

# REPRODUCIBILITY RESEARCH AND OPEN SCIENCE

---

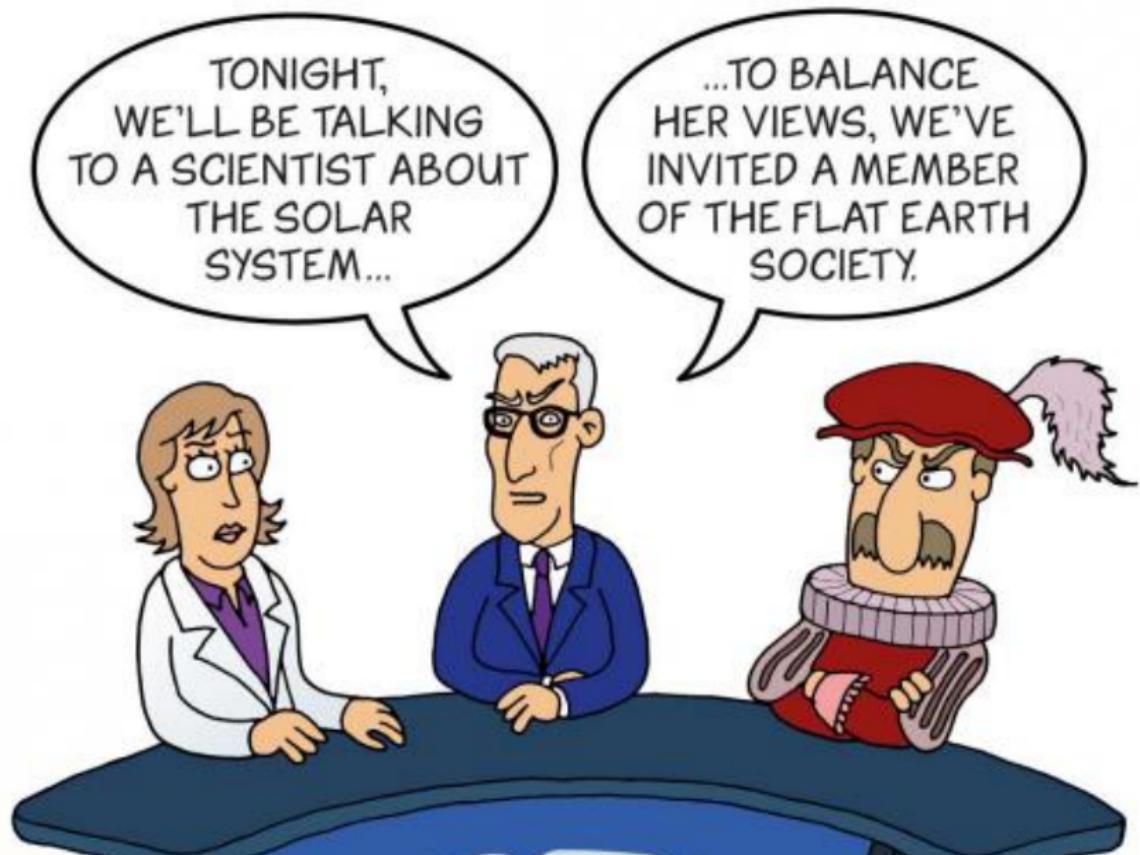
Arnaud Legrand



*Doing a PhD, good practice and pitfalls to avoid*  
December 2025



## SCIENTIFIC CONSENSUS



# NO TRANSPARENCY NO CONSENSUS





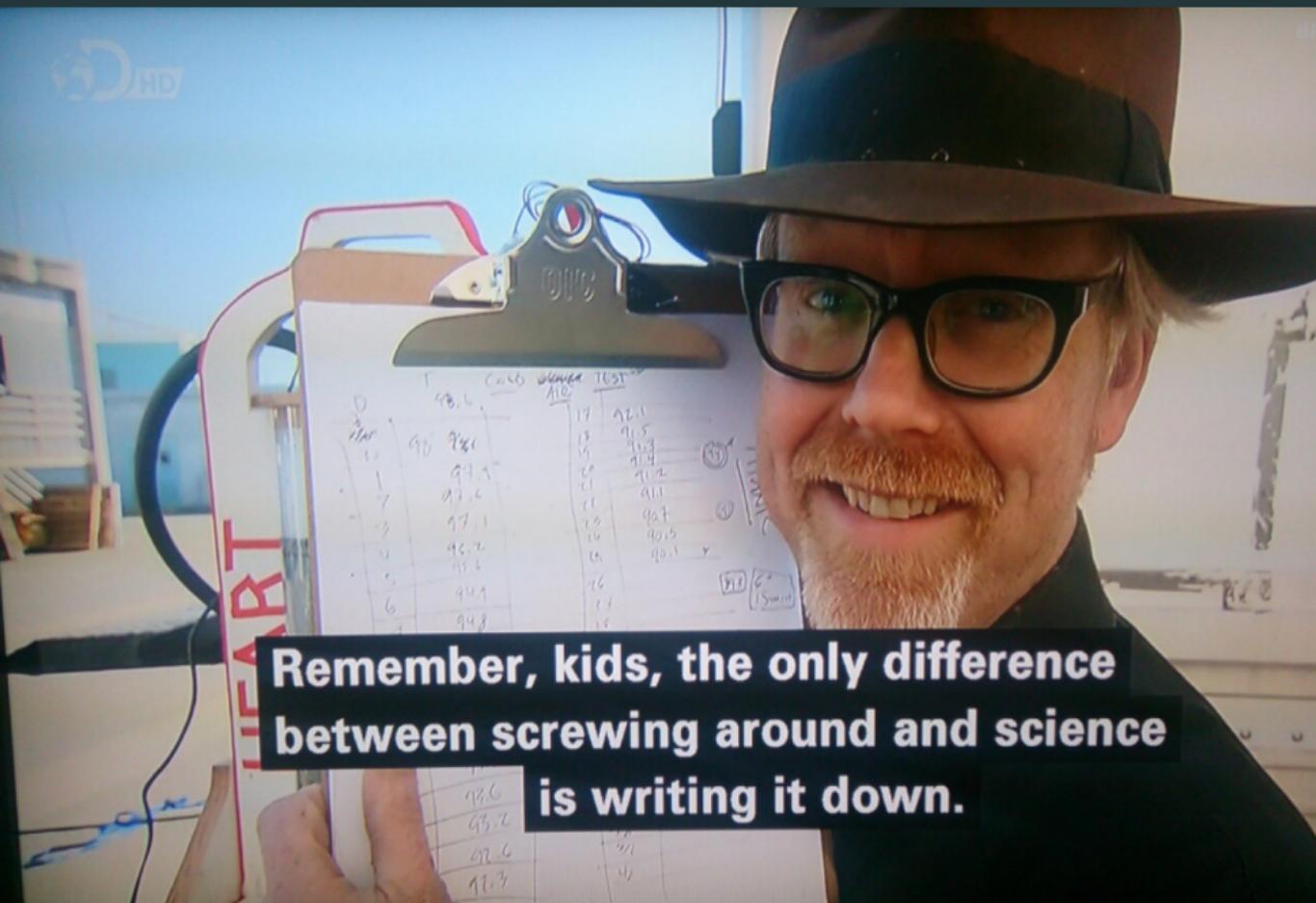
## Author

- I thought I used the same parameters but I'm getting different results!
- The new student wants to compare with the method I proposed last year
- My advisor asked me whether I took care of setting this or this but I can't remember
- The damned fourth reviewer asked for a major revision and wants me to change Figure 3. Which code and which data set did I use?
- It worked yesterday! 6 months later: Why did I do that?

## Reviewer

- As usual, there is no confidence interval, I wonder about the variability and whether the difference is significant or not
- That can't be true, I'm sure they removed some points
- Why is this graph in logscale? How would it look like otherwise? I'm not even sure of what this value means. If only I could access the generation script

# MYTHBUSTERS: SCIENCE VS. SCREWING AROUND



**Remember, kids, the only difference  
between screwing around and science  
is writing it down.**

# COMPUTATIONAL DOCUMENTS...

Document initial dans son environnement

jupyter example\_pi (Python 3)

File Edit View Insert Cell Kernel Widgets Help Hide Code Python 3

# Un document computationnel

Mon ordinateur m'indique que  $\pi$  vaut approximativement

In [1]:

```
from math import *
print(pi)

3,141592653589793
```

Mais calculé avec la `_methods_` des [alg\u00fclles de Buffon](https://fr.wikipedia.org/wiki/Alg\u00fclle_de_Buffon), on obtiendrait comme approximation :

In [2]:

```
import numpy as np
N = 1000000
x = np.random.uniform(size=N, low=0, high=1)
theta = np.random.uniform(size=N, low=0, high=pi/2)
2/(sum((x+np.sin(theta))>1))/N
