# Introduction to Bioinformatics

Esa Pitkänen
esa.pitkanen@cs.helsinki.fi
Autumn 2008, I period
www.cs.helsinki.fi/mbi/courses/08-09/itb

**MBI** MASTER'S DEGREE
PROGRAMME IN BIOINFORMATICS

582606 Introduction to Bioinformatics, Autumn 2008

---

# Introduction to Bioinformatics

Lecture 1:
Administrative issues
MBI Programme, Bioinformatics courses
What is bioinformatics?
Molecular biology primer

---

## How to enrol for the course?

- p Use the registration system of the Computer Science department: https://ilmo.cs.helsinki.fi
  - n You need your user account at the IT department ("cc account")
- p If you cannot register yet, don't worry: attend the lectures and exercises; just register when you are able to do so

3

---

## Teachers

- p Esa Pitkänen, Department of Computer Science, University of Helsinki
- p Elja Arjas, Department of Mathematics and Statistics, University of Helsinki
- p Sami Kaski, Department of Information and Computer Science, Helsinki University of Technology
- p Lauri Eronen, Department of Computer Science, University of Helsinki (exercises)

4

---

## Lectures and exercises

- p Lectures: Tuesday and Friday 14.15-16.00 Exactum C221

- p Exercises: Tuesday 16.15-18.00 Exactum C221
  - n First exercise session on Tue 9 September

5

---

## Status & Prerequisites

- p Advanced level course at the Department of Computer Science, U. Helsinki
- p 4 credits
- p Prerequisites:
  - n Basic mathematics skills (probability calculus, basic statistics)
  - n Familiarity with computers
  - n Basic programming skills recommended
  - n No biology background required

6

## Course contents

- p What is bioinformatics?
- p Molecular biology primer
- p Biological words
- p Sequence assembly
- p Sequence alignment
- p Fast sequence alignment using FASTA and BLAST
- p Genome rearrangements
- p Motif finding (tentative)
- p Phylogenetic trees
- p Gene expression analysis

7

## How to pass the course?

- p Recommended method:
  - n Attend the lectures (not obligatory though)
  - n Do the exercises
  - n Take the course exam
- p Or:
  - n Take a separate exam

8

## How to pass the course?

- p Exercises give you max. 12 points
  - n 0% completed assignments gives you 0 points, 80% gives 12 points, the rest by linear interpolation
  - n "A completed assignment" means that
    - p You are willing to present your solution in the exercise session and
    - p You return notes by e-mail to Lauri Eronen (see course web page for contact info) describing the main phases you took to solve the assignment
  - n Return notes at latest on Tuesdays 16.15
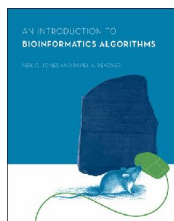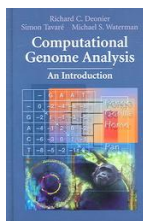
- p Course exam gives you max. 48 points

9

## How to pass the course?

- p Grading: on the scale 0-5
  - n To get the lowest passing grade 1, you need to get at least 30 points out of 60 maximum
- p Course exam: Wed 15 October 16.00-19.00 Exactum A111
- p See course web page for separate exams
- p Note: if you take the first separate exam, the best of the following options will be considered:
  - n Exam gives you 48 points, exercises 12 points
  - n Exam gives you 60 points
- p In second and subsequent separate exams, only the 60 point option is in use
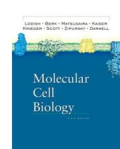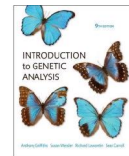
10

## Literature

- p Deonier, Tavaré, Waterman: Computational Genome Analysis, an Introduction. Springer, 2005
- p Jones, Pevzner: An Introduction to Bioinformatics Algorithms. MIT Press, 2004
- p Slides for some lectures will be available on the course web page



11

## Additional literature

- p Gusfield: Algorithms on strings, trees and sequences
- p Griffiths et al: Introduction to genetic analysis
- p Alberts et al.: Molecular biology of the cell
- p Lodish et al.: Molecular cell biology

- p Check the course web site



12

## Questions about administrative & practical stuff?

13

## Master's Degree Programme in Bioinformatics (MBI)

- p Two-year MSc programme
- p Admission for 2009-2010 in January 2009
  - n You need to have your Bachelor's degree ready by August 2009

