

<b>Manuscript Number:</b>	GIGA-D-17-00317R1	
<b>Full Title:</b>	Experimenting with Reproducibility: a case study of Robustness in Bioinformatics	
<b>Article Type:</b>	Review	
<b>Funding Information:</b>	Institut Pasteur (FR)	Not applicable
	H2020 Health	Not applicable
	Institut National des Sciences de l'Univers, Centre National de la Recherche Scientifique (FR)	Not applicable
	Université Paris Diderot (FR)	Not applicable
	Conny-Maeva Charitable Foundation	Not applicable
	Cognacq-Jay Foundation	Not applicable
	Orange	Not applicable
	Fondation pour la Recherche Médicale	Not applicable
	GenMed Labex	Not applicable
	BioPsy Labex	Not applicable
<b>Abstract:</b>	<p>Reproducibility has been shown to be limited in many scientific fields. This question is a fundamental tenet of the scientific activity, but the related issues of reusability of scientific data are poorly documented. Here, we present a case study of our difficulties to reproduce a bioinformatics method [1] although code and data were available. First, we tried to re-run the analysis with the code and data provided by the authors. Second, we reimplemented the whole method in a Python package to avoid dependency on a MATLAB license and ease the execution of the code on HPCC (High-Performance Computing Cluster). Third, we assessed reusability of our reimplementation and the quality of our documentation, testing how easy it would be to start from our implementation to reproduce the results. In a second section, we propose solutions from this case study and other observations to improve reproducibility and research efficiency at the individual and collective level.</p> <p>While finalizing our code, we created case specific documentation and tutorials for the associated Python package StratiPy. Readers are thus invited to experiment our reproducibility case study by generating the two confusion matrices of Fig 3 (see more in 2.2.2).</p> <p>Here we decided to propose two options: 1) a step-by-step process to follow in a Jupyter/IPython notebook; or 2) a Docker container ready to be built and run. Availability: last version of StratiPy (Python) with two examples of reproducibility and dataset are available at GitHub [2] and Zenodo [3].</p>	
<b>Corresponding Author:</b>	Yang-Min KIM Institut Pasteur Paris, Île-de-France FRANCE	
<b>Corresponding Author Secondary Information:</b>		
<b>Corresponding Author's Institution:</b>	Institut Pasteur	
<b>Corresponding Author's Secondary Institution:</b>		
<b>First Author:</b>	Yang-Min KIM	
<b>First Author Secondary Information:</b>		
<b>Order of Authors:</b>	Yang-Min KIM	

	Jean-Baptiste Poline
	Guillaume Dumas
<b>Order of Authors Secondary Information:</b>	
<b>Response to Reviewers:</b>	----- see revision letter -----
<b>Additional Information:</b>	
<b>Question</b>	<b>Response</b>
Are you submitting this manuscript to a special series or article collection?	No
<p><b>Experimental design and statistics</b></p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	No
<p>If not, please give reasons for any omissions below.</p> <p>as follow-up to "<b>Experimental design and statistics</b></p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p> <p>"</p>	<p>This review is especially based on the reproducibility issues. The illustrated method is directly based on the original paper by Hofree et al, 2013. We nevertheless described any specific variations.</p>
<p><b>Resources</b></p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.</p>	Yes

<p>Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	
<p><b>Availability of data and materials</b></p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in <a href="#">publicly available repositories</a> (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	<p>Yes</p>

# Experimenting with Reproducibility: a case study of Robustness in Bioinformatics

## Authors

Yang-Min KIM<sup>1,2,3,4,\*</sup>, Jean-Baptiste Poline<sup>5,6</sup>, Guillaume Dumas<sup>1,2,3,4</sup>

<sup>1</sup>Institut Pasteur, Human Genetics and Cognitive Functions Unit, Paris, France, <sup>2</sup>CNRS UMR 3571 Genes, Synapses and Cognition, Institut Pasteur, Paris, France, <sup>3</sup>University Paris Diderot, Sorbonne Paris Cité, Paris, France, <sup>4</sup>Centre de Bioinformatique, Biostatistique et Biologie Intégrative (C3BI, USR 3756 Institut Pasteur and CNRS), Paris, France, <sup>5</sup>Montreal Neurological Institute, Brain Imaging Center, Ludmer Center, McGill University, <sup>6</sup>Henry H. Wheeler Jr. Brain Imaging Center, Helen Wills Neuroscience Institute, University of California, Berkeley, California, USA

\*To whom correspondence should be addressed

Correspondence: [yang-min.kim@pasteur.fr](mailto:yang-min.kim@pasteur.fr); [jbpoline@gmail.com](mailto:jbpoline@gmail.com); [guillaume.dumas@pasteur.fr](mailto:guillaume.dumas@pasteur.fr)

## Abstract

Reproducibility has been shown to be limited in many scientific fields. This question is a fundamental tenet of the scientific activity, but the related issues of reusability of scientific data are poorly documented. Here, we present a case study of our difficulties to reproduce a bioinformatics method [1] although code and data were available. First, we tried to re-run the analysis with the code and data provided by the authors. Second, we reimplemented the whole method in a Python package to avoid dependency on a MATLAB license and ease the execution of the code on HPCC (High-Performance Computing Cluster). Third, we assessed reusability of our reimplementation and the quality of our documentation, testing how easy it would be to start from our implementation to reproduce the results. In a second section, we propose solutions from this case study and other observations to improve reproducibility and research efficiency at the individual and collective level.

While finalizing our code, we created case specific documentation and tutorials for the associated Python package *StratiPy*. Readers are thus invited to experiment our reproducibility case study by generating the two confusion matrices of Fig 3 (see more in 2.2.2).

Here we decided to propose two options: 1) a step-by-step process to follow in a Jupyter/IPython notebook; or 2) a Docker container ready to be built and run.

**Availability:** last version of *StratiPy* (Python) with two examples of reproducibility and dataset are available at GitHub [2] and Zenodo [3].

## 30 Keywords

- 31 • Reproducibility
- 32 1 • Robustness
- 33 2 • Reusability
- 34 3 • Network Based Stratification (NBS)
- 35 4 • Standard consensus dataset
- 36 5 • Cancer

## 37 1 Background

38 17 The collective endeavor of science depends on researchers being able to replicate the work of others. In a recent survey of 1,576  
39 18 researchers, 70% of them admitted having difficulty in reproducing experiments proposed by other scientists [4]. For 50%, this  
40 19 reproducibility issue even concerns their own experiments. Despite the growing attention on the replication crisis in science [5,6],  
41 20 this controversial subject is far from being new: already in the 17th century, scientists criticized the air pump invented by physicist  
42 21 Robert Boyle because it was too complicated and expensive to build [7].

