# Day 5 sampling - clustering

# SAMPLE POPULATION

**SAMPLING**: IS ESTIMATING THE CHARACTERISTICS OF THE WHOLE POPULATION USING INFORMATION COLLECTED FROM A **SAMPLE** GROUP.

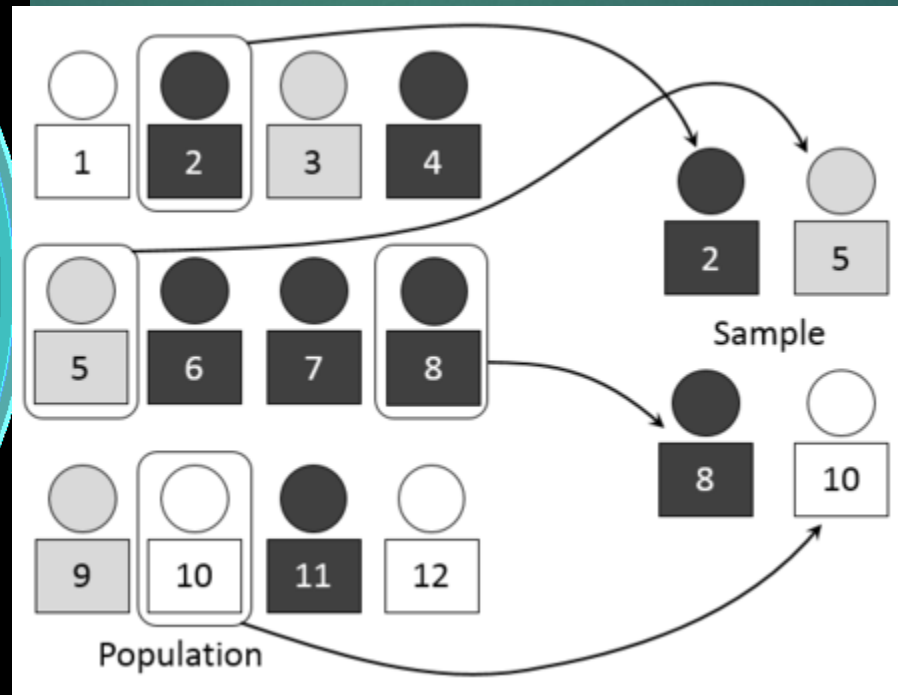The sampling process comprises several stages:

•Defining the population of concern
•Specifying a sampling frame, a set of items or events possible to measure
•Specifying a sampling method for selecting items or events from the frame
•Determining the sample size
•Implementing the sampling plan
•Sampling and data collecting

# Simple random sampling

In a simple random sample (SRS) of a given size, all such subsets of the frame are given an equal probability.

In particular, the variance between individual results within the sample is a good indicator of variance in the overall population, which makes it relatively easy to estimate the accuracy of results.
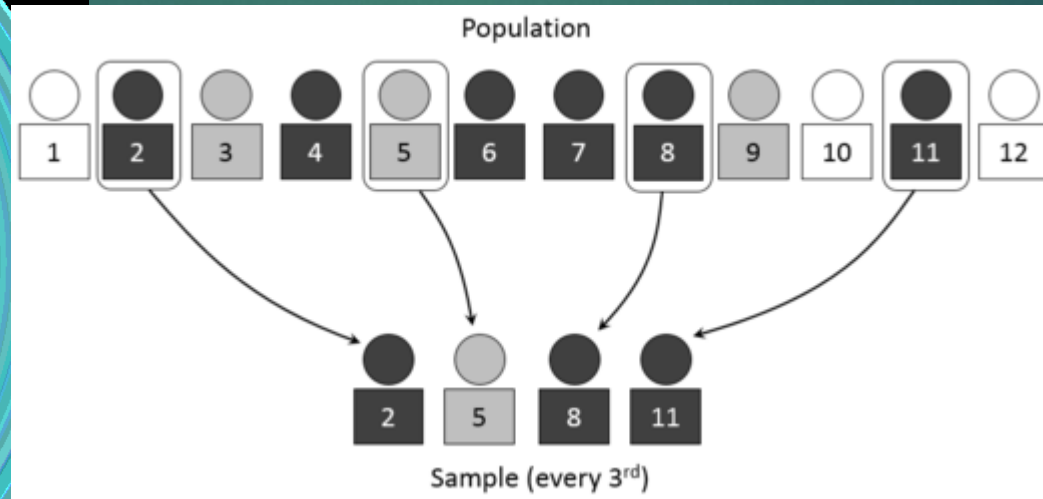
SRS can be vulnerable to sampling error because the randomness of the selection may result in a sample that doesn't reflect the makeup of the population.
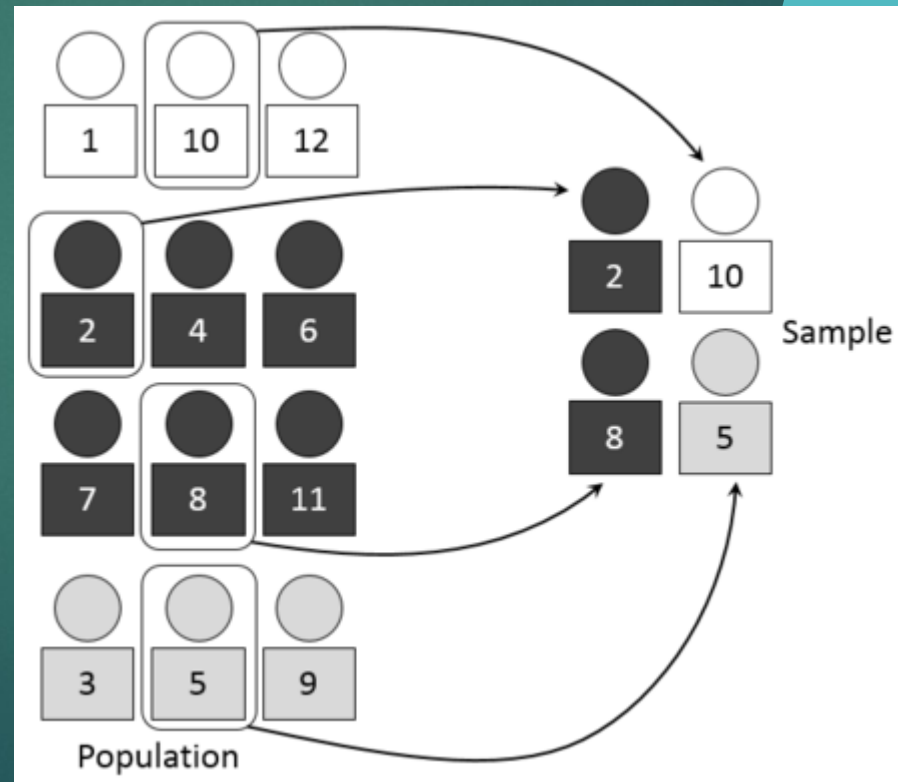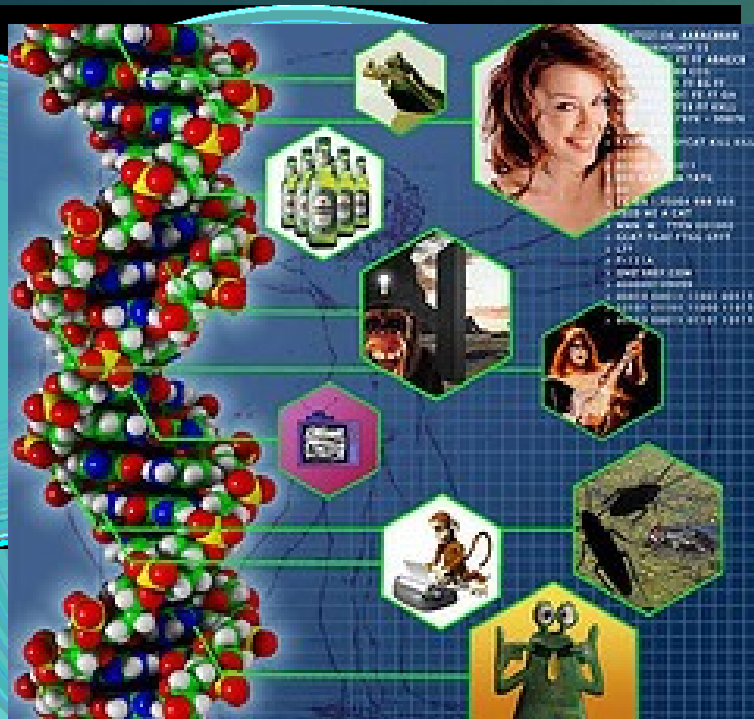
## Systematic sampling

