

From open source and open data to “open computation”

A climate science perspective

Roly Perera

Institute of Computing for Climate Science, University of Cambridge
School of Computer Science, University of Bristol

ICCS overview

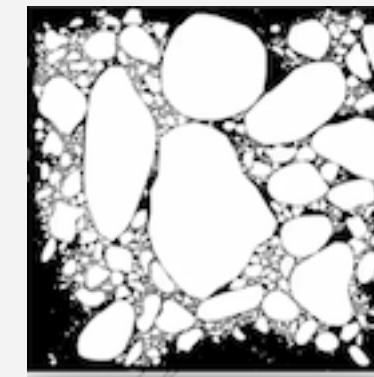
director + 3 co-directors
programme manager
3 researchers
9 RSEs



Institute of
Computing for
Climate Science

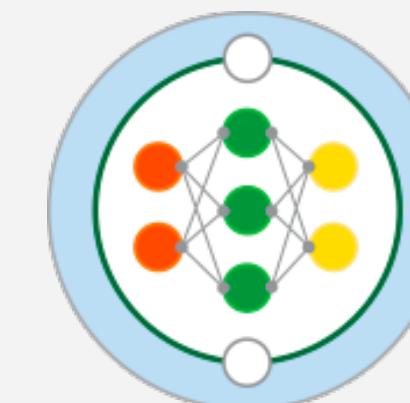
Collaboration/RSE support for 4 VESRIs (Virtual Earth Systems Research Institutes)

DataWave: Collaborative
Gravity Wave Research



SASIP: Scale-Aware
Sea Ice Project

LEMONTREE: Land Ecosystem
Models based On New Theory,
Observations, and Experiments



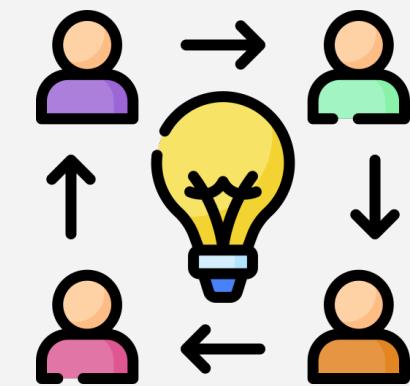
M²LInES: Multiscale Machine
Learning In Coupled Earth
System Modeling

Ingredients of open science

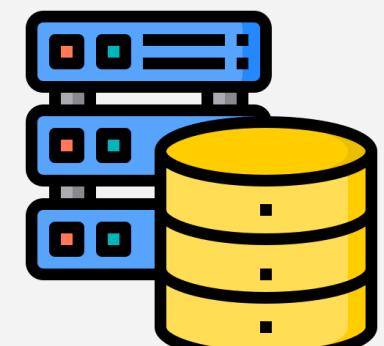
What is open science?