43 22 Several concepts for reproducibility in computational science are closely associated [8,9]. Here we define them as mentioned by  
44 23 K. Whitaker [9]: obtaining the same results using same data and same code is **Reproducibility**; if code is different, it is  
45 24 **Robustness**. If we used different data but with the same code, it is **Replicability**. Lastly, using different data and different code is  
46 25 referred as **Generalizability**. Here we will primarily elaborate on **Reproducibility** and **Robustness**, and acknowledge that new  
47 26 datasets or hardware environment introduce additional hurdles [10]. Reproducibility is a key first step, for instance, among the  
48 27 400 algorithms published during the major artificial intelligence conferences, only 6% offered the code [11]. Even when authors  
49 28 provide data and code, the outcome can vary either marginally or fundamentally [12]. Tackling irreproducibility in bioinformatics  
50 29 thus requires considerable effort beyond code and data availability, an effort that is still poorly recognized in the current  
51 30 publication based research community. In most cases, there is a significant gap between apparent executable work (Fig 1 - i.e.  
52 31 above water portion of iceberg) and necessary effort in practice (Fig 1 - i.e. full iceberg). Such effort is nevertheless necessary to  
53 32 increase the consistency of the literature and *efficiency* of the scientific research process by making research products easily  
54 33 reusable.

## 2 Reproducibility and Robustness in bioinformatics: a case study

### 2.1 Reproducibility: from MATLAB to MATLAB, OS and environment

Our team studies Autism Spectrum Disorders (ASD), a group of neurodevelopmental disorders well known for its heterogeneity. One of the current challenges of our research is to uncover homogeneous subgroups of patients (i.e. stratification) with more precise clinical outcomes, improving their prognosis and treatment [13,14]. An interesting stratification method was recently proposed in the field of cancer research [1], where the authors proposed to combine genetic profiles of patients' tumors with protein-protein interaction networks to uncover meaningful homogeneous subgroups, a method called Network Based Stratification (NBS).

Before using NBS method on our data, we studied the method by reproducing results from the original study. We are very grateful to the main authors who kindly provided online all the related data and code, and gave us invaluable input upon request. The authors of this study did much more to help reproduce their results than is generally done. Despite their help we experienced a number of difficulties that we document here, hoping that this report will help future researchers to improve the reproducibility of results and reusability of research products.

The first step of our project was to execute the original method code with the given data: reproducibility (Table 1). To improve execution speed, the original authors used a library for MATLAB on a Linux platform, using executable compiled code MEX file [15]: MTIMESX [16], a library allowing acceleration of large matrix multiplication. MEX files however are specific to the architecture and have to be recompiled for each Operating System (OS). Since our lab was using Mac OS X Sierra, the compilation of this MEX file into a mac64 binary required a new version of MTIMESX. It was also necessary to install and to configure properly OpenMP [17], a development library for parallel computing. After this, the original MATLAB code was successfully run in our environment.

	Code	Data	Technical issues	Other issues
<b>Reproducibility</b>	Same: MATLAB	Same	<b>OS:</b> MEX file specificity linked to OS (e.g. Linux → OSX)	
<b>Robustness</b>	MATLAB → Python	Same	<p><b>File format:</b> we can load sparse matrices from <i>.mat</i> file but cannot save them into HDF5 using h5py package</p> <p><b>Default parameters:</b> linkage method use for the hierarchical clustering</p> <ul style="list-style-type: none"> <li>• MATLAB (MathWorks): UPGMA (average)</li> <li>• Python (SciPy): single</li> </ul>	<ul style="list-style-type: none"> <li>• Metadata structure</li> <li>• Important parameter value not explained in the original paper</li> <li>• Remaining discarded work ('code ruins') and traces of debugging</li> </ul>
<b>Reproducibility of Robustness</b>	Same: Python	Same	<b>OS:</b> Numpy package and BLAS library compiled for specific OS (e.g. OSX → Linux)	Documentation

## **Table 1: Technical problems encountered during our reproducibility and robustness case study.**

These issues are classic but may not be overcome by researchers with little experience in compilation or installation issues. For these reasons alone, many individuals may turn down the opportunity of reusing code. The next part will focus on code re-implementation, a procedure, which can help understanding the method, but is even more time consuming.

## **2.2 Robustness: from MATLAB to Python, language and organization**

To fully master the method, we developed a complete open source toolkit of genomic stratification in Python [2]. Python is also an interpreted programming language, but contrary to MATLAB is free of use and has a GPL-compatible license [18], which fosters both robustness and generalizability. Recoding in another language in a different environment will lead to be some unavoidable problems such as variation in low level libraries (e.g. glibc): it is likely that the outcomes will vary even if the same algorithm is implemented [19]. In addition, we rely on Python packages to perform visualization or linear algebra computations (e.g. Matplotlib, SciPy, NumPy [20–22]), and results may depend on these packages versions. Python is currently in a transitional period between two major versions 2 and 3. We chose to write the code in Python 3, which is the current recommendation.

### **2.2.1 Metadata and File formats**

Even if the original code could be run, we had to handle several file formats to check and understand the structure of the original data. For instance the data was provided by The Cancer Genome (TCGA) [23] and made available in a MATLAB *.mat* file format as compressed data (sparse matrices). Thanks to SciPy, Python can load all versions through v7.2 MATLAB files, but to read v7.3 *.mat* files, we needed an HDF5 Python library. We decided to continue using Python's h5py package but Scipy's sparse matrices could not be stored in HDF5 format (Table 1). Moreover, the original authors had denoted download dates of patients' data of TCGA, thereby clarifying source of data. But in the absence of structural metadata, it was not always obvious how to interpret dataset variables (e.g. patient ID, gene ID, phenotype). Fig 2 shows an analogy between robustness issues and road transport: driving in a different environment (e.g. OS), we attempt to obtain identical results (i.e. to reach the same location) using the same input data (i.e. gasoline), but with different computational environment (i.e. cars), different implementation of the method (i.e. engine) and different programming languages (i.e. MATLAB and Python roads).

### **2.2.2 Codes and parameters**

Once the environment, file format and data issues were resolved, the code was finally executed. Unfortunately, "unexpected" results were obtained. One cause was the application of the hierarchical clustering step for which we used the clustering tools of

104 SciPy. Both SciPy and MATLAB (MathWorks) functions offer seven linkage methods, however, SciPy’s default option (single  
105 method) [24] differs from MATLAB’s default option (UPGMA or average method) [25], which was used in the original study  
106 (Table 1). Another cause for the variation in results is the value of one of the most important parameters of the method, the graph  
107 regulator factor, which was not clarified in the original paper. From the article, we believed that this factor had a constant value of  
108 1.0 until we found in the original code that its value varies across iterations and converges to an optimal value around 1800.  
109 Therefore, we initially obtained very different results from the original NBS (Fig 3 a) with heterogeneous subgroups. Once the  
110 optimal value was set up, we finally observed homogenous clusters (Fig 3 b). Moreover, during our attempts to run the original  
111 code to understand the causes of the errors, we realized that some parts of the code were not run anymore (e.g. discarded work,  
112 remaining traces of debugging) which made understanding the implementation harder.  
113  
114 To allow others to reproduce our results, we wrote some documentation and tutorials for the Python package *StratiPy* [2]. Readers  
115 are thus invited to experiment with reproducibility too by generating the two confusion matrices of Fig 3. They will use the  
116 different tools described in the following: GitHub, Docker, and Jupyter/IPython notebook.

