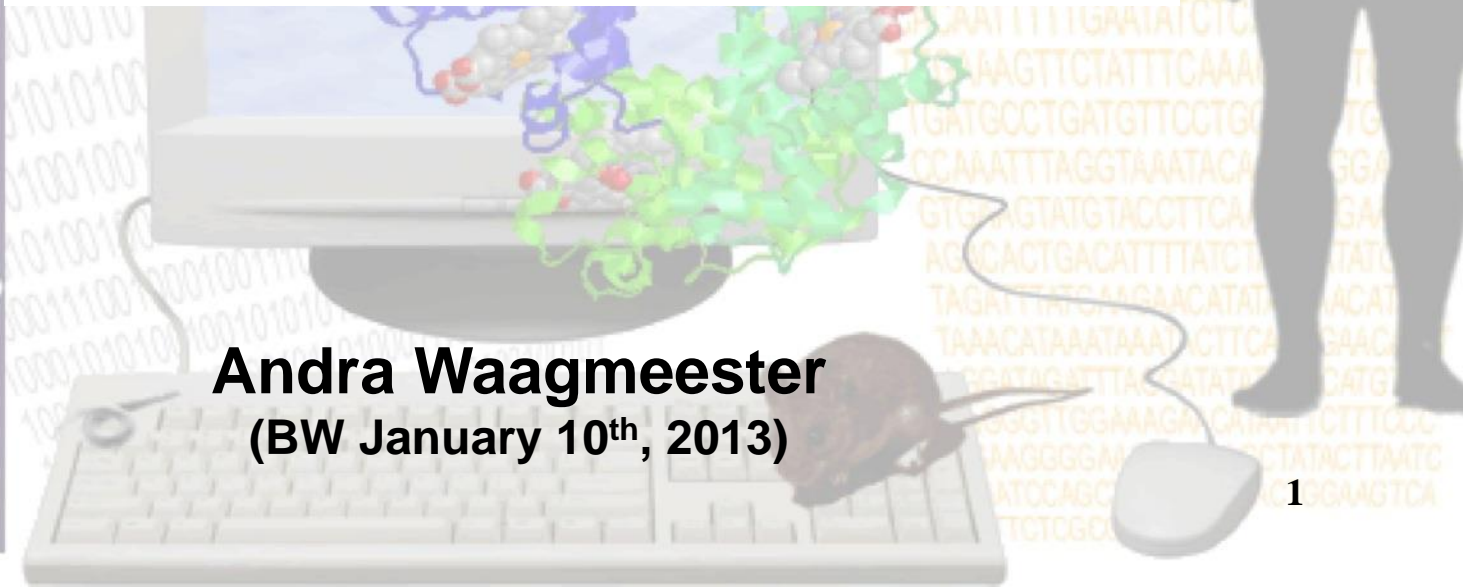# INTRODUCTION BIOINFORMATICS & BIOLOGICAL DATABASES

**Andra Waagmeester**
**(BW January 10th, 2013)**

**1**

# BLOCKCOORDINATOR

- **Dr. Ir. Chris Evelo**        chris.evelo@maastrichtuniversity.nl

# ORGANISATION

- **Dr. Susan Coort (BGK)**    susan.coort@maastrichtuniversity.nl
- **Dr. Lars Eijssen (MLW)**    l.eijssen@maastrichtuniversity.nl
- **Andra Waagmeester (BW)**andra.waagmeester@maastrichtuniversity

# OTHER INSTRUCTORS (left to right)

- **Dr. Egon Willighagen**            egon.willighagen@maastrichtuniversity.nl
- **Stan Gaj, MSc**                s.gaj@maastrichtuniversity.nl
- **Martina Kutmon, MSc**            martina.kutmon@maastrichtuniversity.nl
- **Anwesha Dutta, MSc**            anwesha.dutta@maastrichtuniversity.nl

# Practical information

- Lectures will mostly be in **Dutch**
  Slides and practical exercises are in **English**

- All the practical sessions, should be signed off.

- You are required to study the literature provided for each practical session, <u>before</u> the session starts.
  - ➢ For the first session, the literature can be studied in the session

- The Bioinformatics trajectory (BMW2003) in year 2 is present in periods 1, 3, 4 and 5. Per block a bioinformatics exam will be given. This exam is separate from the block exam. In the end you will get **one grade** for the Bioinformatics trajectory

- **First exam Bioninformatics trajectory**
  - ➢ Tuesday  April 4rd , 13.00-16:00
    <u>Open</u> book exam.

# Course Material

**ELEUM:**

- The *slides* of the lecture will become available after the lecture.
- Per practical session *literature references* are provided.
- The *exercises* are available before the start of each practical session.
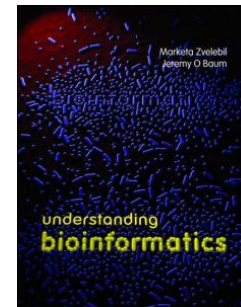- The *answers* to the questions are available a week after each practical session.

# Course Material

**ELEUM:**

- The *slides* of the lecture will become available after the lecture.
- Per practical session *literature references* are provided. We strongly recommend you to read these before the session.
- The *exercises* are available before the start of each practical session.
- The *answers* to the questions are available a week after each practical session.

**BOOKS** (available at the "studielandschap")

- **Understanding bioinformatics**
  *M Zvelebil and J.O. Baum*
- **Bioinformatics:**
  **Sequence and genome analysis**
  *David W. Mount*
- **Bioinformatics and Functional Genomics**
  *J. Pevsner*
- **Learning Perl**
  *Randal L. Schwartz and Tom Phoenix*

# Subjects of bioinformatics track in BW2.3

1. Introduction Bioinformatics and Biological Databases

2. Protein Structures

Next block 2.4 the programme will continue:

Gene Expression data

Pathway analysis – Network analysis

Introduction Programming:

