# Computational and Statistical Methods
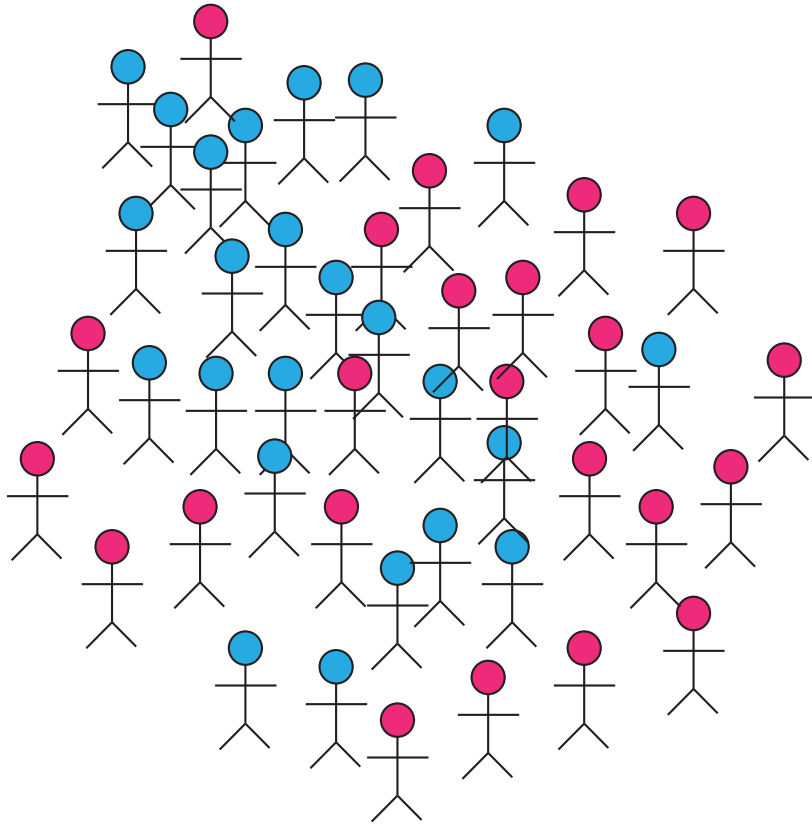# for High-Throughput Genomics

*Héctor Corrada Bravo*
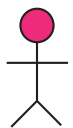Dept. of Computer Science
Center for Bioinformatics and Computational Biology
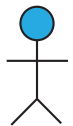University of Maryland