**MBI** MASTER'S DEGREE PROGRAMME IN BIOINFORMATICS          www.cs.helsinki.fi/mbi

| News & events | Programme | Studies | Admission | People | Contact | |

News and Events

## MBI programme organizers



Department of Computer Science, Department of Mathematics and Statistics Faculty of Science, Kumpula Campus, HY

Laboratory of Computer and Information Science, Laboratory of CS and Engineering,TKK

Faculty of Biosciences Faculty of Agriculture and Forestry Viikki Campus, HY

Faculty of Medicine, Meilahti Campus, HY

15

## Four MBI campuses



HY, Viikki

HY, Meilahti

HY, Kumpula

TKK, Otaniemi

## MBI highlights

- p You can take courses from both HY and TKK
- p Two biology courses tailored specifically for MBI
- p Bioinformatics is a new exciting field, with a high demand for experts in job market

- p Go to www.cs.helsinki.fi/mbi/careers to find out what a bioinformatician could do for living
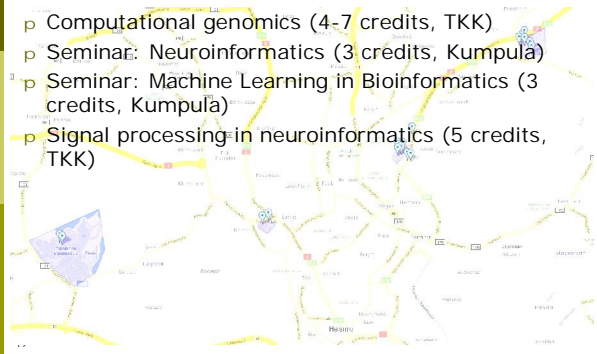
17

## Admission

- p Admission requirements
  - n Bachelor's degree in a suitable field (e.g., computer science, mathematics, statistics, biology or medicine)
  - n At least 60 ECTS credits in total in computer science, mathematics and statistics
  - n Proficiency in English (standardized language test: TOEFL, IELTS)
- p Admission period opens in late Autumn 2009 and closes in 2 February 2009
- p Details on admission will be posted in www.cs.helsinki.fi/mbi during this autumn

18

## Bioinformatics courses in Helsinki region: 1st period

- p Computational genomics (4-7 credits, TKK)
- p Seminar: Neuroinformatics (3 credits, Kumpula)
- p Seminar: Machine Learning in Bioinformatics (3 credits, Kumpula)
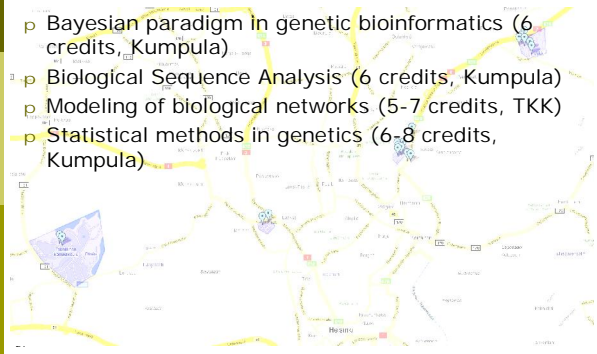- p Signal processing in neuroinformatics (5 credits, TKK)

## A good biology course for computer scientists and mathematicians?

- p Biology for methodological scientists (8 credits, Meilahti)
  - n Course organized by the Faculties of Bioscience and Medicine for the MBI programme
  - n Introduction to basic concepts of microarrays, medical genetics and developmental biology
  - n Study group + book exam in I period (2 cr)
  - n Three lectured modules, 2 cr each
  - n Each module has an individual registration so you can participate even if you missed the first module
  - n www.cs.helsinki.fi/mbi/courses/08-09/bfms/

20

## Bioinformatics courses in Helsinki region: 2nd period

- p Bayesian paradigm in genetic bioinformatics (6 credits, Kumpula)
- p Biological Sequence Analysis (6 credits, Kumpula)
- p Modeling of biological networks (5-7 credits, TKK)
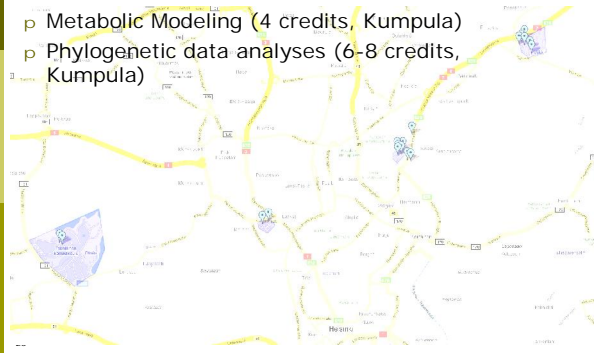- p Statistical methods in genetics (6-8 credits, Kumpula)

