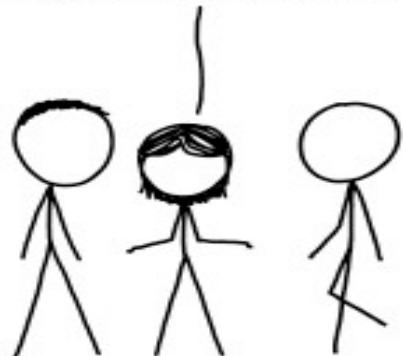


Open data science

OUR FIELD HAS BEEN
STRUGGLING WITH THIS
PROBLEM FOR YEARS.



STRUGGLE NO MORE!
I'M HERE TO SOLVE
IT WITH ALGORITHMS!



SIX MONTHS LATER:
WOW, THIS PROBLEM
IS REALLY HARD.



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.





Statistical and Machine Learning Techniques in Human Microbiome Studies: Contemporary Challenges and Solutions

Isabel Moreno-Indias^{1,2*}, Leo Lahti³, Miroslava Nedyalkova⁴, Ilze Elbere⁵, Gennady

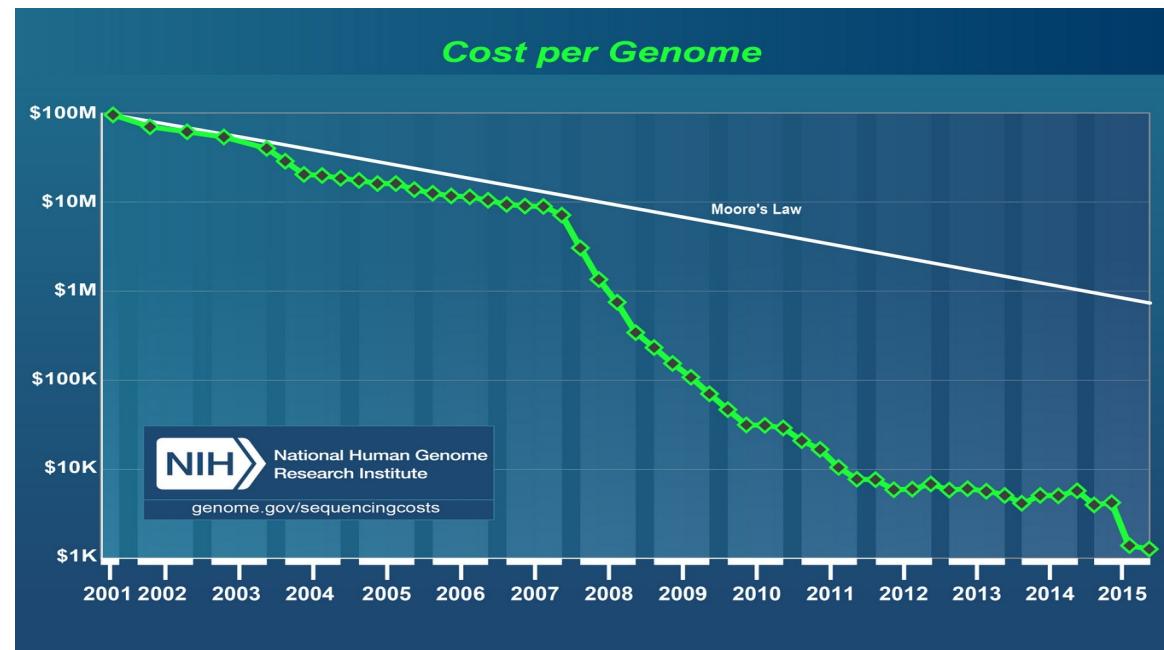
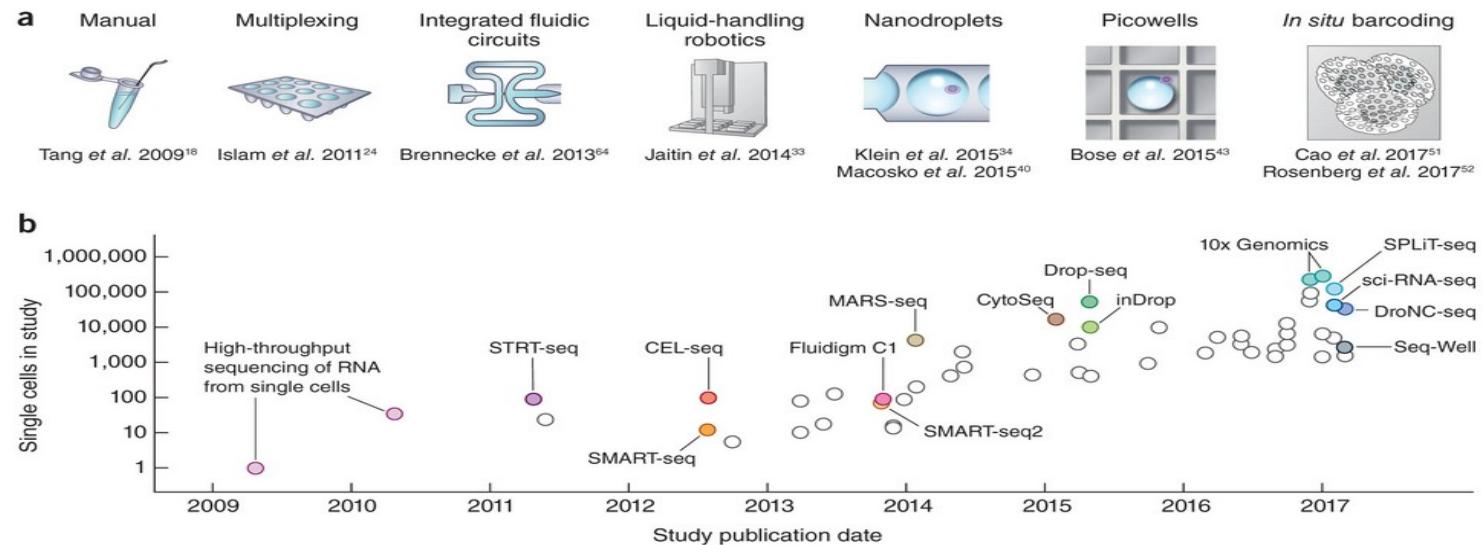


Figure 1: Scaling of scRNA-seq experiments.