# Computational Genomics



- Study the **molecular** basis of *variation* in development and disease

- Using **high-throughput** experimental methods

  - statistical learning

  - visualization
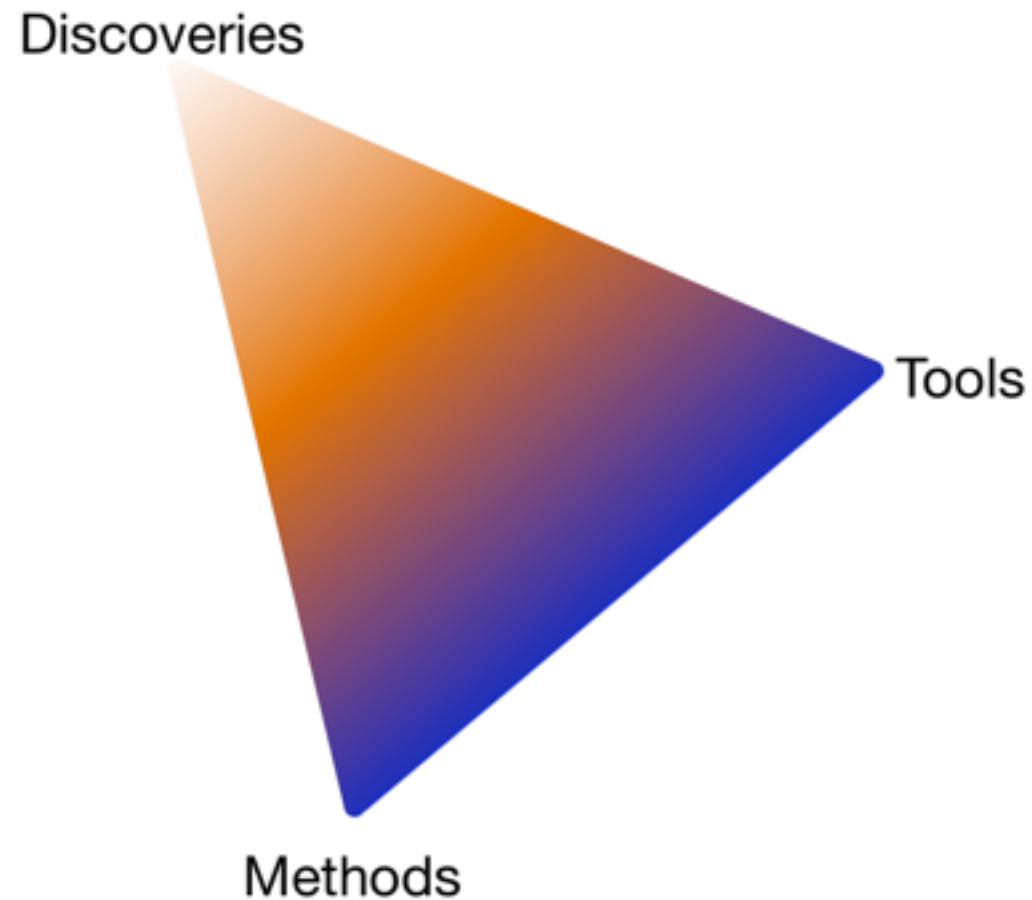
  - algorithms
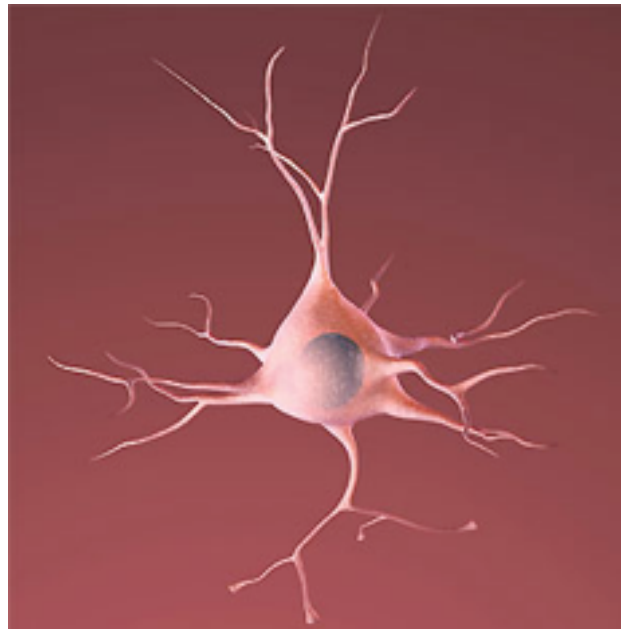
  - data management

