# *Introduction to Comparative Genomics*
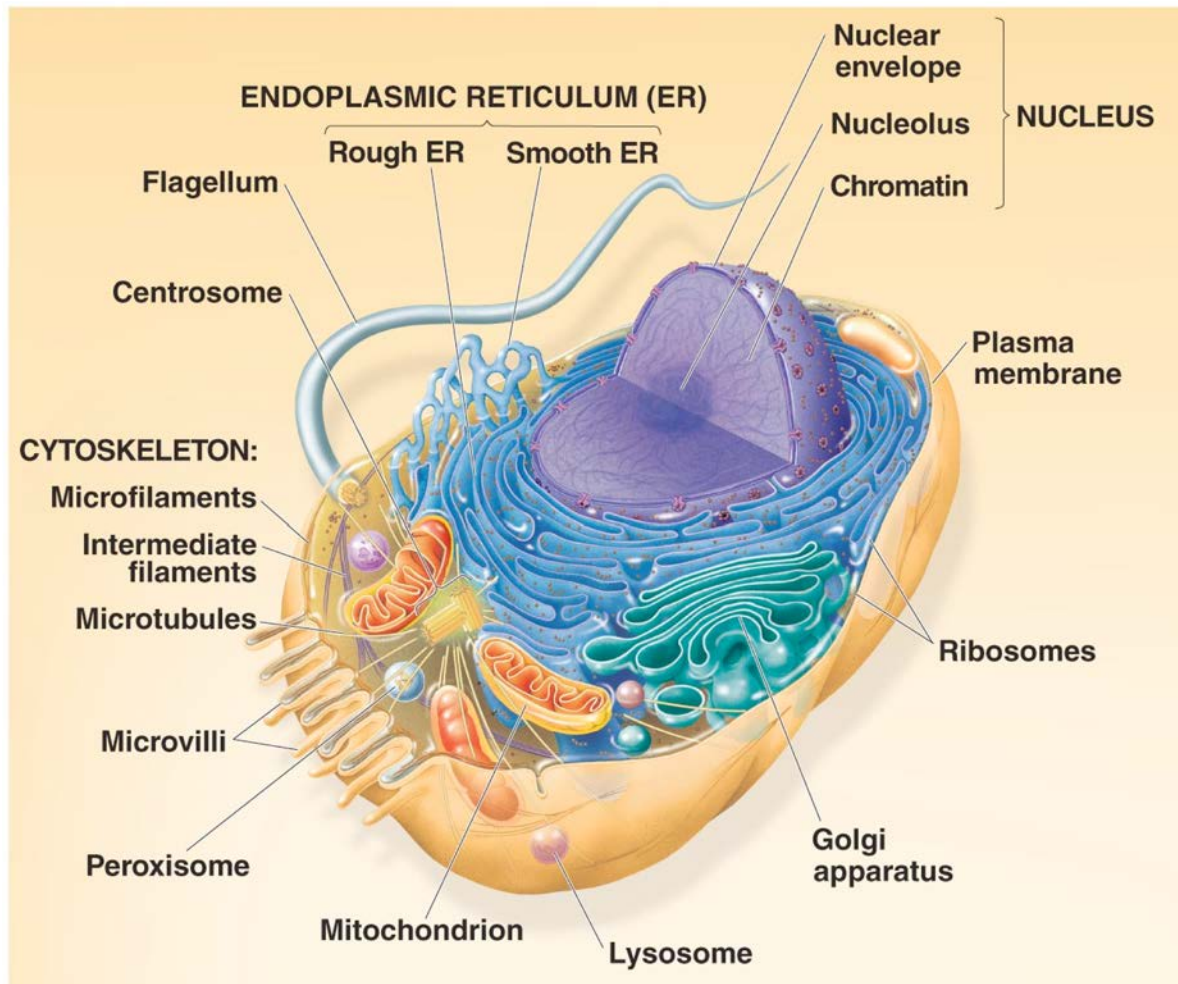
*Evgeny.Zdobnov@unige.ch*

UNIVERSITÉ
DE GENÈVE

FACULTÉ DE MÉDECINE

# Cell Elements



Copyright © 2008 Pearson Education, Inc., publishing as Pearson Benjamin Cummings.

**Cell elemental composition**

Cells are 90% water.

The remaining is approximately:

- 50% **protein**
- 15% **carbohydrate**
- 15% **nucleic acid**
- 10% **lipid**
- 10% **miscellaneous**

# Cell Elements

- Proteins are the main cellular machinery
- All proteins – proteome
- All DNA – genome
- All RNA – transcriptome
- All lipids – lipidome

# Terms

- -omics     ⇔     high <u>throughput</u> data acquisition
      in Molecular <u>Biology</u>

- Bioinformatics ⇔ computational <u>management</u>
      and <u>analysis</u> of biological data

# *Why Genomics?*

# Genome encodes hereditary information

# The dogma

# DNA/RNA sequencing is far ahead

# DNA, chromatin, chromosomes

# *Sequencing cost is decreasing and data are being accumulated fast*



sequencing has been industrialized

# Sequencing «Generations»

# Sequencers: read length and output

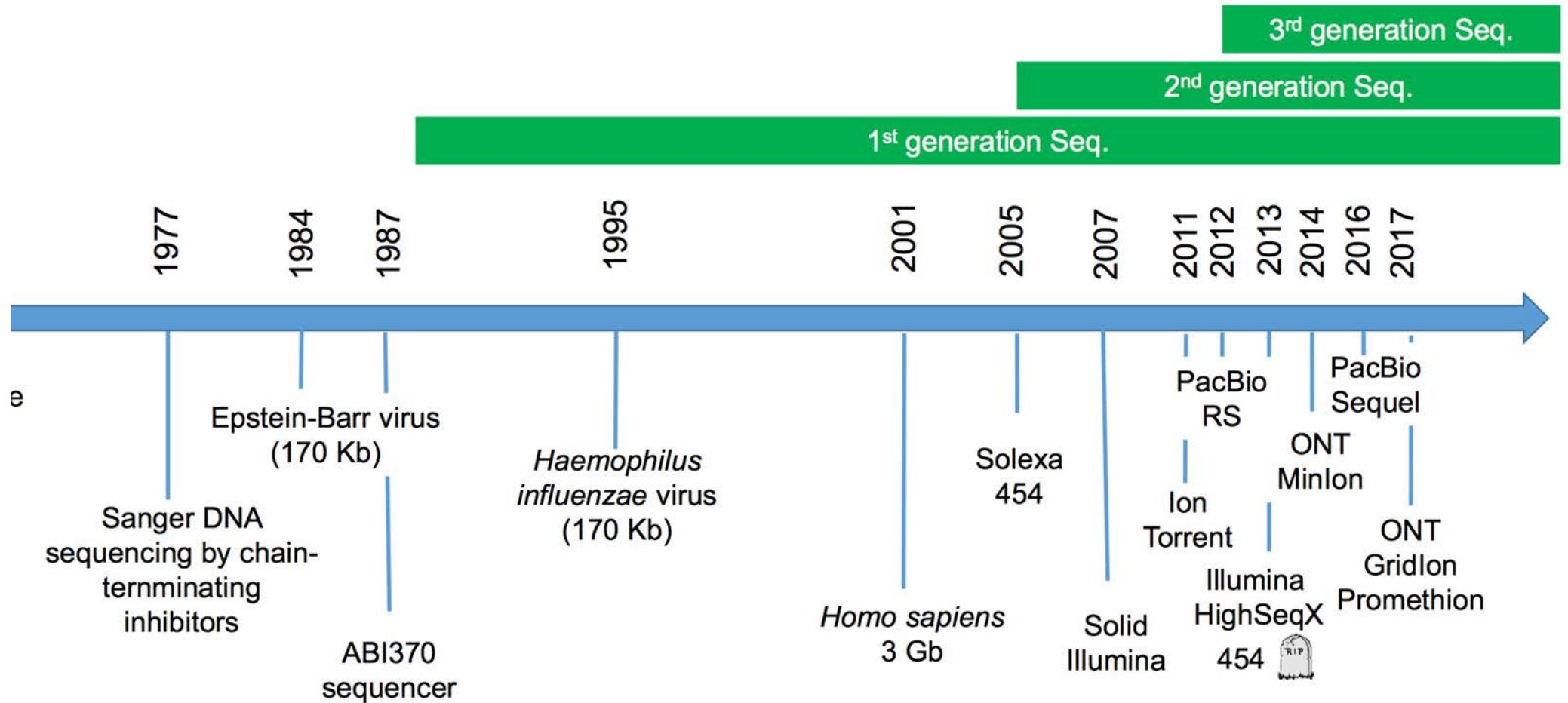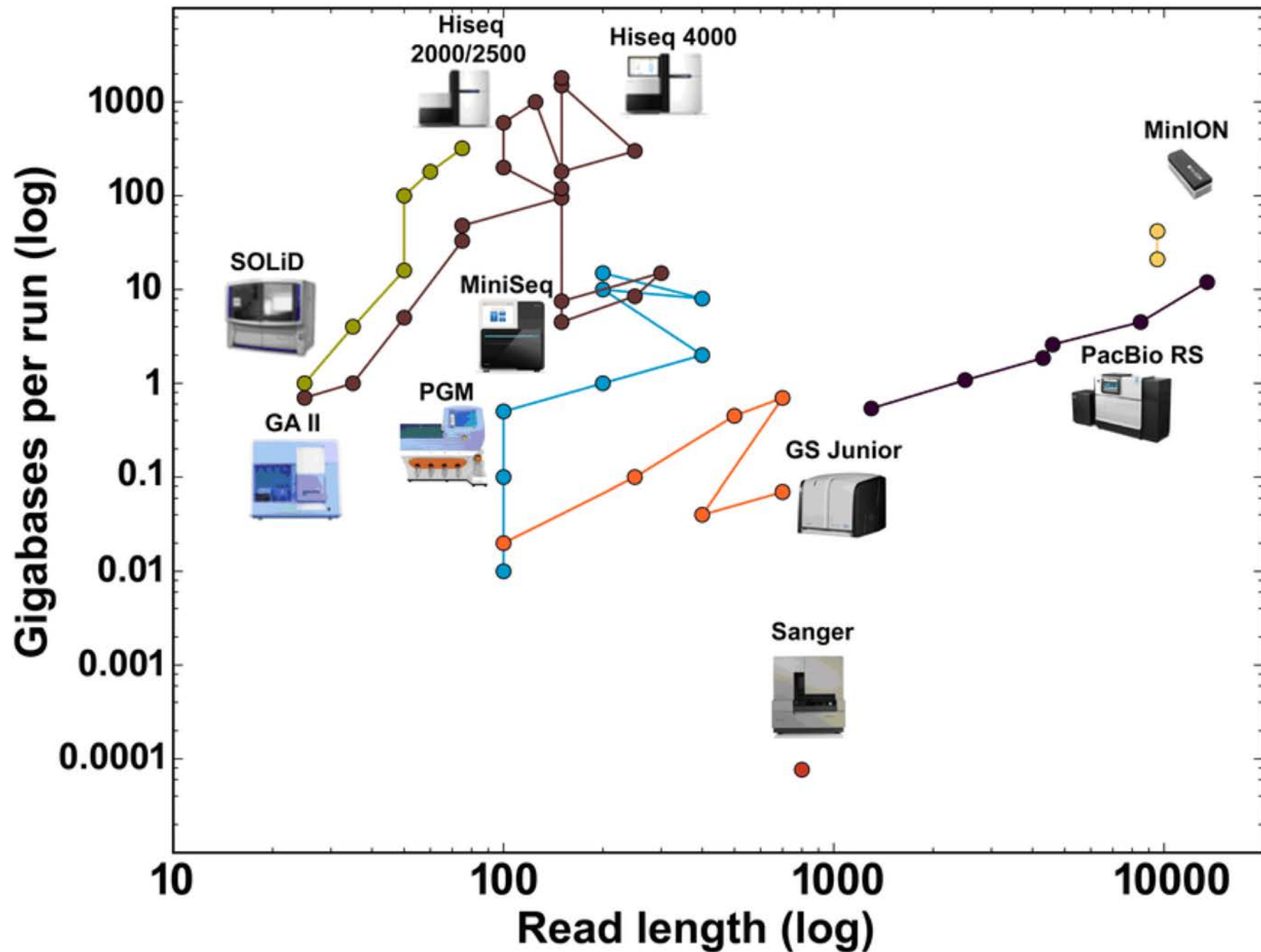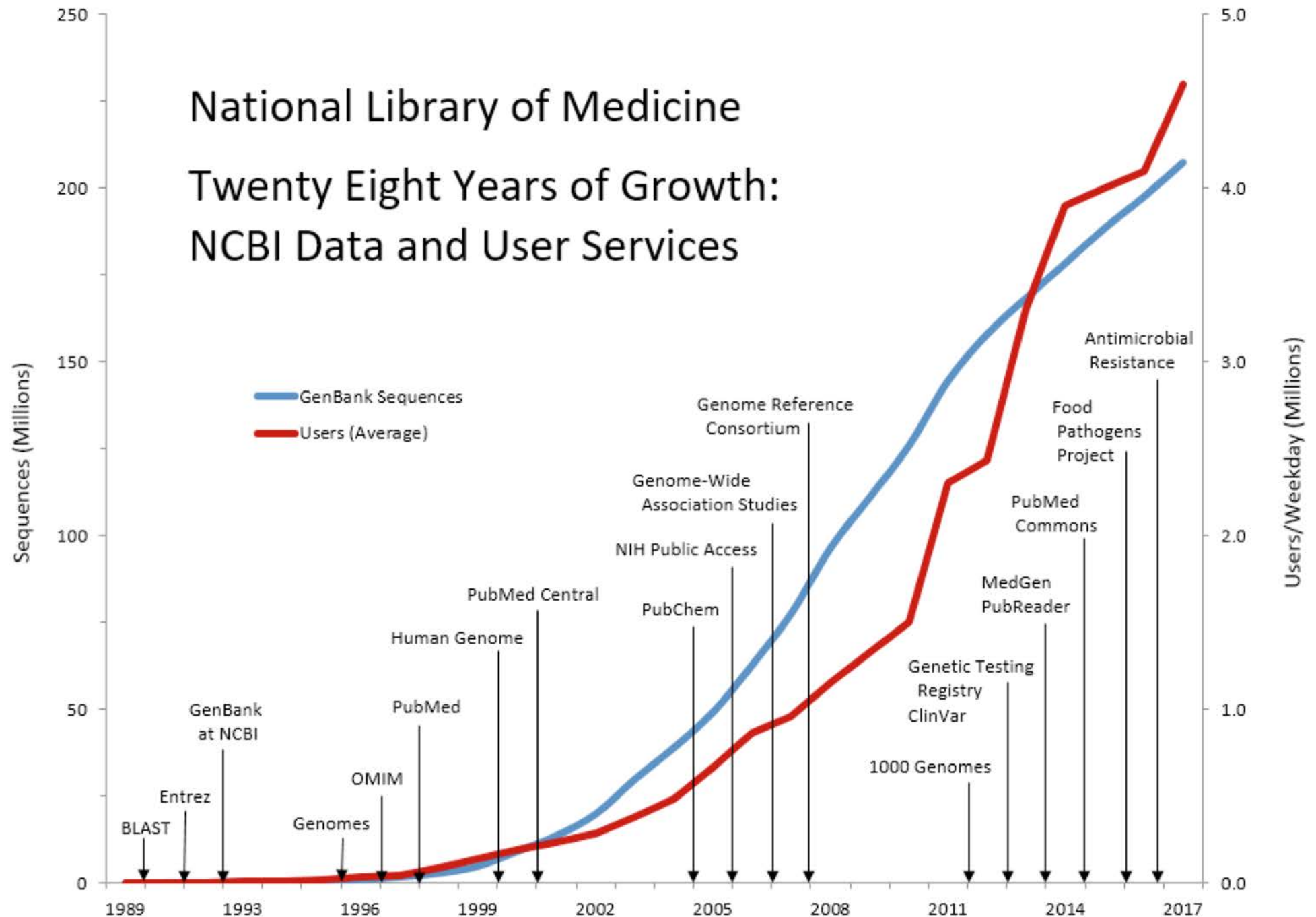National Library of Medicine

