plataforma LINUX es el método más eficiente para analizar distintas pistas de anotaciones. Sin embargo, para aquellos usuarios que no poseen un dominio suficiente para trabajar con soltura en este tipo de entornos, el propio portal de UCSC permite gestionar las consultas enmascarando la complejidad de un gestor de base de datos. La aplicación Table browser (navegador de tablas), accesible desde el enlace situado dentro de la barra de menú superior de color azul, ofrece una interfaz web para interrogar varias pistas simultáneamente. Podemos definir la región genómica donde dirigir nuestras demandas o trabajar con la totalidad de las anotaciones de las pistas. Con esta interfaz de comandos podemos instruir al navegador para recuperar datos a

#### Ved también

Estos conceptos están introducidos en la asignatura *Fundamentos de informática en entornos bioinformáticos*.

gran escala. Por ejemplo, en el marco de la investigación sobre diferentes catálogos de genes, es posible extraer sistemáticamente la región codificante de todos los genes anotados en cualquier especie para llevar a cabo un análisis bioinformático posterior.

Figura 45. El navegador de tablas de UCSC.

Table Browser

Use this program to retrieve the data associated with a track in text format, to calculate intersections between tr DNA sequence covered by a track. For help in using this application see [Using the Table Browser](#) for a descriptor form, the [User's Guide](#) for general information and sample queries, and the [OpenHelix Table Browser tutorial](#) for of the software features and usage. For more complex queries, you may want to use [Galaxy](#) or our [public MySQL s](#) biological function of your set through annotation enrichments, send the data to [GREAT](#). Refer to the [Credits](#) pag contributors and usage restrictions associated with these data.

clade: Mammal genome: Human assembly: Feb. 2009 (GRCh37/hg19)

group: Genes and Gene Prediction Tracks track: RefSeq Genes add custom tracks

table: refGene describe table schema

region: genome position chr2:80530895-80530944 lookup define regions

identifiers (names/accessions): paste list upload list

filter: create

intersection: create

correlation: calculate clear (with: UCSC Genes)

output format: all fields from selected table Send output to Galaxy GREAT

output file: (leave blank to keep output in browser)

file type returned: plain text gzip compressed

get output summary/statistics

To reset all user cart settings (including custom tracks), [click here](#).

Esta aplicación permite ejecutar fácilmente múltiples acciones que involucren varias pistas. Por ejemplo, podemos utilizar esta herramienta para calcular la correlación entre las anotaciones servidas por el proyecto RefSeq y la pista UCSC Genes proporcionada independientemente por el navegador UCSC (ver figura 46). Este tipo de análisis permite establecer el nivel de solapamiento entre pistas de datos similares, evaluando globalmente la precisión de varios sistemas de anotaciones. Con esta aproximación podemos también calcular el grado de similitud entre dos réplicas del mismo experimento de secuenciación.

Figura 46. Cálculo de correlaciones con el navegador de tablas.

Correlate table 'RefSeq Genes' (refGene) with table 'knownGene'

Select a group, track and table to correlate with:

group: Genes and Gene Prediction Tracks track: UCSC Genes

table: knownGene

Limit total data points in result: 40,000,000 Window data to: 1 bases

calculate clear selections return to table browser

position: full genome, 49 chromosomes

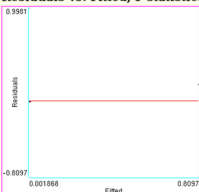
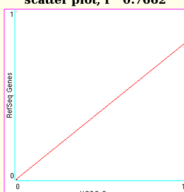
intersecting both tables 'RefSeq Genes' and 'UCSC Genes' with 'knownGene':

warning: reached maximum data points: 40,000,000 before end of data.

