

RECHERCHE REPRODUCTIBLE: BONNES PRATIQUES POUR UNE SCIENCE TRANSPARENTE

Arnaud Legrand



Congrès de la ROADEF

Champs-sur-Marne, Février 2025



WHAT'S EXPECTED FROM SCIENCE?



Goal

- **Describe** and **explain** **Nature** (human organizations and artificial objects as well)

WHAT'S EXPECTED FROM SCIENCE?



Goal

- **Describe** and **explain** **Nature** (human organizations and artificial objects as well)
- Find **solutions** to **society** needs and issues
 - Feed technological development and **industry**: health, energy, communications, ...
 - Enlighten public decisions: educate citizens, politics, regulation, ...

WHAT'S EXPECTED FROM SCIENCE?



Goal

- **Describe** and **explain** **Nature** (human organizations and artificial objects as well)
- Find **solutions** to **society** needs and issues
 - Feed technological development and **industry**: health, energy, communications, ...
 - Enlighten public decisions: educate citizens, politics, regulation, ...

What else distinguishes science from other human activities?

WHAT'S EXPECTED FROM SCIENCE?



Goal

- **Describe** and **explain** **Nature** (human organizations and artificial objects as well)
- Find **solutions** to **society** needs and issues
 - Feed technological development and **industry**: health, energy, communications, ...
 - Enlighten public decisions: educate citizens, politics, regulation, ...

What else distinguishes science from other human activities?

Method

SCIENTIFIC CONSENSUS



NO TRANSPARENCY NO CONSENSUS



NO TRANSPARENCY NO CONSENSUS



Enlightening the society requires moral/methodological/technical warranties

MERTON'S FOUR NORMS OF SCIENCE (1940)

Universalism Scientific **validity** is independent of the sociopolitical status/personal attributes of its participants (origin, gender, sexuality, religion, etc.)

Communism/communalism All scientists should have common ownership of scientific goods (intellectual property), to promote collective collaboration; **secrecy is the opposite of this norm**

- We say “Newton’s law” to remember that Newton made the original discovery, but not because he has any property of this law

Disinterestedness Scientific institutions act for the **benefit of a common scientific enterprise**, rather than for the personal gain of individuals

Organized skepticism Scientific claims should be exposed to **critical scrutiny** (including **reproduction**) before being accepted

MERTON'S FOUR NORMS OF SCIENCE (1940)

Universalism Scientific **validity** is independent of the sociopolitical status/personal attributes of its participants (origin, gender, sexuality, religion, etc.)

Communism/communalism All scientists should have common ownership of scientific goods (intellectual property), to promote collective collaboration; **secrecy is the opposite of this norm**

- We say “Newton’s law” to remember that Newton made the original discovery, but not because he has any property of this law

Disinterestedness Scientific institutions act for the **benefit of a common scientific enterprise**, rather than for the personal gain of individuals

Organized skepticism Scientific claims should be exposed to **critical scrutiny** (including **reproduction**) before being accepted

Merton was not saying “*you should follow these norms*”, but rather
“*this is the behavior belief I am observing in the scientific community*”

TRANSPARENT RESEARCH ~ OPEN SCIENCE ?

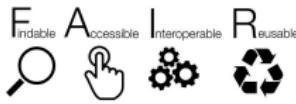
Plan National pour la Science Ouverte: CNRS, Inria, INRAE, Citizen Science

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



(started before 2000)



CC-BY-SA by default at CNRS!



**NO TRANSPARENCY
NO CONSENSUS**



TRANSPARENT RESEARCH ~ OPEN SCIENCE ?

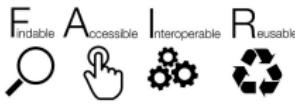
Plan National pour la Science Ouverte: CNRS, Inria, INRAE, Citizen Science

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



(started before 2000)



CC-BY-SA by default at CNRS!



**NO TRANSPARENCY
NO CONSENSUS**



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

E.g., Artifact evaluation and ACM badges



COMMON REPRODUCIBILITY PITFALLS

GO READ THE PAPER BY SMITH ET. AL. 2009



Access through your institution

Purchase PDF

Article preview

Abstract

Introduction

Section snippets

References (30)

Cited by (31)



Future Generation Computer Systems

Volume 25, Issue 3, March 2009, Pages 315-325



Secure on-demand grid computing

M. Smith , M. Schmidt , N. Fallenbeck , T. Dörnemann , C. Schridde , B. Freisleben

Show more

+ Add to Mendeley Share Cite

<https://doi.org/10.1016/j.future.2008.03.002>

[Get rights and content](#)

Abstract

In this paper, a novel approach for enabling Grid users to autonomously install and use custom software on demand using an image creation station is presented, while at the same time offering new security mechanisms to protect both software and data from other Grid users and external attackers. An automated dynamic firewalling mechanism

GO READ THE PAPER BY SMITH ET. AL. 2009

Purchase options

▼ Corporate

For R&D professionals working in corporate organizations.

▲ Academic and personal

For academic or personal use only.

US\$27.95

Local taxes may apply

Online access for 48 hours with the option to save or download the article in PDF format. [Learn more ↗](#)

Add to Cart

Looking for a customized option?

Contact sales for special pricing for your organization

Contact sales ↗

- Use your institution subscription... or Sci-Hub 😊,

GO READ THE PAPER BY SMITH ET. AL. 2009

 Access through your institution [View Open Manuscript](#) [Purchase PDF](#)

[Article preview](#)
[Abstract](#)
[Introduction](#)
[Section snippets](#)
[References \(32\)](#)
[Cited by \(1\)](#)

**Journal of Parallel and Distributed Computing**
Volume 166, August 2022, Pages 111-125 

Simulation-based optimization and sensibility analysis of MPI applications: Variability matters

Tom Cornebize , Arnaud Legrand  

Show more 

+ Add to Mendeley  Share  Cite

<https://doi.org/10.1016/j.jpdc.2022.04.002>  Get rights and content 

Abstract

Finely tuning MPI applications and understanding the influence of key parameters (number of processes, granularity, collective operation algorithms, virtual topology, and

- Use your institution subscription... or Sci-Hub 😊, or HAL/Arxiv

GO READ THE PAPER BY SMITH ET. AL. 2009

HAL

Science ouverte

Chercher un document, un auteur, un mot clef...

Télécharger pour visualiser

Article Dans Une Revue Journal of Parallel and Distributed Computing Année : 2022

Simulation-based Optimization and Sensibility Analysis of MPI Applications: Variability Matters

Tom Cornebize (1, 2), Arnaud Legrand (3, 1)

Afficher plus de détails

1 POLARIS - Performance analysis and optimization of LARge Infrastructures and Systems
2 UGA - Université Grenoble Alpes
3 CNRS - Centre National de la Recherche Scientifique

Résumé en

Finely tuning MPI applications and understanding the influence of key parameters (number of processes, granularity, collective operation algorithms, virtual topology, and process placement) is critical to obtain

Mots clés en

Simulation validation
sensibility analysis SimGrid

The screenshot shows a detailed view of a research paper on the HAL platform. At the top, there's a navigation bar with 'HAL' and language options ('FR'). Below it is a search bar. The main content area has a sidebar on the left with download links and identification details. The main panel displays the article title, authors, and abstract. It also lists funding institutions and provides a summary and keywords section.

- Use your institution subscription... or Sci-Hub 😊, or HAL/Arxiv

Rodriguez et al., CONCUR'15

Unfolding-based Partial Order Reduction*

César Rodríguez¹, Marcelo Sousa², Subodh Sharma³, and
Daniel Kroening⁴

¹ Université Paris 13, Sorbonne Paris Cité, LIPN, CNRS, France

^{2,4} Department of Computer Science, University of Oxford, UK

³ Indian Institute of Technology Delhi, India

Abstract

Partial order reduction (POR) and net unfoldings are two alternative methods to tackle state-space explosion caused by concurrency. In this paper, we propose the combination of both approaches in an effort to combine their strengths. We first define, for an abstract execution model, unfolding semantics parameterized over an arbitrary independence relation. Based on it, our main contribution is a novel stateless POR algorithm that explores at most one execution per Mazurkiewicz trace, and in general, can explore exponentially fewer, thus achieving a form of *super-optimality*. Furthermore, our unfolding-based POR copes with non-terminating executions and incorporates state-caching. Over benchmarks with busy-waits, among others, our experiments show a dramatic reduction in the number of executions when compared to a

JUST COMPARE TO THE ALGORITHM THEY PROPOSED

Rodriguez et al., CONCUR'15

Algorithm 1: An unfolding-based POR exploration algorithm.

```
1 Initially, set  $U := \{\perp\}$ , set  $G := \emptyset$ , and call Explore( $\{\perp\}, \emptyset, \emptyset$ ).
2 Procedure Explore( $C, D, A$ )
3   Extend( $C$ )
4   if  $\text{en}(C) = \emptyset$  return
5   if  $A = \emptyset$ 
6     | Choose  $e$  from  $\text{en}(C)$ 
7   else
8     | Choose  $e$  from  $A \cap \text{en}(C)$ 
9   Explore( $C \cup \{e\}, D, A \setminus \{e\}$ )
10  if  $\exists J \in \text{Alt}(C, D \cup \{e\})$ 
11    | Explore( $C, D \cup \{e\}, J \setminus C$ )
12  Remove( $e, C, D$ )
13 Procedure Extend( $C$ )
14   | Add  $ex(C)$  to  $U$ 
15 Procedure Remove( $e, C, D$ )
16   | Move  $\{e\} \setminus Q_{C,D,U}$  from  $U$  to  $G$ 
17   | foreach  $\hat{e} \in \#_U^i(e)$ 
18     |   | Move  $[\hat{e}] \setminus Q_{C,D,U}$  from  $U$  to  $G$ 
```

- Looks good!

JUST COMPARE TO THE ALGORITHM THEY PROPOSED

Rodriguez et al., CONCUR'15

We give some new definitions. Let C be a configuration of \mathcal{U} . The *extensions* of C , written $ex(C)$, are all those events outside C whose causes are included in C . Formally, $ex(C) := \{e \in E : e \notin C \wedge [e] \subseteq C\}$. We let $en(C)$ denote the set of events *enabled* by C , i.e., those corresponding to the transitions enabled at $state(C)$, formally defined as $en(C) := \{e \in ex(C) : C \cup \{e\} \in conf(\mathcal{U})\}$. All those events in $ex(C)$ which are not in $en(C)$ are the *conflicting extensions*, $cex(C) := \{e \in ex(C) : \exists e' \in C, e \#^i e'\}$. Clearly, sets $en(C)$ and $cex(C)$ partition the set $ex(C)$. Lastly, we define $\#^i(e) := \{e' \in E : e \#^i e'\}$, and $\#_U^i(e) := \#^i(e) \cap U$. The difference between both is that $\#^i(e)$ contains events from *anywhere* in the unfolding structure, while $\#_U^i(e)$ can only *see* events in U .

The algorithm is given in [Alg. 1](#). `Explore`(C, D, A), the main procedure, is given the configuration that is to be explored as the parameter C . The parameter D (for *disabled*) is the set of set of events that have already been explored and prevents that `Explore()` repeats work. It can be seen as a *sleep set* [7]. Set A (for *add*) is occasionally used to guide the direction of the exploration.

Additionally, a global set U stores all events presently known to the algorithm. Whenever some event can safely be discarded from memory, `Remove` will move it from U to G (for *garbage*). Once in G , it can be discarded at any time, or be preserved in G in order to save work when it is re-inserted in U . Set G is thus our *cache memory* of events.

The key intuition in [Alg. 1](#) is as follows. A call to `Explore`(C, D, A) visits all maximal configurations of \mathcal{U} which contain C and do not contain D ; and the first one explored will

- Looks good! Err... **not so simple**. Depending on how you do this, you quickly move from polynomial to exponential.

JUST COMPARE TO THE ALGORITHM THEY PROPOSED

Rodriguez et al., CONCUR'15

POSIX threads. The analyzer accepts deterministic programs, implements a variant of [Alg. 1](#) where the computation of the alternatives is memoized, and supports cutoffs events with the criteria defined in [§ 5](#).

We ran POET on a number of multi-threaded C programs. Most of them are adapted from benchmarks of the Software Verification Competition [17]; others are used in related works [8, 19, 2]. We investigate the characteristics of average program unfoldings (depth, width, etc.) as well as the frequency and impact of cutoffs on the exploration. We also compare POET with NIDHUGG [1], a state-of-the-art stateless model checking for multi-threaded C programs that implements Source-DPOR [2], an efficient but non-optimal DPOR. All experiments were run on an Intel Xeon CPU with 2.4 GHz and 4 GB memory. [Tables 1](#) and [2](#) give our experimental data for programs with acyclic and non-acyclic state spaces, respectively.

For programs with acyclic state spaces ([Table 1](#)), POET with and without cutoffs seems to perform the same exploration when the unfolding has no cutoffs, as expected. Furthermore, the number of explored executions also coincides with NIDHUGG when the latter reports 0 sleep-set blocked executions (cf., [§ 4](#)), providing experimental evidence of POET's optimality.

The unfoldings of most programs in [Table 1](#) do not contain cutoffs. All these programs are deterministic, and many of them highly sequential (STF, SPIN08, FIB), features known to make cutoffs unlikely. CCNF(n) are concurrent programs composed of $n - 1$ threads where thread i and $i + 1$ race on writing one variable, and are independent of all remaining

³ Source code and benchmarks available from: <http://www.cs.ox.ac.uk/people/marcelo.sousa/poet/>.