```

Out[2]:

```
3,143719869498765
```

On peut inclure des formules math\u00e9matiques comme  $\frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$  et des dessins qui n'ont rien \u00e0 voir avec  $\pi$  (si ce n'est une constante de normalisation...).

In [3]:

```
%matplotlib inline
import matplotlib.pyplot as plt

mu, sigma = 100, 15
x = mu + sigma*np.random.randn(100000)

plt.hist(x, 99)
plt.grid(True)
plt.show()
```



Document final

## Un document computationnel

Mon ordinateur m'indique que  $\pi$  vaut approximativement

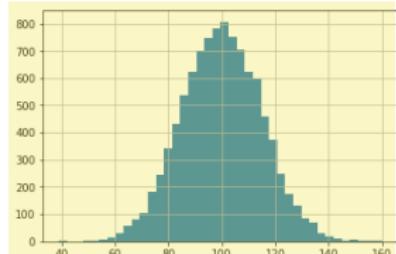
3.141592653589793

Mais calcul\u00e9 avec la [m\u00e9thode des alg\u00fclles de Buffon](#), on obtiendrait comme approximation :

```
import numpy as np
N = 1000000
x = np.random.uniform(size=N, low=0, high=1)
theta = np.random.uniform(size=N, low=0, high=pi/2)
2/(sum((x+np.sin(theta))>1))/N
```

3.143719869498765

On peut inclure des formules math\u00e9matiques comme  $\frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$  et des dessins qui n'ont rien \u00e0 voir avec  $\pi$  (si ce n'est une constante de normalisation...).



Export

# COMPUTATIONAL DOCUMENTS...

Document initial dans son environnement

The screenshot shows a Jupyter Notebook interface with several code cells:

- In [1]:**

```
from math import *  
print(pi)
```

 Output: 3,141592653589793
- In [2]:**

```
import numpy as np  
N = 1000000  
x = np.random.uniform(size=N, low=0, high=1)  
theta = np.random.uniform(size=N, low=0, high=pi/2)  
2/(sum((x+np.sin(theta))>1))/N
```

 Output: 3,1437198694098765
- In [3]:**

```
%matplotlib inline  
import matplotlib.pyplot as plt  
  
mu, sigma = 100, 15  
x = mu + sigma*np.random.randn(10000)  
  
plt.hist(x, 99)  
plt.grid(True)  
plt.show()
```

 Output: A histogram showing a normal distribution centered around 100.



