



Toward a FAIR Reproducible Research

Christophe Bontemps, Valérie Orozco

► To cite this version:

| Christophe Bontemps, Valérie Orozco. Toward a FAIR Reproducible Research. 2021. hal-03018653

HAL Id: hal-03018653

<https://hal.inrae.fr/hal-03018653>

Preprint submitted on 29 Jan 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

November 2020

“Toward a FAIR Reproducible Research”

Christophe Bontemps and Valérie Orozco

Toward a FAIR Reproducible Research

Christophe Bontemps^{1*} | Valérie Orozco^{1*}

¹Toulouse School of Economics, INRAE,
University of Toulouse Capitole, Toulouse,
France

Correspondence

Valérie Orozco
Email: valerie.orozco@inrae.fr

Present address

*Toulouse School of Economics, INRAE,
University of Toulouse Capitole, Toulouse,
France

Funding information

Two major movements are actively at work to change the way research is done, shared and reproduced. The first is the reproducible research (RR) approach, which has never been easier to implement given the current availability of tools and DIY manuals. The second is the FAIR (Findable, Accessible, Interoperable, and Reusable) approach, which aims to support the availability and sharing of research materials. We show here that despite the efforts made by researchers to improve the reproducibility of their research, the initial goals of RR remain mostly unmet. There is great demand, both within the scientific community and from the general public, for greater transparency and for trusted published results. As a scientific community, we need to reorganize the diffusion of all materials used in a study and to rethink the publication process. Researchers and journal reviewers should be able to easily use research materials for reproducibility, replicability or reusability purposes or for exploration of new research paths. Here we present how the research process, from data collection to paper publication, could be reorganized and introduce some already available tools and initiatives. We show that even in cases in which data are confidential, journals and institutions can organize and promote "FAIR-like RR" solutions where not only the published paper but also all related materials can be used by any researcher..

Keywords: Reproducible Research; FAIR; Trusted Science; Confidential Data.

1 | THE NEED FOR REPRODUCIBLE RESEARCH

During the last decade, a great number of papers have been published on the problem of irreproducibility of research (Nature, 2013; Munafò et al., 2017) and on the crisis in science due to errors (Reinhart and Rogoff, 2010) or fraud (Ioannidis, 2005), leading to a lack of trust in published results. One response of this credibility crisis has been a renewal of interest in the "reproducible research" (RR) approach, as defined initially by geologist John Claerbout as the possibility of the "*replication [of a paper] by other scientists*" (Claerbout, 1990).

However, despite an apparent consensus on the general problem, the publication of papers exhorting the scientific community to publish reproducible results, and the dissemination of tools, good practices and courses (Stodden et al., 2013; Gentzkow and Shapiro, 2013; Sandve et al., 2013; Orozco et al., Forthcoming 2021), we still observe considerable weaknesses in both researcher and journal practices, leading to the scarce dissemination of raw scientific materials.¹ Many journals do not have a replication or data and code availability policy. Among those that do, there is a great diversity in practices, with many implementing a simple supplementary materials section on their website. Thus, a clear organization and precise guidelines on how to achieve the initial goals of reproducibility in science are still lacking.

Since science is based on the accumulation of verified results, its growth is limited if researchers cannot easily reproduce one another's results and verify them. To achieve verification, several conditions have to be fulfilled: First, the research has to be done in such a manner that it can easily be reproduced. Second, the materials used to produce the results have to be available to others. Finally, somebody, i.e., a referee for a journal or another researcher, has to reproduce and validate the published results using the materials available. These conditions may seem very strong, but we argue that they are necessary to prove the validity of any research. In mathematics, we would all agree that a paper with a new theorem should, first, be accompanied with its complete proof in a readable format. Second, all the elements of the proof (lemmas, other theorems) must be provided in the paper or easily available, and finally, someone must have checked—reproduced—the theorem using these elements. Our message here is that even in empirical work where we use data and code to produce a result, we have to prove our findings. We follow the idea of LeVeque (2009) that "*constructing a computer program isn't so different from constructing a formal proof*" and claim that reproducing a result issued from a computer program should not be different from reproducing a formal proof.

¹We consider here that the research process starts once the data are collected and in possession of the researcher. We do not address here the issue of reproducibility for data collection in experimental economics or field experiments (Bowers et al., 2017).

This paper is devoted to two practical problems that have received little attention in economics and statistics so far: How, in practice, can we ensure that the results published in a paper have been reproduced and verified? How are all the materials used to produce the results of that paper shared with the community? These are very complex questions that can be even more complex when materials are confidential or when certain technical expertise (on machine learning, big data, etc.) is required for storing and using them. We question here the overall organization of research leading to the publication of a paper. The workflow leading to publication is long, and many obstacles may limit the publication of anything else but the paper. The researcher has a great role in preparing not only the document to be submitted but also all the materials leading to the results. Thus, the focus has been on researchers' practices for a long time. However, other actors, in particular journals as well as data providers and research institutions, also have a role to play in facilitating the diffusion of science, not only through publications but also through the availability of all the materials necessary for science to be truly reproducible and trusted.

The paper is organized as follows. In the next section, we identify and illustrate the current problems that limit the reproducibility of research. We also clarify and illustrate the definitions of the terms and notions used in the RR and FAIR approaches and survey some important initiatives that have been proposed. In section 3, we introduce some recent initiatives and propose new schemes involving researchers, journals and the FAIR and reproducible research community. Some issues arising from the nature of the raw materials (big data) or of the analysis (machine learning) are also addressed in this section. We illustrate the problem of sharing research materials when they are confidential in section 4 and show that some solutions to this problem have already been implemented by some journals. Section 5 lists the incentives and impediments related to the proposed approach and concludes the paper.

2 | REPRODUCIBLE RESEARCH IN PRACTICE

The notion of reproducibility has been discussed lengthily in the literature, sometimes with conflicting terminologies (Barba, 2018; Benureau and Rougier, 2018; Bollen et al., 2015). We follow Barba (2018) who summarize it by the equation "same data + same method = same results".² For Gentleman and Temple Lang (2007), this idea means that the "data+code" *compendium* - or reproducibility package - used in a paper is made available to the readers so that they can first verify (reproduce) the results and second conduct alternative analyses of the work. The notion of

²We will not discuss here the question of the precise meaning of "same results".

reproducibility is thus related to the similar notion of verification or scientific proof.

In complement to RR, the Open Science movement emerged 10 years ago and aims at sharing data used in research as a patrimonial and cumulative goal. This movement has seen the involvement of many stakeholders, communities, and institutions (e.g., the Open Government Partnership comprising 70 countries, the Center for Open Science, the Research Data Alliance).³ These initiatives have focused primarily on the big questions behind open science implementation (e.g., repository certification, reusability, data citation, data management plans). Many tools and platforms emerged from these discussions: e.g., the European Open Science Cloud, the GOFAIR initiative and the Transparency and Openness Promotion (TOP) guidelines (Nosek and coauthors, 2015). The FAIR principles (findable, accessible, interoperable, and reusable) proposed by Wilkinson M. and et al. (2016) recently emerged as a global template for sharing data. FAIR does not mean open but, in brief, requires some accessibility to findable elements (most often datasets or at least metadata). The principles call for materials to be shared in a format that others can use (interoperable) and reuse.

These two movements are very active and influence the way institutions, research centers and national statistical offices (NSOs) construct their infrastructures and datacenters. To illustrate how these movements may affect the research publication landscape, we propose to reduce all the materials needed to produce a paper to only three key elements: the data, the code and the workflow, even if elements such as the documentation and the computing environment are also of great importance for some papers. We use the pictograms presented in Figure 1 throughout the paper. In view of the RR approach as defined by Claerbout (1990), all these elements, not only the research paper, should be shared with the scientific community.



