BIOLOGICAL DATABASES

Dr. Susan Steinbusch

susan.coort@maastrichtuniversity.nl

May 13th 2019

# Content

- Introduction
- Learning goals
- Biological sequence databases
  - Ensembl
  - NCBI
- Human genome project
- ENCODE project
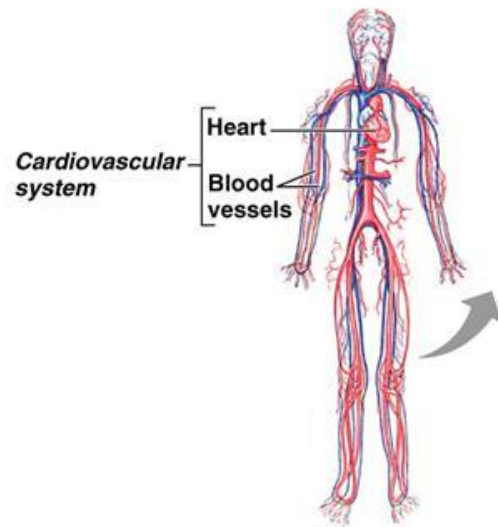- Genetic variation
- Gene Ontology
- WikiPathways

# Introduction

# What happens with the human body when you are running?

# Organ systems work together

- Skeletal system- supports the skeleton
- Muscular system - pulls on the bones to enable you to move
- Respiratory system - makes sure your muscles have enough oxygen for respiration
- Circulatory system- provides oxygen and glucose to the skeletal muscle cells

# Human body structure



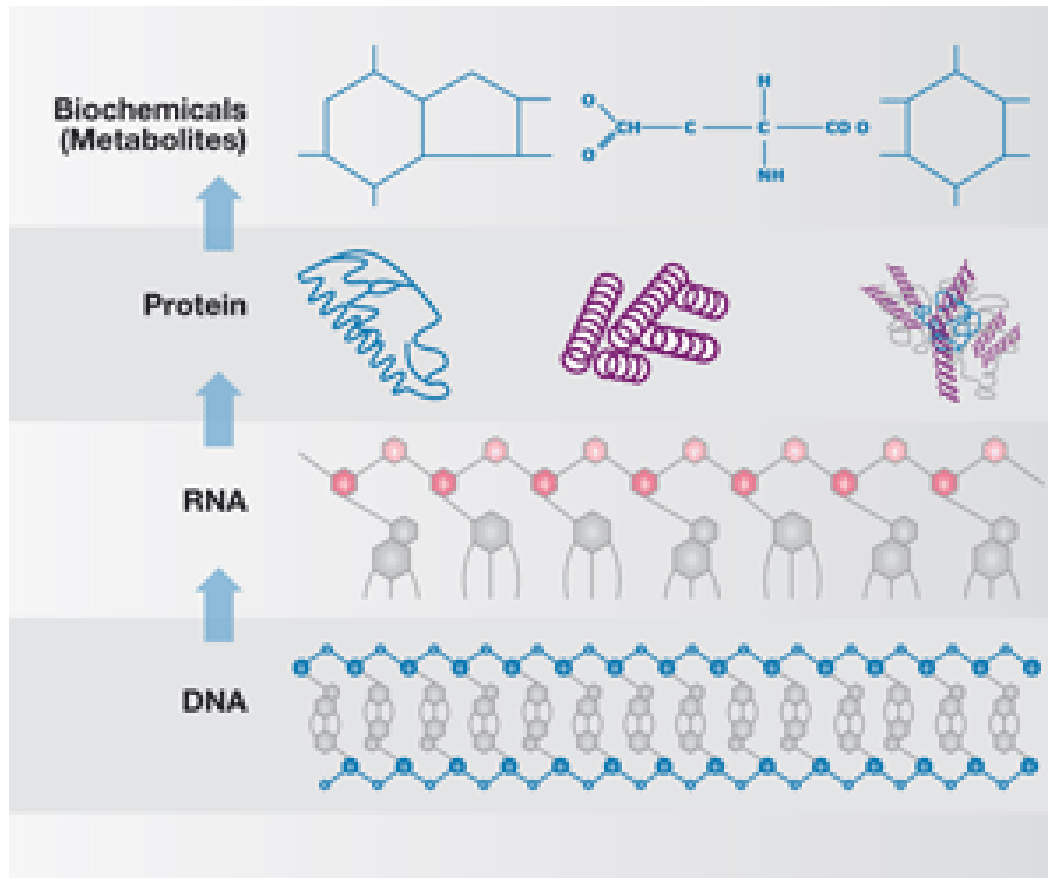Cardiovascular system — Heart, Blood vessels

