

Bioinformatics

Genome information resources

Bioinformatics - lectures

- Introduction
- Information networks
- Protein information resources
- **Genome information resources**
- DNA sequence analysis
- Pairwise sequence alignment
- Multiple sequence alignment
- Secondary database searching
- Analysis packages
- Protein structure modelling

Genome information resources

- primary DNA sequence databases
- specialised DNA sequence databases

Primary DNA sequence databases

- EMBL
- DDBJ
- GenBank
- dbEST
- GSDB

Store DNA sequences and annotations.

Primary protein sequence databases

- **EMBL - European Molecular Biology Laboratory**
 - European Bioinformatics Institute (EBI)
 - collaboration with DDBJ and GenBank - exchange of new entries on daily basis
 - source of sequences: **direct author submissions, genome projects, scientific literature, patents**
 - rate of growth is exponential with doubling time ~9-12 months
 - most entries from model organisms
 - retrieval through SRS

Primary protein sequence databases

■ DDBJ - DNA Data Bank of Japan

- National Institute of Genetics
- collaboration with EMBL and GenBank
- retrieval through DDBGet

■ GenBank

- National Center for Biotechnology Information (NCBI)
- collaboration with DDBJ and EMBL
- data split into **17 divisions**
- retrieval through Entrez

Codes for 17 divisions of GenBank

<i>Division</i>	<i>Sequence subset</i>
PRI	Primate
ROD	Rodent
MAM	Other mammalian
VRT	Other vertebrate
INV	Invertebrate
PLN	Plant, fungal, algal
BCT	Bacterial
RNA	Structural RNA
VRL	Viral
PHG	Bacteriophage
SYN	Synthetic
UNA	Unannotated
EST	EST (Expressed Sequence Tags)
PAT	Patent
STS	STS (Sequence Tagged Sites)
GSS	GSS (Genome Survey Sequences)
HTG	HTG (High Throughput Genomic Sequences)

LOCUS DRODPPC 4001 bp mRNA INV 15-MAR-1990
 DEFINITION D.melanogaster decapentaplegic gene complex (DPP-C), complete cds.
 ACCESSION M30116
 NID g157291
 KEYWORDS .
 SOURCE D.melanogaster, cDNA to mRNA.
 ORGANISM Drosophila melanogaster
 Eukaryotae; mitochondrial eukaryotes; Metazoa; Arthropoda;
 Tracheata; Insecta; Pterygota; Diptera; Brachycera; Muscomorpha;
 Ephydroidea; Drosophilidae; Drosophila.
 REFERENCE 1 (bases 1 to 4001)
 AUTHORS Padgett,R.W., St Johnston,R.D. and Gelbart,W.M.
 TITLE A transcript from a Drosophila pattern gene predicts a protein
 homologous to the transforming growth factor-beta family
 JOURNAL Nature 325, 81-84 (1987)
 MEDLINE 87090408

COMMENT The initiation codon could be at either 1188-1190 or 1587-1589.
 FEATURES Location/Qualifiers
 source 1..4001
 /organism="Drosophila melanogaster"
 /db_xref="taxon:7227"
 mRNA <1..3918
 /gene="dpp"
 /note="decapentaplegic protein mRNA"
 /db_xref="FlyBase:FBgn0000490"
 gene 1..4001
 /note="decapentaplegic"
 /gene="dpp"
 /allele=""
 /db_xref="FlyBase:FBgn0000490"
 CDS 1188..2954
 /gene="dpp"
 /note="decapentaplegic protein (1188 could be 1587)"
 /codon_start=1
 /db_xref="FlyBase:FBgn0000490"
 /db_xref="PID:g157292"
 /translation="MRAWLLLLAVLATFQTIVRVASTEDISQRFIAAIAPVAAHIPLA
 SASGSGSGRSGRSRSGASTSTALAKAFNPFSEPASFSDDSKSHRSKTNKKPSKSDANR

 LGYDAYYCHGKCPFPLADHFNSTNHAVVQTLVNNMNPVKVPKACCVPTQLDSVAMLYL
 NDQSTVVLKKNYQEMTVVGCGCR"

