

PUNCH use case: Compute4PUNCH for data-intensive analysis applications

Matthias Hoefft

Manuel Giffels, Andreas Henkel, Dominik Schwarz,
Kilian Schwarz, Christoph Wissing

for PUNCH4NFDI Task Area 2

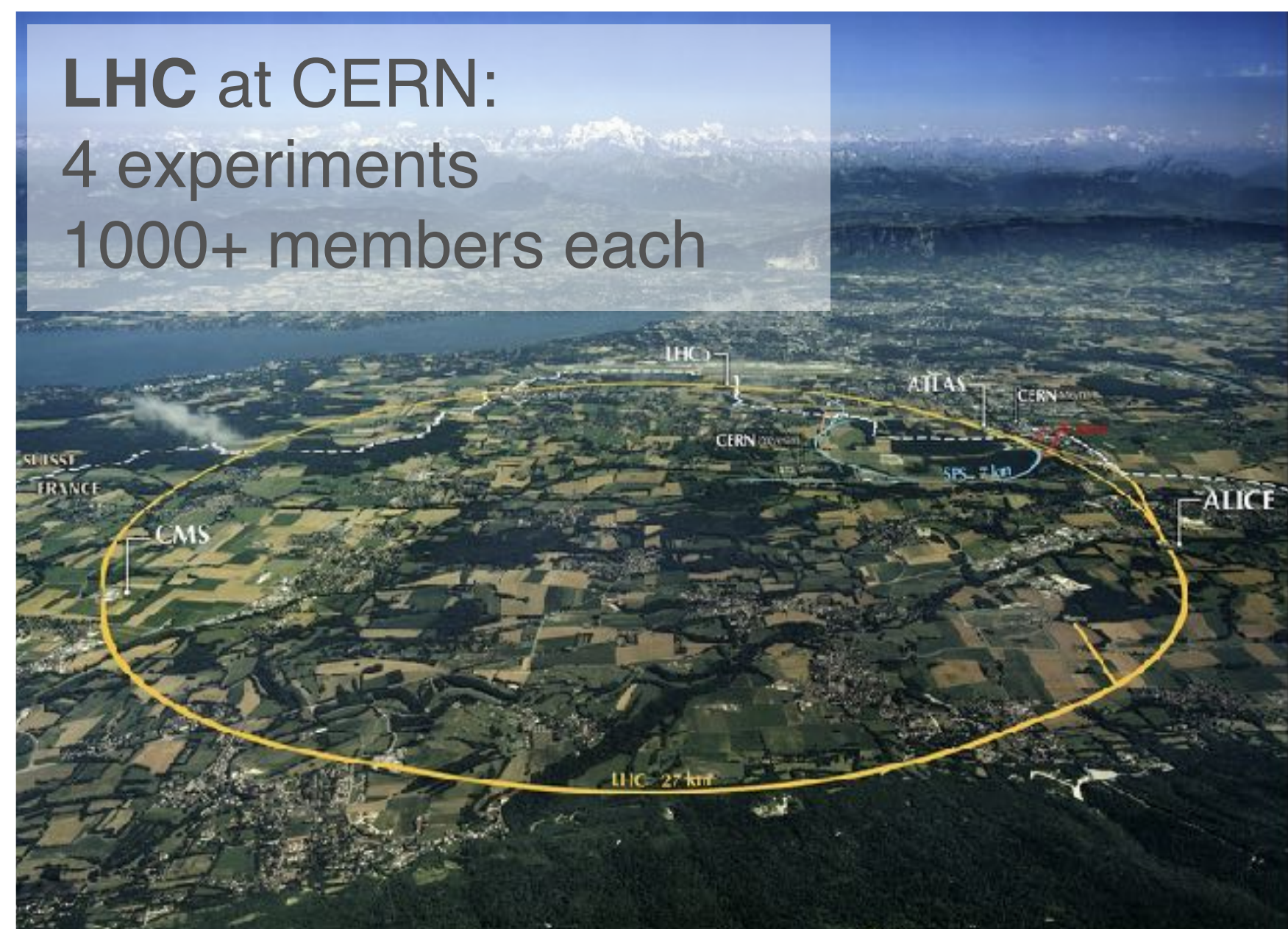


Examples for scientific questions in the PUNCH sciences



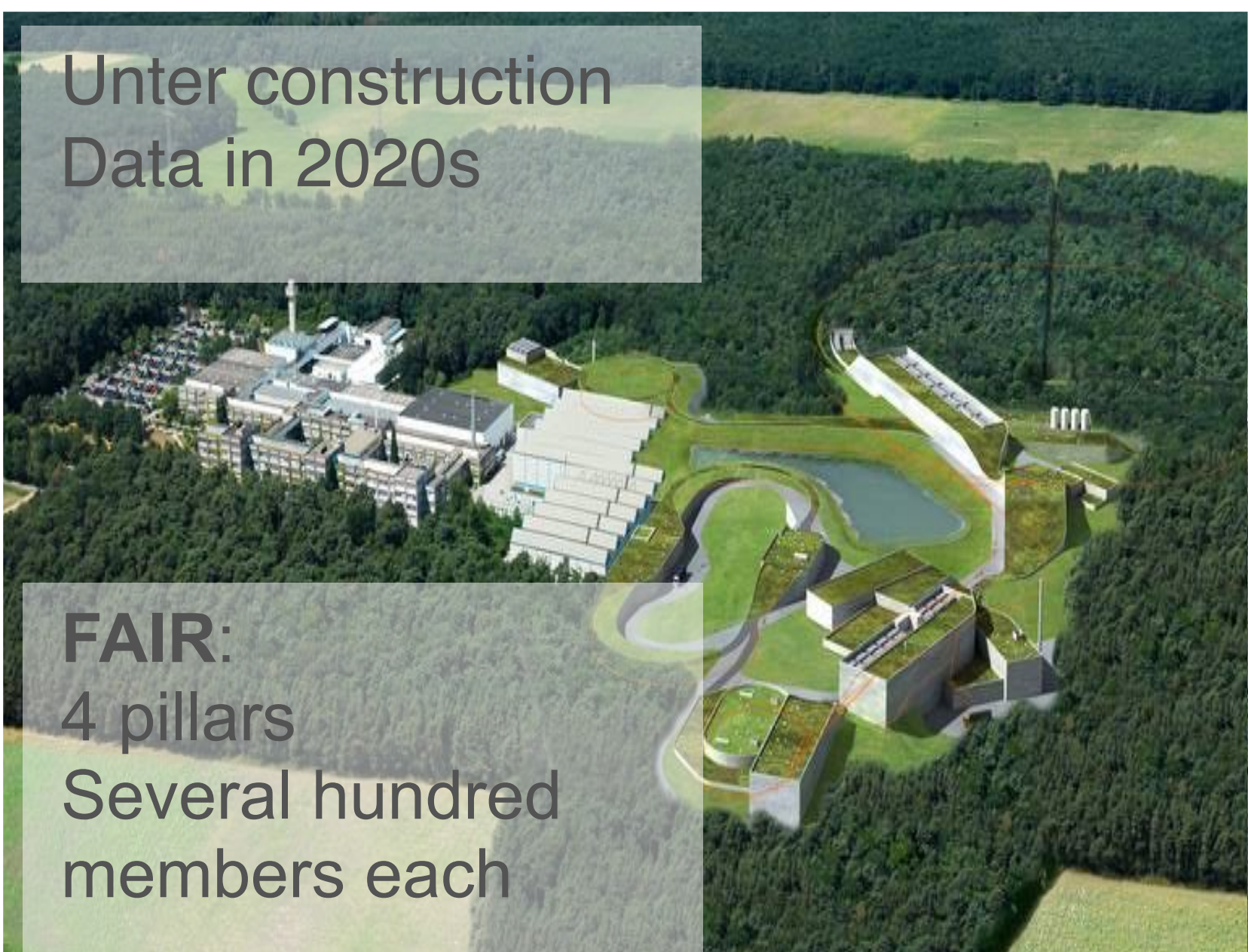
Answering these questions needs large facilities

