

some questions ....

Did you once find data on your or the institute's computer and no one knew any details about this data anymore because forgotten or the person who generated the data or put it there is no longer at the institute?

You once needed data from someone else, e.g. to compare your results, but perhaps also a script or part of model code. You may even have known the person personally .. BUT you had to wait forever for the data, were put off because of lack of time and got the data late or maybe not at all?

Or perhaps you have already been asked for data from a previous project, but then simply did not have the time to look up the data again or provide it in time?

How many of you have several versions of the final report of a project ...

```
.  
..  
DETER_FinalReport_final.doc  
DETER_FinalReport_final2.doc  
DETER_FinalReport_final2_mb.doc  
DETER_FinalReport_final2_mb_ia.doc  
DETER_FinalReport_final2_mb_pr.doc  
DETER_FinalReport_final2_mb_pr_final.doc  
DETER_FinalReport_final2_mb_pr_final2.doc  
DETER_FinalReport_final2_mb_pr_final3.doc
```

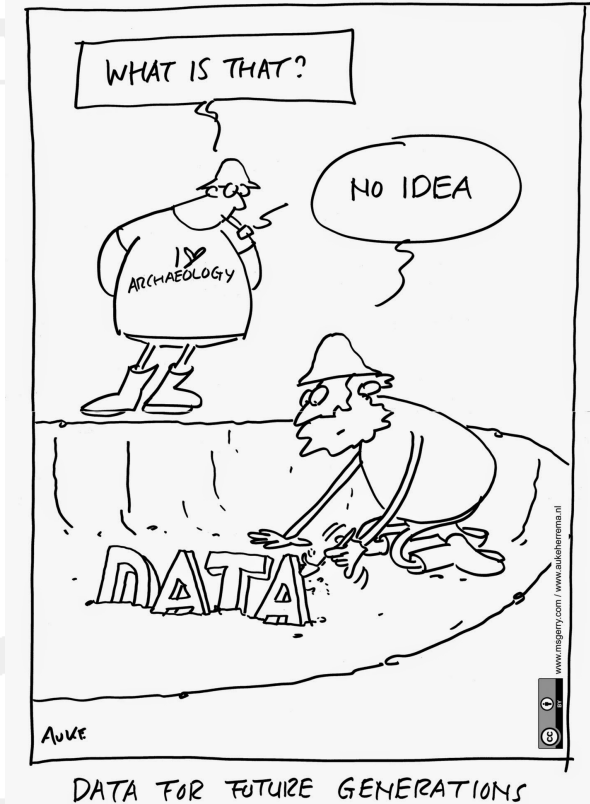
... and don't even really have the final version?

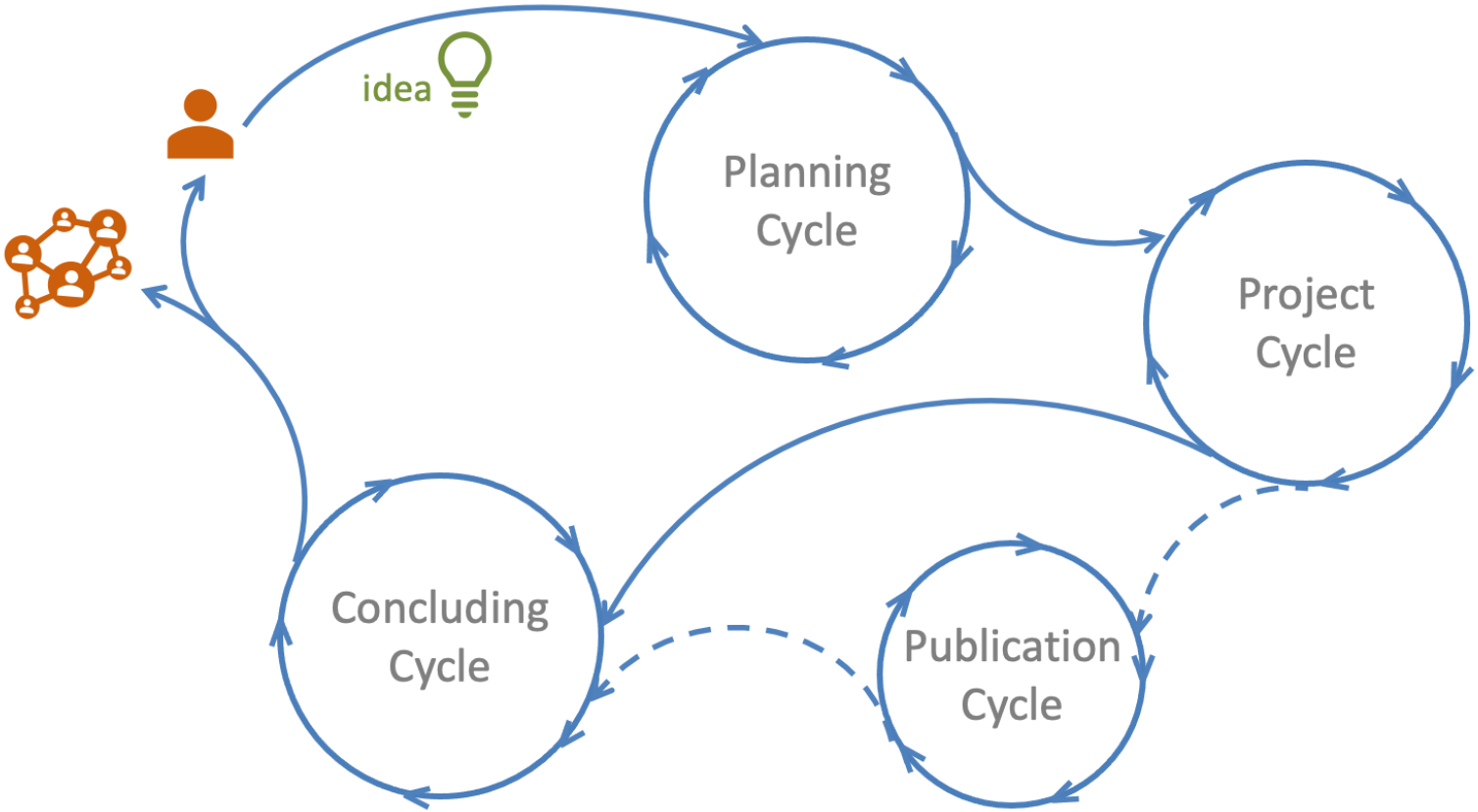
# Research Data Management Services, current Workflows and DMPs

