

# Genome Informatics

Systems Biology and the Omics Cascade (Course 2143)  
Day 3, June 11<sup>th</sup>, 2008

Kiyoko F. Aoki-Kinoshita

# Introduction

- Genome informatics covers the computer-based modeling and data processing of genome-based data.
- This includes databases and resources for genomic analysis.
- You were introduced to KEGG on Day 2.
- Some other useful databases and resources will be covered today.

# But first! Data formats

- It is usually not enough to simply look at the data provided by databases
- To actually use the data for analysis, one often needs to save the retrieved data
- This requires knowledge about the data formats used by each database
- So we will cover the major data formats used in bioinformatics

# Data Formats

## ■ Major data formats:

- GenBank
- EMBL/UniProt
- FASTA
- PDB

## ■ Formats suited for programming:

- ASN.1 (Abstract Syntax Notation One)
- XML (eXtensible Markup Language)

# GenBank format

- Each line starts with a keyword in capital letters.
- Each keyword is followed by a tab and the information corresponding to it.
- Some keywords are hierarchical:

```
LOCUS      SCU49845      5028 bp      DNA      PLN      21-JUN-1999
DEFINITION Saccharomyces cerevisiae TCP1-beta gene, partial cds, and Axl2p
            (AXL2) and Rev7p (REV7) genes, complete cds.
ACCESSION  U49845
VERSION    U49845.1  GI:1293613
KEYWORDS   .
SOURCE     Saccharomyces cerevisiae (baker's yeast)
            ORGANISM  Saccharomyces cerevisiae
            Eukaryota; Fungi; Ascomycota; Saccharomycotina; Saccharomycetes;
            Saccharomycetales; Saccharomycetaceae; Saccharomyces.
REFERENCE  1  (bases 1 to 5028)
            AUTHORS   Torpey,L.E., Gibbs,P.E., Nelson,J. and Lawrence,C.W.
            TITLE     Cloning and sequence of REV7, a gene whose function is required for
            DNA damage-induced mutagenesis in Saccharomyces cerevisiae
            JOURNAL   Yeast 10 (11), 1503-1509 (1994)
            PUBMED    7871890
REFERENCE  2  (bases 1 to 5028)
            AUTHORS   Roemer,T., Madden,K., Chang,J. and Snyder,M.
            TITLE     Selection of axial growth sites in yeast requires Axl2p, a novel
            plasma membrane glycoprotein
            JOURNAL   Genes Dev. 10 (7), 777-793 (1996)
            PUBMED    8846915
REFERENCE  3  (bases 1 to 5028)
            AUTHORS   Roemer,T.
            TITLE     Direct Submission
            JOURNAL   Submitted (22-FEB-1996) Terry Roemer, Biology, Yale University, New
            Haven, CT, USA
FEATURES   Location/Qualifiers
            source     1..5028
                    /organism="Saccharomyces cerevisiae"
                    /db_xref="taxon:4932"
                    /chromosome="IX"
                    /map="9"
            CDS        <1..206
                    /codon_start=3
                    /product="TCP1-beta"
                    /protein_id="AAA98665.1"
                    /db_xref="GI:1293614"
                    /translation="SSIYNGISTSGLDLNNGTIADMRQLGIVESYKLRKRAVSSASEA
                    AEVLLRVDNIIIRARPRTANRQHM"
            gene       687..3158
                    /gene="AXL2"
            CDS        687..3158
                    /gene="AXL2"
```

# EMBL format

- Similar to GenBank, except that keywords are two-letter IDs.
- UniProt's format is similar to this format.

```
ID SC10H5 standard; DNA; PRO; 4870 BP.
XX
AC AL031232;
XX
DE Streptomyces coelicolor cosmid 10H5.
XX
OS Streptomyces coelicolor
OC Eubacteria; Firmicutes; Actinomycetes; Streptomycetes;
OC Streptomycetaceae; Streptomyces.
XX
RN [1]
RP 1-4870
RA Redenbach M., Kieser H.M., Denapaite D., Eichner A.,
RA Cullum J., Kinashi H., Hopwood D.A.;
RT "A set of ordered cosmids and a detailed genetic and physical
RT map for the 8 Mb Streptomyces coelicolor A3(2) chromosome.";
RL Mol. Microbiol. 21(1):77-96(1996).
XX
CC Notes:
CC
CC Streptomyces coelicolor sequencing at The Sanger Centre is funded
CC by the BBSRC.
XX
FH Key Location/Qualifiers
FH
FT source 1..4870
FT /organism="Streptomyces coelicolor"
FT /strain="A3(2)"
FT /clone="cosmid 10H5"
FT CDS complement(350..805)
FT /note="SC10H5.02c, probable integral membrane protein, len:
FT 151 aa; similar to S. coelicolor hypothetical protein
FT TR:054194 (EMBL:AL021411) SC7H1.35 (155 aa), fasta scores;
FT opt: 431 z-score: 749.8 E(): 0, 53.5% identity in 114 aa
FT overlap."
FT /product="putative integral membrane protein"
FT /gene="SC10H5.02c"
FT RBS complement(812..815)
FT /note="possible RBS upstream of SC10H5.02c"
FT CDS complement(837..1301)
FT /note="SC10H5.03c, probable integral membrane protein, len:
FT 154 aa"
FT /product="putative integral membrane protein"
FT /gene="SC10H5.03c"
FT RBS complement(1308..1312)
FT /note="possible RBS upstream of SC10H5.03c"
```

# FASTA sequence format

- > Randseq1 first randomly generated seq  
GGTGGTTACTAACCGTAAGAGATGATGTCGCCGTGGTTCGCGTGGC  
GCCGCGGACCCAGATTGTACTTCTCTGAGTCGTTCTAGATCGACC  
AGTCTTCTAGCTTGCCCGTGAGGTATGGGG  
AGCCGCATATTGCCCAACAAT
- > Randseq2 second randomly generated seq  
GCGACGCGTCTCTACACCAGACGCTTCTGTTGAGGAAGAGTGCCT  
GAGTGCAGGTCCTCGAGAACCCTGGAAGTTGAAGGGCGCGTCT  
CACTGGTCGTGAGAAGGCTCCGTCGATACG  
AAAGTCCATGCCAAGGACAT
- > Randseq3 third randomly generated seq  
GGCGAGTCTGAACTCACAAATATTGCACGAGAGTTTAGTGTATGT  
TCCTCTTAGGCTGATAACAATAGTTTAGTGAGCGGAAATGCAACC  
GCGAGGCGGTCCCCTGCGCTTGTAATGGCC  
ACCTGTTGCCCGTCGGATAT



# Nucleic acid code for FASTA

- A → adenosine
- C → cytidine
- G → guanine
- T → thymidine
- U → uridine
- R → G A (purine)
- Y → T C (pyrimidine)
- K → G T (keto)
- M → A C (amino)
- S → G C (strong)
- W → A T (weak)
- B → G T C
- D → G A T
- H → A C T
- V → G C A
- N → A G C T (any)
- - → gap of indeterminate length



# Amino acid code for FASTA

- A alanine
- B aspartate or asparagine
- C cystine
- D aspartate
- E glutamate
- F phenylalanine
- G glycine
- H histidine
- I isoleucine
- K lysine
- L leucine
- M methionine
- N asparagine
- P proline
- Q glutamine
- R arginine
- S serine
- T threonine
- U selenocysteine
- V valine
- W tryptophan
- Y tyrosine
- Z glutamate or glutamine
- X any
- \* translation stop
- - gap of indeterminate length

# PDB format

- Similar to GenBank, using different keywords
- Includes 3-dimensional coordinates of amino acids

```

HEADER  HYDROLASE(ACID PROTEINASE ZYMOGEN)      03-SEP-91   3PSG   3PSG   2
COMPND  PEPINOGEN                                3PSG   3
SOURCE  PORCINE (SUS ¥$SCROFA)                   3PSG   4
AUTHOR  J.A.HARTSUCK,G.KOELSCH,S.J.REMINGTON    3PSG   5
REVDAT  1 15-JAN-93 3PSG   0                    3PSG   6
SPRSDE  15-JAN-93 3PSG   1PSG                   3PSG   7
JRNL    AUTH  J.A.HARTSUCK,G.KOELSCH,S.J.REMINGTON 3PSG   8
JRNL    TITL  THE HIGH RESOLUTION CRYSTAL STRUCTURE OF PORCINE 3PSG   9
JRNL    TITL 2 PEPINOGEN                                3PSG  10
JRNL    REF  PROTEINS.STRUCT.,FUNCT.,          V. 13      1 1992 3PSG  11
JRNL    REF  2 GENET.                                3PSG  12
JRNL    REFN  ASTM PSFGEY US ISSN 0887-3585                867 3PSG  13
REMARK  1                                           3PSG  14
REMARK  2                                           3PSG  15
REMARK  2 RESOLUTION. 1.65 ANGSTROMS.                3PSG  16
REMARK  3                                           3PSG  17
REMARK  3 REFINEMENT.                                3PSG  18
REMARK  3 PROGRAM                                TNT                3PSG  19
REMARK  3 AUTHORS                                TRONRUD,TEN EYCK & MATHEWS 3PSG  20
REMARK  3 R VALUE                                0.173                  3PSG  21
REMARK  3 RMSD BOND DISTANCES                    0.014 ANGSTROMS       3PSG  22
REMARK  3 RMSD BOND ANGLES                       2.60 DEGREES          3PSG  23
SEQRES  1 370 LEU VAL LYS VAL PRO LEU VAL ARG LYS LYS SER LEU ARG 3PSG  76
SEQRES  2 370 GLN ASN LEU ILE LYS ASP GLY LYS LEU LYS ASP PHE LEU 3PSG  77
SEQRES  3 370 LYS THR HIS LYS HIS ASN PRO ALA SER LYS TYR PHE PRO 3PSG  78
SEQRES  4 370 GLU ALA ALA ALA LEU ILE GLY ASP GLU PRO LEU GLU ASN 3PSG  79
SEQRES  5 370 TYR LEU ASP THR GLU TYR PHE GLY THR ILE GLY ILE GLY 3PSG  80
SEQRES  6 370 THR PRO ALA GLN ASP PHE THR VAL ILE PHE ASP THR GLY 3PSG  81
FTNOTE  1                                           3PSG 105
FTNOTE  1 RESIDUE PRO 23 IS A CIS PROLINE.            3PSG 106
FORMUL  2 HOH *180(H2 O1)                            3PSG 107
CRYST1  106.100  43.700  88.900  90.00  91.40  90.00 C 2          4 3PSG 108
ORIGX1  1.000000  0.000000  0.000000                0.000000                3PSG 109
ORIGX2  0.000000  1.000000  0.000000                0.000000                3PSG 110
ORIGX3  0.000000  0.000000  1.000000                0.000000                3PSG 111
SCALE1  0.009425  0.000000  0.000230                0.000000                3PSG 112
SCALE2  0.000000  0.022883  0.000000                0.000000                3PSG 113
SCALE3  0.000000  0.000000  0.011252                0.000000                3PSG 114
ATOM    1  N   LEU      1P   57.364 -9.595  2.554  1.00  21.58  3PSG 115
ATOM    2  CA  LEU      1P   56.687 -8.586  3.371  1.00  22.68  3PSG 116
ATOM    3  C   LEU      1P   57.052 -8.758  4.847  1.00  16.32  3PSG 117
ATOM    4  O   LEU      1P   57.270 -9.875  5.293  1.00  26.09  3PSG 118
ATOM    5  CB  LEU      1P   55.129 -8.641  3.226  1.00  36.56  3PSG 119
ATOM    6  CG  LEU      1P   54.503 -8.303  1.850  1.00  32.34  3PSG 120
ATOM    7  CD1 LEU      1P   52.977 -8.461  1.892  1.00  37.73  3PSG 121
ATOM    8  CD2 LEU      1P   54.798 -6.861  1.453  1.00  30.91  3PSG 122
ATOM    9  N   VAL      2P   57.149 -7.672  5.627  1.00  17.11  3PSG 123
ATOM   10  CA  VAL      2P   57.472 -7.872  7.047  1.00  15.32  3PSG 124

```

# ASN.1 format

```
{
  name      {givenName "John", initial "P", familyName "Smith"},
  title     "Director",
  number    51,
  dateOfHire "19710917",
  nameOfSpouse {givenName "Mary", initial "T", familyName "Smith"},
  children
  { {name {givenName "Ralph", initial "T", familyName "Smith"} ,
    dateOfBirth "19571111"},
    {name {givenName "Susan", initial "B", familyName "Jones"} ,
    dateOfBirth "19590717" }
  }
}
```

Name:	John P Smith
Title:	Director
Employee Number:	51
Date of Hire:	17 September 1971
Name of Spouse:	Mary T Smith
Number of Children:	2
Child Information	
Name:	Ralph T Smith
Date of Birth	11 November 1957
Child Information	
Name:	Susan B Jones
Date of Birth	17 July 1959

- Hierarchical data format
- Groups are delineated by curly brackets
- Data type names precede the brackets
- Data within a group are separated by commas

# XML format

```
person ::=
  <PersonnelRecord>
    <name>
      <givenName>John</givenName>
      <initial>P</initial>
      <familyName>Smith</familyName>
    </name>
    <title>Director</title>
    <number>51</number>
    <dateOfHire>19710917</dateOfHire>
    <nameOfSpouse>
      <givenName>Mary</givenName>
      <initial>T</initial>
      <familyName>Smith</familyName>
    </nameOfSpouse>
    <children>
      <ChildInformation>
        <name>
          <givenName>Ralph</givenName>
          <initial>T</initial>
          <familyName>Smith</familyName>
        </name>
        <dateOfBirth>19571111</dateOfBirth>
      </ChildInformation>
      <ChildInformation>
        <name>
          <givenName>Susan</givenName>
          <initial>B</initial>
          <familyName>Jones</familyName>
        </name>
        <dateOfBirth>19590717</dateOfBirth>
      </ChildInformation>
    </children>
  </PersonnelRecord>
```

Name:	John P Smith
Title:	Director
Employee Number:	51
Date of Hire:	17 September 1971
Name of Spouse:	Mary T Smith
Number of Children:	2
Child Information	
Name:	Ralph T Smith
Date of Birth	11 November 1957
Child Information	
Name:	Susan B Jones
Date of Birth	17 July 1959

- Hierarchical data format
- Tags define the type of data:
  - Opening tag: <name>
  - Closing tag: </name>
- Data are delineated by opening and closing tags

# Data format converter

## ■ READSEQ

– URL : <http://thr.cit.nih.gov/molbio/readseq/>

### WWW READSEQ Sequence Conversion

**Function:** Converts input DNA/AA sequence to specified format (Input format is determined automatically).

As of Feb. 3, 2006, Sequence Conversion uses the JAVA version of the READSEQ program (Readseq version 2.1.22 (02-May-2005), developed by [Dr. Don Gilbert](#), Information on READSEQ is maintained at the [IUBio Archive](#) site at University of Indiana.

