Slide 1:

Spring 2024

**BIOS 477/877**

*Bioinformatics and Molecular Evolution*

**Lecture 2**

BIOS477/877 L2 - 1

---

Slide 2:

**TODAY'S TOPICS**

➢ **Introduction to Internet Resources (NCBI, databases)**

➢ **Assignment #1**

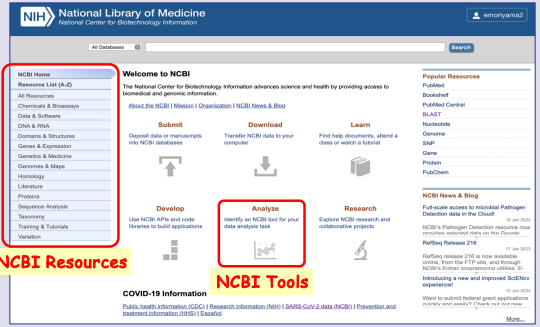BIOS477/877 L2 - 2

---

Slide 3:

**General Web Search**

DO NOT blindly believe what you find on internet (*e.g.,* Google, Wikipedia)!
Lots of information are incorrect. Misinformation propagates rapidly!
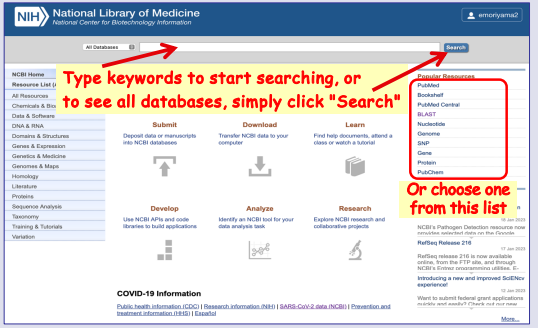Always double-check with the original information source (journal articles).

BIOS477/877 L2 - 3

---

Slide 4:

**National Center for Biotechnology Information (NCBI)**
https://www.ncbi.nlm.nih.gov

NCBI Resources

NCBI Tools

BIOS477/877 L2 - 4

---

Slide 5:

**National Center for Biotechnology Information (NCBI)**
https://www.ncbi.nlm.nih.gov

Type keywords to start searching, or to see all databases, simply click "Search"

Or choose one from this list

BIOS477/877 L2 - 5

---

Slide 6:

**NCBI Databases**
https://www.ncbi.nlm.nih.gov/search

Literature search

Literature

Many online books

Scroll down to find more databases

BIOS477/877 L2 - 6

1

## Slide 7

### NCBI Databases
https://www.ncbi.nlm.nih.gov/search

**Data**

**Genes**
Gene sequences and annotations used as references for the study of orthologs structure, expression, and evolution
**Gene**
Collected information about gene loci
**GEO DataSets**
Functional genomics studies
**GEO Profiles**
Gene expression and molecular abundance profiles
**HomoloGene**
Homologous genes sets for selected organisms
**PopSet**
Sequence sets from phylogenetic and population studies

**GEO (Gene Expression Omnibus): microarray & RNA-Seq data**

**Proteins**
Protein sequences, 3-D structures, and tools for the study of functional protein domains and active sites
**blastp**
**Conserved Domains**
**Protein Family Models**
Models representing homologous proteins with a common function
**Structure**
Experimentally-determined biomolecular structures

**BLAST**
A tool to find regions of similarity between biological sequences
**blastn**
Search nucleotide sequence databases
**tblastn**
Search translated nucleotide databases using a translated

**Genomes**
Genome sequence assemblies, large-scale functional genomics data, and source biological samples
**Assembly**
Genome assembly information
**BioCollections**
Museum, herbaria, and other biorepository collections
**BioProject**
Biological projects providing data to NCBI
**BioSample**
Descriptions of biological source materials
**Genome**
Genome sequencing projects by organism
**Nucleotide**
DNA and RNA sequences
**SRA**
High throughput sequence reads
**Taxonomy**
Taxonomic classification and nomenclature

**Clinical**
Heritable DNA variations, associations with human pathologies, and clinical diagnostics and treatments
**ClinicalTrials.gov**
Privately and publicly funded clinical studies conducted around the world
**ClinVar**
Human variations of clinical significance
**dbGaP**
Genotype/phenotype interaction studies
**dbSNP**
Short genetic variations