### 116 2.2.3 Jupyter/IPython

117 During the re-coding process, we used an enhanced Python interpreter to debug: IPython, an interactive shell supporting both  
118 Python 2 and 3. Since the dataset is large and the execution takes a significant amount of time, we used IPython to re-run  
119 interactively some sub-sections of the script, which is one of the most helpful features. IPython can be integrated in the web  
120 interface Jupyter Notebook, offering an advanced structure for mixing code and documentation. While the Jupyter/IPython  
121 notebook was therefore initially convenient, it does not scale well to large programs and is not well adapted to versioning.  
122 However, ability of mixing code with document text is very useful for tutorials: a user of the code can read documentation  
123 (docstring), text explanations, and see how to run the code, explore parameters and visualize results in the browser. Our work on  
124 NBS, as related here, can be reproduced with a Jupyter/IPython notebook available on our GitHub [2]. One can find more  
125 examples and several helpful links on this “gallery of interesting Jupyter Notebooks” [26], which contains a section about  
126 “Reproducible academic publications”.

## 127 2.3 Reproducibility of Robustness: from Python to Python

128 Besides Jupyter/IPython notebooks, we used versioning tools like the Git code version control system (VCS) to document the  
129 development of our Python code. Git is arguably one of the most powerful VCS, allowing easy development of branches and  
130 helping us to work together as a distributed team (Paris, Berkeley, Montreal) on the same project. This project, *StratiPy*, is hosted  
131 on GitHub, a web-based Git repository hosting service [2]. While the original code was not available on GitHub, the main authors  
132 shared their code on a website. This should be sufficient for our purpose but makes it less easy to collaborate on code. While  
133 working on our GitHub repository, researchers from USA, India, China, and Europe contacted us about our robustness



134 experiment. Not only GitHub supports a better organization of projects, it also facilitates the collaboration of open-source  
135 software projects, thanks to its social network functions [27]. We adopted open source coding standards and learnt how to  
136 1 efficiently use Git and GitHub. Both required considerable training efforts on the short-term but brought clear benefits on the  
137 2 long-term, especially regarding collaboration and debugging.

138 3 We then attempted to re-run and reproduce the results we obtained on another platform. While the Python code was developed  
139 4 under Mac OS X Sierra (10.12), we used an Ubuntu 16.04.1 (Xenial) computer to test the Python implementation. Some  
140 5 additional issues emerged (Table 1). First, our initial documentation did not include the list of the required packages and  
141 6 instructions to launch the code. Second, the code was very slow to the extent that it was impractical to run it on a laptop because  
142 7 the Numpy package had not been compiled with BLAS (Basic Linear Algebra Subprograms) that speeds up low-level routines  
143 8 performing basic vector and matrix operations. Last, there was (initially) no easy way to check whether the results obtained on a  
144 9 different architecture were the expected ones. We added documentation and tests on the results files md5sum to solve this. To  
145 10 summarize, although the reuse and reproducibility of the results of the developed package were improved, these were far from  
146 11 being optimal.

## 147 12 **3 Potential solutions: from local to global**

### 148 13 **3.1 Act locally: simple practices and available tools**

149 14 Given the observed difficulties, in this section we draw some conclusions on this reproducibility case study experiment and  
150 15 suggest some tools and best practices. In addition, we suggest to follow the programming best practices of Wilson et al. such as  
151 16 modularizing and re-using code, unit testing, document design, data management, and project organization [5,28]. Sandve and  
152 17 colleagues [29] also suggest to keep the data provenance with recording all intermediate results.

#### 153 18 **3.1.1 Environment**

154 19 In 1995, Buckheit and Donoho were already thinking about reproducible research in computer science. Their motto was “When  
155 20 we publish articles containing figures which were generated by computer, we also publish the complete software environment  
156 21 which generates the figures” by offering a complete and free package (WaveLab) to reproduce the published output [30].  
157 22 Container and virtual machines technologies such as Docker [31], Vagrant [32], Singularity [33,34] (easily works in cluster  
158 23 environments) are becoming a standard solution to installation issues. These rely however on competencies that we think too few  
159 24 biologists possess today. While a container might encapsulate everything needed for a software execution, it could be hard to  
160 25 develop in a container. For instance, running Jupyter/IPython notebooks in Docker’s container requires knowledge on advanced  
161 26 port forwarding , which can be discouraging for some biologists. Therefore, we decided to propose two options in our example  
162 27 implementation of reproducibility: 1) a step-by-step process to follow in a Jupyter/IPython notebook; or 2) a Docker container

163 ready to be built and run. Nevertheless, mastering Docker –or other container tools– will become an important skill for  
164 computational reproducible researchers.

### 165 2 **3.1.2 Metadata**

166 4 Standard metadata are vital for an efficient documentation of both data and software. In our example, we still lack the standard  
167 6 lexicon to document the data as well as documenting the software. We however aim to follow the recommendations by Stodden *et*  
168 8 *al.* [35]: “*Software metadata should include, at a minimum, the title, authors, version, language, license, Uniform Resource*  
169 9 *Identifier/DOI, software description (including purpose, inputs, outputs, dependencies), and execution requirements*”. The more  
170 12 comprehensive is the metadata description, the more likely the reuse will be both efficient and appropriate [36].

### 171 15 **3.1.3 Write readable code**

172 17 Anyone who has spent time to understand someone else’s code would advise some simple basic rules to help make the code  
173 19 readable and understandable.

174 21 First, the structure of the program should be clear and easily accessible. Second, good concise code documentation and naming  
175 23 convention will help readability. Third, the code should not contain left-overs of previously tested solutions. When a solution  
176 25 takes a long time to compute, an option to store it locally can be proposed. Using standard coding and documentation conventions  
177 27 (e.g. PEP 8 and PEP 257 in Python [37,38]) with detailed comments and references of papers makes the code more accessible.  
178 29 When an algorithm is used, any modification from the original reference should be explained and discussed in the article as well  
179 32 as in the code. We advocate for researchers to write code “for their colleagues”, hence, ask for the opinion and review of co-  
180 34 working or partner laboratories. Furthermore, the collaboration between researchers working on different environments can more  
181 36 easily isolate reproducibility problems. In the future, journals may consider review of code as part of the standard review process  
182 38 [39].