FIGURE 1 Pictograms of inputs and outputs of empirical research (simplified)

We also illustrate the process leading to the publication of a paper by focusing on 3 major actors: the researcher(s), the journal that handles the publication process and the scientific community that should benefit from a new publication. We represent these actors with the pictograms in Figure 2. Later in the paper, we will see that other actors, such as research institutions, data providers and funding agencies, may play an important role in the publication process

³At the European level, one should mention OpenAIRE and in France the "Plan national pour la science ouverte" (<https://www.ouvrirelascience.fr/>).

and its outcome.



FIGURE 2 Main actors in the research and publishing process (simplified)

Since the publication of the seminal paper by [Claerbout \(1990\)](#), several authors have pointed out weaknesses in researchers' day-to-day practices and have proposed tools, solutions and advice. Many of these proposals underline existing ways of improving our practices toward RR, focusing on technical solutions. Applied statisticians and econometricians can now enjoy tools developed in R ([Leisch, 2002](#); [Meredith and Racine, 2009](#); [Gandrud, 2015](#); [Xie, 2015](#); [Xie et al., 2018](#)), Python ([Bilina and Lawford, 2012](#)), Stata ([Gentzkow and Shapiro, 2013](#); [Jann, 2016, 2017](#)), SAS ([Lenth and Højsgaard, 2007](#)) and MATLAB ([Hunt et al., 2014](#); [Woodward, 2012](#)).

Others ([Baiocchi, 2007](#); [Orozco et al., Forthcoming 2021](#)) have identified possible organizational improvements at the researcher level. These authors proposed principles to link a paper to its raw components, data and code through a clear workflow, following the original idea of *reproducible research documents* proposed by [Knuth \(1984, 1992\)](#) and his *literate programming* approach. Other principles include a clear organization of work and files, greater attention to versioning, good documentation of the research workflow, good writing practices for code using layouts and naming conventions. Automating the whole workflow is also recommended and encouraged.

We have also seen the emergence of companion websites and "executable" papers allowing online code editing and execution ([Hurlin et al., 2014](#); [Gorp and Mazanek, 2011](#)). Platforms (e.g., Code Ocean, Exec&Share, SHARE) have been created to allow code to be run online using materials stored by other researchers. With this technology, researchers from around the world are able to rerun the exact same code as the author, change parameter values to see the impact on the results, or even replicate the code with another dataset.

However, despite the efforts observed and all the tools and methods mentioned above, implementing RR has often been a challenge in practice. In a recent survey, [Chang and Li \(2017\)](#) attempted to update the seminal work of [Dewald et al. \(1988\)](#) and successfully replicated only 33% of 67 papers published in economics journals. Other examples in other disciplines exhibit similar features ([Miyakawa, 2020](#); [Asendorpf et al., 2013](#)). This current situation, whereby journals do not ask for or check the raw materials used to produce a research paper, is represented in Figure

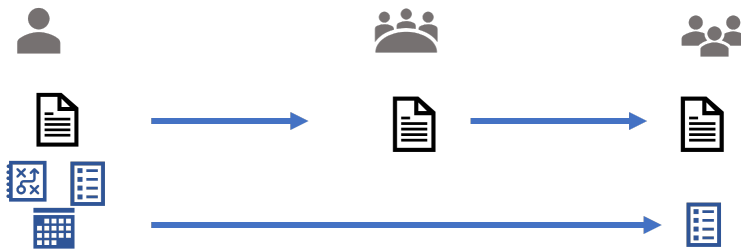


FIGURE 3 The researcher (left) uses code and data following a (written) workflow and submits a (reproducible) paper to a journal (middle) that neither asks for nor checks the raw materials. The researcher then decides to share (or not) some materials (e.g., the code) with the community (right) online or "upon request".

3.

There are many reasons for the overall "irreproducibility" (Nature, 2013) of research. It is true that crafting reproducible papers may require more time and effort than that needed for papers with code that will be used only once. However, when researchers, especially young ones, invest in RR practices, habits and tools, they generally become more efficient in their day-to-day practices (Stodden et al., 2013). They are able to reuse their own materials, reproduce their own results, answer referees' questions more quickly and test various specifications of their models with little effect on the time spent reporting new results in a paper (McCullough, 2009). There is also some evidence that papers that share their research materials are more likely to be cited (Christensen et al., 2019). The case of the software industry has shown that open-source creations are likely to increase the reputation and hence the value of their authors in the labor market (Lerner and Tirole, 2002).

The way researchers apply RR methods and share their research materials may be another important reason for the prevalence of irreproducibility. When they share their materials, researchers do it in an unstructured manner, on their websites or on GitHub, and many do not link the materials to each other or to the associated publication (Baker, 2016). Even if available, the data may not be easily accessible. In their book, Christensen et al. (2019) remark that when shared, submitted data were also frequently an "unlabeled and undocumented mess". Many authors cite only general source information but do not give enough precise information about the data used (Dewald et al., 1986). Many researchers focus on code and share their files without a clear document explaining the order of execution of each piece of code or the overall workflow. The *readme* files that serve as a workflow description are rarely structured to explain how to use the programs and often only describe the materials used and are far from a manual for reproduction. Some materials, mostly code, may be shared after acceptance in the journal's supplementary material section (if any), on the researcher's website or on GitHub. Most of the time, researchers decide on their own initiative to share (or

not) their materials "upon request".

Data availability and data accessibility are thus ubiquitous problems.⁴ These problems are not the sole responsibility of researchers. Many scientists currently do share code and data, but this takes time and effort. [Wacharamanotham et al. \(2019\)](#) recently surveyed more than 1000 authors of papers in computer science, asking if they were sharing their materials and if not why. The results underline some "misunderstandings about the purpose of sharing and reliable hosting". We have to acknowledge that even if some researchers do share their research materials, others simply do not want to.

In other disciplines, data collection represents a substantial effort. Researchers may not want to lose their investment and do not want to share this rent with others after a first publication.

The lack of incentives to share and the lack of sharing solutions are thus two barriers to reproducibility. They exist either because no clear data and code-sharing policies are defined by journals or because the technical solutions proposed are not good ones.

In [Figures 3 and 4](#), we illustrate different unsatisfactory yet frequently observed situations in which the research community has no access to the original raw materials even when the paper was originally designed by the researcher to be reproducible. In [Figure 4](#), only the code is finally shared with the research community, even if the journal has access to all the materials. The materials may be "available upon request", allowing the author the freedom to decide what to share. It is not guaranteed in either case that the paper was crafted to be reproducible nor that the referees were able to reproduce the results when reviewing and then accepting the paper for publication. Several examples of errors, such as the famous one in [Reinhart and Rogoff \(2010\)](#), cast doubt on the role of referees and their ability or willingness to ask for and use submitted materials to reproduce results during the referee process.

This issue is not limited to applied research: even theoretical mathematics papers display the same types of problems, with some analyses failing to be reproduced during the referee process or with authors citing unpublished papers in their work. Recently, a mathematician discovered a basic error (an inversion of an inequality) in the proof of a published paper ([Shchur, 2013](#)) and published the correct version where all the proofs had been checked using an open-source computer proof checker, opening the way to a new type of computer-assisted referee process ([Gouëzel and Shchur, 2019](#)).

For journals, reproducing a submitted paper requires many combined conditions that are still rarely met. First,

⁴Other issues that we do not address directly here include the digital preservation of research data ([Gutmann et al., 2009](#); [Akers and Doty, 2013](#)) or the preservation of software ([Di Cosmo and Zacchioli, 2017](#)).

