

Introduction to multi-omics data analysis & machine learning workshop

H2020/FindingPheno, Jan 13-14, 2022

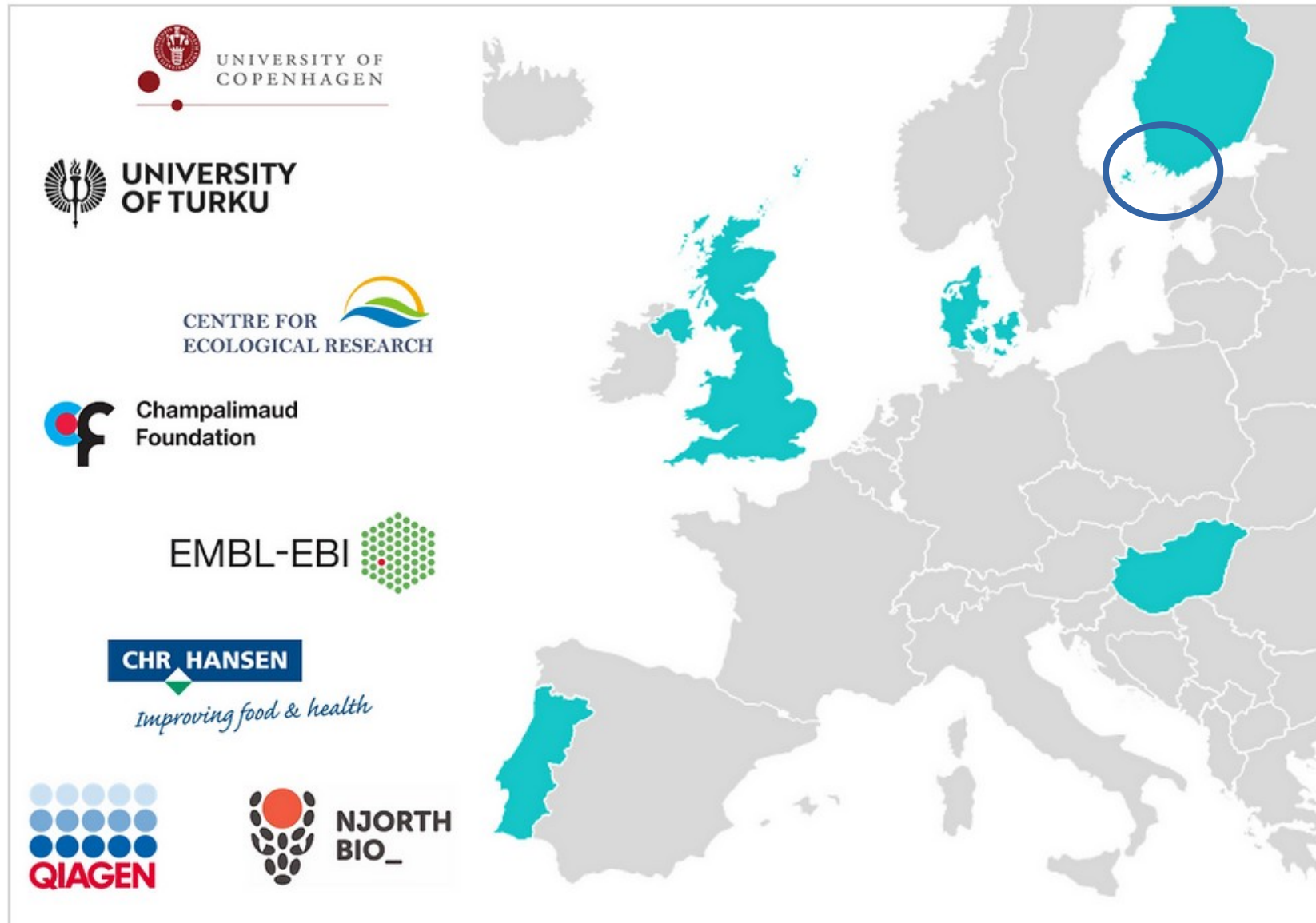


Associate Prof. Leo Lahti | datascience.utu.fi
Department of Computing, University of Turku, Finland



Turun yliopisto
University of Turku

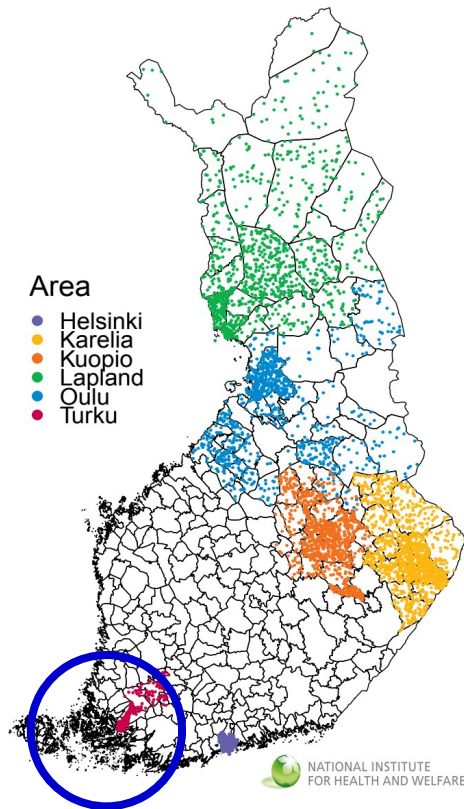
FindingPheno is a truly interdisciplinary project featuring close collaboration between academic and industrial partners. Our combined scientific expertise and experiences include statistics and machine learning, ecology and evolutionary genetics, industrial food production systems experience and development of world class genomics software.



Department of Computing, University of Turku, Finland

datascience.utu.fi

N=7231



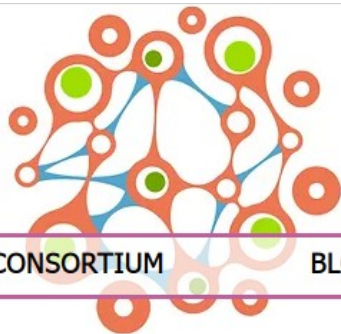
The University of Turku

Turku, Finland

The University of Turku (UTU) was established in 1920 but its roots reach back to the Royal Academy of Turku in the 1650s. Today, UTU has almost 20,000 students and 3,271 staff members.

The UTU researchers participating in FindingPheno come from the Turku Data Science Group within the Department of Computing. This group is headed by Assoc. Prof. Leo Lahti and combines theory and methods of algorithmic data analysis, with a particular focus on probabilistic machine learning, complex systems, high-throughput data analysis and statistical programming.





FindingPheno

[ABOUT](#)[CONSORTIUM](#)[BLOG](#)[EVENTS](#)[DELIVERABLES](#)[CONTACT](#)[Log In](#)

Multi-omics data analysis for genotype-phenotype associations

The generation of large-scale biological data sets is exploding. Genomics, transcriptomics, proteomics, metabolomics of both host and microbiome.. But! The methods to integrate and exploit this data do not keep up.

FindingPheno is answering this call.

Developing new computational tools for untangling host-microbiome interactions in multi-omic data sets

[Find out more](#)

EU-funded research collaboration with eight academic, research and industry partners from across Europe

[Find out more](#)

Validating our tools against commercially relevant test cases for more sustainable food production

[Find out more](#)



This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 952914

Contact us: findingpheno.eu@gmail.com

Day 1 (Times in CET)

Lectures (45 min + 15 min breaks)

9:15-10:00 - **Welcome & introduction** - Leo Lahti, Associate professor (UTU)

10:15-11:00 - **Metagenomics** - Katariina Pärnänen, Postdoctoral researcher (UTU)

11:15-12:00 - **Metabolomics** - Pande Putu Erawijantari, Postdoctoral researcher (UTU)

12:15-13:00 - **Multimomics** - Leo Lahti, Associate professor (UTU)

13:00-14 - **Lunch** break

Practical session

14:15-17:00 - Tuomas Borman and Chouaib Benchraka, Research assistants (UTU)

Data import and data structures

Microbiome data exploration & visualization

Day 2 (Times in CET)

Lectures

9:15-10:00 - **Unsupervised ML**- Matti Ruuskanen, Postdoctoral researcher (UTU)

10:15-11:00 - **Supervised ML** - Matti Ruuskanen

11:15-12:00 - **Individual-based modeling** - Gergely Boza, Research fellow (CER)

12:15-13:00 - **Data integration** - Leo Lahti, Associate professor (UTU)

13:00-14 - **Lunch** break

Practical

14:15-17:00 - Tuomas Borman, Matti Ruuskanen and Chouaib Benchraka (UTU)