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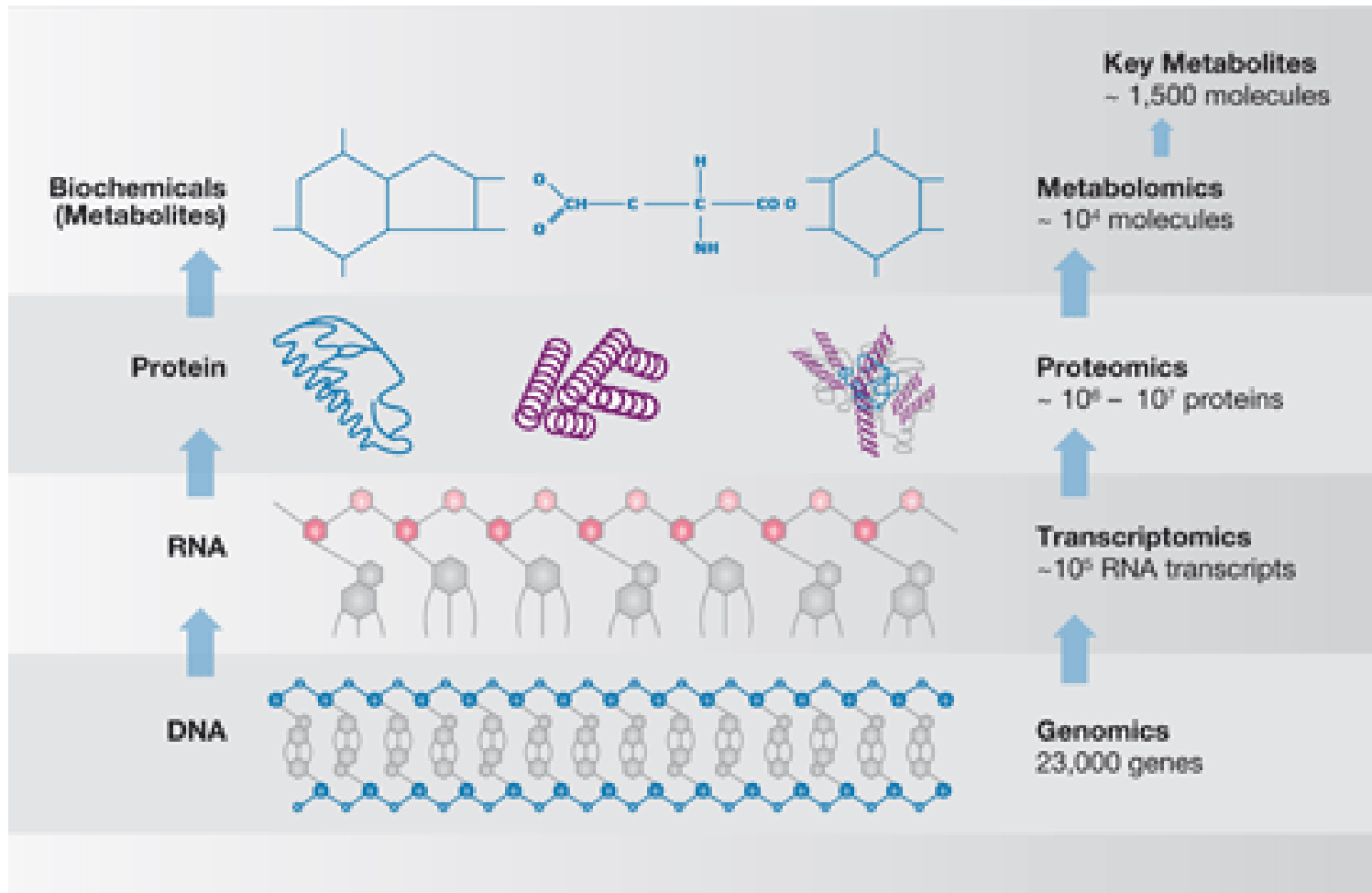The human organism is made up of many organ systems.

