

# Enjeux et défis de la recherche reproductible

Konrad HINSEN

Centre de Biophysique Moléculaire, Orléans, France  
et  
Synchrotron SOLEIL, Saint Aubin, France

23 mai 2019

# Pourquoi tout le monde parle de reproductibilité ?



Royal Society of London for Improving Natural Knowledge (1663)

# Pourquoi tout le monde parle de reproductibilité ?

**Reproductibilité** : preuve de **rigueur** qui inspire **confiance**

Un résultat non-reproductible suggère...

- ... une description incomplète
- ... une maîtrise insuffisante des techniques
- ... une erreur
- ... une fraude

Ou encore...

- ... une prise insuffisante sur les sujets d'étude (souris, étoiles, ...)

- Quelle est le rôle de la reproductibilité pour la méthode scientifique ?

- Quelle est le rôle de la reproductibilité pour la méthode scientifique ?

Ma réponse : c'est un accélérateur important

- Quelle est le rôle de la reproductibilité pour la méthode scientifique ?

Ma réponse : c'est un accélérateur important

- Quel effort est raisonnable pour améliorer la reproductibilité ?
- Quel effort est raisonnable pour vérifier la reproductibilité ?
- Comment motiver les chercheurs à faire plus d'effort ?
- Comment réduire l'effort nécessaire par la technologie ?

# Du technique au scientifique

## Reproductibilité

- Niveau technique
- Est-ce bien fait ?
- Évaluation simple
- Réponse simple : oui/non
- **Vérification**

## Réplicabilité

- Niveau scientifique
- Est-ce la bonne chose à faire ?
- Évaluation laborieuse
- Réponse complexe : si...
- **Validation**

# Du technique au scientifique

## Reproductibilité

- Niveau technique
- Est-ce bien fait ?
- Évaluation simple
- Réponse simple : oui/non
- **Vérification**

## Réplicabilité

- Niveau scientifique
- Est-ce la bonne chose à faire ?
- Évaluation laborieuse
- Réponse complexe : si...
- **Validation**

## En calcul scientifique :

- Même logiciel
- Mêmes paramètres
- Mêmes données
- **Résultat identique ?**

- Nouveau logiciel
- Mêmes paramètres
- Mêmes données (ou pas)
- **Résultat équivalent ?**



# Les enjeux de la reproductibilité en calcul scientifique

- Être sûr de ce qui a été calculé.
- Pouvoir vérifier le calcul.
- Pouvoir adapter le calcul.

Particularité du calcul : **déterminisme**

# Exercices d'échauffement

Pour pouvoir vérifier et adapter un calcul, il faut évidemment :

- 1 Avoir le code source et les données d'entrée.
- 2 Être sûr que c'est bien ce code source qui a été utilisé.  
→ pouvoir recompiler / réinstaller.

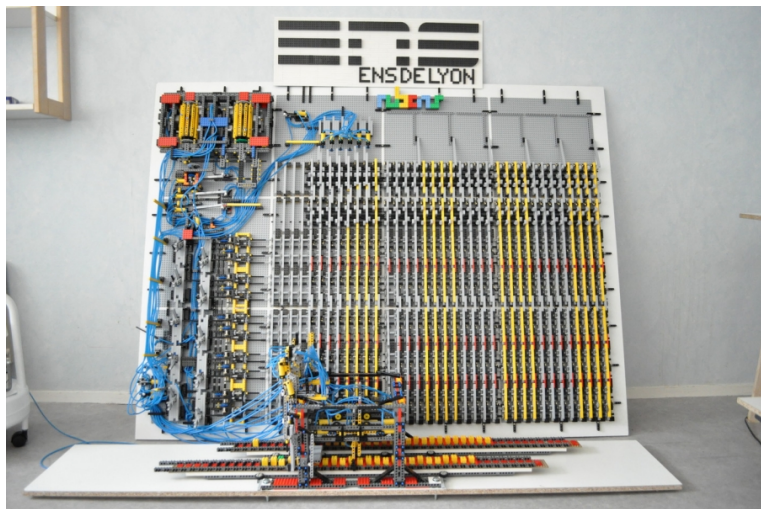
Il est avantageux d'avoir en plus

- 3 L'historique du code source
- 4 Des tests et exemples d'utilisation

## Mots clés :

sauvegardes, archives (Zenodo, Software Heritage, ...),  
gestion de versions, documentation, tests

# C'est quoi un calcul ?



Vue globale de la machine de Turing en Légo du [projet Rubens](#).

