# CG020 Genomika

# Lesson 1
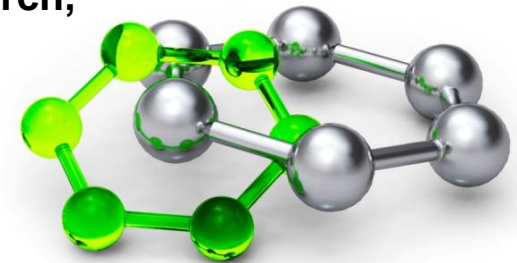## Introduction into Bioinformatics

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

CEITEC

# Course Syllabus

- **Lesson 01**
  - Introduction into Bioinformatics

- **Lesson 02**
  - Identification of Genes

- **Lesson 03**
  - Reverse Genetics Approaches

- **Lesson 04**
  - Forward Genetics Approaches

CEITEC

# 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

CEITEC

# 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

CEITEC

# 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

CEITEC

# Outline

- Syllabus of thecourse

- Definition of Genomics

CEITEC

# 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

CEITEC

# GENOMICS – What is it?
## The role of BIOINFORMATICS in FUNCTIONAL GENOMICS

**Forward („classical") Genetics Approaches**

**Reverse Genetics Approaches**

5'TTATATATATATATTAAAAAATAAAATAA
AAGAACAAAAAAGAAAATAAAATA....3'

BIOINFORMATICS

Insertional mutagenesis

exon
intron
transmembrane region
duplication
En-1 element

1      G1 F G2      4494

FUNCTIONAL GENOMICS

3      :      1

exon
intron
transmembrane region
duplication
En-1 element

1      G1 F G2      4494

**En-1**

gaattcaagtcgtCACTACAAGA/TCTTGTAGTGcgtggagact

1122      1123

...aat tca agt cgt gga gac tac act...
N   S   S   R   G   D   Y   T

?