Unsupervised learning: Beta-diversity and biclustering

Supervised learning: Regression and classification with random forests

Validation and interpretation of black box models



Matti Ruuskanen
Postdoc
Turku, Finland



Pande Putu Erawijantari
Postdoc
Turku, Finland



Katariina Pärnänen
Postdoc
Turku, Finland



Gergely Boza
Research fellow,
Center for Ecological Research,
Hungary



Chouaib Benchraka
Scientific programmer
Turku, Finland



Tuomas Borman
Scientific programmer
Turku, Finland



Leo Lahti
Associate Prof.
Turku, Finland

Learning objectives

motivation & challenges in multi-omics

examples of computational approaches

hands-on experience on R/Bioconductor tools

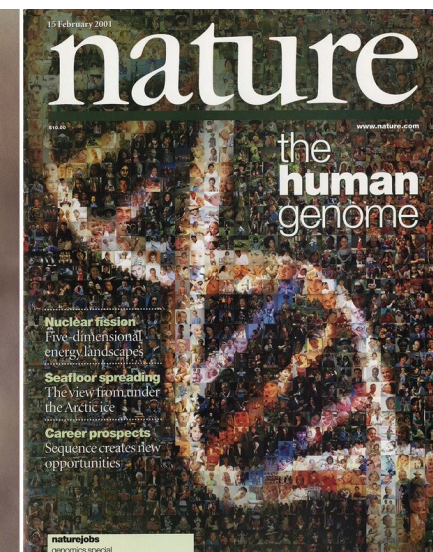
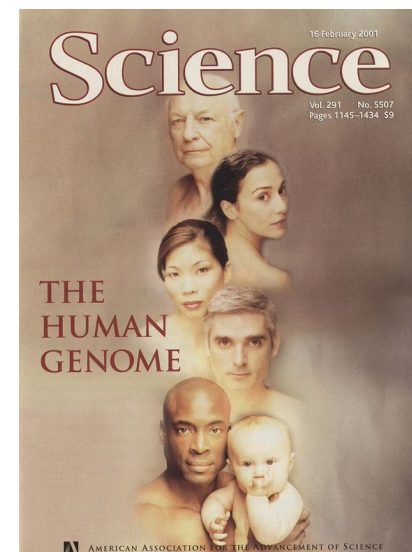
Initial sequencing and analysis of the human genome

~2001

International Human Genome Sequencing Consortium*

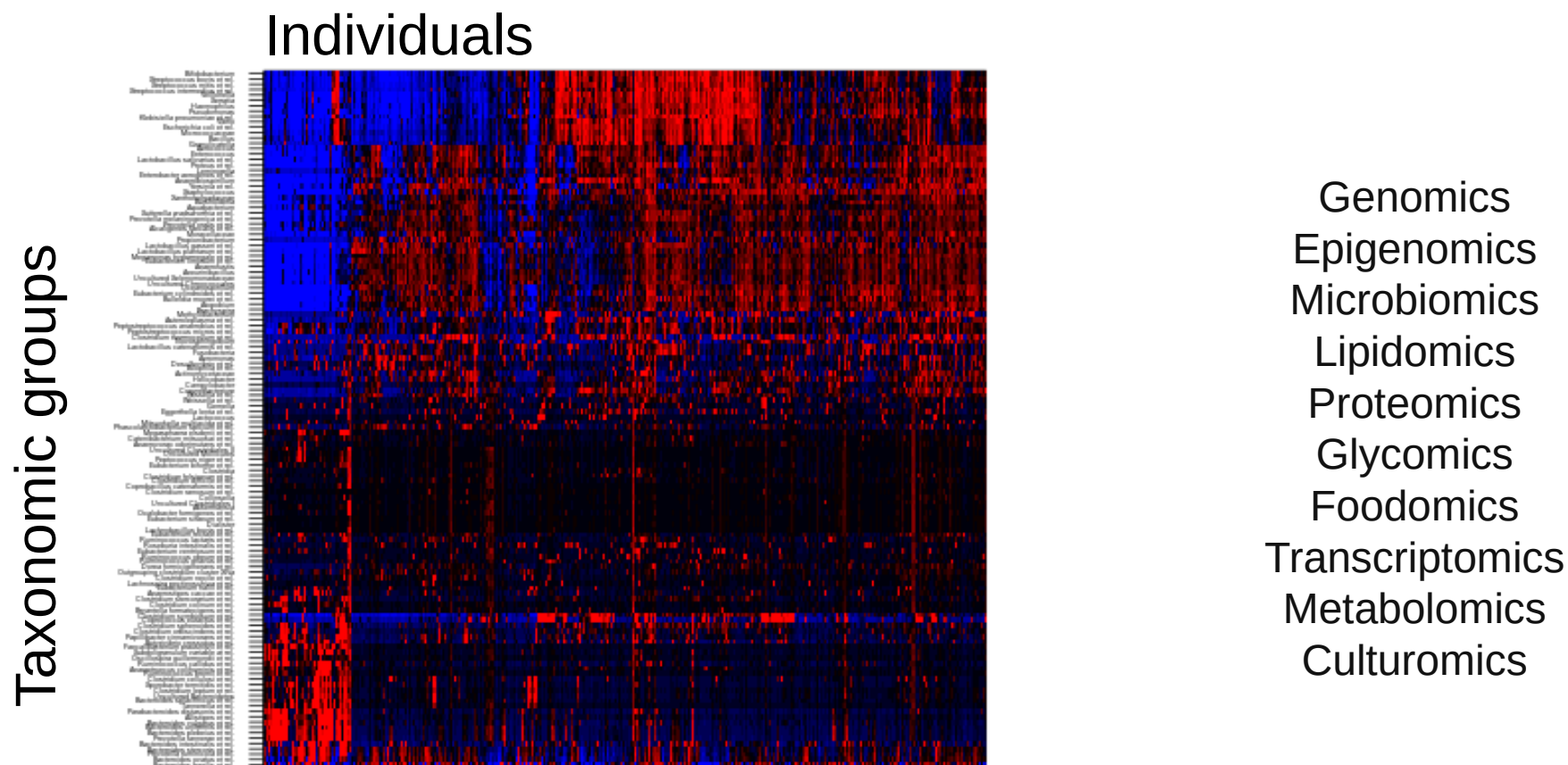
* A partial list of authors appears on the opposite page. Affiliations are listed at the end of the paper.

The human genome holds an extraordinary trove of information about human development, physiology, medicine and evolution. **Here we report the results of an international collaboration to produce and make freely available a draft sequence of the human genome.** We also present an initial analysis of the data, describing some of the insights that can be gleaned from the sequence.



Omics: taxonomic abundance table

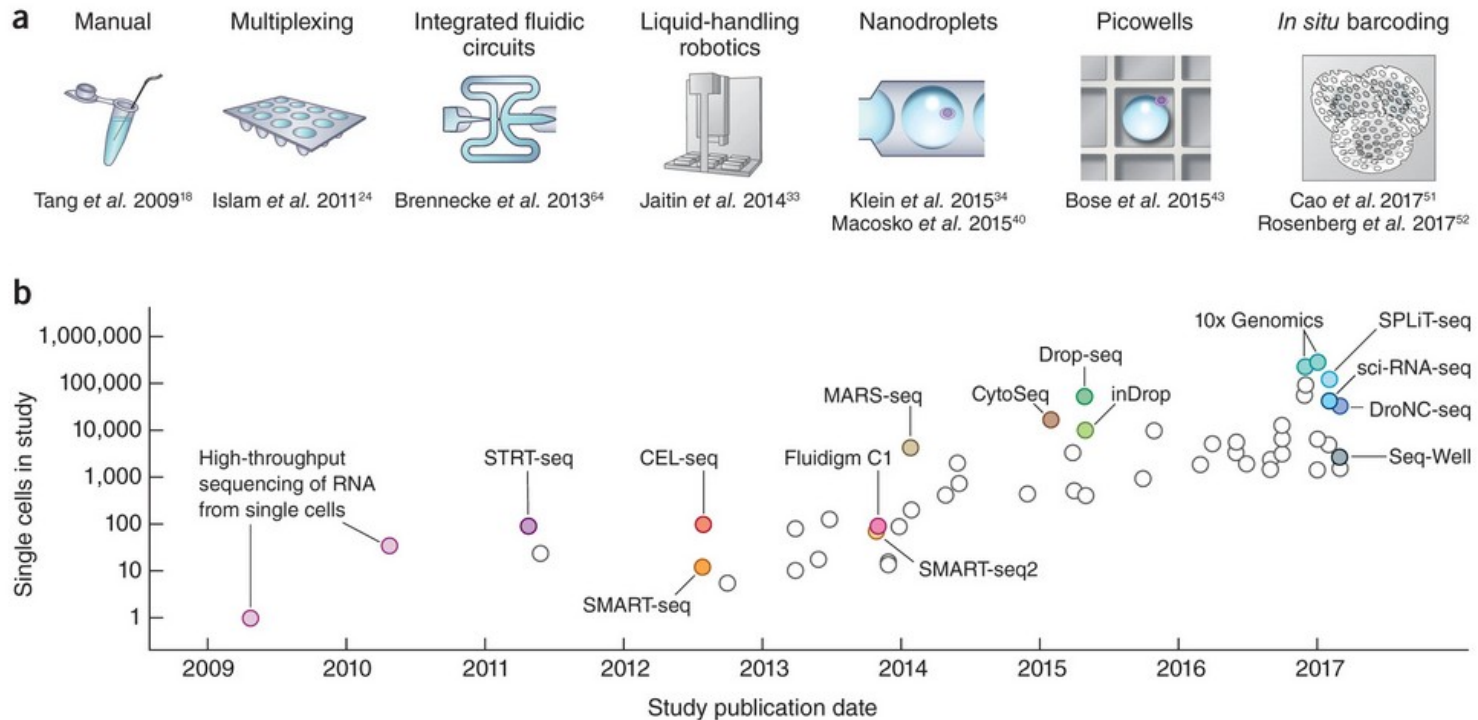
Omics in Oxford English Dictionary:
in cellular and molecular biology,
forming nouns with the sense
"all constituents considered collectively"



Gut microbiota: 1000 western adults
(Lahti *et al.* Nature Comm. 2014)