DKRZ USER WORKSHOP

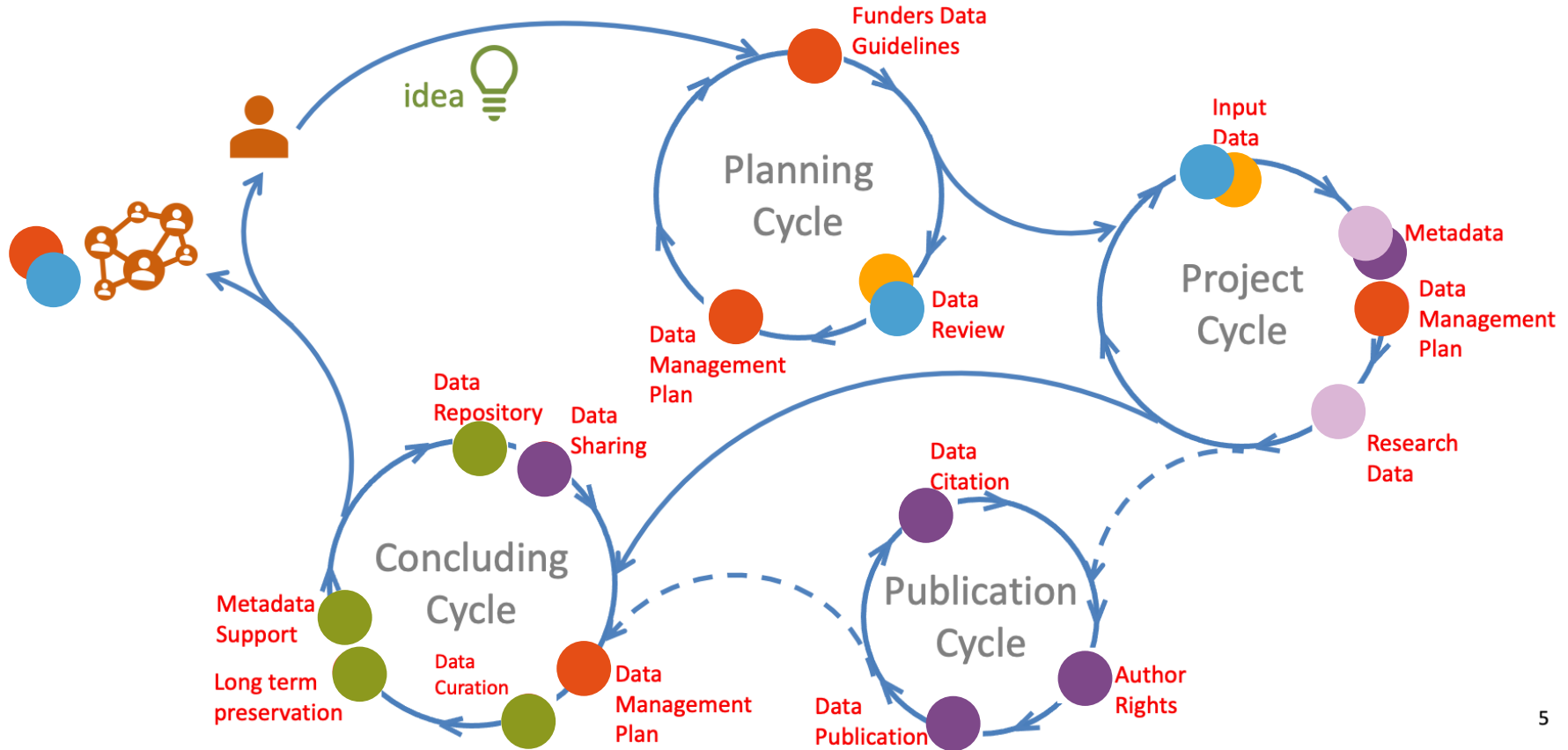
13. 10. 2022