### 183 41 **3.1.4 Test the code**

184 43 To check if the code is yielding a correct answer, software developers associate test suites (unit tests or integration tests) with their  
185 45 software. While we developed only a few tests in this project, we realize that this has a number of advantages, such as checking if  
186 47 the software installation seems correct, check if updates in the code or in the operating system impact the results, etc. In our case,  
187 49 we propose to check for the integrity of the data and for the results of some key processing.

## 188 3.2 Think globally: from education to community standards

### 189 3.2.1 Training the new generation of scientists to digital tools and practices

190 2 The training in coding is still too limited for biologists. Often, it is self-training, from searching answers on Stack Overflow or  
191 3 equivalent. Despite efforts by organizations such as Software [40] or Data Carpentry [41] and the growing demand for ‘data  
192 4 scientists’ in life science, university training on coding practices is still not enough generalized. The difficulty to access and  
193 5 understand code may lead to applying code blindly without checking the validity of the results: often, scientists may prefer to  
194 6 believe that the results are correct because of the time that would be needed to check the validity of the results. Mastering a  
195 7 package such that results are truly understood can take a long time, as it was the case in our experiment.

196 8 Academia could instruct young scientists best practices for reproducibility. For instance, Hothorn and Leisch organized a  
197 9 reproducibility workshop gathering mostly PhD students and young postdocs specialized in bioinformatics and biostatistics. Then  
198 10 they evaluated 100 random sample papers from *Bioinformatics* [6]. Their study revealed how such a workshop can raise young  
199 11 scientists awareness about “*what makes reproduction easy or hard at first hand*”. Indeed, they found out that only a third of the  
200 12 original papers and two-thirds for applications notes had given access to the source code of software used.

### 201 26 3.2.2 Standard consensus dataset and testing ecosystem

202 27 We propose here that bioinformatics methods publications are systematically accompanied with a test dataset, code source and  
203 28 some basic tests. As the method is tested on new datasets, the number of tests of the method would increase in number and cover a  
204 29 wider range of applications. We give a first example with our NBS re-implementation. We develop below how this could  
205 30 generalize and what would be the benefit for the scientific community.

206 31 A schematic overview of a possible testing ecosystem is shown in Fig 4. The core of this system would be a set of standard  
207 32 consensus datasets used to validate methods. For instance in the field of machine learning, standard image databases are widely  
208 33 used for training and testing (e.g. MNIST for handwritten digits [42]). In the case of our proposal, data could be from different  
209 34 categories such as binary, text, image (shown as folders in different colors, Fig 4 b), and sub-categories to introduce criteria such  
210 35 as size, quantitative/qualitative, discrete/continuous using a tagging system. Datasets could be issued from simulations or from  
211 36 acquisition, and would validate a method on a particular component. This testing ecosystem will help scientists that cannot release  
212 37 their data because of privacy issues (Fig 4 a.1) (although these can often be overcome) but also give access to data and tests to a  
213 38 wide community including establishments with weak financial means.

214 39 We divide those who interact with scientific software or analysis code in two broad categories. First, the authors (“A”) who  
215 40 propose a method and need to verify its validity and usefulness with open and/or private – data. Second, the users (“U”, e.g.  
216 41 developers, engineers, bioinformaticians) who need to test and evaluate the proposed methods with other data.

217 When authors propose a new method, this method could have a reproducibility profile, which will progressively be built by  
218 authors and users (Fig 4 b.3, b.4). The information of which method does or does not work with well identified data is crucial for  
219 1 future work. During the optimization of a project, the software code and associated documentation should be accessible to foster  
220 2 collaboration on additional use cases and data. When the work achieve some level of maturity, a full fledge article can be posted  
221 3 on a preprint servers such as bioRxiv [43,44] and be associated with a GitHub repository by digital object identifiers (DOI). With  
222 4 considerable effort, Stodden *et al.* conducted a reproducibility study on 204 random articles of Science: despite some availability  
223 5 of the code, it had often been changed after publication, causing difficulties in replication [45]. In our proposed testing ecosystem,  
224 6 users will be able to launch reproducibility projects more easily thanks to code and article versioning.  
225 7 Users who test and approve reproducibility on original or new data could be credited and recognized by the scientific and  
226 8 developer communities (i.e. Stack Overflow, GitHub). This testing ecosystem could thus facilitate collaborations between  
227 9 methodology development and biological research communities.

## 228 21 **4 Conclusion and perspective**

229 24 In the 19<sup>th</sup> century, Pasteur introduced a detailed "Methods" section in his report: this advanced approach was necessary to  
230 25 reproduce his experiments and became new norms in the philosophy of science [46]. Today with the advent of computational  
231 26 science, the reproducibility issue is seen as a growing concern. Before reusing a published method, we attempted to reproduce the  
232 27 initial results and recoded the method to have a deep understanding of it. The investment in time to verify a previously published  
233 28 method can be as or even more important than the work needed to publish a new paper. Despite the willingness of the authors to  
234 29 share their tool and help us in our work, we have faced computational reproducibility and robustness problems due to  
235 30 compatibility between environments, programming languages and software versions, choice of parameters, data identification, etc.  
236 31 In addition to individual effort to write well documented and readable code, we recommend to use online repositories and tools to  
237 32 help other scientists in their exploration of the method: Docker for environment standardization, GitHub for code version  
238 33 management, and Jupyter notebooks for demonstration and tutorial [26,27,31]. We suggest to adopt such practices, not only for  
239 34 writing code but also manuscripts [5]. At the community level, we should enhance the cooperation between academic education  
240 35 and industry to foster a new generation of well-trained scientists in software development. For instance, Academia-Industry  
241 36 Software Quality & Testing summit (AISTQ) organizes conference in order to encourage collaboration between Academia and  
242 37 Industry [47]. This is also in line with the work of the Software or Data carpentry organizations [40,41]. Here, we propose a  
243 38 testing ecosystem where the community uses standard, well-identified datasets to validate tools and their versions. The scope of a  
244 39 proposed method could be continuously evaluated on new datasets. Eventually, data and software can be versioned and cited to  
245 40 give credit to the individuals who have contributed to these building blocks of Science. This testing ecosystem is not only a  
246 41 reproducibility validation tool; it is an attempt to make research product more reusable using online platforms and open source  
247 42 tools, beyond the publication of a PDF file. Such system could be seen as a generalization of already existing workflow systems

248 such as Galaxy or GATK, integrating data and software provenance [48,49]. Some top-down initiatives already provide some  
249 incentives for such a process i.e. Horizon 2020 (H2020) [50] project of the European Commission (EC) mandates open access of  
250 research data, while respecting security and liability. H2020 supports OpenAIRE [51] a technical infrastructure of the open access,  
251 which allows the interconnection between projects, publications, datasets, and author information across Europe. Thanks to  
252 common guidelines, OpenAIRE interoperates with other web-based *generalist* scientific data repositories such as Zenodo, hosted  
253 by CERN, which allows combining data and GitHub repository using DOI. The Open Science Framework also hosts data and  
254 software for a given project [52]. Respecting standard guidelines to transparently communicate the scientific work is a key step  
255 towards tackling irreproducibility and insures a robust scientific endeavor.