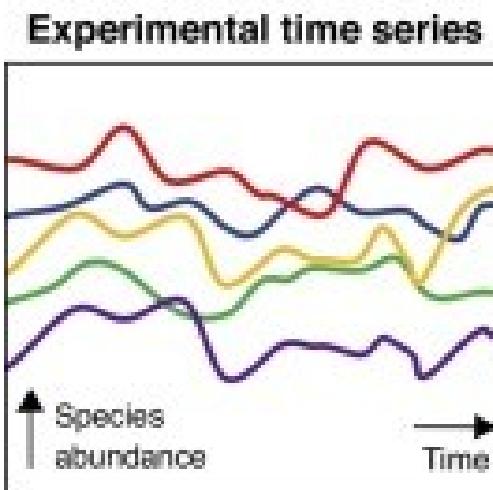
From: [Exponential scaling of single-cell RNA-seq in the past decade](#)



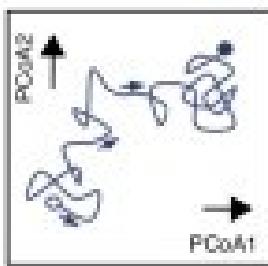
a) Key technologies that have allowed jumps in experimental scale. A jump to ~ 100 cells was enabled by sample multiplexing, and then a jump to $\sim 1,000$ cells was achieved by large-scale studies using integrated fluidic circuits, followed by a jump to several thousands of cells with liquid-handling 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.

Microbial communities as dynamical systems

Didier Gonze ^{1, 2}✉, Katharine Z Coyte ^{3, 4}, Leo Lahti ^{5, 6, 7}, Karoline Faust ⁵✉



Visualization



Mathematical modeling

$$\frac{dX_i}{dt} = X_i \left(b_i + \sum_{j=1}^N a_{ij} X_j \right)$$

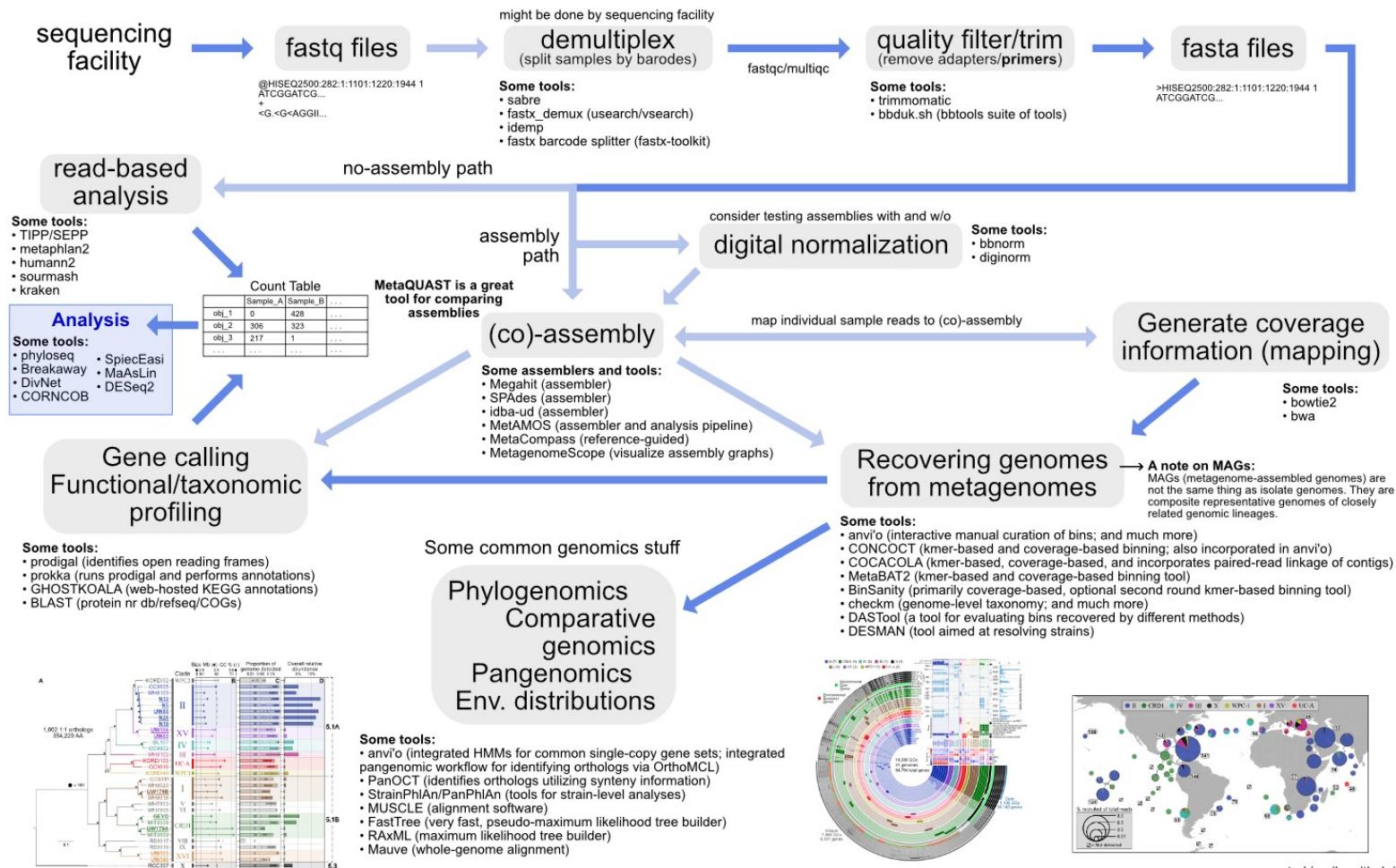
Dynamical properties



- Stability
- Alternative states
- Response to perturbation
- Stochasticity

Overview of generic* metagenomics workflow

* This is generic; specific workflows can vary on the order of steps here and how they are done.