# C'est quoi un calcul ?

## Input

```
100111100001001100110101101100
001010011101010111110001001101
010111101100011110111011110001
001100001110111000100100000111
110101100111001110100000100110
11011110011100001111101101111
111001001011110001100110000101
011100001000010001011110000010
110101110011101111001010100111
111000101110011001101101001001
011001010100101011000001001100
11010011100101111100001011101
011110111110001111011110101101
000001110110011001010101011100
100010110001100000111001100010
000000111011100100100101010111
000010000001100001000010110110
101111101111000111100101110101
100101010100001001110100010001
011110011010100101111011110101
100011000110110001011101100110
110100000100000011011000001101
100000011100100111101101011011
010110010001000101110111001010
```



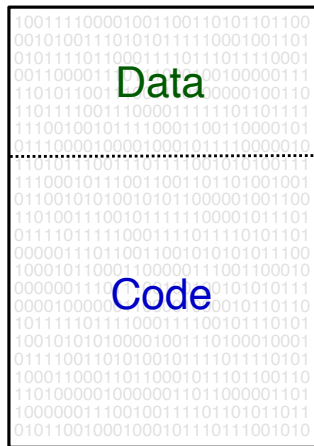
## Output

```
0000001110110010100001001110111
110000101110111011111110101010
000110010101111100101110110100
001110110011010110000101011010
111101111100000100010111010111
111010001000010010111100111001
111001100101000111101000011100
10111110000011011011011110001
100100110111101111000101010100
111110011010111011010011011100
111011100011110101011111000100
010111011010100100011110100011
001111000001111110001011100111
101101100000100011100111110011
001101000010011000110011000011
10101110111101010000011010001
010111100101010010011100011011
0010101011100101000001010000110
100000101001110011010000011100
0011100110001111111111000001100
100100010100000110001011010000
010110010111101001000010100010
101011110001001001010010111000
011000100000010000000011100111
```

Computer by Creative Stall from the Noun Project

# C'est quoi un calcul ?

Input



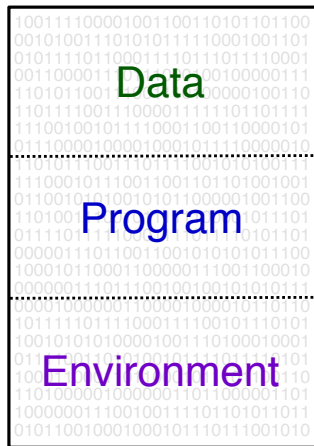
Output



Computer by Creative Stall from the Noun Project

# C'est quoi un calcul ?

Input

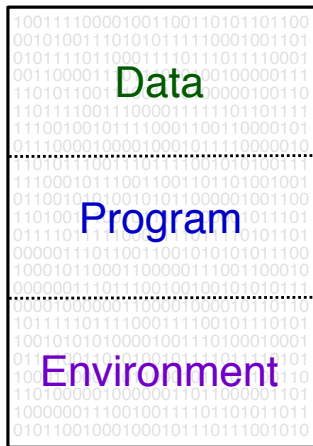


Output



Computer by Creative Stall from the Noun Project

## Input



my research

my colleagues' code

stuff I don't care about

# Que fait ce programme ?

```
data_analysis.py
```

```
from datalib import Dataset
```

```
points = [(1, 1), (-1, 1), (2, 4)]
```

```
data = Dataset()
```

```
for x, y in points:
```

```
    if x > 0:
```

```
        data.add_value(y)
```

```
print(data.average())
```



# Que fait ce programme ?

```
data_analysis.py
```

```
from datalib import Dataset

points = [(1, 1), (-1, 1), (2, 4)]

data = Dataset()
for x, y in points:
    if x > 0:
        data.add_value(y)
print(data.average())
```

Réponse rapide :

Calculer la moyenne de  $y$  pour les points où  $x$  est positif.