CEITEC

# Outline

- Syllabus of this course

- Definition of genomics

- **Role of BIOINFORMATICS in FUNCTIONAL GENOMICS**

CEITEC

# Bioinformatics



- **Definiction 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.**

CEITEC

# 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

CEITEC

# 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

CEITEC

# Outline

- Syllabus of this course

- Definition of genomics

- Role of BIOINFORMATICS in FUNCTIONAL GENOMICS

- Databases
    - Spectre of „on-line" resources

CEITEC

# Spectre of On Line Resources

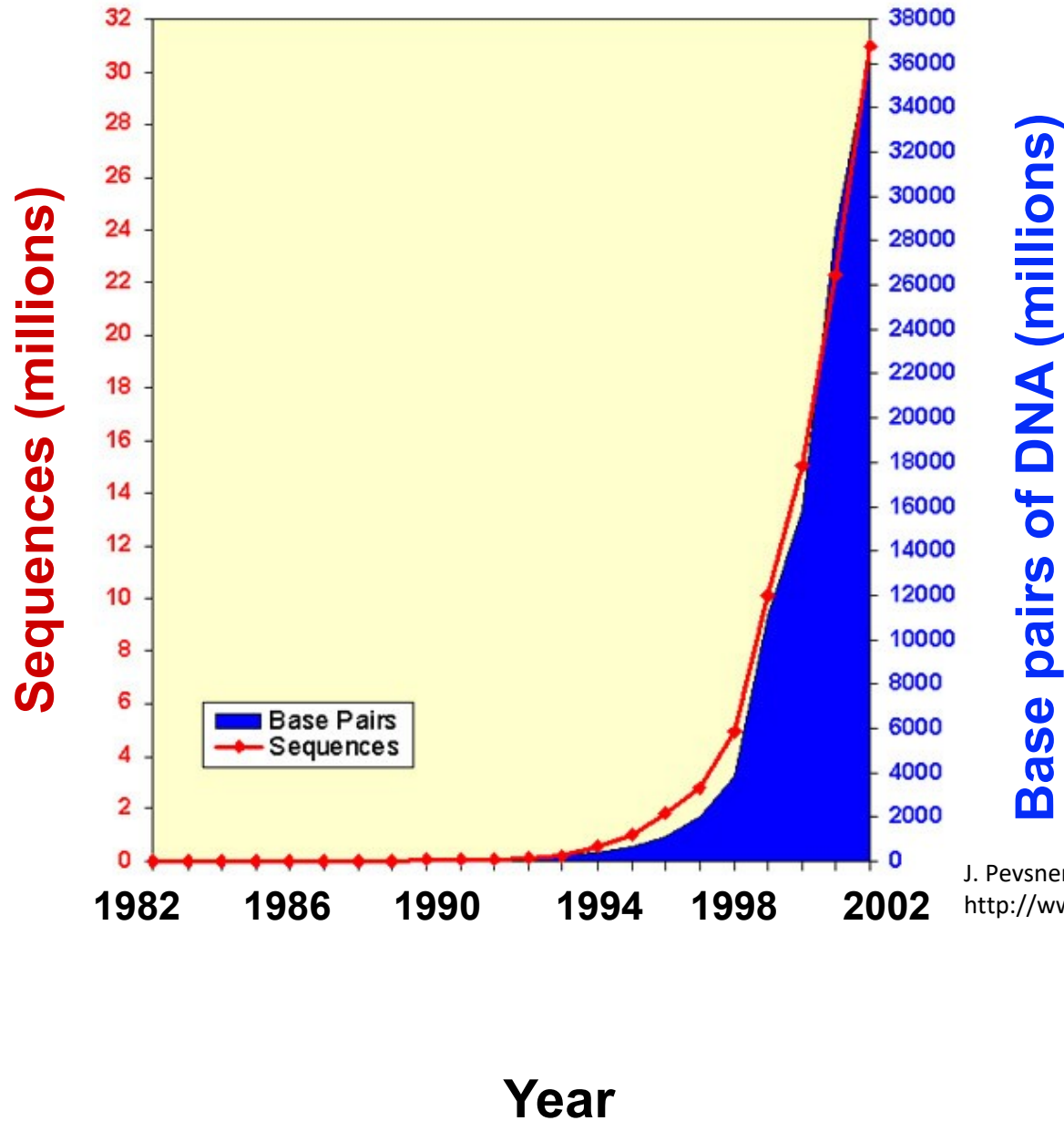**EMBnet National Nodes**

| | | |
|---|---|---|
| Vienna Biocenter | Austria | http://www.at.embnet.org/ |
| BEN | Belgium | http://www.be.embnet.org/ |
| BioBase | Denmark | http://biobase.dk/ |
| CSC | Finland | http://www.fi.embnet.org/ |
| INFOBIOGEN | France | http://www.infobiogen.fr/ |
| GENIUSnet | Germany | http://genome.dkfz-heidelberg.de/biounit/ |
| IMBB | Greece | http://www.imbb.forth.gr/ |
| HEN | Hungary | http://www.hu.embnet.org/ |
| INCBI | Ireland | http://acer.gen.tcd.ie/ |
| INN | Israel | http://dapsas.weizmann.ac.il/bcd/inn.html |
| IEN-ADR | Italy | http://bio-www.ba.cnr.it:8000/BioWWW/Bio-WWW.htm |
| CAOS/CAMM | Netherlands | http://www.caos.kun.nl/ |
| Bio | Norway | http://www.no.embnet.org/ |
| IBB | Poland | http://www.ibb.waw.pl/ |
| IGC | Portugal | http://www.igc.gulbenkian.pt/ |
| GeneBee | Russia | http://www.genebee.msu.su/ |
| CNB-CSIC | Spain | http://www.es.embnet.org/ |
| BMC | Sweden | http://www.embnet.se/ |
| SIB | Switzerland | http://www.ch.embnet.org/ |
| SEQNET | UK | http://www.seqnet.dl.ac.uk/ |

**EMBnet Specialist Nodes**

| | | |
|---|---|---|
| MIPS | Germany | http://www.mips.biochem.mpg.de/ |
| ICGEB | Italy | http://www.icgeb.trieste.it/ |
| Pharmacia Upjohn | Sweden | http://www.pnu.com/ |
| F.Hoffmann-La Roche | Switzerland | http://www.roche.com/ |
| EBI | UK | http://www.ebi.ac.uk/ |
| HGMP-RC | UK | http://www.hgmp.mrc.ac.uk/ |
| Sanger | UK | http://www.sanger.ac.uk/ |
| UMBER | UK | http://www.bioinf.man.ac.uk/dbbrowser |

**EMBnet Associate Nodes**

| | | |
|---|---|---|
| IBBM | Argentina | http://sol.biol.unlp.edu.ar/embnet |
| ANGIS | Australia | http://www.angis.su.oz.au/ |
| CBI | China | http://www.cbi.pku.edu.cn/ |
| CIGB | Cuba | http://bio.cigb.edu.cu/ |
| CDFD | India | http://salarjung.embnet.org.in/ |
| SANBI | South Africa | http://www.sanbi.ac.za |