Matlab

2 sessions in

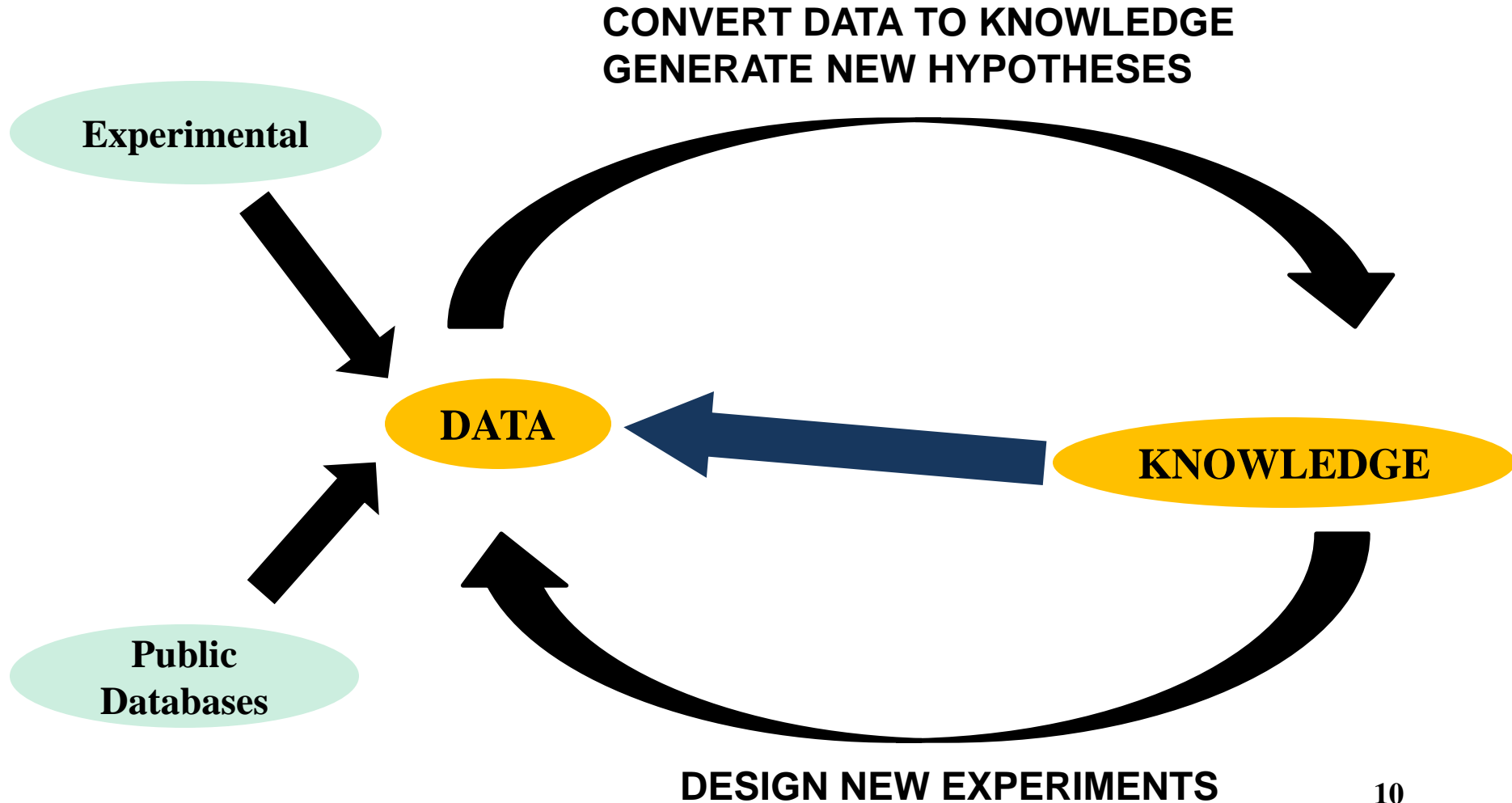# Introduction to Bioinformatics

# Pathways
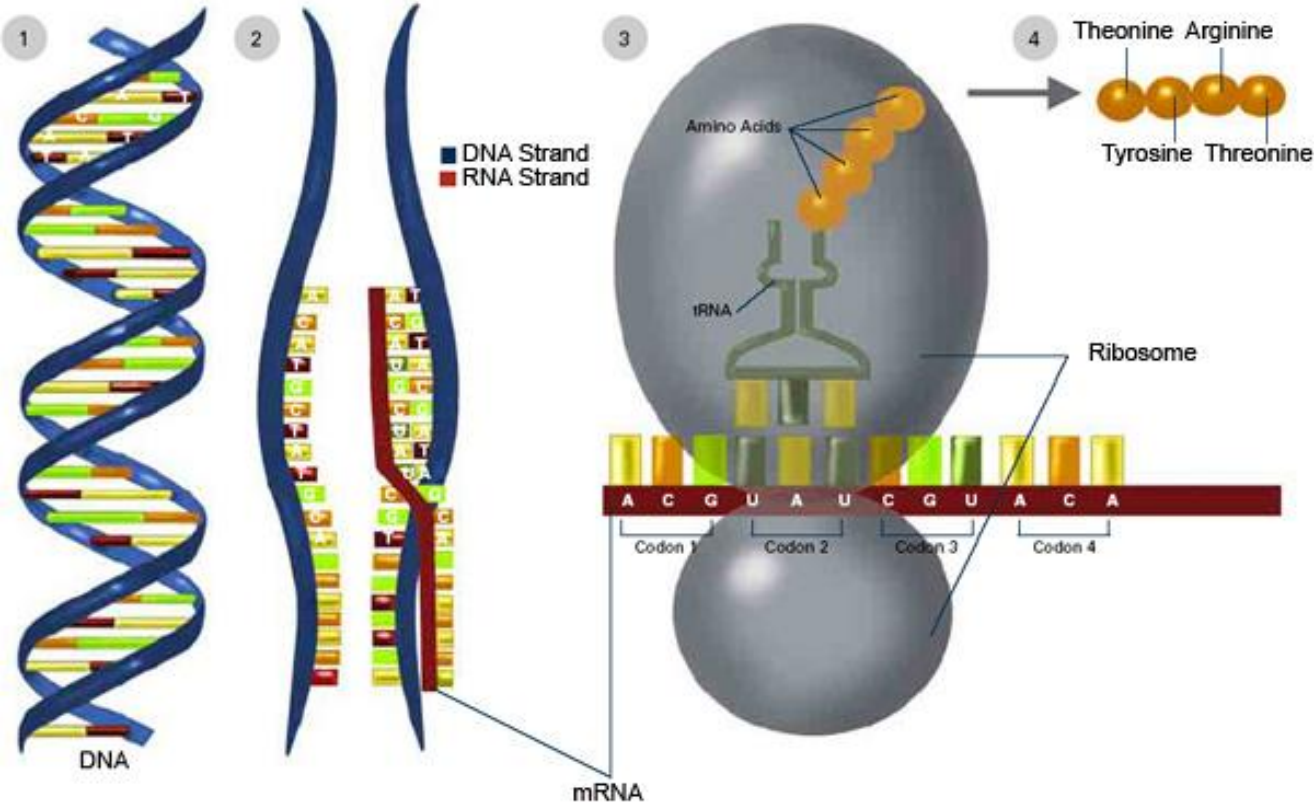
# What is Bioinformatics?



Bioinformatics uses **"informatics" techniques**
(from applied math, computer science, statistics, etc.)
to **understand** and **organize** biological information,
like genes, proteins and molecules on a **large-scale.** 9
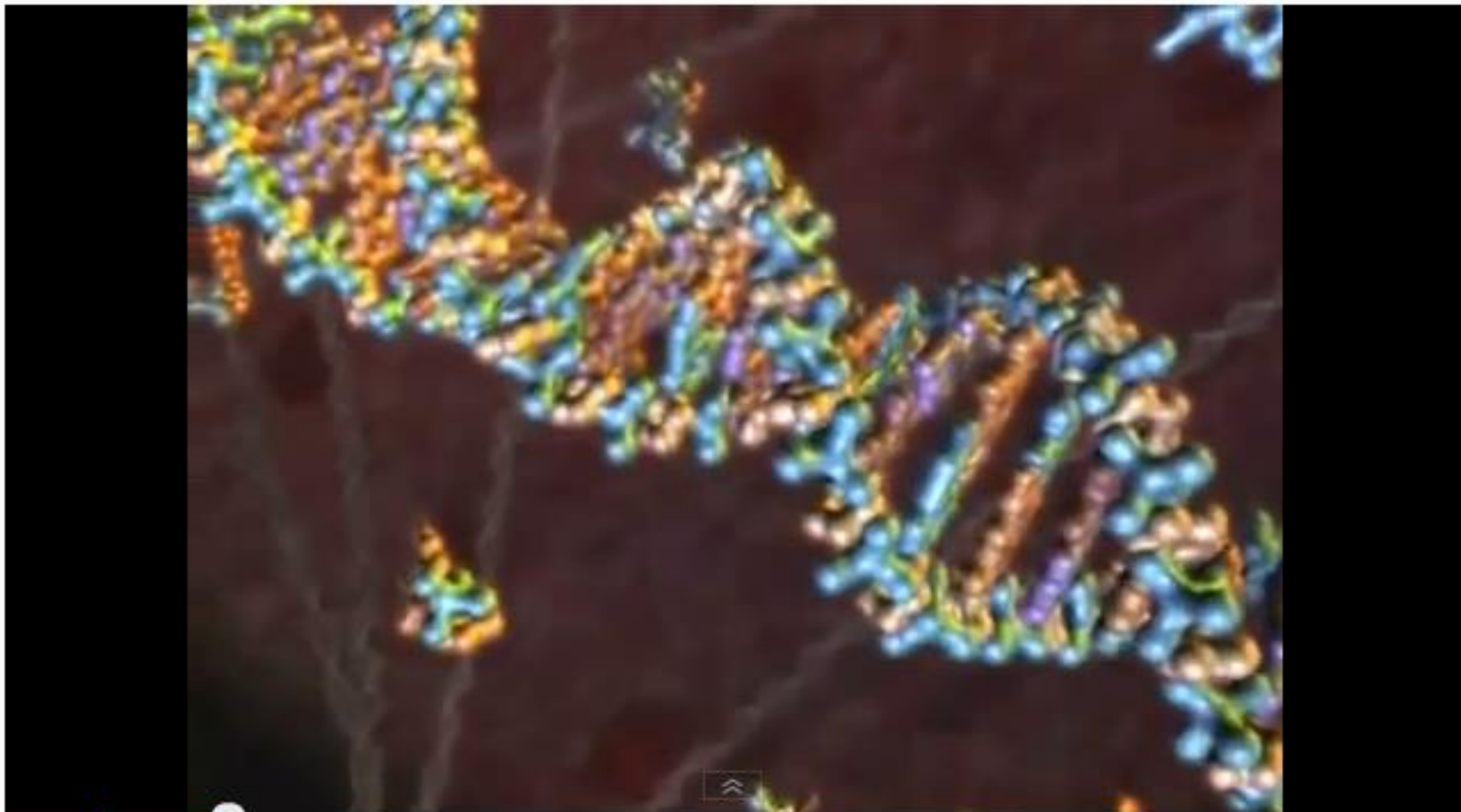
# Why Bioinformatics?



**CONVERT DATA TO KNOWLEDGE**
**GENERATE NEW HYPOTHESES**

Experimental

Public
Databases

DATA

KNOWLEDGE

**DESIGN NEW EXPERIMENTS**

# Central dogma of Molecular Biology



| Gene (DNA) | Transcription ⟹ | mRNA | Translation ⟹ | Protein |

Cells express **different** subset of the genes in different tissues and under different conditions
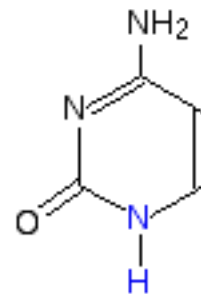
0:26 / 3:02

# DNA Transcription and Protein Assembly

redandbrownpaperbag · 8 video's

Abonneren    187

# DNA



Base pairs

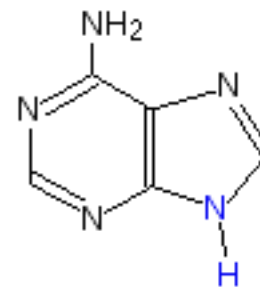Adenine    Thymine

Guanine    Cytosine

Sugar phosphate backbone

cytosine (C)

thymine (T)

adenine (A)

guanine (G)

**13**

# mRNA



Image adapted from: National Human Genome Research Institute.

