# Systems Biology: introduction

**PhD Program, 2023**

*Michele Caselle – University of Torino and INFN*
*caselle@to.infn.it*

# **Plan of the lecture**

1. Introduction

2. The last twenty years:    The "genomic revolution"

3. New tools and ideas: Computational Biology and Systems biology

4. Example 1: Evolutionary models

5. Example 2: Gene Regulation

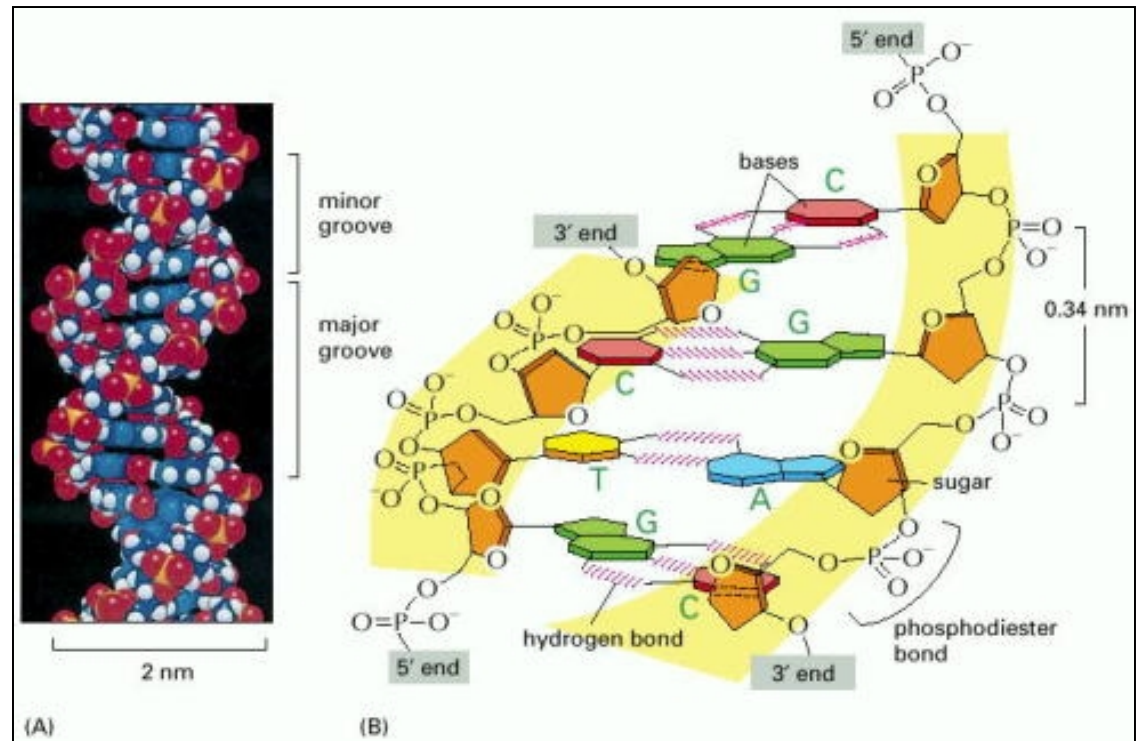6. Example 3: Identification of Cancer Driver Genes

# **References**

1. U. Alon:  An introduction to Systems Biology,
            (second edition)  CRC press

2. E.Klipp, W. Liebermeister, C. Wierling, A Kowald, R. Herwig
   Systems Biology: A Textbook  Wiley ed.

"There is no precise technical definition of a "**complex system**", but most researchers in the field would probably agree that it is a system composed of many interacting parts, such that the collective behavior of those parts together is more than the sum of their individual behaviors. The collective behaviors are sometimes also called "emergent" behaviors, and a complex system can thus be said to be a system of interacting parts that displays emergent behavior."

M.E.Newman,   Complex Systems: a Survey
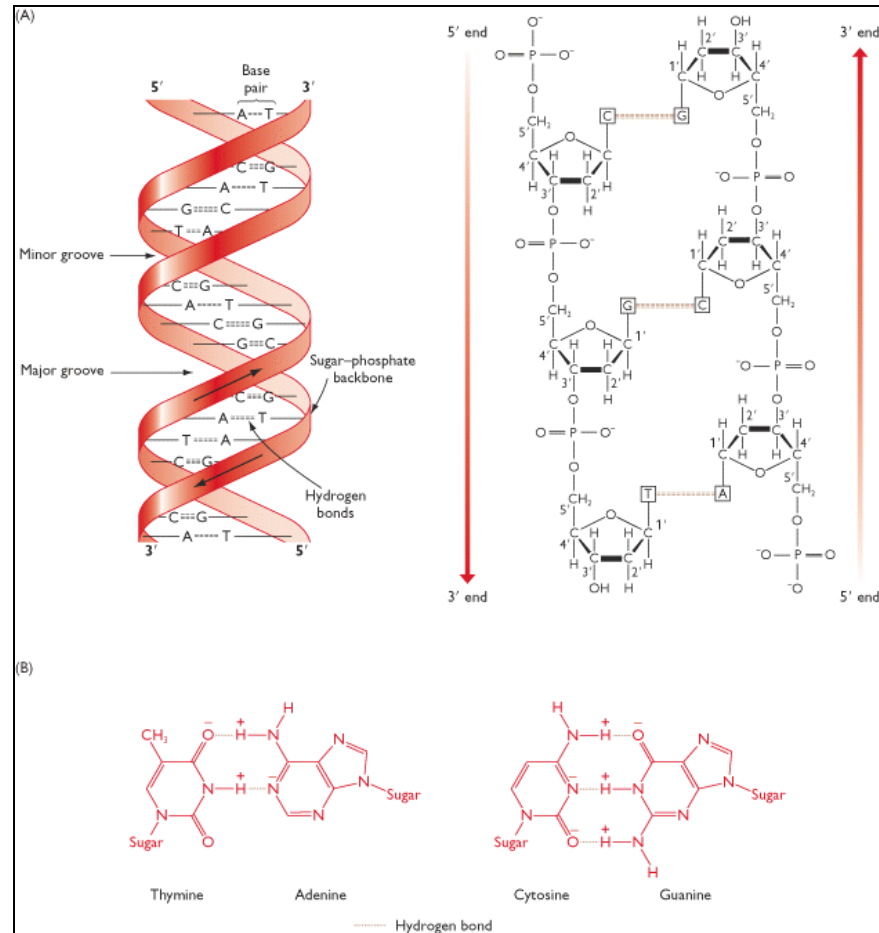http://arxiv.org/abs/1112.1440

# DNA

- Genomic information is encoded in the **DNA** chain.
- In the human case the genome is composed by 3x10^9 base pairs which may take four possible values: **A,C,G,T**

# DNA

The main property of the **DNA** chain is base pairing: (A,T) and (C,G). This allows both DNA replication and the use of the chain as a template for protein production.