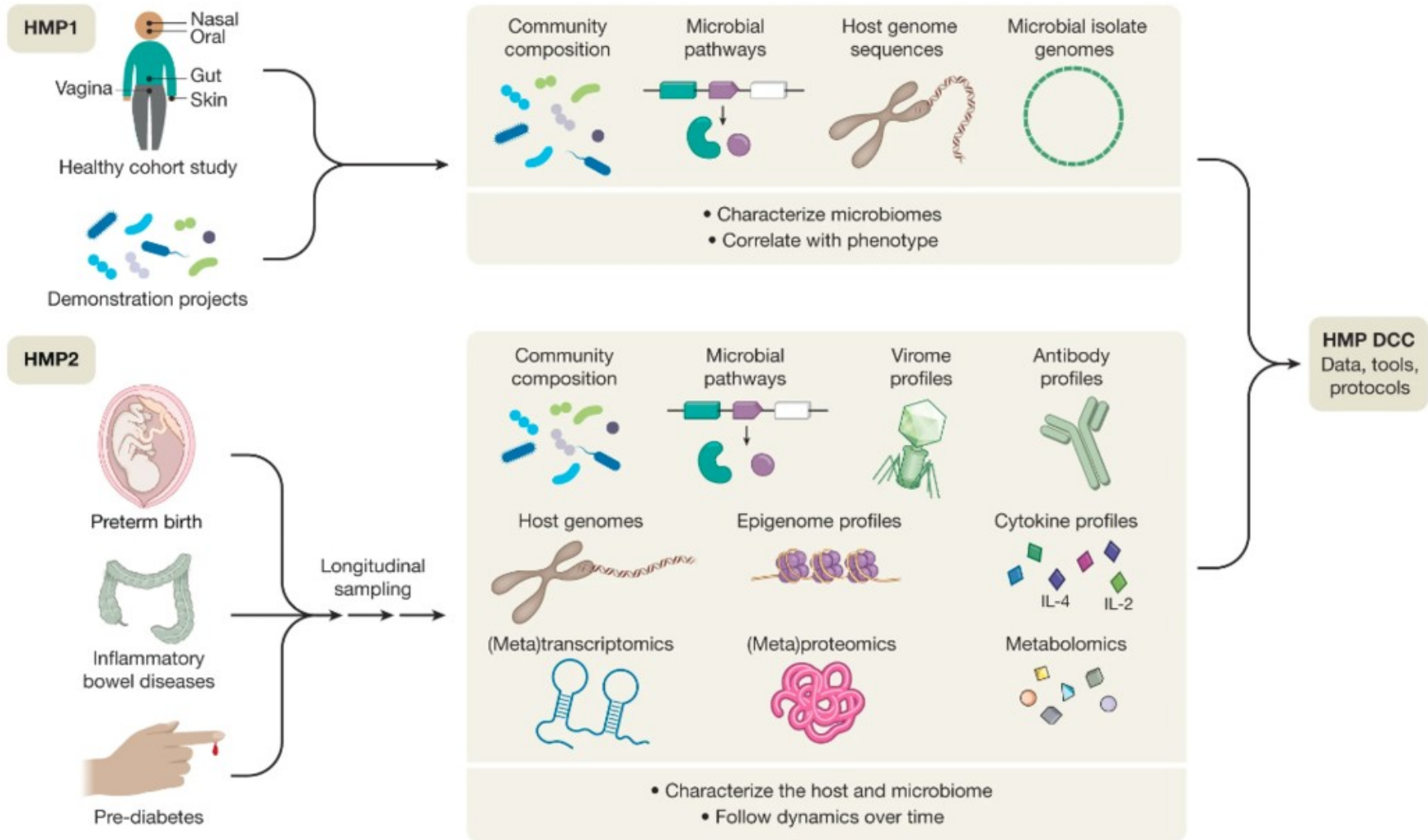
Figure 1: Scaling of scRNA-seq experiments.

From: [Exponential scaling of single-cell RNA-seq in the past decade](#)



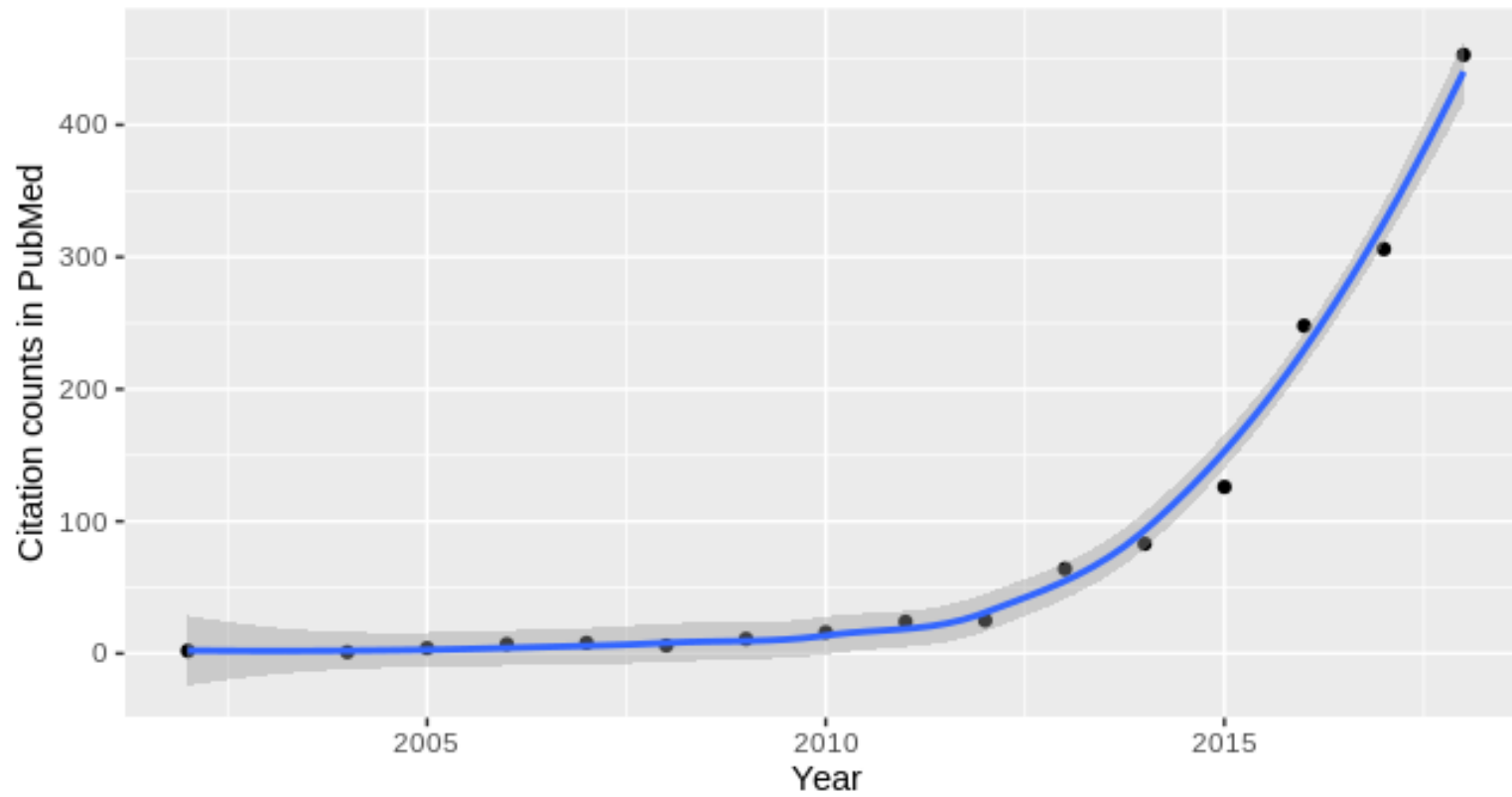
a) Key technologies that have allowed jumps in experimental scale. A jump to ~ 100 cells was enabled by sample multiplexing, and then a jump to $\sim 1,000$ cells was achieved by large-scale studies using integrated fluidic circuits, followed by a jump to several thousands of cells with liquid-handling robotics. Further orders-of-magnitude increases bringing the number of cells assayed into the tens of thousands were enabled by random capture technologies using nanodroplets and picowell technologies. Recent studies have used *in situ* barcoding to inexpensively reach the next order of magnitude of hundreds of thousands of cells. (b) Cell numbers reported in representative publications by publication date. Key technologies are indicated.

