Databases

Objective of this lecture

The goal of this lecture is to introduce the databases that store these data and strategies to extract information from them.

What is a Gene and Genome

- A genome is the complete set of DNA.
- A gene is a section of DNA that makes up the building plans for physical traits.

Genes vary a lot in size:

- > Humans: average 3000 bp.
- Largest 2.4 million bp.
- A base pair (bp) is a unit consisting of two nucleobases bound to each other by hydrogen bonds.
- Genes are separated by sequences with unknown function.

Genomes of 1000s of organisms have been completely sequenced (300,000 species!)

 Publicly available databanks now contain quadrillions (>10¹⁵) of nucleotides of DNA sequence data, soon to be quintillions (>10¹⁸ bases).

What is a DNA Sequence?

- The DNA double helix is made up of a series of chemical bases stung along a sugar backbone.
- There are 4 bases usually represented by the letters A, T, C and G.
- The linear sequence in which these bases occur determines all the instructions for building an organism.

Features of DNA Sequence analysis

• Detecting open reading frames (ORF):

- Initial codon: ATG
- Stop codon: TGA, TAA, TAG
- Features may be used as indicators of potential protein coding regions in DNA:
 - Sufficient ORF length
 - Recognition of flanking Kozak consensus sequence ((gcc) gccRccAUGG)
 - Patterns of codon usage
 - A general preference for G/C over A/T in third base (wobble) position of a codon
 - Ribosome binding sites

• Alignment with homologous protein sequences

An example of DNA sequence

ACCACTTTCACAATCTGCTAGCAAAGGTTATGCAGCGCGTGAACATGATCATGGCAGAATCACCAGGCCT CATCACCATCTGCCTTTTAGGATATCTACTCAGTGCTGAATGTACAGTTTTTCTTGATCATGAAAAACGCC AACAAAATTCTGAATCGGCCAAAGAGGTATAATTCAGGTAAATTGGAAGAGTTTGTTCAAGGGAACCTTG AGAGAGAATGTATGGAAGAAAAGTGTAGTTTTGAAGAAGCACGAGAAGTTTTTGAAAAACACTGAAAGAAC AACTGAATTTTGGAAGCAGTATGTTGATGGAGATCAGTGTGAGTCCAATCCATGTTTAAATGGCGGCAGT TGCAAGGATGACATTAATTCCTATGAATGTTGGTGTCCCTTTGGATTTGAAGGAAAGAACTGTGAATTAG ATGTAACATGTAACATTAAGAATGGCAGATGCGAGCAGTTTTGTAAAAAATAGTGCTGATAACAAGGTGGT TTGCTCCTGTACTGAGGGATATCGACTTGCAGAAAACCAGAAGTCCTGTGAACCAGCAGTGCCATTTCCA TGTGGAAGAGTTTCTGTTTCACAAACTTCTAAGCTCACCCGTGCTGAGACTGTTTTTCCTGATGTGGACT ATGTAAATTCTACTGAAGCTGAAACCATTTTGGATAACATCACTCAAAGCACCCAATCATTTAATGACTT CACTCGGGTTGTTGGTGGAGAAGATGCCAAACCAGGTCAATTCCCTTGGCAGGTTGTTTTGAATGGTAAA GTTGATGCATTCTGTGGAGGCTCTATCGTTAATGAAAAATGGATTGTAACTGCTGCCCACTGTGTTGAAA CTGGTGTTAAAATTACAGTTGTCGCAGGTGAACATAATATTGAGGAGACAGAACATACAGAGCAAAAGCG AAATGTGATTCGAATTATTCCTCACCACAACTACAATGCAGCTATTAATAAGTACAACCATGACATTGCC CTTCTGGAACTGGACGAACCCTTAGTGCTAAACAGCTACGTTACACCTATTTGCATTGCTGACAAGGAAT ACACGAACATCTTCCTCAAATTTGGATCTGGCTATGTAAGTGGCTGGGGGAAGAGTCTTCCACAAAGGGAG ATCAGCTTTAGTTCTTCAGTACCTTAGAGTTCCACTTGTTGACCGAGCCACATGTCTTCGATCTACAAAG TTCACCATCTATAACAACATGTTCTGTGCTGGCTTCCATGAAGGAGGTAGAGATTCATGTCAAGGAGATA GTGGGGGGACCCCATGTTACTGAAGTGGAAGGGACCAGTTTCTTAACTGGAATTATTAGCTGGGGTGAAGA GTGTGCAATGAAAGGCAAATATGGAATATATACCAAGGTATCCCGGTATGTCAACTGGATTAAGGAAAAA ACAAAGCTCACTTAATGAAAGATGGATTTCCAAGGTTAATTCATTGGAATTGAAAATTAACAGGGCCTCT CACTAACTAATCACTTTCCCCATCTTTTGTTAGATTTGAATATATACATTCTATGATCATTGCTTTTTCTC TTTACAGGGGGAGAATTTCATATTTTACCTGAGCAAATTGATTAGAAAATGGAACCACTAGAGGAATATAA TGTGTTAGGAAATTACAGTCATTTCTAAGGGCCCAGCCCTTGACAAAATTGTGAAGTTAAATTCTCCACT CCATCTTCCCGATCTTCTTTGCTTCTCCCAACCAAAACATCAATGTTTATTAGTTCTGTATACAGTACAGG CTAAAACTCATCAAAAACACTACTCCTTTTCCTCTACCCTATTCCTCAATCTTTTACCTTTTCCAAATCC CAATCCCCCAAATCAGTTTTTCTCTTTCTTACTCCCTCTCTCCCCTTTTACCCCTCCATGGTCGTTAAAGGAG AGATGGGGGGGCATCATTCTGTTATACTTCTGTACACAGTTATACATGTCTATCAAACCCCAGACTTGCTTC CGTAGTGGAGACTTGCTTTTCAGAACATAGGGATGAAGTAAGGTGCCTGAAAAGTTTGGGGGGAAAAGTTT TAGTGTGTGTGTATGCGTGTGTGTGTAGACACACACGCATACACACATATAATGGAAGCAATAAGCCATTCT AAGAGCTTGTATGGTTATGGAGGTCTGACTAGGCATGATTTCACGAAGGCAAGATTGGCATATCATTGTA ACTAAAAAAGCTGACATTGACCCAGACATATTGTACTCTTTCTAAAAATAATAATAATAATGCTAACAGA AAGAAGAGAACCGTTCGTTTGCAATCTACAGCTAGTAGAGACTTTGAGGAAGAATTCAACAGTGTGTCTT CAGCAGTGTTCAGAGCCAAGCAAGAAGTTGAAGTTGCCTAGACCAGAGGACATAAGTATCATGTCTCCTT TAACTAGCATACCCCGAAGTGGAGAAGGGTGCAGCAGGCTCAAAGGCATAAGTCATTCCAATCAGCCAAC TAAGTTGTCCTTTTCTGGTTTCGTGTTCACCATGGAACATTTTGATTATAGTTAATCCTTCTATCTTGAA TCTTCTAGAGAGTTGCTGACCAACTGACGTATGTTTCCCTTTGTGAATTAATAAACTGGTGTTCTGGTTC AT