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

6

Figure 1.1
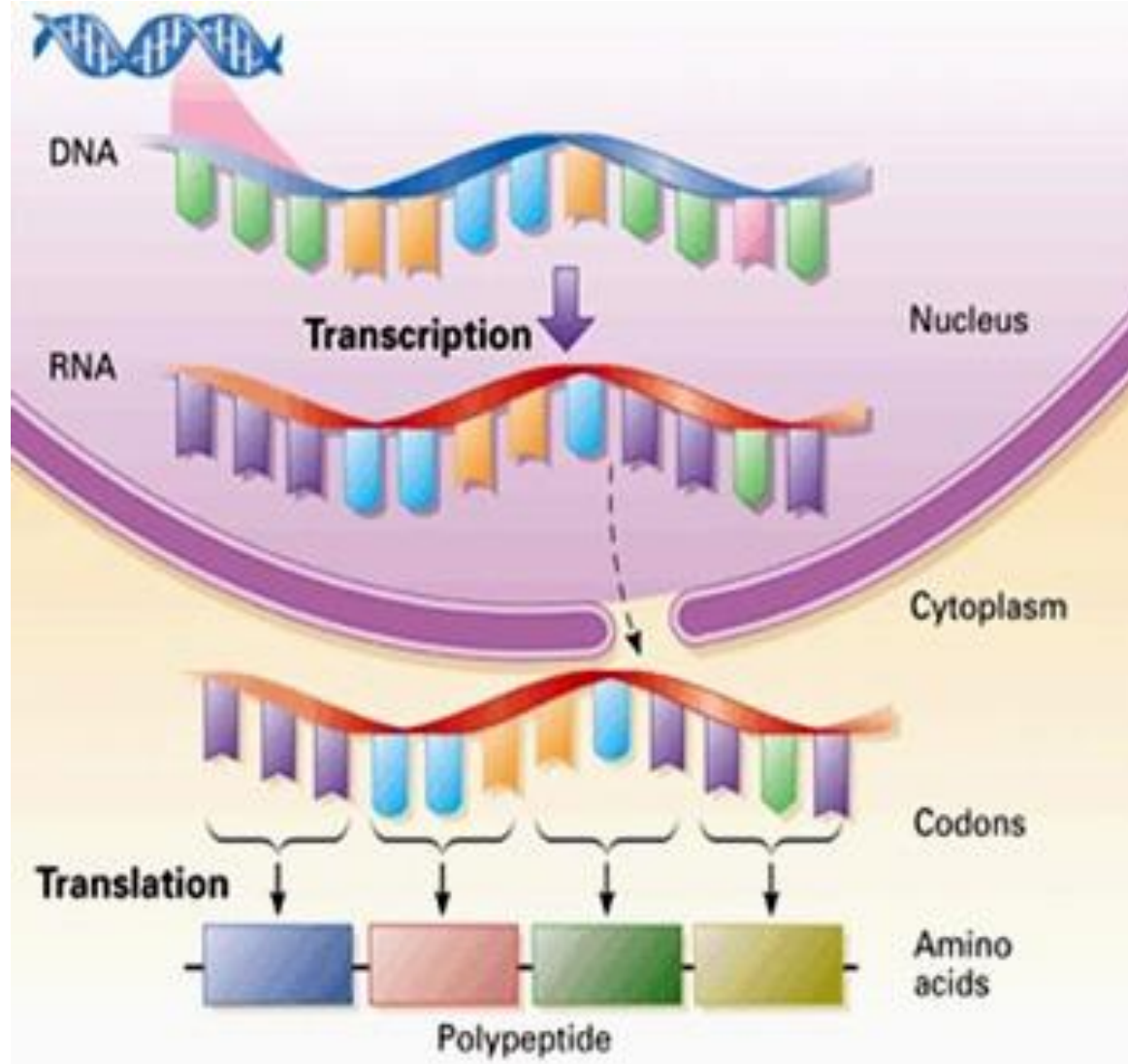
# (Bio)Molecules
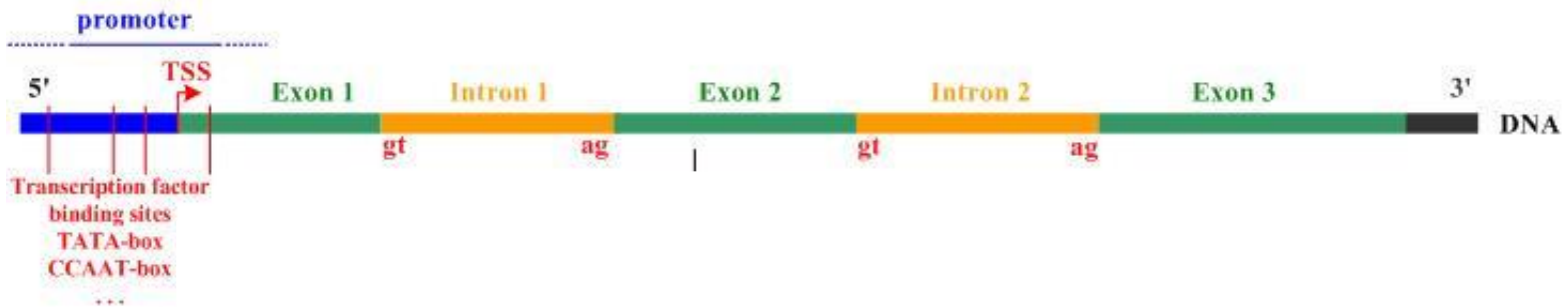# Individual players are important