Systematic sampling (also known as interval sampling) relies on arranging the study population according to some ordering scheme and then selecting elements at regular intervals through that ordered list

Systematic sampling involves a random start and then proceeds with the selection of every $k$th element from then onwards. In this case, $k$=(population size/sample size). It is important that the starting point is not automatically the first in the list, but is instead randomly chosen from within the first to the $k$th element in the list.



Population

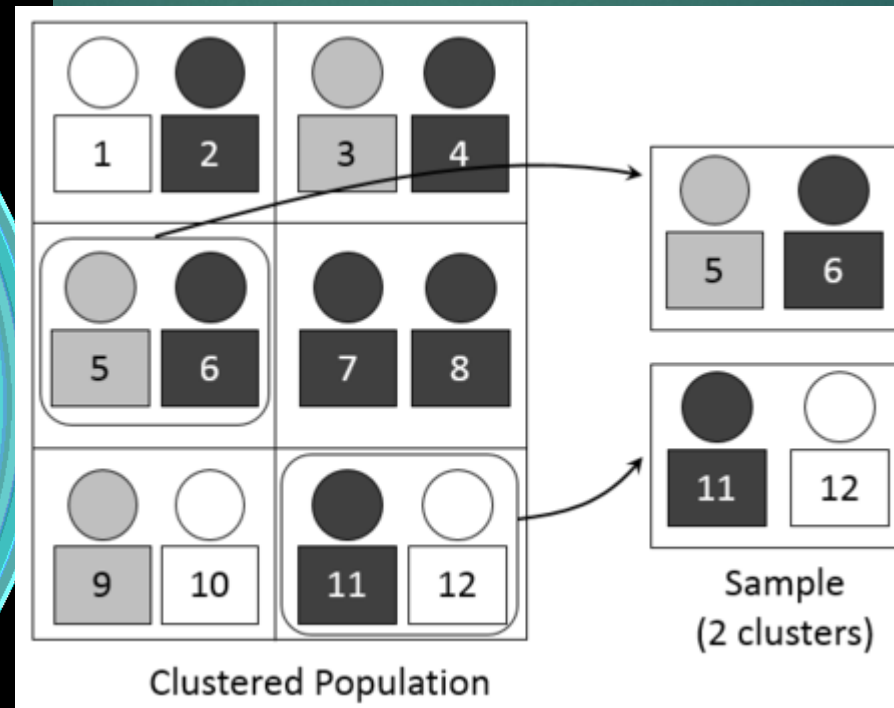1 2 3 4 5 6 7 8 9 10 11 12

2 5 8 11

Sample (every 3rd)

# STRATIFIED SAMPLING

WHEN THE POPULATION EMBRACES A NUMBER OF DISTINCT CATEGORIES, THE FRAME CAN BE ORGANIZED BY THESE CATEGORIES INTO SEPARATE "STRATA." EACH STRATUM IS THEN SAMPLED AS AN INDEPENDENT SUB-POPULATION, OUT OF WHICH INDIVIDUAL ELEMENTS CAN BE RANDOMLY SELECTED
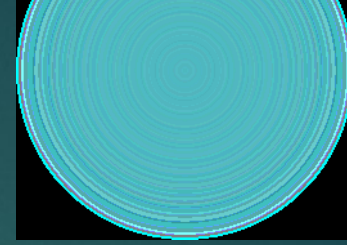
# Cluster sampling

▶ Sometimes it is more cost-effective to select respondents in groups ('clusters')



Clustered Population

Sample (2 clusters)

Quota sampling
Minimax sampling
Accidental sampling
Voluntary Sampling

# Windows Phone

## Offerta Privati

| LG Optimus 7 Ricaricabile | Vantaggi | Mobile Internet | Costo Telefono | |
|---|---|---|---|---|
| Ricaricabile | SIM Vodafone con 5€ di traffico | 3€ settimana 500MB inclusi | 399€ | Avvisami |

| LG Optimus 7 Abbonamento | Vantaggi | Mobile Internet | Costo Telefono | Contributo Mensile | Dettagli |
|---|---|---|---|---|---|
| Stile Libero New | 9 cent al minuto verso tutti | Incluso 2GB al mese | 0€ | 19€ | Avvisami |
| Tutto Facile Small | 50€ per chiamare e inviare SMS | Incluso 2GB al mese | 0€ | 44€ | Avvisami |
| Tutto Facile Medium | 100€ per chiamare e inviare SMS | Incluso 2GB al mese | 0€ | 69€ | Avvisami |
| Tutto Facile Large | 150€ per chiamare e inviare SMS | Incluso 2GB al mese | 0€ | 84€ | Avvisami |
| Tutto Facile Top Club | 200€ per chiamare e inviare SMS | Incluso 2GB al mese | 0€ | 100€ | Avvisami |

SKY

vodafone

## Scegli da tre a cinque GENERI di Mondo.

INTRATTENIMENTO | BAMBINI | DOCUMENTARI | MUSICA | NEWS

## Aggiungi i PACCHETTI che ti interessano.

CINEMA | SPORT | CALCIO

Trova la tua combinazione ideale.

