# NGS sequence application, a few examples

Dr Francois Sabot & Christine Tranchant-Dubreuil

8th of October, 2018

IRD - UMR DIADE

# Analyses

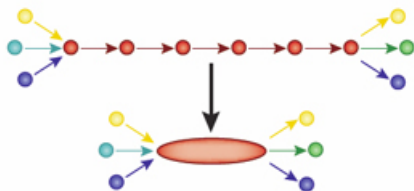1. Fragment DNA and sequence

2. Find overlaps between reads

...AGCCTAGACCTACAGGATGCGCGACACGT
GGATGCGCGACACGTCGCATATCCGGT...

From Baker, 2012

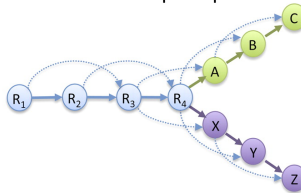3. Assemble overlaps into contigs

4. Assemble contigs into scaffolds
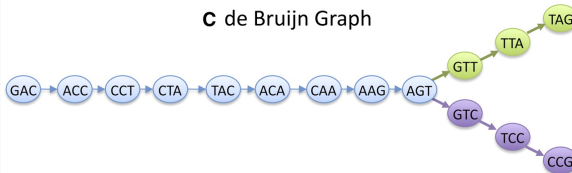
*Michael Schatz, Cold Spring Harbor*

From Baker, 2012

From Schatz, 2010

From Wikipedia

## Generally with **Pair-End** data

## Generally with **Pair-End** data

From Eurofins

From CGFB, Bordeaux, France

- Mainly in RNA sequencing, but also in CNV (Copy Number Variation)
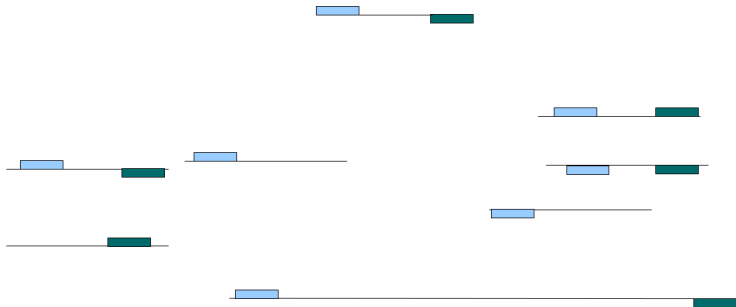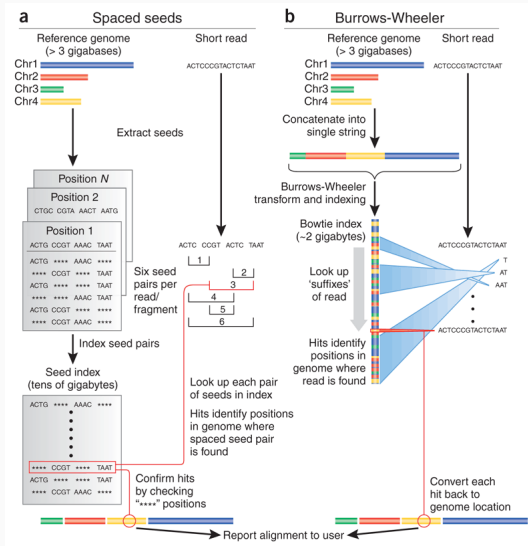
- Mainly in RNA sequencing, but also in CNV (Copy Number Variation)
- Counting the number of reads/bases at each position

- Mainly in RNA sequencing, but also in CNV (Copy Number Variation)
- Counting the number of reads/bases at each position
- More precise than ChiP