# Heaps of knowledge on biomolecules online available.

# Protein synthesis

# Gene structure



*Alternative splicing!*

CDS = Coding DNA Sequence
UTR = UnTranslated region

www.carolguze.com

# Learning goals

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.
- How to find information on genetic diseases
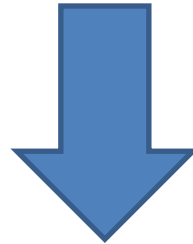- What is gene ontology and WikiPathways?

# Biological sequence databases

# What is a database

[https://www.youtube.com/watch?v=gfT7EGibry0](https://www.youtube.com/watch?v=gfT7EGibry0)

# Genes in stead of persons

| Name | Identifier | Sequence | Synonyms | Chromosomal location | Disease | Many more |
|------|-----------|----------|----------|---------------------|---------|-----------|
| Gene 1 | 2456 | AGTCCCGT | DAH, HSD | 4q12 | Cancer | ..... |
| Gene2 | 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/

**GQuery**

NCBI Global Cross-database Search

### Search NCBI databases

[                                                    ] [ Search ]

#### Literature

**PubMed**: scientific & medical abstracts/citations

**PubMed Central**: full-text journal articles

**NLM Catalog**: books, journals and more in the NLM Collections

**MeSH**: ontology used for PubMed indexing

**Books**: books and reports

**Site Search**: NCBI web and FTP site index

#### Health

**PubMed Health**: clinical effectiveness, disease and drug reports

**MedGen**: medical genetics literature and links

**GTR**: genetic testing registry

**dbGaP**: genotype/phenotype interaction studies

**ClinVar**: human variations of clinical significance

**OMIM**: online mendelian inheritance in man

**OMIA**: online mendelian inheritance in animals

#### Organisms

**Taxonomy**: taxonomic classification and nomenclature catalog

#### Nucleotide Sequences

**Nucleotide**: DNA and RNA sequences

**GSS**: genome survey sequences

**EST**: expressed sequence tag sequences

**SRA**: high-throughput DNA and RNA sequence read archive

**PopSet**: sequence sets from phylogenetic and population studies

**Probe**: sequence-based probes and primers

#### Genomes

**Genome**: genome sequencing projects by organism

**Assembly**: genomic assembly information

**Epigenomics**: epigenomic studies and display tools

**UniSTS**: sequence-tagged sites for genome mapping

**SNP**: short genetic variations

**dbVar**: genome structural variation studies

**BioProject**: biological projects providing data to NCBI

**BioSample**: descriptions of biological source materials

**Clone**: genomic and cDNA clones

20

# Gene (NCBI)
# DHH as example

# Homologene



**Homologue = One of a group of similar DNA sequences that share a common ancestry.**

# PubMed (NCBI)

# Ensembl homepage

# Ensembl
# example DHH (human)

# 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 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.

# Human Genome & ENCODE project

# 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 identification*      GENOME

⇓                                                      ⇓

RNA            *Sequencing and gene expression*      TRANSCRIPTOME

⇓                                                      ⇓

PROTEIN        *Identification and structure determination*      PROTEOME

# 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



[movie](movie)

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

# 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**

# Sanger sequencing (chain termination)

The Sanger sequencing method is the most commonly used analysis technique in genetic diagnostics.

It was also used to sequence the whole human genome.

