

CG020 Genomika

Lesson 1

Introduction into Bioinformatics

Jan Hejátko

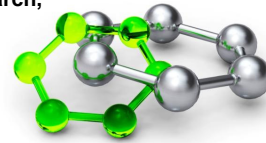
Functional Genomics and Proteomics of Plants,
Mendel Centre for Plant Genomics and Proteomics,
CEITEC - Central European Institute of Technology
and

National Centre for Biomolecular Research,
Faculty of Science,

Masaryk University, Brno

hejatko@sci.muni.cz, www.ceitec.eu

MUNI
SCI



Outline

- Syllabus Of The Course
- Definition Of Genomics
- Role Of Bioinformatics In Functional Genomics
- Databases
 - Spectre Of „On-line“ Resources
 - PRIMARY, SECONDARY and STRUCURAL Databases
 - GENOME Resources
- Analytical Tools
 - Homologies Searching
 - Searching Of Sequence Motifs, Open Reading Frames, Restriction Sites...
 - Other On-line Genome Tools

Course Syllabus

- **Lesson 01**
 - Introduction into Bioinformatics
- **Lesson 02**
 - Identification of Genes
- **Lesson 03**
 - Reverse Genetics Approaches
- **Lesson 04**
 - Forward Genetics Approaches

Course Syllabus

- **Lesson 05**
 - RNA Interference and Genome Editing
- **Lesson 06**
 - Gene Expression and Chemical Genetics
- **Lesson 07**
 - Protein-Protein Interactions And Their Analysis
- **Lesson 08**
 - Recent Approaches in DNA Sequencing

Course Syllabus

- **Lesson 09**
 - Structure of Genomes
- **Lesson 10**
 - Genome evolution
- **Lesson 11**
 - Genomics and Systems Biology
- **Lesson 12**
 - Practical Aspects Of Functional Genomics
 - Model Organisms,
 - PCR

Literature

- Literature resources for **Chapter 01**:
 - **Bioinformatics and Functional Genomics**, 3rd Edition, Jonathan Pevsner, Wiley-Blackwell, 2015
<http://www.bioinfbook.org/php/?q=book3>
 - **Úvod do praktické bioinformatiky**, Fatima Cvrčková, 2006, Academia, Praha
 - **Plant Functional Genomics**, ed. Erich Grotewold, 2003, Humana Press, Totowa, New Jersey

Outline

- Syllabus of the course
- Definition of Genomics

GENOMICS – What is it?

- *Sensu lato* (in the broad sense) – it is interested in **STRUCTURE and FUNCTION** of genomes
 - Necessary prerequisite: knowledge of the genome (sequence) – work with databases
- *Sensu stricto* (in the narrow sense) – it is interested in **FUNCTION** of **INDIVIDUAL GENES** – **FUNCTIONAL GENOMICS**
 - It uses mainly the reverse genetics approaches

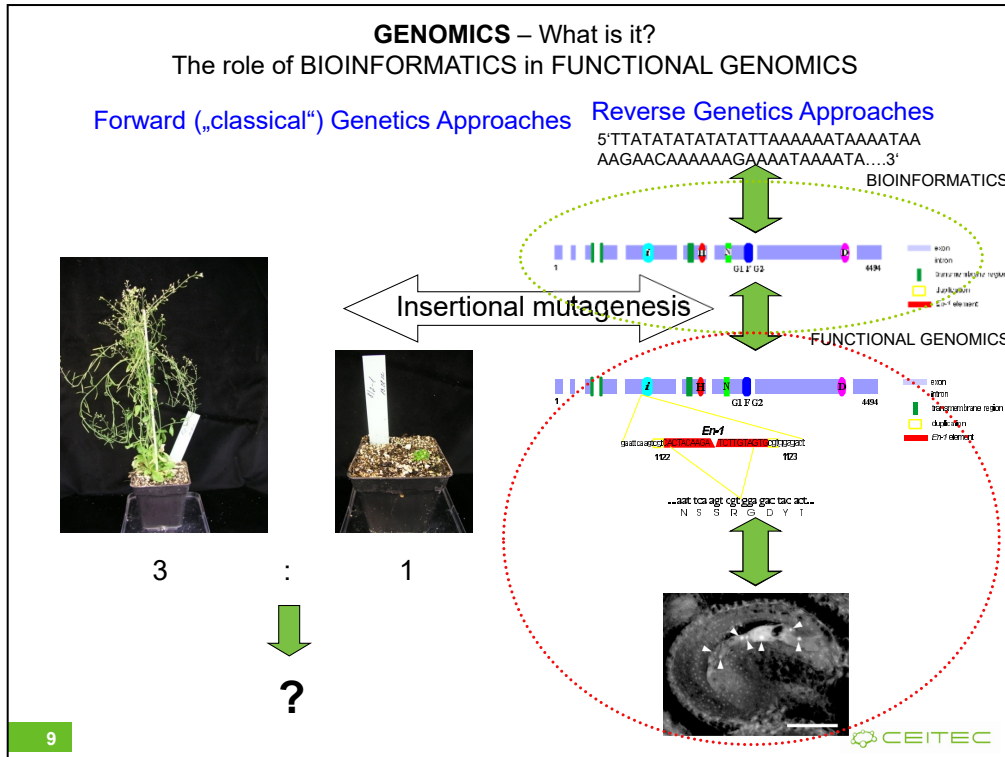
8

 CEITEC

Genomics is a science discipline that is interested in the analysis of genomes. Genome of each organism is a complex of all genes of the respective organism. The genes could be located in cytoplasm (prokaryots) nucleus (in most eukaryotic organisms), mitochondria or chloroplasts (in plants).

The critical prerequisite of genomics is the knowledge of gene sequences.

Functional genomics is interested in function of individual genes.



With the knowledge of gene sequences (or the knowledge of the gene files in the individual organisms, i.e. the knowledge of genomes), **Reverse Genetics** appears that allows study their function.

In comparison to "classical" or **Forward Genetics**, starting with the phenotype, the reverse genetics starts with the sequence identified as a gene in the sequenced genome. The gene identification using approaches of **Bioinformatics** will be described later (see Lesson 02).

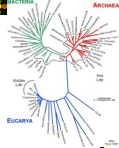
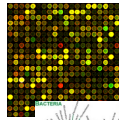
Reverse genetics uses a spectrum of approaches that will be described in the Lesson 03 that allow isolation of sequence-specific mutants and thus their phenotype analysis.

The necessity of having phenotype alterations in the forward genomics approach introduces important difference between those two approaches. Thus, the gene is no longer understood as a factor (*trait*) determining *phenotype*, but rather as a piece of DNA characterized by the unique *string of nucleotides*. i.e. **physical DNA molecule**.

Outline

- Syllabus of this course
- Definition of genomics
- **Role of BIOINFORMATICS in FUNCTIONAL GENOMICS**

Bioinformatics



- **Definition of Bioinformatics** (according to NIH Biomedical Information Science and Technology Initiative Consortium)

Research, development, or application of computational tools and approaches for expanding the use of biological, medical, behavioral or health data, including those to acquire, store, organize, archive, analyze, or visualize such data.

11



NIH WORKING DEFINITION OF BIOINFORMATICS AND COMPUTATIONAL BIOLOGY

July 17, 2000

The following working definition of bioinformatics and computational biology were developed by the BISTIC Definition Committee and released on July 17, 2000. The committee was chaired by Dr. Michael Huerta of the National Institute of Mental Health and consisted of the following members:

Bioinformatics Definition Committee BISTIC Members Expert Members

Michael Huerta (Chair) Gregory Downing
Florence Haseltine Belinda Seto
Yuan Liu

Preamble

Bioinformatics and computational biology are rooted in life sciences as well as computer and information sciences and technologies. Both of these interdisciplinary approaches draw from specific disciplines such as mathematics, physics, computer science and engineering, biology, and behavioral science. Bioinformatics and computational biology each maintain close interactions with life sciences to realize their full potential. Bioinformatics applies principles of information sciences and technologies to make the vast, diverse, and complex life sciences data more understandable and useful. Computational biology uses mathematical and computational approaches to address theoretical and experimental questions in biology. Although bioinformatics and computational biology are distinct, there is also significant overlap and activity at their interface.

Definition

The NIH Biomedical Information Science and Technology Initiative Consortium agreed on the following definitions of bioinformatics and computational biology recognizing that no definition could completely eliminate overlap with other activities or preclude variations in interpretation by different individuals and organizations.

Bioinformatics: Research, development, or application of computational tools and approaches for expanding the use of biological, medical, behavioral or health data, including those to acquire, store, organize, archive, analyze, or visualize such data.

Computational Biology: The development and application of data-analytical and theoretical methods, mathematical modeling and computational simulation techniques to the study of biological, behavioral, and social systems.

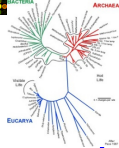
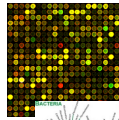
What is bioinformatics?

- **Interface** between the **biology** and **computers**
- **Analysis** of **proteins, genes** and **genomes** using **computer algorithms** and **databases**
- **Genomics** is the **analysis** of **genomes**.

The **tools of bioinformatics** are used to make **sense** of the **billions** of **base pairs** of **DNA** that are sequenced by genomics projects.

