# Bioinformatics
# (statistics/databases/big data)

Bruno Pot

Applied Maths

**Veyrier-du-Lac**
07/04/2016

# What are the challenges in our domain?

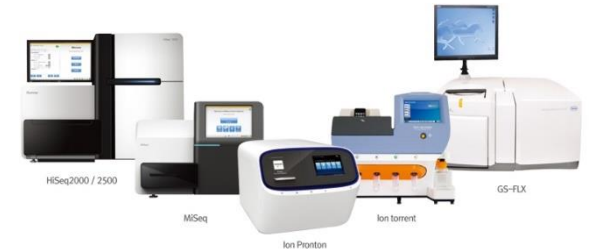- The quality of a *bioinformatics* analysis will depend on:

- Quality
  - of *sample preparation* (soil, stool, sputum, swaps, blood ...)
  - of *wet lab protocols* (DNA preparation, library preparation)
    - → as standardized as possible
  - sequencing *equipment* (quality, read length, sequencing depth,...)

- Can *bioinformatics* help you to judge the robustness and reliability of your sequencing preparations?
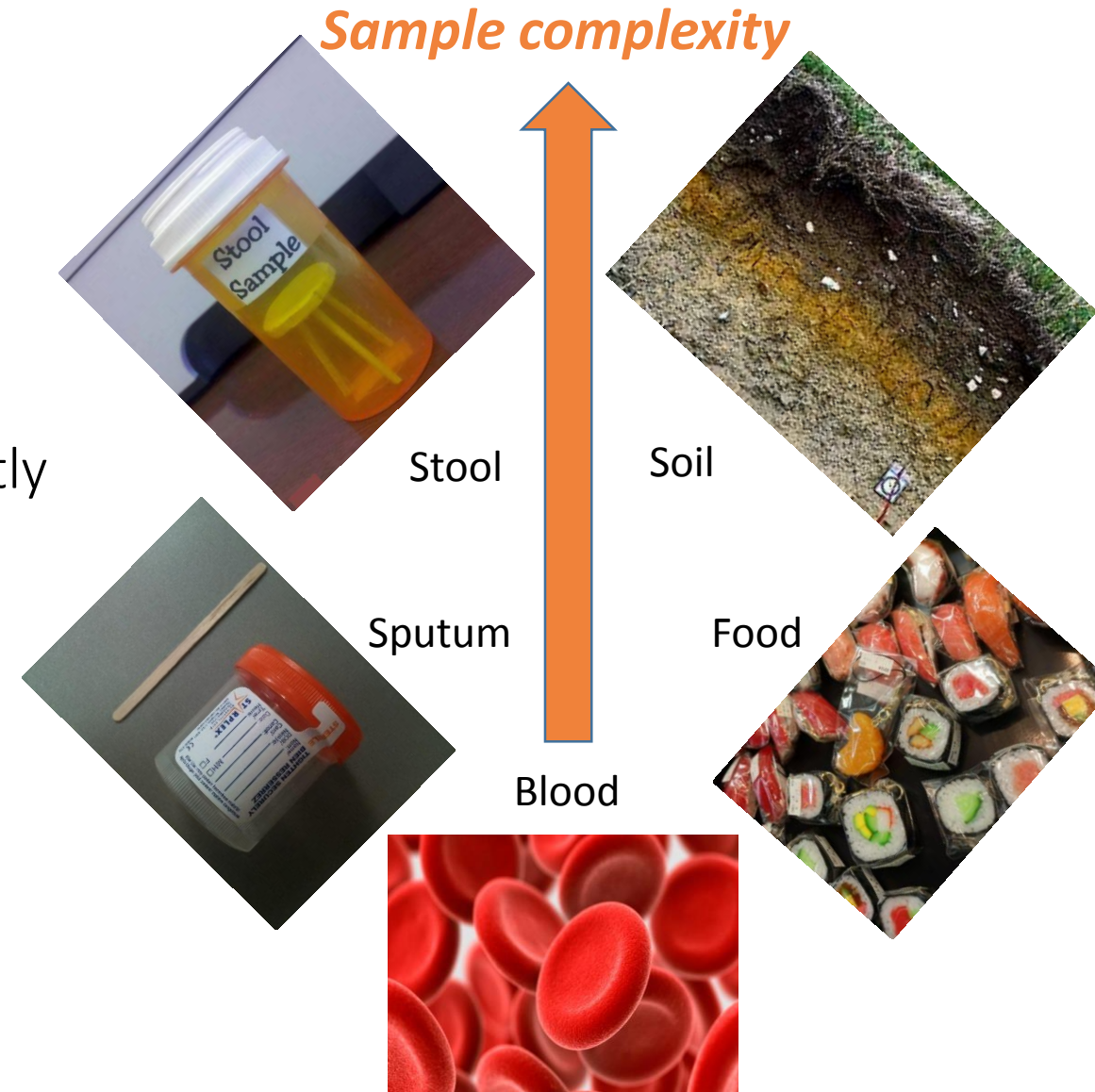

