

# REPRODUCIBLE RESEARCH AND COMPUTER SCIENCE

---

Arnaud Legrand



Master 1 UGA  
December 2024



# WHAT IS SCIENCE ABOUT?

Question: In less than 5 lines give a definition of "Science"

# WHAT IS SCIENCE ABOUT?

**Question:** In less than 5 lines give a definition of "Science"

## Dictionary of science and technology

1. the study of the physical and natural world and phenomena, especially by using systematic observation and experiment
2. a particular area of study or knowledge of the physical world
3. a systematically organized body of knowledge about a particular subject

**New Oxford Dictionary** the intellectual and practical activity encompassing the systematic study of the structure and behavior of the physical and natural world through observation and experiment : the world of science and technology.

1. a particular area of this : veterinary science | the agricultural sciences.
2. a systematically organized body of knowledge on a particular subject : the science of criminology.
3. archaic knowledge of any kind.

# WHAT IS SCIENCE ABOUT?

Question: In less than 5 lines give a definition of "Science"

## Dictionary of science and technology

1. the study of the physical and natural world and phenomena, especially by using systematic observation and experiment
2. a particular area of study or knowledge of the physical world
3. a systematically organized body of knowledge about a particular subject

**New Oxford Dictionary** the intellectual and practical activity encompassing the systematic study of the structure and behavior of the physical and natural world through observation and experiment : the world of science and technology.

1. a particular area of this : veterinary science | the agricultural sciences.
2. a systematically organized body of knowledge on a particular subject : the science of criminology.
3. archaic knowledge of any kind.

**Building Reliable Knowledge**

## SCIENTIFIC CONSENSUS VS. DEMOCRACY AND FREEDOM OF SPEECH



# REPRODUCIBILITY CRISIS

## Is there a reproducibility crisis A Nature survey, 2016

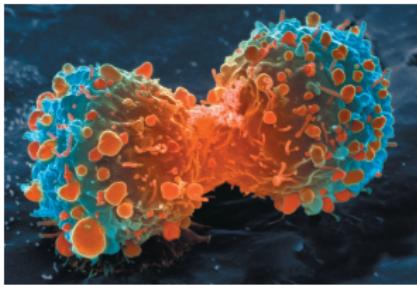
### COMMENT

DATA MINERS Shift expertise to track mutations where they strike a blow

EARTH SYSTEM Past climates give valuable clues to future warming [107](#)

HISTORY OF SCIENCE Don't let us forget better track using Google [149](#)

ENVIRONMENT Wyllie Vida and colleagues stress how much [147](#)



Many landmark findings in preclinical oncology research are not reproducible, in part because of inadequate cell lines and animal models.

## Raise standards for preclinical cancer research

C. Glenn Begley and Lee M. Ellis propose how methods, publications and incentives must change if patients are to benefit.

Efforts over the past decade to characterize the genetic alterations in human cancers have led to a better understanding of disease and driven some successes in therapy. Although this is a cause for optimism, it also has led to a cancer field that hoped this would lead to more effective drugs, historically, our ability to translate cancer research from the lab to the clinic has been miserably low. Sadly, clinical

trials in oncology have the highest failure rate compared with other therapeutic areas. Given the high unmet need in oncology, it is understandable that the barriers to clinical development may be lower than for other diseases. However, a larger number of drugs with suboptimal preclinical validation will contribute to the failure of trials. The success rate is not sustainable or acceptable, and

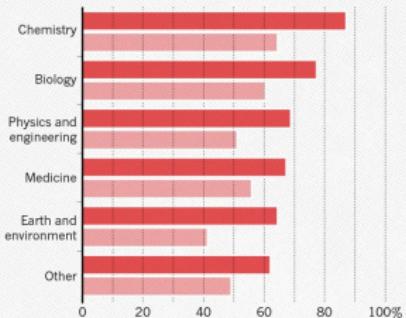
Investigation must focus on those that approach to translating discovery research into greater clinical success and impact.

Many factors are responsible for the high failure rate in oncology, including the inherently difficult nature of this disease. Certainly, the limitations of preclinical tools with regard to cell-line and animal models and mouse models make it difficult to even

## HAVE YOU FAILED TO REPRODUCE AN EXPERIMENT?

Most scientists have experienced failure to reproduce results.

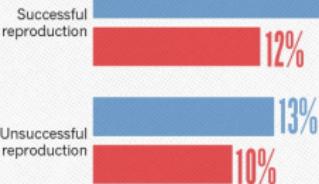
- Someone else's
- My own



## HAVE YOU EVER TRIED TO PUBLISH A REPRODUCTION ATTEMPT?

Although only a small proportion of respondents tried to 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

## Must try harder

Too many sloppy mistakes are creeping into scientific papers at the data — and at themselves.

## Six red flags for suspect work

C. Glenn Begley explains how to recognize the preclinical papers in which the data won't stand up

## Error prone

Biologists must realize the pitfalls of massive amount of data

## Know when your numbers are significant

- Nekrutenko & Taylor, *Nature Genetics* (2012)
- Alsheikh-Ali et al. *PLoS ONE* (2011)
- Begley & Ellis *Nature* (2012)

## COMMON REPRODUCIBILITY PITFALLS

---

# GO READ THE PAPER BY SMITH ET. AL. 2009

Article preview      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. Follenbeck , T. Dörnemann , C. Schridde , B. Freisleben 

Show more 

+ Add to Mendeley  Share 

---

<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  Cite  
<https://doi.org/10.1016/j.jpdc.2022.04.002> 

### 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

The screenshot shows the HAL digital library interface. At the top, there is a search bar with the placeholder "Chercher un document, un auteur, un mot clé...". Below the search bar, there is a button labeled "Télécharger pour visualiser". On the left, there is a sidebar with sections for "Dates et versions" and "Identifiants". The "Dates et versions" section lists two versions: "hal-03141988, version 1 (15-02-2021)" and "hal-03141988, version 2 (06-01-2022)". The "Identifiants" section lists "HAL Id: hal-03141988, version 2", "ARXIV: 2102.07674", and "DOI: 10.1016/j.jpdc.2022.04.002". The main content area displays the title "Simulation-based Optimization and Sensibility Analysis of MPI Applications: Variability Matters" by Tom Cornebize (1, 2), Arnaud Legrand (3, 1). It also lists three institutions: 1 POLARIS - Performance analysis and optimization of LARge Infrastructures and Systems, 2 UGA - Université Grenoble Alpes, and 3 CNRS - Centre National de la Recherche Scientifique. Below the title, there is a summary in French: "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". There are also sections for "Mots clés" (tags) such as "Simulation", "validation", "sensibility analysis", and "SimGrid".

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

Rodriguez et al., CONCUR'15

## Unfolding-based Partial Order Reduction\*

César Rodríguez<sup>1</sup>, Marcelo Sousa<sup>2</sup>, Subodh Sharma<sup>3</sup>, and  
Daniel Kroening<sup>4</sup>

<sup>1</sup> Université Paris 13, Sorbonne Paris Cité, LIPN, CNRS, France

<sup>2,4</sup> Department of Computer Science, University of Oxford, UK

<sup>3</sup> 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.

## 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

---

<sup>3</sup> 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)				
	Name	$ P $	$ I $	$ B $	$t(s)$	$ E $	$ \Omega $	$\langle  U_\Omega  \rangle$	$t(s)$	$ E $	$ E_{\text{cut}} $	$ \Omega $	$\langle  U_\Omega  \rangle$
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

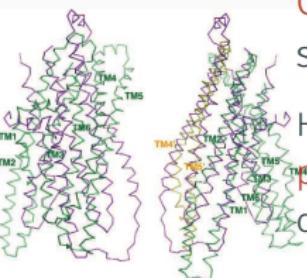
- Looks good! Err... **not so simple**. Depending on how you do this, you quickly move from polynomial to exponential.
- **Experiments!**  $\triangle$ Possible 404, code not found! ahead!!!
- Wait, what's this language? Did this ever run one day?
- Wow, I'll need **CPLEX** and **Gurobi** but all I have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

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

## BLAMING "COMPUTER SCIENCE"

---

# How COMPUTERS BROKE SCIENCE



Geoffrey Chang (Scripps, UCSD) works on crystallography and studies the structure of cell membrane proteins.

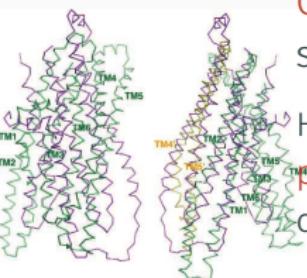
He specialized in structures of **multidrug resistant transporter proteins in bacteria**: MsbA de Escherichia Choli (Science, 2001), Vibrio cholera (Mol. Biology, 2003), Salmonella typhimurium (Science, 2005)

**2006:** Inconsistencies reveal a programming mistake

*A homemade data-analysis program had flipped two columns of data, inverting the electron-density map from which his team had derived the protein structure.*

5 retractions that motivate improved software engineering practices in comp. biology

# How COMPUTERS BROKE SCIENCE



Geoffrey Chang (Scripps, UCSD) works on crystallography and studies the structure of cell membrane proteins.

He specialized in structures of **multidrug resistant transporter proteins in bacteria**: MsbA de Escherichia Choli (Science, 2001), Vibrio cholera (Mol. Biology, 2003), Salmonella typhimurium (Science, 2005)

**2006:** Inconsistencies reveal a programming mistake

*A homemade data-analysis program had flipped two columns of data, inverting the electron-density map from which his team had derived the protein structure.*

5 retractions that motivate improved software engineering practices in comp. biology

# COMPUTERS...

## How computers broke science – and what we can do about it

*Most modern science is so complicated, and most journal articles so brief, it's impossible for the article to include details of many important methods and decisions made by the researcher as he analyzed his data on his computer. How, then, can another researcher judge the reliability of the results, or reproduce the analysis?*



– Ben Marwick,  
The conversation, 2015

**Point-and-click** procedures are rampant but they hinder reproducibility.

**Spreadsheets** are generalized and intensively used in biology:

- **Membrane-Associated Ring Finger (C3HC4) 1,**  
**E3 Ubiquitin Protein Ligase** → **MARCH1** → 2016-03-01 →  
1456786800
- **2310009E13** → 2.31E+19

And more recently, we had the **COVID tracing failure**.



Machine Learning: Trouble at the lab, The Economist 2013

*According to some estimates, three-quarters of published scientific papers in the field of machine learning are bunk because of this "overfitting".*

– Alex "Sandy" Pentland

The Reproducibility Crisis in ML-based science (Princeton workshop 2022)

*Reproducibility failures in ML-based science are systemic. We found 20 reviews across 17 scientific fields (medicine, neuroimaging, autism diagnosis, genomics, computer security, ...) that find errors in a total of 329 papers that use ML-based science and in some cases leading to wildly overoptimistic conclusion. [...] complex ML models don't perform substantively better than decades-old LR models.*

*Data leakage:* spurious relationship between the independent variables and the target variable that arises as an artifact of the data collection, sampling, or pre-processing strategy.

– S. Kapoor and A. Narayanan

## THIS IS ABOUT COMPUTATIONAL SCIENCE. SHOULD MATHEMATICIANS CARE?

Computer Science is young and inherits from Mathematics, Engineering,  
Linguistic, Nat. Sciences, ...

Purely theoretical scientists whose practice is close to mathematics may not be concerned (can't publish a math article without releasing the proofs).

# THIS IS ABOUT COMPUTATIONAL SCIENCE. SHOULD MATHEMATICIANS CARE?

Computer Science is young and inherits from Mathematics, Engineering,  
Linguistic, Nat. Sciences, ...

Purely theoretical scientists whose practice is close to mathematics may not be concerned (can't publish a math article without releasing the proofs).

Yet, incoherencies are common, especially in a fast moving field:

- E.g., definitions/concepts in book/article A and B are *slightly different* and the resulting theorems cannot be mixed
- Have a look at Vladimir Voevodsky's talk in 2014 at Princeton 😊
- ERC Nano bubbles: how, when and why does science fail to correct itself?

Flagging incorrect nucleotide sequence reagents in biomedical papers:

To what extent does the leading publication format impede automatic  
error detection?

(Labbe et al., 2020)

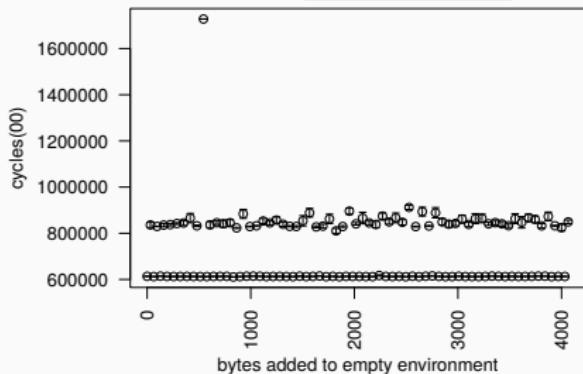
# EXPERIMENTING WITH COMPUTERS

Machines are real!



Brendan Gregg: Shouting in the data center

Machines are complicated



Mytkowicz et al. *Producing wrong data without doing anything obviously wrong!*  
ACM SIGPLAN Not. 44(3), March 2009

Our reality evolves!!! The hardware keeps evolving so most results on old platforms quickly become obsolete (although, we keep building on such results 😊).

We need to regularly revisit and allow others to build on our work!

# COMPUTER PERFORMANCE ? WELL, I DESIGN ALGORITHMS!

- "Real" problems are all NP-hard, Log-APX, etc.
- Real workload = ~~NP-completeness proof widgets~~, regularities and properties (difficult to formally state but that should be exploited)

Algorithms are evaluated on particular **workloads** that impact both their running time and the quality of the solutions

# COMPUTER PERFORMANCE ? WELL, I DESIGN ALGORITHMS!

- "Real" problems are all NP-hard, Log-APX, etc.
- Real workload = ~~NP-completeness proof~~ widgets, regularities and properties (difficult to formally state but that should be exploited)

Algorithms are evaluated on particular **workloads** that impact both their running time and the quality of the solutions

**Image Processing:** **True horror stories**, E. Meinhardt-Llopis, CANUM 2016

- *The proposed multigrid algorithm converges to the solution of the problem in  $O(N)$  using biharmonic functions*
- Surprisingly, our naive multi-scale Gauss-Seidel converges much faster

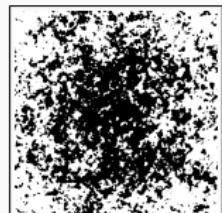
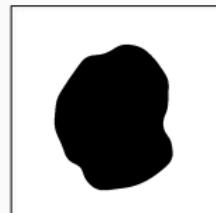
# COMPUTER PERFORMANCE ? WELL, I DESIGN ALGORITHMS!