HD SEMPRE INCLUSA*

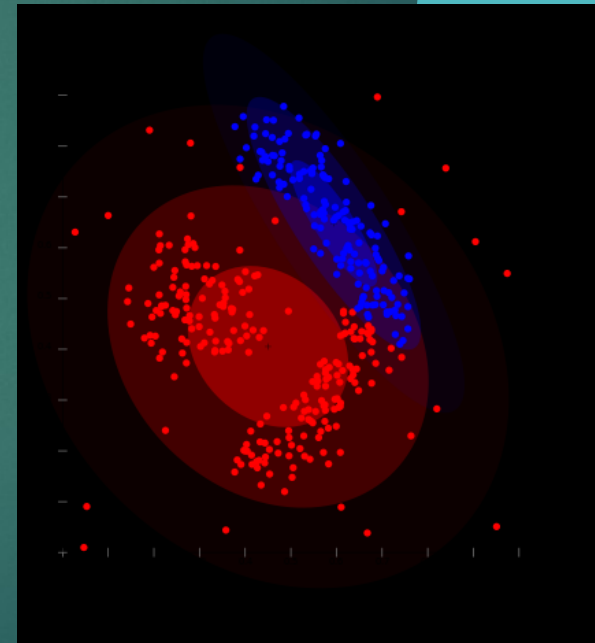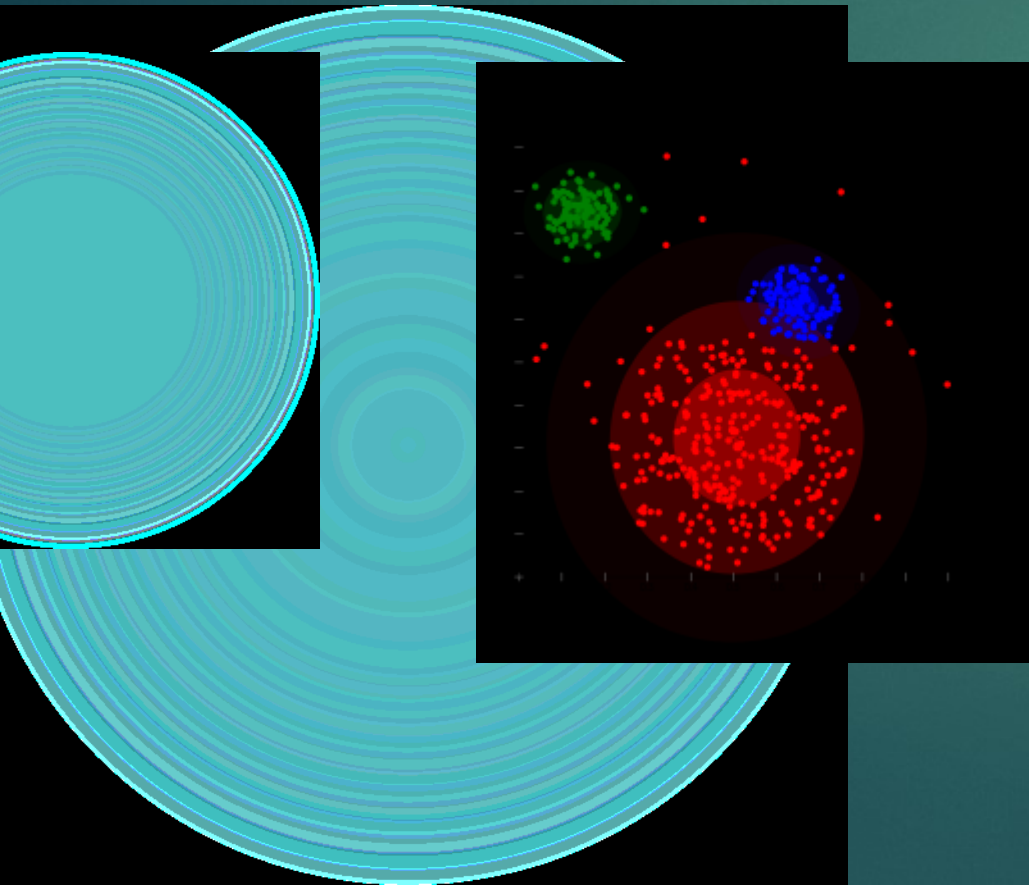| | 3 GENERI | 4 GENERI | 5 GENERI |
|---|---|---|---|
| MONDO | 19.90€ | 24.90€ | 29.90€ |
| MONDO + CINEMA | | 34€ NOVITÀ | 39€ NOVITÀ | 43€ |
| MONDO + SPORT CALCIO (1 PACCHETTO a scelta tra Sport e Calcio) | | 39€ NOVITÀ | 43€ |
| MONDO + (2 PACCHETTI a scelta tra Cinema Sport e Calcio) | | 52€ NOVITÀ | 56€ |
| MONDO + CINEMA SPORT CALCIO | | 65€ NOVITÀ | 69€ |