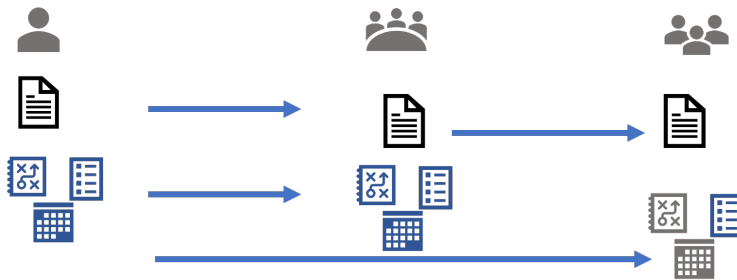


FIGURE 4 The researcher (left) sends all materials to the journal (middle) but in the end shares some material (e.g. code) directly online (e.g. GitHub). There is no indication that the shared materials are the ones used in the paper, that they have been checked, or that they are working.

the journal's policy must require that all materials be sent before or during the referee process. This implies that the submitted paper must have been done in a reproducible manner by the researchers in the first place. Second, the referees must have the skills, the willingness and the incentives to check the empirical proofs. Finally, the paper has to be reproduced, which implies some technical requirements, time and resources. Even if a journal is willing to reproduce a paper, it may not be able to do so and such conditions are still rarely fulfilled. Since replicating others' work can be quite technical and time consuming, this burden could also be outsourced to specialized third-parties agencies (see section 4). The journal then has several options to organize the way that materials are shared or not with the community. In practice, many journals do not impose or specify a way for materials to be shared. The most popular way of sharing consists of a "supplementary materials" section on journal websites. In a poll, [Science \(2011\)](#) reports that less than half of data requests from researchers to authors lead to positive outcomes.

In fact, the situations depicted in [Figure 4](#) in which the researchers are the sole party responsible for the quality

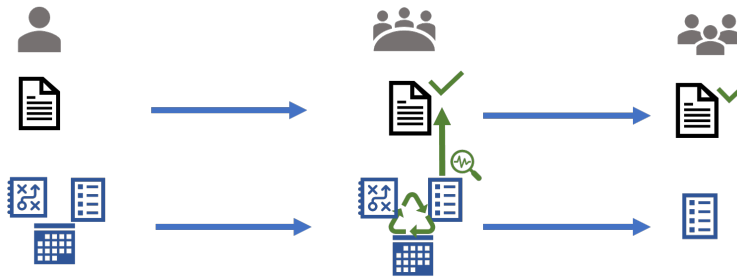


FIGURE 5 The researcher (left) sends all materials to the journal (middle). The journal shares only the code online for the community (right) and signals that the code has been carefully checked, either internally or by a trusted third-party agency.

and reproducibility of the files shared and no one checks what is stored are still the most commonly observed.

In Figures 5 to 7, we illustrate different types of observed situations where the community has access to a verified paper.⁵ These solutions assume that journals have a clear data and code availability policy that is enforced. This requires that the raw materials provided by the researchers closely follow the RR principles. The whole workflow, including the curation of the original (raw) dataset collected either by the researcher or through a data provider, should be written in a readable form for humans and computers, following the literate programming ideas proposed by Knuth (1992).

If these mandatory conditions are fulfilled by researchers, journals may implement two different strategies. They could try to reproduce the results and then signal or certify that the results are correct (Figures 5 and 6). This is what is expected from any scientific journal willing to maintain a reputation as a trusted publication. Alternatively, if a journal lacks sufficient human or financial resources to achieve this task, it should at least organize the diffusion of the materials and leave the verification process to the community, as in Figure 7.⁶ The best solution is illustrated in Figure 6, where the journal or a trusted third party reproduces and checks the paper's results and then organizes the sharing of the materials they have used for the benefit of the community. Such transparent organization may seem difficult to establish in practice, but it is in fact already implemented by some journals.

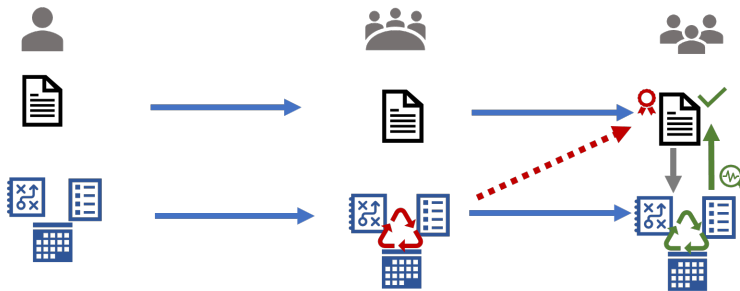


FIGURE 6 The researcher (left) sends all materials to the journal (middle). The journal or a trusted third party certifies that the materials used in the paper actually reproduce the results. The journal also shares all the materials (code-data-workflow) online for the community (right), which can also check, reproduce and reuse the materials.

⁵In these figures, for clarity reasons, we do not illustrate the fact that researchers may share their materials themselves.

⁶In 2003, Hashem Pesaran announced the creation of a new section of the Journal of Applied Econometrics dedicated to the replication of previously published empirical papers (Pesaran, 2003). Since then, some journals have followed this idea leading to an increase in the number of replication papers in economics (Mueller-Langer et al., 2019).

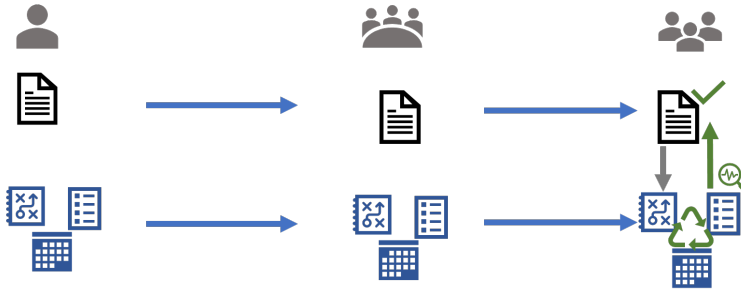


FIGURE 7 The researcher (left) sends all materials to the journal (middle). The journal shares all the materials (code-data-workflow) online for the community (right), which may check the results.

3 | IMPLEMENTING FAIR AND RR PRINCIPLES IN PRACTICE

In an empirical study examining 346 journals in economics and business studies, [Vlaeminck and Herrmann \(2015\)](#) showed that only 20% of the journals have a data policy. In Table 1, we compile the data and code availability policies for statistics and economics journals as published on their websites in 2020. We can observe a great heterogeneity of practices. From this nonexhaustive list, we confirm that many journals still do not publish any policy at all. Even if the publication of erroneous results has probably helped journals improve their referee process, the vast majority only "encourage" authors to share their materials, sometimes without any guidance on how or where to put the materials. Until very recently, very few (e.g., JSS, AJPS) checked that the code runs and reproduces the key results ([Christian et al., 2018](#)). Moreover, [Duvendack et al. \(2017\)](#) showed that in only 28 economic journals (out of 333) do a majority of the empirical papers supply their data and code. This suggests that even when data and code policies are written, their enforcement is lax.

Nevertheless, there are some very good examples that should be inspiring. The JASA, for example, has an "associate editor for reproducibility", responsible for the technical review of manuscripts before, during or after the usual review process ([Fuentes, 2016](#)). The Journal of Statistical Software (JSS) asks for a standalone replication script that must enable reproducibility. The materials are checked and then stored on the journal website, following the scheme of Figure 6. Since July 2019, all journals published by the American Economic Association require that authors share their data and code, which are systematically checked "within reasonable limits of time and computing resources". These verification's can be costly for journals, but we may expect that the costs should decrease over time with the

improvement of researchers' and reviewers' practices.⁷ This process can also be outsourced to specialists.⁸ An example is given by the American Journal of Political Science (AJPS), which contracted the Odum Institute for Research in Social Science to systematically check that research materials confirm the results of submissions (Crabtree, 2011).

