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: <u>https://ilmo.cs.helsinki.fi</u>
 - 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

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)

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

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

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)
- Phylogenetic trees
- p Gene expression analysis

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

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

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

Literature

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





Additional literature

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







Lodish • Berk • Matsudaira • Kaiser Rieger • Scott • Zipursky • Darnell Molecular Cell Biology

Questions about administrative & practical stuff?

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



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

Four MBI campuses



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

Admission

p Admission requirements

- 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

Bioinformatics courses in Helsinki region: 1st period



Tolo Reservano 51 Lennang 5

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/

Bioinformatics courses in Helsinki region: 2nd period

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



Bioinformatics courses in Helsinki region: 3rd period



Bioinformatics courses in Helsinki region: 4th period





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."

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?



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
 - 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

THE CHEMICAL BASIS OF MORPHOGENESIS

By A. M. TURING, F.R.S. University of Manchester

(Received 9 November 1951-Revised 15 March 1952)

It is suggested that a system of chemical substances, called morphogens, reacting together and diffusing through a tissue, is adequate to account for the main phenomena of morphogenesis. Such a system, although it may originally be quite homogeneous, may later develop a pattern or structure due to an instability of the homogeneous equilibrium, which is triggered off by random disturbances. Such reaction-diffusion systems are considered in some detail in the case of an isolated ring of cells, a mathematically convenient, though biologically unusual system. The investigation is chieffy concerned with the onset of instability. It is found that there are six essentially different forms which this may take. In the most interesting form stationary waves appear on the ring. It is suggested that this might account, for instance, for the tentacle patterns on *Hydra* and for whorled leaves. A system of reactions and diffusion on a sphere is also considered. Such a system appears to account for gastrulation. Another reaction system in two dimensions gives rise to patterns reminiscent of dappling. It is also suggested that stationary waves in two dimensions could account for the phenomena of phyllotaxis.

The purpose of this paper is to discuss a possible mechanism by which the genes of a zygote may determine the anatomical structure of the resulting organism. The theory does not make any new hypotheses; it merely suggests that certain well-known physical laws are sufficient to account for many of the facts. The full understanding of the paper requires a good knowledge of mathematics, some biology, and some elementary chemistry. Since readers cannot be expected to be experts in all of these subjects, a number of elementary facts are explained, which can be found in text-books, but whose omission would make the paper difficult reading.

1. A model of the embryo. Morphogens

In this section a mathematical model of the growing embryo will be described. This model will be a simplification and an idealization, and consequently a falsification. It is to be hoped that the features retained for discussion are those of greatest importance in the present state of knowledge.

The model takes two slightly different forms. In one of them the cell theory is recognized but the cells are idealized into geometrical points. In the other the matter of the organism is imagined as continuously distributed. The cells are not, however, completely ignored, for various physical and physico-chemical characteristics of the matter as a whole are assumed to have values appropriate to the cellular matter.

With either of the models one proceeds as with a physical theory and defines an entity called 'the state of the system'. One then describes how that state is to be determined from the state at a moment very shortly before. With either model the description of the state consists of two parts, the mechanical and the chemical. The mechanical part of the state describes the positions, masses, velocities and elastic properties of the cells, and the forces between them. In the continuous form of the theory essentially the same information is given in the form of the stress, velocity, density and elasticity of the matter. The chemical part of the state is given (in the cell form of theory) as the chemical composition of each separate cell; the diffusibility of each substance between each two adjacent cells must also

Vol. 237. B. 641. (Price 8s.) 5 [Published 14 August 1952

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

Related concepts

- p Systems biology
 - n "Biology of networks"
 - 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/



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
- Paradigm shift in biology to utilise bioinformatics in research

Bioinformatician's skill set

p Statistics, data analysis methods

- n Lots of data
- n High noise levels, missing values
- n #attributes >> #data points

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

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!

Communication skills: case 1


Communication skills: case 2



Communication skills: case 2



...biologist/bioinformatician ratio is important!

Communication skills: case 3



Bioinformatician's skill set

p How much biology you should know?

Bioinformatician's skill set



Where would you be in this triangle?

A problem involving bioinformatics?