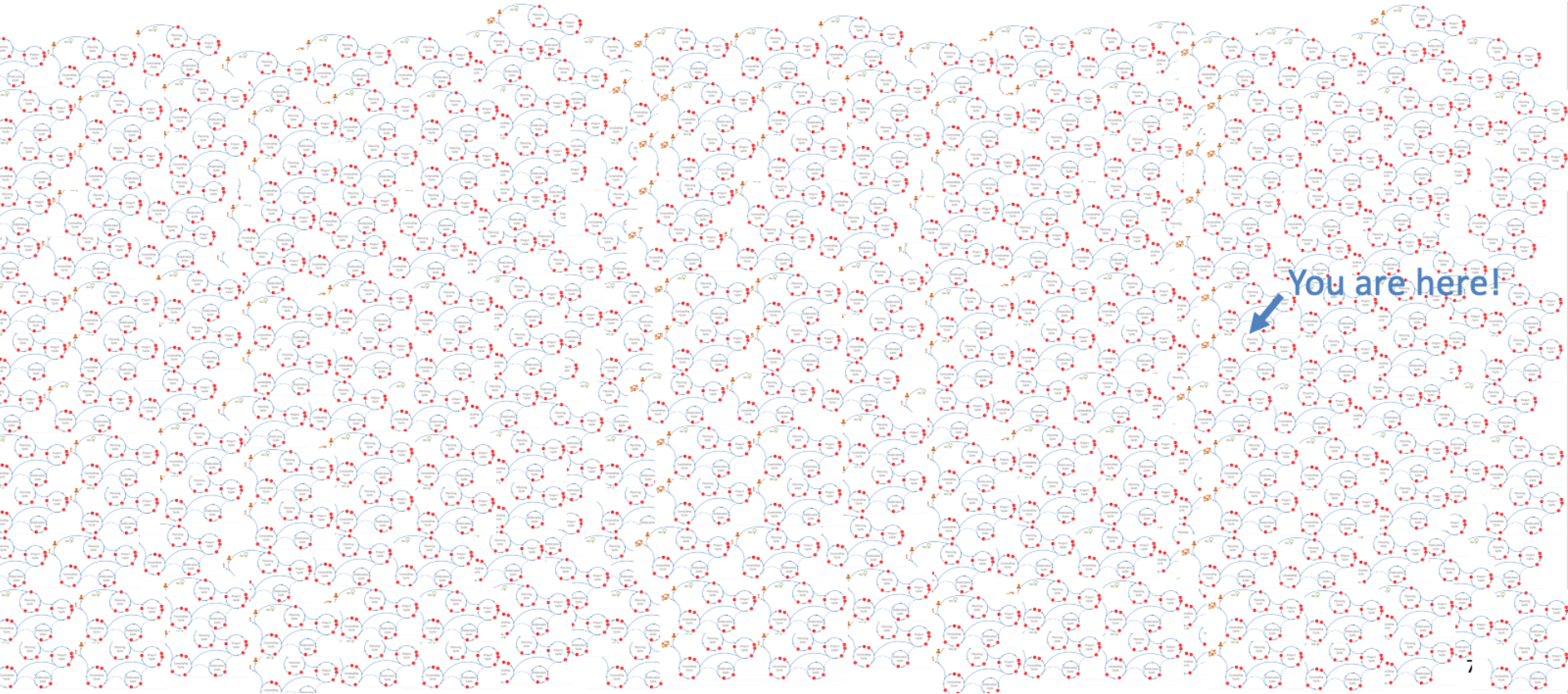
Ivonne Anders, Karsten Peters-von Gehlen  
& ALL@DM