Multi-omics

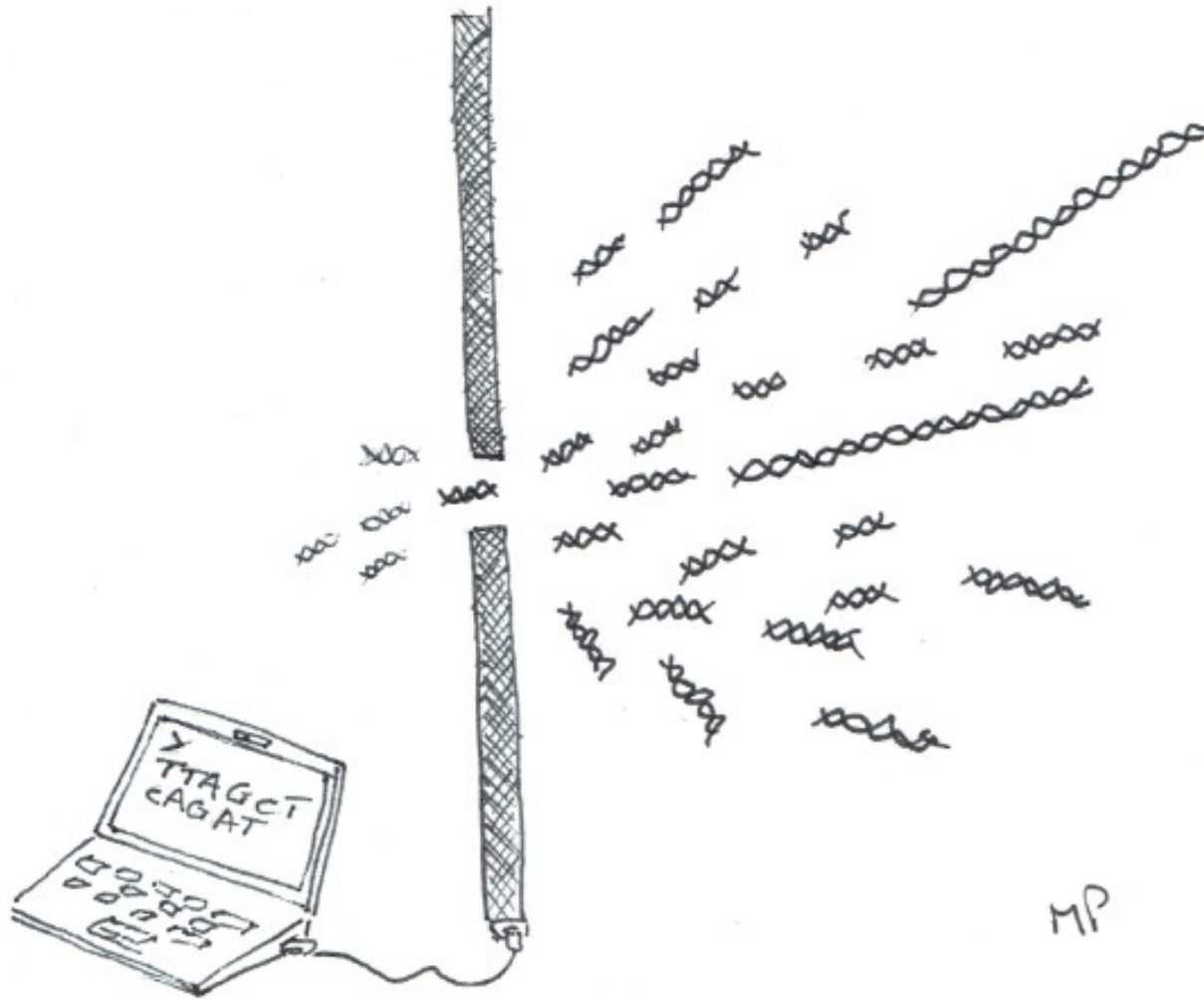


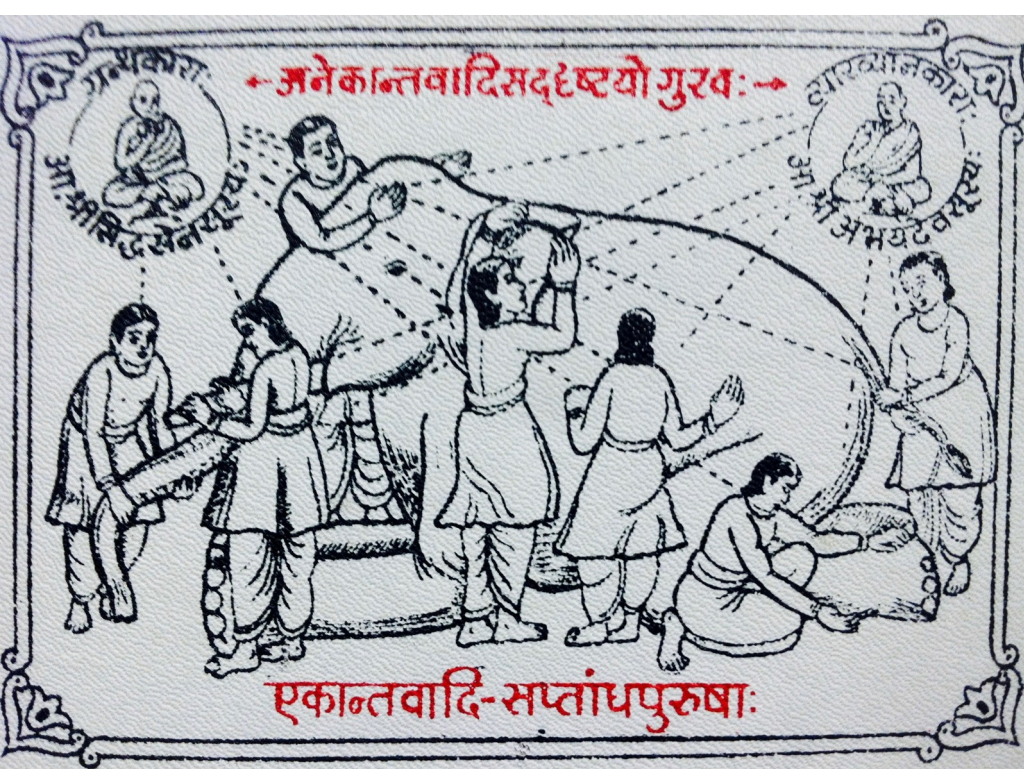
Number of articles related to "multiomics" in PubMed until 2018

(source: Wikipedia)



Limited observations → data integration?



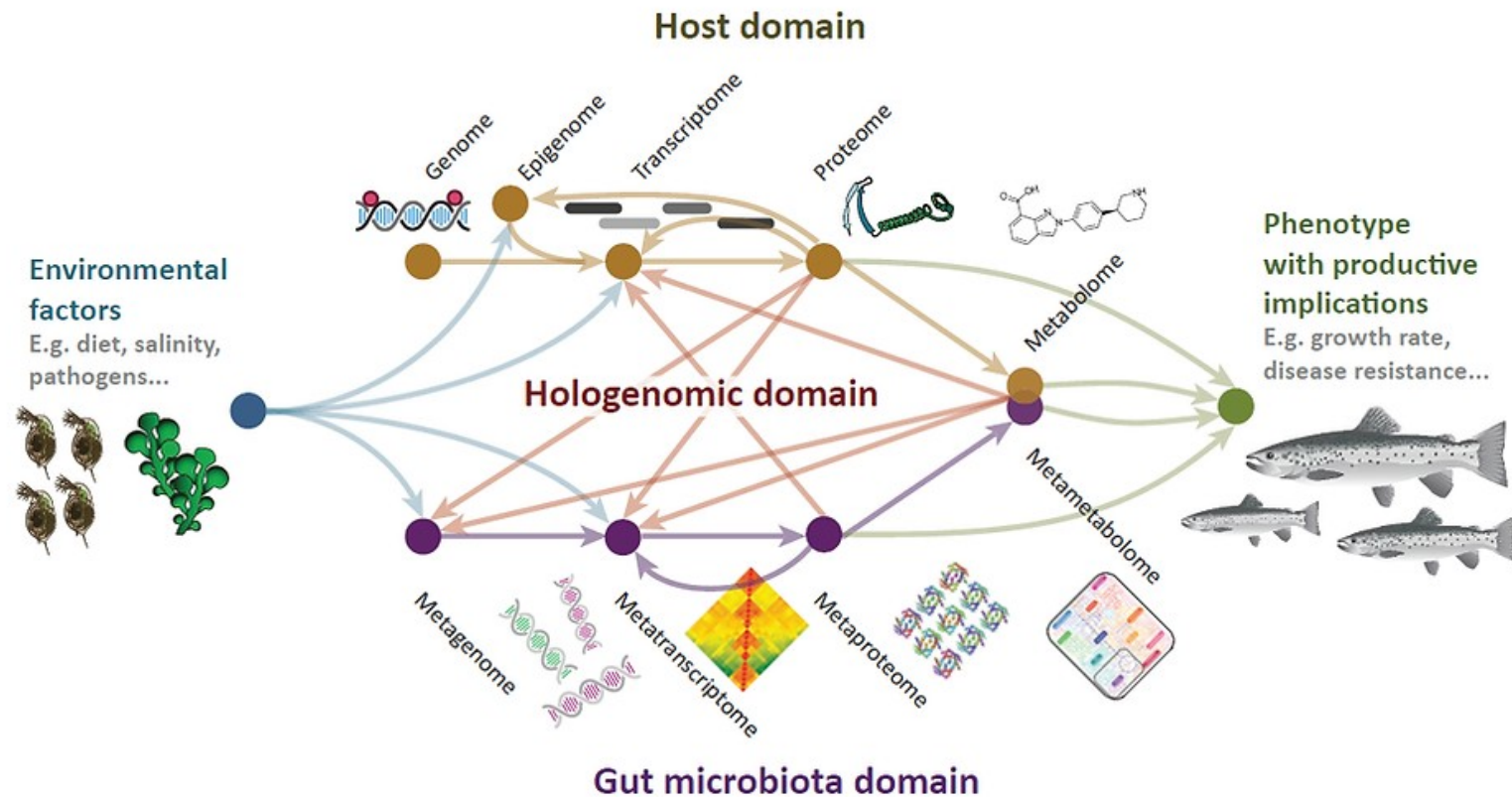


Elephant in the dark

The medieval era Jain texts explain the concepts of anekāntavāda (or "many-sidedness") and syādvāda ("conditioned viewpoints") with the parable of the blind men and an elephant (Andhgajanyāyah), which addresses the manifold nature of truth.

The Buddhist text *Tittha sutta*, *Udāna* 6.4, *Khuddaka Nikaya*, contains one of the earliest versions of the story. The *Tittha sutta* is dated to around c. 500 BCE, although the parable is likely older.

FindingPheno is creating an integrated computational framework for hologenomic big data, providing the tools to better understand how host-microbiome interactions can affect growth and other outcomes.



Understanding the hologenomic domain is a fiendishly difficult problem, with a complex tangle of interactions at many molecular levels both within and between organisms. FindingPheno aims to solve this problem, developing a unified statistical framework for the intelligent integration of multi-omic data from both host and microbiome to understand biological outcomes.

We apply state-of-the-art mathematical and machine learning approaches taken from evolutionary genomics, collective behaviour analysis, ecosystem dynamics, statistical modelling, and applied agricultural research to give us a truly interdisciplinary perspective towards solving this difficult problem. Our project takes a unique two-pronged approach: combining biology-agnostic machine learning methods with biology-informed hierarchical modelling to increase the power and adaptability of our predictive tools.

