# Doing FAIR (software)

# in Environmental and Life Sciences
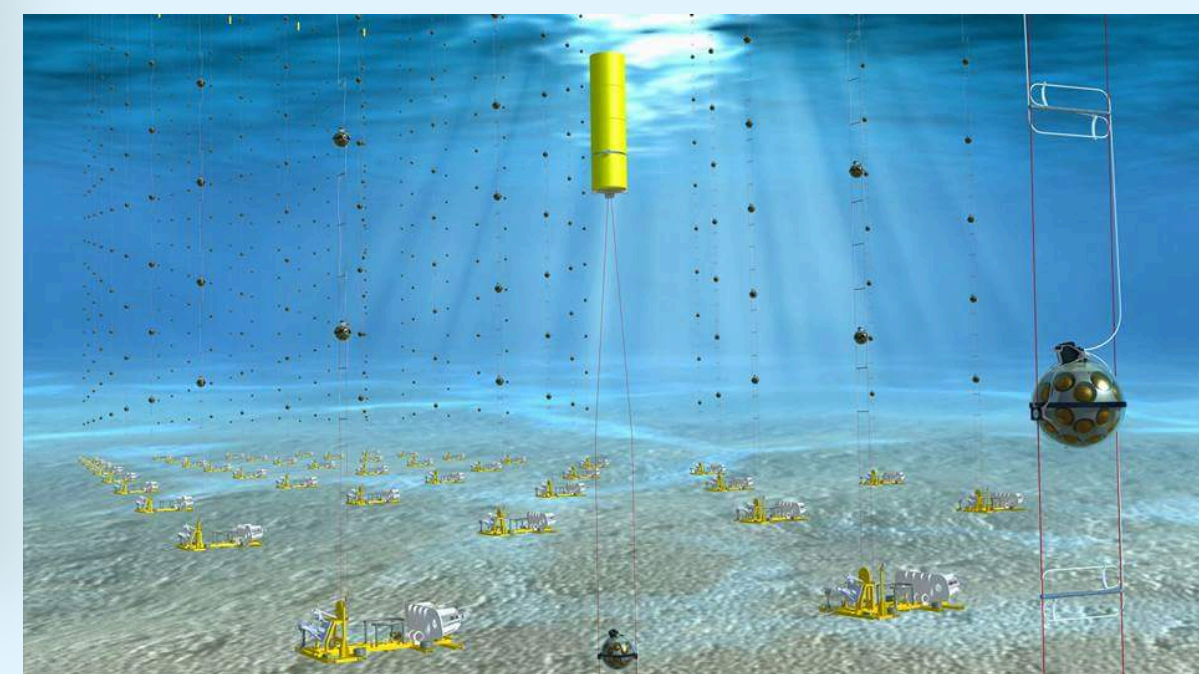
Wageningen, 12-12-2018

# 100+ projects



## Humanities & Social Sciences

incl. SMART cities, text analysis, creative technologies

## Physics & Beyond

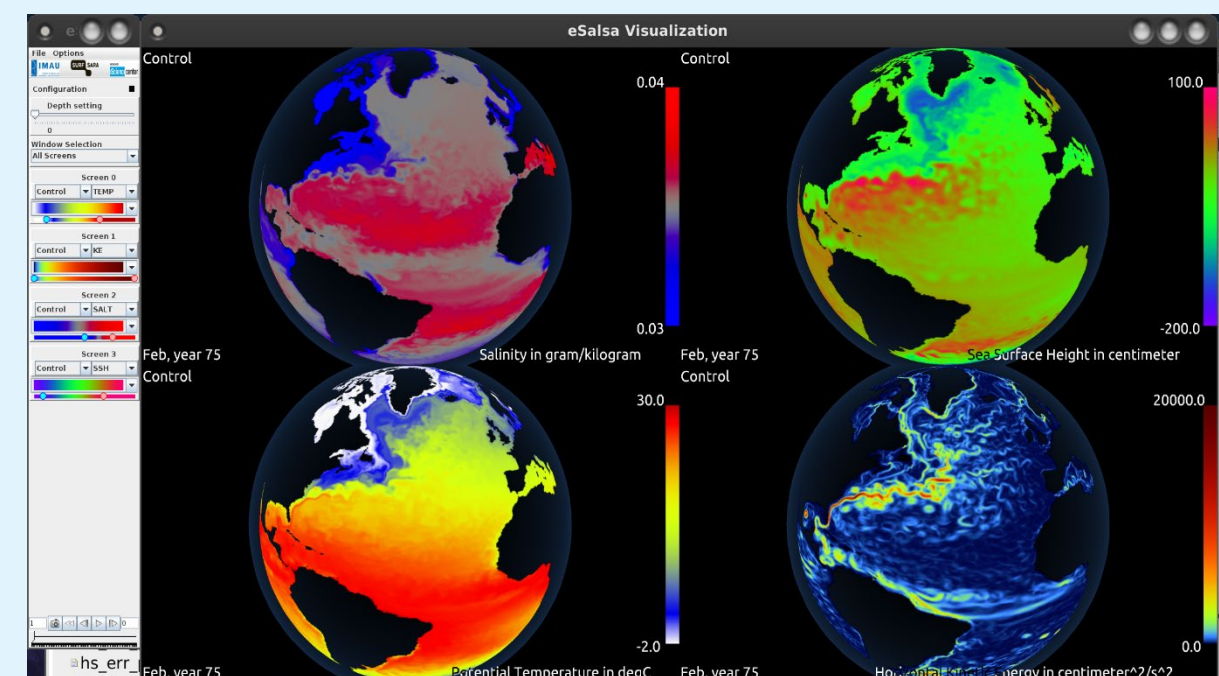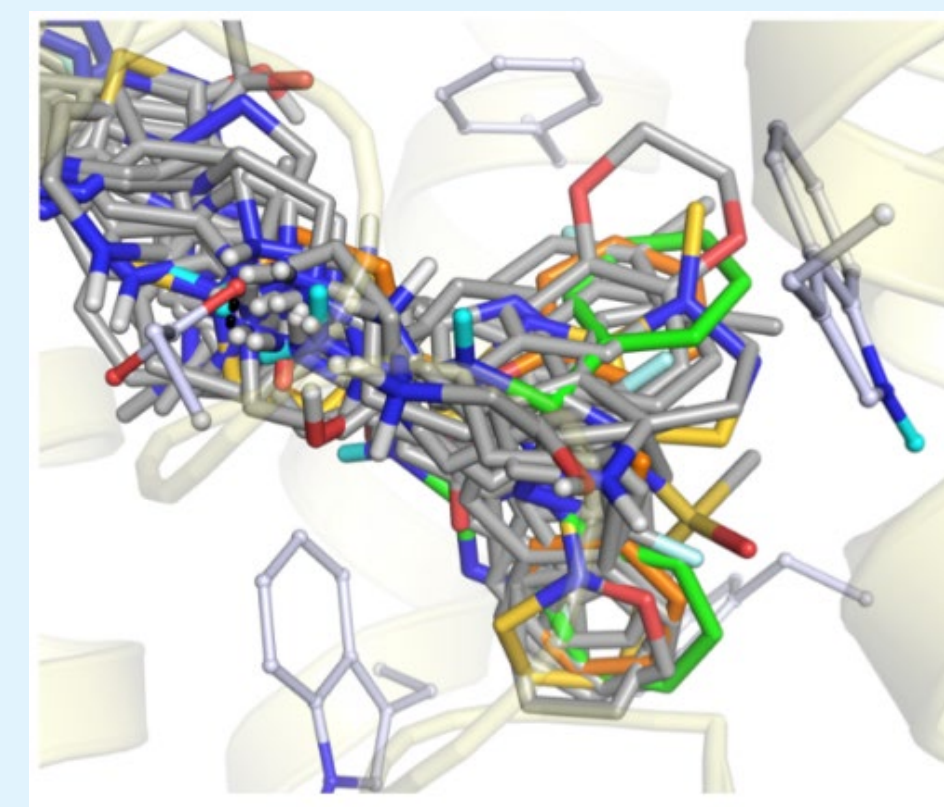incl. astronomy, high-energy physics, advanced materials