- Looks good! Err... **not so simple**. Depending on how you do this, you quickly move from polynomial to exponential.
- **Experiments!**

JUST COMPARE TO THE ALGORITHM THEY PROPOSED

Rodriguez et al., CONCUR'15

■ **Table 1** Programs with acyclic state space. Columns are: $|P|$: nr. of threads; $|I|$: nr. of explored traces; $|B|$: nr. of sleep-set blocked executions; $t(s)$: running time; $|E|$: nr. of events in \mathcal{U} ; $|E_{\text{cut}}|$: nr. of cutoff events; $|\Omega|$: nr. of maximal configurations; $\langle |U_\Omega| \rangle$: avg. nr. of events in U when exploring a maximal configuration. A * marks programs containing bugs. <7K reads as “fewer than 7000”.

| Benchmark | NIDHUGG | POET (without cutoffs) | | | | POET (with cutoffs) | | | | | | | |
|-----------|----------------|------------------------|-------|--------|-------|---------------------|------------------------------|--------|-------|--------------------|------------|------------------------------|--------|
| | | $ I $ | $ B $ | $t(s)$ | $ E $ | $ \Omega $ | $\langle U_\Omega \rangle$ | $t(s)$ | $ E $ | $ E_{\text{cut}} $ | $ \Omega $ | $\langle U_\Omega \rangle$ | $t(s)$ |
| STF | 3 | 6 | 0 | 0.06 | 121 | 6 | 79 | 0.04 | 121 | 0 | 6 | 79 | 0.06 |
| STF* | 3 | - | - | 0.05 | - | - | - | 0.02 | - | - | - | - | 0.03 |
| SPIN08 | 3 | 84 | 0 | 0.08 | 2974 | 84 | 1506 | 2.04 | 2974 | 0 | 84 | 1506 | 2.93 |
| FIB | 3 | 8953 | 0 | 3.36 | <185K | 8953 | 92878 | 305 | <185K | 0 | 8953 | 92878 | 704 |
| FIB* | 3 | - | - | 0.74 | - | - | - | 81.0 | - | - | - | - | 133 |
| CCNF(9) | 9 | 16 | 0 | 0.05 | 49 | 16 | 46 | 0.07 | 49 | 0 | 16 | 46 | 0.06 |
| CCNF(17) | 17 | 256 | 0 | 0.15 | 97 | 256 | 94 | 5.76 | 97 | 0 | 256 | 94 | 6.09 |
| CCNF(19) | 19 | 512 | 0 | 0.28 | 109 | 512 | 106 | 22.5 | 109 | 0 | 512 | 106 | 22.0 |
| SSB | 5 | 4 | 2 | 0.05 | 48 | 4 | 38 | 0.03 | 46 | 1 | 4 | 37 | 0.03 |
| Ssb(1) | 5 | 22 | 14 | 0.06 | 245 | 23 | 143 | 0.11 | 237 | 4 | 23 | 140 | 0.11 |
| SSB(3) | 5 | 169 | 67 | 0.12 | 2798 | 172 | 1410 | 3.51 | 1179 | 48 | 90 | 618 | 0.90 |
| SSB(4) | 5 | 336 | 103 | 0.15 | <7K | 340 | 3333 | 20.3 | 2179 | 74 | 142 | 1139 | 2.07 |
| SSB(8) | 5 | 2014 | 327 | 0.85 | <67K | 2022 | 32782 | 4118 | <12K | 240 | 470 | 6267 | 32.1 |

Lastly, we note that this cutoff approach imposes no liability on what events shall be kept in the prefix, set G can be cleaned at discretion. Also, redefining (7) to use adequate orders [5] is straightforward, cf. App. C.1 (in our proofs we actually assume adequate orders).

- Looks good! Err... **not so simple**. Depending on how you do this, you quickly move from polynomial to exponential.
- **Experiments!** \triangle Possible 404, code not found! ahead!!!

JUST COMPARE TO THE ALGORITHM THEY PROPOSED

Rodriguez et al., CONCUR'15

- Looks good! Err... **not so simple**. Depending on how you do this, you quickly move from polynomial to exponential.
- **Experiments!** \triangle Possible 404, *code not found!* ahead!!!
- Wait, what's this language? Did this ever run one day?

JUST COMPARE TO THE ALGORITHM THEY PROPOSED

Rodriguez et al., CONCUR'15

- Looks good! Err... **not so simple**. Depending on how you do this, you quickly move from polynomial to exponential.
- **Experiments!** \triangle Possible 404, *code not found!* ahead!!!
- Wait, what's this language? Did this ever run one day?
- Wow, I'll need **CPLEX** and **Gurobi** but all I have is **lpsolve** or **glpk**

JUST COMPARE TO THE ALGORITHM THEY PROPOSED

Rodriguez et al., CONCUR'15

- Looks good! Err... **not so simple**. Depending on how you do this, you quickly move from polynomial to exponential.
- **Experiments!** \triangle Possible 404, *code not found!* ahead!!!
- Wait, what's this language? Did this ever run one day?
- Wow, I'll need **CPLEX** and **Gurobi** but all I have is **lpsolve** or **glpk**
- Sweet, they provided a **binary!** Oh, wait, MacOSX in 2015 ?!?

JUST COMPARE TO THE ALGORITHM THEY PROPOSED

Rodriguez et al., CONCUR'15

- Looks good! Err... **not so simple**. Depending on how you do this, you quickly move from polynomial to exponential.
- **Experiments!** \triangle Possible 404, *code not found!* ahead!!!
- Wait, what's this language? Did this ever run one day?
- Wow, I'll need **CPLEX** and **Gurobi** but all I have is **lpsolve** or **glpk**
- Sweet, they provided a **binary**! Oh, wait, MacOSX in 2015 ?!?
- The GitHub webpage says it requires Foo, Bar, and Baz, but none of the **versions** I find appear to work.

JUST COMPARE TO THE ALGORITHM THEY PROPOSED

Rodriguez et al., CONCUR'15

- Looks good! Err... **not so simple**. Depending on how you do this, you quickly move from polynomial to exponential.
- **Experiments!** Possible 404, *code not found!* ahead!!!
- Wait, what's this language? Did this ever run one day?
- Wow, I'll need **CPLEX** and **Gurobi** but all I have is **lpsolve** or **glpk**
- Sweet, they provided a **binary**! Oh, wait, MacOSX in 2015 ?!?
- The GitHub webpage says it requires Foo, Bar, and Baz, but none of the **versions** I find appear to work.
- With which **parameters** and data set do you run this code? And Why?

In the end, **one new thesis** to reproduce/understand this paper... and finally contribute.

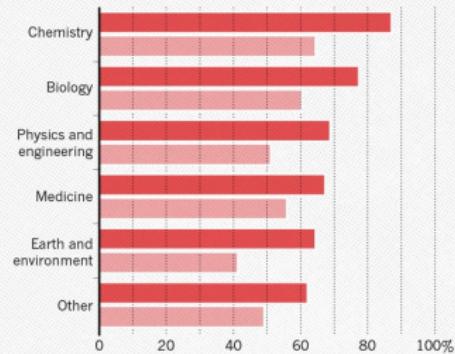
DIFFERENT KINDS OF REPRODUCIBILITY

SOCIO-TECHNICAL CHALLENGES

HAVE YOU FAILED TO REPRODUCE AN EXPERIMENT?

Most scientists have experienced failure to reproduce results.

● Someone else's ● My own



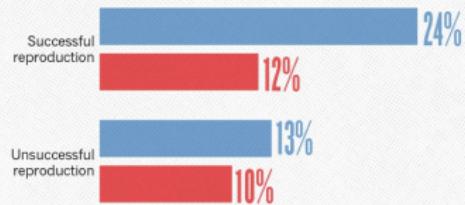
1,500 scientists lift the lid on reproducibility,

Nature, May 2016

HAVE YOU EVER TRIED TO PUBLISH A REPRODUCTION ATTEMPT?

Although only a small proportion of respondents tried to publish replication attempts, many had their papers accepted.

● Published ● Failed to publish



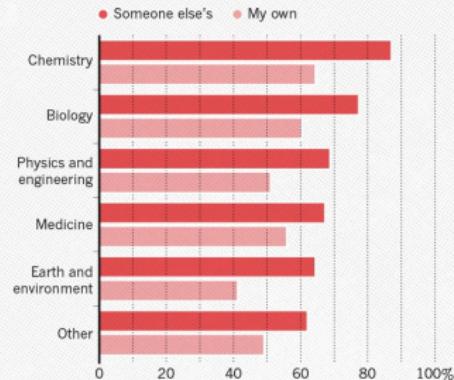
Number of respondents from each discipline:

Biology 703, Chemistry 106, Earth and environmental 95, Medicine 203, Physics and engineering 236, Other 233.

SOCIO-TECHNICAL CHALLENGES

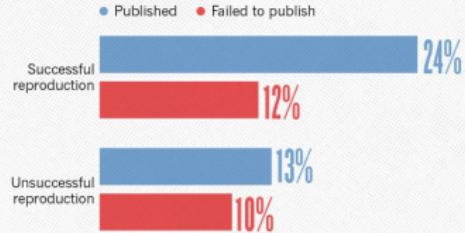
HAVE YOU FAILED TO REPRODUCE AN EXPERIMENT?

Most scientists have experienced failure to reproduce results.



HAVE YOU EVER TRIED TO PUBLISH A REPRODUCTION ATTEMPT?

Although only a small proportion of respondents tried to publish replication attempts, many had their papers accepted.



Number of respondents from each discipline:
Biology 703, Chemistry 106, Earth and environmental 95,
Medicine 203, Physics and engineering 236, Other 233.

1,500 scientists lift the lid on reproducibility,

Nature, May 2016

Social causes

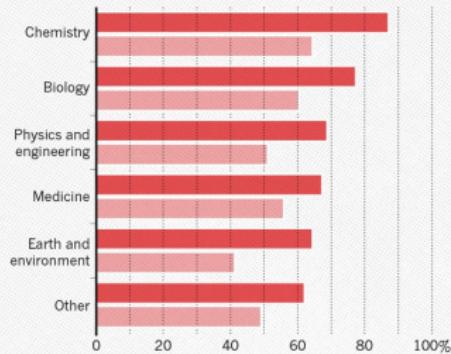
- Fraud, conflict of interest (pharmaceutic, ...)
- No incentive to reproduce/check our own work (afap), nor the work of others (big results!), nor to allow others to check (competition)
- Peer review does not scale: 1M+ articles per year!

SOCIO-TECHNICAL CHALLENGES

HAVE YOU FAILED TO REPRODUCE AN EXPERIMENT?

Most scientists have experienced failure to reproduce results.

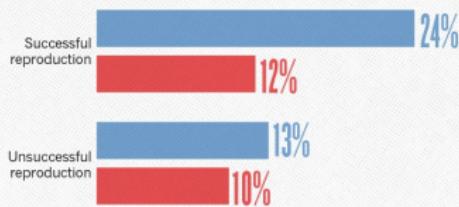
● Someone else's ● My own



HAVE YOU EVER TRIED TO PUBLISH A REPRODUCTION ATTEMPT?

Although only a small proportion of respondents tried to publish replication attempts, many had their papers accepted.

● Published ● Failed to publish



Number of respondents from each discipline:

Chemistry 106, Biology 106, Earth and environmental 95, Medicine 203, Physics and engineering 236, Other 233.

1,500 scientists lift the lid on reproducibility,

Nature, May 2016

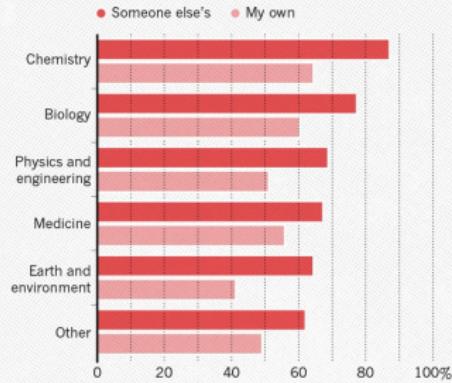
Social causes

- Fraud, conflict of interest (pharmaceutic, ...)
- No incentive to reproduce/check our own work (afap), nor the work of others (big results!), nor to allow others to check (competition)
- Peer review does not scale: 1M+ articles per year!
- Emerging practices: DORA/Plan S/COARA, DMP and FAIR data, artefact evaluation, reproducibility badges, reproducibility challenges, open reviews, ...

SOCIO-TECHNICAL CHALLENGES

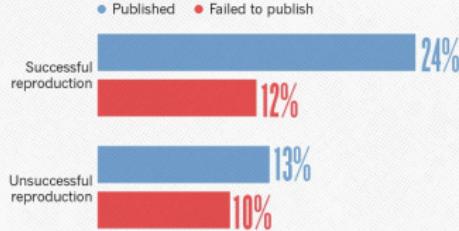
HAVE YOU FAILED TO REPRODUCE AN EXPERIMENT?

Most scientists have experienced failure to reproduce results.



HAVE YOU EVER TRIED TO PUBLISH A REPRODUCTION ATTEMPT?

Although only a small proportion of respondents tried to publish replication attempts, many had their papers accepted.



Number of respondents from each discipline:
Biology 703, Chemistry 106, Earth and environmental 95,
Medicine 203, Physics and engineering 236, Other 233.