DNA extraction


Library preparation


Sequencing depth

# What are the challenges in our domain?

- *Bioinformatics* analysis will depend on the sample origin (read: **complexity**)

- This will depend on the availability of sufficently '**Complete**' and '**Reliable**' reference **databases** for proper identification of (all) (relevant) μorganisms and functionalities, and

- require the necessary **validation** processes (inter- and intra run controls; artificial mixes)

**Sample complexity**

Stool

Soil

Sputum

Food

Blood

# High-throughput data processing !
## (diagnostics, quality control,...)

Requires
- Extended computer power
  (Cloud, Cluster, PC)
- Extended storage capacity
  (fast memory!)
- Faster data transfer

- Advanced (fast) algorithms
- Read analysis strategy
  (assemble-identify/ identify)
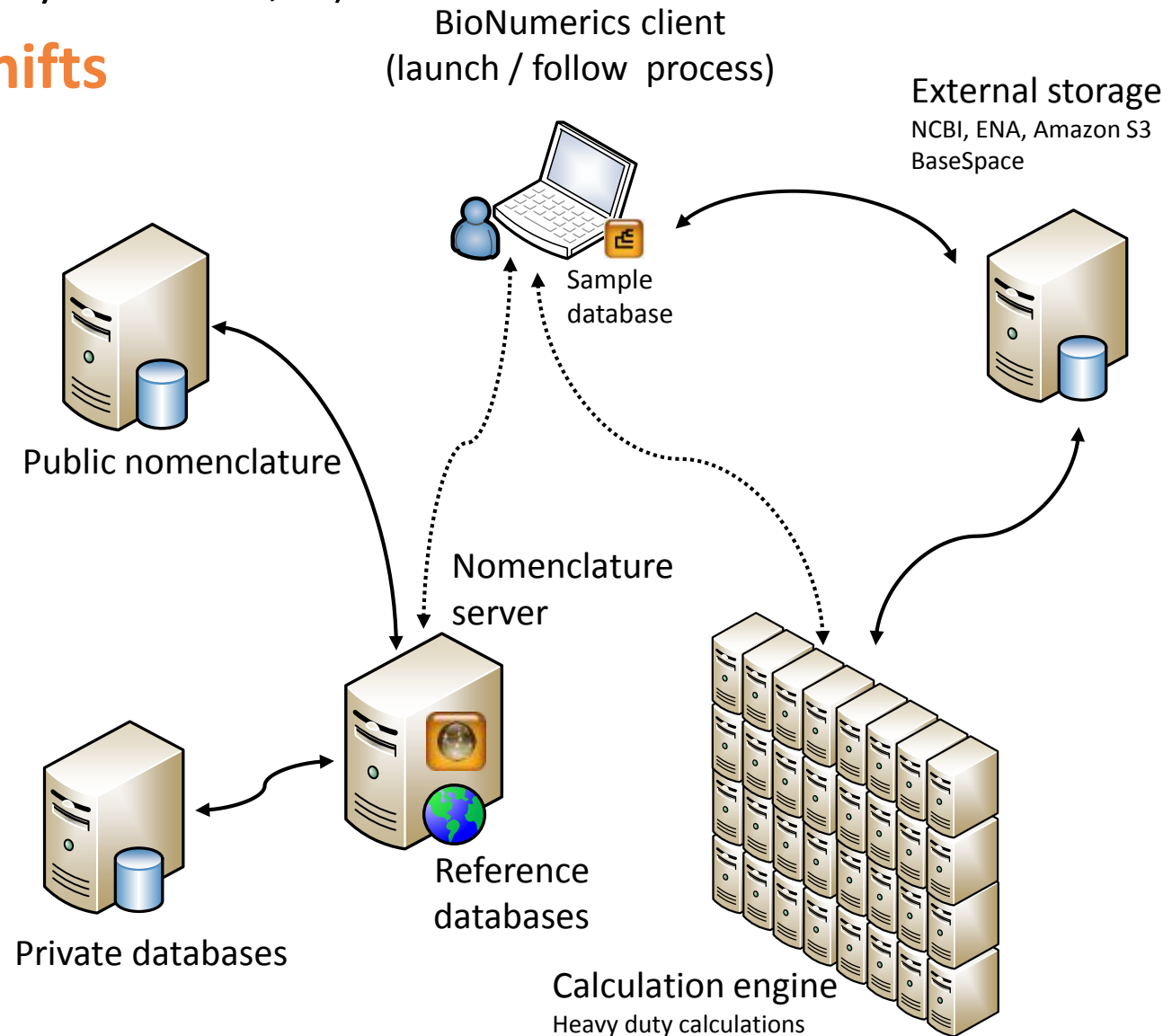- Advanced visualization tools