## Sustainability & Environment

incl. climate, ecology, energy, logistics, water management

## Life Sciences & eHealth

incl. bio-imaging, next generation sequencing, molecules

# DAILY EXPRESS

THE WORLD'S GREATEST NEWSPAPER

## Clooney's amazing mother

## Iran threatens serious action against sailors

# THE BIG CLIMATE CHANGE 'FRAUD'

**We are not to blame says top scientist... It's a con to raise tax**

By John English

THE scientific consensus that mankind has caused climate change was rocked yesterday as a leading academic called it a "load of hot air underpinned by fraud".



## IT'S FURRY NICE TO MEET YOU

- In the context of the sharing of data and methodologies, …. Professor XX's actions were in line with common practice in the climate science community.

- It is not standard practice in climate science to publish the raw data and the computer code in academic papers. However, climate science is a matter of great importance and the quality of the science should be irreproachable. We therefore consider that climate scientists should take steps to make available all the data that support their work (including raw data) and full methodological workings (including the computer codes).
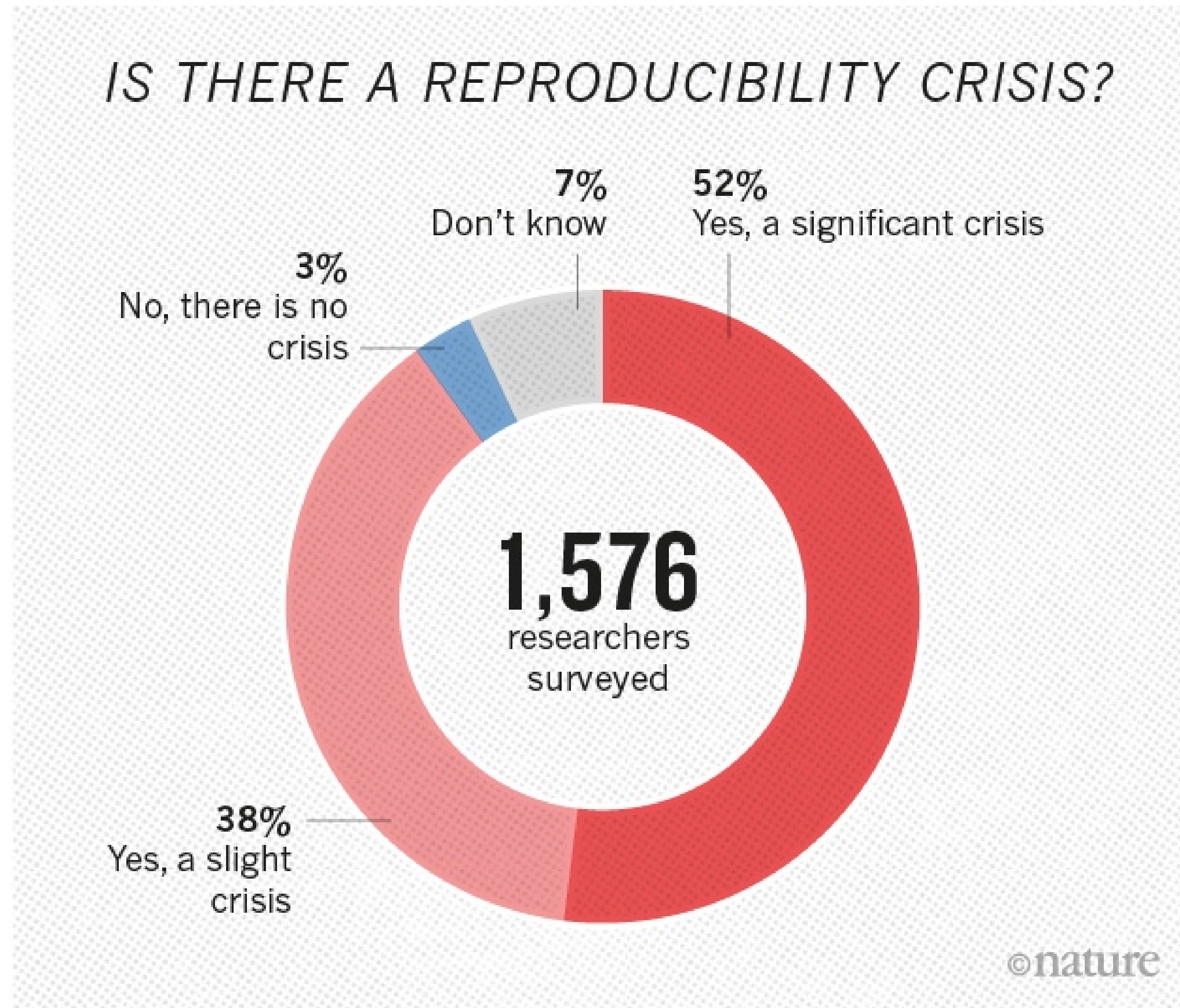
IS THERE A REPRODUCIBILITY CRISIS?

