# THE FRENCH OPEN SCIENCE MONITOR

the French Open Science Monitor

Measure the evolution of open science in France using reliable, open and controlled data.
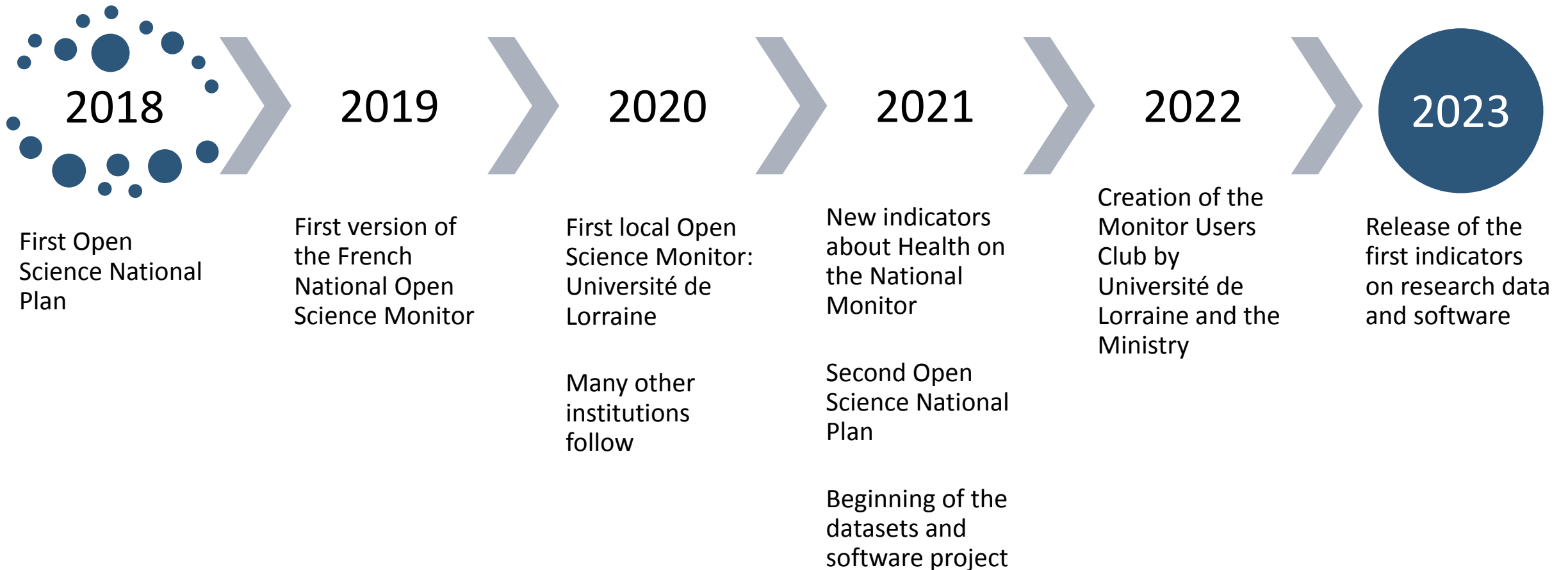
EOSC Tripartite event, 2024/09/12

