Dr. Lars Eijssen

l.eijssen@maastrichtuniversity.nl

# Content

- Introduction
- Background
- Learning goals
- What is a database?
- What are biological sequence databases?
  - NCBI, Ensembl, UCSC
- Identifiers
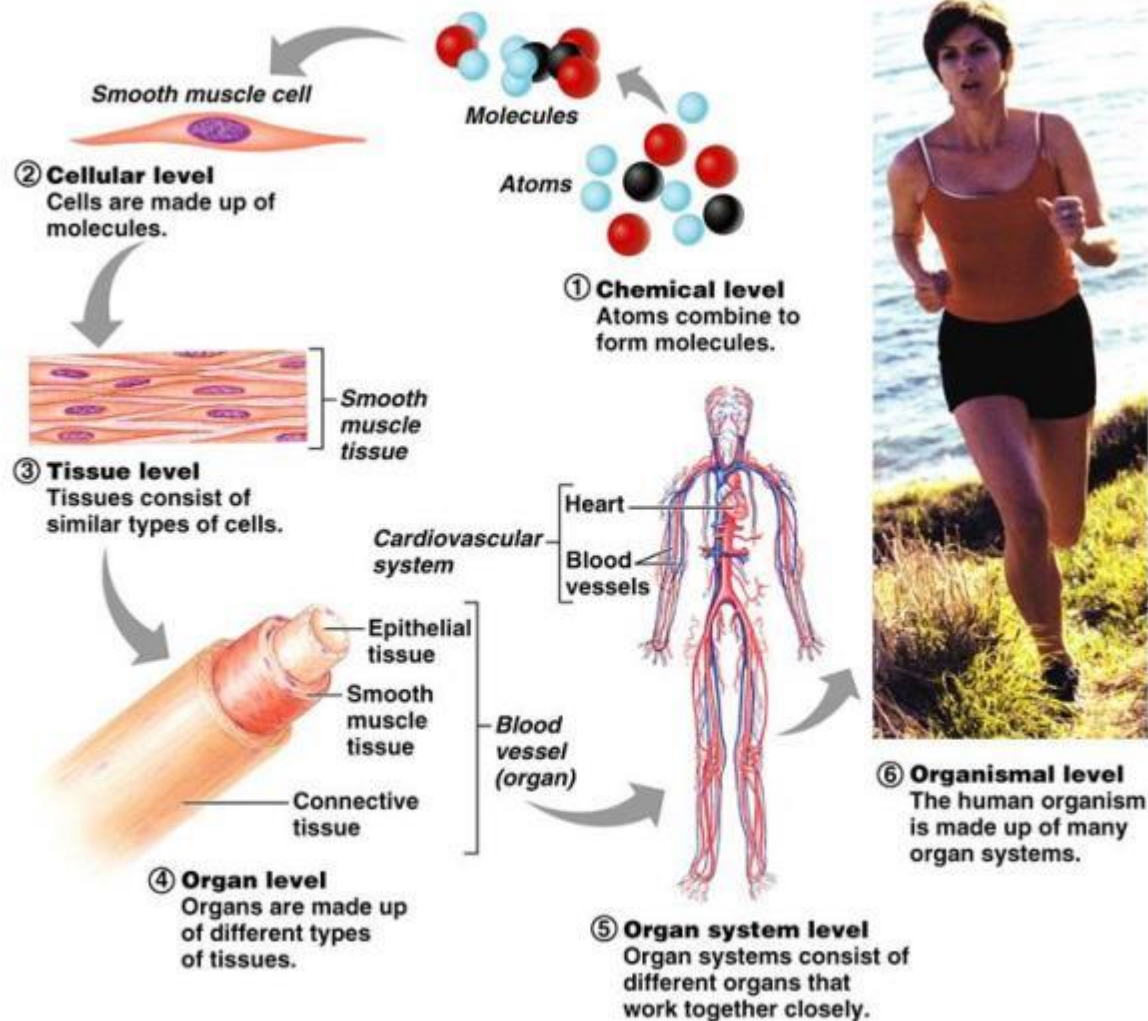- Human genome project
- ENCODE project
- Gene Ontology

# What happens with the human body when you are running?
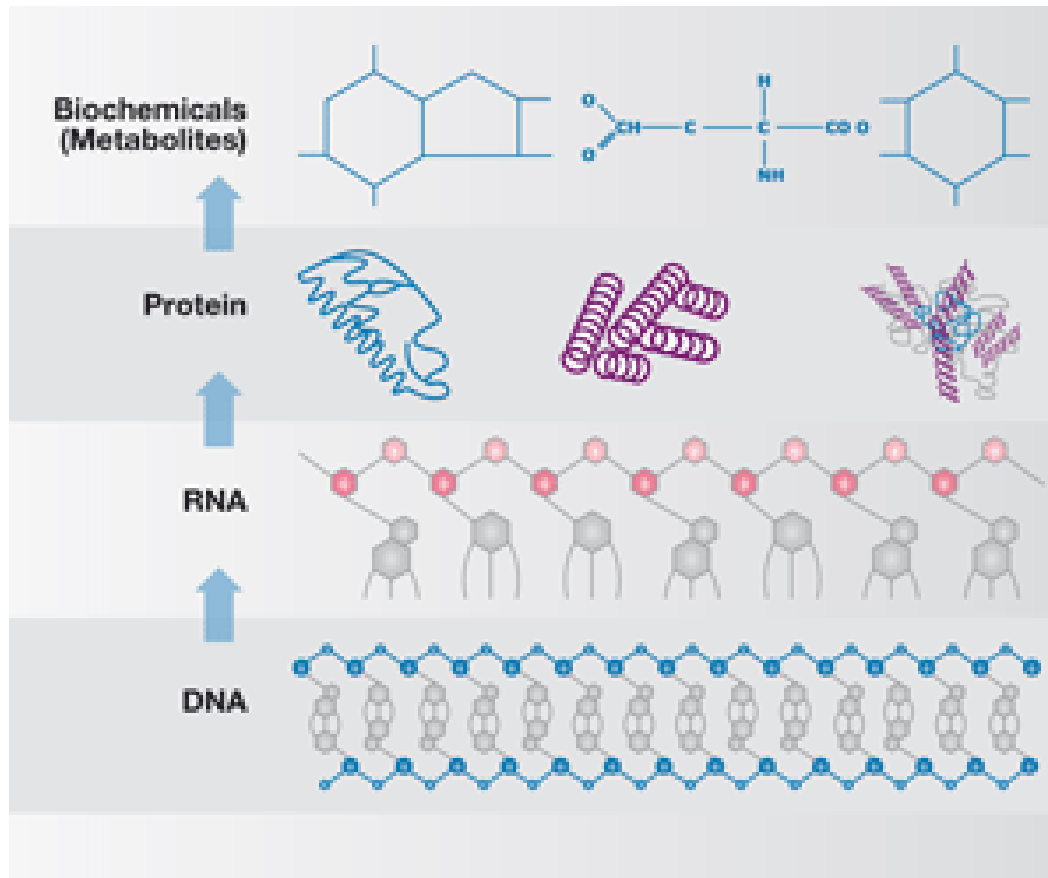
# Organ systems work together