1,500 scientists lift the lid on reproducibility,

Nature, May 2016

Social causes

- Fraud, conflict of interest (pharmaceutic, ...)
- No incentive to reproduce/check our own work (afap), nor the work of others (big results!), nor to allow others to check (competition)
- Peer review does not scale: 1M+ articles per year!
- Emerging practices: DORA/Plan S/COARA, DMP and FAIR data, artefact evaluation, reproducibility badges, reproducibility challenges, open reviews, ...

Methodological/technical causes

- The many biases (apophenia, confirmation, hindsight, experimenter, ...): bad designs
- Selective reporting, weak analysis (statistics, data manipulation mistakes, computational errors)
- Lack of information, code/raw data unavailable

DIFFERENT REPRODUCIBILITY CONCERN IN MODERN SCIENCE

Biology, Oncology sample provenance, clinical trials \rightsquigarrow standardized protocols

DIFFERENT REPRODUCIBILITY CONCERN IN MODERN SCIENCE

Biology, Oncology sample provenance, clinical trials \rightsquigarrow standardized protocols

Psychology, Nutrition HARKING, p-hacking \rightsquigarrow pre-registration

Brian Wansink Professor, Psychological Nutrition, Cornell, 2016

I gave her a data set of a self-funded, failed study which had null results. I said "This cost us a lot of time and our own money to collect. There's got to be something here we can salvage because it's a cool (rich & unique) data set." I told her what the analyses should be. [...] Every day she came back with puzzling new results, and every day we would scratch our heads, ask "Why," and come up with another way to reanalyze the data with yet another set of plausible hypotheses

17 retracted publications

DIFFERENT REPRODUCIBILITY CONCERNS IN MODERN SCIENCE

Biology, Oncology sample provenance, clinical trials \rightsquigarrow standardized protocols

Psychology, Nutrition HARKING, p-hacking \rightsquigarrow pre-registration

Genomics software engineering, computational reproducibility, provenance

Computational fluid dynamics numerical chaos, parallel architectures

DIFFERENT REPRODUCIBILITY CONCERN IN MODERN SCIENCE

Biology, Oncology sample provenance, clinical trials \rightsquigarrow standardized protocols

Psychology, Nutrition HARKING, p-hacking \rightsquigarrow pre-registration

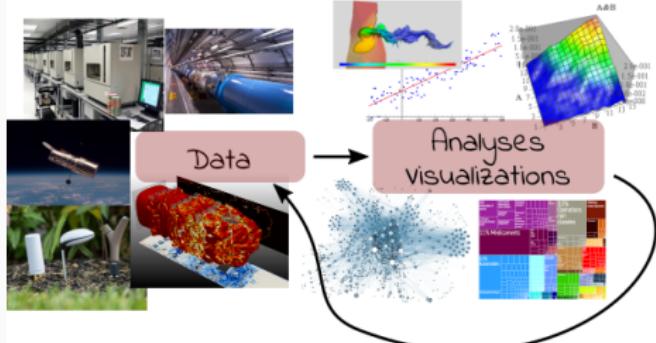
Genomics software engineering, computational reproducibility, provenance

Computational fluid dynamics numerical chaos, parallel architectures

Artificial Intelligence most of the above 😊

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

Authors



DIFFERENT REPRODUCIBILITY CONCERN IN MODERN SCIENCE

Biology, Oncology sample provenance, clinical trials \rightsquigarrow standardized protocols

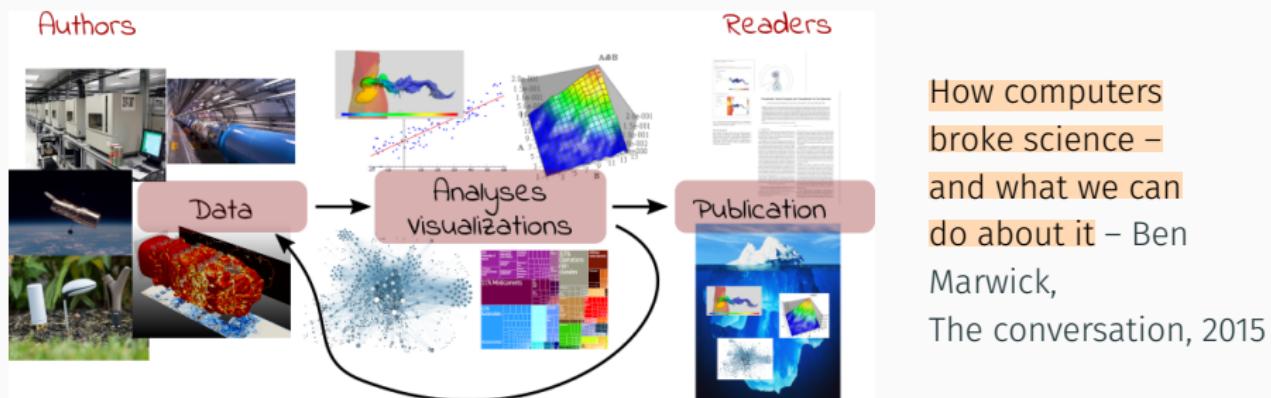
Psychology, Nutrition HARKING, p-hacking \rightsquigarrow pre-registration

Genomics software engineering, computational reproducibility, provenance

Computational fluid dynamics numerical chaos, parallel architectures

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

GOOD PRACTICE #1: TAKING NOTES AND DOCUMENTING

What your research supposedly looks like:

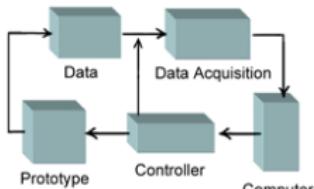


Figure 1. Experimental Diagram

What your research *actually* looks like:

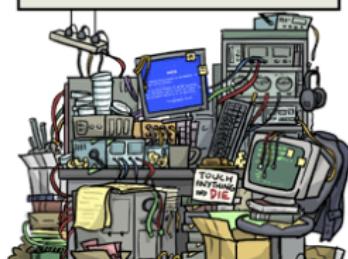


Figure 2. Experimental Mess

MAIN CHALLENGES FOR A COMPUTATIONAL SCIENTIST

```
1 my_code --cfg=magical_param:0.94572 '*.dat' --output foo.csv
```

Tracking code version

- `my_code` is revision 21b95ecfa0911d6ca87668482b11ab9498edd8f3

MAIN CHALLENGES FOR A COMPUTATIONAL SCIENTIST

```
1 my_code --cfg=magical_param:0.94572 '*.dat' --output foo.csv
```

Tracking code version

- `my_code` is revision `21b95ecfa0911d6ca87668482b11ab9498edd8f3`

Tracking software environment

- `my_code` depends on a dozen of libraries, which depend on dozens of libraries
- `my_code` was compiled with `clang 1:9.0-49.1` and
`-O3 -funroll-loops -fno-strict-aliasing -finline-functions ...`

MAIN CHALLENGES FOR A COMPUTATIONAL SCIENTIST

```
1 my_code --cfg=magical_param:0.94572 '*.dat' --output foo.csv
```

Tracking code version

- `my_code` is revision `21b95ecfa0911d6ca87668482b11ab9498edd8f3`

Tracking software environment

- `my_code` depends on a dozen of libraries, which depend on dozens of libraries
- `my_code` was compiled with `clang 1:9.0-49.1` and
`-O3 -funroll-loops -fno-strict-aliasing -finline-functions ...`

Tracking parameters and data

- `*.dat`? Ooh, you ran this in `data/2091293-AJXQ37`?
- Wasn't `mymap.dat` updated since then?
- That was for `foo.csv`. What about `bar.csv`? Is it reproducible?

MAIN CHALLENGES FOR A COMPUTATIONAL SCIENTIST

```
1 my_code --cfg=magical_param:0.94572 '*.dat' --output foo.csv
```

Tracking code version

- `my_code` is revision `21b95ecfa0911d6ca87668482b11ab9498edd8f3`

Tracking software environment

- `my_code` depends on a dozen of libraries, which depend on dozens of libraries
- `my_code` was compiled with `clang 1:9.0-49.1` and
`-O3 -funroll-loops -fno-strict-aliasing -finline-functions ...`

Tracking parameters and data

- `*.dat`? Ooh, you ran this in `data/2091293-AJXQ37`?
- Wasn't `mymap.dat` updated since then?
- That was for `foo.csv`. What about `bar.csv`? Is it reproducible?

Tracking the process (on short/long term)

- Why did I run this? What did I learn from it? I remember doing this but when?

MAIN CHALLENGES FOR A COMPUTATIONAL SCIENTIST

```
1 my_code --cfg=magical_param:0.94572 '*.dat' --output foo.csv
```

Tracking code version

- `my_code` is revision `21b95ecfa0911d6ca87668482b11ab9498edd8f3`

Tracking software environment

- `my_code` depends on a dozen of libraries, which depend on dozens of libraries
- `my_code` was compiled with `clang 1:9.0-49.1` and
`-O3 -funroll-loops -fno-strict-aliasing -finline-functions ...`

Tracking parameters and data

- `*.dat`? Ooh, you ran this in `data/2091293-AJXQ37`?
- Wasn't `mymap.dat` updated since then?
- That was for `foo.csv`. What about `bar.csv`? Is it reproducible?

Tracking the process (on short/long term)

- Why did I run this? What did I learn from it? I remember doing this but when?

Handle complex sequences and reuse results (leverage supercomputers)

TOOL 1: COMPUTATIONAL NOTEBOOKS (LITTERATE PROGRAMMING)

Un document computationnel

Mon ordinateur m'indique que π vaut approximativement

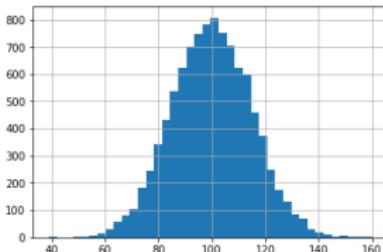
3.141592653589793

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

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 π (si ce n'est une constante de normalisation... ☺).



TOOL 1: COMPUTATIONAL NOTEBOOKS (LITTERATE PROGRAMMING)

Document initial dans son environnement

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

- In [1]:**

```
from math import *
print(pi)
3.141592653589793
```

Mais calculé avec la [méthode des aiguilles de Buffon](#) (https://fr.wikipedia.org/wiki/Aiguille_de_Buffon), on obtient aussi 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]: 3.14371986944998765

On peut inclure des formules mathématiques comme $\sqrt{2/\pi} \exp(-x^2/2)$ et des dessins qui n'ont rien à voir avec π (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(10000)
plt.hist(x,40)
plt.grid(True)
plt.show()
```

Document final

Un document computationnel

Mon ordinateur m'indique que π vaut approximativement

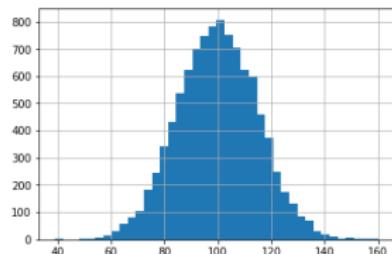
3.141592653589793

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
```

3.14371986944998765

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 π (si ce n'est une constante de normalisation...).



TOOL 1: COMPUTATIONAL NOTEBOOKS (LITTERATE PROGRAMMING)

Document initial dans son environnement

The screenshot shows a Jupyter Notebook interface. At the top, there's a toolbar with File, Edit, View, Insert, Cell, Kernel, Widgets, Help, Hide Code, and Python 3. Below the toolbar, there's a header bar with tabs for example_nb (selected), Cell, Kernel, Widgets, Help, Hide Code, and Python 3. The main area contains a cell titled '# Un document computationnel' with the following text:
Mon ordinateur m'indique que π vaut "approximativement"
In [1]:

```
from math import *
print(pi)
3.141592653589793
```


Out[1]:
Mais calculé avec la [méthode des aiguilles de Buffon](#) (https://fr.wikipedia.org/wiki/Aiguille_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]:
3.14371986944998765
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 π (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(10000)
plt.hist(x,40)
plt.grid(True)
plt.show()
```


A histogram is displayed below the code, showing a bell-shaped curve centered at approximately 100.

Document final

Un document computationnel

Mon ordinateur m'indique que π vaut approximativement

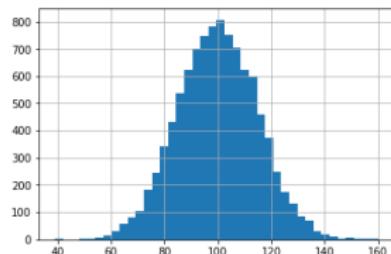
3.141592653589793

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
```

3.14371986944998765

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 π (si ce n'est une constante de normalisation...).



TOOL 1: COMPUTATIONAL NOTEBOOKS (LITTERATE PROGRAMMING)

Document initial dans son environnement

The screenshot shows a Jupyter Notebook interface with three 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,60)
plt.grid(True)
plt.show()` → Output: A histogram showing a normal distribution centered at 100.