# Genome size

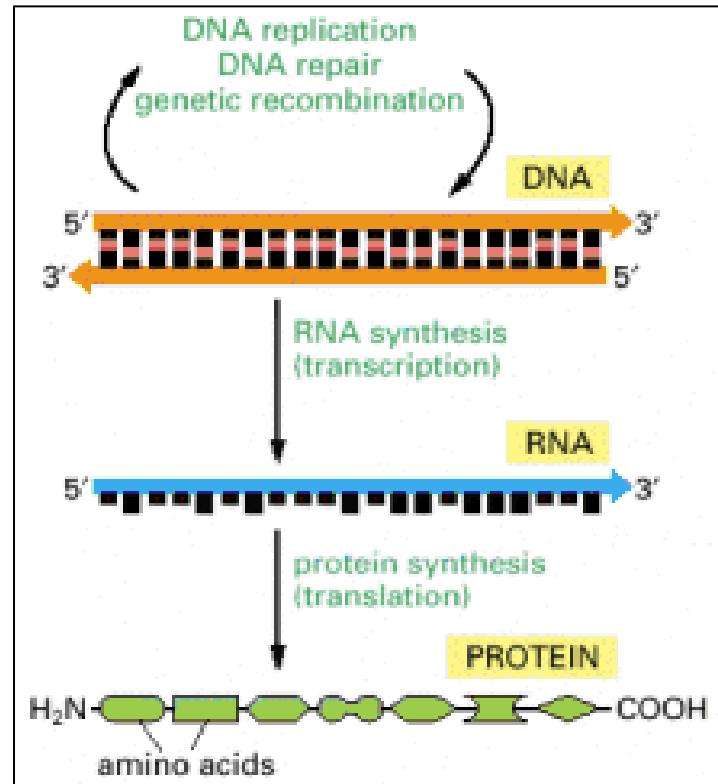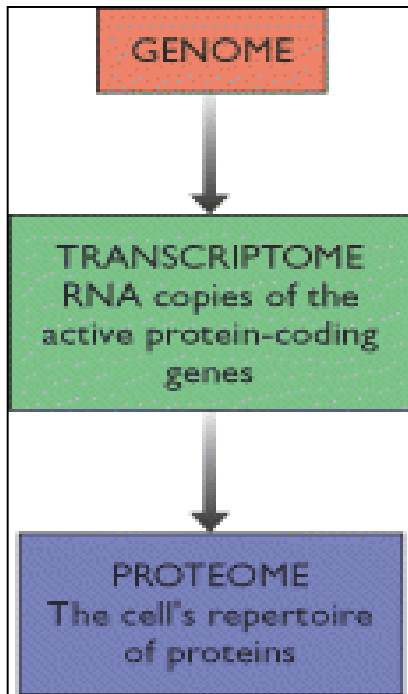| Organism | Genome length (Mbp) | Number of coding genes |
|---|---|---|
| *M. genitalium* | 0.58 | 470 |
| *E. coli* | 4.6 | 4288 |
| *S. Cerevisiae* | 12.2 | 6692 |
| *C. Elegans* | 103 | 20447 |
| *D. Melanogaster* | 144 | 13918 |
| *M. Musculus* | 3500 | 22606 |
| *H. Sapiens* | 3300 | 20300 |

# Genome Organization

- While the genome size increases with the complexity of the organism. The **number of genes is almost costant!**

- The portion of the genome coding for proteins decreases as the complexity of the organism increases. It is very high in procaryotes and yeast but very low in mammalian.

  **97% of the human genome is non-coding!!**

- Most of this non-coding DNA is involved in the **regulation of gene expression**

# "Old Paradigm": Information flow in the cell

"**Central Dogma**" of molecular biology

# Genomic Revolution

The main driving force of the Genomic Revolution was the Human Genome Project (2000)

> *homo_sapiens*
ACTTTTTTACCCTCGTGTGTTGC
AGACTTTTTGCCACTTTTAAAAC
GCTGACAATTCGACCCTTTCCAA
GTGCAAAAAGTGCCAAGATTTA
CGATAAAATTCCCCCGAGAGAC
GTGTGCA………

# New Paradigm:

**- Alternative Splicing:** one gene → many proteins

**- RNA regulatory genes:** miRNA, lncRNA → "RNA World"

**- Retrotransposition:** more than 50 % of the genome is composed by Transposons

**- Cell to cell variability:** several mRNAs are produced in few units: stochastic fluctuations are important

**- Network paradigm:** Complex functions are performed by a complex interplay of several genes

# The Present Revolution:

**- Metagenomics:**     Single bacterial strain→ microbiome

**- Horizontal Transfer:**     Genome     →    Pangenome

**- Single Cell Sequencing:**  Each cell is unique! How can we classify cells?  Thousands of different neurons in the brain!

**- Epigenomics:**     Heritable, reversible genetic information without DNA mutations

**-  Hi-C techniques:**     Tridimensional structure  of Chromatine controls gene expression

**A central role in this revolution was played by physics.**

On the experimental side:
- nanotechnolgy
- microfluidics

# Changes in instrument capacity



Timing of the major sequencing projects

**A central role in this revolution was played by physics.**

On the experimental side:
- nanotechnolgy
- microfluidics

**A central role in this revolution was played by physics.**

Both on the experimental side:
- nanotechnolgy
- microfluidics

And on the theoretical side:
 - new inference methods
 - modeling of complex systems
 - network theory
 - alignment tools

That is:

 **Computational Biology and Systems Biology**

# New Theoretical Tools:

# Systems biology and Computational  Biology

# Computational Biology

With the terms  "Computational Biology"  or "Bioinformatics"
one usually refers to all the data mining tool based
on methods and ideas coming from
mathematics / physics / statistics / computer-science .

Genomic data   (both sequences and annotations)
Can be easily downloaded from  huge "open access" data banks.

These data contain a lot of hidden information.
In general only a fraction of it has been recognized and published
by the authors of the experiments.

Relevant original results can be obtained with no need of new
costly experiments but simply using in a clever way existing data.

# Systems Biology

Network theory: Complex functions, must be described at the network level and not at the level of single genes, proteins or neurons.

Modeling: These networks can be decomposed in elementary circuits. ("network motifs") which may be modeled using differential or stochastic equations.

Ontologies: biological (and medical) information must be organized in a quantitative and standardized way

# Modern Genomics: *networks*

- genes and proteins of a given organism are organized in networks .

- Cells react to external stimuli in a "global" way.

H.Jeong et al.
Nature, 411 (2001) 41

Pajek

# *Network motifs*

Example: SIM (Single Input Module)  (a) experimental realization:
arginine biosynthesis   b) Circuit behaviour: different genes are
activated at different times as a function of their different activation
threshold as the concentration  of X (master regulator) changes in
time   R.Milo et al. Science 298 (2002) 824



Nature Reviews | Genetics

# Modern Genomics:
# *Gene Ontology*

- **Gene Ontology** is an example of standardization of biological data.

- The goal is the construction of a controlled vocabulary to describe:
  - Molecular function
  - Biological process
  - Cellular component
  
  of a given gene.

- The ontologies are organized as hierarchical networks (Directed acyclic graphs)



The G.O. Consortium
Nature Genet. 25 (2000) 25

# Systems Biology: Regulatory Networks

Example : **"Circuitry and Dynamics of Human Transcription Factor Regulatory Network**" Neph et al. CELL (2012) 150, 1274 (ENCODE collaboration).