# Proteins
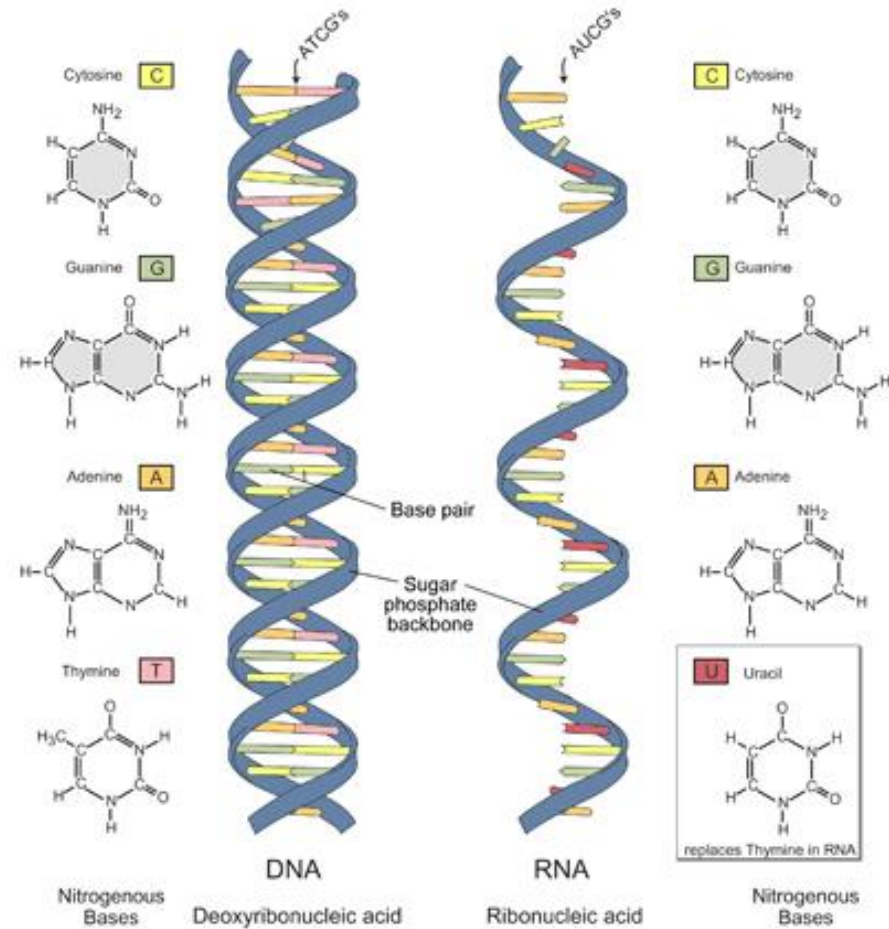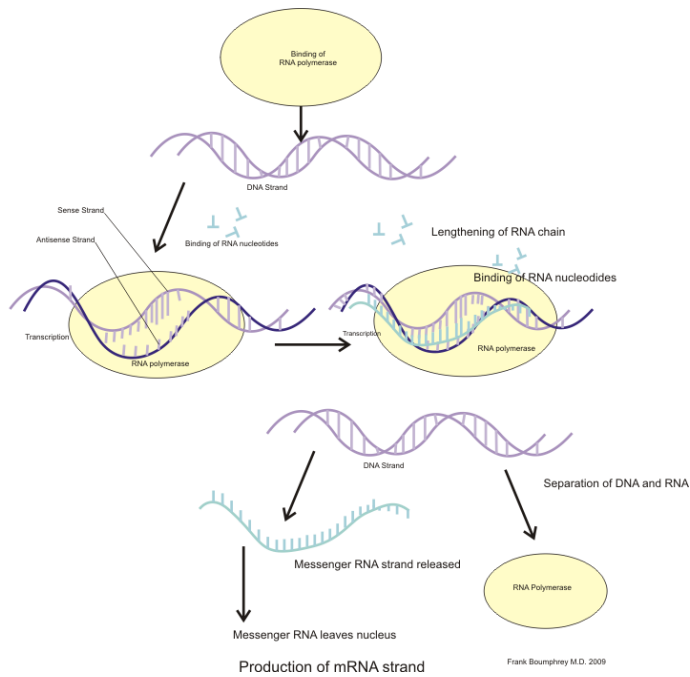
# Genome Sequences /
# the Human Genome Project

AGTCCGCGAATACAGGCTCGGT

# Genomes

A genome is the collection of DNA that comprises an organism.

Today we have assembled the sequence of hundreds of genomes.



Mary S. Gibbs (GNN)

The genome is divided into chromosomes, chromosomes contain genes, and genes are made of DNA.

Each one of earth's organism has its own distinctive genome (except identical twins).

# Genome content
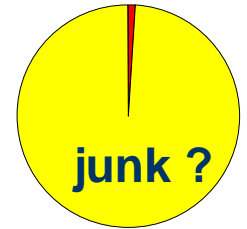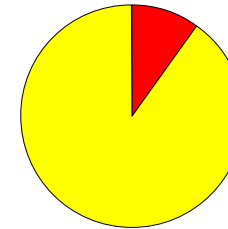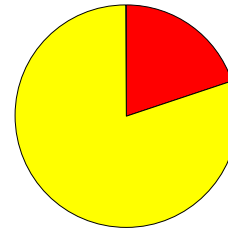
|  | bacteria | yeast | worm | fly | man |
|---|---|---|---|---|---|
| **Size (Mb)** | 2 | 12 | 97 | 137 | 3.500 |

**% genes**

junk ?

| **total genes** | 2.000 | 6.300 | 19.000 | 14.000 | 30.000 ? |

# Functional genomics

**Single genes**

**All genes**

| | | |
|---|---|---|
| **DNA** | *Organisation (sequencing)* | **GENOME** |
| ⇩ | ⬇ | ⇩ |
| **RNA** | *Expression (µ-arrays/sequencing)* | **TRANSCRIPTOME** |
| ⇩ | ⬇ | ⇩ |
| **PROTEIN** | *Synthesis/Structure (2D gels & LCMS/NMR-Xray)* | **PROTEOME** |
| ⇩ | ⬇ | ⇩ |
| **METABOLISM** | *Flux (NMR-kinetics-model)* | **METABOLOME** |
| | ⬇ ⬇ ⬇ | |
| | *FUNCTION* | |