**USA Information Providers**

| | | |
|---|---|---|
| NCBI | USA | http://www.ncbi.nlm.nih.gov/ |
| NLM | USA | http://www.nlm.nih.gov/ |
| NIH | USA | http://www.nih.gov/ |

CEITEC

# Spectre of On Line Resources

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

CEITEC

# Spectre of On Line Resources

☐ NCBI http://www.ncbi.nlm.nih.gov/

CEITEC

# 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**
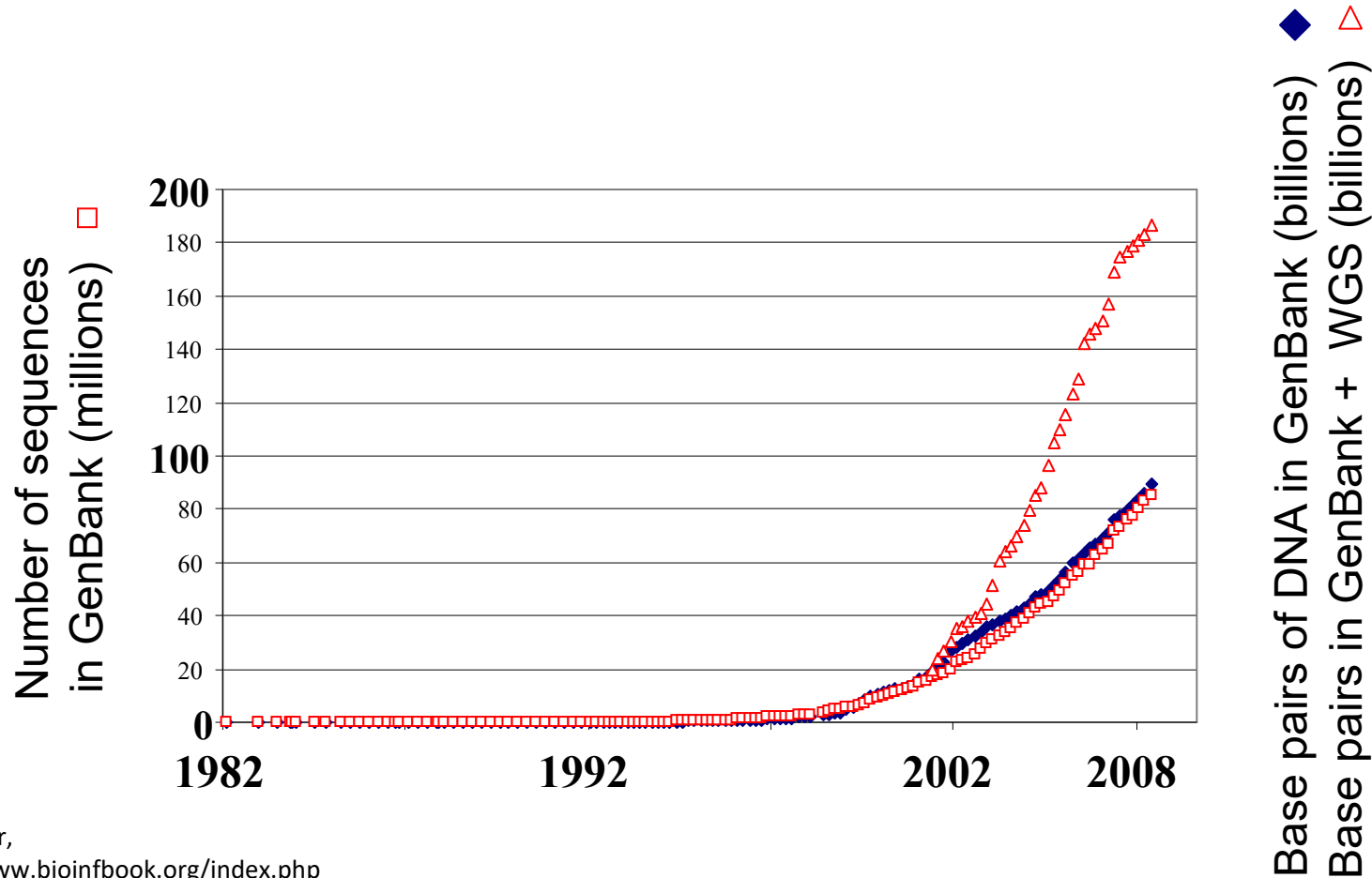
CEITEC

# Primary Databases

- Include primary datasets – <u>DNA</u> and <u>Protein</u> 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 $10^6$ entries (approx. 33 x $10^9$ bp)
  - August 2005 100 x $10^9$ bp from 165.000 organisms