Regulatory network in 41 different human cell lines among 475 TFs using DNAse footprinting

Cell-Specific versus Shared
Regulatory Interactions
in TF Networks of
41 Diverse Cell Types

Number of cell-types that a transcriptional regulatory interaction was observed in

Legend
Cell-specific
2+ cell types
Regulator → Regulated

Legend
Cell-specific
2+ cell types
A75 factors
Regulator → Regulated

**Visceral cells**

Hippocampal Astrocyte
HA-h

Skeletal Myoblast
HSMM

Skeletal Muscle
SKMC

Astrocyte
NH-A

**Cancer**

Neuroblastoma
SK-N-SH_RA

Hepatoblastoma
HepG2

**Embryonic Stem Cells**

Embryonic Stem Cells
H7-hESC

**Regulatory interactions in human ES cells (detail)**

Interaction Types
- Cell-specific
- 2+ cell types
- Constituitive factors
- Pluripotency factors

SOX2
POU5F1
MAX
NFYA
KLF4
CTCF
SP1

NANOG

475 factors

475 factors

Regulator → Regulated

475 transcription factors, ordered from high (SP1) to low degree (ZNF354C) in H7-hES Cells

# Conserved Architecture of Human TF Regulatory Networks

# Three paradigmatic examples

§ Evolutionary models

§ Gene Regulation

§ Identification of Cancer Driver Genes

# Evolutionary models

# Evolution at the Genomic level

There are three different processes which drive sequence evolution and accordingly there are three different **scales** at which the DNA sequence can evolve.

- Single Nucleotide Mutations (SNP)

- Gene duplications

- Whole Genome Duplication (WGD)

# Single Point Mutation



```
                        T
                        T
                        C
                        A
                        A
                        G
                        A
                        C

T → C → A              T              multiple substitutions
T                      T → C          single substitution
C → T                  C → T          parallel substitution
A                      A
A → G → C              A → C          convergent substitution
G → A → G              G              back substitution
A                      A
C                      C

              A T
              T C
              T T
              A A
              C C
              G G
              A A
              C C
```

# Single Point Mutation



**Fig. 1.2** Relative substitution rates between nucleotides under three Markov-chain models of nucleotide substitution: JC69 (Jukes and Cantor 1969), K80 (Kimura 1980), and HKY85 (Hasegawa *et al.* 1985). The thickness of the lines represents the substitution rates while the sizes of the circles represent the steady-state distribution.

**Table 1.1** Substitution-rate matrices for commonly used Markov models of nucleotide substitution

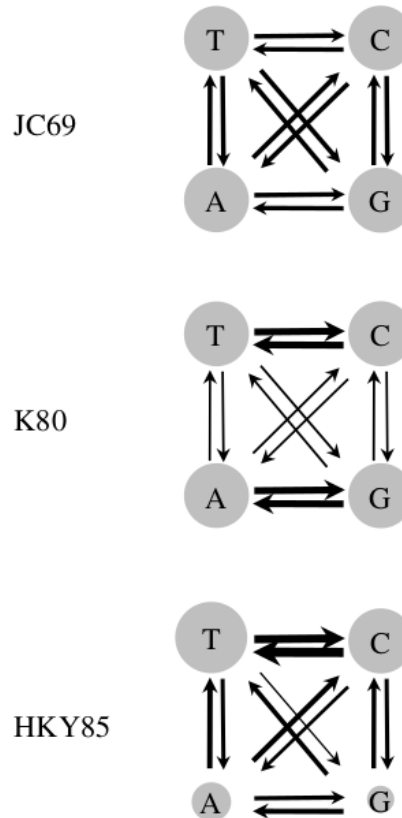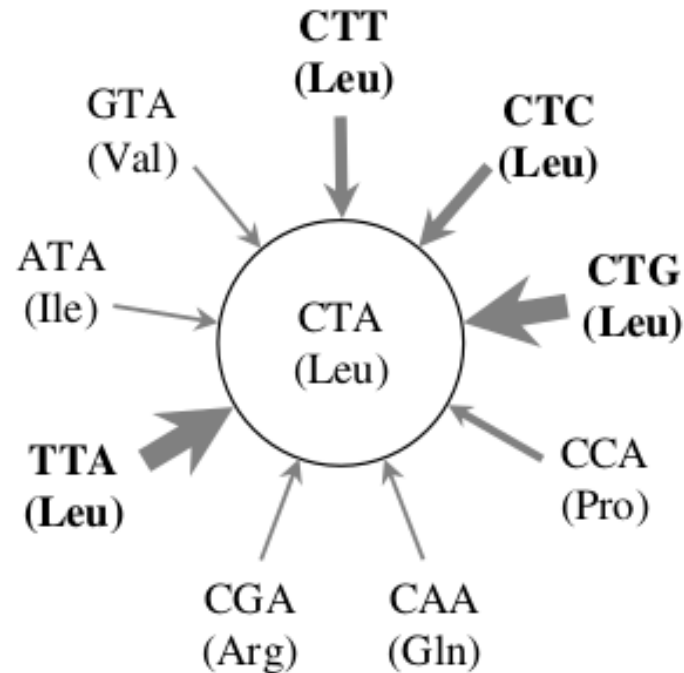| From | | To | | | |
|---|---|---|---|---|---|
| | | T | C | A | G |
| JC69 (Jukes and Cantor 1969) | T | · | $\lambda$ | $\lambda$ | $\lambda$ |
| | C | $\lambda$ | · | $\lambda$ | $\lambda$ |
| | A | $\lambda$ | $\lambda$ | · | $\lambda$ |
| | G | $\lambda$ | $\lambda$ | $\lambda$ | · |
| K80 (Kimura 1980) | T | · | $\alpha$ | $\beta$ | $\beta$ |
| | C | $\alpha$ | · | $\beta$ | $\beta$ |
| | A | $\beta$ | $\beta$ | · | $\alpha$ |
| | G | $\beta$ | $\beta$ | $\alpha$ | · |
| F81 (Felsenstein 1981) | T | · | $\pi_C$ | $\pi_A$ | $\pi_G$ |
| | C | $\pi_T$ | · | $\pi_A$ | $\pi_G$ |
| | A | $\pi_T$ | $\pi_C$ | · | $\pi_G$ |
| | G | $\pi_T$ | $\pi_C$ | $\pi_A$ | · |
| HKY85 (Hasegawa *et al.* 1984, 1985) | T | · | $\alpha\pi_C$ | $\beta\pi_A$ | $\beta\pi_G$ |
| | C | $\alpha\pi_T$ | · | $\beta\pi_A$ | $\beta\pi_G$ |
| | A | $\beta\pi_T$ | $\beta\pi_C$ | · | $\alpha\pi_G$ |
| | G | $\beta\pi_T$ | $\beta\pi_C$ | $\alpha\pi_A$ | · |
| F84 (Felsenstein, DNAML program since 1984) | T | · | $(1+\kappa/\pi_Y)\beta\pi_C$ | $\beta\pi_A$ | $\beta\pi_G$ |
| | C | $(1+\kappa/\pi_Y)\beta\pi_T$ | · | $\beta\pi_A$ | $\beta\pi_G$ |
| | A | $\beta\pi_T$ | $\beta\pi_C$ | · | $(1+\kappa/\pi_R)\beta\pi_G$ |
| | G | $\beta\pi_T$ | $\beta\pi_C$ | $(1+\kappa/\pi_R)\beta\pi_A$ | · |
| TN93 (Tamura and Nei 1993) | T | · | $\alpha_1\pi_C$ | $\beta\pi_A$ | $\beta\pi_G$ |
| | C | $\alpha_1\pi_T$ | · | $\beta\pi_A$ | $\beta\pi_G$ |
| | A | $\beta\pi_T$ | $\beta\pi_C$ | · | $\alpha_2\pi_G$ |
| | G | $\beta\pi_T$ | $\beta\pi_C$ | $\alpha_2\pi_A$ | · |
| GTR (REV) (Tavaré 1986; Yang 1994*b*; Zharkikh 1994) | T | · | $a\pi_C$ | $b\pi_A$ | $c\pi_G$ |
| | C | $a\pi_T$ | · | $d\pi_A$ | $e\pi_G$ |
| | A | $b\pi_T$ | $d\pi_C$ | · | $f\pi_G$ |
| | G | $c\pi_T$ | $e\pi_C$ | $f\pi_A$ | · |
| UNREST (Yang 1994*b*) | T | · | $q_{TC}$ | $q_{TA}$ | $q_{TG}$ |
| | C | $q_{CT}$ | · | $q_{CA}$ | $q_{CG}$ |
| | A | $q_{AT}$ | $q_{AC}$ | · | $q_{AG}$ |
| | G | $q_{GT}$ | $q_{GC}$ | $q_{GA}$ | · |

