# Provenance Data in the Context of Earth System Modelling

**Kerstin Fieg, Luis Kornblueh, Deike Kleberg, Pavan Siligam**

**DKRZ & MPI-M**

1. **Difference Metadata – Provenance Data**
2. **Why Provenance Data?**
3. **How to collect Provenance Data**
4. **Requirements => Luis' talk**

# 1. Metadaten vs. Provenance Data

**Metadata**
- … are structured data products
- … describing the data itself
- … becoming increasingly important with growing amount of data

Possible future demand of:
    … extended machine readability of metadata
    … and / or automatic evaluation procedures of the metadata.

*No data archiving without Metadata*

*e.g. NetCDF as self describing model metadata format*

**Provenance Data**
"The provenance of a piece of data is the process that led to that piece of data", Moreau, 2010

- … describing the origin of the results
- … is more than an "extension" to Metadata
- … includes process details (history / used tools / methods / procedures …..)
- … ideally documenting the complete development process

*No automatic self describing tool available*

## 2. Why Provenance Data ?

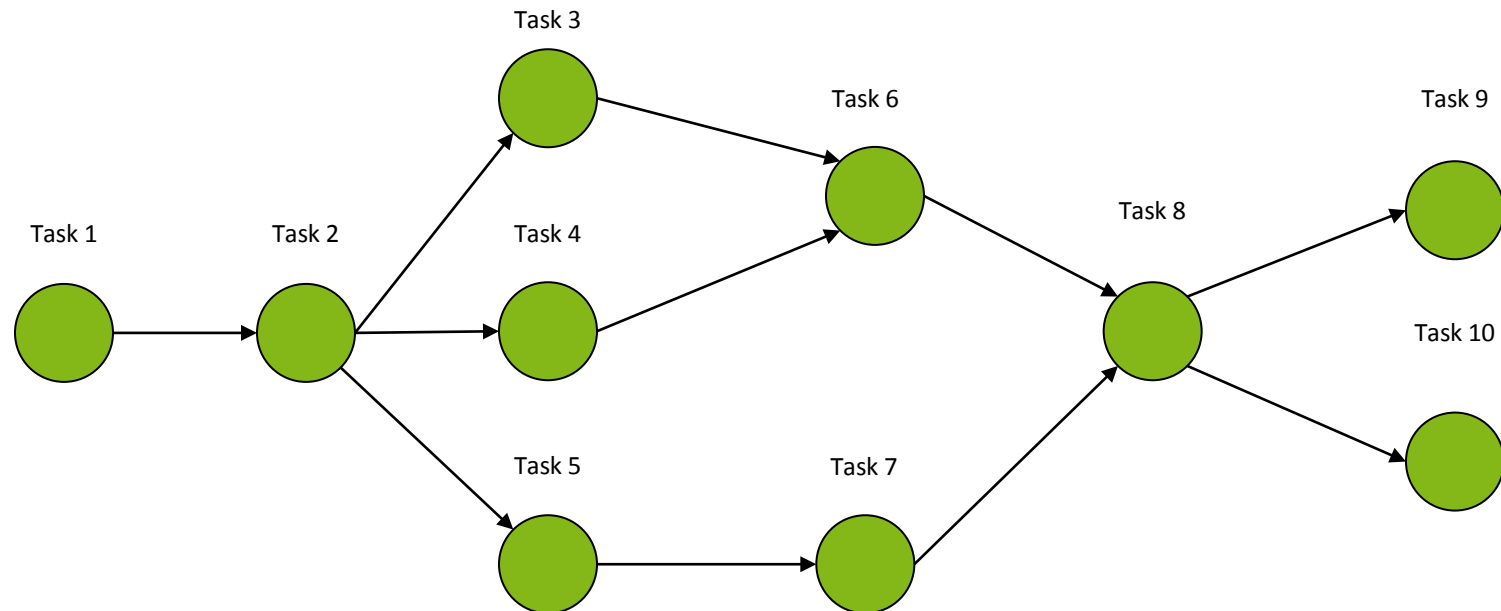➢ we have to - because of "Gute wissenschaftliche Praxis"

### Benefit of Provenance Data – in addition to Metadata

➢ making Provenance data available

… can help to evaluate quality of data

… can document the scientific standard

… enables error – tracking - *even years after data production and leave of the scientist*

➢ degree of reliability is documented and traceable (how / who / why / where /what…)

=> *for possible discussions in future facts are available*!