CEITEC

# Growth of GenBank



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

20

CEITEC

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



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

CEITEC

# Growth of GenBank
## Aug 2016

**Bases**



**Sequences**



- Dec **1982** 680 338 bp, 606 sequences

- Apr **2002** 19 x 10$^9$ bp, 17 x 10$^6$ sequences + WGS 692 x 10$^6$ bp, 172 768 sequences

- Aug **2016** 218 x 10$^9$ bp, 196 x 10$^6$ sequences + WGS 1,6 10$^{12}$ bp, 360 x 10$^6$ sequences

CEITEC

# WGS



Shotgun Sequencing

Fig 1: Genomic DNA is fragmented, ligated into viral DNA and packaged into viral particles to create a library

Computer Aided Sorting

Fig 2: Short fragments of DNA sequence are ordered by overlapping data to recreate the whole genome sequence

Interactive concepts in biochemistry, Rodney Boyer, Wiley, 2002, http://www.wiley.com//college/boyer/0470003790/

CEITEC

# Growth of DNA Sequence in Repositories



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

# Growth of DNA Sequence in Repositories



Legend:
- ◆ GenBank bases
- □ WGS bases
- ▲ Sequence Read Archive bases
- + SRA Open Access bases
- — Total bases sequenced in 2014 at major sequencing centers

A vast amount of sequence data has been generated using next-generation sequencing.

# Growth of DNA Sequence in Repositories



Legend:
- ◆ GenBank bases
- □ WGS bases
- ▲ Sequence Read Archive bases
- + SRA Open Access bases
- — Total bases sequenced in 2014 at major sequencing centers

Y-axis: Bases ($\log_{10}$ scale), from 100,000 to 100,000,000,000,000,000

Right-axis markers: $1 \times 10^{15}$ 1 petabase; $1 \times 10^{12}$ 1 terabase; $1 \times 10^{9}$ 1 gigabase

X-axis: Year — 1985, 1990, 1995, 2000, 2005, 2010, 2015

Perhaps 40 petabases (corresponding to 10 mil. human genomes) of DNA were generated in calendar year 2014 at major sequencing centers.

# 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/

CEITEC

# Primary Databases

- Types of sequences in primary databases

  - Standard nucleotide sequences acquired by high quality sequencing

  - **EST**s (**E**xpressed **S**equence **T**ags)

  - **HGTS** (**H**igh **T**hroughput **G**enome **S**equencing)
    - Results of sequencing projects without annotation

  - **Reference Sequences** of annotated genomes

  - **TPA**s (**T**hird **P**arty **A**nnotation)

    - sequences annotated by third party (by someone else, not the orginal authors)