To succeed in organizing the way materials are shared while preserving the link with the results included in the paper, journals could extend the FAIR principles, developed primarily for datasets, as proposed by Wilkinson M. and et al. (2016). In practice, sharing data and other materials together can be quite complicated. Individual solutions such as GitHub or even Dropbox may seem convenient for researchers but are not adapted or recommended as sharing solutions for published papers (Gandrud, 2015). There are risks that some elements in the code evolve (functions), that the data used are not explicitly identified (increments, versions) or that the links among all the research materials may be fuzzy or lost. Moreover, we believe that researchers should not organize the sharing themselves and that journals should align a strict mandatory policy with a clear organization and resources. In this regard, the Journal of Applied Econometrics (JAE) was a sort of pioneer, implementing its own data archive from the late 1980s. One satisfying technical solution is now offered by the *Dataverse Network*, developed at Harvard University. The network hosts collections of studies, embedding all materials for a paper in a single object, called a *dataverse* (King, 2007; Leeper, 2014). This solution is recommended by some journals (e.g., JASA, AJPS, QJE, PLOS, Nature). Other journals, following the recommendations of research institutions (NSF, ERC), use Figshare, Zenodo, Mendeley or ICPSR (see Table 1).⁹

Following the FAIR principles implies not only that the materials should be shared, or accessible, but also that they should be interoperable and findable. The interoperability principle can be interpreted here as the ability for any reader to have access to materials stored in a readable format.¹⁰ It should also be easy to find the materials without any ambiguity. We agree with the American Sociological Review that "*citing datasets used in published research is just as important as citing journal articles, books, and other sources that contributed to the research*". This means that the data and other materials should be identified in a consistent way. To precisely identify datasets and code, journals or institutions may request that a digital object identifier (DOI) be attached to each element described in the paper.

⁷In a recent analysis, Jacoby William G. (2017) analyzed the AJPS verification policy and reported an average of 8 person-hours per manuscript to replicate the analyses and curate the materials. The publication workflow, involving more rounds and resubmissions, is also much longer.

⁸There are already some useful resources and checklists that should facilitate the process for authors and replicators (see https://social-science-data-editors.github.io/guidance/Verification_guidance.html) or the TOP proposing varying levels of replication policies for journals (Nosek and coauthors, 2015).

⁹A complete list of solutions is detailed in *The Registry of Research Data Repositories* (<http://re3data.org>) a service of DataCite. In addition, CoreTrustSeal provides certification to repositories and lists the certified ones.

¹⁰For datasets, the FAIR interoperability principle suggests the use of open formats such as CSV files instead of proprietary formats (.xls). For code and other materials, open-source software should be preferred to avoid exclusive access (Vilhuber, 2019). The metadata should also follow relevant standards (Dublin core or DDI) and ontologies. References and links to other related data should also be provided (Jones and Grootveld, 2017).

DOIs are the backbones of the findable component of the FAIR principles. Even for data with constrained access (e.g., proprietary or confidential data), a DOI should provide enough elements, at least the metadata, to retrieve enough information describing the data (or other materials) used.¹¹ Several papers explain how to cite data repositories in a paper (Fenner et al., 2017).

In the context of big data or machine learning (ML) analysis, sharing the materials and the analysis are important issues that can be challenging (Crosas et al., 2015). The *Zelig* package, initially developed for R users, may offer a solution. *Zelig* provides a unified framework for statistical analysis allowing exploration, estimation and interpretation of any type of statistical model (Choirat et al., 2016). An important feature in the context of ML is that the *Zelig* package automatically creates a workflow embedding all the procedures and algorithms used in the analysis into a single object that may then be exported and shared. The analysis can then be used by others, even if they do not know R (King, 2007). In a recent paper, Crosas et al. (2015) proposed extending the architecture of *Zelig* to other tools and languages (Python) to provide a framework suitable for big data analysis. The NSF-funded Whole Tale (<https://wholetale.org/>) may also be a scalable solution enabling the creation, publication, and execution of "tales" or executable objects embedding data (potentially big data), code, and the complete software environment used to produce research findings (Chard et al., 2020).

Another difficulty may be due to the length of the publication process, which can be quite extensive. The reviewing process often requires additional tests or modifications. Thus, the code as well as the datasets and even the workflow may change with the evolution of the paper under review. Researchers following the RR approach and writing RR papers are already familiar with version control tools such as GitHub for their code. They should easily integrate their data and workflow in the same spirit. Both researchers and journals must integrate versioning in their organization. They should coordinate to be able to identify the exact version of each element used in the current version of the paper under review.

Another large step would be to question the access policies of publications, which are often paywalled, and to promote open-access publications such as PLOS (<https://plos.org/>) or arXiv (<https://arxiv.org/>) in our fields.

¹¹The DataCite project (Brase, 2009) is a popular resource to locate and precisely identify data through a unique DOI.

4 | CONFIDENTIAL DATA

Confidential or proprietary data are often cited as an obstacle to reproducibility, mainly because of data accessibility restrictions such as those imposed by the European GDPR.¹² In economics, Christensen and Miguel (2018) observed that there has been a small increase in empirical papers (using data) published coinciding with a significant increase in data exemptions. They found that nearly half of the papers using data are not reproducible because the data are not available (see also Vilhuber (2018)). These exemptions may be due to the use of confidential data or to other restrictions limiting data availability.

Restrictions may also come from data providers' preservation policies or national statistical offices (NSOs) data management conditions allowing only remote and strictly controlled access to data through secure virtual terminals.¹³ Journals and reviewers are then unable to access the data and cannot check the results, having access to the paper and the code only.¹⁴ According to Lagoze and Vilhuber (2017), 50% of confidential data used in papers are from NSOs (approximately 60% in 2010 for top economic journals), and each NSO has its own data-sharing restriction and regulation. To handle the complexity of security-level restrictions and to allow third parties or reviewers to access confidential materials, Sweeney et al. (2015) proposed a system of *datatag* repositories. Each *datatag* repository documents the way data and other sensitive materials can be shared, reducing the complexity of the situation to a small number of tags.

Some authors propose altering original raw data into "safe data", potentially accessible by anybody, using blurring or aggregation techniques to remove sensitive details such as individual information (Alter and Gonzalez, 2018). Other methods include adding random noise or swapping individual responses between otherwise similar respondents while maintaining the same likelihood distribution (Boker et al., 2015). In our view, these methods are limited to some specific cases and do not seem compatible with the principles of transparency used in the RR approach. Moreover, the transformation of raw data does not seem to solve the confidentiality issue, since the code used for these transformations must itself remain confidential, leading to a new sharing problem.

Nonetheless, solutions exist to preserve privacy while allowing other researchers to access and replicate results. The principle is based on a trusted third-party having a secure access to confidential data and on an interactive platform

¹²There are many sources of confidential and nonshareable data (Christensen and Miguel, 2018; Lagoze and Vilhuber, 2017).

¹³In France, the CASD (<https://www.casd.eu/>) is a single-access portal to many public data providers (INSEE, ministries, etc.). Researchers are not allowed to copy all the materials locally on their machine, and only some type of outputs can be extracted.

¹⁴The code may also contain some confidential elements. In particular, the code used for the initial data curation may contain, e.g., brand or city names, addresses, etc.

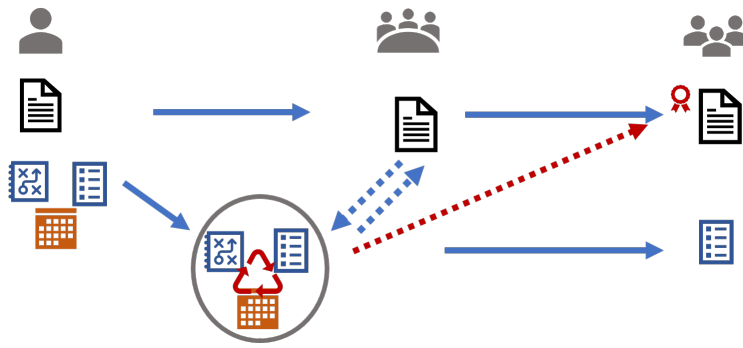


FIGURE 8 The researcher (left) shares all materials (code-data-workflow) privately with the trusted third-party. Access may be contracted with the journal (middle). A certification of reproducibility is sent to the journal together with the paper. The paper is published with the reproducibility certification, but only nonconfidential elements, such as code, are available to the community (right).

for queries and answers (Dwork et al., 2009). The curator model (Crosas et al., 2015), depicted in Figure 8, could be implemented using secured-sharing platforms such as *dataverses* or *datatags*.¹⁵ Referees could have access to materials securely and reproduce and check the validity of the results, even with confidential data. A reproducibility certification could be emitted during or even before the submission process to attest to the full reproducibility of the results.

