

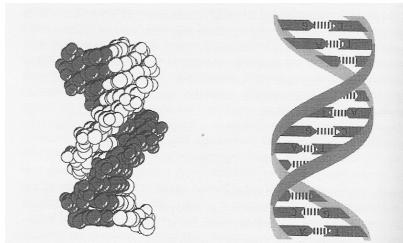
Computational Biology

part II technological developments in molecular biology & databases

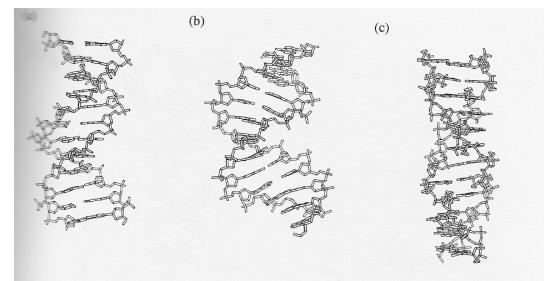
Jaap Kaandorp

	Technology development	Structure determination
1950	49 Edman degradation 54 Isomorphous replacement	51 α -helix model 53 DNA double helix model Insulin primary structure
1960	62 Restriction enzyme	60 Myoglobin tertiary structure 65 tRNA ^{Ala} primary structure
1970	72 DNA cloning 73 DNA sequencing	73 tRNA ^{Asp} tertiary structure 77 4X174 complete genome
1980	84 Pulse field gel electrophoresis 85 Polymerase chain reaction 87 YAC vector	79 Z DNA by single crystal diffraction 86 Protein structure by 2D NMR 88 Human Genome Project
1990	93 DNA chip	95 <i>H. influenzae</i> complete genome
2000		

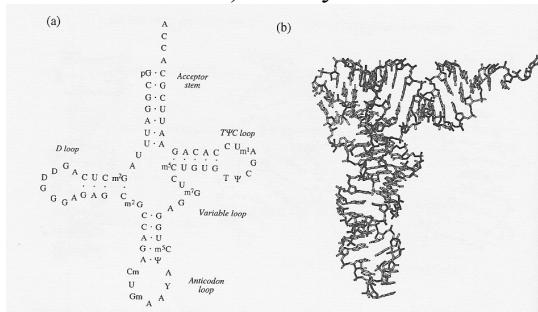
DNA double helix



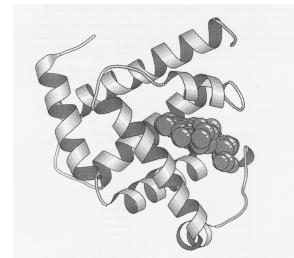
Polymorphic DNA tertiary structures



Transfer RNA a) primary and secondary structure b) tertiary structure



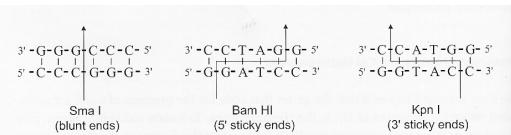
3D protein structure protein: example myoglobin



History of structure determination for nucleic acids and proteins II: restriction enzymes

	Technology development	Structure determination
1950	49 Edman degradation 54 Isomorphous replacement	51 α -helix model 53 DNA double helix model Insulin primary structure
1960	62 Restriction enzyme	60 Myoglobin tertiary structure 65 tRNA ^{Ala} primary structure
1970	72 DNA cloning 75 DNA sequencing	73 tRNA ^{Pro} tertiary structure 77 ϕ X174 complete genome 79 Z-DNA by single crystal diffraction
1980	84 Pulse field gel electrophoresis 85 Polymerase chain reaction 87 YAC vector	86 Protein structure by 2D NMR 88 Human Genome Project
1990	93 DNA chip	95 <i>H. influenzae</i> complete genome
2000		

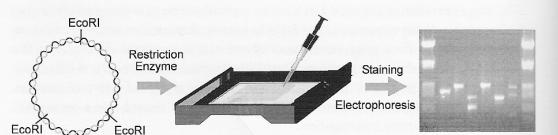
Restriction enzymes



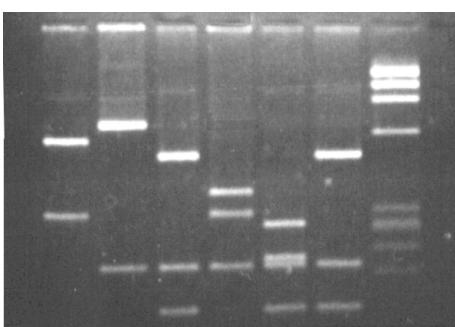
Restriction enzymes, different recognition sequences