**Format for the output** (Use this [form](#) for 'Pretty' format)

**Additional formatting options:**

Mixed Case  Lower Case  Upper Case

Remove Gap Symbols ('-')

Please enter or paste sequence[s] to be converted (most [formats](#) accepted):

Credits: WWW implementation by [BIMAS Staff](#)



# Types of databases

## ■ Data resources of multiple types of data

- EBI (European Bioinformatics Institute)
- NCBI (National Center for Biotechnology Information)
- KEGG (Kyoto Encyclopedia of Genes and Genomes)

## ■ Gene and protein information

- GenBank, UniProt, and PDB
- Species specific: FlyBase, dictyBase, etc.

## ■ Ontological data

- Gene Ontology

## ■ Pathway data

- KEGG PATHWAY, Reactome, BRENDA, etc.

## ■ Protein-protein interaction data

- IntAct, BioGRID, etc.

# EBI

- <http://www.ebi.ac.uk/>
- European base of molecular biology information, including genomic, gene expression, and literature information





### Data Resources & Tools

- EMBL-BANK
- UniProt
- ArrayExpress
- Ensembl
- InterPro
- PDB-EBI
- Genomes
- Nucleotide Sequences
- Protein Sequences
- Macromolecular Structures
- Small Molecules
- Gene Expression
- Molecular Interactions
- Literature
- Taxonomy
- Sequence Similarity & Analysis
- Pattern & Motif Searches

### About the EBI

- Research
- PhD Studies
- Training
- Industry Support
- Group & Team Leaders
- EBI Funders
- User Support
- EBI Missions
- People
- Events and News
- How to Find Us



Search for *HIV* in *All the EBI*

Expand all Collapse all

▶ Genomes	87	▶ Molecular Interactions	3
▶ Nucleotide Sequences	250,495	▶ Reactions & Pathways	36
▶ Protein Sequences	51,222	▶ Protein Families	41
▶ Macromolecular Structures	628	▶ Enzymes	2
▶ Small molecules	3	▶ Literature	183,676
▶ Gene Expression	23	▶ Ontologies	6
		▶ EBI Web Site	34

#### Refine your search:

Search for *HIV* in *All the EBI*

with the following keywords

Refine

# NCBI

- <http://www.ncbi.nlm.nih.gov/>
- Contains public databases of molecular biology information including genomes, microarray gene expression, protein sequence domains, etc.
- Develops software for analyzing genome data, including BLAST
- Provides PubMed, an archive of biomedical and life science journals



## SITE MAP

[Alphabetical List](#)  
[Resource Guide](#)

## About NCBI

An introduction to  
NCBI

## GenBank

Sequence  
submission support  
and software

## Literature databases

[PubMed](#), [OMIM](#),  
[Books](#), and [PubMed  
Central](#)

## Molecular databases

## ▶ What does NCBI do?

Established in 1988 as a national resource for molecular biology information, NCBI creates public databases, conducts research in computational biology, develops software tools for analyzing genome data, and disseminates biomedical information - all for the better understanding of molecular processes affecting human health and disease. [More about NCBI...](#)

## Protein Clusters

The new Protein Clusters database contains Reference Sequence (RefSeq) protein records that are grouped and annotated by sequence and functional similarity. Source sequences come from the complete genomes of prokaryotes, plasmids, and organelles.

## Hot Spots

- ▶ [Assembly Archive](#)
- ▶ [Clusters of orthologous groups](#)
- ▶ [Coffee Break, Genes & Disease, NCBI Handbook](#)
- ▶ [Electronic PCR](#)
- ▶ [Entrez Home](#)
- ▶ [Entrez Tools](#)
- ▶ [Gene expression omnibus \(GEO\)](#)
- ▶ [Human genome](#)

Search

All Databases

for HIV

Go

## SITE MAP

[Alphabetical Resource](#)

## About NCBI

[An introduction to NCBI](#)

## GenBank

[Sequence submission and software](#)

## Literature databases

[PubMed, Books, and PubMed Central](#)

## Molecular databases

[Genomes](#)

All Databases

NCBI Web Site

PubMed

Protein

Nucleotide

EST

GSS

Structure

Genome

Books

CancerChromosomes

Conserved Domains

3D Domains

Gene

Genome Project

dbGaP

GENSAT

GEO Profiles

GEO Datasets

## What does NCBI do?

Established in 1988 as a national resource for biology information, NCBI creates databases, conducts research in computational biology, develops software for analyzing genome data, and disseminates biomedical information - all for the understanding of molecular biology affecting human health and disease. [More about NCBI...](#)

## Clusters

The Protein Clusters database contains Reference Sequence (RefSeq) protein records that are grouped and annotated by sequence and functional similarity. Source sequences come from the complete genomes of prokaryotes, plasmids, and organelles.

## Hot Spots


























































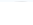
- ▶ Assembly Archive
- ▶ Clusters of orthologous groups
- ▶ Coffee Break, Genes & Disease, NCBI Handbook
- ▶ Electronic PCR
- ▶ Entrez Home
- ▶ Entrez Tools
- ▶ Gene expression omnibus (GEO)
- ▶ Human genome

Search across databases



[Help](#)

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

<b>184871</b>		<b>PubMed:</b> biomedical literature citations and abstracts	
<b>40435</b>		<b>PubMed Central:</b> free, full text journal articles	
<b>106</b>		<b>Site Search:</b> NCBI web and FTP sites	
<b>5015</b>		<b>Books:</b> online books	
<b>222</b>		<b>OMIM:</b> online Mendelian Inheritance in Man	
<b>none</b>		<b>OMIA:</b> online Mendelian Inheritance in Animals	
<b>263187</b>		<b>Nucleotide:</b> Core subset of nucleotide sequence records	
<b>429</b>		<b>EST:</b> Expressed Sequence Tag records	
<b>174358</b>		<b>GSS:</b> Genome Survey Sequence records	
<b>273432</b>		<b>Protein:</b> sequence database	
<b>20</b>		<b>Genome:</b> whole genome sequences	
<b>1013</b>		<b>Structure:</b> three-dimensional macromolecular structures	
<b>5</b>		<b>Taxonomy:</b> organisms in GenBank	
<b>none</b>		<b>SNP:</b> single nucleotide polymorphism	
<b>841</b>		<b>Gene:</b> gene-centered information	
<b>28</b>		<b>HomoloGene:</b> eukaryotic homology groups	
<b>13</b>		<b>GENSAT:</b> gene expression atlas of mouse central nervous system	
<b>14</b>		<b>dbGaP:</b> genotype and phenotype	
<b>219</b>		<b>UniGene:</b> gene-oriented clusters of transcript sequences	
<b>28</b>		<b>CDD:</b> conserved protein domain database	
<b>3894</b>		<b>3D Domains:</b> domains from Entrez Structure	
<b>83</b>		<b>UniSTS:</b> markers and mapping data	
<b>2570</b>		<b>PopSet:</b> population study data sets	
<b>301379</b>		<b>GEO Profiles:</b> expression and molecular abundance profiles	
<b>66</b>		<b>GEO DataSets:</b> experimental sets of GEO data	
<b>1</b>		<b>Cancer Chromosomes:</b> cytogenetic databases	
<b>49</b>		<b>PubChem BioAssay:</b> bioactivity screens of chemical substances	
<b>220</b>		<b>PubChem Compound:</b> unique small molecule chemical structures	
<b>17</b>		<b>PubChem Substance:</b> deposited chemical structures	

# Gene and protein databases

## ■ UniProt

- Universal Protein Resource
- <http://beta.uniprot.org/>
- Consists of three components:
  - The **UniProt Knowledgebase (UniProtKB)** is the central access point for extensive curated protein information, including function, classification, and cross-reference.
  - The **UniProt Reference Clusters (UniRef)** databases combine closely related sequences into a single record to speed searches.
  - The **UniProt Archive (UniParc)** is a comprehensive repository, reflecting the history of all protein sequences.



Notice: This site will be replaced with [beta.uni-prot.org](#)

- [Text Search](#)
- [BLAST](#)
- [FAQ](#)
- [Help Desk](#)
- [Download](#)

## Welcome to UniProt

UniProt (Universal Protein Resource) is a catalog of information on protein sequence and function created from Swiss-Prot, TrEMBL, and PIR.

UniProt has three components: **UniProt Knowledgebase** (an extensive curated protein sequence and cross-reference database), **UniProt Archive** (reflecting the history of all protein sequences), and **UniProt Consortium**.

The sequences and information are available via [BLAST similarity search](#), a



[European Bioinformatics Institute](#)

Search in

Query

Protein Knowledgebase (UniProtKB)

[Fields >](#)

### WELCOME

The mission of UniProt is to provide the scientific community with a comprehensive, high-quality and freely accessible resource of protein sequence and functional information.

### What we provide

UniProtKB	Protein knowledgebase, consists of two sections: <ul style="list-style-type: none"> <li>★ Swiss-Prot, which is manually annotated and reviewed.</li> <li>★ TrEMBL, which is automatically annotated and is <b>not</b> reviewed.</li> </ul>
UniRef	Sequence clusters, used to speed up similarity searches.
UniParc	Sequence archive, used to keep track of sequences and their identifiers.
Supporting data	Literature citations, taxonomy, keywords and more.



This is a beta version of our new web site, we are welcoming your [feedback](#). Please take into account that some functions are incomplete and many help pages are still preliminary.

### NEWS



#### Release 13.4 - May 20, 2008

*Swiss-Prot in the Wonderland of protein names* · [Cross-references to CGD](#)

- > [Statistics for UniProtKB: Swiss-Prot · TrEMBL](#)
- > [Forthcoming changes](#)
- > [News archives](#)

### SITE TOUR



Learn how to make best use of the tools and data on this site.

### PROTEIN SPOTLIGHT

#### The selfish smell May 2008

We are surrounded by smells. Pleasant ones and not so pleasant ones, hard to distinguish ones, mild ones and strong ones. Smells are not part of our everyday life for the simple sake of pleasure. They are there for a purpose...



Search in

Query

Protein Knowledgebase (UniProtKB)

HIV

Search

Clear

Fields >

Search

Blast

Align

Retrieve

ID Mapping \*

1 – 25 of 151,248 results for HIV in UniProtKB sorted by score descending

Browse by taxonomy, keyword, gene ontology, enzyme class or pathway | Reduce sequence redundancy to 100%, 90% or 50% | Customize display

Download...

