# Systems Immunology & Immune Oncology
## A Data-Centric View
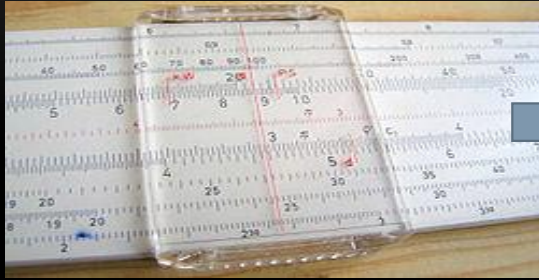
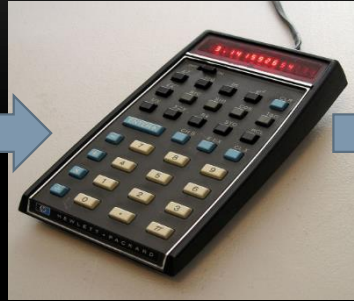**Magnus Fontes**

https://pics-about-space.com/pics-new-world-nasa?p=2

# Disclosures & Affiliations

- General Manager of Institut Roche, France https://institut.roche.com/
- Adjunct professor of mathematics, Lund University, Sweden https://portal.research.lu.se/en/persons/magnus-fontes
- Co-founder of the bioinformatics software company Qlucore www.qlucore.com

# A brief personal history of computing


Slide rule


Hp-35 (1972)


ABC80 (1978)
16 kB RAM


Sparc Station
Beginning of the 90's
Up to 500 MB in RAM



?

Nearly a century-and-a-half ago, Louis Pasteur envisioned a world where "Institut Pasteur" scientists across the globe would share their research and knowledge in service of humanity...

Pasteur Global Health Genomics Center
https://vimeo.com/171747507

...standing united "on the edge of (global) mysteries" and aspiring to "lift the veil" on the origin of diseases.

**Digitalized Biomedicine and Healthcare**
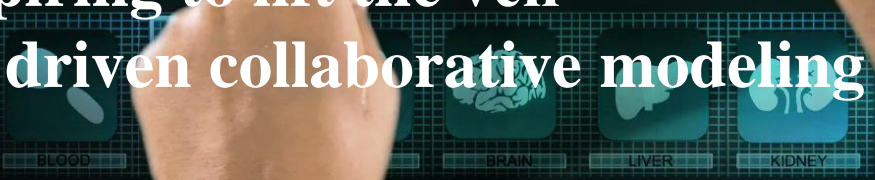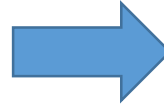
The Human Machine Partnership has placed us in a position to deliver even more on the vision of Louis Pasteur "standing united and aspiring to lift the veil" through Data driven collaborative modeling
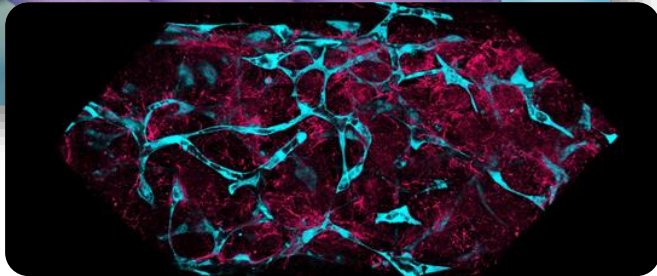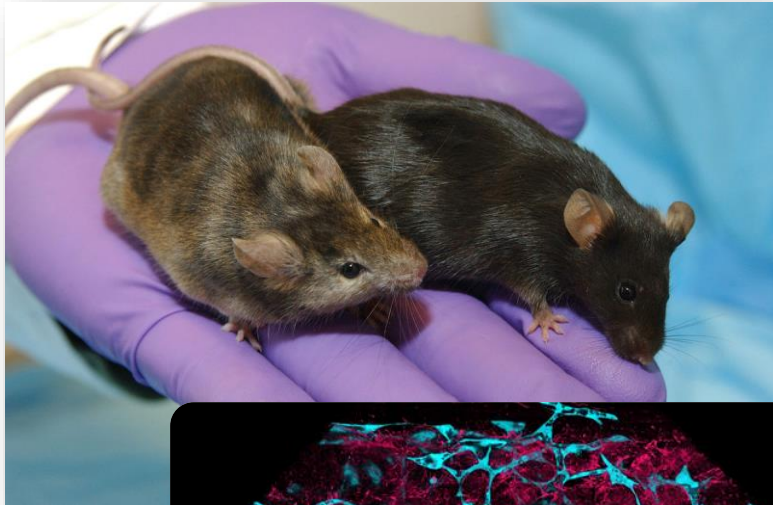
# Multidisciplinary & Collaborative research through Human-Machine Partnerships



- **Individual Human Knowledge → Connected Global Knowledge**
- **Legal & Ethical & Psychological considerations**
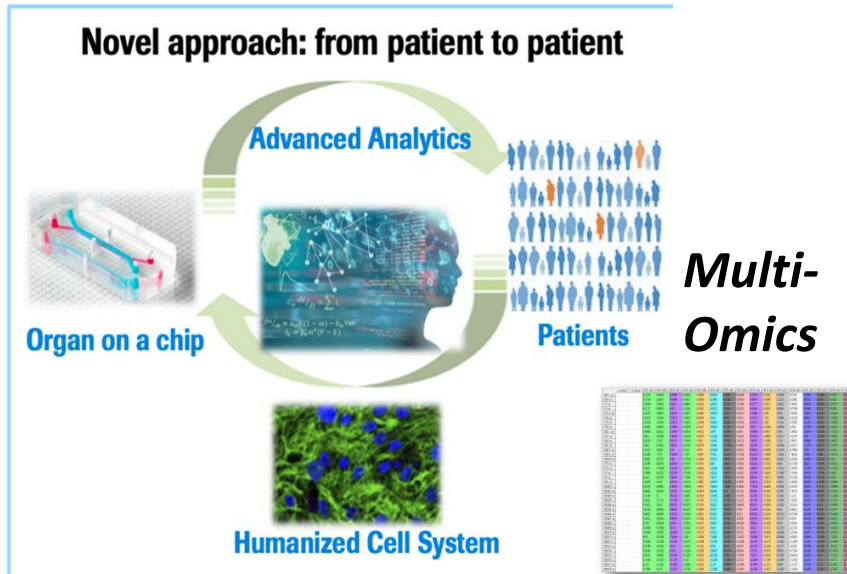- **How do we use our new tools? Repurposing Generative AI?**

# Do you have a model for that?



$$\frac{dx}{dt} = \alpha x - \beta xy,$$
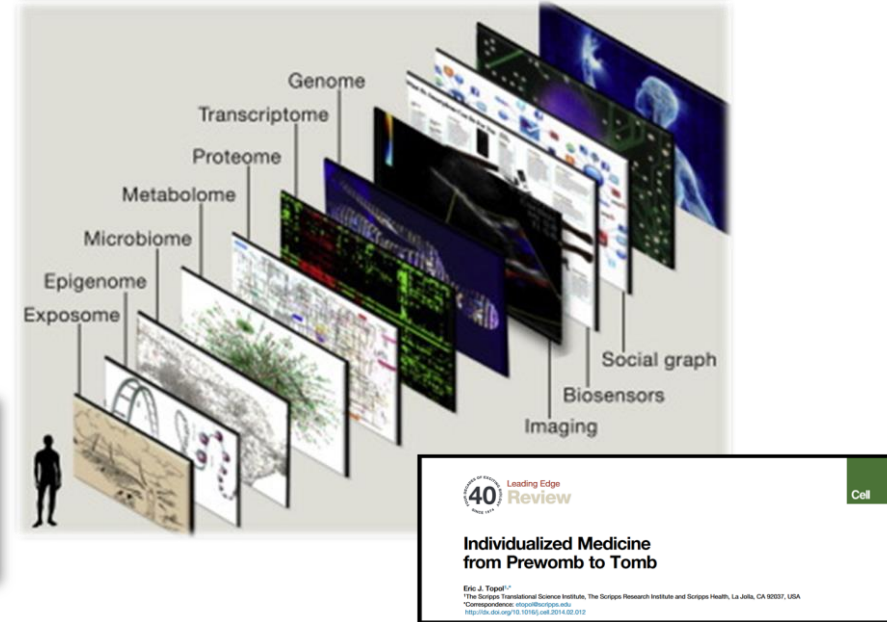$$\frac{dy}{dt} = \delta xy - \gamma y,$$