# CLUSTERING
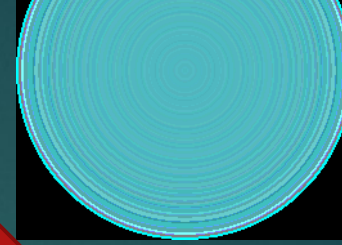
**1D**

**2D**

**3D**

**nD**

Clustering techniques are based on measures of similarity between elements. In many approaches this similarity is conceived in terms of distance in a multidimensional space.
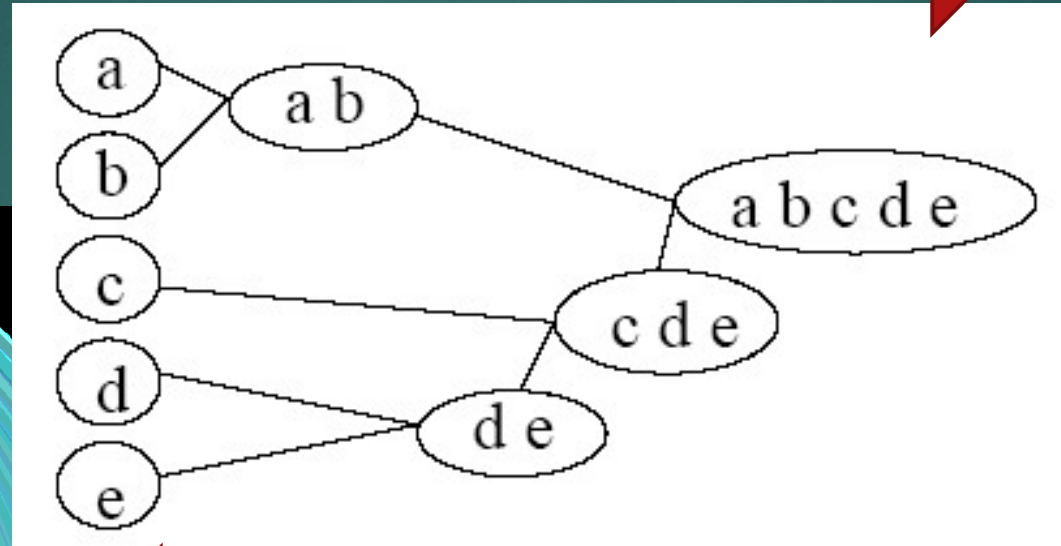
The notion of a "cluster" cannot be precisely defined, which is one of the reasons why there are so many clustering algorithms.There is a common denominator: a group of data objects
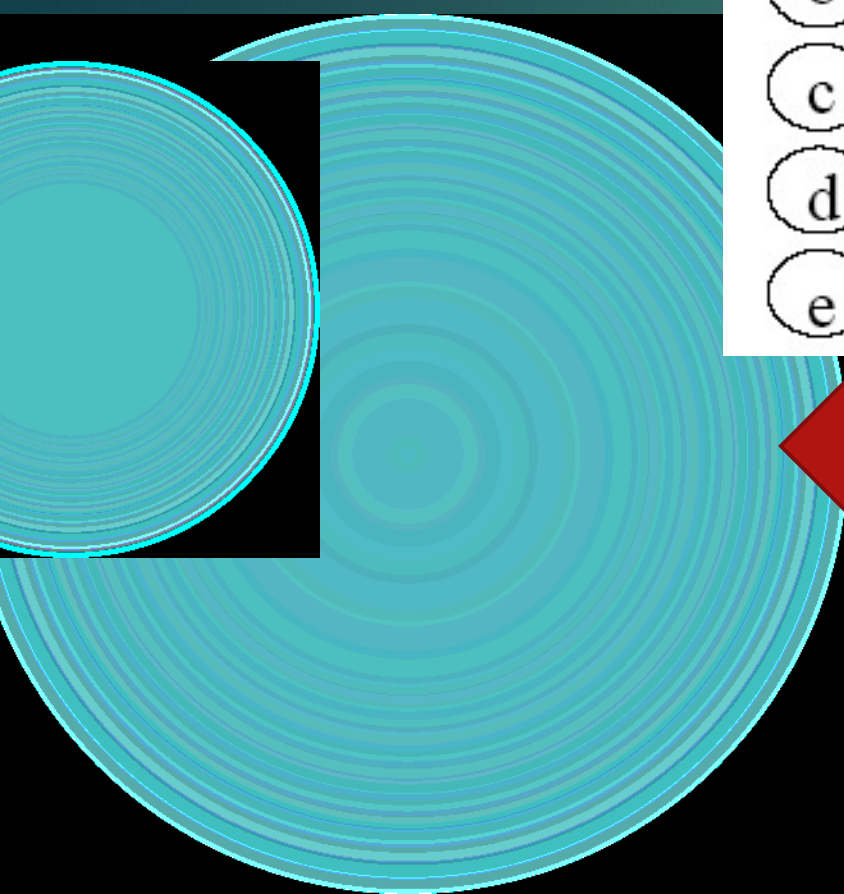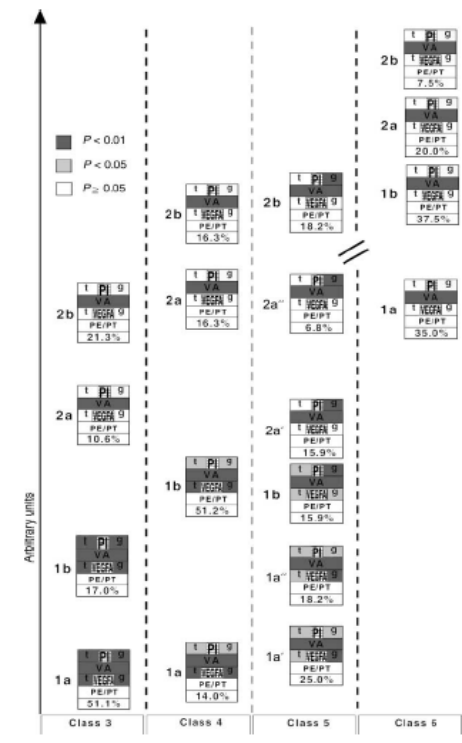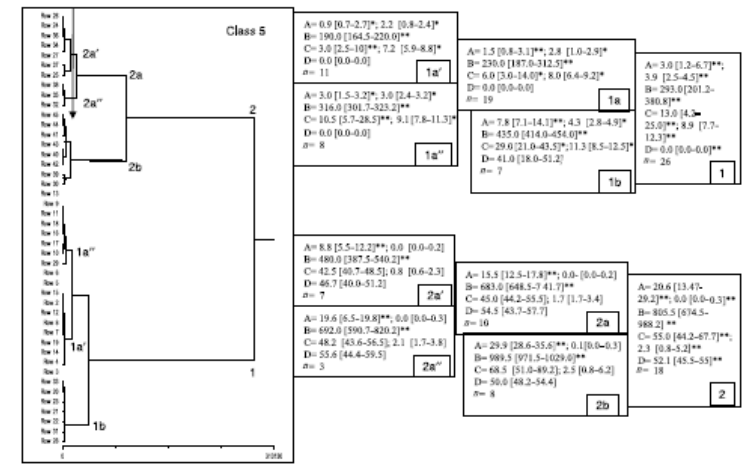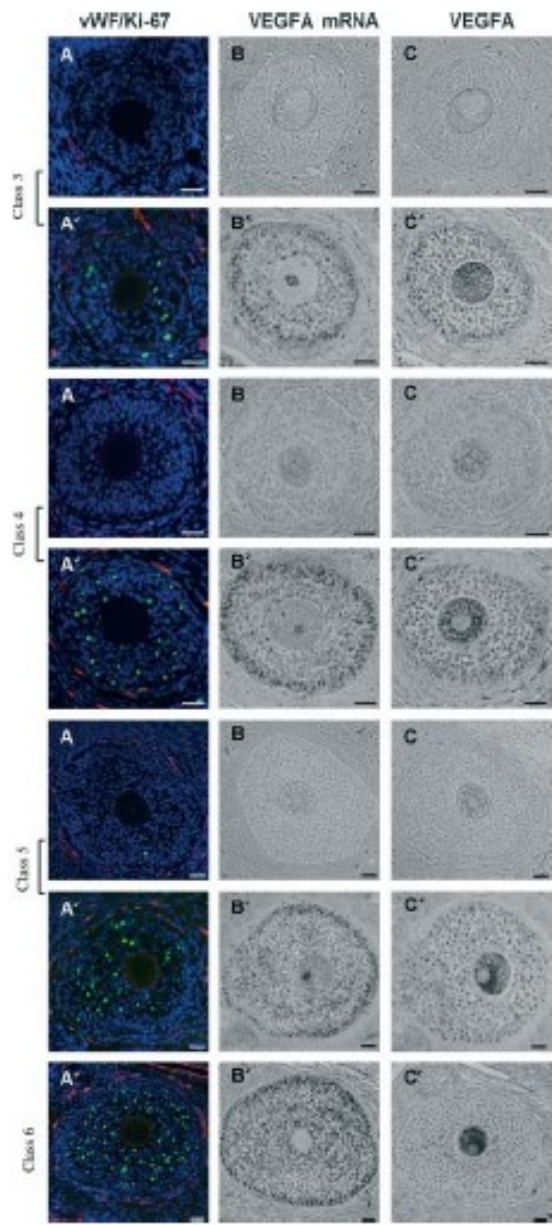
# Clustering



Bottom-up

Top-Down

dendogramma

Figure 5 A general model that describes the distribution of the follicular subpopulations identified inside each follicular class by using cluster analysis. The subpopulations were represented within each class and arranged on the y-axis by considering the distances between the bifurcation obtained in the dendograms. The grey scale represents significant differences recovered for somatic and vascular parameters. The thickness of the dotted line represents the statistical difference among classes. t PI, theca proliferation index; g PI, granulosa proliferation index; VA, vascular area; t VEGFA, theca VEGFA mRNA; g VEGFA, granulosa VEGFA mRNA; PE/PT, proportion of proliferating endothelial cells; % value, percentage of preantral follicles belonging to the subpopulation.