7%
Don't know

52%
Yes, a significant crisis

3%
No, there is no crisis

1,576
researchers surveyed

38%
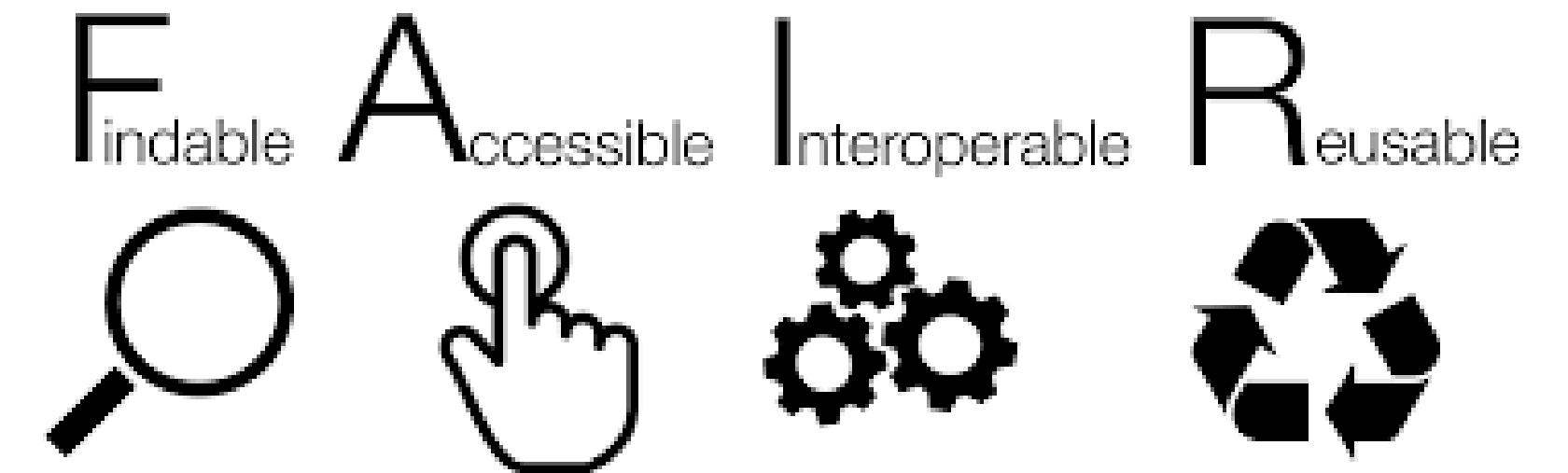Yes, a slight crisis

©nature

Baker, 2016, Nature

# About science and scholarly research

- Increasingly problem -driven on big societally relevant themes

- Increasingly inter -, multi -, trans -disciplinary

- Grand challenges: clean energy, safe societies, healthy societies, etc.

Open Science is about **extending the principles of openness to the whole research cycle**, fostering sharing and collaboration as early as possible thus entailing a systemic change to the way science and research is done

*FAIR data principles and FAIR software principles contribute to open science*

F<sub>indable</sub> A<sub>ccessible</sub> I<sub>nteroperable</sub> R<sub>eusable</sub>

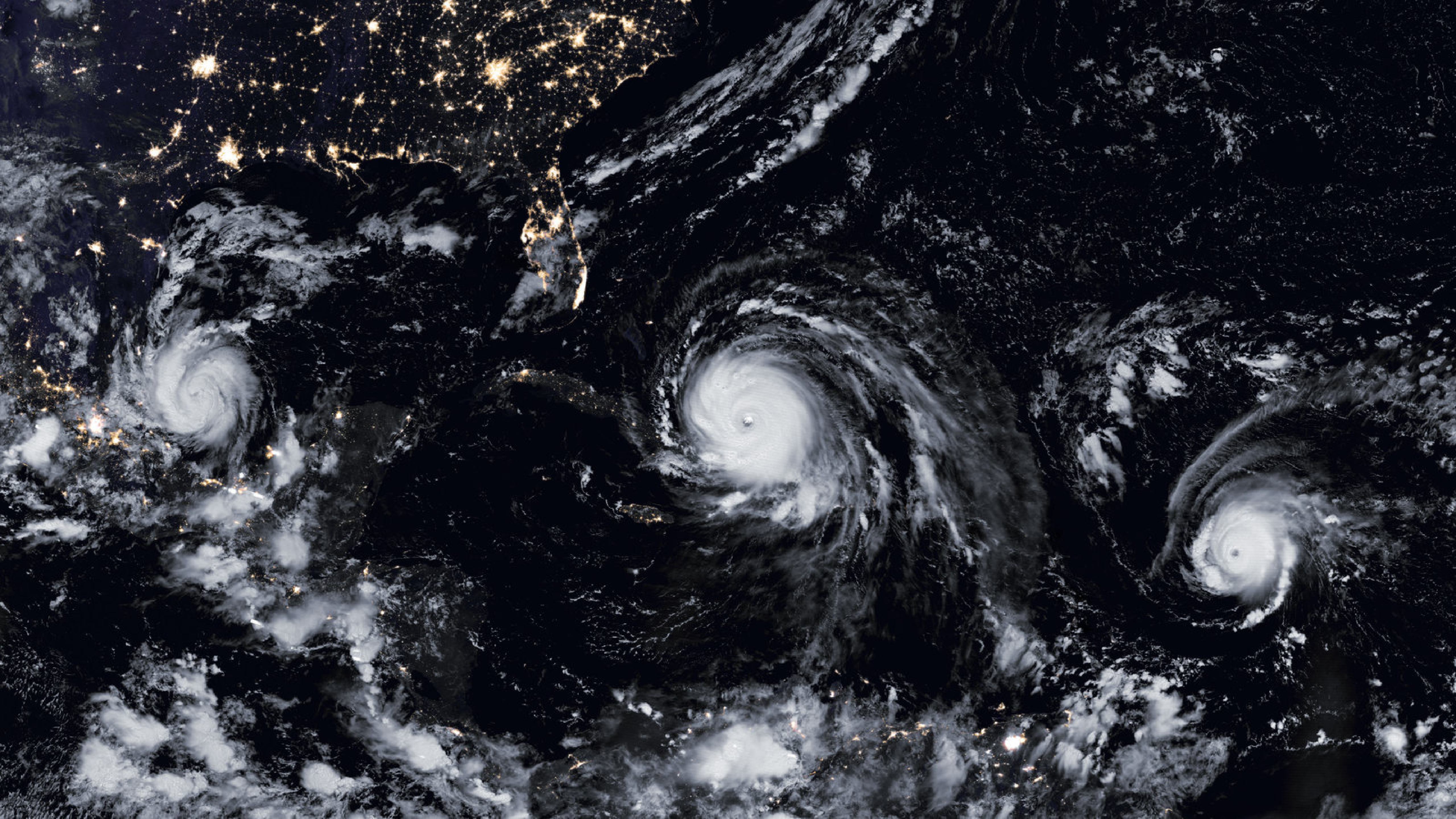**Findable** : sufficiently rich metadata and unique persistent identifier

**Accessible** : metadata is in machine and human readable format

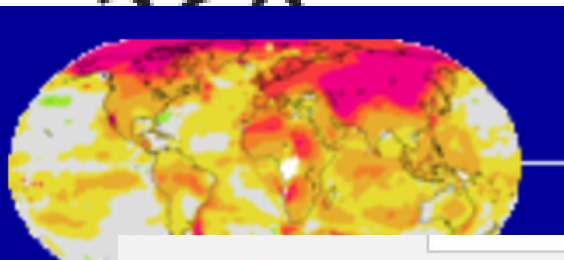Software and metadata is deposited in trusted community approved

repository

**Interoperable** : uses community accepted standards and platforms,

making it possible for users to run the software

**Reusable:** has clear license and documentation

# Examples of FAIR in weather & climate research

# Findable open data in climate research

# Interoperable in climate research

- Netcdf : is a set of software libraries and self-describing, machine-independent data formats that support the creation, access, and sharing of array-oriented scientific data.

