

Representation of sequence – sequence file format

1.) FASTA – simple and clean

> gene_name, (other info)

MASASASKJHKLJLKJLDSDFSF

SSDSASFSD...

Practice / DIY: retrieve sequence in Fasta format and save the file in the local computer.

How to store sequence files

- Pure text format is clean and allows downstream sequence analysis.
- .doc or .rtf allows formatting during annotation – however, extra information are inserted thus NOT suitable for computational analysis.
- **!! No space in file or folder name!!**- Or trouble will find you.

Practice – file types

- Using file Finder (Mac) or Explorer (PC) to view downloaded files.
- Change the “Options” so that the file **extensions (.txt) are revealed.**
- Edit the downloaded sequence file in MS Word, highlight a section of the sequence with Bold font or color and save as .doc
- Open the .doc file in a TextEditor – observe the inserted characters.

Practice – file types (Cont.)

- Load the “Mysequence.doc” file to Webcutter using “Choose file” and then “Upload sequence file”.
 - Notice that the “sequence” in the sequence box are nonsense characters.
- Clear input; Browse and then load the .txt file. Run an analysis.

Keep you sequences as .seq or .fasta file for downstream analysis.

Public Resources for Bioinformatics

- **Databases**

- **Analysis Tools**

Observe: List of databases and service at NCBI, EBI, KEGG, and Ensembl.

Public Resources (I) – Databases and data sources

- **Over 1,000 in the sea of databases.**
- **Content-specific, such as DNA, Protein, Structure, etc.**
- **Species-specific, such as flybase, wormbase, OMIM, etc.**
- **System-specific, such as MetaCyc, AFCS, etc.**

Pet Project:

IL-6, or your favorite gene

What can we know about this gene?

- **Search for “curated” databases.**
- **To prepare for future analysis, save annotated sequence files as `genename.html` (in a target folder).**
- **For downstream sequence analysis, save pure sequence as FASTA format file.**

Where and how much information are available for my gene?

Observe: The information contents and presentation format for the same gene in SwissProt, NCBI protein, NCBI Genes, etc..

Observe/Practice – basic DB search skills

Search for IL6 (or your favorite gene) in the NCBI Gene or Proteins databases.

- Why do we get so many hits?

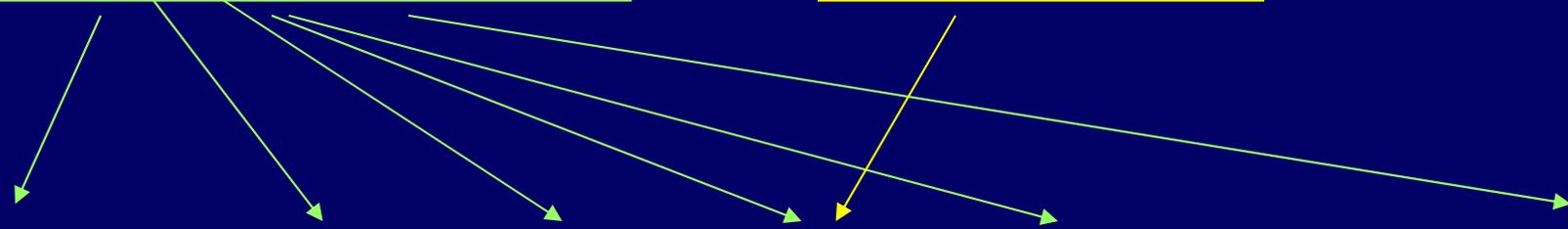
Search for IL6 in the default “All Text” vs. search in the [Gene Name] field only in the Gene database.

Compare results.

Database concept – tables in relational databases

“IL6”=IL6[All Fields]

IL6[Name]



Accession	Organ.	Ref.	Name	Key words	Features	
....	medline1	TNF
....	medline2	P53

Gene table

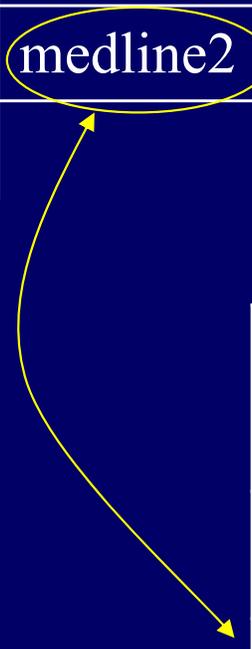
Database concept – relationship between tables allows linkage

Accession	Organ.	Ref.	Name	Key words	Features	
....	medline1	P27
....	medline2	P53

Protein table

ID	title	year	author	abstract	
medline1	1970
medline2	1980

Reference table



Observe/Practice – content specific DBs

Search for your favorite gene at:

- NCBI nucleotide
- NCBI protein
- NCBI gene
- GEO
- Etc.

Observe/Practice – the Gene entry

Observe the links from the IL6 **Gene** page:

- RefSeq
- OMIM
- SNP
- GEO
- Etc.

The “Gene” entry is the pivotal point for many NCBI resources.

Representation of genes and related information

The need to represent associated info with sequence

- Different aspects of the gene (such as protein, nucleotide, structure (PDB), OMIM etc.)
- Specialized databases (such GEO, SNP)
- Complex / customized data structure
 - Object-oriented data representation

Observe

Observe entries involving IL6 (or your gene)
in Reactome.