# CLUSTERING



Comp Clin Pathol

Fig. 2 Multivariate hierarchical cluster analysis of MICRO dogs. The *caption* shows the area of interest where clustering is significant



Fig. 3 Mono-dimensional distribution of distances among the dogs as measured by clustering. The *caption* shows a clear subgroup among the MICRO dogs
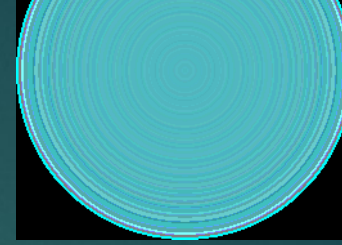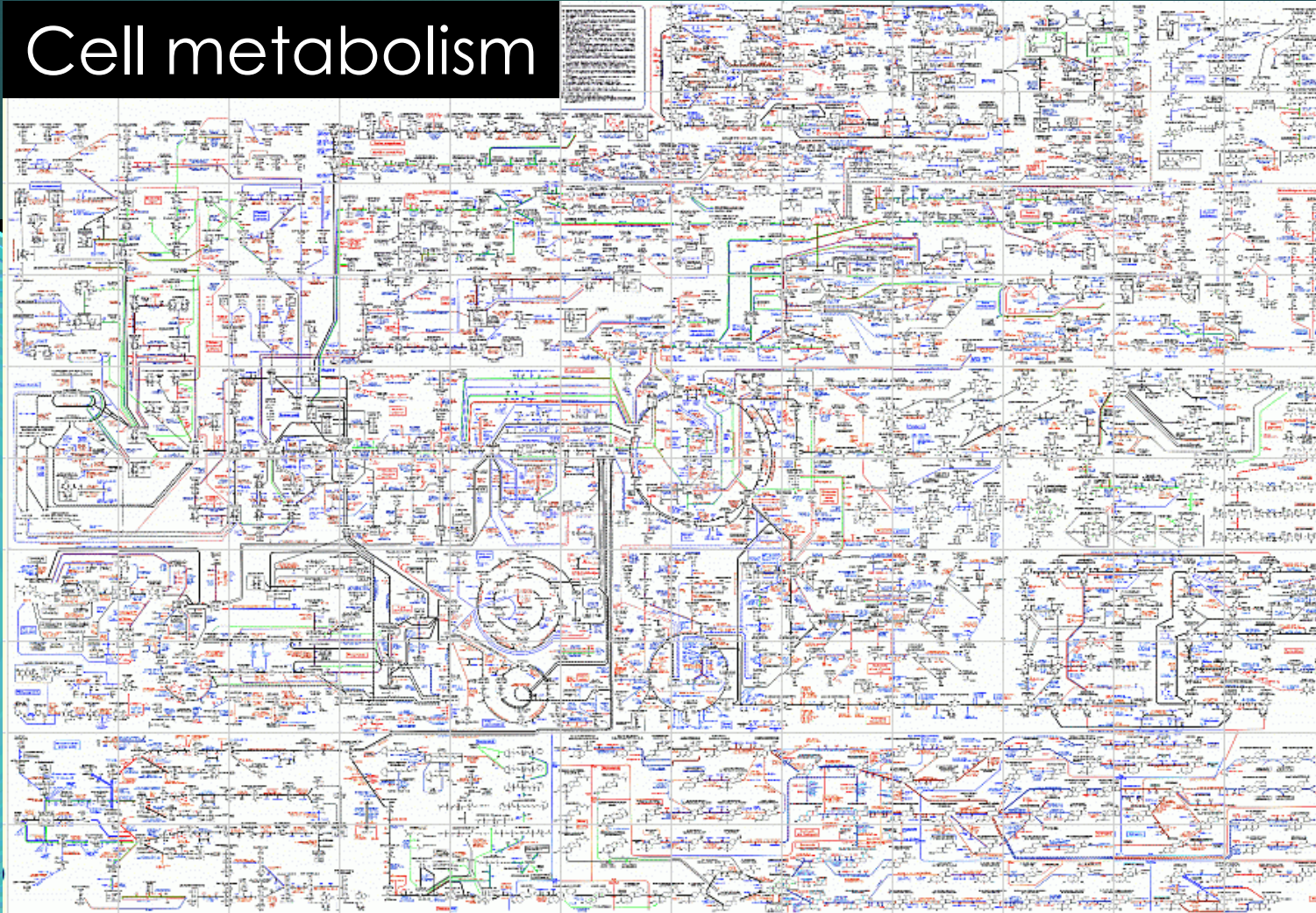
# Inividuality Complexity