Document final

Un document computationnel

Mon ordinateur m'indique que π vaut approximativement

3.141592653589793

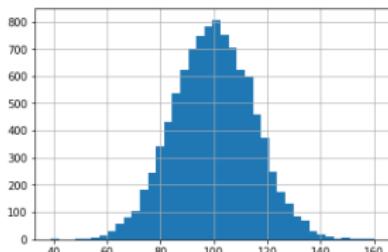
Mais calculé avec la méthode des [aiguilles de Buffon](#), on obtiendrait comme approximation :

```
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

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 π (si ce n'est une constante de normalisation... ☺).



TOOL 1: COMPUTATIONAL NOTEBOOKS (LITTERATE PROGRAMMING)

Document initial dans son environnement

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

- In [1]:** Prints the value of pi (3.141592653589793) and includes a note about calculating pi with Buffon's needle method.
- In [2]:** Generates random points (x, y) and calculates the ratio of points below the unit circle to total points, which approximates pi.
- In [3]:** Plots a histogram of x values from -100 to 100, showing a symmetric bell-shaped distribution centered around 0.

Annotations in red highlight the output of In [1] and In [2], and a large red arrow points from the text "Résultats" to the histogram in In [3].

Document final

Un document computationnel

Mon ordinateur m'indique que π vaut approximativement

3.141592653589793

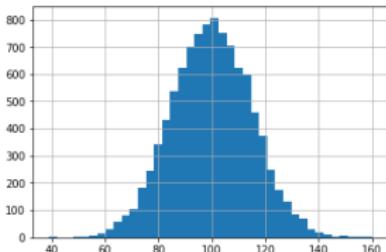
Mais calculé avec la méthode des [aiguilles de Buffon](#), on obtiendrait comme approximation :

```
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
```

3.1437198694908765

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 π (si ce n'est une constante de normalisation... ☺).



TOOL 1: COMPUTATIONAL NOTEBOOKS (LITTERATE PROGRAMMING)

Document initial dans son environnement

Un document computationnel

```
In [1]:  
from math import *  
print(pi)  
3,141592653589793
```

Mais calculé avec la `_methodes_ des éimpulles de Buffon` (https://fr.wikipedia.org/wiki/Algille_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]: 3,1437198694098765
```

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`, avec `matplotlib` (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(10000)  
  
plt.hist(x, 99)  
plt.grid(True)  
plt.show()
```

A histogram plot showing a normal distribution. The x-axis ranges from 40 to 160 with major ticks every 20 units. The y-axis ranges from 0 to 800 with major ticks every 100 units. The distribution is symmetric and centered around 100, with the highest frequency occurring at 100.

Document final

Un document computationnel

Mon ordinateur m'indique que π vaut approximativement

3.141592653589793

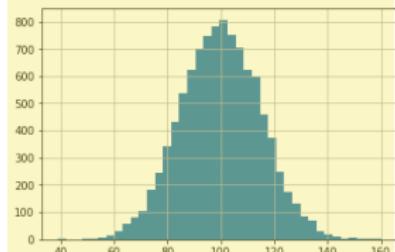
Mais calculé avec la **méthode des aiguilles de Buffon**, on obtiendrait comme approximation :

```
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 π (si ce n'est une constante de normalisation... ☺).



TOOL 1: COMPUTATIONAL NOTEBOOKS (LITTERATE PROGRAMMING)

Document initial dans son environnement

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

- In [1]:** Prints the value of pi.
- In [2]:** Calculates the ratio of points inside a unit square with sine values greater than zero to the total number of points, which approximates pi/2.
- In [3]:** Generates a histogram of 100,000 random numbers between 0 and 1, showing a bell-shaped distribution centered at 0.5.

A red arrow labeled "Export" points from the bottom left towards the final document.

Document final

Un document computationnel

Mon ordinateur m'indique que π vaut approximativement

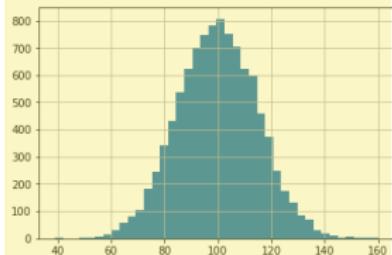
3.141592653589793

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
```

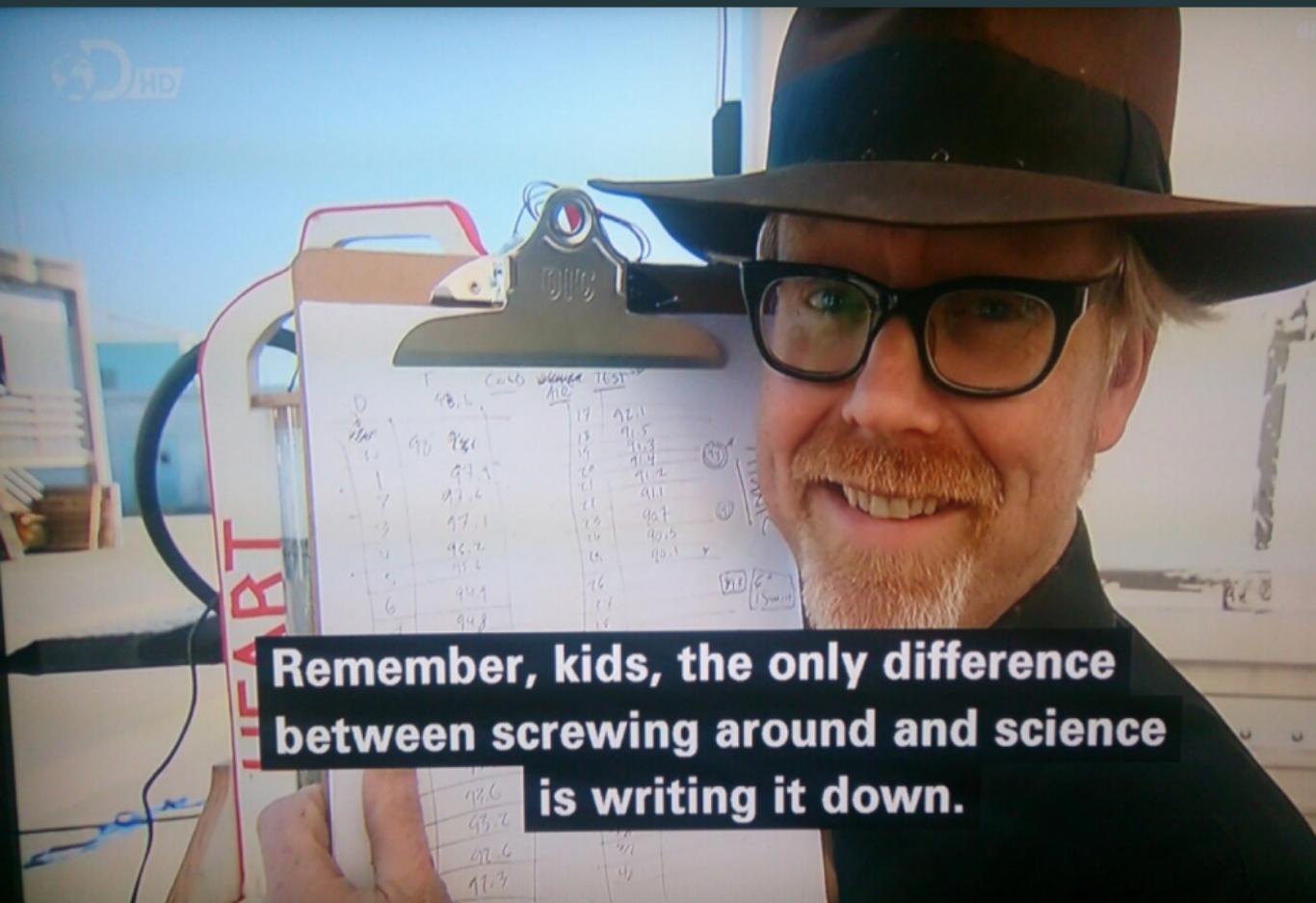
3.1437198694098765

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 π (si ce n'est une constante de normalisation... ☺).



allow to write a reproducible document that links raw data with final figures

TOOL 2: ELECTRONIC NOTEBOOKS



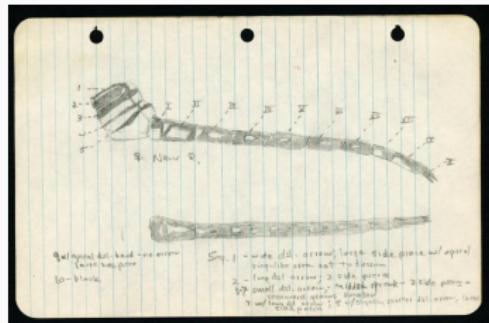
**Remember, kids, the only difference
between screwing around and science
is writing it down.**

FIELD AND LABORATORY NOTEBOOKS

Social Sciences, Ecology, Biology



Chemistry, Physics, Biology



Robert Henry Gibbs, Jr.,
ichthyologist (1929 – 1988)



Emil Heinrich du Bois-Reymond,
electrophysiologist (1818 – 1896)

FIELD AND LABORATORY NOTEBOOKS

Social Sciences, Ecology, Biology



| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 | 32 | 33 | 34 | 35 | 36 | 37 | 38 | 39 | 40 | 41 | 42 | 43 | 44 | 45 | 46 | 47 | 48 | 49 | 50 | 51 | 52 | 53 | 54 | 55 | 56 | 57 | 58 | 59 | 60 | 61 | 62 | 63 | 64 | 65 | 66 | 67 | 68 | 69 | 70 | 71 | 72 | 73 | 74 | 75 | 76 | 77 | 78 | 79 | 80 | 81 | 82 | 83 | 84 | 85 | 86 | 87 | 88 | 89 | 90 | 91 | 92 | 93 | 94 | 95 | 96 | 97 | 98 | 99 | 100 |
|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|-----|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 | 32 | 33 | 34 | 35 | 36 | 37 | 38 | 39 | 40 | 41 | 42 | 43 | 44 | 45 | 46 | 47 | 48 | 49 | 50 | 51 | 52 | 53 | 54 | 55 | 56 | 57 | 58 | 59 | 60 | 61 | 62 | 63 | 64 | 65 | 66 | 67 | 68 | 69 | 70 | 71 | 72 | 73 | 74 | 75 | 76 | 77 | 78 | 79 | 80 | 81 | 82 | 83 | 84 | 85 | 86 | 87 | 88 | 89 | 90 | 91 | 92 | 93 | 94 | 95 | 96 | 97 | 98 | 99 | 100 |

Marie-Claude Quidoz, Centre d'Ecologie
Fonctionnelle et Evolutive
Essentially no evolution throughout the last century

Chemistry, Physics, Biology



Contemporary laboratory notebook
in Neurology

TAKING NOTES

Org-mode and Markdown two simple text formats

- simple formatting and export to more elaborate formats
 - Hyperlinks, images, code, etc.

Journal structure

- My journal/notebook (single org file)
 - Tom Cornebize's journal (single org file + Jupyter notebooks)

Zettelkasten possible structure

- Org-roam, Zettler, Roam, Obsidian, ...
 - Architects, librarians, **gardeners**



TAKING NOTES

Org-mode and Markdown two simple text formats

- simple formatting and export to more elaborate formats
 - Hyperlinks, images, code, etc.

Journal structure

- My journal/notebook (single org file)
 - Tom Cornebize's journal (single org file + Jupyter notebooks)

Zettelkasten possible structure

- Org-roam, Zettler, Roam, Obsidian, ...
 - Architects, librarians, **gardeners**



My recommendations: Do not use a fancy cloud-based proprietary tool

- Simple open source **text-based** format
 - **Control versions** and backup yourself (e.g., using gitlab, github)
 - Single location if possible
 - Annotate (tags in a journal, links in a Zettelkasten)

GOOD PRACTICE #2

CONTROLLING SOFTWARE ENVIRONMENT

NIGHTMARE 1: FIGHTING SOFTWARE ENVIRONMENTS NIGHTMARE

What is hiding behind a simple

```
1 import matplotlib
```

Package: python3-matplotlib
Version: 2.1.1-2
Depends: python3-dateutil, python-matplotlib-data (>= 2.1.1-2),
python3-pyparsing (>= 1.5.6), python3-six (>= 1.10), python3-tz,
libjs-jquery, libjs-jquery-ui, python3-numpy (>= 1:1.13.1),
python3-numpy-abi9, python3 (<< 3.7), python3 (>= 3.6~),
python3-cycler (>= 0.10.0), python3:any (>= 3.3.2-2~), libc6 (>=
2.14), libfreetype6 (>= 2.2.1), libgcc1 (>= 1:3.0), libpng16-16 (>=
1.6.2-1), libstdc++6 (>= 5.2), zlib1g (>= 1:1.1.4)

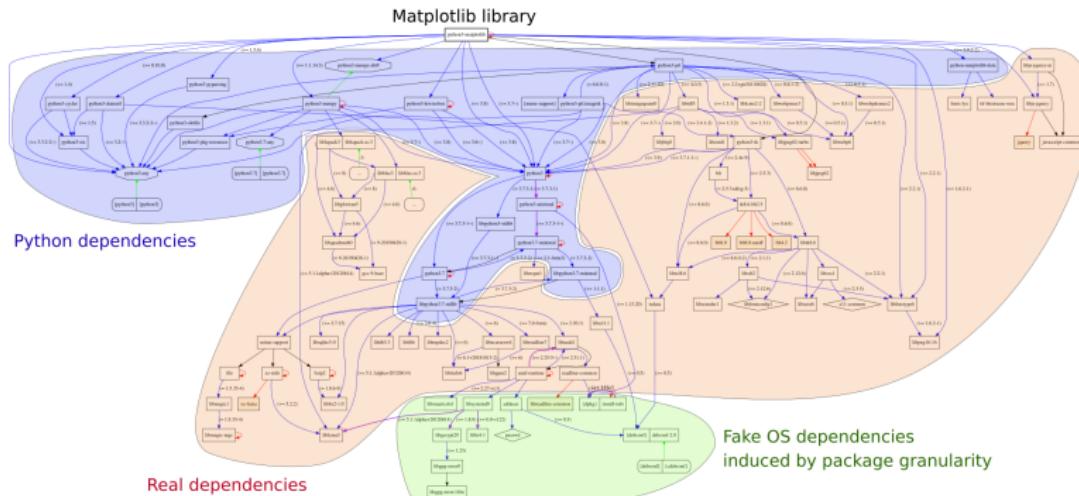
NIGHTMARE 1: FIGHTING SOFTWARE ENVIRONMENTS NIGHTMARE

What is hiding behind a simple

1