- Mainly in RNA sequencing, but also in CNV (Copy Number Variation)
- Counting the number of reads/bases at each position
- More precise than ChiP
- Need to be reproduced
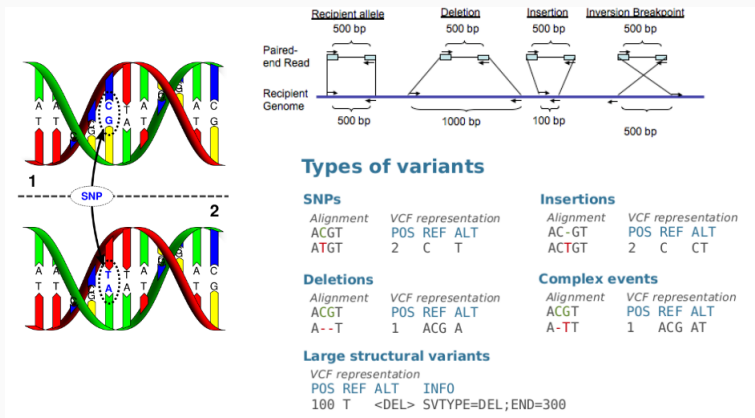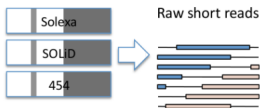
- Mainly in RNA sequencing, but also in CNV (Copy Number Variation)
- Counting the number of reads/bases at each position
- More precise than ChiP
- Need to be reproduced
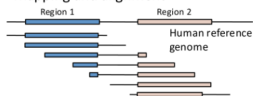- Lots of Statistical models and Controls behind

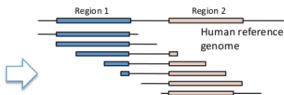From unmapped reads to true genetic variation in next-generation sequencing data

From Korbel et al, 2007

# Common File for all Variations, the VCF



**Example**

VCF header

Body

```
##fileformat=VCFv4.0          ← Mandatory header lines
##fileDate=20100707
##source=VCFtools
##reference=NCBI36
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele">
##INFO=<ID=H2,Number=0,Type=Flag,Description="HapMap2 membership">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality (phred score)">
##FORMAT=<ID=GL,Number=3,Type=Float,Description="Likelihoods for RR,RA,AA genotypes (R=ref,A=alt)">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##ALT=<ID=DEL,Description="Deletion">
##INFO=<ID=SVTYPE,Number=1,Type=String,Description="Type of structural variant">
##INFO=<ID=END,Number=1,Type=Integer,Description="End position of the variant">
#CHROM POS ID    REF  ALT   QUAL FILTER  INFO                   FORMAT    SAMPLE1    SAMPLE2
1      1    .    ACG  A,AT  .    PASS                            GT:DP     1/2:13     0/0:29
1      2    rs1  C    T,CT  .    PASS   H2;AA=T                  GT:GQ     0|1:100    2/2:70
1      5    .    A    G     .    PASS                            GT:GQ     1|0:77     1/1:95
1      100  .    T    <DEL> .    PASS   SVTYPE=DEL;END=300       GT:GQ:DP  1/1:12:3   0/0:20
```

Optional header lines (meta-data about the annotations in the VCF body)

Mandatory header lines

Reference alleles (GT=0)

Alternate alleles (GT>0 is an index to the ALT column)

Deletion

SNP

Large SV

Insertion

Other event

Phased data (G and C above are on the same chromosome)

VCF = Variant Call Format From 1000 Genomes Project

| Application | GS FLX++ | GS Junior | HiSeq 2500 | MiSeq | PacBio RS |
|---|---|---|---|---|---|
| **Genome Sequencing** | | | | | |
| De novo sequencing of bacterial & fungal genomes | ✓✓✓ | | ✓ | ✓✓ | ✓ |
| De novo sequencing of higher eukaryotic genomes | ✓✓ | | ✓✓✓ | | ✓ |
| De novo sequencing of BACs, viruses & plasmids | ✓✓✓ | ✓✓✓ | | | ✓ |
| Resequencing of genomes | | | ✓✓✓ | ✓✓ | |
| **Transcriptome Sequencing** | | | | | |
| De novo Transcriptome sequencing | ✓✓✓ | | ✓✓ | ✓✓ | |
| Expression profiling | | | ✓✓✓ | | |
| Small RNA sequencing | | | ✓✓✓ | ✓✓ | |
| ChIP sequencing | | | ✓✓✓ | ✓✓ | |
| **Resequencing & Amplicons** | | | | | |
| Ultra deep amplicon sequencing | ✓✓✓ | ✓✓✓ | ✓ | ✓ | |
| Resequencing by Sequence Capture | ✓✓ | ✓ | ✓✓✓ | | |

