# Computer genomics lectures

- **Internet-resources for DNA sequences search and analysis**

- **General description of bioinformatics**

- **Genes and genomes databases**

- **Genomic DNA: superposition of many codes**

- **Repeats in DNA, formal description and analysis**

- **Consensus, position weight matrices, sequence alignment methods**

- **Machine learning methods overview: hidden markov models, neural networks, genetic algorithm**

- **Methods of eukaryotic gene structure recognition**

- **Statistical approaches used for different recognition methods accuracy comparison**

- **Investigation of qualitative and quantitative characteristics of transcriptome**

# Internet-resources for DNA sequences search and analysis

## Major bioinformatics databases

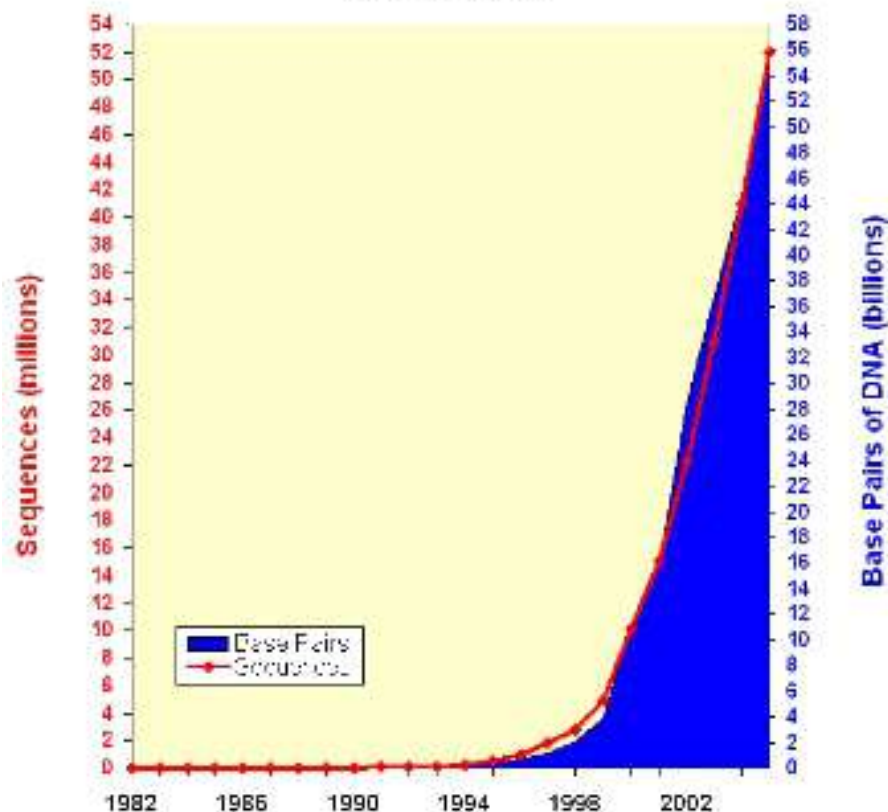| | | |
|---|---|---|
| GenBank<br>EMBL<br>DDBJ | http://www.pubmed.com<br>http://www.ebi.ac.uk/embl.html<br>http://www.ddbj.nig.ac.jp | Nucleotide and Protein Sequences |
| PubMed | http://www.pubmed.com | Bibliographic database |
| EnsEMBL<br>UCSC Genome Browser | http://www.ensembl.org<br>http://genome.ucsc.edu | Genes and genomes analysis and annotation.<br>Genes and genomes search and visualization tool |
| RefSeq | http://www.ncbi.nlm.nih.gov/RefSeq | Non-redundant sequence database of genomes, transcripts and proteins |
| UniGene<br>STACK<br>GeneCards<br>GenAtlas<br>GeneOntology<br>TIGR Gene Indices | http://www.ncbi.nlm.nih.gov/UniGene<br>http://www.sanbi.ac.za/Dbases.html<br>http://www.genecards.org<br>http://www.dsi.univ-paris5.fr/genatlas<br>http://www.geneontology.org<br>http://www.tigr.org/tdb/tgi.shtml | Gene Database |
| SWISSPROT | http://www.expasy.ch | Protein Database |
| EPD | http://www.epd-isb-sib.ch | Eukaryotic Promoter Database |