# Transformation of Biomedical R&D

- Reverse Translational Research
- Exvivo  biological systems – Biological Avatars
- Insilico Systems – Digital Avatars
- Advanced Analytics through AI and Human-Machine partnerships

# Transformation of Healthcare:

- Longitudinal and deep precision sampling
- Early detection and diagnosis
- Treating early disease
- Personalization & Combination therapies
- Clinical decision support via Personalized Biological + Digital Avatars



Novel approach: from patient to patient

Advanced Analytics

Organ on a chip

Humanized Cell System

Patients

*Multi-Omics*



Genome
Transcriptome
Proteome
Metabolome
Microbiome
Epigenome
Exposome
Social graph
Biosensors
Imaging

Individualized Medicine from Prewomb to Tomb

Eric J. Topol

# Biological Avatars – Exvivo Systems → Modeling holistic dynamics



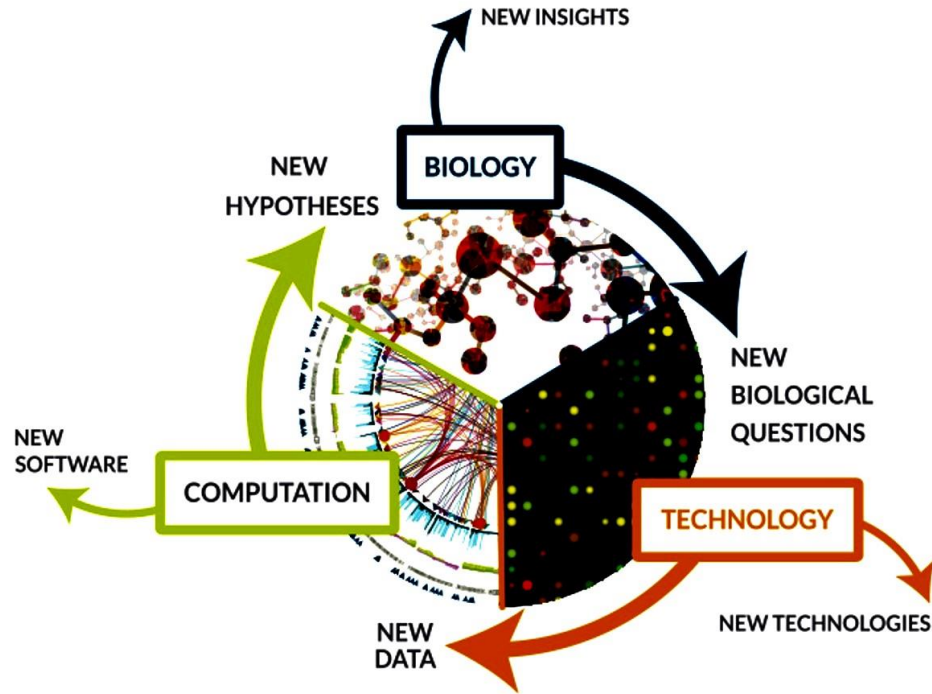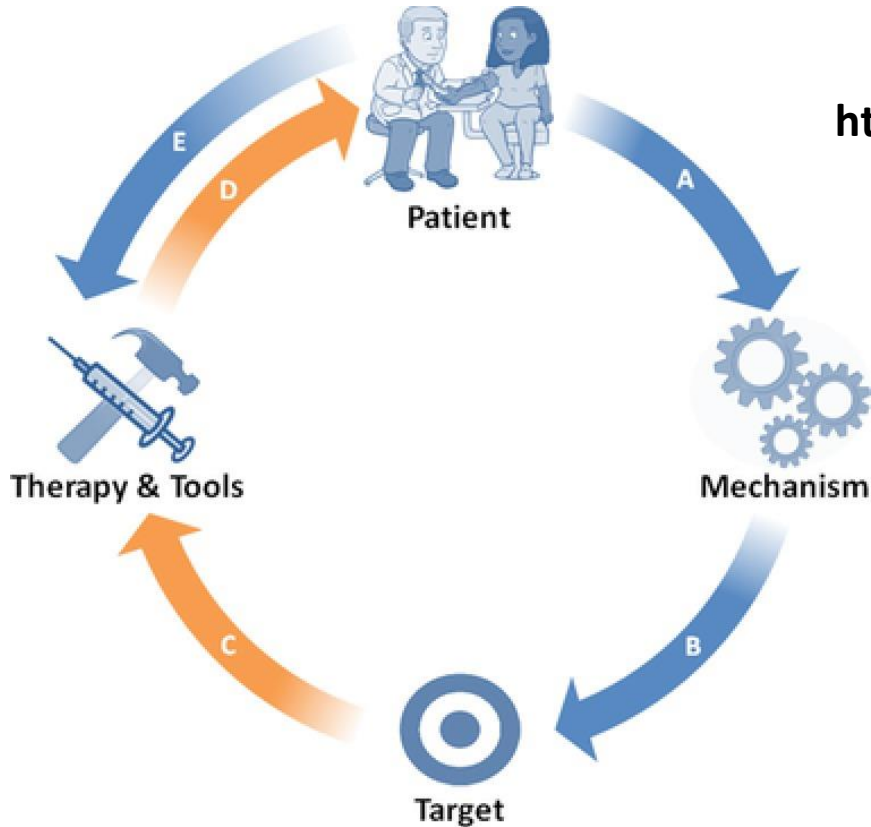| | 2D cell culture | C.elegans | D. melanogaster | D. rerio | M. musculus | PDX | Human organoids |
|---|---|---|---|---|---|---|---|
| Ease of establishing system | ✓/✗ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Ease of maintenance | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Recapitulation of developmental biology | ✗ | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ |
| Duration of experiments | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Genetic manipulation | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ |
| Genome-wide screening | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | ✓ |
| Physiological complexity | ✗ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Relative cost | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Recapitulation of human physiology | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

✓ Best   ✓ Good   ✓ Partly suitable   ✗ Not suitable

Kim, J., Koo, BK. & Knoblich, J.A. Human organoids: model systems for human biology and medicine. Nat Rev Mol Cell Biol 21, 571–584 (2020). https://doi.org/10.1038/s41580-020-0259-3

# Digital Avatars -- Computational Systems → Predicitive modeling

# Reverse Translational Resarch – « *The ground truth* »

**The ultimate model for human health and disease is human health and disease**



## https://www.gene.com/stories/reverse-translation

*"In a clinical trial of one of our cancer immunotherapies, our scientists observed differences in people who responded to the medicine. By digging deeper into the biology of the non-responders, the team discovered that some of them had an up-regulated gene signature associated with a protein called TGF-beta. With this knowledge, they then explored the phenomenon in pre-clinical models. When they combined the cancer immunotherapy with an investigational antibody that blocks TGF-beta, it resulted in improved anti-tumor activity in pre-clinical models that mimic the biology of some non-responders."*