The diagonals of the matrix are determined by the requirement that each row sums to 0. The equilibrium distribution is $\pi = (1/4, 1/4, 1/4, 1/4)$ under JC69 and K80, and $\pi = (\pi_T, \pi_C, \pi_A, \pi_G)$ under F81, F84, HKY85, TN93, and GTR. Under the general unrestricted (UNREST) model, it is given by the equations $\pi Q = 0$ under the constraint $\sum_i \pi_i = 1$.
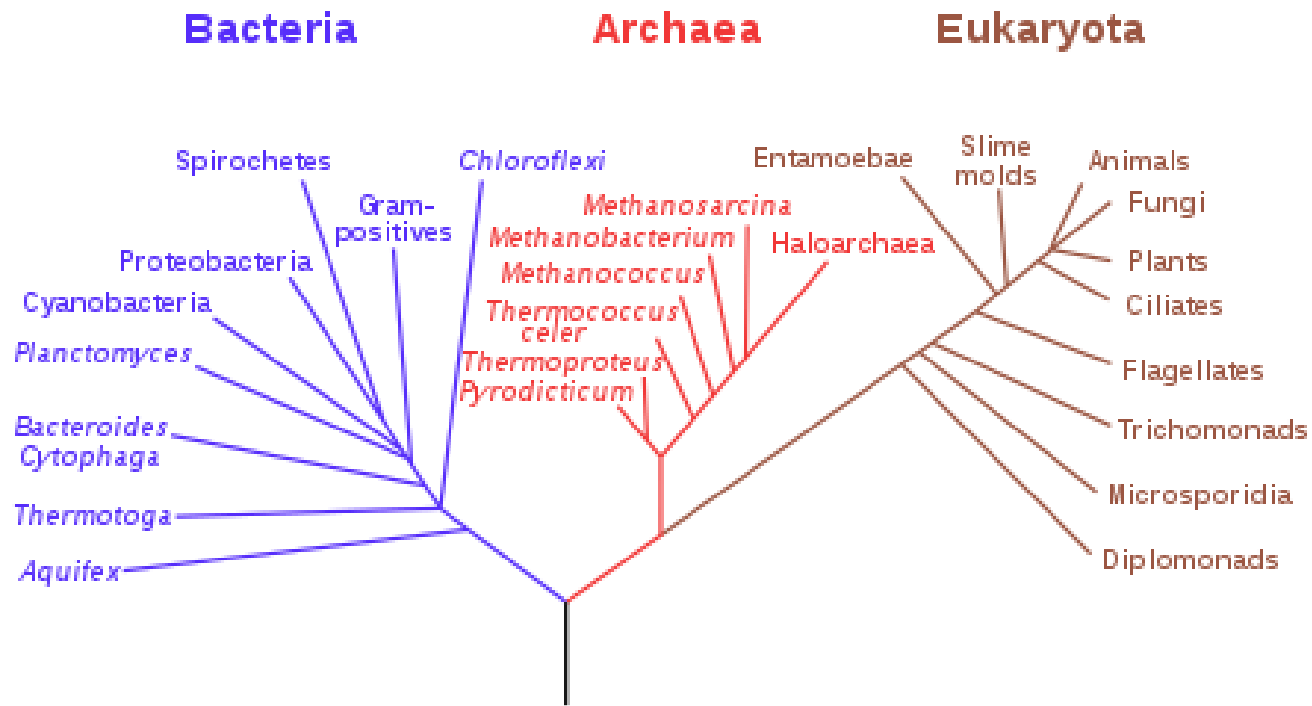
# Mutations in the coding regions

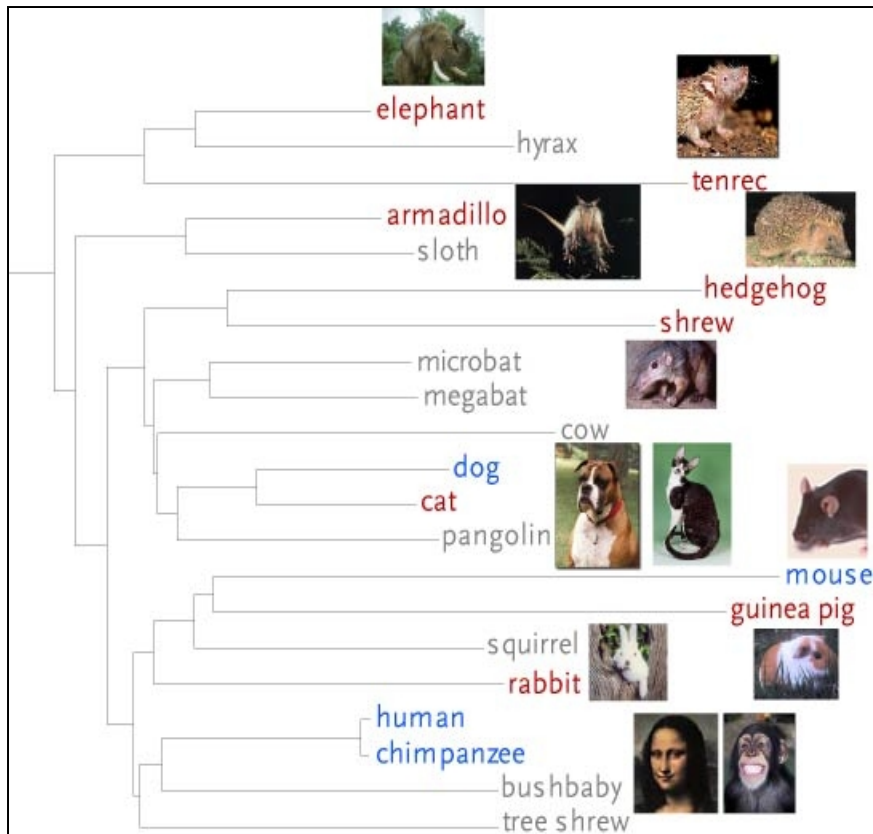Synonymous versus nonsynonymous mutations

# Important observations

- Most of the SNPs are neutral and are not selected by evolution. Since they occur at a fixed rate, by counting their number we can infer the time at which the process started. This the **"Molecular clock!"**

- By comparing the sequences of two species with **"alignment algorthms"** we can infer the number of mutations and thus the time of speciation. This is the **Genomic Tree of Life!**

- By comparing synonymous versus non-synonymous mutations we can measure the effect of selection: "**important regions of DNA are kept conserved under evolution**" . The same holds also for non coding regulatory regions. **Sequence conservation is a hallmark of Biological relevance!**