- "Real" problems are all NP-hard, Log-APX, etc.
- Real workload = ~~NP-completeness proof~~ widgets, regularities and properties (difficult to formally state but that should be exploited)

Algorithms are evaluated on particular **workloads** that impact both their running time and the quality of the solutions

**Image Processing:** True horror stories, E. Meinhardt-Llopis, CANUM 2016

- *The proposed multigrid algorithm converges to the solution of the problem in  $O(N)$  using biharmonic functions*
- Surprisingly, our naive multi-scale Gauss-Seidel converges much faster



## DIFFERENT KINDS OF REPRODUCIBILITY

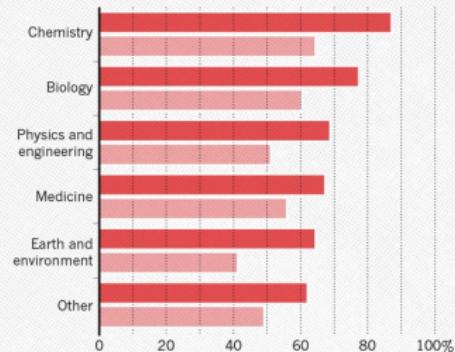
---

# SOCIO-TECHNICAL CHALLENGES

## HAVE YOU FAILED TO REPRODUCE AN EXPERIMENT?

Most scientists have experienced failure to reproduce results.

● Someone else's ● My own



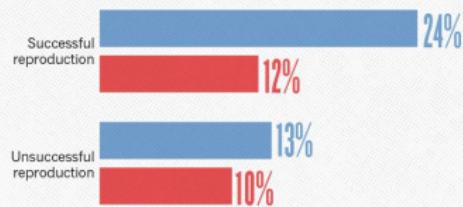
1,500 scientists lift the lid on reproducibility,

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



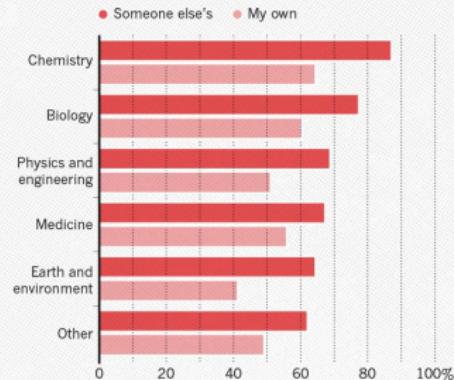
Number of respondents from each discipline:

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

# SOCIO-TECHNICAL CHALLENGES

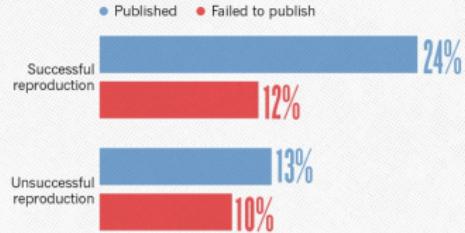
## HAVE YOU FAILED TO REPRODUCE AN EXPERIMENT?

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.



1,500 scientists lift the lid on reproducibility,

Nature, May 2016

## Social causes

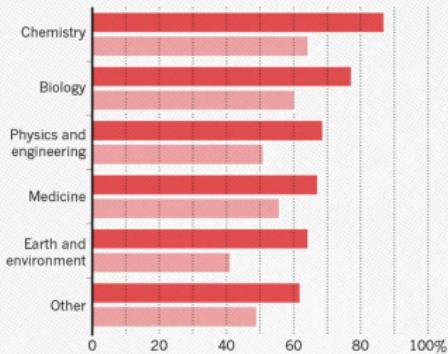
- Fraud, conflict of interest (pharmaceutic, ...)
- No incentive to reproduce/check our own work (afap), nor the work of others (big results!), nor to allow others to check (competition)
- Peer review does not scale: 1M+ articles per year!

# SOCIO-TECHNICAL CHALLENGES

## HAVE YOU FAILED TO REPRODUCE AN EXPERIMENT?

Most scientists have experienced failure to reproduce results.

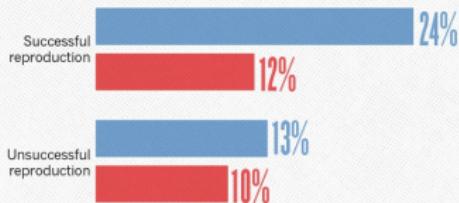
● Someone else's ● My own



## 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

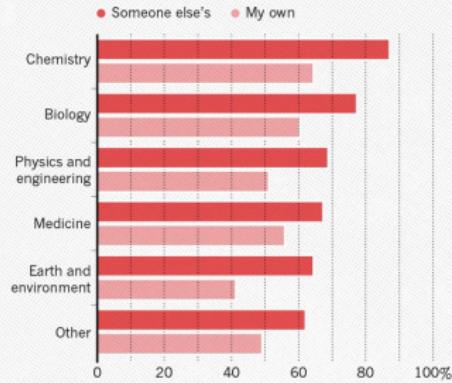
## Social causes

- 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

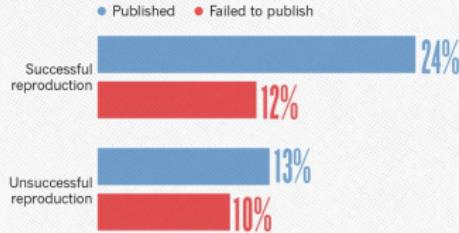
## HAVE YOU FAILED TO REPRODUCE AN EXPERIMENT?

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.



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

# NO TRANSPARENCY NO CONSENSUS



# DIFFERENT REPRODUCIBILITY CONCERNS IN MODERN SCIENCE

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

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

# 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

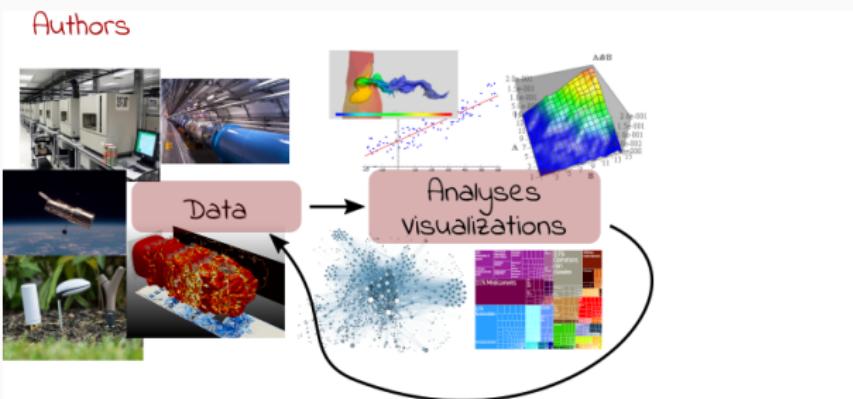
**Psychology, Nutrition** HARKING, p-hacking  $\rightsquigarrow$  pre-registration

**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*



# DIFFERENT REPRODUCIBILITY CONCERN IN MODERN SCIENCE

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

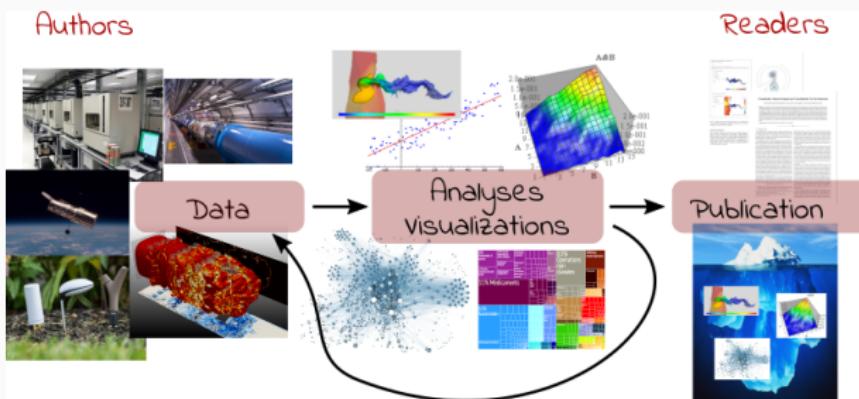
**Psychology, Nutrition** HARKING, p-hacking  $\rightsquigarrow$  pre-registration

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

**Computational fluid dynamics** numerical chaos, parallel architectures

**Artificial Intelligence** most of the above 😊

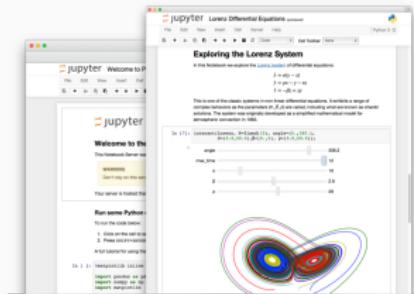
*The processing steps between raw observations and findings have gotten increasingly numerous and complex*



Reproducible Research = Bridging the Gap by working Transparently

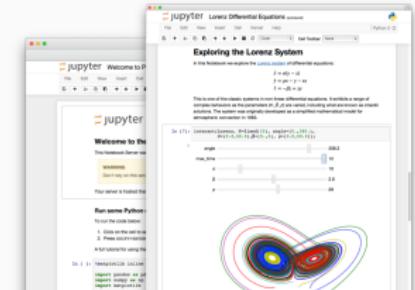
# REPRODUCIBILITY ISSUES RELATED TO THE USE OF COMPUTERS

## Computation provenance: notebooks and workflows

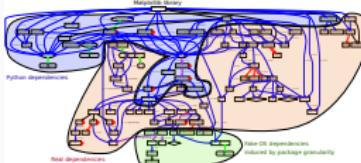


# REPRODUCIBILITY ISSUES RELATED TO THE USE OF COMPUTERS

## Computation provenance: notebooks and workflows



## Software environments

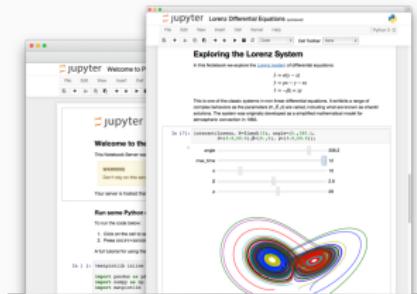


**ReproZip**

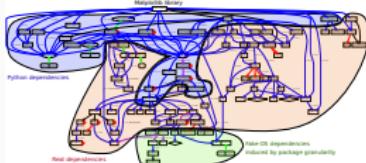


# REPRODUCIBILITY ISSUES RELATED TO THE USE OF COMPUTERS

## Computation provenance: notebooks and workflows



## Software environments



## Sharing and Archiving



## GOOD PRACTICE #1

### TAKING NOTES AND DOCUMENTING

---

# FRUSTRATION AS AN AUTHOR/REVIEWER



## 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

# TOOL 1: COMPUTATIONAL NOTEBOOKS/LITTERATE PROGRAMMING

## Un document computationnel

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

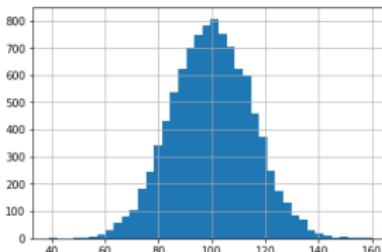
3.141592653589793

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

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... ☺).



# TOOL 1: COMPUTATIONAL NOTEBOOKS/LITTERATE PROGRAMMING

Document initial dans son environnement

The screenshot shows a Jupyter Notebook interface with the following content:

**# Un document computationnel**

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

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) ([https://fr.wikipedia.org/wiki/Aiguille\\_de\\_Buffon](https://fr.wikipedia.org/wiki/Aiguille_de_Buffon)), on obtientrait 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=np.pi/2)
2/(sum((x+np.sin(theta))>1))/N
```

Out[2]: 3.14371986944998765

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

Document final

## Un document computationnel

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

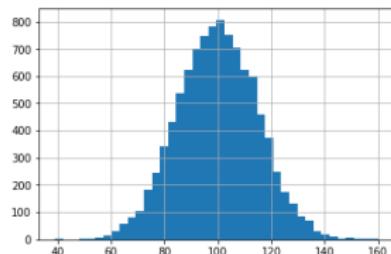
3.141592653589793

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=np.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  $\pi$  (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 [méthode des aiguilles de Buffon](#) ([https://fr.wikipedia.org/wiki/Aiguille\\_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  $\frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$  et des dessins qui n'ont rien à voir avec  $\pi$  (si ce n'est une constante de normalisation...).

```
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  $\pi$  vaut approximativement

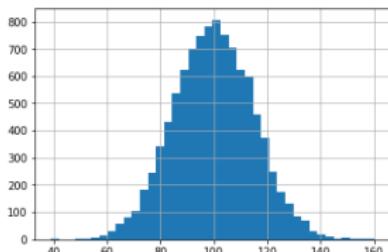
3.141592653589793

Mais calculé avec la [méthode des aiguilles de Buffon](#), on obtient 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  $\pi$  (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. A red arrow points from this cell to the word "Code" in the "Document final" section.
- In [2]:** Calculates the area under a curve using the Monte Carlo method. It imports numpy, generates random numbers, calculates theta, and sums the results. The output is 3,1437198694098765. Another red arrow points from this cell to the word "Code".
- In [3]:** Plots a histogram of 100,000 random numbers. The plot shows a bell-shaped distribution centered around 100, with the x-axis ranging from 40 to 160 and the y-axis from 0 to 800. A red arrow points from this cell to the histogram itself.

Document final

## Un document computationnel

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

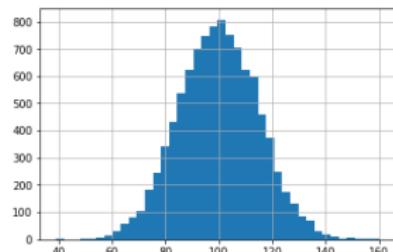
3.141592653589793

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

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... ☺).