open source



working in the open

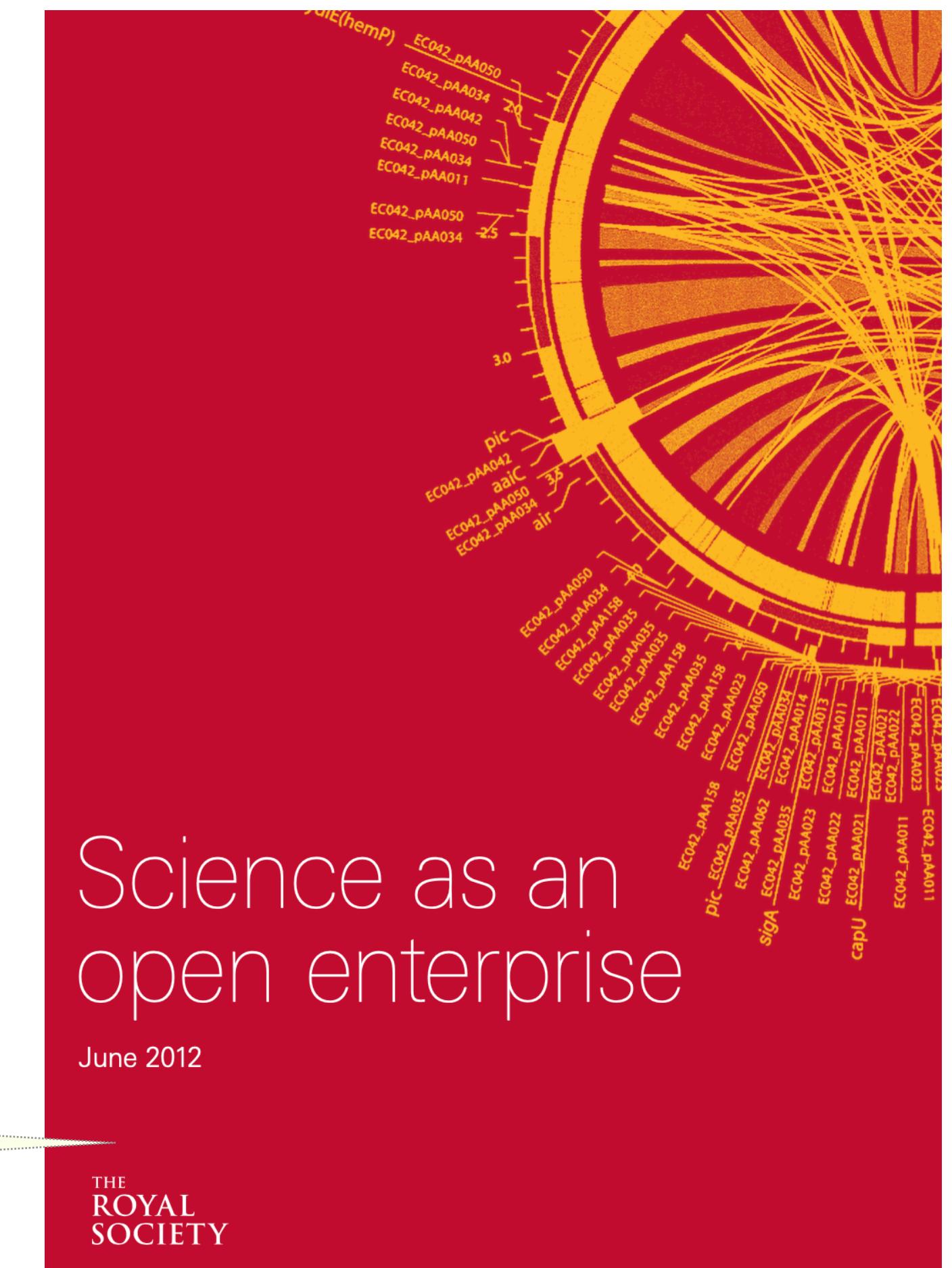


open data

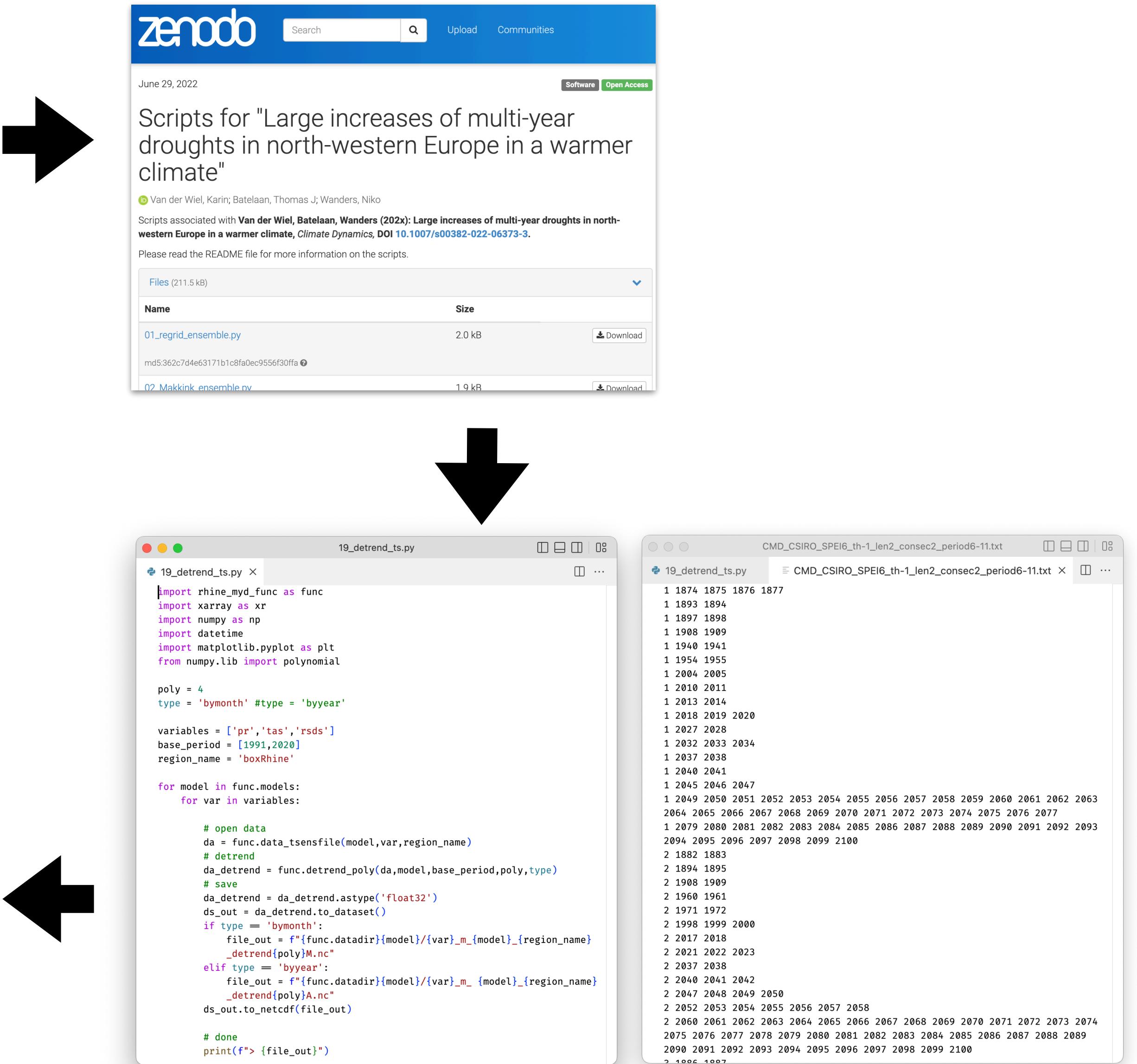
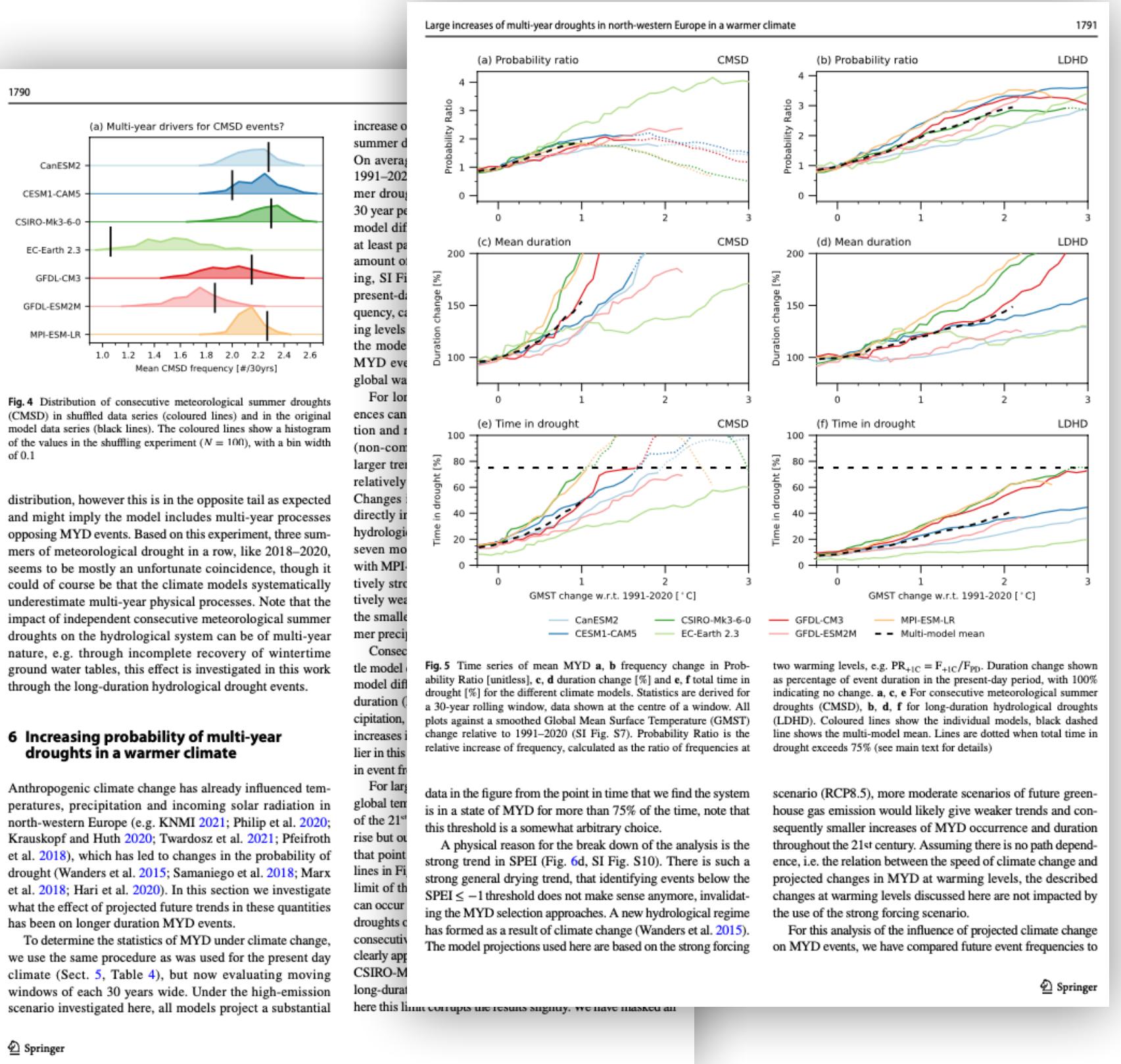
But what else?

Not a magic sauce that makes outcomes transparent and intelligible

Complex relationships between data and outcomes need to be **accessible**, **intelligible**, and easy to **evaluate** and **reuse**

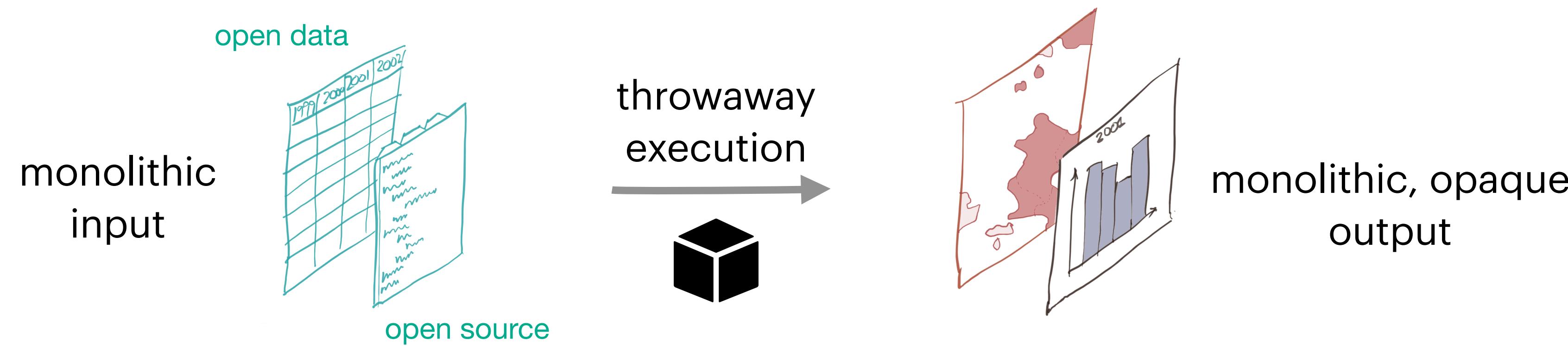


Open science UX today



From open source and open data to “open computation”

We can and should continue to improve open science practices, but there is a problem:



Crucial **relationships** between data and outputs established during execution are lost and are very hard to reconstruct after the fact

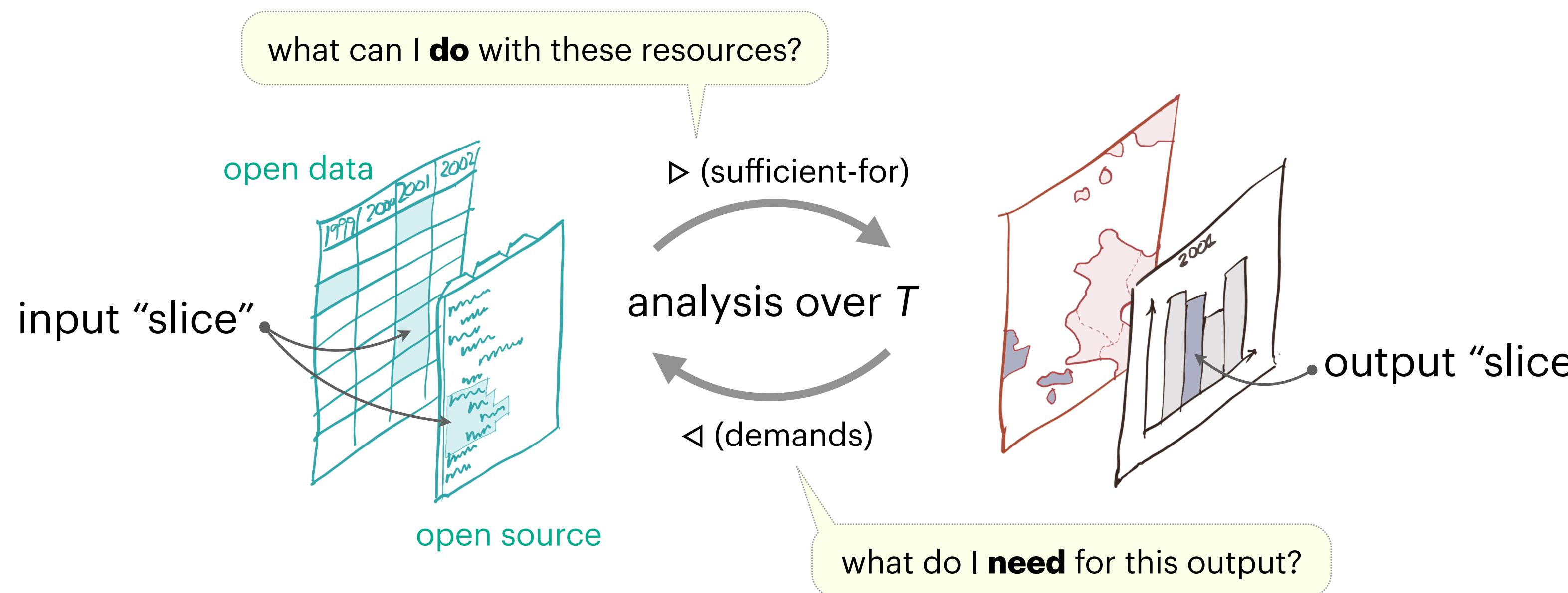
My research: can techniques from **programming languages** help us rethink this picture?