BASE COUNT 1170 a 1078 c 956 g 797 t
 ORIGIN

```

1 gtcgttcaac agcgctgata gagtttaaat ctataccgaa atgagcggcg gaaagtgagc
61 cacttggcgt gaacccaaag ctttcgagga aaattctcgg acccccatat acaaatatcg
121 gaaaaagtat cgaacagttt cgcgacgcga agcgtaaga tcgccaaaag atctccgtgc
181 ggaaacaaag aaattgaggc actattaaga gattgttgtt gtgcgcgagt gtgtgtcttc
241 agctgggtgt gtggaatgtc aactgacggg ttgtaaaggg aaaccctgaa atccgaacgg
301 ccagccaaag caataaagc tgtgaatacg aattaagtac acaaacagt tactgaaaca
361 gatacagatt cggattcgaa tagagaaaca gatactggag atgccccag aaacaattca
421 attgcaaata tagtgcgttg cgcgagtgcc agtggaaaaa tatgtggatt acctgcgaac
481 cgtccgcca aggagccgcc gggtgacagg tgtatcccc aggataccaa cccgagcca
541 gaccgagatc cacatccaga tcccgaccgc agggtgccag tgtgtcatgt gccgcggcat
601 accgaccgca gccacatcta ccgaccaggt gcgcctcgaa tgcggaaca caattttcaa
.....
3841 aactgtataa acaaaacgta tgccctataa atatatgaat aactatctac atcgttatgc
3901 gttctaagct aagctcgaat aaatccgtac acgttaatta atctagaatc gtaagaccta
3961 acgcgtaagc tcagcatggt ggataaatta atagaaacga g
  
```


Primary protein sequence databases

■ dbEST

- National Center for Biotechnology Information (NCBI)
- maintains only **Expressed Sequence Tag (EST)** data

■ GSDB - Genome Sequence DataBase

- National Center for Genome Resources
- complete collections of DNA sequence for **genome-sequencing** laboratories
- on-line submission of large-scale data
- quality checks
- format consistent with GenBank + **GSDBID**

Specialised DNA sequence databases

- SGD
- UniGene
- TDB
- ACeDB

Store species-specific and technique-specific DNA sequences.

Specialised DNA sequence databases

■ SGD - *Saccharomyces* Genome Database

- molecular biology and genetics of *S. cerevisiae*
- complete genome, genes, proteins, phenotypes
- first eukaryotic genome sequenced (1998)
- sequence analysis, register of genes, 3D structural data, primer sequences for cloning

■ UniGene

- collection of genes encoding proteins (**transcript map**)
- non-redundant; derived from GenBank
- data organised in clusters (1 cluster = 1 unique gene)
- gene-mapping projects and gene expression analysis

Specialised DNA sequence databases

■ TDB - TIGR Database

- suite of databases: DNA and protein sequences, gene expression, protein families, taxonomic data
- links: TIGR microbial genome sequencing projects, parasite databases, gene index projects, *A. thaliana* database, human genomic dataset

■ ACeDB - A *Cernorhabditis elegans* DataBase

- *C. elegans* genome project
- restriction maps, gene structural information, cosmid maps, sequence data, bibliographic information
- software to organise data ACEDB: CGI script and perl

Search:

In Class: **Ready**

Genetic_map	Sequence_map	Author
Locus	Clone	Paper
Gene	Clone_Grid	
Other_Locus	Sequence	Expr_pattern
Rearrangement		Cell
Allele		Pathway
Strain		
Gene_Class	KeySet	Model

Global Search:

Main KeySet

Show As.:

Query..

- 2056 items 1 selected
- | | | | | |
|---------|-------|--------|-------|--------|
| abl-1 | aex-3 | arl-1 | bli-3 | cat-4 |
| ace-1 | aex-4 | arl-2 | bli-4 | cat-5 |
| ace-2 | aex-5 | arl-3 | bli-5 | cat-6 |
| ace-3 | aex-6 | arl-4 | bli-6 | cct-1 |
| acr-2 | age-1 | arl-5 | bor-1 | cct-2 |
| act-2 | ale-1 | arp-1 | brp-1 | cct-4 |
| act-3 | ali-1 | art-1 | cad-1 | cct-5 |
| act-4 | ana-1 | atn-1 | cah-1 | cct-6 |
| act-5 | ana-2 | avr-15 | cal-1 | cdc-42 |
| act-123 | anc-1 | bac-1 | cal-5 | cdh-1 |
| adr-1 | aph-1 | bar-1 | can-1 | cdh-2 |
| adr-1 | aph-2 | bas-1 | cap-1 | cdh-3 |
| adr-2 | apl-1 | bas-2 | cap-2 | cdk-5 |
| adp-1 | apx-1 | ben-1 | cat-1 | ced-1 |
| aex-1 | arf-1 | bli-1 | cat-2 | ced-2 |
| aex-2 | arf-3 | bli-2 | cat-3 | ced-3 |