> Show only reviewed (UniProtKB/Swiss-Prot) or unreviewed (UniProtKB/TrEMBL) entries

> Restrict term "hiv" to author, gene name, protein name, organism, strain, taxonomy, web resource

Page 1 of 6,050 | Next >

All	Accession	Entry Name	Status	Protein Names	Genes	Organism	Length
<input type="checkbox"/>	<a href="#">P04585</a>	POL_HV1H2	★	<b>Gag-Pol polyprotein</b> (Pr180Gag-Pol) [Cleaved into: Matrix protein p17 (MA); Capsid protein p24 (CA); Spacer peptide p2; Nucleocapsid protein p7 (NC); Transframe peptide (TF); p6-pol (p6*); Protease (EC 3.4.23.16) (Retropesin) (PR); Reverse transcriptase/ribonuclease H (EC 2.7.7.49) (EC 2.7.7.7) (EC 3.1.26.4) (p66 RT); p51 RT; p15; Integrase (IN)]	<b>gag-pol</b>	Human immunodeficiency virus type 1 (isolate HXB2 group M subtype B) (HIV-1)	1,435
<input type="checkbox"/>	<a href="#">P03366</a>	POL_HV1B1	★	<b>Gag-Pol polyprotein</b> (Pr180Gag-Pol) [Cleaved into: Matrix protein p17 (MA); Capsid protein p24 (CA); Spacer peptide p2; Nucleocapsid protein p7 (NC); Transframe peptide (TF); p6-pol (p6*); Protease (EC 3.4.23.16) (Retropesin) (PR); Reverse transcriptase/ribonuclease H (EC 2.7.7.49) (EC 2.7.7.7) (EC 3.1.26.4) (p66 RT); p51 RT; p15; Integrase (IN)]	<b>gag-pol</b>	Human immunodeficiency virus type 1 (isolate BH10 group M subtype B) (HIV-1)	1,447
<input type="checkbox"/>	<a href="#">P03367</a>	POL_HV1BR	★	<b>Gag-Pol polyprotein</b> (Pr180Gag-Pol) [Cleaved into: Matrix protein p17 (MA); Capsid protein p24 (CA); Spacer peptide p2; Nucleocapsid protein p7 (NC); Transframe peptide (TF); p6-pol (p6*); Protease (EC 3.4.23.16) (Retropesin) (PR); Reverse transcriptase/ribonuclease H (EC 2.7.7.49) (EC 2.7.7.7) (EC 3.1.26.4) (p66 RT); p51 RT; p15; Integrase (IN)]	<b>gag-pol</b>	Human immunodeficiency virus type 1 (isolate BRU/LAI group M subtype B) (HIV-1)	1,447

▪ **Data Resources & Tools**

- **EMBL-BANK**
- **UniProt**
- **ArrayExpress**
- **Ensembl**
- **InterPro**
- **PDB-EBI**
- Genomes
- Nucleotide Sequences
- Protein Sequences
- Macromolecular Structures
- Small Molecules
- Gene Expression
- Molecular Interactions
- Literature
- Taxonomy
- Sequence Similarity & Analysis
- Pattern & Motif Searches

▪ **About the EBI**

- **Research**
- **PhD Studies**
- **Training**
- **Industry Support**
- **Group & Team Leaders**
- **EBI Funders**
- User Support
- EBI Mission
- People
- Events and News
- How to Find Us



Search for *HIV* in *All the EBI*

<input type="button" value="▶"/> <a href="#">Genomes</a>	87	<input type="button" value="▶"/> <a href="#">Molecular Interactions</a>	3
<input type="button" value="▶"/> <a href="#">Nucleotide Sequences</a>	250,495	<input type="button" value="▶"/> <a href="#">Reactions &amp; Pathways</a>	36
<input type="button" value="▶"/> <a href="#">Protein Sequences</a>	51,222	<input type="button" value="▶"/> <a href="#">Protein Families</a>	41
<input type="button" value="▶"/> <a href="#">Macromolecular Structures</a>	628	<input type="button" value="▶"/> <a href="#">Enzymes</a>	2
<input type="button" value="▶"/> <a href="#">Small molecules</a>	3	<input type="button" value="▶"/> <a href="#">Literature</a>	183,676
<input type="button" value="▶"/> <a href="#">Gene Expression</a>	23	<input type="button" value="▶"/> <a href="#">Ontologies</a>	6
		<input type="button" value="▶"/> <a href="#">EBI Web Site</a>	34

**Refine your search:**

Search for *HIV* in *All the EBI*  
 with the following keywords

## 51,108 results found in UniProt KB

### [HTSF1\\_HUMAN](#)

O43719, Q59G06, Q99730

HIV Tat-specific factor 1 (Tat-SF1).

View: [▶ in UniProt format](#) [▶ in SRS](#) [▶ in UniSave](#) [▶ in Interpro matches](#)

References: [▶ EMBL-Bank](#) [▶ Ensembl](#) [▶ InterPro](#) [▶ Medline](#) [▶ Taxonomy](#) [▶ MSD/PDB](#) [▶ EMBL-Bank \(Coding Sequence\)](#) [▶ GO](#)

### [A2R057\\_ASPNG](#)

A2R057

Function: N. elliposporum CV-N has anti-HIV activity.

View: [▶ in UniProt format](#) [▶ in SRS](#) [▶ in UniSave](#) [▶ in Interpro matches](#)

References: [▶ Taxonomy](#) [▶ EMBL-Bank](#) [▶ EMBL-Bank \(Coding Sequence\)](#) [▶ InterPro](#) [▶ Medline](#)

### [Q9UDI3\\_HUMAN](#)

Q9UDI3

26 S protease subunit 7, MSS1=MODULATOR of HIV TAT-mediatedtransactivation (Fragments).

View: [▶ in UniProt format](#) [▶ in SRS](#) [▶ in UniSave](#) [▶ in Interpro matches](#)

References: [▶ Taxonomy](#) [▶ Medline](#)

[more...](#)

## 51 results found in UniRef100

### [UNIREF100\\_O43719](#)

Cluster: HIV Tat-specific factor 1

View: [▶ in UniRef format](#) [▶ in SRS](#)

References: [▶ Taxonomy](#) [▶ UniParc](#) [▶ UniRef50](#) [▶ UniRef90](#) [▶ UniProt KB](#)

### [UNIREF100\\_UPI0000DBF22B](#)

Cluster: HIV TAT specific factor 1

View: [▶ in UniRef format](#) [▶ in SRS](#)

References: [▶ Taxonomy](#) [▶ UniParc](#) [▶ UniRef50](#) [▶ UniRef90](#)

### [UNIREF100\\_UPI00015A4CD3](#)

Cluster: HIV TAT specific factor 1

View: [▶ in UniRef format](#) [▶ in SRS](#)

References: [▶ Taxonomy](#) [▶ UniParc](#) [▶ UniRef50](#) [▶ UniRef90](#)

[more...](#)

The UniRef databases provide clustered sets of sequences from UniProt (including splice variants and isoforms) and selected UniParc records, in order to obtain complete coverage of sequence space at several resolutions while hiding redundant sequences (but not their descriptions) from view. The UniRef100 database combines identical sequences and sub-fragments with 11 or more residues (from any organism) into a single UniRef entry, displaying the sequence of a representative protein, the accession numbers of all the merged entries, and links to the corresponding UniProtKB and UniParc records.

# Types of databases

## ■ Data resources of multiple types of data

- EBI (European Bioinformatics Institute)
- NCBI (National Center for Biotechnology Information)
- KEGG (Kyoto Encyclopedia of Genes and Genomes)

## ■ Gene and protein information

- GenBank, UniProt, and PDB
- Species specific: FlyBase, dictyBase, etc.

## ■ Ontological data

- Gene Ontology

## ■ Pathway data

- KEGG PATHWAY, Reactome, BRENDA, etc.

## ■ Protein-protein interaction data

- IntAct, BioGRID, etc.

# Gene and protein databases

## ■ GenBank



























































- A part of NCBI
- <http://www.ncbi.nlm.nih.gov/>
- Search can be performed through the Entrez interface, which searches for the query in all databases available at the NCBI

Search across databases



[Help](#)

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

<b>184871</b>		<b>PubMed:</b> biomedical literature citations and abstracts	
<b>40435</b>		<b>PubMed Central:</b> free, full text journal articles	
<b>106</b>		<b>Site Search:</b> NCBI web and FTP sites	
<b>5015</b>		<b>Books:</b> online books	
<b>222</b>		<b>OMIM:</b> online Mendelian Inheritance in Man	
<b>none</b>		<b>OMIA:</b> online Mendelian Inheritance in Animals	
<b>263187</b>		<b>Nucleotide:</b> Core subset of nucleotide sequence records	
<b>429</b>		<b>EST:</b> Expressed Sequence Tag records	
<b>174358</b>		<b>GSS:</b> Genome Survey Sequence records	
<b>273432</b>		<b>Protein:</b> sequence database	
<b>20</b>		<b>Genome:</b> whole genome sequences	
<b>1013</b>		<b>Structure:</b> three-dimensional macromolecular structures	
<b>5</b>		<b>Taxonomy:</b> organisms in GenBank	
<b>none</b>		<b>SNP:</b> single nucleotide polymorphism	
<b>841</b>		<b>Gene:</b> gene-centered information	
<b>28</b>		<b>HomoloGene:</b> eukaryotic homology groups	
<b>13</b>		<b>GENSAT:</b> gene expression atlas of mouse central nervous system	
<b>14</b>		<b>dbGaP:</b> genotype and phenotype	
<b>219</b>		<b>UniGene:</b> gene-oriented clusters of transcript sequences	
<b>28</b>		<b>CDD:</b> conserved protein domain database	
<b>3894</b>		<b>3D Domains:</b> domains from Entrez Structure	
<b>83</b>		<b>UniSTS:</b> markers and mapping data	
<b>2570</b>		<b>PopSet:</b> population study data sets	
<b>301379</b>		<b>GEO Profiles:</b> expression and molecular abundance profiles	
<b>66</b>		<b>GEO DataSets:</b> experimental sets of GEO data	
<b>1</b>		<b>Cancer Chromosomes:</b> cytogenetic databases	
<b>49</b>		<b>PubChem BioAssay:</b> bioactivity screens of chemical substances	
<b>220</b>		<b>PubChem Compound:</b> unique small molecule chemical structures	
<b>17</b>		<b>PubChem Substance:</b> deposited chemical structures	



Search Protein for HIV    Go    Clear    Save Search

Limits    Preview/Index    History    Clipboard    Details

Display Summary    Show 20    Sort by Relevance    Send to

All: 274399    Bacteria: 249    RefSeq: 1520    Related Structures: 232331

Items 1 - 20 of 274399    Page 1 of 13720 Next

This search in Gene shows 750 results, including:

- [Hiv](#) (*Drosophila melanogaster*): Haploinviable
- [env](#) (*Human immunodeficiency virus 1*): gp160; envelope glycoprotein
- [gag](#) (*Human immunodeficiency virus 1*): Pr55(Gag)

▼ Top Organisms [Tree]

- Human immunodeficiency virus 1 (242817)
- Simian immunodeficiency virus
- Human immunodeficiency virus 2 (2921)
- Hepatitis C virus (2563)
- Homo sapiens (2407)
- All other taxa (8848)

More...

- 1:** [1HSI B](#)    Reports    BLink, Conserved Domains, Links  
 Chain B, Crystal Structure At 1.9 Angstroms Resolution Of Human Immunodeficiency Virus (Hiv) Ii Protease Complexed With L- 735,524, An Orally Bioavailable Inhibitor Of The Hiv Proteases  
 gi|1421442|pdb|1HSI|B[1421442]
- 2:** [1HSI A](#)    Reports    BLink, Conserved Domains, Links  
 Chain A, Crystal Structure At 1.9 Angstroms Resolution Of Human Immunodeficiency Virus (Hiv) Ii Protease Complexed With L- 735,524, An Orally Bioavailable Inhibitor Of The Hiv Proteases  
 gi|1421441|pdb|1HSI|A[1421441]
- 3:** [1HSH D](#)    Reports    BLink, Conserved Domains, Links  
 Chain D, Crystal Structure At 1.9 Angstroms Resolution Of Human Immunodeficiency Virus (Hiv) Ii Protease Complexed With L- 735,524, An Orally Bioavailable Inhibitor Of The Hiv Proteases  
 gi|1421436|pdb|1HSH|D[1421436]
- 4:** [1HSH C](#)    Reports    BLink, Conserved Domains, Links  
 Chain C, Crystal Structure At 1.9 Angstroms Resolution Of Human Immunodeficiency Virus (Hiv) Ii Protease Complexed With L- 735,524, An Orally Bioavailable Inhibitor Of The Hiv Proteases  
 gi|1421435|pdb|1HSH|C[1421435]
- 5:** [1HSH B](#)    Reports    BLink, Conserved Domains, Links  
 Chain B, Crystal Structure At 1.9 Angstroms Resolution Of Human Immunodeficiency Virus (Hiv) Ii Protease Complexed With L- 735,524, An Orally Bioavailable Inhibitor Of The Hiv Proteases  
 gi|1421434|pdb|1HSH|B[1421434]
- 6:** [1HSH A](#)    Reports    BLink, Conserved Domains, Links  
 Chain A, Crystal Structure At 1.9 Angstroms Resolution Of Human Immunodeficiency Virus (Hiv) Ii Protease Complexed With L- 735,524, An Orally Bioavailable Inhibitor Of The Hiv Proteases  
 gi|1421433|pdb|1HSH|A[1421433]



Search  for

Display  Show  Send to

Range: from  to  Features:  CDD

1: [1HSI\\_B](#) Reports Chain B, Crystal ...[gi:1421442]

[Comment](#) [Features](#) [Sequence](#)

LOCUS 1HSI\_B 99 aa linear VRL 01-OCT-2007  
 DEFINITION Chain B, Crystal Structure At 1.9 Angstroms Resolution Of Human Immunodeficiency Virus (Hiv) Ii Protease Complexed With L- 735,524, An Orally Bioavailable Inhibitor Of The Hiv Proteases.  
 ACCESSION 1HSI\_B  
 VERSION 1HSI\_B GI:1421442  
 DBSOURCE pdb: molecule 1HSI, chain 66, release Aug 27, 2007;  
 deposition: Aug 27, 2007;  
 class: Hydrolase (Acid Proteinase);  
 source: Mol\_id: 1; Organism\_scientific: Human Immunodeficiency Virus 2; Gene: Hiv-2 Protease From The Rod Isolate;  
 Expression\_system: Escherichia Coli;  
 Exp. method: X-Ray Diffraction.  
 KEYWORDS .  
 SOURCE Human immunodeficiency virus 2 (HIV-2)  
 ORGANISM [Human immunodeficiency virus 2](#)  
 Viruses; Retro-transcribing viruses; Retroviridae;  
 Orthoretrovirinae; Lentivirus; Primate lentivirus group.  
 REFERENCE 1 (residues 1 to 99)  
 AUTHORS Chen,Z., Li,Y., Chen,E., Hall,D.L., Darke,P.L., Culberson,C., Shafer,J.A. and Kuo,L.C.  
 TITLE Crystal structure at 1.9-A resolution of human immunodeficiency virus (HIV) II protease complexed with L-735,524, an orally bioavailable inhibitor of the HIV proteases

# Gene and protein databases

## ■ PDB: Protein DataBank

- Contains 3-dimensional protein structures
- <http://www.rcsb.org>
- Data submitted by individual researchers

Advanced Keyword Query for: HIV

1 2 3 4 5 .. 46 ↩

1j9v



Solution structure of a lactam analogue (DabD) of HIV gp41 600-612 loop.