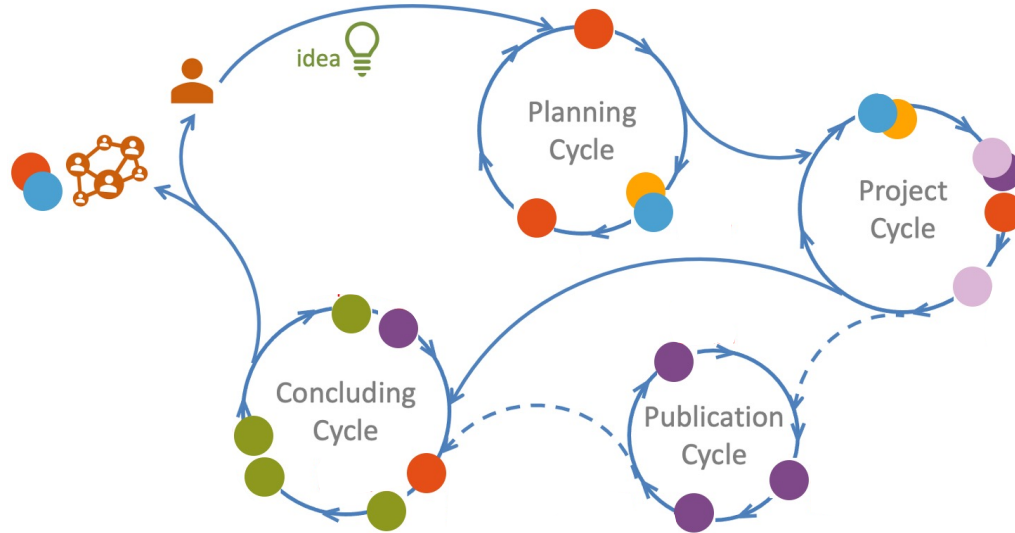




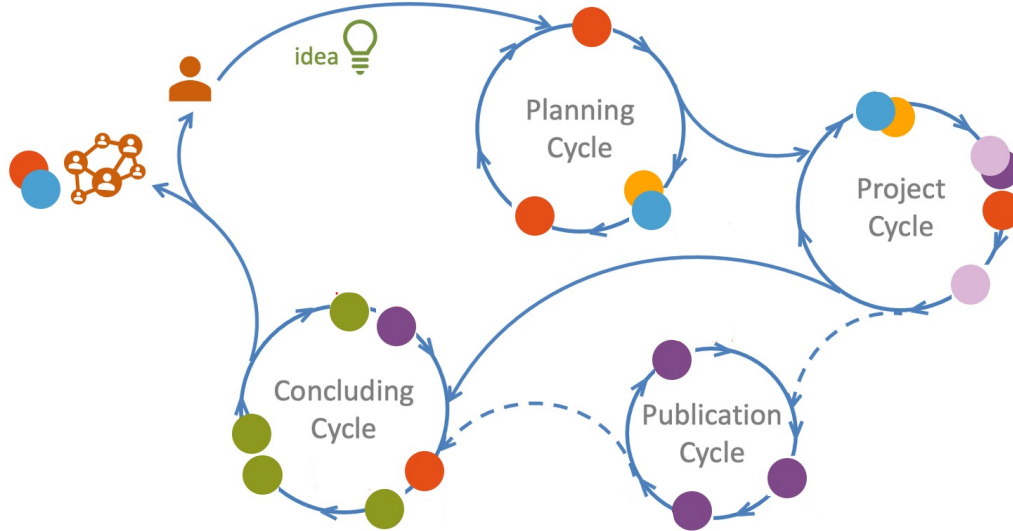


... current developments in Earth System Modeling demand for efficient workflow strategies to cope with ever increasing data amounts, complexity of analyses, server-side processing and distributed databases