- Userfriendliness for non-specialized users
  - software installation level
  - software usage level
  - software maintenance level
  - user management and confidentiality management

**Mindshifts**

**HARDware**

**SOFTware**

**BI&IT support**

BioNumerics client
(launch / follow  process)

External storage
NCBI, ENA, Amazon S3
BaseSpace

Sample database

Public nomenclature

Nomenclature server

Private databases

Reference databases

Calculation engine
Heavy duty calculations
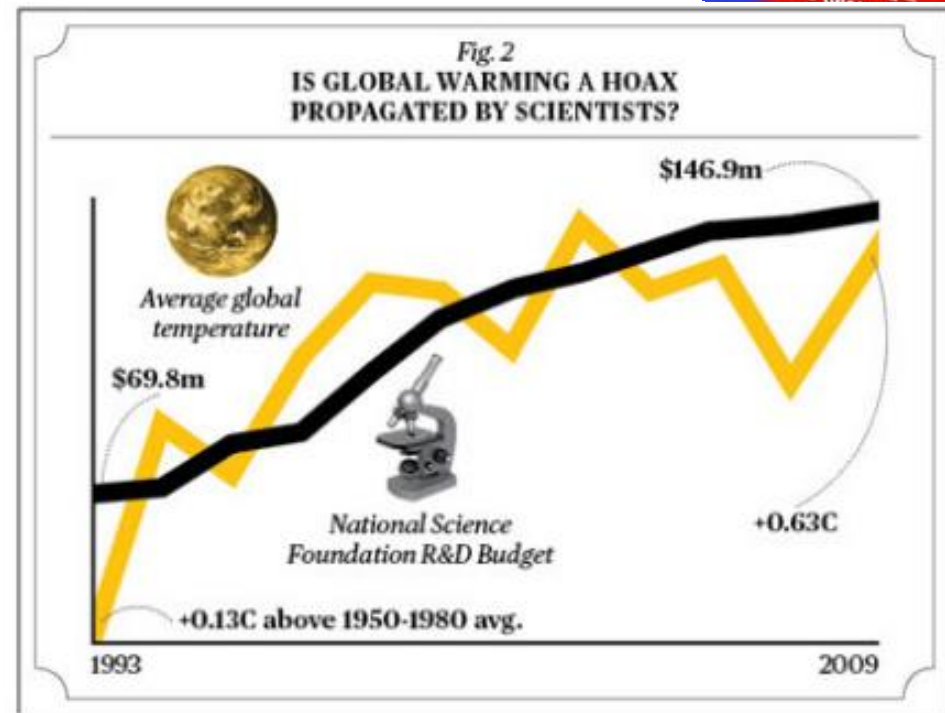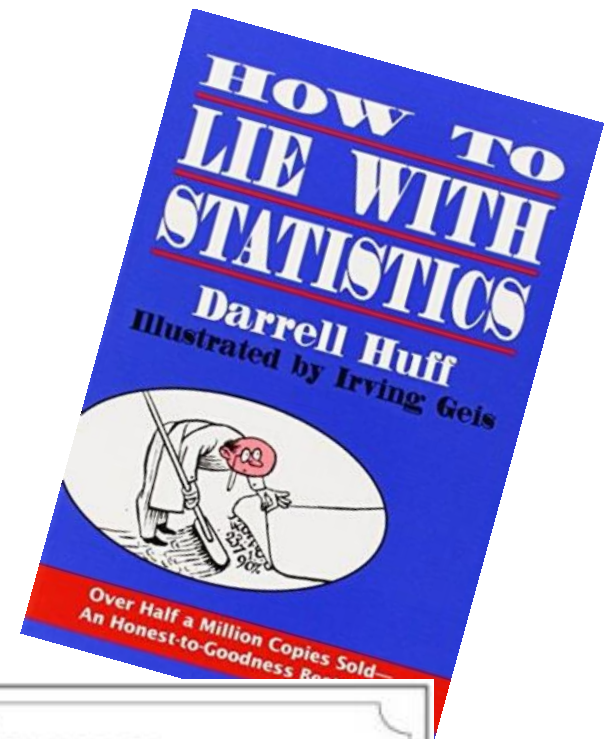
# Statistics tell it &ll…

**Read statistics** (length, numbers, quality,..):     Q U A L I T Y    evaluation

- Filtering (singletons, chimaeras, barcoding…)
- Assembly statistics
- Matching statistics

**Analytical statistics**

R E L I A B I L I T Y    evaluation

- Diversity (indices)
- Distribution / geography (plots)
- Timeframe analyses (plots)
- Cluster algorithms
- PCA and other visualization algorithms
- Correlation versus association versus cause