Position and # of data points in intersection	Correlation coefficient r	r <sup>2</sup>	Track	Minimum	Maximum	Mean	Variance	Standard deviation	Regression line y = m*x + b
chr1:1-40,000,000 40,000,000 data points	0.8753	0.7662	RefSeq Genes	0	1	0.0428	0.04097	0.2024	0.8078
			UCSC Genes	0	1	0.05067	0.0481	0.2193	0.001868

scatter plot, r<sup>2</sup> 0.7662

Residuals vs. Fitted, F statistic: 1.311e+08



## 4. El navegador genómico ENSEMBL

ENSEMBL es otro navegador ampliamente extendido que está gestionado por el Instituto Europeo de Bioinformática (EBI) en Hinxton (Reino Unido). ENSEMBL también sirve las anotaciones en forma de pistas que pueden mostrarse u ocultarse dentro del visor gráfico. El usuario puede cargar sus propias pistas en ENSEMBL para integrarlas junto con la información convencional. La plataforma ENSEMBL produce también desde sus orígenes su propia anotación génica de referencia, colaborando intensamente con varios consorcios para la curación manual de estos datos. Como se aprecia en la figura 47, para acceder a la anotación de un gen el usuario debe seleccionar primero la distribución apropiada del genoma e introducir posteriormente su nombre en la caja de búsqueda.

### Lectura complementaria

T. Hubbard y otros (2002). "The Ensembl genome database project". *Nucleic Acids Research* (núm. 30, págs. 38-41).

Figura 47. Accediendo a las anotaciones en ENSEMBL.



**e!Ensembl** | Login | Register | BLAST/BLAT | BioMart | Tools | Downloads | Help  
Home

Search: **All species** for  
 **Go**

e.g. **human gene BRCA2** or **rat X:100000..200000** or **coronary heart disease**

### Browse a Genome

The Ensembl project produces genome databases for vertebrates and other eukaryotic species, and makes this information freely available online.  
Click on a link below to go to the species' home page.

**Popular genomes** ([Log in to customize this list](#))

-  **Human**  
GRCh37
-  **Mouse**  
NCBIM37
-  **Zebrafish**  
Zv8

**All genomes**

-- Select a species --

[View full list of all Ensembl species](#)

Other species are available in [Ensembl Pre!](#) and [EnsemblGenomes](#)

Human (GRCh37)

About this species

Description

Genome Statistics

Assembly and Genebuild

Top 40 InterPro hits

Top 500 InterPro hits

What's New

Sample entry points

Karyotype

Location (6:133017695-133017695-133017695)

Gene (BRCA2)

Transcript (FOXP2-203)

Variation (rs1333049)

Regulation (ENSR00000133049)

Configure this page

Manage your data

Export data

Bookmark this page

Search Ensembl Human

Search for:

LRRTM1

Go

e.g. gene BRCA2 or 6:133017695-133161157 or osteoarthritis

Description Assembly and Genebuild »

Human (*Homo sapiens*)

Assembly

This site provides a data set based on the February 2009 *Homo sapiens* high coverage assembly from the [Genome Reference Consortium](#). The data set consists of gene models built from the genewise alignments of the human proteome as well as from alignments of human cDNAs using the cDNA2genome model of exonerate.

This release of the assembly has the following properties:

- 27478 contigs.
- contig length total 3.2 Gb.
- chromosome length total 3.1 Gb.