19

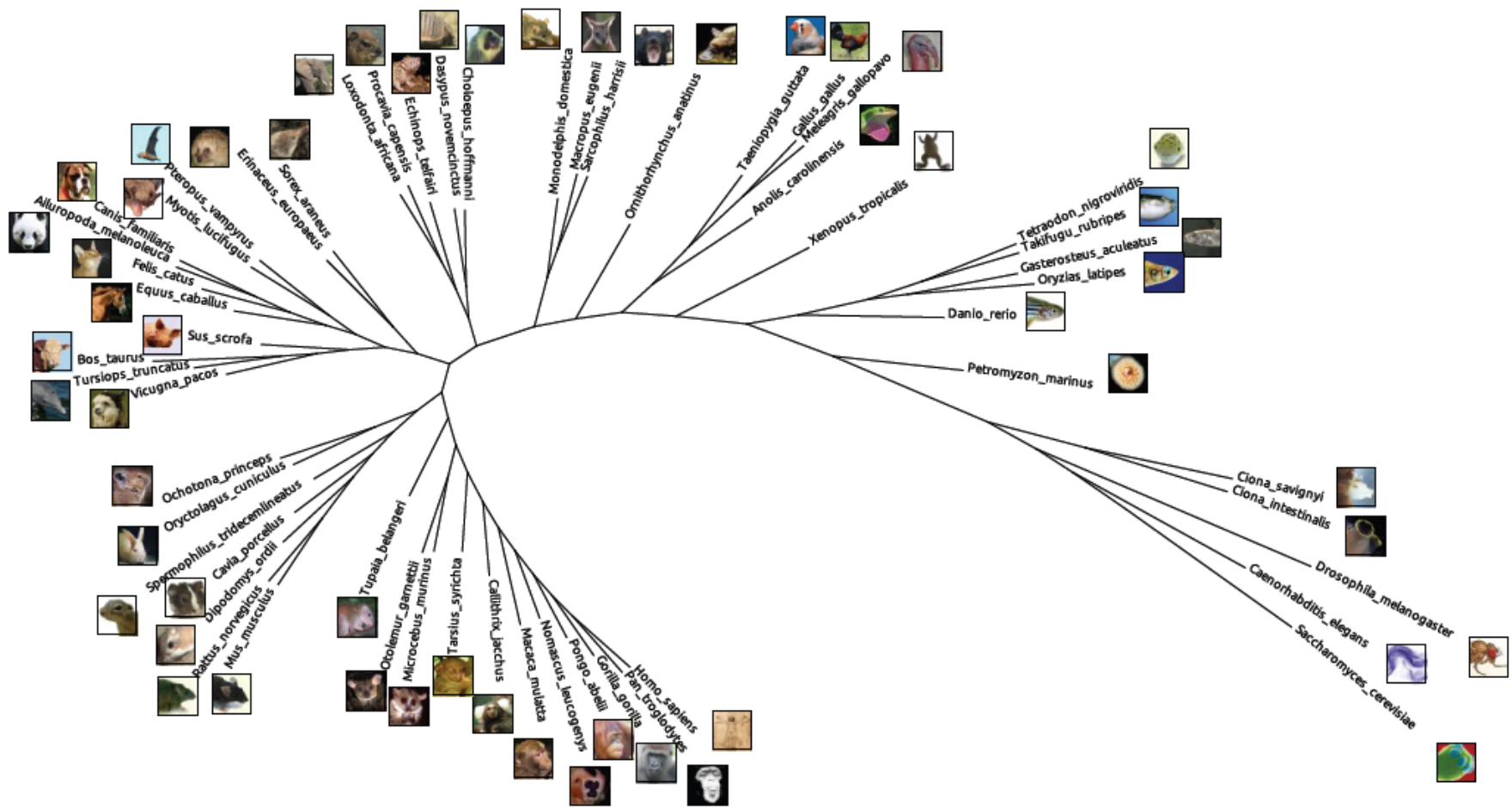# The vertebrate genomes available in Ensembl



Image obtained using Dendroscope (D.H. Huson et al. "Dendroscope- An interactive viewer for large phylogenetic trees", BMC Bioinformatics 8:460, 2007)

# Human Genome project

- ## [Introduction video](http://www.youtube.com/watch?v=N4i6lYfYQzY&list=PLF0701633C91835BF&index=1)
  ( http://www.youtube.com/watch?v=N4i6lYfYQzY&list=PLF0701633C91835BF&index=1)

- # Strategies

- # Conclusions

  **International Human Genome Sequencing Consortium, Finishing the euchromatic sequence of the human genome. Nature 431, 931-945 (21 October 2004).**
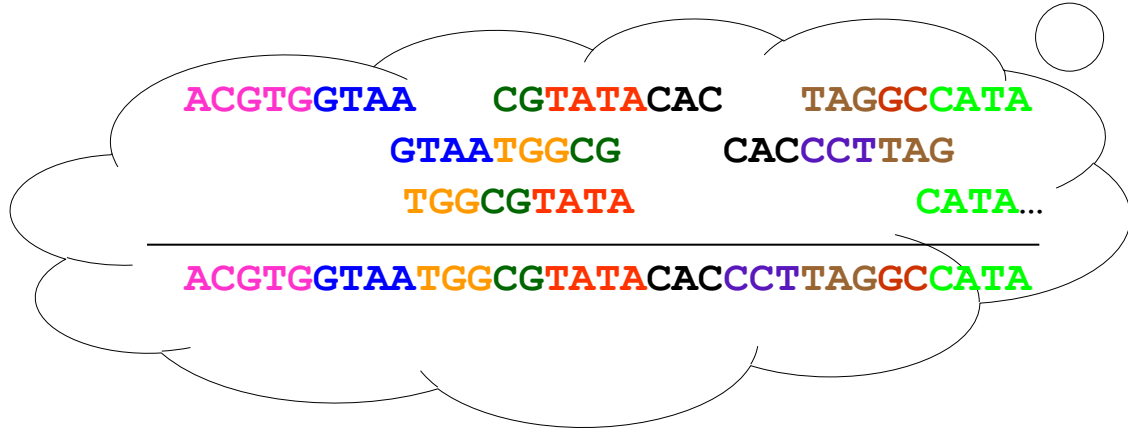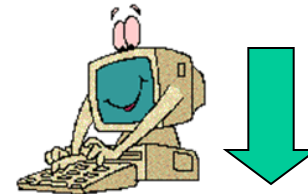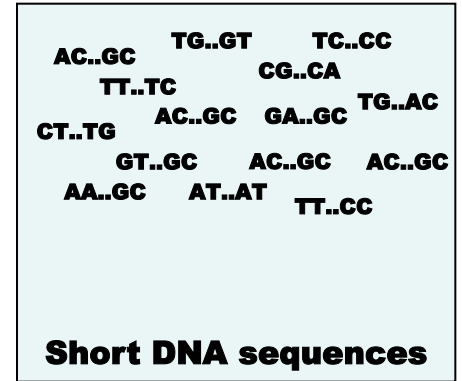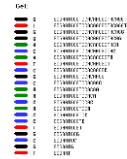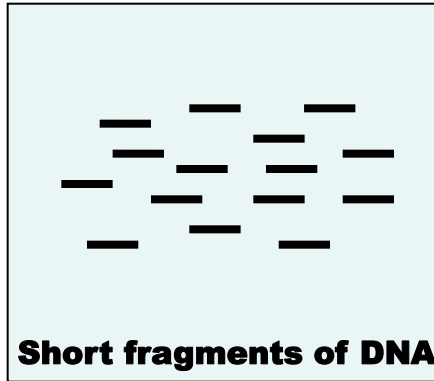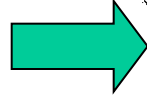
# Overview of genome analysis