- Muscular system - pulls on the bones to enable you to move

- Respiratory system - makes sure your muscles have enough oxygen for respiration

- Cardiovascular system- provides oxygen and glucose to the skeletal muscle cells

- Nervous system – controls your movements and heart rate
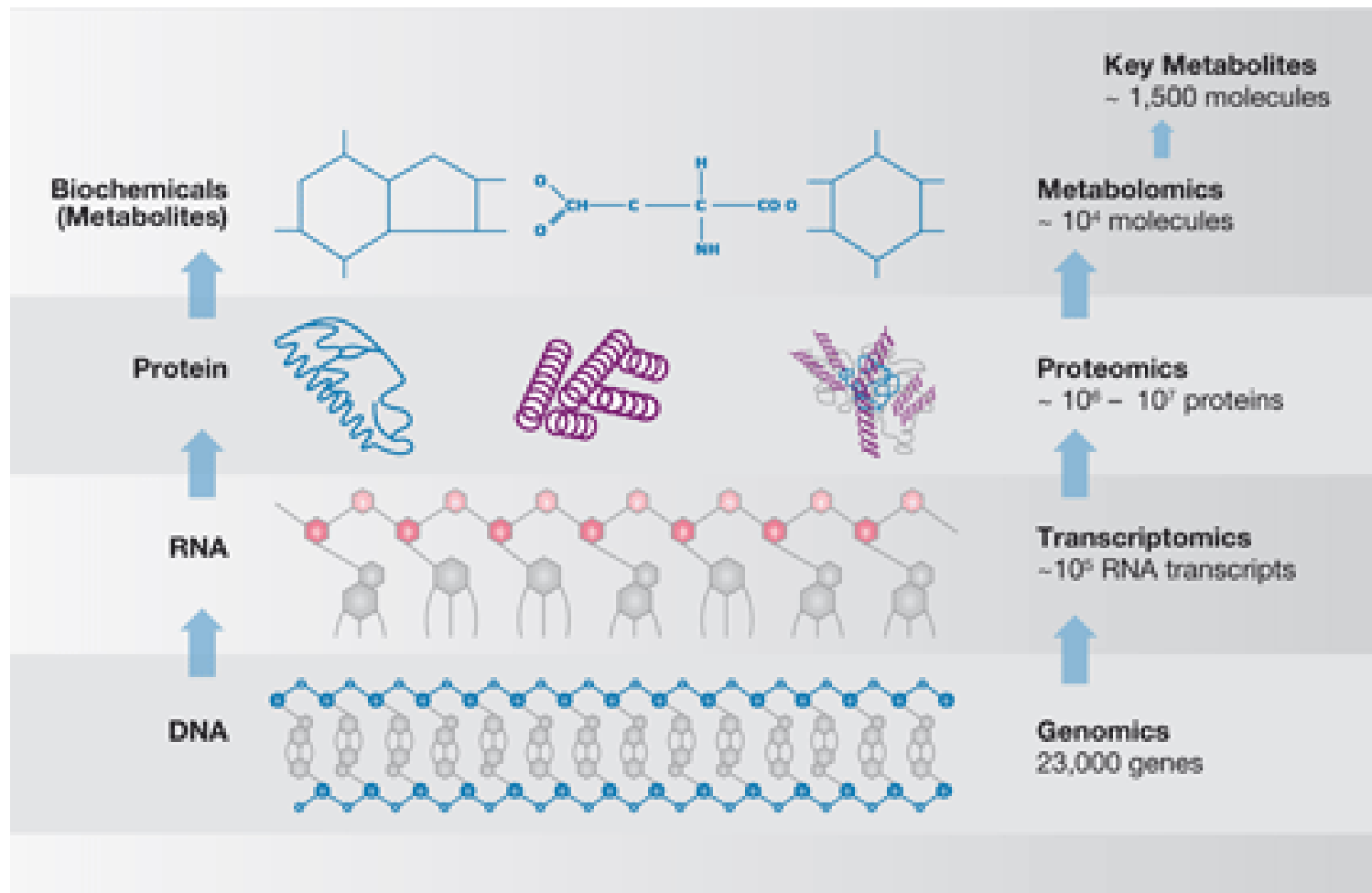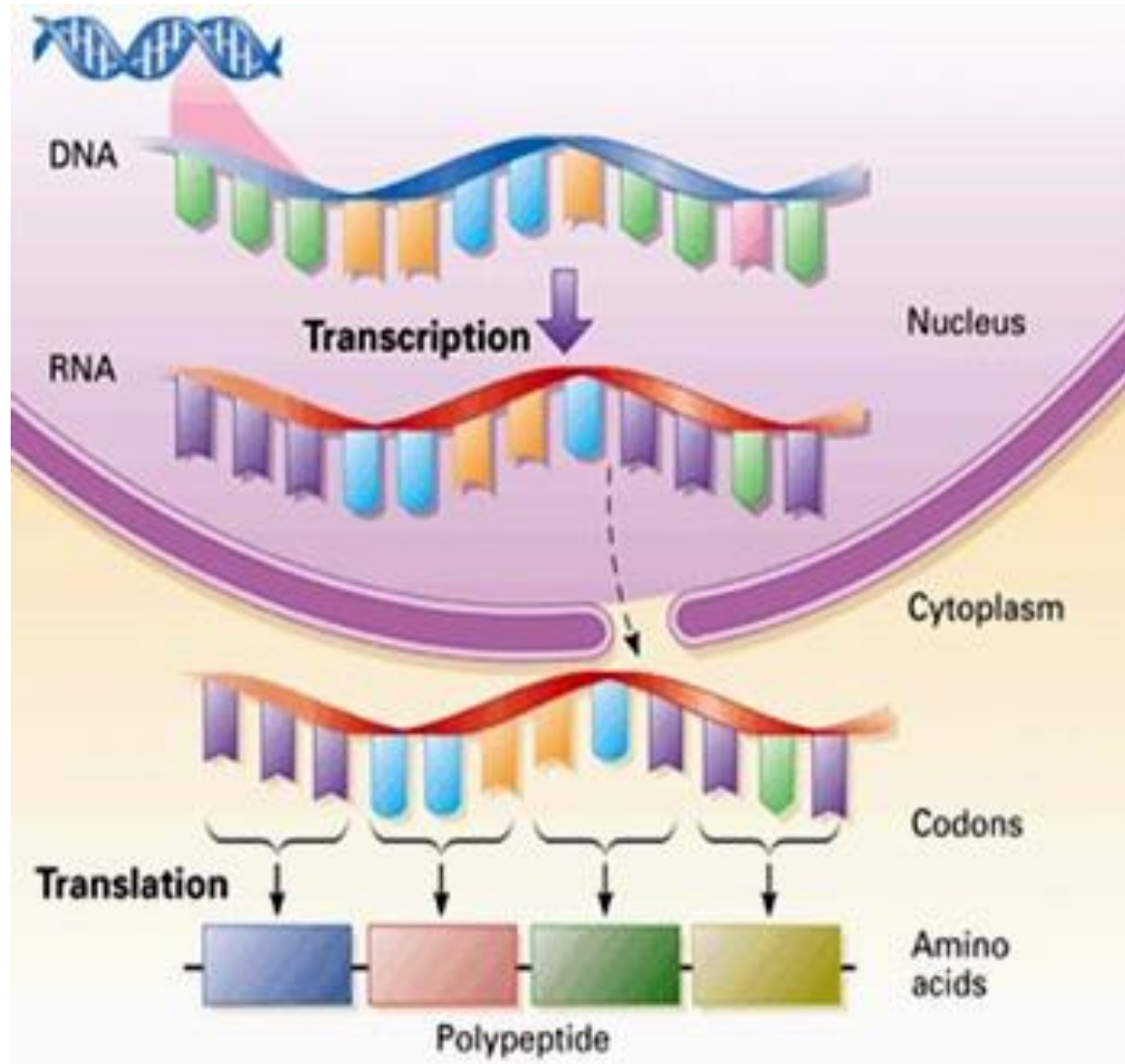
# Human body structure