LHC at CERN:
4 experiments
1000+ members each

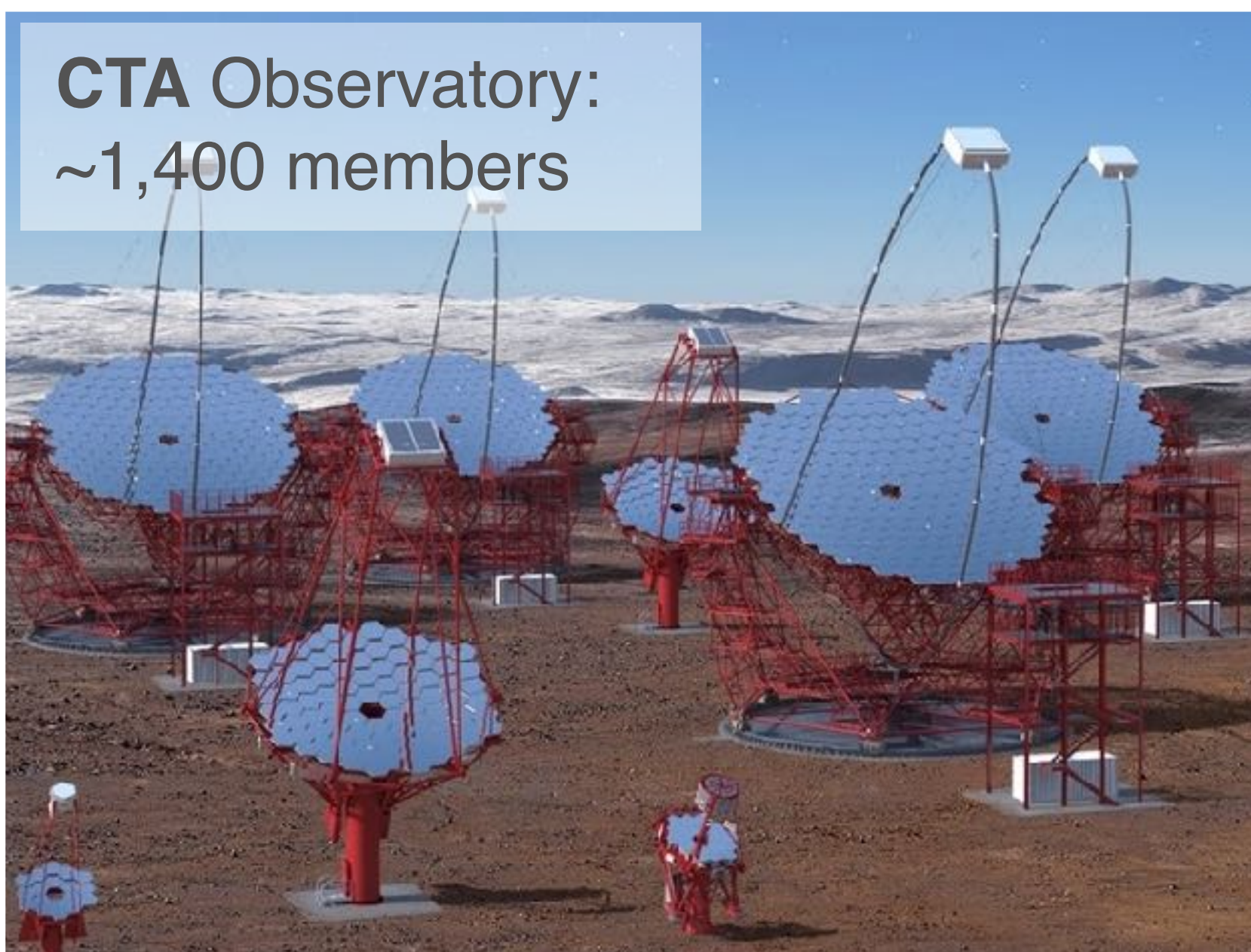


Unter construction
Data in 2020s

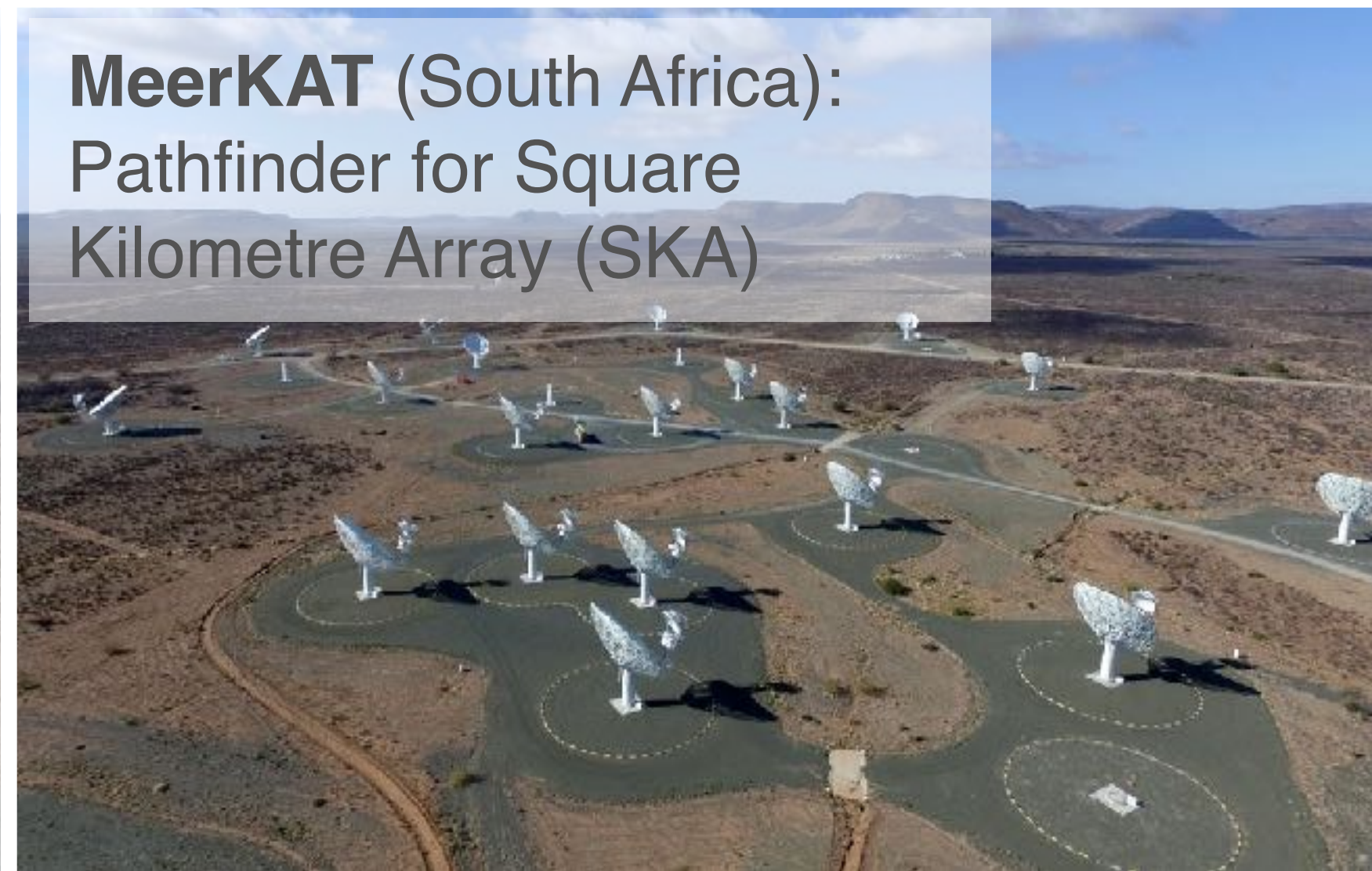
FAIR:
4 pillars
Several hundred
members each