Lee, (2019). Happy Belly Bioinformatics: an open-source resource dedicated to helping biologists utilize bioinformatics. Journal of Open Source Education, 4(41), 53, <https://doi.org/10.21105/jose.00053>

Happy Belly Bioinformatics
JOSE 10.21105/jose.00053

 AstroBioMike
 Orcid: 0000-0001-7750-9145

The influence of hidden researcher decisions in applied microeconomics

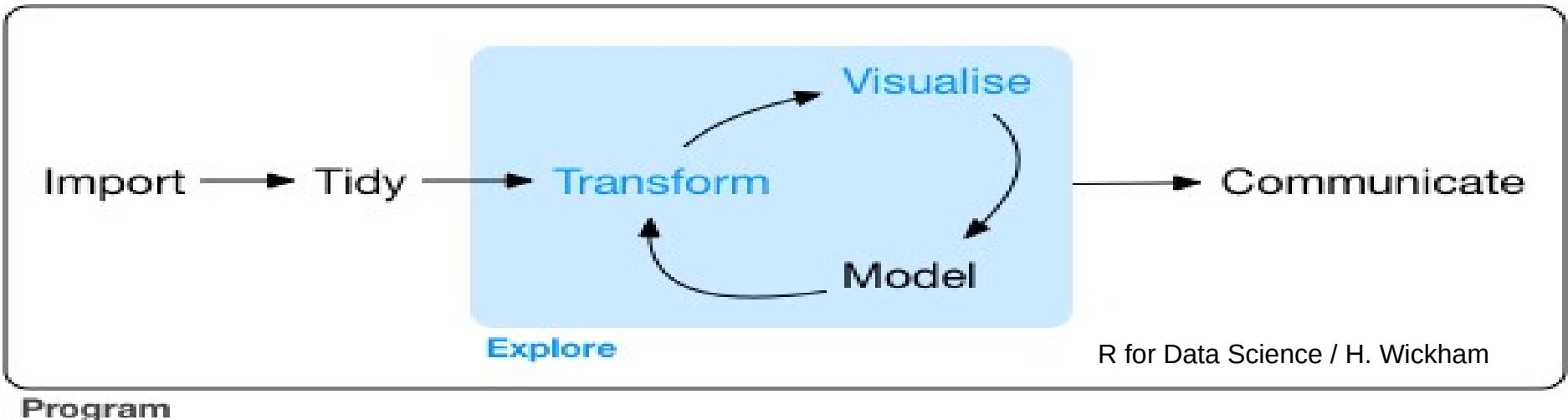
Nick Huntington-Klein ✉, Andreu Arenas, Emily Beam, Marco Bertoni, Jeffrey R. Bloem, Pralhad Burli,
Naibin Chen, Paul Grieco, Godwin Ekpe, Todd Pugatch, Martin Saavedra, Yaniv Stopnitzky

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.

Data science workflow



REVISED Bioconductor Workflow for Microbiome Data Analysis: from raw reads to community analyses [version 2; peer review: 3 approved]

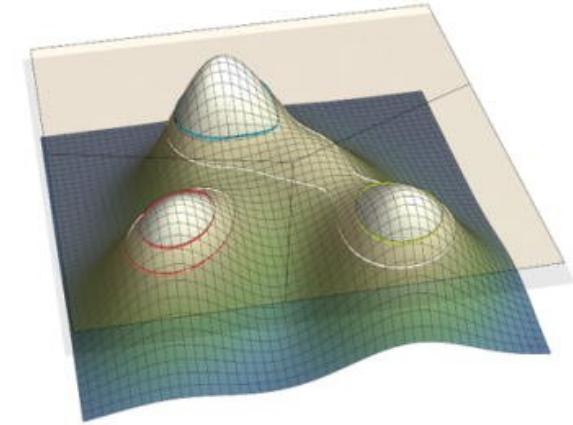
Ben J. Callahan¹, Kris Sankaran¹, Julia A. Fukuyama¹, Paul J. McMurdie²,  Susan P. Holmes



This article is included in the [Bioconductor](#) gateway.

How to choose a correct model?