# 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 `_methode_ des épingles de Buffon` ([https://fr.wikipedia.org/wiki/Algille\\_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  $\pi$  (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()
```

Document final

## Un document computationnel

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

3.141592653589793

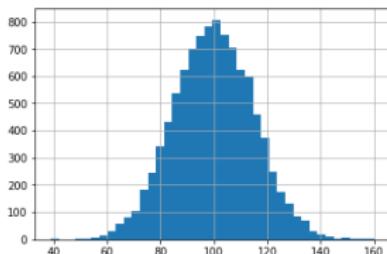
Mais calculé avec la méthode des [algilles de Buffon](#), on obtiendrait comme approximation :

```
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  $\pi$  (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 the Buffon needle method.
- In [2]:** Generates random points (x, theta) and calculates the ratio of points where sin(theta) > x, which approximates pi/2.
- In [3]:** Plots a histogram of x values from 40 to 160, showing a bell-shaped distribution centered around 100.

Document final

## Un document computationnel

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

3.141592653589793

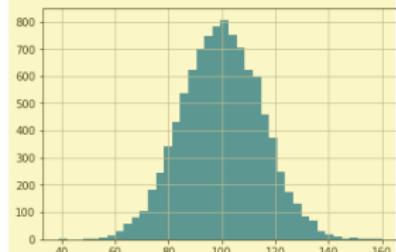
Mais calculé avec la méthode des [aiguilles de Buffon](#), on obtiendrait comme approximation :

```
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... ☺).



# 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 it with the Buffon's needle method.
- In [2]:** Generates random points and calculates the ratio of points below a line to the total number of points, which is used to approximate  $\pi$ .
- In [3]:** Plots a histogram of 100,000 random numbers between 0 and 1, showing a bell-shaped distribution centered at 0.5.

Document final

## Un document computationnel

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

3.141592653589793

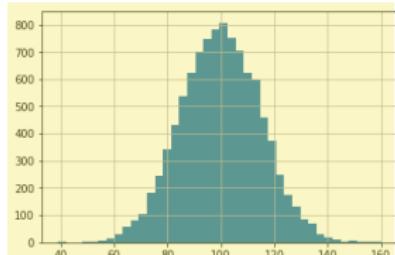
Mais calculé avec la méthode des [aiguilles de Buffon](#), on obtiendrait comme approximation :

```
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)
```

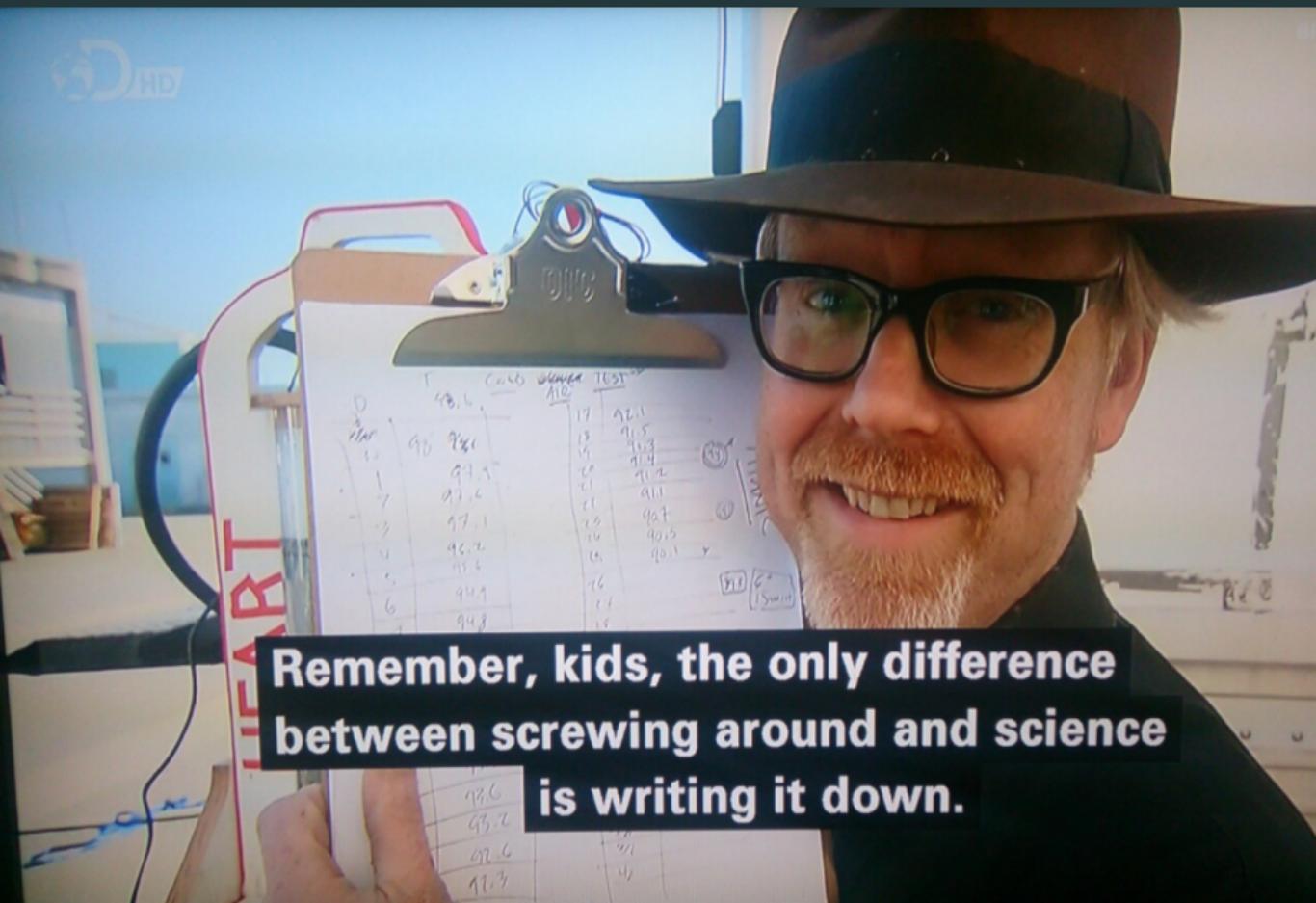
3.1437198694098765

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... ☺).



## TOOL 1 BIS: LABORATORY NOTEBOOKS, COMPUTATIONAL DOCUMENTS



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

# TOOL 1 TER: WORKFLOWS

Notebooks are no panacea and do not help developing clean code

The screenshot shows a Jupyter Notebook interface with the title bar "jupyter example\_pi.ipynb". The notebook contains several cells:

- In [1]:** A code cell containing:

```
# Un document computationnel

# Mon ordinateur n'indique que j'ai 15 chiffres "approximativement"

In [1]:
```

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

Annotations above this cell say "Hide Prompt" and "Hide Code". Annotations below the output say "Mais calculé avec la \_\_method\_\_ des (ajoutées de Buffet) `math.pi_as_double()`, on obtientrait comme approximation...".
- In [2]:** A code cell containing:

```
import numpy as np
n = 1000000
x = np.random.uniform(0, low=0, high=1)
theta = np.random.uniform(0, low=0, high=np.pi/2)
if (x**2 + np.sin(theta)**2) < 1/n
```

Annotations above this cell say "Hide Prompt" and "Hide Code". Annotations below the output say "On peut inclure des formules mathématiques comme `Sqrt(2)/pi` ou `(4*pi/3)*sqrt(1/(pi+2))` dans les cellules de code et elles seront automatiquement dessinées ou n'ont rien à voir avec l'output (si ce n'est une constante de normalisation...)."
- In [3]:** A code cell containing:

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

n = 1000000
x = np.random.uniform(0, 1)
y = np.sqrt(1-x**2)
plt.hist(x, 40)
plt.title("Pi")
plt.show()
```

Annotations above this cell say "Hide Prompt" and "Hide Code". Annotations below the output say "Hide Output".

The output of this cell is a histogram showing a distribution of points within a unit circle, centered at the origin, with the x-axis ranging from 0 to 1 and the y-axis ranging from 0 to 1. The distribution is roughly triangular, peaking at x=0.5 and y=0.5, with a total count of approximately 100,000 points.

# TOOL 1 TER: WORKFLOWS

Notebooks are no panacea and do not help developing clean code

The screenshot shows a Jupyter Notebook interface with the title "analyse-syndrome-grippal". The notebook contains several code cells and text sections. One section is titled "Incidence du syndrome grippal". It includes a snippet of R code for calculating incidence from a dataset. Another section discusses the calculation of the number of consultations per day and provides a table with data. A third section contains a histogram visualizing the data.

```
In [1]: #!/usr/bin/rscript -e
library(tidyverse)
library(lubridate)
library(httr)
library(dplyr)
```

```
In [2]: # Les données de l'Institut de Santé Publique peuvent être obtenues à l'adresse suivante : https://www.insee.fr/fr/statistiques/4556946. Nous les disposons sous forme d'un fichier 'syndrome_grippal.csv'. Pour consulter ce fichier, nous devons le télécharger et pour accéder à son contenu utiliser la commande 'head()' ou 'tail()'. Le premier ligne de tête ('row 1') est un nom des colonnes et ne doit pas être traité.
```

```
In [3]: # Voici l'incidence des consultations établie sur le site datagouv
```

```
In [4]: # Fonction pour calculer l'incidence
incidence<-function(df) {
  df %>%
    mutate(date=ymd(hm(strftime(df$date,"%Y-%m-%d %H:%M")))) %>%
    filter(date>=ymd("2016-01-01")) %>%
    group_by(date) %>%
    summarise(n_consultations=n(), n_consultants=n_distinct(consultant))
}
```

```
In [5]: # Calcul de l'incidence à partir de l'ensemble des données
incidence(df)>%>
  mutate(incidence=(n_consultations/n_consultants)*1000000) %>%
  select(-n_consultants,-n_consultations,-date)
```

```
In [6]: # Il existe également une fonction de l'Insee permettant de calculer l'incidence pour 100 000 personnes
# par département. Voici le lien pour accéder à cette fonction : https://www.insee.fr/fr/statistiques/4556946
# Dans le cas de la ligne suivante nous avons utilisé la fonction 'syndrome_grippal' pour obtenir la liste des colonnes
# dont nous nous serviront plus tard pour calculer l'incidence par département.
```

```
In [7]: # Fonction pour calculer l'incidence à partir de l'ensemble des données
incidence_insee<-function(df) {
  df %>%
    group_by(departement) %>%
    summarise(n_consultants=n(), n_consultations=n())
}
```

```
In [8]: # Calcul de l'incidence à partir de l'ensemble des données
incidence_insee(df)
```

departement	n_consultants	n_consultations
0	2993382	2993382
1	2217636	2217636
2	3120439	3120439
3	2218843	2218843
4	2250662	2250662
5	310996	310996
6	481635	481635
7	32506	32506
8	91541	91541
9	523957	523957
10	14485	14485

```
In [9]: # Ensuite nous devons bien sûr faire quelques tests, qui fournissent environ 50% de la population française, soit environ 500 000 personnes au total
```

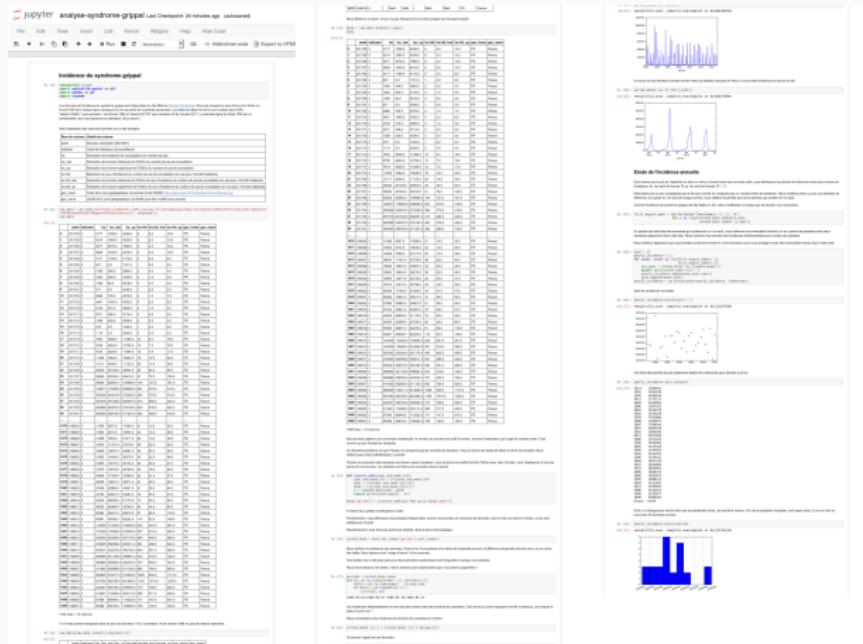
```
In [10]: #parcours_insee %>% head(50000)
```

```
In [11]: # Analyse de la variance des valeurs. Ainsi, on peut voir que les deux dernières colonnes sont fortement corrélées.
```

```
In [12]: ggplot(data=syndrome_grippal, aes(x=date, y=incidence)) +
```

# TOOL 1 TER: WORKFLOWS

Notebooks are no panacea and do not help developing clean code



# TOOL 1 TER: WORKFLOWS

Notebooks are no panacea and do not help developing clean code

The image displays a 4x3 grid of Jupyter Notebook screenshots, each illustrating a different aspect of data analysis or machine learning. The notebooks include code snippets, explanatory text, and visualizations such as scatter plots and heatmaps.