Document final

## Un document computationnel

Mon ordinateur m'indique que  $\pi$  vaut approximativement

3.141592653589793

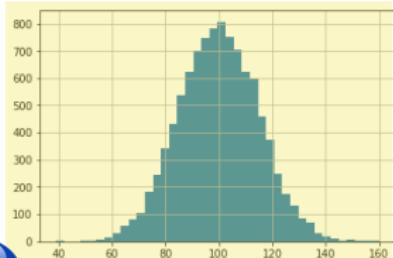
Mais calculé avec la [méthode des aiguilles de Buffon](#), on obtiendrait comme approximation :

```
import numpy as np  
N = 1000000  
x = np.random.uniform(size=N, low=0, high=1)  
theta = np.random.uniform(size=N, low=0, high=pi/2)  
2/(sum((x+np.sin(theta))>1))/N
```

3.1437198694098765

Export

On peut inclure des formules mathématiques comme  $\frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$  et des dessins qui n'ont rien à voir avec  $\pi$  (si ce n'est une constante de normalisation... ☺).



## ... AND LABORATORY NOTEBOOKS

Subject: Fin des pièces jointes dans Mattermost (Inria)

- > Mattermost est le principal outil de communication et d'échanges dans certaines équipes. Nous avons un **canal par projet**, un **canal avec chaque doctorant**, etc. Ce ne sont pas de simples chats "informels"... ce sont de vraiment outils de **travail dans la durée** !
- > Comme XXX, nous aussi l'utilisons beaucoup pour **échanger des fichiers**, des PDF de **papiers**, faire la \*biblio\*...
- Nous avons positionné Mattermost comme un service de discussion instantanée (ou asynchrone, mais pour des messages à court terme).
- Sa finalité [...] n'est ni la gestion documentaire, ni l'archivage de documents, ni le suivi des expérimentations (au sens carnet de laboratoire), même si, finalement, on peut l'utiliser de cette manière.
- Si des documents importants pour votre activité au sein d'Inria sont stockés dans Mattermost, et conservés seulement ici, **c'est un risque pour Inria (et pour vous)**: que se passe-t-il en cas de **départ des agents**? **remplacement du service** par un autre ? **évolution d'une équipe de recherche** ? Qui est le **propriétaire d'un document** partagé ici ?

## COMMON HORROR STORIES 2/4: ARGH... DAMNED COMPUTERS

- Hey! Here is my code. It's on GitHub so feel free to play with it!  
I'm doing open science 😊
  - Alice: I got 3.123123      Bob: I got segfault      Cal: I got 3.123125
- Damned! It used to work!!! Whenever I upgrade my computer,  
things break so I try to stay away from this 😞
- Whenever trying the code of my colleague, I had to install Foo  
but I broke everything and now neither his code nor mine works!  
😞

Seriously ? It's 21st century. 😊 How come all this is so painful ?

# CONTAINERS AND PACKAGE MANAGERS

The good



The bad



The ugly



Automatic tracking

# CONTAINERS AND PACKAGE MANAGERS

The good



The bad



The ugly



Automatic tracking

Containers

- Pros: Lightweight, Good isolation, Easy to use
  - Running as easy as `docker run <img> <cmd>`
  - Building images: `docker build -f <Dockerfile>`
  - Sharing through the Docker Hub: `docker pull/push <img>`

# CONTAINERS AND PACKAGE MANAGERS

The good



The bad



The ugly



Automatic tracking

## Containers

- **Pros:** Lightweight, Good isolation, Easy to use
- **Cons:** Opaque, Container build is generally not reproducible
  - Recipes rarely follow *reproducible good practices*

---

```
1  FROM ubuntu:20.04
2  RUN apt-get update
3      && apt-get upgrade -y
4      && apt-get install -y ...
```

- 
- Choose a stable image (and the smallest possible)
  - Include only the necessary libraries (e.g. no graphics libs)
  - Avoid system updates (instead freeze sources)

# CONTAINERS AND PACKAGE MANAGERS

The good



The bad



The ugly



Automatic tracking

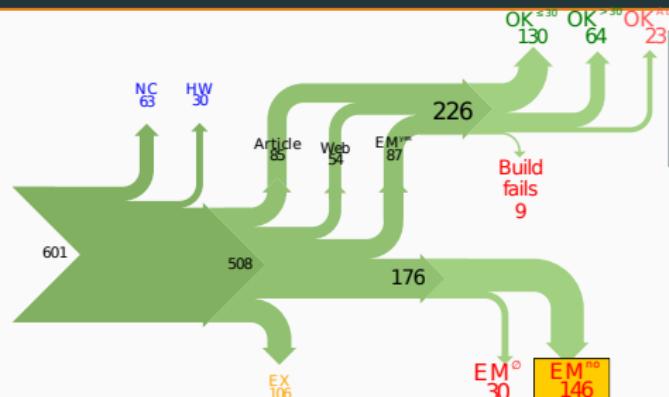
Containers

- Pros: Lightweight, Good isolation, Easy to use
- Cons: Opaque, Container build is generally not reproducible

Package managers (the ugly and the good)

- Language specific: `pip/pipenv/virtualenv`, `conda`, `CRAN/Bioconductor`
  - Limits: version management, durability, permeable, language centric
- **GUIX/NiX = Full-fledged functional package manager**
  - Native support for environment (*à la git*)
  - Isolation through `--pure`
  - Recompile from source (cache recommended)