Galois slicing (i)

Replace run-once execution by **persistent and fine-grained bidirectional relationship**

program produces a **trace T**
as well as an output

pair of functions $\triangleleft \dashv \triangleright$ relating
input and output “slices” called a
Galois connection



Bake this metadata into
computed artefact as a kind
of **built-in transparency**

Demo: “data-transparent” visualisations

Research prototype called **Fluid**

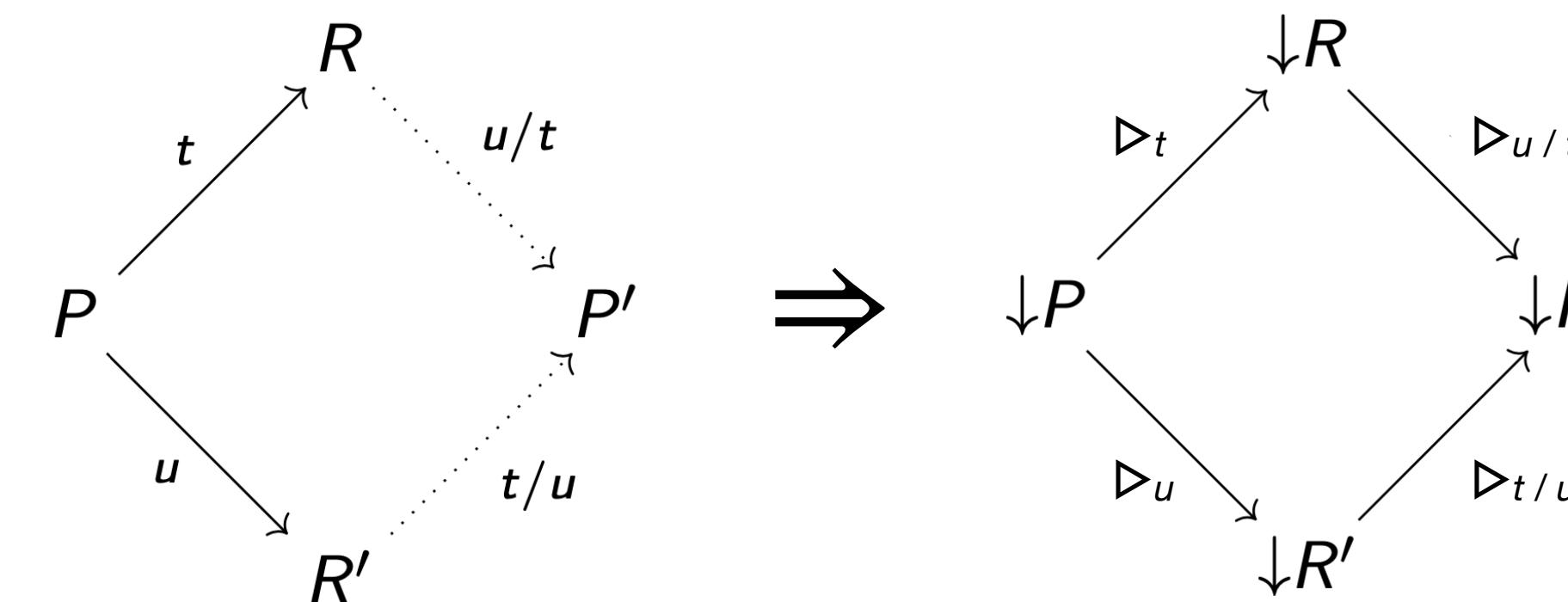
Galois slicing (ii)

Later work extended to mutable state, exceptions, and **concurrency**

To define \triangleleft and \triangleright for a **specific execution**, need to distinguish:

- **observable** non-determinism (e.g. choice of who to communicate with)
- **unobservable** non-determinism (interleaving of independent steps)

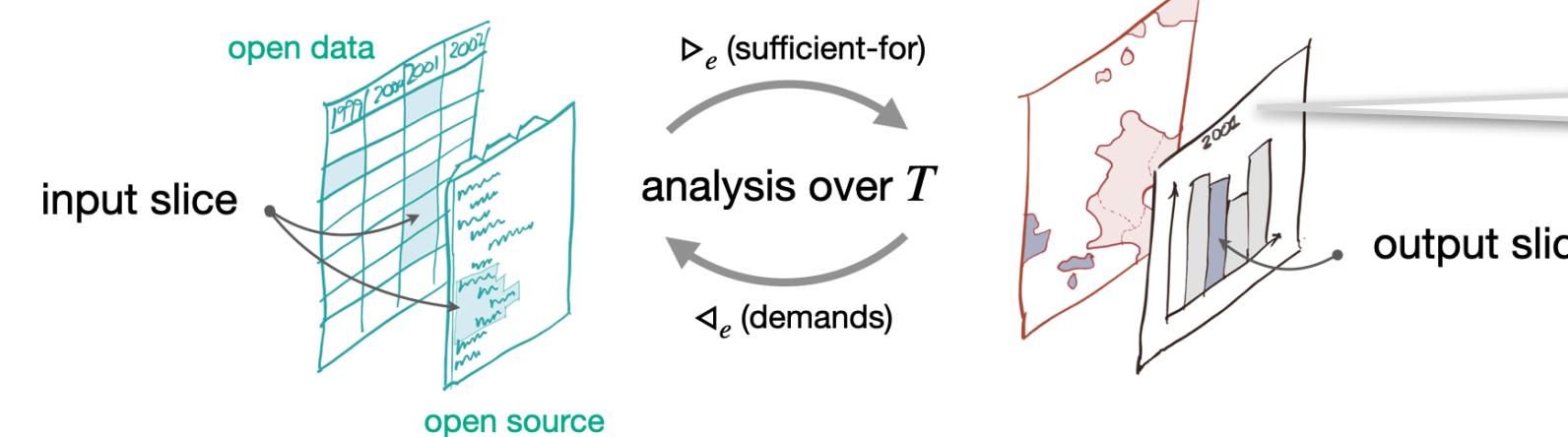
Executions are
causally equivalent iff
they differ only in
unobservable choices



Causally equivalent
executions induce
**equivalent slicing
operators** \triangleleft and \triangleright

Allows slicing to exploit concurrency whilst producing consistent results

Linked outputs (i)



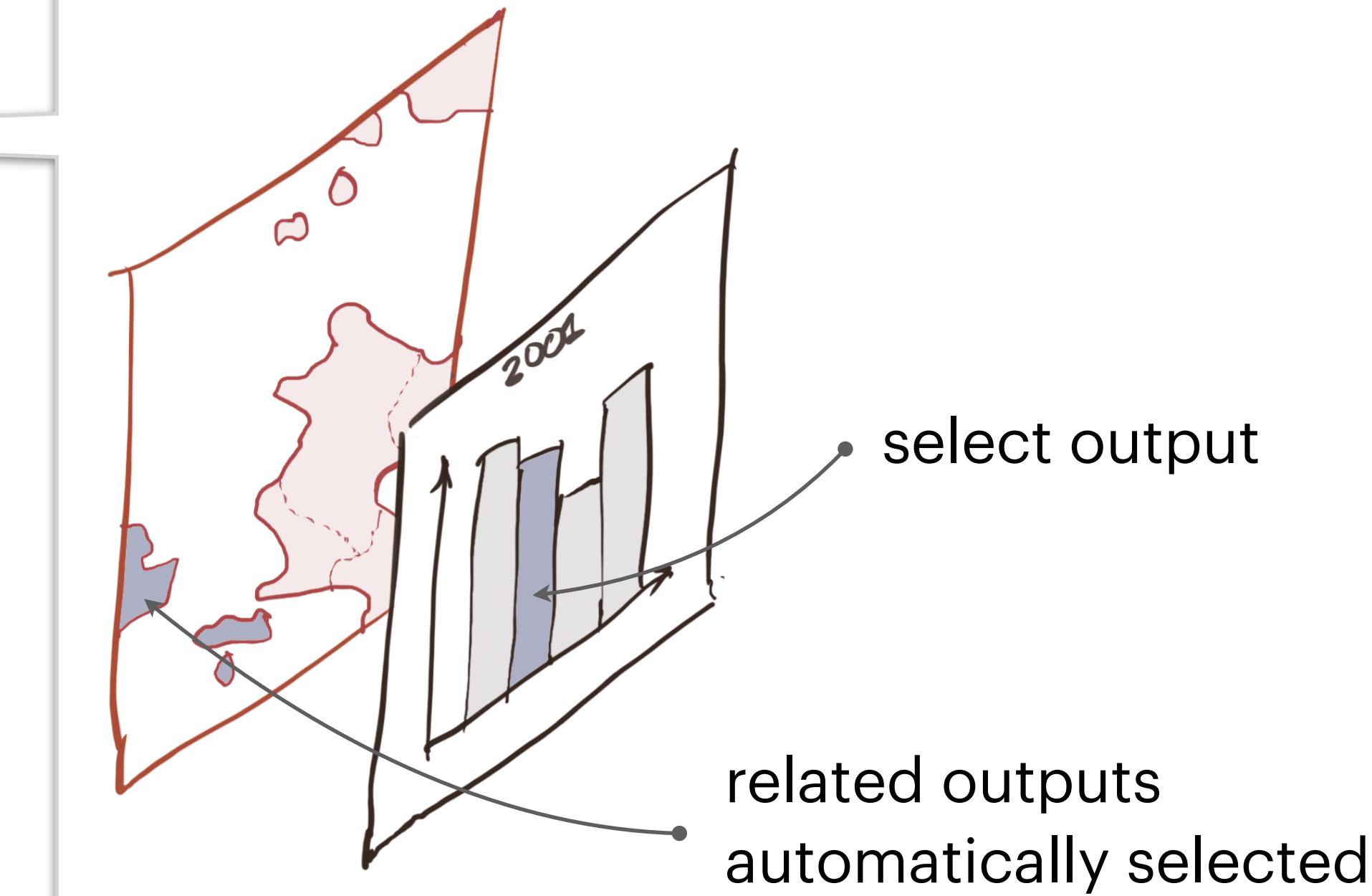
What else might support intelligibility?
What about how **outputs relate to other outputs**?