- Cell 1: Estimating Color Names by Web Image**

Code to extract color names from a web image. Includes a scatter plot titled "Chromaticity distribution of training data".
- Cell 2: Dimensionality reduction**

Code to reduce dimensionality of a dataset. Includes a heatmap titled "Dimensionality plane and linear model results, dimensionality filtered data".
- Cell 3: Dimensionality reduction**

Code to reduce dimensionality of a dataset. Includes a heatmap titled "Dimensionality plane and linear model results, dimensionality filtered data".
- Cell 4: Dimensionality reduction**

Code to reduce dimensionality of a dataset. Includes a heatmap titled "Dimensionality plane and linear model results, dimensionality filtered data".
- Cell 5: Dimensionality reduction**

Code to reduce dimensionality of a dataset. Includes a heatmap titled "Dimensionality plane and linear model results, dimensionality filtered data".
- Cell 6: Dimensionality reduction**

Code to reduce dimensionality of a dataset. Includes a heatmap titled "Dimensionality plane and linear model results, dimensionality filtered data".
- Cell 7: Dimensionality reduction**

Code to reduce dimensionality of a dataset. Includes a heatmap titled "Dimensionality plane and linear model results, dimensionality filtered data".
- Cell 8: Dimensionality reduction**

Code to reduce dimensionality of a dataset. Includes a heatmap titled "Dimensionality plane and linear model results, dimensionality filtered data".
- Cell 9: Dimensionality reduction**

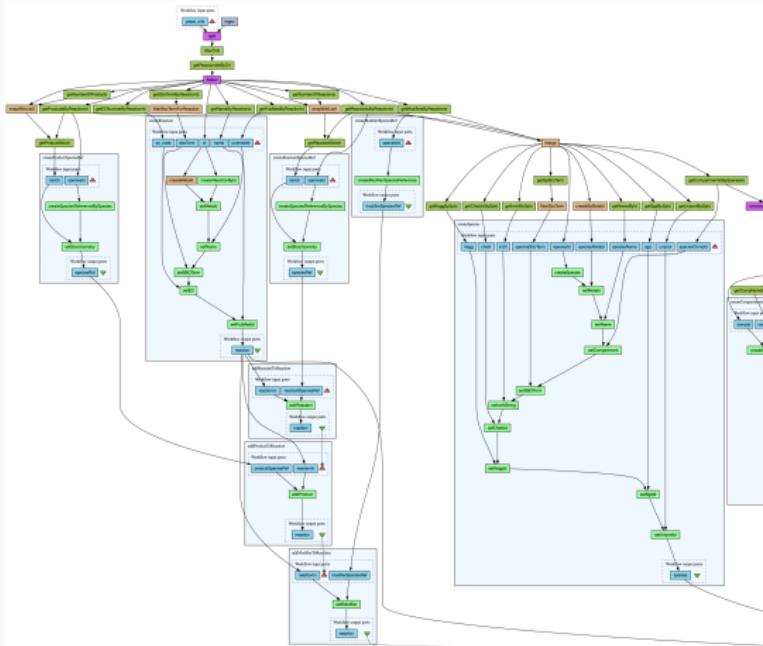
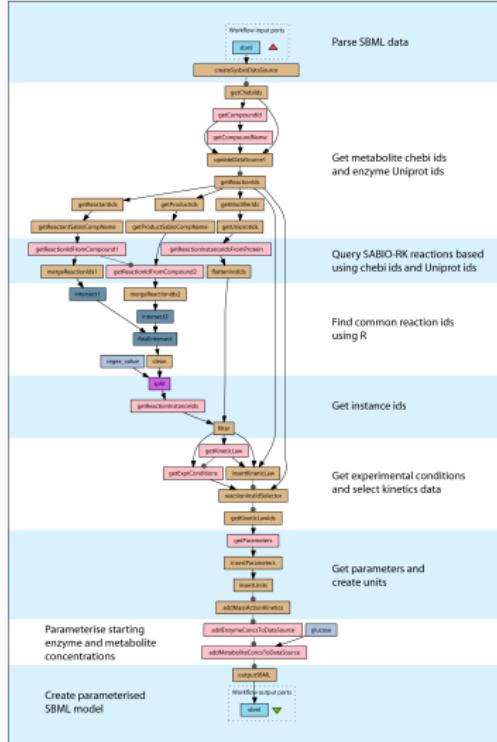
Code to reduce dimensionality of a dataset. Includes a heatmap titled "Dimensionality plane and linear model results, dimensionality filtered data".
- Cell 10: Dimensionality reduction**

Code to reduce dimensionality of a dataset. Includes a heatmap titled "Dimensionality plane and linear model results, dimensionality filtered data".
- Cell 11: Dimensionality reduction**

Code to reduce dimensionality of a dataset. Includes a heatmap titled "Dimensionality plane and linear model results, dimensionality filtered data".
- Cell 12: Dimensionality reduction**

Code to reduce dimensionality of a dataset. Includes a heatmap titled "Dimensionality plane and linear model results, dimensionality filtered data".

# TOOL 1 TER: 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, ...

## GOOD PRACTICE #2

### CONTROLLING SOFTWARE ENVIRONMENT

---

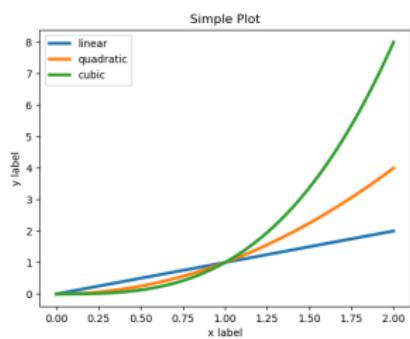
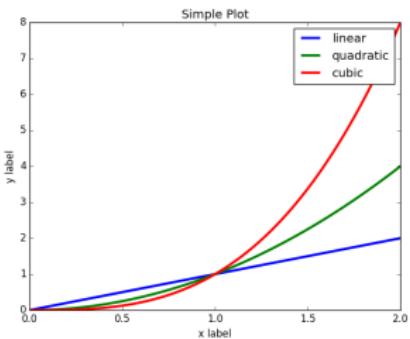
# ARGH... DAMNED COMPUTERS

- Alice: I got 3.123123 Bob: I got segfault
- 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 `libFoo-1.5c` and `pip install blah` but I broke everything and now neither his code nor mine works! 😞
- But hey! Here is my code. It's on GitHub so feel free to play with it! I'm doing open science 😊
  1. No one will ever run/use your code if it isn't easy to install
  2. No one will ever manage to run your code if you don't document how to run it
  3. Others (even you) are unlikely to get the same results unless you control and share your software environment

## SOFTWARE DEPENDENCIES: HORROR STORIES

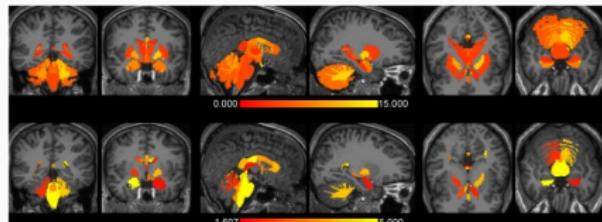
# SOFTWARE DEPENDENCIES: HORROR STORIES

- Software environment evolution



# SOFTWARE DEPENDENCIES: HORROR STORIES

- Software environment evolution
- OS heterogeneity



The Effects of FreeSurfer Version, Workstation Type, and Macintosh Operating System Version on Anatomical Volume and Cortical Thickness Measurements (PLOS ONE, 2012)

*Significant differences in volume and cortical thickness were revealed across FreeSurfer versions:*

- volume:  $8.8 \pm 6.6\%$  (range 1.3-**64.0%**)
- cortical thickness:  $2.8 \pm 1.3\%$  (range 1.1-7.7%)

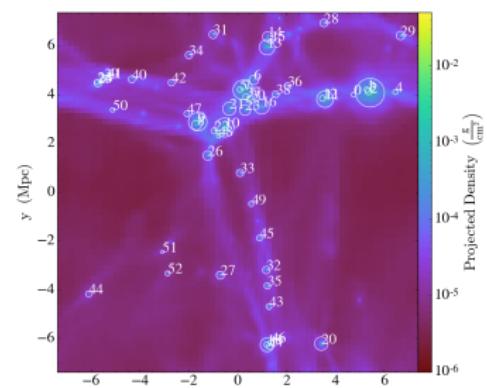
*About a factor two smaller differences were found between the Mac and HP workstations and between Mac OSX 10.5 and OSX 10.6.*

*In the context of an ongoing study, users are discouraged to update to a new major release of either FreeSurfer or operating system.*

*Formal assessment of the accuracy of FreeSurfer is desirable.*

# SOFTWARE DEPENDENCIES: HORROR STORIES

- Software environment evolution
- OS heterogeneity
- Impact of the compiler

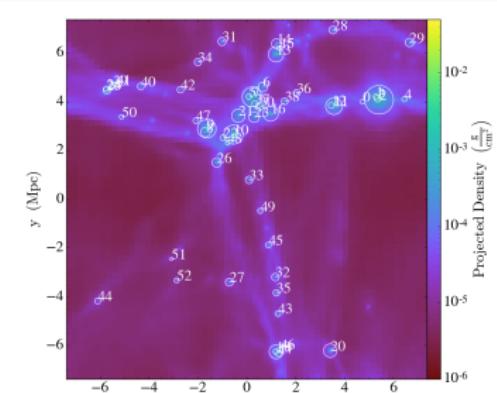


Assessing Reproducibility: An Astrophysical Example of Computational Uncertainty in the HPC Context (ResCuE-HPC, 2018)

Compiler	Optim.	Largest Halo Avg Mass.	Std. Err	Walltime
gcc@6.2.0	None	2.273E 46	1.069E 44	22h

# SOFTWARE DEPENDENCIES: HORROR STORIES

- Software environment evolution
- OS heterogeneity
- Impact of the compiler

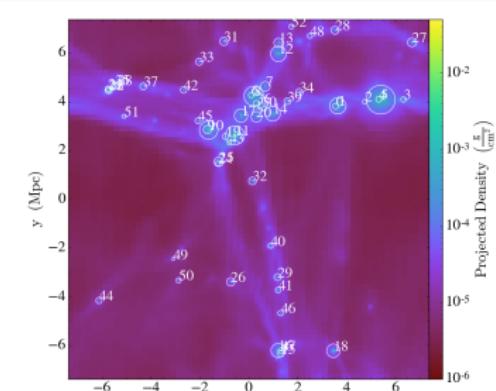


Assessing Reproducibility: An Astrophysical Example of Computational Uncertainty in the HPC Context (ResCuE-HPC, 2018)

Compiler	Optim.	Largest Halo Avg Mass.	Std. Err	Walltime
gcc@6.2.0	None	2.273E 46	1.069E 44	22h

# SOFTWARE DEPENDENCIES: HORROR STORIES

- Software environment evolution
- OS heterogeneity
- Impact of the compiler

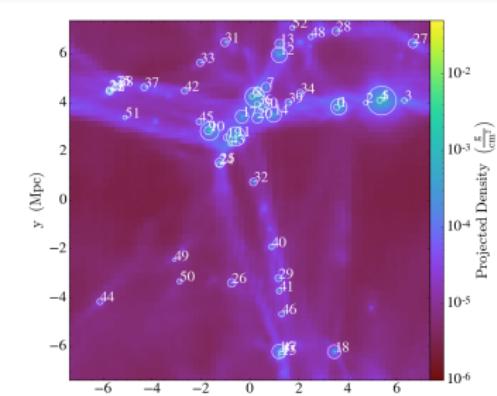


Assessing Reproducibility: An Astrophysical Example of Computational Uncertainty in the HPC Context (ResCuE-HPC, 2018)

Compiler	Optim.	Largest Halo		Walltime
		Avg Mass.	Std. Err	
gcc@6.2.0	None	2.273E 46	1.069E 44	22h
gcc@6.2.0	Normal	2.266E 46	1.218E 44	10h
gcc@6.2.0	High	2.275E 46	1.199E 44	9h

# SOFTWARE DEPENDENCIES: HORROR STORIES

- Software environment evolution
- OS heterogeneity
- Impact of the compiler



Assessing Reproducibility: An Astrophysical Example of Computational Uncertainty in the HPC Context (ResCuE-HPC, 2018)

Compiler	Optim.	Largest Halo		Walltime
		Avg Mass.	Std. Err	
gcc@6.2.0	None	2.273E 46	1.069E 44	22h
gcc@6.2.0	Normal	2.266E 46	1.218E 44	10h
gcc@6.2.0	High	2.275E 46	1.199E 44	9h
intel@16.0.3	None	<b>22.71</b> E 46	1.587E 44	39h
intel@16.0.3	Normal	<b>43.30</b> E 46	1.248E 44	7h
intel@16.0.3	High	2.268E 46	1.414E 44	6h
cce@8.5.5	Low	<b>43.11</b> E 46	1.353E 44	16h
cce@8.5.5	Normal	2.271E 46	1.261E 44	6h
cce@8.5.5	High	2.272E 46	1.341E 44	5h

# COMPLEX ECOSYSTEMS

```
1 import matplotlib  
2 print(matplotlib.__version__)
```

3.5.1

# COMPLEX ECOSYSTEMS

```
1 import matplotlib  
2 print(matplotlib.__version__)
```

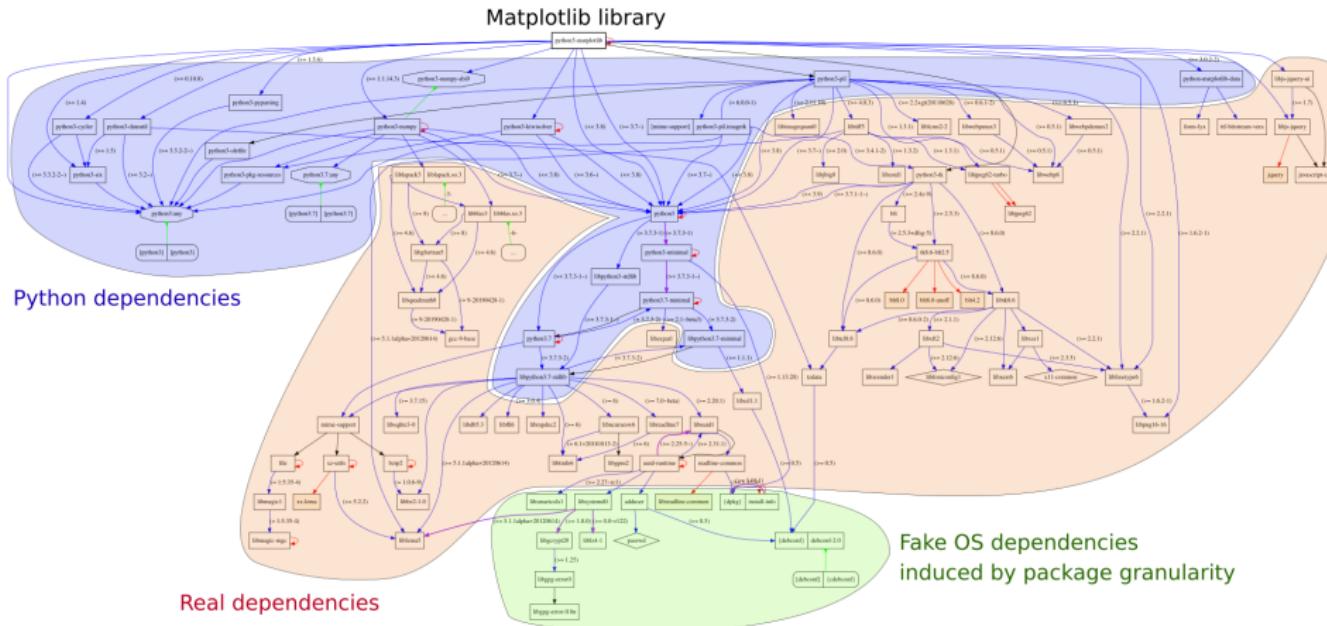
3.5.1

```
1 apt show python3-matplotlib
```