## COMMON HORROR STORIES 3/4: PLEASE HOLD ON



- Versionning Problems

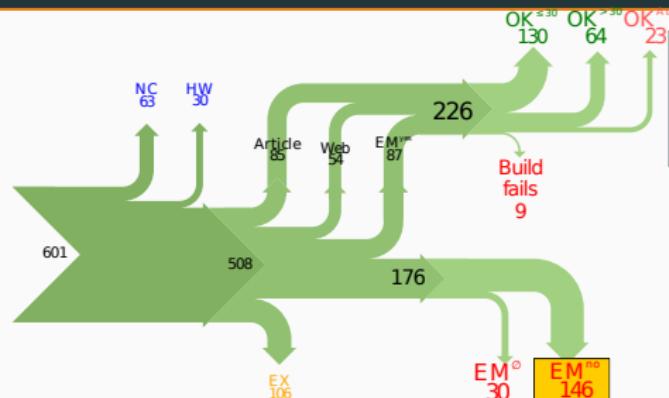
Collberg, Christian et Al., Measuring Reproducibility in Computer Systems Research, <http://reproducibility.cs.arizona.edu/> 2013

- 8 ACM conferences (ASPLOS'12, CCS'12, OOPSLA'12, OSDI'12, PLDI'12, SIGMOD'12, SOSP'11, VLDB'12) and 5 journals
- EM<sup>no</sup> = the code cannot be provided

Thanks for your interest in the implementation of our paper. The good news is that I was able to find some code. I am just hoping that it is a stable working version of the code, and matches the implementation we finally used for the paper. Unfortunately, I have lost some data when my laptop was stolen last year. The bad news is that the code is not commented and/or clean.

Attached is the <system> source code of our algorithm. I'm not very sure whether it is the final version of the code used in our paper, but it should be at least 99% close. Hope it will help.

## COMMON HORROR STORIES 3/4: PLEASE HOLD ON



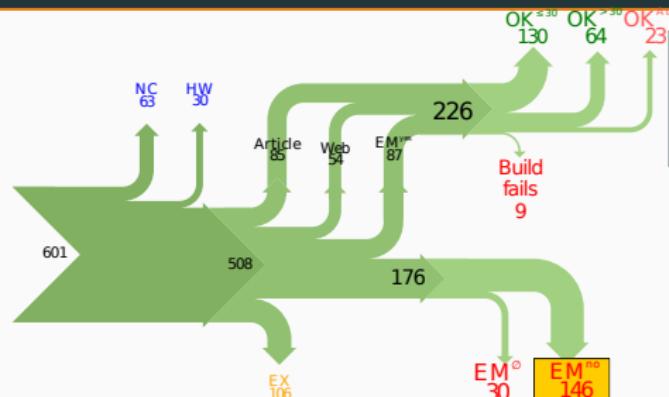
- Versionning Problems
- Bad Backup Practices

Collberg, Christian et Al., *Measuring Reproducibility in Computer Systems Research*,  
<http://reproducibility.cs.arizona.edu/> 2013

- 8 ACM conferences (ASPLOS'12, CCS'12, OOPSLA'12, OSDI'12, PLDI'12, SIGMOD'12, SOSP'11, VLDB'12) and 5 journals
- EM<sup>no</sup>= the code cannot be provided

*Unfortunately, the server in which my implementation was stored had a disk crash in April and three disks crashed simultaneously. While the help desk made significant effort to save the data, my entire implementation for this paper was not found.*

## COMMON HORROR STORIES 3/4: PLEASE HOLD ON



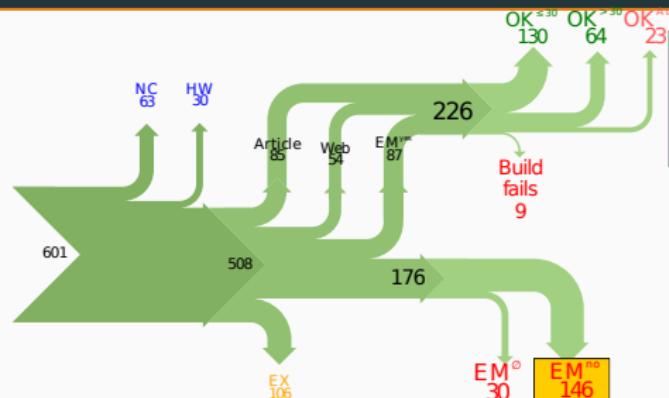
- Versionning Problems
- Bad Backup Practices
- Code Will be Available Soon

Collberg, Christian et Al., Measuring Reproducibility in Computer Systems Research, <http://reproducibility.cs.arizona.edu/> 2013

- 8 ACM conferences (ASPLOS'12, CCS'12, OOPSLA'12, OSDI'12, PLDI'12, SIGMOD'12, SOSP'11, VLDB'12) and 5 journals
- EM<sup>no</sup>= the code cannot be provided

Unfortunately the current system is *not mature enough at the moment*, so it's not yet publicly available. We are actively working on a number of extensions and *things are somewhat volatile*. However, once things stabilize we plan to release it to outside users. At that point, we would be happy to send you a copy.