Name	Organism	Recognition sequence
Alu I	<i>Anthrobacter luteus</i>	5' A G ↑ C T 3'
Bam H1	<i>Bacillus amyloliquefaciens</i>	5' G ↑ G A T C C 3'
Bgl II	<i>Bacillus globigii</i>	5' A ↑ G AT C T 3'
Eco RI	<i>Escherichia coli</i>	5' G T A A T T C 3'
Eco RV	<i>Escherichia coli</i>	5' G A T ↑ A T C 3'
Hind III	<i>Haemophilus influenzae</i>	5' A T A G C T T 3'
Kpn I	<i>Klebsiella pneumonia</i>	5' G G T A C ↑ C 3'
Pme I	<i>Pseudomonas mendocina</i>	5' G T T T ↑ A A A C
Sau96 I	<i>Staphylococcus aureus</i>	5' G ↑ G C C 3'
Sma I	<i>Serratia marcescens</i>	5' C C C ↑ G G G 3'
Xba I	<i>Xanthomonas malvacearum</i>	5' C ↑ C C G G G 3'

Gel electrophoresis of DNA restriction fragments



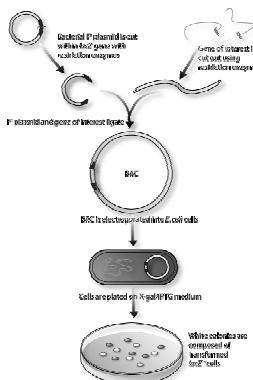
Gel electrophoresis of DNA restriction fragments



History of structure determination for nucleic acids and proteins III: DNA cloning

	Technology development	Structure determination
1950	49 Edman degradation 54 Isomorphous replacement	51 α -helix model 53 DNA double helix model Insulin primary structure
1960	62 Restriction enzyme	60 Myoglobin tertiary structure 65 tRNA ^{Ala} primary structure
1970	72 DNA cloning 75 DNA sequencing	73 tRNA ^{Pro} tertiary structure 77 ϕ X174 complete genome 79 Z-DNA by single crystal diffraction
1980	84 Pulse field gel electrophoresis 85 Polymerase chain reaction 87 YAC vector	86 Protein structure by 2D NMR 88 Human Genome Project
1990	93 DNA chip	95 <i>H. influenzae</i> complete genome
2000		

Cloning using bacteria (Shizuya et al, PNAS, 1992)



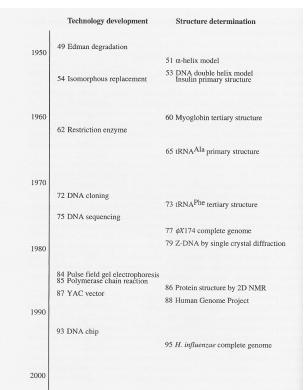
DNA libraries

- Genomic DNA libraries, large number of fragments is cloned in bacteria. Used in for example DNA chips
- Complementary DNA (cDNA) libraries, use mRNA pool of cell or tissue of interest as starting point.. Used to study proteins which are expressed in small amounts

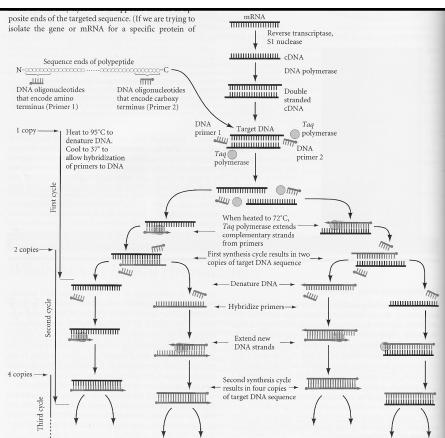
cDNA libraries differ from genomic libraries:

- Complementary DNA
- Only contain coding regions
- Tissue specific, snapshot of current gene expression
- Frequency of clones in library indicates the level of gene expression