Eric JEANGIRARD, French Ministry of Higher Education and Research
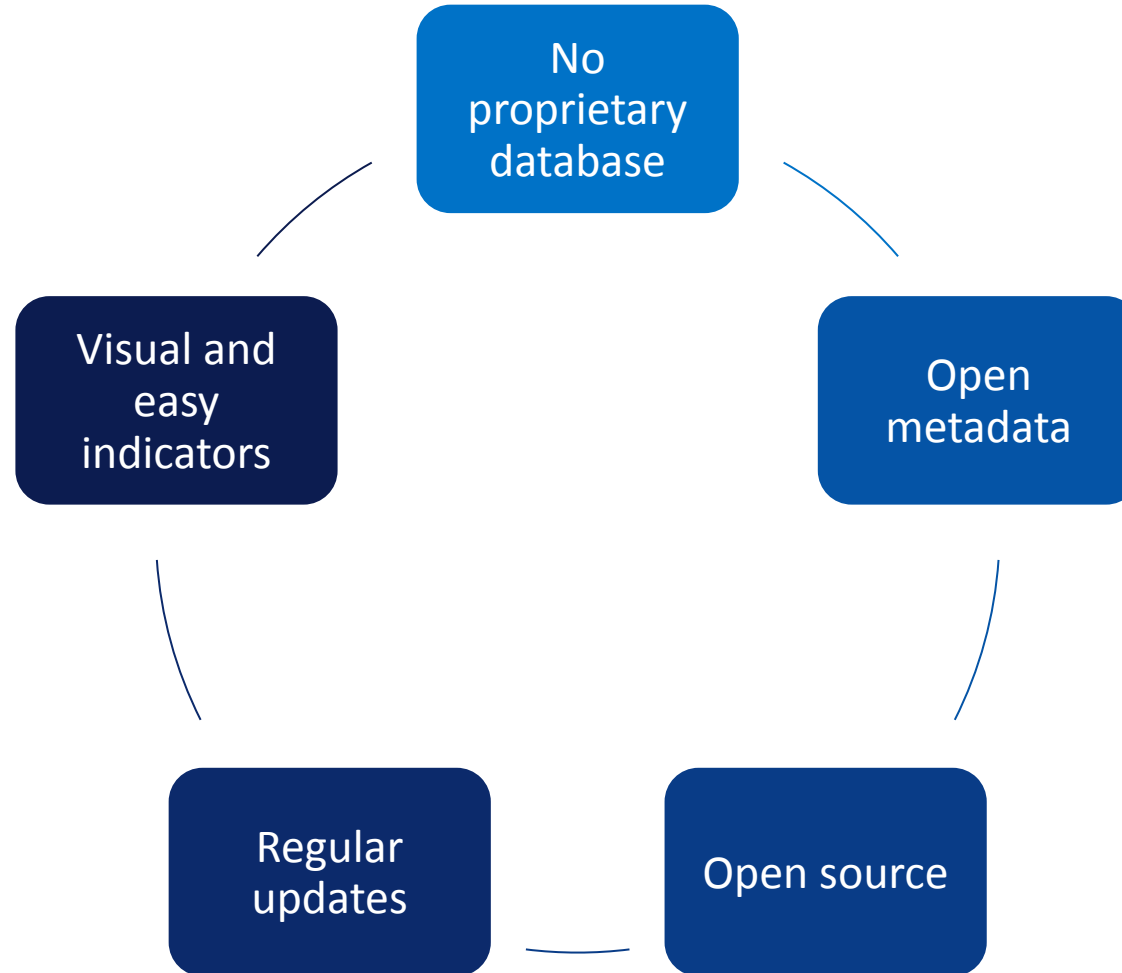
# FROM MONITORING OPEN ACCESS ...