## Project Cycle

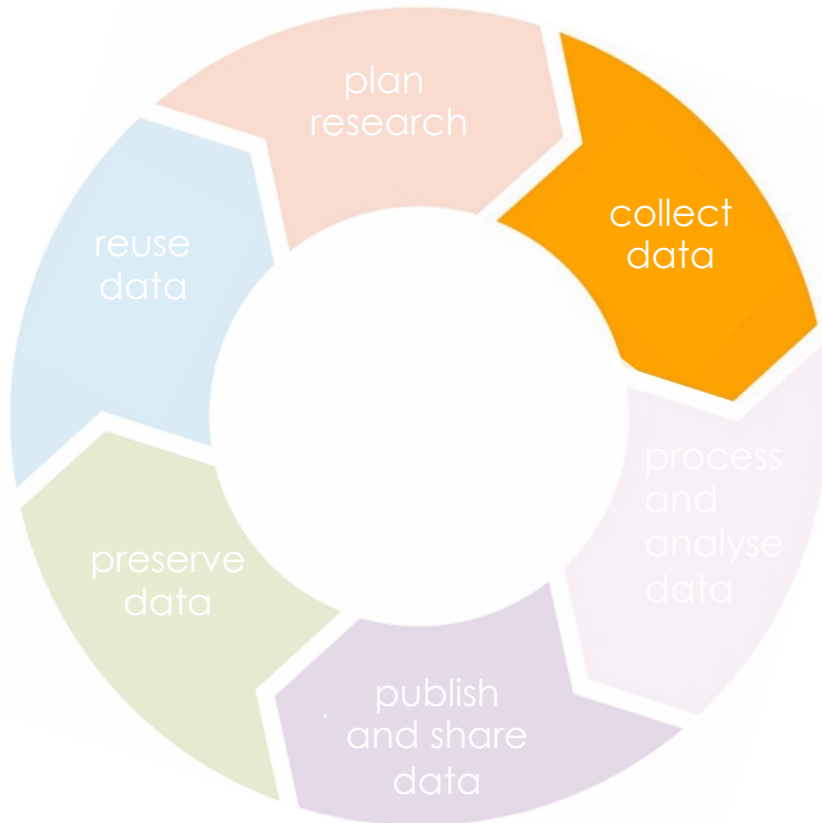


## Project Cycle



## Data Life Cycle





- Model simulation output
  - Everything related to Earth System Science
- Data collections hosted by DKRZ
  - CMIP3/5/6
  - CORDEX
  - ERA5, ERA-Interim
  - ...
- Archived data
  - Project-based tape archive
  - World Data Center for Climate (WDCC)
  - DOKU



<https://www.wdc-climate.de/>



- Discipline-specific and centrally maintained software stack
  - Dedicated support available
  - Expanded if required

- Intake catalogs

- CMIP5/6
- CORDEX
- ERA5
- DYAMOND
- nextGEMS



- Experimental STAC implementation



- Jupyterhub @DKRZ

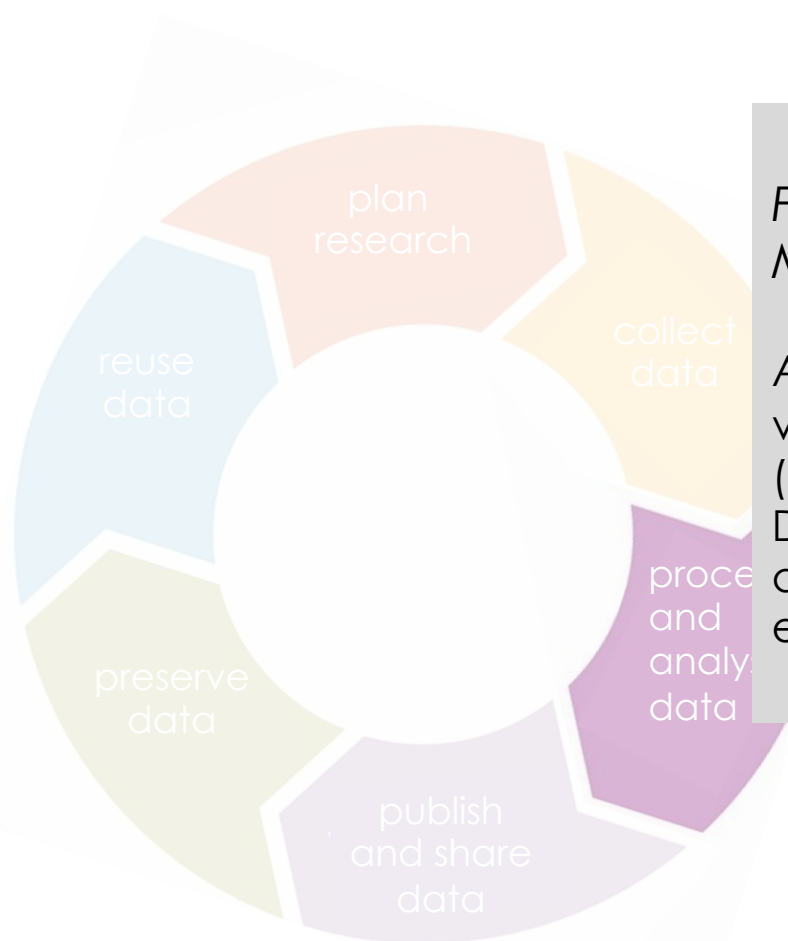
- Direct access to HPC resources (compute and disk)

<https://jupyterhub.dkrz.de/>



- Data standardization support

- Community- or project-specific (depending on user requirements)