History of structure determination for nucleic acids and proteins IV:PCR



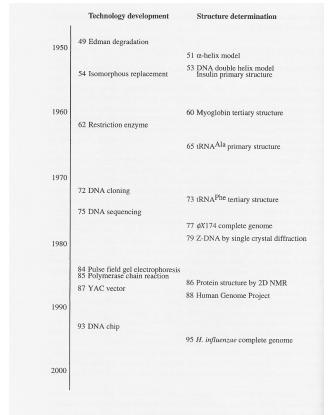
Polymerase Chain Reaction (Gilbert, 2001)



PCR

- Has revolutionized molecular genetics
- You can clone genes directly from genomic DNA (without DNA libraries)
- You can clone mRNA with this technique
- Very sensitive technique (method of choice in forensic studies)
- Real time PCR can provide quantitative information

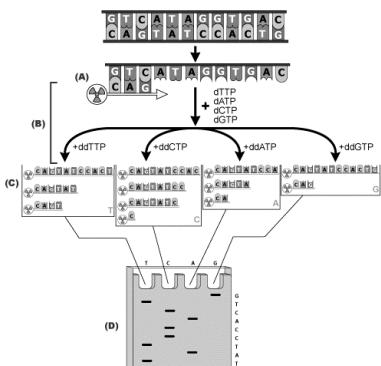
History of structure determination for nucleic acids and proteins V: DNA sequencing



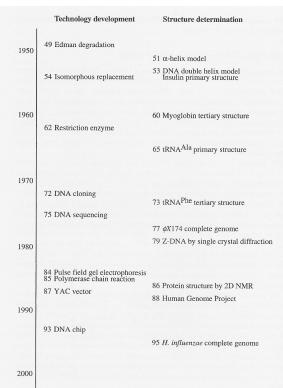
DNA-sequencing

- Use same procedure like in PCR to produce copies of target sequences
- Instead of synthesizing copies of the complete sequence, copies are forced to terminate before reaching the end
- Chain termination method by adding dideoxynucleotides (ddNTPs), chemical modified nucleotid nothing can be added after adding a ddNTP
- Resulting copies are separated by gel electrophoresis
- Four different reactions are visualized on 4 different lanes on the gel

DNA-sequencing



History of structure determination for nucleic acids and proteins V: DNA chips (microarrays)



DNA chips (microarrays) I

- Method for high-throughput gene expression analysis
- Monitor the expression of several thousand of genes in one experiment, not a single gene
- You need a DNA library

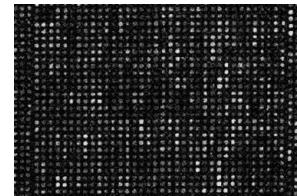
DNA chips (microarrays) II

- Construction of a chip from a DNA library
- Amplify individual genes by PCR
- individual genes are spotted onto glass slides

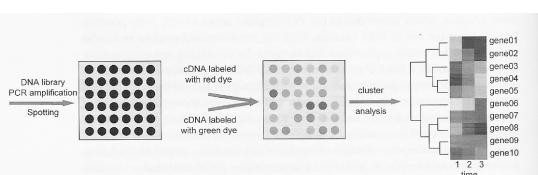
DNA chips (microarrays) III

- Extract mRNA from 2 samples you want to compare
- Use reverse transcriptase mRNA->cDNA
- Label cDNA with fluorescent colours
- cDNA is hybridized with DNA chip

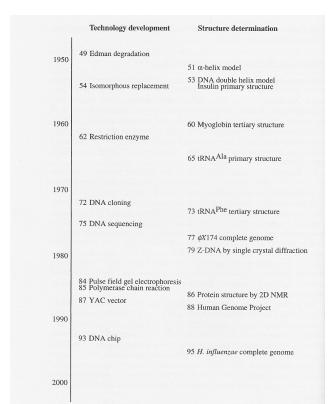
Microarrays III



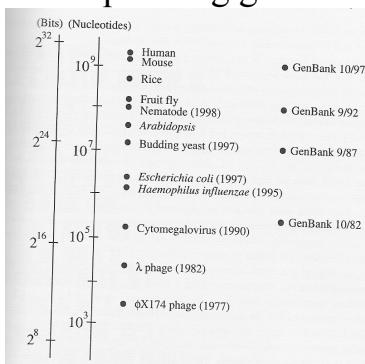
Microarrays IV



History of structure determination for nucleic acids and proteins V: sequenced genomes (human genome 2001)