Our insight:

Outputs x and y are related if they
compete for inputs i.e. $\triangleleft x$ and $\triangleleft y$ overlap

however: not easy to compute related
outputs using only \triangleleft !

brushing and linking



In data visualisation, **brushing and linking** is an important comprehension tool, but ad hoc and poorly understood

Linked outputs (ii)

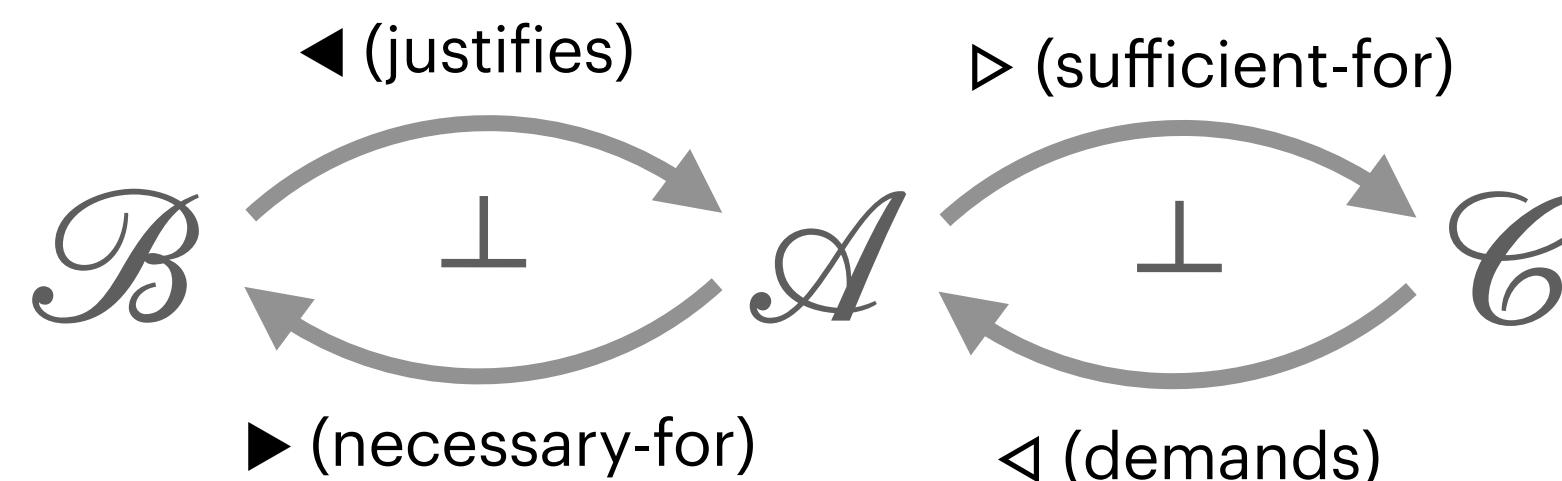
Showed how to define **related-outputs** operator $\lhd\lhd$
that picks out “competing outputs” by allowing slices to have **complements**

Previous notion of slice

- value with “holes”
- no notion of complement
- lattice

New notion (“selection”)

- subset of paths in a value
- complements
- Boolean lattice



Complements mean \lhd and \rhd
have De Morgan duals

$\rhd := \neg \lhd \neg$ **(necessary-for)**
 $\lhd := \neg \rhd \neg$ **(justifies)**

Demo 2: Linked visualisations

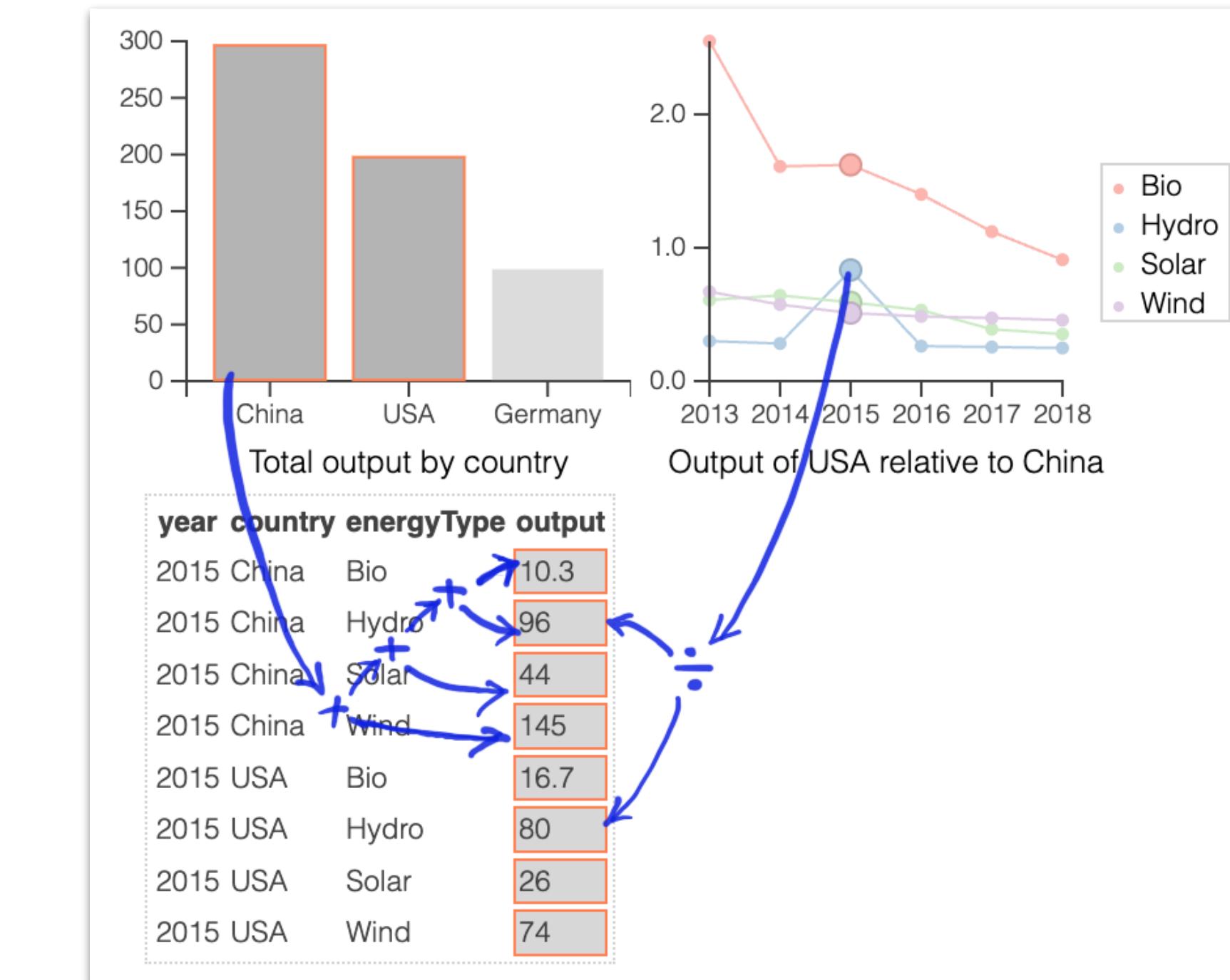
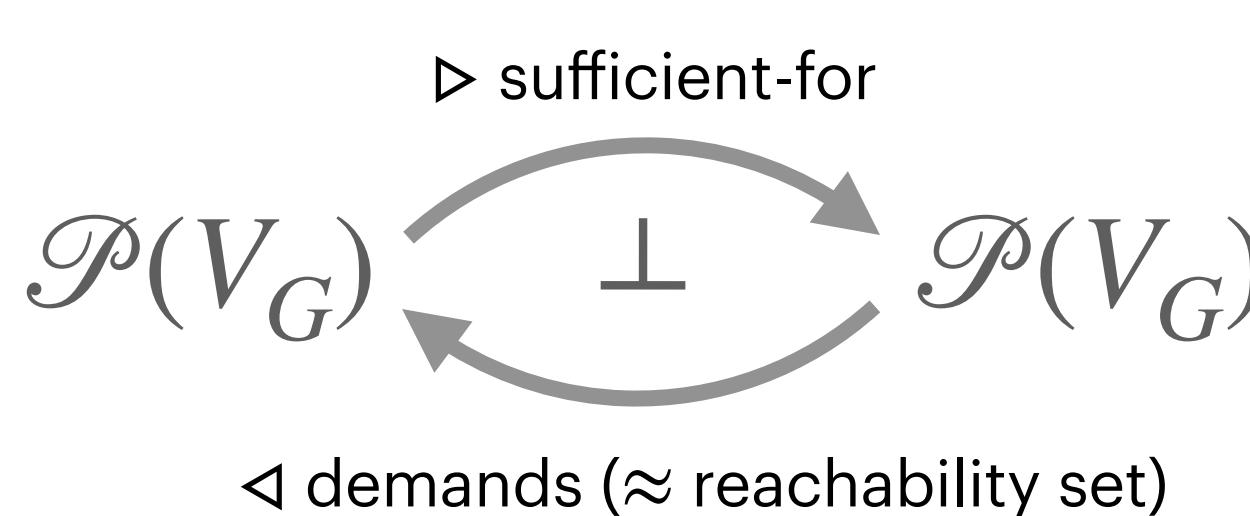
Graphical Galois slicing

with Joe Bond, Minh Nguyen, Cristina David

Dynamic dependence graphs (DDGs)

Used for program slicing (Field and Tip 1990)