# Computational Mathematical Modeling through Compartmentalized Integrative Grey Box Modeling
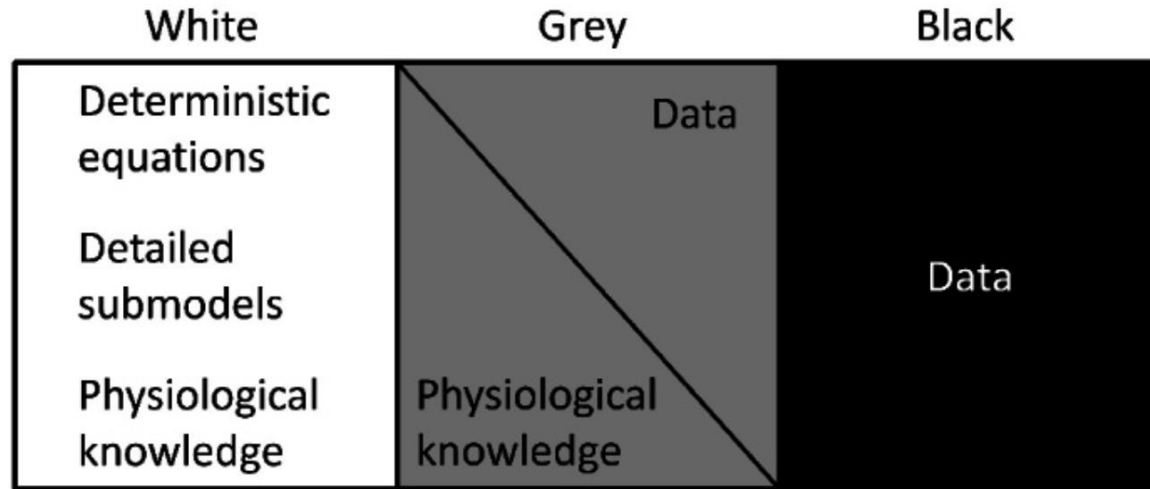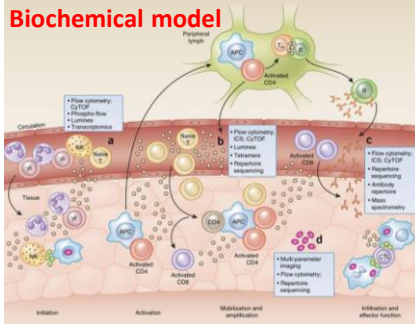


Illustration of the concept of grey-box modeling. White-box models are based mainly on knowledge about the system. Black-box models are built on statistical information from the data. Grey-box modeling combines the two approaches.

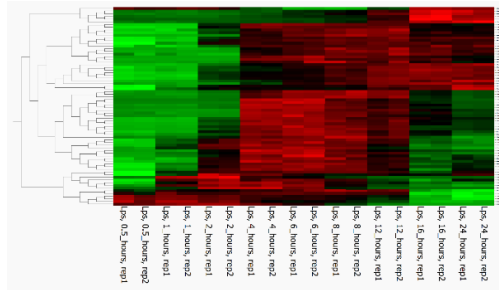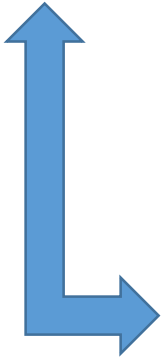The Grey Box Modeling Loop
Connecting Data & Models

Biochemical model
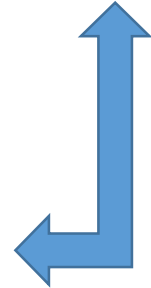
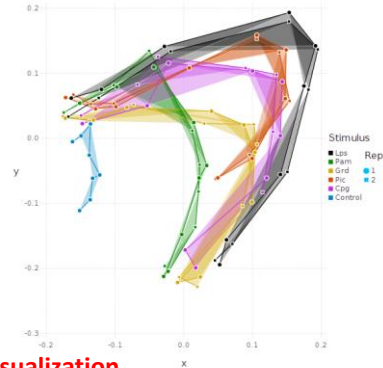Mechanistic mathematical model

$$\frac{dx}{dt} = \alpha x - \beta x y,$$

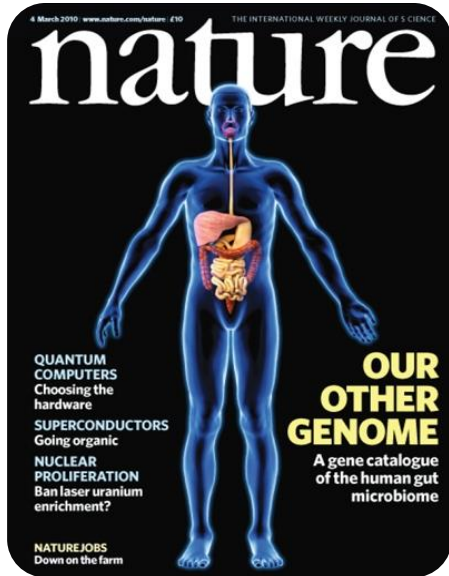$$\frac{dy}{dt} = \delta x y - \gamma y,$$

"A Systems Biology Workbench"
A holistic engine for scientific discovery and innovation

Statistical learning and visualization

# The human model system



**Nature March 2010**

**Nature June 2012**

- $3 \times 10^{13}$ human cells
- As many bacterial cells
- Around 10 times as many viruses
- Around $10^{10}$ proteins per human cell

A very dynamical system…

Check for updates

# The distribution of cellular turnover in the human body

Ron Sender and Ron Milo ✉

We integrated ubiquity, mass and lifespan of all major cell types to achieve a comprehensive quantitative description of cellular turnover. We found a total cellular mass turnover of $80 \pm 20$ grams per day, dominated by blood cells and gut epithelial cells. In terms of cell numbers, close to 90% of the $(0.33 \pm 0.02) \times 10^{12}$ cells per day turnover was blood cells.

To better understand the function of the human body in health and disease, it is of major interest to quantify its cellular compo-

cells comprising the human body[3] or ones with an especially fast turnover of $\tau < 10$ d.
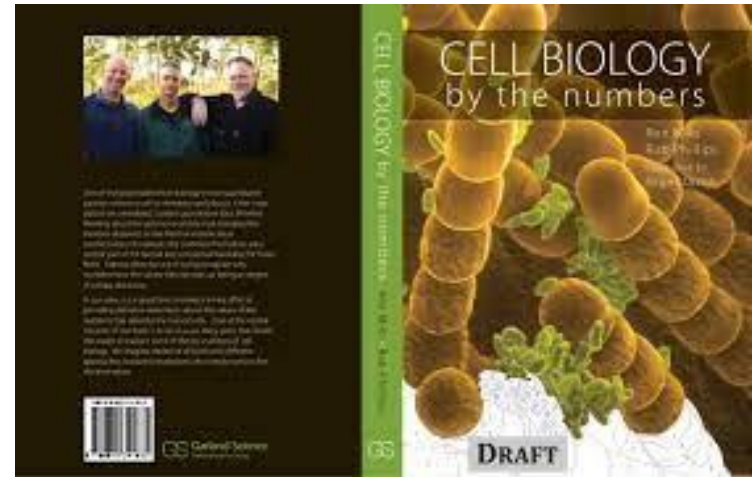
We analyzed many of the tissues thought to be relevant and found them to make a negligible contribution in terms of both number and mass (Supplementary Tables 1–4); for example, sperm cells, kidney cells and osteocytes. For cell types with a short lifespan, we revised earlier estimates[1,3] of the total cell number, as documented in Methods. Figure 1a presents for each of these cell types, the number

**Approximately 4 million human cells die per second in an average human**

| cell type | turnover time | BNID |
|---|---|---|
| small intestine epithelium | 2-4 days | 107812, 109231 |
| stomach | 2-9 days | 101940 |
| blood Neutrophils | 1-5 days | 101940 |
| white blood cells Eosinophils | 2-5 days | 109901, 109902 |
| gastrointestinal colon crypt cells | 3-4 days | 107812 |
| cervix | 6 days | 110321 |
| lungs alveoli | 8 days | 101940 |
| tongue taste buds (rat) | 10 days | 111427 |
| platelets | 10 days | 111407,111408 |
| bone osteoclasts | 2 weeks | 109906 |
| intestine Paneth cells | 20 days | 107812 |
| skin epidermis cells | 10-30 days | 109214, 109215 |
| pancreas beta cells (rat) | 20-50 days | 109228 |
| blood B cells (mouse) | 4-7 weeks | 107910 |
| trachea | 1-2 months | 101940 |
| hematopoietic stem cells | 2 months | 109232 |
| sperm (male gametes) | 2 months | 110319, 110320 |
| bone osteoblasts | 3 months | 109907 |
| red blood cells | 4 months | 101706, 107875 |
| liver hepatocyte cells | 0.5-1 year | 109233 |
| fat cells | 8 years | 103455 |
| cardiomyocytes | 0.5-10% per year | 107076, 107077, 107078 |
| central nervous system | life time | 101940 |
| skeleton | 10% per year | 109908 |
| lens cells | life time | 109840 |
| oocytes (female gametes) | life time | 111451 |