From Eurofins

From my own personal Experience:

**Assembly** : Nanopore, PacBio, Illumina (MySeq + HiSeq, various libraries)

From my own personal Experience:

**Assembly** : Nanopore, PacBio, Illumina (MySeq + HiSeq, various libraries)

**SNP detection** : Illumina

From my own personal Experience:

**Assembly** : Nanopore, PacBio, Illumina (MySeq + HiSeq, various libraries)

**SNP detection** : Illumina

**SV Variation** : Nanopore, PacBio, Illumina, IonTorrent

From my own personal Experience:

**Assembly** : Nanopore, PacBio, Illumina (MySeq + HiSeq, various libraries)

**SNP detection** : Illumina

**SV Variation** : Nanopore, PacBio, Illumina, IonTorrent

**Quantification** : Illumina

- Amount of original samples

- Amount of original samples
- Size of sequenced unit

- Amount of original samples
- Size of sequenced unit
- Error rate

- Amount of original samples

- Size of sequenced unit

- Error rate

- Volume of Outputted data

- Amount of original samples

- Size of sequenced unit

- Error rate

- Volume of Outputted data

All linked to technical constraints

- Cleaning data level

- Cleaning data level
- Mapping Conditions

- Cleaning data level
- Mapping Conditions
- Mapping Cleaning Conditions

- Cleaning data level
- Mapping Conditions
- Mapping Cleaning Conditions
- Variation Calling level

- Cleaning data level
- Mapping Conditions
- Mapping Cleaning Conditions
- Variation Calling level

All linked to the Specificity/Sensitivity Informatics Paradox

- Availability of Sample

- Availability of Sample
- Choice of Sample

- Availability of Sample
- Choice of Sample
- Amount of Sample

- Availability of Sample
- Choice of Sample
- Amount of Sample
- Purity of Sample

- Availability of Sample
- Choice of Sample
- Amount of Sample
- Purity of Sample
- Size of sample (for Assembly/Mapping essentially)
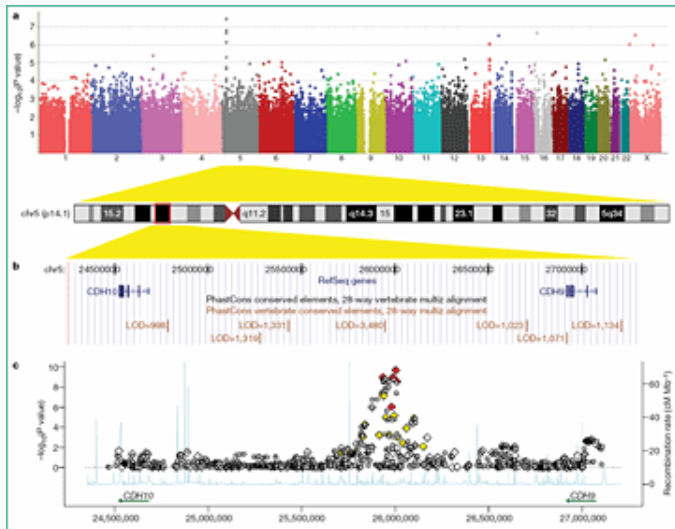
# Applications

- Gene discovery/GWAs

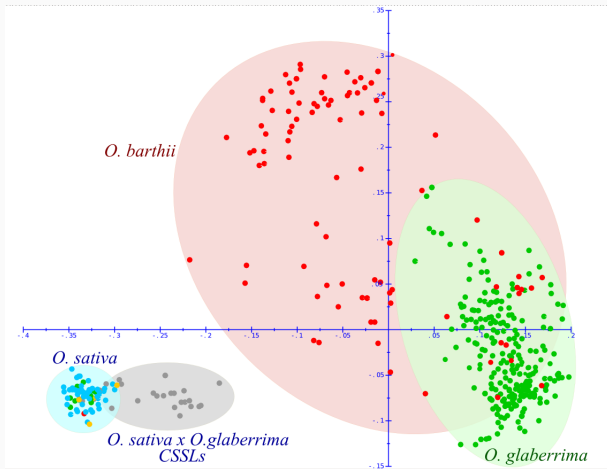- Gene discovery/GWAs
- Species Definition