Smooth muscle cell

② **Cellular level**
Cells are made up of molecules.

**Molecules**

**Atoms**

① **Chemical level**
Atoms combine to form molecules.

Smooth muscle tissue

③ **Tissue level**
Tissues consist of similar types of cells.

*Cardiovascular system*

Heart

Blood vessels

Epithelial tissue

Smooth muscle tissue

*Blood vessel (organ)*

Connective tissue

④ **Organ level**
Organs are made up of different types of tissues.

⑤ **Organ system level**
Organ systems consist of different organs that work together closely.

⑥ **Organismal level**
The human organism is made up of many organ systems.

5

Figure 1.1

# (Bio)Molecules
# Individual players are important

# Heaps of knowledge on biomolecules online available.

# Protein synthesis

# Gene structure



*Alternative splicing!*

CDS = Coding DNA Sequence
UTR = UnTranslated region

# GOAL

To understand biological sequence databases

- Which biological sequence databases are available?
- How can you find information in these databases?
- What is the content of the databases?
- Two projects aimed at deciphering the content of the human genome, the human genome project & ENCODE.
- What is Gene Ontology?

# What is a database

https://www.youtube.com/watch?v=gfT7EGibry0

(till 2:58)

# Genes in stead of persons

| Name | Identifier | Sequence | Synonyms | Chromosomal location | Disease | Many more |
|---|---|---|---|---|---|---|
| Gene 1 | 2456 | AGTCCCGT | DAH, HSD | 4q12 | Cancer | ..... |
| Gene 2 | 4333 | CGGTAACT | HGR | 7p10 | Diabetes | ....... |
| Gene 3 | 6799 | AGTCGGCGGG | | | | |
| etc | | | | | | |

## All the available information is stored in databases!

# Biological sequence databases

Originally – just a storage place for sequences.

Currently – the databases are bioinformatics work bench which provide many tools for retrieving, comparing and analyzing sequences.

1. Global nucleotide/protein sequence storage databases:
   – GenBank of NCBI (National Center for Biotechnology Information)
   – The European Molecular Biology Laboratory (EMBL) database
   – The DNA Data Bank of Japan (DDBJ)

2. Genome-centered databases
   – NCBI genomes
   – Ensembl Genome Browser
   – UCSC Genome Bioinformatics Site

3. Protein Databases
   – UniProt                    Lecture protein structures

# NCBI nucleotide databases

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

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

# Growth of GenBank



Growth of GenBank (1971-2013)