## COMMON HORROR STORIES 3/4: PLEASE HOLD ON

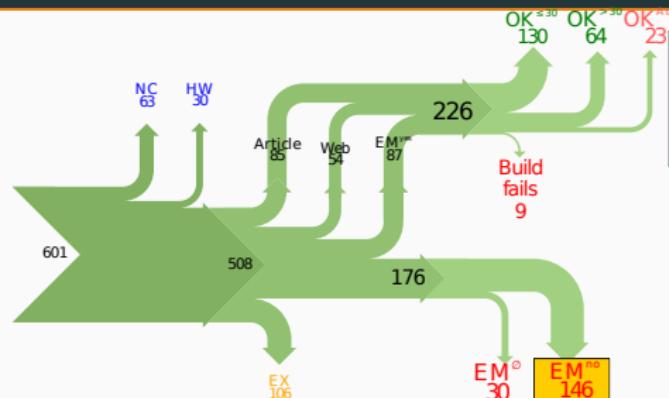


Collberg, Christian et Al., *Measuring Reproducibility in Computer Systems Research*, <http://reproducibility.cs.arizona.edu/> 2013

- Versionning Problems
- Bad Backup Practices
- Code Will be Available Soon
- No Intention to Release

*I am afraid that the source code was never released. The code was never intended to be released so is not in any shape for general use.*

## COMMON HORROR STORIES 3/4: PLEASE HOLD ON



- Versionning Problems
- Bad Backup Practices
- Code Will be Available Soon
- No Intention to Release

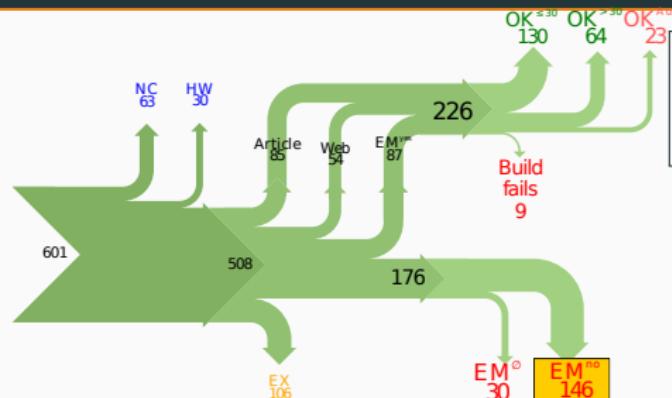
Collberg, Christian et Al., *Measuring Reproducibility in Computer Systems Research*, <http://reproducibility.cs.arizona.edu/> 2013

- 8 ACM conferences (ASPLOS'12, CCS'12, OOPSLA'12, OSDI'12, PLDI'12, SIGMOD'12, SOSP'11, VLDB'12) and 5 journals
- EM<sup>no</sup>= the code cannot be provided
  - Programmer Left

*(STUDENT)* was a graduate student in our program but *he left a while back* so I am responding instead. For the paper we used a prototype that included many moving pieces that only *(STUDENT)* knew how to operate and we did not have the time to integrate them in a ready-to-share implementation before he left. Still, I hope you can build on the ideas/technique of the paper.

Unfortunately, the author who has done most of the coding for this paper has *passed away* and the code is no longer maintained.

## COMMON HORROR STORIES 3/4: PLEASE HOLD ON



- Versionning Problems
- Bad Backup Practices
- Code Will be Available Soon
- No Intention to Release

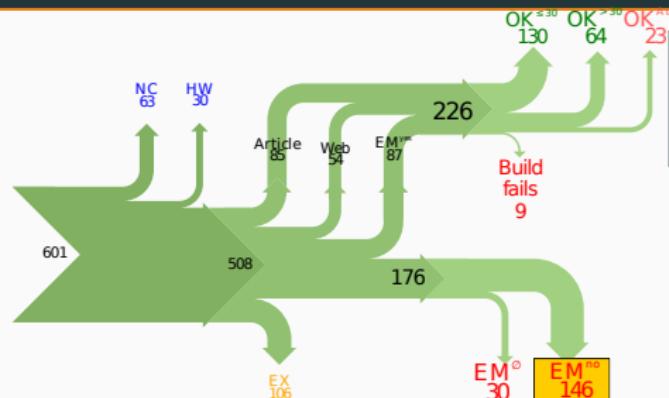
Collberg, Christian et Al., *Measuring Reproducibility in Computer Systems Research*, <http://reproducibility.cs.arizona.edu/> 2013

- 8 ACM conferences (ASPLOS'12, CCS'12, OOPSLA'12, OSDI'12, PLDI'12, SIGMOD'12, SOSP'11, VLDB'12) and 5 journals
- EM<sup>no</sup>= **the code cannot be provided**
  - Programmer Left
  - Commercial Code

Since this work has been done at **(COMPANY)** we don't open-source code unless there is a compelling business reason to do so. So unfortunately I don't think we'll be able to share it with you.

The code owned by **(COMPANY)**, and AFAIK the code is not open-source. Your best bet is to reimplement :( Sorry.

## COMMON HORROR STORIES 3/4: PLEASE HOLD ON



- Versionning Problems
- Bad Backup Practices
- Code Will be Available Soon
- No Intention to Release

Collberg, Christian et Al., *Measuring Reproducibility in Computer Systems Research*, <http://reproducibility.cs.arizona.edu/> 2013

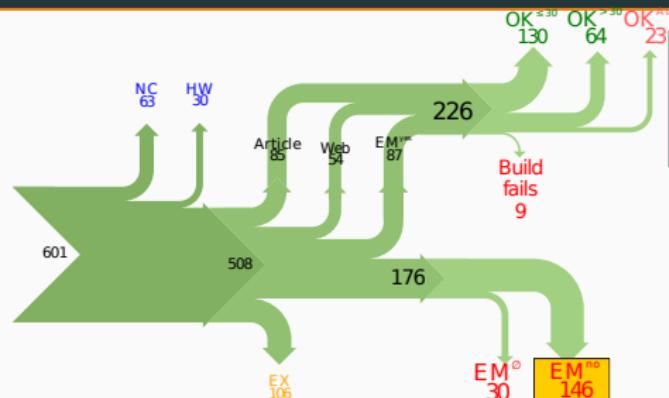
- 8 ACM conferences (ASPLOS'12, CCS'12, OOPSLA'12, OSDI'12, PLDI'12, SIGMOD'12, SOSP'11, VLDB'12) and 5 journals
- EM<sup>no</sup> = **the code cannot be provided**
  - Programmer Left
  - Commercial Code
  - Proprietary Academic Code

Unfortunately, the *(SYSTEM)* sources are not meant to be opensource (the code is partially *property of (UNIVERSITY 1), (UNIVERSITY 2) and (UNIVERSITY 3)*.)