© www.matteofossati.it

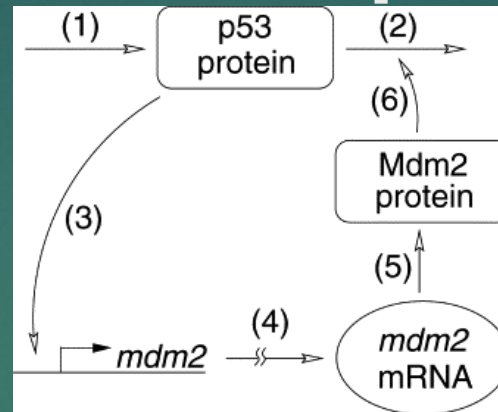# Complicated vs. complex

# fertility

# Biological complexity
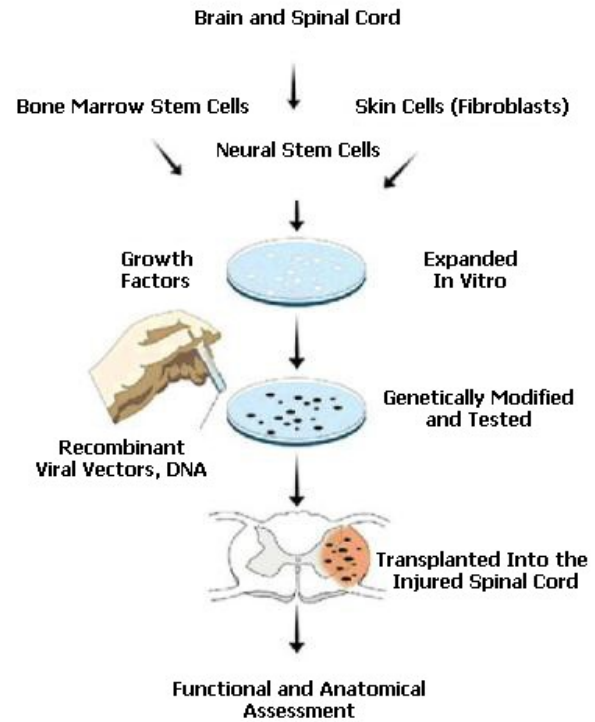
Cell metabolism

# Complexity:
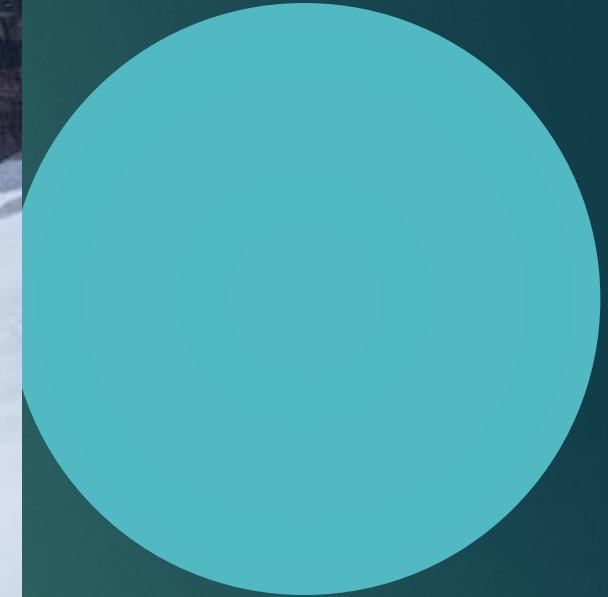
non-linearity of interactions:
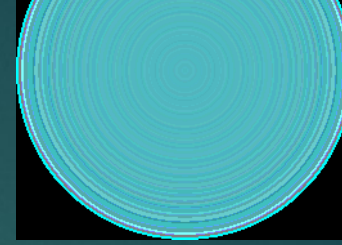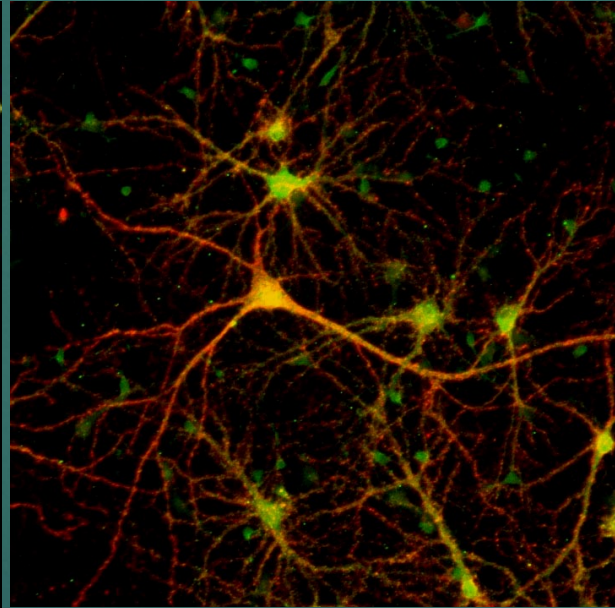
# An example…
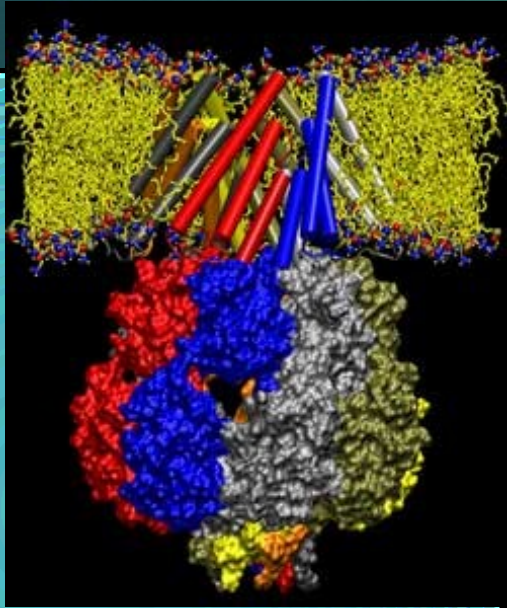
un

# unpredictability



Strategies of spinal cord transplantation and gene therapy
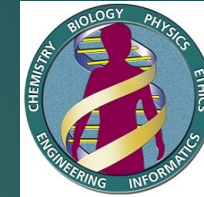
# Butterfly effect