An example of DNA sequence

ACCACTTTCACAATCTGCTAGCAAAGGTTATGCAGCGCGTGAACATGATCATGGCAGAATCACCAGGCCT CATCACCATCTGCCTTTTAGGATATCTACTCAGTGCTGAATGTACAGTTTTTCTTGATCATGAAAAACGCC AACAAAATTCTGAATCGGCCAAAGAGGTATAATTCAGGTAAATTGGAAGAGTTTGTTCAAGGGAACCTTG AGAGAGAATGTATGGAAGAAAAGTGTAGTTTTGAAGAAGCACGAGAAGTTTTTGAAAAACACTGAAAGAAC AACTGAATTTTGGAAGCAGTATGTTGATGGAGAGCAGTGTGAGTCCAATCCATGTTTAAATGGCGGCAGT TGCAAGGATGACATTAATTCCTATGAATGTTGGTGTCCCTTTGGATTTGAAGGAAAGAACTGTGAATTAG ATGTAACATGTAACATTAAGAATGGCAGATGCGAGCAGTTTTGTAAAAATAGTGCTGATAACAAGGTGGT TTGCTCCTGTACTGAGGGATATCGACTTGCAGAAAACCAGAAGTCCTGTGAACCAGCAGTGCCATTTCCA TGTGGAAGAGTTTCTGTTTCACAAACTTCTAAGCTCACCCGTGCTGAGACTGTTTTTCCTGATGTGGACT ATGTAAATTCTACTGAAGCTGAAACCATTTTGGATAACATCACTCAAAGCACCCCAATCATTTAATGACTT CACTCGGGTTGTTGGTGGAGAAGATGCCAAACCAGGTCAATTCCCTTGGCAGGTTGTTTTGAATGGTAAA GTTGATGCATTCTGTGGAGGCTCTATCGTTAATGAAAAATGGATTGTAACTGCTGCCCACTGTGTTGAAA CTGGTGTTAAAATTACAGTTGTCGCAGGTGAACATAATATTGAGGAGACAGAACATACAGAGCAAAAGCG AAATGTGATTCGAATTATTCCTCACCACAACTACAATGCAGCTATTAATAAGTACAACCATGACATTGCC CTTCTGGAACTGGACGAACCCTTAGTGCTAAACAGCTACGTTACACCTATTTGCATTGCTGACAAGGAAT ACACGAACATCTTCCTCAAATTTGGATCTGGCTATGTAAGTGGCTGGGGGAAGAGTCTTCCACAAAGGGAG ATCAGCTTTAGTTCTTCAGTACCTTAGAGTTCCACTTGTTGACCGAGCCACATGTCTTCGATCTACAAAG TTCACCATCTATAACAACATGTTCTGTGCTGCCTGCCATGAAGGAGGTAGAGATTCATGTCAAGGAGATA GTGGGGGGACCCCA TTACTGAAGTGGAAGGGACCAGTTTCTTAACTGGAATTATTAGCTGGGGTGAAGA GTGTGCAATGAAAGGCAAATATGGAATATATACCAAGGTATCCCGGTATGTCAACTGGATTAAGGAAAAA ACAAAGCTCACT TAA GAAAGATGGATTTCCAAGGTTAATTCATTGGAAATTGAAAATTAACAGGGCCTCT CACTAACTAATCACTTTCCCCATCTTTTGTTAGATTTGAATATATACATTCTATGATCATTGCTTTTTCTC TTTACAGGGGAGAATTTCATATTTTACCTGAGCAAATTGATTAGAAAATGGAACCACTAGAGGAATATAA TGTGTTAGGAAATTACAGTCATTTCTAAGGGCCCCAGCCCTTGACAAATTGTGAAGTTAAATTCTCCACT CCATCTTCCCGATCTTCTTTGCTTCTCCCAACCAAAACATCAATGTTTATTAGTTCTGTATACAGTACAGG CTAAAACTCATCAAAAACACTACTCCTTTTCCTCTACCCTATTCCTCAATCTTTTACCTTTTCCAAATCC CAATCCCCCAAATCAGTTTTTCTCTTTCTTTACTCCCTCTCTCCCCTTTTACCCCTCCATGGTCGTTAAAGGAG AGATGGGGGGGCATCATTCTGTTATACTTCTGTACACAGTTATACATGTCTATCAAACCCCAGACTTGCTTC CGTAGTGGAGACTTGCTTTTCAGAACATAGGGATGAAGTAAGGTGCCTGAAAAGTTTGGGGGGAAAAGTTT TAGTGTGTGTGTATGCGTGTGTGTGTGGAGACACACACGCATACACACATATAATGGAAGCAATAAGCCATTCT AAGAGCTTGTATGGTTATGGAGGTCTGACTAGGCATGATTTCACGAAGGCAAGATTGGCATATCATTGTA ACTAAAAAAGCTGACATTGACCCCAGACATATTGTACTCTTTCTAAAAATAATAATAATAATGCTAACAGA AAGAAGAGAACCGTTCGTTTGCAATCTACAGCTAGTAGAGACTTTGAGGAAGAATTCAACAGTGTGTCTT CAGCAGTGTTCAGAGCCAAGCAAGAAGTTGAAGTTGCCTAGACCAGAGGACATAAGTATCATGTCTCCTT TAACTAGCATACCCCGAAGTGGAGAAGGGTGCAGCAGGCTCAAAGGCATAAGTCATTCCAATCAGCCAAC TAAGTTGTCCTTTTCTGGTTTCGTGTTCACCATGGAACATTTTGATTATAGTTAATCCTTCTATCTTGAA TCTTCTAGAGAGTTGCTGACCAACTGACGTATGTTTCCCCTTTGTGAATTAATAAACTGGTGTTCTGGTTC AT

An example of DNA sequence

ATG 🕼 AG CGC GTG AAC ATG ATC ATG GCA GAA TCA CCA GGC CTC ATC ACC ATC TGC CTT TTA GGA TAT CTA CTC AGT GCT GAA TGT ACA GTT TTT CTT GAT CAT GAA AAC GCC AAC AAA ATT CTG AAT CGG CCA AAG AGG TAT AAT TCA GGT GAA GTT TTT GAA AAC ACT GAA AGA ACA ACT GAA TTT TGG AAG CAG TAT GTT GAT GGA GAT CAG TGT GAG TCC AAT CCA TGT TTA AAT GGC GGC AGT TGC AAG GAT GAC ATT AAT TCC TAT GAA TGT TGG TGT CCC TTT GGA TTT GAA GGA AAG AAC TGT GAA TTA GAT GTA ACA TGT AAC ATT AAG AAT GGC AGA TGC GAG CAG TTT TGT AAA AAT AGT GCT GAT AAC AAG GTG GTT TGC TCC TGT ACT GAG GGA TAT CGA CTT GCA GAA AAC CAG AAG TCC TGT GAA CCA GCA GTG CCA TTT CCA TGT GGA AGA GTT TCT GTT TCA CAA ACT TCT AAG CTC ACC CGT GCT GAG ACT GTT TTT CCT GAT GTG GAC TAT GTA AAT TCT ACT GAA GCT GAA ACC ATT TTG GAT AAC ATC ACT CAA AGC ACC CAA TCA TTT AAT GAC TTC ACT CGG GTT GTT GGT GGA GAA GAT GCC AAA CCA GGT CAA TTC CCT TGG CAG GTT GTT TTG AAT GGT AAA GTT GAT GCA TTC TGT GGA GGC TCT ATC GTT AAT GAA AAA TGG ATT GTA ACT GCT GCC CAC TGT GTT GAA ACT GGT GTT AAA ATT ACA GTT GTC GCA GGT GAA CAT AAT ATT GAG GAG ACA GAA CAT ACA GAG CAA AAG CGA AAT GTG ATT CGA ATT ATT CCT CAC CAC AAC TAC AAT GCA GCT ATT AAT AAG TAC AAC CAT GAC ATT GCC CTT CTG GAA CTG GAC GAA CCC TTA GTG CTA AAC AGC TAC GTT ACA CCT ATT TGC ATT GCT GAC AAG GAA TAC ACG AAC ATC TTC CTC AAA TTT GGA TCT GGC TAT GTA AGT GGC TGG GGA AGA GTC TTC CAC AAA GGG AGA TCA GCT TTA GTT CTT CAG TAC CTT AGA GTT CCA CTT GTT GAC CGA GCC ACA TGT CTT CGA TCT ACA AAG TTC ACC ATC TAT AAC AAC ATG TTC TGT GCT GGC TTC CAT GAA GGA GGT AGA GAT TCA TGT CAA GGA GAT AGT GCC GGA CCC CAT GTT ACT GAA GTG GAA GGG ACC AGT TTC TTA ACT GGA ATT ATT AGC TGG GGT GAA GAG TGT GCA ATG AAA GGC AAA TAT GGA ATA TAT ACC AAG GTA TCC CGG TAT GTC AAC TGG ATT AAG GAA AAA ACA AAG CTC ACT TAA

Figure 1.5 Practical Bioinformatics (© Garland Science 2013)

ACCACTTTCACAATCTGCTAGCAAAGGTT

What is a Protein Sequence?

- Proteins are complex molecules which control most aspects of cell biology.
- Constructed of small subunits called amino acids.
- There are 20 main types of amino acids that form proteins.
- Assembled by 'reading' (or translating) the DNA sequence.
- Every set of 3 bases (e.g. ATG) corresponds to an amino acid.
- So a protein is built up one amino acid at a time according to the DNA blueprint.

Sequence analysis

•Sequence analysis is the process of subjecting a DNA, RNA or peptide sequence to any of a wide range of analytical methods to understand its features, function, structure, or evolution.

Applications

Sequence analysis can be used in following fields of molecular biology:

- The comparison of sequences in order to find similarity, often to infer if they are related (homologous)
- Identification of intrinsic features of the sequence such as active sites, post translational modification sites, genestructures, reading frames, distributions of introns and exons and regulatory elements.
- Identification of sequence differences and variations such as point mutations and single nucleotide polymorphism (SNP) in order to get the genetic marker.
- Revealing the evolution and genetic diversity of sequences and organisms
- Identification of molecular structure from sequence alone

Uses of sequence analysis

- 1. Proteomics (determining protein structure).
- 2. Sequence Alignment.
- 3. Genomics.
- 4. Evolutionary Biology (Phylogenetics).
- 5. Systems biology.

Looking at DNA sequences

- Analysis of DNA or protein sequences is a frequent requirement of research.
 - Locating genes within a sequence.
 - Comparing two sequences for similarity.
 - Searching for similar genes (orthologues) in other organisms.

Using a DNA Sequence In



Centralized Databases Store DNA Sequences

• We begin with three main sites that have been responsible for storing nucleotide sequence data from 1982 to the present



FIGURE 2.1 The nucleotide collections of GenBank at NCBI, EMBL-Bank at the European Bioinformatics Institute, and DDBJ at the DNA Data Bank of Japan are all coordinated by the International Nucleotide Sequence Database Collaboration (INSDC).



What is database?

- It is a <u>collection</u> of data from different sources.
- Each stored data is called <u>entry</u>.
- Used to <u>retrieve</u> data (or entrez).
- Searchable with keyword (called query).