- Gene discovery/GWAs
- Species Definition
- Subspecies/specific subgroup definition

- Gene discovery/GWAs
- Species Definition
- Subspecies/specific subgroup definition
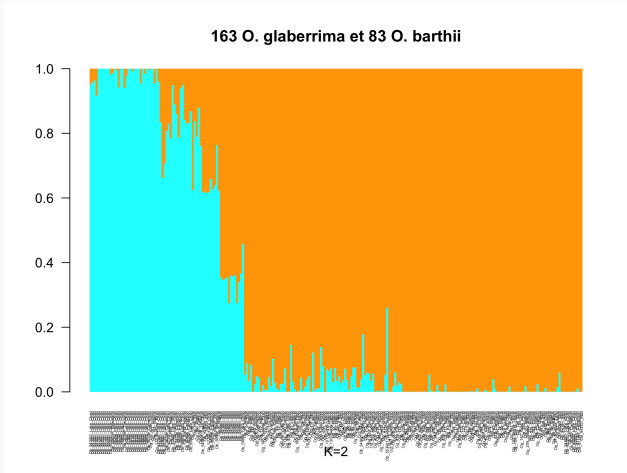- Global genotyping (for breeding in agriculture e.g.)

- Gene discovery/GWAs
- Species Definition
- Subspecies/specific subgroup definition
- Global genotyping (for breeding in agriculture e.g.)
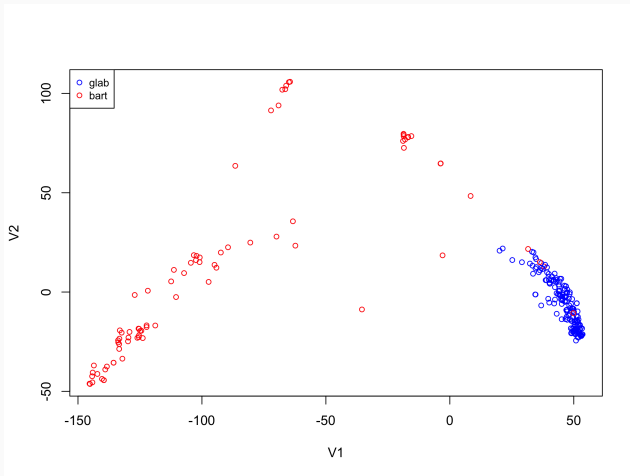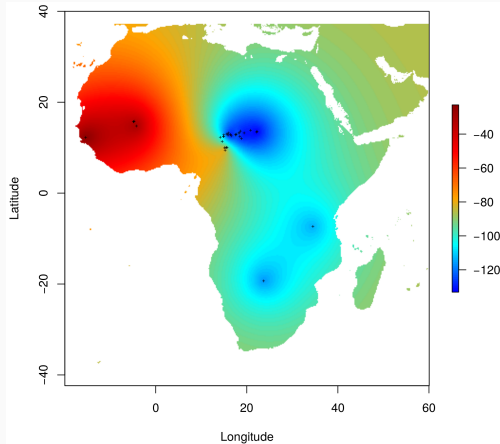- Genomic Ecology (Transposable elements, etc...)
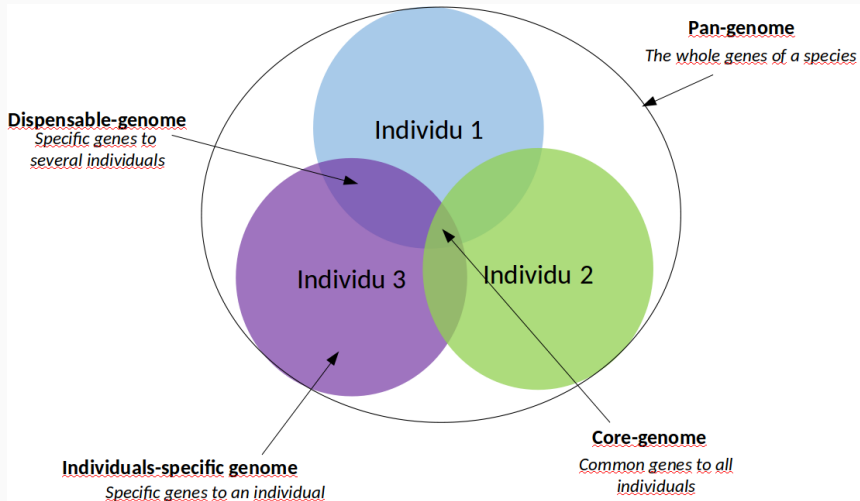
From Orjuela et al, 2014