J. Pevsner,
<http://www.bioinfbook.org/index.php>

Bioinformatics



- **Bioinformatics in functional genomics**
 - **Processing and analysis of sequencing data**
 - Identification of reference sequences
 - Identification of genes
 - Identification of homologues, orthologues and paralogues
 - Correlative analysis of genomes and phenotypes (incl. human)
 - **Processing and analysis of transcriptional data**
 - Transcriptional profiling using DNA chips or next-gen sequencing
 - **Evaluation of experimental data and prediction of new regulations in systems biology approaches**
 - Mathematical modelling of gene regulatory networks

Outline

- Syllabus of this course
- Definition of genomics
- Role of BIOINFORMATICS in FUNCTIONAL GENOMICS
- Databases
 - Spectre of „on-line“ resources

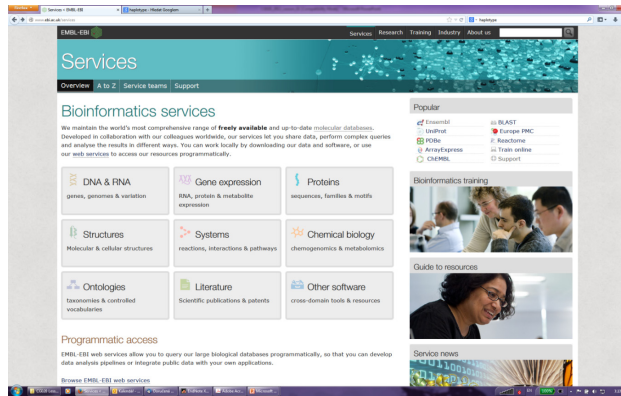
Spectre of On Line Resources

EMBLnet National Nodes		
Vienna BioCenter	Austria	http://www.at.emblnet.org/
EMBL	Belgium	http://www.be.emblnet.org/
EMBLbase	Denmark	http://db.emblbase.dk/
CSC	Finland	http://www.fi.emblnet.org/
INFOBIOGEN	France	http://www.infobiogen.fr/
GENESNET	Germany	http://genomem.dfb-haldenberg.de/biocont/
IMSB	Greece	http://www.imsb.forth.gr/
EMBL	Hungary	http://www.hu.emblnet.org/
INCEB	Ireland	http://www.gen.tcd.ie/
EMBL	Israel	http://dgpas.welamann.ac.il/foef/fin.html
EMBL-ABR	Italy	http://dbi-www.ba.cnr.it/2000/BioWWW/Bo-WWW.htm
KAOS/CANN	Netherlands	http://www.kaos.kun.nl/
EMBL	Norway	http://www.no.emblnet.org/
EMBL	Poland	http://www.ils.waw.pl/
EMBL	Portugal	http://www.lgc.gubportugal.pt/
EMBLnet	Russia	http://www.genetec.msk.ru/
EMBL-CSC	Spain	http://www.es.emblnet.org/
EMBL	Sweden	http://www.se.emblnet.org/
EMBL	Switzerland	http://www.ch.emblnet.org/
EMBLNET	UK	http://www.seqnet.rl.ac.uk/
EMBLnet Specialist Nodes		
EMBL	Germany	http://www.mips.biochem.mpg.de/
EMBL	Italy	http://www.igibn.biochem.it/
Pharmacia Upjohn	Sweden	http://www.gnu.com/
F. Hoffmann La Roche	Switzerland	http://www.rhbc.com/
EMBL	UK	http://www.abi.ac.uk/
EMBL-RC	UK	http://www.bgmp.mrc.ac.uk/
EMBL	UK	http://www.sanger.ac.uk/
EMBL	UK	http://www.biolinf.man.ac.uk/ibbswester
EMBLnet Associate Nodes		
EMBL	Argentina	http://fsl.biol.unla.edu.ar/emblnet
EMBL	Australia	http://www.emph.usc.edu.au/
EMBL	China	http://www.cbl.pku.edu.cn/
EMBL	Cuba	http://dbi.cigb.edu.cu/
EMBL	India	http://sakarjung.emblnet.org.in/
EMBL	South Africa	http://www.sarbi.ac.za
EMBL Information Providers		
EMBL	USA	http://www.ncbi.nlm.nih.gov/
EMBL	USA	http://www.nlm.nih.gov/
EMBL	USA	http://www.nih.gov/

There are many of on-line resources that could be used.

Spectre of On Line Resources

- EBI <http://www.ebi.ac.uk/services>



16

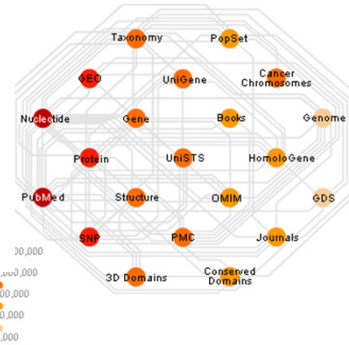


Nowadays, the resources are interconnected and could be accessed via dedicated web pages. Among the best and mostly used www resources integrating plenty of database resources belong www portal of European Bioinformatics Institute (EBI) in Europe (Germany) and National Center of Biotechnology Information (NCBI) in the USA (

Spectre of On Line Resources

NCBI <http://www.ncbi.nlm.nih.gov/>

The screenshot shows the NCBI homepage with a search bar at the top. On the left is a navigation menu with categories like 'All Resources', 'Chemicals & Bioassays', 'Data & Software', 'DNA & RNA', 'Domains & Structures', 'Genes & Expression', 'Genetics & Medicine', 'Genomes & Maps', 'Homology', 'Literature', 'Proteins', 'Sequence Analysis', 'Taxonomy', 'Training & Tutorials', and 'Variation'. The main content area includes a 'Welcome to NCBI' message, a 'Get Started' section with links to 'Tools', 'Downloads', 'How To's', and 'Submissions', a 'NCBI YouTube channel' link, and a 'NCBI Announcer' section with news about a new version of Genes, a new integrated download, a new July Newsletter, and a new Microbial BLAST.



17

CEITEC

Nowadays, the resources are interconnected and could be accessed via dedicated web pages.

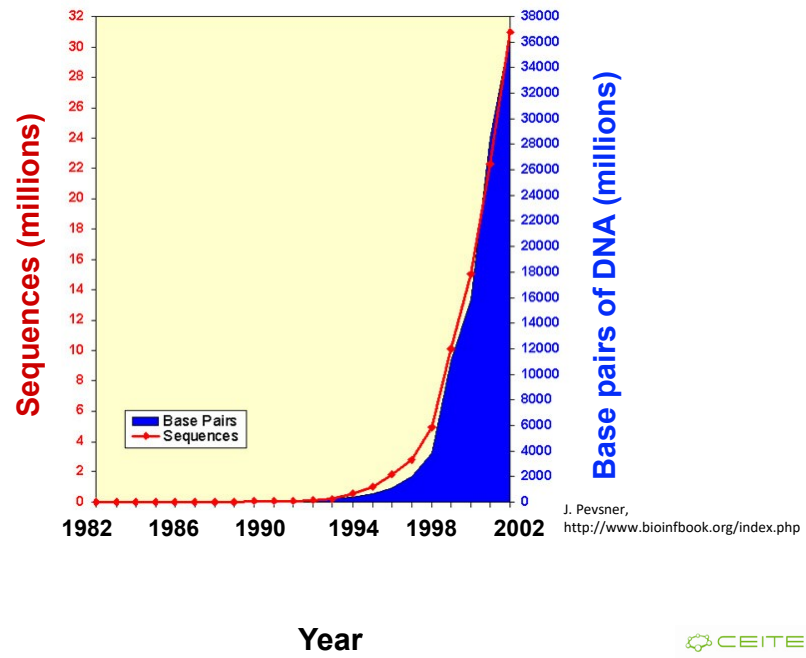
Outline

- Syllabus of this course
- Definition of genomics
- Role of BIOINFORMATICS in FUNCTIONAL GENOMICS
- Databases
 - Spectre of „on-line“ resources
 - PRIMARY, SECONDARY and STRUCURAL databases

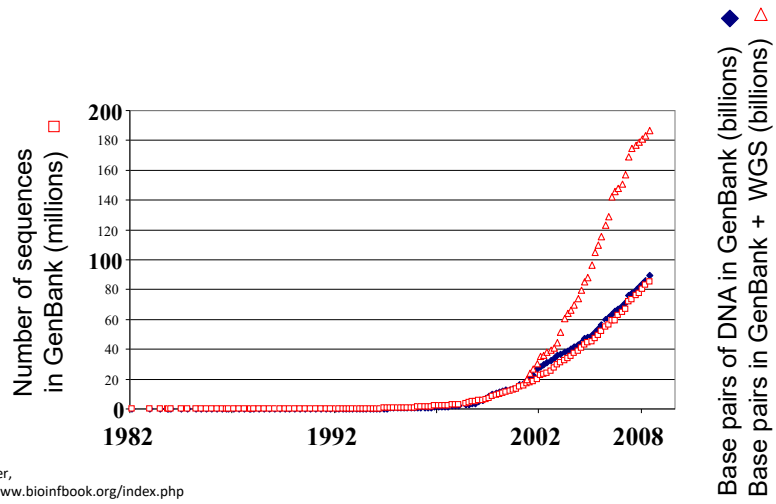
Primary Databases

- Include primary datasets – DNA and Protein sequences
 - Sequences in databases of „The Big Three“:
 - **EMBL**
 - <http://www.ebi.ac.uk/embl/>
 - **GenBank**
 - <http://www.ncbi.nih.gov/Genbank/GenbankSearch.html>
 - **DDBJ**
 - <http://www.ddbj.nig.ac.jp>
 - Daily mutual exchange and backup of data
 - Works with large amount of data (capacity and software requirements)
 - September 2003 27,2 x 10⁶ entries (approx. 33 x 10⁹ bp)
 - August 2005 100 x 10⁹ bp from 165.000 organisms