colon cancer

not

Hector Corrada Bravo

UNIVERSITY OF MARYLAND

# Computational Genomics

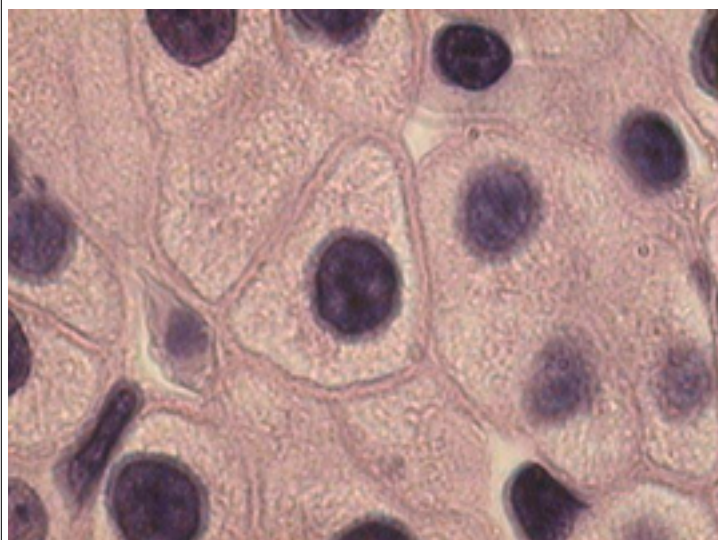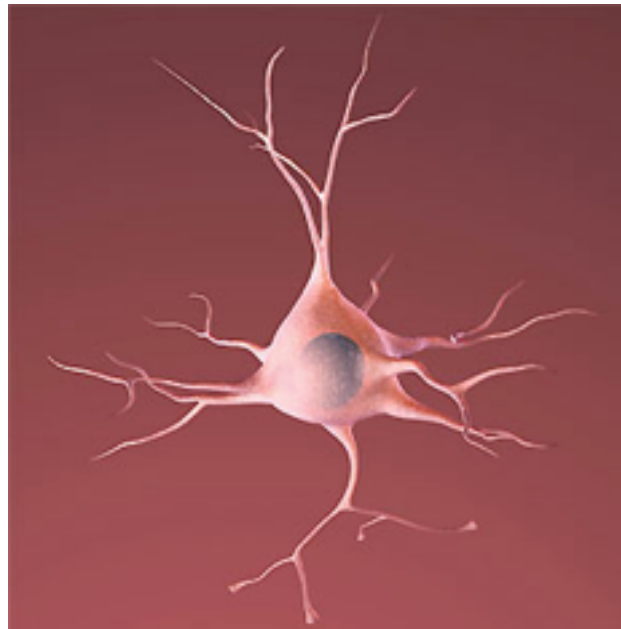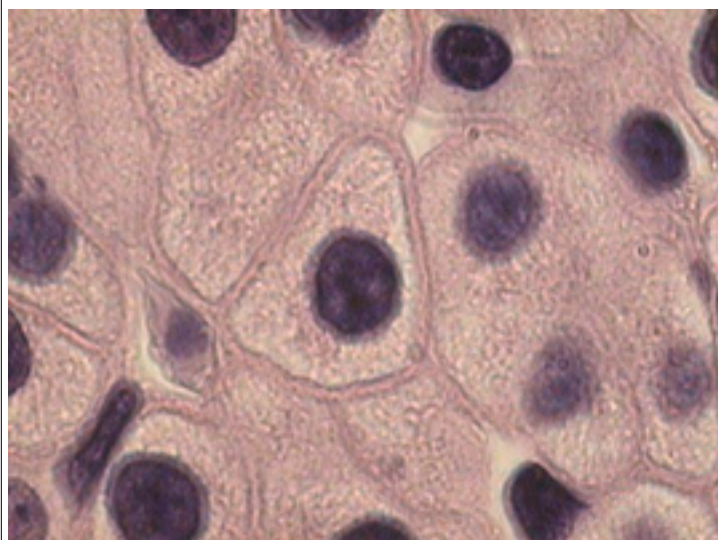Hector Corrada Bravo

# Differentiation

Different genes are **expressed** during different **stages** and in different **tissues**