CEITEC

# Primary Databases

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

CEITEC

# Primary Databases

# Primary Databases

# Primary Databases

# Primary Databases

```
/translation="MNGRYSPTRQDFKTGAKPWSILALIVAAMIFAFMAVASWQDNAT
TQAILSQLRSINADSASLQRDVLRAHTGTVANYRPIISRLGALRKNLEDLKQLPRQSH
IVSESNAAQLLRQLEVSLNSADAAVAAPGAQNVRLQDSLASPTRALSSLPGKASTDQT
LEKPTELASMMLQFLRQPSPAISFEISLELERLQKQRGLDEAPVRILAREGPIILSLL
PQVKDLVNMIQTSDTAEIAEMLQRECLEVYSLKNVEERSARIFLGSASVGLCLYIITL
VYRLRKKTDWLARRLDYEELIKEIGVCPBGEAATTSSAQAALRIIQRPFDADTCALAL
VDHDRRWAVETPGAKHPKPVWDDSVLREIVSRTKADERATVFRIISSKKIVHLPLEIP
GLSILLAHKSTDKLIAVCSLGYQSYRPRPCQGEIQLLELATACLCHYIDVRRKQTRCD
VLARRLEHAQRLEAVGTLAGGIAHEFNNILGSILGHAELAQNSVSRTSVTRRYIDYII
SSGDRAMLIIDQILTLSRKQERMIKPPSVSELVTEIAPLLRMALPPNIELSPRPDQMQ
SVIBGSPLELQQVLINICKNASQAMTANGQIDIIISQAFLPVKKILAHGVMPPGDYVL
LSISDNGGGIPEAVLPHIFEPPFTTRARNGGTGLGLASVHGHISAPAGYIDVSSTVGH
GTRPDIYLPPSSKEPVNPDSFFGRNKAPRGNGEIVALVEPDDLLREAYEDKIAALGYE
PVGPRTFNEIRDWISKGNEADLVMVDQASLPEDQSPNSVDLVLKTASIIIGGNDLKMT
LSREDVTRDLYLPKPISSRTMAHAILTKIKT"
ORIGIN
        1 atgaacggaa gatattcacc gacgcggcag gattttaaga caggcgcgaa gccttggtct
       61 atattggccc ttatcgttgc tgcaatgatt ttcgcgttca tggcgggttgc gtcctggcag
      121 gacaatgcga ctacccaggc aatcctcagc caactacgat cgattaacgc cgacagcgcc
      181 tcactgcagc gcgatgtact ccgcgctcac acgggcaccg tggcgaacta ccgccccatt
      241 atctccaggc tgggagctct gcggaagaat ctggaagatt tgaagcaatt atttagacaa
      301 tctcatattg taagtgagag caatgctgct caactgctac gccagctaga agtgtctcta
      361 aattcggctg acgcggcggt cgccgccttt ggtgcgcaaa atgtacgcct gcaagattcg
      421 ctggccagtt tcactcgtgc tttgagcagt cttccaggaa aagcctcaac cgatcagact
      481 ttagaaaaac caacagaatt ggctagcatg atgctccaat ttcttcggca accaagcccg
      541 gctatttcat tcgagatcag ccttgaacta gagaggctcc aaaaacaacg cggtcttgat
      601 gaagctcccg tgcgcatact tgcacgtgaa ggtcccatta tcttatcgct tttgccacag
      661 gtgaaagatc tggtgaacat gattcagacg tctgacaccg cagaaattgc ggagatgctg
      721 cagcgcgagt gtttggaggt ctatagcttg aaaaatgtag aggagcggag cgcacgtatc
      781 tttcttgggt ccgcttcagt gggtctttgc ctctacatca tcaccttagt ctataggcta
      841 cgcaaaaaaa ccgattggtt agcgcggcgt ttagattacg aagagctaat caaagagatc
      901 ggagtatgtt ttgaaggtga ggcggccacc acgtcgtccg cgcaagctgc acttcgtatt
      961 attcagcgct tctttgatgc cgatacgtgc gcgttagctc tagtggacca tgaccgtaga
     1021 tgggctgtcg aaacattcgg tgcgaaacac ccaaaacctg tgtgggacga cagcgtgcta
     1081 cgcgaaatag tctctcgtac caaagcggac gaacgggcga cggtattccg catcatatcg
     1141 tcgaaaaaaa tcgtacattt gcctctcgaa attccaggtc tctcgatact actggctcac
     1201 aaatccacag ataaactaat tgcggtttgt tcactgggtt accaaagcta tcgccctcga
     1261 ccttgccaag gcgaaattca gcttcttgaa ctcgccaccg cctgcctctg tcactatatc
     1321 gatgttcggc gtaagcagac cgaatgcgac gttttggcca gacgattgga gcatgcgcaa
     1381 cgccttgagg cagttggtac acttgccggc ggaatagcac atgaatttaa taacattttg
     1441 tcgggcacgc agaattagca caaaactcgg tgtctcgaac atctgtcacc
     1501 cgaagatata ttgactatat catttcgtca ggcgacagag ccatgctcat tatcgatcag
     1561 atcttgacgc tgagccgaaa acaggagcgc atgatcaagc catttagtgt ctcagagctt
     1621 gtgaccgaaa tcgctccctt gctacgtatg gctcttccgc caaacatcga gcttagtttc
     1681 agatttgatc aaatgcagag cgtgatcgaa ggaagcccgc ttgaacttca acaggtacta
     1741 attaacatct gcaagaatgc ttcccaagcc atgactgcaa atggtcaaat cgacatcatc
     1801 atcagccaag ctttttttacc agttaagaaa attctggcgc atggtgttat gccacctggc
     1861 gactatgttc tcctatctat tagcgacaat ggtggaggca ttcccgaggc tgtgttaccc
     1921 cacatttttg aacccttctt tacgacacga gctcgcaacg gtggaacggg tctcggcctt
     1981 gcttctgtgc atggtcatat cagcgcgttt gcgggttaca tcgacgttag ttcaactgtt
     2041 gggcatggga cgcgctttga catttatctc cctccgtctt ctaaggaacc cgtaaatcca
     2101 gacagttttt tcggccgcaa taaggcaccg cgtggaaacg gggagattgt ggcacttgtt
     2161 gagcccgatg acctcctgcg ggagcgcgtat gaagacaaga tcgccgctct aggatatgag
     2221 ccggtcggct ttcgtacctt taatgaaatt cgcgattgga tttcaaaagg caatgaagcc
     2281 gatctggtca tggtcgacca agcgtctctt cctgaagatc aaagtcctaa ttccgtggat
     2341 ttagtgctca agaccgcctc catcatcatt gcgcggaaatg atctcaaaat gacccttttca
```