**Characteristics Classification**

**Release Date:** 01-Jul-2003 **Exp. Method:** NMR 49 Structures  
**Viral Protein**

**Compound**

**Polymer:** 1 **Molecule:** DabD (Ace)IWG(DAB)SGKLIDTTA ANALOGUE OF HIV GP41 **Chains:** A

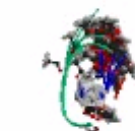
**Authors**

**Phan Chan Du, A., Limal, D., Semetey, V., Dali, H., Jolivet, M., Desgranges, C., Cung, M.T., Briand, J.P., Petit, M.C., Muller, S.**

1jaa



Solution structure of lactam analogue (DapE) of HIV gp41 600-612 loop.



**Characteristics Classification**

**Release Date:** 01-Jul-2003 **Exp. Method:** NMR 50 Structures  
**Viral Protein**

**Compound**

**Polymer:** 1 **Molecule:** DapE : (Ace)IWG(Dap)SGKLIETTA ANALOGUE OF HIV GP41 **Chains:** A

**Authors**

**Phan Chan Du, A., Limal, D., Semetey, V., Dali, H., Jolivet, M., Desgranges, C., Cung, M.T., Briand, J.P., Petit, M.C., Muller, S.**

1jar



Solution structure of lactam analogue (DDab) of HIV gp41 600-612 loop.



**Characteristics Classification**

**Release Date:** 01-Jul-2003 **Exp. Method:** NMR 49 Structures  
**Viral Protein**

**Compound**

**Polymer:** 1 **Molecule:** DDab: (ACE)IWGDSGKLI(DAB)TTA ANALOGUE OF HIV GP41 **Chains:** A

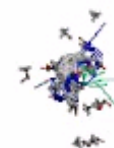
**Authors**

**Phan Chan Du, A., Limal, D., Semetey, V., Dali, H., Jolivet, M., Desgranges, C., Cung, M.T., Briand, J.P., Petit, M.C., Muller, S.**

1jc8



Solution structure of lactam analogue (DDap) of gp41 600-612 loop of HIV



**Characteristics Classification**

**Release Date:** 01-Jul-2003 **Exp. Method:** NMR 20 Structures  
**Viral Protein**

**Compound**

**Polymer:** 1 **Molecule:** DDap: (ACE)IWGDSGKLI(DNP)TTA ANALOGUE OF HIV GP41 **Chains:** A **Other Details:** ATOM OD2

Are you missing data?  
<ftp://ftp.wwpdb.org>.  
For more information click

## Welcome to the

The RCSB PDB provides a resource for studying the structures of proteins and their relationships to disease.

The RCSB is a member of the wwPDB to ensure that the PDB archive is a resource with uniform data.

This site offers tools for browsing and reporting that utilize the data in the archive to create a more complete archive.

Information about compatibility is [here](#).

A [narrated tutorial](#) will illustrate how to navigate, browse, generate reports, and download structures using this new [player](#) download.]

Comments? [info@rcsb.org](mailto:info@rcsb.org)

Molecule of the Month: Pri



CONTACT US | HELP | PRINT PAGE

PDB ID or keyword  Author

Home Search

- Home
- Getting Started
- ▶ Download Files
- ▶ Deposit and Validate
- ▶ Structural Genomics
- ▶ Dictionaries & File Formats
- ▶ Software Tools
- ▶ General Education
- ▶ Site Tutorials
- BioSync
- ▶ General Information
- Acknowledgements
- Frequently Asked Questions
- ✉ Report Bugs/Comments

### Quick Tips:



Having trouble with the web site? Try the tutorial: click [here](#)



Are you missing data updates? The PDB archive has moved to <ftp://ftp.wwpdb.org>.  
For more information click [here](#).

2dit

DOI 10.2210/pdb2dit/pdb

Red - Derived Information

Title	Solution structure of the RRM_1 domain of HIV TAT specific factor 1 variant
Authors	Dang, W., Muto, Y., Inoue, M., Kigawa, T., Shirouzu, M., Terada, T., Yokoyama, S., RIKEN Structural Genomics/Proteomics Initiative (RSGI)
Primary Citation	Dang, W., Muto, Y., Inoue, M., Kigawa, T., Shirouzu, M., Terada, T., Yokoyama, S. Solution structure of the RRM_1 domain of HIV TAT specific factor 1 variant. <i>To be Published</i>
History	Deposition 2006-03-30 Release 2006-09-30
Experimental Method	Type NMR, 20 STRUCTURES Data N/A
NMR Ensemble	Conformers Calculated 100 Conformers Submitted 20 Selection Criteria target function, structures with the lowest energy, structures with the least restraint violations
NMR Refine	Method NMR, 20 STRUCTURES

Molecular Description Asymmetric Unit  
 Polymer: 1 Molecule: HIV TAT specific factor 1 variant Fragment: RRM\_1 domain Chains: A

Classification RNA Binding Protein

Source Polymer: 1 Scientific Name: **Homo sapiens** Common Name: **Human** Expression system: **Cell free synthesis**

GO Terms  
 Polymer: HIV TAT specific factor 1 variant (2DIT:A)  
 Molecular Function: none  
 Biological Process: none  
 Cellular Component: none

Images and Visualization

Asymmetric Unit



Display Options

- KiNG
- Jmol**
- WebMol

- MBT Protein Workshop
- QuickPDB
- All Images

Quick Tips: [Navigation icons]

When exploring a structure, select Structure

CONTACT US | HELP | PRINT PAGE

PDB ID or keyword  Author

Site Search

Advanced Search

Home Search Structure Results

Queries

Are you missing data updates? The PDB archive has moved to <ftp://ftp.wwpdb.org>. For more information click [here](#).

Jmol

2dit

2DIT

Download Files

FASTA Sequence

Download Original Files

Display Files

Display Molecule

Image Gallery

KING Viewer

Jmol Viewer

WebMol Viewer

Protein Workshop

FirstGlance

Rasmol Viewer  
(Plugin required)

Swiss-PDB Viewer  
(Plugin required)

Molecular Viewers Help

KING Help

Jmol Help

WebMol Help

Protein Workshop Help

QuickPDB

Asymmetric Unit

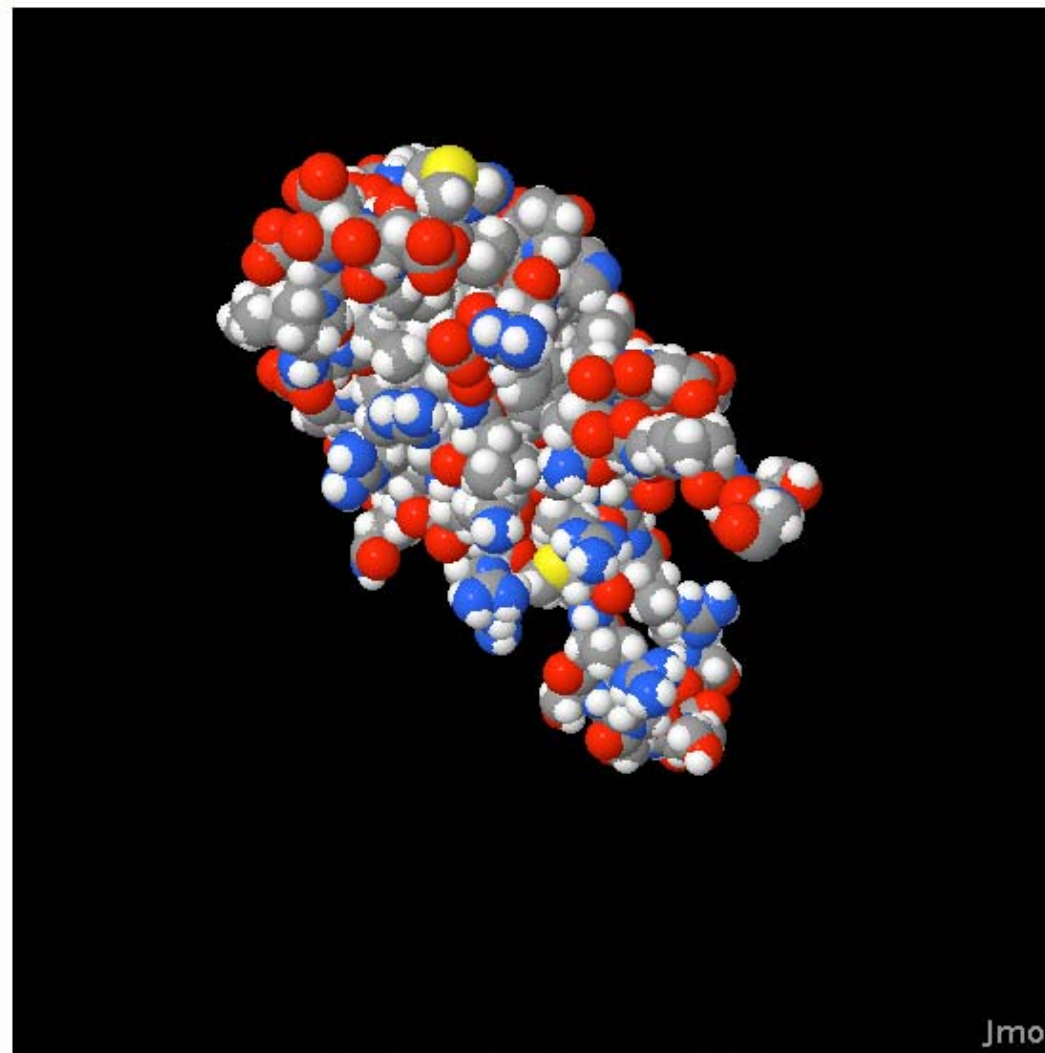
Assumed Biological Molecule 1

Structural Reports

External Links

Structure Analysis

Help



For help select one of the options below:

# Databases of Model Organisms

- Model organisms are those whose genome has been extensively studied such that the workings of biological phenomena for more complex organisms can be inferred.
- For example, the mouse has been most studied to understand other mammalian systems.
- For plant species, *Arabidopsis thaliana* is most often used as a model organism.



# Databases of model organisms

## ■ *Arabidopsis* (mustard plant)

- TAIR: The Arabidopsis Information Resource
  - <http://www.arabidopsis.org/>
- The Carnegie Institution of Washington, the National Center for Genome Resources

## ■ *C. elegans*

- WormBase
  - <http://www.wormbase.org/>
- Cold Spring Harbor Laboratory

## ■ *Dictyostelium* (slime mold)

- dictyBase
  - <http://dictybase.org/>
- Northwestern University



# Databases of model organisms

## ■ *Drosophila* (Fruit fly)

- FlyBase

  - <http://www.flybase.org/>

- Indiana University

## ■ Mouse

- Mouse Genome Informatics (MGI)

  - <http://www.informatics.jax.org/>

- the Jackson Laboratory

## ■ Rat

- Rat Genome Database (RGD)

  - <http://rgd.mcw.edu>

- The Medical College of Wisconsin

# Databases of model organisms

- *Saccharomyces* (Yeast)
  - *Saccharomyces* Genome Database (SGD)
    - <http://www.yeastgenome.org/>
  - Stanford University
- *Xenopus* (African clawed frog)
  - Xenbase
    - <http://www.xenbase.org/>
  - The University of Calgary, Alberta
- Zebrafish
  - ZFIN
    - <http://zfin.org/>
  - The University of Oregon

# Types of databases

- **Data resources of multiple types of data**
  - EBI (European Bioinformatics Institute)
  - NCBI (National Center for Biotechnology Information)
  - KEGG (Kyoto Encyclopedia of Genes and Genomes)
- **Gene and protein information**
  - GenBank, UniProt, and PDB
  - Species specific: FlyBase, dictyBase, etc.
- **Ontological data**
  - Gene Ontology
- **Pathway data**
  - KEGG PATHWAY, Reactome, BRENDA, etc.
- **Protein-protein interaction data**
  - IntAct, BioGRID, etc.

Break time!

# The Gene Ontology

- The Gene Ontology (GO) project is a collaborative effort to address the need for consistent descriptions of gene products in different databases.
- <http://www.geneontology.org/>
- In annotating gene and protein function, scientists from different backgrounds tend to use different terminology to describe the same event.
  - Eg. Looking for all the gene products that are involved in bacterial protein synthesis, if one database describes these molecules as being involved in 'translation', whereas another uses the phrase 'protein synthesis', it will be difficult for you - and even harder for a computer - to find functionally equivalent terms.

# The GO Hierarchy

- The GO project has developed three structured controlled vocabularies (ontologies) that describe gene products in terms of their associated *biological processes, cellular components and molecular functions* in a species-independent manner.
- There are three separate aspects to this effort:
  1. the development and maintenance of the ontologies themselves
  2. the annotation of gene products, which entails making associations between the ontologies and the genes and gene products in the collaborating databases
  3. the development of tools that facilitate the creation, maintenance and use of ontologies.