163 O. glaberrima et 83 O. barthii

K=2

From Cubry et al, 2018

From Cubry et al, 2018

From Cubry et al, 2018

From C. Monat

*O. glaberrima*
current

*O. barthii*
current

86,44 %
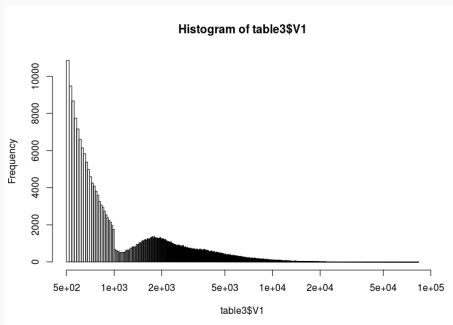
98,15 %

From Monat et al, 2016

# Pangenomic



Chromosome Chr2

From Monat et al, 2018

Legend
O. glaberrima dispensable gene
O. barthii dispensable genome
O. glaberrima core gene
O. barthii core genome

Percentage of genes per window of 10000 bases

Position on the chromosome reported to 1 Mb

# Some really recent results...



Histogram of table3$V1

**Table 3**
Micro-Collinearity Statistics for CG14 vs. TOG5681

| | Valid Scaffolds | Not Valid Scaffolds | Not Referenced Scaffolds |
|---|---|---|---|
| Number of sequences | 48223 | 16672 | 93 |
| Minimal size | 200 | 201 | 202 |
| Maximal size | 86103 | 90835 | 3041 |
| Mean size | 4110 | 6087 | 447 |
| Median size | 1942 | 2592 | 320 |
| Number of functionally annotated gene model | 10685 | 2147 | 2 |
| Number of GO | 23634 | 4817 | 4 |

Sizes are given in bp.

From Monat et al, 2017

34

- Level of expression in different conditions or in different individuals

- Level of expression in different conditions or in different individuals
- Variation in sequences

- Level of expression in different conditions or in different individuals
- Variation in sequences
- Variation of splicing