# General description of bioinformatics

**Basic directions of bioinformatics development**

- Homology search, multiple alignment

- Statistical analysis of genetic texts, genome segmentation

- Recognition of coding sequences and open reading frames

- Recognition of DNA functional sites

- Analysis of RNA secondary structure

- Analysis of protein sequences, protein secondary structure prediction, recognition of functional sites in proteins

- Phylogenetic analysis

- DNA-chips, DNA microarray: expression analysis

- Database surfing: manipulation with a large amount of data

# Genes and genomes databases

## Genome annotation progress

### Growth of GenBank
(1982 - 2005)



### Map Viewer - genome annotation updates:

| Species | Build | Map Viewer Release |
|---|---|---|
| Rattus norvegicus | RGSC v3.4 | July 6, 2006 |
| Macaca mulatta (rhesus macaque) | 1.1 | June 23, 2006 |
| Caenorhabditis elegans | WS150 | May 11, 2006 |
| Mus musculus | 36.1 | May 8, 2006 |
| Drosophila melanogaster | 4.3 | April 19, 2006 |
| Tribolium castaneum (red flour beetle) | 1.1 | April 18, 2006 |
| Homo sapiens | 36.1 | March 9, 2006 |
| Dictyostelium discoideum | 1.1 | November 22, 2005 |
| Arabidopsis thaliana | TAIR6.0 | November 21, 2005 |
| Bos taurus (cow) | 2.1 | October 12, 2005 |
| Canis familiaris (dog) | 2.1 | September 8, 2005 |
| Strongylocentrotus purpuratus (sea urchin) | 1.1 | August 17, 2005 |
| Danio rerio (zebrafish) | Zv4 | July 5, 2005 |
| Anopheles gambiae (mosquito) | 2.2 | June 30, 2005 |
| Apis mellifera (bee) | 2.1 | May 31, 2005 |
| Pan troglodytes (chimpanzee) | 1.1 | November 23, 2004 |
| Gallus gallus (chicken) | 1.1 | August 11, 2004 |

▸ The Human Genome

**The Human Genome**
The Human Genome Project generated an unprecedented amount of knowledge about human genetics. Explore human genome resources, browse the human genome sequence using the Map Viewer.

**Organism-Specific**
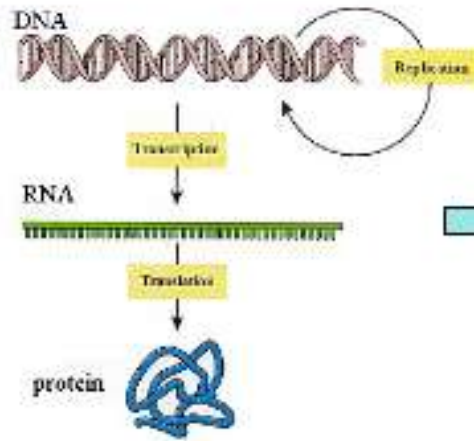- Genome Resources
- BLAST
- Map Viewer
- Genome Project DB

- ▸ Arabidopsis
- ▸ Aspergillus
- ▸ Bee
- ▸ Beetle NEW
- ▸ Cat
- ▸ Chicken
- ▸ Chimpanzee
- ▸ Cow
- ▸ Dictyostelium
- ▸ Dog
- ▸ Frog
- ▸ Fruit Fly
- ▸ Human
- ▸ Malaria
- ▸ Mosquito
- ▸ Mouse
- ▸ Nematode
- ▸ Pig
- ▸ Rabbit
- ▸ Rat
- ▸ Rhesus macaque NEW
- ▸ Sea Urchin
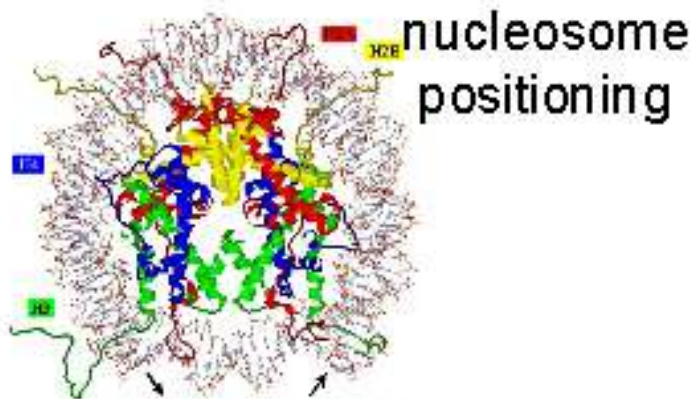- ▸ Sheep
- ▸ Yeast (Saccharomyces)
- ▸ Zebrafish