## 256 Key points

- 257 • Main barrier for reproducibility is in the lack of compatibility between environments, programming languages, software  
258 versions, etc.
- 259 • At the individual level, the key is in research practices such as well written, tested and documented code, and well  
260 curated data and the use of online repositories and collaborative tools.
- 261 • At the community level, we propose a testing ecosystem where standard consensus datasets are used to validate new  
262 methods and foster their generalizability.

## 263 Declarations

### 264 • **Ethics approval and consent to participate**

265 We used the uterine endometrial carcinoma data as they were downloaded on January 1<sup>st</sup>, 2013 from the The Cancer Genome  
266 Atlas (TCGA) portal by Hofree and colleagues [1].

### 267 • **Consent for publication**

268 Not applicable

### 269 • **Availability of data and material**

270 Last version of StratiPy (Python) with two examples of reproducibility and dataset are available at GitHub [2].

271 Zenodo DOI: <https://doi.org/10.5281/zenodo.1042546> [3]

### 272 • **Competing interests**

273 The authors declare that they have no competing interests.

274 • **Funding**

275 This work was supported by:

- 276 1 ○ Institut Pasteur (<http://dx.doi.org/10.13039/501100003762>)  
2
- 277 3 ○ H2020 Societal Challenges (<http://dx.doi.org/10.13039/100010676>)  
4
- 278 5 ○ Centre National de la Recherche Scientifique (<http://dx.doi.org/10.13039/501100004794>)  
6
- 279 7 ○ Université Paris Diderot (<http://dx.doi.org/10.13039/501100005736>)  
8
- 280 9 ○ Conny-Maeva Charitable Foundation  
10
- 281 11 ○ Cognacq-Jay Foundation  
12
- 282 13 ○ Orange (<http://dx.doi.org/10.13039/501100003951>)  
14
- 283 15 ○ Fondation pour la Recherche Médicale (<http://dx.doi.org/10.13039/501100002915>)  
16
- 284 17 ○ GenMed Labex  
18
- 285 19 ○ BioPsy Labex.  
20  
21

286 • **Authors' contributions**

287 25 Y-M. K., J-B. P., and G.D. wrote the manuscript, Y-M. K. and G.D. developed the StratiPy module.  
26

288 27 All authors read and approved the final manuscript.  
28  
29

289 • **Acknowledgements**

290 32 We thank Thomas Rolland and Freddy Cliquet for sharing their advices and comments.  
33  
34

291 • **Authors' information**

292 38 **Yang-Min KIM**<sup>1,2,3,4,\*</sup> is a PhD student at Human Genetics and Cognitive Functions unit at the neuroscience department of the  
39  
293 40 Institut Pasteur in Paris. Her research on next-generation sequencing data and biological networks is focused on stratification of  
41  
294 42 patients with Autism.  
43

295 44 *Keywords:* Autism Spectrum Disorder; Network; Protein-Protein Interaction; Personalize Medicine; Network-Based Stratification;  
45  
296 46 Computational Biology.  
47

297 48  
49  
298 50 **Jean-Baptiste Poline**<sup>5,6</sup> is a researcher at Henry H. Wheeler Jr. Brain Imaging Center, Helen Wills Neuroscience Institute,  
51  
299 52 University of California, Berkeley, California, USA. His research focuses on neuroimaging methods, imaging-genetic biostatistics  
53  
300 54 and neuroinformatics.  
55

301 56 *Keywords:* Neuroinformatics; Statistical methods; Brain imaging; Imaging genomics.  
57  
302 58  
59  
60  
61  
62  
63  
64  
65

303 **Guillaume Dumas**<sup>1,2,3,4</sup> is research fellow of the Human Genetics and Cognitive Functions unit at the neuroscience department of  
304 the Institut Pasteur in Paris. His interdisciplinary work is at the cross-road of social psychology, cognitive neuroscience, and  
305 system biology.

306 *Keywords:* Open Science; Complex Systems; Computational Biology; Cognitive Science; Social Neuroscience; Autism Spectrum  
307 Disorder.