➢ Provenance Data

• … enable reusebility of data  *for a longer time*

• … enable runtime recovery

✓ repeat / reproduce model runs

We have a defined model **workflow**
  … which describes a chain of chronological and logical connected <u>data</u>
  steps or <u>processing</u> steps - or both in combination



To perform the workflow, we need a model **framework**
  … which is a transparent runtime environment, that <u>drives</u> the
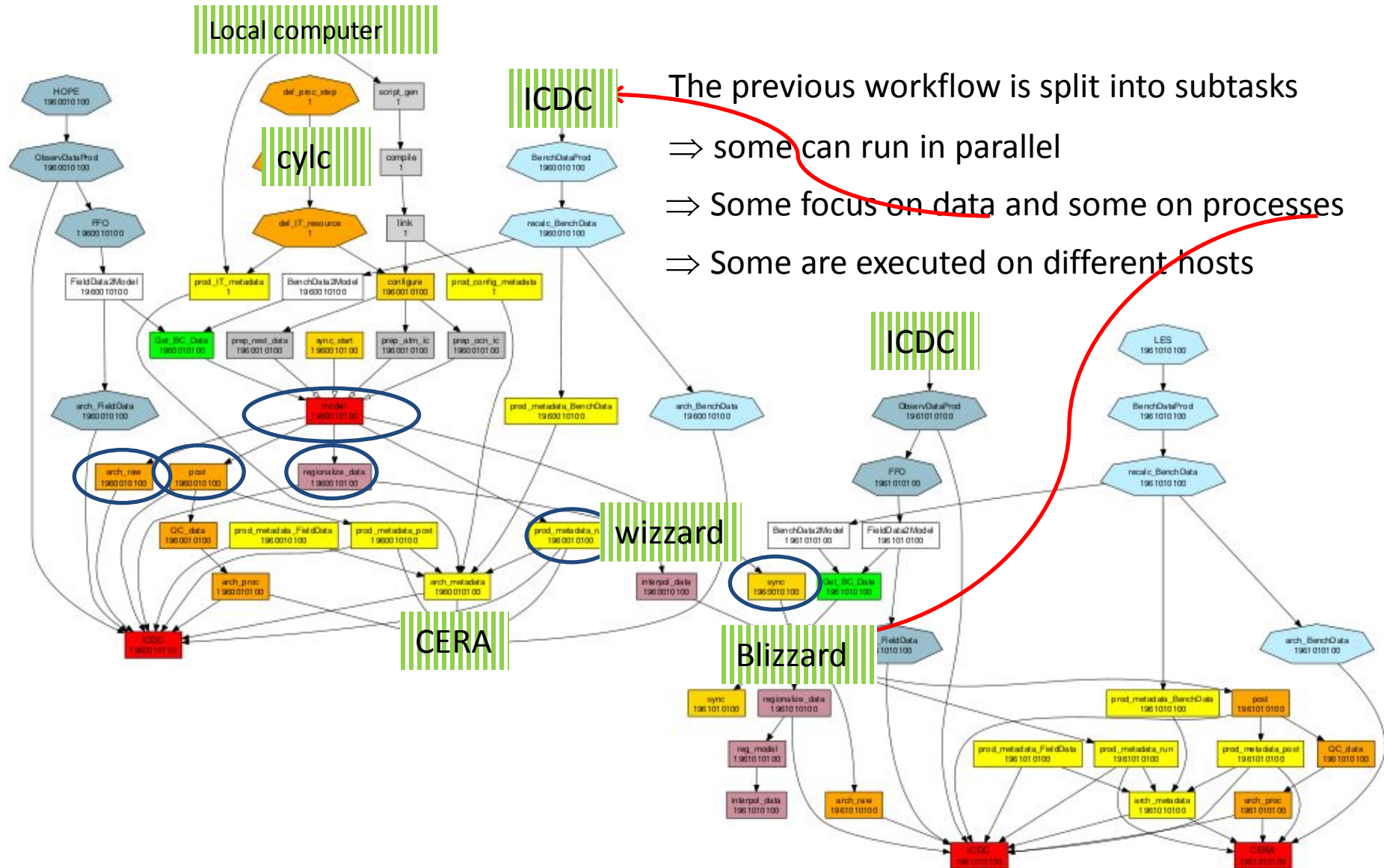  workflow and <u>supports</u> it during the entire lifetime of a model run
  or project