The following are mixed in a test tube:

1. DNA template
2. One primer
3. DNA polymerase
4. dNTPs: the DNA building blocks A, C, G, T,
   -> a mixture of normal nucleotides
5. ddNTPs: Modified nucleotides with
   fluorescent markers.
   -> do not allow the chain to lengthen, so they **stop the reaction**.

# Sanger sequencing

Step 1: The primer recognizes and binds to a complementary piece of DNA.

Step 2: DNA polymerase transcribes the code using letters that are freely 'swimming around'.

Step 3: A fluorescent letter is inserted at random in a specific place -> transcription stops.

Step 4:
These fragments are arranged in order of length and separated . The fluorescent signals (which are different for each of the 4 nucleotides) are successively received by the sequencer. In this way the original code is 'assembled' by the computer.

dCTP      dGTP      dATP      dTTP
  +          +          +          +
ldCTP    ddGTP    ddATP    ddTTP

3'  A T C G A T C G A T  5'
                              Templat

5'  | | | |  3'
    T A G C   Primer

# Next generation sequencing (=Massive parallel sequencing)

- Subsequently the different samples - each with their unique bar code - are pooled

- Then every individual DNA piece from the DNA-library is replicated using its adaptor on a glass slide (flowcell). This process is called **'clonal' amplification**.



Fragments    Add adaptors    Attach to flowcell

Bind to primer    PCR extension    Dissociation    Cluster formation

Next step

# Next generation sequencing (=Massive parallel sequencing)

Finally the sequence is determined whereby all replicated (amplified) DNA fragments from the DNA-library are sequenced using sequencing primers, a polymerase enzyme and the simultaneous addition of the 4 fluorescent labelled DNA building blocks (sequencing by synthesis).



Cluster formation

Sequencing

Signal scanning

# Sequencing the Human Genome

| | | |
|---|---|---|
| **$3,000,000,000** | **2003** Human Genome Project | |
| **$20,000,000** | **2006** 1st individual genome | |
| **$2,000,000** | **2007** 1st NGS Genome | |
| **$200,000** | **2008** 1st 30x genome | |
| **$10,000** | **2010** 1st sub-10K genome | |
| **$1,000** | **2014** 1st $1,000 genome | |
| **$100** | **2017** 1st $100 genome | |

**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**

# 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

41

Pray, L. (2008) Eukaryotic genome complexity. Nature Education 1(1)

# Organization of the human genome



Human Genome Structure
from Strachan & Read, Human Molecular Genetics 2E, Wiley-Liss, 1999

# 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.

# 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.

# ENCODE data in Ensembl

# Genetic Variation

# Genetic variations

- In human beings, 99.9 percent of the bases are the same.

- Remaining 0.1 percent makes a person unique.
  - Different attributes / characteristics / traits
    - how a person looks
    - diseases he or she develops

- Most of those variations are in non-coding regions
  - This does not mean they have no effect!



Exon

Intron

Gene

Exon

# Consequences of genetic variations

- Variations can be:

    - Harmless (change in phenotype)

    - Harmful (diabetes, cancer, heart disease, Huntington's disease, and hemophilia )

    - Latent (variations found in coding and regulatory regions that are not harmful on their own, and the change in each gene only becomes apparent under certain conditions, *e.g.* susceptibility to lung cancer)

# Types of genetic variation

SNPs

Deletions

Insertions

Translocations

# Single Nucleotide Polymorphisms (SNP)

- A SNP (single nucleotide polymorphism) is defined as a single base change in a DNA sequence *that occurs in a significant proportion* (more than 1 percent) of a large population

- Currently (2017), dbSNP at NCBI (build 151) has > 100 million validated human SNPs
  - The minimal frequency criterion is <u>not</u> used

# SNP facts

- SNPs are found in
  - coding and (mostly) non-coding regions.