Twenty Eight Years of Growth:
NCBI Data and User Services

*Without interpretation*
*(by comparisons)*
*DNA is unintelligible*

*=> sequence analysis required!*

UNIVERSITÉ
DE GENÈVE
FACULTÉ DE MÉDECINE

14

# *Some assembly required*

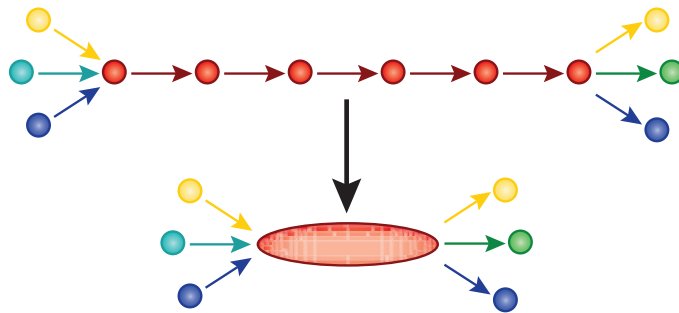1. Fragment DNA and sequence

2. Find overlaps between reads

…AGCCTAGACCTACAGGATGCGCGACACGT
                  GGATGCGCGACACGTCGCATATCCGGT..

3. Assemble overlaps into contigs

4. Assemble contigs into scaffolds

Genome assembly stitches together a genome
from short sequenced pieces of DNA.

UNIVERSITÉ
DE GENÈVE
FACULTÉ DE MÉDECINE

15

# Genomics is unthinkable without computer data analysis



just one genome

our computer

computers can only execute human intelligence

UNIVERSITÉ DE GENÈVE
FACULTÉ DE MÉDECINE

# *The promise: i.e. why we are here*

**Musings**

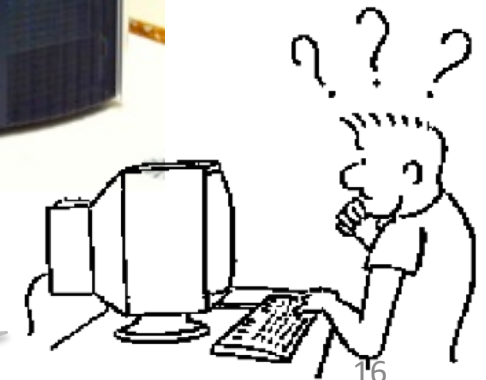**Highly accessed**

## The $1,000 genome, the $100,000 analysis?

**Elaine R Mardis**
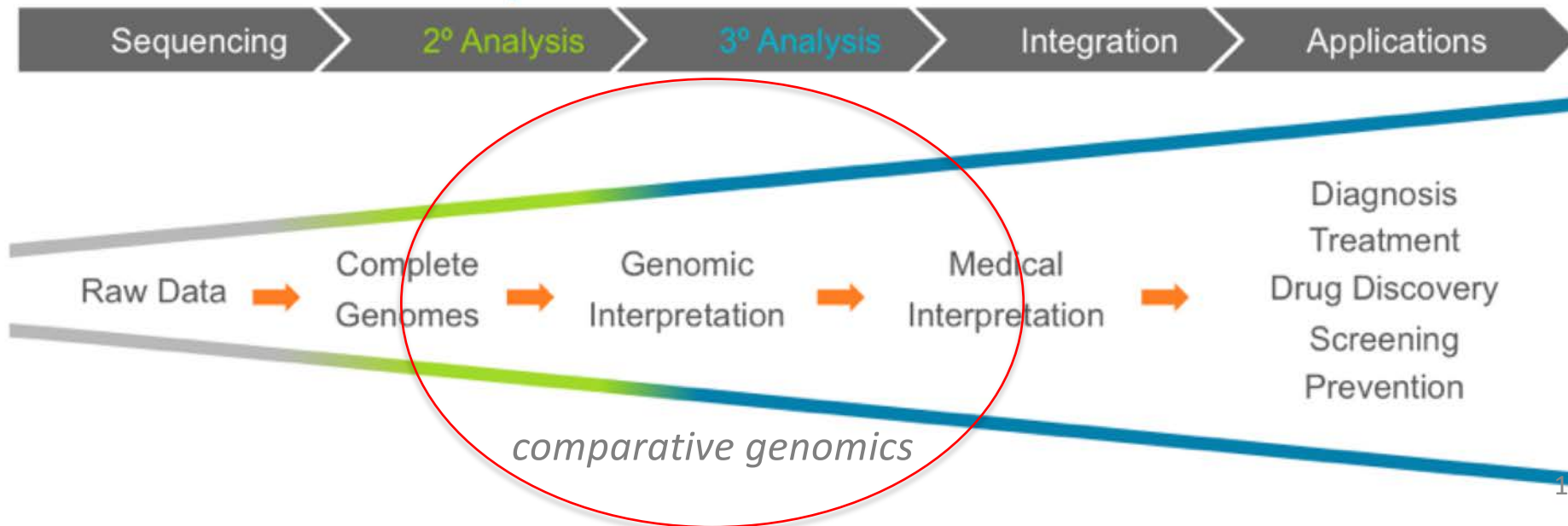
Correspondence: Elaine R Mardis emardis@wustl.edu  ⊻ Author Affiliations

The Genome Center at Washington University School of Medicine, 4444 Forest Park Blvd, St Louis, MO 63108, USA

*Genome Medicine* 2010, **2**:84  doi:10.1186/gm205

## Genomics Landscape      for future medicine

| Sequencing | 2° Analysis | 3° Analysis | Integration | Applications |

Raw Data ➡ Complete Genomes ➡ Genomic Interpretation ➡ Medical Interpretation ➡ Diagnosis / Treatment / Drug Discovery / Screening / Prevention