## 309 Reference

- 310<sup>13</sup> 1. Hofree M, Shen JP, Carter H, Gross A, Ideker T. Network-based stratification of tumor mutations. *Nat Methods* [Internet].  
311<sup>14</sup> 2013 [cited 2015 Aug 12];10. Available from: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3866081/>  
312<sup>16</sup> 2. StratiPy: Graph regularized nonnegative matrix factorization (GNMF) in Python [Internet]. GHFC; 2017. Available from:  
313<sup>17</sup> <https://github.com/GHFC/StratiPy>  
314<sup>19</sup> 3. Kim Yang-Min, Poline Jean-Baptiste, Dumas Guillaume. StratiPy [Internet]. Zenodo; 2017. Available from:  
315<sup>20</sup> <https://zenodo.org/record/1042546>  
316<sup>22</sup> 4. Baker M. 1,500 scientists lift the lid on reproducibility. *Nat News*. 2016;533:452.  
317<sup>24</sup> 5. Wilson G, Bryan J, Cranston K, Kitzes J, Nederbragt L, Teal TK. Good enough practices in scientific computing. *PLOS*  
318<sup>25</sup> *Comput Biol*. 2017;13:e1005510.  
319<sup>27</sup> 6. Hothorn T, Leisch F. Case studies in reproducibility. *Brief Bioinform*. 2011;12:288–300.  
320<sup>29</sup> 7. Shapin S, Schaffer S. *Leviathan and the Air-Pump: Hobbes, Boyle, and the Experimental Life* (New in Paper). Princeton  
321<sup>31</sup> University Press; 2011.  
322<sup>33</sup> 8. Peng RD. Reproducible research in computational science. *Science*. 2011;334:1226–7.  
323<sup>35</sup> 9. Whitaker K. Showing your working: a how to guide to reproducible research [Internet]. 2017. Available from:  
324<sup>36</sup> [https://figshare.com/articles/Showing\\_your\\_working\\_a\\_how\\_to\\_guide\\_to\\_reproducible\\_research/5443201](https://figshare.com/articles/Showing_your_working_a_how_to_guide_to_reproducible_research/5443201)  
325<sup>38</sup> 10. Nekrutenko A, Taylor J. Next-generation sequencing data interpretation: enhancing reproducibility and accessibility. *Nat Rev*  
326<sup>39</sup> *Genet*. 2012;13:667–72.  
327<sup>41</sup> 11. HutsonFeb. 15 M, 2018, Pm 12:30. Missing data hinder replication of artificial intelligence studies [Internet]. *Sci. AAAS*.  
328<sup>42</sup> 2018 [cited 2018 Mar 13]. Available from: [http://www.sciencemag.org/news/2018/02/missing-data-hinder-replication-artificial-](http://www.sciencemag.org/news/2018/02/missing-data-hinder-replication-artificial-intelligence-studies)  
329<sup>43</sup> [intelligence-studies](http://www.sciencemag.org/news/2018/02/missing-data-hinder-replication-artificial-intelligence-studies)  
330<sup>45</sup> 12. Herndon T, Ash M, Pollin R. Does high public debt consistently stifle economic growth? A critique of Reinhart and Rogoff.  
331<sup>46</sup> *Camb J Econ*. 2014;38:257–79.  
332<sup>48</sup> 13. Bourgeron T. From the genetic architecture to synaptic plasticity in autism spectrum disorder. *Nat Rev Neurosci*.  
333<sup>49</sup> 2015;16:551–63.  
334<sup>51</sup> 14. Loth E, Spooren W, Ham LM, Isaac MB, Auriche-Benichou C, Banaschewski T, et al. Identification and validation of  
335<sup>52</sup> biomarkers for autism spectrum disorders. *Nat Rev Drug Discov*. 2016;15:70–73.  
336<sup>54</sup> 15. Introducing MEX Files - MATLAB & Simulink - MathWorks France [Internet]. [cited 2017 Aug 18]. Available from:  
337<sup>55</sup> [https://fr.mathworks.com/help/matlab/matlab\\_external/introducing-mex-files.html?requestedDomain=www.mathworks.com](https://fr.mathworks.com/help/matlab/matlab_external/introducing-mex-files.html?requestedDomain=www.mathworks.com)  
338<sup>57</sup> 16. Tursa. MTIMESX - Fast Matrix Multiply with Multi-Dimensional Support - File Exchange - MATLAB Central [Internet].  
339<sup>58</sup> 2009 [cited 2017 Apr 24]. Available from: [http://fr.mathworks.com/matlabcentral/fileexchange/25977-mtimesx-fast-matrix-](http://fr.mathworks.com/matlabcentral/fileexchange/25977-mtimesx-fast-matrix-multiply-with-multi-dimensional-support)  
340<sup>60</sup> [multiply-with-multi-dimensional-support](http://fr.mathworks.com/matlabcentral/fileexchange/25977-mtimesx-fast-matrix-multiply-with-multi-dimensional-support)

- 341 17. tim.lewis. Specifications [Internet]. OpenMP. [cited 2017 Aug 18]. Available from: <http://www.openmp.org/specifications/>
- 342 18. Python Software Foundation. History and License — Python 3.6.1 documentation [Internet]. 2017 [cited 2017 Apr 24].  
343 Available from: <https://docs.python.org/3/license.html#licenses-and-acknowledgements-for-incorporated-software>
- 1  
344 2 19. Glatard T, Lewis LB, Ferreira da Silva R, Adalat R, Beck N, Lepage C, et al. Reproducibility of neuroimaging analyses across  
345 3 operating systems. *Front Neuroinformatics*. 2015;9:12.
- 4  
346 5 20. Michael Droettboom, Thomas A Caswell, John Hunter, Eric Firing, Jens Hedegaard Nielsen, Antony Lee, et al.  
347 6 matplotlib/matplotlib v2.2.2 [Internet]. Zenodo; 2018. Available from: <https://zenodo.org/record/1202077>
- 7  
348 8 21. Pauli Virtanen, Ralf Gommers, Evgeni Burovski, Travis E. Oliphant, David Cournapeau, Warren Weckesser, et al.  
349 9 scipy/scipy: SciPy 1.0.1 [Internet]. Zenodo; 2018. Available from: <https://zenodo.org/record/1206941>
- 10  
350 11 22. NumPy — NumPy [Internet]. [cited 2018 Apr 3]. Available from: <http://www.numpy.org/>
- 12  
351 13 23. TCGA [Internet]. Cancer Genome Atlas - Natl. Cancer Inst. [cited 2017 Apr 24]. Available from:  
352 14 <https://cancergenome.nih.gov/>
- 15  
353 16 24. Eads. Hierarchical clustering (scipy.cluster.hierarchy) — SciPy v0.19.0 Reference Guide [Internet]. 2007 [cited 2017 Apr 24].  
354 17 Available from: <https://docs.scipy.org/doc/scipy/reference/cluster.hierarchy.html>
- 18  
355 19 25. Hierarchical Clustering - MATLAB & Simulink - MathWorks France [Internet]. [cited 2017 Apr 24]. Available from:  
356 20 <https://fr.mathworks.com/help/stats/hierarchical-clustering-12.html>
- 21  
357 22 26. A gallery of interesting Jupyter Notebooks · jupyter/jupyter Wiki [Internet]. [cited 2017 Aug 18]. Available from:  
358 23 <https://github.com/jupyter/jupyter/wiki/A-gallery-of-interesting-Jupyter-Notebooks>
- 24  
359 25 27. Blischak JD, Davenport ER, Wilson G. A Quick Introduction to Version Control with Git and GitHub. *PLoS Comput Biol*.  
360 26 2016;12:e1004668.
- 27  
361 28 28. Wilson G, Aruliah DA, Brown CT, Chue Hong NP, Davis M, Guy RT, et al. Best Practices for Scientific Computing. Eisen  
362 29 JA, editor. *PLoS Biol*. 2014;12:e1001745.
- 30  
363 31 29. Sandve GK, Nekrutenko A, Taylor J, Hovig E. Ten Simple Rules for Reproducible Computational Research. Bourne PE,  
364 32 editor. *PLoS Comput Biol*. 2013;9:e1003285-4.
- 34  
365 35 30. Buckheit JB, Donoho DL. WaveLab and Reproducible Research. *Wavelets Stat* [Internet]. Springer, New York, NY; 1995  
366 36 [cited 2018 Mar 13]. p. 55–81. Available from: [https://link.springer.com/chapter/10.1007/978-1-4612-2544-7\\_5](https://link.springer.com/chapter/10.1007/978-1-4612-2544-7_5)
- 37  
367 38 31. Boettger C. An introduction to Docker for reproducible research, with examples from the R environment. *ACM SIGOPS*  
368 39 *Oper Syst Rev*. 2015;49:71–9.
- 40  
369 41 32. Introduction [Internet]. Vagrant HashiCorp. [cited 2017 Oct 13]. Available from: <https://www.vagrantup.com/intro/index.html>
- 42  
370 43 33. Singularity | Singularity [Internet]. [cited 2017 Oct 13]. Available from: <http://singularity.lbl.gov/>
- 44  
371 45 34. Kurtzer GM, Sochat V, Bauer MW. Singularity: Scientific containers for mobility of compute. *PLOS ONE*.  
372 46 2017;12:e0177459.
- 47  
373 48 35. Stodden V, McNutt M, Bailey DH, Deelman E, Gil Y, Hanson B, et al. Enhancing reproducibility for computational methods.  
374 49 *Science*. 2016;354:1240–1.
- 50  
375 51 36. Hill SL. How do we know what we know? Discovering neuroscience data sets through minimal metadata. *Nat Rev Neurosci*.  
376 52 2016;17:735–6.
- 53  
377 54 37. PEP 8 -- Style Guide for Python Code [Internet]. Python.org. [cited 2017 Aug 21]. Available from:  
378 55 <https://www.python.org/dev/peps/pep-0008/>
- 56  
379 57 38. PEP 257 -- Docstring Conventions [Internet]. Python.org. [cited 2017 Aug 21]. Available from:  
380 58 <https://www.python.org/dev/peps/pep-0257/>
- 59  
60  
61  
62  
63  
64  
65