→ a community typing example



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

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

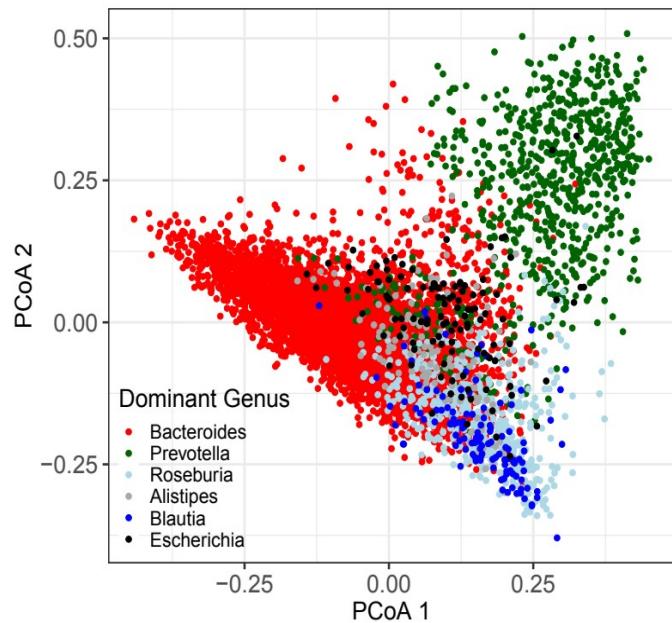
Clustering method

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

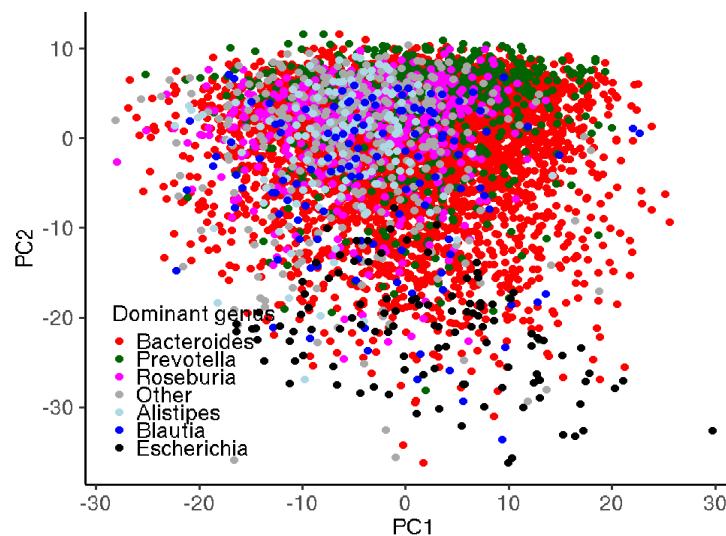
Regulation

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

PCoA + Bray-Curtis



PCA + Aitchison



Reproducible Research: Enterotype Example

Susan Holmes and Joey McMurdie

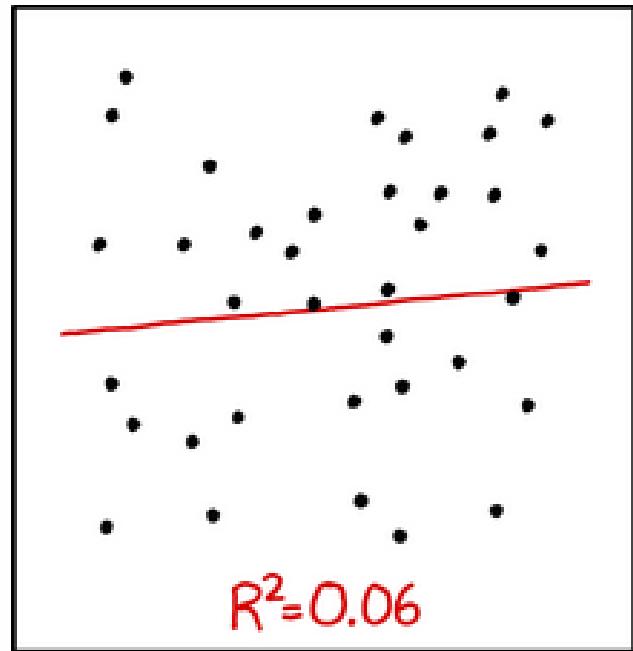
<http://statweb.stanford.edu/~susan/papers/EnterotypeRR.html>

[Comment on this paper](#)

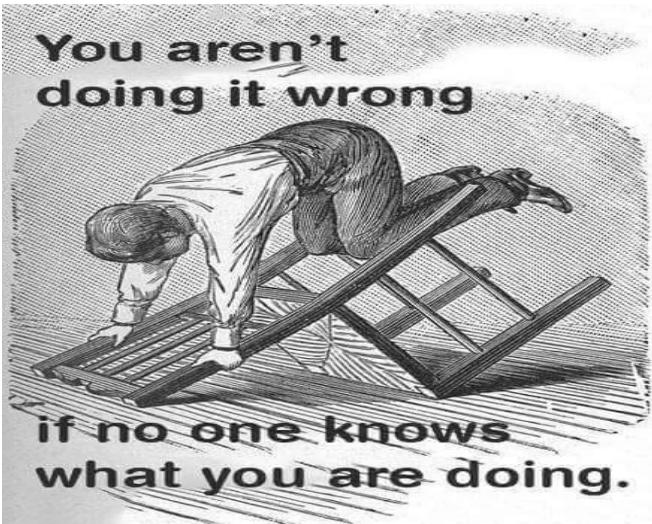
Taxonomic Signatures of Long-Term Mortality Risk in Human Gut Microbiota

Aaro Salosensaari, Ville Laitinen, Aki Havulinna, Guillaume Meric, Susan Cheng, Markus Perola, Liisa Valsta, Georg Alfthan, Michael Inouye, Jeremie D. Watrous, Tao Long, Rodolfo Salido, Karenina Sanders, Caitriona 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>

How we choose which model to apply?



I DON'T TRUST LINEAR REGRESSIONS WHEN IT'S HARDER
TO GUESS THE DIRECTION OF THE CORRELATION FROM THE
SCATTER PLOT THAN TO FIND NEW CONSTELLATIONS ON IT.



"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.

```
int getRandomNumber()
{
    return 4; // chosen by fair dice roll.
              // guaranteed to be random.
}
```

<http://web.stanford.edu/class/cs109l/unrestricted/images/>

RESEARCH PRIORITIES

Shining Light into Black Boxes

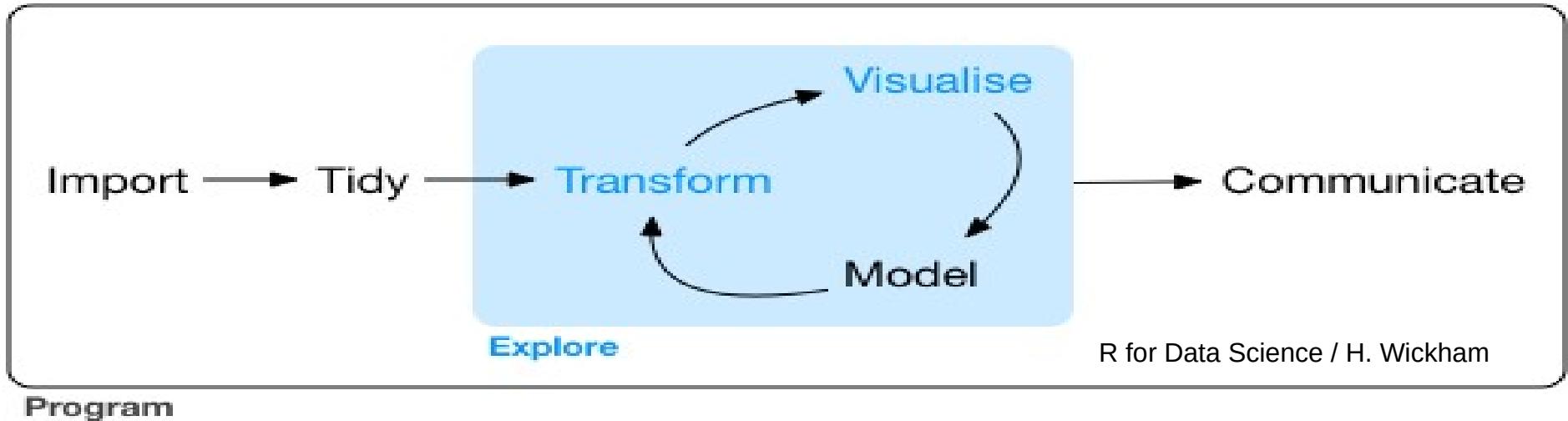
A. Morin¹, J. Urban², P. D. Adams³, I. Foster⁴, A. Sali⁵, D. Baker⁶, P. Sliz^{1,*}