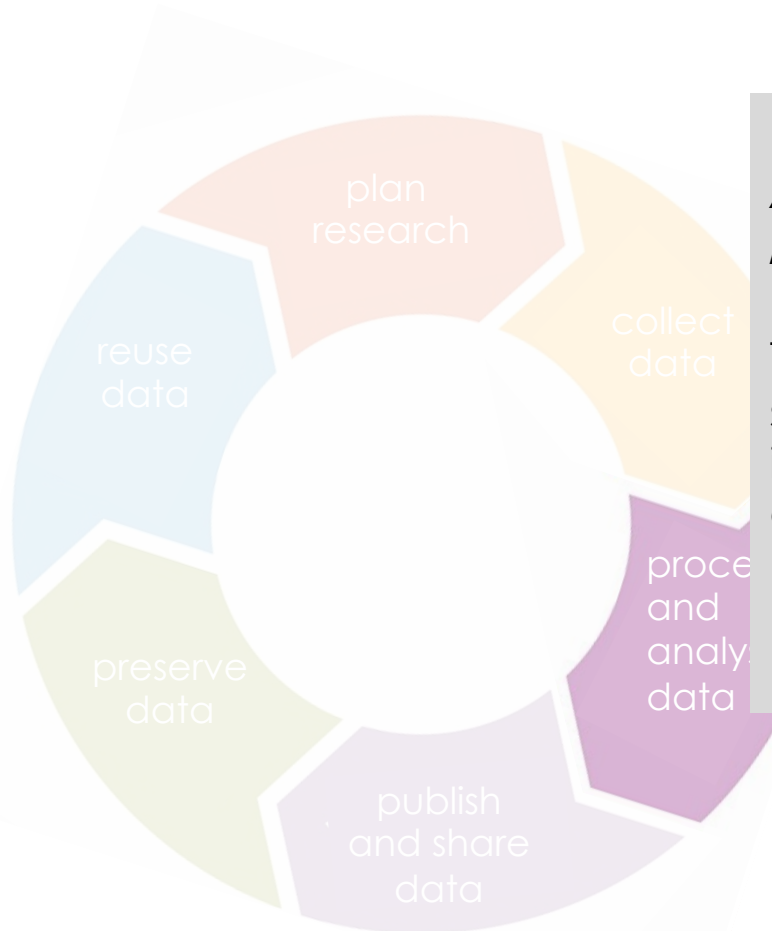
*Fabian Wachsmann and Marco Kulüke*

Availability and accessibility of large-volume datasets hosted at DKRZ (CMIP3/5/6, ERA5, CORDEX, DYAMOND,...) using metadata-driven data access (catalogs), including an extensive hands-on session.



- Direct access to HPC resources (compute and disk)  
<https://jupyterhub.dkrz.de/>
- Data standardization support
  - Community- or project-specific (depending on user requirements)





*Angelika Heil and  
Martin Schupfner*

The basics and benefits of data standardization including an introduction of tools and services offered at DKRZ + an extensive hands-on session



- Direct access to HPC resources (compute and disk)  
<https://jupyterhub.dkrz.de/>
- Data standardization support
  - Community- or project-specific (depending on user requirements)



**SWIFT**

- Sharing via DKRZ cloud (swift object store, 1.5PB current capacity) <https://swiftbrowser.dkrz.de/>

- Publication in WDCC or DOKU (long-term archival services)

- DataCite DOI for WDCC



<https://www.wdc-climate.de/>

- Dissemination via ESGF



- Jupyter-Notebooks

 **jupyterhub**  
<https://jupyterhub.dkrz.de/>

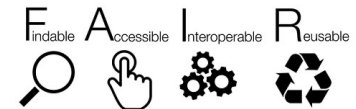
- Project-related support for organizing data publication/sharing/dissemination



- Tape archive for everyday use **STRONGLINK**
  - More than 300 PB available
  - Metadata-driven search and access
- DOKU
  - Archival of project-related reference data
  - PIDs
  - Available to DKRZ users
- WDC (World Data Center for Climate)
  - CoreTrustSeal certified domain-specific repository
  - FAIR compliant ([Peters-v. G. et al., 2022](#))
  - DOIs available
  - Available to the global community



<https://www.wdc-climate.de/>



This Photo by Unknown Author is licensed under [CC BY](#)



*Daniel Heydebreck and Andrej Fast*

Using the DKRZ tape archive (StrongLink) and its associated command-line tool `slk` in your everyday workflow with a focus on its metadata-enabled data access framework; also including hands-on

STRONGLINK

access

reference data

(Climate)  
in-specific

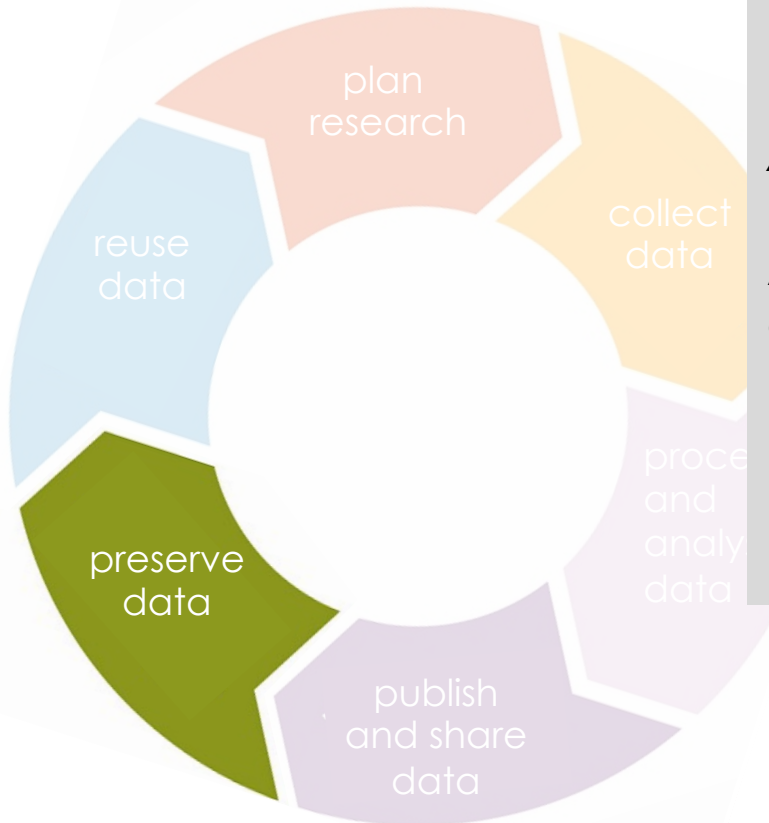
- FAIR compliant ([Petris v. O. et al., 2022](#))
- DOIs available
- Available to the global community