# Emergence of proprieties

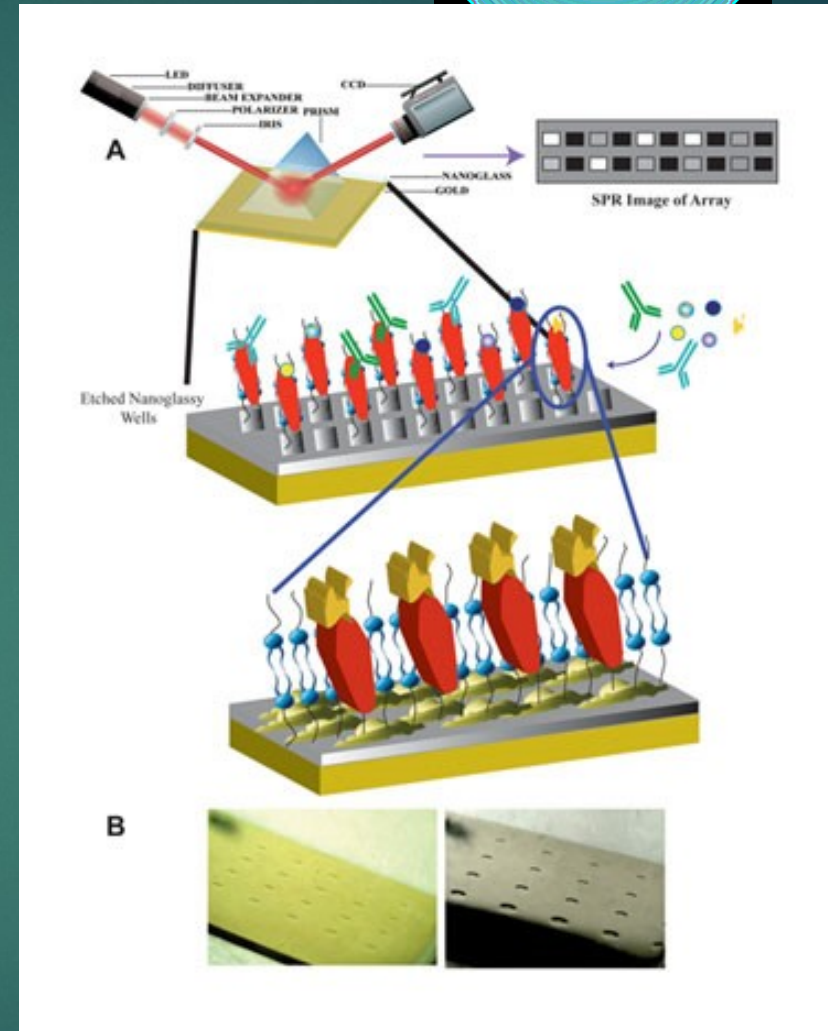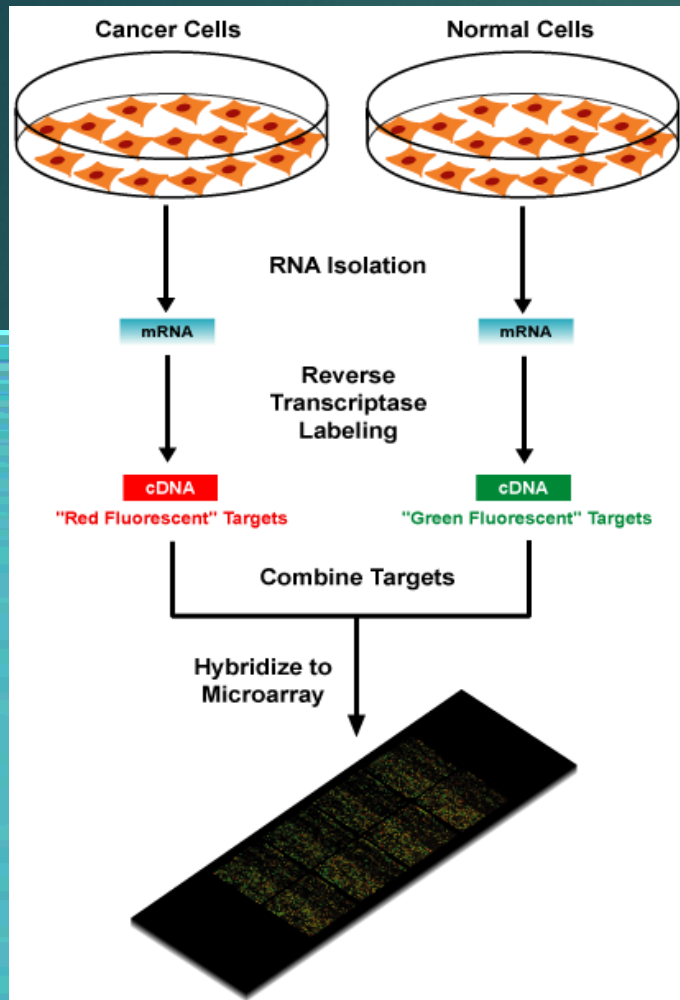# The whole is more (different) than the sum of the individual components
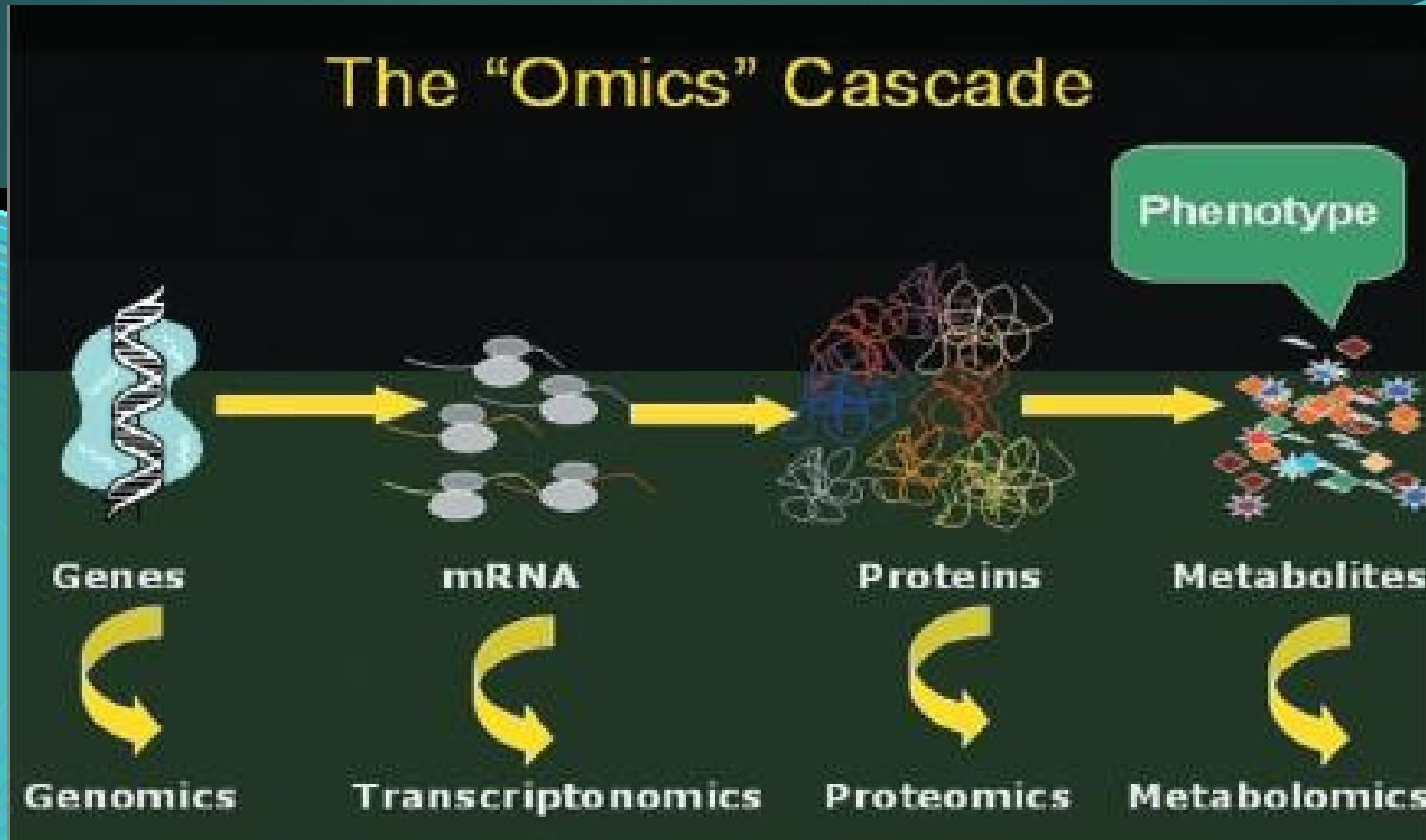
# Human Genome Project



Begun formally in 1990, the U.S. Human Genome Project was a 13-year effort coordinated by the U.S. Department of Energy and the National Institutes of Health. The project originally was planned to last 15 years, but rapid technological advances accelerated the completion date to 2003. Project goals

- *identify* all the approximately 20,000-25,000 genes in human DNA,

- *determine* the sequences of the 3 billion chemical base pairs that make up human DNA,