- Climate and Forecasting conventions

- CMOR: climate model output writer



Figure 1. The structure and syntax of a CDL (ASCII) equivalent to a NetCDF file.

Smartphone

air $k_e$ $k_b$ body

$T_e$ $T_b$

$P_p$ $T_p$

London, VAL:
ME = −0.28 °C
MAE = 1.45 °C
CV = 0.12
$\rho^2$ = 0.65

- Smartphone, CAL
- Smartphone, VAL
- WMO nr. 037683
- Battery temperature

Overeem et al GRL 2013

weeralarm

# Downscaling

**Daily forecasts**
**WRF3.5 + urban module (SLUCM)**
**48 hour runs, 24 hour spin-up**

**Domain 1: 12.5km**
**default setup**

**Domain 2: 2.5km**
**default setup**

**Domain 3: 500m**
**hi-res landuse,**
**Rijkswaterstaat river temperatures**

**Domain 4: 100m**
**Rijkswaterstaat river temperatures,**
**TOP10NL, satellite imagery, AHN2**
**(height map), CBS data**



**Attema et al, IEEE eScience, 2015**

WAGENINGEN UR
For quality of life

b.

# Flexible steering, execution of models and data handling

● ● ●

```python
from ewatercycle.models import PcrGlobWB
from ewatercycle.forcings import Gfs
from ewatercycle.plotting import geo_plot, timeseries_plot
```

```python
parameterset = PcrGlobWB.parametersets['RhineMeuse30min']
# Or generate a parameterset for a region
parameterset = PcrGlobWB.parameterset_from_region(latmin=4, latmax=10, lonmin=45, lonmax=55)
```