<https://www.wdc-climate.de/>



This Photo by Unknown Author is licensed under [CC BY](#)



*Andrea Lammert*

An overview of DKRZ long-term archiving services WDCCL and DOKU

STRONG LINK

access

reference data

(Climate)  
in-specific

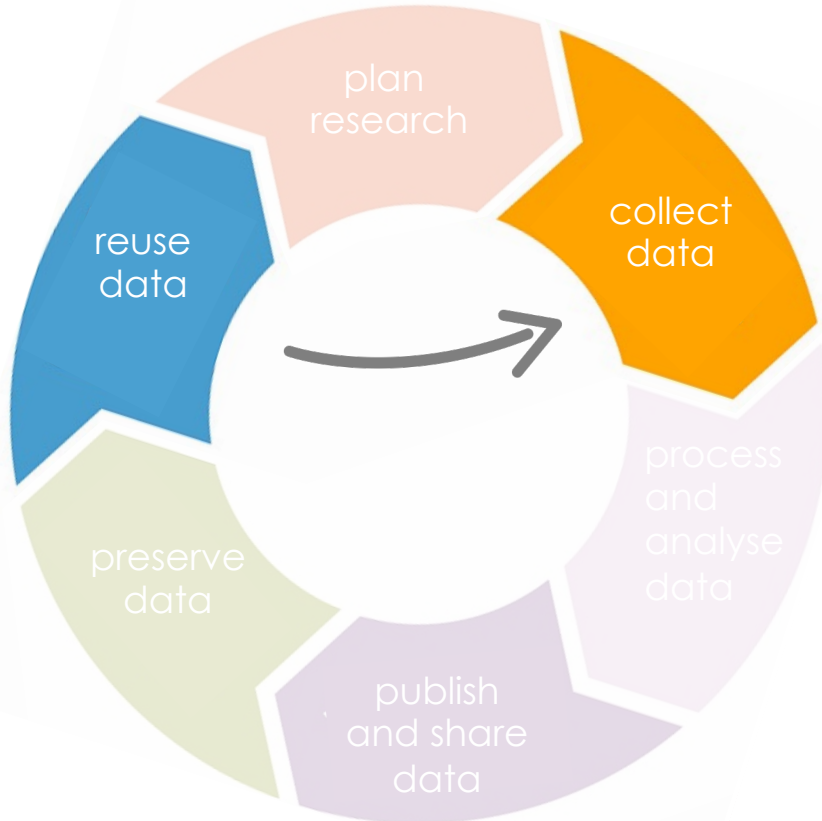
- FAIR compliant ([Peters v. O. et al., 2022](#))
- DOIs available
- Available to the global community



<https://www.wdc-climate.de/>



This Photo by Unknown Author is licensed under [CC BY](#)



[www.digitalbevaring.dk](http://www.digitalbevaring.dk)



We always recommend a good planning!

See a Data Management Plan as your “Project’s Brain” and a therefore a living document.



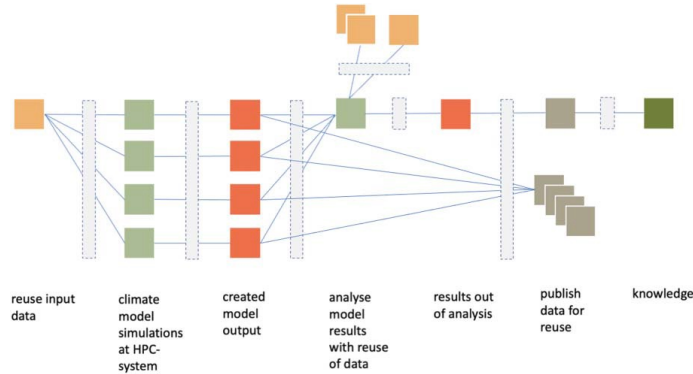
Created by Teewara soontom  
from Noun Project

We support e.g. large projects.  
Get in contact with us!