- "I found a fruit fly that is immune to all diseases!"
- "It was one of these"

Molecular biology primer



Molecular Biology Primer by Angela Brooks, Raymond Brown, Calvin Chen, Mike Daly, Hoa Dinh, Erinn Hama, Robert Hinman, Julio Ng, Michael Sneddon, Hoa Troung, Jerry Wang, Che Fung Yung

Edited for Introduction to Bioinformatics (Autumn 2007, Summer 2008, Autumn 2008) by Esa Pitkänen

Molecular biology primer

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

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

Cells

- p Fundamental working units of every living system.
- Every organism is composed of one of two radically different types of cells:
 - n prokaryotic cells or
 - n eukaryotic cells.
- 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.

Two types of cells: Prokaryotes and Eukaryotes





Prokaryotes and Eukaryotes

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



Lecture: Phylogenetic trees

All Cells have common Cycles



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



James Watson and Francis Crick



"Photo 51"



Rosalind Franklin

DNA, continued



- 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



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

Human chromosomes

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

K	\$	X		5		L
Ķ	and the second s	(ر ۱	(r	10	?)) 12
儿 13	JL.	15){ 16	71 17	18
	19	20		21	11 22	Si

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

Length of DNA and number of chromosomes

Organism	#base pairs	#chromosomes (germline)
Prokayotic Escherichia coli (bacterium)	4x10 ⁶	1
Eukaryotic		
Saccharomyces cerevisia (yeast)	1.35x10 ⁷	17
Drosophila melanogaster (insect)	1.65x10 ⁸	4
Homo sapiens (human)	2.9x10 ⁹	23
Zea mays (corn / maize)	5.0x10 ⁹	10

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	cgaggggccc	ggcggccaag	tttggaggac
241	actccggccc	gaagggttga	gagtacccca	gagggaggaa	gccacacgga	gtagaacaga
301	gaaatcacct	ccagaggacc	ccttcagcga	acagagagcg	catcgcgaga	gggagtagac
361	catagcgata	ggaggggatg	ctaggagttg	ggggagaccg	aagcgaggag	gaaagcaaag
421	agagcagcgg	ggctagcagg	tgggtgttcc	gccccccgag	aggggacgag	tgaggcttat
481	cccggggaac	tcgacttatc	gtccccacat	agcagactcc	cggaccccct	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	ccgaagagga	aagaaggacg	cgagacgcaa	acctgcgagt	ggaaacccgc	tttattcact
961	ggggtcgaca	actctgggga	gaggagggag	ggtcggctgg	gaagagtata	tcctatggga
1021	atccctggct	tccccttatg	tccagtccct	ccccggtccg	agtaaagggg	gactccggga
1081	ctccttgcat	gctggggacg	aagccgcccc	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	atcctcccct	ggaaggcctc	ttcctaggtc	cggagtctac	ttccatctgg	tccgttcggg
1381	ccctcttcgc	cggggggagcc	ccctctccat	ccttatcttt	ctttccgaga	attcctttga
1441	tgtttcccag	ccagggatgt	tcatcctcaa	gtttcttgat	tttcttctta	accttccgga
1501	ggtctctctc	gagttcctct	aacttctttc	ttccgctcac	ccactgctcg	agaacctctt
1561	ctctcccccc	gcggtttttc	cttccttcgg	gccggctcat	cttcgactag	aggcgacggt
1621	cctcagtact	cttactcttt	tctgtaaaga	ggagactgct	ggccctgtcg	cccaagttcg
1681	ag					

RNA

- P RNA is similar to DNA chemically. It is usually only a single strand. T(hyamine) is replaced by U(racil)
- p 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

DNA, RNA, and the Flow of Information



Denis Noble: The principles of Systems Biology illustrated using the virtual heart http://velblod.videolectures.net/2007/pascal/eccs07_dresden/noble_denis/eccs07_noble_psb_01.ppt