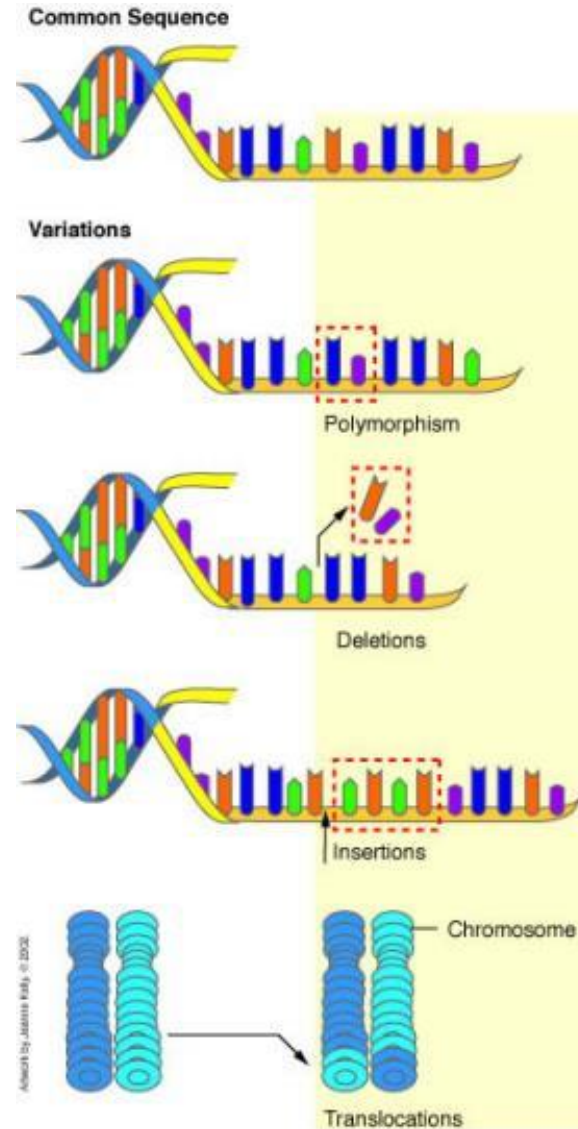- Occur with a very high frequency
  - about 1 in 1000 bases to 1 in 100 to 300 bases.

- The abundance of SNPs and the ease with which they can be measured make these genetic variations significant.

- SNPs in coding regions alter the protein sequence made by that coding region:
  - **Synonymous** SNP: no protein sequence alteration
  - **Non-synonymous** SNP: protein sequence alteration -> also known as **missense** mutation
    - Special case: a truncating SNP: premature end of protein -> also known as **nonsense** mutation

# Types of SNPs in a gene

# Inheritance of single-gene disorders

- Errors in DNA sequences

    - Autosomal dominant

    - Autosomal recessive

    - X-linked recessive
    - X-linked dominant
    - Y-linked (holandric)

Male

# NCBI - OMIM
# Online Mendelian Inheritance in Man

\* 605423

## DESERT HEDGEHOG; DHH

*HGNC Approved Gene Symbol: DHH*

*Cytogenetic location: 12q13.12* Genomic coordinates (GRCh38): *12:49,086,655–49,094,818* (from NCBI)

### Gene-Phenotype Relationships

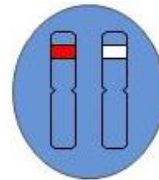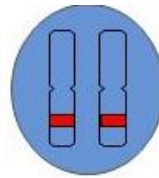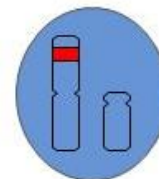| Location | Phenotype | Phenotype MIM number | Inheritance | Phenotype mapping key |
|---|---|---|---|---|
| 12q13.12 | 46XY partial gonadal dysgenesis, with minifascicular neuropathy | 607080 | | 3 |
| | 46XY sex reversal 7 | 233420 | AR | 3 |

### TEXT

#### ▼ Description

The hedgehog gene family encodes signaling molecules that play an important role in regulating morphogenesis. Mammalian hedgehog genes share striking homology to the Drosophila segment polarity gene hedgehog, a key regulator of pattern formation in the embryonic and adult fly.

#### ▼ Cloning and Expression

Tate et al. (2000) found that the human DHH gene encodes a 396-amino acid polypeptide (GenBank AB010994). ⊕

Bitgood and McMahon (1995) and Parmantier et al. (1999) showed that during development in the mouse, Dhh mRNA shows a very restricted distribution, being expressed primarily in Sertoli cells of developing testes and in Schwann cells of peripheral nerves. ⊕

**▼ External Links**

► Genome

► DNA

► Protein

► Gene Info

► Clinical Resources

▼ Variation

  1000 Genome
  ClinVar
  ExAC Beta
  GWAS Catalog
  GWAS Central
  HGMD
  HGVS
  NHLBI EVS
  PharmGKB

► Animal Models

► Cellular Pathways

# OMIM Content: Scope of Phenotypes

- Single-gene mendelian disease/disorders/phenotypes
(including: cystic fibrosis, sickle cell anemia, achondroplasia, phenotypic traits such as hair and eye color, susceptibility to drug reaction as in malignant hyperthermia and warfarin sensitivity, altered reaction to infection such as herpes simplex encephalitis and progression to AIDS in HIV infection, germline susceptibilities to cancer such as BRCA1 and breast/ovarian cancer, etc.)