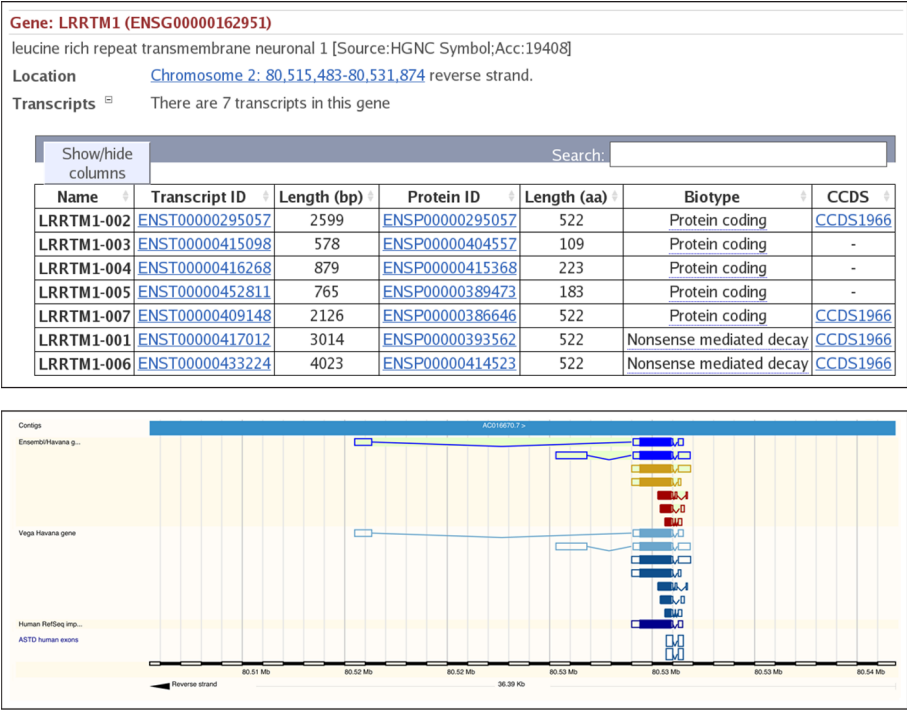
ENSEMBL fue pionero en presentar la información sobre los genes y los transcritos de forma estructurada. Los genes, como unidades que agrupan múltiples ARN mensajeros, reciben un identificador con las iniciales ENSG (ENSEMBL gene), mientras que los diferentes productos de este, transcritos y proteínas, poseen sus propios códigos ENST y ENSP, respectivamente. Es posible acceder a toda la información recopilada sobre cada gen o visualizar la secuencia de éste en pantalla (ver figura 48). Para modificar el inventario de pistas o realizar cambios en su volcado por pantalla, el usuario debe presionar el botón *Configure this page* (configurar esta página). A diferencia del servidor genómico UCSC, que genera imágenes estáticas de las pistas que permiten posteriormente acceder a más información, ENSEMBL reconstruye dinámicamente el visor genómico en el momento de la actualización. A medida que desplazamos el puntero de nuestro ratón sobre ciertas secciones de la imagen, aparece un resumen de la información asociada a cada pista (incluyendo enlaces hacia otras bases de datos). ENSEMBL genera igualmente imágenes para impresión en formato de alta calidad.

Los usuarios de ENSEMBL tienen a su disposición numerosos tutoriales, páginas de ayuda e incluso películas con ejemplos ilustrativos (con su propio portal de acceso en YouTube). Para cada opción también existe específicamente un manual de asistencia propio.

#### Anotaciones comunes

Tanto UCSC como ENSEMBL sirven sus propias anotaciones junto con un conjunto de pistas comunes proporcionadas por otros recursos. La anotación génica de cada servidor puede encontrarse como una pista más de datos en el otro programa.

Figura 48. Anotación del gen *LRRTM1* en ENSEMBL.



BIOMART es una de las utilidades más potentes de ENSEMBL. Esta aplicación combina registros de numerosas fuentes con el objetivo de extraer nuevo conocimiento sobre los genes y otras pistas de anotaciones almacenadas en su base de datos. Como se aprecia en la figura 49, la pantalla de BIOMART está dividida principalmente en dos áreas claramente diferenciadas: el menú (a la izquierda) y la zona de trabajo (ocupando la mayoría de la pantalla). El usuario debe en primer lugar establecer un genoma de referencia para trabajar sobre éste. La unidad de trabajo en BIOMART es el conjunto de genes de dicho genoma (denominado Dataset en el menú, figura 49). Inicialmente, podremos realizar operaciones sobre todos los transcritos anotados por ENSEMBL en el genoma seleccionado (51737 ARN mensajeros en la distribución GRCh37 del genoma humano).

Nota

Dado que las acciones que podemos hacer empleando BIOMART o el navegador de tablas de UCSC son ligeramente diferentes, recomendamos que el estudiante realice prácticas intensivas con ambos.

