# Data Analysis and Visualization with R

Carlos Jaime Barrios Hernández, PhD

Gilberto Javier Diaz Toro, MSc

MAESTRIA EN MICROBIOLOGIA 2017

# HPC Data Reduction

- Large Data Sets
- Data Reduction
- Collaboration
- High Resolution
- Real Time (almost)

**Big Data + Large Size Data = Great Problem**

"… In the last 5 years, ONLY the astronomy research has produce more of 200PBytes of data/day… we have a lot of data to process… more than all the history of the humanity…"

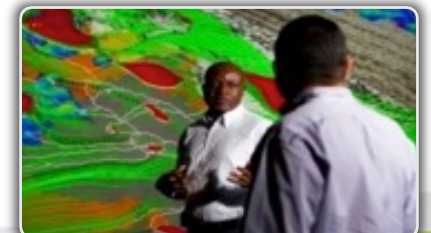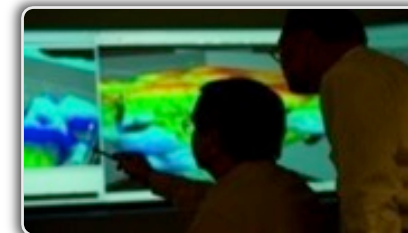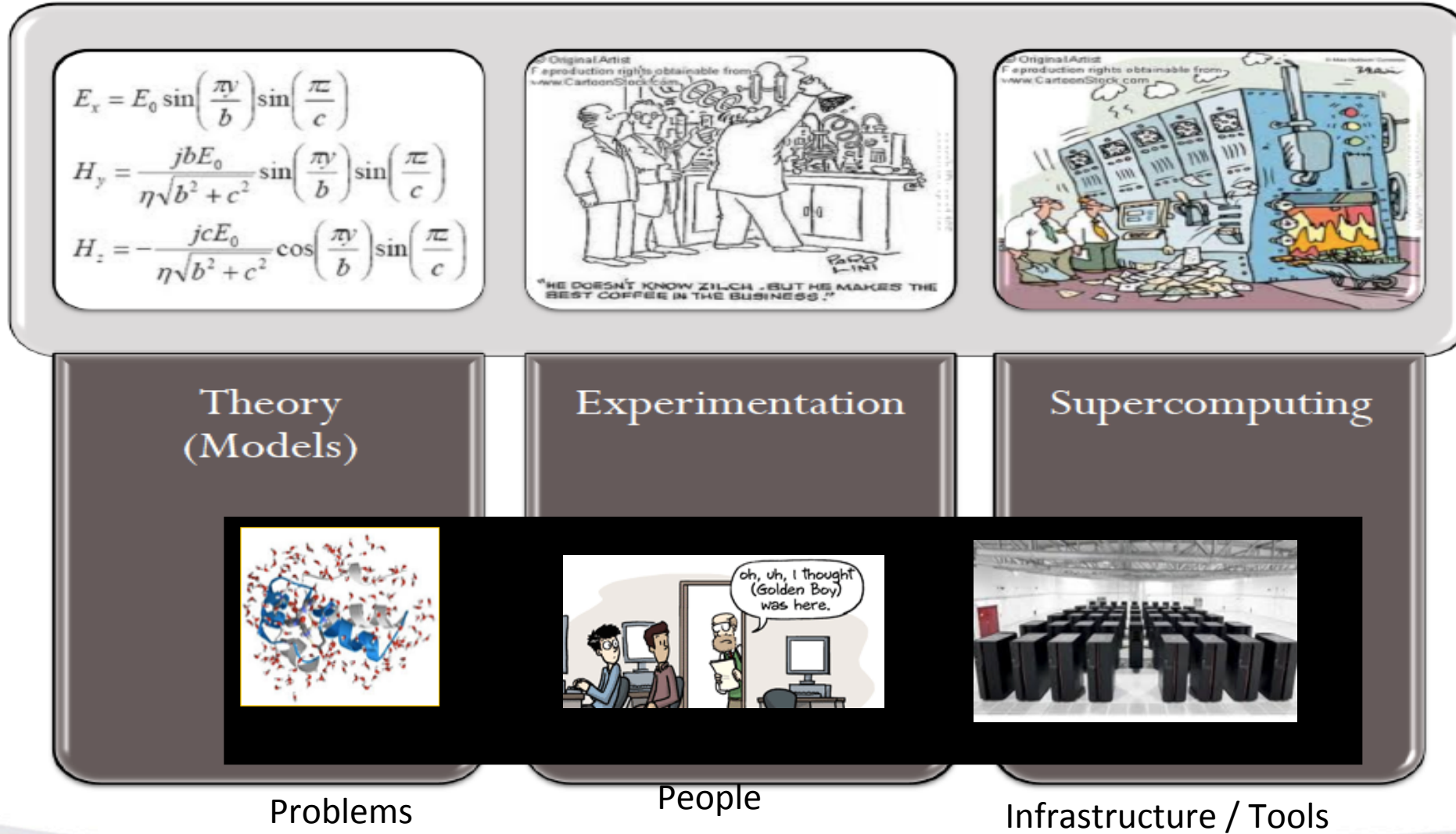"… aja y esto es una oportunidad".