# 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 | |
| NT_030059 | Genomic contig | **DNA** |
| 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 | |
| AAC02945 | GenBank protein | **Protein** |
| Q28369 | SwissProt protein | |
| 1KT7 | Protein Data Bank structure record | |

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

34

CEITEC

# 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

CEITEC

# RefSeq

# 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

CEITEC

# Primary Databases

# Primary Databases

# Secondary Databases

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

- **PROSITE**, http://www.expasy.org/prosite/

# Secondary Databases

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

- **PROSITE**, http://www.expasy.org/prosite/

```
>PDOC00003 PS00003 SULFATION Tyrosine sulfation site [rule] [Warning: rule with a high probability of occurrence].

    571 - 585   nkeesstYeteisns

>PDOC00004 PS00004 CAMP_PHOSPHO_SITE cAMP- and cGMP-dependent protein kinase phosphorylation site [pattern] [Warning: pattern with a high probability of occurrence].

    744 - 747   RRvT
    814 - 817   KRrS

>PDOC00005 PS00005 PKC_PHOSPHO_SITE Protein kinase C phosphorylation site [pattern] [Warning: pattern with a high probability of occurrence].

     148 -  150   SsR
     164 -  166   TgR
     171 -  173   StK
     219 -  221   SkK
     369 -  371   TrR
     460 -  462   SgK
     513 -  515   SgR
     585 -  587   SiR
     602 -  604   TgK
     652 -  654   TdK
     716 -  718   SpR
     726 -  728   SpK
     747 -  749   TeK
     794 -  796   SsR
     854 -  856   ScK
     864 -  866   StR
     868 -  870   SeR
     921 -  923   SpK
     957 -  959   SvR
     960 -  962   TgR
     974 -  976   TsK
     997 -  999   SrK
    1002 - 1004   TgK
    1018 - 1020   SgK
    1031 - 1033   TgR
    1119 - 1121   SkR
```

# Secondary Databases

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

- **PROSITE**, http://www.expasy.org/prosite/

>PDOC50109 PS50109 **HIS_KIN** Histidine kinase domain [profile].

```
402 - 671   NASHDIRGALAGMKGLIDICRDGVKPGSDVDTTLNQVNVCAKDLVALLNSVLDMSKIESG
            KMQLVKEDPNLSKLLEDVIDFYHPVAMKKGVDVVLDPHDgsvfKPSNVRGDSGRLKQILN
            NLVSNAVKFTVD--GHIAVRAWAQrpgsnssvvlasypkgvskfvksmfcknkeesstye
            teisnsirnnanTMEFVPEVDDTGKGIPMEMRKSVPENYVQVREtAQGHQGTGLGLGIVQ
            SLVRLMGGEIRITDKAMGekGTCPQPNVLLTT
```

>PDOC50110 PS50110 **RESPONSE_REGULATORY** Response regulatory domain [profile].

```
987 - 1085  RVLVVDDNPISRKVATGKLKKMGVSeVEQCDSGKEALRLVTEGLtqreeqgsvdklpFDY
            IFMDCQMPEMDGYEATREIRkvekSYGVRTPIIAVSGHD--------------------
            --------------
```

**Graphical summary of hits** *(java applet)*

Click on items to see a description. Drag the two red cursors to select a zoom region.          About   Prefs

[graphical applet area: HIS_KIN ... RESPONSE_REGU]

0   Zoom   Back   Reset   1123   |||||||: 100 residues

**98 hits with 12 PROSITE entries**

| ExPASy Home page | Site Map | Search ExPASy | Contact us | Swiss-Prot | PROSITE | Proteomics tools |

CEITEC

# 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 compendium of protein **fingerprints**. A fingerprint is a group of conserved motifs used to characterise a protein family; its diagnostic power is refined by iterative scanning of a *SWISS-PROT/TrEMBL* composite. Usually the motifs do not overlap, but are separated along a sequence, though they may be contiguous in 3D-space. Fingerprints can encode protein folds and functionalities more flexibly and powerfully than can single motifs, full diagnostic potency deriving from the mutual context provided by motif neighbours. **References**