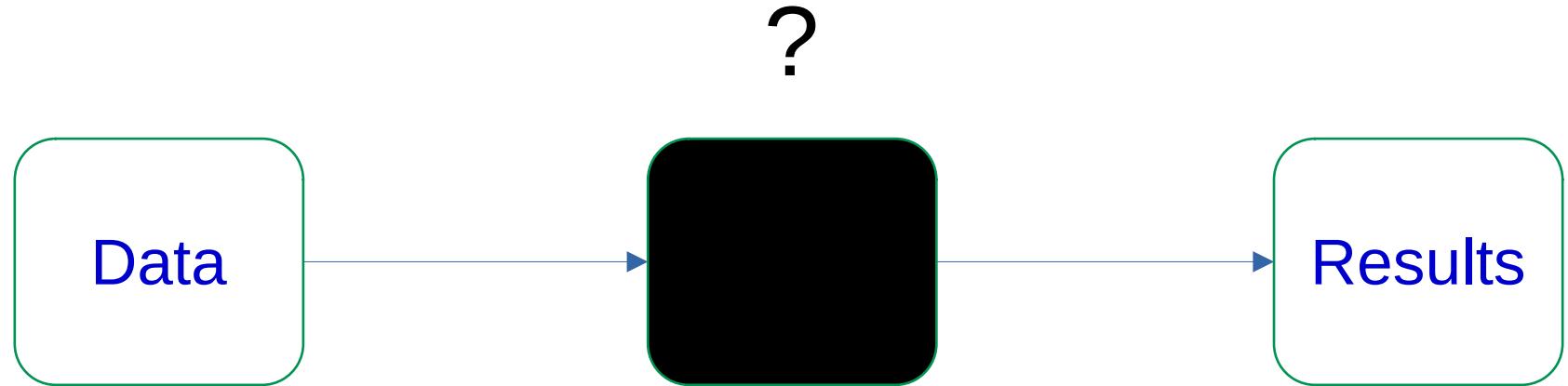


Data silo



Data science workflow





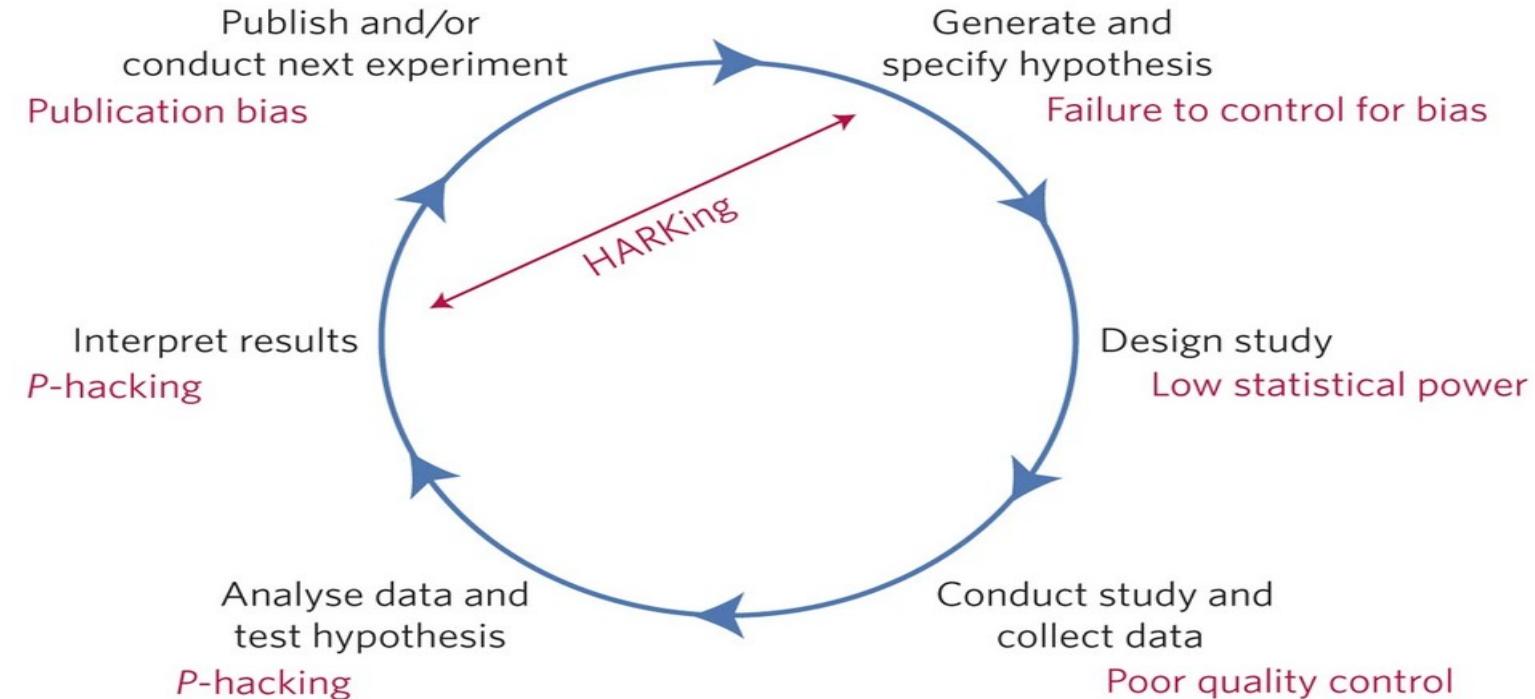
RESEARCH PRIORITIES

Shining Light into Black Boxes

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

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.

OPEN ACCESS

ESSAY

898,944

VIEWS

1,119

CITATIONS

Why Most Published Research Findings Are False

John P. A. Ioannidis

Published: Aug 30, 2005 • DOI: 10.1371/journal.pmed.0020124

OPEN ACCESS

ESSAY

898,944

VIEWS

Why Most Published Research Findings Are False

John P. A. Ioannidis

Published: Aug 30, 2005 • DOI: 10.1371/journal.pmed.0020124

How to Make More Published Research True

John P. A. Ioannidis 

Published: October 21, 2014 • DOI: 10.1371/journal.pmed.1001747

Open reporting and communication were part of academic culture since the early days

BJHS 45(2): 165–188, June 2012. © British Society for the History of Science 2012
doi:10.1017/S0007087412000064 First published online 20 March 2012

Openness versus secrecy? Historical and historiographical remarks

KOEN VERMEIR*



Source: Wikimedia Commons / Public domain

Alchemy & algorithms: perspectives on the philosophy and history of open science

Research Ideas and Outcomes 3:e13593, 2017

▼ Leo Lahti, Filipe da Silva, Markus Petteri Laine, Viivi Lähteenaja, Mikko Tolonen

Beyond Open Access - The Changing Culture of Producing and Disseminating Scientific Knowledge

Heidi Laine

Leo Lahti

Anne Lehto



Home > Open Science > UNESCO Recommendation on Open Science

Open Science



UNESCO Recommendation on Open Science

The UNESCO Recommendation on Open Science was adopted by the General Conference of UNESCO at its 41st session, in November 2021.