Le résultat est 2.5.

# Que fait ce programme ?

```
data_analysis.py
```

```
from datalib import Dataset
```

```
points = [(1, 1), (-1, 1), (2, 4)]
```

```
data = Dataset()
```

```
for x, y in points:
```

```
    if x > 0:
```

```
        data.add_value(y)
```

```
print(data.average())
```

Réponse correcte :

Ça dépend de `datalib`

# Il faut bien connaître ses bibliothèques

```
data.lib.py
```

```
class Dataset(object):  
  
    def __init__(self):  
        self.values = []  
  
    def add_value(self, value):  
        self.values = [value]  
  
    def average(self):  
        return sum(self.values, 0)/len(self.values)
```

Quelle surprise ! `add_value` ne garde que la dernière valeur !  
Le résultat de `data_analysis.py` est donc 4.

# Il faut bien connaître ses bibliothèques *et langages*

```
data.lib.py
```

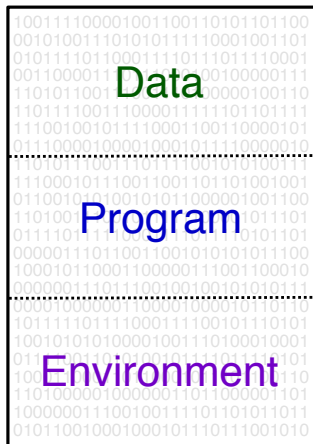
```
class Dataset(object):  
  
    def __init__(self):  
        self.values = []  
  
    def add_value(self, value):  
        self.values = [value]  
  
    def average(self):  
        return sum(self.values, 0)/len(self.values)
```

Quelle surprise ! `add_value` ne garde que la dernière valeur !

Plus précisément : 4 en Python 2 mais 4.0 en Python 3.

# Donner un sens aux bits

## Input



Data

zeros and ones

Program

interpretation of the data

Environment

interpretation of the program

TURING AWARD LECTURE

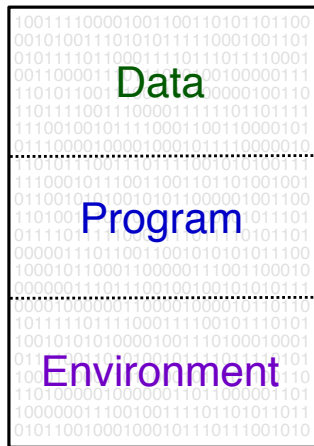
## Reflections on Trusting Trust

*To what extent should one trust a statement that a program is free of Trojan horses? Perhaps it is more important to trust the people who wrote the software.*