# Genome-centered databases

UCSC

NCBI

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

http://genome.ucsc.edu/

Ensembl

http://www.ensembl.org/

# NCBI homepage

# NCBI Global Cross-database search
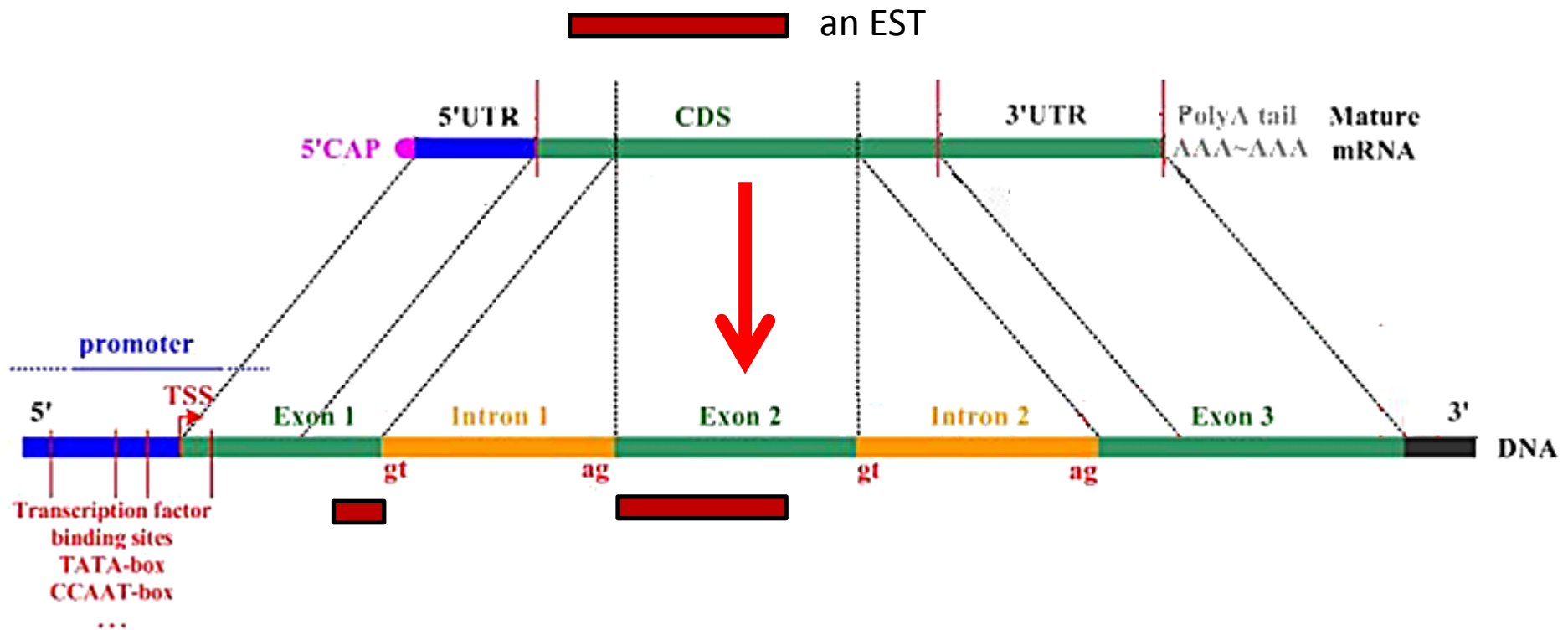## http://www.ncbi.nlm.nih.gov/gquery/

# UniGene

- EST (=expressed sequence tag):
  - DNA sequence corresponding to mRNA from expressed gene
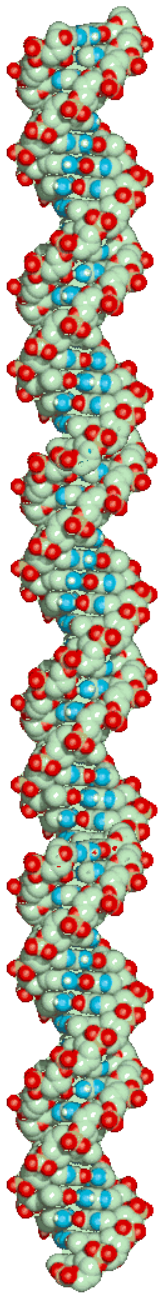  - ~500 base pairs long
  - Sequenced from a cDNA library