Package: python3-matplotlib  
Version: 3.5.1-2+b1  
Source: matplotlib (3.5.1-2)  
Maintainer: Sandro Tosi <morph@debian.org>  
Installed-Size: 27.6 MB  
Depends: libjs-jquery, libjs-jquery-ui, python-matplotlib-data (>= 3.5.1),  
 python3-dateutil, python3-pil.imagetk, python3-pyparsing (>= 1.5.6),  
 python3-six (>= 1.4), python3-numpy (>= 1:1.20.0), python3-numpy-abi9,  
 python3 (<< 3.11), python3 (>= 3.9~), python3-cycler (>= 0.10.0),  
 python3-fonttools, python3-kiwisolver, python3-packaging, python3-pil,  
 python3:any, libc6 (>= 2.29), libfreetype6 (>= 2.2.1),  
 libgcc-s1 (>= 3.3.1), libqhull-r8.0 (>= 2020.1), libstdc++6 (>= 11)  
Recommends: python3-tk  
Suggests: dvipng, ffmpeg, fonts-staypuft, ghostscript, gir1.2-gtk-3.0, inkscape,  
 ipython3, librsvg2-common, python-matplotlib-doc, python3-cairoffi,  
 python3-gi, python3-gi-cairo, python3-gobject, python3-pyqt5,  
 python3-scipy, python3-sip, python3-tornado, texlive-extra-utils,  
 texlive-latex-extra  
Enhances: ipython3

# COMPLEX ECOSYSTEMS

```
1 import matplotlib  
2 print(matplotlib.__version__)
```



python3-scipy, python3-sip, python3-tornado, texlive-extra-utils,  
texlive-latex-extra

Enhances: ipython3

# POTENTIAL SOLUTIONS: CONTAINERS AND PACKAGE MANAGERS

The good



The bad



The ugly



Automatic tracking

# POTENTIAL SOLUTIONS: CONTAINERS AND PACKAGE MANAGERS

The good



Guix



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>`

# POTENTIAL SOLUTIONS: 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)

# POTENTIAL SOLUTIONS: 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



# DEBIAN DEPENDENCIES

```
1 dpkg --status python3-matplotlib
```

```
Package: python3-matplotlib
Version: 3.6.3-1+b1
Source: matplotlib (3.6.3-1)
Depends: libjs-jquery, libjs-jquery-ui, python-matplotlib-data (>= 3.6.3),
          python3-dateutil, python3-pil.imagetk, python3-pyparsing (>= 1.5.6),
          python3-six (>= 1.4), python3-numpy (>= 1:1.22.0), python3-contourpy,
          python3 (<< 3.12), python3 (>= 3.11~), python3-numpy-abi9,
          python3-cycler (>= 0.10.0), python3-fonttools, python3-kiwisolver,
          python3-packaging, python3-pil, python3:any, libc6 (>= 2.34),
          libfreetype6 (>= 2.2.1), libgcc-s1 (>= 3.3.1),
          libqhull-r8.0 (>= 2020.1), libstdc++6 (>= 11)
```

# DEBIAN DEPENDENCIES

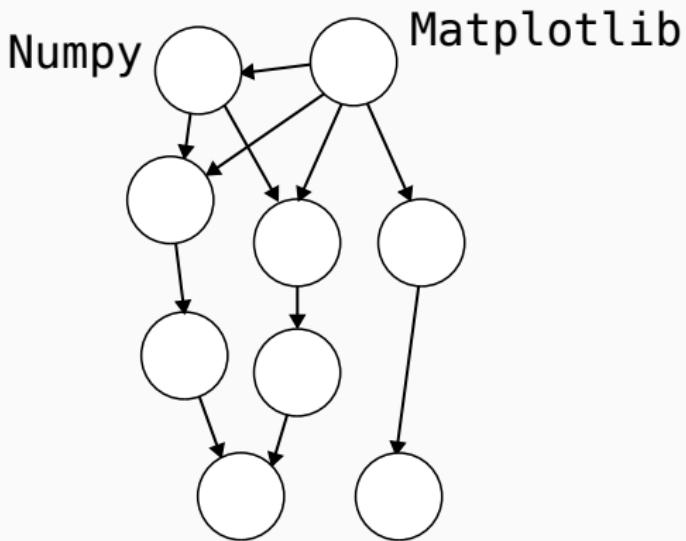
```
1 dpkg --status python3-matplotlib
```

```
Package: python3-matplotlib
Version: 3.6.3-1+b1
Source: matplotlib (3.6.3-1)
Depends: libjs-jquery, libjs-jquery-ui, python-matplotlib-data (>= 3.6.3),
          python3-dateutil, python3-pil.imagetk, python3-pyparsing (>= 1.5.6),
          python3-six (>= 1.4), python3-numpy (>= 1:1.22.0), python3-contourpy,
          python3 (<< 3.12), python3 (>= 3.11~), python3-numpy-abi9,
          python3-cycler (>= 0.10.0), python3-fonttools, python3-kiwisolver,
          python3-packaging, python3-pil, python3:any, libc6 (>= 2.34),
          libfreetype6 (>= 2.2.1), libgcc-s1 (>= 3.3.1),
          libqhull-r8.0 (>= 2020.1), libstdc++6 (>= 11)
```

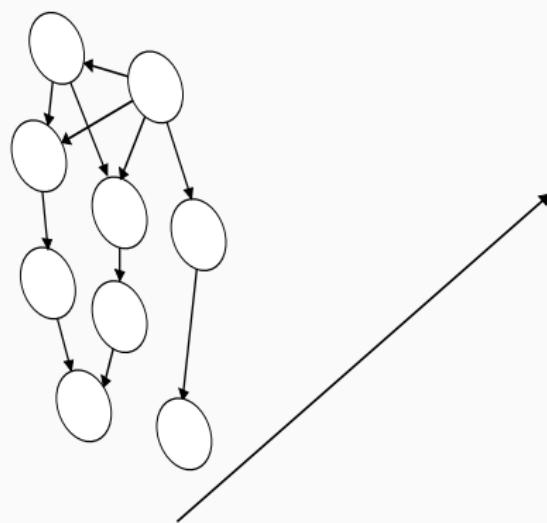
On a given day:

- Several versions of each package are available on the server
- Installing the latest version of a package may require upgrading some other packages

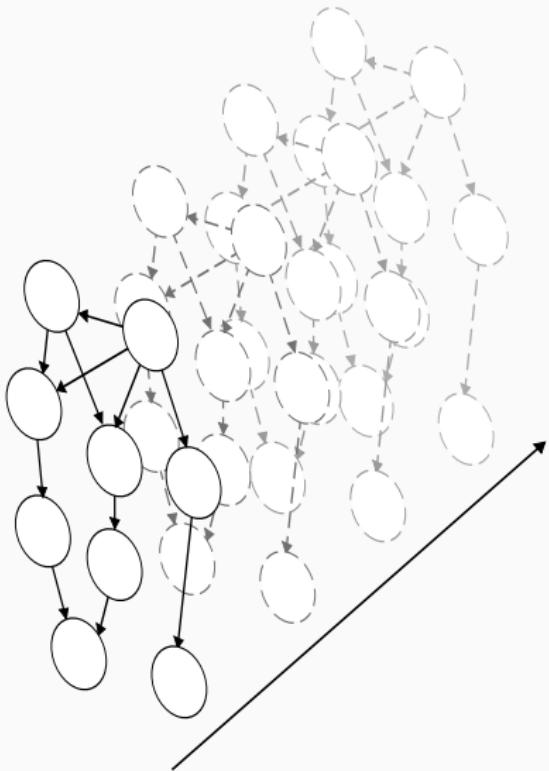
## LOOSE VS. STRICT DEPENDENCIES IN PICTURE



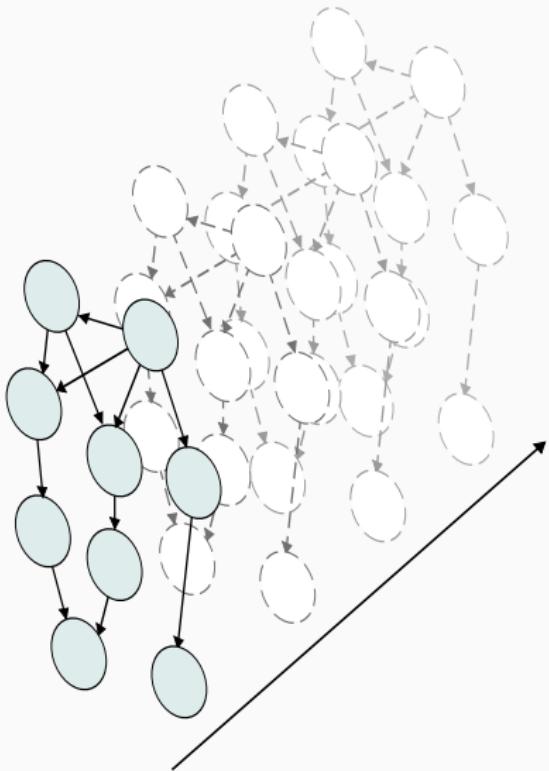
## LOOSE VS. STRICT DEPENDENCIES IN PICTURE



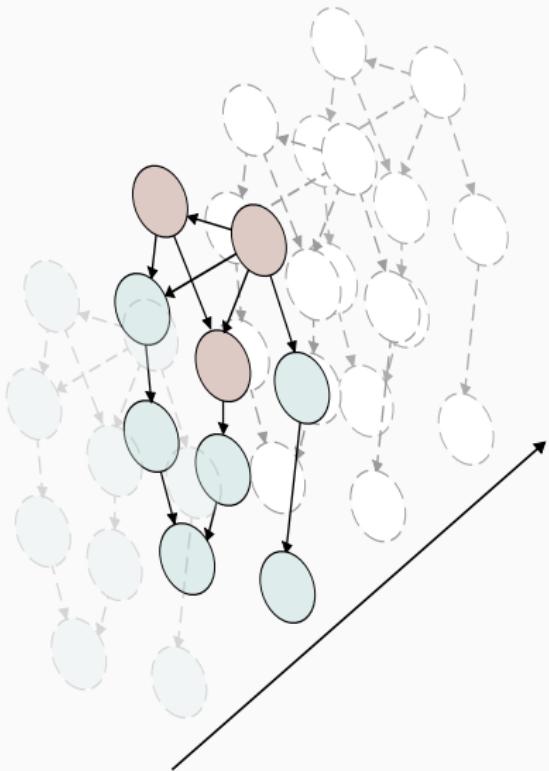
## LOOSE VS. STRICT DEPENDENCIES IN PICTURE



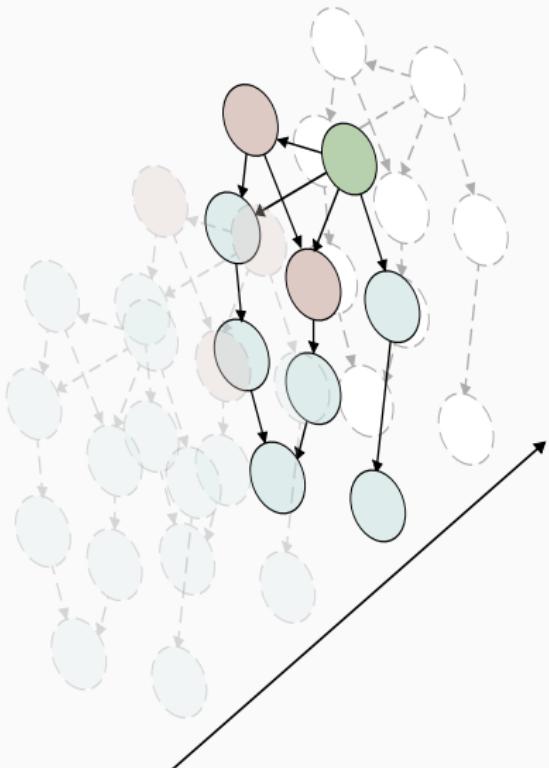
## LOOSE VS. STRICT DEPENDENCIES IN PICTURE



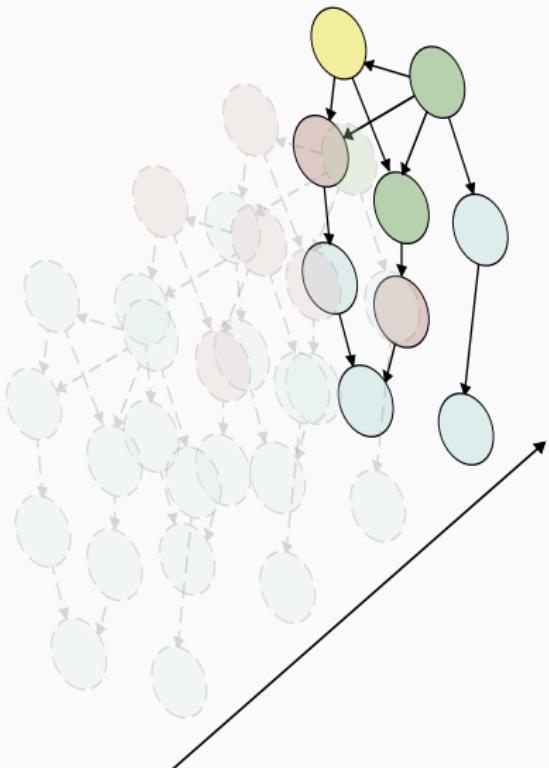
## LOOSE VS. STRICT DEPENDENCIES IN PICTURE



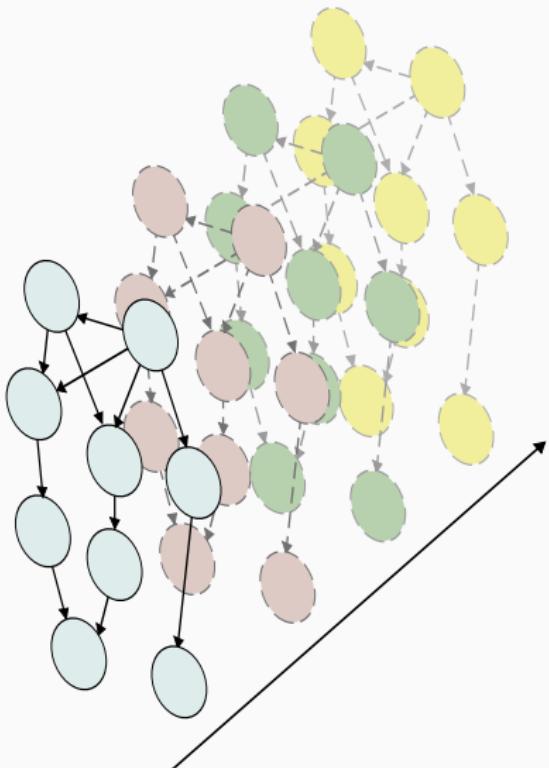
## LOOSE VS. STRICT DEPENDENCIES IN PICTURE



## LOOSE VS. STRICT DEPENDENCIES IN PICTURE



## LOOSE VS. STRICT DEPENDENCIES IN PICTURE



## GUIX IN A NUTSHELL (1/2)

GUIX is not a ~~container technology~~, it is meant for **package management!**

## GUIX IN A NUTSHELL (1/2)

GUIX is not a ~~container technology~~, it is meant for **package management!**

- A GUIX **recipe** (aka `manifest.scm`)

```
1 (specifications->manifest
 2   (list "hello" "coreutils"))
```

```
1 guix shell -C -m manifest.scm -- ls /    # -C = --container
```