If this will change I will let you know, albeit I do not think there is an intention to make the *(SYSTEM)* sources opensource in the near future.

If you're interested in obtaining the code, we only ask for a description of the research project that the code will be used in (*which may lead to some joint research*), and we also have a software license agreement that the University would need to sign.

## COMMON HORROR STORIES 3/4: PLEASE HOLD ON



- Versionning Problems
- Bad Backup Practices
- Code Will be Available Soon
- No Intention to Release

Collberg, Christian et Al., *Measuring Reproducibility in Computer Systems Research*, <http://reproducibility.cs.arizona.edu/> 2013

- 8 ACM conferences (ASPLOS'12, CCS'12, OOPSLA'12, OSDI'12, PLDI'12, SIGMOD'12, SOSP'11, VLDB'12) and 5 journals
- EM<sup>no</sup> = the code cannot be provided
  - Programmer Left
  - Commercial Code
  - Proprietary Academic Code
  - Research vs. Sharing

*In the past when we attempted to share it, we found ourselves spending more time getting outsiders up to speed than on our own research. So I finally had to establish the policy that we will not provide the source code outside the group.*

# CHANGING RESEARCH PRACTICES

## Soft. Engineering, Statistics, and Reproducible Research in the curricula

**Manifesto:** "*I solemnly pledge*" ([WSSSPE](#), [Lorena Barba](#), [FAIR](#))

1. I will teach my graduate students about reproducibility
2. All our research code (and writing) is under version control
3. We will always carry out verification and validation
4. We will share data, plotting script & figure under CC-BY
5. We will upload the preprint to arXiv at the time of submission of a paper
6. We will release code at the time of submission of a paper
7. We will add a "Reproducibility" declaration at the end of each paper
8. I will keep an up-to-date web presence



Learn and Teach using online resources like

- [Software Carpentry](#), [The Turing Way](#), ...

# CHANGING PUBLISHING PRACTICES

## Artifact evaluation and ACM badges



## Major conferences

- Supercomputing: Artifact Description (AD) mandatory, Artifact Evaluation (AE) still optional, Double blind vs. RR
- NeurIPS, ICLR: open reviews, reproducibility challenge



Joelle Pineau @ NeurIPS'18

- ACM SIGMOD 2015-2019, Most Reproducible Paper Award...

Mentalities are evolving people care, make stuff available, errors are found and fixed

## HORROR STORIES 4/4: FIGHTING INFORMATION LOSS WITH ARCHIVES



or



= awesome collaborations ( $\neq$  archive)

- D. Spinellis. *The Decay and Failures of URL References*. CACM, 46(1), 2003  
*The half-life of a referenced URL is approximately 4 years*
- P. Habibzadeh. *Decay of References to Web sites in Articles Published in General Medical Journals: Mainstream vs Small Journals*. Applied Clinical Informatics. 4 (4), 2013  
*half life ranged from 2.2 years in EMHJ to 5.3 years in BMJ*
- Discontinued forges: Code Space, Gitorious, Google code, Inria Gforge

## HORROR STORIES 4/4: FIGHTING INFORMATION LOSS WITH ARCHIVES



or



= awesome collaborations ( $\neq$  archive)

- D. Spinellis. *The Decay and Failures of URL References*. CACM, 46(1), 2003  
*The half-life of a referenced URL is approximately 4 years*
- P. Habibzadeh. *Decay of References to Web sites in Articles Published in General Medical Journals: Mainstream vs Small Journals*. Applied Clinical Informatics. 4 (4), 2013  
*half life ranged from 2.2 years in EMHJ to 5.3 years in BMJ*
- Discontinued forges: Code Space, Gitorious, Google code, Inria Gforge

Article archives

arXiv.org

**HAL**  
archives-ouvertes.fr

Data archives



**figshare**

**zenodo**



Software Archive



**Software Heritage**

Collect/Preserve/Share

Plan for disaster with **git** and **git-annex** (not **git LFS!**)

Separation between articles, code, and data is not so simple though

# DIFFERENT REPRODUCIBILITY CONCERN IN MODERN SCIENCE

**Social Sciences, Oncology, ...** methodology, statistics, pre-registration

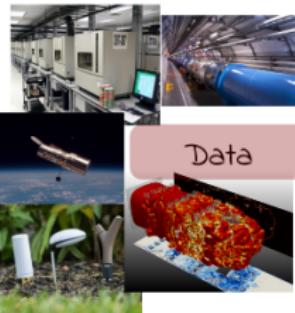
**Genomics** software engineering, computational reproducibility, provenance

**Computational fluid dynamics** numerical issues

**Artificial Intelligence** most of the above

*The processing steps between raw observations and findings have gotten increasingly numerous and complex*