http://www.ncbi.nih.gov/Genbank/genbankstats.html    http://www.ncbi.nlm.nih.gov/Genomes/

# Genomic DNA: superposition of many codes



= **genomic DNA codes?**

nucleosome positioning

transcription regulation

state 1

state 2          TSS

→ product 1

→ product 2

higher order packaging
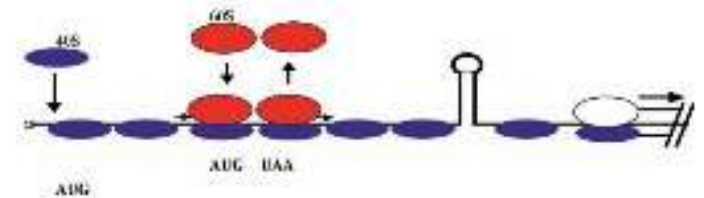
DNA conformation

translation regulation

etc...

# Repeats in DNA, formal description and analysis

## Repeats: basic types

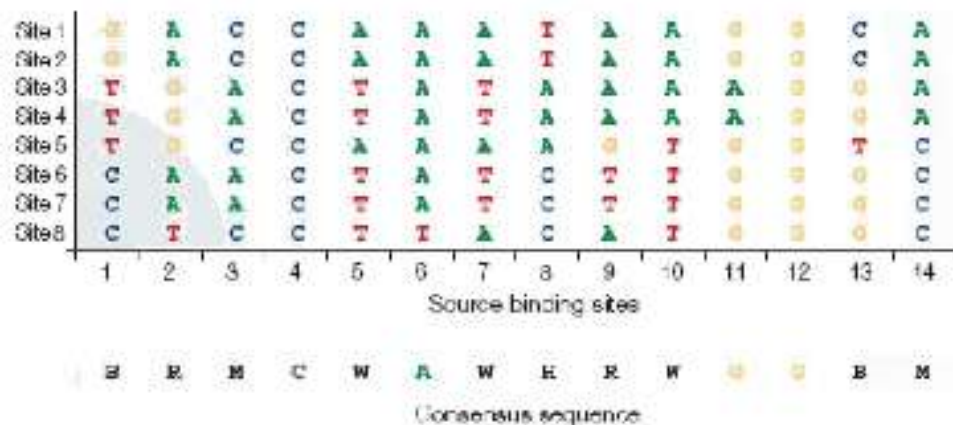| Type | Example | | Direction | Complementarity |
|---|---|---|---|---|
| Direct | AGCTTT<br>TCGAAA | AGCTTT<br>TCGAAA | Forward | No |
| Invert | AGCTTT<br>TCGAAA | AAAGCT<br>TTTCGA | Reverse | Yes |
| Symmetric | AGCTTT<br>TCGAAA | TTTCGA<br>AAAGCT | Reverse | No |
| Direct complementary | AGCTTT<br>TCGAAA | TCGAAA<br>AGCTTT | Forward | Yes |
| Palindrome | AAGCCGAA<br>TTCGGCTT | | Reverse | No |
| Complementary palindrome | AAGCGCTT<br>TTCGCGAA | | Reverse | No |

## Repeats: possible mutual positioning

| Type | Example |
|---|---|
| Dispersed | …AGTTC…..AGTTC… |
| Tandem | …AGTTCAGTTC… |
| Overlapped | …AGTTCAGTTCAGTTC… |

# Consensus, position weight matrices, sequence alignment methods

## Position weight matrix (PWM) model



Source binding sites

Consensus sequence

Position frequency matrix (PFM)