CTA Observatory:
~1,400 members



MeerKAT (South Africa):
Pathfinder for Square
Kilometre Array (SKA)

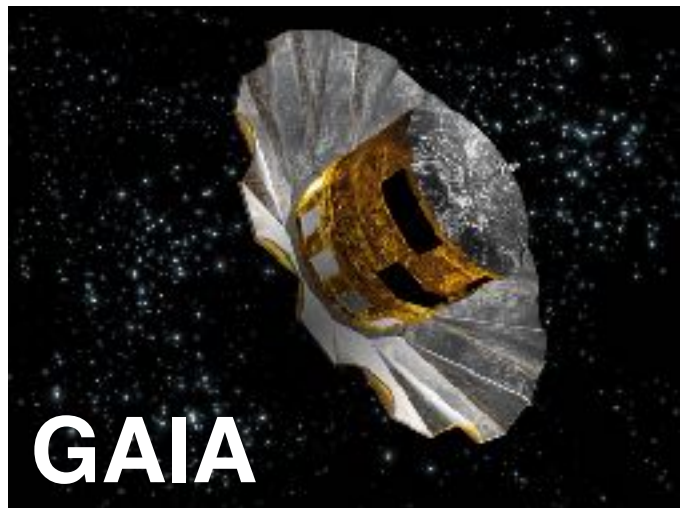


Supercomputing:
Theoretical physics
Cosmology



And many more
examples

GAIA



S-DALINAC



(Some) Challenges

PUNCH communities deal with **increasingly large datasets**
for instance:

High-Luminosity LHC and the Square Kilometre Array

Advanced studies often require
combination of very different observations, experiments and simulations
for instance:

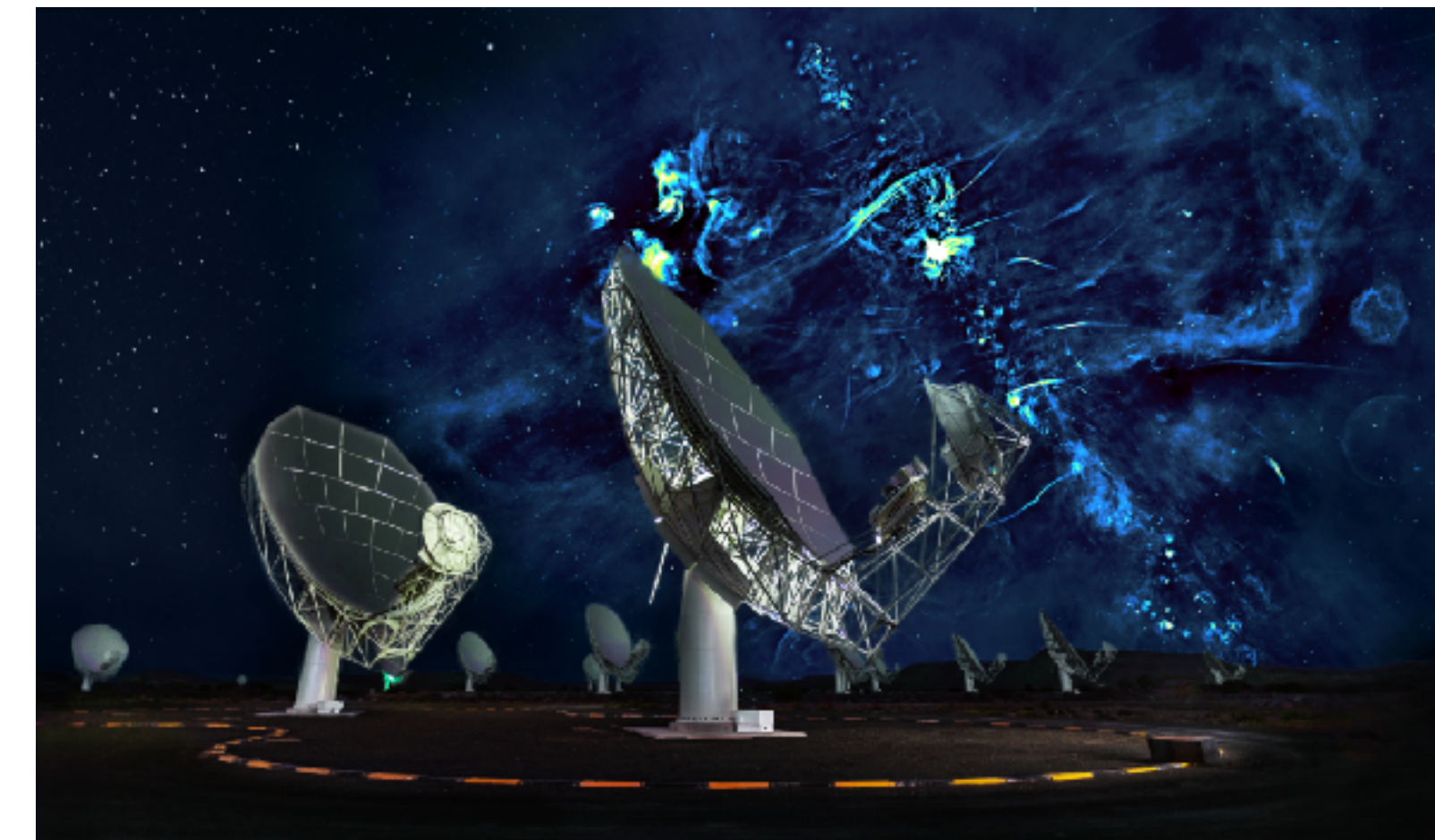
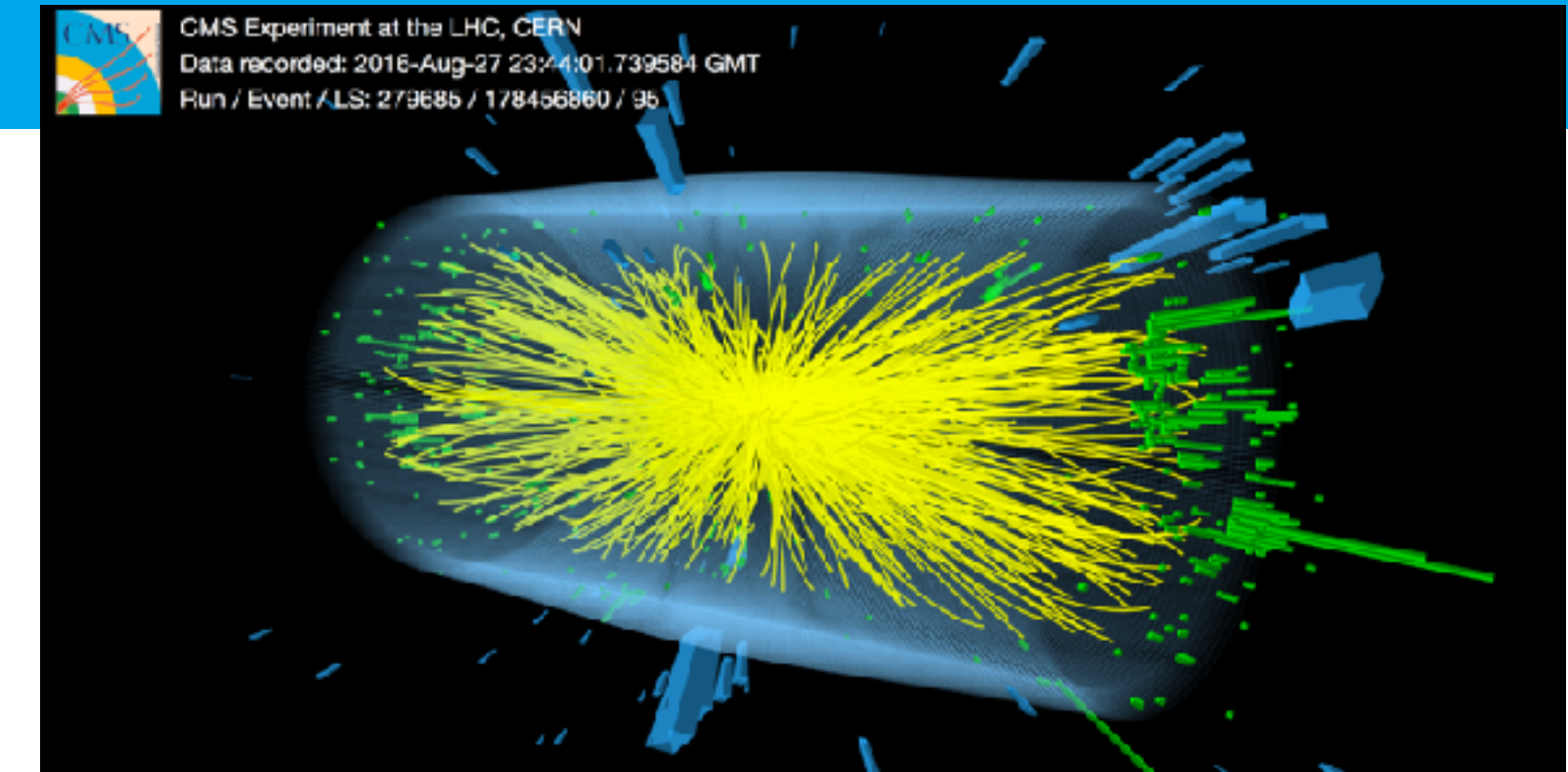
constraining dark matter properties from
astronomical observations and particle experiments