- Predict genes based on ESTs.
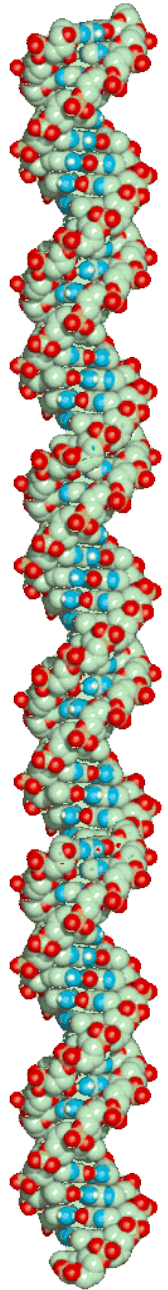
- Cluster ESTs from many cDNA libraries to predict distinct genes
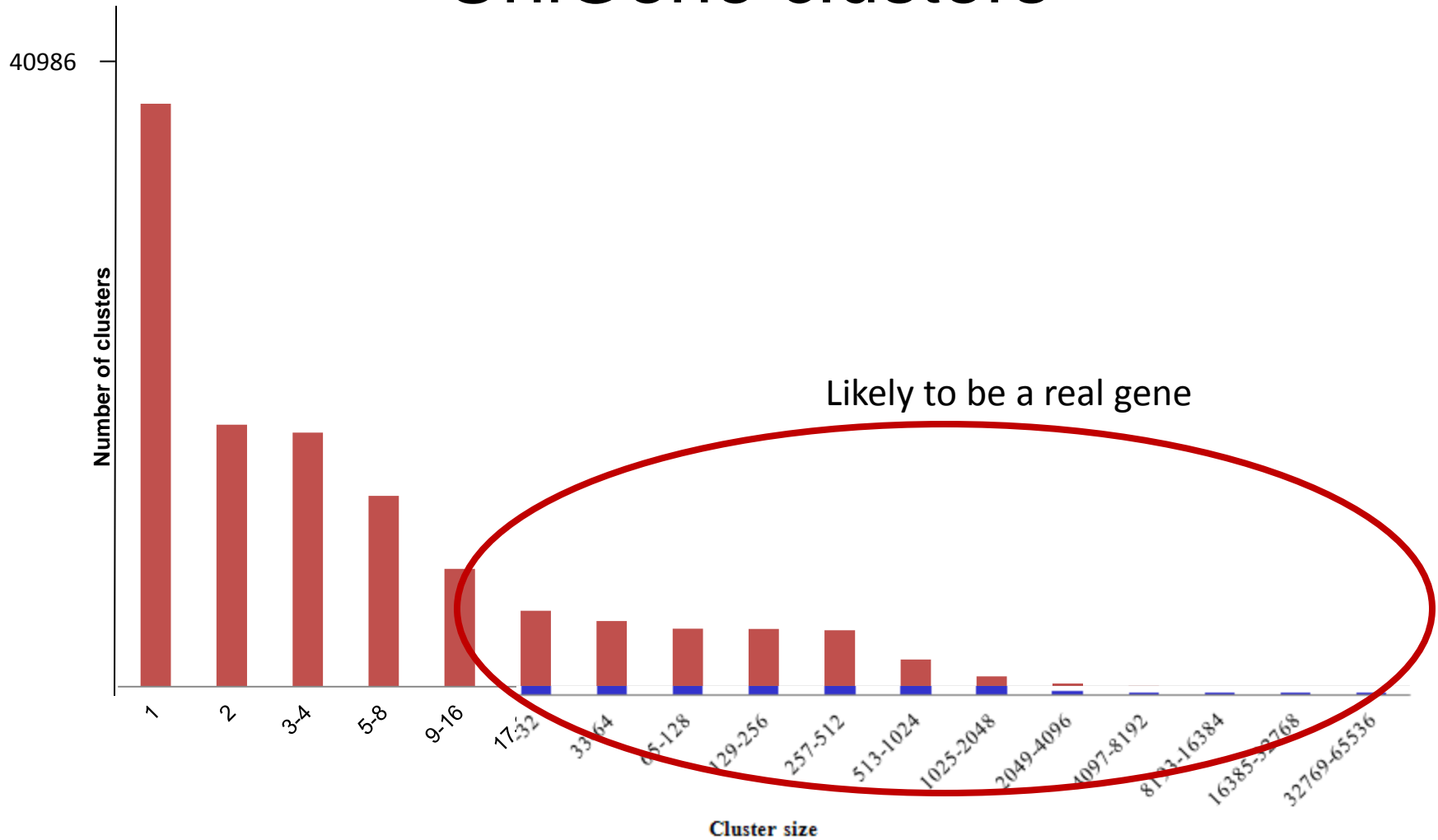
# Map mRNA (EST) back to DNA

# 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

# Gene (NCBI)
# DHH as example

# OMIM (NCBI)
# Online Mendelian Inheritance in Man



An Online Catalog of Human Genes
and Genetic Disorders

# Homology

- Homologous protein or DNA sequences share common ancestry.

- Homology need not imply similar function.

- A pair of sequences is either homologous or not homologous.



a, b homologous

a, b NOT homologous

# HomoloGene (NCBI)

# PubMed (NCBI)

# Ensembl homepage

# Ensembl
# example DHH (human)

# UCSC homepage

# UCSC: Entry page (DHH)

# Search for genomic information using identifiers

How can you store genes with a unique name?

➢ Regular gene names are not suited

- Structured identifiers
- These are different for different databases

# NCBI identifiers

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

- Genbank:
  - Many types of IDs

- NCBI (Entrez) gene ID:
  - Number

- OMIM ID:
  - Number

- Pubmed ID:
  - Number

# Ensembl identifiers

- ENSG###          Ensembl Gene ID
- ENST###          Ensembl Transcript ID
- ENSP###          Ensembl Peptide ID
- ENSE###          Ensembl Exon ID


- For other species than human a suffix is added:

  MUS (*Mus musculus*) for mouse: ENSMUSG###
  DAR (*Danio rerio*) for zebrafish: ENSDARG###, etc.

# Where does all this information come from?

- Submissions (e.g. Sequences)
- Literature
- Curators and contributors
- Automated generation by computer tools
- High-throughput lab screenings
- Individual contributions and large scale contributions

# Functional genomics

**Single biomolecules**                    **High throughput**

DNA                    *Sequencing and gene*        GENOME
⇩                      *identification*             ⇩

RNA                    *Sequencing and gene*        TRANSCRIPTOME
⇩                      *expression*                 ⇩

PROTEIN                *Identification and*         PROTEOME
                       *structure determination*

## Gezondheid

Gepubliceerd: 6 september 2012 18:42
Laatste update: 6 september 2012 18:59

Deel:

# 'Wegenkaart' menselijk DNA gepubliceerd

AMSTERDAM – Een gecoördineerde massapublicatie van 30 wetenschappelijke artikelen, waarvan zes in Nature, doet deze week vrijwel alle functies van het menselijk DNA uit de doeken.

Elk van onze cellen bevat bijna drie meter aan minutieus opgevouwen DNA. Slechts één procent daarvan doet dienst als gen. Lange tijd was dan ook de vraag: wat is het nut van al het overige, zogenaamde junk-DNA?

Het antwoord daarop wordt deze week gegeven door ENCODE (Encyclopedia of DNA Elements), een internationaal samenwerkingsverband tussen 440 onderzoekers uit 32 laboratoria.

Foto: ANP

## Junk-DNA

De belangrijkste vondst is dat in het menselijk 'junk-DNA' maar liefst vier miljoen genetische schakelaars liggen besloten. Deze schakelaars bepalen of een gen meer of minder actief wordt, zoals de dimmer op een schemerlamp. Het systeem van genetische schakelaars blijkt extreem complex. De computerberekeningen om de data te analyseren duurden bij elkaar opgeteld meer dan 300 jaar.

## Human Genome Project

ENCODE is een vervolg op het Human Genome Project, één van de

dimmer op een schemerlamp. Het systeem van genetische schakelaars blijkt extreem complex. De computerberekeningen om de data te analyseren duurden bij elkaar opgeteld meer dan 300 jaar.