Sequencing genomes



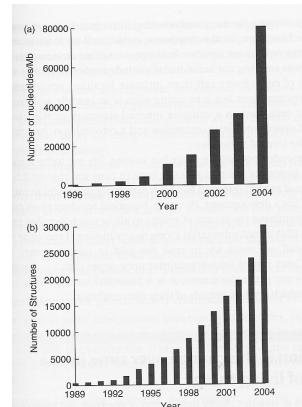
Genome sequences

- Viruses genome size 10^3 - 10^5 nucleotids
- bacteria 10^6 nucleotids
- Saccharomyces cerevisiae (yeast) 10^7 nucleotids
- Caenorhabditis elegans 10^8 nucleotids
- free-living organisms 10^6 - 10^9 nucleotids

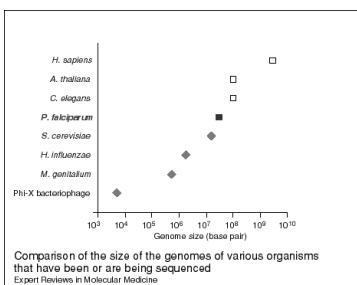
Genome sequences: increasing genome size
-> increasing fraction of non-coding regions

- Bacteria and archeabacteria genes every 1000 nucleotids
- *Saccharomyces cerevisiae* (yeast) genes every 2000 nucleotids
- *Caenorhabditis elegans* genes every 5000 nucleotids

Number of nucleotids / proteins in databases



Genomes different organisms I



Arabidopsis thaliana



Genomes in different species II

TABLE I		
Nuclear genome size in different species		
Common name	Scientific name	Nuclear genome size (1)
Wheat	<i>Triticum aestivum</i>	15,966
Onion	<i>Allium cepa</i>	15,290
Garden pea	<i>Pisum sativum</i>	3,947
Corn	<i>Zea mays</i>	2,292
Asparagus	<i>Asparagus officinalis</i>	1,308
Tomato	<i>Lycopersicum esculentum</i>	907
Sugarbeet	<i>Beta vulgaris</i>	758
Apple	<i>Malus X domestica</i>	743
Common bean	<i>Phaseolus vulgaris</i>	637
Cantaloupe	<i>Cucumis melo</i>	454
Grape	<i>Vitis vinifera</i>	483
Man	<i>Homo sapiens</i>	2,910

1: Expressed in Megabases (1 Mb: 1,000,000 bases)

Genomes different organisms III

Organism	Number of chromosomes (haploid genome)	Genome size (base pairs; genes)
<i>Mycoplasma genitalium</i> (prokaryote)	1 circular chromosome	$580 \cdot 10^3$ bp; 480 genes
<i>Escherichia coli</i> (prokaryote)	1 circular chromosome	$4.6 \cdot 10^6$ bp; 4,290 genes
<i>Saccharomyces cerevisiae</i> (budding yeast; eukaryote)	16 chromosomes	$12.5 \cdot 10^6$ bp; 6,186 genes
<i>Arabidopsis thaliana</i> (flowering plant; eukaryote)	5 chromosomes	$100 \cdot 10^6$ bp; ~25,000 genes
<i>Drosophila melanogaster</i> (fruit fly, eukaryote)	4 chromosomes	$180 \cdot 10^6$ bp; ~14,000 genes
<i>Mus musculus</i> (mouse, eukaryote)	20 chromosomes	$2.5 \cdot 10^9$ bp; ~30,000 genes
<i>Homo sapiens</i> (human, eukaryote)	23 chromosomes	$2.9 \cdot 10^9$ bp; ~30,000 genes

Rice genome

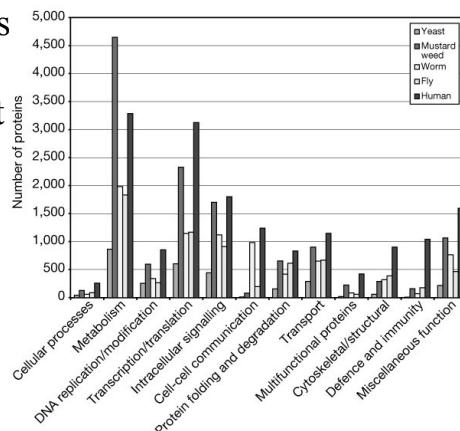
- Goff, S. A. et al. A draft sequence of the rice genome (*Oryza sativa* L. ssp *japonica*). *Science* 296, 92 - 100 (2002).