# Bionumbers, e.g. Lifespan of a cell:

**https://bionumbers.hms.harvard.edu/search.aspx**



**http://book.bionumbers.org/**

# Some (log) length scales of Life



http://www.bates.edu/gould-research-lab

# A view of our hierarchically organized system
## Organizational levels of life:

- …
- **Molecules** ←→
- **Cells** ←→
- **Tissues** ←→
- **Organs** ←→
- **Organ complexes** ←→
- **Organisms**
- **Populations**
- ….



**The focus depends on « *the target application* »**

# Focusing on Cells – The basic units of Life

On June 15th 2023 Aviv Regev received the L'Oreal-Unesco for women in science award *"for her pioneering work applying mathematics and computer science to revolutionize cell biology."*

**https://www.humancellatlas.org/**

**https://data.humancellatlas.org/**

# Cells are specialized and adapted to their neighborhood

# The Human Immune System

# The Immunity State Space

*The Immunity State Space* is constructed from « precision measurements » of molecular, cellular and higher order biological activities connected with the Immune System.
The goal is to first identify domains of health and non-health and then find interventions to push patients back to healthy states.

# Computational modeling experience tells us that « *Life operates in low dimension* »

*A typical example: Invitro stimulated mouse dendritic cells 0-0.5-2-4-6-8-12-16- 24h bulk mRNA expression for all 10716 genes (two samples per time point) result in « clear » 1 dimensional trajectories embedded in low (<7) dimensional space:*

# Where to find and deposit biomedical data?



https://www.nature.com/sdata/policies/repositories

# https://www.ebi.ac.uk/



# https://www.ebi.ac.uk/services

# https://www.ncbi.nlm.nih.gov/

Search

Advanced

Help

## SRA - Now available on the cloud

Sequence Read Archive (SRA) data, available through multiple cloud providers and NCBI servers, is the largest publicly available repository of high throughput sequencing data. The archive accepts data from all branches of life as well as metagenomic and environmental surveys. SRA stores raw sequencing data and alignment information to enhance reproducibility and facilitate new discoveries through data analysis.

### Getting Started

How to Submit

How to search and download

How to use SRA in the cloud

Submit to SRA

### Tools and Software

Download SRA Toolkit

SRA Toolkit Documentation

SRA-BLAST

SRA Run Browser

SRA Run Selector

### Related Resources

Submission Portal

dbGaP Home

BioProject

BioSample

**https://www.ncbi.nlm.nih.gov/sra**

Search for [                    ] Search  Clear  Show All  Advanced Search

Page size 20

4348 DataSet records                                                      Page 1 of 218 > >>

| DataSet | Title | Organism(s) | Platform | Series | ▶ Samples |
|---|---|---|---|---|---|
| GDS6063 | Influenza A effect on plasmacytoid dendritic cells | *Homo sapiens* | GPL10558 | GSE68849 | 10 |
| GDS6010 | Influenza virus H5N1 infection of U251 astrocyte cell line: time course | *Homo sapiens* | GPL6480 | GSE66597 | 18 |
| GDS5879 | Pulmonary CDC11c+ cells from young and middle-age animals | *Mus musculus* | GPL6885 | GSE71868 | 8 |
| GDS5826 | Multiple myeloma cell lines with acquired resistance to chemotherapeutic agent carfilzomib | *Homo sapiens* | GPL570 | GSE69078 | 12 |
| GDS5825 | Interleukin-1α deficiency effect on injured spinal cord | *Mus musculus* | GPL6246 | GSE70302 | 12 |
| GDS5881 | Nebulin deficiency effect on the soleus | *Mus musculus* | GPL6246 | GSE70213 | 12 |
| GDS5880 | Nebulin deficiency effect on the quadriceps | *Mus musculus* | GPL6246 | GSE70213 | 12 |
| GDS5913 | SRPIN803 small molecule inhibitor of SRPK1 effect on retinal pigment epithelial cell line | *Homo sapiens* | GPL570 | GSE62947 | 6 |
| GDS5665 | Pathogen-associated molecular-pattern curdlan effect on interleukin-2 deficient GM-CSF myeloid dendritic cells | *Mus musculus* | GPL6246 | GSE58120 | 12 |
| GDS5662 | Histone demethylase KDM3A-deficiency effect on estrogen-stimulated breast cancer cells in vitro | *Homo sapiens* | GPL10558 | GSE68918 | 11 |

**DataSet Record GDS6063:** (Expression Profiles) (Data Analysis Tools) (Sample Subsets)

| | | |
|---|---|---|
| **Title:** | Influenza A effect on plasmacytoid dendritic cells | |
| **Summary:** | Analysis of primary plasmacytoid dendritic cells (pDC) exposed to influenza A for 8 hours ex vivo. pDCs are vital to antiviral defense, directing immune responses via secretion of interferon-alpha. Results provide insight into the regulation of the response of pDC to viral pathogens. | |
| **Organism:** | *Homo sapiens* | |
| **Platform:** | GPL10558: Illumina HumanHT-12 V4.0 expression beadchip | |
| **Citation:** | Bajwa G, DeBerardinis RJ, Shao B, Hall B et al. Cutting Edge: Critical Role of Glycolysis in Human Plasmacytoid Dendritic Cell Antiviral Responses. *J Immunol* 2016 Mar 1;196(5):2004-9. PMID: 26826244 | |
| **Reference Series:** | GSE68849 | **Sample count:** 10 |
| **Value type:** | count | **Series published:** 2016/02/01 |

Cluster Analysis

Download
- DataSet full SOFT file
- DataSet SOFT file
- Series family SOFT file
- Series family MINiML file
- Annotation SOFT file

NLM  NIH  GEO Help  Disclaimer  Accessibility

# https://www.ncbi.nlm.nih.gov/geo/

# https://www.immport.org/home



https://www.immport.org/shared/home

# https://www.broadinstitute.org/



**https://www.broadinstitute.org/data-sciences-platform**

**https://terra.bio/** - **A joint effort between the Broad Institute & Microsoft & Verily**

# https://www.ga4gh.org/



Formed in 2013, the Global Alliance for Genomics and Health (GA4GH) unites an international community dedicated to advancing human health through genomic data. We build technical standards and policy frameworks and tools that will expand responsible, voluntary, and secure use of genomic and other related health data.