# Since 2018, the French OSM is a monitoring tool for the Open Science public policy

**2018**

First Open Science National Plan

**2019**

First version of the French National Open Science Monitor

**2020**

First local Open Science Monitor: Université de Lorraine

Many other institutions follow

**2021**

New indicators about Health on the National Monitor

Second Open Science National Plan

Beginning of the datasets and software project

**2022**

Creation of the Monitor Users Club by Université de Lorraine and the Ministry

**2023**

Release of the first indicators on research data and software
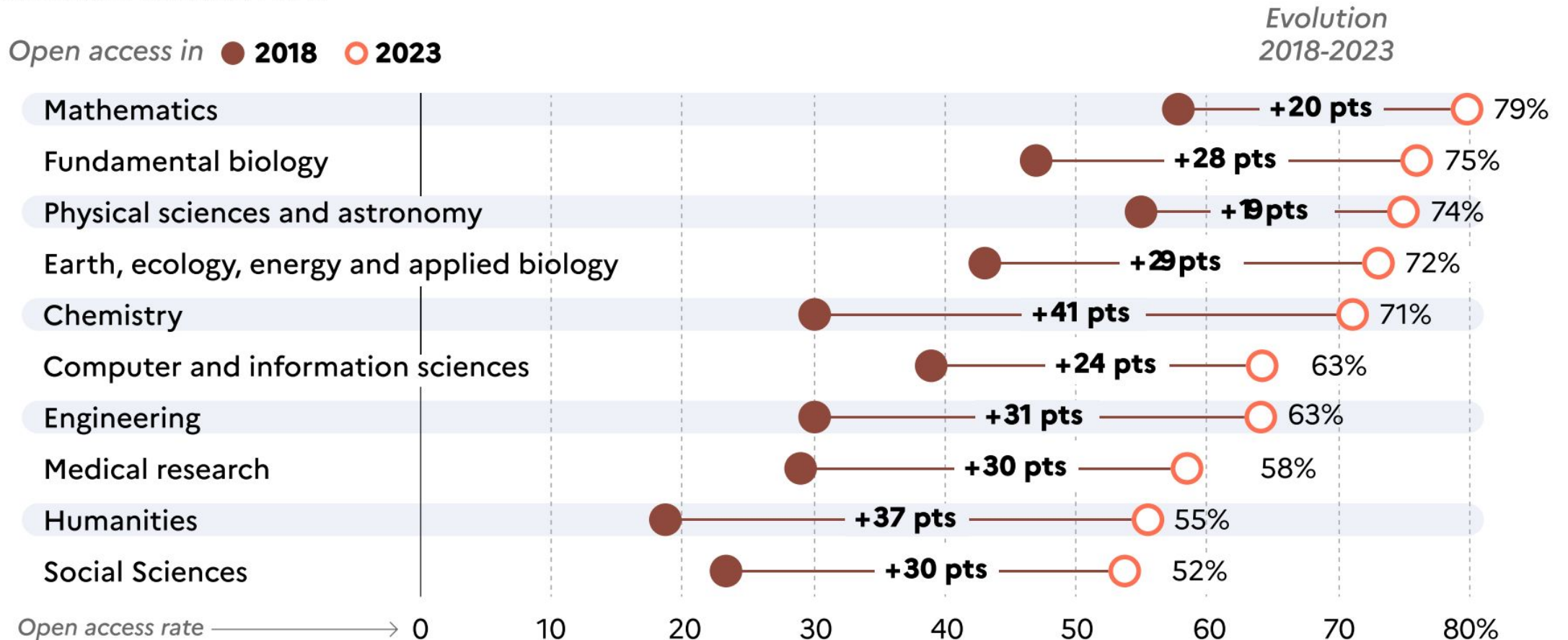
# Two principles: usefulness to the community and openness of data and processing to create a virtuous circle

- As a monitoring tool, useful for decision makers at different levels
  - national
  - institutions
  - laboratories
  - libraries
  - disciplines
  - …

- **Its openness enables the creation of new services based on the French OSM data**, making the tool even more relevant and useful for the community

# Quick overview: the open access rate in France has strongly improved since 2018

**Rate of open access publications in France, for each discipline between 2018 and 2023**

Open access in ● 2018  ○ 2023

Evolution 2018-2023

| Discipline | Evolution | 2023 |
|---|---|---|
| Mathematics | +20 pts | 79% |
| Fundamental biology | +28 pts | 75% |
| Physical sciences and astronomy | +19 pts | 74% |
| Earth, ecology, energy and applied biology | +29 pts | 72% |
| Chemistry | +41 pts | 71% |
| Computer and information sciences | +24 pts | 63% |
| Engineering | +31 pts | 63% |
| Medical research | +30 pts | 58% |
| Humanities | +37 pts | 55% |
| Social Sciences | +30 pts | 52% |

Open access rate → 0   10   20   30   40   50   60   70   80%

# Changing scale: from national to local monitoring (institution level)

- The national monitoring tool can be derived to monitor local OS progress in any French institutions