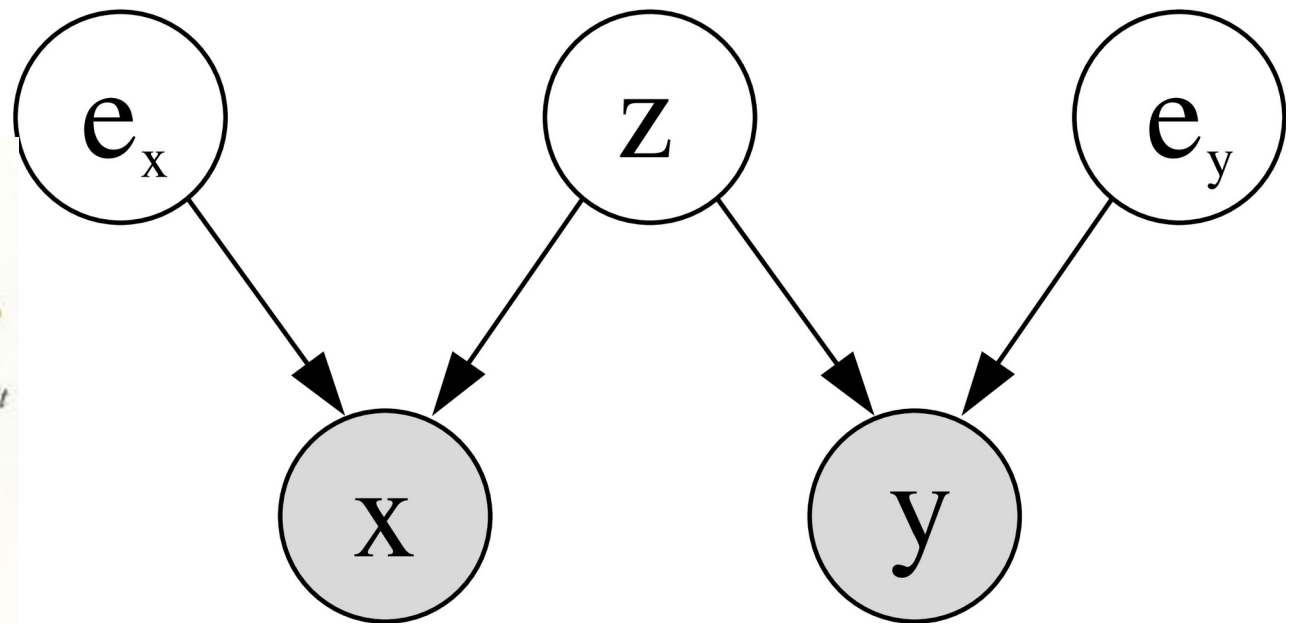
The tools created in FindingPheno are expected to significantly improve how we understand and utilise the functions provided by microbiomes in combating human diseases as well as the way we produce sustainable food for future generations.

$$\begin{cases} X = W_x \mathbf{z} + \varepsilon_x \\ Y = W_y \mathbf{z} + \varepsilon_y \end{cases}$$

Multi-view learning



Mon dessin ne représentait pas un chapeau. Il représentait un serpent boa qui digérait un éléphant

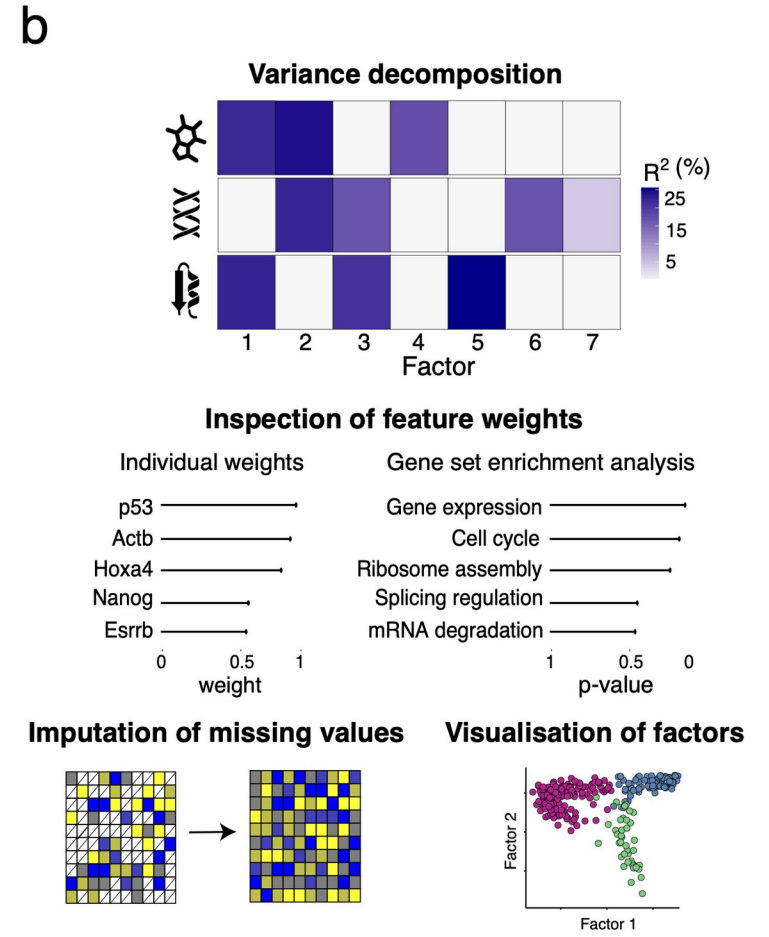
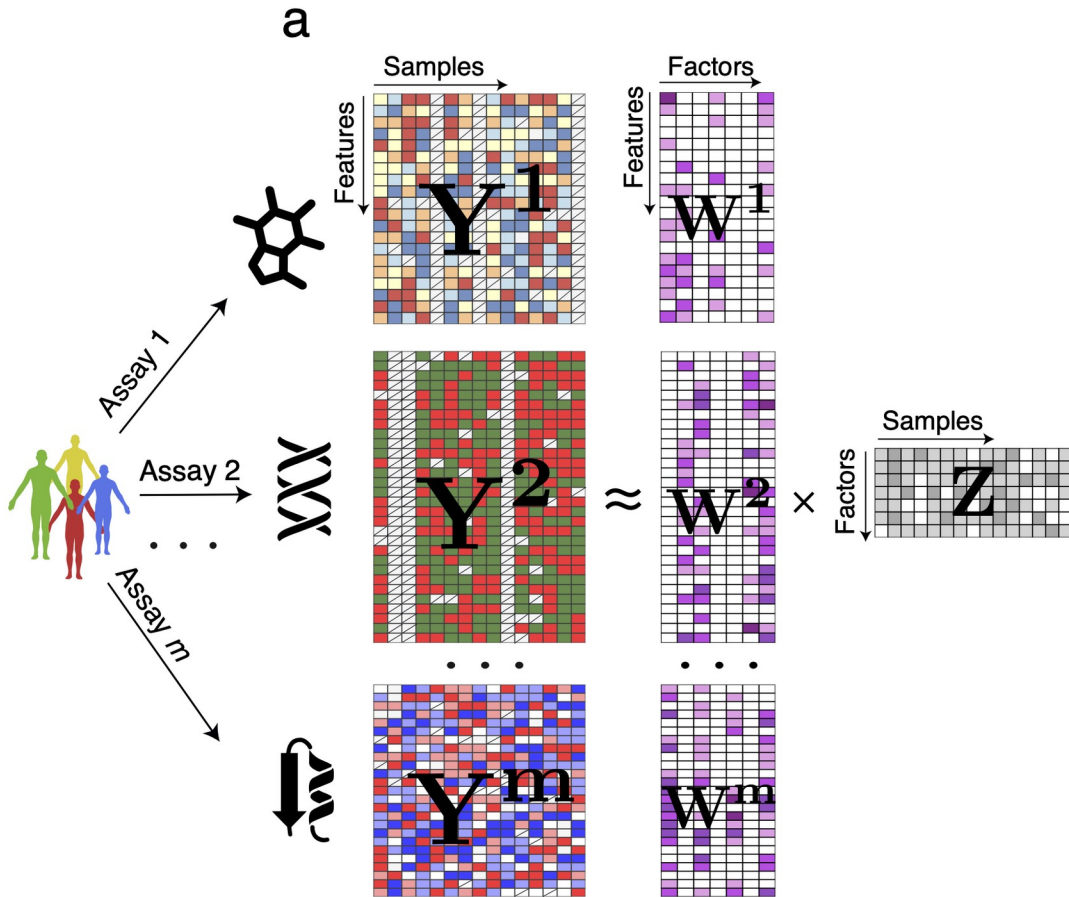


Multi-Omics Factor Analysis—a framework for unsupervised integration of multi-omics data sets

Ricard Argelaguet, Britta Velten, Damien Arno, Sascha Dietrich, Thorsten Zenz, John C. Marioni, Florian Buettner, Wolfgang Huber, Oliver Stegle