# Terms in the GO

- GO entries, or terms, are given IDs in the form GO:nnnnnnnn and a term name, such as “signal transduction”
- Each term is assigned to one of the three ontologies (molecular function, cellular component, or biological process)

# Term-Term relationships

- GO terms can be linked by five types of relationships: `is_a`, `part_of`, `regulates`, `positively_regulates` and `negatively_regulates`.
- A `is_a` B: A is a subclass of B
  - E.g. nuclear chromosome `is_a` chromosome

```
GO:0043232 : intracellular non-membrane-bound organelle
[i] GO:0005694 : chromosome
---[i] GO:0000228 : nuclear chromosome
```

# Term-Term relationships

- C part\_of D: C is always a part of D if C is present
  - E.g. periplasmic flagellum part\_of periplasmic space

```
GO:0044464 : cell part
[i] GO:0042995 : cell projection
---[i] GO:0019861 : flagellum
-----[i] GO:0009288 : flagellin-based flagellum
-----[i] GO:0055040 : periplasmic flagellum
[i] GO:0042597 : periplasmic space
---[p] GO:0055040 : periplasmic flagellum
```

# Term-Term relationships

- The **regulates**, **positively\_regulates** and **negatively\_regulates** relationships describe interactions between biological processes and other biological processes, molecular functions or biological qualities.
- E regulates F: E modulates the occurrence of F. If F is a biological quality, then E modulates the value of F.
  - The term **regulation of transcription**. When **regulation of transcription** occurs, it always alters the rate, extent or frequency at which a gene is transcribed.

```
GO:0009987 : cellular process
[i] GO:0010467 : gene expression
---[r] GO:0010468 : regulation of gene expression
-----[i] GO:0045449 : regulation of transcription
---[p] GO:0006350 : transcription
-----[r] GO:0045449 : regulation of transcription
```

[Open menus](#)[Home](#)[FAQ](#)[Downloads](#)[Tools](#)[Documentation](#)[About GO](#)[Projects](#)[Contact GO](#)[Site Map](#)

## Gene Ontology Home

The Gene Ontology project provides a controlled vocabulary to describe gene and gene product attributes in any organism. [Read more about the Gene Ontology...](#)

### Search the Gene Ontology Database

Search for genes, proteins or GO terms using [AmiGO](#):

gene or protein name  GO term or ID

[AmiGO](#) is the official GO browser and search engine. [Browse the Gene Ontology with AmiGO.](#)

### GO website

- ◆ The latest news and views in the [GO newsletter](#)
- ◆ [GO downloads](#), including [ontology files](#), [annotations](#) and the [GO database](#)
- ◆ [Tools](#) for using GO, including [OBO-Edit downloads](#), [AmiGO](#), and the [GO Online SQL Environment](#).
- ◆ [Request new terms or ontology changes](#) or [get help with new term submission](#)
- ◆ [Documentation](#) on all aspects of the GO project and the [GO FAQ](#)
- ◆ Projects within the GO consortium, including [Reference Genomes](#) and [immune system annotation](#)
- ◆ [Gene Ontology mailing lists](#) and [contact details](#)

14 results for **lectin** in terms fields **term accession, term name, synonyms**

▼ Filter search results ?

Ontology

- All
- biological process
- cellular component
- molecular function

Set filters

Remove all filters

Select all

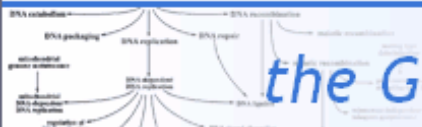
Clear all

Perform an action with the selected terms...

Go!

rel ↓	Accession , Term		Ontology
<input type="checkbox"/>	GO:0001862 : <a href="#">col</a> <b>lectin</b> binding <a href="#">[show def]</a>	0 gene products <a href="#">view in tree</a>	<a href="#">molecular function</a>
<input type="checkbox"/>	GO:0043698 : <a href="#">iridosome</a> <a href="#">[show def]</a> Query matches synonym "ref <b>lectin</b> g platelet" [exact synonym]	0 gene products <a href="#">view in tree</a>	<a href="#">cellular component</a>
<input type="checkbox"/>	GO:0030246 : <a href="#">carbohydrate binding</a> <a href="#">[show def]</a> Query matches synonym "se <b>lectin</b> " [related synonym]	692 gene products <a href="#">view in tree</a>	<a href="#">molecular function</a>
<input type="checkbox"/>	GO:0001863 : <a href="#">col</a> <b>lectin</b> receptor activity <a href="#">[show def]</a>	0 gene products <a href="#">view in tree</a>	<a href="#">molecular function</a>
<input type="checkbox"/>	GO:0046703 : <a href="#">natural killer cell</a> <b>lectin</b> -like receptor binding <a href="#">[show def]</a> Query matches synonym "NK cell <b>lectin</b> -like receptor binding" [exact synonym]	10 gene products <a href="#">view in tree</a>	<a href="#">molecular function</a>
<input type="checkbox"/>	GO:0001867 : <a href="#">complement activation, <b>lectin</b> pathway</a> <a href="#">[show def]</a>	16 gene products <a href="#">view in tree</a>	<a href="#">biological process</a>
<input type="checkbox"/>	GO:0001868 : <a href="#">regulation of complement activation, <b>lectin</b> pathway</a> <a href="#">[show def]</a>	8 gene products <a href="#">view in tree</a>	<a href="#">biological process</a>
<input type="checkbox"/>	GO:0002772 : <a href="#">inhibitory C-type <b>lectin</b> receptor signaling pathway</a> <a href="#">[show def]</a>	0 gene products <a href="#">view in tree</a>	<a href="#">biological process</a>
<input type="checkbox"/>	GO:0002223 : <a href="#">stimulatory C-type <b>lectin</b> receptor signaling pathway</a> <a href="#">[show def]</a>	4 gene products <a href="#">view in tree</a>	<a href="#">biological process</a>





Search GO   terms  genes or proteins  exact match

## Tree Browser

▼ Filter tree view ?

Filter by ontology

Ontology

All  
 biological process  
 cellular component  
 molecular function

Filter Gene Product Counts

Data source	Species
All	All
CGD	Anaplasma phagocy...
dictyBase	Arabidopsis thaliana
FlyBase	Bacillus anthraci...

View Options

Tree view  Full  Compact

- ▣ all : all [250413 gene products]
  - ▣ ⓘ GO:0003674 : molecular\_function [168550 gene products]
    - ▣ ⓘ GO:0005488 : binding [46693 gene products]
      - ▣ ⓘ GO:0030246 : carbohydrate binding [692 gene products]

- Actions...
- Last action: Reset the tree
  - Graphical View
  - Permalink
  - Download...
  - OBO
  - RDF-XML
  - GraphViz dot

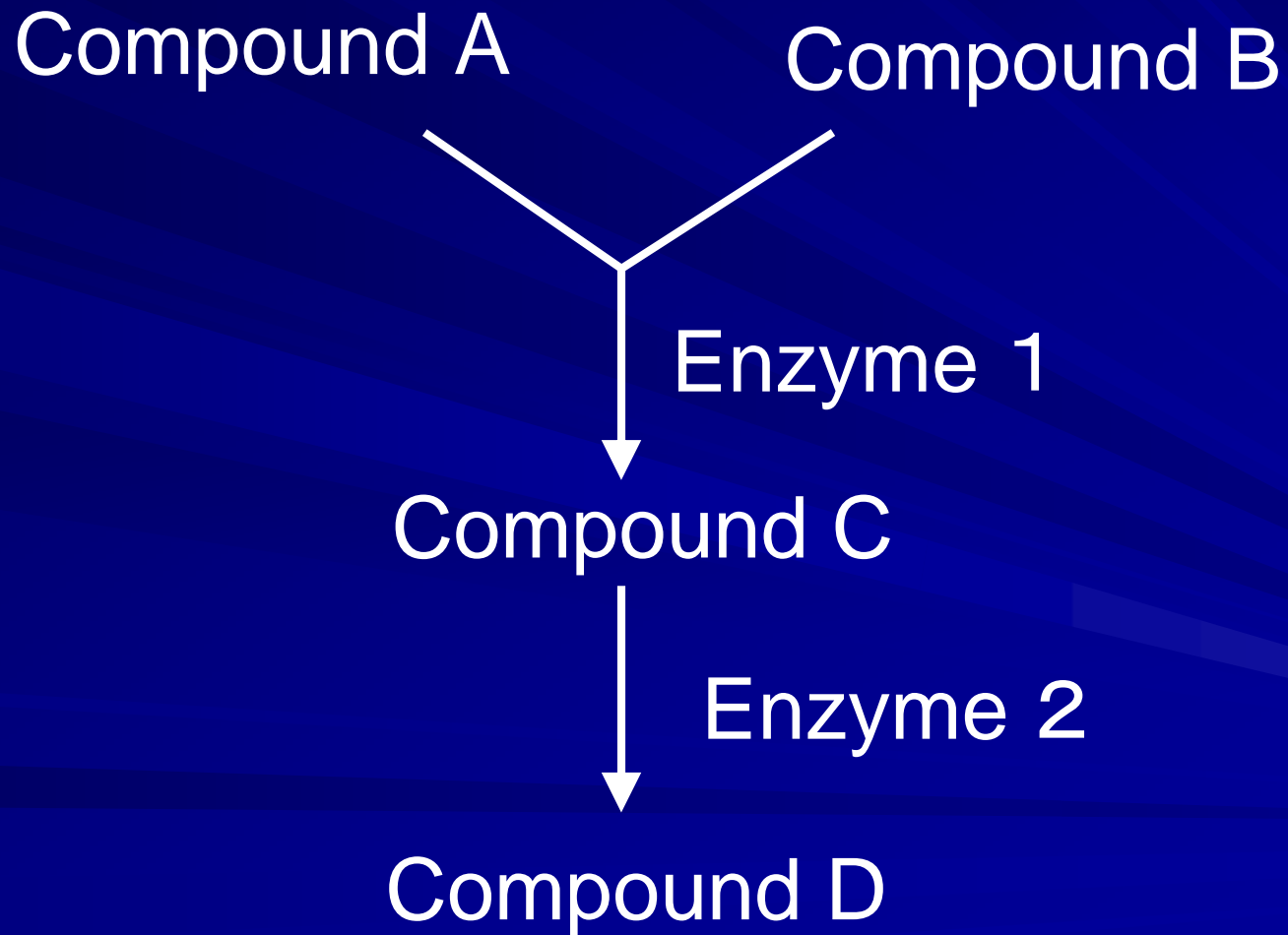
# Types of databases

- **Data resources of multiple types of data**
  - EBI (European Bioinformatics Institute)
  - NCBI (National Center for Biotechnology Information)
  - KEGG (Kyoto Encyclopedia of Genes and Genomes)
- **Gene and protein information**
  - GenBank, UniProt, and PDB
  - Species specific: FlyBase, dictyBase, etc.
- **Ontological data**
  - Gene Ontology
- **Pathway data**
  - KEGG PATHWAY, Reactome, BRENDA, etc.
- **Protein-protein interaction data**
  - IntAct, BioGRID, etc.

# What are pathways?

- Within the body, pathways consist of signal transduction, metabolism and transcription regulation, where information of biological entities interact with one another, forming a map.
- Since the genomic era began, an abundant amount of such pathway information has accumulated, and pathway analysis has become important for the understanding of biological systems.
  - For example, given a particular protein known to express in a certain disease, we want to find the other proteins with which it interacts.
  - As another example, we can build a network of multiple proteins known to interact with one another to glean insights into how it affects the system.

# A simple pathway



# Pathway Databases

## ■ KEGG

- Bioinformatics Center, Kyoto University
- <http://www.genome.jp>

## ■ BRENDA

- Braunschweig Dept. of Bioinformatics
- <http://www.biocyc.org>





## ■ Reactome

- EBI (European Bioinformatics Institute), Cold Spring Harbor Laboratory and the Gene Ontology Consortium
- <http://www.reactome.org>

# BRENDA

- Database of enzyme information, linked with KEGG pathways
- Data on enzyme function are extracted directly from the primary literature by scientists holding a degree in Biology or Chemistry.
- Formal and consistency checks are done by computer programs, each data set on a classified enzyme is checked manually by at least one biologist and one chemist.