*comparative genomics*

17

# Genome sizes

# Genomics "Holy Grail": predicting phenotypic variability from genetic variability

# *Gene expression*

- Proxy to cell functions (via proteins)
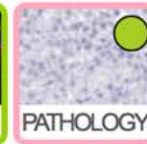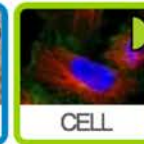- Not all genes expressed
- Highly uneven expression levels

# Gene expression



www.proteinatlas.org

# Not only genome can be sequenced

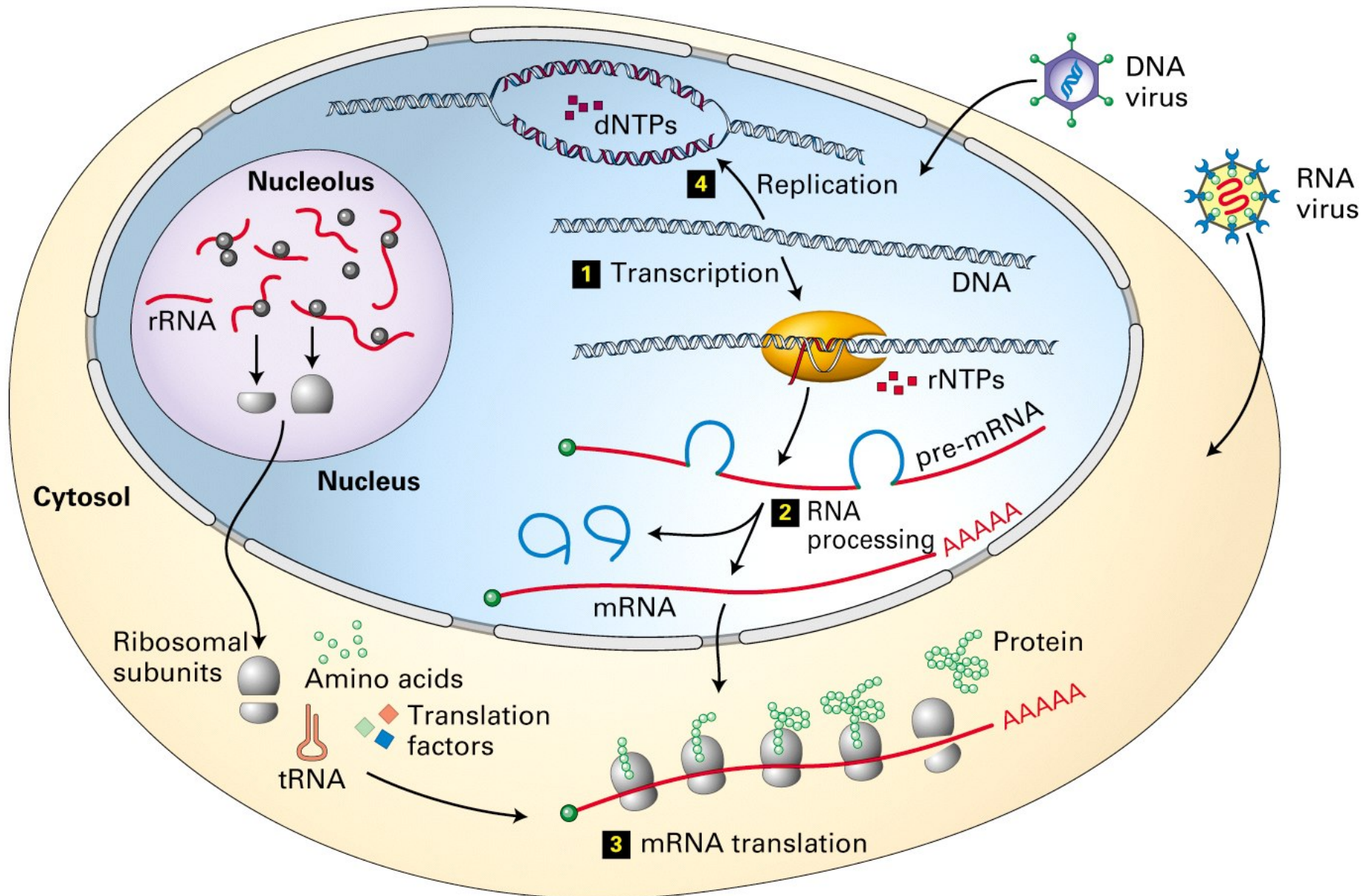# Genomics keywords

- **DNA-Seq** is sequencing DNA in the sample

- **RNA-Seq** is sequencing RNA in the sample

- **ChIP-Seq** is sequencing DNA sites interacting with specific protein

- **Meta-genomics** is sequencing many organisms in one sample

# Why genomics?

+ "Complete" cellular DNA/RNA snapshot,
+ Protein/NA & NA/NA interactions
+ relative abundance of "reads"
+ **wealth of data**


- it doesn't tell you about biology;
  proteins, interactions, metabolites,  etc.;
  not even which sequences are meaningful
  and which not.

# Biological systems



- **Ecosystem**
- **Population**
- Organism
- Organ
- Tissue
- Cell
- Complexes/networks
- Molecules

# Metagenomics:
# direct sequencing of total DNA/RNA

a mix of intestinal bacteria
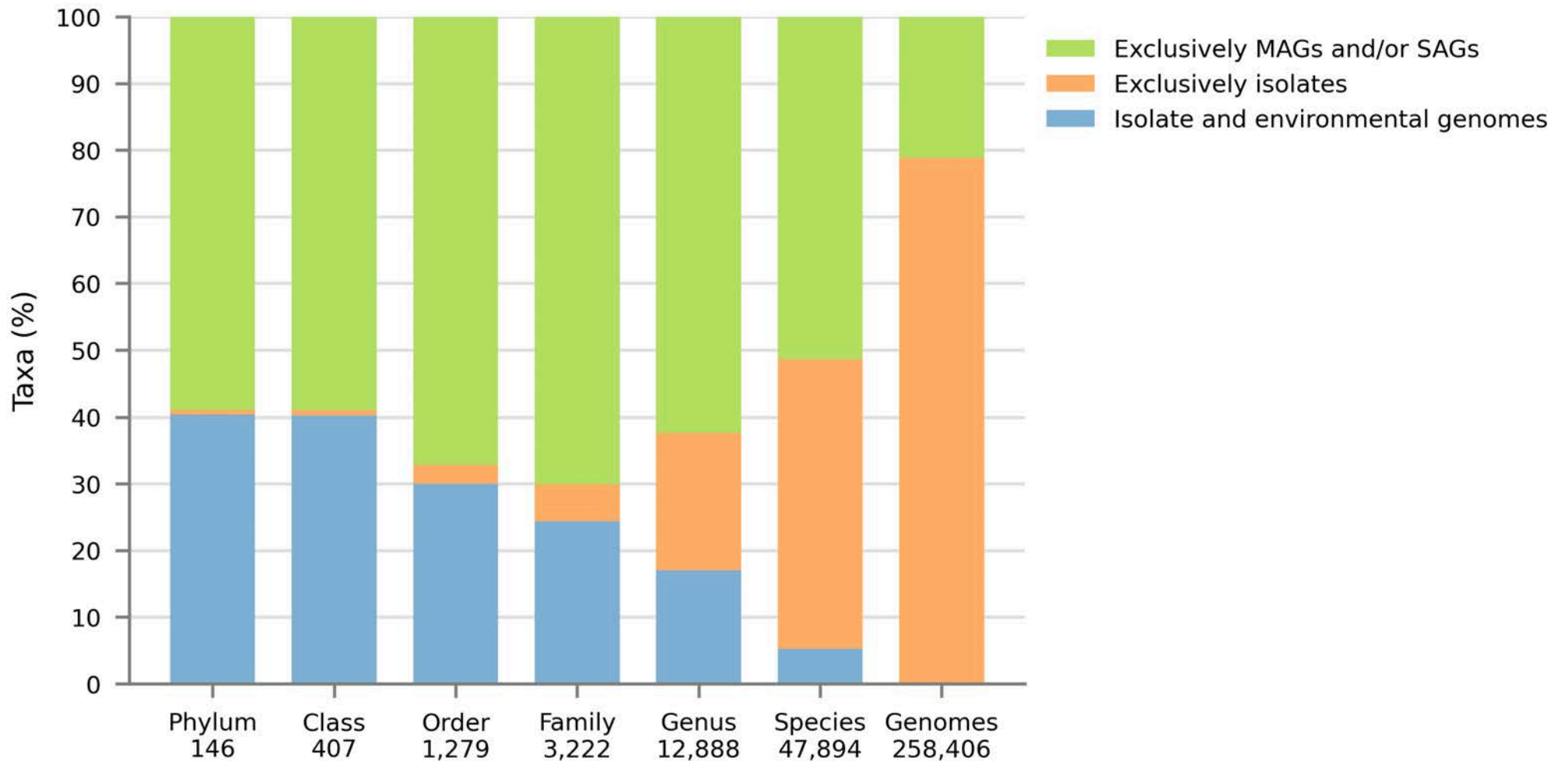
can be sequenced without culturing



*www.micronaut.ch*

# Scaling-up and mixing the puzzles

- Who is there?
  *or* What they can do?
- How many?

# Bacterial and archaeal genomes



*gtdb.ecogenomic.org/stats*

# Recent trends in genomics

- cancer / clinical human variations
- metagenomics
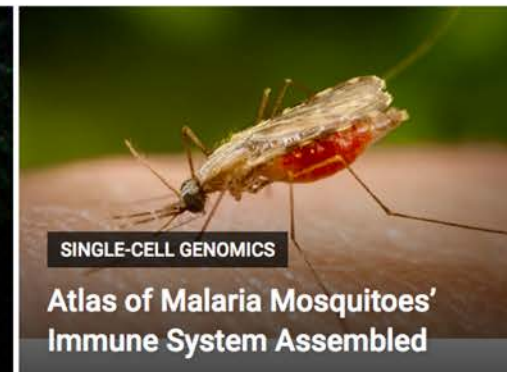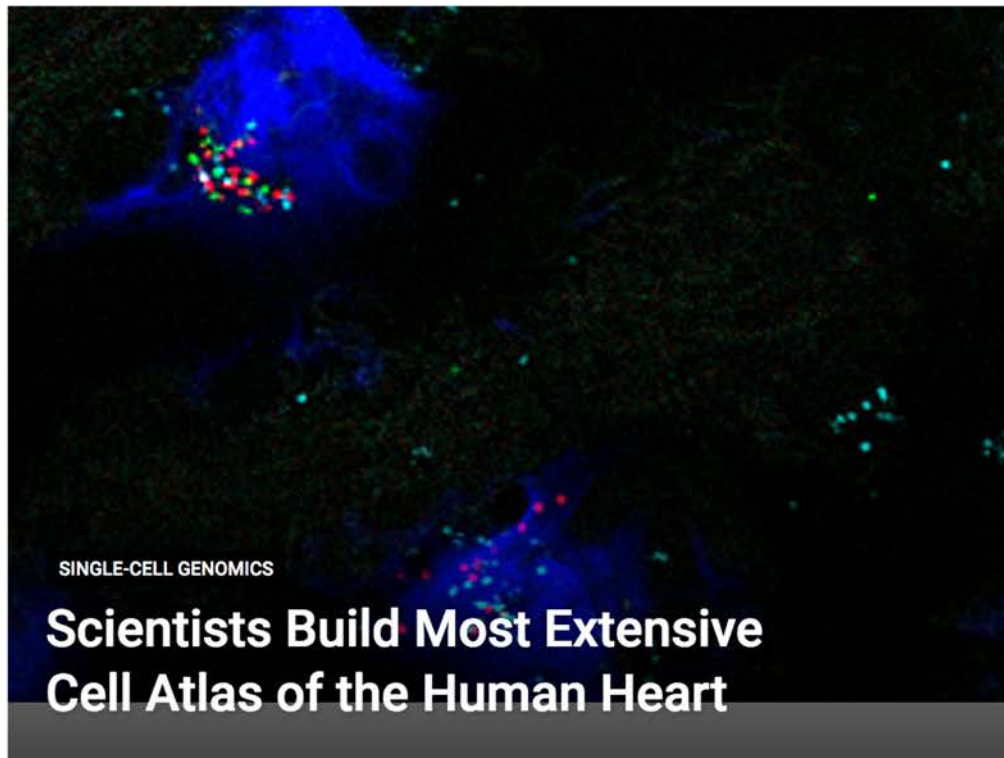- single-cell and spatial transcriptomics