## GUIX IN A NUTSHELL (1/2)

GUIX is not a ~~container technology~~, it is meant for **package management!**

- A GUIX **recipe** (aka **manifest.scm**)

```
1 (specifications->manifest
2   (list "hello" "coreutils"))

1 guix shell -C -m manifest.scm -- ls /  # -C = --container
```

- The description of the **versions** is in the **manifest.scm**

```
1 guix describe -f channels > channels.scm
```

```
1 (list (channel
2       (name 'guix)
3       (url "https://git.savannah.gnu.org/git/guix.git")
4       (branch "master")
5       (commit
6         "d09a4cc7c739f4201821623f62c69bcef9c20f52")
7       (introduction
8         (make-channel-introduction
9           "9edb3f66fd807b096b48283debdcccfea34bad"
10          (openpgp-fingerprint
11            "BBB0 2DDF 2CEA F6A8 0D1D E643 A2A0 6DF2 A33A 54FA")))))
```

# GUIX IN A NUTSHELL (1/2)

GUIX is not a ~~container technology~~, it is meant for **package management!**

- A GUIX **recipe** (aka `manifest.scm`)

```
1 (specifications->manifest
2   (list "hello" "coreutils"))

1 guix shell -C -m manifest.scm -- ls / # -C = --container
```

- The description of the **versions** is in the `manifest.scm`

```
1 guix describe -f channels > channels.scm
```

```
1 (list (channel
2       (name 'guix)
3       (url "https://git.savannah.gnu.org/git/guix.git")
4       (branch "master")
5       (commit
6         "d09a4cc7c739f4201821623f62c69bcef9c20f52")
7       (introduction
8         (make-channel-introduction
9           "9edb3f66fd807b096b48283debdcccfea34bad"
10          (openpgp-fingerprint
11            "BBB0 2DDF 2CEA F6A8 0D1D E643 A2A0 6DF2 A33A 54FA")))))
```

- A **time-machine**

```
1 guix time-machine -C channels.scm -- shell -m manifest.scm -- ls /
```

## GUIX IN A NUTSHELL (2/2)

- Under the hood:
  - A **deamon** compiles everything in a collection of directories  
`/gnu/store/8fpk2cja3f07xls48jfnpgrzrljpqivr-coreutils-8.32/`
  - All the directories are assembled (with symlinks) in a  
`/gnu/store/j5964hh821p2h5mcadpvj16l1m9330gv-profile/` dir
  - Environment variables (**PATH**, **LD\_LIBRARY\_PATH**, ...) are updated accordingly

## GUIX IN A NUTSHELL (2/2)

- Under the hood:
  - A **deamon** compiles everything in a collection of directories  
`/gnu/store/8fpk2cja3f07xls48jfnpgrzrljpqivr-coreutils-8.32/`
  - All the directories are assembled (with symlinks) in a  
`/gnu/store/j5964hh821p2h5mcadpvj16l1m9330gv-profile/` dir
  - Environment variables (**PATH**, **LD\_LIBRARY\_PATH**, ...) are updated accordingly
- Several **containerization** options

```
1 guix shell --container coreutils -- ls
2 guix shell coreutils -- ls  # Fully permeable: expends the env, ...
```

## GUIX IN A NUTSHELL (2/2)

- Under the hood:
  - A **deamon** compiles everything in a collection of directories  
`/gnu/store/8fpk2cja3f07xls48jfnpgrzrljpqivr-coreutils-8.32/`
  - All the directories are assembled (with symlinks) in a  
`/gnu/store/j5964hh821p2h5mcadpvj16l1m9330gv-profile/` dir
  - Environment variables (**PATH**, **LD\_LIBRARY\_PATH**, ...) are updated accordingly
- Several **containerization** options

```
1 guix shell --container coreutils -- ls
2 guix shell coreutils -- ls # Fully permeable: expends the env, ...
```

- Various **export** formats (**docker**, **squashfs**, **debian**, **tarball**, **module**, **relocatable**...)

```
1 guix pack --format=docker --save-provenance -m manifest.scm
```

Allows to carefully control/nest environments and adapt to your context

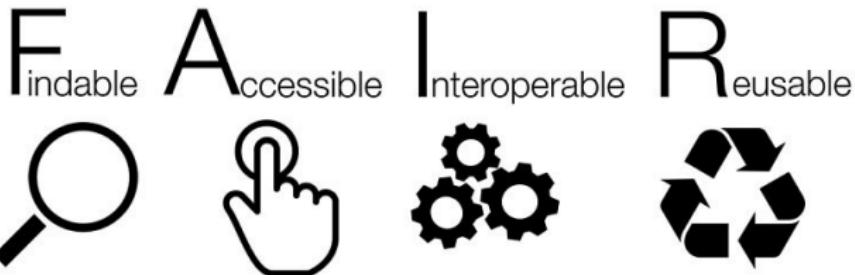
- For more information: <https://hpc.guix.info/>
  - Lastest Guix Workshop on HPC was held in Bordeaux on Nov. 7, 2024, after the JCAD
  - Check out the **Café Guix**

## GOOD PRACTICE #3

## VERSION CONTROL AND ARCHIVING

---

# FAIR PRINCIPLES



<https://www.go-fair.org/fair-principles/>

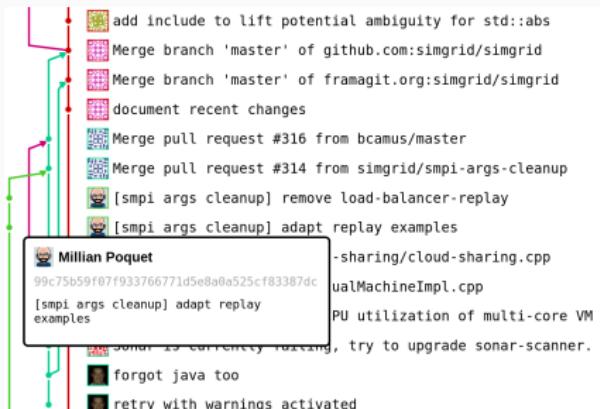
- "*Open as much as possible and close as much as necessary*"
- Management, publication, annotation (metadata), archiving
- Source code = specific data with specific consideration

Let's go beyond general principles!

# TOOL 3: VERSION CONTROL AND FORGE

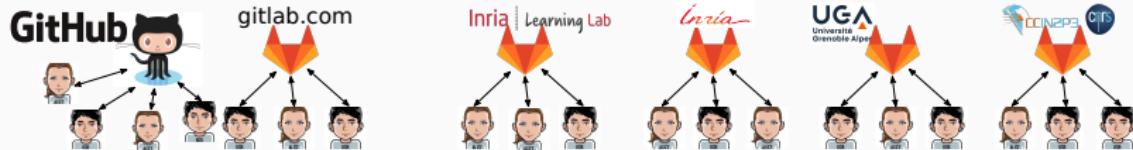
## Git = version control

- Developed in 2005 by Linus Torvalds for the kernel development
- Local and efficient rollbacks
- Distributed: everyone has a full copy of the history



## GitHub, GitLab, and Co

- Free hosting of public projects, social network



## Limitation

- Managing large data: **Git LFS**   **Git Annex** (or DataLad)

## TOOL 3BIS: 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 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*
- Discontinued forges: Code Space, Gitorious, Google code, Inria Gforge

## TOOL 3BIS: 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 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*
- Discontinued forges: Code Space, Gitorious, Google code, Inria Gforge

Article archives



Data archives



figshare



Software Archive



Software Heritage

Collect/Preserve/Share

CONTROLLING THE WHOLE  
SOFTWARE/COMPIILING STACK IS NOT  
SUFFICIENT

---

# 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]

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

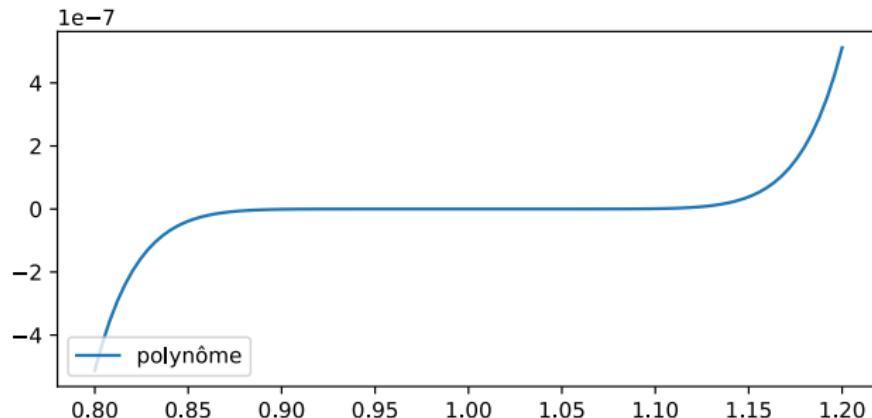
Evaluation of  $f(77617.0, 33096.0)$

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

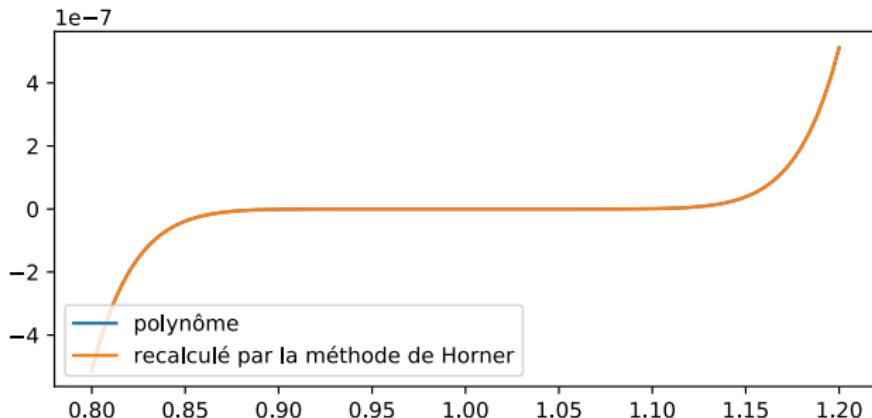
At scale (complex code + non-determinism), all this can become particularly harmful and painful.

# ALL I CARE ABOUT IS THE ALGORITHM OUTPUT (FP)



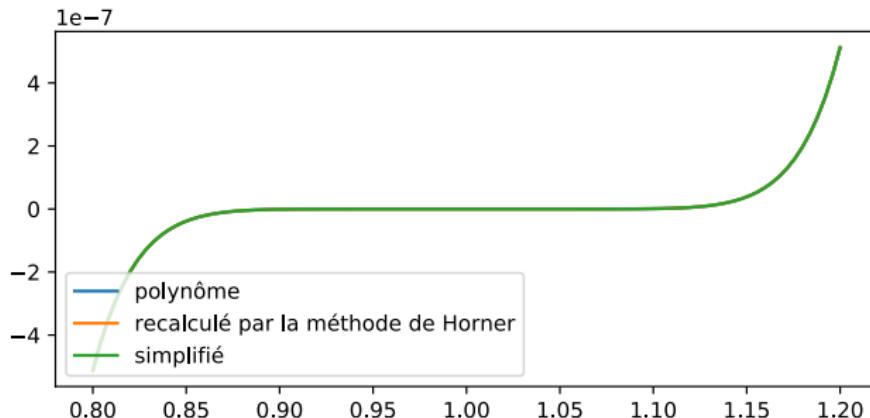
```
1 def polynome(x):  
2     return x**9 - 9.*x**8 + 36.*x**7 - 84.*x**6 + 126.*x**5 \  
3         - 126.*x**4 + 84.*x**3 - 36.*x**2 + 9.*x - 1.
```

# FLOATING-POINT ARITHMETIC



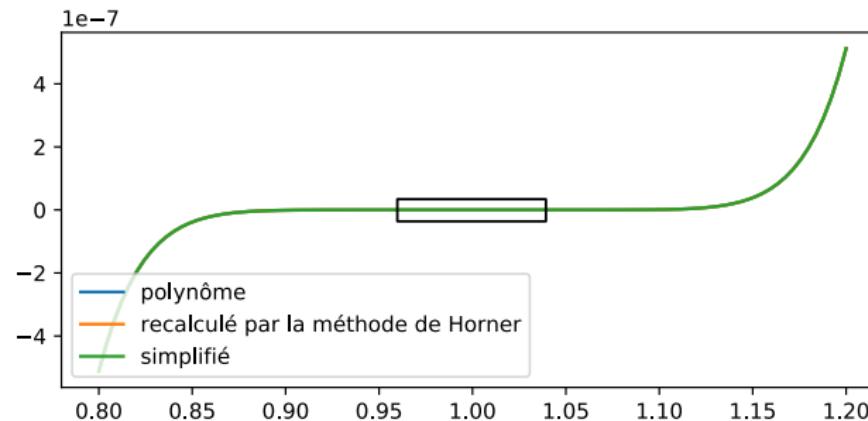
```
1 def horner(x):  
2     return x*(x*(x*(x*(x*(x*(x - 9.) + 36.) - 84.) + 126.) \  
3             - 126.) + 84.) - 36.) + 9.) - 1.
```