```
import matplotlib
```

Package: python3-matplotlib



NIGHTMARE 1: FIGHTING SOFTWARE ENVIRONMENTS NIGHTMARE

Python and its rapidly evolving environment

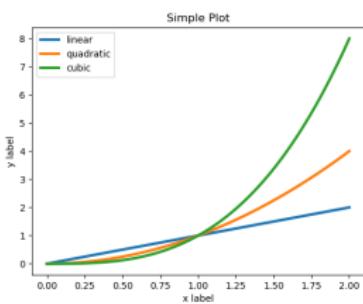
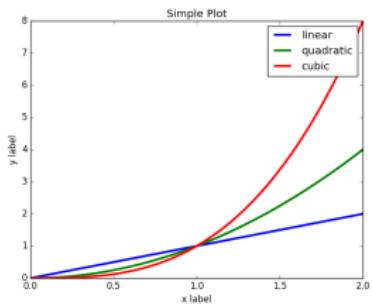
```
1 python2 -c "print(10/3)"  
2 python3 -c "print(10/3)"
```

```
3  
3.333333333333335
```

NIGHTMARE 1: FIGHTING SOFTWARE ENVIRONMENTS NIGHTMARE

Python and its rapidly evolving environment

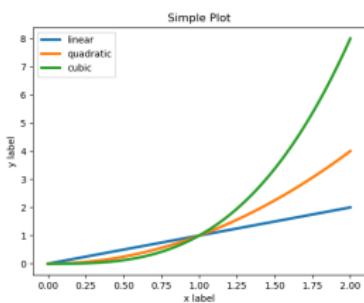
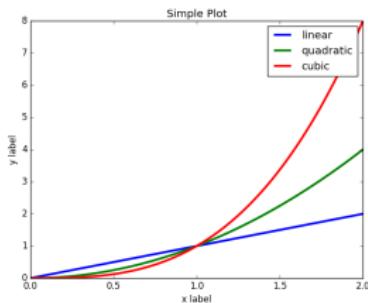
```
1 python2 -c "print(10/3)"  
2 python3 -c "print(10/3)"
```



NIGHTMARE 1: FIGHTING SOFTWARE ENVIRONMENTS NIGHTMARE

Python and its rapidly evolving environment

```
1 python2 -c "print(10/3)"  
2 python3 -c "print(10/3)"
```



Cortical Thickness Measurements (PLOS ONE, June 2012) in FreeSurfer
About a *factor two smaller differences* were found *between*
the Mac and HP workstations and *between Mac OSX 10.5 and*
OSX 10.6.

TOOL 3: CONTAINERS AND PACKAGE MANAGERS

The good



The bad



The ugly



Automatic tracking

TOOL 3: 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 <cmd>`
- Building images: `docker build -f <Dockerfile>`
- Sharing through the Docker Hub: `docker pull/push `

TOOL 3: 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 *reproducibility 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)

TOOL 3: CONTAINERS AND PACKAGE MANAGERS

The good



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)

The bad



The ugly

OTHER COMPUTER RELATED NIGHTMARES

NIGHTMARE 2: FIGHTING INFORMATION LOSS WITH ARCHIVES

D. Spinellis. *The Decay and Failures of URL References*. CACM, 46(1), Jan 2003.

The half-life of a referenced URL is approximately 4 years from its publication date.

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

NIGHTMARE 2: FIGHTING INFORMATION LOSS WITH ARCHIVES

D. Spinellis. *The Decay and Failures of URL References*. CACM, 46(1), Jan 2003.

The half-life of a referenced URL is approximately 4 years from its publication date.

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

Article archives arXiv.org HAL
archives-ouvertes.fr

Data archives figshare zenodo

Software Archive Software Heritage

 or  = awesome but:

- Inadequate for **large** data
- ≠ Archive

TOOL 5: AUTOMATING COMPUTATIONS WITH WORKFLOWS

Notebooks are no panacea and do not help developing clean code

The screenshot shows a Jupyter Notebook interface with several code and output cells. The first cell (In [1]) contains a single line of Python code: `from math import *; print(pi)`. The output (Out[1]) is the value `3.141592653589793`. A note in the cell states: "Mon ordinateur m'indique que `pi` vaut *approximativement*". The second cell (In [2]) uses NumPy to generate a random sample from a uniform distribution between 0 and pi/2, and calculates the mean. The output (Out[2]) is `0.143719686499785`. A note in the cell states: "On peut inclure des formules mathématiques comme `\frac{1}{\sqrt{\pi}}` via `\frac{1}{\sqrt{\pi}}`". The third cell (In [3]) contains a line of Matplotlib code to plot a histogram. The output (Out[3]) is a histogram showing a bell-shaped curve centered around 0.5.

```
# Un document computationnel
# Hide Prompt □ Hide Code □ Hide Outputs □
Mon ordinateur m'indique que pi vaut *approximativement*
In [1]: Hide Prompt □ Hide Code □ Hide Outputs □
from math import *
print(pi)
3.141592653589793

Mais calculé avec la __méthode__ des [assemblés de Buffets]
https://christian.sirang.org/wiki/Assemblee_de_Buffets, on obtiendrait comme
approximation : □ Hide Prompt □ Hide Code □ Hide Outputs □

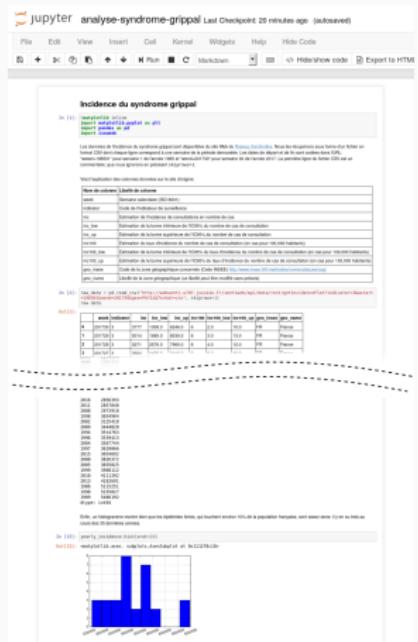
In [2]: Hide Prompt □ Hide Code □ Hide Outputs □
import numpy as np
n = 1000000
x = np.random.uniform(0, low=0, high=pi)
theta = np.random.uniform(0, low=0, high=pi/2)
l = len([x for x in theta if x >= pi/2])
Out[2]: 0.143719686499785

On peut inclure des formules mathématiques comme \frac{1}{\sqrt{\pi}} via \frac{1}{\sqrt{\pi}}
ce qui peut être intéressant pour les présentations. Notez que l'ordre des termes dans une
division qui n'est rien à voir avec latex (si ce n'est une constante de
normalisation...). □ Hide Prompt □ Hide Code □ Hide Outputs □

In [3]: Hide Prompt □ Hide Code □ Hide Outputs □
%matplotlib inline
import matplotlib.pyplot as plt
n = 1000000
x = np.random.uniform(0, high=pi)
plt.hist(x, 40)
plt.title("Histogramme")
plt.show()
```

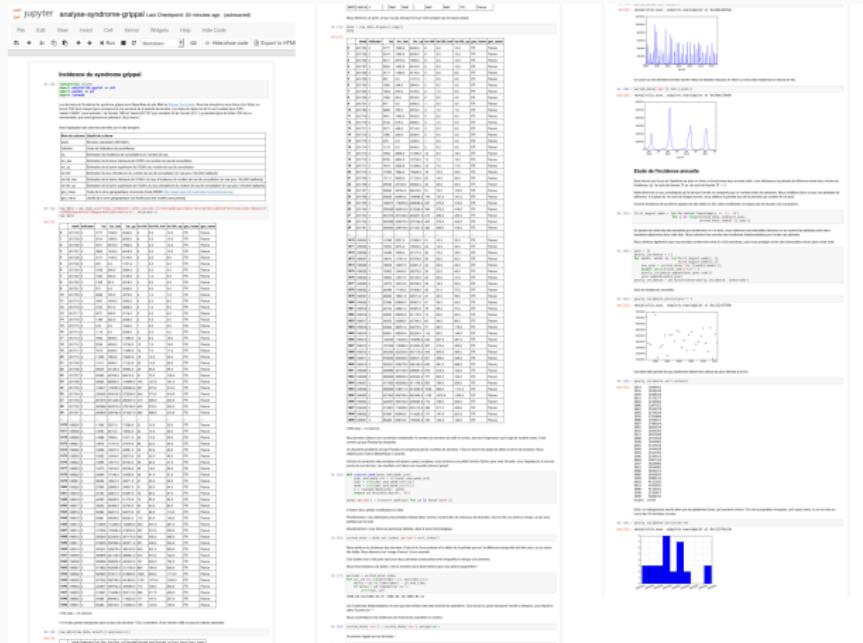
TOOL 5: AUTOMATING COMPUTATIONS WITH WORKFLOWS

Notebooks are no panacea and do not help developing clean code



TOOL 5: AUTOMATING COMPUTATIONS WITH WORKFLOWS

Notebooks are no panacea and do not help developing clean code



TOOL 5: AUTOMATING COMPUTATIONS WITH WORKFLOWS

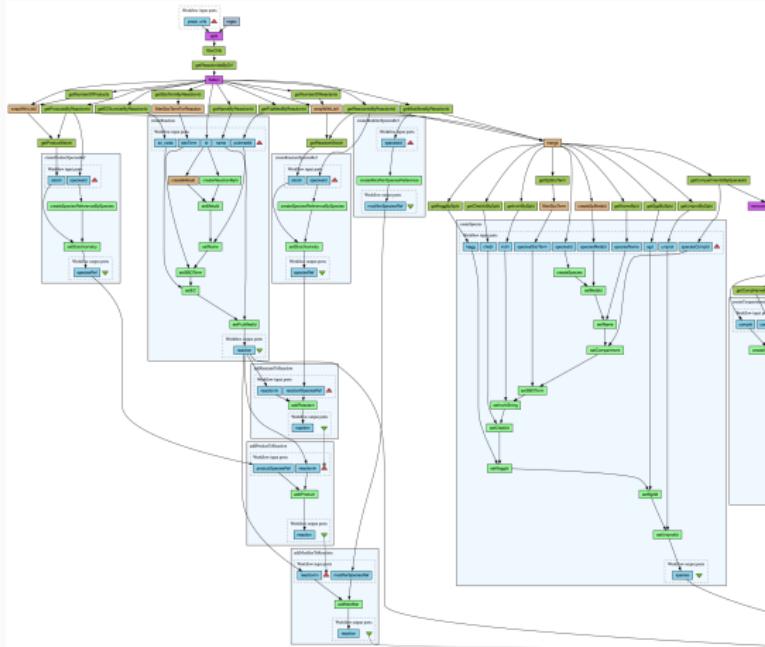
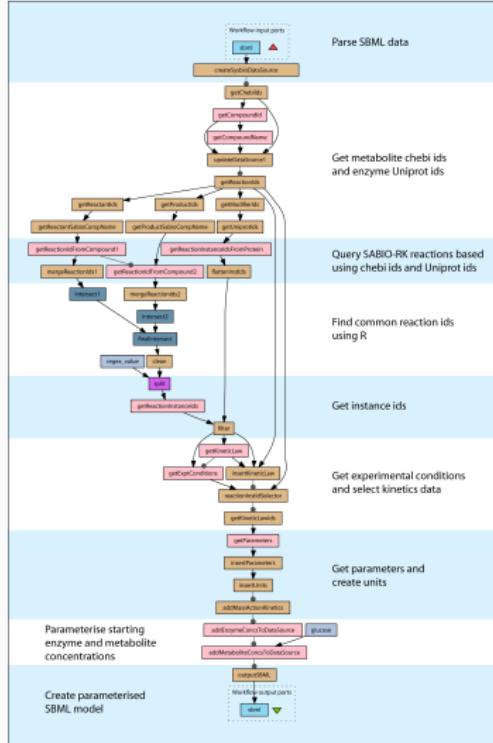
Notebooks are no panacea and do not help developing clean code

The collage consists of ten screenshots from Jupyter Notebooks, each showing a different step in a data science process:

- Screenshot 1: Extracting Color Names by Web Image Search.** Shows code for extracting color names from a web image using the `colornames` library.
- Screenshot 2: Generating a color palette.** Shows code for generating a color palette based on a color name and a number of colors.
- Screenshot 3: Generating a color palette by color name.** Shows code for generating a color palette based on a color name and a number of colors.
- Screenshot 4: Generating a color palette by color name.** Shows code for generating a color palette based on a color name and a number of colors.
- Screenshot 5: Generating a color palette by color name.** Shows code for generating a color palette based on a color name and a number of colors.
- Screenshot 6: Generating a color palette by color name.** Shows code for generating a color palette based on a color name and a number of colors.
- Screenshot 7: Generating a color palette by color name.** Shows code for generating a color palette based on a color name and a number of colors.
- Screenshot 8: Generating a color palette by color name.** Shows code for generating a color palette based on a color name and a number of colors.
- Screenshot 9: Generating a color palette by color name.** Shows code for generating a color palette based on a color name and a number of colors.
- Screenshot 10: Generating a color palette by color name.** Shows code for generating a color palette based on a color name and a number of colors.

Each screenshot includes a snippet of Python code, a corresponding visualization (such as a heatmap or scatter plot), and a detailed explanatory text block below it. The visualizations include a color palette, a heatmap of color distribution, and a scatter plot of training data.

TOOL 5: AUTOMATING COMPUTATIONS WITH WORKFLOWS



TOOL 5: AUTOMATING COMPUTATIONS WITH WORKFLOWS

Workflows:

- Clearer high-level view
- **Explicit** composition of codes and data movement
- Safer sharing, reusing, and execution
- Notebooks are a variant that is both impoverished and richer
 - No simple/mature path from a notebook to a workflow

Examples:

- Galaxy, Kepler, Taverna, Pegasus, Collective Knowledge, VisTrails
- Light-weight: `make`, dask, drake, swift, `snakemake`, ...
- Hybrids: SOS-notebook, ...

FLOATING POINTS ROUNDING: THE OTHER ROOT OF ALL EVIL ?

- Every operation includes implicit rounding.