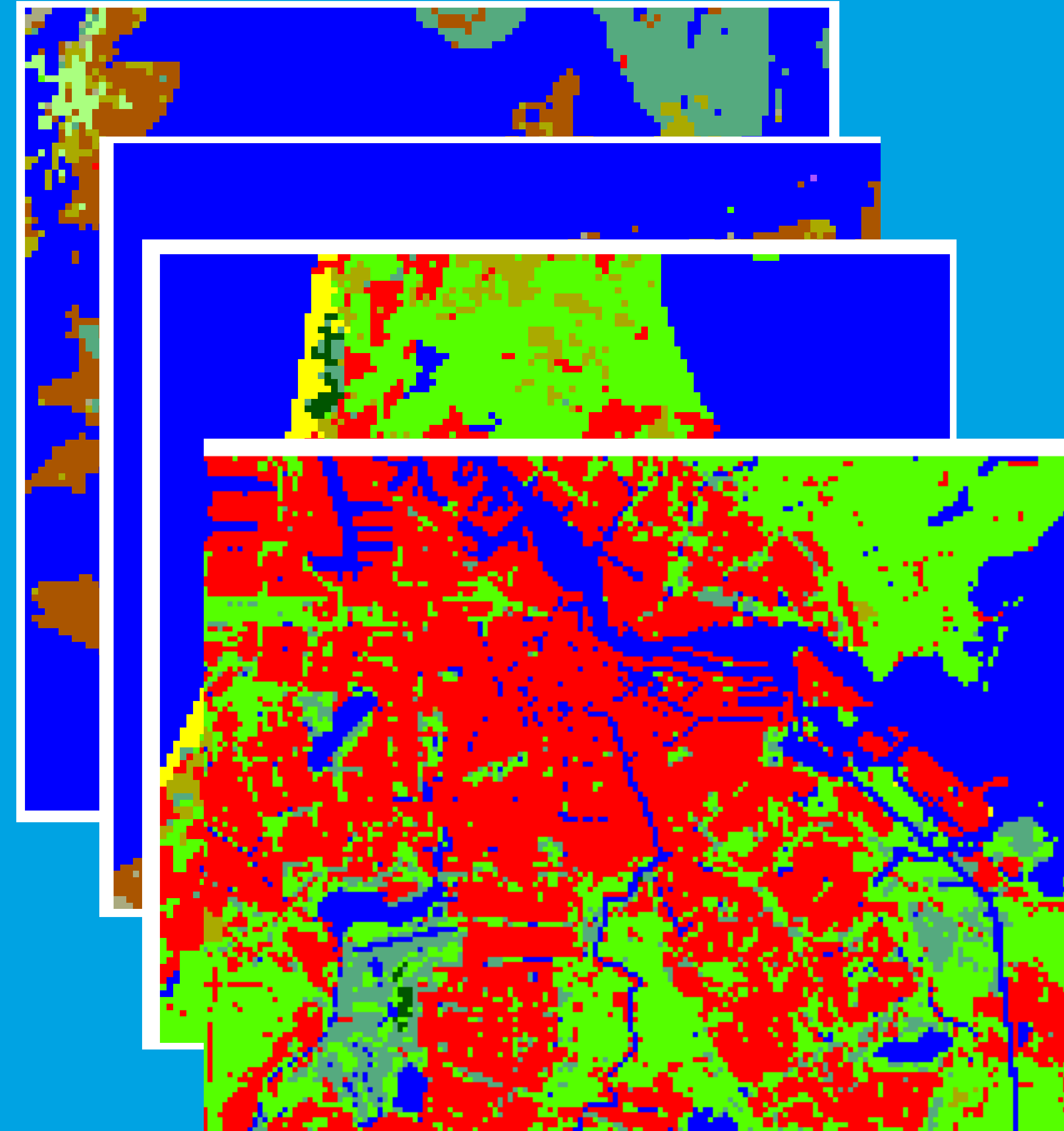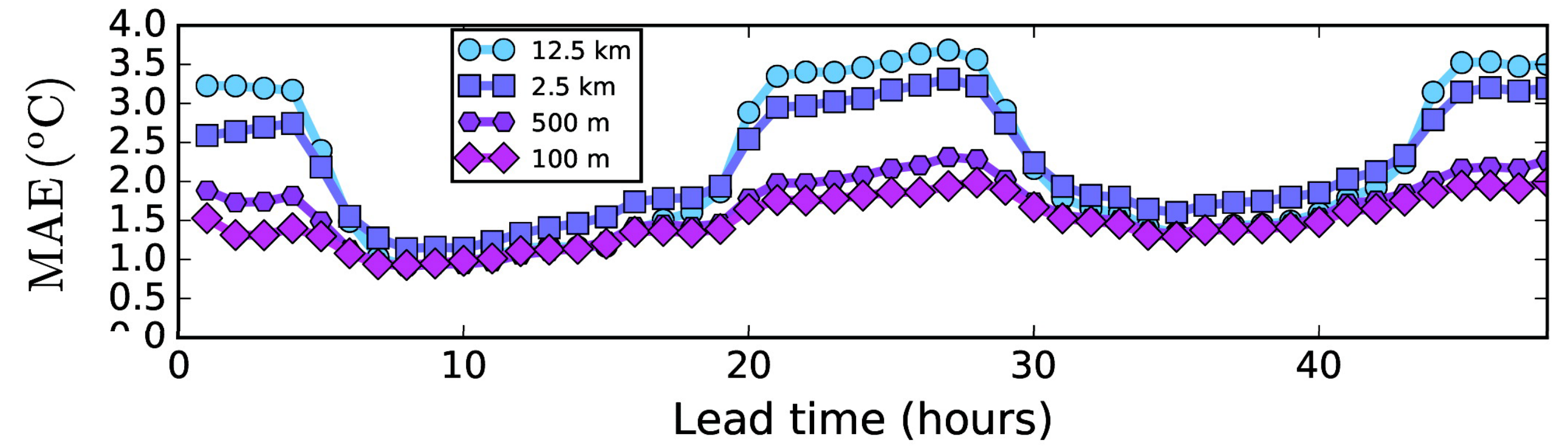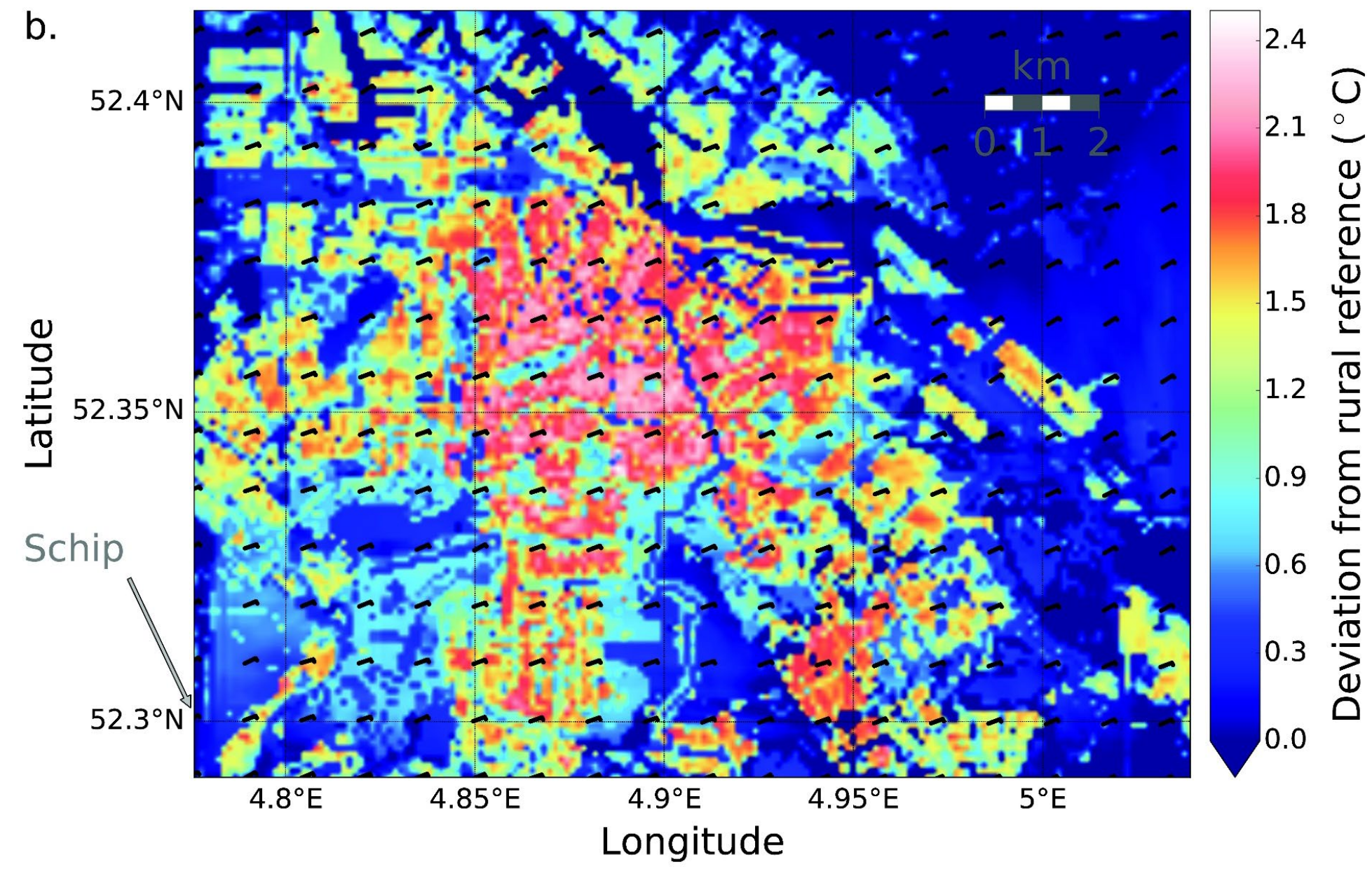
```python
forcing = Gfs()
```

```python
start = '1999-01-01T00:00:00Z'
end = '2010-31-12T23:59:59Z'
```

```python
model = PcrGlobWB(parameterset=parameterset,
                  forcing=forcing,
                  start=start,
                  end=end,
                  )
```
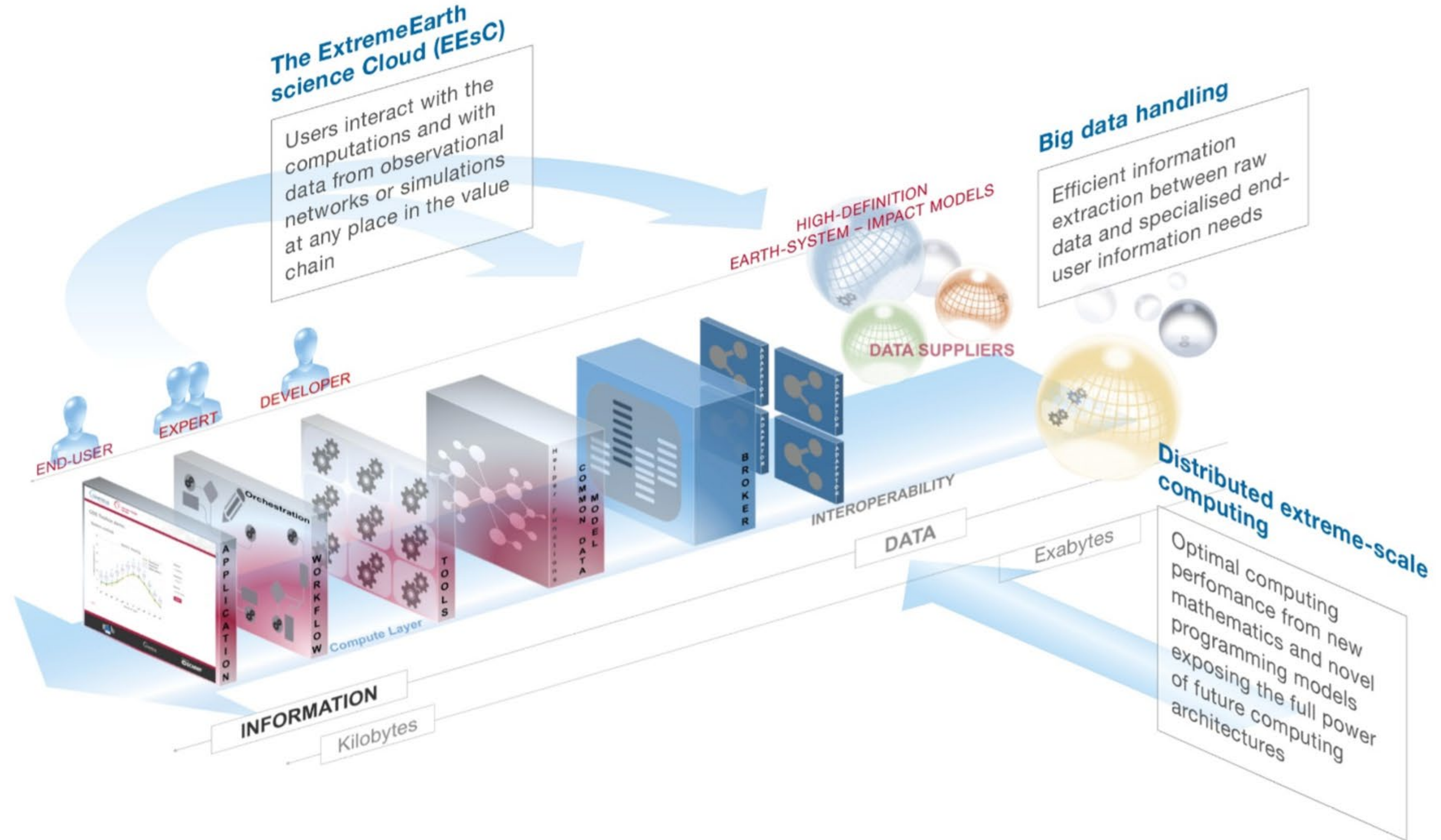
```python
discharge_over_time = []
while model.current_time < model.end_time:
    model.update()
    discharge_over_time.append(model.discharge)
```

```python
# Plot discharge of last time step
geo_plot(model.discharge)
```