# Human microbiomes

# Earth Microbiome Project: Mapping the microbiome of... everything

by University of California - San Diego

# SINGLE-CELL GENOMICS

SINGLE-CELL GENOMICS

**Scientists Build Most Extensive Cell Atlas of the Human Heart**

SINGLE-CELL GENOMICS

**Atlas of Malaria Mosquitoes' Immune System Assembled**

SINGLE-CELL GENOMICS

**Single-Cell Company Berkeley Lights Sets Terms for $153M IPO**

SINGLE-CELL GENOMICS

**Single-Cell RNA-Seq Reveals Progenitor Glioblastoma Stem Cell**

SINGLE-CELL GENOMICS

**New Single-Cell Technique Reveals Genetic Diversity of Cancer Tumors**

UNIVERSITÉ DE GENÈVE

FACULTÉ DE MÉDECINE

**a**

Unspliced mRNA → Spliced mRNA

Ratio → RNA velocity

**b**

Differentiated cell types

Early progenitors

©nature

# NCBI SARS-CoV-2 Resources

## Quick Navigation Guide

Sequence Submission

Literature

Sequence-Related Resources

Clinical Resources

Other Websites

## SARS-CoV-2 Data

**2,973,768**
**SRA runs**

**3,646,037**
**Nucleotide records**

**3,215**
**ClinicalTrials.gov**

**223,652**
**PubMed**

**275,714**
**PMC**

UNIVERSITÉ
DE GENÈVE
FACULTÉ DE MÉDECINE

# *Comparative approaches*

# It all started with microbes

Research requires

*a tool to see, and*

*bookkeeping the comparisons*

THE INFECTIOUS WORLD OF GERMS AND MICROBES

BSL-4 — BSL-4
BSL-3 — BSL-3
COLD — EBOLA — FLU — RABIES

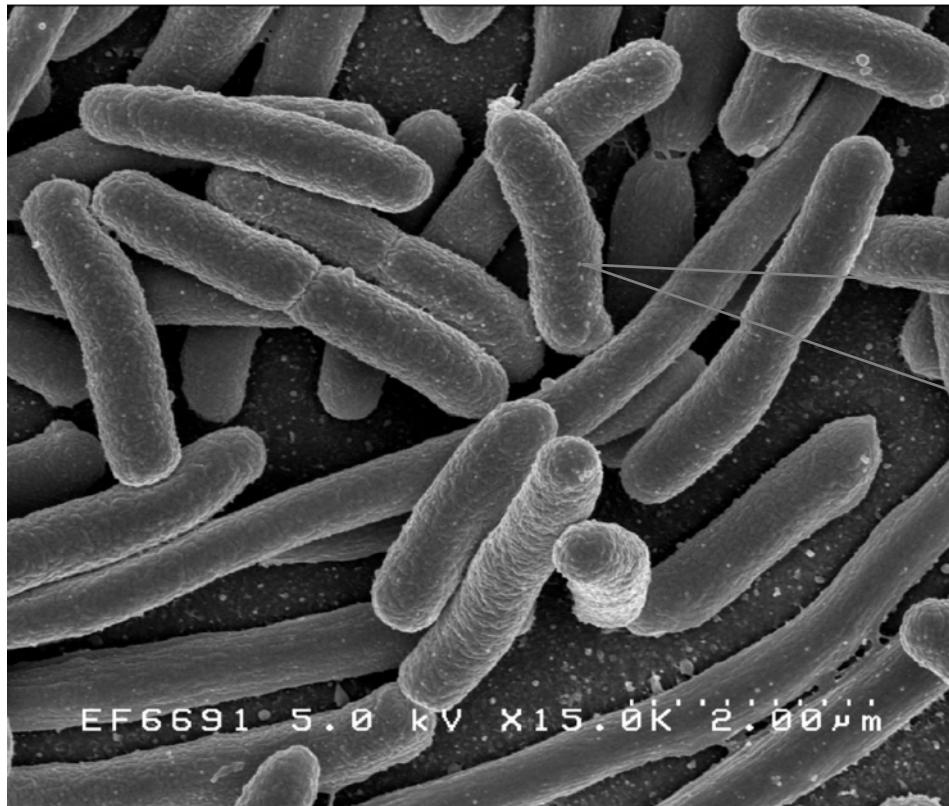BY JENNIFER GARDY · ILLUSTRATED BY JOSH HOLINATY

## and model organisms

UNIVERSITÉ DE GENÈVE
FACULTÉ DE MÉDECINE

# *Looking same but some are pathogenic, requiring molecular-level investigations*



Three strains
of *E. coli*
can have only
40% genes
in common..

Welch, R.A. (2002). Extensive mosaic structure revealed by the complete genome sequence of uropathogenic Escherichia coli. Proceedings of the National Academy of Sciences, 99(26), 17020-17024.

# *Comparative genomics is about comparing the genomic features of different organisms.*

*An example*

Initial impact of the sequencing of the human genome

**Eric S. Lander**

*Nature* **470**, 187–197 (10 February 2011) | doi:10.1038/nature09792

The sequence of the human genome has dramatically accelerated biomedical research.

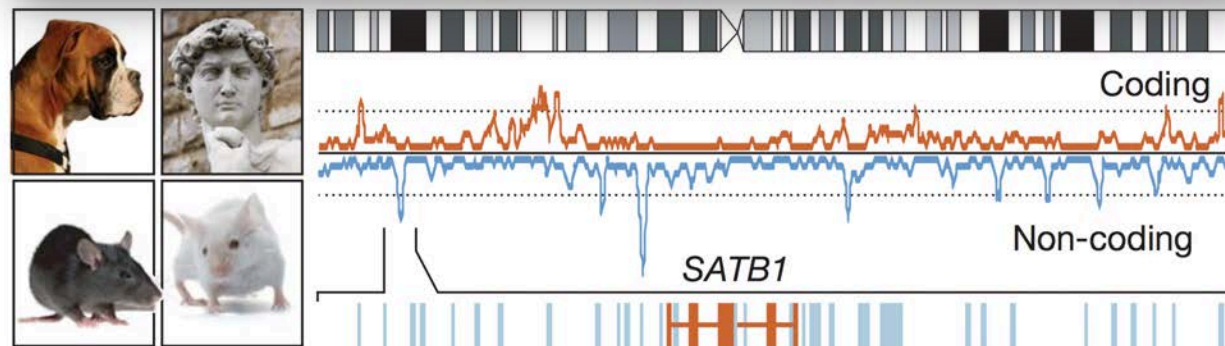Coding

Non-coding

SATB1

**Figure 1 | Evolutionary conservation maps.** Comparison among the human, mouse, rat and dog genomes helps identify functional elements in the genome.
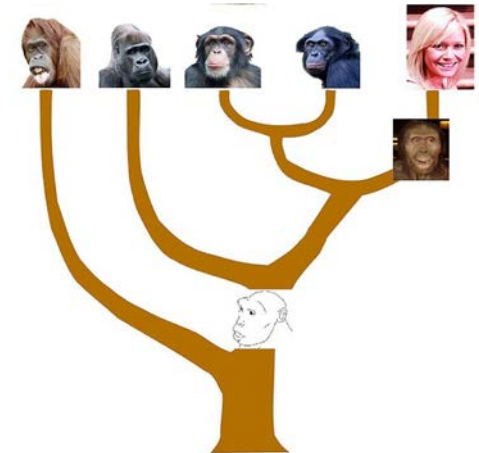
FACULTÉ DE MÉDECINE

*The aims are to:*

- *Enable knowledge transfer, e.g. from models to human*

- *Interpret Nature's molecular experimentation*

38

# How genomes evolve

- Accumulation of mutations ⇔ divergence

- Vertical descent by speciation



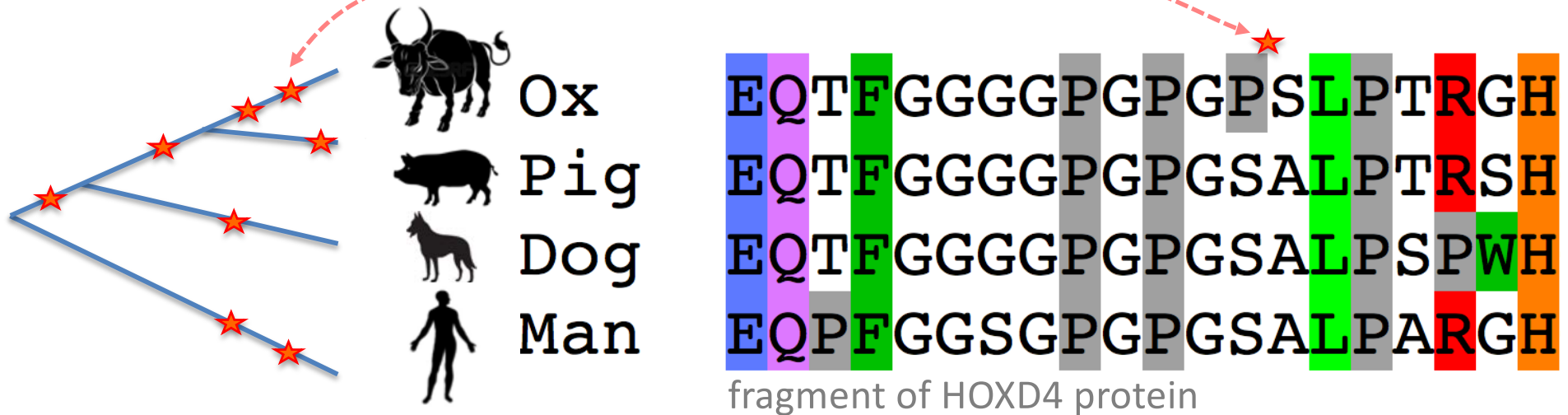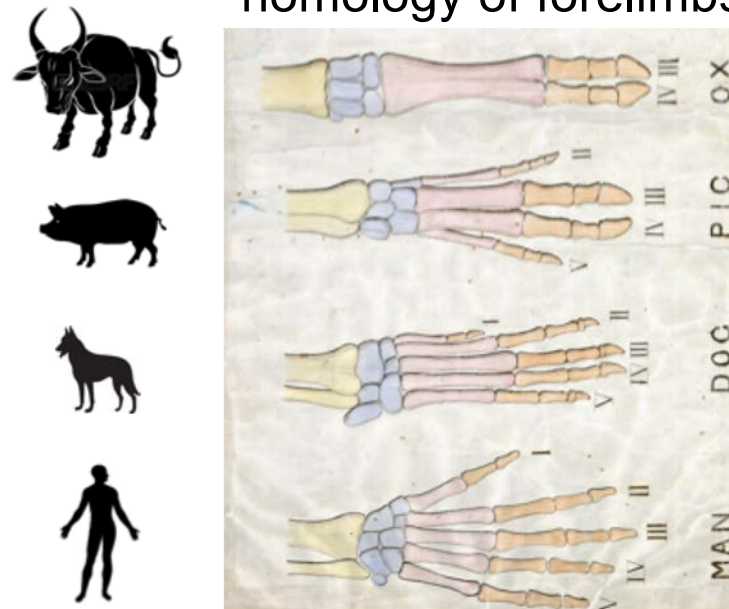UNIVERSITÉ
DE GENÈVE
FACULTÉ DE MÉDECINE