```
1 print(2.1-2.0 == 0.1)
```

```
False
```

FLOATING POINTS ROUNDING: THE OTHER ROOT OF ALL EVIL ?

- Every operation includes implicit rounding.

```
1 print(2.1-2.0 == 0.1)
```

False

- Unfortunately: `round(round(a+b)+c) ≠ round(a+round(b+c))`

Hence, operation order matters. For a reproducible computation,

operation order should be preserved!!! Which order is more relevant is an other debate 😊

FLOATING POINTS ROUNDING: THE OTHER ROOT OF ALL EVIL ?

- Every operation includes implicit rounding.

```
1 print(2.1-2.0 == 0.1)
```

```
False
```

- Unfortunately: `round(round(a+b)+c) ≠ round(a+round(b+c))`
Hence, operation order matters. For a reproducible computation,
operation order should be preserved!!! Which order is more relevant is
an other debate 😊
- Numerical **instability** may be closer than you think [Rump, 1988]

$$f(x,y) = 333.75y^6 + x^2(11x^2y^2 - y^6 - 121y^4 - 2)2 + 5.5y^8 + \frac{x}{2y}$$

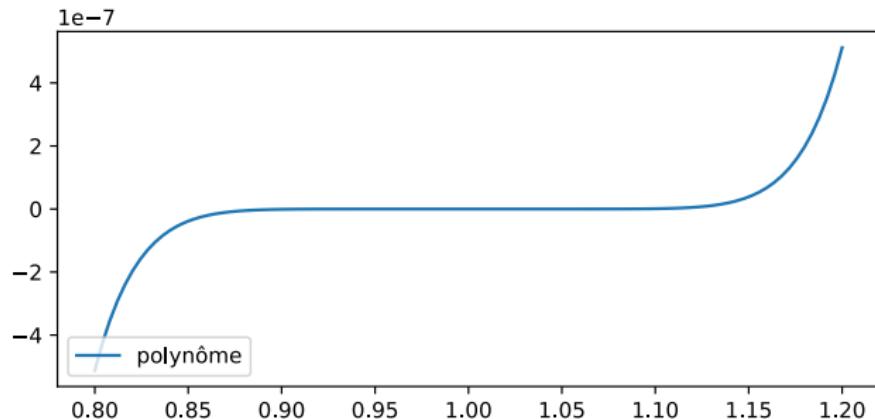
Evaluation of $f(77617.0, 33096.0)$

| | |
|------------------------------------------|----------------------------------|
| Single precision | 1.172603 |
| Double precision | 1.1726039400531 |
| Extended precision | 1.172603940053178 |
| MPFI | [-0.827396059946821368141165...] |
| (multiple precision interval arithmetic) | -0.827396059946821368141165...] |

Courtesy of Christophe Denis

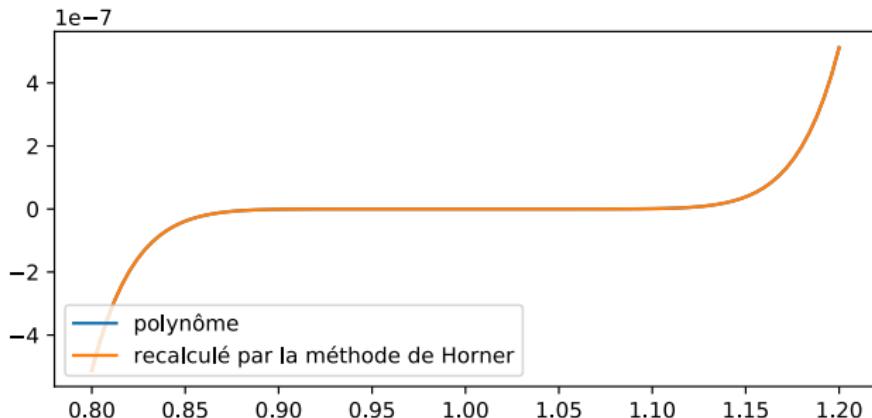
At scale (complex code + non-determinism), all this can become particularly harmful and painful.

ALL I CARE ABOUT IS THE ALGORITHM OUTPUT (FP)



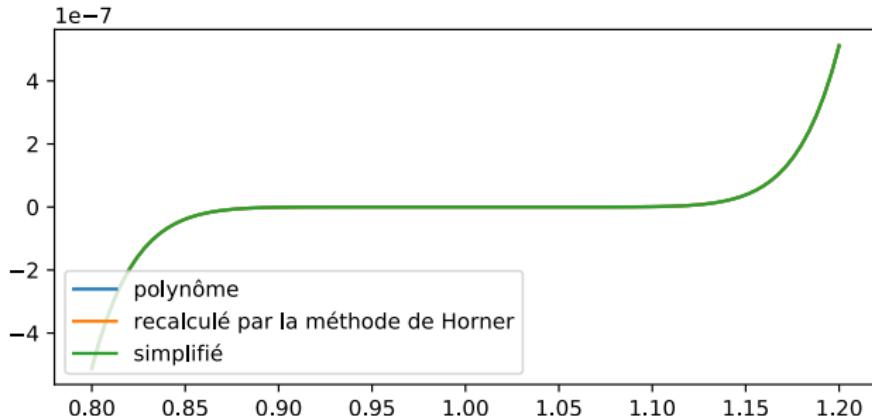
```
1 def polynome(x):  
2     return x**9 - 9.*x**8 + 36.*x**7 - 84.*x**6 + 126.*x**5 \  
3         - 126.*x**4 + 84.*x**3 - 36.*x**2 + 9.*x - 1.
```

FLOATING-POINT ARITHMETIC



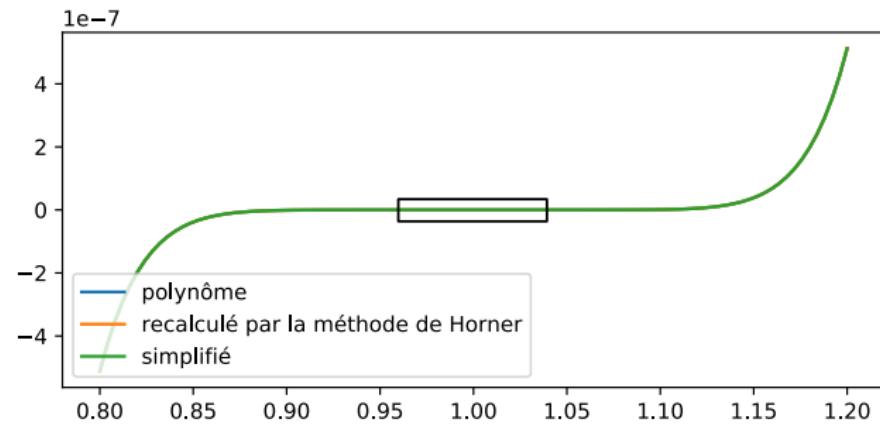
```
1 def horner(x):  
2     return x*(x*(x*(x*(x*(x*(x - 9.) + 36.) - 84.) + 126.) \  
3             - 126.) + 84.) - 36.) + 9.) - 1.
```

FLOATING-POINT ARITHMETIC

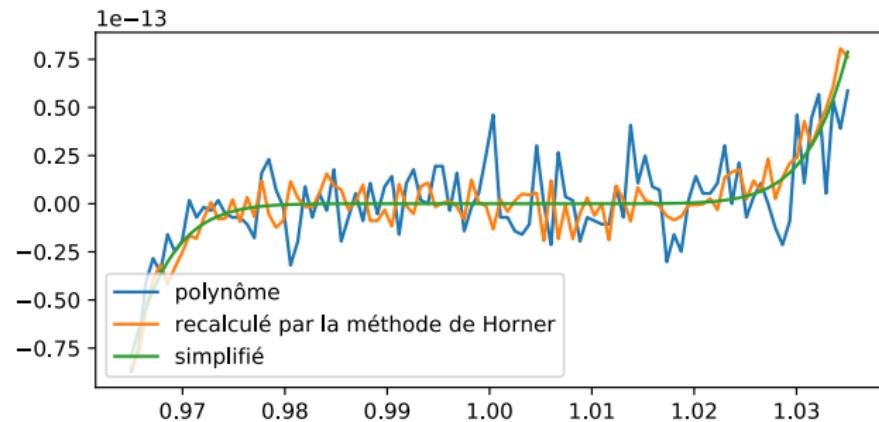