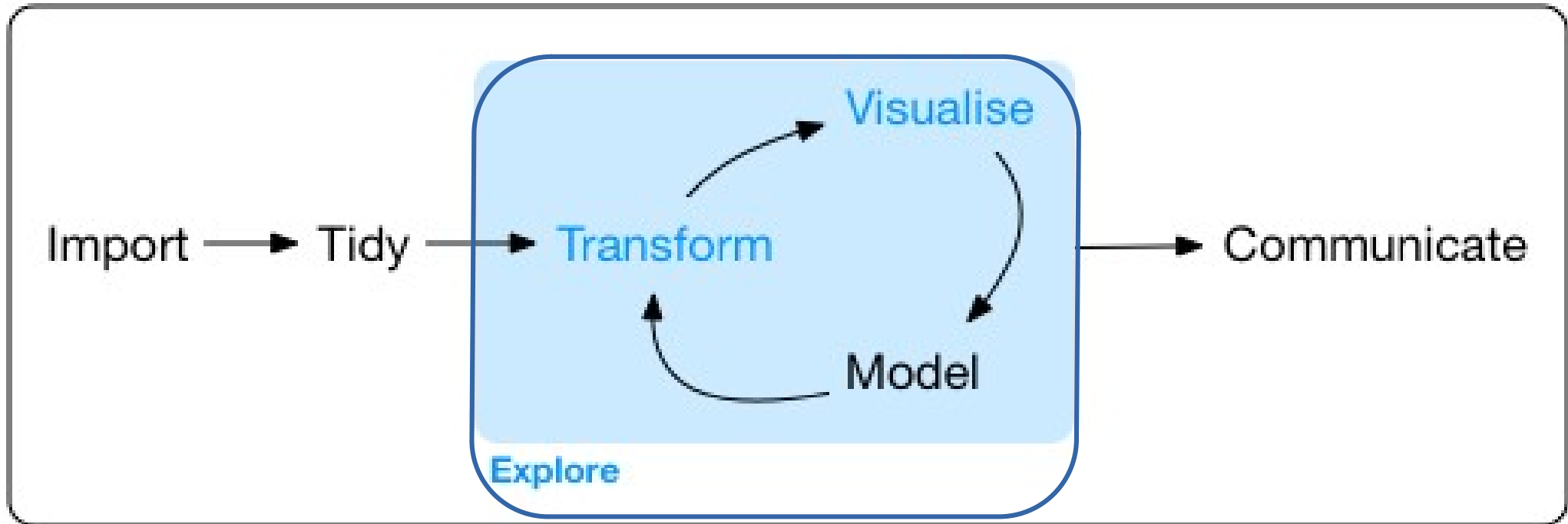
Author Information

Molecular Systems Biology (2018) 14: e8124 | <https://doi.org/10.15252/msb.20178124>



open data science frameworks

Computational workflows have an increasingly central role in research!



Program

R for Data Science / H. Wickham

The influence of hidden researcher decisions in applied microeconomics

Nick Huntington-Klein ✉, Andreu Arenas, Emily Beam, Marco Bertoni, Jeffrey R. Bloem, Pralhad Burli, Naibin Chen, Paul Grieco, Godwin Ekpe, Todd Pugatch, Martin Saavedra, Yaniv Stopnitzky

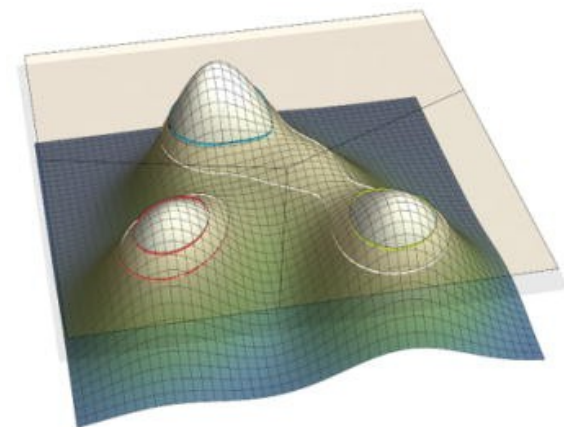
First published: 22 March 2021

<https://doi.org/10.1111/ecin.12992>

Researchers make hundreds of decisions about data collection, preparation, and analysis in their research. We use a many-analysts approach to measure the extent and impact of these decisions. Two published causal empirical results are replicated by seven replicators each. We find large differences in data preparation and analysis decisions, many of which would not likely be reported in a publication. No two replicators reported the same sample size. Statistical significance varied across replications, and for one of the studies the effect's sign varied as well. The standard deviation of estimates across replications was 3–4 times the mean reported standard error.

How to choose a correct model?

→ a community typing example



Enterotypes in the landscape of gut microbial community composition. Costea *et al.* Nature 2018.

$$2 \times 6^6 = 93312$$

Taxonomic level

- Phylum
- Family
- Order
- Genus
- Species
- Strain..

Filtering

- None
- Prevalent
- Core
- Excl. outliers
- High variance
- Custom

Normalization

- None
- TSS
- CSS
- ILR/ALR/CLR
- phILR
- Hellinger

(Dis)similarity

- Eulidean
- Aitchison
- Bray-Curtis
- Jaccard
- weighted Unifrac
- unweighted Unifrac

Clustering method

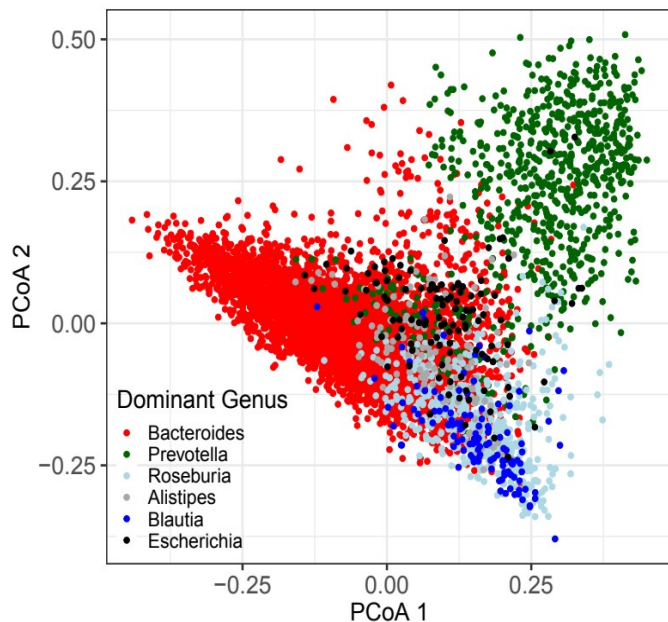
- Hierarchical / Ward
- Hierarchical / Complete
- Gaussian mixture
- DMM
- PAMR
- K-means

Regulation

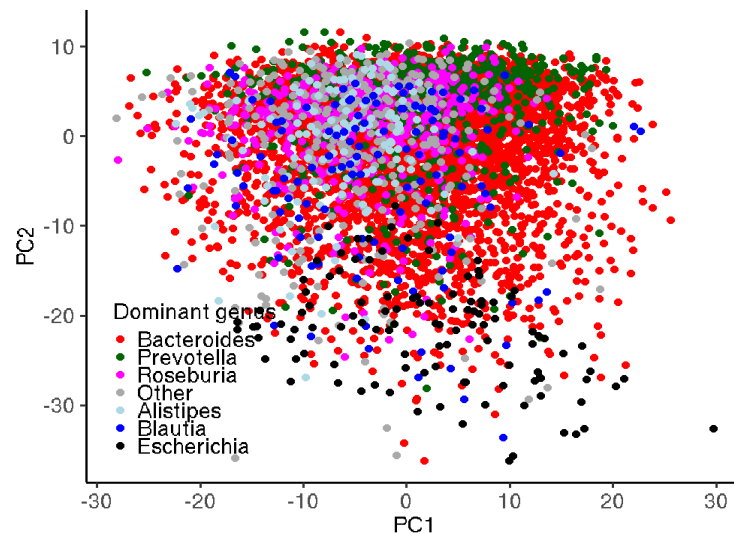
- Calinski-Harabasz
- Dirichlet Process
- Silhouette Index
- AIC
- BIC
- DIC

Walk-through example in R/Bioc by Holmes & McMurdie
<http://statweb.stanford.edu/~susan/papers/EnterotypeRR.html>

PCoA + Bray-Curtis



PCA + Aitchison



Reproducible Research: Enterotype Example

Susan Holmes and Joey McMurdie

<http://statweb.stanford.edu/~susan/papers/EnterotypeRR.html>

[Comment on this paper](#)

Taxonomic Signatures of Long-Term Mortality Risk in Human Gut Microbiota

• Aaro Salosensaari, • Ville Laitinen, • Aki Havulinna, Guillaume Meric, • Susan Cheng, • Markus Perola, Liisa Valsta, • Georg Alfthan, • Michael Inouye, Jeremie D. Watrous, Tao Long, • Rodolfo Saldo, Karenina Sanders, Cairiona Brennan, Gregory C. Humphrey, Jon G. Sanders, • Mohit Jain, Pekka Jousilahti, • Veikko Salomaa, • Rob Knight, • Leo Lahti, • Teemu Niiranen
doi: <https://doi.org/10.1101/2019.12.30.19015842>