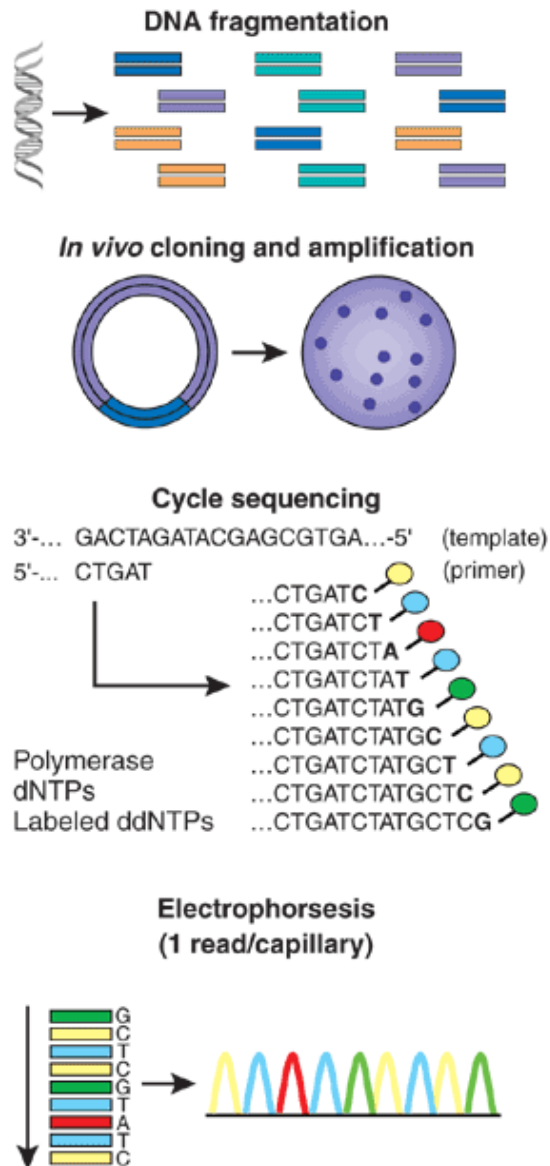There are two main strategies for sequencing genomes

An approach used to decode an organism's genome by shredding it into smaller fragments of DNA which can be sequenced individually. The sequences of these fragments are then ordered, based on overlaps in the genetic code, and finally reassembled into the complete sequence.

The 'whole genome shotgun' (WGS) method is applied to the entire genome all at once, while the 'hierarchical shotgun' method is applied to large, overlapping DNA fragments of known location in the genome.