## Authors



## Data

# DIFFERENT REPRODUCIBILITY CONCERN IN MODERN SCIENCE

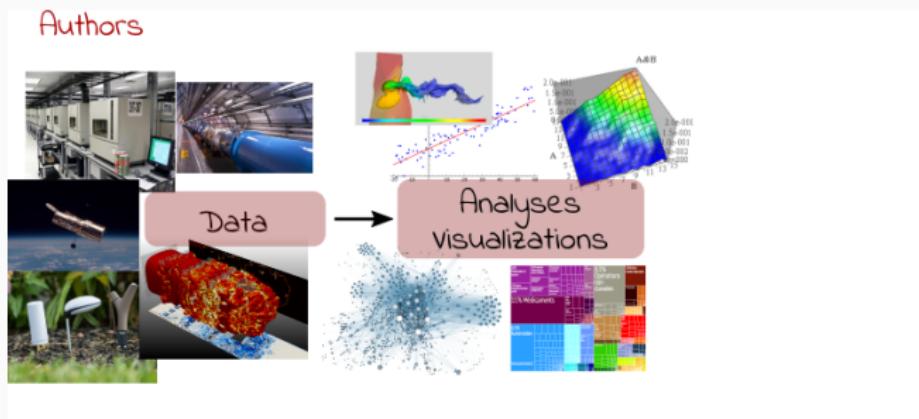
**Social Sciences, Oncology, ...** methodology, statistics, pre-registration

**Genomics** software engineering, computational reproducibility, provenance

**Computational fluid dynamics** numerical issues

**Artificial Intelligence** most of the above

*The processing steps between raw observations and findings have gotten increasingly numerous and complex*



# DIFFERENT REPRODUCIBILITY CONCERN IN MODERN SCIENCE

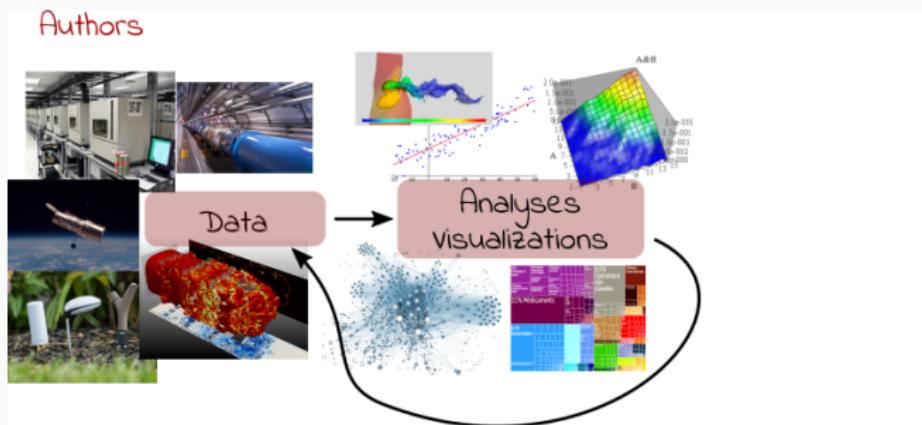
**Social Sciences, Oncology, ...** methodology, statistics, pre-registration

**Genomics** software engineering, computational reproducibility, provenance

**Computational fluid dynamics** numerical issues

**Artificial Intelligence** most of the above

*The processing steps between raw observations and findings have gotten increasingly numerous and complex*



# DIFFERENT REPRODUCIBILITY CONCERN IN MODERN SCIENCE

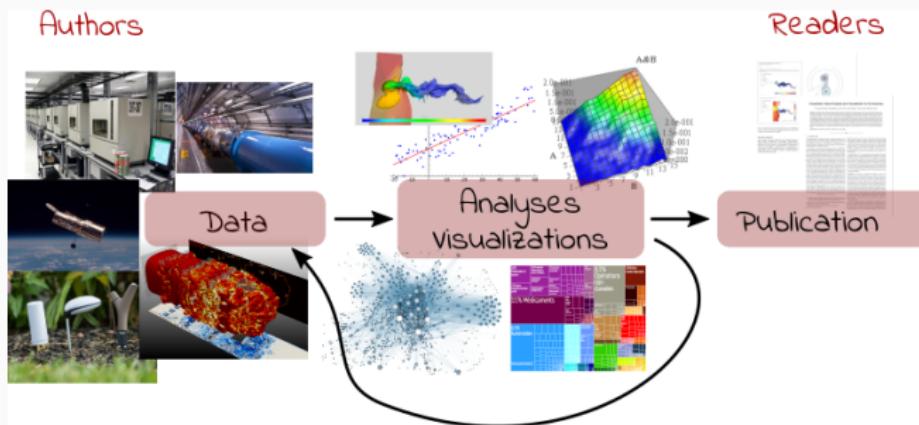
**Social Sciences, Oncology, ...** methodology, statistics, pre-registration

**Genomics** software engineering, computational reproducibility, provenance

**Computational fluid dynamics** numerical issues

**Artificial Intelligence** most of the above

*The processing steps between raw observations and findings have gotten increasingly numerous and complex*



# DIFFERENT REPRODUCIBILITY CONCERN IN MODERN SCIENCE

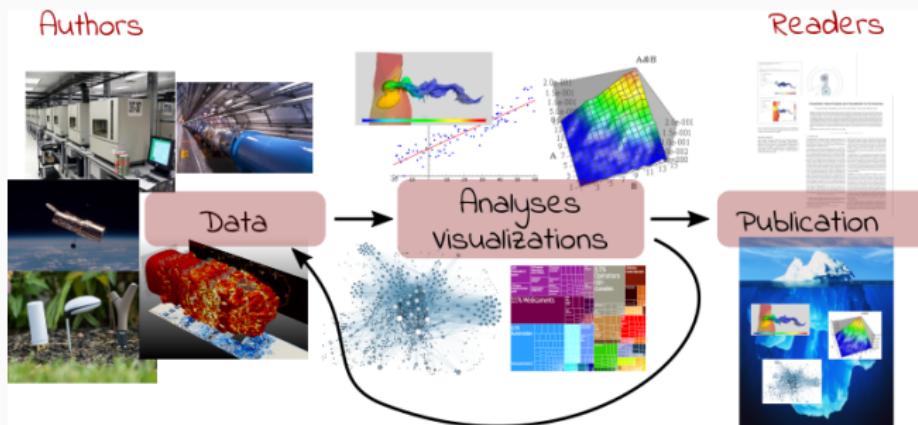
**Social Sciences, Oncology, ...** methodology, statistics, pre-registration