**New:**

- SPRINT - *Search PRINTS-S (relational PRINTS)*
- prePRINTS - *Search PRINTS' automatic supplement*
- InterPro - *Search the integrated InterPro family database*

**Direct PRINTS access:**

- By accession number
- By PRINTS code
- By database code
- By text
- By sequence
- By title
- By number of motifs
- By author
- By query language

**PRINTS search:**

- Search PRINTS with NEW FingerPRINTScan
- FPScan
- GRAPHScan
- MULScan
- FingerPRINTScan binaries and source are available: contact scordis@bioinf.man.ac.uk

43

CEITEC

# Secondary Databases

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



Scaffold/Matrix Attached Region transaction Database

CEITEC

# Structural Databases

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

CEITEC

# Structural Databases

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

CEITEC

# Structural Databases

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



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

47

# 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**

CEITEC

# Genome Resources

❑ **Human Genome Browser** http://genome.ucsc.edu/cgi-bin/hgGateway

# Genome Resources

☐ **Human Genome Browser** http://genome.ucsc.edu/cgi-bin/hgGateway



50

# Genome Resources

☐ **Human Genome Browser** http://genome.ucsc.edu/cgi-bin/hgGateway

CEITEC

# Genome Resources

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

# Genome Resources

☐ **Human Genome Browser** http://genome.ucsc.edu/cgi-bin/hgGateway

# Genome Resources

❑ The Arabidopsis Information Resource (TAIR) http://www.arabidopsis.org

CEITEC

# Genome Resources

☐ **TAIR, The Arabidopsis Information Resource**, http://www.arabidopsis.org

# 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**

CEITEC

# Analytical Tools

□ **Global** versus **Local** alignment



```
Globální přiřazení
SLAV-----------APATNIK-------PIQNYR-I------AKSETQRYMVIE
SLAVYTYIEFVRANAPATNIKSECVRAAPIQNYRRVEHVRATAKSETQRYMVIE

Lokální přiřazení
SLAVYTYIEFVRANAPATNIKSECVRAAPIQNYRRVEHVRATAKSETQRYMVIE
--------------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

CEITEC

# 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 treshold

# 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 neccessary to do a local alignment

- Dotplot can be obtained using BLAST2 (see further in the presentation)

CEITEC

# Analytical Tools

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

# BLAST

**B**asic **L**ocal **A**lignment **S**earch **T**ool

- 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

Matice PAM 250



Cvrčková, Úvod do praktické bioinformatiky

CEITEC

# BLAST

### Basic Local Alignment Search Tool



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

E= expectancy value

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

```
Query: 156  gtcgacaacggctctggcatgtgcaaggccggatttgccggagacgatgctccccgcgcc 215
            |||||||||||||| ||||||||||||||||||| ||| |||||||| |||||| |||
Sbjct: 101  gtcgacaacggctccggcatgtgcaaggccggcttcgcgggcgacgatgccccccgggcc 160

Query: 216  gtcttcccatcgattgtgggacgtccccgtcaccagggtgtgatggtcggcatgggccag 275
            ||||||| || || |||| || ||| | ||||||||| ||||||||| |||||||| |||
Sbjct: 161  gtcttcccctccatcgtggggcgccccaggcaccagggcgtgatggtgggcatgggtcag 220

Query: 276  aaggactcgtacgtgggtgatgaggcgcagagcaagcgtggtatcctcaccctgaagtac 335
            ||||| || || ||||| | ||| ||||||||||||| | ||||||||||||||||||||
Sbjct: 221  aaggattcctatgtgggcgacgaggcccagagcaagagaggcatcctcaccctgaagtac 280

Query: 336  cccattgagcacggtatcgtgaccaactgggacgatatggagaagatctggcaccacacc 395
            ||||| ||||||||| ||||| |||||||||||||| ||||||| |||||||||||||||
Sbjct: 281  cccatcgagcacggcatcgtcaccaactgggacgacatggagaaaatctggcaccacacc 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 whe  searching in the database consisiting of randomly assembled sequences

- the results shows fraction of identical and in case of proteins also similar sequence positions and/or inserted spaces

CEITEC

# Primary Databases

# BLAST

**B**asic **L**ocal **A**lignment **S**earch **T**ool

# 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.

CEITEC

# BLAST

## Specialized Versions

- ☐ Currently there exist a lot of specialized versions of BLAST

  - ▪ **PSI-BLAST** (**P**osition-**S**pecific **I**terated Blast)
    - • 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 (Position Specific Substitution Matrix)

    - • 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.