## Bioinformatics courses in Helsinki region: 3rd period

- p Evolution and the theory of games (5 credits, Kumpula)
- p Genome-wide association mapping (6-8 credits, Kumpula)
- p High-Throughput Bioinformatics (5-7 credits, TKK)
- p Image Analysis in Neuroinformatics (5 credits, TKK)
- p Practical Course in Biodatabases (4-5 credits, Kumpula)
- p Seminar: Computational systems biology (3 credits, Kumpula)
- p Spatial models in ecology and evolution (8 credits, Kumpula)
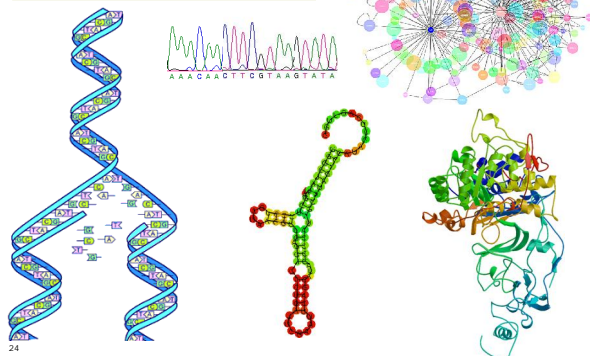- p Special course in bioinformatics I (3-7 credits, TKK)

## Bioinformatics courses in Helsinki region: 4th period

- p Metabolic Modeling (4 credits, Kumpula)
- p Phylogenetic data analyses (6-8 credits, Kumpula)

# 1. What is bioinformatics?

24

## What is bioinformatics?

p Bioinformatics, *n.* The science of information and information flow in biological systems, esp. of the use of computational methods in genetics and genomics. (Oxford English Dictionary)

p "The mathematical, statistical and computing methods that aim to solve biological problems using DNA and amino acid sequences and related information."                    -- Fredj Tekaia

## What is bioinformatics?

p "I do not think all biological computing is bioinformatics, e.g. mathematical modelling is not bioinformatics, even when connected with biology-related problems. In my opinion, bioinformatics has to do with management and the subsequent use of biological information, particular genetic information."
-- Richard Durbin

## What is *not* bioinformatics?

p Biologically-inspired computation, e.g., genetic algorithms and neural networks
p However, application of neural networks to solve some biological problem, could be called bioinformatics
p What about DNA computing?



*http://www.wisdom.weizmann.ac.il/~lbn/new_pages/Visual_Presentation.html*

## Computational biology

p Application of computing to biology (broad definition)
p Often used interchangeably with bioinformatics
p Or: *Biology* that is done with computational means

## Biometry & biophysics

p Biometry: the statistical analysis of biological data
  n Sometimes also the field of identification of individuals using biological traits (a more recent definition)
p Biophysics: "an interdisciplinary field which applies techniques from the physical sciences to understanding biological structure and function"   -- British Biophysical Society

## Mathematical biology

p Mathematical biology "tackles biological problems, but the methods it uses to tackle them need not be numerical and need not be implemented in software or hardware."
-- Damian Counsell

*Alan Turing*

## Turing on biological complexity

p "It must be admitted that the biological examples which it has been possible to give in the present paper are very limited.

This can be ascribed quite simply to the fact that biological phenomena are usually very complicated. Taking this in combination with the relatively elementary mathematics used in this paper one could hardly expect to find that many observed biological phenomena would be covered.
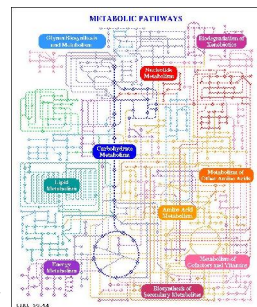
It is thought, however, that the imaginary biological systems which have been treated, and the principles which have been discussed, should be of some help in interpreting real biological forms."