Sobre el catálogo de genes de trabajo podemos realizar tres tipos de acciones:

- **Filtros.** Reducir el conjunto de genes de trabajo, focalizando nuestro interés sobre aquellos que posean determinadas características. Podemos definir filtros sobre cromosomas, regiones genómicas, homología entre especies o funciones biológicas. También es posible introducir nuestra propia lista de genes (por ejemplo, los resultados de un experimento de expresión).
- **Atributos.** Seleccionar el tipo de información que nos gustaría obtener del conjunto actual de genes. Podemos establecer los atributos generales de cada gen o definir la información sobre homología con otras especies que deseamos recuperar. Existe la opción de extraer sistemáticamente distintas partes de la estructura exónica de los genes.
- **Dataset.** Combinar las anotaciones actuales con otro genoma distribuido por ENSEMBL para reducir el conjunto inicial de resultados.

Para actualizar la cifra total de genes de trabajo, debemos presionar el botón *Count* tras definir un nuevo filtro sobre los datos. Cuando hemos finalizado con el filtrado, mantenemos en el área de trabajo el grupo de genes que cumplen con todos los criterios especificados. En ese instante, justo antes de generar los resultados, es necesario realizar la selección de los atributos. Una vez terminados los procesos de filtrado y formateado de los resultados, debemos utilizar el botón *Results* para recibir el archivo final. BIOMART exhibe siempre una vista preliminar más reducida de la información solicitada. De este modo, el usuario puede introducir cambios en las preguntas realizadas sobre el catálogo de genes sin necesidad de descargar completamente en cada ocasión el fichero final de resultados (figura 49).



## Resumen

En este capítulo hemos introducido los conceptos fundamentales sobre la navegación genómica. Trabajando con el servidor genómico UCSC hemos aprendido a configurar el programa, manipular las pistas de anotaciones e integrar nuestros propios datos para ser visualizados dentro del mismo entorno gráfico. También nos hemos familiarizado con los diferentes ficheros que constituyen la base de cualquier distribución genómica. Finalmente, como alternativa a los sistemas de análisis basados en la plataforma LINUX, hemos mostrado cómo manipular estas informaciones fácilmente para extraer nuevo conocimiento, utilizando herramientas web como el navegador de tablas de UCSC o BIOMART.



## Actividades

1. Buscad en la Red la revista científica *Nucleic Acids Research*. En su interior, explorad el contenido las últimas ediciones especiales sobre bases de datos y servidores web.
2. Seleccionad un gen humano cualquiera con el navegador UCSC. Después, abrid el navegador ENSEMBL para visualizar el mismo gen. A continuación, explorad en paralelo las anotaciones disponibles para este gen en ambos servidores.
3. La herramienta Liftover de UCSC es extremadamente útil para los bioinformáticos. Explorad qué operaciones realiza, cómo funciona y probad su eficacia sobre la anotación concreta de un gen humano cuyas coordenadas hayan variado entre dos versiones de este genoma.
4. Aprended a crear una pista propia en formato BED para mostrar una caja artificial de 1.000 nucleótidos de longitud en color rojo, concretamente en la posición 1000000 del cromosoma 21 del genoma humano (versión hg19). Ahora cread una pista WIG en las mismas coordenadas para representar un elemento artificial que posea un valor de diez unidades durante los primeros 300 nucleótidos, un valor de veinte unidades en los siguientes 400 nucleótidos y un valor de diez unidades nuevamente en los últimos 300 nucleótidos.
5. Activad las pistas de conservación entre vertebrados, elegid un gen humano constituido por cinco o más exones y estudiad el patrón de conservación filogenética en las diferentes regiones de este gen (exones, intrones, regiones codificantes, regiones no traducibles).
6. Dentro de la versión de UCSC para servir datos del proyecto ENCODE, activad una superpista que contenga múltiples modificaciones de histonas y visitad este paisaje genómico a lo largo de un cromosoma humano concreto.
7. Inspeccionad el navegador de modENCODE e intentad reproducir la figura 41.
8. Con UCSC, realizad la descarga de la última distribución del genoma humano y cread un DVD que contenga la secuencia de sus cromosomas y las anotaciones más relevantes.
9. En el marco de la última distribución del genoma de la mosca de la fruta, abrid el navegador de tablas de UCSC y calculad la correlación existente entre las anotaciones servidas por los consorcios FlyBase y RefSeq. Visualizad casos concretos de anotaciones divergentes y valorad los resultados obtenidos.
10. Utilizad el navegador de tablas de UCSC para extraer las regiones promotoras de todos los genes humanos (longitud: 500 bps). Emplead el programa BLAT sobre un gen en cada hebra de ADN para validar que la descarga haya funcionado correctamente.
11. Utilizad la herramienta BIOMART de ENSEMBL para recuperar los identificadores de aquellos genes cuya forma ortóloga entre humano y ratón está documentada.

