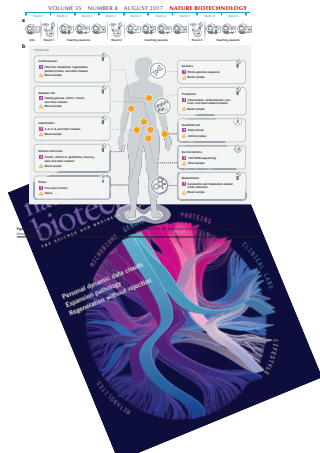


# Statistical Genomics: Master of Science in Bioinformatics and Master of Science in Statistical Data Analysis

Lieven Clement  
*Ghent University, Belgium*

# Scientific Integrity and Reproducible Research

## Bio-informatics research is based on empirical data



NATURE METHODS | VOL.12 903 | JULY 2015 | 485

Mass spectrometrists should search only for peptides they care about

William Stafford Noble



DATA

The open future is brighter than ever for the pharmaceutical Clinical-Pathways Office

1

2

3

4

5

Big biological impacts from big data

by Mike May | Jan 15, 2014, 9:40 AM

ONE APP.  
THOUSANDS OF JOBS

📱 📊 🔄 🌐

# Scientific Integrity and Reproducible Research

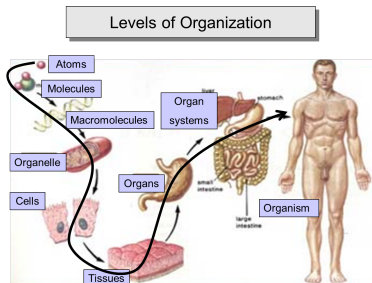
Bio-informatics research is based on empirical data



- Number of observations  $\lll$  number of features
- Need for statistics to distinguish real patterns from random patterns in high dimensional data

# Genomics

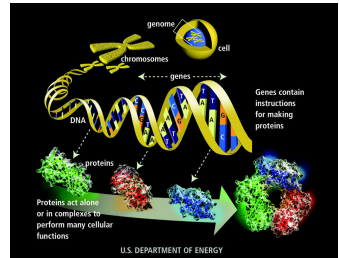
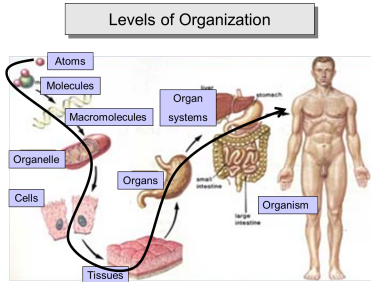
- The genome is entire hereditary information of an organism
- Contains all info needed for each function of an organism
- Most of the functions are carried out by proteins
- **Gene** is genomic region that directs synthesis of a **protein**
- **Genomics** studies all genetic information of an organism together: specific code, effects, functions and interactions



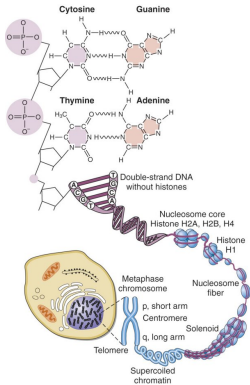


# Genomics

- The genome is entire hereditary information of an organism
- Contains all info needed for each function of an organism
- Most of the functions are carried out by proteins
- **Gene** is genomic region that directs synthesis of a **protein**
- **Genomics** studies all genetic information of an organism together: specific code, effects, functions and interactions

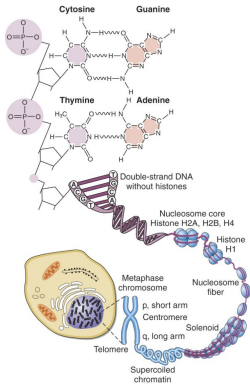


# Genome - DNA



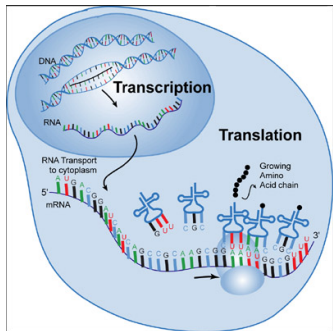
- Is stored in DNA (DeoxyriboNucleic Acid)(for many types of viruses in RNA)
- A code of 4 nucleotides
  - purines: adenine (A) and guanine (G)
  - pyrimidines: thymine (T) and cytosine (C)
  - a phosphate group;
  - a deoxyribose sugar;
- Double helix structure (2-3 hydrogen bounds)
- Organized in chromosomes
- Most of it is in the nucleus, also a part in mitochondrion (energy organelle of the cell)

# DNA structure



- Polynucleotide chains are directional molecules with slightly different ends: 3' end and 5' end.
- 3' and 5' refers to carbon atom numbering in the sugar ring. (3' hydroxyl group, 5' phosphate group)
- Complementary DNA strands are antiparallel (i.e, 5' to 3' ends for each strand are opposite)
- Most of it is coiled and condensed: very stable