**PubMed**
Repository of chemical information, molecular pathways, and tools for bioactivity screening
**BioAssays**
Bioactivity screening studies
**Compounds**
Chemical information with structures, information and links
**Pathways**
**Substances**
Deposited substance and chemical information

**SRA (Sequence Read Archive): next-gen sequencing data (454, Illumina, PacBio, Nanopore, etc.)**

**Genome & DNA sequences**

BIOS477/877 L2 - 7

---

## Slide 8

### NCBI Databases
https://www.ncbi.nlm.nih.gov/search

**Data**

**Genes**
Gene sequences and annotations used as references for the study of orthologs structure, expression, and evolution
**Gene**
Collected information about gene loci
**GEO DataSets**
Functional genomics studies
**GEO Profiles**
Gene expression and molecular abundance profiles
**HomoloGene**
Homologous genes sets for selected organisms
**PopSet**
Sequence sets from phylogenetic and population studies

**Proteins**
Protein sequences, 3-D structures, and tools for the study of functional protein domains and active sites
**Conserved Domains**
Find conserved protein domains
**Identical Protein Groups**
Protein sequences grouped by identity
**Protein**
Protein sequences
**Protein Family Models**
Models representing homologous proteins with a common function
**Structure**
Experimentally-determined biomolecular structures

**Protein sequences & structures**

**BLAST**
A tool to find regions of similarity between biological sequences
**blastn**
**tblastn**
Search translated nucleotide databases using a translated
**Primer-BLAST**
Find primers specific to your PCR template

**Genomes**
Genome sequence assemblies, large-scale functional genomics data, and source biological samples
**Assembly**
Genome assembly information
**BioCollections**
Museum, herbaria, and other biorepository collections
**BioProject**
Biological projects providing data to NCBI
**BioSample**
Descriptions of biological source materials
**Genome**
Genome sequencing projects by organism
**Nucleotide**
DNA and RNA sequences
**SRA**
High throughput sequence reads
**Taxonomy**
Taxonomic classification and nomenclature

**Clinical**
Heritable DNA variations, associations with human pathologies, and clinical diagnostics and treatments
**ClinicalTrials.gov**
Privately and publicly funded clinical studies conducted around the world
**ClinVar**
Human variations of clinical significance
**dbGaP**
Genotype/phenotype interaction studies
**dbSNP**
Short genetic variations
**dbVar**
Genome structural variation studies
**GTR**
Genetic testing registry
**MedGen**
Medical genetics literature
**OMIM**
Online mendelian inheritance in man

**PubMed**
Repository of chemical information, molecular pathways, and tools for bioactivity screening
**BioAssays**
Bioactivity screening studies
**Compounds**
Chemical information with structures, information and links
**Pathways**
Molecular pathways with links to genes, proteins and chemicals
**Substances**
Deposited substance and chemical information

**Biological systems & pathways**

BIOS477/877 L2 - 8

---

## Slide 9

### PubMed: Literature Search
https://pubmed.ncbi.nlm.nih.gov

**NIH** National Library of Medicine
National Center for Biotechnology Information

Log in

**PubMed.gov**

bioinformatics ← **Type a query**    **Start searching** → **Search**
Advanced

PubMed® comprises more than 35 million citations for biomedical literature from MEDLINE, life science journals, and online books. Citations may include links to full text content from PubMed Central and publisher web sites.

BIOS477/877 L2 - 9

---

## Slide 10

### PubMed: Literature Search
https://pubmed.ncbi.nlm.nih.gov

**PubMed®**    bioinformatics    ✕    **Search**
Advanced Create alert Create RSS    User Guide

Save    Email    Send to    Sort by: Best match    Display options ⚙

**Advanced search**

MY NCBI FILTERS    ‹‹  ‹  Page 1 of 50,880  ›  ››

RESULTS BY YEAR
1958 — 2024

**Bioinformatics** in translational drug discovery.
Wooller SK, Benstead-Hume G, Chen X, Ali Y, Pearl FMG.
Biosci Rep. 2017 Jul 7;37(4):BSR20160180. doi: 10.1042/BSR20160180. Print 2017 Aug 31.
PMID: 28487472    Free PMC article.    Review.
**Bioinformatics** approaches are becoming ever more essential in translational drug discovery both in academia and within the pharmaceutical industry. ...Here, we highlight some of the areas in which **bioinformatics** resources and methods are being developed to support t ...