# The tree of life

# The tree of life: Zoom
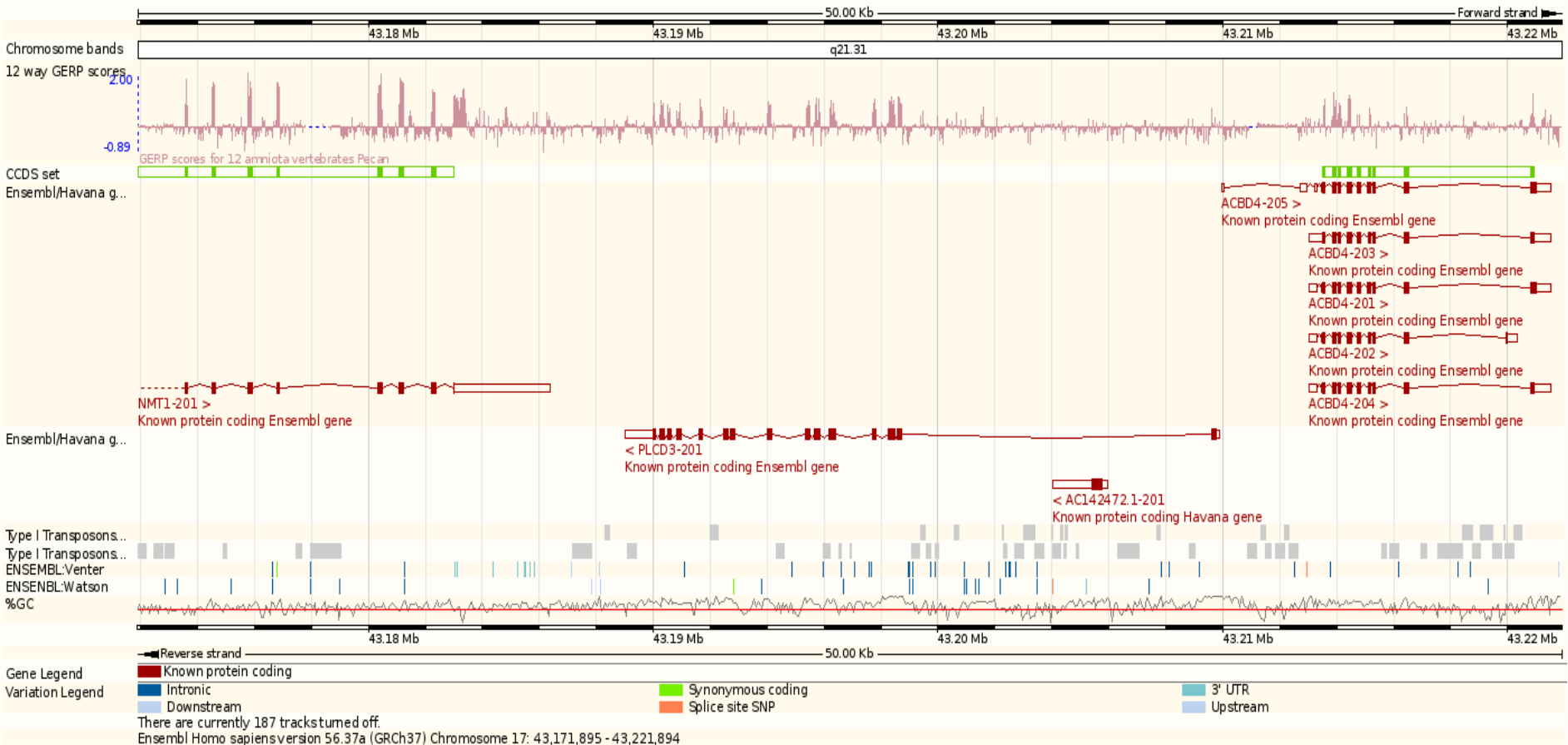
# Taxonomic versus Genomic trees



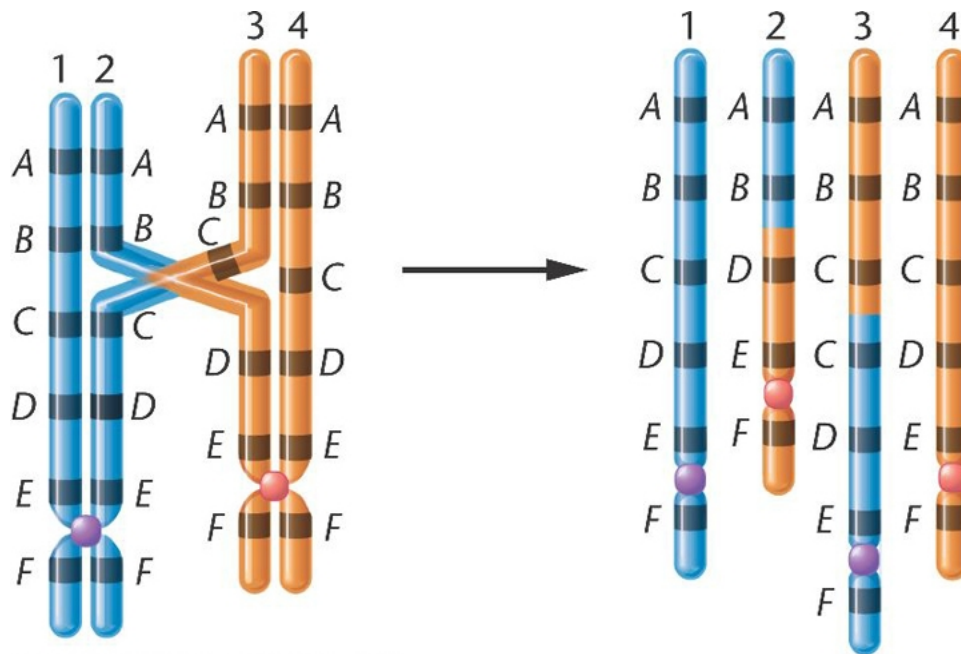Genomic trees may be obtained using alignment algorithms. They are impressively similar to the taxonomic trees!!