M K

# Three Pillars and Three Actors



Theory (Models)

Experimentation

Supercomputing

Problems

People

Infrastructure / Tools

# The transition to precision medicine



Degree of personalization

| Intuitive<br>Based on trial and error | Identification of probabilistic patterns | Decisions and actions based on knowledge |

**Intuitive Medicine**  **Empirical Medicine**  **Precision Medicine**

Today  Tomorrow

Molecular biomarkers

Genomic biomarkers

**Precision medicine** is based on a better knowledge of phenotype-genotype relationships. That is the knowledge od **disease** and **drug action mechanisms**

With the introduction of **molecular biomarkers** we are living now the **transition** from **intuitive** to **empirical** medicine

From Dopazo

# And how do we identify patterns?
## Using single-gene biomarkers

Most "personalized" therapies are based on this type of biomarkers

From Dopazo

# In spite of its simplicity, Empirical medicine based on biomarkers really works: Increasingly personalized treatments increase patient survival

Super Computación y
Cálculo Científico UIS | High Performance Knowledge

Guatiguará
PARQUE TECNOLÓGICO

Universidad
Industrial de
Santander

# New sequencing technologies change the rules of the game.

*http://www.genome.gov/sequencingcosts/*



With exome sequencing costs ranging 600-800€ and clinical panels below 300€ the use of **NGS for routine diagnosis** matches the price of many other clinical tests.

While **cost** is in continuous **reduction, data volume and complexity increases**

Medicine becomes more and more computational

From Dopazo

# The Spanish "1000 genomes"
## Initiative to sequence rare disease patients

**Sample providers**

**Sequencing platforms**

http://www.gbpa.es/

**Data analysis**

Diseases with
- Unknown genes
- No mutations in known genes

Search for:
- New genes
- Known genes with unknown modifier genes
- Susceptibility genes

**A total of 1044 patients (including 300 controls) of more than 30 diseases were sequenced**

From Dopazo

Data management, analysis and storage = knowledge increase

From Dopazo

# Discovery of new disease genes by whole exome sequencing



**Impressing pace of discovery:**
Only during the last year, the CIBERER initiative has discovered 13 new disease genes and 36 new causative mutations in known genes. An enormous increase of the diagnostic portfolio

From Dopazo

Scaling-up the developments to support the rare diseases initiative to the development of a comprehensive system for diagnosis and gene discovery

From Dopazo

# Implementation of a system for genomic data management in the supercomputing center IT4I (Czech Republic)

This pilot project has been set up in the IT4I supercomputing center, where the genomic data analysis of the country will be centralized.
Obviously not yet challenged with 10M genomes

# Genomic medicine bridges the gap between empirical and precision medicine

**Empirical medicine based on simple biomarkers**

Biomarker

Therapy 1

Therapy 2

Therapy 3

**Genomic Medicine**

Biomarker

Therapy 1

Therapy 2

Therapy 3

Ensayo clínico

Resultado

Systematic genomic analysis enables the **association** of patient **biomarkers** to **therapy results** allowing the immediate application the new knowledge generated, **saving time and costs** and incrementing the **success** of **treatments**.

feedback

From Dopazo

**Big Data: Astronomical or Genomical?**
Stephens ZD. et al., PLOS Biology, July 2015

ABSTRACT - Genomics is a Big Data science and is going to get much bigger, very soon, but it is not known whether the needs of genomics will exceed other Big Data domains. **Projecting to the year 2025**, we compared genomics with three other major generators of Big Data: **astronomy, YouTube, and Twitter**. Our estimates show that **genomics is a "four-headed beast"—it is either on par with or the most demanding of the domains analyzed here in terms of data acquisition, storage, distribution, and analysis**



From Dopazo

# NGS genomic variation data, big and complex

**Logical view** of genomic variation data, real data comes in **different VCF files**.

Each cell represents one specific genotype for one mutation in one sample

Hundreds of millions of mutations, some meta data needed: **Variant annotation**
●Clinical info
●Consequence type
●Conservation scores
●Population frequencies
●...

**Genomics England** project:
200M variants x 100K samples.
About **20 trillion** points with a lot of meta data. About **500-1000TB** to be indexed.

**Samples**

**Genomic Variants**

| | | | | | | |
|---|---|---|---|---|---|---|
| var_1 | A/T | A/A | A/T | T/T | A/A | A/T |
| var_2 | C/C | C/G | C/C | C/G | C/C | G/G |
| .. | .. | .. | .. | .. | .. | .. |
| .. | .. | .. | .. | .. | .. | .. |
| .. | .. | .. | .. | .. | .. | .. |
| .. | .. | .. | .. | .. | .. | .. |
| .. | .. | .. | .. | .. | .. | .. |
| .. | .. | .. | .. | .. | .. | .. |
| .. | .. | .. | .. | .. | .. | .. |
| .. | .. | .. | .. | .. | .. | .. |
| .. | .. | .. | .. | .. | .. | .. |
| .. | .. | .. | .. | .. | .. | .. |
| var_n | .. | .. | .. | .. | .. | .. |

Meta data: **Sample annotation**
●Phenotype
●Family pedigree, Population
●Clinical variables
●...