## Ejercicios de autoevaluación

1. Enumerad qué dos tipos de información definen cada distribución de un genoma.
2. Definid cuál es la principal misión de un servidor genómico.
3. Enumerad qué informaciones pueden emplearse para localizar una región del genoma.
4. ¿Qué es una pista de datos en el contexto de un navegador genómico?
5. Describid el procedimiento habitual para acceder a una región genómica con un servidor de genomas como UCSC o ENSEMBL.
6. Enumerad las opciones básicas para configurar la visualización de un navegador.
7. ¿Qué opciones posee habitualmente el usuario para mostrar una pista?
8. Enumerad cinco bloques de opciones en los que se agrupan las pistas de datos en el navegador genómico UCSC.
9. Definid qué es el recurso RefSeq.
10. ¿Gráficamente, cómo se representan los exones u otros elementos genómicos anotados sobre una pista en cualquier navegador genómico?

11. Describid el contenido de un fichero FASTA.
12. ¿Cuál es la función del programa BLAT?
13. Enumerad qué informaciones necesitamos obtener para crear nuestra propia pista.
14. Explicad la principal diferencia entre los formatos BED y WIG.
15. Definid el concepto de superpista (existente en el navegador UCSC para ENCODE).
16. Definid las funciones principales del navegador de tablas de UCSC.
17. BIOMART: definid el protocolo de uso de esta herramienta para extraer anotaciones.

## Solucionario

1. La secuencia del genoma y las anotaciones a lo largo de los cromosomas de éste.
2. Un servidor genómico proporciona herramientas para navegar a través de las anotaciones de los genomas.
3. Las coordenadas dentro de un determinado cromosoma, el nombre de un gen, el nombre de una proteína, una secuencia similar.
4. Una pista contiene una serie de anotaciones sobre cierto elemento biológico ubicado en determinadas regiones del genoma.
5. Primero, introducimos una determinada palabra clave que identifique nuestro objeto de búsqueda. Segundo, el servidor nos muestra una lista de posibles coincidencias en varias pistas de datos. Tercero, escogemos una pista para que el servidor construya una representación gráfica de la región seleccionada. Cuarto, utilizamos los botones de desplazamiento y diferentes enfoques para configurar esta vista. Quinto y último, visitamos los detalles de las anotaciones pinchando con el ratón sobre las pistas.
6. Podemos desplazarnos por una región, acercarnos o alejarnos de ésta o invertir la secuencia junto con sus anotaciones en la otra hebra de la molécula de ADN.
7. El conjunto de anotaciones proporcionadas por una pista de datos puede mostrarse compactado en una única línea o desplegarse en varias líneas del visor.
8. Los bloques de opciones más comunes son Mapping and Sequencing Tracks, Phenotype and Disease Associations, Genes and Gene Prediction Tracks, mRNA and EST Tracks, Expression, Regulation, Comparative Genomics, Variation and Repeats.
9. RefSeq es un consorcio público de anotación de productos génicos que fundamentalmente produce anotaciones de calidad evitando información redundante o errónea.
10. Los exones son representados en forma de cajas, mientras que los intrones se indican con una línea que une dos exones contiguos. Otras anotaciones pueden mostrarse también de este modo o exhibiendo una cierta distribución continua de valores a lo largo del genoma.
11. Una secuencia FASTA posee dos tipos de información:
  - La cabecera, que posee el símbolo >, contiene información para clasificar el fichero.
  - La secuencia, a continuación, en forma de líneas de longitud fija (generalmente, 60 caracteres).
12. El programa BLAT busca secuencias prácticamente idénticas a una suministrada por el usuario, generalmente con el objetivo de identificar su ubicación exacta en el genoma.
13. Para crear una *custom track* recomendamos dotar a la pista de instrucciones para el navegador sobre la ubicación de la región anotada, instrucciones sobre la configuración de la visualización de la pista y las coordenadas de nuestras anotaciones en el genoma.
14. El formato BED (*browser extensible data*) es útil para representar anotaciones en el genoma que poseen un inicio y final concretos (por ejemplo, exones). El formato WIG (WIGgle), por contra, es más adecuado para mostrar anotaciones continuas de una cierta característica del genoma.
15. Una superpista realiza una superposición gráfica de varias pistas de datos (generalmente en formato WIG), integrándolas en una única línea del visor genómico.
16. El navegador de tablas permite comparar los datos de diferentes pistas para calcular intersecciones, correlaciones o extraer información adicional a partir de estas.
17. Con BIOMART podemos aplicar sucesivamente varios tipos de operaciones (filtrar, seleccionar atributos y combinar con otros conjuntos de datos) para delimitar claramente la fracción de las anotaciones que deseamos recuperar.