The recent and promising Certification Agency for Scientific Code and Data (CASCAD), supported by the French National Science Foundation (CNRS), proposes the solution of prepublication reproducibility certification (Pérignon et al., 2019). Authors can ask CASCAD to review their materials, even when confidential data are used, and to certify the reproducibility of their paper's result (Figure 8). The replication process can be performed before publication, which facilitates the journal review process. Alternatively, journals may outsource this task to the certification agency during the referee process. Ex ante contracts between the trusted third party, the data provider and the journal may facilitate the process. The code used and the certification report may then be hosted on an open-access repository, such as Zenodo, and be accessible by the research community.

It is possible to attest that a paper using confidential data has been reproduced so that no doubt remains over the verifiability of the results. However, only a few people, namely, the third-party agency or some reviewers, had access to the research materials, as only the certified code may be finally shared with the community.

The reasons to grant researchers (the community) access to the research materials used to prepare papers dealing

¹⁵ Some data providers, in particular national statistical offices (NSOs), already perform RR on their confidential data, controlling output files and code, to check if they satisfy confidentiality restrictions before publication (Lagoze and Vilhuber, 2017).

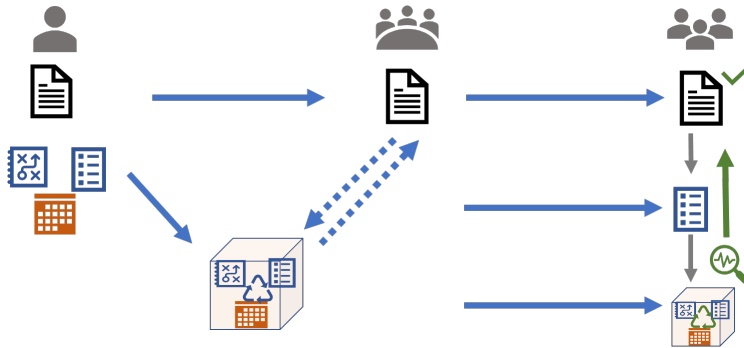


FIGURE 9 The researcher (left) shares all materials (code-data-workflow), preferably as a *package* or *tale* that is securely accessible by the journal's reviewers (middle) within a secured platform. The paper is published with access to a secured platform ("data enclave", such as the ICPSR). The community (right) has access to the secured query platform and can use the materials (e.g., code) without having access to the confidential elements.

with confidential data are numerous. As already stated, independent checks of results are important for science to be trusted, even if those results are based on confidential data. The recent Covid-19 pandemic and the use of individual databases have taught us that *"careful data-management practices should govern both data collection and data processing"* (Ienca and Vayena, 2020). If papers using confidential data are not trustable, the consequences could be dramatic for policymakers and for researchers. We follow Ienca and Vayena (2020) and believe that *"transparent public communication about (confidential) data processing for the common good should be pursued"*. Second, reproducing another researcher's paper can lead to new ideas that can be tested using shared materials. Finally, using each other's material is part of science, and sharing should not be an option.

However, working with confidential data does not necessarily mean having access to the dataset. Technical solutions, such as the one illustrated in Figure 9 and based on the idea of "data enclaves", exist. The Inter-university Consortium for Political and Social Research (ICPSR) has developed virtual and physical "data enclaves", allowing on-line data analysis with strict restrictions on queries, data access and downloads (Dunn and Austin, 1998). At the time of its creation long ago, physical access to ICPSR resources and all the outputs created were humanly controlled. Currently, an online data analysis tool (survey documentation and analysis (SDA) software, created by the University of California, Berkeley) is used to evaluate output for disclosure risk prior to displaying for the end user. This is again a curator model accepting queries from any researcher and providing only approved analysis output, including the elements (estimations, tables, figures) published in the paper.

5 | CONCLUSION

To many, the path toward findable, accessible, interoperable, and reproducible research may seem paved with obstacles. We argue here that this is the path to the future, considering the great challenges that we, as a scientific community, should overcome. Erroneous publications, unavailable research materials, and long and sometimes archaic publication processes have generated a research crisis within the research community and within society itself. Some publications ([Science, 2011](#)) advocate for a great change in the individual and collective practices of scientists, journals, funders, institutions, and societies, acknowledging *Claerbout's principles* that "*an article (· · ·) in a scientific publication is not the scholarship itself, (· · ·) the actual scholarship is the complete set of instructions which generated the figures*" ([de Leeuw, 2001](#)).

We review here some existing publication workflows and identify their weaknesses, observing that many lead to the unavailability of research materials for the research community, even when the researcher has made an effort to prepare the paper to be reproducible. We highlight some promising organizations dedicated to a better transmission of research materials to allow the scientific community to reproduce and verify published results. **XXXXX here something on paywalled journals and open access movement.** We argue that our community and the public at large will greatly benefit from a change toward greater transparency and better-organized research.

This change will only occur if all research actors agree to adhere to FAIR and reproducible research principles. Most of these principles can be gradually implemented as a growing process toward more reproducible practices.

Researchers should initiate these changes. They now have all the resources needed to improve their individual practices to create more reproducible papers by embedding code and data in a documented and written workflow understandable by others, including researchers' own "future selves" ([Gentzkow and Shapiro, 2013](#); [Orozco et al., Forthcoming 2021](#)). They are probably in the process of changing their habits already under the pressure of certain journals and research institutions such as the ERC and NSF.

Researchers are also reviewers, some of them even journal editors, and thus can promote many valuable actions. First, journals should reorganize the review process and the way the results of submitted papers are checked, including when some materials are confidential. This is probably one of the most challenging issues, requiring new skills for reviewers, additional resources and a clear internal setup for sharing the submitted paper's materials. Certified trusted third parties (e.g., CASCAD, Odum Institute, CISER R-Squared) already exist, if that process has to be outsourced. Having an "associate editor for reproducibility", as the JASA does, could also be a good idea to solve many practical

questions such as when to check the validity of the submitted materials (before, during, or after the classical review) or to organize the verification and the relations with the author. Second, journals should have a clear data and code availability policy with a proper check of materials by reviewers. Journals should also organize the way research materials are shared.¹⁶ Contracting with public repositories such as Zenodo, Mendeley, or Figshare would reduce the constellation of individual and self-maintained repositories and the fragmentation of arbitrarily different, incompatible standards (Sansone et al., 2019). These platforms also guarantee perennial access to the core materials of science. Finally, if implemented, this process will provide strong incentives for researchers to only produce RR papers. Some journals (e.g., AJPS, Biostatistics) propose badges tagging reproducible research. Such a practice seems to increase the proportion of papers using open practices and to improve the preservation of research materials (Kidwell et al., 2016; Rowhani-Farid and Barnett, 2018).

Data providers, whether they are private or public, could also facilitate the changes that we call for. Quite often, they are aware of these problems and have implemented some processes for their internal publications that could be inspiring (Lagoze and Vilhuber, 2017). In the near future, they may pay increasing attention to the use of the data they provide, checking published results either to publicize their activity or to criticize a misuse.¹⁷ In the case of confidential data, providers and, in particular, NSOs may also find some interest in promoting and organizing the way their data are findable and accessible (Pérignon et al., 2019). Thus, data providers should encourage researchers, institutions and journals to produce more reproducible and reproduced papers. It is therefore likely that partnerships among journals, data providers and private or public third parties will increase in the future.