A manifesto for reproducible science

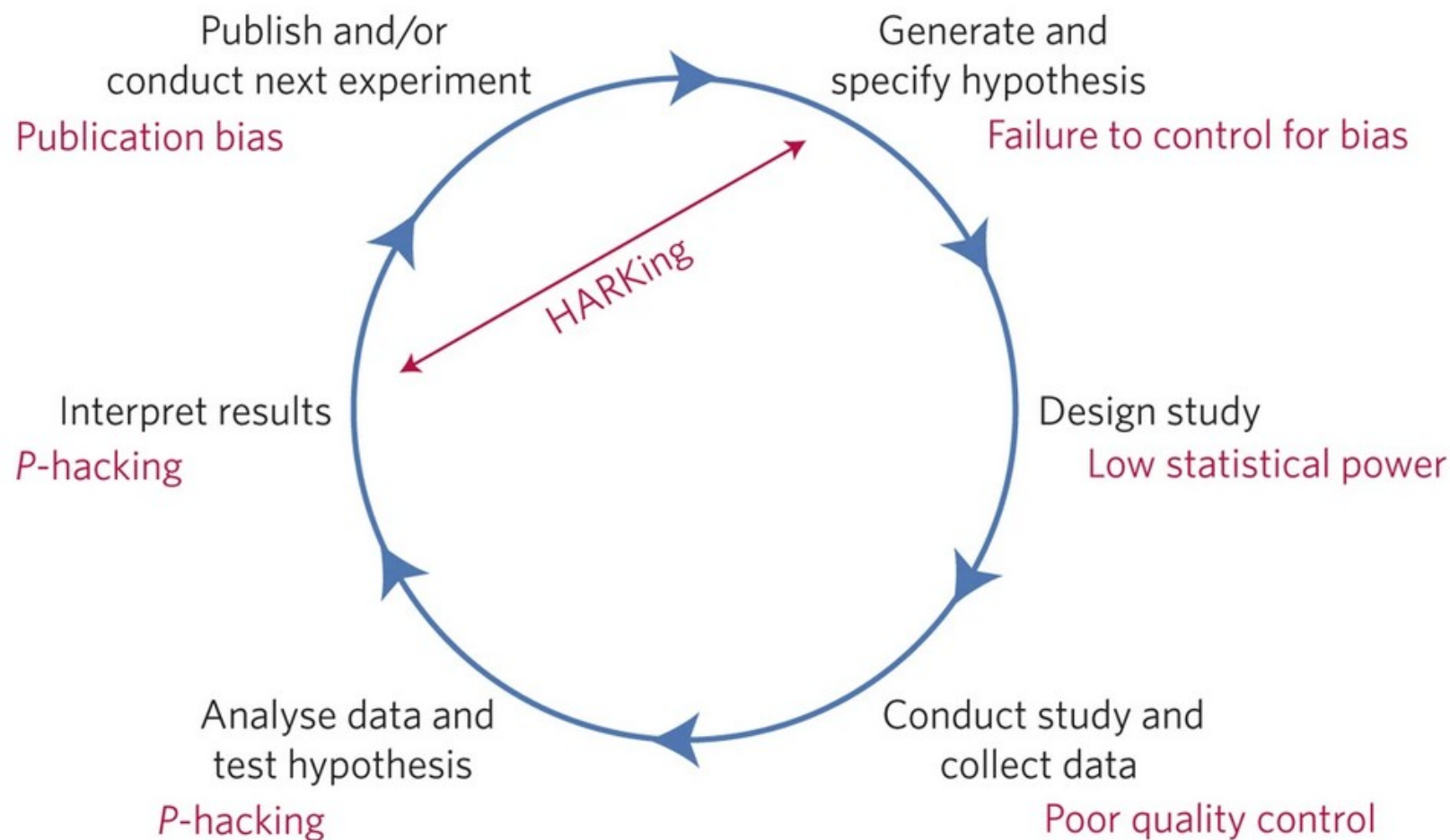
Marcus R. Munafò , Brian A. Nosek, Dorothy V. M. Bishop, Katherine S. Button, Christopher D. Chambers, Nathalie Percie du Sert, Uri Simonsohn, Eric-Jan Wagenmakers, Jennifer J. Ware & John P. A. Ioannidis

Nature Human Behaviour 1, Article number: 0021 (2017) | [Cite this article](#)

204k Accesses | 963 Citations | 2579 Altmetric | [Metrics](#)

Figure 1: Threats to reproducible science.

From: [A manifesto for reproducible science](#)



An idealized version of the hypothetico-deductive model of the scientific method is shown. Various potential threats to this model exist (indicated in red), including lack of replication⁵, hypothesizing after the results are known (HARKing)⁷, poor study design, low statistical power², analytical flexibility⁵¹, P-hacking⁴, publication bias³ and lack of data sharing⁶. Together these will serve to undermine the robustness of published research, and may also impact on the ability of science to self-correct.

You aren't
doing it wrong



if no one knows
what you are doing.

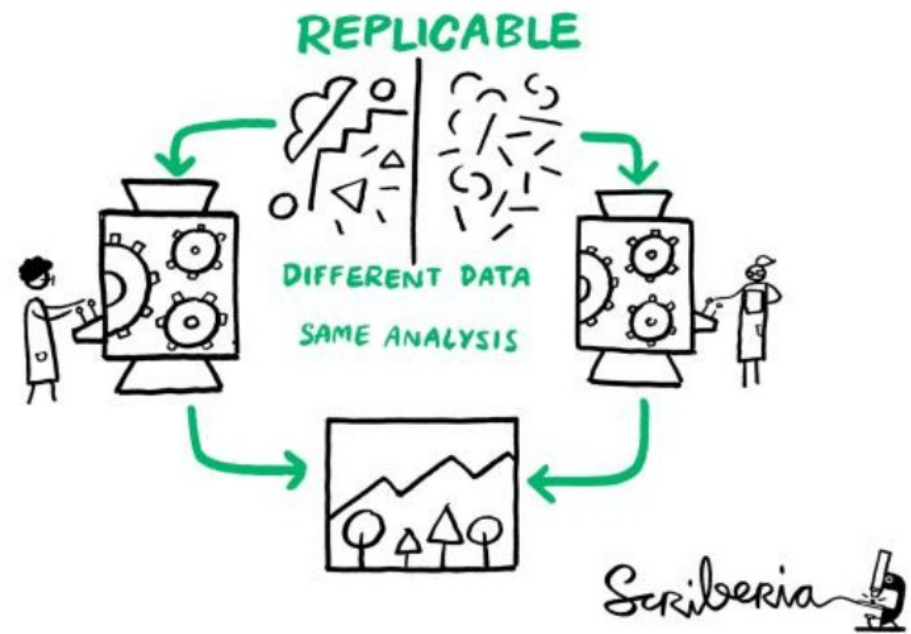
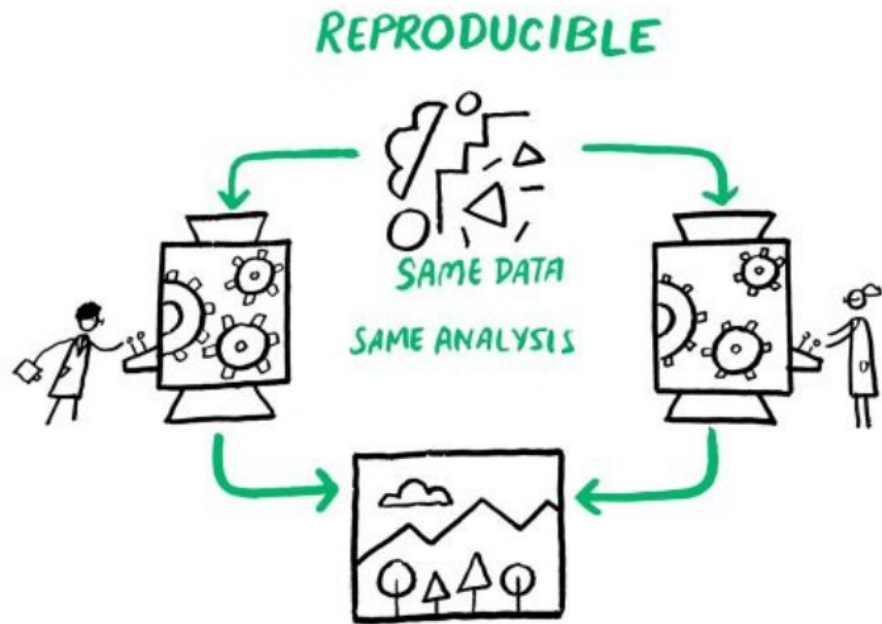
RESEARCH PRIORITIES

Shining Light into Black Boxes

A. Morin¹, J. Urban², P. D. Adams³, I. Foster⁴, A. Sali⁵, D. Baker⁶, P. Sliz^{1,*}

“I have begun to think that no one ought to publish biometric results, without lodging a well arranged and well bound manuscript copy of all his data, in some place where it should be accessible, under reasonable restrictions, to those who desire to verify his work.”

Francis Galton (1901), *Biometrika* 1:1, pp. 7-10.



open microbiome analysis frameworks



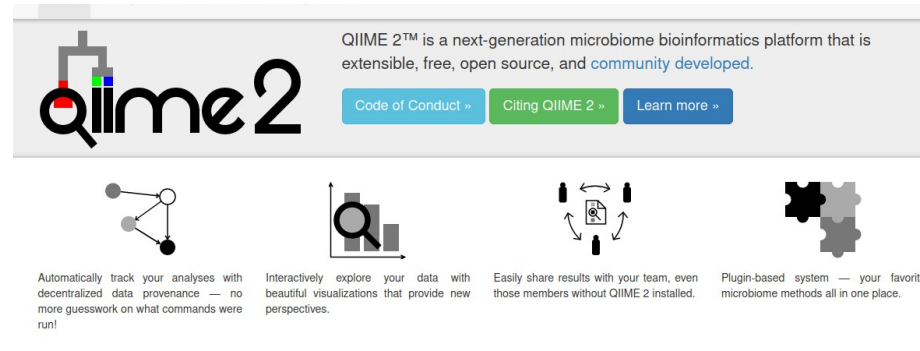
mothur
Download Wiki Forum Blog R bloggers Facebook

Welcome to the website for the mothur project, edited by Dr. Patrick Schloss and the research group in the Department of Microbiology & Immunology at The University of Michigan. This project seeks to develop a single piece of open-source, reusable software to fit the bioinformatics needs of the microbial ecology community. In February 2009 we released the first version of mothur, which had accelerated versions of the popular DOTUR and SONS programs. mothur has gone on to become one of the most cited bioinformatics tools for analyzing 16S rRNA gene sequences. Stop inside the wiki and user forum and learn how you can use mothur to process data generated by Sanger, PacBio, IonTorrent, 454, and Illumina (MiSeq/HiSeq). If you would like to contribute code to the project feel free to download the source code and make your own improvements. Alternatively, if you have an idea or a need, but lack the programming expertise, let us know through the forum and we'll add it to the queue of features we would like to add.