Niels Drost, pers. Comm, NLeSC/TUD/UU/WUR/Deltares eWatercycle II project

The ExtremeEarth science Cloud (EEsC)

Users interact with the computations and with data from observational networks or simulations at any place in the value chain

Big data handling

Efficient information extraction between raw data and specialised end-user information needs

HIGH-DEFINITION EARTH-SYSTEM – IMPACT MODELS

DATA SUPPLIERS

DEVELOPER

EXPERT

END-USER

APPLICATION

WORKFLOW

TOOLS

Helper functions

COMMON DATA MODEL

BROKER

INTEROPERABILITY

DATA

Exabytes

Distributed extreme-scale computing

Optimal computing performance from new mathematics and novel programming models exposing the full power of future computing architectures

Orchestration

Compute Layer

INFORMATION

Kilobytes
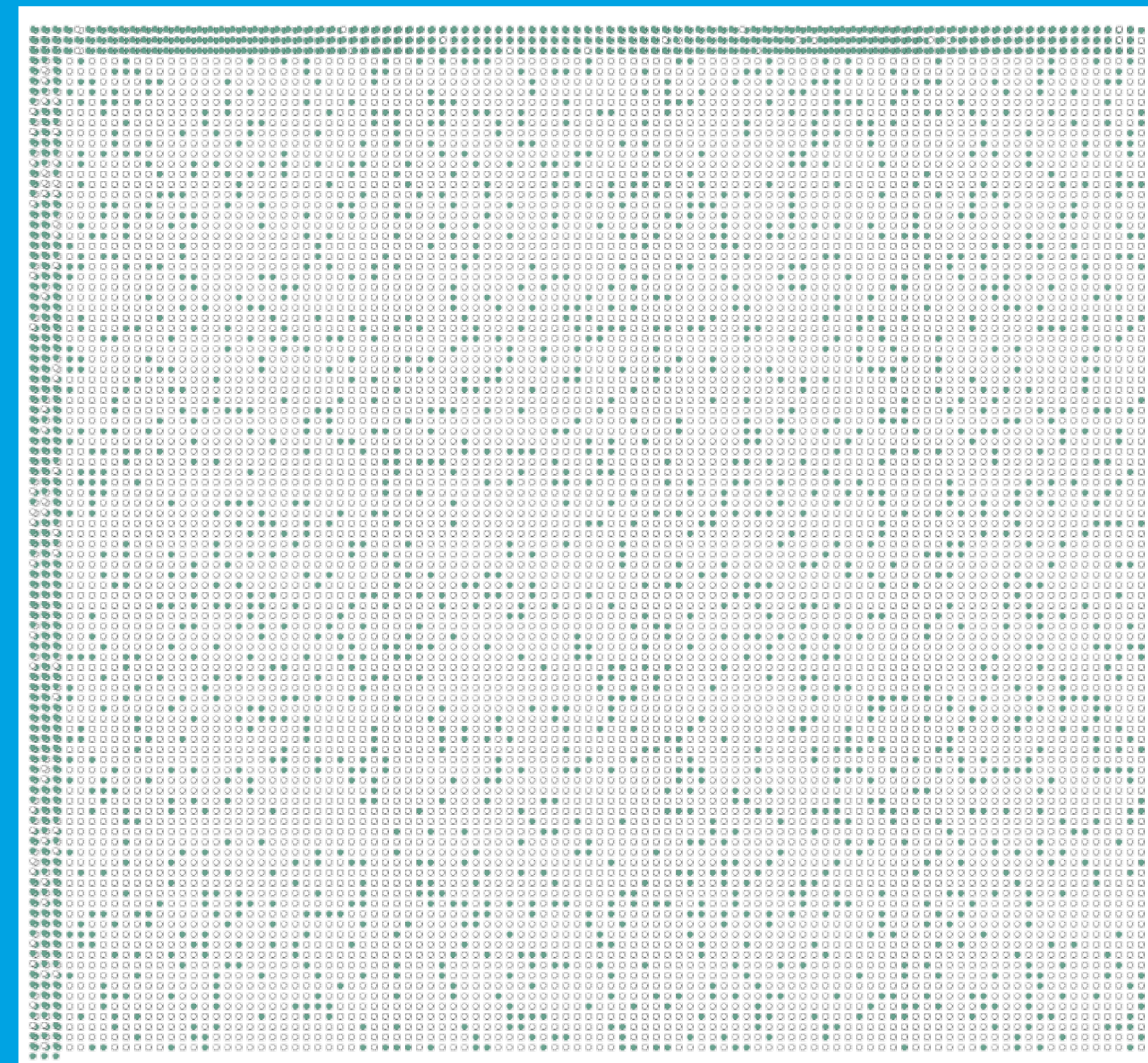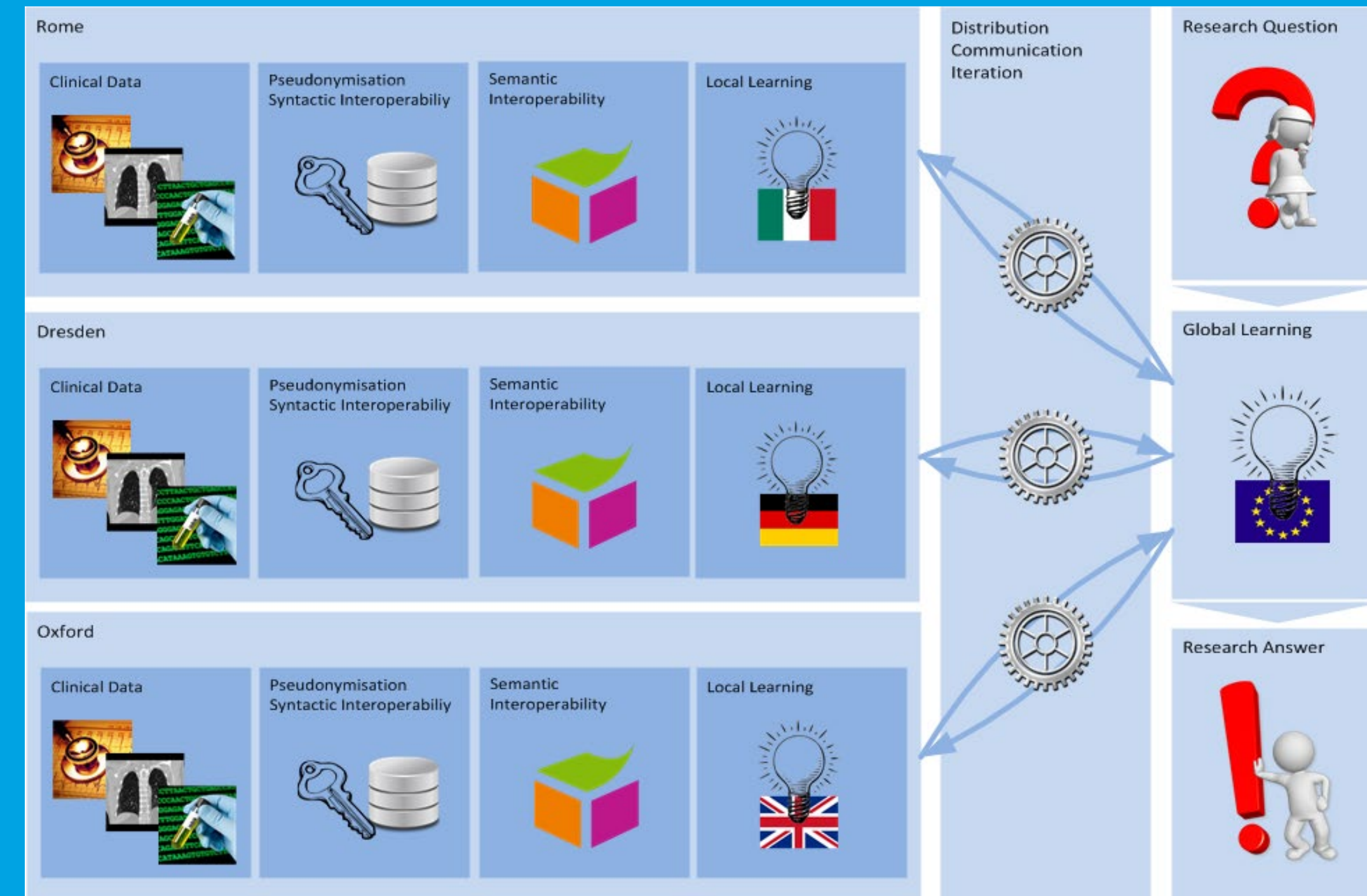
# Examples of FAIR in life sciences

# Big data landscape in health care

- **Clinical research**
  - 3% of patients
  - 100% of features
  - 5% missing
  - 285 data points

- **Clinical registries**
  - 100% of patients
  - 3% of features
  - 20% missing
  - 240 data points

- **Clinical routine**
  - 100% of patients
  - 100% of features
  - 80% missing
  - <u>2000</u> data points