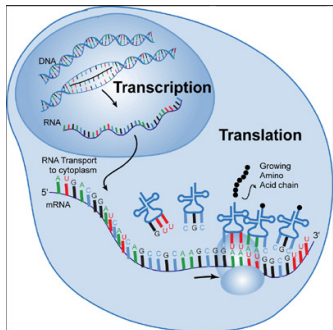
# Transcription-Translation



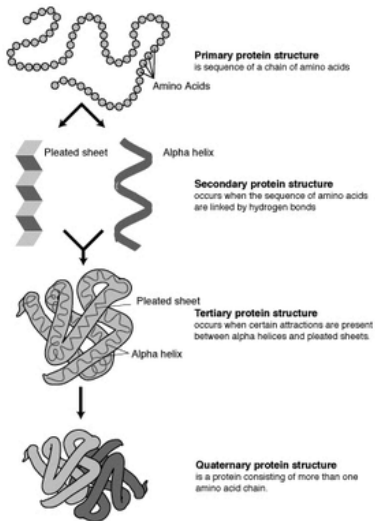
- Genome/DNA all genetic info in each cell: “Hard Drive”. Four letter code: A, C, T, G
- Transcriptome/RNA: genetic info actively used by cell: “RAM”
- Transcription
  - Unwinding of DNA
  - RNA polymerase
  - DNA template: antisense strand
  - Single complementary RNA strand
  - Splicing

# Transcription-Translation

- Genome/DNA: “Hard Drive”
- Transcriptome/RNA: active genetic info in cell: “RAM”
- Proteome
- Translation RNA→Protein
  - At ribosomes: factories of the cell
  - 24 amino acids (aa)
  - 3 consecutive bases codon
  - tRNA: with antisense codon, carries one type of aa
  - several codons exist for same aa
  - start codon AUG (methionine, often removed)
  - stop codon UAG, UAA, UGA
- Post-translational modification+protein folding



# Proteins



# The human genome



- Humans:  $2 \times 3$  billion base pairs
- 2 meters of DNA
- $\pm 20,000$  protein coding genes (500-4000/ chromosome)
- 99.9% in common with each-other
- Only 2% is protein coding



# The human genome



- Humans:  $2 \times 3$  billion base pairs
- 2 meters of DNA
- $\pm 20,000$  protein coding genes (500-4000/ chromosome)
- 99.9% in common with each-other
- Only 2% is protein coding
- 96% in common with chimp





# The human genome



- Humans:  $2 \times 3$  billion base pairs
- 2 meters of DNA
- $\pm 20,000$  protein coding genes (500-4000/ chromosome)
- 99.9% in common with each-other
- Only 2% is protein coding
- 96% in common with chimp
- 50% in common with banana



# The human genome



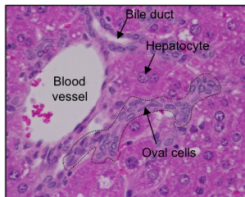
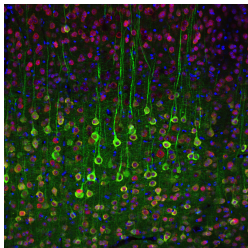
- Humans:  $2 \times 3$  billion base pairs
- 2 meters of DNA
- $\pm 20,000$  protein coding genes (500-4000/ chromosome)
- 99.9% in common with each-other
- Only 2% is protein coding
- 96% in common with chimp
- 50% in common with banana
- Organized in 23 pairs of chromosomes
  - 22 autosomal pairs
  - One sex chromosome pair: XX for females and XY for males
  - In each pair, one paternally other maternally inherited (cf. meiosis)

All cells of organism have same genome: still huge differences between different cells and over time?

Brain

vs

liver cell

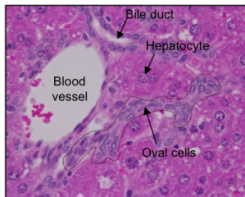
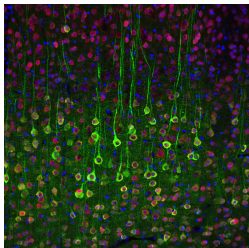


All cells of organism have same genome: still huge differences between different cells and over time?

Brain

vs

liver cell



Development of butterfly



# Differential Gene Expression

- Different genes are expressed in different cells and at different times
- Genes are expressed at different levels in different cells and over time

## Human

Tissue/Cell	Number of genes*	Fraction of genes*	Ensembl genes†
Skeletal muscle <sup>‡</sup>	11,276	0.61	11,953
Liver <sup>‡,§</sup>	11,392	0.61	12,191
BT474 <sup>¶</sup>	11,844	0.64	12,808
MB435 <sup>¶</sup>	11,847	0.64	12,726
HME <sup>‡</sup>	12,084	0.65	12,920
T47D <sup>¶</sup>	12,205	0.66	12,983
Heart	12,209	0.66	13,159
MCF7 <sup>¶</sup>	12,281	0.66	13,216
Adipose tissue	12,553	0.68	13,503
Colon	13,016	0.70	14,052
Cerebellum <sup>‡,§</sup>	13,132	0.70	14,043
Kidney	13,235	0.71	14,177
Brain <sup>‡</sup>	13,298	0.71	14,107
Breast	13,406	0.72	14,537
Lymph node	13,534	0.73	14,686
Testes	15,518	0.84	16,869

\*annotations from RefSeq, protein-coding genes.