# Interacting with Data - Software tools



https://www.tiobe.com/tiobe-index/

**I have gene expression data on 100 patients. What is the first thing to do to start looking for patterns in my data?**

If you have gene expression data on 100 patients and want to start looking for patterns in the data, here are some general steps you could follow:

1. Normalize the data: Gene expression data is typically subject to variation due to technical and biological factors. Normalization is important to remove systematic variations that are not related to biological variation. There are several methods for normalization, including quantile normalization, RMA normalization, and Z-score normalization.
2. Visualize the data: Visualization is an important first step in exploring gene expression data. You can use scatter plots, heatmaps, and boxplots to visualize the distribution of gene expression values across samples and identify patterns.
3. Perform exploratory data analysis: Explore the distribution of gene expression values, identify outliers, and check for correlations between genes.
4. Identify differentially expressed genes: Use statistical tests such as t-tests, ANOVA, or non-parametric tests to identify genes that are differentially expressed between groups of interest, such as patients with and without a disease.
5. Conduct pathway analysis: Pathway analysis can help identify biological pathways that are enriched for differentially expressed genes. This can provide insights into the underlying biology of the disease or condition of interest.
6. Use machine learning algorithms: Machine learning algorithms such as clustering, principal component analysis, and dimensionality reduction can help identify patterns in the data that may not be immediately apparent from visual inspection.
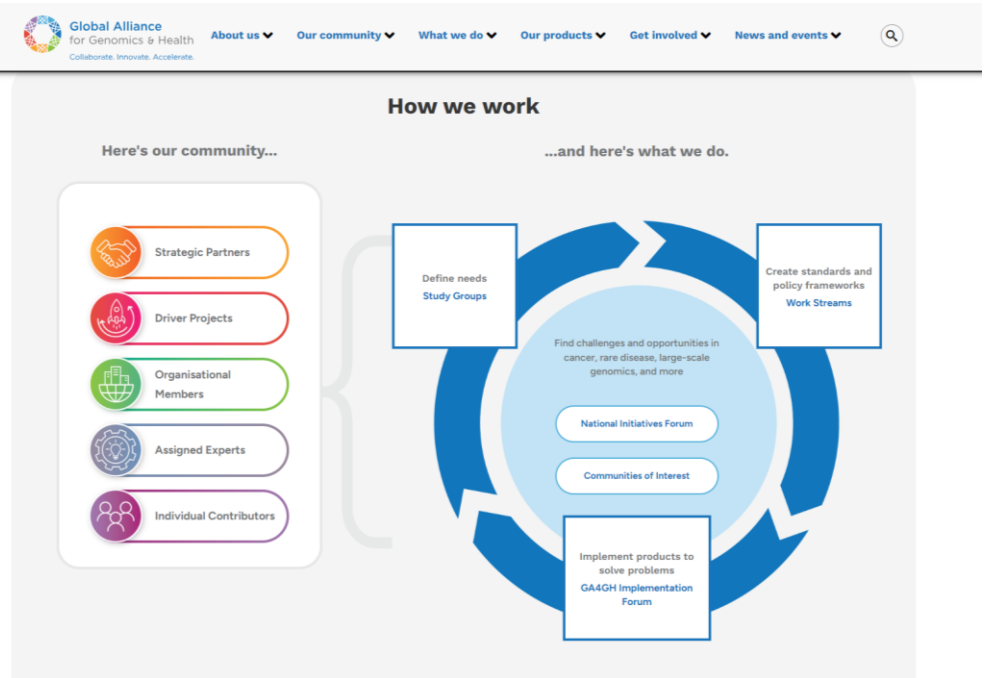7. Interpret the results: Once you have identified patterns in the data, it is important to interpret the results in the context of existing knowledge and hypotheses. This may involve consulting the literature, conducting follow-up experiments, or developing new hypotheses.

# *What tools do I use?*
## → *Python* and *R are* still dominating in bioinformatics

**https://www.r-project.org/**

### The R Project for Statistical Computing

Getting Started

R is a free software environment for statistical computing and graphics. It compiles and runs on a wide variety of UNIX platforms, Windows and MacOS. To download R, please choose your preferred CRAN mirror.

If you have questions about R like how to download and install the software, or what the license terms are, please read our answers to frequently asked questions before you send an email.

[Home]

**Download**
CRAN

**R Project**
About R

**posit** | PRODUCTS ⌄   SOLUTIONS ⌄   LEARN & SUPPORT ⌄   EXPLORE MORE ⌄   PRICING    🔍   **DOWNLOAD RSTUDIO**

RSTUDIO IDE

# The most trusted IDE for open source data science

RStudio is an integrated development environment (IDE) for R and Python. It includes a console, syntax-highlighting editor that supports direct code execution, and tools for plotting, history, debugging, and workspace management. RStudio is available in open source and commercial editions and runs on the desktop (Windows, Mac, and Linux).

# Bioconductor
OPEN SOURCE SOFTWARE FOR BIOINFORMATICS

## About *Bioconductor*

The mission of the *Bioconductor* project is to develop, support, and disseminate free open source software that facilitates rigorous and reproducible analysis of data from current and emerging biological assays. We are dedicated to building a diverse, collaborative, and welcoming community of developers and data scientists.

*Bioconductor* uses the R statistical programming language, and is open source and open development. It has two releases each year, and an active user community. *Bioconductor* is also available as Docker images.

**https://www.bioconductor.org/**

### Install »

- Discover 2083 software packages available in *Bioconductor* release 3.14.

Get started with *Bioconductor*

- Install *Bioconductor*
- Get support
- Latest newsletter
- Follow us on twitter
- Install R

### Learn »

Master *Bioconductor* tools

- Courses
- Support site
- Package vignettes
- Literature citations
- Common work flows
- FAQ
- Community resources
- Videos

### Use »

Create bioinformatic solutions with *Bioconductor*

- Software, Annotation, and

### Develop »

Contribute to *Bioconductor*

- Developer resources
- Use Bioc 'devel'

Please provide R code for principal component analysis

Sure! Here is an example of how to perform principal component analysis (PCA) in R:

```r
# Load data
data <- read.csv("data.csv")

# Perform PCA
pca <- prcomp(data, scale = TRUE)

# Print summary of results
summary(pca)

# Plot results
plot(pca, type = "l")
```

In this example, `data` is a dataframe containing the data we want to perform PCA on. We use the `prcomp()` function to perform the PCA, with the `scale = TRUE` argument indicating that we want to scale the data before performing the analysis. We then use `summary()` to print a summary of the results, and `plot()` to visualize the results.

Note that in practice, you may want to customize the PCA analysis to fit your specific needs. For example, you may want to specify the number of principal components you want to extract, or you may want to use a different scaling method. The `prcomp()` function has several optional arguments that allow you to customize the analysis, so be sure to consult the documentation for more information.

# Interactive data analysis using e.g. RShiny web-applications

**https://gdc.cancer.gov/**



**https://gdc.cancer.gov/analyze-data**

# https://cri-iatlas.org/



# https://isb-cgc.shinyapps.io/iatlas/

# Cell

# Optimal-Transport Analysis of Single-Cell Gene Expression Identifies Developmental Trajectories in Reprogramming

## Graphical Abstract

## Authors

Geoffrey Schiebinger, Jian Shu, Marcin Tabaka, ..., Rudolf Jaenisch, Aviv Regev, Eric S. Lander

## Correspondence

jianshu@broadinstitute.org (J.S.), aregev@broadinstitute.org (A.R.), lander@broadinstitute.org (E.S.L.)

## In Brief

Application of a new analytical approach to examine developmental trajectories of single cells offers insight into how paracrine interactions shape reprogramming.

**https://broadinstitute.github.io/wot/**

# Principal Moment Analysis

</> View Principal Moment Analysis on GitHub

Welcome to the Principal Moment Analysis home page!

Principal Moment Analysis is described in this paper:

> " " *Fontes, M., & Henningsson, R. (2020). Principal Moment Analysis. arXiv arXiv:2003.04208.*

There is a Julia implementation of (Simplex) Principal Moment Analysis. You can also use the Principal Moment Analysis App.

# PMA for dimension reduction and visualizations of State Spaces

## https://principalmomentanalysis.github.io/

PMA is a framework that incorporates the sample distribution and lower dimensional (noise reduced) approximations of it.

# Take home messages concerning PMA

- Fast (*sample based PMA* is as fast as corresponding PCA)
- Robust (equivalent to local infinite pseudo sample bootstrap)
- Statistically and conceptually sound ("optimal" approximation of underlying probability measure respecting intrinsic local dimensionality and sampling density, quality etc)
- Possible to supervise using annotation information and expert knowledge
- Immediate: *PMA projection score* generalizes.
- Immediate: Kernel PMA

# STATISTICAL LEARNING & VISUALIZATION → ACQUIRING INSIGHTS AROUND PATTERNS IN DATA RELATED TO BIOLOGICAL VARIATION



Design and Photo: Emilia Fontes 2016

# Exploratory Data Analysis vs Confirmatory Data Analysis



John Wilder Tukey (1915-2000)
Inventor of the FFT, the Box plot and the word "bit".