Research institutions have already started to impose some conditions on funded projects or grants by requiring researchers to follow a strict RR approach, by promoting the dissemination of the FAIR and RR approaches, and by financing public infrastructures hosting FAIR research materials repositories.¹⁸

If "*science is organized knowledge*" (Spencer, 1854), then we should all work for better organization for better science. We believe that the FAIR and reproducible research movements are there to jointly provide organized resources, tools and practices. Changing the publication workflow and our habits may be a long and probably costly journey. Not changing could be even more costly.

¹⁶Alter and Gonzalez (2018) suggested that to "protect" researchers who want to use their data first (before sharing), journals can propose an "embargo".

¹⁷A recent lawsuit involving the popular training program CrossFit showed that a paper by Smith et al. (2013) erroneously showed an increased risk for injuries for its users. Although the paper was retracted later, the impacts on the researcher's career were severe (for details, see <https://retractionwatch.com/>).

¹⁸The European Research Council (ERC) recommends "*to all its funded researchers that they follow best practice by retaining files of all the research data they have used during the course of their work and that they be prepared to share this data with other researchers*" (see https://erc.europa.eu/sites/default/files/document/file/ERC_Open_Access_Guidelines-revised_2013.pdf).

references

- Akers, K. G. and Doty, J. (2013) Disciplinary differences in faculty research data management practices and perspectives. *International Journal of Digital Curation*, **8**, 5–26.
- Alter, G. and Gonzalez, R. (2018) Responsible practices for data sharing. *American Psychologist*, **73**, 146–156.
- Asendorpf, J. B., Conner, M., De Fruyt, F., De Houwer, J., Denissen, J. J., Fiedler, K., Fiedler, S., Funder, D. C., Kliegl, R., Nosek, B. A. et al. (2013) Recommendations for increasing replicability in psychology. *European Journal of Personality*, **27**, 108–119.
- Baiocchi, G. (2007) Reproducible research in computational economics: guidelines, integrated approaches, and open source software. *Computational Economics*, **30**, 19–40. URL: <http://ideas.repec.org/a/kap/compec/v30y2007i1p19-40.html>.
- Baker, M. (2016) Why scientists must share their research code. *Nature News*.
- Barba, L. A. (2018) Terminologies for reproducible research. *arXiv preprint arXiv:1802.03311*.
- Benureau, F. C. Y. and Rougier, N. P. (2018) Re-run, repeat, reproduce, reuse, replicate: Transforming code into scientific contributions. *Frontiers in Neuroinformatics*, **11**, 69. URL: <https://www.frontiersin.org/article/10.3389/fninf.2017.00069>.
- Bilina, R. and Lawford, S. (2012) Python for Unified Research in Econometrics and Statistics. *Econometric Reviews*, **31**, 558–591. URL: <https://ideas.repec.org/a/taf/emetr/v31y2012i5p558-591.html>.
- Boker, S. M., Brick, T. R., Pritikin, J. N., Wang, Y., Oertzen, T. v., Brown, D., Lach, J., Estabrook, R., Hunter, M. D., Maes, H. H. et al. (2015) Maintained individual data distributed likelihood estimation (middle). *Multivariate behavioral research*, **50**, 706–720.
- Bollen, K., Cacioppo, J., Kaplan, R., Krosnick, J. and Olds, J. (2015) Social, behavioral, and economic sciences perspectives on robust and reliable science: Report of the subcommittee on replicability in science, advisory committee to the national science foundation directorate for social, behavioral, and economic sciences. *Advisory Committee to the National Science Foundation Directorate for Social, Behavioral, and Economic Sciences*.
- Bowers, J., Higgins, N., Karlan, D., Tulman, S. and Zinman, J. (2017) Challenges to replication and iteration in field experiments: Evidence from two direct mail shots. *American Economic Review*, **107**, 462–65.
- Brase, J. (2009) Datacite - a global registration agency for research data. In *2009 Fourth International Conference on Cooperation and Promotion of Information Resources in Science and Technology*, 257–261.
- Chang, A. C. and Li, P. (2017) A preanalysis plan to replicate sixty economics research papers that worked half of the time. *American Economic Review*, **107**, 60–64. URL: <http://www.aeaweb.org/articles?id=10.1257/aer.p20171034>.
- Chard, K., Gaffney, N., Hategan, M., Kowalik, K., Ascher, B. L., Mcphillips, T., Nabrzyski, J., Stodden, V., Taylor, I., Thelen, T. and Others (2020) Toward enabling reproducibility for data-intensive research using the whole tale platform. *Advances in Parallel Computing*, **36**, 766–778.
- Choirat, C., Honaker, J., Imai, K., King, G. and Lau, O. (2016) Zelig: everyone's statistical software. URL <http://zeligproject.org/Version>.
- Christensen, G., Freese, J. and Miguel, E. (2019) *Transparent and Reproducible Social Science Research: How to Do Open Science*. University of California press.
- Christensen, G. and Miguel, E. (2018) Transparency, reproducibility, and the credibility of economics research. *Journal of Economic Literature*, **56**, 920–80. URL: <http://www.aeaweb.org/articles?id=10.1257/jel.20171350>.

- Christian, T.-M., Lafferty-Hess, S., Jacoby, W. and Carsey, T. (2018) Operationalizing the replication standard: A case study of the data curation and verification workflow for scholarly journals. *International Journal of Digital Curation*, **13**, 114–124.
- Claerbout, J. (1990) Active documents and reproducible results. *SEP*, **67**, 139–144. URL: <http://sepwww.stanford.edu/data/media/public/docs/sep67/jon2.pdf>.
- Crabtree, J. D. (2011) Odum institute user study: Exploring the applicability of the dataverse network.
- Crosas, M., King, G., Honaker, J. and Sweeney, L. (2015) Automating open science for big data. *ANNALS of the American Academy of Political and Social Science*, **659**, 260–273. URL: <http://ann.sagepub.com.ezp-prod1.hul.harvard.edu/content/659/1/260.full.pdf+html>.
- Dewald, W. G., Thursby, J. and Anderson, R. (1986) Replication in empirical economics: The journal of money, credit and banking project. *American Economic Review*, **76**, 587–603. URL: <https://EconPapers.repec.org/RePEc:aea:aecrev:v:76:y:1986:i:4:p:587-603>.
- Dewald, W. G., Thursby, J. G. and Anderson, R. G. (1988) Replication in empirical economics: The journal of money, credit and banking project: Reply. *American Economic Review*, **78**, 1162–1163. URL: <http://www.jstor.org/stable/1807177>.
- Di Cosmo, R. and Zacchiroli, S. (2017) Software heritage: Why and how to preserve software source code. URL: <https://hal.archives-ouvertes.fr/hal-01590958/document>.
- Dunn, C. S. and Austin, E. W. (1998) Protecting confidentiality in archival data resources. *IASSIST Quarterly*, **22**, 16–16.
- Duvendack, M., Palmer-Jones, R. and Reed, W. R. (2017) What is meant by “replication” and why does it encounter resistance in economics? *American Economic Review*, **107**, 46–51.
- Dwork, C., Naor, M., Reingold, O., Rothblum, G. N. and Vadhan, S. (2009) On the complexity of differentially private data release: efficient algorithms and hardness results. In *Proceedings of the forty-first annual ACM symposium on Theory of computing*, 381–390. URL: <https://www.microsoft.com/en-us/research/wp-content/uploads/2009/05/dnrrv09.pdf>.
- Fenner, M., Crosas, M., Grethe, J., Kennedy, D., Hermjakob, H., Rocca-Serra, P., Durand, G., Berjon, R., Karcher, S., Martone, M. and Clark, T. (2017) A data citation roadmap for scholarly data repositories. *bioRxiv*. URL: <https://www.biorxiv.org/content/early/2017/10/09/097196>.
- Fuentes, M. (2016) Reproducible research in jasa. *AMSTAT news: the membership magazine of the American Statistical Association*, **17**. URL: <https://magazine.amstat.org/blog/2016/07/01/jasa-reproducible16/>.
- Gandrud, C. (2015) *Reproducible Research with R and RStudio Second Edition*. Chapman & Hall/CRC The R Series.
- Gentleman, R. and Temple Lang, D. (2007) Statistical analyses and reproducible research. *Journal of Computational and Graphical Statistics*, **16**, 1–23. URL: <http://dx.doi.org/10.1198/106186007X178663>.
- Gentzkow, M. and Shapiro, J. (2013) Nuts and bolts: Computing with large data. In *Summer Institute 2013 Econometric Methods for High-Dimensional Data*. Summer Institute 2013 Econometric Methods for High-Dimensional Data. URL: http://www.nber.org/econometrics_minicourse_2013/.
- Corp, P. V. and Mazanek, S. (2011) Share: a web portal for creating and sharing executable research papers. *Procedia Computer Science*, **4**, 589–597. URL: <http://www.sciencedirect.com/science/article/pii/S1877050911001207>.
- Gouëzel, S. and Shchur, V. (2019) A corrected quantitative version of the morse lemma. *Journal of Functional Analysis*, **277**, 1258–1268.
- Gutmann, M. P., Abrahamson, M., Adams, M. O., Altman, M., Arms, C., Bollen, K., Carlson, M., Crabtree, J., Donakowski, D., King, G. et al. (2009) From preserving the past to preserving the future: The data-pass project and the challenges of preserving digital social science data. *Library Trends*, **57**, 315–337.