- UNESCO Recommendation on Open Science
English | Français

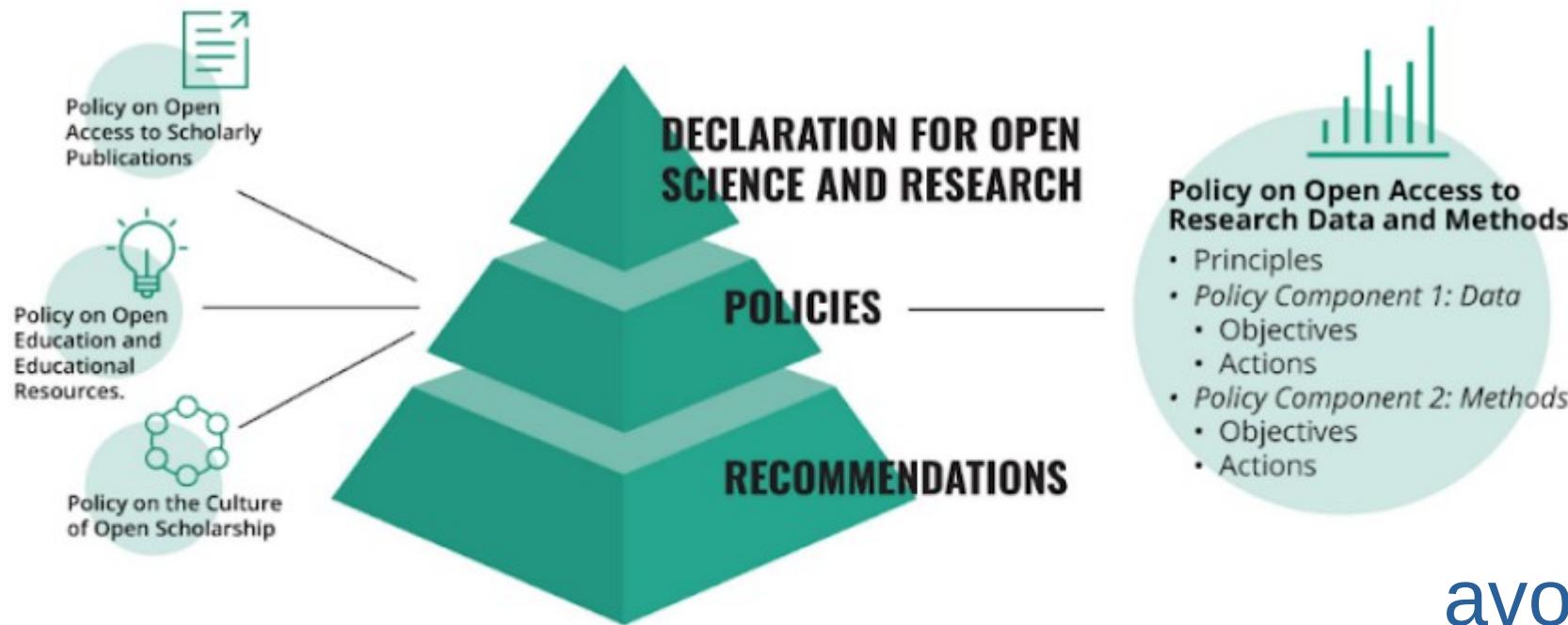
[Home](#)[UNESCO Recommendation on Open Science](#)[Multistakeholder Consultations on Open Science](#)



Avoin tiede

Picture 1. The policy in relation to other national open science documents.

National Policy on Open Science National policy and executive plan by the higher education and research community for 2021-2025



avointiede.fi

**PRINCIPLE 1: RESEARCH DATA AND METHODS
SHALL BE MANAGED, OPENED AND USED
RESPONSIBLY AND APPROPRIATELY.**

**PRINCIPLE 2: RESEARCHERS HAVE ACCESS TO
INFRASTRUCTURES AND SERVICES THAT ENABLE
RESPONSIBLE DATA MANAGEMENT, AND THESE
ARE DEVELOPED FURTHER IN AN ECONOMICALLY
SUSTAINABLE WAY, TAKING INTO ACCOUNT THE
RESEARCHERS' NEEDS.**

**PRINCIPLE 3: THE RESEARCHER'S MERITS IN
THE PROMOTION OF GOOD DATA MANAGEMENT,
WORK RELATED TO RESEARCH DATA AND
METHODS, AND THE APPROPRIATE OPENING OF
RESEARCH DATA AND METHODS ARE VALUED
AND CAN SUPPORT THE RESEARCHER'S CAREER.**

open data science ecosystems

mothur

Download Wiki Forum Blog Riggemores Facebook

Welcome to the website for the mothur project, initiated by Dr. Patrick Schloss and his research group at the Department of Microbiology & Immunology at the University of Michigan. This project seeks to develop a single piece of open-source, expandable software to fit the bioinformatics needs of the microbial community. In February 2009 we released the first version of mothur, which had accelerated versions of the most common bioinformatics programs. mothur has gone on to become one of the most cited bioinformatics tools for analyzing 16S rRNA gene experiments. They now have a mailing list, forum and team how you can use mothur to process data from samples from Bangladesh, India, and Bhutan (MothurGang). 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 of a need, but lack the programming expertise, let us know through the forums and we will 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
© 2008-2019
The University of Michigan

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 >](#)

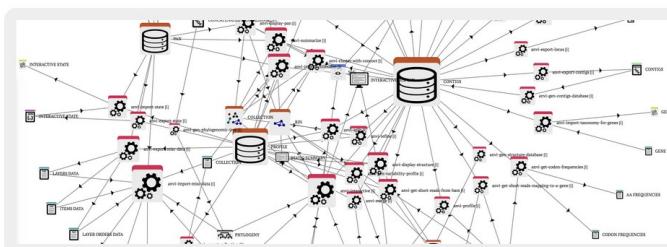
Anvi'o: an advanced analysis and visualization platform for 'omics data

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.

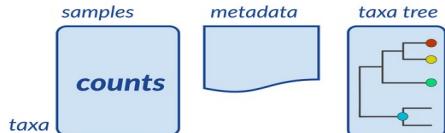
Example workflow - microbiome.github.io

Figure by Domenick Braccia (EuroBioC 2020)

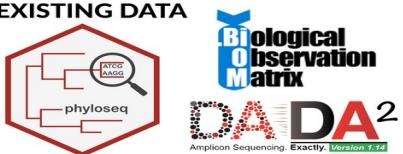
Import Data

This workflow starts with either raw data directly from relative abundance estimation or taxonomic classification OR pre-existing data objects from widely used software.