## Human Genome Project

ENCODE is een vervolg op het Human Genome Project, één van de grootste wetenschappelijke projecten uit de geschiedenis. Hiermee werd in 2003 het bijna volledige menselijke DNA uitgelezen. ENCODE ging vervolgens op zoek naar alle functionele elementen daarin. Ze vonden dat ten minste 80 procent van ons DNA een biologische functie vervult.

De resultaten vormen een doorbraak in de biologie en wellicht ook de geneeskunde. Experts vergelijken het met de wegenkaart van het menselijk DNA. Het schept enorme potentie voor de ontwikkeling van nieuwe medicatie voor een veelvoud aan ziektes. Al moet daar, gezien de complexiteit, nog wel een slag om de arm worden gehouden.

Door: NU.nl/Kevin Janssen

37

# HGP and ENCODE

- We will now discuss these two major projects that contributed a lot of data

- The Humane Genome Project (1990-2003)
  - Sequencing of the human genome
  - Characterizing the genes on the DNA sequence

- The ENCODE project (2003-2012)
  - Focuses on regulatory elements on the DNA

# the Human Genome Project

AGTCCGCGAATACAGGCTCGGT

[movie](#)

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

# The human genome project

**HGP aim:** sequence the entire human genome and provide the data free to the world.

First major global collaboration of its kind and the largest biological research project ever undertaken, involving thousands of staff in institutes across the globe.

By assigning different portions of the genome to different research groups in a coordinated and efficient way, the HGP researchers were able to overcome this challenge.

**Africans**
1 Bantu
2 Mandenka
3 Yoruba
4 San
5 Mbuti pygmy
6 Biaka
7 Mozabite

**Europeans**
8 Orcadian
9 Adygei
10 Russian
11 Basque
12 French
13 North Italian
14 Sardinian
15 Tuscan

**Western Asians**
16 Bedouin
17 Druze
18 Palestinian

**Central and Southern Asians**
19 Balochi
20 Brahui
21 Makrani
22 Sindhi
23 Pathan
24 Burusho
25 Hazara
26 Uygur
27 Kalash

**Eastern Asians**
28 Han (S. China)
29 Han (N. China)
30 Dai
31 Daur
32 Hezhen
33 Lahu
34 Miao
35 Orogen
36 She
37 Tujia
38 Tu
39 Xibo
40 Yi
41 Mongola
42 Naxi
43 Cambodian
44 Japanese
45 Yakut

**Oceanians**
46 Melanesian
47 Papuan

**Native Americans**
48 Karitiana
49 Surui
50 Colombian
51 Maya
52 Pima

# Genome sequencing: general principle

**Genome**

**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

ACGTGACCGGTACTGGTAACGTACA
CCTACGTGACCGGTACTGGTAACGT
ACGCCTACGTGACCGGTACTGGTAA
CGTATACACGTGACCGGTACTGGTA
ACGTACACCTACGTGACCGGTACTG
GTAACGTACGCCTACGTGACCGGTA
CTGGTAACGTATACCTCT...

**Sequenced genome**

# How was the human genome sequenced?

# Costs of sequencing the human genome



2001: *Human Genome Project*
2.7G$, 11 years

2007: *454*
1M$, 3 months

2001: *Celera*
100M$, 3 years

2008: *ABI SOLiD*
60K$, 2 weeks

2010: *5K$, a few days*

2009: *Illumina, Helicos*
40-50K$

2017: *100$, <24 hrs?*

Log$_{10}$(price)

Year

# When has a genome been fully sequenced?

- *N*-fold coverage
  - A typical goal is to obtain five to ten-fold coverage.
  - With next-generation sequencing typically even more, like 30-fold coverage
  - Mostly both strands are sequenced
- Finished sequence
  - Usually no gaps in the sequence
  - High quality standard; error rate <0.01%.

# Genome sizes in nucleotide base pairs (log scale)



plasmids

viruses

bacteria

fungi

plants

algae

insects

mollusks

bony fish

amphibians

reptiles

birds

mammals

The size of the human genome is ~ 3 X $10^9$ bp; almost all of its complexity is in single-copy DNA.

The human genome is thought to contain ~20,000-30,000 genes.

$10^4$ $10^5$ $10^6$ $10^7$ $10^8$ $10^9$ $10^{10}$ $10^{11}$

46

http://www3.kumc.edu/jcalvet/PowerPoint/bioc801b.ppt

Species List - Mozilla Firefox

File   Edit   View   History   Bookmarks   Tools   Help

Species List   +

www.ensembl.org/info/about/species.html   Google

Most Visited   Aan de slag   Laatste nieuws   http://ftp.bigcat.unim...   https://webmail.maas...   Log in   Kääntäjä

**Collared flycatcher** (preview - assembly only)
*Ficedula albicollis*
FicAlb_1.4

**Cow**
*Bos taurus*
UMD3.1

**Dog**
*Canis lupus familiaris*
CanFam3.1

**Dolphin**
*Tursiops truncatus*
turTru1

**Duck** (preview - assembly only)
*Anas platyrhynchos*
duck1

**Elephant**
*Loxodonta africana*
loxAfr3

**Ferret**
*Mustela putorius furo*
MusPutFur1.0

**Fruitfly**
*Drosophila melanogaster*
BDGP5

**Mouse**
*Mus musculus*
GRCm38

**Mouse Lemur**
*Microcebus murinus*
micMur1

**Opossum**
*Monodelphis domestica*
BROADO5

**Orangutan**
*Pongo abelii*
PPYG2

**Painted Turtle** (preview - assembly only)
*Chrysemys picta bellii*
ChrPicBel3.0.1

**Panda**
*Ailuropoda melanoleuca*
ailMel1

**Pig**
*Sus scrofa*
Sscrofa10.2

**Pig FPC_map** (preview - assembly only)
*Sus scrofa map*
MAP

**Tilapia**
*Oreochromis niloticus*
Orenil1.0

**Tree Shrew**
*Tupaia belangeri*
TREESHREW

**Turkey**
*Meleagris gallopavo*
UMD2

**Wallaby**
*Macropus eugenii*
Meug_1.0

**Xenopus**
*Xenopus tropicalis*
JGI_4.2

**Zebra Finch**
*Taeniopygia guttata*
taeGut3.2.4

**Zebrafish**
*Danio rerio*
Zv9

Credits page for species images

**Other Metazoa**

Additional metazoan genomes (initially insect vectors and nematodes) are available from EnsemblMetazoa