Proteins

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



Urease enzyme from Helicobacter pylori

Amino acids

Н	Н	Н	Н	Н
0	0, 1	<u>م ا</u>		<u> </u>
H ₃ N ⁺ - ^a C - C ⊖	H ₃ N⁺ - °C - Cle	H ₃ N⁺ - °C - C (⊖	H ₃ N ⁺ - °C - C (⊖	H ₃ N ⁺ - ^a C - Ci⊖
(CH _a),	CH.	CH.	CH.	L L V
4/5	0.12			
NH	$\dot{C}H_2$			
C-bill		\sim		H
C=NH ₂	U=U 1			
NH,	I NH _a	Phenvlalanine	Tvrosine	Tryptophan
Aroinine	Glutamine	(Phe / F)	(Tyr / Y)	(Trp, W)
(Arg/R)	(Gln / O)			
	· · · ·	Н		
	Н	HNT-C.CA		
H₂N⁺ - ℃ - C(⊖	0, 1			
	H₃N⁺ - °C - C⊖	ĊH ₃	/ CH2	ĊH2
$(CH_2)_4$	0″		нң ји	
	H	Alanina	Histidine	OH
Lveine	(Glv/G)	(Ala/A)	(His / H)	(Ser / S)
(Lys / K)	Н	H	Н	Н
H _a	 1 .0	 	0, 1	
, C	H_3N^+ - ${}^{o}C$ - $C \stackrel{\circ}{\leftarrow} \Theta$	H₃N⁺ -ªC - C€	H₃N⁺ - °C - C (⊖	H₃N⁺ - ºC - C(€
H ₁ C CH ₂	0	I `0	0 1	<u> </u>
	CH_2	CH ₂	H = C = OH	
$H_2N^* \cdot {}^{\alpha}C \cdot C \overleftrightarrow{\Theta}$	і СН	і Соон	СН	। सरु
Proline	1			
(Pro / P)	COOH			
Н	Glutamic Acid	Aspartic Acid	Threonine	Cysteine
0	(Glu / E)	(Asp / D)	(Thr / T)	(Cys / C)
H₃N⁺ - °C - C (⊖	Н	Н	Н	Н
		0, 1	0, 1	
	H₃N⁺ - ℃ - C 🤴	$H_3N^+ - C - C \in$	H₃N⁺ -°C - C(⊖	H ₃ N ⁺ - ^u C - C e
ĊH,				сн с
S	ĊH	Ċ = O	ĊH2	CH ₃ CH ₃
 ~u		I		
UTI3	UH3 UH3	NH_2	CH3	
Methionine (Met / Mix	Leucine	Asparagine	Isoleucine	Valine
(INICE / IVI)	(Leu/L)	(Asn / N)	(IIe / I)	(Val / V)

How DNA/RNA codes for protein?

- p DNA alphabet contains four letters but must specify protein, or polypeptide sequence of 20 letters.
- Dinucleotides are not enough: 4² = 16 possible dinucleotides
- p Trinucleotides (triplets)
 allow 4³ = 64 possible
 trinucleotides
- p Triplets are also called codons



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
- How to find a gene? Look
 for start and stop codons
 (not that easy though)



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 essential work 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.
- Proteins work together with other proteins or nucleic acids as "molecular machines"
 - n structures that fit together and function in highly specific, lockand-key ways.

Lecture 8: Proteomics

Genes

- "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



FoldIt: Protein folding game



Genes & alleles

A gene can have different variants
 The variants of the same gene are called alleles



Genes can be found on both strands







Alternative splicing

Different splice variants may be generated



Where does the variation in genomes come from?

- Prokaryotes are typically haploid: they have a single (circular) chromosome
- DNA is usually inherited vertically (parent to daughter)
- p Inheritance is clonal
 - Descendants are faithful copies of an ancestral DNA
 - Nariation 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
- 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

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
- Meiosis: reduction of chromosome number from 2N to N during reproductive cycle
 - 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

- 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
- Allele combinations in offspring are usually different from combinations found in parents
- P Recombination errors lead into additional variations



Mitosis



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

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
 - Published data is submitted to databases
- p General vs specialised databases
- P This topic is discussed extensively in *Practical course in biodatabases* (III period)



10 most important biodatabases... according to "Bioinformatics for dummies"

- р
- Ensembl р
- PubMed р
- NR р
- UniProt D
- InterPro р
- OMIM р
- Enzymes р
- PDB р
- KEGG D

GenBank/DDJB/EMBL www.ncbi.nlm.nih.gov www.ensembl.org www.ncbi.nlm.nih.gov www.ncbi.nlm.nih.gov www.expasy.org www.ebi.ac.uk www.ncbi.nlm.nih.gov www.expasy.org www.rcsb.org/pdb/ www.genome.ad.jp

Nucleotide sequences Human/mouse genome Literature references Protein sequences Protein sequences Protein domains Genetic diseases Enzymes Protein structures Metabolic pathways

FASTA format

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

Header line, begins with >

>Hepatitis delta virus, complete genome