1977

First question: What should it mean to be similar?
-> Choice of similarity measure or distance function
*Example genuine metrics:*

A **metric space** is a set $M$ together with a fixed **distance function** or **metric** $d : M \times M \longrightarrow [0, \infty)$ such that for all $x, y$ and $z$ in $M$ we have

$$d(x, y) \geq 0 \quad \text{with equality if and only if (iff) } x = y \qquad (2.1)$$

$$d(x, y) = d(y, x) \quad (symmetry) \qquad (2.2)$$

$$d(x, y) \leq d(x, z) + d(z, y) \quad (the\ triangle\ inequality) \qquad (2.3)$$

Examples:
- **Different types of edit distances** for e.g. sequencing data
- L2 or *Euclidean distance* for quantitative data
- L1 or *Manhattan (or Taxicab) distance*

# Effects of different similarity measures or distance functions



The "unit sphere" in the $L^1$ norm. A constraint formulated with this norm favors *sparse vectors*. The idea behind the *Lasso*, see Tibshirani, R. (1996). *Regression shrinkage and selection via the lasso.* J. Royal. Statist. Soc B., Vol. 58, No. 1, pages 267-288).

The normal "unit sphere" in the standard $L^2$ norm

The "unit sphere" in the $L^\infty$ norm

# **Feature extraction and selection** to "understand" phenotypic variation



Measurements

Sample space population

Data



Data Analysis

→ **clustering & classification** → **Biological Insights**

# How to explore MULTI-OMICS data?



- N=10-10^5 Samples
- P=10-10^10 Variables

## How to find relevant structure?

## Problems:
- Noise
- Artifacts: technical, batch effects,..
- High Dimensions
- Often few samples and many variables

*Visualize and Analyze* the data
- Choice of "*similarity measure*"
- Choice of dimension/model-reduction method = choice of "*objective function*"

*Classification of pediatric acute lymphoblastic leukemia by gene expression profiling*
Mary E. Ross et al.
Blood 2003 102:2951-2959; doi:10.1182/blood-2003-01-0338
**Data available at https://www.stjuderesearch.org/site/data/ALL3/**

# Example Data: Affymetrix chip 22282 ProbeIDs and 132 Samples



NEOPLASIA

Classification of pediatric acute lymphoblastic leukemia by gene expression profiling

Mary E. Ross, Xiaodong Zhou, Guangchun Song, Sheila A. Shurtleff, Kevin Girtman, W. Kent Williams, Hsi-Che Liu, Rami Mahfouz, Susana C. Raimondi, Noel Lenny, Anami Patel, and James R. Downing

Contemporary treatment of pediatric acute lymphoblastic leukemia (ALL) requires the assignment of patients to specific risk groups. We have recently demonstrated that expression profiling of leukemic blasts can accurately identify the known prognostic subtypes of ALL, including T-cell lineage ALL (T-ALL), *E2A-PBX1*, *TEL-AML1*, *MLL* rearrangements, *BCR-ABL*, and hyperdiploid karyotypes with more than 50 chromosomes. As the next step toward developing this methodology into a frontline diagnostic tool, we have now analyzed leukemic blasts from 132 diagnostic samples using higher density oligonucleotide arrays that allow the interrogation of most of the identified genes in the human genome. Nearly 60% of the newly identified subtype discriminating genes are novel markers not identified in our previous study, and thus should provide new insights into the altered biology underlying these leukemias. Moreover, a proportion of the newly selected genes are highly ranked as class discriminators, and when incorporated into class-predicting algorithms resulted in an overall diagnostic accuracy of 97%. The performance of an array containing the identified discriminating genes should now be assessed in frontline clinical trials in order to determine the accuracy, practicality, and cost effectiveness of this methodology in the clinical setting. (Blood. 2003; 102:2951-2959)

© 2003 by The American Society of Hematology

# Dimension reduction techniques: *Classical PCA*



**Legend:**
- BCR-ABL (black)
- E2A-PBX1 (red)
- Hyperdiploid (>5 (green)
- MLL (blue)
- Other (cyan)
- T-ALL (magenta)
- TEL-AML1 (yellow)

ALL (Leukemia) samples
mRNA bulk

22282 variables 132 samples

**PCA based on the correlation matrix**;
Mean centering and normalizing all variables to unit variance

**Distances**: Euclidean      **Objective function**: Total Variance

Scree Plot

132/132 Samples, 22282/22282 Variables

We can use a ***Scree plot*** to obtain a rough estimate of the global dimension of our data set

**BMC Bioinformatics**

# The projection score – an evaluation criterion for variable subset selection in PCA visualization

Magnus Fontes[*] and Charlotte Soneson

$$\alpha_2(\mathbf{\Lambda_X}, S) = \frac{\sum_{k \in S} \lambda_k^2}{\sum_{k=1}^{r} \lambda_k^2}.$$

$$\tau(\phi_m(\mathbf{X}), S, \mathcal{P}_{\phi_m(\mathbf{X})}) = \left(\alpha_2(\mathbf{\Lambda}_{\phi_m(\mathbf{X})}, S)\right)^{1/2}$$
$$- \mathbb{E}_{\mathcal{P}_{\phi_m(\mathbf{X})}} \left[ \left(\alpha_2(\mathbf{\Lambda}_{\phi_m(\mathbf{X})}, S)\right)^{1/2} \right].$$

Ross DT, Scherf U, Eisen MB, Perou CM, Rees C, Spellman P, Iyer V, Jeffrey SS, Van de Rijn M, Waltham M, Pergamenschikov A, Lee JC, Lashkari D, Shalon D, Myers TG, Weinstein JN, Botstein D, Brown PO: Systematic variation in gene expression patterns in human cancer cell lines. *Nat Genet* 2000, **24**:227-235.

Duality of PCA → *Bi-plot* after optimization over *projection score*

# A Global Geometric Framework for Nonlinear Dimensionality Reduction

Joshua B. Tenenbaum,[1]* Vin de Silva,[2] John C. Langford[3]

## isomap

### Isometric Feature Mapping Ordination

The function performs isometric feature mapping which consists of three simple steps: (1) retain only some of the shortest dissimilarities among objects, (2) estimate all dissimilarities as shortest path distances, and (3) perform metric scaling (Tenenbaum et al. 2000).

Keywords    multivariate

### Usage

```
isomap(dist, ndim=10, ...)
isomapdist(dist, epsilon, k, path = "shortest", fragmentedOK =FALSE, ...)
"summary"(object, axes = 4, ...)
"plot"(x, net = TRUE, n.col = "gray", type = "points", ...)
```
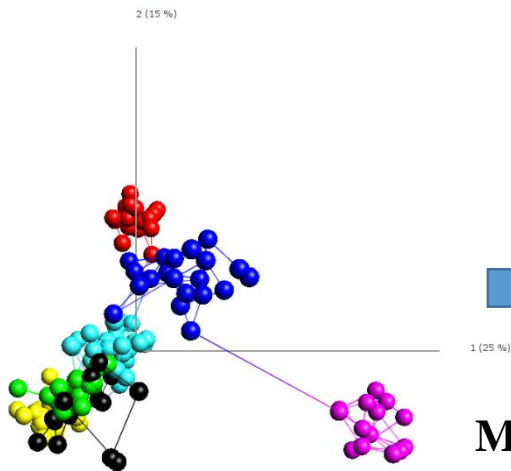
### Arguments

| | |
|---|---|
| dist | Dissimilarities. |
| ndim | Number of axes in metric scaling (argument `k` in `cmdscale`). |
| epsilon | Shortest dissimilarity retained. |