- ① **Persistent execution** provided by dependence graph G instead of trace

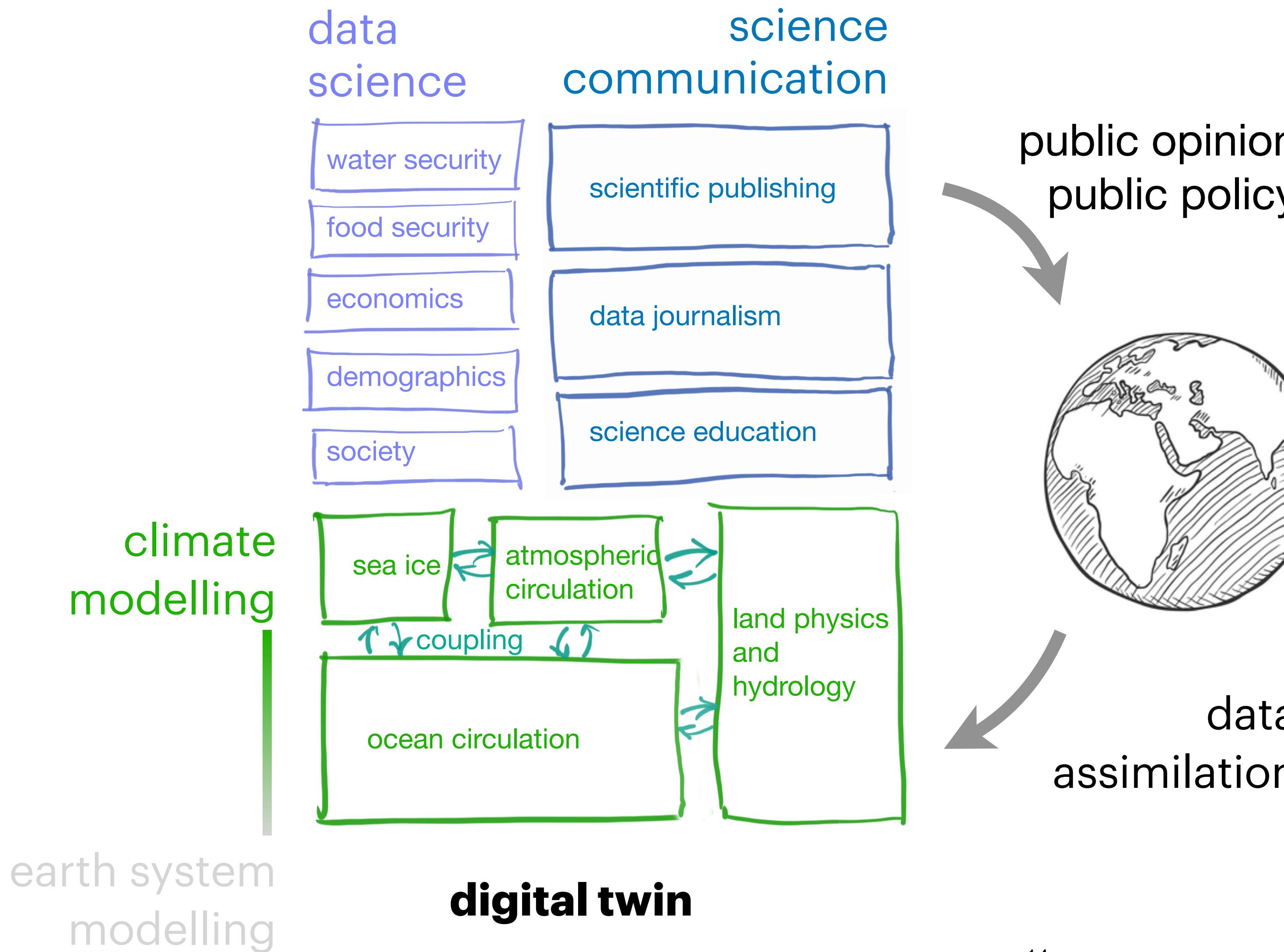


- ② **Fine-grained analyses** now defined over G (language-agnostic)

- G directly represents **causal equivalence class** so no rewinding of irrelevant steps
- asymptotically faster in common parallel cases
- **stages** execute concurrently

Towards transparent computing for climate science

Computing for climate science



Transparent software

would improve

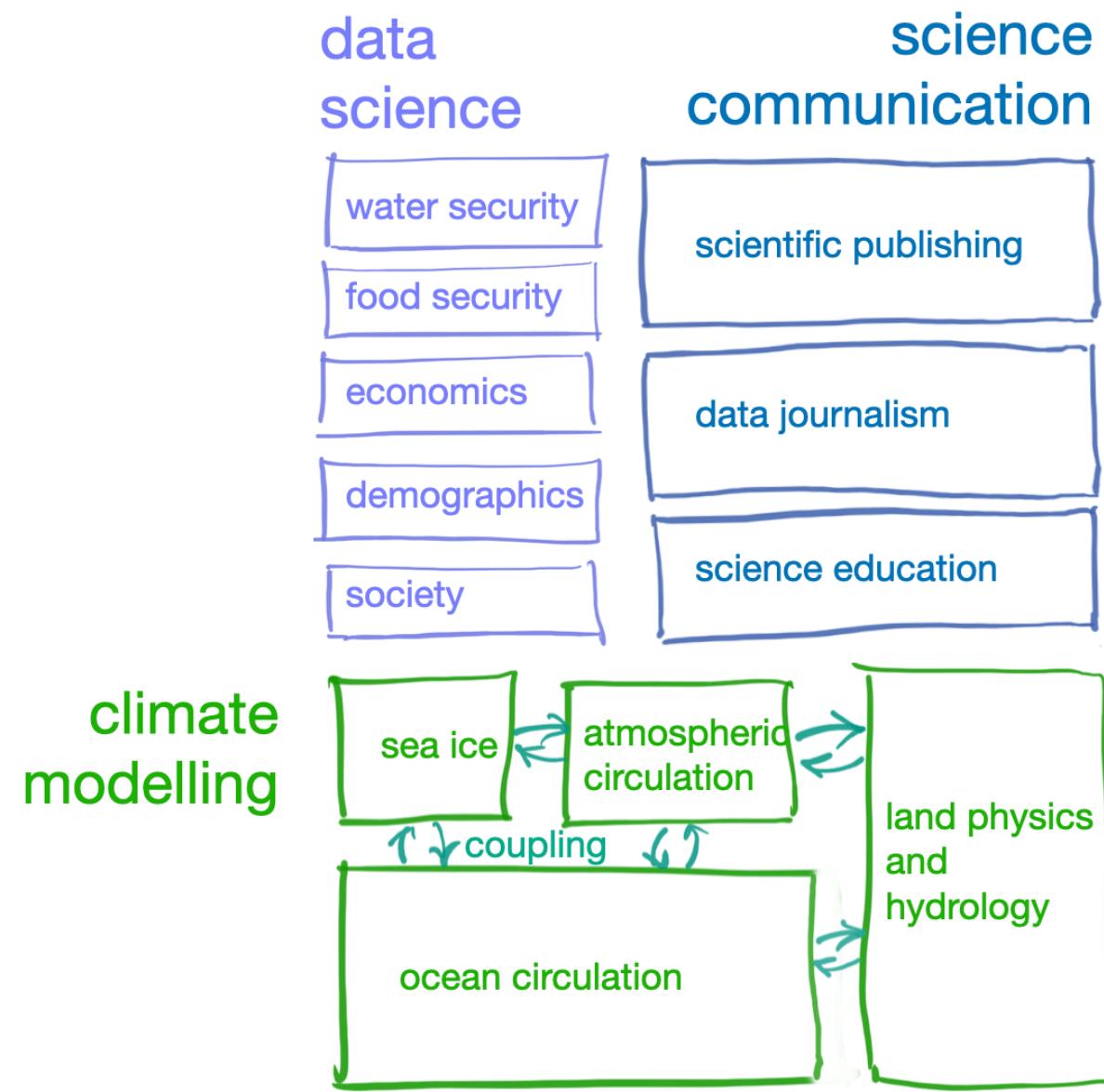
- peer review
- reproducibility
- reuse
- diagnostics

More **efficient** science

- inclusivity
- participation
- accountability
- trustworthiness

More **responsible** science

However: some challenges



Value proposition of pervasive transparency is clear

But hard to achieve in **heterogeneous, computationally demanding** systems

multi-paradigm

numerical simulation, probabilistic, ML, data analysis, streaming data, big data, visualisation, ...