-  [BRENDA home](#)
-  [login](#)
-  [history](#)
-  [All enzymes](#)

### SEARCH-Navigator

- close all  open all
- Nomenclature
- Reaction & Specificity
- Functional Parameters
- Organism related Information
- Enzyme Structure
- Isolation & Preparation
- Stability
- Disease & References
- Application & Engineering

-  [Quick search](#)
-  [Fulltext search](#)
-  [Advanced search](#)
-  [Substructure search](#)
-  [TaxTree Explorer](#)
-  [EC Explorer](#)
-  [Sequence Search](#)
-  [Genome Explorer](#)
-  [Ontology Explorer](#)
-  [Download](#)

- [Introduction/References](#)
- [News](#)
- [Contact/Team/Errors](#)
- [Jobs](#)
- [Copyright](#)
- [Related Links](#)
- [Help](#)

<b>EC-Number</b>	<b>Enzyme Name</b>	<b>Organism</b>	<b>Protein</b>	<b>Full text</b>	<b>Advanced Search</b>
<input style="width: 100%; height: 20px;" type="text"/> <input type="button" value="Search"/> Display <input style="width: 30px;" type="text" value="10"/> entries					

Latest BRENDA update 12/2007

Nomenclature	Reaction & Specificity	Functional Parameters
Enzyme Names EC Number Common/ Recommended Name Systematic Name Synonyms CAS Registry Number	<b>Pathway</b> Catalysed Reaction Reaction Type Natural Substrates and Products Substrates and Products Substrates Natural Substrate Products Natural Product Inhibitors Cofactors Metals/Ions Activating Compounds Ligands	Km Value Ki Value pI Value Turnover Number Specific Activity pH Optimum pH Range Temperature Optimum Temperature Range
<b>Isolation &amp; Preparation</b>  Purification Cloned Renatured Crystallization	<b>Enzyme Structure</b>  Sequence/ SwissProt link 3D-Structure/ PDB link Molecular Weight Subunits Posttranslational Modification	<b>Organism-related information</b>  Organism Source Tissue Localization Protein-Specific Search
<b>Stability</b>  pH Stability Temperature Stability General Stability Organic Solvent Stability Oxidation Stability Storage Stability	<b>Disease &amp; References</b>  Disease References	<b>Application &amp; Engineering</b>  Engineering Application

✉ Webmaster: **Maurice Scheer**  
n.scheer@tu-bs.de

**SEARCH-Navigator**

- close all  open all  
 Nomenclature  
 Reaction & Specificity  
 Functional Parameters  
 Organism related Information  
 Enzyme Structure  
 Isolation & Preparation  
 Stability  
 Disease & References  
 Application & Engineering

- [Quick search](#)  
[Fulltext search](#)  
[Advanced search](#)  
[Substructure search](#)  
[TaxTree Explorer](#)  
[EC Explorer](#)  
[Sequence Search](#)  
[Genome Explorer](#)  
[Ontology Explorer](#)  
[Download](#)

- [Introduction/References](#)  
[News](#)  
[Contact/Team/Errors](#)  
[Jobs](#)  
[Copyright](#)  
[Related Links](#)  
[Help](#)

[Any question? -> Use the BRENDA Discussion groups](#)

**Search Pathway**

contains  show  results [clear](#)

Recommended Name:  contains   
 EC Number:  contains   
 KEGG Link:  contains

:= amino acid sequences := comprehensive online version := show the catalyzed reaction

Results 1 - 10 of 23



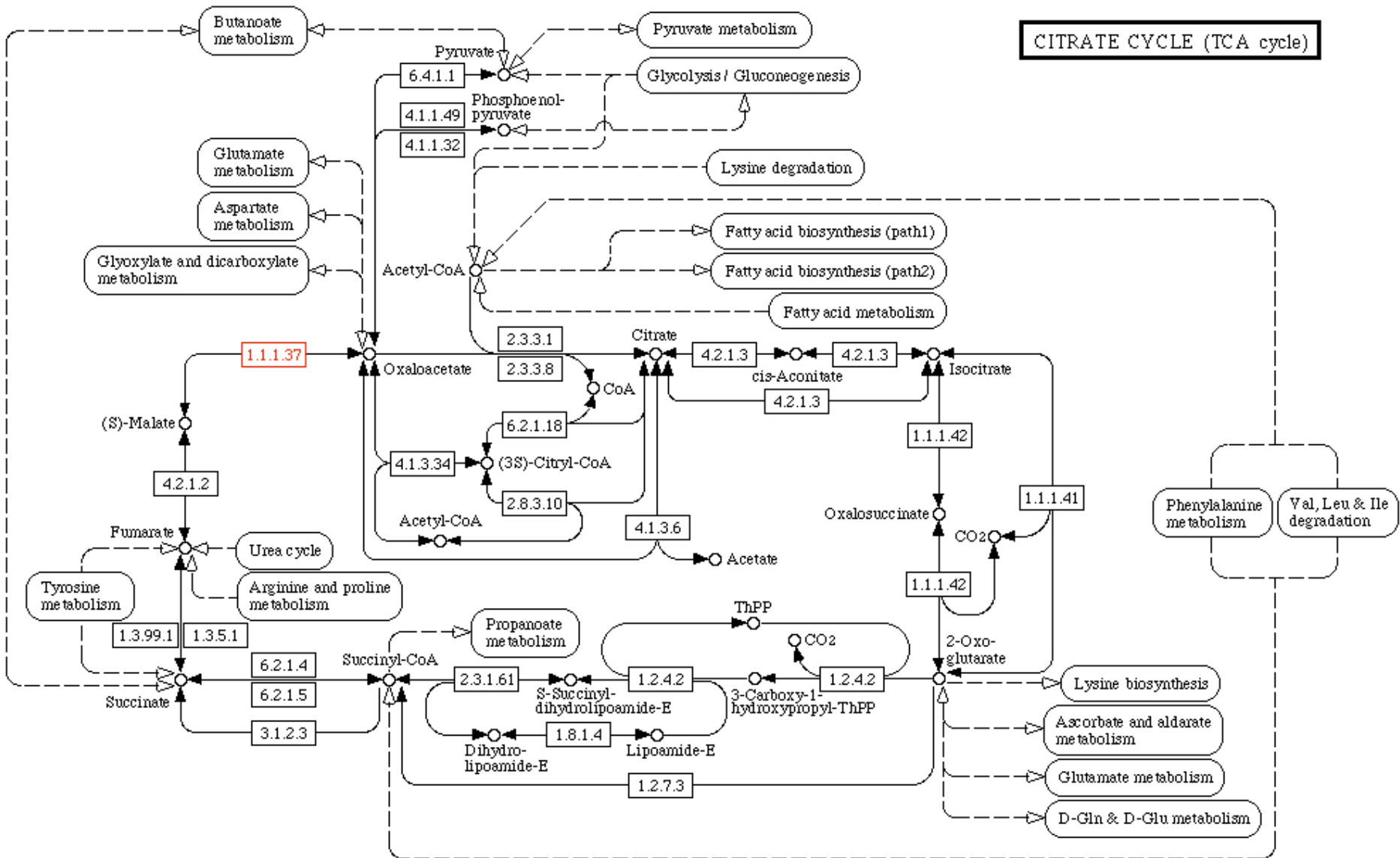
[download this result as tab stop separated values \(Excel,OpenOffice\) format](#)

EC Number	Recommended Name	Pathway	KEGG Link
<a href="#">1.1.1.37</a>	malate dehydrogenase	Citrate cycle (TCA cycle)	<input checked="" type="checkbox"/>
<a href="#">1.1.1.41</a>	isocitrate dehydrogenase (NAD <sup>+</sup> )	Citrate cycle (TCA cycle)	<input checked="" type="checkbox"/>
<a href="#">1.1.1.42</a>	isocitrate dehydrogenase (NADP <sup>+</sup> )	Citrate cycle (TCA cycle)	<input checked="" type="checkbox"/>
<a href="#">1.2.4.2</a>	oxoglutarate dehydrogenase (succinyl-transferring)	Citrate cycle (TCA cycle)	<input checked="" type="checkbox"/>
<a href="#">1.2.7.3</a>	2-oxoglutarate synthase	Citrate cycle (TCA cycle)	<input checked="" type="checkbox"/>
<a href="#">1.3.5.1</a>	succinate dehydrogenase (ubiquinone)	Citrate cycle (TCA cycle)	<input checked="" type="checkbox"/>
<a href="#">1.3.99.1</a>	succinate dehydrogenase	Citrate cycle (TCA cycle)	<input checked="" type="checkbox"/>
<a href="#">1.8.1.4</a>	dihydrolipoyl dehydrogenase	Citrate cycle (TCA cycle)	<input checked="" type="checkbox"/>
<a href="#">2.3.1.61</a>	dihydrolipoyllysine-residue succinyltransferase	Citrate cycle (TCA cycle)	<input checked="" type="checkbox"/>
<a href="#">2.3.3.1</a>	citrate (Si)-synthase	Citrate cycle (TCA cycle)	<input checked="" type="checkbox"/>

[ Pathway menu | Ortholog table ]

Reference pathway Go Current selection Select

**CITRATE CYCLE (TCA cycle)**



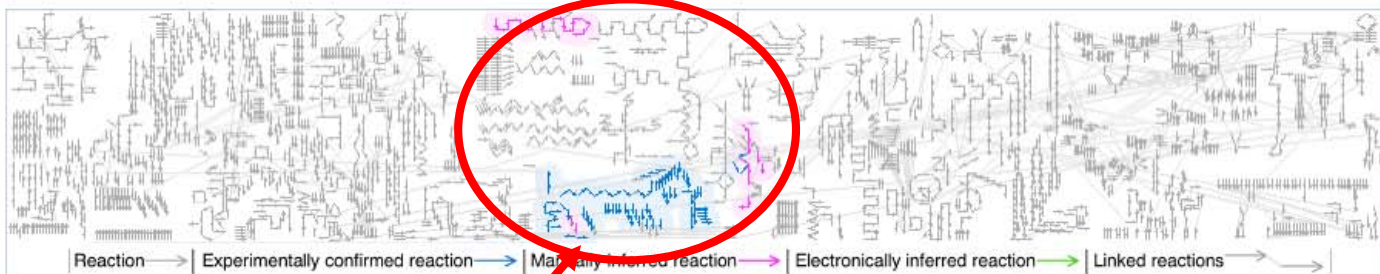
# Reactome

- A curated resource of core pathways and reactions in human biology.
- The information in this database is authored by biological researchers with expertise in their fields, and cross-referenced with [NCBI](#), [UniProt](#), [KEGG \(Gene and Compound\)](#), [PubMed](#), [GO](#) and others.
- In addition to curated human events, [inferred orthologous events](#) in 22 non-human species including mouse, rat, chicken, puffer fish, worm, fly, yeast, and *E.coli* are also available.



# Reactome - a curated knowledgebase of biological pathways

The data displayed is for Homo sapiens. Use the menu to change the species. Check for cross-species comparison.




Reaction → Experimentally confirmed reaction → Manually inferred reaction → Electronically inferred reaction → Linked reactions →

Apoptosis	Botulinum neurotoxicity	Cell Cycle Checkpoints	Cell Cycle, Mitotic
DNA Repair	DNA Replication	Electron Transport Chain	Gap junction trafficking and regulation
Gene Expression	<b>HIV Infection</b>	Hemostasis	Influenza Infection
Integration of energy metabolism	Lipid and lipoprotein metabolism	Membrane Trafficking	Metabolism of amino acids
Metabolism of carbohydrates	Metabolism of nitric oxide	Metabolism of non-coding RNA	Metabolism of vitamins and cofactors
Metabolism of xenobiotics	Nucleotide metabolism	Porphyryn metabolism	Pyruvate metabolism and TCA cycle
Post-translational protein modification	Regulatory RNA pathways	Signaling by BMP	Signaling by EGFR
Signaling by FGFR	Signaling in Immune system	Signaling by Insulin receptor	Signalling by NGF
Signaling by Notch	Signaling by Rho GTPases	Signaling by TGF beta	Signaling by VEGF
Signaling by Wnt	Telomere Maintenance	Transcription	Translation
mRNA Processing			

## Reactome is seeking expert help for the curation of new modules

The Reactome knowledgebase relies on collaborations with research biologists to construct expert consensus views of key biological processes, and to integrate these with other processes already in Reactome. We are seeking new author-collaborators. If you're interested, or would like more information about our data acquisition process, please contact us at [editorial@reactome.org](mailto:editorial@reactome.org). Click here to view/hide a list of high-priority projects now being developed.

### About Reactome

 The **Reactome** project is a collaboration among Cold Spring Harbor Laboratory, The European Bioinformatics Institute, and The Gene Ontology Consortium to develop a curated resource of core pathways and reactions in human biology. The information in this database is authored by biological researchers with expertise in their fields, maintained by the Reactome editorial staff, and cross-referenced with the sequence databases at NCBI, Ensembl and UniProt, the UCSC Genome Browser, HapMap, KEGG (Gene and Compound), ChEBI, PubMed and GO. In addition to curated human events, *inferred orthologous events* in 22 non-human species including mouse, rat, chicken, puffer fish, worm, fly, yeast, two plants and E.coli are also available. A description of Reactome has been published in [Genome Biology](#).

Reactome is a free on-line resource, and Reactome software is open-source. However, please take note of our [disclaimer](#).

### News and Notes

#### • March 12, 2008 Version 24 Released

New additions to curated content include the pathway topics: *Signaling by VEGF*, *Metabolism of nitric oxide*, and *MicroRNA biogenesis*. The pathway topics updated with new curated events include: **Apoptosis - Cleavage reactions**, **Gene expression - Generic transcription pathway**, **Hemostasis - Tie2 and PECAM events**, **Immune signaling - TCR cascade**, **Energy metabolism - PDC regulation**, **Membrane Trafficking - COP II mediated events**, **Post translational modifications - Hypusine synthesis**, **EGFR signaling - PLC-gamma and Grb2 mediated downstream events**. K Schulze-Osthoff, and L Castagnoli are our external expert annotators, and S Ranganathan, LP Freedman, J Trowsdale, J Lippincott-Schwartz, G Enikolopov, F Karginov, GJ Hannon, CH Heldin, and L Claesson-Welsh are our external reviewers, in this release. As usual, *protein-protein interaction datasets*, *Statistics* and the *Editorial Calendar* are available. Reactome team thanks the users for their comments on the *pathway visualization tool*.

Click on [contact](#) to reach us, on [editorial](#) to contribute to Reactome content, and on [subscribe](#) to receive Reactome announcements.