Why do we need databases?

- To arrange these large number of sequences, i.e., to make a library and catalogue of genes, RNA and proteins
- To remove redundant sequences. So, we have only unique sequences per gene.
- To annotate (name) the sequences. So, we know what the protein does or where the gene starts and ends.

Scales of DNA Base Pairs

TABLE 2.1 Scales of DNA base pairs.

Base pairs	Unit	Abbreviation	Example
1	1 base pair	1 bp	
1000	1 kilobase pair	1 kb	Size of a typical coding region of a gene
1,000,000	1 megabase pair	1 Mb	Size of a typical bacterial genome
10 ⁹	1 gigabase pair	1 Gb	The human genome is 3 billion base pairs
10 ¹²	1 terabase pair	1 Tb	
10 ¹⁵	1 petabase pair	1 Pb	

GenBank

- A database consisting of DNA and protein sequences.
- Contains bibliographic and biological annotation.
- Data are available free of charge from NCBI.
- Over 310,000 different species!

GenBank

- Over 1000 new species added per month
- Over the past 30 years the number of bases in GenBank has doubled approximately every 18 months
- GenBank received submissions since 1982, including sequences from thousands of individual submitters.

Organisms in GenBank

TABLE 2.3 Taxa represented in GenBank.

Ranks	Higher taxa	Genus	Species	Lower taxa	Total
Archaea	143	140	525	0	808
Bacteria	1,370	2,611	13,331	819	18,131
Eukaryota	20,443	67,606	297,207	22,608	407,864
Fungi	1,550	4,620	29,450	1,128	36,748
Metazoa	14,670	45,517	145,044	11,428	216,659
Viridiplantae	2,622	14,680	113,529	9,789	140,620
Viruses	618	442	2,349	0	3,409
All taxa	22,603	70,806	313,443	23,427	430,279

Source: GenBank, NCBI, (#) http://www.ncbi.nlm.nih.gov/Taxonomy/txstat.cgi.

NCBI and EBI

Two of the main centralized bioinformatics hubs:

- The National Center for Biotechnology Information (NCBI)
 - The European Bioinformatics Institute (EBI).

In many cases those sites begin with similar raw data and then provide distinct ways of organizing, analyzing, and displaying data across a broad range of bioinformatics applications.

NCBI

- Creates public databases.
- Conducts research in computational biology
- Develops software tools for analyzing genome data
- Disseminates biomedical information

NCBI Resources

- Entrez: integrates scientific literature, DNA, and protein sequence databases, 3D protein structure data, population study datasets, and assemblies of complete genomes into a tightly coupled system
- **PubMed**: search service from the National Library of Medicine (NLM) that provides access to over 24 million citations
- **BLAST**: (Basic Local Alignment Search Tool) is NCBI's sequence similarity search
- **OMIM**: Online Mendelian Inheritance in Man (OMIM) is a catalog of human genes and genetic disorders
- Books: NCBI offers about 200 books online
- o other

NCBI Entrez

Entrez:

- Accessed from the home page of NCBI
- Provides links to results from 40 different NCBI databases
 - scientific literature, DNA, and protein sequence databases, 3D protein structures, population study datasets, and assemblies of complete genomes



Entrez Main Screen

Search across databases			
Welcome to th		GO Clear Help	
Welcome to th	ne Entrez cr	oss-database search page	
PubMed: biomedical literature citations and abstracts	0	Books: online books	0
PubMed Central: free, full text journal articles	۵	OMIM: online Mendelian Inheritance in Man	0
Site Search: NCBI web and FTP sites	۲	OMIA: online Mendelian Inheritance in Animals	0
CoreNucleotide: Core subset of nucleotide sequence records	۲	dbGaP: genotype and phenotype	0
EST: Expressed Sequence Tag records	0	UniGene: gene-oriented clusters of transcript sequences	0
GSS: Genome Survey Sequence records	0	CDD: conserved protein domain database	0
Protein: sequence database	0	3D Domains: domains from Entrez Structure	0
Genome: whole genome sequences	0	UniSTS: markers and mapping data	0
Structure: three-dimensional macromolecular structures	0	PopSet: population study data sets	0
Taxonomy: organisms in GenBank	0	GEO Profiles: expression and molecular abundance profiles	0
SNP: single nucleotide polymorphism	0	GEO DataSets: experimental sets of GEO data	0
Gene: gene-centered information	0	Cancer Chromosomes: cytogenetic databases	0
HomoloGene: eukaryotic homology groups	0	PubChem BioAssay: bioactivity screens of chemical substances	0
GENSAT: gene expression atlas of mouse central nervous system	0	PubChem Compound: unique small molecule chemical structures	Ø
Probe: sequence-specific reagents	0	PubChem Substance: deposited chemical substance records	0

۲

0

 $\ensuremath{\textbf{Journals:}}$ detailed information about the journals indexed in PubMed and other Entrez databases

NLM Catalog: catalog of books, journals, and audiovisuals in the NLM

U.

Č

collections

۲



NCBI PubMed

PubMed: provides access to over 24 million citations.

PubMed citations come from

- 1) MEDLINE indexed journals (largest subset of PubMed)
- 2) Journals/manuscripts deposited in PMC, and
- 3) NCBI Bookshelf.



PubMed Quick Start Guide

PubMed Mobile

🗧 NCBI 🛛 Resources 🗹 How	v To 🗹	
Publed.gov Pu	ubMed 🗘	
US National Library of Medicine National Institutes of Health	Advanced	

Format: Abstract -

Send to

J Infect Public Health. 2016 Sep 15. pii: S1876-0341(16)30146-0. doi: 10.1016/j.jiph.2016.09.007. [Epub ahead of print]

Building predictive models for MERS-CoV infections using data mining techniques.

<u>Al-Turaiki I¹, Alshahrani M², Almutairi T³.</u>

Author information

Abstract

BACKGROUND: Recently, the outbreak of MERS-CoV infections caused worldwide attention to Saudi Arabia. The novel virus belongs to the coronaviruses family, which is responsible for causing mild to moderate colds. The control and command center of Saudi Ministry of Health issues a daily report on MERS-CoV infection cases. The infection with MERS-CoV can lead to fatal complications, however little information is known about this novel virus. In this paper, we apply two data mining techniques in order to better understand the stability and the possibility of recovery from MERS-CoV infections.

METHOD: The Naive Bayes classifier and J48 decision tree algorithm were used to build our models. The dataset used consists of 1082 records of cases reported between 2013 and 2015. In order to build our prediction models, we split the dataset into two groups. The first group combined recovery and death records. A new attribute was created to indicate the record type, such that the dataset can be used to predict the recovery from MERS-CoV. The second group contained the new case records to be used to predict the stability of the infection <u>ba</u>sed on the current status attribute.

What is the difference between pubmed and medline?

- Pubmed is an interface used to search Medline, as well as additional biomedical content. Medline is an interface for searching only Medline content.
- Pubmed is more user-friendly and allows you to search through more content than Ovid Medline. However, Medline allows you to perform a more focused search. You will get slightly different results by searching in each database.
- In Pubmed, in addition to Medline articles, you will have access to PubMedCentral papers, which are full text articles deposited to promote open access, and articles that are "in press" that is, prior to being indexed with MeSH terms, and articles submitted by publishers, "ahead of print." This is why if you search for the same term in Medline and in Pubmed, you may obtain as many as ten thousand more articles in Pubmed.

Pubmed database

S NCBI Resources 🖸 How To 🗹		Sign in to NCBI
US National Library of Medicine National Institutes of Health	Advanced	Search
Pub Pub MEDL full-tex	Med ed comprises more than 25 million citations for biomedical literature from INE, life science journals, and online books. Citations may include links to tt content from PubMed Central and publisher web sites.	PubMed Commons PubMed Commons Featured comment - Jan 22 A Messori highlights literature on frequently documented drug- drug interactions in patients with cancer. <u>1.usa.gov/1kSAsDd</u>
Using PubMed	PubMed Tools	More Resources
PubMed Quick Start Guide	PubMed Mobile	MeSH Database
Full Text Articles	Single Citation Matcher	Journals in NCBI Databases
PubMed FAQs	Batch Citation Matcher	Clinical Trials
PubMed Tutorials	Clinical Queries	E-Utilities (API)
New and Noteworthy	Topic-Specific Queries	LinkOut

You are here: NCBI > Literature > PubMed

GETTING STARTED
NCBI Education
NCBI Help Manual
NCBI Handbook
Training & Tutorials
Submit Data

RESOURCES Chemicals & Bioassays Data & Software DNA & RNA Domains & Structures Genes & Expression Genetics & Medicine Genomes & Maps Homology Literature Proteins Sequence Analysis Taxonomy Variation

POPULAR PubMed Bookshelf PubMed Central PubMed Health BLAST Nucleotide Genome SNP Gene Protein PubChem

FEATURED

Genetic Testing Registry PubMed Health GenBank Reference Sequences Gene Expression Omnibus Map Viewer Human Genome Mouse Genome Influenza Virus Primer-BLAST Sequence Read Archive