# *Inheritance of sequence and function*

mutations happen



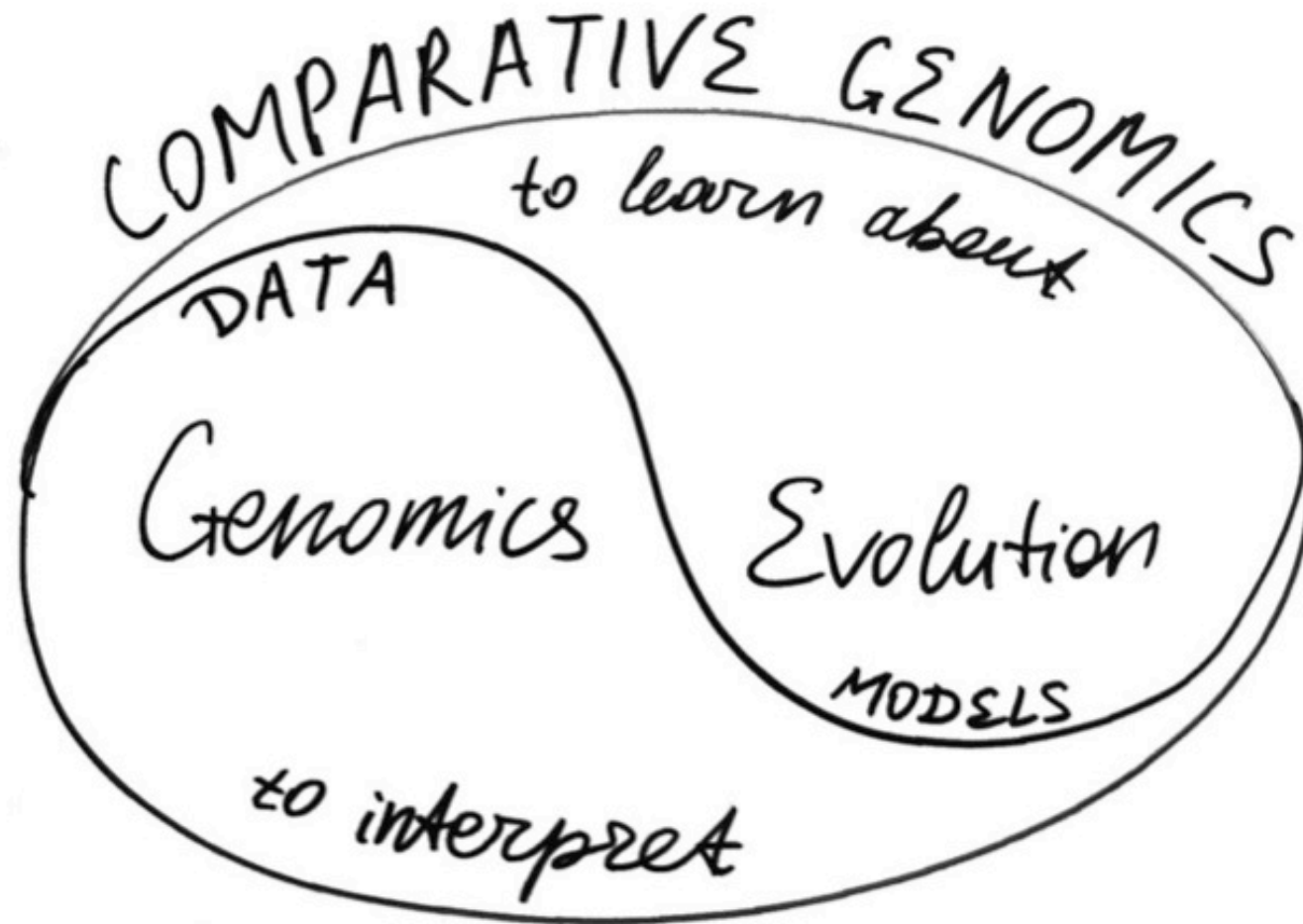| | | |
|---|---|---|
| Ox | EQTFGGGGPGPGPSLPTRGH | |
| Pig | EQTFGGGGPGPGSALPTRSH | |
| Dog | EQTFGGGGPGPGSALPSPWH | |
| Man | EQPFGGSGPGPGSALPARGH | |

fragment of HOXD4 protein

time

homology of forelimbs



*functional selection
accepts or rejects mutations*

# How to get there:
## employing knowledge to interpret genomes and using genomes to further our knowledge
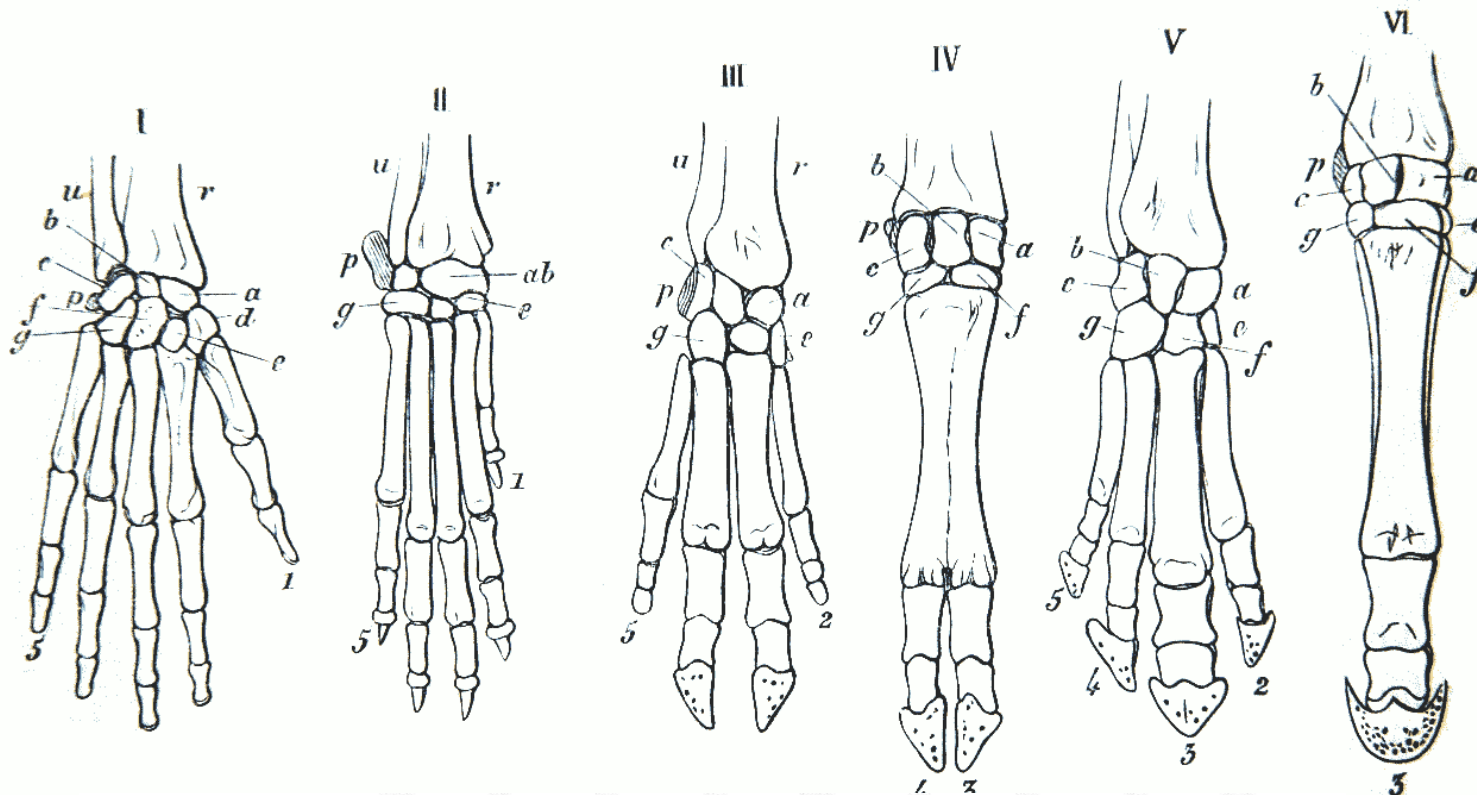
# *General aims*

- **Similarities** allow to transfer our knowledge
  from well studied model organisms
  to the newly sequenced ones

- **Differences** may shed light on unique
  species adaptation processes

UNIVERSITÉ
DE GENÈVE
FACULTÉ DE MÉDECINE

- Similarity

  *vs*

- Homology

  *vs*

- Orthology

# How would you compare?



44

**Homology**, in biology, similarity of the structure, physiology, or development of different species of organisms based upon their descent from a common evolutionary ancestor. ...

www.britannica.com › science › homology-evolution

# *Sequence alignment*
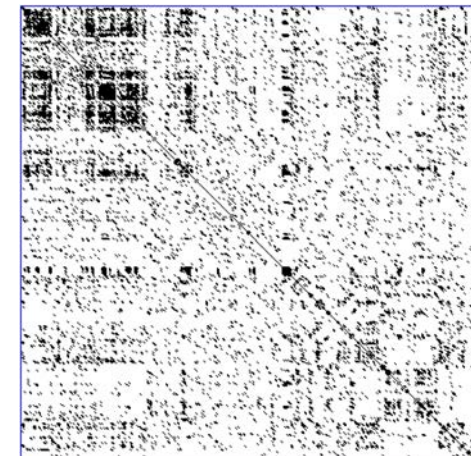
```
AAB24882   TYHMCQFHCRYVNNHSGEKLYECNERSKAFSCPSHLQCHKRRQIGEKTHEHNQCGKAFPT 60
AAB24881   -------------------YECNQCGKAFAQHSSLKCHYRTHIGEKPYECNQCGKAFSK 40
                              ****: .***:   * *:** * :****.:* *******..

AAB24882   PSHLQYHERTHTGEKPYECHQCGQAFKKCSLLQRHKRTHTGEKPYE-CNQCGKAFAQ- 116
AAB24881   HSHLQCHKRTHTGEKPYECNQCGKAFSQHGLLQRHKRTHTGEKPYMNVINMVKPLHNS 98
              ****  *:**********:***:**.:  .****************    :  *.: :
```

\* - identical

: - conserved substitutions (same colour group)

. - semi-conserved substitution (similar shapes).

```
Global   FTFTALILLAVAV
         F--TAL-LLA-AV


Local    FTFTALILL-AVAV
         --FTAL-LLAAV--
```
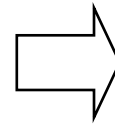
46

# *The significance of similarity scores*

Extreme-value distribution of <u>random sequence</u> alignment scores



**E-value (expected value):** the expected number of *random hits* with the same score expected by chance in *this* database.