- Complex diseases with significant single gene contribution (
          such as: complement factor H and age related macular degeneration)

- Descriptions of recurrent deletion and duplication syndromes
(e.g., Potocki-Shaffer syndrome, and chromosome 10q26 deletion syndrome)

# How to name a SNP? – SNP identifiers

- A standard ID for SNPs is the dbSNP ID
  - also called "rs number"
  - example: rs4986852
  - Standardised, unique, stable

- An alternative for disease related SNPs is the OMIM variation ID
  - example: 113705.0011        (this is: gene_number.SNP_number)
  - Standardised, unique, stable

- A final possibility is the
  - For non-coding or coding SNPs: variation
    - Example: BRCA1, 2978G>A
  - For coding SNPs (also): mutation
    - Example: BRCA1, SER1040ASN
  - Easier to interpret, but not stable

# Gene Ontology

# Gene Ontology

- Built for a very specific purpose:

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

- Applicable to all species

# 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*
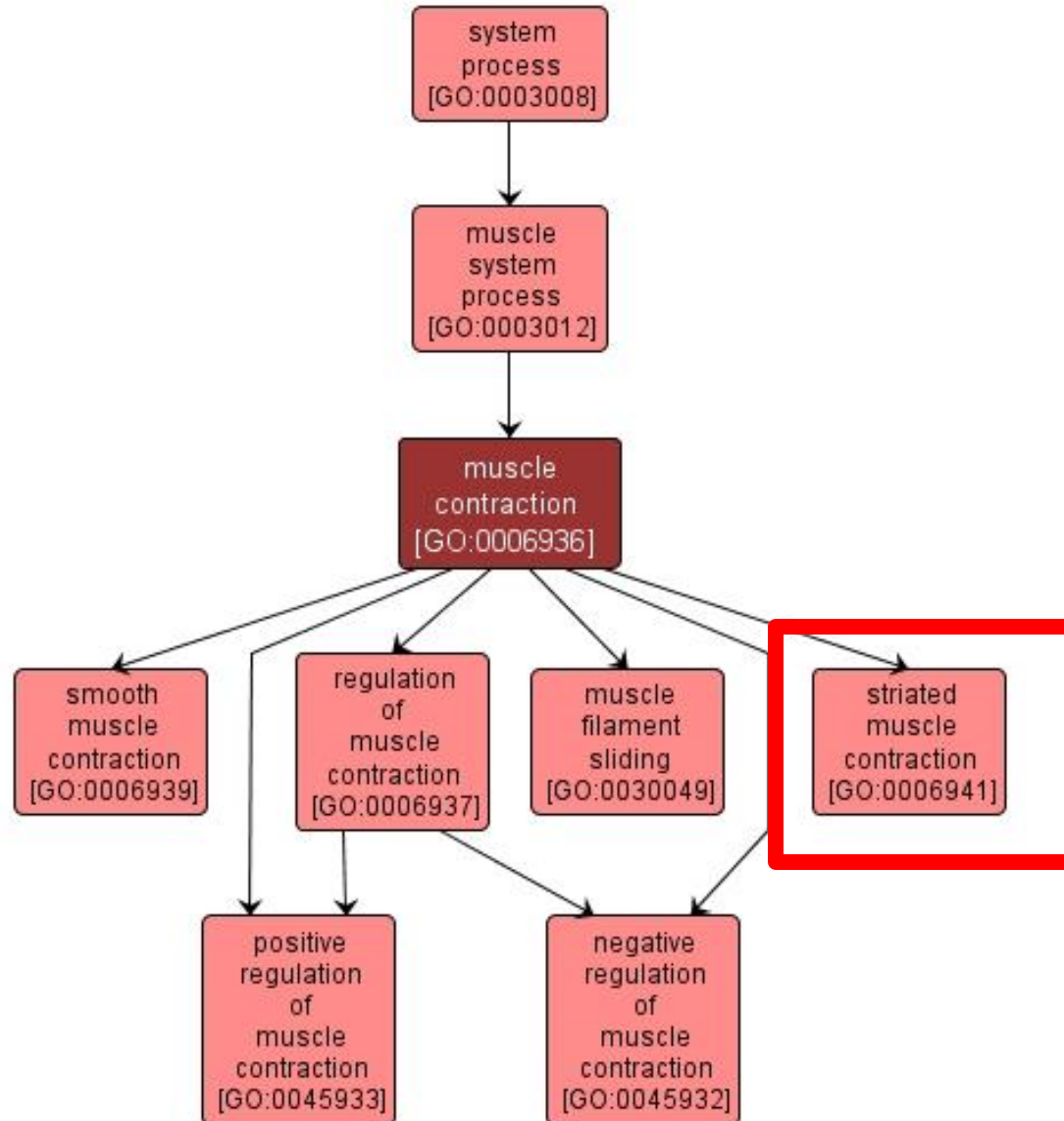
# GO muscle contraction – tree view