**KEN THOMPSON**

# C'est quoi un calcul ?

Input



Output

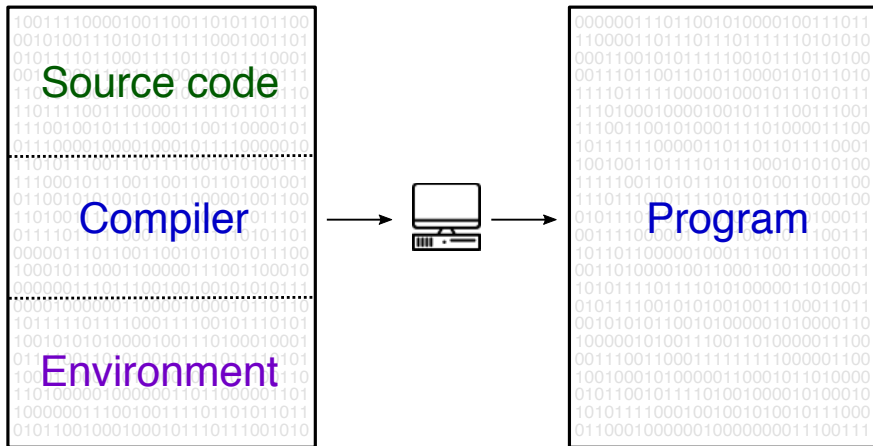


Computer by Creative Stall from the Noun Project

# D'où vient le programme ?

Input

Output

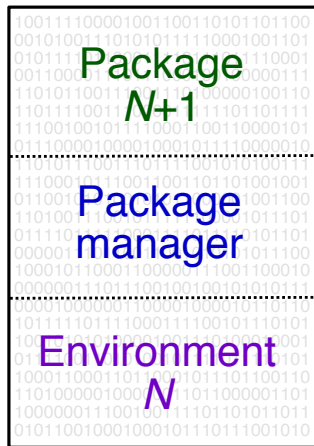


Computer by Creative Stall from the Noun Project

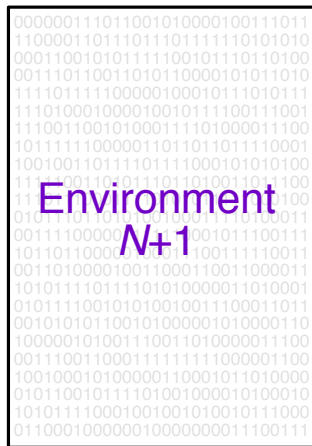


# Et d'où vient l'environnement ?

Input



Output

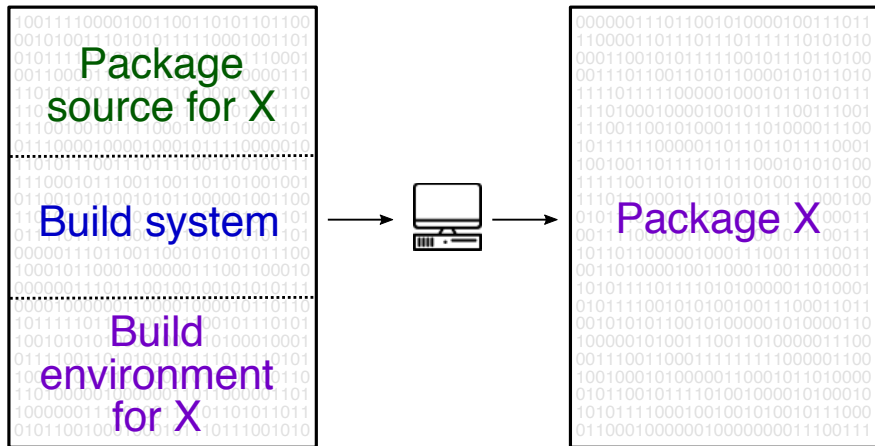


Computer by Creative Stall from the Noun Project

# Bon, alors, d'où viennent les paquets ?

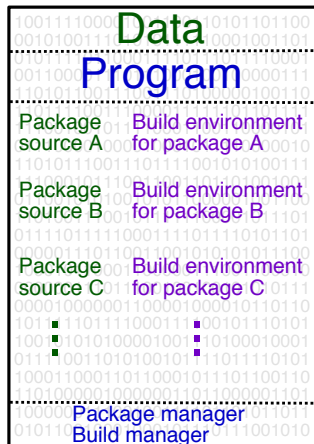
Input

Output



Computer by Creative Stall from the Noun Project

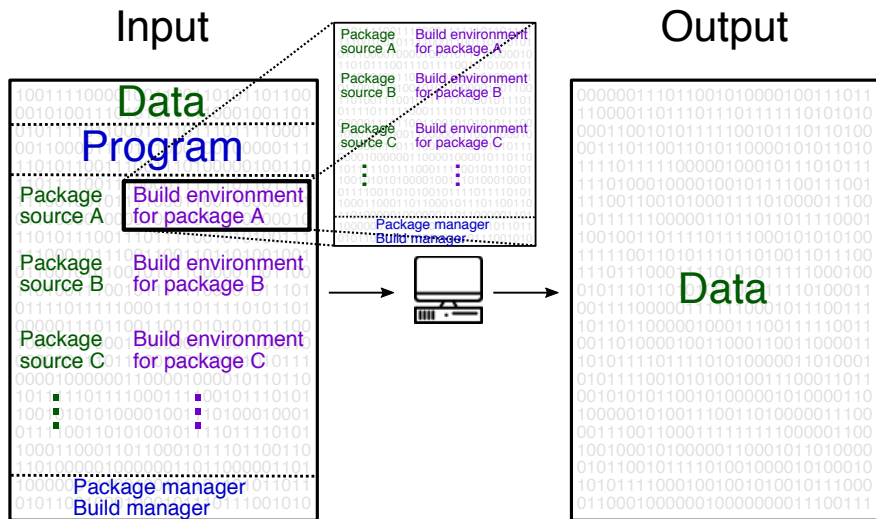
## Input



## Output

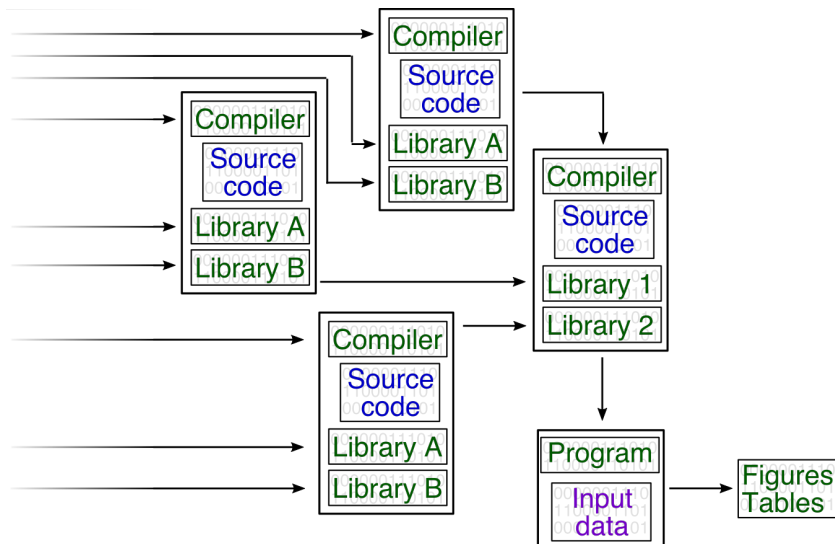


Computer by Creative Stall from the Noun Project



Computer by Creative Stall from the Noun Project

# Un autre point de vue



# La provenance des données numériques

## Observations, expériences

- Source : instruments scientifiques
- Enregistrement avec des méta-données (date, instrument, ...)

## Apports humains

- Modèles, codes source, paramètres, ...
- Traçables par la gestion de versions

## Résultats de calcul

- Données scientifiques mais aussi logiciels en binaire
- Reproductibles à partir des entrées

# Les gestionnaires de paquets classiques

Linux : dpkg/apt, RPM, pacman, ...

macOS : Homebrew, MacPorts, Fink

Multi-plateforme : conda, Spack, ...

- L'environnement est découpé en “paquets” qui sont assemblés sur un serveur puis téléchargés pour l'installation.
- On perd la trace des environnements de construction des paquets individuels.

L'environnement n'est pas reproductible à partir du code source  
...sans mesures complémentaires

Machines virtuelles : QEMU, VirtualBox, VMWare, JVM, ...

Conteneurs : Docker, Singularity, ...

- Paquets gros-grain
- On perd la trace des environnements de construction.
- Il faut assurer la reproductibilité par d'autres moyens.