RAW DATA

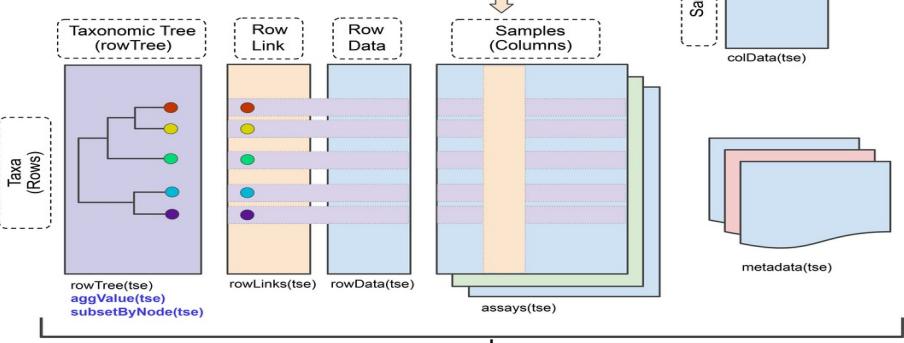


EXISTING DATA



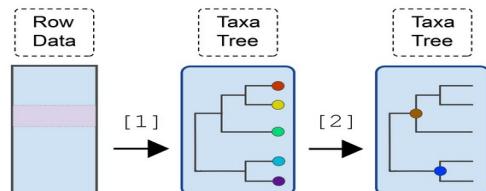
The TreeSE object

The tse object is uniquely positioned to support the next generation of microbiome data manipulation and visualization.



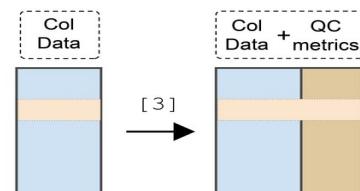
The mia Pipeline

Accessing Taxonomic Info.



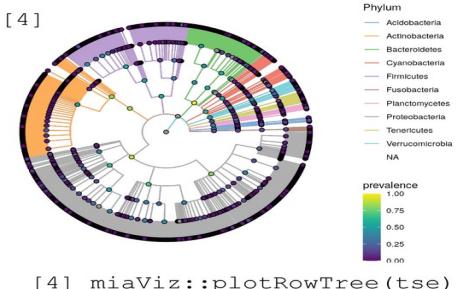
```
[1] mia::addTaxonomyTree(tse)  
[2] TreeSE::aggValue(tse)
```