```
1 def simple(x):  
2     return (x-1.)**9  
3 # Easy! ;)
```

FLOATING-POINT ARITHMETIC



FLOATING-POINT ARITHMETIC

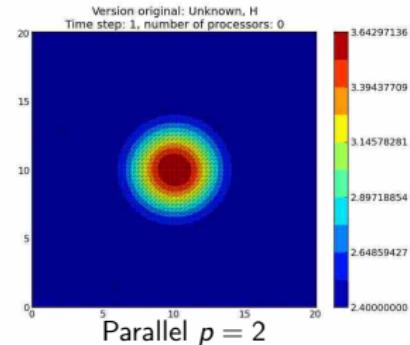
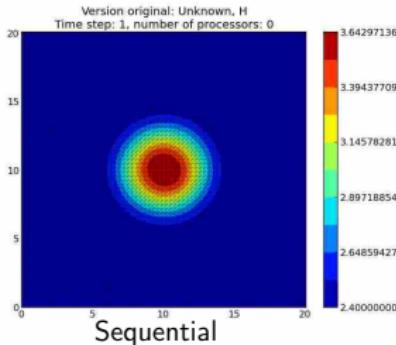


Telemac2D: the simplest gouttedo simulation

The gouttedo test case

- 2D-simulation of a water drop fall in a square bassin
- Unknown: water depth for a 0.2 sec time step
- Triangular mesh: 8978 elements and 4624 nodes

Expected numerical reproducibility (time step = 1, 2, ...)



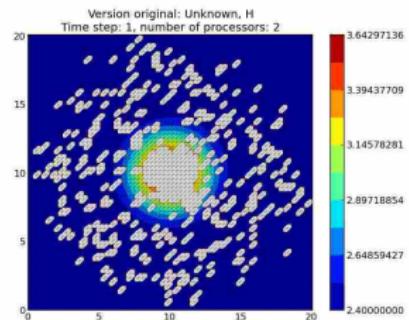
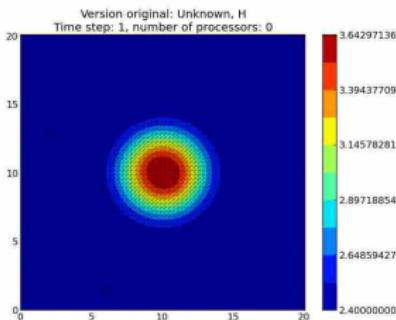
Courtesy of P. Langlois and R. Nheili
13 / 54

DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 1



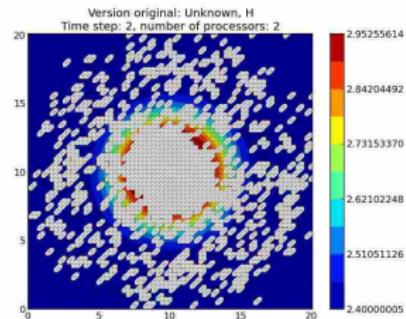
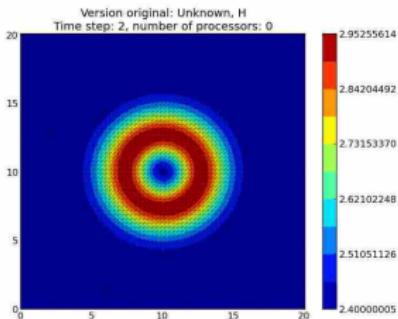
Courtesy of P. Langlois and R. Nheili
14 / 54

DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 2



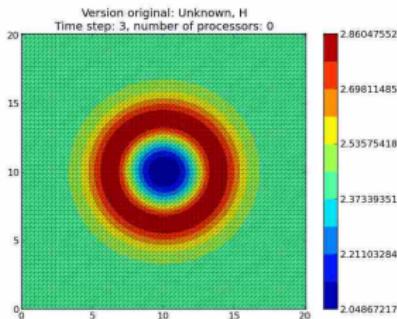
Courtesy of P. Langlois and R. Nheili
14 / 54

DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

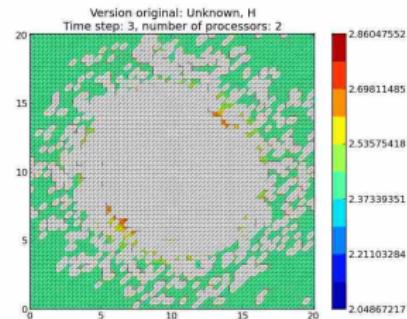
A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 3



Sequential



Parallel $p = 2$

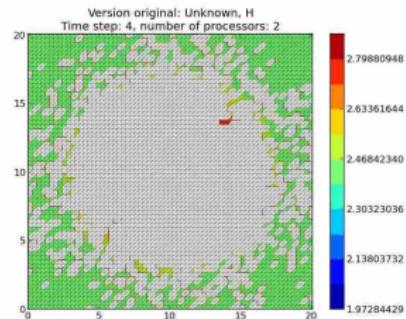
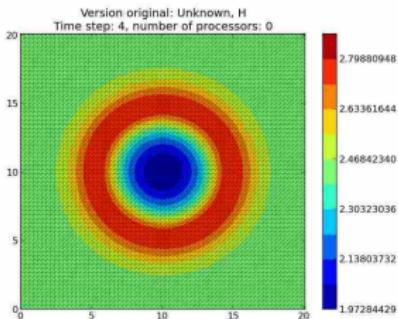
Courtesy of P. Langlois and R. Nheili
14 / 54

DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 4



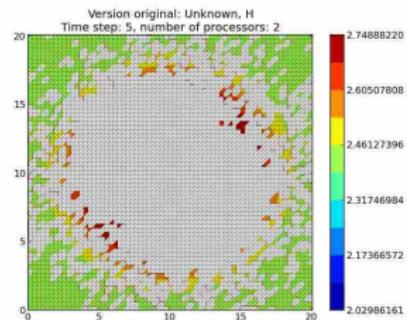
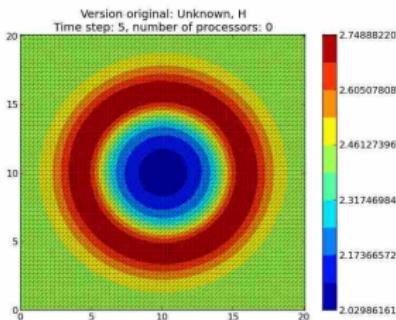
Courtesy of P. Langlois and R. Nheili
14 / 54

DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 5



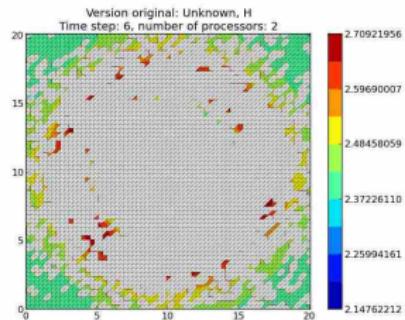
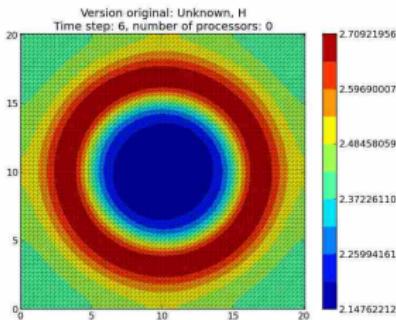
Courtesy of P. Langlois and R. Nheili
14 / 54

DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 6



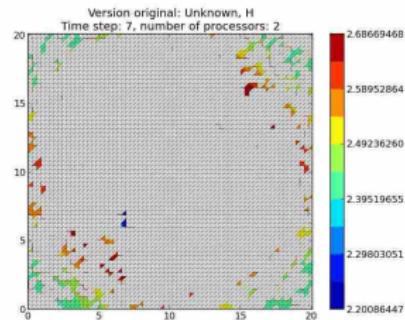
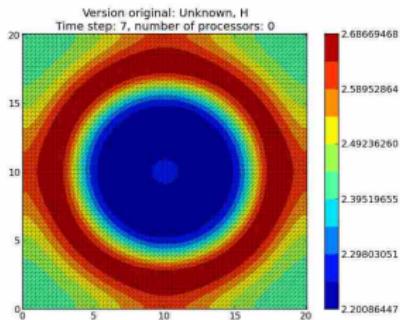
Courtesy of P. Langlois and R. Nheili
14 / 54

DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 7



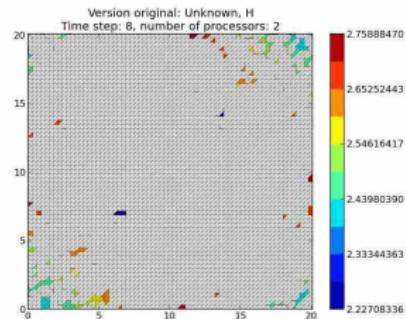
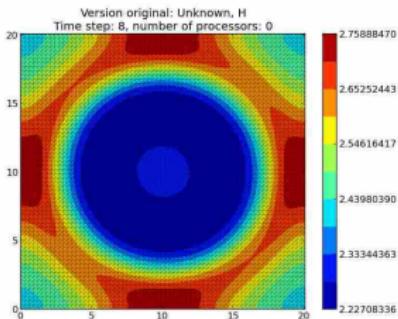
Courtesy of P. Langlois and R. Nheili
14 / 54

DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 8



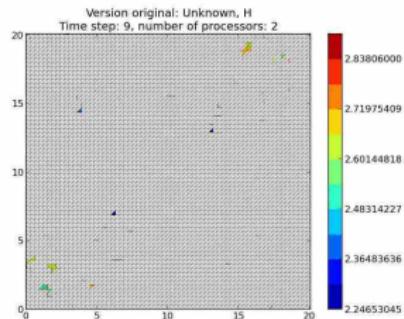
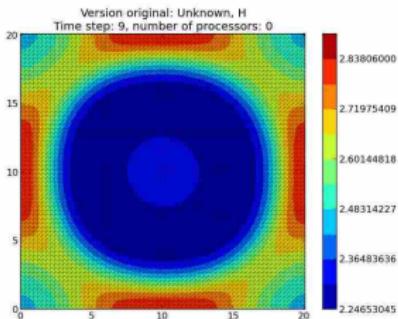
Courtesy of P. Langlois and R. Nheili
14 / 54

DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 9



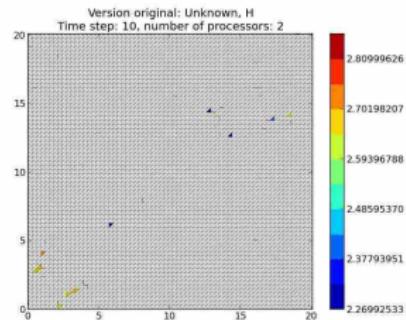
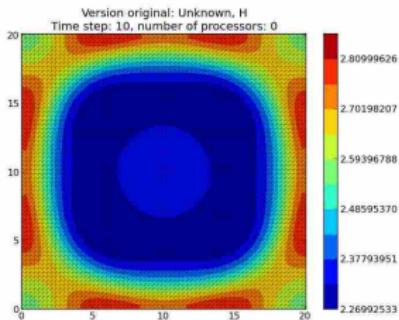
Courtesy of P. Langlois and R. Nheili
14 / 54

DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 10



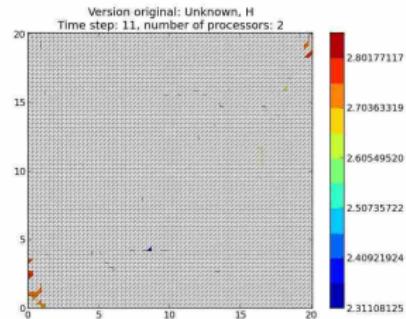
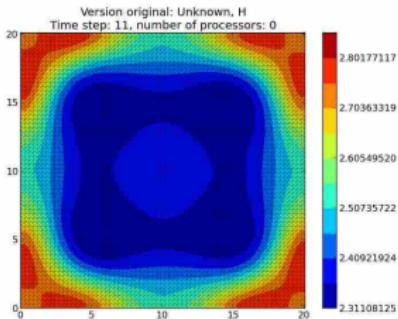
Courtesy of P. Langlois and R. Nheili
14 / 54

DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 11



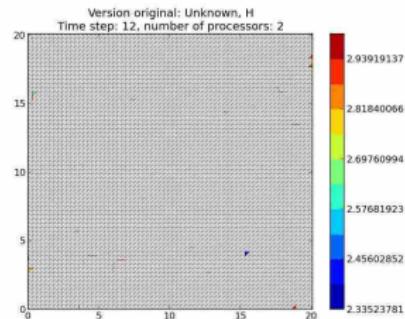
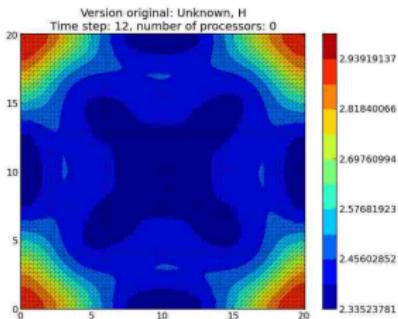
Courtesy of P. Langlois and R. Nheili
14 / 54

DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 12



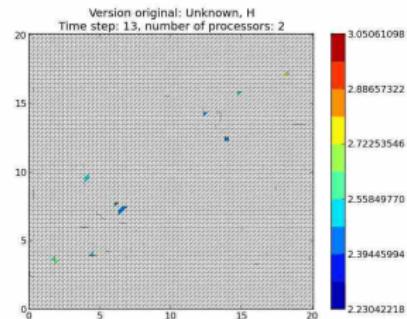
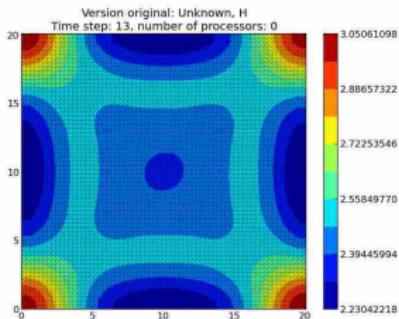
Courtesy of P. Langlois and R. Nheili
14 / 54

DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 13



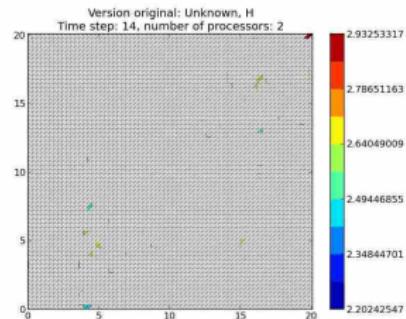
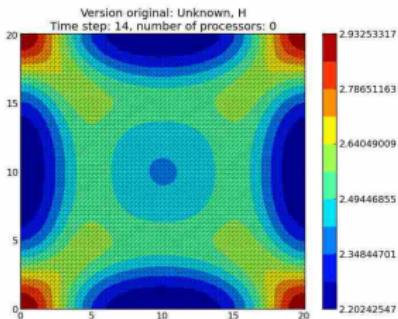
Courtesy of P. Langlois and R. Nheili
14 / 54

DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 14



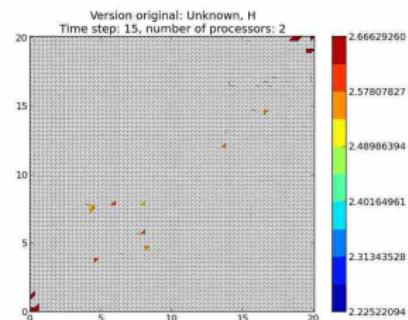
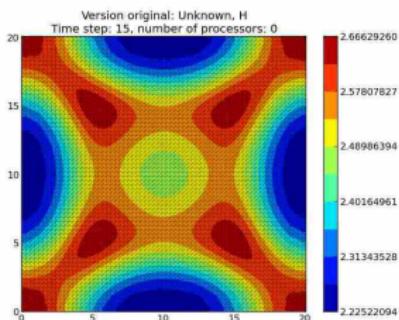
Courtesy of P. Langlois and R. Nheili
14 / 54

DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

NO numerical reproducibility!

time step = 15



Courtesy of P. Langlois and R. Nheili
14 / 54

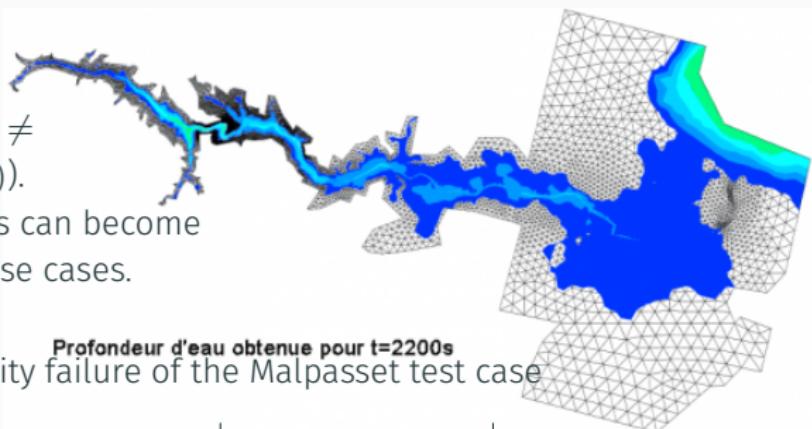
DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

$\text{round}(\text{round}(a + b) + c) \neq$
 $\text{round}(a + \text{round}(b + c)).$

These numerical issues can become quite harmful in real use cases.

Profondeur d'eau obtenue pour t=2200s

TABLE 1.1: Reproducibility failure of the Malpasset test case



| | The sequential run | a 64 procs run | a 128 procs run |
|------------|--------------------|----------------|-----------------|
| depth H | 0.3500122E-01 | 0.2748817E-01 | 0.1327634E-01 |
| velocity U | 0.4029747E-02 | 0.4935279E-02 | 0.4512116E-02 |
| velocity V | 0.7570773E-02 | 0.3422730E-02 | 0.7545233E-02 |

Numerical reproducibility: Approximations in the model, in the algorithm, in its implementation, in its execution.

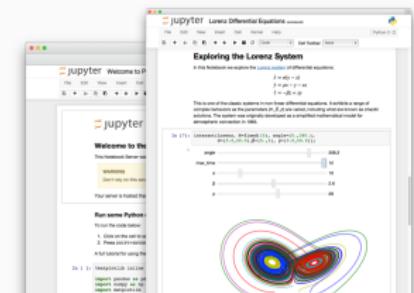
The whole chain needs to be revisited.

Courtesy of P. Langlois and R. Nheili

CONCLUSION AND ADVERTISING

GOOD RESEARCH REQUIRES TIME AND RESOURCES

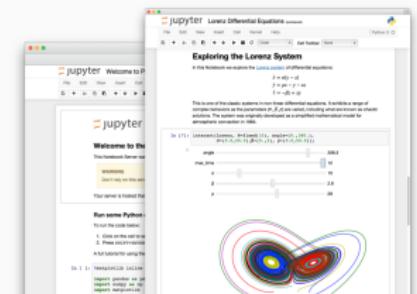
Computation provenance: notebooks



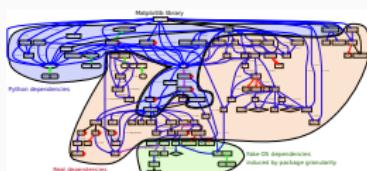
MOOC RR 1: Methodological
principles for a transparent science
3rd Edition: March 2020 – ... (22,800+)

GOOD RESEARCH REQUIRES TIME AND RESOURCES

Computation provenance:
notebooks and workflows



Software environments



Sharing and Archiving



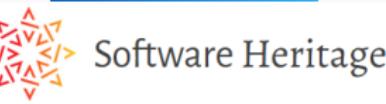
GNU Make



snake~~make~~

MOOC RR 1: Methodological
principles for a transparent science
3rd Edition: March 2020 – ... (22,800+)

MOOC RR 2: Practices and tools for
managing computations and data
1st Edition: May 2024–Oct 2024



WHAT'S THE POINT ?



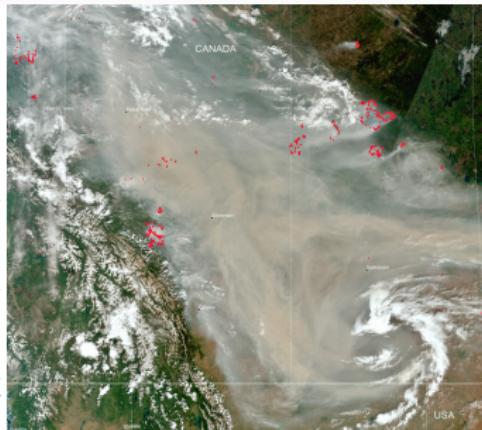
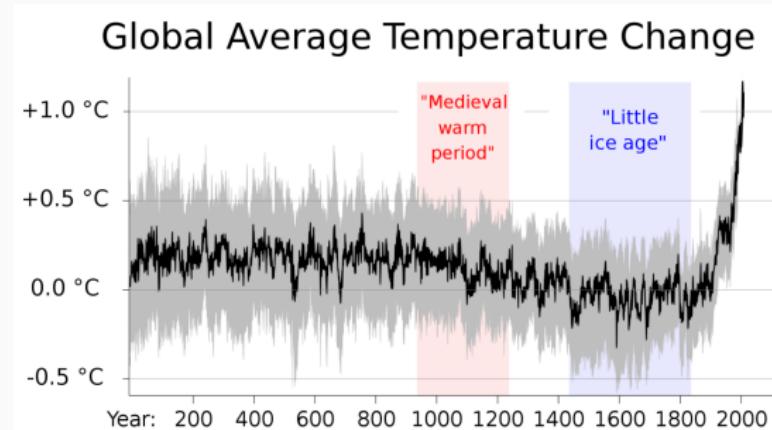
THE SCIENCE IS CLEAR

Why are we
ignoring it?