- (M)anova analyses
- Significance tests
- …

# Conclusions

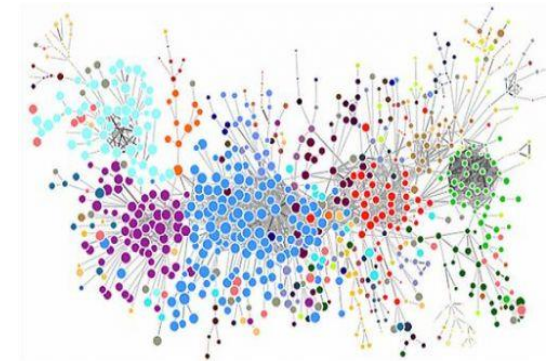## Reliability issues

- Different wet lab protocols will give different outcomes on the same sample (e.g. library preparation)
- Different pipelines will yield different outcomes on the same dataset
  - → Artifical *in silico* data are difficult to reconstruct with existing pipelines (due to need for filtering, etc…)
- Different statistical settings will give different interpretations on the same datasets

## Complexity issues

- Amplicon-based metagenomics approach verus full shot gun metagenomics
- Causalilty questions
- Interpretation of results in non-perfect conditions (geography, time analysis, different populations, climates, disease types, …)

## Practical hurdles

- Amount of data (storage, analysis, visualization,…)
- Complexity of the bioinformatics analysis (several quality dependent steps)
- Manage complexity of the IT infrastructure required
- Lack of adapted databases for 'new' subjects