– Alan Turing, The Chemical Basis of Morphogenesis, 1952

31

## Related concepts

p Systems biology
- n "Biology of networks"
- n Integrating different levels of information to understand how biological systems work

p Computational systems biology



*Overview of metabolic pathways in KEGG database, www.genome.jp/kegg/*

32

## Why is bioinformatics important?

p New measurement techniques produce huge quantities of biological data
- n Advanced data analysis methods are needed to make sense of the data
- n Typical data sources produce noisy data with a lot of missing values

p Paradigm shift in biology to utilise bioinformatics in research

33

## Bioinformatician's skill set

p Statistics, data analysis methods
- n Lots of data
- n High noise levels, missing values
- n #attributes >> #data points

p Programming languages
- n Scripting languages: Python, Perl, Ruby, …
- n Extensive use of text file formats: need parsers
- n Integration of both data and tools

p Data structures, databases

34

## Bioinformatician's skill set

p Modelling
- n Discrete vs continuous domains
- n -> Systems biology

p Scientific computation packages
- n R, Matlab/Octave, …

p Communication skills!

35

## Communication skills: case 1



?

Biologist presents a problem to computer scientists / mathematicians

"I am interested in finding what affects the regulation gene x during condition y and how that relates to the organism's phenotype."

"Define input and output of the problem."

36

## Communication skills: case 2

Bioinformatician is a part of a group that consists mostly of biologists.

## Communication skills: case 2

…biologist/bioinformatician ratio is important!

## Communication skills: case 3

A group of bioinformaticians offers their services to more than one group

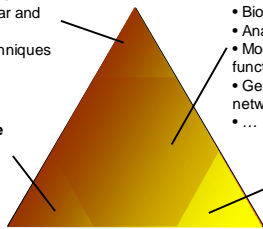## Bioinformatician's skill set

p How much biology you should know?

## Bioinformatician's skill set

**Biology & Medicine**
• Basics in molecular and cell biology
• Measurement techniques

**Computer Science**
• Programming
• Databases
• Algorithmics

**Bioinformatics**
• Biological sequence analysis
• Biological databases
• Analysis of gene expression
• Modeling protein structure and function
• Gene, protein and metabolic networks
• …

**Mathematics and statistics**
• Calculus
• Probability calculus
• Linear algebra

*Where would you be in this triangle?*

*Prof. Juho Rousu, 2006*

## A problem involving bioinformatics?

- "I found a fruit fly that is immune to all diseases!"

- "It was one of these"

*Pertti Jarla, http://www.hs.fi/fingerpori/*

## Molecular biology primer

43

## Molecular biology primer

p Part 1: What is life made of?
p Part 2: Where does the variation in genomes come from?

44

## Life begins with Cell



p A cell is a smallest structural unit of an organism that is capable of independent functioning
p All cells have some common features

45

## Cells

p Fundamental working units of every living system.
p Every organism is composed of one of two radically different types of cells:
  n prokaryotic cells or
  n eukaryotic cells.
p Prokaryotes and Eukaryotes are descended from the same primitive cell.
  n All prokaryotic and eukaryotic cells are the result of a total of 3.5 billion years of evolution.

46

## Two types of cells: Prokaryotes and Eukaryotes



47

## Prokaryotes and Eukaryotes

p According to the most recent evidence, there are three main branches to the tree of life
p Prokaryotes include Archaea ("ancient ones") and bacteria
p Eukaryotes are kingdom Eukarya and includes plants, animals, fungi and certain algae



► Lecture: Phylogenetic trees

48

## All Cells have common Cycles



p Born, eat, replicate, and die

## Common features of organisms

p Chemical energy is stored in ATP
p Genetic information is encoded by DNA
p Information is transcribed into RNA
p There is a common triplet genetic code
p Translation into proteins involves ribosomes
p Shared metabolic pathways
p Similar proteins among diverse groups of organisms

## All Life depends on 3 critical molecules

p DNAs (Deoxyribonucleic acid)
  n Hold information on how cell works
p RNAs (Ribonucleic acid)
  n Act to transfer short pieces of information to different parts of cell
  n Provide templates to synthesize into protein
p Proteins
  n Form enzymes that send signals to other cells and regulate gene activity
  n Form body's major components (e.g. hair, skin, etc.)
  n "Workhorses" of the cell

## DNA: The Code of Life



p The structure and the four genomic letters code for all living organisms
p Adenine, Guanine, Thymine, and Cytosine which pair A-T and C-G on complimentary strands.