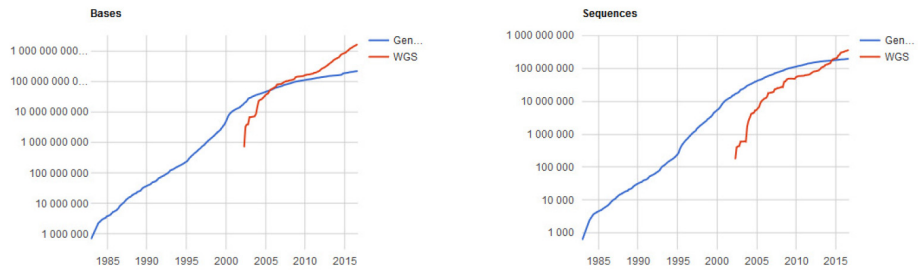
Growth of GenBank



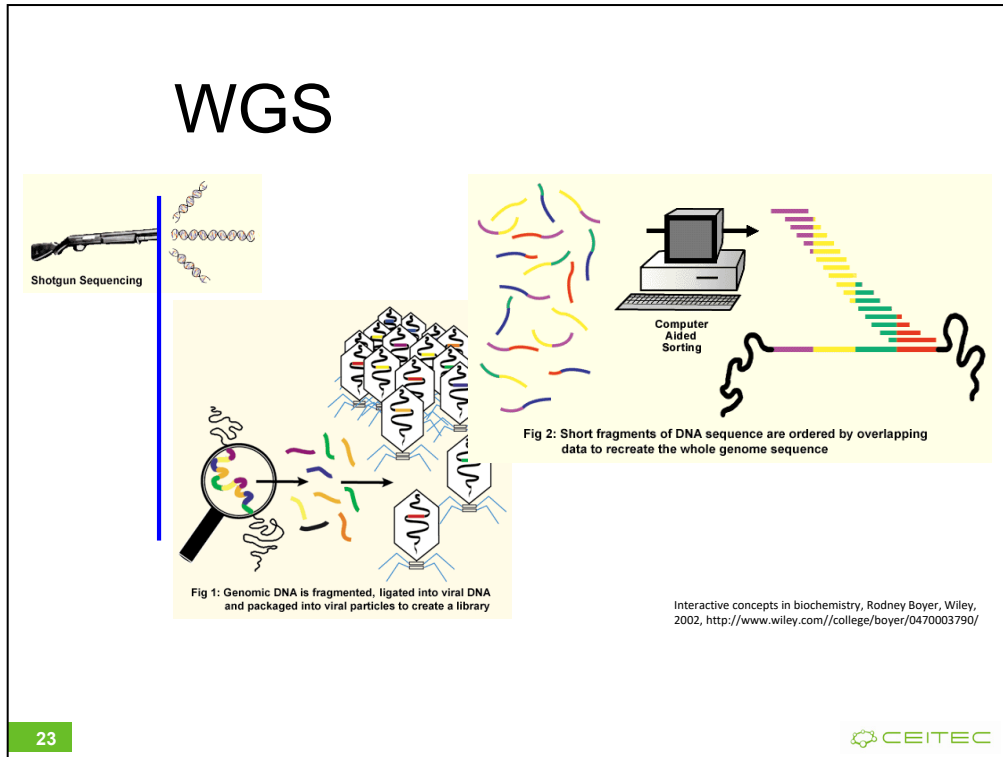
Growth of GenBank + Whole Genome Shotgun (1982-November 2008): we reached **0.2 terabases**



Growth of GenBank Aug 2016



- Dec 1982 680 338 bp, 606 sequences
- Apr 2002 19×10^9 bp, 17×10^6 sequences + WGS 692×10^6 bp, 172 768 sequences
- Aug 2016 218×10^9 bp, 196×10^6 sequences + WGS $1,6 \times 10^{12}$ bp, 360×10^6 sequences

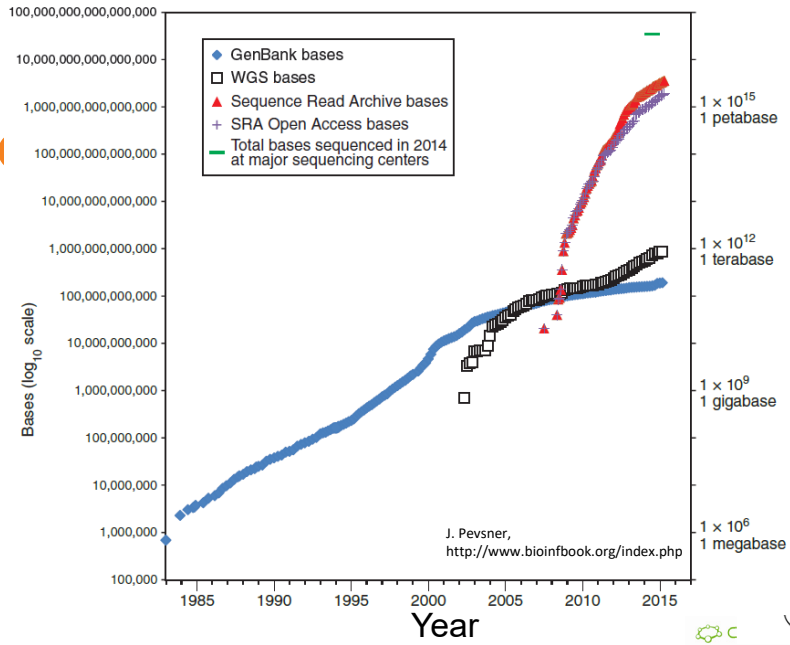


Shotgun sequencing allows a scientist to rapidly determine the sequence of very long stretches of DNA. The key to this process is fragmenting of the genome into smaller pieces that are then sequenced side by side, rather than trying to read the entire genome in order from beginning to end. The genomic DNA is usually first divided into its individual chromosomes. Each chromosome is then randomly broken into small strands of hundreds to several thousand base pairs, usually accomplished by mechanical shearing of the purified genetic material. Each of the short DNA pieces is then inserted into a DNA vector (a viral genome), resulting in a viral particle containing "cloned" genomic DNA (Fig. 1).

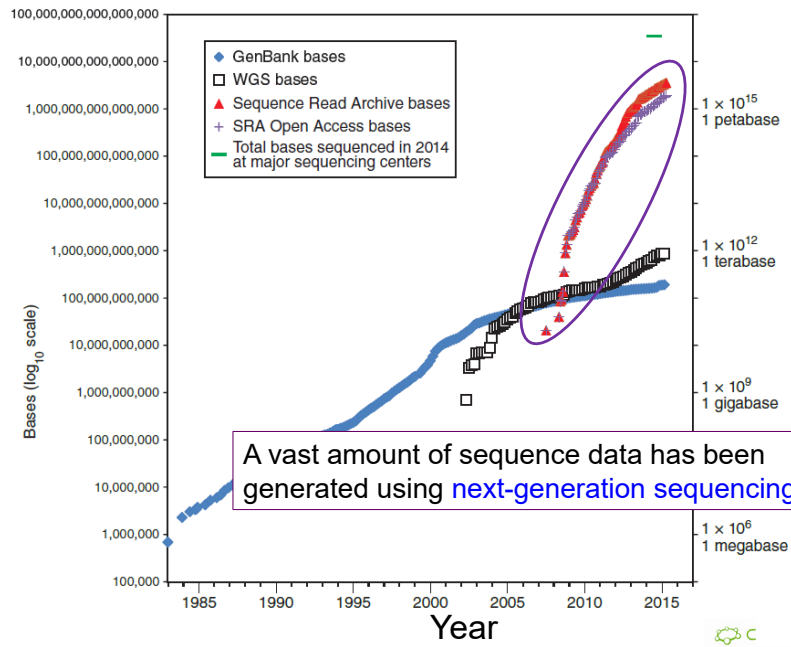
The collection of all the viral particles with all the different genomic DNA pieces is referred to as a library. Just as a library consists of a set of books that together make up all of human knowledge, a genomic library consists of a set of DNA pieces that together make up the entire genome sequence. Placing the genomic DNA within the viral genome allows bacteria infected with the virus to faithfully replicate the genomic DNA pieces. Additionally, since a little bit of known sequence is needed to start the sequencing reaction, the reaction can be primed off the known flanking viral DNA.