• [More...](#)



# HIV Infection [Homo sapiens]

+ Reactionmap

- Details

[open to selected event](#) [open all](#) [close all](#)

HIV Infection [Homo sapiens]

HIV Life Cycle

Host Interactions of HIV factors

Host Immune responses to HIV infection

## HIV Infection

**Stable identifier** REACT\_6185.3

**Authored** Bukrinsky, M, D'Eustachio, P, Gillespie, ME, Gopinathrao, G, Iordanskiy, S, Morrow, MP, Matthews, L, Rice, AP

The global pandemic of Human Immunodeficiency Virus (HIV) infection has resulted in tens of millions of people infected by the virus and millions more affected. UNAIDS estimates around 40 million HIV/AIDS patients worldwide with 75% of them living in sub-Saharan Africa. The primary method of HIV infection is by sexual exposure while nonsexual HIV transmission also can occur through transfusion with contaminated blood products, injection drug use, occupational exposure, accidental needlesticks or mother-to-child transmission. HIV damages the immune system, leaving the infected person vulnerable to a variety of "opportunistic" infections arising from host immune impairment (Hare, 2004).

HIV-1 and the less common HIV-2 belong to the family of retroviruses. HIV-1 contains a single-stranded RNA genome that is 9 kilobases in length and contains 9 genes that encode 15 different proteins. These proteins are classified as: structural proteins (Gag, Pol, and Env), regulatory proteins (Tat and Rev), and accessory proteins (Vpu, Vpr, Vif, and Nef) (Frankel and Young, 1998).

**HIV infection cycle** can be divided into two phases:

1. An **Early phase** consisting of early events occurring after HIV infection of a susceptible target cell and a
2. **Late phase** comprising the later events in the HIV-infected cell resulting in the assembly of new infectious virions. The section titled **HIV lifecycle** consists of annotations of events in these two phases.

The virus has developed various molecular strategies to suppress the antiviral immune responses (innate, cellular and humoral) of the host. HIV-1 viral auxiliary proteins (Tat, Rev, Nef, Vif, Vpr and Vpu) play important roles in these host-pathogen interactions (Li et al., 2005). The section titled **Host interactions of HIV factors** will highlight these complex post-infection processes and the annotations will be released in near future.

[Hare 2004, Frankel & Young 1998, Li et al 2005]

**Organism** Homo sapiens  
Human immunodeficiency virus 1

**Cellular compartment** cell [GO](#)

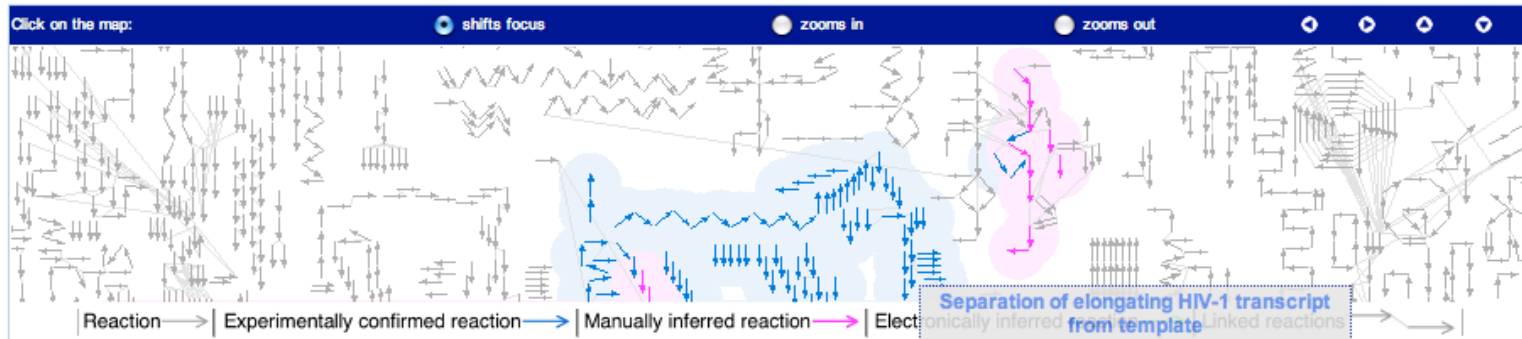
**Represents GO biological process** viral life cycle [GO](#)

### Participating molecules

- 26S protease regulatory subunit 4 [cytosol] [UK](#) [E](#) [O](#) [G](#) [R](#) [H](#) [U](#) [R](#)
- 26S protease regulatory subunit 6A [cytosol] [UK](#) [E](#) [O](#) [G](#) [R](#) [H](#) [U](#) [R](#)
- 26S protease regulatory subunit 6B [cytosol] [UK](#) [E](#) [O](#) [G](#) [R](#) [H](#) [U](#) [R](#)
- 26S protease regulatory subunit 7 [cytosol] [UK](#) [E](#) [O](#) [G](#) [R](#) [H](#) [U](#) [R](#)
- 26S protease regulatory subunit S10B [cytosol] [UK](#) [E](#) [O](#) [G](#) [R](#) [H](#) [U](#) [R](#)
- 26S proteasome non-ATPase regulatory subunit 1 [cytosol] [UK](#) [E](#) [O](#) [G](#) [R](#) [H](#) [U](#) [R](#)
- 26S proteasome non-ATPase regulatory subunit 10 [cytosol] [UK](#) [E](#) [O](#) [G](#) [R](#) [H](#) [U](#) [R](#)
- 26S proteasome non-ATPase regulatory subunit 11 [cytosol] [UK](#) [E](#) [O](#) [G](#) [R](#) [H](#) [U](#) [R](#)
- 26S proteasome non-ATPase regulatory subunit 12 [cytosol] [UK](#) [E](#) [O](#) [G](#) [R](#) [H](#) [U](#) [R](#)
- 26S proteasome non-ATPase regulatory subunit 13 [cytosol] [UK](#) [E](#) [O](#) [G](#) [R](#) [H](#) [U](#) [R](#)

# HIV Infection [Homo sapiens]

## Reactionmap



## Details

[open to selected event](#) [open all](#) [close all](#)

### HIV Infection [Homo sapiens]

- [-] HIV Life Cycle
  - [-] Early Phase of HIV Life Cycle
    - [+] Binding and entry of HIV virion
    - [+] Uncoating of the HIV Virion
      - [+] Formation of RTC (Reverse Transcriptase)
      - [+] Annealing of 3'-end of unwound transcript
    - [+] Reverse Transcription of HIV RNA
      - [+] Removal of plus-strand flap and gap closure
    - [+] Integration of provirus
  - [+] Late Phase of HIV Life Cycle
- [+] Host Interactions of HIV factors
- [+] Host immune responses to HIV infection

HIV Infection	
<b>Stable identifier</b>	REACT_6185.3
<b>Authored</b>	Bukrinsky, M, D'Eustachio, P, Gillespie, ME, Gopinathrao, G, Iordanskiy, S, Morrow, MP, Matthews, L, Rice, AP
	<p>The global pandemic of Human Immunodeficiency Virus (HIV) infection has resulted in tens of millions of people infected by the virus and millions more affected. UNAIDS estimates around 40 million HIV/AIDS patients worldwide with 75% of them living in sub-Saharan Africa. The primary method of HIV infection is by sexual exposure while nonsexual HIV transmission also can occur through transfusion with contaminated blood products, injection drug use, occupational exposure, accidental needlesticks or mother-to-child transmission. HIV damages the immune system, leaving the infected person vulnerable to a variety of "opportunistic" infections arising from host immune impairment (Hare, 2004).</p> <p>HIV-1 and the less common HIV-2 belong to the family of retroviruses. HIV-1 contains a single-stranded RNA genome that is 9 kilobases in length and contains 9 genes that encode 15 different proteins. These proteins are classified as: structural proteins (Gag, Pol, and Env), regulatory proteins (Tat and Rev), and accessory proteins (Vpu, Vpr, Vif, and Nef) (Frankel and Young, 1998).</p> <p><b>HIV infection cycle</b> can be divided into two phases:</p> <ol style="list-style-type: none"> <li>1. An <b>Early phase</b> consisting of early events occurring after HIV infection of a susceptible target cell and a</li> <li>2. <b>Late phase</b> comprising the later events in the HIV-infected cell resulting in the assembly of new infectious virions. The section titled <b>HIV lifecycle</b> consists of annotations of events in these two phases.</li> </ol> <p>The virus has developed various molecular strategies to suppress the antiviral immune responses (innate, cellular and humoral) of the host. HIV-1 viral auxiliary proteins (Tat, Rev, Nef, Vif, Vpr and Vpu) play important roles in these host-pathogen interactions (Li et al., 2005). The section titled <b>Host interactions of HIV factors</b> will highlight these complex post-infection processes and the annotations will be released in near future.</p> <p>[Hare 2004, Frankel &amp; Young 1998, Li et al 2005]</p>
<b>Organism</b>	Homo sapiens Human immunodeficiency virus 1
<b>Cellular compartment</b>	cell <a href="#">GO</a>
<b>Represents GO biological process</b>	viral life cycle <a href="#">GO</a>
<b>Participating molecules</b>	<ul style="list-style-type: none"> <li>▪ 26S protease regulatory subunit 4 [cytosol] <a href="#">UKEOGRHUR</a></li> <li>▪ 26S protease regulatory subunit 6A [cytosol] <a href="#">UKEOGRHUR</a></li> <li>▪ 26S protease regulatory subunit 6B [cytosol] <a href="#">UKEOGRHUR</a></li> <li>▪ 26S protease regulatory subunit 7 [cytosol] <a href="#">UKEOGRHUR</a></li> <li>▪ 26S protease regulatory subunit S10B [cytosol] <a href="#">UKEOGRHUR</a></li> <li>▪ 26S proteasome non-ATPase regulatory subunit 1 [cytosol] <a href="#">UKEGRHUR</a></li> <li>▪ 26S proteasome non-ATPase regulatory subunit 10 [cytosol] <a href="#">UKEOGRHUR</a></li> <li>▪ 26S proteasome non-ATPase regulatory subunit 11 [cytosol] <a href="#">UKEOGRHUR</a></li> <li>▪ 26S proteasome non-ATPase regulatory subunit 12 [cytosol] <a href="#">UKEOGRHUR</a></li> <li>▪ 26S proteasome non-ATPase regulatory subunit 13 [cytosol] <a href="#">UKEOGRHUR</a></li> <li>▪ ...</li> </ul> <p><a href="#">List all 1191 participating molecules</a></p>



EBI &gt; Databases &gt; QuickGO

[Home](#) [Help](#) [Downloads](#)  [Your selection \(0 terms\)](#) 

Search

**GO:0016032 viral reproduction**[Back](#)

The process by which a virus reproduces. Usually, this is by infection of a host cell, replication of the viral genome, and assembly of progeny virus particles. In some cases the viral genetic material may integrate into the host genome and only subsequently, under particular circumstances, 'complete' its life cycle.

[Term Information](#)[Ancestor chart](#)[Ancestor table](#)[Child Terms](#)[Protein Annotation](#)[Statistics](#)

<a href="#">i</a> <b>ID</b>	GO:0016032
<a href="#">i</a> <b>Name</b>	viral reproduction
<a href="#">i</a> <b>Definition</b>	The process by which a virus reproduces. Usually, this is by infection of a host cell, replication of the viral genome, and assembly of progeny virus particles. In some cases the viral genetic material may integrate into the host genome and only subsequently, under particular circumstances, 'complete' its life cycle.
<a href="#">i</a> <b>Comment</b>	See also the biological process terms 'viral infectious cycle ; GO:0019058' and 'lysogeny ; GO:0030069'.

[i](#) **Synonyms**

Type	Synonym
related	viral infection
related	virulence
exact	viral replication cycle
exact	viral life cycle

# Types of databases

- **Data resources of multiple types of data**
  - EBI (European Bioinformatics Institute)
  - NCBI (National Center for Biotechnology Information)
  - KEGG (Kyoto Encyclopedia of Genes and Genomes)
- **Gene and protein information**
  - GenBank, UniProt, and PDB
  - Species specific: FlyBase, dictyBase, etc.
- **Ontological data**
  - Gene Ontology
- **Pathway data**
  - KEGG PATHWAY, Reactome, BRENDA, etc.
- **Protein-protein interaction data**
  - IntAct, BioGRID, etc.



# Protein-protein interactions

- Pathways are a type of protein-protein interaction (PPI)
- High-throughput methods of determining PPI has produced a large amount of data
  - Two-hybrid
  - Microarray gene expression

# Protein-protein interaction (PPI) data

## ■ IntAct

- Proteomic database provided by the EBI
- <http://www.ebi.ac.uk/intact/>
- Interaction data is derived from the literature or via direct user submissions
- Interaction data can be searched, analyzed, and downloaded



# Database Search

Search Results for {ac=%HIV%;shortLabel=%HIV%;description=%HIV%;xref=%HIV%;disjunction=true}, (short labels of search criteria matches are **highlighted**)

Please click on any name to view more detail.



Search IntAct

**Search** Clear

- IntAct Home
- Advanced Search
- Tools
- Data Submission
- Downloads
- Documentation
  - FAQ
  - User manual
  - Annotation manual
  - Publications
  - Statistics
- Developer Resources
  - Development Site
- Contact IntAct
- Printer Friendly View

News **RSS**

2008/02/27  
**Dataset of the month**  
 You can read more about this dataset in the [Sanger Press Release](#).