# GO muscle contraction – tree view

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

**striated muscle contraction**

Term associations ↓  Term information ➡  Term lineage ➡  External references ➡

## Gene Product Associations to striated muscle contraction ; GO:0006941 and children

Download all association information in: 🗋 gene association format  🗋 RDF-XML

▼ Filter associations displayed ❓

Filter by Gene Product
Gene Product Type | Data source | Species
All | All | All
complex | ASAP | Arabidopsis thaliana
gene | AspGD | Aspergillus fumig...
gene product | CGD | Aspergillus fumig...

Filter by Association
Evidence Code
All
IBA
IKR
IRD

View associations
◉ All  ○ Direct associations

[Set filters]
[Remove all filters]

1 2 3 4 5 6 7 8 9 ... 17 View all results

**striated muscle contraction** ; **GO:0006941**  [show def]  [view in tree]

| | Symbol, full name | Information | Qualifier | Evidence | Reference | Assigned by |
|---|---|---|---|---|---|---|
| ☐ | Aldoa<br>aldolase A, fructose-bisphosphate | 15 associations **protein** from *Mus musculus* | | ISO<br>With UniProtKB:P04075 | MGI:MGI:4834177 | MGI |
| ☐ | Aldoa<br>aldolase A, fructose-bisphosphate | 27 associations **gene** from *Rattus norvegicus*<br>BLAST | | ISO<br>With RGD:735815 | RGD:1624291 | RGD |
| ☐ | ALDOA<br>Fructose-bisphosphate aldolase | 12 associations **protein** from *Bos taurus*<br>BLAST | | IEA<br>With Ensembl:ENSP00000378669 | GO_REF:0000019 | Ensembl (via UniProtKB) |
| ☐ | ALDOA<br>Fructose-bisphosphate aldolase A | 29 associations **protein** from *Homo sapiens*<br>BLAST | | IMP | PMID:14615364 | BHF-UCL (via UniProtKB) |
| ☐ | Arg2<br>arginase 2 | 35 associations **gene** from *Rattus norvegicus*<br>BLAST | | IEA<br>With Ensembl:ENSMUSP00000021550<br><br>ISO<br>With RGD:736823 | RGD:1600115<br><br>RGD:1624291 | Ensembl (via RGD)<br><br>RGD |
| ☐ | Arg2<br>arginase type II | 13 associations **protein** from *Mus musculus*<br>BLAST | | IMP | PMID:16537391 | MGI |

65

# Searching and Browsing GO

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

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

# WikiPathways

# WikiPathways

- Biological pathway database
  www.wikipathways.org

- Founded in 2008 by Gladstone Institutes and the Department of Bioinformatics in Maastricht

- **Wiki**Pathways - What is a wiki?
  "A wiki is an application, typically a web application, which allows <u>collaborative</u> modification, extension, or deletion of its content and structure."

# WikiPathways

- A Wikipedia for pathways
  - Collection and curation of knowledge
  - Community curated
    - Everybody can contribute pathways
    - Everybody can edit and curate pathways
    - Everybody can use the pathway collections
  - Tools
    - Not just images but fully annotated models
    - Interactive pathway viewer
    - Full pathway editor and analysis software: PathVisio
  - New findings can be added immediately - fast!

# Pathway pages

# Questions

# Practical session

– Ensembl tutorials

– Ensembl genome browser


– Several NCBI databases

  - Gene

  - OMIM


– WikiPathways

# QUIZ at GoSoapBox

- Go to **app.gosoapbox.com** on your own computer, tablet, or smartphone.

- Type in **233-291-104** in the Access Code field.

- Enter your name prior to joining.