► Lecture: Genome sequencing and assembly

## Discovery of the structure of DNA

p 1952-1953 James D. Watson and Francis H. C. Crick deduced the double helical structure of DNA from X-ray diffraction images by Rosalind Franklin and data on amounts of nucleotides in DNA
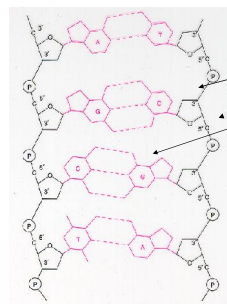


"Photo 51"

Rosalind Franklin

James Watson and Francis Crick

## DNA, continued



p DNA has a double helix structure which is composed of
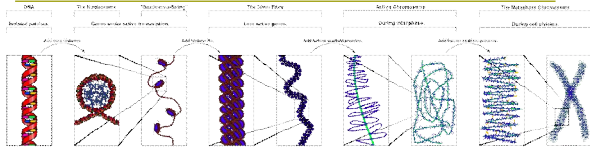  n sugar molecule
  n phosphate group
  n and a base (A,C,G,T)

p By convention, we read DNA strings in direction of transcription: from 5' end to 3' end
  5' ATTTAGGCC 3'
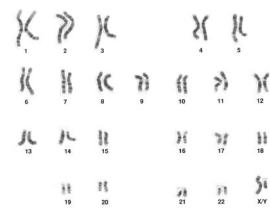  3' TAAATCCGG 5'

## DNA is contained in chromosomes



- In eukaryotes, DNA is packed into *chromatids*
  - In metaphase, the "X" structure consists of two identical chromatids
- In prokaryotes, DNA is usually contained in a single, circular chromosome

---

## Human chromosomes

- Somatic cells in humans have 2 pairs of 22 chromosomes + XX (female) or XY (male) = total of 46 chromosomes
- Germline cells have 22 chromosomes + either X or Y = total of 23 chromosomes



Karyogram of human male using Giemsa staining
(http://en.wikipedia.org/wiki/Karyotype)

56

---

### Length of DNA and number of chromosomes

| Organism | #base pairs | #chromosomes (germline) |
|---|---|---|
| **Prokayotic** | | |
| Escherichia coli (bacterium) | $4 \times 10^6$ | 1 |
| | | |
| **Eukaryotic** | | |
| Saccharomyces cerevisia (yeast) | $1.35 \times 10^7$ | 17 |
| Drosophila melanogaster (insect) | $1.65 \times 10^8$ | 4 |
| Homo sapiens (human) | $2.9 \times 10^9$ | 23 |
| Zea mays (corn / maize) | $5.0 \times 10^9$ | 10 |