- Hunt, B. R., Lipsman, R. L. and Rosenberg, J. M. (2014) *A Guide to MATLAB: For Beginners and Experienced Users Third Edition*. Cambridge University Press.
- Hurlin, C., Pérignon, C. and Stodden, V. (2014) Runmycode.org: a novel dissemination and collaboration platform for executing published computational results. *Open Science Framework*. URL: <https://osf.io/39eq2/>.
- Ienca, M. and Vayena, E. (2020) On the responsible use of digital data to tackle the covid-19 pandemic. *Nature Medicine*, 1–2.
- Ioannidis, J. P. (2005) Why most published research findings are false. *PLoS Med*, **2**, e124. URL: <http://statweb.stanford.edu/~tibs/sta306bfiles/ioan.pdf>.
- Jacoby William G., Sophia Lafferty-Hess, T.-M. C. (2017) Should journals be responsible for reproducibility? URL: <https://www.insidehighered.com/blogs/rethinking-research/should-journals-be-responsible-reproducibility>.
- Jann, B. (2016) Creating latex documents from within stata using texdoc. *Stata Journal*, **16**, 245–263. URL: <http://www.stata-journal.com/article.html?article=pr0062>.
- (2017) Creating html or markdown documents from within stata using webdoc. *Stata Journal*, **17**, 3–38. URL: <http://www.stata-journal.com/article.html?article=pr0065>.
- Jones, S. and Grootveld, M. (2017) How fair are your data? URL: <https://doi.org/10.5281/zenodo.1065991>.
- Kidwell, M. C., Lazarević, L. B., Baranski, E., Hardwicke, T. E., Piechowski, S., Falkenberg, L.-S., Kennett, C., Slowik, A., Sonnleitner, C., Hess-Holden, C. et al. (2016) Badges to acknowledge open practices: A simple, low-cost, effective method for increasing transparency. *PLoS biology*, **14**.
- King, G. (2007) An introduction to the dataverse network as an infrastructure for data sharing. *Sociological Methods & Research*, **36**, 173–199. URL: <https://doi.org/10.1177/0049124107306660>.
- Knuth, D. E. (1984) Literate programming. *The Computer Journal*, **27**, 97–111.
- (1992) *Literate Programming*. Center for the Study of Language and Information.
- Lagoze, C. and Vilhuber, L. (2017) O privacy, where art thou? making confidential data part of reproducible research. *CHANCE*, **30**, 68–72. URL: <https://doi.org/10.1080/09332480.2017.1383118>.
- Leeper, T. J. (2014) Archiving reproducible research with r and dataverse. *R Journal*, **6**. URL: <https://journal.r-project.org/archive/2014/RJ-2014-015/RJ-2014-015.pdf>.
- de Leeuw, J. (2001) Reproducible research. the bottom line. URL: <http://www.escholarship.org/uc/item/9050x4r4>.
- Leisch, F. (2002) Sweave: Dynamic generation of statistical reports using literate data analysis. *Compstat 2002 - Proceedings in Computational Statistics*, 575–580.
- Lenth, R. V. and Højsgaard, S. (2007) Sasweave: Literate programming using sas. *Journal of Statistical Software*, **19**, 1–20.
- Lerner, J. and Tirole, J. (2002) Some simple economics of open source. *The journal of industrial economics*, **50**, 197–234.
- LeVeque, R. J. (2009) Python tools for reproducible research on hyperbolic problems. In *Special issue on Reproducible Research*, 19–27. Computing in Science and Engineering (CiSE).
- McCullough, B. D. (2009) Open access economics journals and the market for reproducible economic research. *Economic Analysis and Policy*, **39**, 117–126. URL: <https://www.sciencedirect.com/science/article/pii/S0313592609500471>.
- Meredith, E. and Racine, J. S. (2009) Towards reproducible econometric research: The sweave framework. *Journal of Applied Econometrics*, **24**, 366–374. URL: <http://www.jstor.org/stable/40206278>.