Position weight matrix (PWM)

Site scoring

$\Sigma = 5.23, 78\%$ of maximum

**A set of aligned binding sites**

**Consensus model**

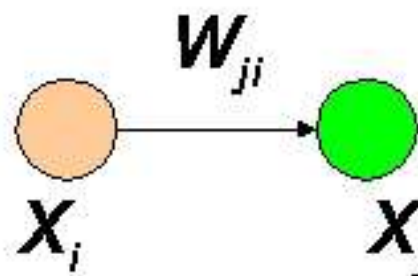**Position frequency matrix (PFM): the count of observed nucleotides at each position**

**PFM is converted to a position weight matrix (PWM) using a special formula**

**Using a PWM model, a score for any DNA sequence can be calculated by summation over all positions**

• Wasserman WW, Sandelin A. Applied bioinformatics for the identification of regulatory elements. Nat Rev Genet. 2004; 276-287.

# Machine learning methods overview: neural networks, genetic algorithm

## Artificial neural network: oriented multigrapf of artificial neurons with weighted connections

Input  Hidden  Output

weights

node

Information →

$W_{ji}$

$X_i$  $X_j$

$$W_{ij} \leftarrow W_{ij} + \Delta W_{ij}$$

**Network parameter**

weight of connection between the neurons $x_i$ and $x_j$.

**Network learning**

Weight modification according to the learning paradigm (supervised, unsupervised, Hebbian, reinforcement...)