57

---

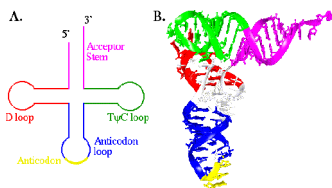Hepatitis delta virus, complete genome

```
   1 atgagccaag ttccgaacaa ggattcgcgg ggaggataga tcagcgcccg agaggggtga
  61 gtcggtaaag agcattggaa cgtcggagat acaactccca agaaggaaaa aagagaaagc
 121 aagaagcgga tgaatttccc cataacgcca gtgaaactct aggaagggga aagagggaag
 181 gtggaagaga aggaggcggg cctcccgatc cgagggggccc ggcggccaag tttggaggac
 241 actccggccc gaaggggttga gagtacccca gagggaggaa gccacacgga gtagaacaga
 301 gaaatcacct ccagaggacc ccttcagcga acagagagcg catcgcgaga gggagtagac
 361 catagcgata ggaggggatg ctaggagttg ggggagaccg aagcgaggag gaaagcaaag
 421 agagcagcgg ggctagcagg tgggtgttcc gcccccgag aggggacgag tgaggcttat
 481 cccggggaac tcgacttatc gtccccacat agcagactcc cggacccccct ttcaaagtga
 541 ccgagggggg tgactttgaa cattggggac cagtggagcc atgggatgct cctcccgatt
 601 ccgcccaagc tccttccccc caagggtcgc ccaggaatgg cgggacccca ctctgcaggg
 661 tccgcgttcc atcctttctt acctgatggc cggcatggtc ccagcctcct cgctggcgcc
 721 ggctgggcaa cattccgagg ggaccgtccc ctcggtaatg gcgaatggga cccacaaatc
 781 tctctagctt cccagagaga agcgagagaa aagtggctct cccttagcca tccgagtgga
 841 cgtgcgtcct ccttcggatg cccaggtcgg accgcgagga ggtggagatg ccatgccgac
 901 cgcaagaga aagaaggacg cgagacgcaa acctgcgagt ggaaacccgc tttattcact
 961 ggggtcgaca actctgggga gaggagggag ggtcggctgg gaagagtata tcctatggga
1021 atccctggct tccccttatg tccagtccct ccccggtccg agtaaagggg gactccggga
1081 ctccttgcat gctgggggacg aagccgccccc cgggcgctcc cctcgttcca ccttcgaggg
1141 ggttcacacc cccaacctgc gggccggcta ttcttctttc ccttctctcg tcttcctcgg
1201 tcaacctcct aagttcctct tcctcctcct tgctgaggtt ctttcccccc gccgatagct
1261 gctttctctt gttctcgagg gccttccttc gtcggtgatc ctgcctctcc ttgtcggtga
1321 atcctccccct ggaaggcctc ttcctaggtc cggagtctac ttccatctgg tccgttcggg
1381 ccctcttcgc cgggggagcc ccctctccat ccttatcttt cttccgaga attcctttga
1441 tgtttcccag ccagggatgt tcatcctcaa gtttcttgat tttcttctta accttccgga
1501 ggtctctctc gagttcctct aacttctttc ttccgctcac ccactgctcg agaacctctt
1561 ctctccccccc gcggttttttc cttcctcgg gccggctcat cttcgactag aggcgacggt
1621 cctcagtact cttactcttt tctgtaaaga ggagactgct ggccctgtcg cccaagttcg
1681 ag
```
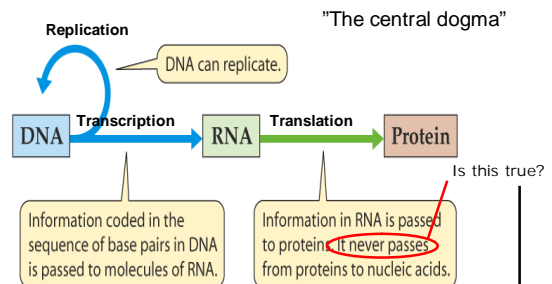
---

## RNA

- RNA is similar to DNA chemically. It is usually only a single strand. T(hyamine) is replaced by U(racil)
- Several types of RNA exist for different functions in the cell.



tRNA linear and 3D view:        http://www.cgl.ucsf.edu/home/glasfeld/tutorial/trna/trna.gif

59

---

## DNA, RNA, and the Flow of Information



"The central dogma"
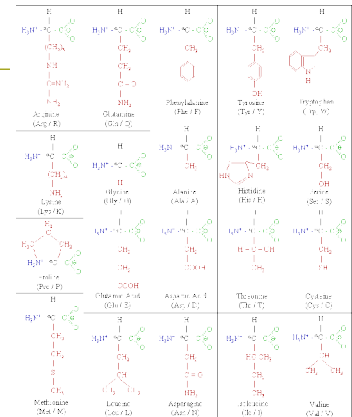
Is this true?

---

## Proteins

- p Proteins are polypeptides (strings of amino acid residues)
- p Represented using strings of letters from an alphabet of 20: AEGLV...WKKLAG
- p Typical length 50...1000 residues



*Urease enzyme from Helicobacter pylori*

61

---

## Amino acids

62

---

## How DNA/RNA codes for protein?

- p DNA alphabet contains four letters but must specify protein, or polypeptide sequence of 20 letters.
- p Dinucleotides are not enough: $4^2 = 16$ possible dinucleotides
- p Trinucleotides (triplets) allow $4^3 = 64$ possible trinucleotides
- p Triplets are also called *codons*



63

---

## How DNA/RNA codes for protein?

- p Three of the possible triplets specify "stop translation"
- p Translation usually starts at triplet AUG (this codes for methionine)
- p Most amino acids may be specified by more than triplet
- p How to find a gene? Look for start and stop codons (not that easy though)



64

---

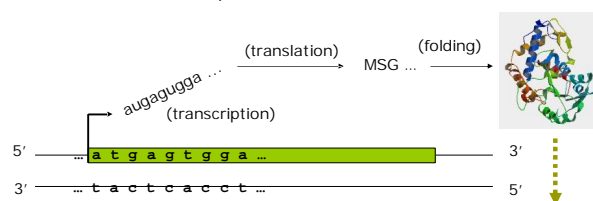## Proteins: Workhorses of the Cell

- p 20 different amino acids
    - n different chemical properties cause the protein chains to fold up into specific three-dimensional structures that define their particular functions in the cell.
- p Proteins do all <u>essential work</u> for the cell
    - n build cellular structures
    - n digest nutrients
    - n execute metabolic functions
    - n mediate information flow within a cell and among cellular communities.
- p Proteins work together with other proteins or nucleic acids as "molecular machines"
    - n structures that fit together and function in highly specific, lock-and-key ways.