- *store* this information in databases,

- *improve* tools for data analysis,

- *transfer* related technologies to the private sector, and

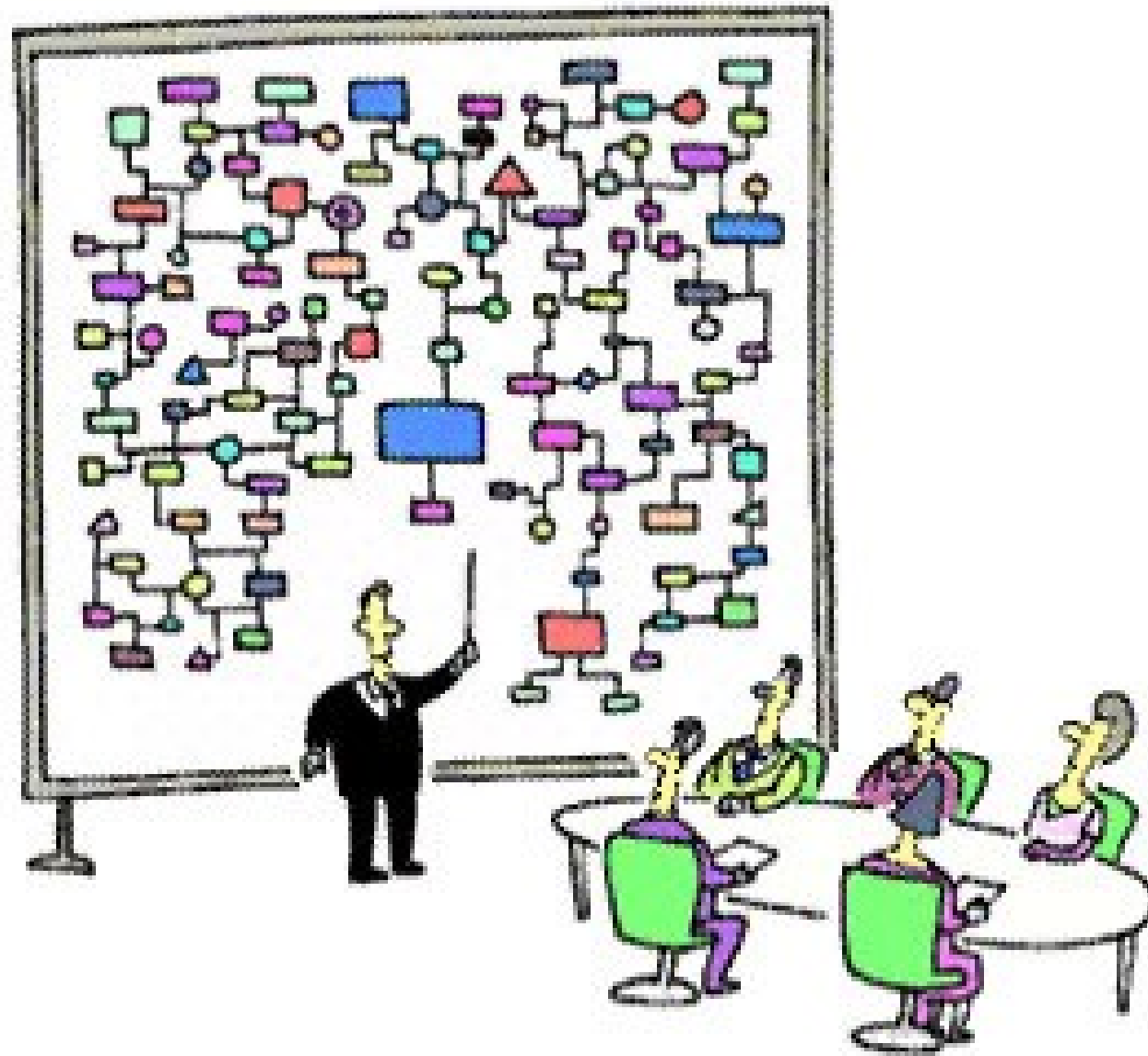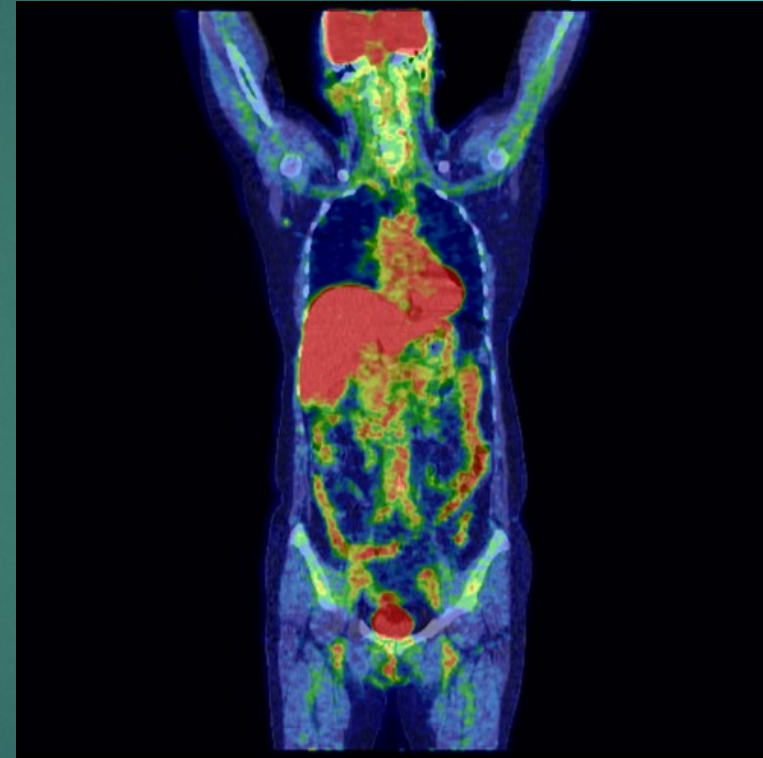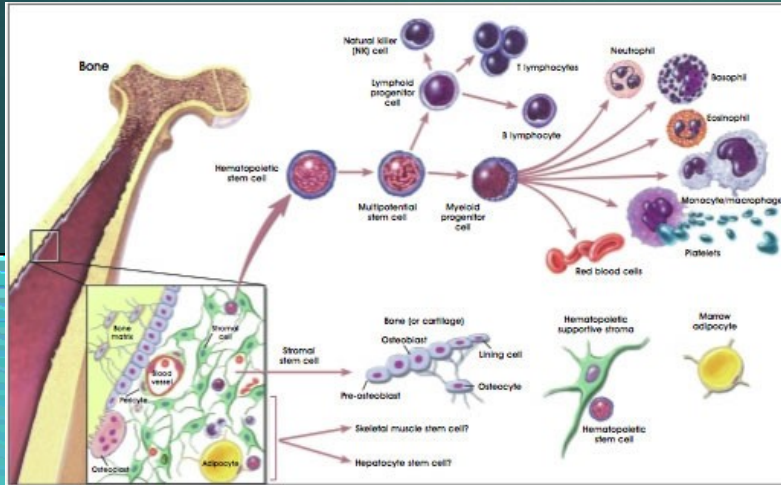- *address* the ethical, legal, and social issues (ELSI) that may arise from the project.

# - OMICS



The "Omics" Cascade

"And that's why we need a computer."

# Computational models in biology and medicine