†number of protein-coding genes, annotations from Ensembl.

‡number of genes detected in mouse: skeletal muscle 11,799; liver 11,201; brain 13,626.

§standard deviation for samples from different individuals: 106.

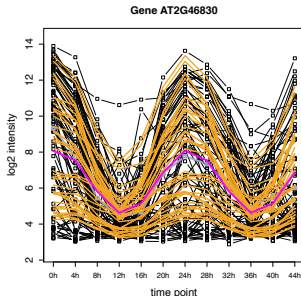
¶mean number for different individuals.

‡breast cancer cell line.

‡human mammary epithelial cell line.

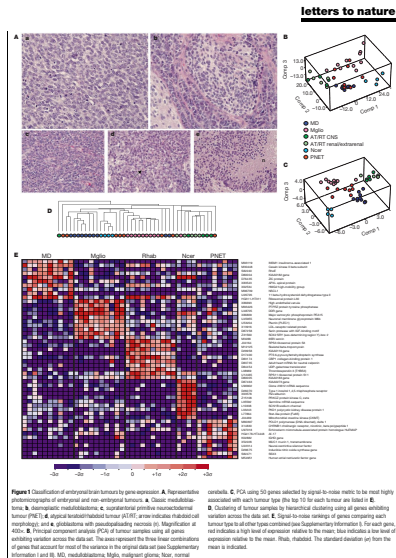
doi:10.1371/journal.pcbi.1000598.t002

## Arabidopsis Clock Gene

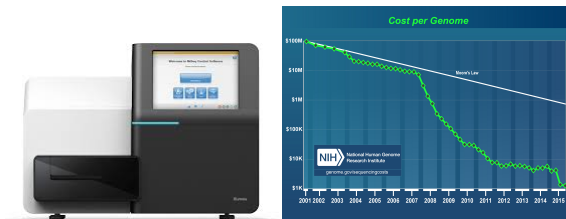


Ramsköld D et al. (2009) PLoS Comput Biol

De Beuf et al. (2012) BMC bioinformatics



# 'omics profiling



- Study all of the genome simultaneously by high throughput 'omics profiling
- Huge number of variables/features for every sample ( $p$  features)
- Number of observations  $n \lll p$
- Statistics is key to distinguish real patterns from random patterns that are observed because of we look in high dimensional data
- We can now profile gene expression at the level of individual cells!! scRNA-seq

# Topics

## Module I: Quantitative Proteomics

- 1 Identification and quantification of peptides and proteins
- 2 Data exploration and quality control using plots
- 3 Preprocessing: log-transformation, Filtering, Normalization, Summarization
- 4 Dealing with batch effects and other confounders
- 5 Statistical Concepts
  - 1 Linear models/Linear mixed models
  - 2 Trade-off between biological relevance/effect size vs statistical significance
  - 3 Empirical Bayes Methods
  - 4 Multiple testing



## Module II: Next generation sequencing (NGS, Transcriptomics)

- 1 NGS Data exploration
- 2 Preprocessing/normalization
- 3 Additional Statistical Concepts
  - 1 Generalized linear models (GLM) for binary data
  - 2 GLM for count data
  - 3 Overdispersion

# Organisation

- 1 Theory and Tutorials are blended
  - Module I: week 1-5
  - Module II: week 6-10
  - Project: week 1-10 via small assignments + week 11-12
- 2 Communication and submission of projects via Ufora
- 3 All tutorials from week 2 onwards are based on R/Bioconductor
  - via R-studio
  - Scripts are made in R/markdown: a file format to combine text, R code and R output.

→ This makes it very easy to document your analysis and to distribute them in a way which is reproducible.

# Organisation

## ④ Project

- Projects: 10/20
- Written Exam: 10/20.
  - Open book
  - Deep insight expected
  - Critical assessment of R-output,

## Projects + Master thesis

- Project 201415, Master thesis 201516: Genome biology  
<https://genomebiology.biomedcentral.com/articles/10.1186/s13059-018-1406-4>
- Project 201516: Frontiers in Neuroscience  
<https://www.frontiersin.org/articles/10.3389/fnins.2018.00136/full?report=>
- Project 201516: Analytical Chemistry  
<https://pubs.acs.org/doi/10.1021/acs.analchem.9b04375>
- Master thesis 201516: Nature Methods  
<https://www.nature.com/articles/nmeth.4338>
- Design Project 201718: Pitfalls in re-analysis of observational omics studies: a post-mortem of the human pathology atlas. submitted to Science.  
<https://www.biorxiv.org/content/10.1101/2020.03.16.994038v1>
- Master thesis 201617:  
<https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-017-1951-y>
- Master thesis 201819: Scalable differential transcript usage analysis for single-cell applications (paper in preparation, talk and poster at euroBioC meeting)

## Projects + Master thesis

- Project 201920: Fast analysis of scRNA-seq data using quasi-likelihood regression. paper in preparation
- Continuing on statistical genomics project for thesis is possible.