► Lecture 8: Proteomics

65

---

## Genes

- p "A gene is a union of genomic sequences encoding a coherent set of potentially overlapping functional products" --Gerstein et al.
- p A DNA segment whose information is expressed either as an RNA molecule or protein
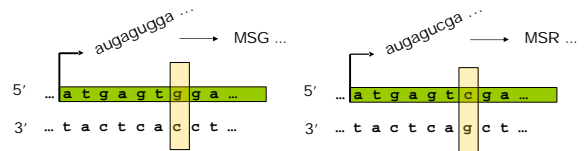


66

http://fold.it
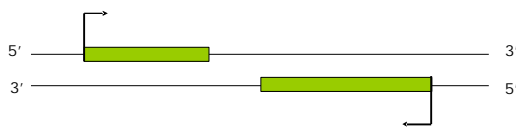
---

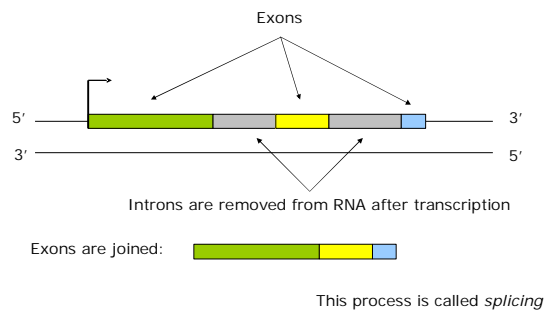## FoldIt: Protein folding game

http://fold.it

## Genes & alleles

p A gene can have different variants

p The variants of the same gene are called *alleles*
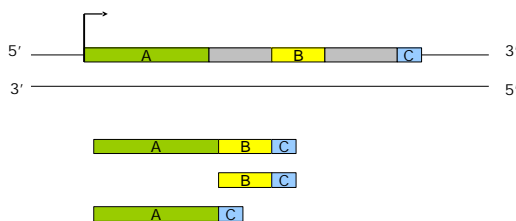


## Genes can be found on both strands



## Exons and introns & splicing

Exons

Introns are removed from RNA after transcription
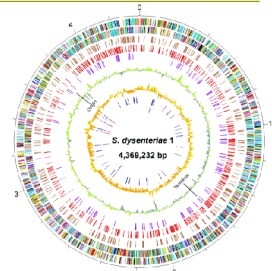
Exons are joined:

This process is called *splicing*

## Alternative splicing

Different *splice variants* may be generated



## Where does the variation in genomes come from?

p Prokaryotes are typically haploid: they have a single (circular) chromosome

p DNA is usually inherited vertically (parent to daughter)

p Inheritance is clonal
  n Descendants are faithful copies of an ancestral DNA
  n Variation is introduced via mutations, transposable elements, and horizontal transfer of DNA

Chromosome map of *S. dysenteriae*, the nine rings describe different properties of the genome
http://www.mgc.ac.cn/ShiBASE/circular_Sd197.htm

## Causes of variation

- p Mistakes in DNA replication
- p Environmental agents (radiation, chemical agents)
- p Transposable elements (transposons)
  - n A part of DNA is moved or copied to another location in genome
- p Horizontal transfer of DNA
  - n Organism obtains genetic material from another organism that is not its parent
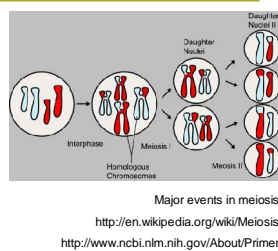  - n Utilized in genetic engineering

## Biological string manipulation

- p Point mutation: substitution of a base
  - n ...ACGGCT... => ...ACGCCT...
- p Deletion: removal of one or more contiguous bases (substring)
  - n ...TTGATCA... => ...TTTCA...
- p Insertion: insertion of a substring
  - n ...GGCTAG... => ...GGTCAACTAG...