L'environnement n'est pas reproductible à partir du code source  
...sans mesures complémentaires



# Comment faire mieux ?

- Nix, Guix → Ludovic Courtès
- **Reproducible Builds**  
“a set of software development practices that create an independently-verifiable path from source to binary code”

# Les trois défis de la reproductibilité numérique

- L'arithmétique à virgule flottante
- Le calcul parallèle
- Comment réconcilier la manipulation interactive avec la reproductibilité ?

# L'arithmétique à virgule flottante

- Standardisé en 1985 norme IEEE 754
- Universellement acceptée aujourd'hui.
- Ses opérations sont précisément spécifiées...
- ...et parfaitement déterministes.
- Donc : aucune particularité pour la reproductibilité !

# L'arithmétique à virgule flottante

- Standardisé en 1985 norme IEEE 754
- Universellement acceptée aujourd'hui.
- Ses opérations sont précisément spécifiées...
- ...et parfaitement déterministes.
- ~~Dont il n'y a aucune particularité pour la reproductibilité~~
- **Aucun langage de programmation donne un accès direct aux opérations IEEE 754.**
- **Les optimisations modifient les résultats.**
- **Le programmeur n'a pas le contrôle complet sur les résultats.**
- **Il devient obligatoire d'intégrer la compilation dans le protocole de reproductibilité : même compilateur, mêmes options.**

# Le calcul parallèle

Pacte avec le Diable :

- calcul plus rapide
- perte de contrôle sur l'ordre des opérations  
→ perte de reproductibilité

Pas de vraie solution,  
on ne peut que réduire l'impact des aléas.

Exemple : <https://bebop.cs.berkeley.edu/reproblas/>

# Chaque bit compte !

- Reproductibilité = résultat **égal**.
- Pas d'exception pour l'arithmétique à virgule flottante.
- Au niveau des bits, le “presque” n'existe pas.
- Le “suffisamment proche” relève de la validation / répliquabilité : c'est un jugement scientifique



Source : [PinClipart.com](https://www.pinciptart.com)

# L'avenir du calcul reproductible

La reproductibilité du calcul sera assurée automatiquement par l'ordinateur.

(sauf pour le calcul parallèle et le travail interactif)

Technologies qui rendent les données numériques traçables :

- content-addressed storage
- blockchain

Déjà utilisés par :

- Git, Software Heritage
- Nix, Guix
- IPFS / IPLD





- La reproductibilité est une source de confiance dans un résultat.

- La reproductibilité est une source de confiance dans un résultat.
- Il y en a d'autres, complémentaires.

- La reproductibilité est une source de confiance dans un résultat.
- Il y en a d'autres, complémentaires.
- Pour le calcul, elle sera assurée automatiquement...

# Résumé

- La reproductibilité est une source de confiance dans un résultat.
- Il y en a d'autres, complémentaires.
- Pour le calcul, elle sera assurée automatiquement...
- ... si nous faisons l'effort pour y arriver.

- La reproductibilité est une source de confiance dans un résultat.
- Il y en a d'autres, complémentaires.
- Pour le calcul, elle sera assurée automatiquement...
- ... si nous faisons l'effort pour y arriver.
- Deux défis se posent par la suite :
  - le travail interactif
  - la compréhensibilité

- La reproductibilité est une source de confiance dans un résultat.
- Il y en a d'autres, complémentaires.
- Pour le calcul, elle sera assurée automatiquement...
- ... si nous faisons l'effort pour y arriver.
- Deux défis se posent par la suite :
  - le travail interactif
  - la compréhensibilité
- Chercher le salut dans les artefacts binaires est sans espoir,

- La reproductibilité est une source de confiance dans un résultat.
- Il y en a d'autres, complémentaires.
- Pour le calcul, elle sera assurée automatiquement...
- ... si nous faisons l'effort pour y arriver.
- Deux défis se posent par la suite :
  - le travail interactif
  - la compréhensibilité
- Chercher le salut dans les artefacts binaires est sans espoir,
- ...mais ils peuvent être des outils efficaces dans le cadre d'une approche plus large.