Write to the Help Desk

NCBI INFORMATION About NCBI Research at NCBI NCBI News NCBI FTP Site NCBI on Facebook NCBI on Tactebook NCBI on YouTube



National Center for Biotechnology Information, U.S. National Library of Medicine 8600 Rockville Pike, Bethesda MD, 20894 USA Policies and Guidelines | Contact

Results of Pubmed

S NCBI Resources 🗹	How To 🕑	Sign in to NCBI
Public dec.gov US National Library of Medicine National Institutes of Health	PubMed (bacterial[Title]) AND promoters[Title] Create RSS Create alert Advanced	Search Help
Article types Clinical Trial Review Customize Text availability Abstract	Summary • 20 per page • Sort by Most Recent • Send to: • Search results Items: 1 to 20 of 56 << First < Prev Page 1 of 3 Next > Last >>	Filters: <u>Manage Filters</u> New feature Try the new Display Settings option - Sort by Relevance
Free full text Full text PubMed Commons Reader comments Trending articles	 The bacterial cell cycle regulator GcrA is a σ70 cofactor that drives gene expression from a subset of methylated promoters. Haakonsen DL, Yuan AH, Laub MT. Genes Dev. 2015 Nov 1;29(21):2272-86. doi: 10.1101/gad.270660.115. PMID: 26545812 Similar articles 	Titles with your search terms Bacterial genetics by flow cytometry: rapid isolation of Salmonella typhi [Mol Microbiol. 1996] Regulation at complex bacterial promoters: how bacteria use different [Curr Opin Microbiol. 2004]
Publication dates 5 years 10 years Custom range Species Humans Other Animals	 <u>32 Electrostatic properties of T7-like phages promoters for host bacterial and native viral RNA</u> <u>polymerases.</u> Osypov AA, Kamzolova SG. J Biomol Struct Dyn. 2015;33 Suppl 1:19-20. doi: 10.1080/07391102.2015.1032581. No abstract available. PMID: 26103243 <u>Similar articles</u> 	Identification of an UP element consensus sequence for ba [Proc Natl Acad Sci U S A. 1998] See more Find related data
<u>Clear all</u> Show additional filters	 Base flipping in open complex formation at bacterial promoters. Karpen ME, deHaseth PL. Biomolecules. 2015 Apr 28;5(2):668-78. doi: 10.3390/biom5020668. Review. PMID: 25927327 Free PMC Article 	Find items

Using Pubmed

Go to http://www.ncbi.nlm.nih.gov/pubmed



Using pubmed



Article types	Summary - 20 per page - Sort by Most Recent -
Clinical Trial	
Review	
Customize	Search results
Text availability	Items: 1 to 20 of 56

56 articles after using advanced option
NCBI

- **BLAST:** sequence similarity search
- **OMIM:** a catalog of human genes and genetic disorders.
- Books: NCBI offers about 200 books online.

EBI

• The EBI website is comparable to NCBI in its scope and mission, and it represents a complementary, independent resource.



Examples: blast, keratin, bfl1, Janet Thornton ...

EBI

• EBI features six core molecular databases:

- EMBL-Bank is the repository of DNA and RNA sequences that is complementary to GenBank and DDBJ
- Swiss-Prot and TrEMBL are two protein databases
- MSD is a protein structure database
- Ensembl is one of the main genome browsers
- ArrayExpress is for gene expression data

Access to Information via Gene Resource at NCBI

• To illustrate the use of NCBI Gene we search for human **beta globin**.

Result of a search for "beta globin" in NCBI Gene (via an

i 🔒 https://www.ncbi.nlm.nih.gov/gene/?term=beta+globin 🛛 🕻 🔍 Search 🖉 🖓 Search								
3 NCBI Resources 🖂	How To 🖸						<u>Sign in to</u>	NC
Gene	Gene ÷ b C	eta globin Greate RSS Create		Search		He		
A NCBI is currently redire	cting web traffic to HTT	PS. <u>Read more</u> abo	ut our https testing.					
Sene sources Tabular - 20 per page - Sort by Relevance - Send to: - Genomic Filters: Manage Filters								ar >
ategories Iternatively spliced Innotated genes Ion-coding Protein-coding	Search results Items: 1 to 20 of See also 8 dis	s f 130 continued or replace	<< First	< Prev Page 1 of 7 Nex	t > Last >>	Results by taxon Top Organisms [Tree] Homo sapiens (40)		
seudogene	Name/Gene ID	Description	Location	Aliases	MIM	Danio rerio (6)		
equence ontent CDS Insembl tefSeg	□ <u>HBB</u> ID: 3043	hemoglobin subunit beta [<i>Homo sapiens</i> (human)]	Chromosome 11, NC_000011.10 (52254665227071, complement)	CD113t-C, beta-globin	141900	Salmo salar (6) Rattus norvegicus (5) All other taxa (45) More		
tefSeqGene Itatus clear	□ <u>hbg1</u> ID: 394453	hemoglobin subunit gamma 1 [<i>Xenopus</i> <i>tropicalis</i>	Chromosome 9, NC_030685.1 (3697888136980393, complement)	beta-globin, hbb1, hbga, hbgr, hsggl1		Find related data Database: Select	÷	-

NCBI Gene entry for human beta globin

HBB hemoglobin subunit beta [Homo sapiens (human)]

Gene ID: 3043, updated on 2	^{22-Sep-2016} Information is provided on the gene structure
	and chromosomal location, as well as a
Summary	
	summary of the protein's function.
Official Symbol	HBB provided by
Official Full Name	hemoglobin subunit beta provided by HGNC
Primary source	HGNC:HGNC:4827
See related	Ensembl:ENSG00000244734 HPRD:00786; MIM:141900; Vega:OTTHUMG00000066678
Gene type	protein coding
RefSeq status	REVIEWED
Organism	Homo sapiens
Lineage	Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Euarchontoglires;
	Primates; Haplorrhini; Catarrhini; Hominidae; Homo
Also known as	CD113t-C; beta-globin
Summary	The alpha (HBA) and beta (HBB) loci determine the structure of the 2 types of polypeptide chains in adult
	hemoglobin, Hb A. The normal adult hemoglobin tetramer consists of two alpha chains and two beta chains.
	Mutant beta globin causes sickle cell anemia. Absence of beta chain causes beta-zero-thalassemia. Reduced
	amounts of detectable beta globin causes beta-plus-thalassemia. The order of the genes in the beta-globin
	cluster is 5'-epsilon gamma-G gamma-A delta beta3'. [provided by RefSeq, Jul 2008]
Orthologs	all

Access to Information via Gene Resource at NCBI

Genomic regions, transcripts, and products

Go to reference sequence details

EAGTA

Graphice

Go to nucleotido:

<!?

ConBonk

Genomic Sequence: NC_000011.10 Chromosome 11 Reference GRCh38.p7 Primary Assembly

							Go to nucl				
5	NC_000011.10	: 5.2M5.2M (2.	1Kbp) C 🗸 🛛 🤇			+	ATG	🔀 Tools 🗸	🚠 🛛 🔅 Track	s • ಿ 🖗	
00	5,227,200	5,227 K	5,226,800	5,226,600	5,226,400	5,226,200	5,226 K	5,225,800	5,225,600	5,225,400)
Gene	es, NCBI Hom	o sapiens A	nnotation	Release 10	8, 2 📧						*
						нвв			OB?X	5221061	.5231
Con	NM_000518.4		,				P		Gene annota	tions provide	ed by
Gene	es, Ensembr		\rightarrow	}			~	>	* *		~
dbSl	NP Build 147	(Homo sapi	ens Annota	tion Relea:	se 107) al:	l data					*
Clir	nVar Short V	ariations b	ased on db	SNP Build 3	147 📧						×
Cite	7 2 7 3 ed Variants, 6 1 7 4	+ + + 6 1 dbSNP Buil + + + 6 1	+ + + + + d 147 (Hom + + + + +	<pre>+ + o sapiens i + + 1 2</pre>	Anno া	1	3 1 1	1 5 + +	+ + 9	1	×
RNA-	-seq exon co	verage, agg	regate (fi	ltered), NG	CBI Homo sa	apiens Annot	ation Rele	ase 108 - 1	log base 2 s	caled	×
				2108818 4036 0							

Sequence format in databases

- A sequence format defines the permitted layout and content of text in a file.
- The <u>FASTA format</u> is a very widely used format. It consists of a header line starting with a ">" character followed by a code identifying the sequence and, very often, some text describing the sequence. The header line is followed by one or more lines containing the sequence itself.
- FASTA files may contain one or more sequences

FASTA Format

Homo sapiens chromosome 11, GRCh38.p7 Primary Assembly

description

sequence

NCBI Reference Sequence: NC_000011.10

GenBank Graphics

Di|568815587:c5227071-5225466 Homo sapiens chromosome 11, GRCh38.p7 Primary Assembly

FASTA is both an alignment program and a commonly used sequence format

FASTA files may contain one or more sequences

>crab_anapl ALPHA CRYSTALLIN B CHAIN (ALPHA(B)-CRYSTALLIN). MDITIHNPLIRRPLFSWLAPSRIFDQIFGEHLQESELLPASPSLSPFLMR SPIFRMPSWLETGLSEMRLEKDKFSVNLDVKHFSPEELKVKVLGDMVEIH GKHEERQDEHGFIAREFNRKYRIPADVDPLTITSSLSLDGVLTVSAPRKQ SDVPERSIPITREEKPAIAGAQRK