# Genome sequencing

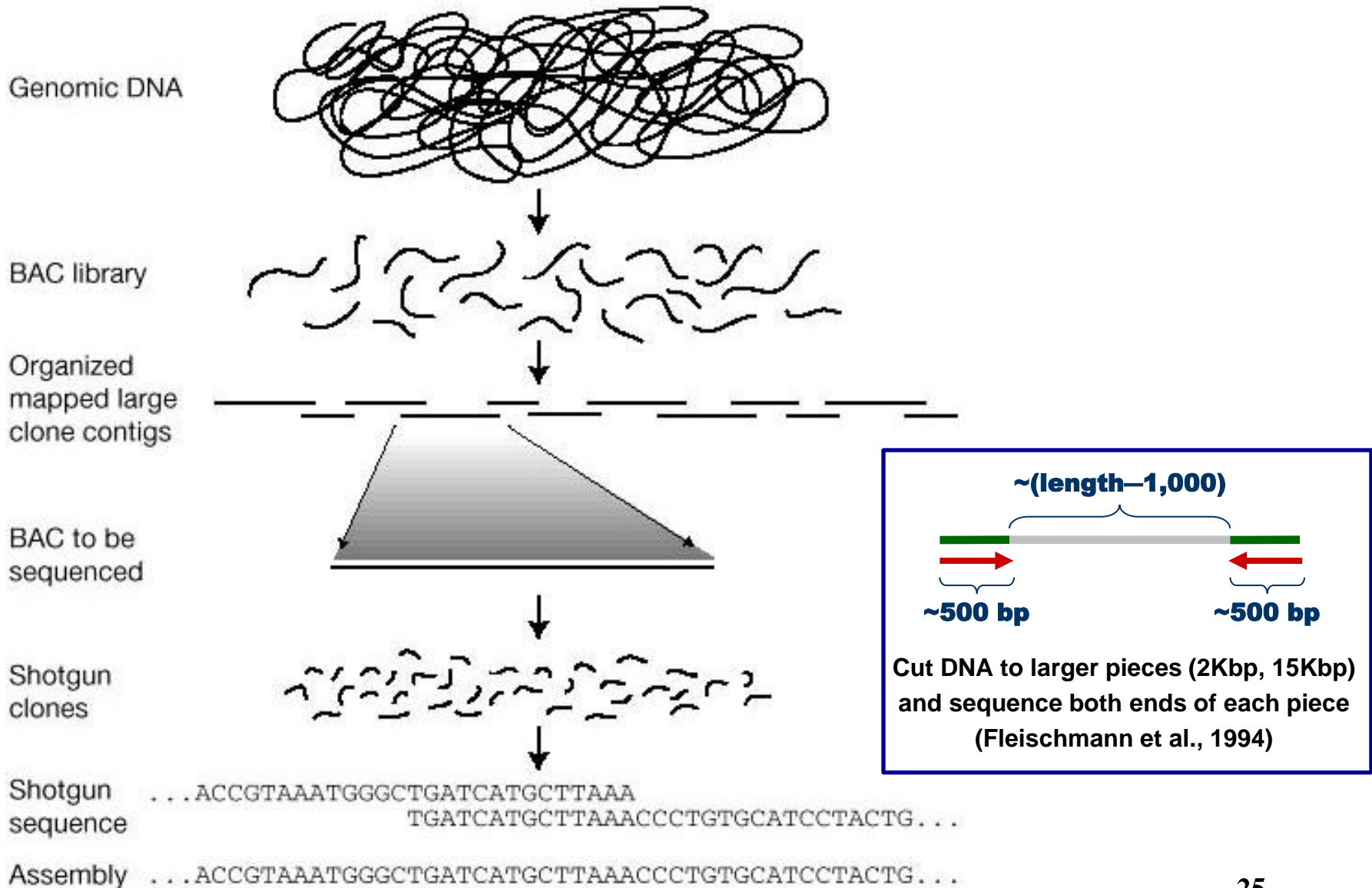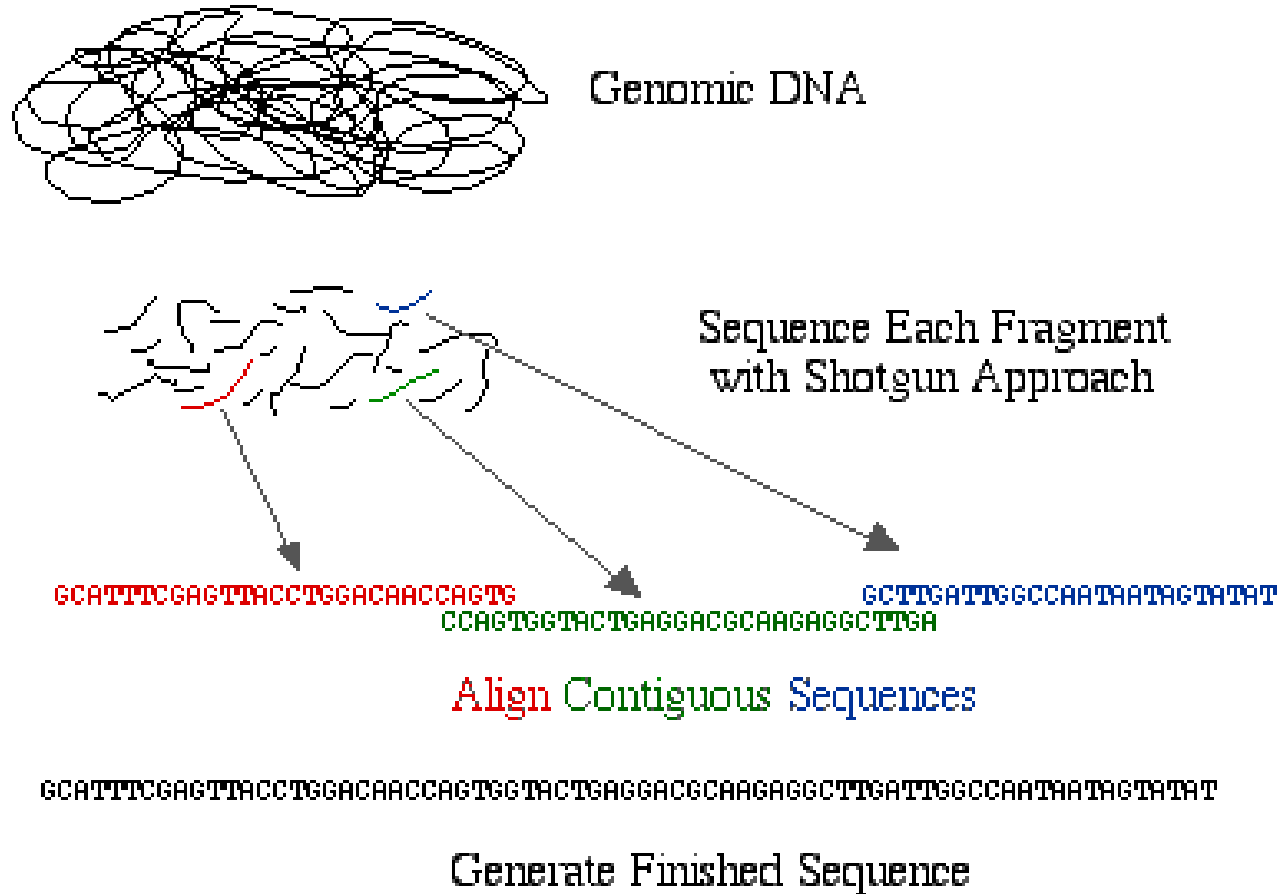**Genome**

**Short fragments of DNA**

**Short DNA sequences**

AC..GC  TG..GT  TC..CC
TT..TC  CG..CA
CT..TG  AC..GC  GA..GC  TG..AC
GT..GC  AC..GC  AC..GC
AA..GC  AT..AT  TT..CC

ACGTGGTAA  CGTATACAC  TAGGCCATA

GTAATGGCG  CACCCTTAG

TGGCGTATA  CATA...

ACGTGGTAATGGCGTATACACCCTTAGGCCATA

ACGTGACCGGTACTGGTAACGTACA
CCTACGTGACCGGTACTGGTAACGT
ACGCCTACGTGACCGGTACTGGTAA
CGTATACACGTGACCGGTACTGGTA
ACGTACACCTACGTGACCGGTACTG
GTAACGTACGCCTACGTGACCGGTA
CTGGTAACGTATACCTCT...

**Sequenced genome**

# Workflow of Sanger sequencing



**DNA fragmentation**

*In vivo* cloning and amplification

**Cycle sequencing**

3'-... GACTAGATACGAGCGTGA...-5' (template)
5'-... CTGAT (primer)

...CTGATC
...CTGATCT
...CTGATCTA
...CTGATCTAT
...CTGATCTATG
...CTGATCTATGC
...CTGATCTATGCT
...CTGATCTATGCTC
...CTGATCTATGCTCG

Polymerase
dNTPs
Labeled ddNTPs

**Electrophorsesis**
(1 read/capillary)

# Hierarchical shotgun sequencing



Genomic DNA

BAC library

Organized mapped large clone contigs

BAC to be sequenced

Shotgun clones

Shotgun sequence ...ACCGTAAATGGGCTGATCATGCTTAAA
TGATCATGCTTAAACCCTGTGCATCCTACTG...

Assembly ...ACCGTAAATGGGCTGATCATGCTTAAACCCTGTGCATCCTACTG...

~(length—1,000)

~500 bp          ~500 bp

**Cut DNA to larger pieces (2Kbp, 15Kbp)
and sequence both ends of each piece
(Fleischmann et al., 1994)**

**25**

**Source: IHGSC (2001)**

# Whole genome shotgun sequencing



Genomic DNA

Sequence Each Fragment
with Shotgun Approach

GCATTTCGAGTTACCTGGACAACCAGTG
CCAGTGGTACTGAGGACGCAAGAGGCTTGA
GCTTGATTGGCCAATAATAGTATAT

Align Contiguous Sequences

GCATTTCGAGTTACCTGGACAACCAGTGGTACTGAGGACGCAAGAGGCTTGATTGGCCAATAATAGTATAT

Generate Finished Sequence

26

**Source: IHGSC (2001)**

# When has a genome been fully sequenced?

A typical goal is to obtain five to ten-fold coverage.