In order to read all the nucleotides of one organism, millions of individual clones are sequenced. The data is sorted by computer, which compares the sequences of all the small DNA pieces at once (in a "shotgun" approach) and places them in order by virtue of their overlapping sequences to generate the full-length sequence of the genome (Fig. 2). To statistically ensure that the whole genome sequence is acquired by this method, an amount of DNA equal to five to ten times the length of the genome must be sequenced. (Interactive concepts in biochemistry, Rodney Boyer, Wiley, 2002, <http://www.wiley.com/college/boyer/0470003790/>)

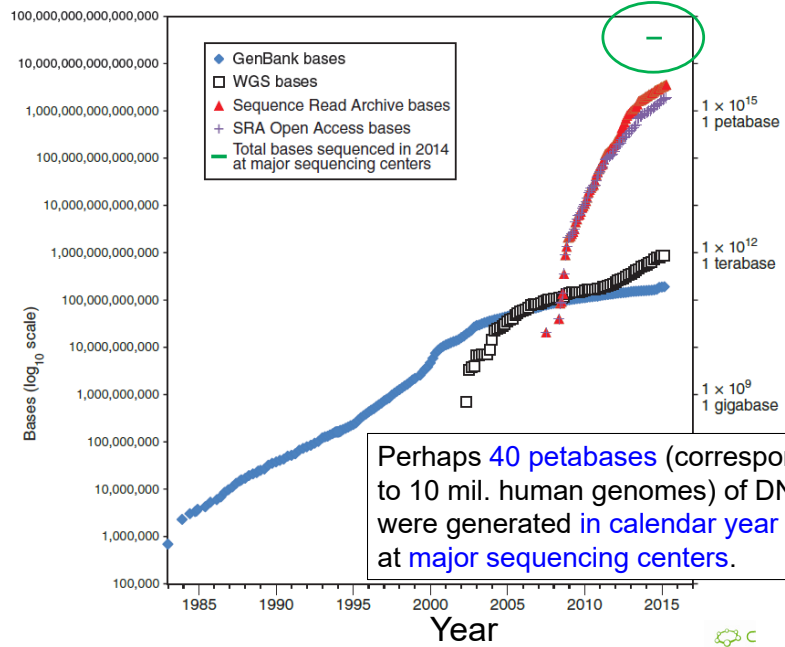
Growth of DNA Sequence in Repositories



Growth of DNA Sequence in Repositories



Growth of DNA Sequence in Repositories



B&FG 3e

Fig. 2-3

23 22

CC

80

Primary Databases

- They include sets of primary data – [DNA](#) and [Protein](#) sequences
 - Protein sequences:
 - **PIR**, <http://pir.georgetown.edu/>
 - **MIPS**, <http://www.mips.biochem.mpg.de>
 - **SWISS-PROT**, <http://www.expasy.org/sprot/>

Primary Databases

- Types of sequences in primary databases
 - **Standard nucleotide sequences** acquired by high quality sequencing
 - **ESTs** (**E**xpressed **S**equences **T**ags)
 - **HGTS** (**H**igh **T**hroughput **G**enome **S**equencing)
 - Results of sequencing projects without annotation
 - **Reference Sequences** of annotated genomes
 - **TPAs** (**T**hird **P**arty **A**nnotation)
 - sequences annotated by third party (by someone else, not the original authors)

Primary Databases

GenBank (NCBI) <http://www.ncbi.nlm.nih.gov/>



The screenshot shows the NCBI homepage with a search bar at the top, a navigation menu on the left, and a main content area with sections for 'Welcome to NCBI', 'Get Started', 'Popular Resources', and 'NCBI YouTube channel'.

NCBI Home

- Resource List (A-Z)
- All Resources
- Chemicals & Bioassays
- Data & Software
- DNA & RNA
- Domains & Structures
- Genes & Expression
- Genetics & Medicine
- Genomes & Maps
- Homology
- Literature
- Proteins
- Sequence Analysis
- Taxonomy
- Training & Tutorials
- Variation

Welcome to NCBI

The National Center for Biotechnology Information advances science and health by providing access to biomedical and genomic information.

[About the NCBI](#) | [Mission](#) | [Organization](#) | [Research](#) | [RSS Feeds](#)

Get Started

- Tools:** Analyze data using NCBI software
- Downloads:** Get NCBI data or software
- How-To's:** Learn how to accomplish specific tasks at NCBI
- Submissions:** Submit data to GenBank or other NCBI databases

NCBI YouTube channel

Learn how to get the most out of NCBI tools and databases with video tutorials on the NCBI YouTube Channel. [GO](#)

Popular Resources

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

NCBI Announcements

New version of GenBank available

An integrated, download for viewing and analyzing NCBI's July Newsletter

Primary Databases

The screenshot displays the NCBI Gene database entry for gene NC_022777.1. The page is organized into several sections:

- Gene symbol:** *adh*
- Gene description:** non-component VMA-like sensor kinase
- Gene tag:** *adh*
- Gene type:** protein coding
- RefSeq name:** *ADH02277*
- Organism:** *Agrobacterium tumefaciens* subsp. *discreta* (strain: ATCC 31624)
- lineage:** Bacteria; Proteobacteria; Alphaproteobacteria; Rhizobiales; Rhizobiales; Rhizobiales/Agrobacterium group; Agrobacterium; Agrobacterium tumefaciens complex

Genomic context:

- Location:** chr1:141,432-141,433
- Sequence:** NC_022777.1 (141854-141855)

Genomic regions, transcripts, and products:

Genomic Sequence: NC_022777.1

Related articles:

1. Sequence analysis of the *adh* gene from *Agrobacterium tumefaciens* strain T-3092. *Biochimica et Biophysica Acta* 1994;124:1-10.
2. The *adh* gene is a homolog of *AdhA* in *Agrobacterium tumefaciens*. *Plant Cell* 1993;5:1105-1110.
3. Characterization of the *adh* locus of *Agrobacterium tumefaciens*. A transcriptional regulator and host gene orthologues. *Genes & Dev* 1997;11:197-209.
4. Analysis of the complete nucleotide sequence of the *Agrobacterium tumefaciens* cell genome. *Temperature* 1997;13:1-10.

GeneID's Gene Reference Site Function: [What's New?](#)

Submit: [View Details](#) [Citations](#)

Primary Databases

The screenshot displays a web browser window showing a GenBank entry for the gene **NP_059797.1**. The browser address bar shows the URL www.ncbi.nlm.nih.gov/nuccore/145694. The main content area shows a genomic map with a scale from 145,400 to 147,600. A red bar represents the gene, with a tooltip providing details:

- NP_059797.1**
- NP_059797.1: two-component VirA-like sensor kinase
- total range: NC_002377.1 (145,694..148,183)
- total length: 2,490
- strand: plus
- protein product length: 829

Below the tooltip, there is a section titled **Links & Tools** with the following links:

- GenBank View: [NC_002377.1 \(145,694..148,183\)](#), [NP_059797](#)
- FASTA View: [NC_002377.1 \(145,694..148,183\)](#), [NP_059797](#)
- BLAST Genomic: [NC_002377.1 \(145,694..148,183\)](#)
- Graphical View: [NP_059797.1](#)
- BLAST Protein: [NP_059797.1](#)
- BLINK Results: [NP_059797.1](#)

Below the links, there is a section titled **Bibliography** and a sub-section titled **Related articles in PubMed**.

Primary Databases

The screenshot displays the NCBI Nucleotide database search results for the accession number BC_022377.1. The interface includes a search bar at the top with the text "Search Nucleotide" and a "Go" button. Below the search bar, the accession number "BC_022377.1" is highlighted with a red circle and labeled "Accession number". The GeneBank Identifier "BC_022377.1" is also circled in red and labeled "GeneBank Identifier". The main content area shows the following details:

```
LOCUS       BC_022377                2490 bp    DNA    linear    M17 29-DEC-2003
DEFINITION  Agrobacterium tumefaciens entomobion plasmid Ti, complete sequence.
ACCESSION   BC_022377
VERSION    BC_022377.1
KEYWORDS   Agrobacterium tumefaciens (Rhizobium radiobacter)
SOURCE     Agrobacterium tumefaciens (Rhizobium radiobacter)
  FAYARD, S.K.,
  TITLE     Oncopine-type Ti plasmid sequence
  JOURNAL   Unpublished
  REFERENCE 2 (Issue 1 to 2490)
  AUTHORS   Zhu, J., Oyer, P.M., Schrammeyer, R., Rooykas, P.J., Parraud, E.K. and
  Wisniewski, S.C.
  TITLE     Direct Submission
  JOURNAL   Submitted (07-MAR-2004) Microbiology, Cornell University, Wing
  Hall, Ithaca, NY 14853, USA
  COMMENT   PROTEIN: 822aa. This record has not yet been subject to final
  NCBI review. The reference sequence was derived from 822aa.
  FEATURES     Location: 0..2490
  source       1..2490
               /organism="Agrobacterium tumefaciens"
               /mol_type="genomic DNA"
               /db_xref="taxon:314"
               /plasmid="Ti"
               /note="antitumorigenic"
  gene         1..2490
               /gene="viraA"
               /db_xref="GeneID:1224316"
  CDS         1..2490
               /gene="viraA"
               /note="viraA component regulator of vira replication; ViraA is a
               transmembrane histidine kinase"
               /codon_start=1
               /transl_start=1
               /product="viraA"
               /protein_id="F003837.1"
               /db_xref="GI:10983141"
```


What is an **Accession Number**?

An accession number is label that used to identify a sequence. It is a string of letters and/or numbers that corresponds to a molecular sequence.

Examples (all for retinol-binding protein, RBP4):

X02775	GenBank genomic DNA sequence	DNA
NT_030059	Genomic contig	
Rs7079946	dbSNP (single nucleotide polymorphism)	

N91759.1	An expressed sequence tag (1 of 170)	RNA
NM_006744	RefSeq DNA sequence (from a transcript)	

NP_007635	RefSeq protein	Protein
AAC02945	GenBank protein	
Q28369	SwissProt protein	
1KT7	Protein Data Bank structure record	

J. Pevsner,
<http://www.bioinfbook.org/index.php>

NCBI's important **RefSeq** project: best **representative sequences**

RefSeq (accessible via the main page of NCBI) provides an **expertly curated accession number** that corresponds to **the most stable, agreed-upon "reference" version of a sequence**.