- More than 80 French institution use this tool now, including large research organisms and universities, funders, but also labs

- The infrastructure and development is handled at the national level and benefit to the whole French ecosystem, providing open data and open services

# ...TO MONITORING OPEN SCIENCE

# Monitoring the pillars of Open Science?

# Monitoring research outputs like software or data is harder than publications

| Technical | Factual |
|---|---|
| • No global database for research data and software<br><br>• Poor (if any) metadata for monitoring: affiliations, topics…<br><br>• Variety of PIDs | • Low awareness from researchers on the value of these research products<br><br>• Low recognition in the individual assessment process |

# Multiples approaches to monitor monitor multiple practices

Since 2021

## Using publications

- Downloading the PDF documents of French publications
- Detecting and characterising mentions to datasets and software (GROBID, Softcite, DataStet)
- Computing indicators (ex : proportion of publications that share software or code)

## Using repositories for datasets

- Dump of DataCite
- Identifying "French" DOIs using affiliations, as well as other metadata elements (publisher, clientId)
- Enrichment
- Computing indicators

# Mining full-text to detect software mentions

- **Innovative approach** based upon the use and development of machine learning tools
  - GROBID: full-text structuring
  - Softcite: **software mention detection**
  - DataStet: **data set mention detection**

- Automatic characterisation of mentions: **usage / production or creation / sharing**

- Another challenge: **downloading massive amounts of full-texts**

s c i e n c e - m i n e r

# Software is used across all the disciplines

Proportion of publications in France that mention the use of code or software by discipline

| Discipline | Proportion |
|---|---|
| Fundamental biology (total = 19.3 k) | 59 % |
| Computer and information sciences (total = 6.2 k) | 55 % |
| Earth, ecology, energy and applied biology (total = 11.0 k) | 55 % |
| Engineering (total = 7.2 k) | 47 % |
| Chemistry (total = 6.4 k) | 45 % |
| Physical sciences and astronomy (total = 9.0 k) | 41 % |
| Medical research (total = 27.6 k) | 33 % |
| Mathematics (total = 4.0 k) | 21 % |
| Humanities (total = 5.7 k) | 19 % |
| Social Sciences (total = 5.4 k) | 18 % |

- **Almost half** of French publication **mention the use of software**

- **All disciplines are involved**

- Mining PDFs is costly and difficult to scale because of
  - PDF accessibility
  - computation costs

  ⇒ **There is a need for global cooperation on this topic**

# International initiatives to scale up

- The French OSM is the first national Open Science Monitor to tackle software monitoring using software mention detections on a large scale (more than 1 millions PDFs analyzed)

- Working on a large-scale infrastructure to monitor software use through the scholarly publications remains a challenge. International coordination is key
  - to **build a consensus** on the detection techniques
  - to push for **open source software to be used for these detections** (relying on proprietary tools would create new dependances that we could avoid)
  - to build a large scale infrastructure

- The French Ministry of Higher Education and Research, the Université de Lorraine, Inria and Unesco organized a workshop on the subject in December 2023

- It led to a first draft for Principles of Open Science Monitoring, currently reviewed internationally

- This is the starting point for OSMI, the Open Science Monitoring Initiative

# PERSPECTIVES

# Upcoming challenges

• **International initiative** on open science monitoring

• Dedicated **infrastructure** to analyse software through publications

• A complementary approach, based on software metadata directly (not the mentions in the publications) is also important to **build software catalog with a high quality of metadata**

# THANK YOU!

✉ ERIC.JEANGIRARD@RECHERCHE.GOUV.FR

✹ HTTPS://FRENCHOPENSCIENCEMONITOR.ESR.GOUV.FR/

# CREDITS

Caution: Image by memyselfaneye from Pixabay
Green statistics: Storyset by Freepik
Telescope: Everypixel by Arnaud Papa
Thank you: Image by Ryan McGuire from Pixabay