# FLOATING-POINT ARITHMETIC

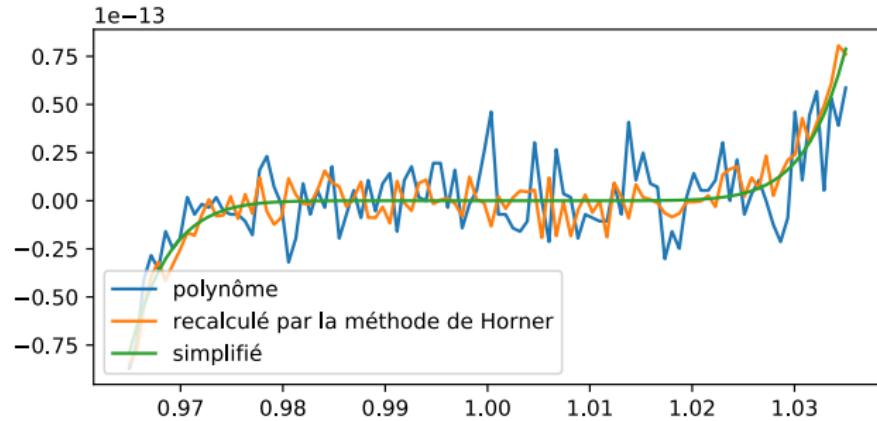


```
1 def simple(x):  
2     return (x-1.)**9  
3 # Easy! ;)
```

# FLOATING-POINT ARITHMETIC



# FLOATING-POINT ARITHMETIC



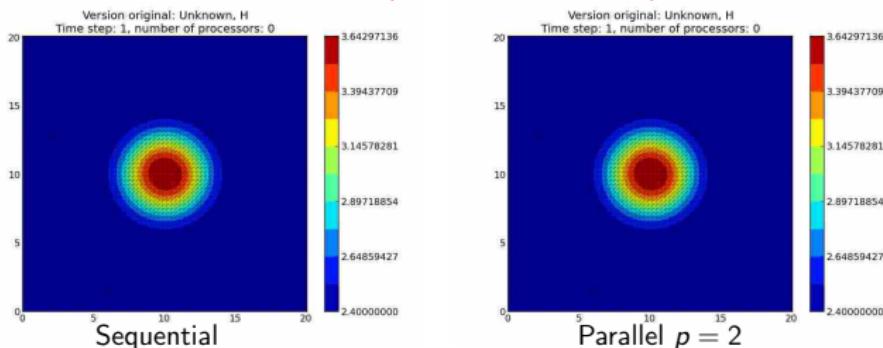
# DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

## Telemac2D: the simplest gouttedeo simulation

### The gouttedeo 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, ...)



13 / 64

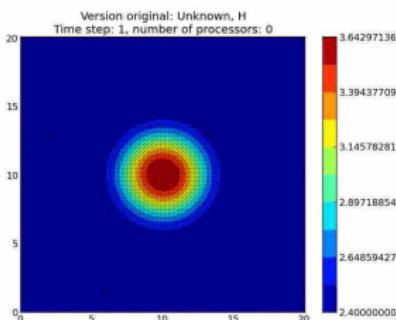
Courtesy of P. Langlois and R. Nheili

# DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

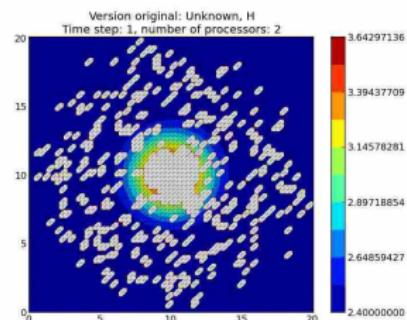
A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 1



Sequential



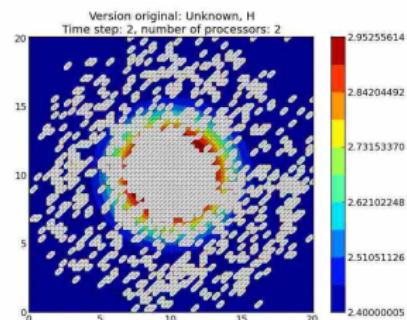
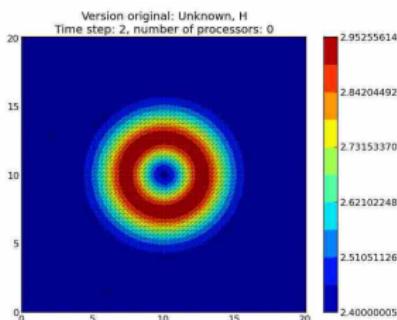
Parallel  $p = 2$

# DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 2

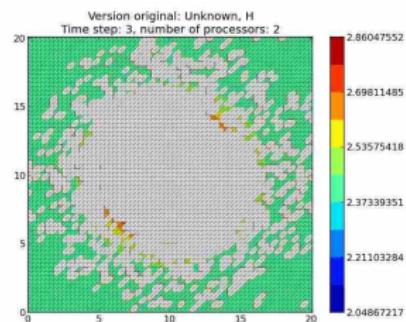
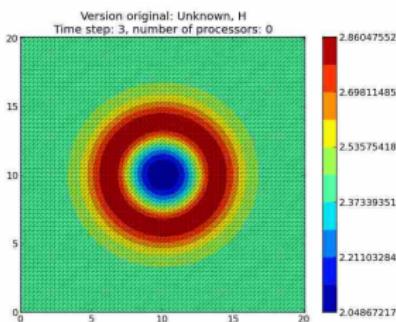


# DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 3

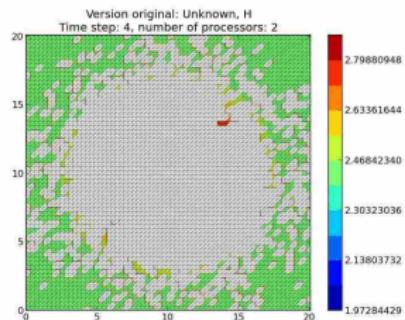
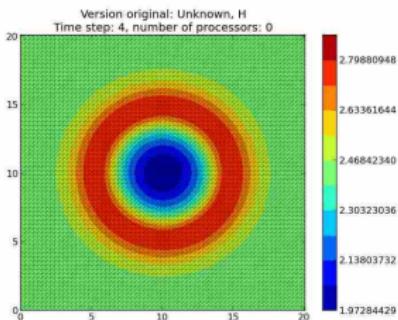


# DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 4

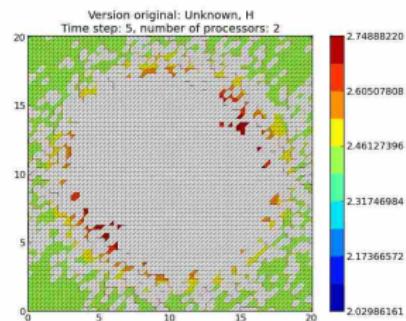
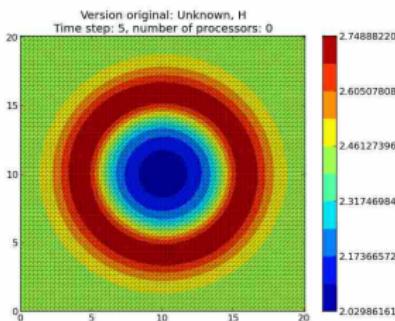


# DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 5

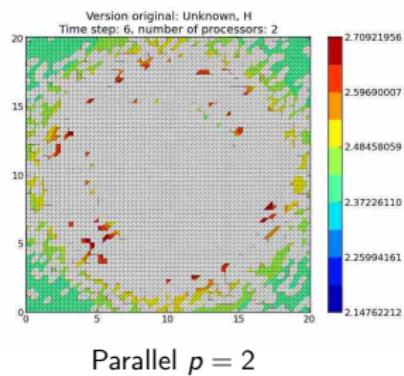
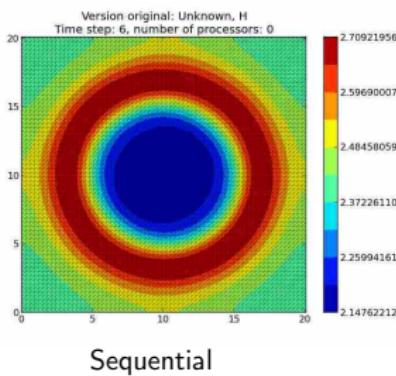


# DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 6

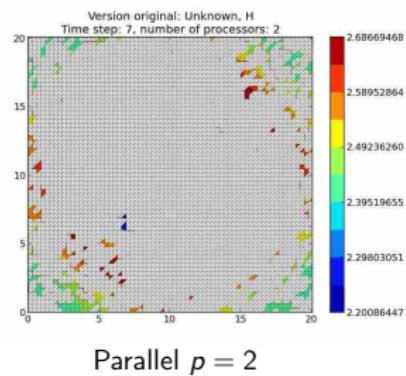
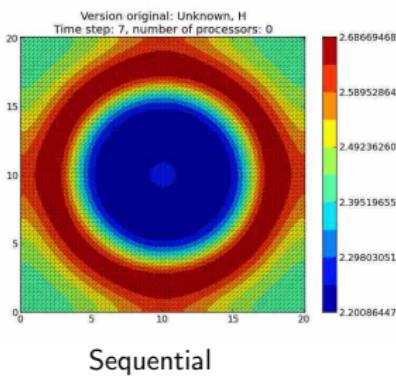


# DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 7

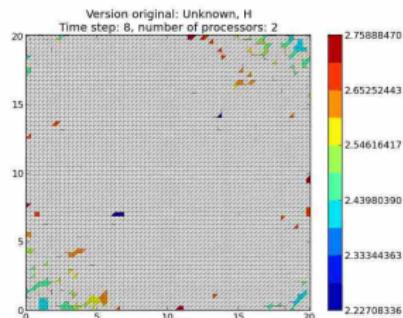
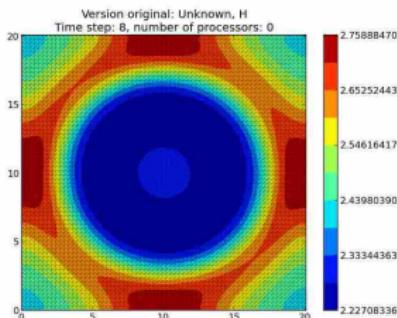


# DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 8

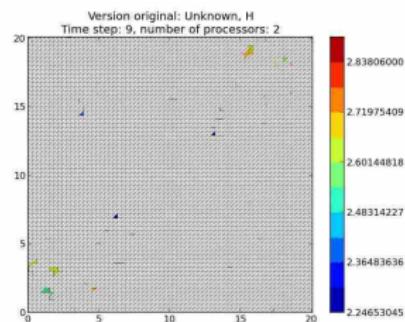
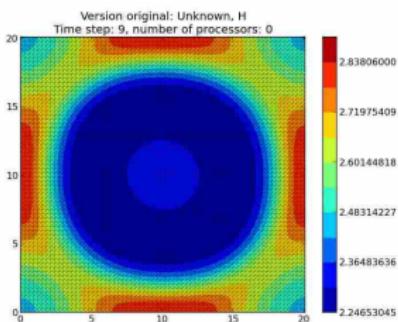


# DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 9

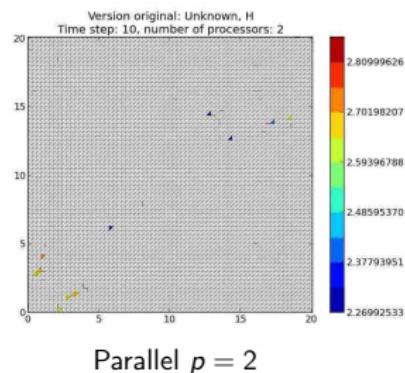
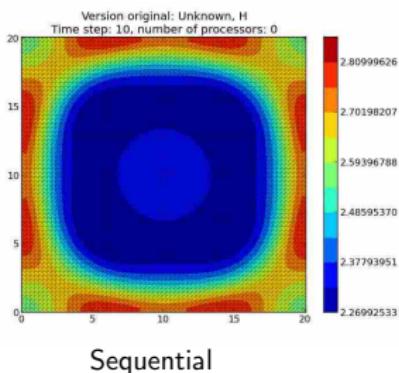


# DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 10

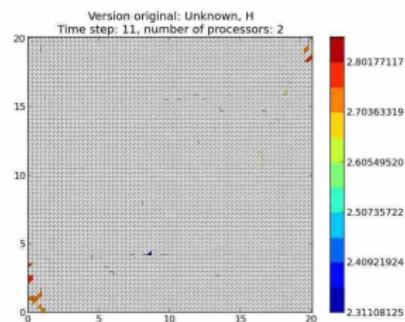
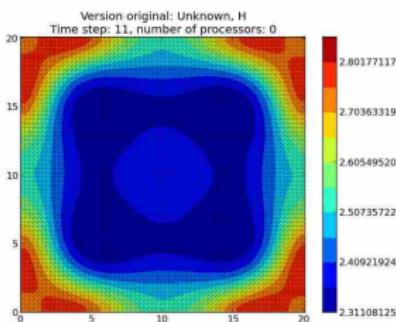


# DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 11

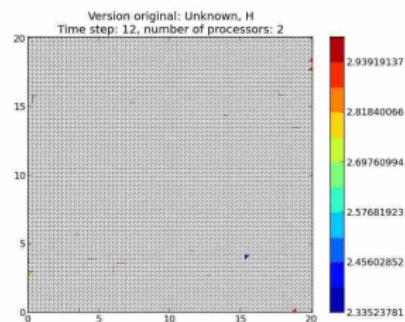
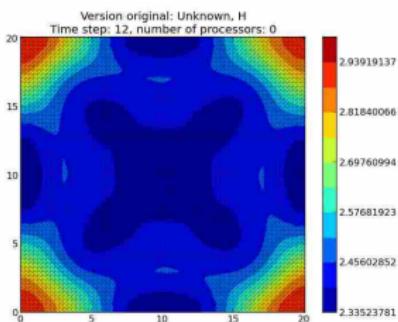


# DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 12

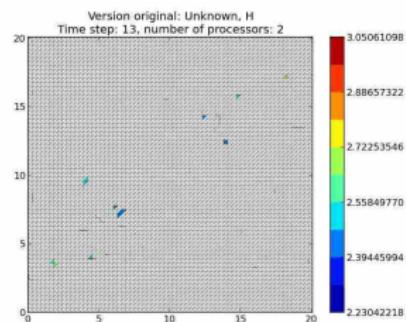
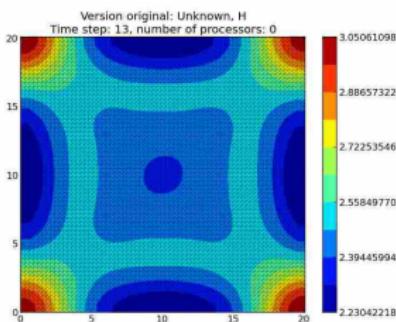


# DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 13

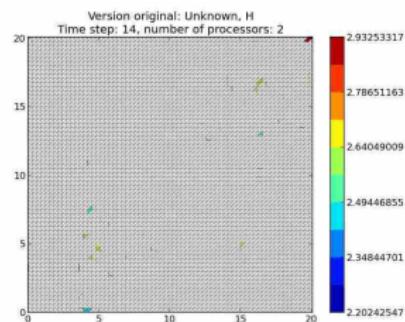
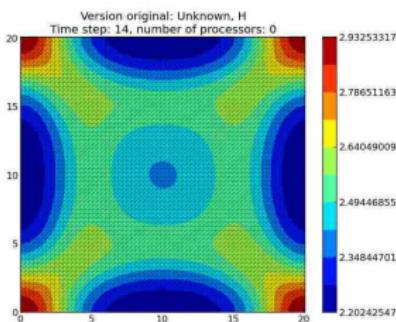


# DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

Numerical reproducibility?

time step = 14

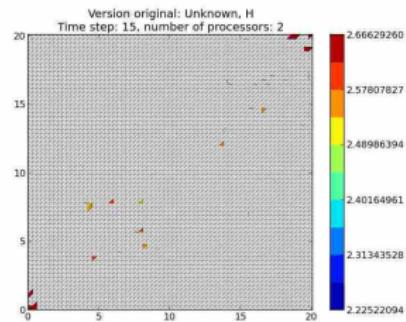
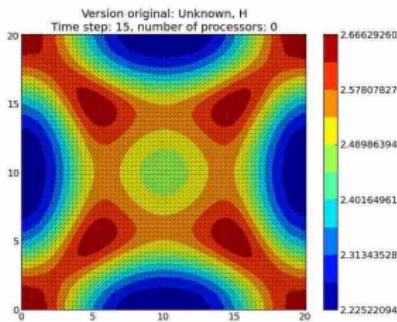


# DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

A white plot displays a non-reproducible value

NO numerical reproducibility!

time step = 15

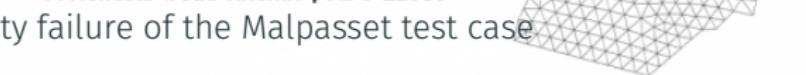


# DID I MENTION WE HAVE PARALLEL MACHINES NOWADAYS?

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

ADVERTISEMENT

---

# MOOCs on REPRODUCIBLE RESEARCH:

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

- Konrad Hinsen, Christophe Pouzat
- 3rd Edition: March 2020 – ... (16,800+)
- Notebooks, version control, simple data formats



# MOOCs on REPRODUCIBLE RESEARCH:

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

- Konrad Hinsen, Christophe Pouzat
- 3rd Edition: March 2020 – ... (16,800+)
- Notebooks, version control, simple data formats



MOOC Reproducible Research II: Practices and tools for managing computations and data (May-Sep 2024, ≈ 2,000)

- Managing data (`FITS/HDF5, git annex`)
- Software environment control (`docker, singularity, guix`)
- Scientific workflow (`make, snakemake`)
- Statistics, Numerical Chaos

WHAT'S THE POINT ?

---

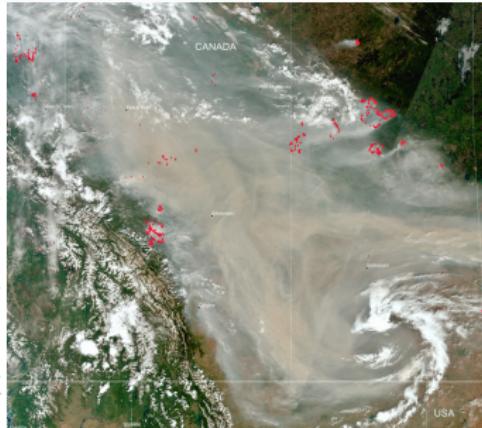
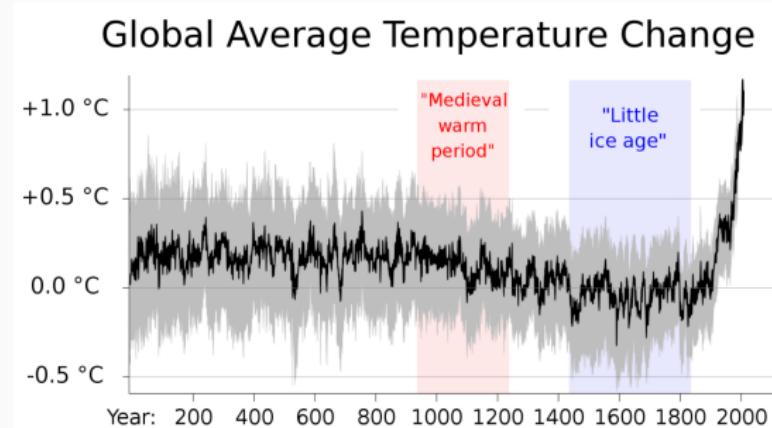
## THE SCIENCE IS CLEAR

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](https://en.wikipedia.org/wiki/Global_temperature_record)

2023 Alberta wildfires (> 1 Mha)

## THE SCIENCE IS CLEAR

scientist rebellion

Why are we ignoring it?

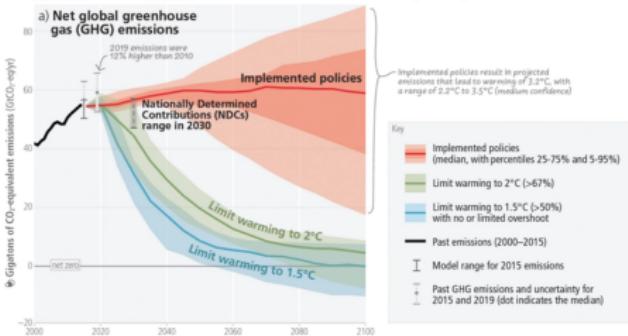


IPCC, IPBES, <https://climate.nasa.gov/>

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

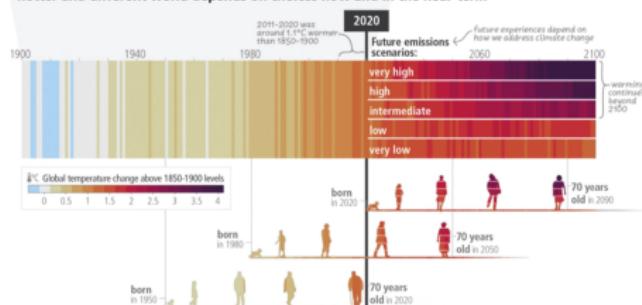
Limiting warming to **1.5°C** and **2°C** involves rapid, deep and in most cases immediate greenhouse gas emission reductions

Net zero: CO<sub>2</sub> and net zero GHG emissions can be achieved through strong reductions across all sectors



Paris Agreement'15 ~ Net Zero by 2050

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

39/44

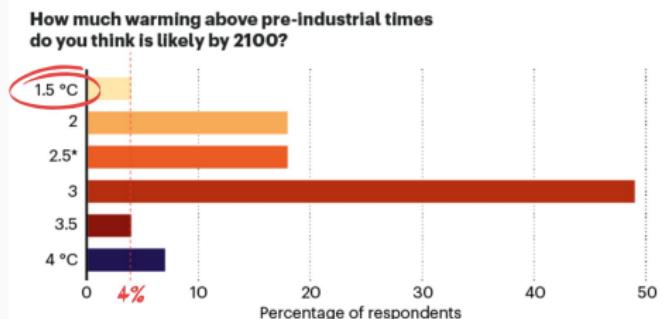
## THE SCIENCE IS CLEAR

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**



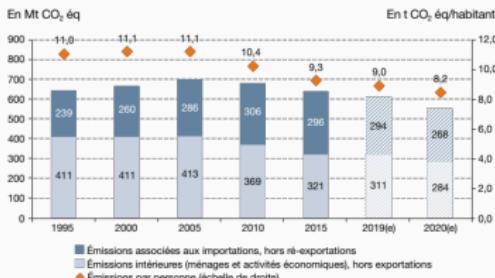
\*Includes 2 responses between 2.7 °C and 2.75 °C; 2.5 °C and 3.5 °C were write-in answers.

@natu Nature survey, Nov. 2021

# THE ELEPHANT IN THE ROOM: CLIMATE CHANGE

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

ÉVOLUTION DE L'EMPREINTE CARBONE DE LA FRANCE



(e) = estimations.

Note : l'empreinte carbone porte sur les trois principaux gaz à effet de serre (CO<sub>2</sub>, CH<sub>4</sub>, N<sub>2</sub>O). En 2021, la méthodologie a été ajustée afin de mieux tenir compte de l'évolution des coûts du pétrole brut, du gaz et du charbon. L'ensemble de la série a ainsi été révisé, l'essentiel des ajustements portant sur les émissions importées de CH<sub>4</sub>.

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

Sources : Citepa ; AIE ; FAO ; Douanes ; Eurostat ; Insee. Traitement : SDES, 2021

Empreinte carbone moyenne en France  
10 tonnes de CO<sub>2</sub>e/an/pers.



÷2  
d'ici  
2030

<2t CO<sub>2</sub>e

Objectif d'ici 2050

- de 2 t de CO<sub>2</sub>e/an/pers.

+ Faire plus d'activités bas carbone !

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

Bref, inventer nos vies bas carbone désirables !

Par exemple :

0,5 t CO<sub>2</sub>e/Annee : À la maison : préférence pour les produits végétaliens

0,5 t CO<sub>2</sub>e/Annee : Transport : 2000 km en voiture (3000 km de fabrication amortie sur 30 ans, risques de pollution et de dégradation de l'environnement dans les transports en commun)

0,5 t CO<sub>2</sub>e/Annee : Consommation : Utiliser rien de neuf, réutiliser, recycler, faire diverses expérimentations dans les achats et les commandes en ligne

0,2 t CO<sub>2</sub>e/Annee : L'agriculture : Choisir bio sur un îlot (PCP ou paupier, 100% bio) et d'un agriculteur local et durable, privilier la châtaigne ou solaire thermique

0,2 t CO<sub>2</sub>e/Annee : Services publics : faire enseignement, éducation, recherche, culture, sport, etc.

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

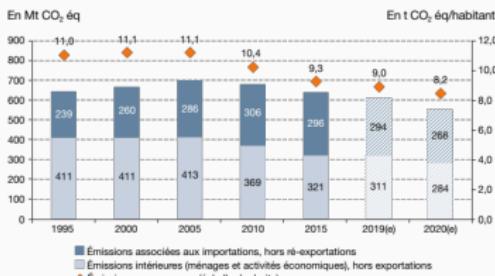
INVENTONS  
NOS VIES  
BAS CARBONE

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...

ÉVOLUTION DE L'EMPREINTE CARBONE DE LA FRANCE



Empreinte carbone moyenne en France  
10 tonnes de CO<sub>2</sub>e/an/pers.



÷2  
d'ici  
2030

<2t CO<sub>2</sub>e

Objectif d'ici 2050

- de 2 t de CO<sub>2</sub>e/an/pers.

+ Faire plus d'activités bas carbone !

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

Bref, inventer nos vies bas carbone désirables !

Par exemple :

1. Émissions Alimentation : À tendance régulière, l'empreinte se rapproche progressivement

2. Émissions Transport : 2000 km en voiture routière (80%) de fabrication amortie sur 30 ans, importations et transports maritimes (15%) et transports en commun (15%).

3. Émissions Construction : Basé sur des matériaux et technologies de construction, faire évoluer les équipements dédiés au chauffage et à l'éclairage, améliorer dans les bâtiments et les infrastructures.

4. Émissions Gestion : Chaque ligne sur un îlot (PPC) en paix, 10% de l'effort d'un logement basé isolé dans une zone de campagne ou à la périphérie d'une ville ou d'une agglomération.

5. Émissions Lagement : Chaque ligne sur un îlot (PPC) en paix, 10% de l'effort d'un logement basé isolé dans une zone de campagne ou à la périphérie d'une ville ou d'une agglomération.

6. Émissions Services publics : Bureaux, enseignement, culture, sport, etc.

7. Émissions Ménages : Consommation de biens et services.

8. Émissions Importations : Importations de biens et services.

9. Émissions Exportations : Exportations de biens et services.

10. Émissions Services publics : Bureaux, enseignement, culture, sport, etc.

11. Émissions Ménages : Consommation de biens et services.

12. Émissions Importations : Importations de biens et services.

13. Émissions Exportations : Exportations de biens et services.

14. Émissions Services publics : Bureaux, enseignement,



<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

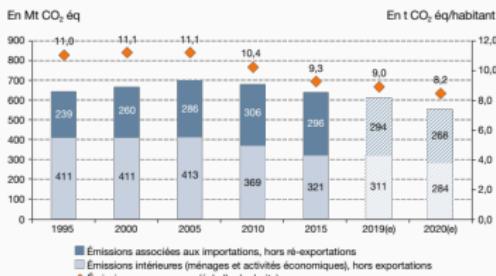
- Verdissement de l'industrie: « pause » sur les normes environnementales
- Loi de programmation militaire (+41%)
- Nous devons préparer la France à une élévation de la température de 4 °C
- Academia ? PEPR 5G, Cloud, NUMPEX, Quantique, IA, Agroécologie et numérique



# THE ELEPHANT IN THE ROOM: CLIMATE CHANGE

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

ÉVOLUTION DE L'EMPREINTE CARBONE DE LA FRANCE



Empreinte carbone moyenne en France  
10 tonnes de CO<sub>2</sub>e/an/pers.



÷2  
d'ici  
2030

Objectif d'ici 2050  
- de 2 t de CO<sub>2</sub>e/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

## Several scenarios on the table

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