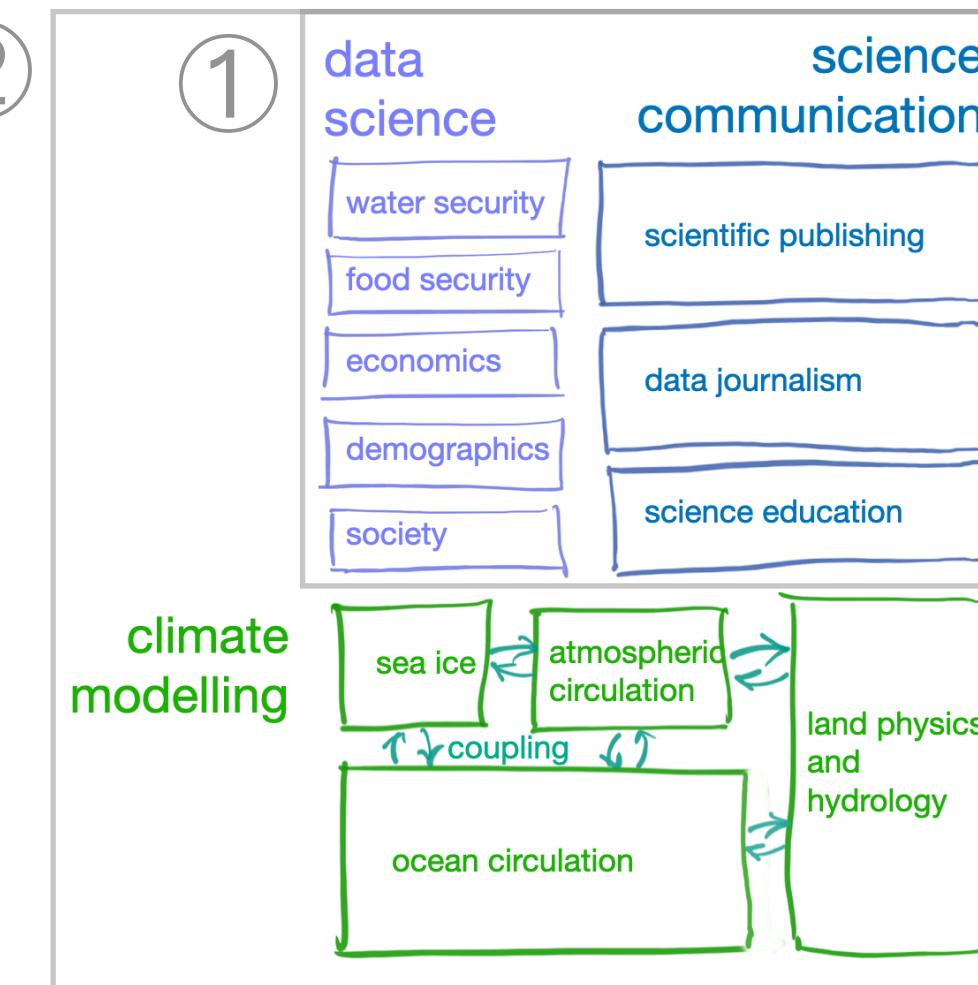
multi-language

C/C++, Fortran, R, Python, Julia, JavaScript, SQL, ...

So, what to do? Two research strategies

① Focus efforts where such features are most valuable and least costly:
climate science communication

- transparency at a premium
- computationally simpler
- less established tools



② More generally: design explicitly for **gradual adoption** and **scalability**

- techniques for interoperable, partial & approximate transparency in heterogeneous systems

Transparent research artefacts (i)

Today: papers, news articles and educational tools are opaque

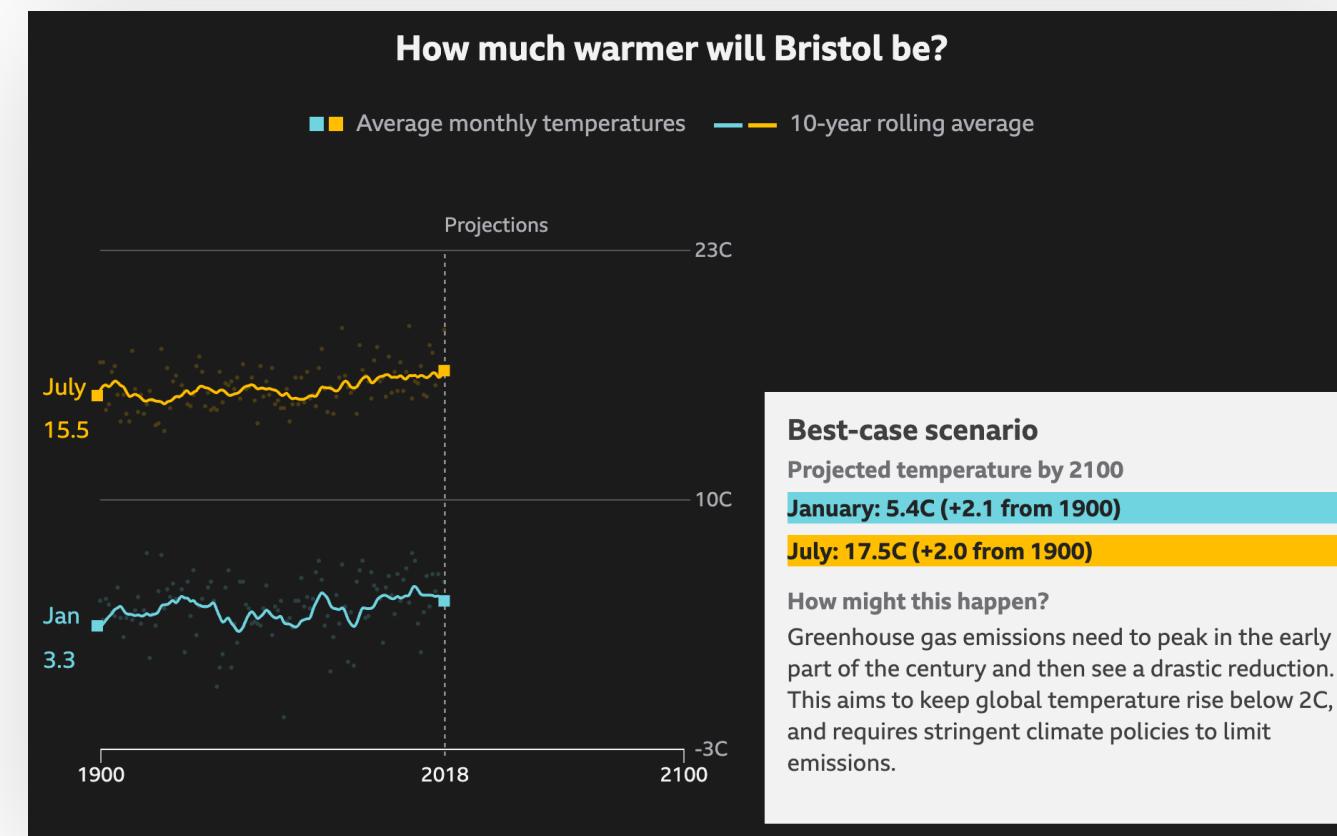
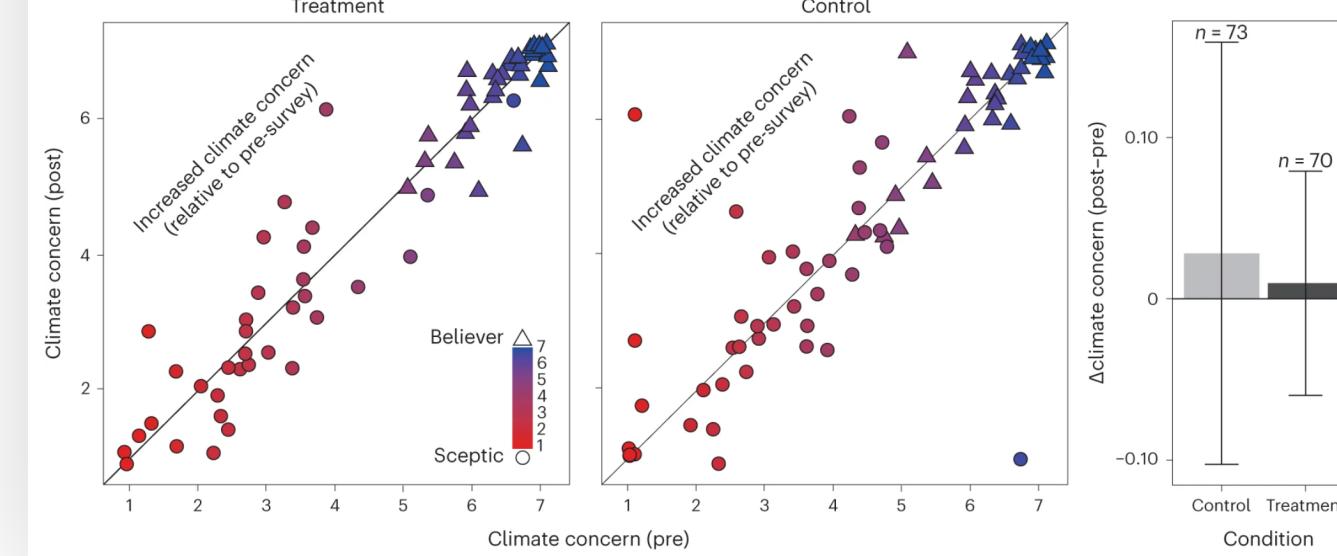
Goal: online artefacts which are **transparent and self-explaining** via rich automated interactions, building on Fluid prototype