Peak usage of resources

for instance:
reprocessing of archival data with improved analysis scheme

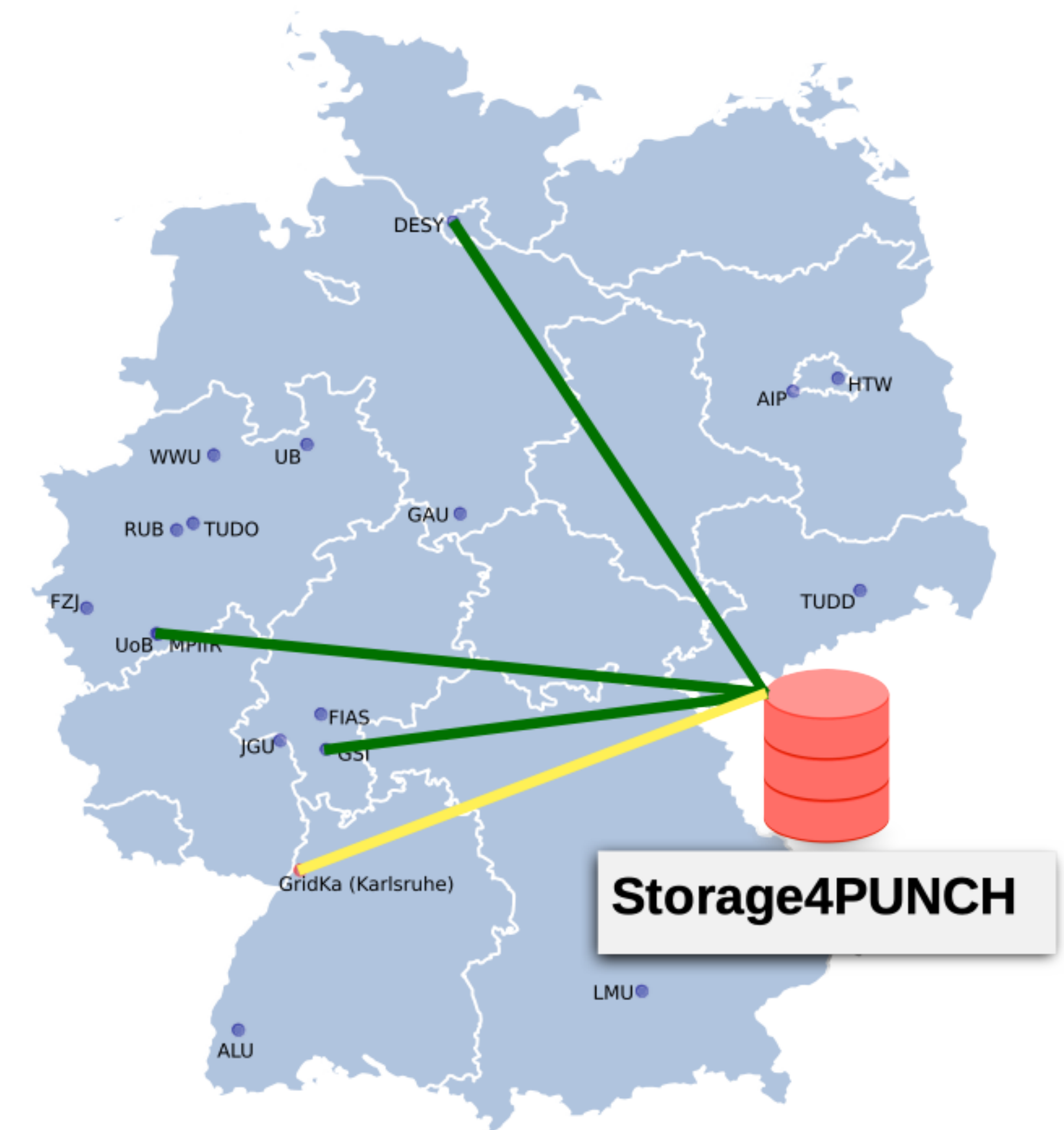
Our Goals

- Federate existing storage and compute infrastructure
- **Provide 'easy-to-use' entry points to storage and computing resources**
- Provide an easy access to a large variety of data collections