TEXT AVAILABILITY
☐ Abstract
☐ Free full text
☐ Full text

Is "**bioinformatics**" dead?
Bourne PE.
PLoS Biol. 2021 Mar 18;19(3):e3001165. doi: 10.1371/journal.pbio.3001165. eCollection 2021 Mar.
PMID: 33736179    Free PMC article.
Why would a **computational** biologist with 40 years of research experience say **bioinformatics** is dead? ...

ARTICLE ATTRIBUTE
☐ Associated data

ARTICLE TYPE
☐ Books and Documents
☐ Clinical Trial
☐ Meta-Analysis
☐ Randomized Controlled Trial
☑ Review
☐ Systematic Review

**Aptamer Bioinformatics**.
Kinghorn AB, Fraser LA, Lang S, Shiu SCC, Tanner JA.
Int J Mol Sci. 2017 Nov 24;18(12):2516. doi: 10.3390/ijms18122516.
PMID: 29186809    Free PMC article.    Review.
They are isolated via SELEX (Systematic Evolution of Ligands by Exponential Enrichment), an evolutionary process that involves iterative rounds of selection and amplification before sequencing and aptamer characterization. As aptamers are genetic in nature, **bioinformatic**s a ...

**To see only reviews**

...ence and bioengineering: Recent advances, ...ctives.
...K, Kojima T, Tsugawa H, Toda Y, Horinouchi T.
J Biosci Bioeng. 2022 Nov;134(5):363–373. doi: 10.1016/j.jbiosc.2022.08.004. Epub 2022 Sep 17.

BIOS477/877 L2 - 10

---

## Slide 11

### PubMed: Literature Search

**PubMed Advanced Search Builder**

**PubMed®**
User Guide

Add terms to the query box

All Fields ⬥  [Enter a search term]    **ADD** ▾
                                      Show Index

**Type new query(s) here**

Query box
[Enter / edit your search query here]    **Search** ▾

**History and Search Details**    ⬇ Download    🗑 Delete

| Search | Actions | Details | Query | Results | Time |
|--------|---------|---------|-------|---------|------|
| #1 | ••• | › | Search: **bioinformatics** | 508,793 | 04:42:16 |

Showing 1 to 1 of 1 entries

**Click the number to see the result**

https://pubmed.ncbi.nlm.nih.gov/

BIOS477/877 L2 - 11

---

## Slide 12

### PubMed: Literature Search

**PubMed Advanced Search Builder**

**PubMed®**
User Guide

Add terms to the query box

All Fields ⬥  gene prediction    **ADD** ▾
                                Show Index

**new query**    **Start searching**

Query box
[Enter / edit your search query here]    **Search** ▾
                                         Add to History

**Or show only the number of hits without doing the actual search**

**History and Search Details**    ⬇ Download    🗑 Delete

| Search | Actions | Details | Query | Results | Time |
|--------|---------|---------|-------|---------|------|
| #2 | ••• | › | Search: **gene prediction** | 268,771 | 04:45:23 |
| #1 | ••• | › | Search: **bioinformatics** | 508,793 | 04:42:16 |