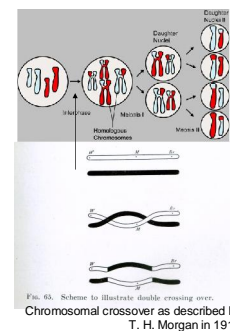► Lecture: Sequence alignment
Lecture: Genome rearrangements

## Meiosis

- p Sexual organisms are usually diploid
  - n Germline cells (gametes) contain N chromosomes
  - n Somatic (body) cells have 2N chromosomes
- p Meiosis: reduction of chromosome number from 2N to N during reproductive cycle
  - n One chromosome doubling is followed by two cell divisions



Major events in meiosis

http://en.wikipedia.org/wiki/Meiosis

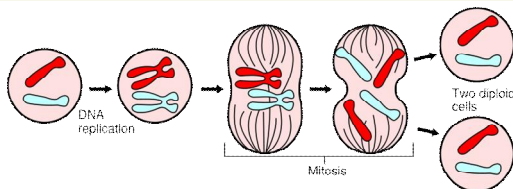http://www.ncbi.nlm.nih.gov/About/Primer

## Recombination and variation

- p Recap: Allele is a viable DNA coding occupying a given locus (position in the genome)
- p In recombination, alleles from parents become suffled in offspring individuals via chromosomal crossover over
- p Allele combinations in offspring are usually different from combinations found in parents
- p Recombination errors lead into additional variations



Fig. 65. Scheme to illustrate double crossing over.

Chromosomal crossover as described by T. H. Morgan in 1916

## Mitosis



- p Mitosis: growth and development of the organism
  - n One chromosome doubling is followed by one cell division

http://en.wikipedia.org/wiki/Image:Major_events_in_mitosis.svg

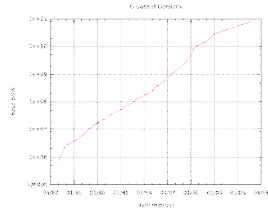## Recombination frequency and linked genes

- p Genetic marker: some DNA sequence of interest (e.g., gene or a part of a gene)

- p Recombination is more likely to separate two distant markers than two close ones

- p Linked markers: "tend" to be inherited together

- p Marker distances measured in centimorgans: 1 centimorgan corresponds to 1% chance that two markers are separated in recombination

## Biological databases

p Exponential growth of biological data

  n New measurement techniques

  n Before we are able to use the data, we need to store it efficiently -> biological databases

  n Published data is submitted to databases

p General vs specialised databases

p This topic is discussed extensively in *Practical course in biodatabases* (III period)

79

## 10 most important biodatabases… according to "Bioinformatics for dummies"

| | | | |
|---|---|---|---|
| p | GenBank/DDJB/EMBL | www.ncbi.nlm.nih.gov | Nucleotide sequences |
| p | Ensembl | www.ensembl.org | Human/mouse genome |
| p | PubMed | www.ncbi.nlm.nih.gov | Literature references |
| p | NR | www.ncbi.nlm.nih.gov | Protein sequences |
| p | UniProt | www.expasy.org | Protein sequences |
| p | InterPro | www.ebi.ac.uk | Protein domains |
| p | OMIM | www.ncbi.nlm.nih.gov | Genetic diseases |
| p | Enzymes | www.expasy.org | Enzymes |
| p | PDB | www.rcsb.org/pdb/ | Protein structures |
| p | KEGG | www.genome.ad.jp | Metabolic pathways |

80
*Sophia Kossida, Introduction to Bioinformatics, Summer 2008*

## FASTA format

p A simple format for DNA and protein sequence data is FASTA

Header line, begins with >

>Hepatitis delta virus, complete genome

```
atgagccaagttccgaacaaggattcgcggggaggatagatcagcgcccgagaggggtga
gtcggtaaagagcattggaacgtcggagatacaactcccaagaaggaaaaaagagaaagc
aagaagcggatgaatttccccataacgccagtgaaactctaggaaggggaaagagggaag
gtggaagagaaggaggcgggcctcccgatccgaggggcccggcggccaagtttggaggac
actccggcccgaagggttgagagtaccccagaggggaggaagccacacggagtagaacaga
gaaatcacctccagaggacccctttcagcgaacagagagcgcatcgcgagagggagtagac
catagcgataggaggggatgctaggagttgggggagaccgaagcgaggaggaaagcaaag
agagcagcggggctagcaggtgggtgttccgcccccgagagggggacgagtgaggcttat
cccggggaactcgacttatcgtccccacatagcagactcccggacccccctttcaaagtga
…
```

81