Climate science communication today

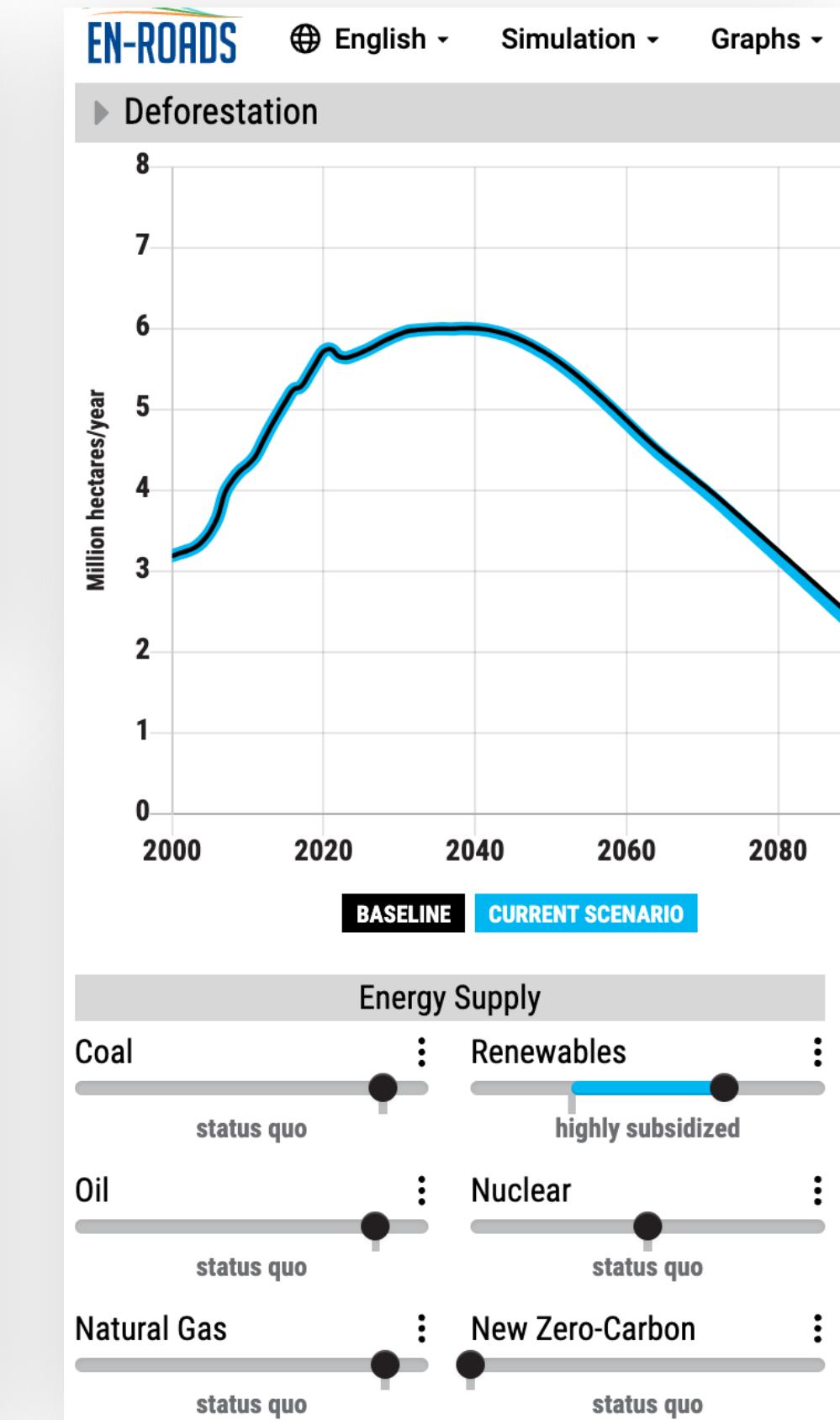
research papers

Fig. 2: Distributions of climate beliefs before and after participating in the climate market.

From: [Participating in a climate prediction market increases concern about global warming](#)



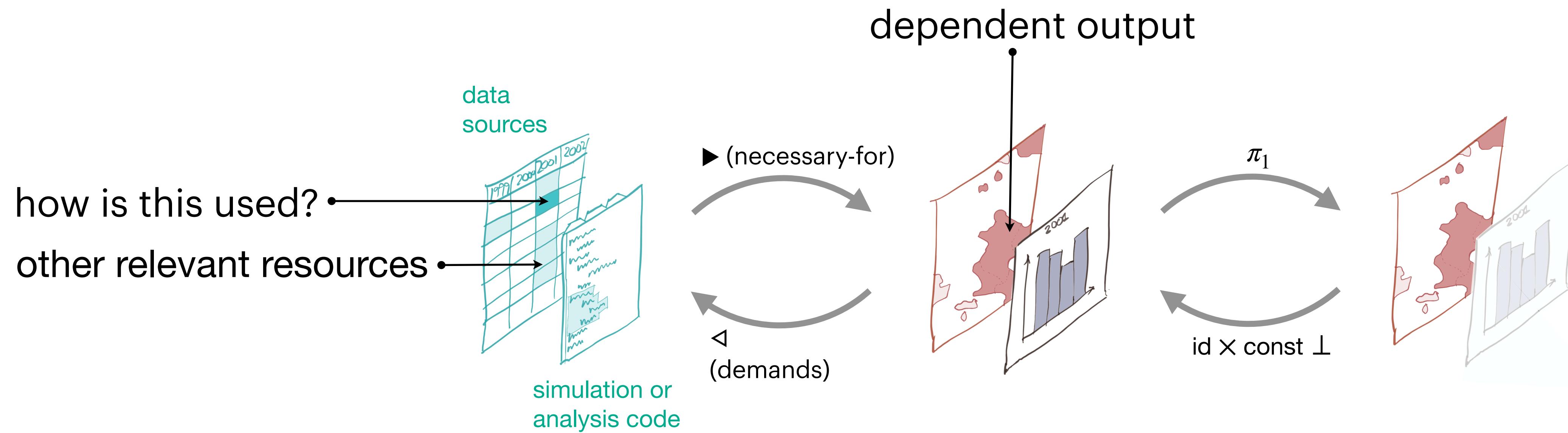
visual journalism



interactive simulation

Transparent research artefacts (ii)

Automated interactions to support **comprehension** and **reuse**



related-inputs := $\blacktriangleright\triangleleft (\text{demands} \circ \text{necessary-for})$

Mostly engineering/UX work

Other adjoint operators for focusing views
on **geographical regions** or **time slices**

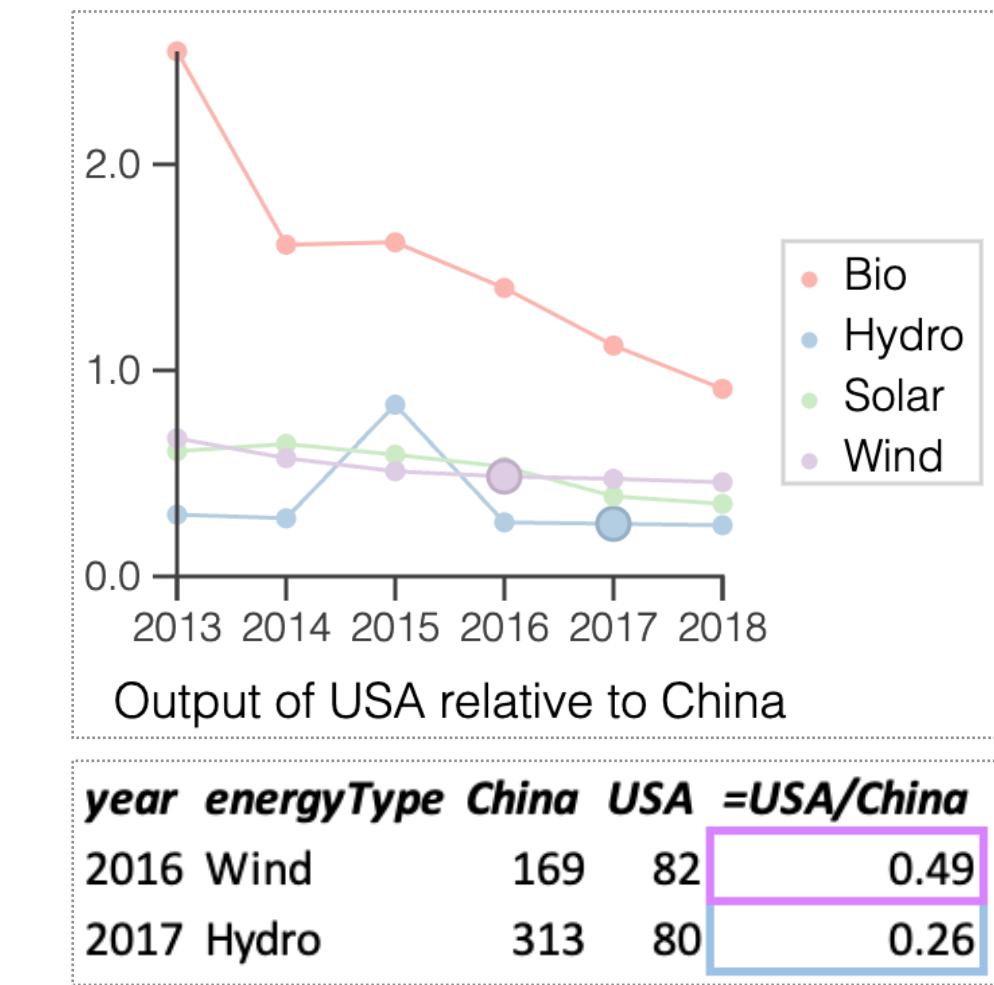
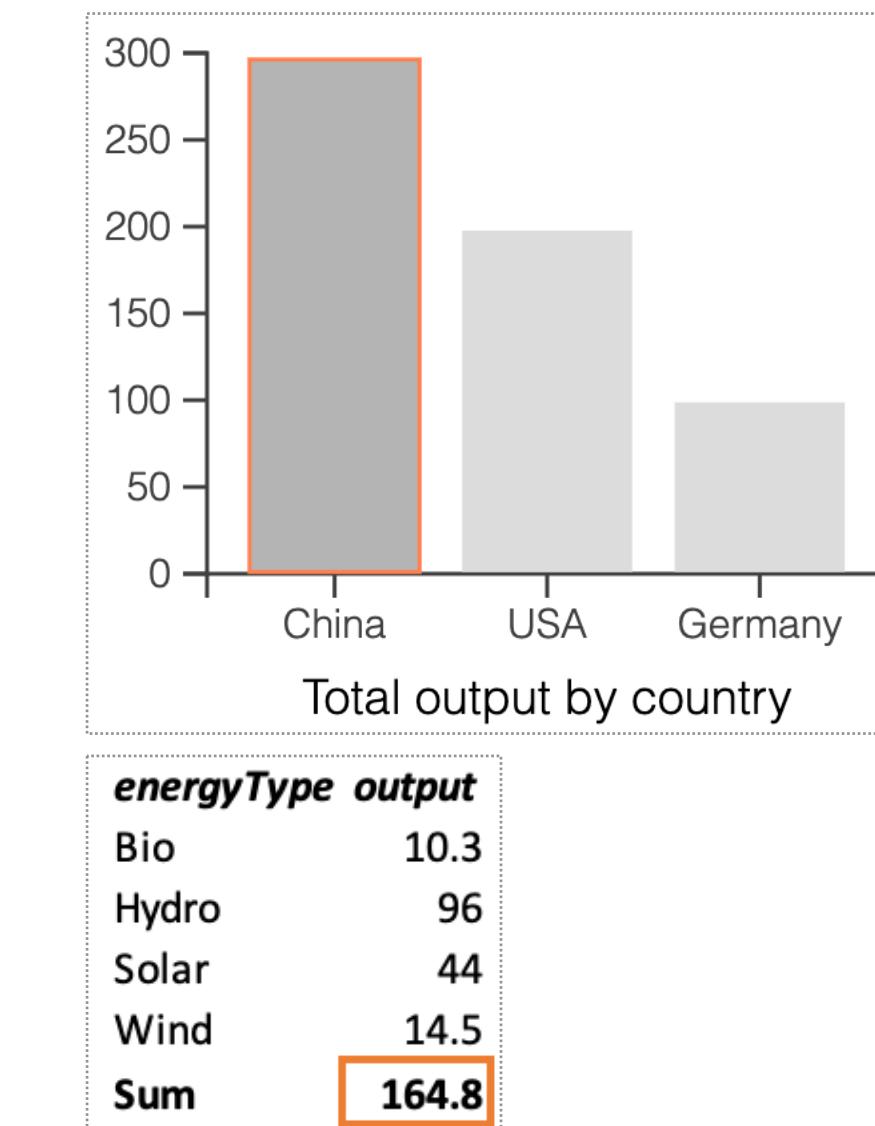
Transparent research artefacts (iii)