>crab_bovin ALPHA CRYSTALLIN B CHAIN (ALPHA(B)-CRYSTALLIN). MDIAIHHPWIRRPFFPFHSPSRLFDQFFGEHLLESDLFPASTSLSPFYLR PPSFLRAPSWIDTGLSEMRLEKDRFSVNLDVKHFSPEELKVKVLGDVIEV HGKHEERQDEHGFISREFHRKYRIPADVDPLAITSSLSSDGVLTVNGPRK QASGPERTIPITREEKPAVTAAPKK

>crab_chick ALPHA CRYSTALLIN B CHAIN (ALPHA(B)-CRYSTALLIN). MDITIHNPLVRRPLFSWLTPSRIFDQIFGEHLQESELLPTSPSLSPFLMR SPFFRMPSWLETGLSEMRLEKDKFSVNLDVKHFSPEELKVKVLGDMIEIH GKHEERQDEHGFIAREFSRKYRIPADVDPLTITSSLSLDGVLTVSAPRKQ SDVPERSIPITREEKPAIAGSQRK

Some identifiers used in fasta description

Database Name	Abbreviations
GenBank	gb
EMBL Data Library	emb
DNA Database of Japan	dbj
SWISS-PROT	sp
Protein Data Bank	pdb
NCBI Reference Sequence	ref

Other formats

- Beyond FASTA, the most widespread sequence formats are those used by the major sequence databases:
- EMBL<u>http://www.ebi.ac.uk/embl/Docume</u> ntation/User_manual/format.html
- o GenBank<u>http://www.ncbi.nlm.nih.gov/Ge</u> nbank/GenbankOverview.html
- SwissProt<u>http://ca.expasy.org/sprot/userm</u> an.html#whatis

HVR sequence

A hypervariable region (HVR) is a location within nuclear DNA or the D-loop of mitochondrial DNA in which base pairs of nucleotides repeat (in the case of nuclear DNA) or have substitutions (in the case of mitochondrial DNA). Changes or repeats in the hypervariable region are highly polymorphic.

Example of HVR sequence

YOUR MITOCHONDRIAL HVR I SEQUENCE 16126C, 16147T, 16183C, 16189C, 16294T, 16296T, 16297C, 16304C, 16519C

COMMAND-LINE ACCESS TO DATA AT NCBI

- The websites of NCBI, EBI, Ensembl, and other bioinformatics sites offer convenient access to resources through a web browser
- An alternative is to use command-line tools.

ACCESS TO INFORMATION GENOME BROWSERS

Genome browsers are databases with a graphical interface that presents a representation of sequence information and other data as a function of position across the chromosomes.

• Three principal genome browsers (Ensembl, UCSC, and NCBI)

ACCESS TO INFORMATION GENOME BROWSERS

Genome build refers to an assembly in which DNA sequence is collected and arranged to reflect the sequence along each chromosome. For a given organism's genome, a build is released only occasionally (typically every few years)

ACCESS TO INFORMATION GENOME BROWSERS

Annotation is the assignment of information such as the start and stop position of genes, exons, repetitive DNA elements, or other features.

- Best to use the most recent available build.
- Earlier builds have richer annotation

Annotations (Features) • It is a type of metadata (sideway database)

 It is the complete information required to be known about that particular sequence (could be DNA, RNA or protein). In case of nucleotides, annotation is the process of:

 Identifying the locations of genes, coding regions and other specific locations that are of importance in a DNA sequence or genome.

 Associating relevant information with those locations (e.g. determining what the identified genes do).

In case of proteins, annotation is the process of:

 Describing regions or sites of interest in the protein sequence, such as post-translational modifications, binding sites, enzyme active sites, local secondary structure or other characteristics.

- Provides graphical views of chromosomal locations at various levels of resolution (from several base pairs up to hundreds of millions of base pairs spanning an entire chromosome).
- Each chromosomal view is accompanied by horizontally oriented annotation tracks.

- There are hundreds of available userselected tracks in categories such as mapping and sequencing, phenotype and disease associations, genes, expression, comparative genomics, and genomic variation.
- These annotation tracks offer the Genome Browser tremendous depth and flexibility.





Our tools

- Genome Browser interactively visualize genomic data
- BLAT rapidly align sequences to the genome
- Table Browser download data from the Genome Browser database
- Variant Annotation Integrator get functional effect predictions for variant calls