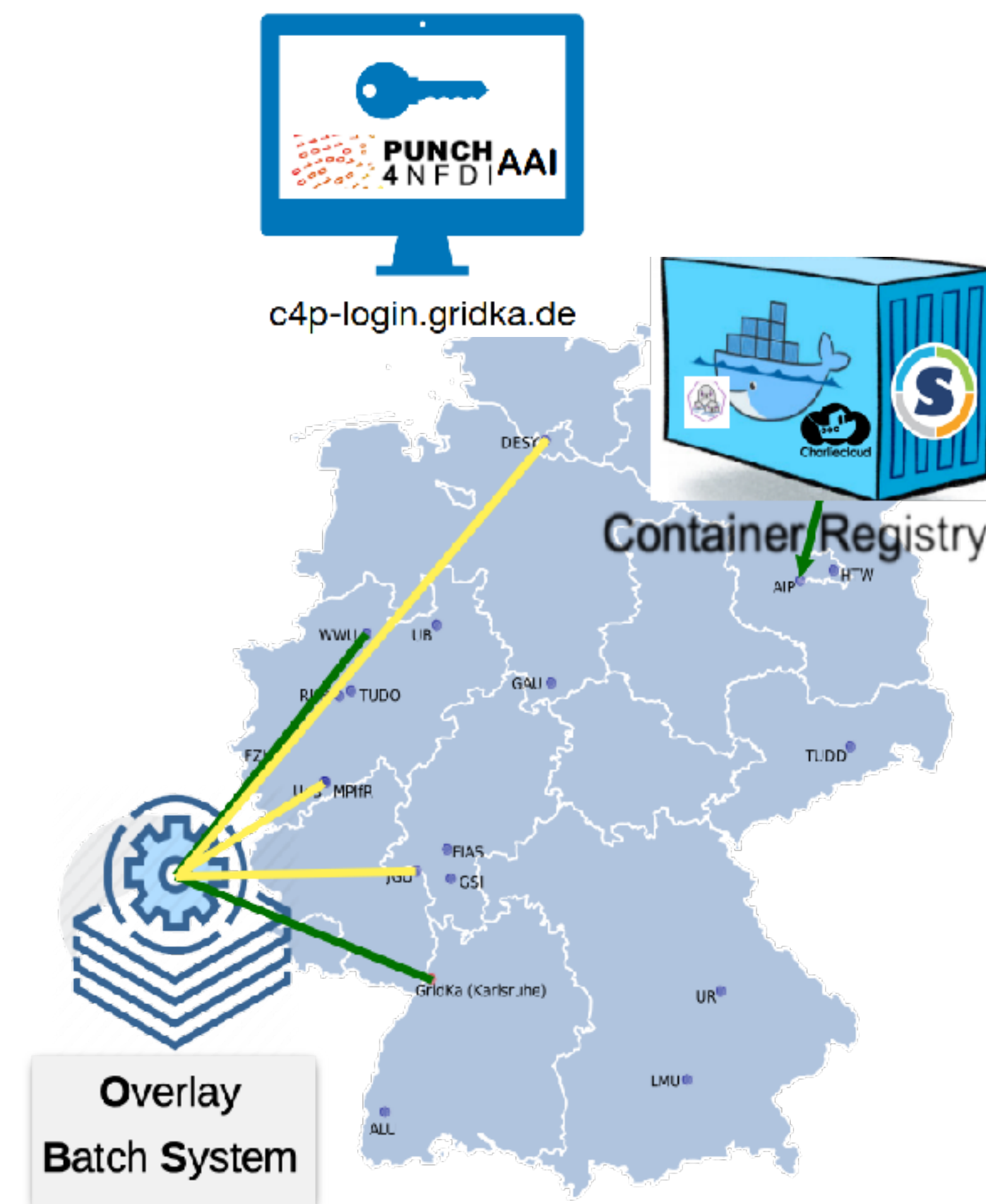


Demonstrator Workflows

- Based on two storage technologies
 - dCache (Instance at DESY & KIT)
 - XrootD (Test instances at U Bonn & GSI)
- Token based access using PUNCH AAI
- Supported protocols: WebDAV & XrootD
- Further integration & evaluation:
 - Systems for file & replica catalogs
 - Rucio – Common DM tool in HEP