This is a highly non trivial test of Evolution theory.
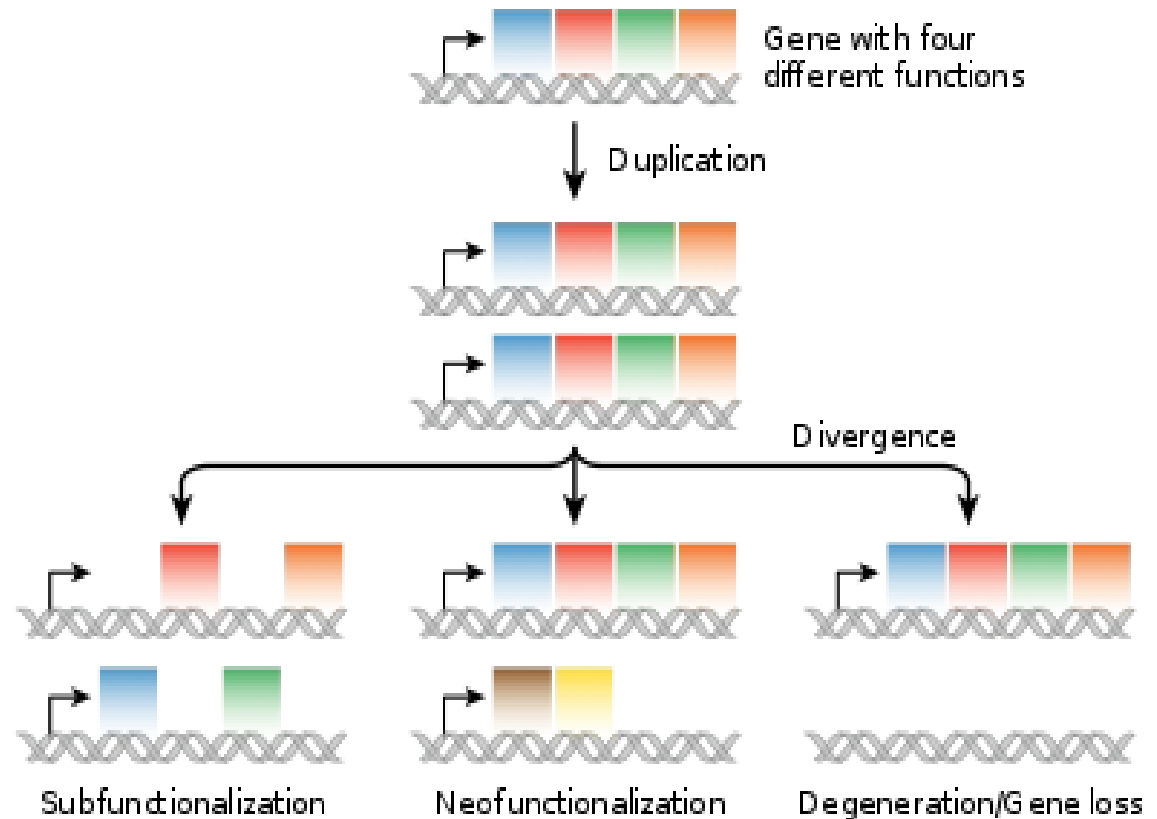
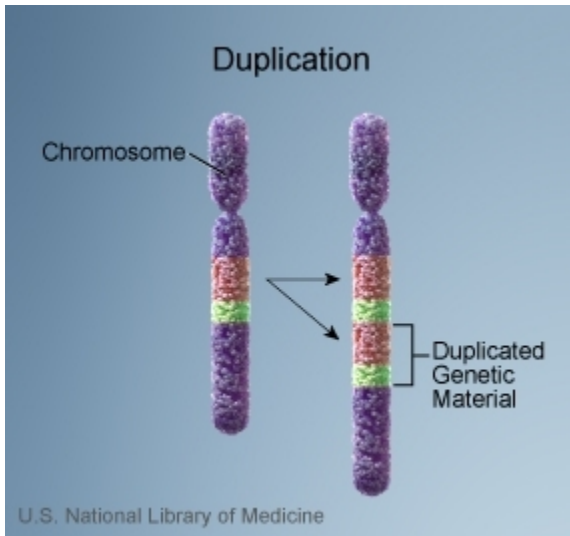# Coding and regulatory regions are conserved!

# Gene duplication via unequal crossingover
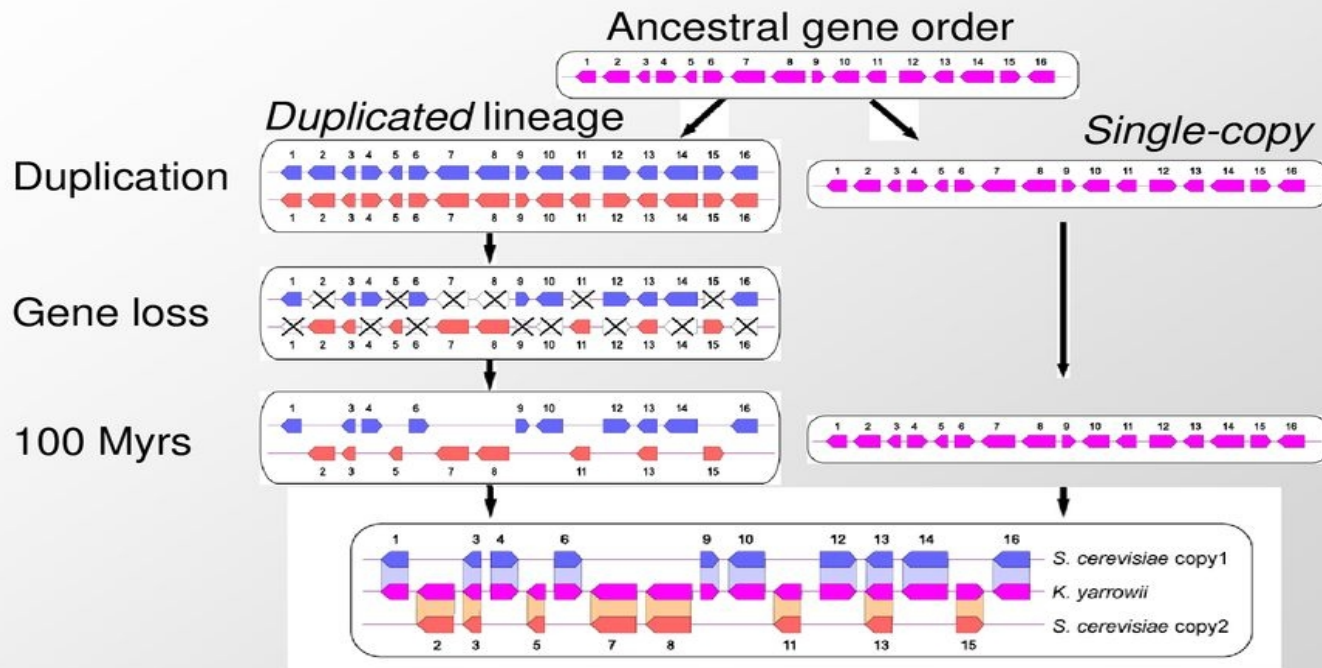


This tetrad is mispaired at meiotic synapsis.

The result, after crossing over, is two unequal chromosomes: one with a duplication (3) and one with a deletion (2).