# Differentiation and disease

### Behavior of these genes in age and age-related disease

# Computational Epigenetics



- DNA methylation

  - replicates after cell division

  - modified in **disease**

# Computational Epigenetics

1. Novel computational analysis methods
2. First project to measure and characterize genome-wide in cancer

*Big takeaway:*
   *Hyper-variability* of specific genes involved in *differentiation* is a stable characteristic of cancer
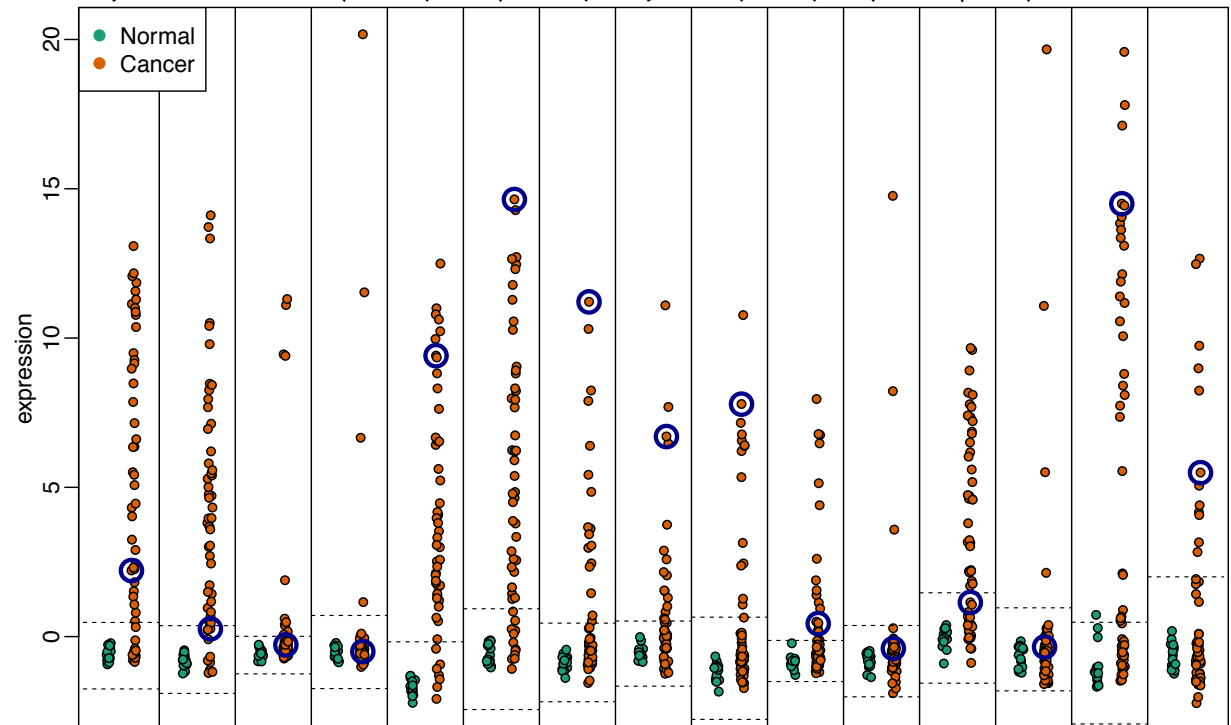
# Anti-profiles

1. Understanding *hyper-variability*
2. How do we predict using *hyper-variability*?
   - Anomaly classification
   - Robustness of predictors

Héctor Corrada Bravo[1,*], Vasyl Pihur[2], Matthew McCall[3], Rafael A Irizarry[2] and Jeffrey T Leek[2]

* Corresponding author: Héctor Corrada Bravo
hcorrada@umiacs.umd.edu
• Author Affiliations

[1] Department of Computer Science, Center for Bioinformatics and Computational Biology, University of Maryland, College Park, MD, USA

[2] Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, USA

[3] Department of Biostatistics and Computational Biology, University of Rochester Medical Center, Rochester, NY, USA

For all author emails, please log on.

JOURNAL
## Cancer Informatics
Journal Analytics

Contents | About | Call for Papers | Editor in Chief | Editorial Board

## Gene Expression Signatures Based on Variability can Robustly Predict Tumor Progr...

Wikum Dina...

### Nucleic Acids Research
ABOUT THIS JOURNAL | CONTACT THIS JOURNAL | SUBSCRIPTIONS | CUR...

Oxford Journals › Life Sciences › Nucleic Acids Research › Advance Access › ...

## Determinants of expression variability

Elfalem Y. Alemu, Joseph W. Carl Jr, Héctor Corrada Bravo* and Sridhar Hannenhalli*
+ Author Affiliations

*To whom correspondence should be addressed. Tel: +1 301 405 8219; Fax: +1 301 314 1410; Email: sridhar@umiacs.umd.edu

# Epigenomic Heterogeneity



**Heterogeneous Celltype Composition**

**Celltype-Specific Methylation Patterns**

● Methylated CpG          ○ Unmethylated CpG

CG  CG  CG  CG  CG  CG  CG  CG  CG

**Aligned Reads**

**Region Graph**

no. reads assigned

33
25
16
8
0

CG  CG  CG  CG  CG  CG  CG  CG  CG

# Interactive Visualization



http://epiviz.cbcb.umd.edu

# Host-Pathogen Systems

**Microbiome:** *environmental sequencing*

- Detecting microbes associated with disease
- Modeling and detecting differences in microbial community composition in health and disease

Differential abundance analysis for microbial marker-gene surveys

Joseph N Paulson, O C

Genome **Biology** IMPACT FACTOR 10.8  Highly accessed  Open Access

**Diarrhea in young children from low-income countries leads to large-scale alterations in intestinal microbiota composition**