**Plants and Fungi**

47

# Number of genes

| Species and Common Name | Estimated Total Size of Genome (bp)* | Estimated Number of Protein-Encoding Genes* |
|---|---|---|
| *Saccharomyces cerevisiae* (unicellular budding yeast) | 12 million | **6,000** |
| *Trichomonas vaginalis* | 160 million | **60,000** |
| *Plasmodium falciparum* (unicellular malaria parasite) | 23 million | **5,000** |
| *Caenorhabditis elegans* (worm) | 95.5 million | **18,000** |
| *Drosophila melanogaster* (fruit fly) | 170 million | **14,000** |
| *Arabidopsis thaliana* (mustard; thale cress) | 125 million | **25,000** |
| *Oryza sativa* (rice) | 470 million | **51,000** |
| *Gallus gallus* (chicken) | 1 billion | **20,000-23,000** |
| *Canis familiaris* (domestic dog) | 2.4 billion | **19,000** |
| *Mus musculus* (laboratory mouse) | 2.5 billion | **30,000** |
| *Homo sapiens* (human) | 2.9 billion | **20,000-25,000** |

Plants and amphibians with huge genomes (not in table) do not have huge amounts of genes

48

# Organization of the human genome

# Non-Protein coding DNA



NONPROTEIN-CODING SEQUENCES make up only a small fraction of the DNA of prokaryotes. Among eukaryotes, as their complexity increases, generally so, too, does the proportion of their DNA that does not code for protein. The noncoding sequences have been considered junk, but perhaps it actually helps to explain organisms' complexity.

www.carolguze.com

# The ENCODE Project: ENCyclopedia Of DNA Elements
## A public research consortium



Launched: September 2003, upgraded to the entire genome September 2007.
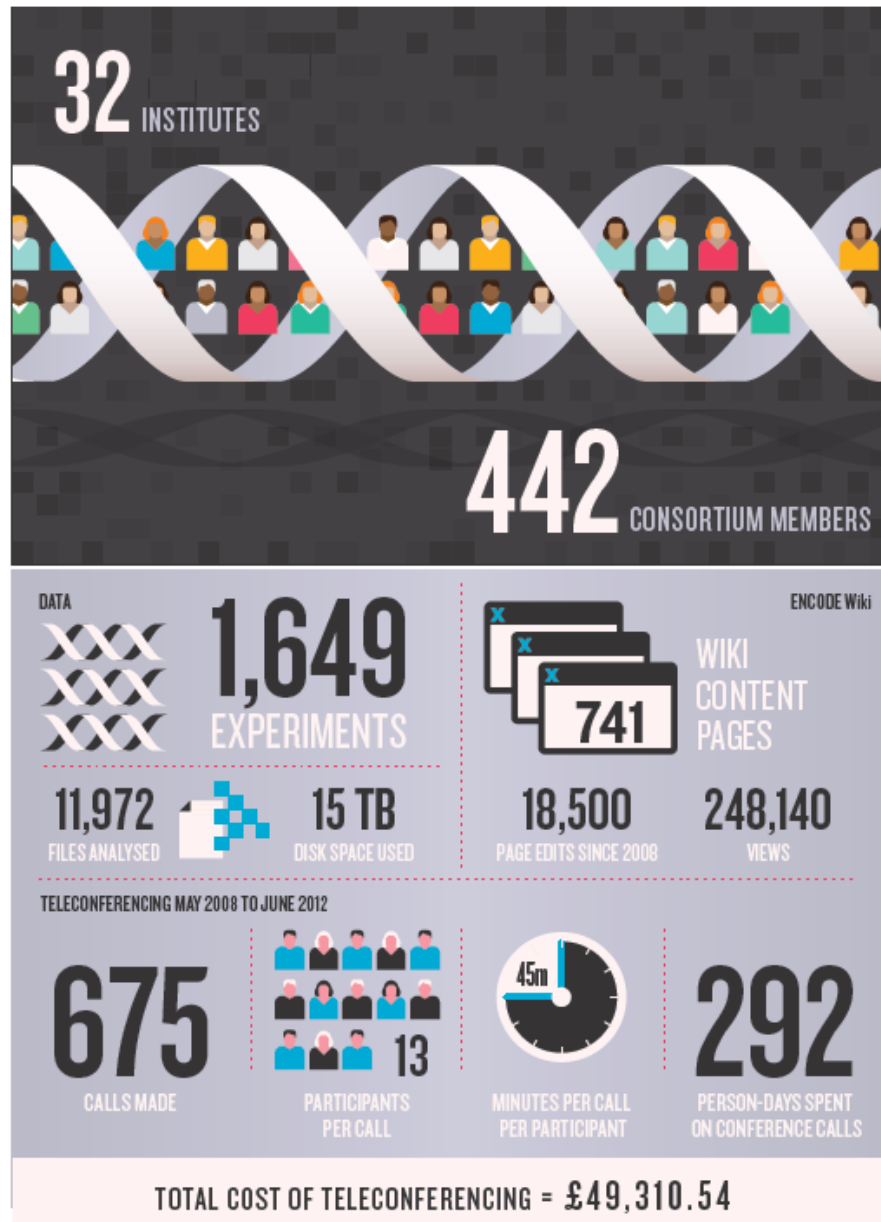
Goal: to carry out a project to identify all the functional elements in the human genome sequence.

**BY THE NUMBERS** The ENCODE project involved hundreds of people from around the world, and a lot of editing, disk space and phone calls.