# Seq. similarity identification tools:

- SSAHA
- Blat
- BLAST
- Smith-Waterman (Paralign)
- PSI-Blast / RPS-Blast
- CS-Blast
- HHMER
- HHsearch

# *Orthology definition*

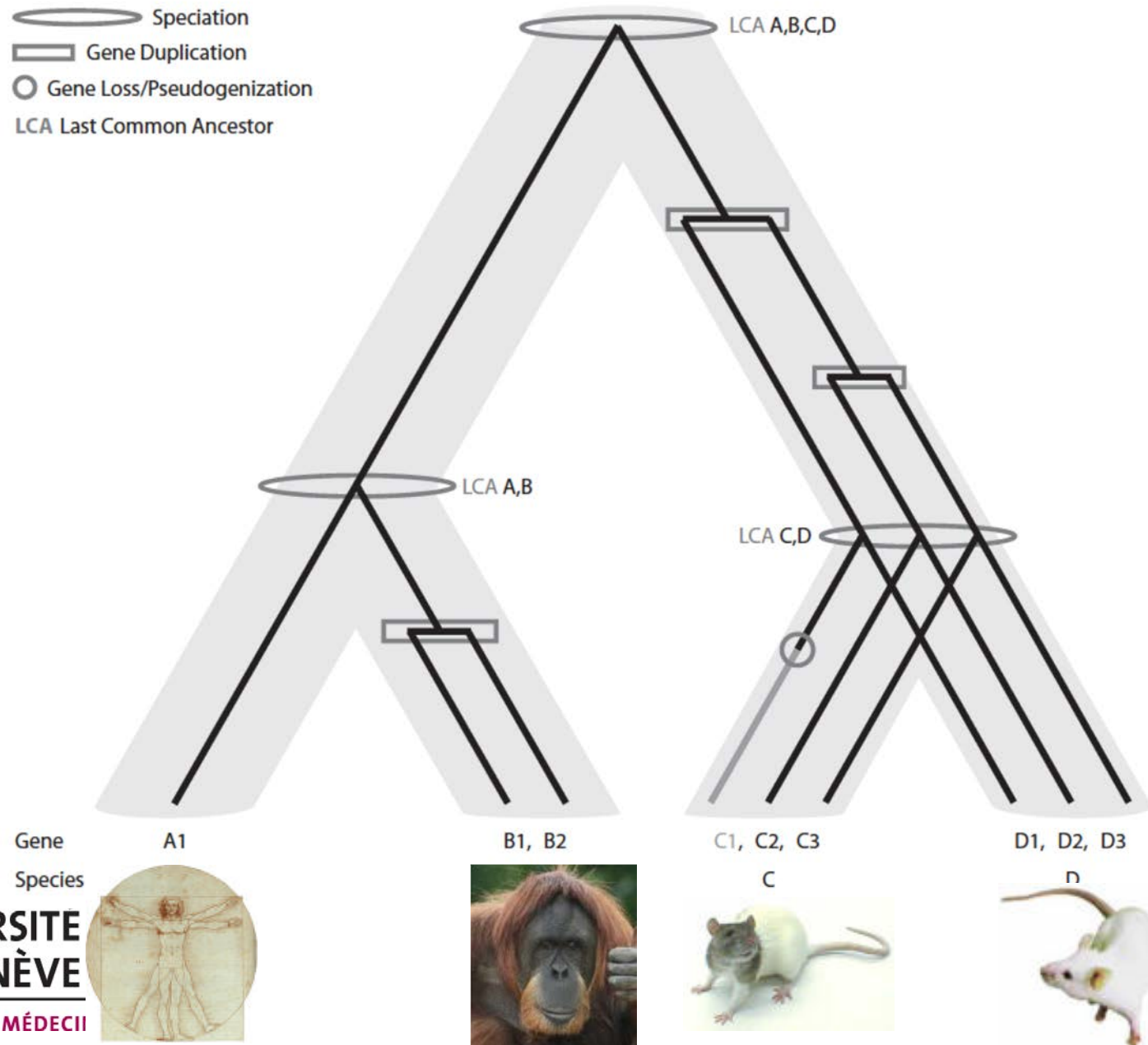Originally the term was introduced in 1970 by Walter Fitch

*Two homologous genes in two different species that derive from a single gene in the last common ancestor of the species*

and better rephrased by Koonin in 2005:
*Genes originating from a single ancestral gene in the last common ancestor of the compared genomes.*

UNIVERSITÉ
DE GENÈVE
FACULTÉ DE MÉDECINE

# Orthologs



Speciation
Gene Duplication
Gene Loss/Pseudogenization
LCA Last Common Ancestor

LCA A,B,C,D

LCA A,B

LCA C,D

Gene    A1          B1,  B2        C1,  C2,  C3        D1,  D2,  D3
Species                            C                   D

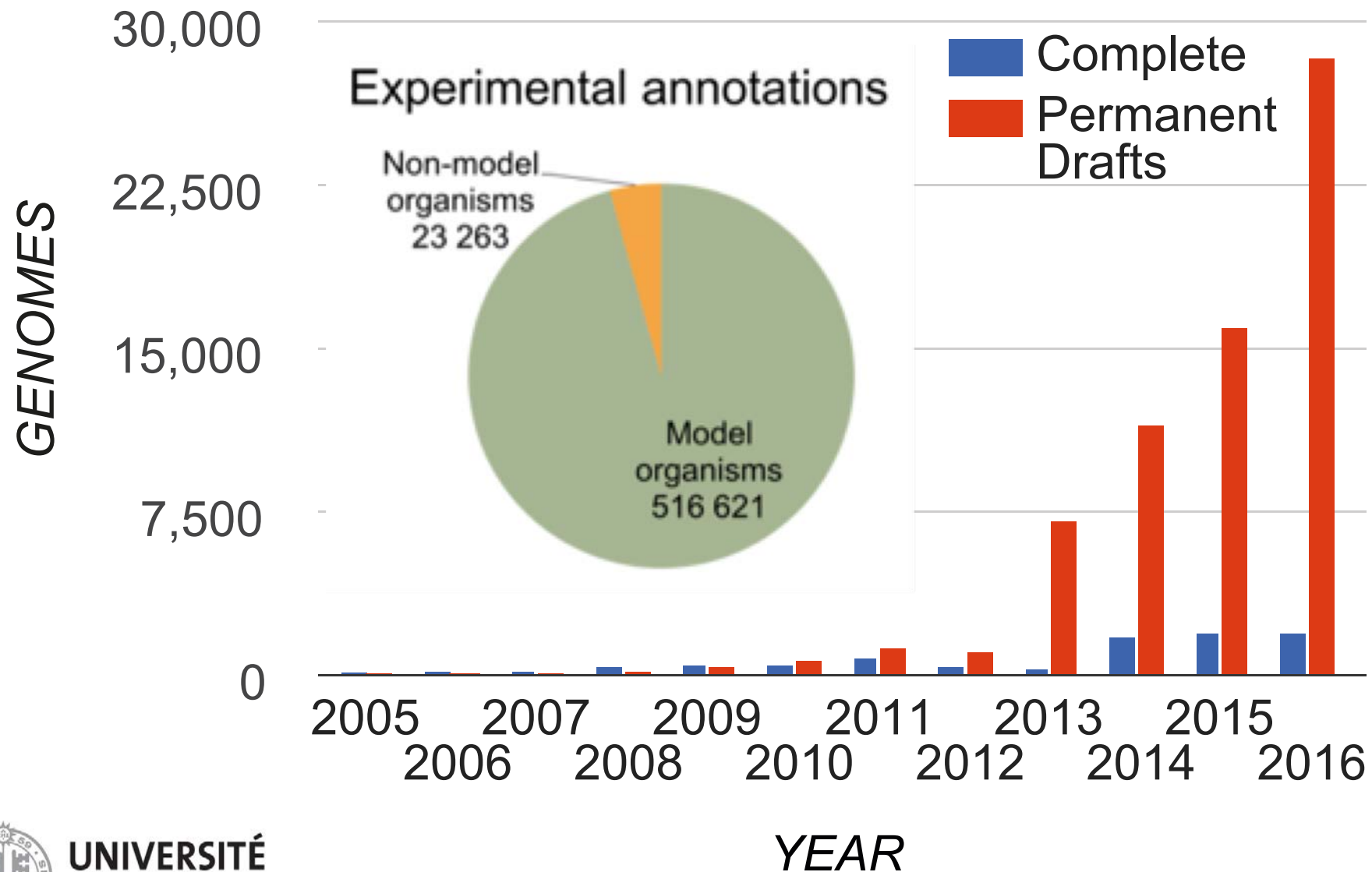UNIVERSITE DE GENÈVE
FACULTÉ DE MÉDECINE

# *Please note*

- Similarity
  *could be*

- Homology
  *could be*

- Orthology

*i.e. all orthologs are homologs and look similar;*
*not all similar looking sequences are homologs,*
*and not all homologs are orthologs;*
*and there is no 'function' in these definitions.*