Mihai Pop[1], Alan W Walker[2], Joseph Paulson[1], Brianna Lindsay[3], Martin Antonio[4], M Anowar Hossain[5], Joseph Oundo[6], Boubou Tamboura[7], Volker Mai[8], Irina Astrovskaya[1], Hector Corrada Bravo[1], Richard Rance[2], Mark Stares[2], Myron M Levine[3], Sandra Panchalingam[3], Karen Kotloff[3], Usman N Ikumapayi[4], Chinelo Ebruke[4], Mitchell Adeyemi[5], Dilruba Ahmed[5], Firoz Ahmed[5], Meer Taifur Alam[5], Ruhul Amin[5], Sabbir Siddiqui[5], John B Ochieng[6], Emmanuel Ouma[6], Jane Juma[6], Euince Mailu[6], Richard Omore[6], J Glenn Morris[8], [illegible] Breiman[?], Debasish Saha[4], Julian Parkhill[2], James P Nataro[10] and O Colin Stine[3]

***Joint transcriptome profiling:*** Understanding host-pathogen interaction *throughout* course of infection

**Transcriptomic profiling of gene expression and RNA processing during *Leishmania major* differentiation**

Laura A. L. Dillon[1,2], Kwame Okrah[3], V. Keith Hughitt[1,2], Rahul Suresh[1], Yuan Li[1,2], Maria Cecilia Fernandes[1,2], A. Trey Belew[1,2], Hector Corrada Bravo[2,4], David M. Mosser[1] and Najib M. El-Sayed[1,2,*]
+ Author Affiliations

# NHGRI strategic plan

- What does the NIH think genomics should be for the next 10 years?

## PERSPECTIVE

# Charting a course for genomic medicine from base pairs to bedside

Eric D. Green[1], Mark S. Guyer[1] & National Human Genome Research Institute*

There has been much progress in genomics in the ten years since a draft sequence of the human genome was published. Opportunities for understanding health and disease are now unprecedented, as advances in genomics are harnessed to obtain robust foundational knowledge about the structure and function of the human genome and about the genetic contributions to human health and disease. Here we articulate a 2011 vision for the future of genomics research and describe the path towards an era of genomic medicine.

[Nature, Feb. 2011]

# Where do we fit in?

- The major bottleneck in genome sequencing is no longer data generation—the computational challenges around data analysis, display and integration are now rate limiting. New approaches and methods are required to meet these challenges.

- **Data analysis**
  - Computational tools are quickly becoming inadequate for analysing the amount of genomic data that can now be generated, and this mismatch will worsen. Innovative approaches to analysis, involving close coupling with data production, are essential.

- **Data integration**
  - Genomics projects increasingly produce disparate data types (for example, molecular, phenotypic, environmental and clinical), so computational approaches must not only keep pace with the volume of genomic data, but also their complexity. New integrative methods for analysis and for building predictive models are needed.

- **Visualization**
  - In the past, visualizing genomic data involved indexing to the one-dimensional representation of a genome. New visualization tools will need to accommodate the multidimensional data from studies of molecular phenotypes in different cells and tissues, physiological states and developmental time. Such tools must also incorporate non-molecular data, such as phenotypes and environmental exposures. The new tools will need to accommodate the scale of the data to deliver information rapidly and efficiently.

- **Computational tools and infrastructure**
  - Generally applicable tools are needed in the form of robust, well-engineered software that meets the distinct needs of genomic and non-genomic scientists. Adequate computational infrastructure is also needed, including sufficient storage and processing capacity to accommodate and analyse large, complex data sets (including metadata) deposited in stable and accessible repositories, and to provide consolidated views of many data types, all within a framework that addresses privacy concerns. Ideally, multiple solutions should be developed[105].

# Where do we fit in?

- Meeting the computational challenges for genomics requires scientists with expertise in biology as well as in informatics, computer science, mathematics, statistics and/or engineering.

- A new generation of investigators who are proficient in two or more of these fields must be trained and supported.

# Education and Training

- PhDs: Computer Science, Applied Statistics, Scientific Computation, Computational Biology

  - MIT/Broad, Harvard/Dana Farber Cancer Center, U. of Chicago, Genentech, Johns Hopkins Medicine, Dow Jones Data Science

- New UG courses in Bioinformatics and Data Science