## Complete Workflow in climate research context (= Reinhards graph)

1. Prepare model (dynamic, physics, optimization)
2. Set up components (site, machine, compiler)
3. Define experiment
4. Build model
5. Compute host directory structure and input file setup
6. Restart and time dependent input data preparation
7. Model run step
8. Restart and log handling
9. Archiving of raw data and restart
10. Post-processing
11. In-house quality control
12. Archiving of derived data products
13. Data publishing process (CIM, ESGF data publication)
14. ESGF (data and index nodes)
15. Data search and data access interface
16. Community analysis

# What we try to cover (processing & data)

1. Prepare model (dynamic, physics, optimization)
2. Set up components (site, machine, compiler, libs)
3. Define experiment
4. Build model
5. Compute host directory structure and input file setup
6. Restart and time dependent input data preparation
7. Model run step
8. Restart and log handling
9. Archiving of raw data and restart
10. Post-processing
11. In-house quality control
12. Archiving of derived data products
13. Data publishing process (CIM, ESGF data publication)
14. ESGF (data and index nodes)
15. Data search and data access interface
16. Community analysis

=> No metadata, but a protocol of **all actions** leading to the final data is collected and stored

The previous workflow is split into subtasks

$\Rightarrow$ some can run in parallel

$\Rightarrow$ Some focus on data and some on processes

$\Rightarrow$ Some are executed on different hosts

General technical requirements:
- allow modularity / flexibility / portability
- allow user interaction / experiment configuration
- enable program portability
- allow general reproducibility of a task / result
- bundle / integrate
    - heterogeneous, distributed services / software tools
    - heterogeneous and distributed data management
- collect provenance data and / or metadata

**Interface of data and process workflow**

Our **frameworks** tool should …
- schedule the individual subtasks of the process chain
- be platforms independent
- enable monitoring of processes
- support testing and quality checking (QC)
- ease failure handling
- enable restart  / controlled repetition of an experiment
- deliver  / produce provenance data( = "Gute wissenschaftliche Praxis …")

# What are we using?

⇒ **We selected CYLC as our Frameworks tool**

What is CYLC?
- Cylc (*"silk"*) is a **suite engine** and ~~~~~~~~~~~~~~**duler**, designed to manage suites of cycling tasks in weather and ~~~~~~~~~~~~casting
- Using Python => platfor~~~~~~~~~~~~~ent
- Enables running of ~~~~~~~~~~es and / or processes in parallel
- Developed by ~~~~~~~~~~~~~NIWA, NZ) under GNU Licence

CYLC
⇒ ... will be our new runtime environment for climate
   research (e.g. for projects like MIKLIP & HD(CP)2)
⇒ ... furthermore, will manage the collection of the
   Provenance Data

⇒ Data will be collected in a separate Provenance Data Base

H. Oliver will be in Hamburg for the ISENES2 Workshop on workflows, June 3 – 5 2014

**Concept of a data base schema for provenance data
by Deike Kleberg (MPI-M)**

Target:

Develop a data base scheme to store all kind of Provenance data
- in an abstract and simple way
- flexible and easy expandable for possible future requirements

# Concept of a data base schema for provenance data
## by Deike Kleberg (MPI-M)



MPI-ESM Runtime Environment, MiKlip, database ER diagram, Deike Kleberg, 29.12.2013, Draft v0047