## Bibliografía

**Abril, J. F.; Guigó, R.** (2000). "gff2ps: visualizing genomic annotations". *Bioinformatics* (núm. 8, págs. 743-744).

**Barski, A. y otros** (2007). "High-resolution profiling of histone methylations in the human genome". *Cell* (núm. 129, págs. 823-837).

**Benson, D. A. y otros** (2008). "GenBank". *Nucleic Acids Research* (núm. 36, págs. D25-30).

**Berger, S. L.** (2007). "The complex language of chromatin regulation during transcription". *Nature* (núm. 447, págs. 407-412).

**Celniker, S. E. y otros** (2009). Unlocking the secrets of the genome. *Nature* (núm. 459, págs. 927-930).

**Cline, M. S.; Kent, W. J.** (2009). "Understanding genome browsing". *Nature Biotechnology* (núm. 27, págs. 153-155).

**Flicek, P. y otros** (2010). "Ensembl's 10th year". *Nucleic Acids Research* (núm. 38, págs. D557-D562).

**Harrow, J. y otros** (2006). "GENCODE: producing a reference annotation for ENCODE". *Genome Biology* (supl. 1 núm. S4, págs. 1-9).

**Hubbard, T. y otros** (2002). "The Ensembl genome database project". *Nucleic Acids Research* (núm. 30, págs. 38-41).

**Kent, W. J.** (2002). "BLAT: the BLAST-like alignment tool". *Genome Research* (núm. 12, págs. 656-664).

**Kent, W. J. y otros** (2002). "The human genome browser at UCSC". *Genome Research* (núm. 12, págs. 996-1006).

**Mayor, C. y otros** (2000). "VISTA: visualizing global DNA sequence alignments of arbitrary length". *Bioinformatics* (núm. 16, págs. 1046-1047).

**Mouse Genome Database Group** (2008). "The Mouse Genome Database (MGD): mouse biology and model systems". *Nucleic Acids Research* (núm. 36, págs. D724-728).

**Ouzounis, C. A.; Valencia, A.** (2003). "Early bioinformatics: the birth of a discipline - a personal view". *Bioinformatics* (núm. 19, págs. 2176-2190).

**Rosenbloom, K. R. y otros** (2010). "ENCODE whole-genome data in the UCSC Genome Browser". *Nucleic Acids Research* (núm. 38, págs. D620-D625).

**Stein, L. D. y otros** (2002). "The generic genome browser: a building block for a model organism system database". *Genome Research* (núm. 12, págs. 1599-1610).

**The Encode Project Consortium** (2007). "Identification and analysis of functional elements in 1% of the human genome by the ENCODE pilot project". *Nature* (núm. 447, págs. 799-816).

**The Encode Project Consortium** (2010). "Identification of functional elements and regulatory circuits by Drosophila modENCODE". *Science* (núm. 330, págs. 1787-1797).

**The Flybase Consortium** (2009). "FlyBase: enhancing Drosophila Gene Ontology annotations". *Nucleic Acids Research* (núm. 37, págs. D555-D559).

**Wheeler, D. L. y otros** (2004). "Database resources of the NCBI: update". *Nucleic Acids Research* (núm. 32, págs. D35-D40).