Finished sequence: a clone insert is contiguously sequenced with high quality standard of error rate 0.01%. There are usually no gaps in the sequence.

Draft sequence: clone sequences may contain several regions separated by gaps. The true order and orientation of the pieces may not be known.

# Main conclusions of human genome project (1)

1. We have about the same number of genes as fish and plants, and not that many more genes than worms and flies – 20,000-30,000

2. The human proteome is far more complex than the set of proteins encoded by invertebrate genomes – alternative splicing

3. Hundreds of human genes were acquired from bacteria by lateral gene transfer

4. 98% of the genome does not code for genes and >50% of the genome consists of repetitive DNA

5. Segmental duplication is a frequent occurrence in the human genome

# Main conclusions of human genome project (2)

6. There are 300,000 Alu repeats in the human genome
    - These are about 300 base pairs and contain an AluI restriction enzyme site.
    - They occupy 3% of the genome and may confer some benefit

7. The mutation rate is about twice as high in male meiosis than female meiosis; most mutation probably occurs in males

8. 1.5 – 2 million single base pair changes or single nucleotide polymorphisms (SNPs) were originally identified.
    - Currently, dbSNP at NCBI over 10 million human SNPs
    - Half of these have been validated
    - A SNP occurs every 100 to 300 base pairs
    - Fewer than 1% of SNPs alter protein sequence

9. Noncoding RNAs are also important (for example miRNAs)

# Sequencing the Human Genome



**Cost per Raw Megabase of DNA Sequence**

Moore's Law

National Human
Genome Research
Institute

genome.gov/sequencingcosts

# Sequencing the Human Genome



Cost per Genome — National Human Genome Research Institute — genome.gov/sequencingcosts

# Nowadays: Work flow second-generation sequencing



DNA fragmentation

*In vitro* adaptor ligation

Generation of polony array

Cyclic array sequencing
(>$10^6$ reads/array)

Cycle 1    Cycle 2    Cycle 3

What is base 1?    What is base 2?    What is base 3?

**32**

# Gateways to the genome sequences

# Genome browsers: gateways to the genome sequences

- Over the last few decades a gigantic amount of information on DNA sequences, gene locations, gene transcripts, protein functions and so on has been gathered

- Now we will discuss several websites that provide all this information collection, and that you will use in the afternoon session

- They all contain essentially the same information, but have a differente interface, look-and-feel, viewing options

# Genome Browsers

**UCSC\***

**NCBI** http://www.ncbi.nlm.nih.gov/sites/genome

http://genome.uc

**Ensembl** http://www.ensembl.org/

**\* We will use the UCSC browser later during the course**

# Nucleotide databases

**The underlying raw DNA sequences are identical**

**EMBL** ⟷ **GenBank** ⟷ **DDBJ**

**Housed
at EBI**

**Housed
at NCBI**

**Housed
in Japan**

**European
Bioinformatics
Institute**

**National
Center for
Biotechnology
Information**

**DNA Data Bank
of Japan**

**www.ddbj.nig.ac.jp/**

**www.ebi.ac.uk/embl/**

**www.ncbi.nlm.nih.gov/Genbank/**

**Hundreds of thousands of species are represented**

# Growth of GenBank (1982-2008)

# NCBI nucleotide databases

- GenBank
  - Individual submissions
  - Bulk submissions (Genome centers)
    - High throughput sequencing (DNA)
    - Expressed Sequence Tags (mRNA)

- RefSeq
  - Curated subset of GenBank
  - "Reference" sequence
  - Single sequence per locus / molecule

# Protein sequence databases

- NCBI
  - **RefSeq** and **Protein**

- EBI
  - **Swiss-Prot** and **TrEMBL → UniProt**

- **Translated** from nucleotide sequence
- **Curated**
- **Combined**

# Accession numbers (Identifiers)

Label to unambiguously identify a sequence

Examples (all for retinol-binding protein, RBP4):

| | | |
|---|---|---|
| DNA | X02775 | GenBank genomic DNA sequence |
| | NT_030059 | Genomic contig |
| | Rs7079946 | dbSNP (single nucleotide polymorphism) |
| | | |
| RNA | N91759.1 | An expressed sequence tag (1 of 170) |
| | NM_006744 | RefSeq DNA sequence (from a transcript) |
| | | |
| protein | NP_007635 | RefSeq protein |
| | AAC02945 | GenBank protein |
| | Q28369 | UniProt protein |
| | 1KT7 | Protein Data Bank structure record |

# From Sequence to Genes: where are the genes?

- Gene prediction
  - Extrinsic
    - Search for genes based on observed mRNA / Protein sequences
    - UniGene

  - Ab initio
    - Predict genes based on genomic sequence alone
    - Promoter sequence
    - Poly(A) tail binding sites, CG content (higher in genes), splicing sites

# UniGene

- Predict genes based on ESTs
- EST:
  – DNA sequence corresponding to mRNA from expressed gene
  – ~500 base pairs long
  – Sequenced from a cDNA library

- Cluster ESTs from many cDNA libraries to predict distinct genes

# EST clusters

This is a gene with
1 EST associated;
the cluster size is 1

This is a gene with
10 ESTs associated;
the cluster size is 10

# UniGene clusters



**Number of clusters**

40986

18424  17855

13411

8288

Likely to be a real gene

5332
4607  4075  4052  3958
1902
710  210  57  17  6  1

1  2  3-4  5-8  9-16  17-3.  12,-256  257-512  33-64  65-128  513-1024  1025-2048  2049-4096  4097-8192  8193-16384  16385-32768  32769-65536

**Cluster size**

**44**

# Ensembl website (1)

# Ensembl website (2)

# Ensembl identifiers

- Gene: ENSG…

- Transcript: ENST…

- Protein: ENSP…

# NCBI website (1)

# NCBI website (2)

# NCBI identifiers

- RefSeq:
  - Chromosome: NC_
  - mRNA: NM_
  - Protein: NP_

- Genbank:
  - Many types of IDs

- Entrez gene ID:
  - Number

- OMIM ID:
  - Number

- Pubmed ID:
  - Number

- UniGene ID:
  - Abbreviation of species.number (e.g. Hs.50223)