2007/12/11  
**Export of protein**

EBI > Databases > Proteomic Databases

[IntAct Home](#)

## Search IntAct

To perform a search in the IntAct database, please use the following examples:

- ◆ Gene name: [BRCA2](#)
- ◆ UniProtKB Ac: [Q06609](#)
- ◆ UniProtKB Id: [dmc1](#)
- ◆ Pubmed Id: [10831611](#)

## Introduction

IntAct provides a freely available, open access resource for protein interaction data. All interactions are direct user submissions and are freely available.

## Dataset of the month: May

- ◆ **Protein-protein interactions**  
**Myxococcus xanthus.**

Whitworth et al.

Go to [Archive](#).

## License

All software, available on this site, is available under the [Creative Commons Attribution License](#). This license allows the use of all records from the IntAct database.

## Acknowledgements



IntAct is funded by the European Union (FP6-021902 (R))

[Graph](#) [Path](#) [InterPro](#) [Select all](#) [Clear all](#)

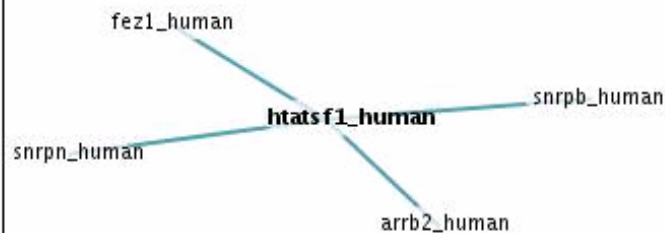
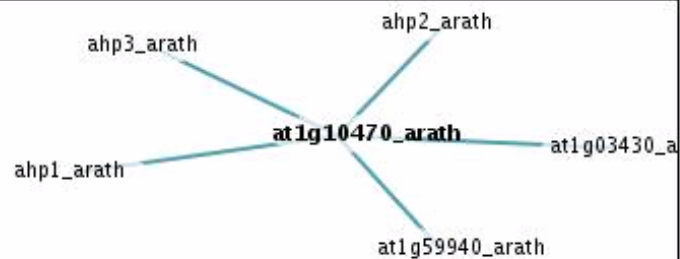
Proteins <sup>2</sup>	Name	Ac	Gene-Name	Description	Interactions	Interactors
<input type="checkbox"/>	<a href="#">q72497_9hiv1</a>	<a href="#">EBI-1036263</a>	gag	Gag polyprotein	6	2
<input type="checkbox"/>	<a href="#">q2hiv3_arath</a>	<a href="#">EBI-1100868</a>	At1g10470	At1g10470	7	5
<input type="checkbox"/>	<a href="#">q97429_9hiv1</a>	<a href="#">EBI-1183896</a>	pol	Reverse transcriptase	0	0
<input type="checkbox"/>	<a href="#">q97705_9hiv1</a>	<a href="#">EBI-1562682</a>	-	Matrix protein P17	0	0
<input type="checkbox"/>	<a href="#">p88053_9hiv1</a>	<a href="#">EBI-352318</a>	env	Envelope glycoprotein, C2-V5 region	0	0
<input type="checkbox"/>	<a href="#">q75888_9hiv1</a>	<a href="#">EBI-519201</a>	gag	Gag protein	0	0
<input type="checkbox"/>	<a href="#">o10752_9hiv1</a>	<a href="#">EBI-538834</a>	env	Envelope glycoprotein	0	0
<input type="checkbox"/>	<a href="#">q97762_9hiv1</a>	<a href="#">EBI-626690</a>	env	Envelope glycoprotein	0	0
<input type="checkbox"/>	<a href="#">q9qjd5_9hiv1</a>	<a href="#">EBI-638685</a>	env	Envelope protein	0	0
<input type="checkbox"/>	<a href="#">htsf1_human</a>	<a href="#">EBI-720468</a>	HTATSF1	HIV Tat-specific factor 1	5	265
<input type="checkbox"/>	<a href="#">q78560_9hiv1</a>	<a href="#">EBI-995253</a>	env	Envelope protein	0	0

[Graph](#) [Path](#) [InterPro](#) [Select all](#) [Clear all](#)

# HierarchView?

Query: [htatsf1\\_human](#), [gag\\_9hiv1](#), [at1g10470\\_arath](#)

Result: 13 molecules, 10 interactions.



Annotation for molecules in network. Click to highlight molecules. [?](#)

MoleculeType (1)	Confidence (0)	Publication (5)	
GO (23)	Interpro (21)	Role (5)	Species (3)

Description	Show	Count	ID
protein binding		10	<a href="#">GO:0005515</a>
RNA splicing		2	<a href="#">GO:0008380</a>
small nuclear ribonucleoprotein complex		2	<a href="#">GO:0030532</a>
nucleus		2	<a href="#">GO:0005634</a>
spliceosome		2	<a href="#">GO:0005681</a>
cytoplasm		2	<a href="#">GO:0005737</a>
transcription elongation regulator activity		1	<a href="#">GO:0003711</a>
nucleic acid binding		1	<a href="#">GO:0003676</a>
plasma membrane		1	<a href="#">GO:0005886</a>
unfolded protein binding		1	<a href="#">GO:0051082</a>
virion binding		1	<a href="#">GO:0046790</a>
identical protein binding		1	<a href="#">GO:0042802</a>
regulation of viral genome replication		1	<a href="#">GO:0045069</a>
structural molecule activity		1	<a href="#">GO:0005198</a>
RNA polymerase II transcription factor activity		1	<a href="#">GO:0003702</a>
cytokinin mediated signaling		1	<a href="#">GO:0009736</a>
cell adhesion		1	<a href="#">GO:0007155</a>
regulation of			

# Protein-protein interaction (PPI) data

## ■ BioGRID

- Biological General Repository for Interaction Datasets
- Contains protein and genetic interactions from major model organism species
- <http://www.thebiogrid.org/>
- Currently includes a virtually complete set of interactions reported in the literature for *S. cerevisiae* and *S. pombe*.

## Search the BioGRID

Examples: Genbank ID's, Entrez-Gene ID's, SGD ID's, Gene Names [\[more\]](#)

Organism:

**Submit Your Search**

Having Problems  
Searching?

Download  
*Osprey*

Osprey is a software platform for visualization of complex interaction networks. Osprey builds data-rich graphical representations from Gene Ontology (GO) annotated interaction data maintained by the BioGRID.

<http://biodata.mshri.on.ca/osprey>

## Interaction Statistics

Total Raw	208685
Total Raw Physical	143038
Total Raw Genetic	65647
Total Non-Redundant	136376
Non-Redundant Physical	93551
Non-Redundant Genetic	42825

## Database Statistics

Proteins	529018
Publications	22314
Organisms	22

## Latest News

[XML](#)

### ➤ [BioGRID version 2.0.40 release \( 3938 physical and genetic interactions added \)](#)

May. 1st, 2008 @ 02:44:20

The BioGRID's curated set of physical and genetic interactions has been updated to include an additional 3938 interactions. These additions bring our total number of non-redundant interactions to 136,376 and raw interactions to 208,685. New interactions will be added in curation updates on a monthly basis. Please let us know if we have missed or incorrectly reported any interactions by sending an e-mail to [biogridadmin@gmail.com](mailto:biogridadmin@gmail.com).

### ➤ [BioGRID version 2.0.39 release \( 980 physical and genetic interactions added \)](#)

Apr. 1st, 2008 @ 03:57:20

The BioGRID's curated set of physical and genetic interactions has been updated to include an additional 980 interactions. These additions bring our total number of non-redundant interactions to 133,595 and raw interactions to 204,747. New interactions will be added in curation updates on a monthly basis. Please let us know if we have missed or incorrectly reported any interactions by sending an e-mail to [biogridadmin@gmail.com](mailto:biogridadmin@gmail.com).

### ➤ [BioGRID version 2.0.38 release \( 354 physical and genetic interactions added \)](#)

Mar. 1st, 2008 @ 02:37:43

Your search for **"HIV"** produced the following 1 results with associations:

Click on your gene of interest for a detailed interaction summary

### Homo sapiens ( Human )



**HIVEP3**

Organism: *Homo sapiens*

# of Associations: 3

*human immunodeficiency virus type 1 enhancer-binding protein 3*|*kappabinding protein-1*;  
*human immunodeficiency virus type 1 enhancer binding protein 3* [\[more\]](#)

Your search for **"HIV"** produced the following 19 results without associations:

Click on your gene of interest for a detailed interaction summary

### Canis familiaris ( Dog )



**HIVEP2**

Organism: *Canis familiaris*

# of Associations: 0

*human immunodeficiency virus type 1 enhancer binding protein 2* [\[more\]](#)

### Gallus gallus ( Chicken )



**HIVEP1**

Organism: *Gallus gallus*

# of Associations: 0

*human immunodeficiency virus type 1 enhancer binding protein 1*; *chicken alphaA-CRYBP1*  
[\[more\]](#)



**HIVEP2**

Organism: *Gallus gallus*

# of Associations: 0

*human immunodeficiency virus type 1 enhancer binding protein 2* [\[more\]](#)



**HIVEP3**

Organism: *Gallus gallus*

# of Associations: 0

*human immunodeficiency virus type 1 enhancer binding protein 3* [\[more\]](#)

### Homo sapiens ( Human )

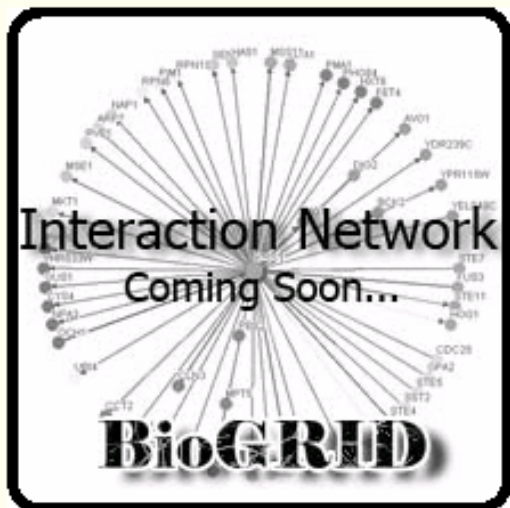


**HIVEP2**

Organism: *Homo sapiens*

# of Associations: 0

*MHC binding protein-2*|*OTTHUMP00000040180*|*Schnurri-2*|*c-myc intron-binding protein 1*|*human immunodeficiency virus type 1 enhancer-binding protein 2*; *human immunodeficiency virus* [\[more\]](#)



[\[JPG\]](#) [\[PNG\]](#) [\[SVG\]](#)

## HIVEP3

*Homo sapiens*

**Aliases:** KBP-1, KBP1, KIAA1555, KRC, ZAS3, FLJ16752, EG59269

**Description:** human immunodeficiency virus type I enhancer-binding protein 3|kappabinding protein-1; human immunodeficiency virus type I enhancer binding protein 3

**Links:** [\[Entrez-Gene\]](#) [\[OMIM\]](#) [\[HGNC\]](#) [\[HPRD\]](#) [\[Ensembl\]](#)

**Gene Ontology:** 4 total GO mappings for this record. [\[view list\]](#)

[Download interactions associated with HIVEP3](#)

### HIVEP3 was identified with 5 protein interactions

Name	Aliases	Description	Evidence Code(s)	Role	Source(s)
<a href="#">TRAF2</a>	<ul style="list-style-type: none"> <li><i>MGC:45012</i></li> <li><i>TRAP</i></li> <li><i>TRAP3</i></li> <li><i>EG7186</i></li> </ul>	OTTHUMP00000064745 tumor necrosis factor type 2 receptor associated protein 3; TNF receptor-associated factor 2 Gene Ontology: <a href="#">[View List]</a>	Invivo	hit	<a href="#">Oukka M (2002)</a>
			Two-hybrid	hit	<a href="#">Oukka M (2002)</a>
<a href="#">TRAF1</a>	<ul style="list-style-type: none"> <li><i>EBI6</i></li> <li><i>MGC:10353</i></li> <li><i>EG7185</i></li> </ul>	Epstein-Bar virus-induced protein 6; TNF receptor-associated factor 1 Gene Ontology: <a href="#">[View List]</a>	Invivo	hit	<a href="#">Oukka M (2002)</a>
			Two-hybrid	hit	<a href="#">Oukka M (2002)</a>
<a href="#">TNFRSF1B</a>	<ul style="list-style-type: none"> <li><i>CD120b</i></li> <li><i>TBPII</i></li> <li><i>TNF-R-II</i></li> <li><i>EG7133</i></li> </ul>	p75 TNF receptor tumor necrosis factor beta receptor tumor necrosis factor binding protein 2 tumor necrosis factor receptor 2; tumor necrosis factor receptor superfamily, member 1B Gene Ontology: <a href="#">[View List]</a>	Invivo	hit	<a href="#">Oukka M (2002)</a>



# Summary (1)

- Many major databases, with many useful information
- For a specific purpose, a single database may suffice
- For more comprehensive analyses, computer programs may be developed to access multiple databases and process the data to retrieve the desired information
  - Application Programming Interfaces (APIs), Web services, etc.

# Summary (2)

- **APIs**: programs (libraries) that can be downloaded such that users can develop computer programs that use the provided libraries to retrieve data via the Internet
- **Web services**: programs that can be executed over the Internet by users' programs to retrieve information
  - **Workflows** can be developed from these web services, such that the retrieved data from DB 1 can be automatically sent as an input query to DB 2, for example.