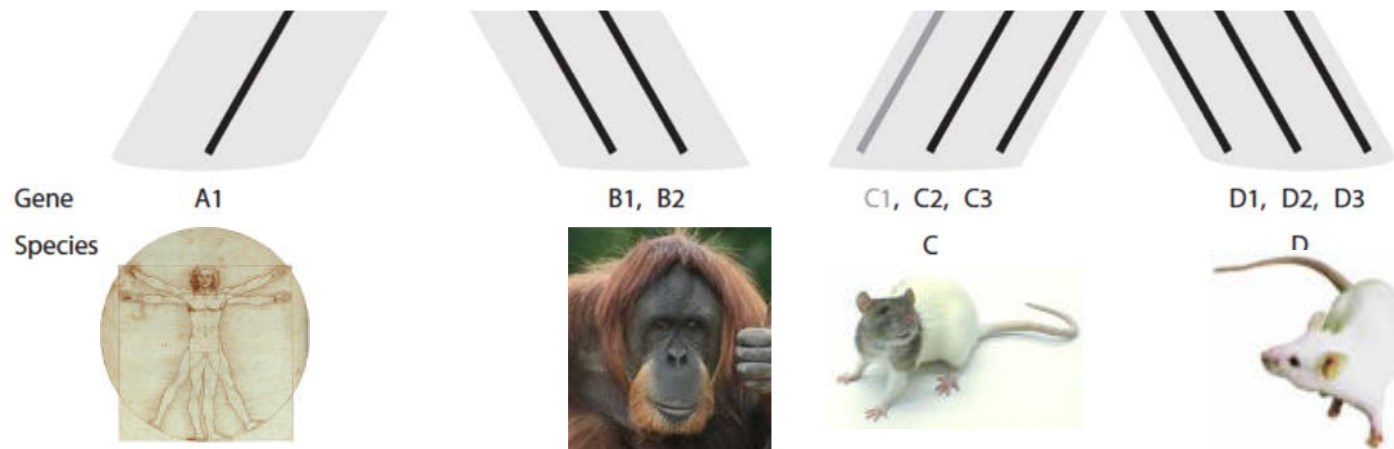The growth of the need: linking genomics data to gene function knowledge

Source: GOLD database

# How to identify orthologs

## ?



Gene      A1            B1, B2      C1, C2, C3      D1, D2, D3

Species                               C          D

# 1. Tree-reconciliation (complicated)

# 2. Best-Reciprocal-Hits (BRH)

- *what are these?*



- *why do they work?*

# #1 Why BRH is indicative of orthology?



Nobody knows the exact history

rat genes

human genes

# #2 lets BLAST
# a rat gene to all human genes



gene duplication

Mammalian ancestor

*the longer evolutionary distance*
*the worse similarity score*

2nd hit

1st hit

UNIVERSITÉ
DE GENÈVE
FACULTÉ DE MÉDECINE

# #3 lets BLAST in reverse
# the best human gene hit to all rat genes



gene duplication

Mammalian ancestor

the longer evolutionary distance
the worse similarity score

2nd hit       1st hit

# #4 there is a reciprocally best matching pair of genes between rat and human

- *BRH joins a pair of genes via a single last-common-ancestor gene*

☞ *orthologs by definition*



gene duplication

Mammalian ancestor

UNIVERSITÉ DE GENÈVE
FACULTÉ DE MÉDECINE

# a real-life BRH graph

☞ *likely orthologs*

UNIVERSITÉ
DE GENÈVE
FACULTÉ DE MÉDECINE

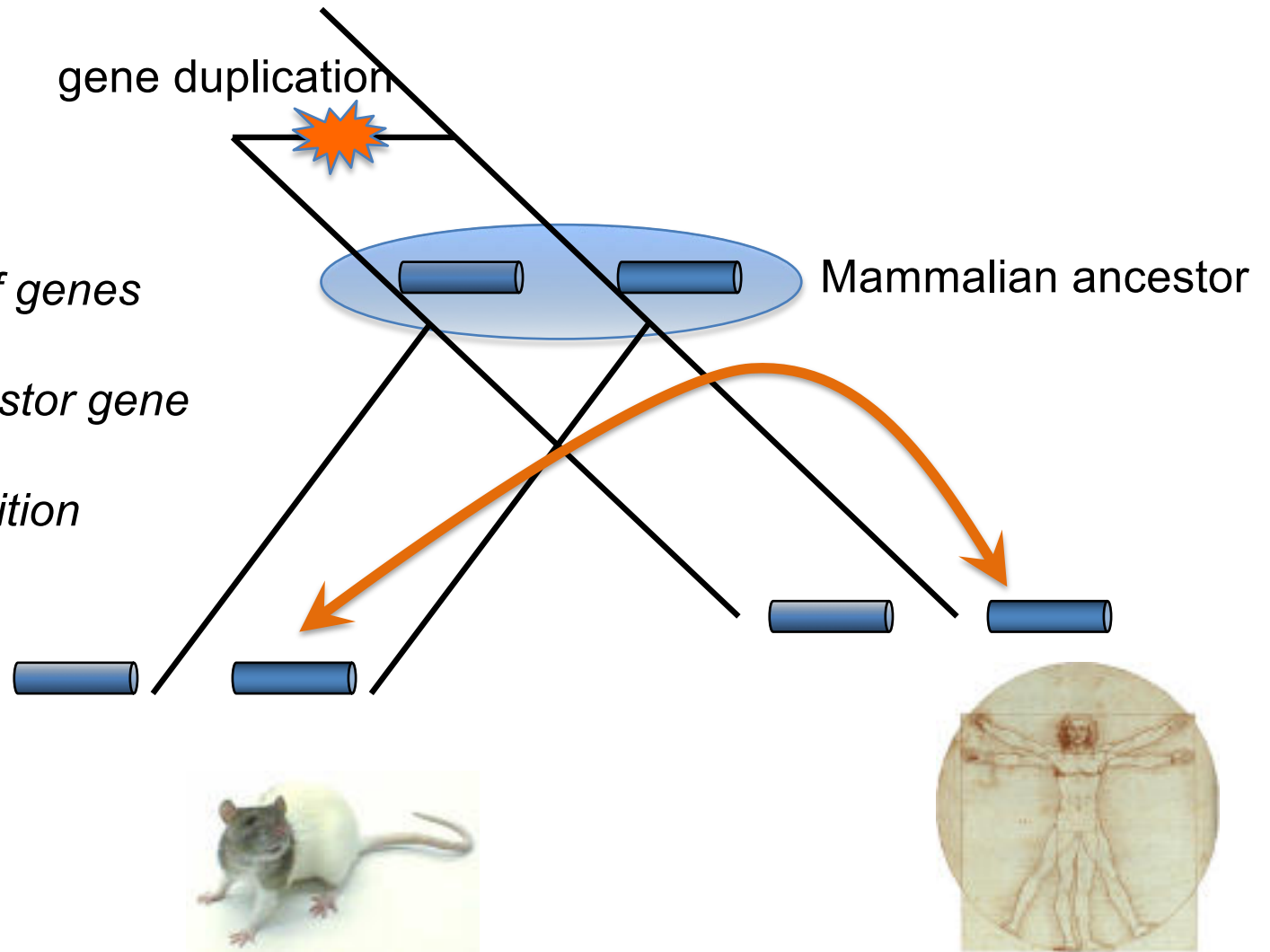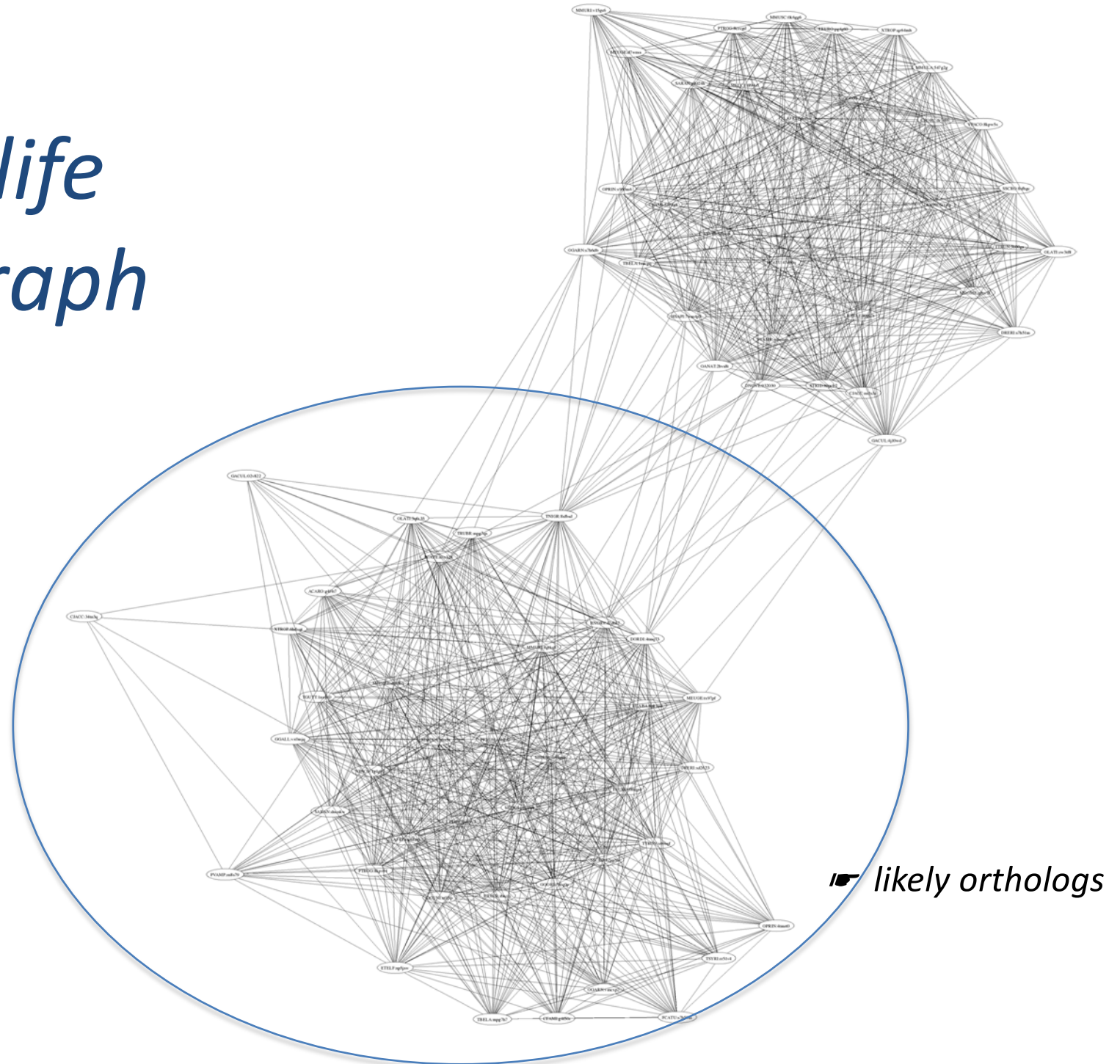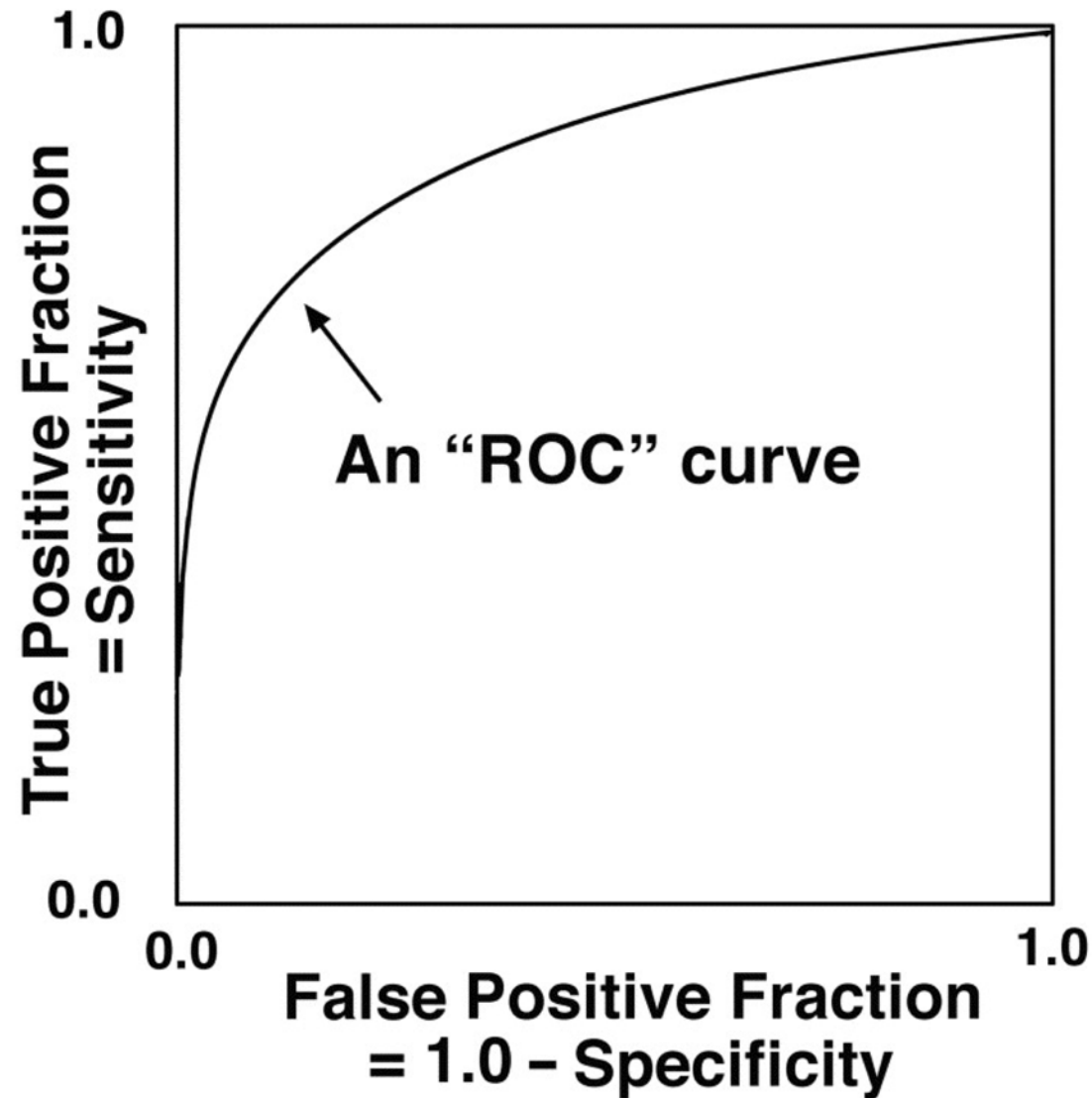# Orthology identification: DIY is more error-prone