- Level of expression in different conditions or in different individuals
- Variation in sequences
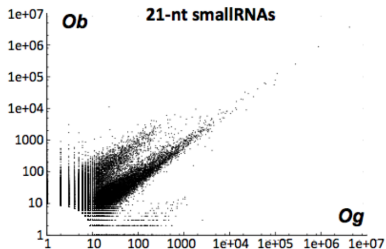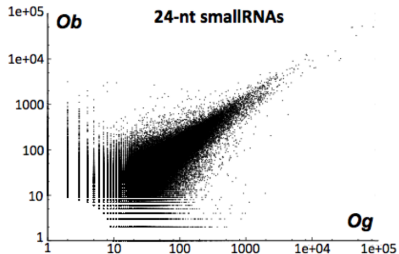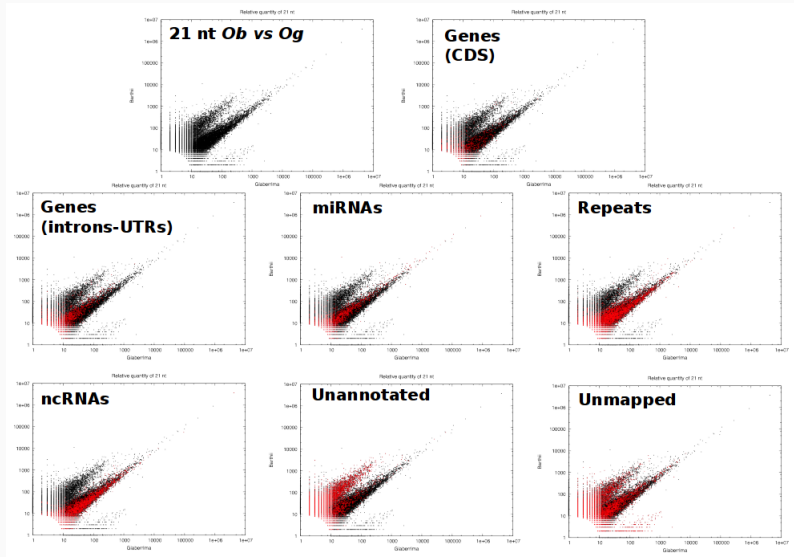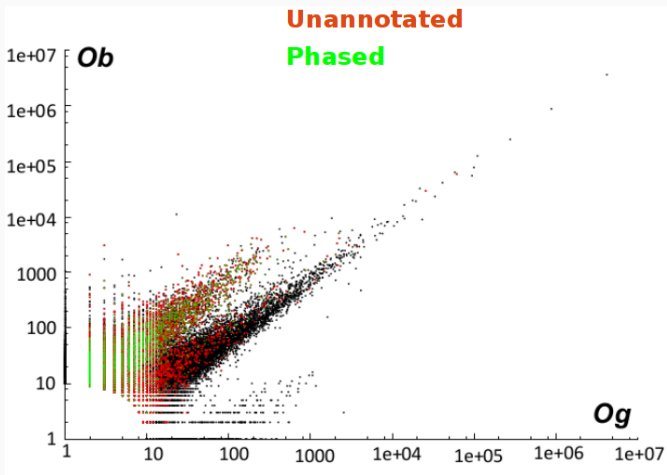- Variation of splicing
- Variation of editing

- Level of expression in different conditions or in different individuals
- Variation in sequences
- Variation of splicing
- Variation of editing
- Detection of putative coding/active sequence

- Level of expression in different conditions or in different individuals

- Level of expression in different conditions or in different individuals
- Variation in sequences

- Level of expression in different conditions or in different individuals
- Variation in sequences
- Variation in specific forms

- Level of expression in different conditions or in different individuals
- Variation in sequences
- Variation in specific forms
- Detection of new forms

From Ta et al, 2015

From Ta et al, 2015

- Pre-diagnostic (Genetic illness, putative resistance)

- Pre-diagnostic (Genetic illness, putative resistance)
- Tumor sequencing

- Pre-diagnostic (Genetic illness, putative resistance)
- Tumor sequencing
- Viral sequencing

- Pre-diagnostic (Genetic illness, putative resistance)
- Tumor sequencing
- Viral sequencing
- Risk Assessement

- Pre-diagnostic (Genetic illness, putative resistance)
- Tumor sequencing
- Viral sequencing
- Risk Assessement
- Epidemiological Studies

## THE METAGENOMICS PROCESS



**Extract all DNA from microbial community in sampled environment**

**DETERMINE WHAT THE GENES ARE**
**(Sequence-based metagenomics)**
- Identify genes and metabolic pathways
- Compare to other communities
- and more…

**DETERMINE WHAT THE GENES DO**
**(Function-based metagenomics)**
- Screen to identify functions of interest, such as vitamin or antibiotic production
- Find the genes that code for functions of interest
- and more…

# Large Metagenomic assays



From Tara Ocean website

From Dinsdale et al, 2008

A Mixture Model for Shotgun Metagenomics

Genome Relative Abundance for Known Genomes
(from Estimated Mixing Parameters)

GRAMMy

Probabilistic Assignment of Reads
(Approximating Component Distributions)

$G$ Reference Genome Set
$q_1$
$q_2$
...
$q_{m-1}$

$r_1$ $R$ Read Set

$r_2$ $r_n$ ... $r_n$

$q_m$

Reference Genome Sequencing

Whole Genome Shotgun Sequencing
(iid sampling from the mixture)