https://pubmed.ncbi.nlm.nih.gov/

BIOS477/877 L2 - 12

13



14



15



16



17



18

3

19



20



21



22



23



24

25



26



27



28



29



30

5

**Nucleotide Database Search**

31



**GenBank Entry**

32



**GenBank Entry**

33



**GenBank Entry**

34



**GenBank Entry**

35



**GenBank Entry**

36

## GenBank Entry

```
                Pterygota; Neoptera; Endopterygota; Diptera; Brachycera.
                Muscomorpha; Ephydroidea; Drosophilidae; Drosophila; Sophophora.
REFERENCE       1
  AUTHORS       Aguade,M.
  TITLE         Nucleotide and copy-number polymorphism at the odorant receptor
                genes Or22a and Or22b in Drosophila melanogaster
  JOURNAL       Mol. Biol. Evol. 26 (1), 61-70 (2009)
  PUBMED        18922763
REFERENCE       2  (bases 1 to 4518)
  AUTHORS       Aguade,M.
  TITLE         Direct Submission
  JOURNAL       Submitted (22-AUG-2008) Aguade M., Genetics, Universitat de
                Barcelona, Diagonal 645, 08028, SPAIN
FEATURES             Location/Qualifiers
     source          1..4518
                     /organism="Drosophila simulans"
                     /mol_type="genomic DNA"
                     /strain="MO"
                     /db_xref="taxon:7240"
     gene            <364..>1742
                     /gene="or22a"
     CDS             join(364..491,552..1190,1257..1614,1674..1742)
                     /gene="or22a"
                     /codon_start=1
                     /product="odorant receptor 22a"
                     /protein_id="CAR79115.1"
                     /db_xref="GOA:B4Q7S7"
                     /db_xref="InterPro:IPR004117"
                     /db_xref="UniProtKB/TrEMBL:B4Q7S7"
                     /translation="MLCKFPPHISKKPLSERVHSRDAFIYLDRVMNSFGWTEPENRHW
                     VLPYNLWFALVHIVHLILLPISMSIEYLRHFKTFSAGEFLSSLEIGVNIYGSSFKCAF
                     TMMGFKKRQEAKVLLDQLDKRCVRDEERSTVHRYVAWGNFFDILYHIFYSSFVNHFP
                     YFLLQRHHAWRNYFPYIDPEKQFFISSIABCFLMTEVIYMDLCTDVCPLISMLKARCH
                     ISLLKQRLRNLRSEPGRTEDEYLEELTECIRDHKLILDYVDALRPVFSGTIFVQFLLI
                     GTVLGLSHINLMFFSTFWTGVATCLFMFDVSMETFPFCYLCNMIIDDCQEMADSLFQS
```

Find the corresponding protein sequence

Related information
Gene
Protein
PubMed
PubMed (Weighted)
Taxonomy

BIOS477/877 L2 - 37

37

---

## GenPept Entry (tranlasted GenBank)

GenPept

odorant receptor 22b [Drosophila simulans]
GenBank: CAR79116.1
Identical Proteins   FASTA   Graphics

```
LOCUS       CAR79116           397 aa     linear   INV 17-DEC-2008
DEFINITION  odorant receptor 22b [Drosophila simulans].
ACCESSION   CAR79116
VERSION     CAR79116.1
DBSOURCE    embl accession FM212179.1
KEYWORDS    .
SOURCE      Drosophila simulans
  ORGANISM  Drosophila simulans
            Eukaryota; Metazoa; Ecdysozoa; Arthropoda; Hexapoda; Insecta;
            Pterygota; Neoptera; Holometabola; Diptera; Brachycera;
            Muscomorpha; Ephydroidea; Drosophilidae; Drosophila; Sophophora.
REFERENCE   1
  AUTHORS   Aguade,M.
  TITLE     Nucleotide and copy-number polymorphism at the odorant receptor
            genes Or22a and Or22b in Drosophila melanogaster
  JOURNAL   Mol. Biol. Evol. 26 (1), 61-70 (2009)
  PUBMED    18922763
REFERENCE   2  (residues 1 to 397)
  AUTHORS   Aguade,M.
  TITLE     Direct Submission
  JOURNAL   Submitted (22-AUG-2008) Aguade M., Genetics, Universitat de
            Barcelona, Diagonal 645, 08028, SPAIN
FEATURES             Location/Qualifiers
     source          1..397
                     /organism="Drosophila simulans"
                     /strain="MO"
                     /db_xref="taxon:7240"
     Protein         1..397
                     /product="odorant receptor 22b"
     Region          81..381
                     /region_name="7tm_6"
                     /note="7tm Odorant receptor; pfam02949"
                     /db_xref="CDD:251636"
     CDS             1..397
                     /gene="or22b"
                     /coded_by="join(FM212179.1:2521..2648,...
                     FM212179.1:2700..3347,FM212179.1:3438..3795,...
```

Find the corresponding DNA sequence

Find Conserved Domains

Analyze this sequence
Run BLAST
Identify Conserved Domains
Highlight Sequence Features
Find in this Sequence

Related information
Nucleotide
PubMed
Taxonomy

CDD Search Results
Conserved Domains (Concise)
Conserved Domains (Full)