**Genomics** software engineering, computational reproducibility, provenance

**Computational fluid dynamics** numerical issues

**Artificial Intelligence** most of the above

*The processing steps between raw observations and findings have gotten increasingly numerous and complex*



# DIFFERENT REPRODUCIBILITY CONCERN IN MODERN SCIENCE

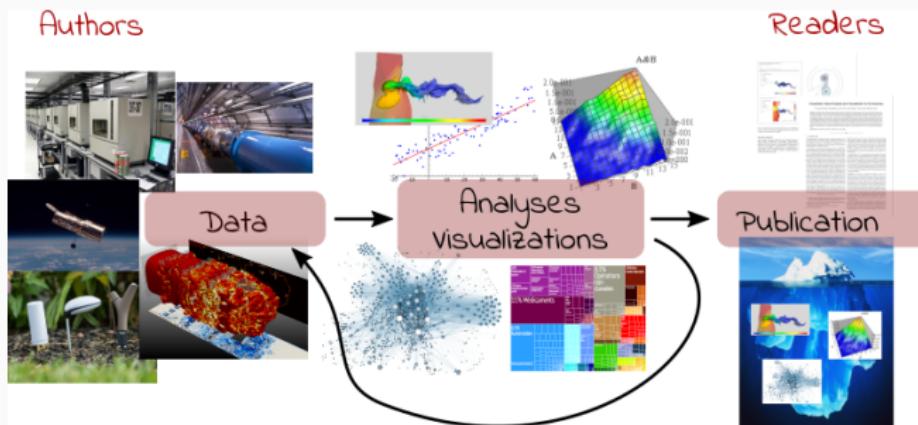
**Social Sciences, Oncology, ...** methodology, statistics, pre-registration

**Genomics** software engineering, computational reproducibility, provenance

**Computational fluid dynamics** numerical issues

**Artificial Intelligence** most of the above

*The processing steps between raw observations and findings have gotten increasingly numerous and complex*



Reproducible Research = Bridging the Gap by working Transparently 12/14

# REPRODUCIBLE RESEARCH ~ OPEN SCIENCE ?

Plan National pour la Science Ouverte (BSN ~ CoSO)

- CNRS, Inria, INRAE, ...
- Many flavors: *Citizen Science*

## Main pillars:

1. Open access
2. Open data
3. Open source
  - Open hardware
4. Open methodology (**Reproducible Research**)
  - Open-notebook science
  - Open science infrastructures
5. Open peer review (avoid **collusion**)
6. Open educational resources



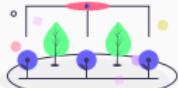
**NO TRANSPARENCY  
NO CONSENSUS**



Obviously **making code/data available for the reproduction of results from published papers has become the new norm**

# RESOURCES AND ACKNOWLEDGMENTS

Vers une recherche  
reproductible  
Faire évoluer ses pratiques

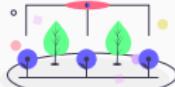


Lion Dangelier, Sabrina Drange, Boris Huguenin,  
Arnaud Legendre, Pascal Pervot, Nicolas Rouger

A non-technical introduction to reproducibility issues

# RESOURCES AND ACKNOWLEDGMENTS

Vers une recherche  
reproductible  
Faire évoluer ses pratiques

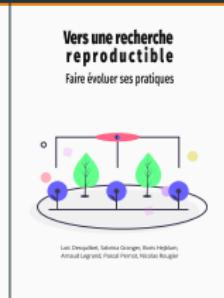


Lion Dangelier, Sabrina Grange, Boris Juhani,  
Amaury Legendre, Pascal Pervit, Nicolas Rougerie

## A non-technical introduction to reproducibility issues



# RESOURCES AND ACKNOWLEDGMENTS



## A non-technical introduction to reproducibility issues



## MOOC Reproducible Research: Methodological principles for a transparent science, Learning Lab

- Konrad Hinsen, Christophe Pouzat
- Laurence Farhi, Madeline Montigny, ...
- 2018 – ... (25,000+)



# RESOURCES AND ACKNOWLEDGMENTS

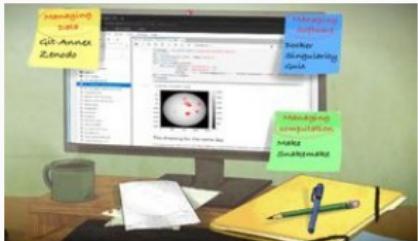


## A non-technical introduction to reproducibility issues



## MOOC Reproducible Research: Methodological principles for a transparent science, Learning Lab

- Konrad Hinsen, Christophe Pouzat
- Laurence Farhi, Madeline Montigny, ...
- 2018 – ... (25,000+)



## MOOC RR2: Practices and tools for managing computations and data 2024 – ... (3,000+)

- Managing data
- Software environment control
- Scientific workflow

([git annex](#), Zenodo, SWH)  
([docker](#), [singularity](#), [guix](#))  
([make](#), [snakemake](#))

THAT'S ALL FOLKS!

---