Subscribe to the mothur mailing list
email address
Subscribe

Department of Microbiology & Immunology
The University of Michigan Medical School
The University of Michigan

This site is maintained by Pat Schloss
© 2008-2015



QIIME 2™ is a next-generation microbiome bioinformatics platform that is extensible, free, open source, and community developed.

Code of Conduct » Citing QIIME 2 » Learn more »

Automatically track your analyses with decentralized data provenance — no more guesswork on what commands were run!

Interactively explore your data with beautiful visualizations that provide new perspectives.

Easily share results with your team, even those members without QIIME 2 installed.

Plugin-based system — your favorite microbiome methods all in one place.



PeerJ >

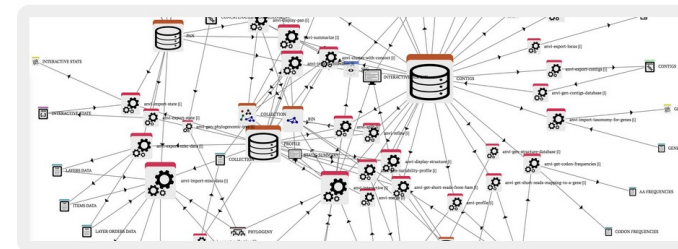
Anvi'o: an advanced analysis and visualization platform for 'omics data

Research article | Bioinformatics | Biotechnology | Computational Biology | Genomics | Microbiology

A. Murat Eren^{1,2}, Özcan C. Esen¹, Christopher Quince³, Joseph H. Vineis¹, Hilary G. Morrison¹, Mitchell L. Sogin¹, Tom O. Delmont¹

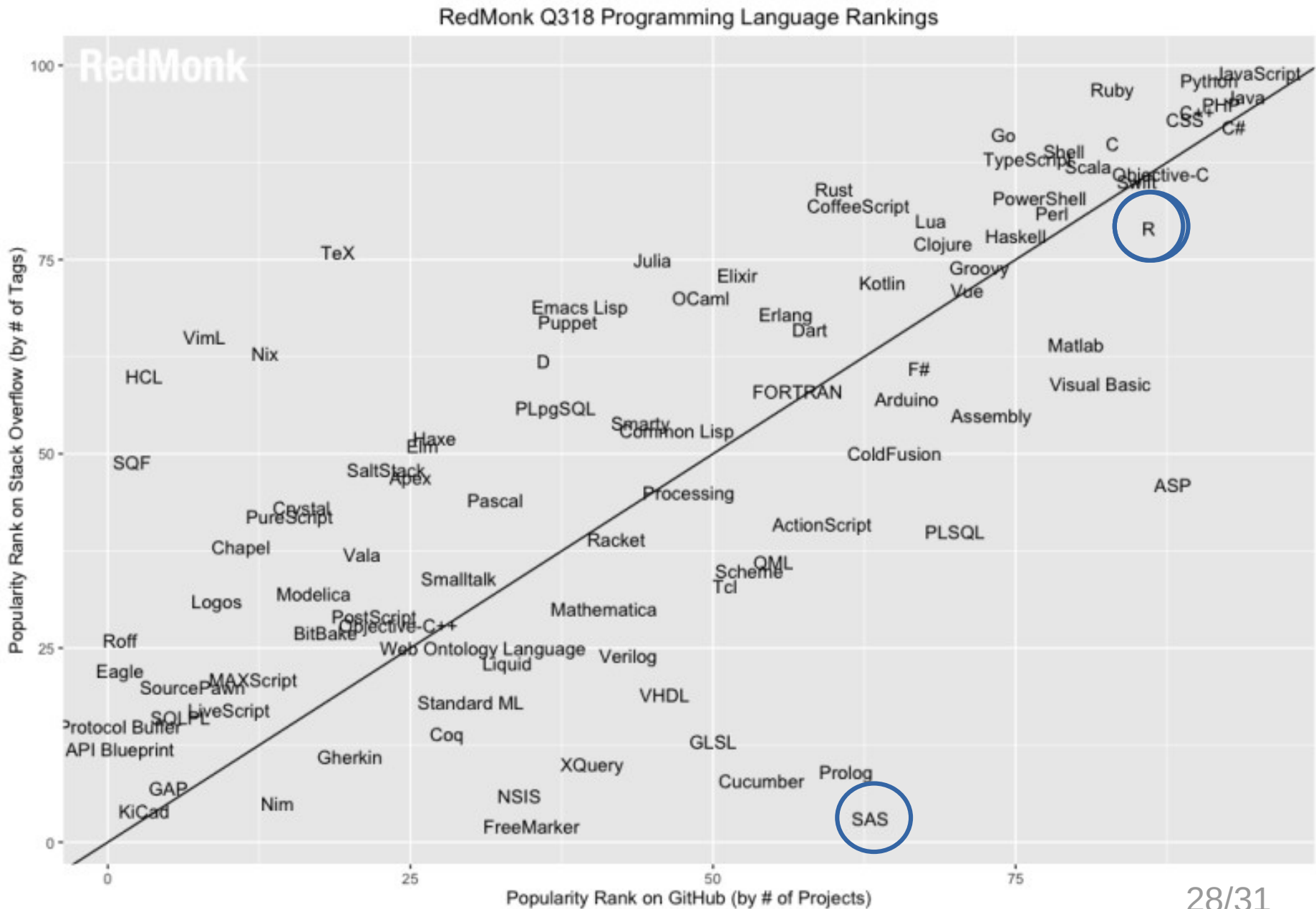
Published October 8, 2015

Anvi'o in a nutshell



Anvi'o is an open-source, community-driven analysis and visualization platform for 'omics data.

Cultures of open data science collaboration



Day 1 (Times in CET)

Lectures (45 min + 15 min breaks)

9:15-10:00 - **Welcome & introduction** - Leo Lahti, Associate professor (UTU)

10:15-11:00 - **Metagenomics** - Katariina Pärnänen, Postdoctoral researcher (UTU)

11:15-12:00 - **Metabolomics** - Pande Putu Erawijantari, Postdoctoral researcher (UTU)

12:15-13:00 - **Multimomics** - Leo Lahti, Associate professor (UTU)

13:00-14 - **Lunch** break

Practical session

14:15-17:00 - Tuomas Borman and Chouaib Benchraka, Research assistants (UTU)

Data import and data structures

Microbiome data exploration & visualization

1 Overview

1.1 Introduction

1.2 Learning goals

1.3 Acknowledgments

2 Program

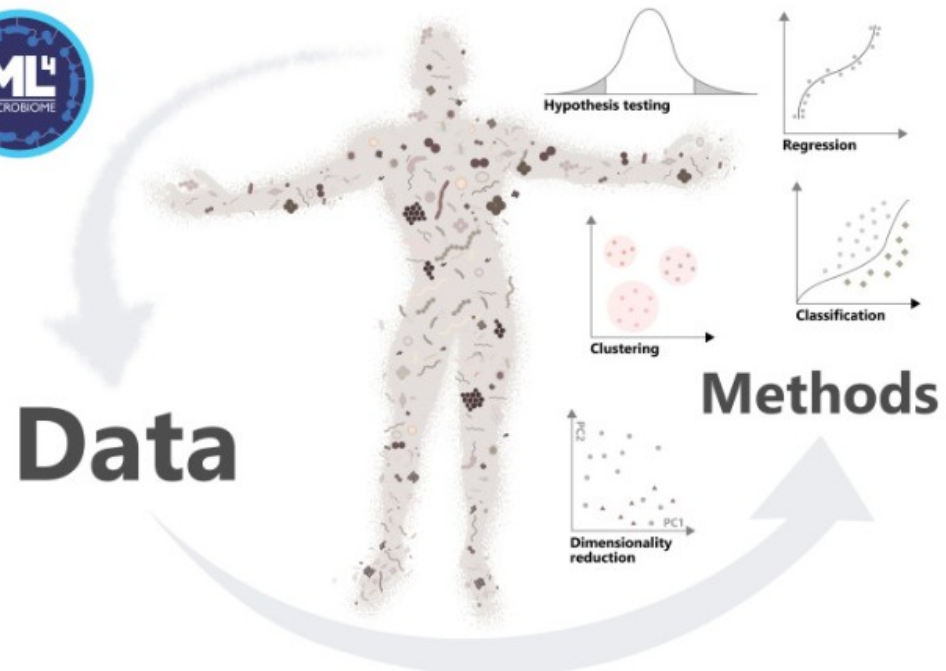
2.1 Day 1

2.2 Day 2

3 Getting started

3.1 Checklist (before the workshop)

Welcome to the multi-omics data analysis workshop



Chapter 3 Getting started

3.1 Checklist (before the workshop)

Install the following software in advance in order to avoid unnecessary delays and leaving more time for the workshop contents.

- [R \(version >4.1.0\)](#)
- [RStudio](#); choose “Rstudio Desktop” to download the latest version. Optional but preferred. For further details, check the [Rstudio home page](#).
- For Windows users: [Rtools](#); Follow the instructions to install the toolkit. This might be required to compile some of the packages required on this course.
- Install and load the required R packages

3.2 Support and resources

For online support on installation and other matters, you can join us at:

- Users: [miaverse Gitter channel](#)
- Developers: [Bioconductor Slack](#) #microbiomeexperiment channel (ask for an invitation)

3.3 Installing and loading the required R packages

This section shows how to install and load all required packages into the R session. Only uninstalled packages are installed.