**BCR-ABL**
**E2A-PBX1**
**Hyperdiploid (>50)**
**MLL**
**Other**
**T-ALL**
**TEL-AML1**

2 (15 %)

1 (25 %)

2 (8 %)

1 (83 %)

**MDS based on graph distances**

# ISOMAP

# Stochastic Neighbor Embedding (SNE)

Original distances

$$p_{j|i} = \frac{\exp\left(-\|x_i - x_j\|^2 / 2\sigma_i^2\right)}{\sum_{k \neq i} \exp\left(-\|x_i - x_k\|^2 / 2\sigma_i^2\right)},$$

Original SNE

Sam Roweis and Geoffrey Hinton

$$q_{j|i} = \frac{\exp\left(-\|y_i - y_j\|^2\right)}{\sum_{k \neq i} \exp\left(-\|y_i - y_k\|^2\right)}.$$

Distances in reduced space

$$q_{ij} = \frac{\left(1 + \|y_i - y_j\|^2\right)^{-1}}{\sum_{k \neq l} \left(1 + \|y_k - y_l\|^2\right)^{-1}}.$$

tSNE

Van der Maaten

The cost function $C$ is given by

$$C = \sum_i KL(P_i \| Q_i) = \sum_i \sum_j p_{j|i} \log \frac{p_{j|i}}{q_{j|i}},$$

Kullback Leibler
Relative Entropy

Stochastic Neighbor Embedding (SNE) and tSNE

**Package 'Rtsne'**

June 30, 2016

**Type** Package

**Title** T-Distributed Stochastic Neighbor Embedding using a Barnes-Hut Implementation

**Version** 0.11

**Description** An R wrapper around the fast T-distributed Stochastic Neighbor Embedding implementation by Van der Maaten.

**License** BSD_3_clause + file LICENSE

**URL** https://github.com/jkrijthe/Rtsne

**Imports** Rcpp (>= 0.11.0)

Journal of Machine Learning Research 1 (2008) 1-48          Submitted 4/00; Publ

**Visualizing Data using t-SNE**

**Laurens van der Maaten**                    L.VANDERMAATEN@MICC.UNI
*MICC-IKAT*
*Maastricht University*
*P.O. Box 616, 6200 MD Maastricht, The Netherlands*

**Geoffrey Hinton**                           HINTON@CS.TORONTO.EDU
*Department of Computer Science*
*University of Toronto*
*6 King's College Road, M5S 3G4 Toronto, ON, Canada*

latest

Search docs

Read the Docs          v: latest ▾

# UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction

Uniform Manifold Approximation and Projection (UMAP) is a dimension reduction technique that can be used for visualisation similarly to t-SNE, but also for general non-linear dimension reduction. The algorithm is founded on three assumptions about the data

1. The data is uniformly distributed on Riemannian manifold;
2. The Riemannian metric is locally constant (or can be approximated as such);
3. The manifold is locally connected.

From these assumptions it is possible to model the manifold with a fuzzy topological structure. The embedding is found by searching for a low dimensional projection of the data that has the closest possible equivalent fuzzy topological structure.
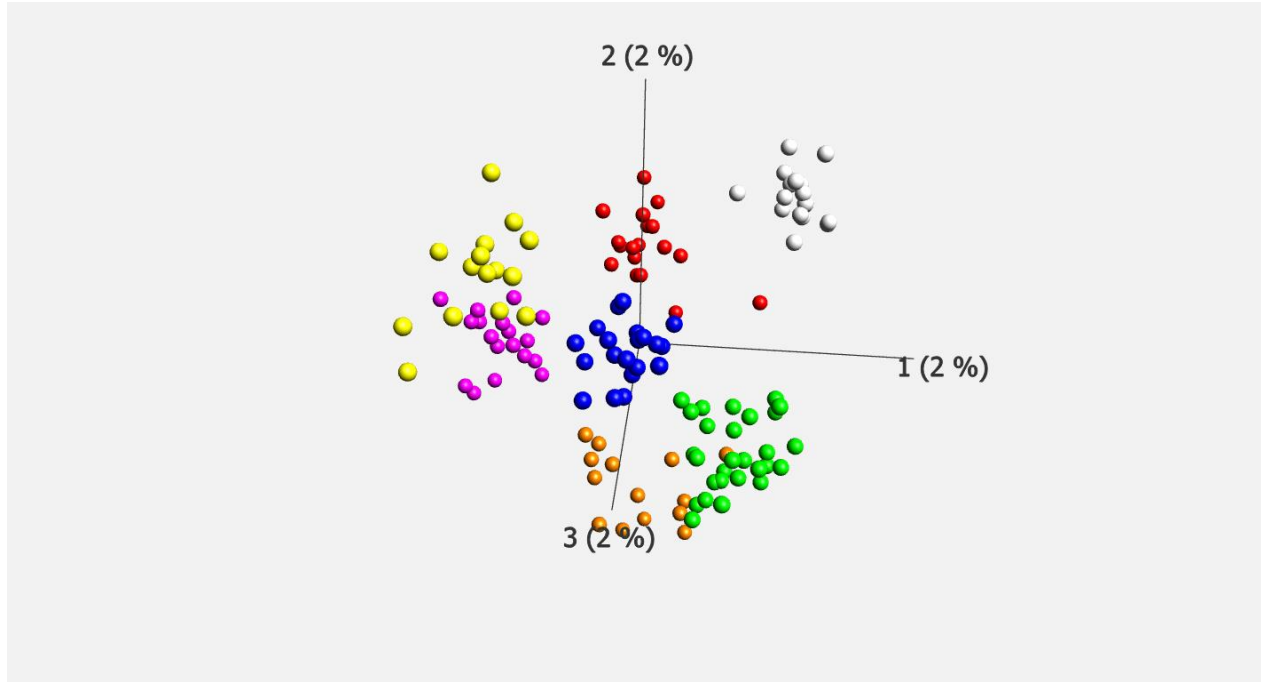
The details for the underlying mathematics can be found in our paper on ArXiv:

McInnes, L, Healy, J, *UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction*, ArXiv e-prints 1802.03426, 2018
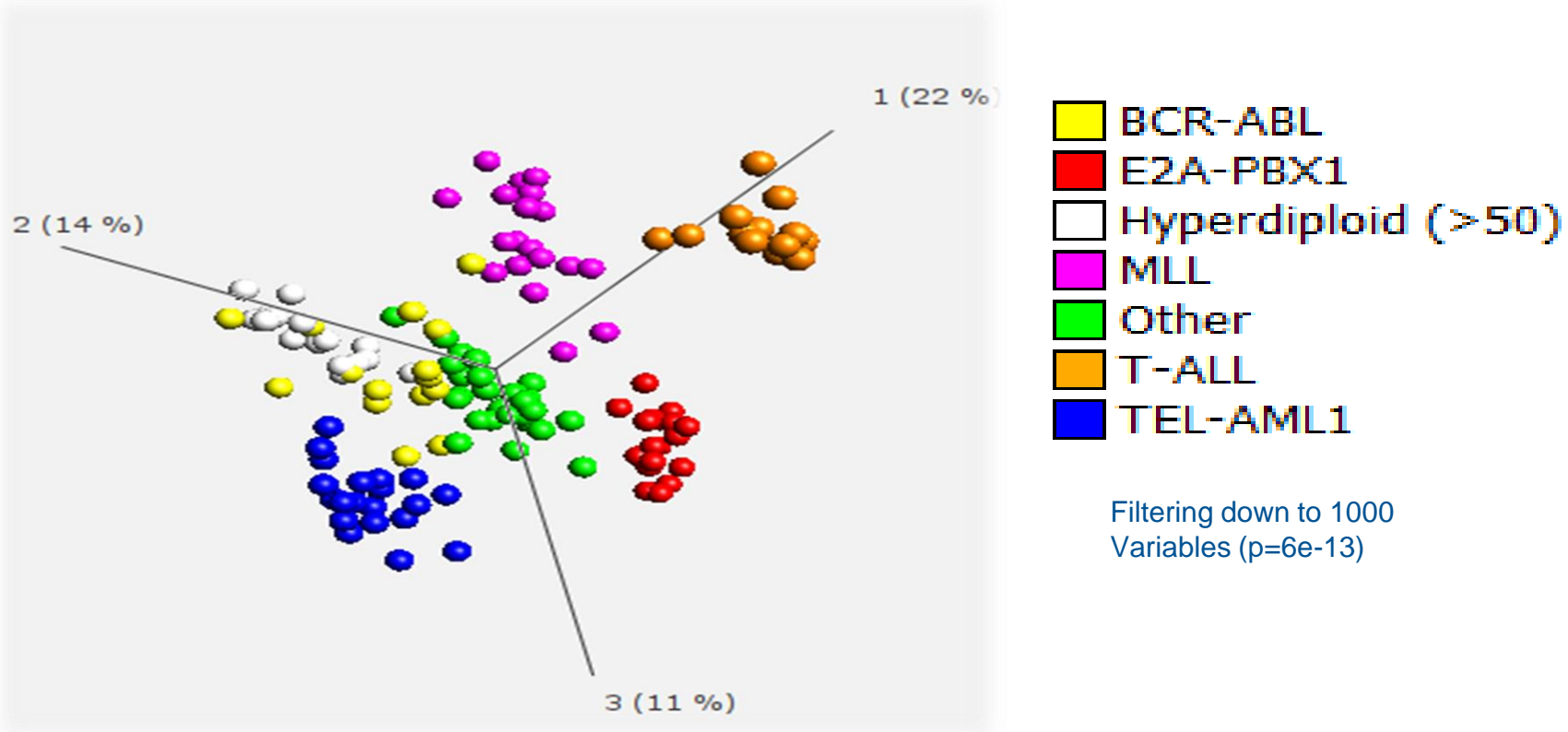
You can find the software on github.

## Installation

# Dangers with exploratory analyses: ANOVA on Random data



Filtering using ANOVA on a 22282*132 dataset with significance threshold p=0.05 on random data resulting in 1108 "*significant discoveries*" + PCA visualization. *Note that 0.05*22282 =1114.1*
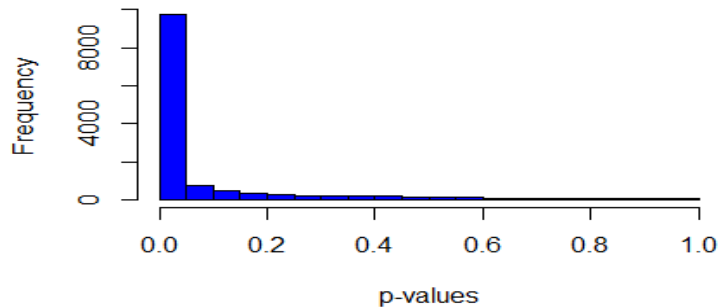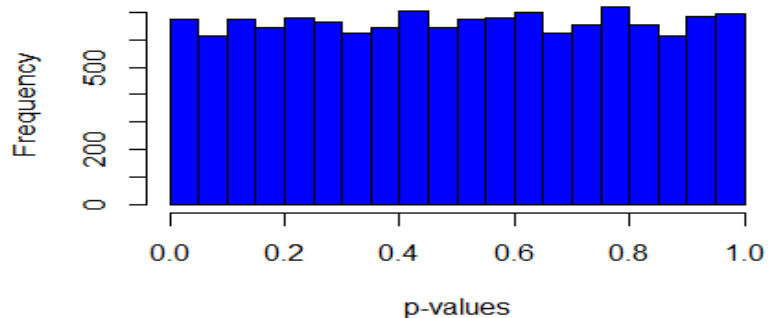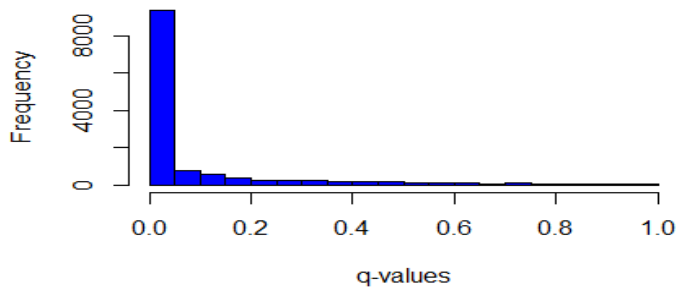
# ANOVA on ALL dataset



Filtering down to 1000
Variables (p=6e-13)

# The p-value distributions

# Recap: General directions in Translational research

Emerging Technology & Data & Insights for

- Longitudinal Sampling:
  Precision Diagnostics & Medicine: Mapping of « Patient journeys » &
  « Trajectories in Biomedical State Space » integrating different data modalities
  (Omics, Imaging, RWD, ….)
- Deep phenotyping sampling:
  Integrative and Holistic analyses using emerging bio-technologies



Multidisciplinary collaborative data & advanced analytics driven biomedical research in order to:

- Better understanding the Dynamics & Complexity:
  Static snapshots will be complemented with Dynamic systems control approaches leading to dynamic & composite biomarkers for response, resistance, safety, quality of life …
- Better Health State Monitoring & Early Detection and dynamical control:
  Controlling and maintaining health in order to prevent disease will lead to treatment paradigm shifts.

# General Systems Immunology for Biomedicine: Predicting and controlling dynamics under perturbations



**Immunospace**

**Immune trajectory**
- Individual 1
- Individual 2
- Individual 3
- Individual 4

- "Human immune systems are relatively stable within individuals over the course of weeks to months, but incredibly variable between individuals"
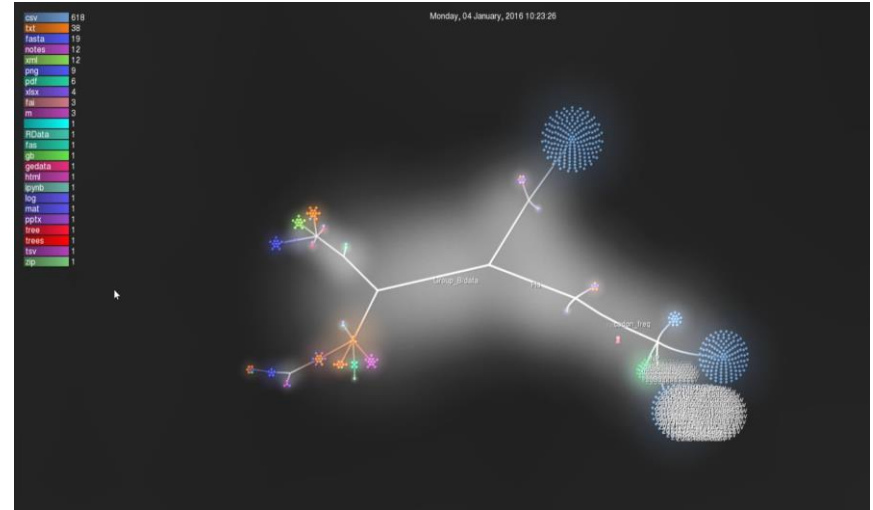
- "…induced responses to pathogens differ markedly among different age groups and… these differences are unique to different kinds of stimuli"

- Functional gene expression responses of human blood cells to common pathogens differ broadly across age groups.

- Divergent immune cell composition with advanced age is associated with chromatin changes that are induced by environmental influences over the course of life.

# Infrastructure connecting the dots through Human-Machine partnerships
## Biomedicine& Data & Advanced Analytics → New Insights

- Connecting data bases

- Connecting IT platforms & tools

- Connecting collaborative research infrastructure (wiki, fora, ...)

- Connecting scientists from bench to computational mathematics

- Incentivize sharing through provenance & tracking



**Drive cultural change** ⟷ **Collaboration & Sharing**