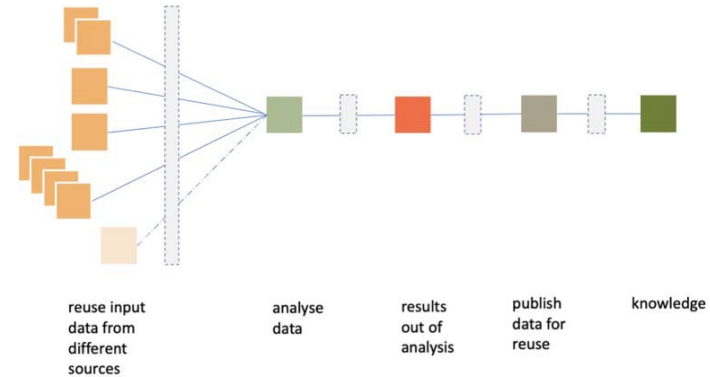
... current developments in Earth System Modeling demand for efficient workflow strategies to cope with ever increasing data amounts, complexity of analyses, server-side processing and distributed databases



## The modeller



## The data analyst / reuser



## Workflow elements



Anders et al., 2022

“efficient access to and discovery and processing of (very) large in-house and remote datasets”

“end-to-end ESM workflow solution including automated output data handling”

Keep model output routines untouched

Publication-ready data package in-line with FAIR principles

Cataloguing of performed ESM simulation setups to enable efficient reuse of large datasets for different scientific questions

Automated setup of model runs given research question (or reuse of existing data)

Online diagnostics and data compression

Comprehensive metadata cataloguing and metadata-driven access across storage tiers

Possibility of very specific database queries

Data formats allowing for fast access and processing

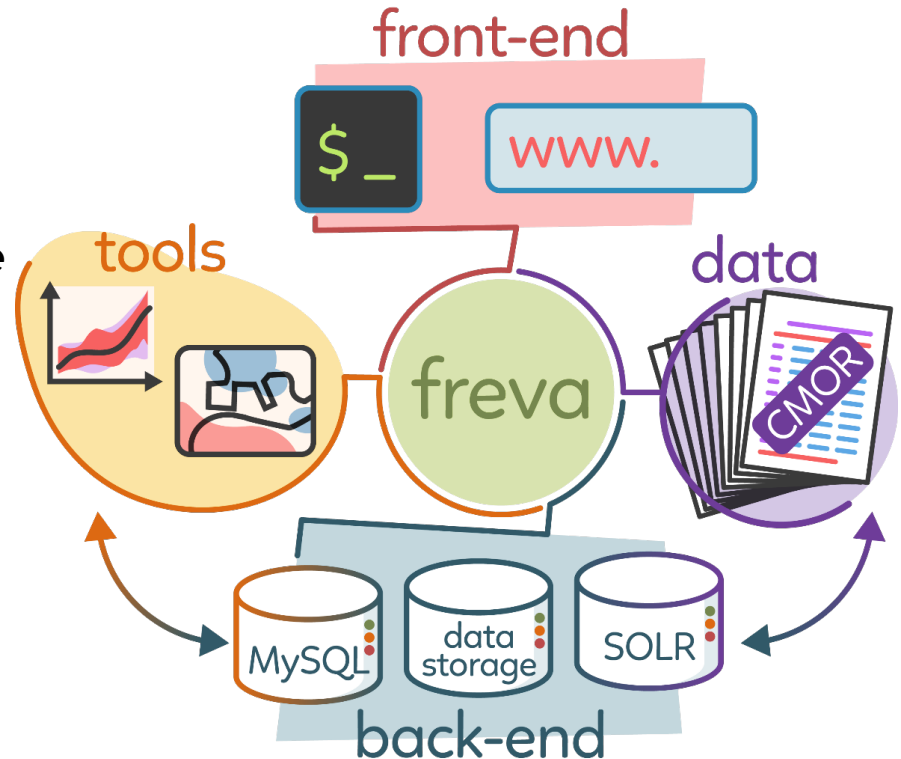
Stock of available and reusable processing scripts/workflows

Interfaces coupled directly to model output stream to ingest data in catalogued database

## freva

Software infrastructure for standardized **data and tool solutions** in Earth system science. Freva runs on **high performance computers, it** to handle **customizable** evaluation systems of research projects, institutes or universities.

See [Kadow et al. \(2021\)](#) or [DKRZ Documentation](#) for more details



## freva

Software infrastructure for standard **data and tool solutions** in Earth system science. Freva runs on **high performance computers** to handle **customizable** evaluation systems of research projects at institutes or universities.

See [Kadow et al. \(2021\)](#) or [DKRZ Documentation](#) for more details.

Martin Bergemann and  
Etor E. Lucio Eceiza

Get your hands on Freva, the Free Evaluation System for Earth System Modeling, and explore the possibilities available for your workflow of the future.