**32** INSTITUTES

**442** CONSORTIUM MEMBERS

DATA

**1,649** EXPERIMENTS

ENCODE Wiki

**741** WIKI CONTENT PAGES

**11,972** FILES ANALYSED

**15 TB** DISK SPACE USED

**18,500** PAGE EDITS SINCE 2008

**248,140** VIEWS

TELECONFERENCING MAY 2008 TO JUNE 2012

**675** CALLS MADE

**13** PARTICIPANTS PER CALL

**45m** MINUTES PER CALL PER PARTICIPANT

**292** PERSON-DAYS SPENT ON CONFERENCE CALLS
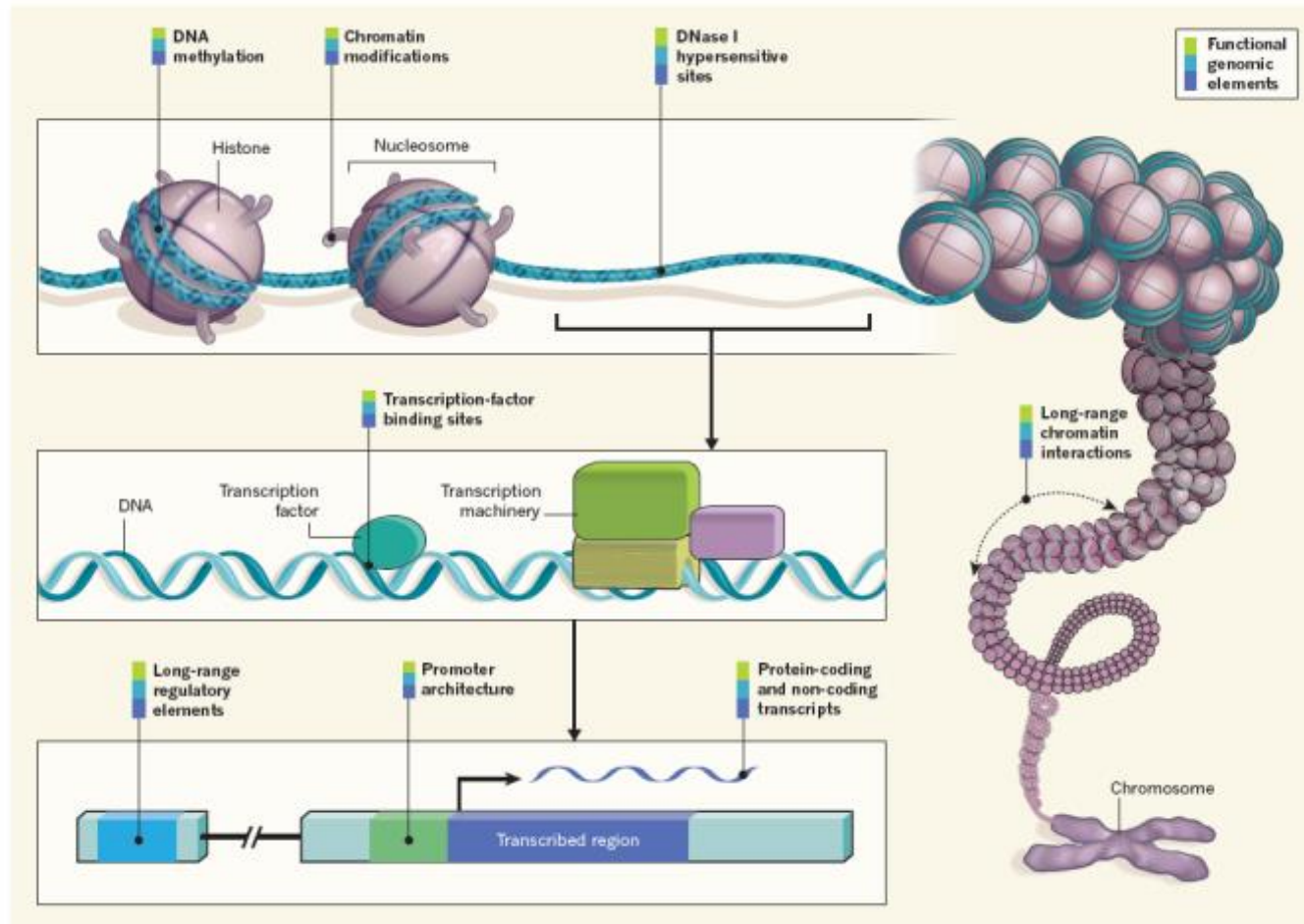
TOTAL COST OF TELECONFERENCING = £49,310.54

Understanding of the human genome is far from complete. We are missing knowledge on:
1. non-coding RNA
2. Alternatively spliced transcripts
3. Regulatory sequences

The making of ENCODE: Lessons for big-data projects. Birney E.
Nature. 2012 Sep 6;489(7414):49-51

# Data retrieved from ENCODE project



Genomics: ENCODE explained. Ecker JR, Bickmore WA, Barroso I, Pritchard JK, Gilad Y, Segal E Nature. 2012 Sep 6;489(7414):52-5.

53

# ENCODE data in Ensembl

# Gene Ontology

- Built for a very specific purpose:

"annotation of genes and proteins in genomic and protein databases"

- Applicable to all species



- GO covers 'normal' functions and processes
  - No pathological processes
  - No experimental conditions

# The 3 Gene Ontologies

- **Molecular Function** = elemental activity/task

  - the tasks performed by individual gene products; examples are *carbohydrate binding* and *ATPase activity*

- **Biological Process** = biological goal or objective

  - broad biological goals, such as *mitosis* or *purine metabolism*, that are accomplished by ordered assemblies of molecular functions

- **Cellular Component** = location or complex

  - subcellular structures, locations, and macromolecular complexes; examples include *nucleus*, *telomere*, and *RNA polymerase II holoenzyme*

# GO muscle contraction – tree view

# GO muscle contraction – tree view

# Gene products - Striated muscle contraction (GO:0006941)

# Searching and Browsing GO

- Gene Ontology consortium: http://geneontology.org/

- AmiGO 2 http://amigo.geneontology.org/amigo

# Practical session

- – Ensembl tutorials
- – Ensembl genome browser

- – Several NCBI databases
  - Gene
  - OMIM

- – Gene Ontology

# Questions