scientist rebellion

IPCC, IPBES, <https://climate.nasa.gov/>

1. Global climate change is not a future problem



https://en.wikipedia.org/wiki/Global_temperature_record

2023 Alberta wildfires (> 1 Mha)

THE SCIENCE IS CLEAR

scientist rebellion

Why are we ignoring it?

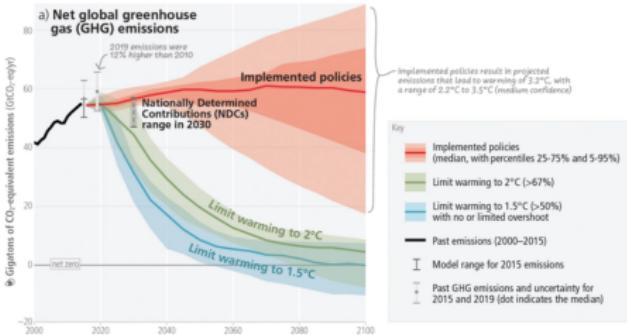


IPCC, IPBES, <https://climate.nasa.gov/>

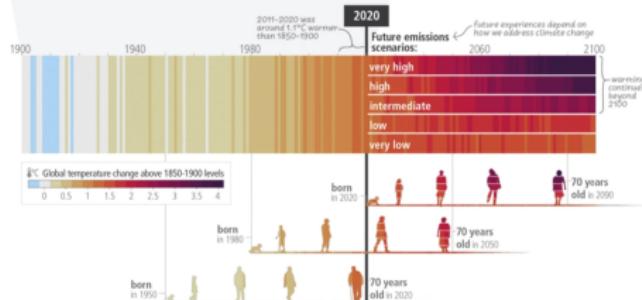
1. Global climate change is **not** a future problem
2. It is **entirely** due to human activity

Limiting warming to **1.5°C** and **2°C** involves rapid, deep and in most cases immediate greenhouse gas emission reductions

Net zero: CO₂ and net zero GHG emissions can be achieved through strong reductions across all sectors



c) The extent to which current and future generations will experience a hotter and different world depends on choices now and in the near-term



Latest IPCC report

27/28

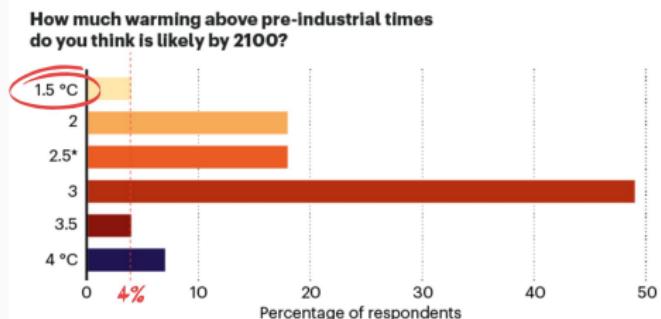
THE SCIENCE IS CLEAR

Why are we
ignoring it?

scientist rebellion

IPCC, IPBES, <https://climate.nasa.gov/>

1. Global climate change is **not** a future problem
2. It is **entirely** due to human activity
3. **9 out of 10 IPCC scientists believe overshoot is likely**



@natu Nature survey, Nov. 2021

THE ELEPHANT IN THE ROOM: CLIMATE CHANGE

Put aside biodiversity loss, pollution, freshwater, land system change...



(e) = estimations.

Note : l'empreinte carbone porte sur les trois principaux gaz à effet de serre (CO₂, CH₄, N₂O).

Champ : périmètre « Kyoto » (métropole et outre-mer appartenant à l'UE).

Sources : Citepa ; AIE ; FAO ; Douanes ; Eurostat ; Insee. Traitement : SDES, 2023

Empreinte carbone moyenne en France
10 tonnes de CO₂e/an/pers.



÷2
d'ici
2030

<2t CO₂e
<https://www.nosviesbascarbonne.org/>

Objectif d'ici 2050

- de 2 t de CO₂e/an/pers.

+ Faire plus d'activités bas carbone !

Danser, chanter, jardinier, rêver, écrire, lire, courir, randonner, planter des arbres, discuter, marcher en forêt, méditer, passer du temps avec ceux qu'on aime, rire...

Bref, inventer nos vies bas carbone désirables !

Par exemple :

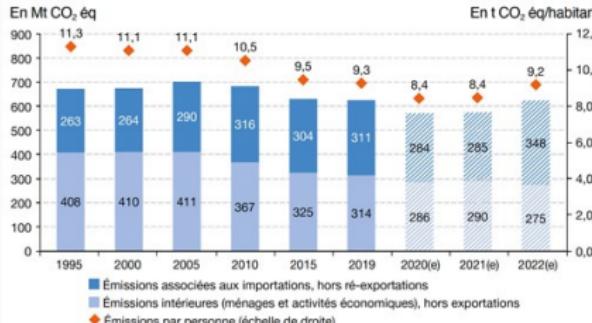
| | |
|-----|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| 0,5 | 1 CO2e Alimentation : J'achète légumes/legumineuses sans produits chimiques |
| 0,5 | 1 CO2e Transport : 200km en auto-petite voiture (BIOdeg, de fabrication anecdotique sur 30 ans, moins de 100g de CO2/km, sans émissions de méthane dans les transports en commun). |
| 0,5 | 1 CO2e Consommation : Quand on est devant, j'éteins la télévision, ferme toutes les installations électriques et informatiques, utilise une chaise métallique, je prends un bain au lieu d'une douche à chaleur ou solaire thermique. |
| 0,2 | 1 CO2e Services publics : Santé, enseignement, culture, etc. |
| 0,2 | 1 CO2e Services privés : Vacances, loisirs, etc. |

INVENTONS
NOS VIES
BAS CARBONE

Sources : Kit Inventons nos vies bas carbone (Fév. 2021), Rapport sur l'état de l'environnement en France (Déc. 2020)

THE ELEPHANT IN THE ROOM: CLIMATE CHANGE

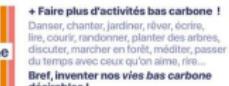
Put aside biodiversity loss, pollution, freshwater, land system change...



Empreinte carbone moyenne en France
10 tonnes de CO₂e/an/pers.



Objectif d'ici 2050
- de 2 t de CO₂e/an/pers.



Par exemple :



<https://www.nosviesbascarbonne.org/>

Sources : Kit Inventons nos vies bas carbone (Fév. 2021), Rapport sur l'état de l'environnement en France (Déc. 2020)

French government response

- Verdissement de l'industrie: « pause » sur les normes environnementales
- Loi de programmation militaire (+41%)
- Nous devons préparer la France à une élévation de la température de 4 °C
- Academia ? PEPR 5G, Cloud, NUMPEX, Quantique, IA, Agroécologie et numérique

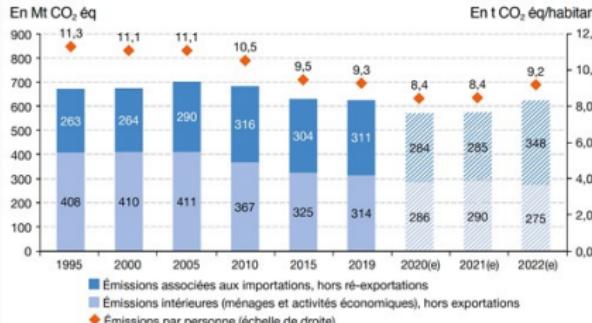


INVENTONS
NOS VIES
BAS CARBONE

THE ELEPHANT IN THE ROOM: CLIMATE CHANGE

2/2

Put aside biodiversity loss, pollution, freshwater, land system change...



(e) = estimations.

Note : l'empreinte carbone porte sur les trois principaux gaz à effet de serre (CO₂, CH₄, N₂O).

Champ : périmètre = Kyoto + (Île-de-France et outre-mer appartenant à l'UE).

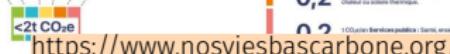
Sources : Citepa ; AIE ; FAO ; Douanes ; Eurostat ; Insee. Traitement : SODES, 2023



Empreinte carbone moyenne en France
10 tonnes de CO₂e/an/pers.



÷2
d'ici
2030



Sources : Kit Inventons nos vies bas carbone (Fév. 2021), Rapport sur l'état de l'environnement en France (Déc. 2020)

Objectif d'ici 2050

- de 2 t de CO₂e/an/pers.

+ Faire plus d'activités bas carbone !

Danser, chanter, jardinier, rêver, écrire, lire, courir, randonner, planter des arbres, discuter, marcher en forêt, méditer, passer du temps avec ceux qu'on aime, rire...

Bref, inventer nos vies bas carbone désirables !

Par exemple :



INVENTONS
NOS VIES
BAS CARBONE

French government response

- Verdissement de l'industrie: « pause » sur les normes environnementales
- Loi de programmation militaire (+41%)
- Nous devons préparer la France à une élévation de la température de 4 °C
- Academia ? PEPR 5G, Cloud, NUMPEX, Quantique, IA, Agroécologie et numérique

Several scenarios on the table

- Energy optimization/saving ≠ sobriety and frugality
- What will research/CS look like/be used for in such a world?