- 381 39. Eglén SJ, Marwick B, Halchenko YO, Hanke M, Sufi S, Gleeson P, et al. Toward standard practices for sharing computer  
382 code and programs in neuroscience. *Nat Neurosci.* 2017;20:770–3.
- 383 40. Software Carpentry [Internet]. *Softw. Carpentry.* [cited 2017 Aug 22]. Available from: [http://software-  
384 carpentry.org/index.html](http://software-<br/>384 carpentry.org/index.html)
- 385 41. Data Carpentry [Internet]. *Data Carpentry.* [cited 2017 Aug 22]. Available from: <http://www.datacarpentry.org/>
- 386 42. MNIST handwritten digit database, Yann LeCun, Corinna Cortes and Chris Burges [Internet]. [cited 2017 Aug 23]. Available  
387 from: <http://yann.lecun.com/exdb/mnist/>
- 388 43. Bourne PE, Polka JK, Vale RD, Kiley R. Ten simple rules to consider regarding preprint submission. *PLOS Comput Biol.*  
389 2017;13:e1005473.
- 390 44. Preprints in biology. *Nat Methods.* 2016;13:277–277.
- 391 45. Stodden V, Seiler J, Ma Z. An empirical analysis of journal policy effectiveness for computational reproducibility. *Proc Natl  
392 Acad Sci.* 2018;115:2584–9.
- 393 46. Day RA, Gastel B. *Historical Perspectives. Write Publ Sci Pap Seventh Ed. ABC-CLIO; 2011. p. 6–8.*
- 394 47. Academia – Industry Software Quality & Testing summit - ISTQB® International Software Testing Qualifications Board  
395 [Internet]. [cited 2017 Aug 23]. Available from: [http://www.istqb.org/special-initiatives/istqb-conference-network-2istqb-  
396 conference-network-academia/academia-%E2%80%93industry-software-quality-testing-summit.html](http://www.istqb.org/special-initiatives/istqb-conference-network-2istqb-<br/>396 conference-network-academia/academia-%E2%80%93industry-software-quality-testing-summit.html)
- 397 48. Kanwal S, Khan FZ, Lonie A, Sinnott RO. Investigating reproducibility and tracking provenance – A genomic workflow case  
398 study. *BMC Bioinformatics.* 2017;18:337.
- 399 49. Karim MR, Michel A, Zappa A, Baranov P, Sahay R, Rebholz-Schuhmann D. Improving data workflow systems with cloud  
400 services and use of open data for bioinformatics research. *Brief Bioinform* [Internet]. [cited 2017 Jul 31]; Available from:  
401 <https://academic.oup.com/bib/article/doi/10.1093/bib/bbx039/3737318/Improving-data-workflow-systems-with-cloud>
- 402 50. Open Research Data in Horizon 2020 [Internet]. [cited 2017 Aug 23]. Available from:  
403 [https://ec.europa.eu/research/press/2016/pdf/opendata-infographic\\_072016.pdf](https://ec.europa.eu/research/press/2016/pdf/opendata-infographic_072016.pdf)
- 404 51. Open Access in Horizon 2020 - EC funded projects [Internet]. [cited 2017 Aug 23]. Available from:  
405 <https://www.openaire.eu/edocman?id=749&task=document.viewdoc>
- 406 52. Foster ED, Deardorff A. Open Science Framework (OSF). *J Med Libr Assoc JMLA.* 2017;105:203–6.