RefSeq identifiers include the following formats:

Complete genome	NC_#####
Complete chromosome	NC_#####
Genomic contig	NT_#####
mRNA (DNA format)	NM_##### e.g. NM_006744
Protein	NP_##### e.g. NP_006735

J. Pevsner,
<http://www.bioinfbook.org/index.php>

RefSeq

two-component ViaA-like sensor kinase

NCBI Reference Sequences (RefSeq)

Genome Annotation

The following sections contain reference sequences that belong to a specific genome build. [Explain](#)

Reference assembly

Genomic

1. **NC_003065.3**

Range 190831..193332

Download [GenBank](#) [FASTA](#) [Sequence Viewer \(Graphics\)](#)

mRNA and Protein(s)

1. **NP_396486.1** two component sensor kinase [Agrobacterium tumefaciens str. C58]

UniProtKB/Swiss-Prot [P18540](#)

Conserved Domains (3) [summary](#)

cd08025	HATPass_c: Histidine kinase-like ATPases. This family includes several ATP-binding proteins for example: histidine kinase, DNA gyrase B, topoisomerase, heat shock protein HSP90, phytochrome-like ATPases and DNA mismatch repair proteins.
cd08032	HtkA, Histidine Kinase A (dimerization/phosphoreceptor) domain: Histidine Kinase
Location:468 - 530	A dimer is formed through parallel association of 2 domains creating 4-helix bundles; usually these domains contain a conserved His residue and are activated via...
PRK13837	PRK13837: two-component ViaA-like sensor kinase; Provisional
Location:18 - 833	Blast Score: 2944

Related Sequences

NCBI's RefSeq project: many accession number formats for genomic, mRNA, protein sequences

Accession	Molecule	Method	Note
AC_123456	Genomic	Mixed	Alternate complete genomic
AP_123456	Protein	Mixed	Protein products; alternate
NC_123456	Genomic	Mixed	Complete genomic molecules
NG_123456	Genomic	Mixed	Incomplete genomic regions
NM_123456	mRNA	Mixed	Transcript products; mRNA
NM_123456789	mRNA	Mixed	Transcript products; 9-digit
NP_123456	Protein	Mixed	Protein products;
NP_123456789	Protein	Curation	Protein products; 9-digit
NR_123456	RNA	Mixed	Non-coding transcripts
NT_123456	Genomic	Automated	Genomic assemblies
NW_123456	Genomic	Automated	Genomic assemblies
NZ_ABCD12345678	Genomic	Automated	Whole genome shotgun data
XM_123456	mRNA	Automated	Transcript products
XP_123456	Protein	Automated	Protein products
XR_123456	RNA	Automated	Transcript products
YP_123456	Protein	Auto. & Curated	Protein products
ZP_12345678	Protein	Automated	Protein products

J. Pevsner,
<http://www.bioinfbook.org/index.php>

Primary Databases

NC_002377.1: 145K..148K (2.9Kbp)

Genes

NP_059797.1

NP_059797.1: two-component VirA-like sensor kinase
total range: NC_002377.1 (145,694..148,183)
total length: 2,490
strand: plus
protein product length: 829

Links & Tools

GenBank View: [NC_002377.1 \(145,694..148,183\)](#), [NP_059797.1 \(145,694..148,183\)](#)
FASTA View: [NC_002377.1 \(145,694..148,183\)](#), [NP_059797.1 \(145,694..148,183\)](#)
BLAST Genomic: [NC_002377.1 \(145,694..148,183\)](#)
Graphical View: [NP_059797.1](#)
BLAST Protein: [NP_059797.1](#)
BLINK Results: [NP_059797.1](#)

Bibliography

Related articles in PubMed

Primary Databases

The screenshot shows the NCBI GenBank entry for Agrobacterium tumefaciens plasmid Ti. The main content is the DNA sequence, which is displayed in a monospaced font. The sequence is preceded by a header: "Agrobacterium tumefaciens plasmid Ti, complete sequence" and "NCBI Reference Sequence: NC_002377.1". To the right of the sequence, there is a sidebar with various options and tools, including "Change region shown", "Analyze this sequence", "Run BLAST", "Find in this Sequence", and "Related Information". The "Change region shown" section is currently set to "Whole sequence" from position 145694 to 148183. The "Analyze this sequence" section includes options for "Run BLAST", "Find Primers", and "Highlight Sequence Features". The "Related Information" section includes "RefSeq", "Full text in PMC", "Gene", "Genome", "Identical GenBank Sequences", "Protein", "Protein Clusters", "PubMed", "PubMed (Weighted)", and "Taxonomy". The "Recent activity" section shows a list of recent searches, including "Agrobacterium tumefaciens plasmid Ti, complete sequence", "vIA [Agrobacterium tumefaciens]", and "vIA [Agrobacterium tumefaciens str. CS1]".

Secondary Databases

- Databases of **functional** or **structural motifs**, acquired by **primary data** (sequences) **comparison**
- **PROSITE**, <http://www.expasy.org/prosite/>

```
->PDOC0003 PS00002 SULFATION Tyrosine sulfation site [rule] [Warning: rule with a high probability of occurrence].
571 - 585 nbnsnatYetsiana

->PDOC0004 PS00004 CAMP_PHOSPHO_SITE cAMP- and cGMP-dependent protein kinase phosphorylation site [pattern] [Warning: pattern with a high probability of occurrence].
744 - 747 DQAT
814 - 817 DEEE

->PDOC0005 PS00005 PKC_PHOSPHO_SITE Protein kinase C phosphorylation site [pattern] [Warning: pattern with a high probability of occurrence].
148 - 150 GPH
144 - 146 TYP
172 - 173 AAK
219 - 221 DSK
369 - 371 TIK
400 - 402 EGP
412 - 416 EGP
485 - 487 ELS
492 - 494 TYP
492 - 494 TYP
716 - 718 EGP
726 - 728 EGP
747 - 749 DAK
794 - 796 EAP
804 - 806 ESK
804 - 806 ESK
860 - 870 ESK
860 - 870 EGP
860 - 870 ESK
860 - 870 ESK
874 - 876 TAP
897 - 899 EAP
1002 - 1004 TYP
1018 - 1020 EGP
1031 - 1033 TYP
1128 - 1131 DAK
```


Secondary Databases

- Databases of **functional** or **structural motifs**, acquired by **primary data** (sequences) comparison

□ **PRINTS**, <http://www.bioinf.man.ac.uk/dbbrowser/PRINTS/>



PRINTS is a compilation of protein fingerprints. A fingerprint is a group of conserved motifs used to characterize a protein family; its diagnostic power is refined by iterative scanning of a PRINTS-PROFFPRINTS component. Usually the motifs do not overlap, but are separated along a sequence, though they may be contiguous in 3D space. Fingerprints can encode protein links and functionalities more flexibly and powerfully than can single motifs, full diagnostic potency deriving from the mutual context provided by motif neighbors. [References](#)

News:

- [SPINCE](#) - Search PRINTS and manual PRINTS
- [mpPRINTS](#) - Search PRINTS automatic exploration
- [PRINTS](#) - Search the improved latest protein database

Direct PRINTS access:

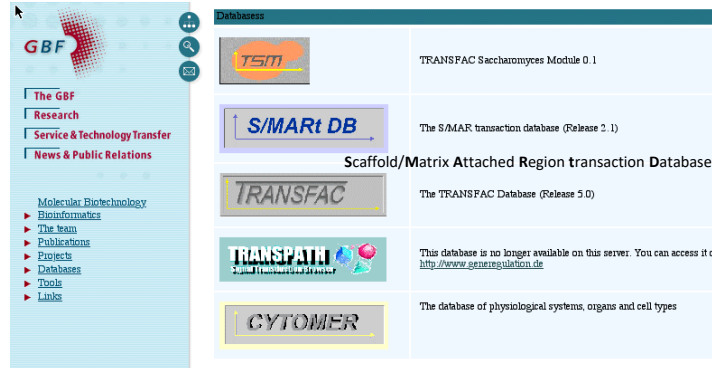
- [By sequence number](#)
- [By PRINTS code](#)
- [By UniProt code](#)
- [By name](#)
- [By domain](#)
- [By PDB](#)
- [By number of motifs](#)
- [By motif](#)
- [By query sequence](#)

PRINTS search:

- Search PRINTS with [NEW LinearPRINTScan](#)
- [PRINTS](#)
- [CEPIScan](#)
- [MILScan](#)
- [FingerPRINTScan](#) binaries and source are available: printscan@bioinf.man.ac.uk

Secondary Databases

- **TRANSFAC** <http://www.gene-regulation.com/>



The screenshot shows the TRANSFAC website interface. On the left is a navigation menu for GBF (German Biotechnology Foundation) with categories like 'The GBF', 'Research', 'Service & Technology Transfer', and 'News & Public Relations'. Below these are sub-categories: 'Molecular Biotechnology', 'Bioinformatics', 'The team', 'Publications', 'Projects', 'Databases', 'Tools', and 'Links'. The main content area is titled 'Databases' and lists several databases:

Database Logo	Description
TSM	TRANSFAC Saccharomyces Module 0.1
S/MARt DB	The S/MAR transaction database (Release 2.1) Scaffold/Matrix Attached Region transaction Database
TRANSFAC	The TRANSFAC Database (Release 5.0)
TRANSPATH	This database is no longer available on this server. You can access it on http://www.gene-regulation.de
CYTOMER	The database of physiological systems, organs and cell types

44

CEITEC

S/MARt DB (saffold/matrix attached region transaction database). This database collects information about S/MARs and the nuclear matrix proteins that are supposed be involved in the interaction of these elements with the nuclear matrix. <http://transfac.gbf.de/SMARTDB/index.html>)

Structural Databases

- PDB <http://www.rcsb.org/pdb/>

DEPOSIT data
DOWNLOAD files
BROWSE LINKS
BETA TEST new features
BETA mirror files

Current Holdings
19623 Structures
Last Update: 30-Dec-2002
PDB Statistics

Molecule of the Month:
Cytochrome c

The Protein Data Bank (PDB) is operated by Rutgers, The State University of New Jersey; the San Diego Supercomputer Center at the University of California, San Diego; and the National Institute of Standards and Technology -- three members of the **Research Collaboratory for Structural Bioinformatics (RCSB)**. The PDB is supported by funds from the **National Science Foundation**, the **Department of Energy**, and two units of the **National Institutes of Health**: the

PROTEIN DATA BANK

Welcome to the PDB, the single worldwide repository for the processing and distribution of 3-D biological macromolecular structure data.