Recent activity
Turn Off  Clear
odorant receptor 22b [Drosophila ...
Protein Links for Nucleotide (Select 218317908) (2)
Drosophila simulans or22a gene for odorant receptor 22a and or22...
odorant receptor 22b [Drosophila melanogaster]

https://www.ncbi.nlm.nih.gov/search/

BIOS477/877 L2 - 38

38

---

## Sequence Formats

➢ **FASTA: a standard format used in many programs**

```
>AAK58030.1 glycoprotein coupled receptor [Human betaherpesvirus 5]
MTPTTTTAELTTEFDYDEDATPCVFTDVLNQSKPVTLFLYGVVFLFGSIGNFLVIFTITWRRRIQCSGDV
YFINLAAADLLFVCTLPLWMQYLLDHNSLASVPCTLLTACFYVAMFASLCFITEIALDRYYAIVYMRYRP
VKQACLFSIFWWIFAVIIAIPHFMVVTKKDNQCMTDYDYLEVSYPIILNVELMLGAFVIPLSVISYCYYR
ISRIVAVSQSRHKGRIVRVLIAVVLVFIIFWLPYHLTLFVDTLKLLKWISSSCEFERSLKRALILTESLA
FCHCCLNPLLYVFVGTKFRQELHCLLAEFRQRLFSRDVSWYHSMSFSRRSSPSRRETSSDTLSDEVCRVS
QIIP
```

• **A sequence begins with a single-line description, followed by lines of sequence data**

• **Up to the first space is usually considered as the sequence id**

• **The description line is distinguished from the sequence data by a greater-than (>) symbol**

BIOS477/877 L2 - 39

39

---

## FASTA Formats

**[Single sequence FASTA file]**
```
>P29627
FGNLSSAQAIMGNPRIRAHGKKVLTSLGLAVQNMDNL
KETFAHLSELHCDKLHVDPENFKLLGNVLVIVLSTHF
AKEFTPEVQAAWQKLVAGVANALSHKYH
```

**[Multiple sequence FASTA file]**
```
>HBB_HORSE
VQLSGEEKAAVLALWDKVNEEEVGGEALGRLLVVYPWT
QRFFDSFGDLSNPGAVMGNPKVKAHGKKVLHSFGEGVH
HLDNLKGTFAALSELHCDKLHVDPENFRLLGNVLVVVL
ARHFGKDFTPELQASYQKVVAGVANALAHKYH
>HBA_HUMAN
VLSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPT
TKTYFPHFDLSHGSAQVKGHGKKVADALTNAVAHVDDM
PNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHLP
AEFTPAVHASLDKFLASVSTVLTSKYR
>HBA_HORSE
VLSAADKTNVKAAWSKVGGHAGEYGAEALERMFLGFPT
TKTYFPHFDLSHGSAQVKAHGKKVGDALTLAVGHLDDL
PGALSNLSDLHAHKLRVDPVNFKLLSHCLLSTLAVHLP
NDFTPAVHASLDKFLSSVSTVLTSKYR
>GLB5_PETMA
PIVDTGSVAPLSAAEKTKIRSAWAPVYSTYETSGVDIL
VKFFTSTPAAQEFFPKFKGLTTADQLKKSADVRWHAER
IINAVNDAVASMDDTEKMSMKLRDLSGKHAKSFQVDPQ
YFKVLAAVIADTVAAGDAGFEKLMSMICILLRSAY
```

**[Alignment FASTA file]**
```
>HBB_HORSE
---------VQLSGEEKAAVLALWDKVN--EEEVGGEALGRLLV
VYPWTQRFFDSFGDLSNPGAVMGNPKVKAHGKKVLHSFGEGVH
HLDN-----LKGTFAALSELHCDKLHVDPENFRLLGNVLVVVL
ARHFGKDFTPELQASYQKVVAGVANALAHKYH------
>HBA_HUMAN
---------VLSPADKTNVKAAWGKVGAHAGEYGAEALERMFL
SFPTTKTYFPHF-DLS------HGSAQVKGHGKKVADALTNAVA
HVDD-----MPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTL
AAHLPAEFTPAVHASLDKFLASVSTVLTSKYR------
>HBA_HORSE
---------VLSAADKTNVKAAWSKVGGHAGEYGAEALERMFL
GFPTTKTYFPHF-DLS------HGSAQVKAHGKKVGDALTLAVG
HLDD-----LPGALSNLSDLHAHKLRVDPVNFKLLSHCLLSTL
AVHLPNDFTPAVHASLDKFLSSVSTVLTSKYR------
>GLB5_PETMA
PIVDTGSVAPLSAAEKTKIRSAWAPVYSTYETSGVDILVKFFT
STPAAQEFFPKFKGLTTADQLKKSADVRWHAERIINAVNDAVA
SMDDT--EKMSMKLRDLSGKHAKSFQVDPQYFRVLAAVIADTV
AAG---------DAGFEKLMSMICILLRSAY-------
```