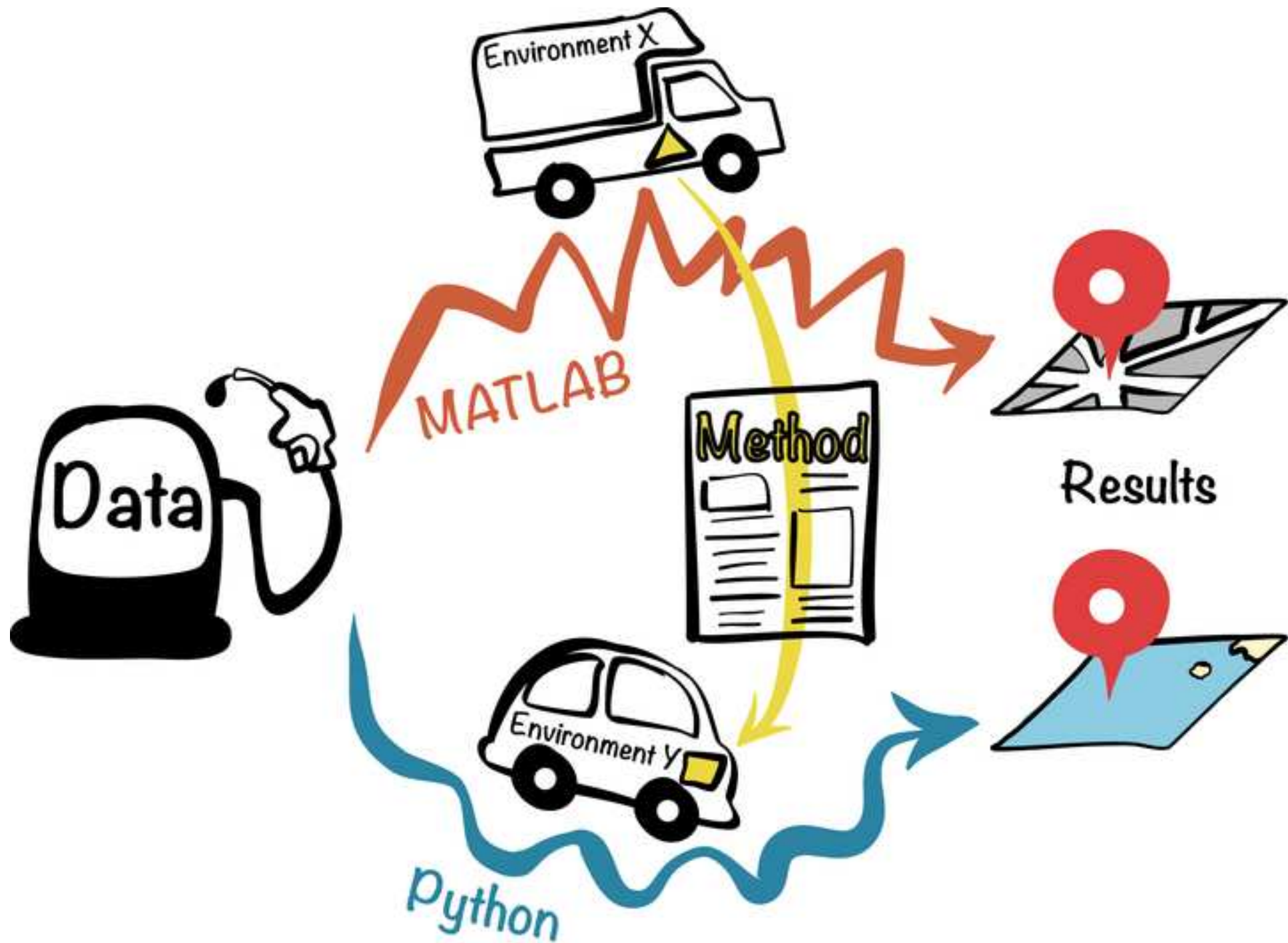
## Figure legends

**Figure 1: Hidden reproducibility issues like underwater iceberg.** Scientific journals readers have the impression that they can almost see the full work of method. But in reality, articles do not take into account adjustment and configuration for significant replication in most cases. Therefore, there is a significant gap between apparent executable work (i.e. above water portion of iceberg) and necessary effort in practice (i.e. full iceberg).

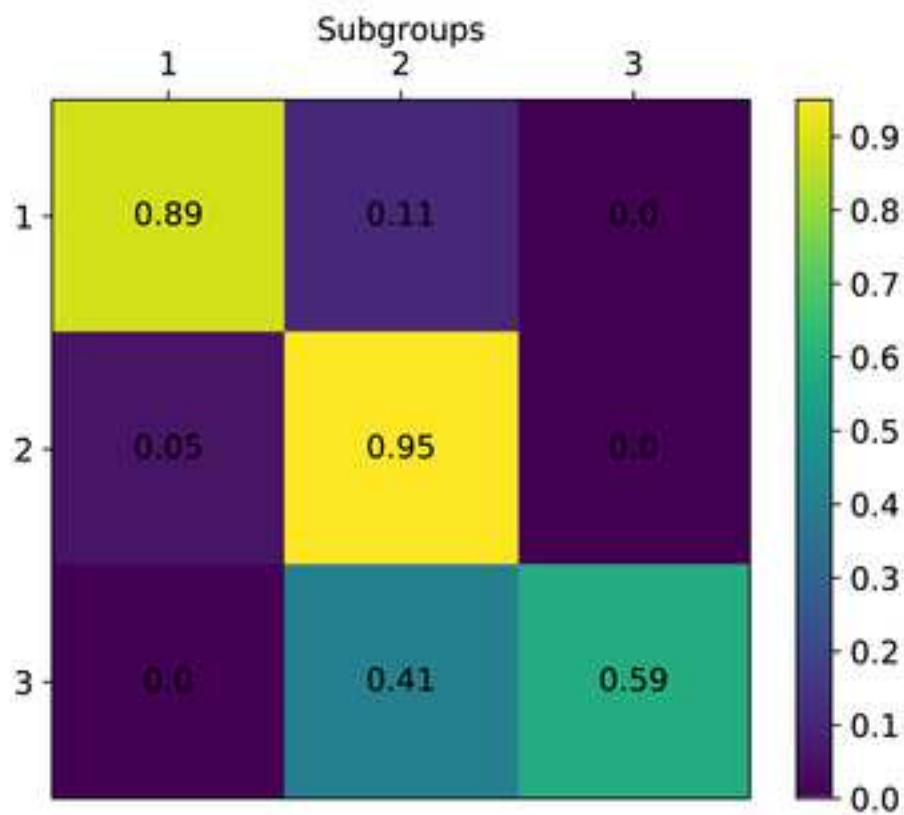
**Figure 2: Analogy between robustness issues and road transport.** The aim is to achieve same output (i.e. to reach the same location) using published methods (i.e. engine). Despite the same input data (i.e. gasoline), we obtained different results due to different programming languages —e.g. MATLAB and Python— (i.e. different roadways) and environments (i.e. different vehicles).

419 **Figure 3: Normalized confusion matrices between original and replicated results.** Before (a) and after (b) applying  
420 1 appropriate value of graph regularization factor on NBS method. Each row or column corresponds to a subgroup of patients (here  
421 2 three subgroups). The diagonal elements show the frequency of correct classifications for each subgroup: a high value indicates a  
422 3 correct prediction.  
423 4  
424 5  
425 6

426 7  
427 8  
428 9 **Figure 4: Working principles of testing ecosystem with private data.** Figure 4a shows a classical case: (a.1) Authors take  
429 10 private data (e.g. blue data) then publish their method and corresponding results; (a.2) Users having their own data (e.g. orange  
430 11 data) find a relevant paper but will be lost in the labyrinth of reproducibility. Figure 4b shows testing ecosystem with standard  
431 12 consensus dataset: (b.1) If authors work with their own data, they must identify corresponding standard data tag(s) (e.g. blue data);  
432 13 (b.2) Authors initiate to develop their method with corresponding standard data and reproducibility profile will be progressively  
433 14 built. Bar length on iceberg corresponds to progression of replication test; (b.3) Users can test proposed method with other  
434 15 standard data (e.g. orange and green data) and thus participate to enhancement of the reproducibility profile; (b.4) Thanks to the  
435 16 collective work on testing, the method could be optimized and authors can upgrade their initial paper (versioning).  
436 17  
437 18  
438 19  
439 20  
440 21  
441 22  
442 23  
443 24  
444 25  
445 26  
446 27  
447 28  
448 29  
449 30  
450 31  
451 32  
452 33  
453 34  
454 35  
455 36  
456 37  
457 38  
458 39  
459 40  
460 41  
461 42  
462 43  
463 44  
464 45  
465 46  
466 47  
467 48  
468 49  
469 50  
470 51  
471 52  
472 53  
473 54  
474 55  
475 56  
476 57  
477 58  
478 59  
479 60  
480 61  
481 62  
482 63  
483 64  
484 65



**a** Confusion matrix with reported tuning parameter value ( $\lambda = 1$ )



**b** Confusion matrix with actually used tuning parameter value ( $\lambda = 1800$ )