ABOUT PDB | DATA UNIFORMITY | RECENT FEATURES | USER GUIDES | FILE FORMATS | EDUCATION | STRUCTURAL GENOMICS | PUBLICATIONS | SOFTWARE

Search the Archive

Enter a PDB ID or keyword [Query Tutorial](#)

query by PDB id only match exact word
 remove sequence homologues

[Search tips](#) keyword search form with examples
[Search fields](#) customizable search form
[Status Search](#) find entries awaiting release

News [Complete News](#) [pdb-L Archive](#)
[Newslister](#) [Subscribe](#)

23-Dec-2002
Happy Holidays from the PDB! The PDB staff wish to extend our best wishes to the community for a happy holiday season and a wonderful new year!

PDB Mirrors

"Please bookmark a mirror site!"

- San Diego Supercomputer Center*
- Rutgers University*
- National Institute of Standards and Technology*
- Cambridge Crystallographic Data Centre, UK
- National University of Singapore
- Osaka University, Japan
- Universidade Federal de Minas Gerais, Brazil
- Max Delbrück Center for Molecular Medicine, Germany

OTHER SITES

Structural Databases


- **PDB** <http://www.rcsb.org/pdb/>

Structure Explorer - 1P5Y

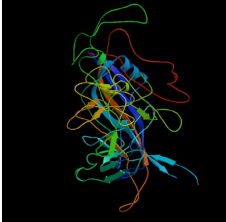
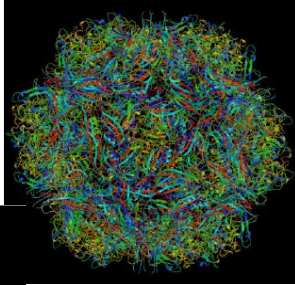
PDB
PROTEIN DATA BANK

Structure Explorer - 1P5Y

Title The Structures Of Heat Range Controlling Regions Of The Capsids Of Canine and Feline Parvoviruses and Mutants
Classification Virus/Viral Protein
Compound Mol. Bt. C. Molecular Coat Protein Vp2, Chain: A; Fragment: Sequence Database Residues 190-237; Engineered: Yes; Mutations: Yes
Exp. Method X-ray Diffraction

 **View Structure**

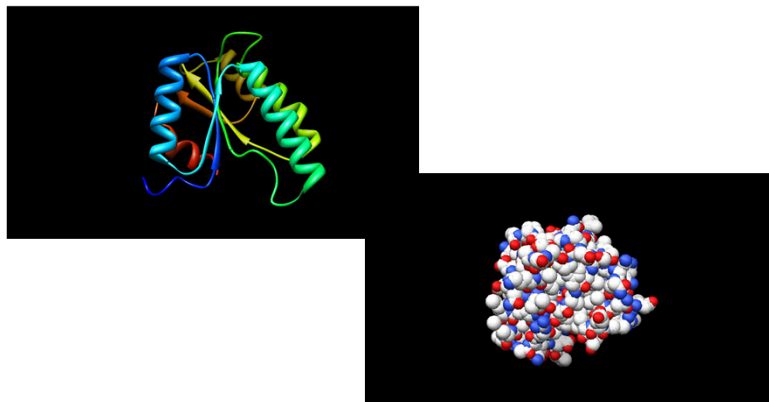
[Summary Information](#)
[View Structure](#)
[Download Display File](#)
[Structural Neighbors](#)
[Geometry](#)
[Other Sources](#)
[Sequence Details](#)



<http://www.rcsb.org/pdb/cgi/structure.cgi?job=graphics&pdb=1P5Y;page=2&id=173561064249344&bio=1&opt=show&size=500> 12/29/2003

Structural Databases

- PDB <http://www.rcsb.org/pdb/>



Pekárová et al., *Plant Journal* (2011)

Outline

- Syllabus Of The Course
- Definition Of Genomics
- Role Of Bioinformatics In Functional Genomics
- Databases
 - Spectre of „on-line“ Resources
 - PRIMARY, SECONDARY And STRUCURAL Databases
 - **GENOME Resources**

Genome Resources

□ **Human Genome Browser** <http://genome.ucsc.edu/cgi-bin/hgGateway>

The screenshot shows the Human Genome Browser interface. At the top, there are navigation tabs for 'Genomes', 'Genome Browser', 'Tools', 'Services', 'Data', and 'About Us'. Below the navigation is a search bar with 'Human (primate assembly) Genome Browser Gateway' as the title. The search bar includes fields for 'chr', 'genome', 'assembly', and 'position', along with a 'Search' button. Below the search bar, there is a section titled 'Human Genome Browser - hg19 assembly (sequences)'. This section includes a paragraph about the 1.4 February 2009 human reference sequence (GRCh37) and a 'Sample position queries' section. The 'Sample position queries' section lists various queries and their corresponding 'Genome Browser Response'.

Request	Genome Browser Response
chr7	Displays all of chromosome 7
chrX_4000212	Displays all of the repeated (copy) 4000212
2p13	Displays region for band p13 on chr 2
chr3_100000	Displays 100,000 bases of chr 3, counting from p-arm telomere
chr3 100000-2000	Displays a region of chr3 that spans 2000 bases, starting with position 100000
RH40061 RH40175	Displays region between genome landmarks, such as the STS markers RH40061 and RH40175, or chromosome bands: 15q11 to 15q13, or SNPs rs104252 and rs1803370. This syntax may also be used for other range queries, such as between uniquely determined ESTs, mRNAs, miRNAs, etc.
15q11 15q13	
rs104252/rs1803370	
D100304	Displays region around D10 (marker D10S2040) from the Chromosome 10 field map. It includes 100,000 bases on each side as well
A020474	Displays region of EST with GenBank accession A020474 in a BRCA1 cancer gene on chr 17
AC00103	Displays region of clone with GenBank accession AC00103
A020311	Displays region of mRNA with GenBank accession number A020311
PRNP	Displays region of genome with HUGO Gene Nomenclature Committee identifier PRNP
MIM_211414	Displays the region of genome with MIM identifier MIM_211414
NP_059110	Displays the region of genome with protein accession number NP_059110
proinsulin mRNA	Lists transcribed products, but not cDNAs
hemoglobin cluster	Lists mRNAs for cluster hemoglobin genes
zinc finger	Lists many zinc finger mRNAs
haptoglobin	Lists only repeated haptoglobin genes
huntingtin	Lists candidate genes associated with Huntington's Disease
zNF97	Lists mRNAs deposited by scientist named Zhan
Evans_J_E	Lists mRNAs deposited by co-author J.E. Evans

Genome Resources

□ Human Genome Browser <http://genome.ucsc.edu/cgi-bin/hgGateway>

The screenshot displays the UCSC Genome Browser interface for the Human Feb. 2009 (GRCh37/hg19) Assembly. The main view shows a genomic region on chromosome 17, with various tracks including RefSeq, Repeat Masker, and other annotations. A green arrow points to the top track, which displays gene models and exons. Below it are tracks for RefSeq, Repeat Masker, and other annotations. The interface includes search bars, zoom controls, and a track configuration panel on the right side.

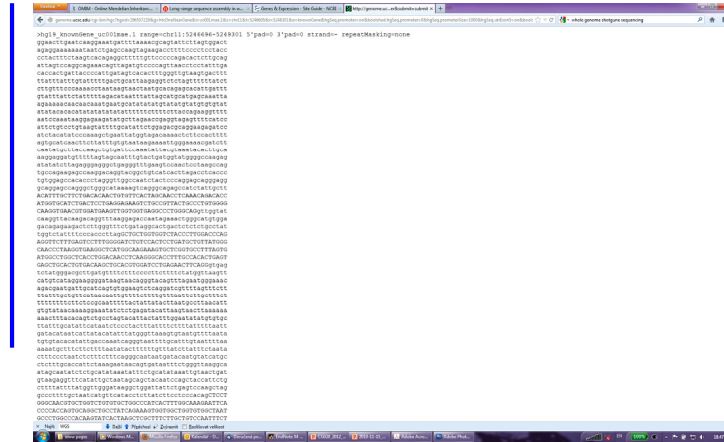
Genome Resources

□ **Human Genome Browser** <http://genome.ucsc.edu/cgi-bin/hgGateway>

The screenshot shows the 'Get Genomic Sequence Near Gene' interface of the Human Genome Browser. The page title is 'Genomic Sequence Near Gene'. Below the title, there is a note: 'Note: if you would prefer to get DNA for more than one feature of this track at a time, try the [Table Browser](#) using the output format sequence.' The main section is titled 'Sequence Retrieval Region Options:' and contains several checkboxes and input fields for configuring the sequence retrieval. The options include: 'Flanking upstream by 1000 bases' (checked), '5' UTR Exons' (checked), 'CDS Exons' (checked), '3' UTR Exons' (checked), 'Introns' (checked), 'Downstream by 1000 bases' (checked), and 'One FASTA record per region (exon, intron, etc.) with 5' extra bases upstream (5') and 3' extra downstream (3')' (checked). Below these options, there is a note: 'Note: if a feature is close to the beginning or end of a chromosome and upstream/downstream bases are added, they may be truncated in order to avoid extending past the edge of the chromosome.' The 'Sequence Formatting Options:' section includes: 'Exons in upper case, everything else in lower case' (checked), 'CDS in upper case, UTR in lower case' (checked), 'All upper case' (unchecked), 'All lower case' (unchecked), and 'Mask repeats: # to lower case | to N' (unchecked). The page is displayed in a web browser window with a Windows taskbar at the bottom.