UCSC Genome Browser on Human Dec. 2013 (GRCh38/hg38) Assembly. mve < >>>> Zoom in 1.5x 3x 10x base Zoom out 1.5x 3x 10x 100x chr11:5,225,464-5,227,071 1,608 bp. enter position, gene symbol or search terms ge mr11 (p15.4) 05.4 05.6 0	Â	Genor	nes Genome	Browser	Tools	Mirrors	Downloads	My Data	View	Help	About Us	
move << >>>>>>>>>>>>>>>>>>>>>>>>>>>>		UCSC Genome Browser on Human Dec. 2013 (GRCh38/hg38) Assembly										
chr11:5,225,464-5,227,071 1,608 bp. enter position, gene symbol or search terms chr11: (p15:4) ising and ising anoreal and ising and ising and ising and ising and ising and ising	m	move <<< << > >> >>> zoom in 1.5x 3x 10x base zoom out 1.5x 3x 10x 100x										
Chr11 (p15.4) 1235# 15.1 D13 11012 11.2 13.4 11014-3 022 22.1 022.3 <th02.3< th=""> 022.</th02.3<>			chr11:5,225,46	4-5,227,07	71 1,608	bp. enter po	sition, gene symbol	or search terms			go	
Scale chr11: 500 bases (h11): 5,226,800 (ENCODE v24 Comprehensive Transcript Set (only Basic displayed by default) 5,227,000 HEE 141900.0356 141900.0356 141900.0367 141900.0356 141900.0366 141900.0356 141900.0366 141900.0356 141900.0366 141900.0356 141900.0366 141900.0356 141900.0356 141900.0356 141900.0356 141900.0356 141900.0356 141900.0356 141900.0356 141900.0356 141900.0356 141900.0356 141900.0357 141900.0356 141900.0356 141900.0356 141900.0356 141900.0357 141900.0356 141900.0356 141900.0356 141900.0356 141900.0357 141900.0356 141900.0356 141900.0356 141900.0356 141900.0351 141900.0356 141900.0351 141900.0356 141900.0356 141900.0351 141900.0351 141900.0351 141900.0356 141900.0356 141900.0351 141900.0351 141900.0351 141900.0356 141900.0356 141900.0351 141900.0351 141900.0351 141900.0356 141900.0356 141900.0355 141900.0355		chr	11 (p15.4) 11 p15	5.4 15.1	p13	11p12 11.2		13,4 11q14,1	q21 22,1	q22.3 q25	.3 q25	
Chrili: 1 5,225,900 5,225,900 5,225,900 CENCODE V24 Comprehensive Transcript Set (only Basic displayed by default) 5,227,000 RefSeq Genes RefSeq Genes 141900.0556 141900.0551 141900.0556 141900.0556 141900.0556 141900.0551 141900.0556 141900.0556 141900.0556 141900.0556 141900.0551 141900.0556 141900.0556 141900.0556 141900.0556 141900.0551 141900.0556 141900.0556 141900.0556 141900.0556 141900.0551 141900.0556 141900.0551 141900.0551 141900.0551 141900.0551 141900.0551 141900.0551 141900.0551 141900.0551 141900.0551 141900.0551 141900.0551 141900.0551 141900.0551 141900.04051 141900.0551 141900.0551 141900.0551 141900.0551 141900.04051 141900.0551 141900.0551 141900.0551 141900.0551 141900.0456 141900.0551 141900.0551 141900.0551 141900.0551 141900.0456	1	Sca le			s	00 bases			hg38			
RefSeq Genes 141990.0336 141900.0357 0MIM Allelic Variants 141900.0426 141900.0682 141900.0682 141900.0585 141900.0366 141900.0561 141900.0366 141900.0366 141900.0365 141900.0365 141900.0365 141900.0365 141900.0365 141900.0355 <t< td=""><th></th><td>Chr11:</td><td>, </td><td></td><td>GENCO</td><td>DE V24 Compret</td><td>nensive Transcript</td><td>set (only Basi</td><td>ic displayed by</td><td>default)</td><td>5,2 •</td><td>227,000</td></t<>		Chr11:	, 		GENCO	DE V24 Compret	nensive Transcript	set (only Basi	ic displayed by	default)	5,2 •	227,000
141900.4333 141900.4112 141900.4357 141900.4356 141900.4356 141900.4356 141900.4356 141900.4355 141900.4356	1	нвв					RefSec	Genes				
141900.0303 141900.0305/ 141900.0305/ 141900.0305/ 141900.0022/ 141900.0022/ 141900.0417 141900.0305/ 141900.0305/ 141900.0305/ 141900.0305/ 141900.0305/ 141900.0302 141900.0007/ 141900.0305/ 141900.0305/ 141900.0305/ 141900.0305/ 141900.0305/ 141900.0302 141900.0007/ 141900.0305/ 141900.0305/ 141900.0223/ 141900.0233/ 141900.0355/ 141900.0302 141900.0407/ 141900.0231/ 141900.0231/ 141900.0355/ 141900.0355/ 141900.04051 141900.04051 141900.0231/ 141900.0231/ 141900.0355/ 141900.0355/ 141900.04051 141900.04051 141900.04051 141900.0231/ 141900.0355/ 141900.0355/ 141900.04051 141900.04051 141900.0405 141900.0232/ 141900.0345/ 141900.0345/ 141900.04051 141900.0405 141900.0405 141900.0405 141900.0405/ 141900.0405/ 141900.04051 141900.0405 141900.0405 141900.0405/ 141900.0405/ 141900.0405/ 141900.04050 141900.0405 141900.0405 141900.0405/	İ						OMIM Allel	ic Variants				
141998.0417 141998.0421 141998.0368 141998.0942 141998.0942 141998.0368 141988.0399 141998.0829 141998.0823 141998.0823 141998.0825 141998.0829 141998.0827 141998.0823 141998.0825 141998.0825 141998.0829 141998.0826 141998.0825 141998.0825 141998.0825 141998.081 141998.0826 141998.0825 141998.0825 141998.0825 141998.0851 141998.0851 141998.0825 141998.0825 141998.0825 141998.0855 141998.0851 141998.0825 141998.0825 141998.0835 141998.0855 141998.0851 141998.0825 141998.0825 141998.0835 141998.0855 141998.0851 141998.0825 141998.0825 141998.0835 141998.0855 141998.0825 141998.0851 141998.0855 141998.0855 141998.0855 141998.0825 141998.0851 141998.0856 141998.0851 141998.0855 141998.0435 141998.0485 141998.0485 141998.0851 141998.0856 141998.0851 141998.0856 141998.0851 141998.0851 141998.0851 <t< th=""><th>14</th><th>1900.0305</th><th>141900.0501</th><th>141966.63</th><th>907 18,8366 </th><th></th><th></th><th></th><th>141900.0420</th><th>141900.0002</th><th>141900.0355</th><th>141900.0387</th></t<>	14	1900.0305	141900.0501	141966.63	907 18,8366				141900.0420	141900.0002	141900.0355	141900.0387
141980.8399 141900.0201 141900.0223 141900.0223 141900.0255 141900.0332 141900.0279 141900.0279 141900.0251 141900.0251 141900.0251 141900.0510 141900.0511 141900.0551 141900.0551 141900.0551 141900.0551 141900.0556 141900.0551 141900.0511 141900.0511 141900.0511 141900.0511 141900.0556 141900.0551 141900.0511 141900.0511 141900.0511 141900.0511 141900.0556 141900.0551 141900.0511 141900.0511 141900.0511 141900.0511 141900.0556 141900.0511 141900.0511 141900.0511 141900.0511 141900.0511 141900.0556 141900.0511 141900.0511 141900.0511 141900.0511 141900.0511 141900.0551 141900.0515 141900.0515 141900.0515 141900.0515 141900.0515 141900.0517 141900.0513 141900.0513 141900.0515 141900.0515 141900.0515 141900.0515 141900.0526 141900.0221 141900.0513 141900.0513 141900.0515 141900.0526 141900.0515 141900.0526	14	1900.0417	141900.0442		141900.0368				141900.0042	141900.0098	141900.0357	
141900.0332 141900.0279 141900.0467 141900.0467 141900.0467 141900.0279 141900.021 141900.0467 141900.0203 141900.0350 141900.0651 141900.0251 141900.0251 141900.0251 141900.0251 141900.0279 141900.0256 141900.0251 141900.0252 141900.0272 141900.0332 141900.0279 141900.0256 141900.0256 141900.0272 141900.0345 141900.0272 141900.0345 141900.0279 141900.0256 141900.0256 141900.0256 141900.0256 141900.0256 141900.0256 141900.0256 141900.0256 141900.0256 141900.0256 141900.0256 141900.0256 141900.0256 141900.0256 141900.0256 141900.0256 141900.0256 141900.0256 141900.0256 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055 141900.0056 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055 1419	14	1900.0399	141900.0007						141900.0427	141900.0223	141900.0358	
141900.02/3 141900.02/3 141900.0501 141900.0501 141900.0051 141900.051 141900.0531 141900.0531 141900.0051 141900.0235 141900.0231 141900.0350 141900.0055 141900.0235 141900.0235 141900.0340 141900.0056 141900.0235 141900.0272 141900.0347 141900.0235 141900.0272 141900.0347 141900.0235 141900.0272 141900.0347 141900.0235 141900.0272 141900.0347 141900.0235 141900.0272 141900.0347 141900.0235 141900.0272 141900.0347 141900.0235 141900.0272 141900.0347 141900.0256 141900.0256 141900.0408 141900.0251 141900.028 141900.0256 141900.0175 141900.0251 141900.0255 141900.0022 141900.0255 141900.0255 141900.0022 141900.0015 141900.0255 141900.0022 141900.0025 141900.0255 141900.0022 141900.0015 141900.0255 141900.0022 141900.0015 141900.0015	14:	1900.0382	141900.0200						141900.0467	141900.0173	141900.0359	
141900.0510 141900.051 141900.051 141900.051 141900.0051 141900.0519 141900.0512 141900.0325 141900.0051 141900.0519 141900.0212 141900.0345 141900.0251 141900.0212 141900.0345 141900.0212 141900.0345 141900.0257 141900.0212 141900.0345 141900.0212 141900.0345 141900.0257 141900.0257 141900.0272 141900.0345 141900.0325 141900.0205 141900.0205 141900.0235 141900.0235 141900.0235 141900.0175 141900.0251 141900.0255 141900.0255 141900.0255 141900.0221 141900.0205 141900.0205 141900.0255 141900.0265 141900.0057 141900.0057 141900.0055 141900.0055 141900.0055 141900.0055 141900.0057 141900.0057 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055 141900.0055			141900.0279						141900.0201	141900.0401	141900.0349	
141900.0051 141900.0519 141900.0212 141900.0346 141900.0056 141900.0272 141900.0347 141900.0279 141900.0236 141900.0272 141900.0347 141900.0205 141900.0236 141900.0272 141900.0347 141900.0205 141900.0236 141900.0238 141900.0347 141900.0205 141900.0238 141900.0248 141900.0248 141900.0101 141900.0259 141900.0248 141900.0556 141900.0469 141900.0251 141900.0551 141900.0556 141900.0475 141900.0252 141900.0551 141900.0268 141900.0622 141900.0551 141900.0268 141900.0268 141900.0622 141900.0628 141900.0265 141900.0268 141900.0622 141900.0268 141900.0268 141900.0268 141900.0628 141900.0283 141900.0268 141900.0248 141900.0027 141900.0028 141900.0015 141900.0248 141900.0028 141900.0015 141900.0024 141900.0024 141900.0028 141900.0035 141900.0025 141900.0025 141900.002			141900.0510						141900.0436	141900.053	141900.0392	
141900.0056 141900.0272 141900.0272 141900.0271 141900.0279 141900.0272 141900.0272 141900.0271 141900.0279 141900.0236 141900.0272 141900.0272 141900.0205 141900.0236 141900.0236 141900.0236 141900.0259 141900.0239 141900.0236 141900.0256 141900.0232 141900.0231 141900.0253 141900.0256 141900.0232 141900.0231 141900.0253 141900.0256 141900.0232 141900.0231 141900.0253 141900.0256 141900.0257 141900.0253 141900.0256 141900.0256 141900.0052 141900.0253 141900.0256 141900.0256 141900.0052 141900.0253 141900.0256 141900.0256 141900.0057 141900.0253 141900.0256 141900.0256 141900.0057 141900.0055 141900.0255 141900.0255 141900.0055 141900.0055 141900.0055 141900.0054 141900.0055 141900.0055 141900.0056 141900.0056 141900.0055 141900.0055 141900.0056 141900.0056			141900.0051						141900.0519	141900.021	2 141900.0346	i
141900.0279 141900.021 141900.0408 141900.014 141900.0305 141900.0145 141900.0208 141900.015 141900.0489 141900.025 141900.0268 141900.025 141900.0489 141900.025 141900.025 141900.025 141900.0175 141900.025 141900.025 141900.0268 141900.0022 141900.0225 141900.025 141900.0268 141900.0057 141900.0225 141900.0265 141900.0268 141900.0057 141900.0228 141900.0225 141900.0265 141900.0057 141900.0228 141900.0285 141900.0242 141900.0057 141900.0288 141900.0295 141900.0242 141900.0057 141900.0295 141900.0295 141900.0295 141900.0028 141900.0091 141900.0295 141900.0091 141900.0028 141900.0091 141900.0094 141900.0091 141900.0028 141900.0095 141900.0094 141900.0094 141900.0028 141900.0095 141900.0094 141900.0094 141900.0028 141900.0095 141900.0094 141900.0094			141900.0056						141900.0236	141900.027	2 141900.0347	l
141988,8385 141988,8385 141988,8385 141988,8355 141988,8385 141988,8455 141988,8455 141988,8555 141988,8385 141988,8455 141988,8555 141988,8555 141988,8385 141988,8555 141988,8555 141988,8655 141988,8355 141988,8555 141988,8555 141988,8655 141988,8555 141988,8555 141988,8555 141988,8655 141988,8555 141988,8555 141988,8142 141988,8142 141988,8555 141988,8555 141988,8142 141988,8142			141900.0279						141900.0021	141900.040	8 141900.0144	
141900,0439 141900,0251 141900,01051 141900,0203 141900,0175 141900,021 141900,0251 141900,0203 141900,022 141900,0226 141900,02051 141900,0205 141900,0022 141900,0226 141900,04021 141900,0402 141900,0027 141900,02051 141900,0402 141900,0402 141900,0027 141900,02051 141900,0402 141900,0402 141900,0027 141900,02051 141900,0402 141900,0402 141900,0027 141900,02051 141900,04051 141900,0402 141900,0028 141900,0028 141900,04051 141900,00342 141900,0028 141900,00295 141900,00395 141900,00342 141900,0028 141900,00395 141900,00395 141900,00342 141900,0038 141900,00395 141900,00395 141900,00395 141900,0045 141900,00395 141900,00395 141900,0149 141900,0178 141900,0178 141900,0179 141900,0149 141900,0178 141900,0178 141900,0179 141900,0149 141900,0178 141900,0178 141900,0179 141900,016			141900.0305						141900.0145	141900.029	8 141900.0161 8 141000 0506	
141900.0175 141900.0006 141900.0513 141900.0268 141900.022 141900.0255 141900.0265 141900.0265 141900.022 141900.0251 141900.0265 141900.0402 141900.022 141900.0251 141900.0421 141900.0420 141900.0057 141900.025 141900.0051 141900.0420 141900.0051 141900.0055 141900.0055 141900.0054 141900.0065 141900.0055 141900.0055 141900.0054 141900.0055 141900.0055 141900.0055 141900.0054 141900.0055 141900.0055 141900.0054 141900.0054 141900.0055 141900.0055 141900.0055 141900.0054 141900.0055 141900.0055 141900.0055 141900.0054 141900.0155 141900.0170 141900.0170 141900.0164 141900.0515 141900.0176 141900.0170 141900.0165 141900.0176 141900.0176 141900.0170 141900.0165			141900.04891						141900.0225	141900.015	51 141900.0093	
141900.0232 141900.0231 141900.0205 141900.0402 141900.0022 141900.0024 141900.0421 141900.0430 141900.0027 141900.0028 141900.0421 141900.0430 141900.0028 141900.0028 141900.0015 141900.0430 141900.0031 141900.0015 141900.0045 141900.0044 141900.0038 141900.0035 141900.0034 141900.0034 141900.0038 141900.0035 141900.0034 141900.0034 141900.0045 141900.0035 141900.0034 141900.0034 141900.0055 141900.0035 141900.0034 141900.0034 141900.0155 141900.0156 141900.0170 141900.0194 141900.0176 141900.0170 141900.0170 141900.0144 141900.0176 141900.0170 141900.0165 141900.0170			141900.0175						141900.0006	141900.05	3 141900.0268	
141960.0622 141960.6242 141900.6421 141900.6480 141960.0657 141900.6281 141900.6051 141900.6042 141960.0628 141900.6055 141900.6042 141900.6044 141900.0608 141900.6045 141900.6094 141900.6094 141900.0486 141900.6045 141900.6035 141900.6094 141900.0486 141900.6045 141900.60395 141900.60295 141900.0486 141900.6045 141900.6030 141900.6049 141900.0176 141900.6045 141900.6030 141900.6049			141900.0232						141900.0031	141900.02	05 141900.0402	2
141900.0057 141900.0028 141900.0015 141900.0042 141900.0021 141900.0028 141900.0091 141900.0094 141900.0086 141900.0045 141900.0091 141900.0094 141900.00486 141900.0045 141900.0091 141900.0094 141900.00486 141900.0045 141900.0094 141900.0094 141900.0155 141900.0045 141900.0094 141900.0094 141900.0155 141900.0045 141900.0094 141900.0094 141900.0155 141900.0045 141900.0094 141900.0094 141900.0178 141900.0045 141900.0045 141900.0046			141900.0022						141900.0242	141900.04	21 141900.048	9]
141968,0211 141968,0095 141968,0091 141968,0091 141968,0091 141968,0086 141968,0195 141968,0091 141968,0091 141968,0091 141968,0195 141968,0195 141968,0091 141968,0091 141968,0091 141968,0195 141968,0195 141968,0091 141968,0091 141968,0091 141968,0195 141968,0195 141968,0195 141968,0191 141968,0191 141968,0176 141968,0176 141968,0176 141968,0194 141968,0194			141900.0057						141900.0028	141900.00	15 141900.0342	21
141900,0455 141900,055 141900,055 141900,055 141900,055 141900,055 141900,0149 141900,0515 141900,0176 141900,0176 141900,0176 141900,0176 141900,0164 141900,055 141			141900.0211						141900,0195	141900.00	91 141900.009	+
141988.8515 141988.0114 141988.0178 141988.0104 141988.0178			141900.0488						141900,0046	141900.00	80 141900.014	al l
141988,8178 141988,8142 141988,8465			141900.0515						141900.0114	141900.01	70 141900.010	4
			141988.8178						141900.0230	141900.01	42 141900.046	<u>61 </u>