CEITEC

# BLAST

## Specialized Versions

- Currently there exists a lot of specialized versions of BLAST

  - **PHI-BLAST** (**P**attern-**H**it **I**nitiated 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

CEITEC

# BLAST

Specialized Versions

- ☐ Example of search by PHI-BLAST

```
>gi|4758958|ref|NP_004148.1| Human cAMP-dependent protein kinase
MSHIQIPPGLTELLQGYTVEVLRQQPPDLVEFAVEYFTRLREARAPASVLPAATPRQSLGHPPPEPGPDR
VADAKGDSESEEDEDLEVPVPSRFNRRVSVCAETYNPDEEEDTDPRVIHPKTDEQRCRLQEACKDILLF
KNLDQEQLSQVLDAMFERIVKADEHVIDQGDDGDNFYVIERGTYDILVTKDNQTRSVGQYDNRGSFGELA
LMYNTPRAATIVATSEGSLWGLDRVTFRRIIVKNNAKKRKMFESFIESVPLLKSLEVSERMKIVDVIGEK
IYKDGERIITQGEKADSFYIIESGEVSILIRSRTKSNKDGGNQEVEIARCHKGQYFGELALVTNKPRAAS
AYAVGDVKCLVMDVQAFERLLGPCMDIMKRNISHYEEQLVKMFGSSVDLGNLGQ

[LIVMF]-G-E-x-[GAS]-[LIVM]-x(5,11)-R-[STAQ]-A-x-[LIVMA]-x-[STACV].
```

CEITEC

# 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…**

CEITEC

# 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.

CEITEC

# SnapGene

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

CEITEC

# 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…
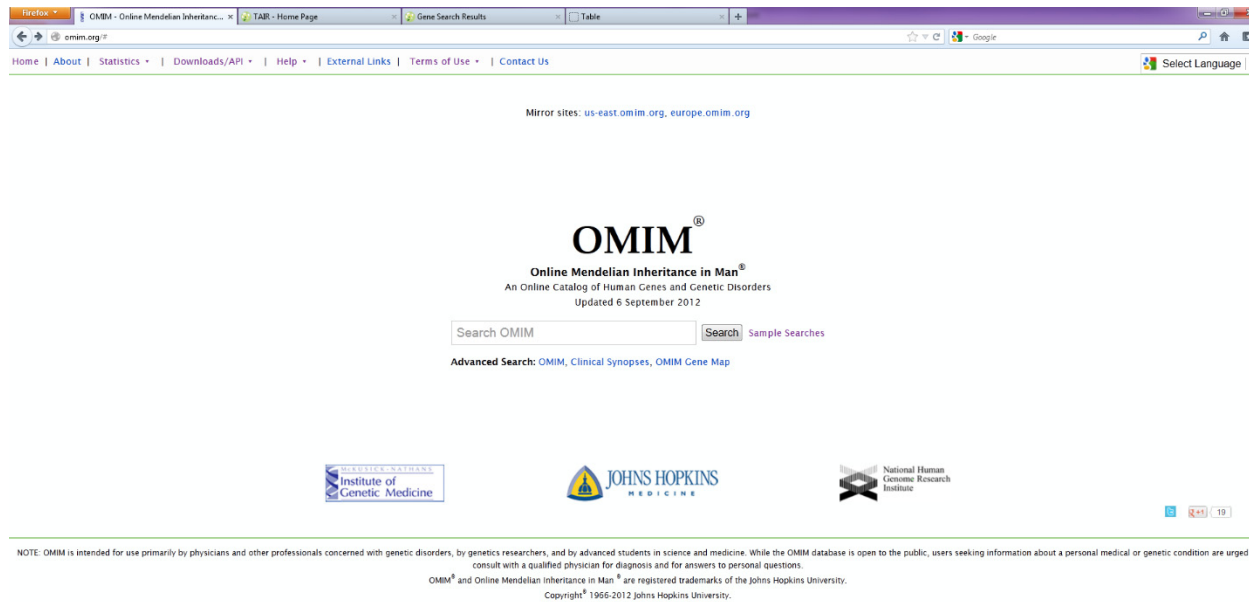    - **Other On-line Genome Tools**

CEITEC

# Other On-Line Genome Resources

- **TIGR** (**T**he **I**nstitute for **G**enomic **R**esearch, http://www.tigr.org/software/)
  - Recently part of the J. Craig Venter Institute

# Other On-Line Genome Resources

- **O**nline **M**endelian **I**nheritance in **M**an (**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

CEITEC

# Discussion

CEITEC