Genome Resources

□ **Human Genome Browser** <http://genome.ucsc.edu/cgi-bin/hgGateway>



Genome Resources

- The Arabidopsis Information Resource (TAIR) <http://www.arabidopsis.org>



Genome Resources

- TAIR, The Arabidopsis Information Resource, <http://www.arabidopsis.org>



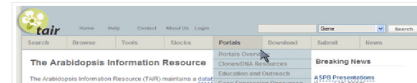
The Arabidopsis Information Resource

The Arabidopsis Information Resource (TAIR) maintains a database of genetic and molecular biology data for the model higher plant *Arabidopsis thaliana*. Data available from TAIR includes the complete genome sequence along with gene structure, gene product information, metabolism, gene expression, DNA and seed stocks, genome maps, genetic and physical markers, publications, and information about the Arabidopsis research community. Gene product function data is updated every two weeks from the latest published research literature and community data submissions. Gene structures are updated 1-2 times per year using computational and manual methods as well as community submissions of new and updated genes. TAIR also provides extensive linkouts from our data pages to other Arabidopsis resources.

The Arabidopsis Biological Resource Center at The Ohio State University collects, reproduces, preserves and distributes seed and DNA resources of *Arabidopsis thaliana* and related species. Stock information and ordering for the ABRC are fully integrated into TAIR.

The NEW arabidopsis.org

We've added new dropdown headers and left navigation bars and reorganized our web pages to make it easier to locate information and resources in TAIR. Please contact us if you experience any problems with our new site.



Breaking News

Data Updates Suspended
[October 15, 2006]
Some TAIR data updates, including loading of new ABRC stocks, will be suspended from Oct 20-Nov 17 while we move our servers.

New Phenotype Search Option
[October 15, 2006]
Search for genes, germplasm and polymorphisms using associated phenotype, and see improved phenotype data display in results and detail pages.

ASPB Presentations
[August 15, 2006]
Following heavy demand, the TAIR workshop presentations given at the ASPB meeting in Boston have been made available from the TAIR website for download.

Outline

- Syllabus Of The Course
- Definition Of Genomics
- Role Of Bioinformatics In Functional Genomics
- Databases
 - Spectre Of „On-line“ Resources
 - PRIMARY, SECONDARY And STRUCURAL Databases
 - GENOME Resources
- Analytical Tools
 - Homology Searching

Analytical Tools

□ Global versus Local alignment

```
Globalní přiřazení
SLAV-----APATNIK-----PIQNYR-I-----AKSETQRYMVE
SLAVYTYIEFVRANAPATNIKSECVRAAPIQNYRRVEHVRATAKSETQRYMVE

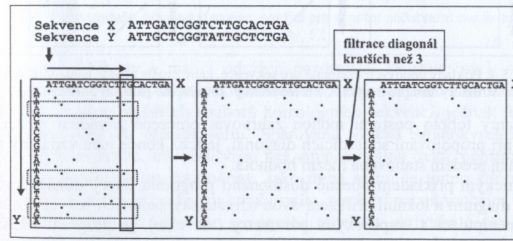
Lokální přiřazení
SLAVYTYIEFVRANAPATNIKSECVRAAPIQNYRRVEHVRATAKSETQRYMVE
-----NAPATNIKSECVRA-PIQNYRRVEHVRA-----
```

Cvrčková, Úvod do praktické bioinformatiky

- **Global Alignment:** only for sequences, which are **similar** and of a **similar length** (BUT can insert spaces into one or both sequences)
- **Global Alignment** is used mainly in case of **multiple alignment** (CLUSTALW, further in the presentation)
- **Local Alignment** provides identification and comparison even in case of alignment of **regions of sequences with high similarity**, e.g. even in case of **change of order** of **protein domains** during evolution

Analytical Tools

- Choosing the right type of alignment using dotplot

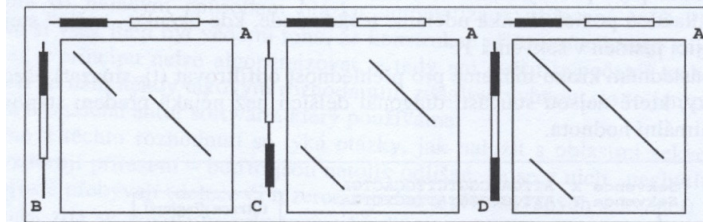


Cvrčková, Úvod do praktické bioinformatiky

- Plotting the sequences against each other (x and y axis)
- Identification of identity in „dot“ of specific size (e.g. 2 bp)
- Filtering the diagonals of lengths lower than a threshold

Analytical Tools

- Examples of sequence alignment using dotplot



Cvrčková, Úvod do praktické bioinformatiky

- **Global Alignment:** possible **only** for **sequences A and B**
- The rest of the sequences underwent change of order of protein domains and therefore it is necessary to do a local alignment
- **Dotplot** can be obtained using **BLAST2** (see further in the presentation)

Analytical Tools

- **BLAST** <http://ncbi.nlm.nih.gov/BLAST/>

The screenshot shows the NCBI Nucleotide BLAST search page. At the top, the NCBI logo is on the left, and the text "nucleotide-nucleotide BLAST" is on the right. Below this, there are tabs for "Nucleotide", "Protein", and "Translations", with "Nucleotide" selected. A link "Retrieve results for an RID" is also present. The main search area contains a text input field with the following sequence: "aaouacucgc cattatcaco atcgttttgg ggcgatgttg tgtggttoca gqgtattaat ataattaatt tattccacat gagatatgat atgatatact atgtatTTTT ttatttgtaa acotttaata taacaagaac tacaaaaaat gaaaa". Below the input field are links for "Set subsequence" and "Choose database", followed by "From:" and "To:" input fields. At the bottom, there are buttons for "BLAST!" (highlighted in blue), "Reset query", and "Reset all".

BLAST

Basic Local Alignment Search Tool

- Word size: 10-11 bp or 2-3 aa
 - Primary similarities (seed matches)
 - Expanding the homology regions to the left and to the right
- Scoring the homology with matrices PAM (Point Accepted Mutation) or BLOSUM (BLOCKS Substitution Matrix)
- Showing the results

	A	T	G	C
A	1	0	0	0
T	0	1	0	0
G	0	0	1	0
C	0	0	0	1

Annotation: hodnota nepáru G-A points to the (G,A) cell (0)

Annotation: hodnota páru G-G points to the (G,G) cell (1)

Cvrčková, Úvod do praktické bioinformatiky

Matrice PAM 250

S	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
T	-2	1	3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
P	-3	1	0	6	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
A	-2	1	1	1	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
G	-3	1	0	-1	1	5	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
D	-4	1	0	-1	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	
E	-5	0	0	-1	0	0	1	2	4	0	0	0	0	0	0	0	0	0	0	0	
N	-5	-1	0	0	0	-1	2	2	4	0	0	0	0	0	0	0	0	0	0	0	
R	-4	0	-1	0	-2	3	0	-1	1	2	6	0	0	0	0	0	0	0	0	0	
K	-5	0	0	-1	-1	-2	1	0	0	1	0	3	5	0	0	0	0	0	0	0	
M	-5	-2	-1	-2	-1	-3	-2	-3	-2	-1	-2	0	0	6	0	0	0	0	0	0	
I	-2	-1	0	-2	-1	-3	-2	-2	-2	-2	-2	-2	-2	2	5	0	0	0	0	0	
L	-6	-3	-2	-3	-2	-4	-3	-4	-3	-2	-3	-3	-4	2	6	0	0	0	0	0	
V	-2	-1	0	-1	0	-1	-2	-2	-2	-2	-2	-2	-2	2	4	4	0	0	0	0	
F	-4	-3	-3	-5	-4	-5	-4	-5	-5	-2	-4	-5	0	1	2	-3	9	0	0	0	
Y	0	-3	-3	-5	-5	-4	-4	-4	-4	-4	-4	-4	-2	-3	-1	-2	7	10	0	0	
W	-8	-2	-9	-8	-8	-7	-8	-7	-8	-7	-8	-7	-3	-2	-3	-4	-5	-2	8	0	17

BLAST

Basic Local Alignment Search Tool

>gi|5016088|ref|NM_001101.2| Length = 1793 E= expectancy value actin, beta (ACTB), mRNA

Score = 1110 bits (560), Expect = 0.0
Identities = 965/1100 (87%)
Strand = Plus / Plus

Query: 156 gtcgacaacggctctggcatgtgcaaggccggatttgcggagacgatgctccccggccc 215
Sbjct: 101 gtcgacaacggctccggcatgtgcaaggccggcttcgggggacgatgccccggggcc 160

Query: 216 gtcttcccatcgattgtgggaogtcccogtcaccaggggtgtgatggctggcag 275
Sbjct: 161 gtcttccctccatcggtggggcggccaggcaccaggggtgtgatggctggcag 220

Query: 276 aaggactcgtacgtgggtgatgaggccagagcaagcgtggtatcctcacccctgaagtac 335
Sbjct: 221 aaggatcctatgtggggacgaggccagagcaagagaggaatcctcacccctgaagtac 280

Query: 336 cccattgagcacggtatcgtgaccaactgggacgatggagaagatctggcaccacacc 395
Sbjct: 281 cccatcgagcacggcatcgtcaccactgggacgatggagaagatctggcaccacacc 340

ds..S=1213 E=0.0
>=200
250 1500

- „expectancy value“ provides the number of expected sequence number with the same or higher similarity when searching in the database consisting of randomly assembled sequences
- the results shows fraction of identical and in case of proteins also similar sequence positions and/or inserted spaces

Primary Databases

The screenshot shows a web browser displaying a GenBank record for NC_002377.1 (145K..148K (2.9Kbp)). A tooltip is open over the gene NP_059797.1, providing the following information:

- NP_059797.1: two-component VirA-like sensor kinase
- total range: NC_002377.1 (145,694..148,183)
- total length: 2,490
- strand: plus
- protein product length: 829

Links & Tools

- GenBank View: [NC_002377.1 \(145,694..148,183\)](#), [NP_059797](#)
- FASTA View: [NC_002377.1 \(145,694..148,183\)](#), [NP_059797](#)
- BLAST Genomic: [NC_002377.1 \(145,694..148,183\)](#)
- Graphical View: [NP_059797.1](#)
- BLAST Protein: [NP_059797.1](#)
- BLINK Results: [NP_059797.1](#)

Below the tooltip, there are sections for **Bibliography** and **Related articles in PubMed**.