Quality Control



```
[3] scatter::addPerCellQC(tse)
```

Visualizing with miaViz

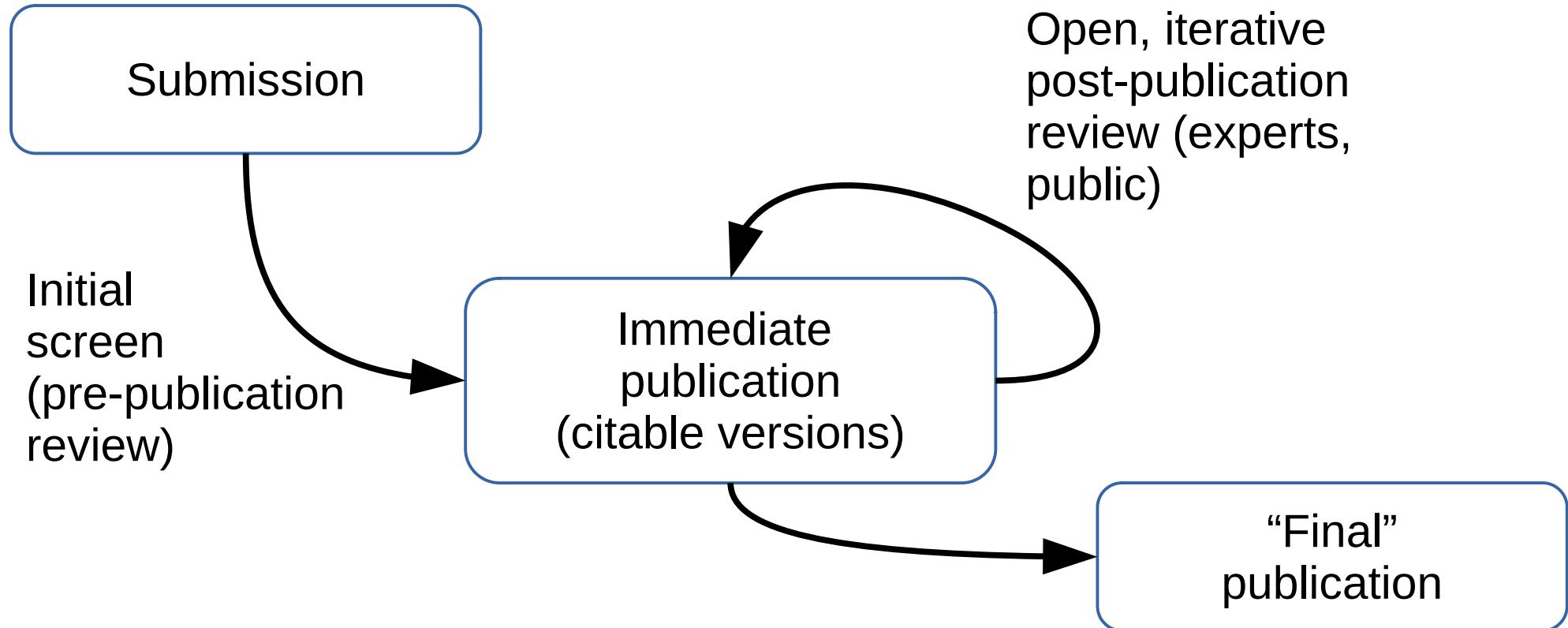


European Bioconductor Meeting 2020

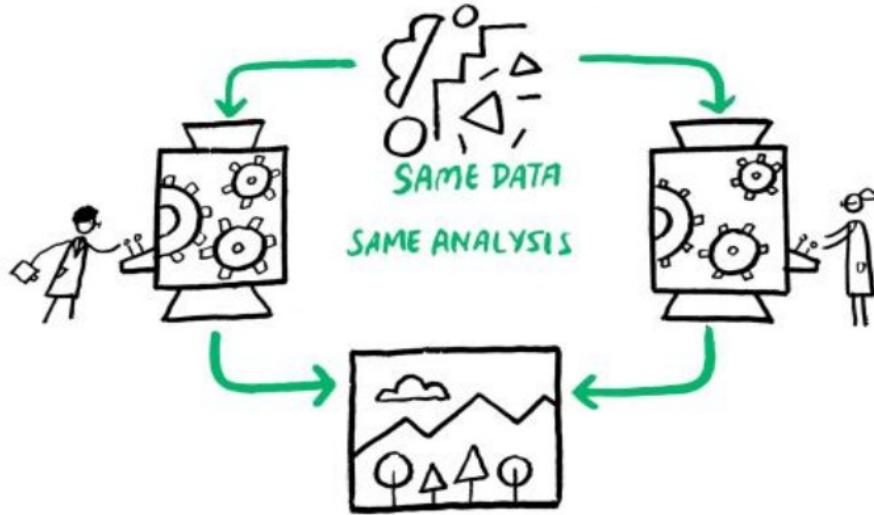
- Where: Virtual Conference
- When: 14-18 December 2020
- On twitter: #EuroBioC2020



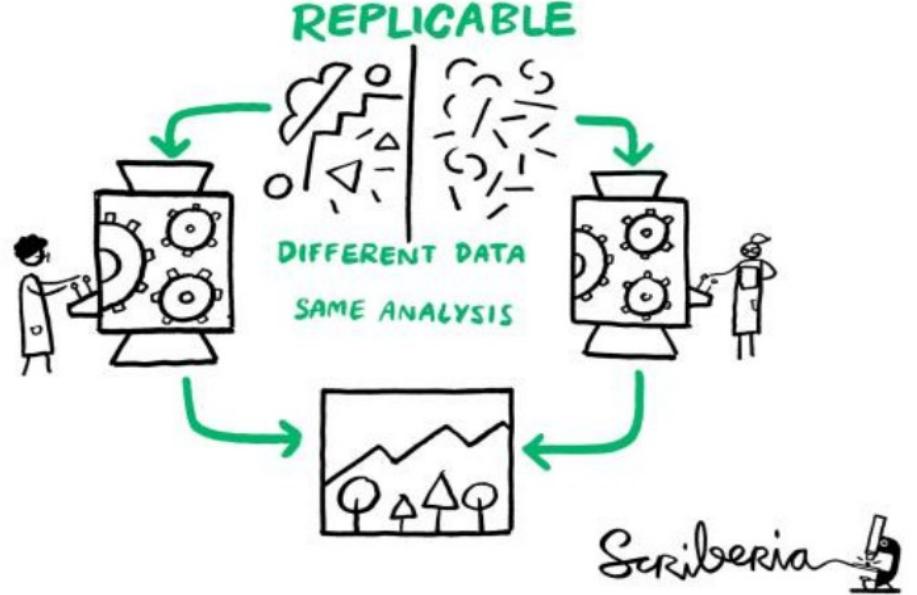
Transparency: pre- vs. post-publication review



REPRODUCIBLE



REPLICABLE



Reusability

PLOS COMPUTATIONAL BIOLOGY

OPEN ACCESS

EDUCATION

A Quick Guide to Software Licensing for the Scientist-Programmer

Andrew Morin, Jennifer Urban, Piotr Sliz

Published: July 26, 2012 • <https://doi.org/10.1371/journal.pcbi.1002598>



Software citation principles

Arfon M. Smith^{1,*}, Daniel S. Katz^{2,*}, Kyle E. Niemeyer^{3,*}
FORCE11 Software Citation Working Group

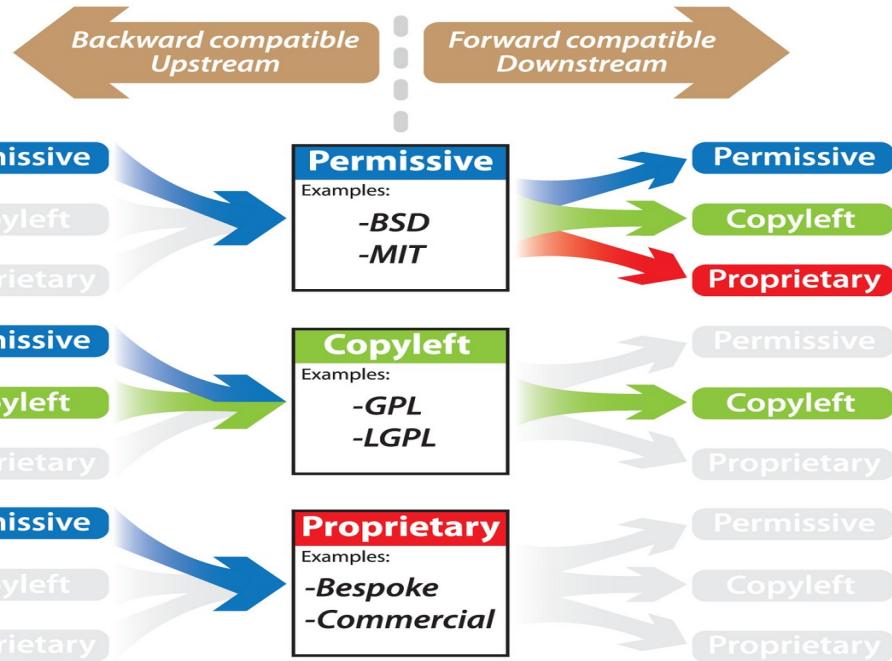
¹ GitHub, Inc., San Francisco, California, United States

² National Center for Supercomputing Applications & Electrical and Computer Department & School of Information Sciences, University of Illinois at Urbana-Champaign, Urbana, Illinois, United States

³ School of Mechanical, Industrial, and Manufacturing Engineering, Oregon State University, Corvallis, Oregon, United States

* These authors contributed equally to this work.

MIT License



Copyright (c) <year> <copyright holders>

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the "Software"), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

Varying cultures of open collaboration



The demise of alchemy provides further evidence, if further evidence were needed, that what marks out modern science is not the conduct of experiments (alchemists conducted plenty of experiments), but the formation of a *critical community capable of assessing discoveries and replicating results*. Alchemy, as a clandestine enterprise, could never develop a community of the right sort. Popper was right to think that science can flourish only in an open society.

The Invention of Science: A New History of the Scientific Revolution, by David Wootton



A family of alchemists at work, an engraving by Philip Galle, after a painting by Pieter Bruegel the Elder, published by Hieronymus Cock, c.1558.

transparency

replicability & reproducibility

reusability

collaborative research