Self-explaining computation: equip output selections with intensional explanations (how, not just what)

Related work:
executable slicing, expression provenance

Some challenges

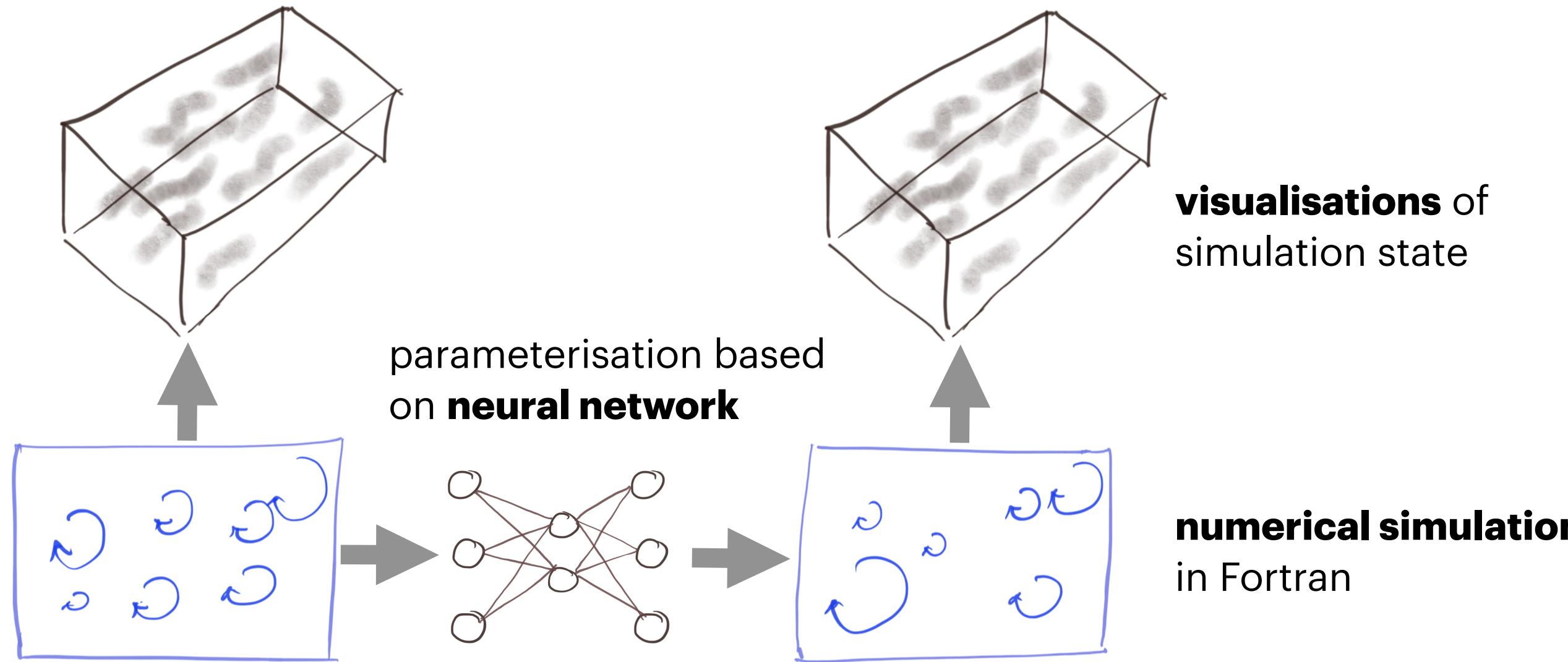
- naming and reference
- domain-specific explanations



- pivoting
- column-wise operations
- rollup along dimension

Use to implement climate articles and explorable simulations like En-ROADS that are **transparent and self-explanatory**

Gradual, interoperable transparency (i)



Back to research direction ②

- complex, heterogeneous systems
- transparency solutions need to **interoperate** and **scale**

Research approach

- treat gradual, interoperable explainability as first-class problem

Example

Subgrid-scale “parameterisations”: pluggable components of climate simulations

How to integrate XAI for neural network with other explainable components?

Gradual, interoperable transparency (ii)

Partial and interoperable transparency

Interoperable

- Compositionality of explainability notions (e.g. quantitative with symbolic)
- Round-tripping properties
- Language interoperability (FFIs)

Partial and gradual

- Flexible instrumentation: disable in production code, retain for debugging or teaching
- Approximate (trade precision for performance)
- Partial (trade intensional information)

Partial transparency

Given dependence graph $G = (V, E)$
if we elide internal nodes but preserve
connectivity, resulting graph will yield the same
Galois connection $\triangleleft \dashv \triangleright : \mathcal{P}(V) \rightarrow \mathcal{P}(V)$

Points to a way to give up intensional
explainability for performance

ExPLan: Explainable
Programming Languages
Project lead: James Cheney
Submitted to ERC

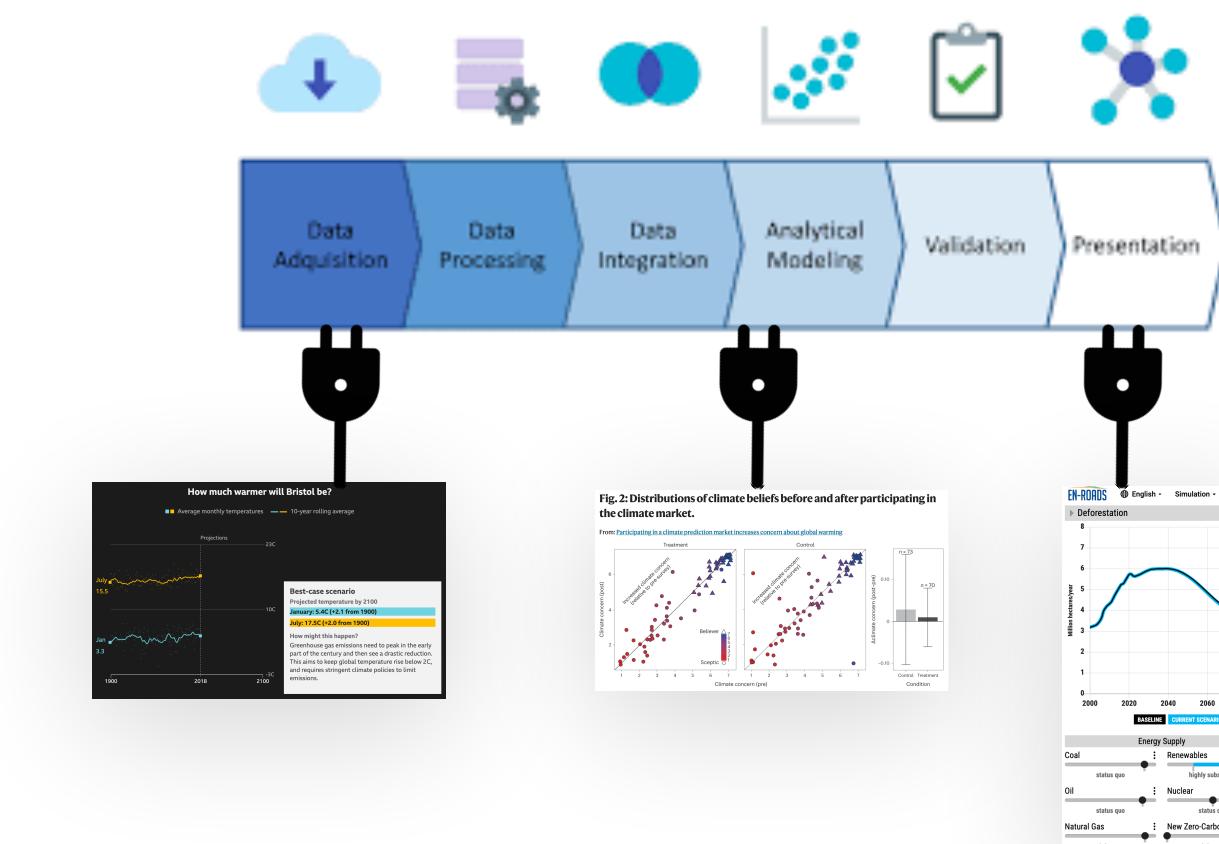
Conclusion: open climate science needs “open computation”

Accessible, intelligible, assessable and **(re)usable** computational science needs better practices but also **programming languages that are more transparent**

Climate science presents specific challenges to this idea (but also the perfect test bed)

More than anywhere else we want:

- engaged citizens
- transparent governance
- accountable policy makers
- intelligible science



With developments in PL infrastructure we can gradually replace **opaque snapshots** with **transparent views of actual workflows** (different views for different stakeholders)

Thank you 😍

Roly Perera
Institute of Computing for Climate Science, University of Cambridge
School of Computer Science, University of Bristol



Institute of
Computing for
Climate Science