# Gene duplication: moving in the phenotype space

# Whole Genome duplication: jumping in the phenotype space!
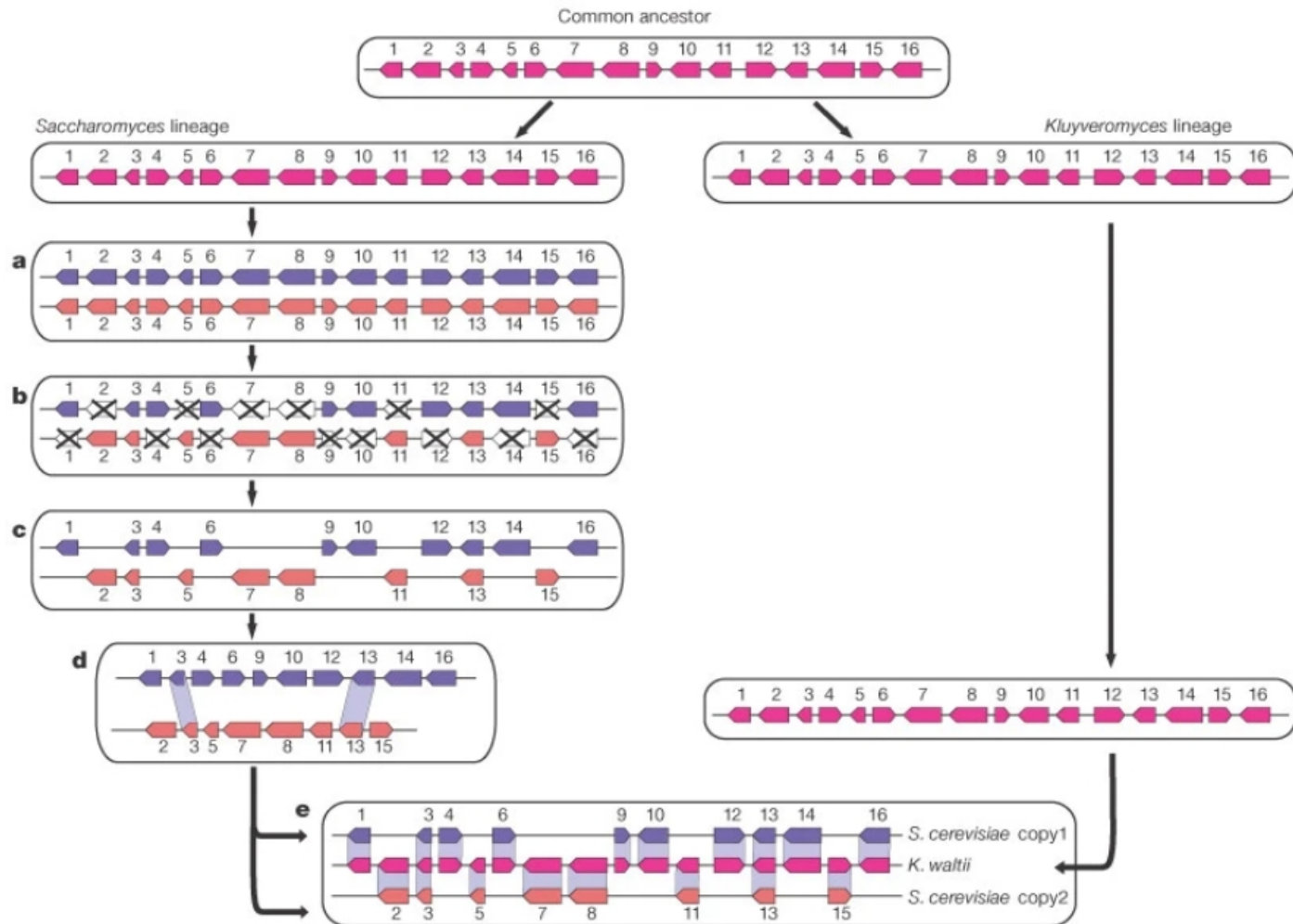


**Evolution by whole-genome duplication**

Ancestral gene order

*Duplicated* lineage

*Single-copy*

Duplication

Gene loss

100 Myrs

S. cerevisiae copy1
K. yarrowii
S. cerevisiae copy2

**Yeast Genome Duplication**
Kellis *et al*. <u>Nature</u>, Apr 8, 2004

**Vertebrate Genome Duplication**
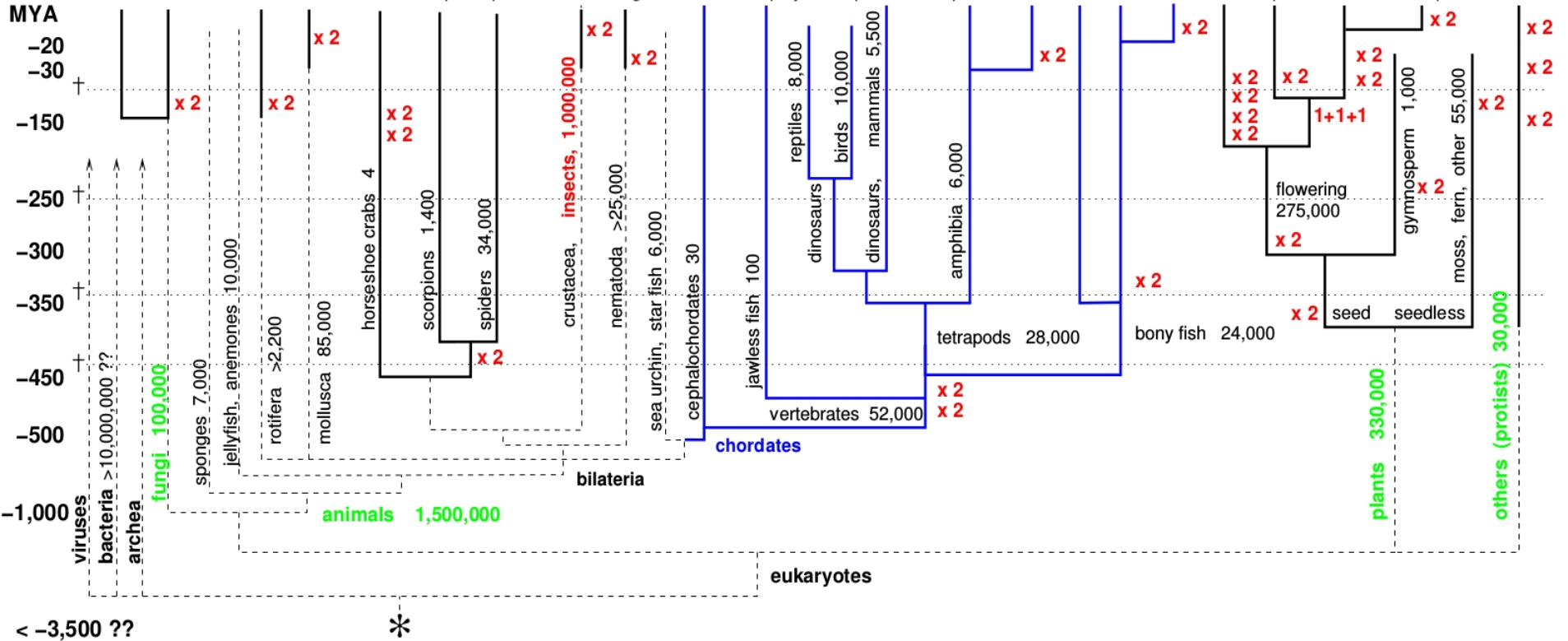Jaillon *et al*. <u>Nature</u>, Oct 21, 2004

# Whole Genome duplication: the yeast case

# Two rounds of WGD at the beginning of the vertebrate lineage!

# **Two examples of applications**

§ Human/Chimp Comparison

§ Evolution of the protein-protein interaction network: the copying model

# Evolution and gene regulation

- **Genomic conservation can be used as an indicator of the functional importance of a given sequence !**

The so called "Ultraconserved regions" have been protecetd by mutations for hundreds of millions of years and they are almost the same in all vertebrates.
Most likely they have a crucial functional role !

Most of them are involved in **gene regulation.** Their identification is the first step toward the construction of a regulatory voculabory in higher eukaryotes.

# Humans and Chimps



96% of the human genome coincides with the chimp's one! Most of the differences are random non-coding  SNPs!

# Humans and Chimps

Big Problem: Human and Chimps are too similar from a genomic point of view!!