Learned network with the fixed weights presents the knowledge about the world

## Program GenScan scheme

N intergenic region
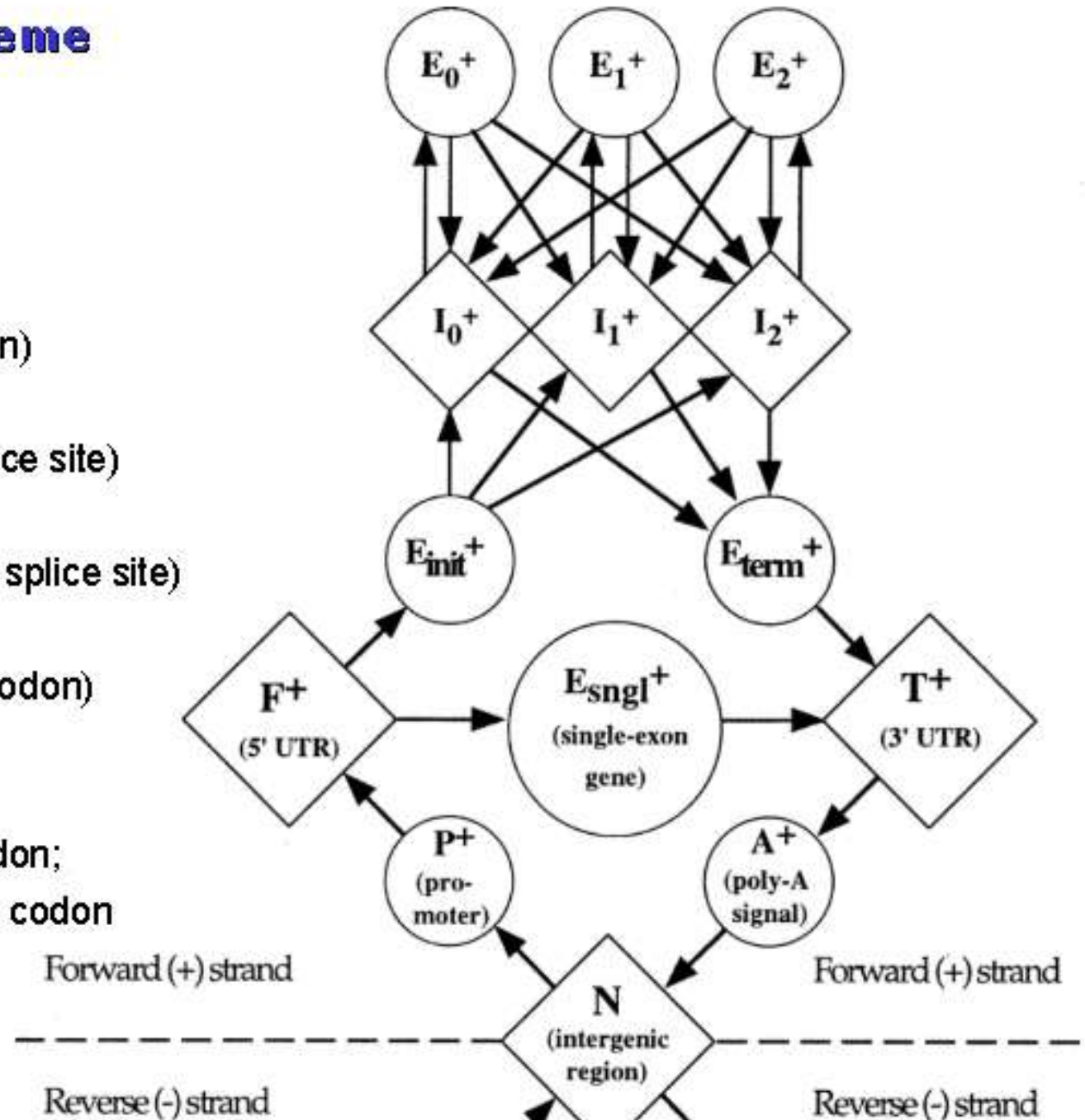
P promoter

F 5'-untranslated region

$E_{sngl}$ single exon (intronless)
(translation start –> stop codon)

$E_{init}$ initial exon
(translation start –> donor splice site)

$E_k$ phase k internal exon
(acceptor splice site –> donor splice site)

$E_{term}$ terminal exon
(acceptor splice site –> stop codon)

$I_k$ phase k intron:
0 – between codons;
1 – after the first base of a codon;
2 – after the second base of a codon

$E_0^+$   $E_1^+$   $E_2^+$

$I_0^+$   $I_1^+$   $I_2^+$

$E_{init}^+$   $E_{term}^+$

$F^+$
(5' UTR)

$E_{sngl}^+$
(single-exon gene)

$T^+$
(3' UTR)

$P^+$
(pro-moter)

$A^+$
(poly-A signal)

Forward (+) strand

$N$
(intergenic region)

Forward (+) strand

Reverse (-) strand

Reverse (-) strand

# Statistical approaches used for different recognition methods accuracy comparison

**Predictions**

Probability

50% — Threshold

YES (real sites)
NO (real non-sites)

Negatives:

Positives:

40%

| True negatives | True positives |
| False negatives | False positives |

30%

false negatives

false positives

20%

true negatives

true positives

10%

0%

**Recognition function score**

**Contingency table**

| | | Prediction | |
| --- | --- | --- | --- |
| | | Non-sites | Sites |
| Reality | Sites | FN, false negatives Not predicted real sites | TP, true positives Correctly predicted real sites |
| | Non-sites | TN, true negatives Correctly predicted non-sites | FP, false positives Real non-sites predicted as sites |

# Investigation of qualitative and quantitative characteristics of transcriptome

## Comparision of methods for transcripts detection and abundance estimation

| Method | Relative/absolute measurability, compatibility | Genes amount | Sensitivity & dynamic range | High throughput capacity |
|---|---|---|---|---|
| **Direct mRNA detection via hybridization of transcripts with ssDNA or RNA probes** | | | | |
| Nothern blot hybridization | + | (1-5) x (5- 20) | - - | - - |
| Ribonuclease protection | + + | (10-15) x (5- 20) | + | - - |
| **Detection of cDNA made by reverse transcription from mRNA** | | | | |
| Quantitative RT-PCA, Real time RT-PCA | + + | tens | + + (!) | + + |
| Differential display | + | tens | + + (!) | + + |
| Oligonucleotide/cDNA microarrays | + + (!) | thousands | + + + | + + + |
| **Computational analysis of cDNA reads, «in silico hybridization» of transcripts** | | | | |
| SAGE | + + + | thousands | + + + | + + + |
| MPSS | + + + + | thousands | + + + + | + + + + |
| EST | + | thousands | + + + (!) | + + + |

(!) Caution about possibility for nonlinearly distorted transcripts abundance estimations