Collective Unknown Genome

$q_{m-1}$
$q_2$
$q_1$
$q_1$
$q_2$
$q_1$
$q_2$
$q_{m-1}$
$q_1$
$q_1$

Environmental Sample
(A Mixture of Genomes)

44

- Real-time Transcriptomics

- Real-time Transcriptomics
- Single-Cell Genomics -> DONE in 2014

- Real-time Transcriptomics
- Single-Cell Genomics -> DONE in 2014
- Single-Cells Transcriptomics (and smallRNA) -> DONE in 2015

**Possibilities in the next 5-10 years (From a presentation in 2013)**

- Real-time Transcriptomics
- Single-Cell Genomics -> DONE in 2014
- Single-Cells Transcriptomics (and smallRNA) -> DONE in 2015
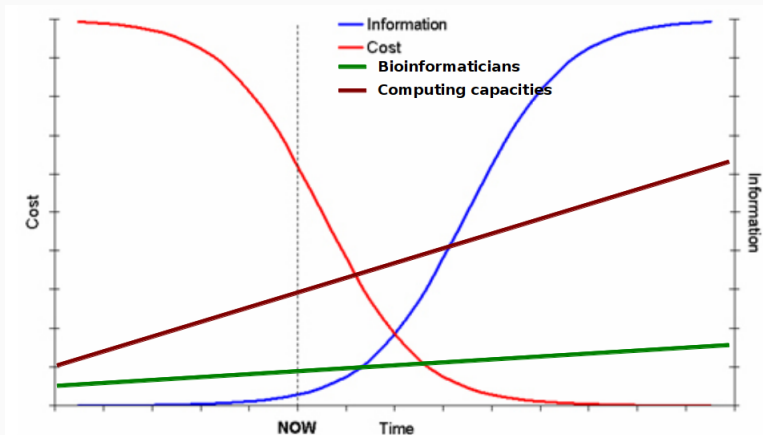- Personal Genomics medicine (ethical problems...) -> Available

- Real-time Transcriptomics
- Single-Cell Genomics -> DONE in 2014
- Single-Cells Transcriptomics (and smallRNA) -> DONE in 2015
- Personal Genomics medicine (ethical problems...) -> Available
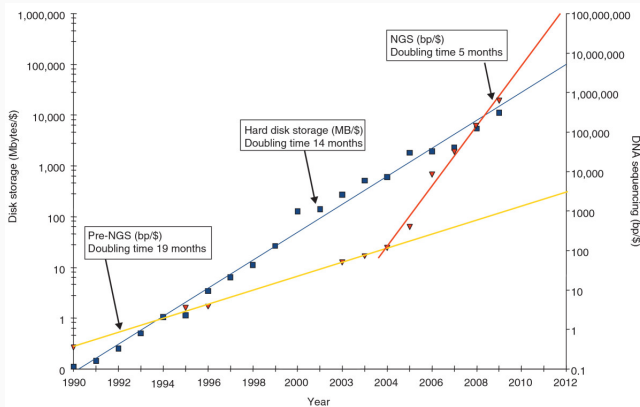- And any new ideas you will have...

- NGS technologies change the way of abording Biology

- NGS technologies change the way of abording Biology
- A lot of Possibilities, a lot of limits

- NGS technologies change the way of abording Biology
- A lot of Possibilities, a lot of limits
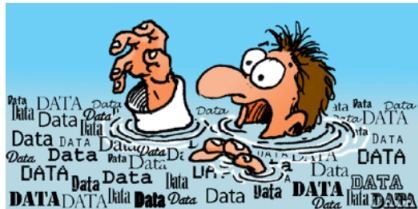- The main limit is no more Sequence, but Sample acquisition and Data treatment
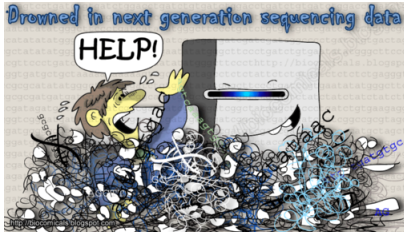
From L. Stein, 2010

# Thanks for your attention