That means rice probably has more genes than we do, despite having only one-seventh as much DNA. One reason for this difference is that plants' genes seem to be on average much shorter than mammals'. A rice gene is usually about 4,500 DNA letters long. The average human gene probably stretches to over 30,000 letters.

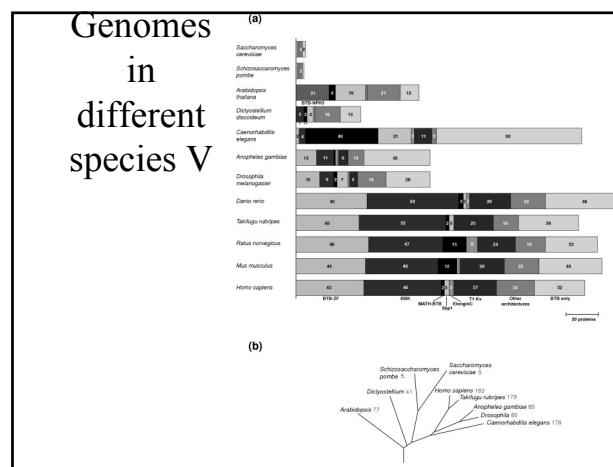
Rice genome II

- Plants are also prone to copying genes, chromosomes, and sometimes their entire genome. This might be a way for them to generate genetic diversity. Mammals can rearrange individual genes to make several different sorts of protein. Plants show much less of this, seeming to rely on a larger number of shorter genes.

Genomes in different species IV



Genomes in different species V



Sequenced genomes

- The genomes of more than 180 organisms have been sequenced since 1995. The Quick Guide includes descriptions of these organisms and has links to sequencing centers and scientific abstracts

http://www.genomenewsnetwork.org/resources/sequenced_genomes/genome_guide_p1.shtml

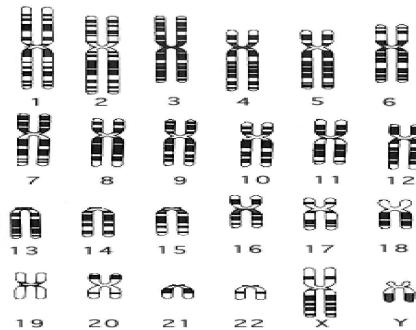
Sequence genomes II

Aeropyrum pernix
Agrobacterium tumefaciens
Anabaena
Anopheles gambiae
Arabidopsis thaliana
Arabidopsis lyrata
Arabidopsis suecica
Arabidopsis thaliana
Archaeoglobus fulgidus
Ashbya gossypii
Bacillus subtilis
Bacillus cereus
Bacillus halodurans
Bacillus licheniformis
Bacillus megaterium
Bacillus subtilis
Bacillus thuringiensis
Bacteroides fragilis
Bacteroides thetaiotaomicron
Bartonella henselae
Bartonella quintana
Bartonella bacilliformis
Bifidobacterium longum
Blochmannia floridanus
Bordetella bronchiseptica
Bordetella parapertussis
Bordetella pertussis
Borrelia burgdorferi
Bradyrhizobium japonicum
Brucella melitensis
Brucella suis
Buchnera aphidicola

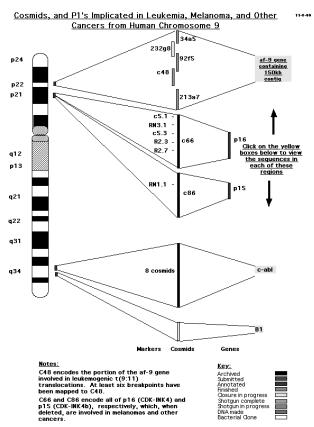
Human Genome project

- Since initiation Human Genome project (1988, finalized in 2001) computational biology (bioinformatics) has become an integral part of molecular biology
- Make biological sense out of sequence data, data produced by DNA chips etc.

Human chromosomes



Human chromosomes



Genetic disorders

- X-linked recessive disorders
- Autosomal dominant disorders
- Autosomal recessive disorders

Computational Biology
end part II technological developments in
molecular biology & databases
Jaap Kaandorp