63



BLINK is a link to the pre-computed BLAST search results for the respective sequence (see the next slide).

BLAST

Basic Local Alignment Search Tool

Pre-computed BLAST results for: [gi|16119781|ref|NP_396486.1](#) two component sensor kinase [Agrobacterium tumefaciens str. C58]

Matching gis: [15162425.20141871.101960](#)

Total (score > 100) : 147068 hits in 148754 proteins in 6309 species

Selected: 147068 hits in 148754 proteins in 6309 species Filter: Min Score: 100 |

Other views (Reports): [Taxonomy report](#) | [Multiple Alignment](#) | [Blast](#)

[Reset all filters](#)

Choose Display Options

1203 Archaea 138095 Bacteria 13 Metazoa 1348 Fungi 594 Plants 6 Viruses 5676 The Others [reset selection](#)

Results: 1 - 100 [Next Page](#) [Last](#)

Rank	% hit	Score	Accession	Length	Protein Description
1	100	41.66	AA090927	833	two component sensor kinase [Agrobacterium tumefaciens str. C58]
2	100	41.66	F18549	833	hectane: full-wide host range vira protein Short-WIR vira
3	100	41.66	AA079282	833	vira [Pisamid pPIC58]
4	100	41.59	NP_053330	833	hypothetical protein pTI-SAK09A_p142 [Agrobacterium tumefaciens]
5	100	41.59	AA017465	833	tiorf140 [Agrobacterium tumefaciens]
6	100	41.53	AA032109	833	vira [Pisamid T1]
7	100	41.53	gi1737127	833	vira protein
8	100	41.53	CAA14727	833	91.2 kDa protein [Agrobacterium tumefaciens]
9	100	39.00	CAA32269	829	vira [Agrobacterium shizogense]
10	100	37.18	gi1227240	849	vira gene
11	100	33.48	AA086413	829	vira [Pisamid T1]

BLAST

Specialized Versions

- Currently there exists a lot of specialized versions of [BLAST](#)
 - Searching according to source (organism) of sequences, e.g. known genomes of [microorganisms](#)
 - **BLASTP**
 - Given the [protein query](#), it returns the most similar protein sequences from the [protein database](#).
 - **BLASTN**
 - Given the [DNA query](#), it returns the most similar DNA sequences from the [DNA database](#).
 - Other variants, e.g. [MEGABLAST](#), for identification of identical or [very similar sequences](#) (searches [long similar regions](#) of nucleotide sequences)
 - **BLASTX**
 - Compares the all possible [six-frame translation products](#) of a [nucleotide query sequence](#) (both strands) against a [protein sequence database](#).

BLAST

Specialized Versions

- Currently there exists a lot of specialized versions of BLAST
 - **TBLASTN**
 - Compares a [protein query](#) against the [all six reading frames](#) of a [nucleotide sequence database](#).
 - **TBLASTX**
 - Translates the [query nucleotide sequence](#) in [all six possible frames](#) and [compares](#) it against the [six-frame translations](#) of a [nucleotide sequence database](#).

BLAST

Specialized Versions

- Currently there exist a lot of [specialized versions](#) of BLAST
 - [PSI-BLAST](#) ([P](#)osition-[S](#)pecific [I](#)terated [B](#)last)
 - [First step](#): [standard BLAST](#), during which PSI-BLAST identifies a [list of similar sequences](#) with [E value better than minimal value](#) (standard = 0,005)
 - For every alignment, PSI-BLAST creates so-called [PSSM](#) ([P](#)osition [S](#)pecific [S](#)ubstitution [M](#)atrix)
 - [PSSM](#) takes into account [relative frequency of specific aminoacid residue in a specific position](#) within sequences identified as similar in first step, which can mean functional conservation.

BLAST

Specialized Versions

- Currently there exists a lot of specialized versions of BLAST
 - **PHI-BLAST** (Pattern-Hit Initiated BLAST)
 - For identification of **specific sequence**, e.g. motif (pattern) in sequence of similar protein sequences
 - Sequence of motif must be inserted using **special syntax**:
 - [LVIMF] means either Leu, Val, Ile, Met or Phe
 - - is spacer (means nothing)
 - x(5) means 5 positions in which any residue is allowed
 - x(3, 5) means 3 to 5 positions where any residue is allowed

BLAST

Specialized Versions

□ Example of search by PHI-BLAST

```
>gi|4758958|ref|NP_004148.1| Human cAMP-dependent protein kinase  
MSHIQIPPGLTPELLQGYTVBVLRQQPPDLVEFAVEYFTRLREARAPASVLPAAATPRQSLGHPPPPEPGPDR  
VADAKGDSSEEDLEVPVPSRFNRRVSVCAETYNPDEEEEDTDPRVIHPKIDEQRCRLQBACKDILLF  
KNLDQEQLSQVLDAMFERIVKADEHVIDQGDDGDNFYVIERGTYDILVTKDNQTRSVGQYDNRGSRGELA  
LMYNTPRAAITVA TSEGSLWGLDRVTFRRIIVKNNAKRKMFPESPIESVPLKSLVSERMKIVDVIgek  
IYKdGERIITQGEKADSFYIIESGevSILIRSRtKSNKdGNgQEvE IARCHKqQYFGELALVtNKpRAAS  
AYAVGDVVKCLVMDVQAFERLLGpCMDIMKRNISHYEEQLVKMFGSSVDLGNLgQ  
  
[LIVMF]-G-E-x-[GAS]-[LIVM]-x(5,11)-R-[STAQ]-A-x-[LIVMA]-x-[STACV].
```

Outline

- Syllabus Of The Course
- Definition Of Genomics
- Role Of Bioinformatics In Functional Genomics
- Databases
 - Spectre Of „On-line“ Resources
 - PRIMARY, SECONDARY And STRUCURAL Databases
 - GENOME Resources
- Analytical Tools
 - Homologies Searching
 - Searching Of Sequence Motifs, Open Reading Frames, Restriction Sites...

Analytical Tools

<https://blog.addgene.org/free-online-molecular-biology-tools>

Early Career Researcher Toolbox: Free Online Molecular Biology Tools

By Beth Kenkel



Beth Kenkel
September 12, 2023

Share this article



Primer design. Plasmid mapping. DNA sequence analysis. We all have our favorite tools for tackling these particular tasks, but they tend to be scattered about the internet. To help you keep your virtual molecular biology toolbox organized, today's post features a list of free online molecular biology tools all in one place.

Plasmid mapping

These tools are for viewing, editing or making plasmid maps, but can also analyze and annotate any DNA sequence.

- **SnapGene Viewer**: The free SnapGene Viewer is great for looking at plasmid maps and viewing sequencing traces, while the paid version provides more tools for plasmid mapping and design (Figure 1).
- **Benchling**: While you might think of Benchling as an electronic lab notebook, it also has a suite of molecular biology tools and can make plasmid maps. Free for academic users.
- **Serial Cloner**: Free desktop-based software for plasmid design and mapping.
- **ApE (A plasmid Editor)**: A free, donation-based plasmid analysis tool including editing, annotating, creating maps, and more. This tool is maintained by M. Wayne Davis from the University of Utah.

SnapGene

<https://www.snapgene.com/snapgene-viewer/download>



<https://www.youtube.com/watch?v=0sQh2s182WQ>

Outline

- Syllabus Of The Course
- Definition Of Genomics
- Role Of Bioinformatics In Functional Genomics
- Databases
 - Spectre Of „On-line“ Resources
 - PRIMARY, SECONDARY And STRUCURAL Databases
 - GENOME Resources
- Analytical Tools
 - Homologies Searching
 - Searching Of Sequence Motifs, Open Reading Frames, Restriction Sites...
 - **Other On-line Genome Tools**

Other On-Line Genome Resources

- Online Mendelian Inheritance in Man (OMIM)



Summary

- Syllabus Of The Course
- Definition Of Genomics
- Role Of Bioinformatics In Functional Genomics
- Databases
 - Spectre Of „On-line“ Resources
 - PRIMARY, SECONDARY and STRUCURAL Databases
 - GENOME Resources
- Analytical Tools
 - Homologies Searching
 - Searching Of Sequence Motifs, Open Reading Frames, Restriction Sites...
 - Other On-line Genome Tools

Discussion