**Heterogeneous data analysis and algorithms**, different technologies and solutions required:
●Search and filter using data and meta data
●Data mining, correlation
●Statistic tests
●Machine learning
●Interactive analysis
●Network-based analysis
●Visualization
●Encryption
●...

Applications:
●Personalized medicine
●...

From Dopazo

# Simulating drug inhibition

"Ideal" KO of EGFR affects 7 pathways

Real inhibition with Afanatib affects 11 pathways

Inhibition with broader spectrum Trastuzumab affects 13 pathways

From Dopazo

# Omics views of genomes

Genetic variation                Epigenetic variation      Expression variation



SNPs, loss-of-heterozygosity          DNA methylation              RNA expression
Copy number variants                     Chromatin                    Gene structure

# HPC Data Analysis, Simulation.. And Visualization

- Scientists need visualize simulations and data
- HPC platforms are necessary
  - Terascale/Petascale/Exascale
    - Energy Efficiency
    - Data Movement
    - Programmability
  - Workflows
  - Hybrid Computing
- Relevance in Simulations
  - Smart
  - Ultrascale

ORNL Jaguar

Julich JUGene

SC3UIS GUANE-1

# Bioinformatics Tools – Open source

Programming languages

R, Perl, Bash scripting (Linux), MySQL, Apache, PHP, Python, Java, …

Software, e.g. Bioconductor, BioPerl, Ensembl Perl API, Bowtie, BWA, Velvet, Varscan, Rmap, …

Alignment, analysis of next-generation sequencing and microarray data

Web browsers, e.g. UCSC, Ensembl

visualize data in relation to genome features

Gene Ontology, e.g. DAVID

functional annotation and enrichment

# And... why R?

R is powerfull, worldwide used and open source.

We can exploit easy HPC architectures (almost).

This course introduces some relatively new additions to the R programming language: advanced reduction and visualization.  R packagess provide a powerful toolkit to make the process of manipulating and visualising data easy and intuitive, in this case for microbiology.

# Schedule and Topics

1. Introduction Course (Today)
2. Some Topics about Algorithms and Platform uses
   1. Personal Installation
   2. HPC UIS Platform Use
3. Introduction to R – Data Structures
4. Writing Analysis workflows with R
5. Summarizing and Combining Data
6. Plotting and Visualization

# Goals of the Course

Participants will gain practical experience and skills to be able to:

- Introduce to scripting computational languages, in this case R
- Meet the challenges of data handling and reduction;
- Introduce to the use of R syntax, functions and packages;
- Understand best practices for scientific computational work.
- Introduce to use visualization tools
- Introduce to use HPC platforms in collaborative environments

# Information on line and Bibliography

- www.sc3.uis.edu.co
  - http://wiki.sc3.uis.edu.co/
    - http://wiki.sc3.uis.edu.co/index.php/An%C3%A1lisis_y_Visualizaci%C3%B3n_de_Datos_con_R
- https://www.r-project.org/
- https://journal.r-project.org/

- https://bioinformatics.ca/workshops/2017/introduction-r-2017
- http://bioinformatics-core-shared-training.github.io/
- https://training.csx.cam.ac.uk/bioinformatics/course/bioinfo-intR
- http://bioinfotraining.bio.cam.ac.uk/postgraduate/specialized/bioinfo-intR
- … and other to find in the first link.

Super Computación y
Cálculo Científico UIS | High Performance Knowledge

Guatiguará
PARQUE TECNOLÓGICO

Universidad
Industrial de
Santander

## Something about US



**Carlos Jaime Barrios Hernandez, PhD. (@carlosjaimebh cbarrios@uis.edu.co )**
PhD in Computer Science (Nice, France) , MSc in Applied Mathematics and Computer Science (Grenoble, France) and Systems Engineering (Bucaramanga, Colombia)
Director of SC3UIS, Assistant Professor UIS.



**Gilberto Javier Diaz Toro ( @gilbertodiazt gilberto.diaz@uis.edu.co )**
MSc in Computer Science (Mérida, Venezuela) Systems and Computing Engineering (Mérida, Venezuela).
CTO of SC3UIS, International Instructor and Professor on HPC and Scientific Computing.

# GUANE-1 and Yajé

- **GUANE-1 Reload**
  - High Density and Green HPE-HPC Platform
  - 128 NVIDIA M2070 TESLA GPUs
  - 32 Intel Xeon E5645 2.4 GHz Processors
  - 3 High Bandwth Networks
  - 1.6 TB RAM
  - General Purpose Platform
- **YAJE**
  - HPE ML150 G9 Development Platform
  - Intel Xeon E52609 1.9 GHz Processor
  - 64GB RAM
  - NVIDIA GRID K2
    - 1536 GPU Cores
    - 2 x 8 GB Memory

Parnertship

« Quien no computa, no compite »

Mateo Valero, PhD.
Director BSC-CNS, Spain

Soon Colombia Advanced Computer Center
Follow us: @sc3uis
Or visit us: www.sc3.uis.edu.co
Thanks!