The Ensembl Genome Browser

• To many users, it is comparable in scope and importance to the UCSC Genome Browser, and it is often useful for new users to visit both sites.



			Login/Regist
	e Ensembles _E	BLAST/BLAT BioMart Tools Downl	loads Help & Documentation More 🔻 🎆 👻 Search Human
ŀ	Human (GRCh38.p7) 🔻 Lo	ocation: 11:5,225,464-5,229,395 Gene:	нвв
G	ene-based displays		
İ	Summary	Gene: HBB ENSG000002	44734
	 Splice variants 		
	 Transcript comparison 	Description	hemoglobin subunit beta [Source:HGNC Symbol;Acc: <u>HGNC:4827</u> &]
	└─ Gene alleles	Synonyms	heta-alohin HBD CD113t-C
	Sequence	Cynonyms	bela-globin, Hbb, CD1131-C
		Location	Chromosome 11: 5,225,464-5,229,395 reverse strand.
Ξ			GRCh38:CM000673.2
		About this gene	This gene has 5 transcripts (splice variants) 136 orthologues, 9 paralogues, is a member of
	- Gene gain/loss tree	30.00	1 Ensembl protein family and is associated with 28 phenotypes.
	- Orthologues	Troposinto	
	- Paralogues	Transcripts	Show transcript table
	Ensembl protein families		
þ	Ontologies	Cump man un A	
	GO: Cellular component	Summary 🕼	
	 GO: Biological process 		
	GO: Molecular function	Name	HBB & (HGNC Symbol)
Ľ	- Phenotypes	CCDS	This gene is a member of the Human CCDS set: CCDS7753.1 🗗
닏	Genetic Variation	UniProtKB	This was been wetains that companyed to the following UniDentificant Decord a
	variant table	UNIFICIND	I his dene has proteins that correspond to the following UniProtkB identifiers: P688/1 A

Havana gene. OT THOMADOODOODOO

Go to Region in Detail for more tracks and navigation options (e.g. zooming)

🌣 🛃 < 🖽 🖬 🎨 🕏

Drag/Select: \leftrightarrow 🏢



(i) https://www.ncbi.nlm.nih.gov/mapview/

Home GenBank BLAST

Map Viewer Home

NCBI

->

4

The Map Viewer provides a wide variety of genome mapping and seq

C

Q

Search		
<u>S</u> earch	Homo sapiens	*
for:	hbb	
		Go
Tools L	egend	
0 B C B G	Search or Browse the Genome BLAST Clone Finder Go to region on a chromosome Genome Resources page	
News		

Vertebrates	
Mammals	
Primates	
Scientific name	Common name
Callithrix jacchus	white-tufted-ear marmoset
Chlorocebus sabaeus	green monkey
Gorilla gorilla	western gorilla
Homo sapiens	human
Macaca fascicularis	crab-eating macaque
Macaca mulatta	rhesus macaque
Nomascus leucogenys	northern white-cheeked gibbon
Otolemur garnettii	small-eared galago
Pan paniscus	pygmy chimpanzee







LARGE-SCALE QUERIES OF REGIONS AND FEATURES

- In many cases we are interested in a single gene.
- In many other cases we want to know about large collections of genes, proteins, or indeed any other element.
- Example:
 - What is the complete set of human globin genes?
 - To which chromosomes are they assigned?
 - How many exons are on chromosome 11, and how many repeat elements occur in each exon?

TABLE 2.10 File formats for custom tracks used at Ensembl and/or UCSC. Two definitions of GTF (from Ensembl and UCSC) are given.