- Miyakawa, T. (2020) No raw data, no science: another possible source of the reproducibility crisis. URL: <https://molecularbrain.biomedcentral.com/articles/10.1186/s13041-020-0552-2#citeas>.
- Mueller-Langer, F., Fecher, B., Harhoff, D. and Wagner, G. G. (2019) Replication studies in economics—how many and which papers are chosen for replication, and why? *Research Policy*, **48**, 62 – 83. URL: <http://www.sciencedirect.com/science/article/pii/S0048733318301847>.
- Munafò, M. R., Nosek, B. A., Bishop, D. V., Button, K. S., Chambers, C. D., Du Sert, N. P., Simonsohn, U., Wagenmakers, E.-J., Ware, J. J. and Ioannidis, J. P. (2017) A manifesto for reproducible science. *Nature human behaviour*, **1**, 1–9.
- Nature, E. (2013) Reducing our irreproducibility. *Nature*, **496**, 398. URL: "http://www.nature.com/polopoly_fs/1.12852!/menu/main/topColumns/topLeftColumn/pdf/496398a.pdf".
- Nosek, B. A. and coauthors (2015) Promoting an open research culture. *Science*, **348**, 1422–1425. URL: <https://science.sciencemag.org/content/348/6242/1422>.
- Orozco, V., Bontemps, C., Maigne, E., Piguët, V., Hofstetter, A., Lacroix, A. M., Levert, F. and Rousselle, J.-M. (Forthcoming 2021) How to make a pie: Reproducible research for empirical economics & econometrics. *Journal Of Economic Surveys*.
- Pérignon, C., Gadouche, K., Hurlin, C., Silberman, R. and Debonnel, E. (2019) Certify reproducibility with confidential data. *Science*, **365**, 127–128. URL: <https://science.sciencemag.org/content/365/6449/127>.
- Pesaran, H. (2003) Introducing a replication section. *Journal of Applied Econometrics*, **18**, 111. URL: <http://www.jstor.org/stable/30035191>.
- Reinhart, C. M. and Rogoff, K. S. (2010) Growth in a time of debt. *American Economic Review*, **100**, 573–78. URL: <http://www.aeaweb.org/articles.php?doi=10.1257/aer.100.2.573>.
- Rowhani-Farid, A. and Barnett, A. (2018) Badges for sharing data and code at biostatistics: an observational study [version 2; peer review: 2 approved]. *F1000Research*, **7**.
- Sandve, G. K., Nekrutenko, A., Taylor, J. and Hovig, E. (2013) Ten simple rules for reproducible computational research. *PLoS Comput Biol*, **9**, 1–4. URL: <http://www.ploscompbiol.org/article/info%3Adoi%2F10.1371%2Fjournal.pcbi.1003285>.
- Sansone, S.-A., McQuilton, P., Rocca-Serra, P., Gonzalez-Beltran, A., Izzo, M., Lister, A. L. and Thurston, M. (2019) Fairsharing as a community approach to standards, repositories and policies. *Nature biotechnology*, **37**, 358–367.
- Science, S. (2011) Challenges and opportunities. *Science*, **331**, 692–693. URL: <https://science.sciencemag.org/content/331/6018/692>.
- Shchur, V. (2013) A quantitative version of the morse lemma and quasi-isometries fixing the ideal boundary. *Journal of Functional Analysis*, **264**, 815 – 836. URL: <http://www.sciencedirect.com/science/article/pii/S002212361200434X>.
- Smith, M. M., Sommer, A. J., Starkoff, B. E., Devor, S. T. et al. (2013) Crossfit-based high-intensity power training improves maximal aerobic fitness and body composition. *J Strength Cond Res*, **27**, 3159–3172.
- Spencer, H. (1854) The art of education. URL: <https://books.google.fr/books?id=-qLQAAAAAAAJ>.
- Stodden, V., Bailey, D., Borwein, J., LeVeque, R., Rider, W. and Stein, W. (2013) Setting the default to reproducible: Reproducibility in computational and experimental mathematics. URL: http://stodden.net/icerm_report.pdf.
- Sweeney, L., Crosas, M. and Bar-Sinai, M. (2015) Sharing sensitive data with confidence: The datatags system. *Technology Science*. URL: <http://techscience.org/a/2015101601/>.
- Vilhuber, L. (2018) Reproducibility and replicability in economics.

- (2019) Report by the aea data editor. *AEA Papers and Proceedings*, **109**, 718–729. URL: <http://www.aeaweb.org/articles?id=10.1257/pandp.109.718>.
- Vlaeminck, S. and Herrmann, L.-K. (2015) Data policies and data archives: A new paradigm for academic publishing in economic sciences? In *New Avenues for Electronic Publishing in the Age of Infinite Collections and Citizen Science* (eds. B. Schmidt and M. Dobrev), 145–155. IOS Press.
- Wacharamanotham, C., Eisenring, L., Haroz, S. and Ehtler, F. (2019) Transparency of chi research artifacts: Results of a self-reported survey. URL: osf.io/3bu6t.
- Wilkinson M., Dumontier M., A. I. and et al. (2016) The fair guiding principles for scientific data management and stewardship. *Scientific Data*, **3**.
- Woodward, M. (2012) Automating communications measurement. *Mathworks newsletters*. URL: http://www.mathworks.com/tagteam/73996_92060v00_automating-communications-measurement.pdf.
- Xie, Y. (2015) *Dynamic Documents with R and knitr, Second Edition*. Chapman & Hall/CRC The R Series.
- Xie, Y., Allaire, J. and Golemund, G. (2018) *R Markdown: The Definitive Guide*. Chapman & Hall/CRC The R Series.
-

Acknowledgments

The authors wish to thank the participants of the Banco de Portugal Reproducible Research Workshop in Porto (2019) for the stimulating discussions, which are at the origin of this paper. We are grateful to Virginie Piguet as well as the two anonymous referees for their careful reading and inspiring comments and suggestions.

Overview of statistics and economics journal policies in 2020

We use here a nonexhaustive list of journals. In statistics, we select the most important ones according to the Web of Science index. In economics, we use the journals listed in McCullough (2009). We choose to eliminate some specialized journals or journals publishing mainly theoretical work.¹⁹

TABLE 1 Overview of statistics and economics journal policies in 2020

Journal	Policy	Platform Used
	<i>(In bold mandatory policies, in plain sharing is encouraged)</i>	
Statistics journals		
Annals of Statistics	–	–
Annals of Applied Statistics	(data+code) sharing	– (Supplementary materials archive)
Biometrics	(data+code) sharing	multiple (OSF, Dataverse)
Biometrika	–	–
Comput Stat & Data Analysis	(data+code) sharing	multiple (Mendeley, DRYAD, ICPSR, RunMyCode, ...)
Electronic Journal of Statistics	code sharing	Statlib
J of Business & Econ. Stat.	(data+code) sharing	Figshare
J of Comp. and Graph. Stat.	(data+code) sharing	Figshare
J of Multivariate Analysis	data citation	–
JASA	(data+code sharing)	JASA Dataverse, JASA GitHub
J of the Royal Statistical Society	(data+code) sharing	–
J of Statistical Software	(data+code) sharing	– (Supplementary materials archive)
Stat. Methods in Medical Research	(data+code) sharing	– (Supplementary materials online)
Statistics & Probability Letters	data citation	–
Statistics and Computing	data sharing	multiple (Figshare, Dryad, openICPSR, Dataverse)
Stoch. Proc. & their Applications	data citation	–
CSBIG	data sharing	–
Economics journals		
Am Econ Review	(data + code) sharing	OpenICPSR (AEA Data and code repository)
J Finance	code, data sharing	–
Q J Economics	(data + code) sharing	data repository (Dataverse) linked to the QJE website
Econometrica	(data + code) sharing	– (Supplementary materials webpage)
J Financial Econ	data + code sharing	multiple (Mendeley, DRYAD, ICPSR, RunMyCode)
J Political Econ	(data + code) sharing	– (JPE website)
Rev Financial Stud	–	–
J Econ Theory	(data + code) sharing	multiple (Mendeley, DRYAD, ICPSR, RunMyCode)
Rev Econ Studies	(data + code) sharing	Oxford Journals Review Archive
J Econometrics	(data + code) sharing	multiple (Mendeley, DRYAD, ICPSR, RunMyCode)
J Econ Literature	(data + code) sharing ^{AEA}	OpenICPSR ^{AEA}
J Monetary Econ	(data + code) sharing	multiple (Mendeley, DRYAD, ICPSR, RunMyCode)
J Econ Perspectives	(data + code) sharing ^{AEA}	OpenICPSR ^{AEA}
Rev Econ & Stat	(data + code) sharing	–
Eur Econ Review	(data + code) sharing	– (EER website)
Int Econ Review	–	–
J Int Econ	data, code sharing	Mendeley repository
Economic Journal	–	–
J Public Econ	(data + code) sharing	multiple (Mendeley, DRYAD, ICPSR, RunMyCode)
Game Econ Behav	(data + code) sharing	multiple (Mendeley, DRYAD, ICPSR, RunMyCode)
RAND J Economics	–	–
J Money Credit Bank	(data + code) sharing	web data archives
Economic Theory	data sharing	multiple repositories
J Bus & Econ Stat	(data + code) sharing	– (Supplementary online materials)
Economics Letters	data, code sharing	multiple (Mendeley, DRYAD, ICPSR, RunMyCode, ...)
J Appl Econometrics	data, code sharing	JAIE data archive
A J Political Science	(data+code) sharing	AJPS Dataverse

^{AEA} indicates that the journal follows the strict AEA Data and Code Availability Policy
(see <https://www.aeaweb.org/journals/policies/data-code/>).

¹⁹ See also Orozco et al. (Forthcoming 2021) for the evolution of economic journal replication policies (2003-2019).