**Andre Dekker**

# A Global Distributed Routine Data Registry

- Keep data locally

- Standardize it according to an ontology

- Make and send around learning and quality indicators

- Share the results & quality indicators – not the data!!
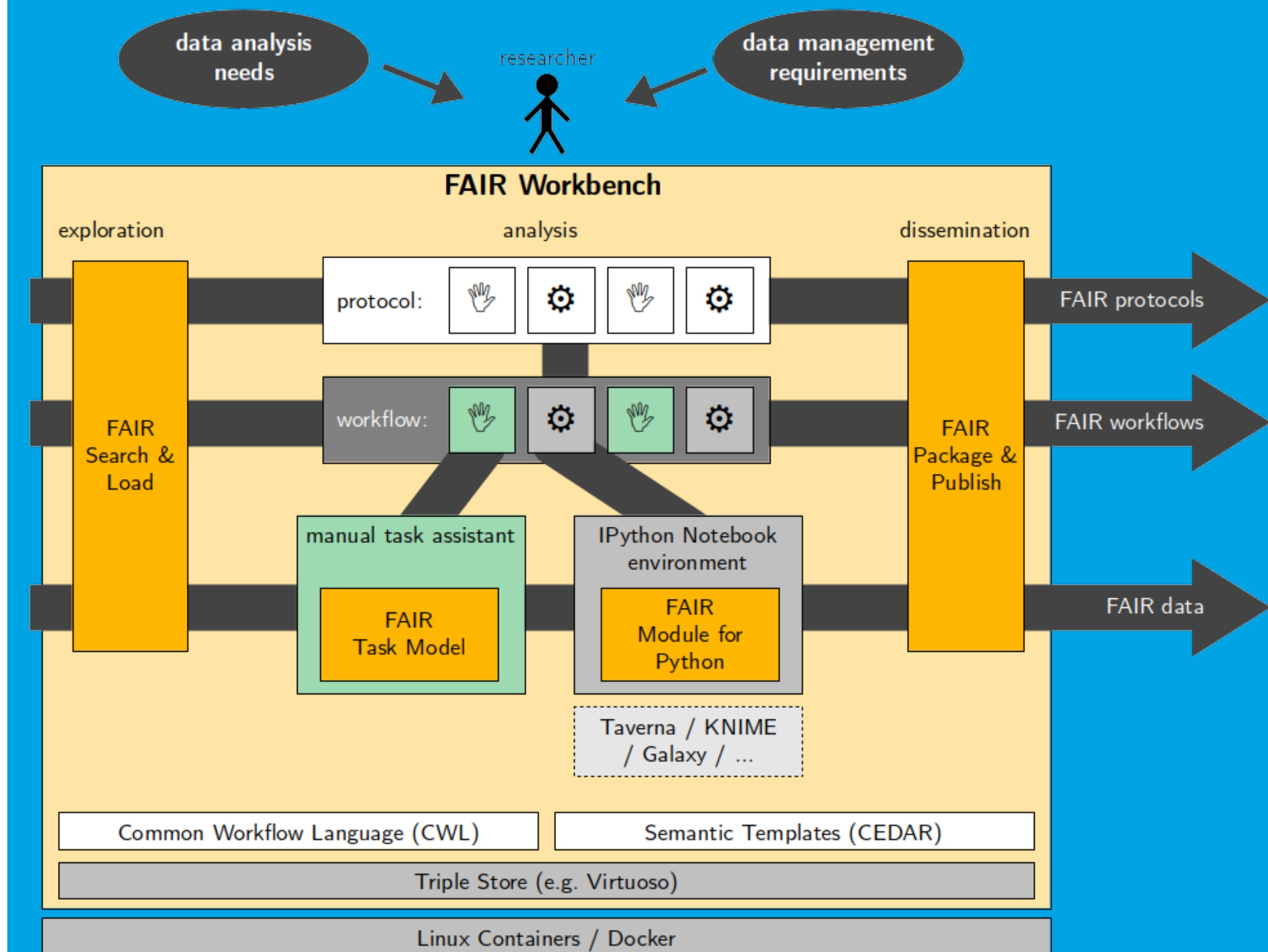
# Reproducible science

Requires not only **data** to be FAIR but also

Software:

- Research Software Directory (F + A)

- Use standard file formats, Docker, API's, etc. (I)

- NLeSC guide: https://guide.esciencecenter.nl/ (R)

Workflows:

- Common Workflow Language

  - platform independent workflow definition and execution



**FAIR Workflows project**
**Collaboration with Tobias Kuhn, Michel Dumontier**

# Linked data platform to relate genes with traits

**Trait (e.g. color)**



**Unstructured QTL data**



*Text/Table mining*

**FAIR data**

**Linked data platform**

VIRTUOSO UNIVERSAL SERVER

| RDF |
| OWL |
| R2RML |
| SPARQL |

**standards**

W3C®

*Analytics/ Ranking*

**Candidate genes for traits**

**Structured biomolecular data**

Genes  Proteins  Networks  Metabolites

**Biomolecular DBs**

*Genome annotations*

*Session F: Arnold Kuzniar, Richard Finkers, Richard Visser*
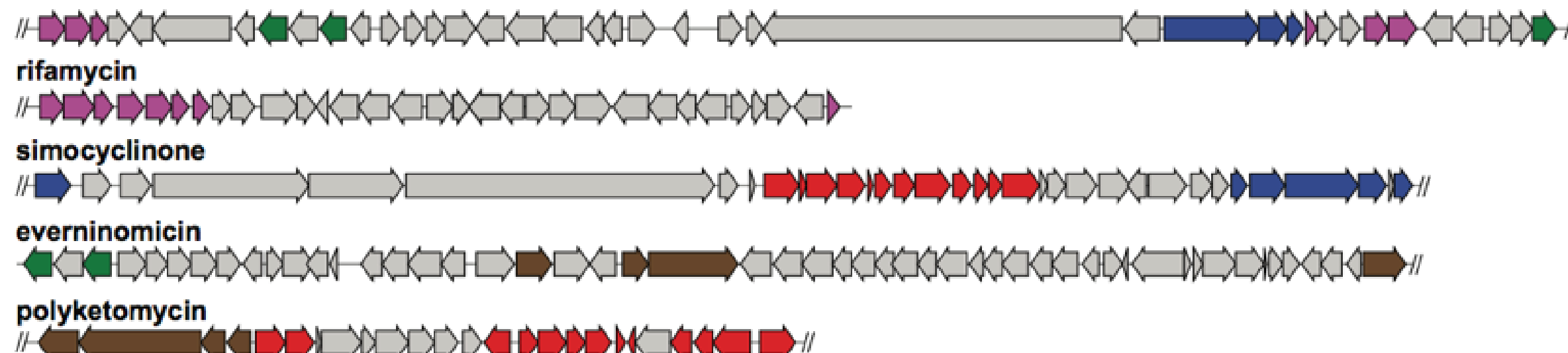
**https://github.com/candYgene**

**Microbiome (e.g. in the human gut)**

**Metabolomes (Natural products)**

**Genomes (Biosynthetic gene clusters)**

*Session F: iOMEGA project*
*Justin van der Hooft, M. Medema*
*S. Verhoeven, F. Huber, L. Ridder*

## FAIR: Just do it!

- Supports working across domains of research
- Requires domain knowledge, digital competences & digital infrastructures, hence an collaborative work environment!
- Absolutely necessary for evidence-based (and transparent) decision making
- Not only data, but also software, worfklows , methods..

Thank you