BIOS477/877 L2 - 40

40

---

## Sequence Databases: DNA

➢ **INSDC: International Nucleotide Sequence Database Collaboration**

• **GenBank: NIH genetic sequence database**
  https://www.ncbi.nlm.nih.gov/genbank/

• **ENA: European Nucleotide Archive**
  https://www.ebi.ac.uk/ena/browser/home

• **DDBJ: DNA Data Bank of Japan**
  https://www.ddbj.nig.ac.jp/index-e.html

➜ **Data are shared among the three databases**

➜ **They share the accession numbers!**

➜ **They use different formats**

BIOS477/877 L2 - 41

41

---

## Sequence Databases: DNA

### EMBL Format

```
ID   X64178; SV 1; linear; mRNA; STD; ROD; 429 BP.
XX
AC   X64178;
XX
DT   26-APR-1992 (Rel. 31, Created)
DT   18-APR-2005 (Rel. 83, Last updated, Version 3)
XX
```

### GenBank (DDBJ) Format

```
LOCUS       X64178       429 bp    mRNA    linear    ROD 26-JUL-2016
DEFINITION  H.auretus mRNA for beta-like x-globin gene.
ACCESSION   X64178
VERSION
```

**ID vs. LOCUS:**
**not required to be the same.**
**Accession numbers (AC, ACCESSION):**
**are the same among databases.**

BIOS477/877 L2 - 42

42

---

7

## Sequence Databases: DNA

**EMBL Format**                    **GenBank (DDBJ) Format**



ID or LOCUS may change in the future.
Accession numbers are permanent.

43

---

## Sequence Databases: Protein

➢ **Protein Databases**
- **UniProt** the Universal Protein Resource
  **https://www.uniprot.org/**
  ➔ **Collaboration between European Bioinformatics Institute (EMBL-EBI), Swiss Institute of Bioinformatics (SIB), and Protein information Resource (PIR, https://proteininformationresource.org/)**
  ➔ **UniProtKB** Protein Knowledgebase
    **Two sections:**
    ◆ **Swiss-Prot** (manually annotated and reviewed)
    ◆ **TrEMBL** (automatically annotated and NOT reviewed)

Protein databases have more information
(functions, domains, superfamily, cross-references, etc.)

44

---

## Sequence Databases: Protein

**UniProt Format** (similar to EMBL format; used by all protein databases)



Protein family information

Cross-references to other databases
(DNA, domain, etc.)

45

---

## UniProt graphic view



Accession #    Entry name

Swiss-Prot entries are manually annotated (high quality)

TrEMBL entries are only computationally annotated.
After careful review, TrEMBL entries will be moved to Swiss-Prot.

**https://www.uniprot.org/uniprotkb/P69905/entry**

46

---

## READSEQ: Sequence Format Converter

- **Readseq @ phylogeny.fr**
  **http://www.phylogeny.fr/one_task.cgi?task_type=readseq**
- **Sequence Format Conversion @ EBI**
  **https://www.ebi.ac.uk/Tools/sfc/emboss_seqret/**
- **Readseq @ mafft website**
  **https://mafft.cbrc.jp/alignment/server/cgi-bin/readseq.txt**
- **Format Converter @ lanl.gov**
  **https://www.hiv.lanl.gov/content/sequence/FORMAT_CONVERSION/form.html**

  [More on the course website, Links page]

47

---

## Sequence Analysis Tools on the Web

➢ **NCBI Sequence Analysis Tools (BLAST, etc.)**
  **https://www.ncbi.nlm.nih.gov/guide/sequence-analysis/**
➢ **EBI Bioinformatics services**
  **https://www.ebi.ac.uk/services**
➢ **ExPASy: SIB Bioinformatics Resource Portal**
  **https://www.expasy.org**
➢ **Max-Planck Institute Bioinformatics Toolkit**
  **http://toolkit.tuebingen.mpg.de/**

- **Many more links are available on the Course website, Links page**

48

## Assignment #1 – Feb 1 due*

➢ **Download the assignment file from Canvas**
- Go to "Assignments" page
- Open "Assignment 1"
- Download the file "Assignment1.doc"

➢ **Submit the file with your answers to Canvas**
- Go to "Assignments" page
- Open "Assignment 1"
- Submit your Assignment 1 file with your answers

**WARNING!!**
Once you click on "Submit" button, you cannot delete/change/add file.
But you can resubmit your assignment file(s) multiple times.

*Each assignment is due at 11:59 pm on the specified date.

49