File Format	Definition	Typical file size
BAM		Any size; often millionsof rows
BED	Browser extensible data	Any size; often dozens to thousands or millions of rows
BedGraph		Any size
bigBed		
GFF/GTF	General feature format, General transfer format Gene transfer format	Any size
MAF		
PSL		Any size
WIG	Wiggle	Any size
BAM	Binary alignment/map	Very large
BigWig		Very large
VCF	Variant call format	Very large

LARGE-SCALE QUERIES OF REGIONS AND FEATURES

Collect information one gene at a time?
 tedious, inefficient, and error-prone

• There are many bioinformatics tools that allow us to collect genome-wide information.

• UCSC Table Browser

Q Search ☆自 https://genome.ucsc.edu C UNIVERSITY OF CALIFORNIA **Genome Browser** Genomes **Genome Browser** Tools Mirrors Downloads My Data Help About Us **Our tools** Genome Browser interactively visualize genomic data BLAT rapidly align sequences to the genome Table Browser download data from the Genome Browser database Variant Annotation Integrator get functional effect predictions for variant calls
Conclusion

- Bioinformatics is an emerging field whose defining feature is the accumulation of biological information in databases.
- The three major traditional DNA databases GenBank, EMBL-Bank, and DDBJ – are adding several million new sequences each year as well as billions of nucleotides.
- Next-generation sequencing technology is producing vastly greater amounts of DNA.

• NCBI's sequence databases accept genome data from sequencing projects from around the world and serve as the cornerstone of bioinformatics research.

GenBank:

 An annotated collection of all publicly available nucleotide and amino acid sequences.

EST database:

• A collection of expressed sequence tags, or short, single-pass sequence reads from mRNA (cDNA).

GSS database:

A database of genome survey sequences, or short, single-pass genomic sequences.

HomoloGene:

 A gene homology tool that compares nucleotide sequences between pairs of organisms in order to identify putative orthologs.

HTG database:

• A collection of high-throughput genome sequences from large-scale genome sequencing centers, including unfinished and finished sequences.

SNPs database:

• A central repository for both single-base nucleotide substitutions and short deletion and insertion polymorphisms.

<u>RefSeq:</u>

 A database of non-redundant reference sequences standards, including genomic DNA contigs, mRNAs, and proteins for known genes. Multiple collaborations, both within NCBI and with external groups, supports data-gathering efforts.

STS database:

• A database of sequence tagged sites, or short sequences that are operationally unique in the genome.

<u>UniSTS:</u>

• A unified, non-redundant view of sequence tagged sites (STSs).

UniGene:

• A collection of ESTs and full-length mRNA sequences organized into clusters, each representing a unique known or putative human gene annotated with mapping and expression information and cross-references to other sources.

UniGene computationally identifies transcripts from the same locus; analyzes expression by tissue, age, and health status; and reports related proteins (protEST) and clone resources.

Single Nucleotide Polymorphism (SNP) database

- The SNP Database (also known as dbSNP) is an archive for genetic variation within and across different species developed and hosted by NCBI in collaboration with the <u>National Human Genome Research</u> Institute (NHGRI).
 - Polymorphism in biology occurs when two or more clearly different phenotypes exist in the same population of a species: related to <u>biodiversity</u>, <u>genetic</u> <u>variation</u> and <u>adaptation</u>
- The dbSNP accepts apparently neutral polymorphisms, polymorphisms corresponding to known phenotypes, and regions of no variation.
- It was created in September 1998 to supplement GenBank (NCBI's nucleic acid and protein sequences)

Single Nucleotide Polymorphism (SNP) database

- Its goal is to act as a single database that contains all identified genetic variation, which can be used to investigate a wide variety of genetically based natural phenomenon. Specifically, access to the molecular variation cataloged within dbSNP aids basic research such as physical mapping, population genetics, investigations into evolutionary relationships, as well as being able to quickly and easily quantify the amount of variation at a given site of interest.
- Applied research, genetic engineering, drug discovery, etc.

SNCBI Resources 🖸 How To 🖸

NCBI All Databases V Search Conserved Domains National Center for dbGaP Biotechnology Information dbVar Epigenomics b NCBI NCBI Home **Popular Resources** EST Resource List (A-Z) Gene PubMed ter for Biotechnology Information advances science and health by providing access to biomedical Genome rmation. Bookshelf All Resources GEO DataSets GEO Profiles PubMed Central Chemicals & Bioassays II | Mission | Organization | Research | NCBI News GSS PubMed Health Data & Software GTR HomoloGene BLAST DNA & RNA MedGen Nucleotide MeSH Domains & Structures alyze data using NCBI software NCBI Web Site Genome Genes & Expression : Get NCBI data or software NLM Catalog SNP Learn how to accomplish specific tasks at NCBI Nucleotide Genetics & Medicine OMIM ons: Submit data to GenBank or other NCBI databases Gene Genomes & Maps PMC Protein PopSet Homology Drohe PubChem Literature NCBI Twitter feed Proteins NCBI Announcements Keep up-to-date on data updates, resource Sequence Analysis announcements, and other information about NCBI's next webinar is The Statistics of Taxonomy what is going on at the NCBI. GO Local Pairwise Sequence Alignment, Training & Tutorials Parts 1 and 2 Jan 13, 201 2 3 4 5 6 7 8 Variation 11 1 On Thursday, January 22nd, Stanbar

Sign in to NCB



NCBI abbreviations for nucleotide

Category	Description
NC_#######	DNA
NM_#######	mRNA
XM_#######	Predicted mRNA

NCBI abbreviations for protein

Category	Description	
NP_#######	Protein	
XP_#######	Predicted Protein	

Predicted & Known Proteins

• Predicted proteins are predicted by gene identification tools (software) and they share sequence similarity to wellcharacterized proteins.

• Known proteins are experimentally proven to exist and are linked to a known gene.

Information Hyperlinked Over Proteins (iHOP)

- iHOP provides fast, accurate, comprehensive, and up-to-date summary information on more than 80,000 biological molecules by automatically extracting key sentences from millions of PubMed documents.
- It is an online text-mining service that provides a gene-guided network to access PubMed abstracts.
- iHOP (Information Hyperlinked over Proteins) allows researchers to explore a network of gene and protein interactions based on published scientific literature. For each gene search, iHOP reports sentences from abstracts associating it with other genes, links out to full abstracts, and reports experimental evidence for the interactions, if available. You can also select sentences to create and visualize your own gene model

http://www.ihop-net.org/UniPub/iHOP/



Symbol	Name	Synonym/ DB-reference	o Organism	Results
		Life cycles of suc	cessful genes rends in Genetics	₽
WWC1	WW and C2 domain containing 1	KIBRA	Homo sapiens
Wwc1	WW, C2 and coiled-coil domain containing 1	KIBRA	Mus musculus	e e 🛛 i

Figure 2.7a Practical Bioinformatics (© Garland Science 2013)

Constant No.			S	Ormaniam	
WWC1 WW	and C2 domain containing 1		FLJ10865, FLJ23369, HBeAg- protein 3, HBEBP3, HBEBP36 KIBRA, Kidney and brain prote KIBRA, WW domain-containing	binding Homo sapiens , KIAA0869, bin, Protein g protein 1	
WikiGenes UniProt IntAct PDB Structure OMIM NCBI Gene NCBI RefSeq NCBI RefSeq NCBI UniGene NCBI Accession	edit this page new Q8IX03, Q7Z4G8, Q6MZX4 Q8IX03 2Z0U 610533 23286 NP_001155133, NP_001155134 NM_001161662, NM_001161661 23286 DR001014, AK296323	more than 2,800 organ alwa	isms, 110,000 genes, 23 ays up to date – every d	.4 million sentences. lay.	
Homologues of	WWC1				
Interaction info	rmation for WWC1 🕎				
Most recent inf	ormation for WWC1 🔯				
Enhanced Publ	Med/Google query				
WARNING: Please ke	eep in mind that gene detection is done automa	tically and can exhibit a certain error. Read	more about synonym ambiguity and the iH0	OP confidence value 삶삶삶.	
				Find in this Page	
Sentences in th symbol 🗐 - Re For a summary	nis view contain definitions for W ad more. overview of the information in the	WC1 - Definitions are available	whenever you see this	Show all Order by relevance	
We also found that KIBRA 2-DLC1 2 interaction is mandatory for the recruitment and transactivation functions of ER 2 or DLC1 2 to E					
Finally we found that KIBRA interacts with histone H3 via its glutamic acid [?]-rich region and that such interaction might play a mechanistic role in conferring an optimal ER interactivation function as well as the proliferation of ligand-stimulated breast cancer cells. [2006]					
KIAA0513 🏠 int and cytoskeletal	t eracts with KIBRA ☆, <u>HAX-1</u> ☆, a regulation. [2006]	nd INTS4 (a), which also interact	with proteins involved in <u>neurop</u>	lasticity, apoptosis, 📓 土	

Figure 2.7b Practical Bioinformatics (© Garland Science 2013)

Conclusion

- Many databases and resources are available, some as websites and some (such as R packages or NCBI E-Utilities) via programing languages.
- There is no single correct way to find information; many approaches are possible.
- Resources are closely interrelated, providing links between the databases.
- Bioinformatics databases are evolving extremely rapidly.
- Each January, the first issue of the journal Nucleic Acids Research includes nearly 100 brief articles on databases.