- **Idea: use evolutionary conservation in the other way around:** **Look for a region ultraconserved in all vertebrates but mutated in human!**

**One of these regions is found within the coding sequence of  the FOXP2 gene**

# FOXP2 !!



FIGURE 25.27. The *FOXP2* gene is the first gene identified that carries a mutation that causes a specific language deficit in humans. The silent and replacement nucleotide substitutions in this gene as mapped on a primate phylogeny are shown. (*Red bars*) Amino acid changes; (*blue tick marks*) nucleotide changes. Data suggest that the *FOXP2* gene has been the target of selection during recent human evolution after the separation of the human lineage from the common ancestor with the chimpanzee. Numbers show how many nonsynonymous/synonymous changes have occurred along each branch.

25.27, adapted from Enard W. et al., *Nature* **418:** 869–872, © 2002 Macmillan, www.nature.com

*Evolution* © 2007 Cold Spring Harbor Laboratory Press

# FOXP2 !!

Mutations (SNPs) in the FOXP2 gene are associated to
deep alterations in speaking ability.

The human version of FOXP2 could be at the origin of the ability of H. Sapiens to articulate words and thus organize complex chasing strategies...

# Modern Genomics: *networks*

- genes and proteins of a given organism are organized in networks .

- Cells react to external stimuli in a "global" way.

H.Jeong et al.
Nature, 411 (2001) 41

# Interaction Networks

- All the interaction networks in biological systems are "**heterogeneous**", with a "fat-tailed" connettivity distribution: few hubs and several peripheral links.

- Standard explanations such as "**preferential attachment**" models used to describe WWW topology cannot work for biological systems.

- Evolutionary models based on gene duplication, the so called "**copying models**" lead to a biologically rooted explanation which fits the available data very well

# Gene Regulation

Gene expression is tightly regulated. All cells in the body carry the full set of genes, but only express about 20% of them at any particular time. Different proteins are expressed in different cells (neurons, muscle cells....) according to the different functions of the cell.

Among the various regulatory steps the most important ones are:

▪ transcriptional control,
  by **Transcription Factors.**

▪ post-transcriptional control,
  by **microRNAs.**



Alberts, *Molecular Biology of the Cell*

# Transcription Factors and miRNAs

**Transcription Factors (TFs)**: proteins binding to specific recognition **motifs (TFBSs)** usually short (5-10 bp) and located **upstream** of the coding region of the regulated gene.



*Wassermann*, Nat. Rev. Genetics

**MicroRNAs (miRNAs)** are a family of small RNAs (typically **21 - 25** nucleotide long) that **negatively regulate gene expression at the posttranscriptional level**, (usually) thanks to the "seed" region in 3'-UTR regions.

# Transcription Factors

# MicroRNAs

# Regulatory Networks 1

**Key 1 -->** **TFs** are themselves proteins produced by other genes, and they act in a combinatorial way, resulting in a complex network of interactions between genes and their products.
**--> Transcriptional Network**

**miRNAs** also act in a combinatorial and one-to-many way, and, moreover, <u>are transcribed from same POL-II promotes of TFs</u>.
**--> Post-Transcriptional Network**

# Regulatory Networks 2

**Key 2 -->** Biological functions are performed by groups of genes which act in an interdependent and synergic way. A complex network can be divided into simpler, distinct regulatory patterns called **network motifs**, typically composed by 3 or 4 interacting components which are able to perform elementary signal processing functions.

# *Network motifs*

Example: SIM (Single Input Module)  (a) experimental realization: arginine biosynthesis   b) Circuit behaviour: different genes are activated at different times as a function of their different activation threshold as the concentration  of X (master regulator) changes in time   R.Milo et al. Science 298 (2002) 824



Nature Reviews | Genetics

# Network Motifs II

Network motifs can be studied using standard tools of theoretical physics:
- Ordinary differential equations
- Stochastic equations
- Montecarlo (Gillespie) simulations.

- Goal: understand the functional role of the motif and why it was selected by evolution

- Example 1: incoherent feedforward loops can reduce the noise in the amount of produced proteins.

A) FFL

B) Open circuit

Time (seconds)

FFL

Open circuit

C) Probability Density

Protein Number

# Cancer Driver Genes

Cancer is the result of a pathological alteration of the regulatory network induced by somatic mutations (to be distinguished from the germinal mutations) and chromosomal alterations (Copy Number Variations)

Main outcome of the recent genomic studies: each tumour is unique!
→ new therapeutic approach: "Personalized Medicine / Precision Medicine"

Can we identify on purely computational basis the drivers of this disregulation process?

Problem: in a typical cancer cell we find thousands of altered genes. Can we identify the real drivers?

# Somatic mutation frequencies observed in exomes from 3,083 tumour–normal pairs.

nature

# Cancer driver genes

Goal: integrate different sources of regulatory information using Network Theory to identify driver genes in cancer

**TOOLS:**

- Multiplex Theory
- Mutual Information (ARACNE)
- Filtering Methods ( Disparity Filter)
- Community detection algorithms (Infomap, OSLOM, Label propagation, Louvain and Modularity optimization via simulated annealing)
- Consensus Clustering

# Proposal: Multi-Network Integration of Regulatory Information

**Cantini L., Medico E., Fortunato S. and Caselle M.**,
*"Detection of gene communities in multi-networks reveals cancer drivers."*
Scientific Reports (2015) **5**: 17386

# Proposal: Multi-Network Integration of Regulatory Information



L1: expression

L2: miRNA

L3: TF

L4: PPI

Multi-network communities reconstruction:

1. Community detection within each network layer
2. Consensus clustering across the four layers.

Open source Community detection algorithms:

- Infomap,
- OSLOM,
- Label propagation,
- Louvain
- Modularity optimization via simulated annealing.

→ The rationale behind this choice is that gene coexpression and protein-protein interactions require a tight coregulation of the partners and that such a fine tuned regulation can be obtained only combining both the transcriptional and post-transcriptional layers of regulation. .

→ Our procedure is valid in principle for any pathology but is particularly suited for cancer, due to role that disregulation plays in cancer. We studied  in particular  gastric, lung, pancreas and colorectal cancer

→ To extract the relevant biological information we constructed for each tumor two multiplex networks: one using expression data for the normal tissue and one for the tumor and then compared their partition into *communities.*

A **community** is a group of nodes that are densely connected to each other, but sparsely connected to the other nodes of the network.

# **Results: Chromosomal Locations**

Three major outcomes:

1) Out of the hundreds of genes contained in each enriched chromosomal location, with the Multiplex recosntruction  *only the few genes which are involved in a common co-regulatory scheme are selected.* and thus are likely to be the real drivers of the cancer.
2) In the communities one find also *genes outside the enriched chromosomal locus, related to them not only by a coexpression link but also by regulatory relations* and this suggests that they could be part of a common biological pathway which is disregulated in the tumour.

3) *In some cases the community is also characterized by a GO or KEGG enriched category* which may give some hint to identify the above pathway.

# Enriched Chromosomal locations

**Locus 1q32** enriched in 43rd community of Pancreatic Cancer Data
- Genes in the intersection: ATF3, BTG2, CD46, IRF6 and PPP1R15B.
- 43rd community also enriched in DREAM pathway

# Conclusions:
## a set of open questions

- Which is the role of **Non coding DNA**?

- Which is the genomic origin of the difference between **Human and Chimps**?

- Which is the genomic origin of the difference between **different human beings**?

- Can we identify the genomic source of the impressive **complexity of multicellular organisms**?

- Can we identify the disregulated pathways leading to **complex deseases, in particular cancer**?