- COG/KOG
- InParanoid
- eggNOG
- OrthoMCL
- OrthoFinder
- OrthoDB (Orthologer)

UNIVERSITÉ
DE GENÈVE
FACULTÉ DE MÉDECINE

# *Sensitivity vs. Specificity tradeoff*

# To keep in mind

- Sensitivity *vs.* Specificity tradeoff

- Cost ($ or time) *vs* Accuracy tradeoff

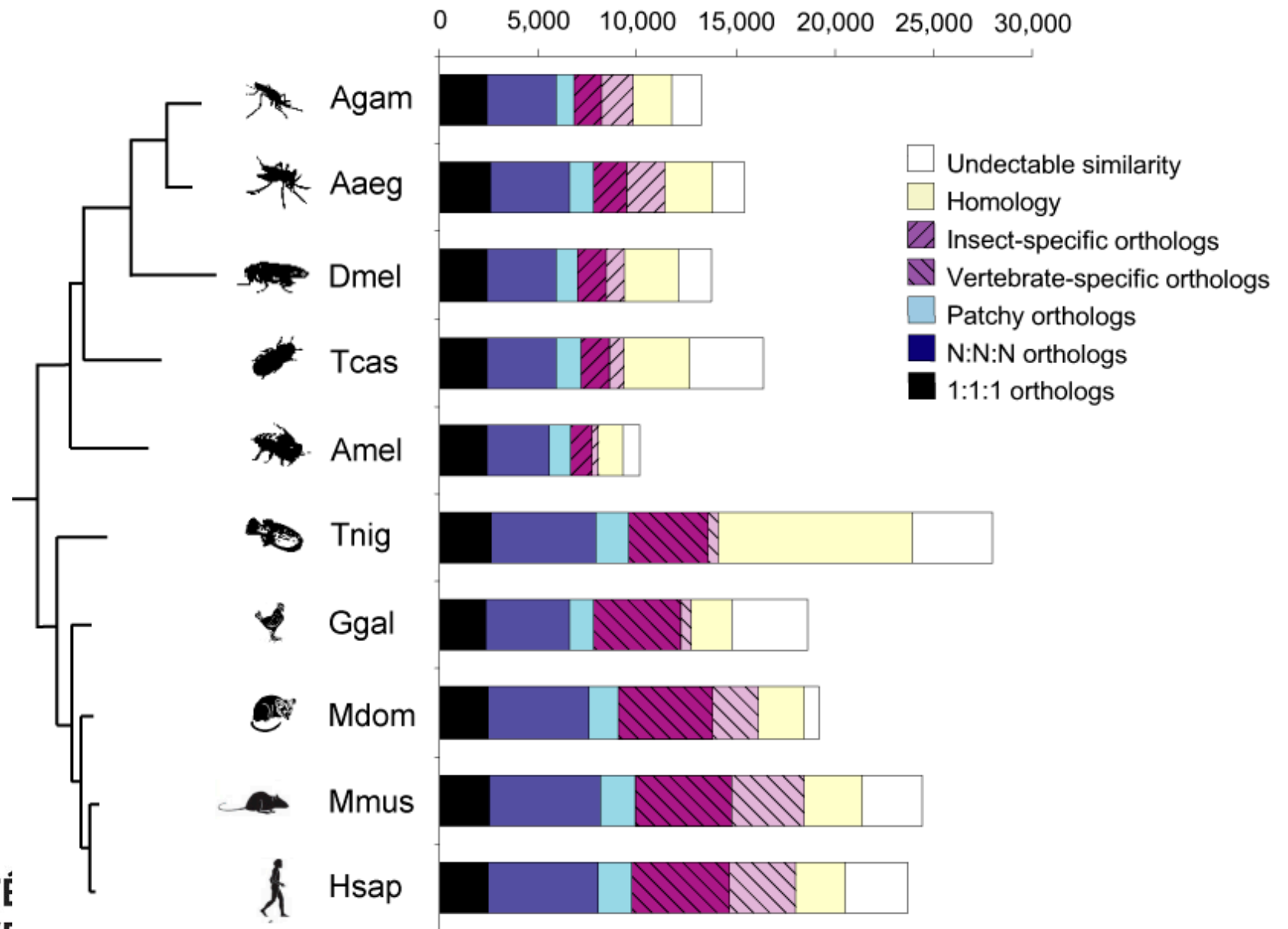☛ find or do benchmarking

UNIVERSITÉ
DE GENÈVE
FACULTÉ DE MÉDECINE

# *Examples of comparative study*

UNIVERSITÉ
DE GENÈVE

FACULTÉ DE MÉDECINE

# Counting genes
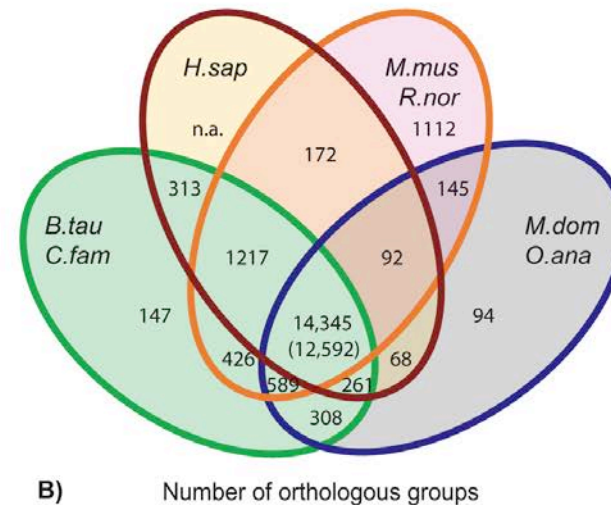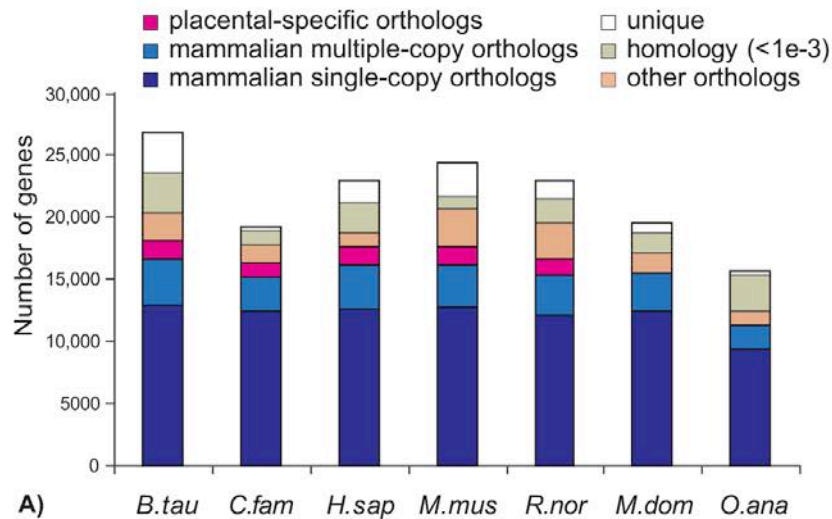
# Cow is molecularly closed to human than mouse



A) Number of genes for B.tau, C.fam, H.sap, M.mus, R.nor, M.dom, O.ana

Legend:
- placental-specific orthologs
- mammalian multiple-copy orthologs
- mammalian single-copy orthologs
- unique
- homology (<1e-3)
- other orthologs

B) Number of orthologous groups

C) Number of single-copy orthologs vs Percent a.a. identity to H.sap

D) Phylogenetic tree: B.tau, C.fam, H.sap, M.mus, R.nor, M.dom, O.ana with divergence times ~80-100 MYA, ~75-95 MYA, ~12-21 MYA, ~170-190 MYA

*Pathway perspective*

# Nasonia problem with amino acid metabolism



UNIVERSITÉ
DE GENÈVE
FACULTÉ DE MÉDECINE