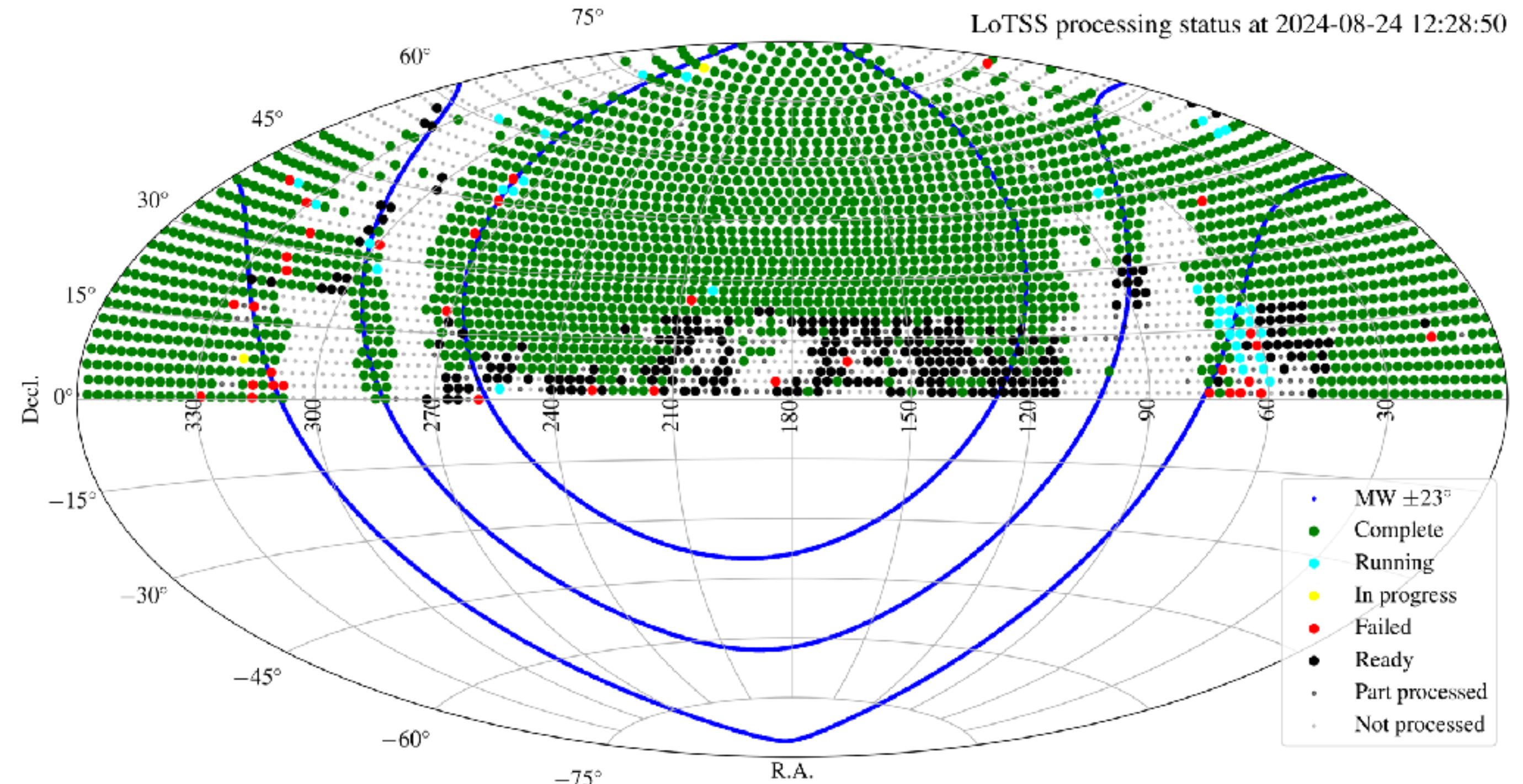


- Prototype of federated Compute4PUNCH infrastructure is set up
- AAI based login node of all PUNCH4NFDI members
- Container registry & Gitlab CI/CD workflows
- Dynamic integration of several sites
U Bonn, KIT, LMU, DESY, FZJ, U Mainz, U Göttingen
- Container distribution via CVMFS
(CERN Virtual Machine File System)

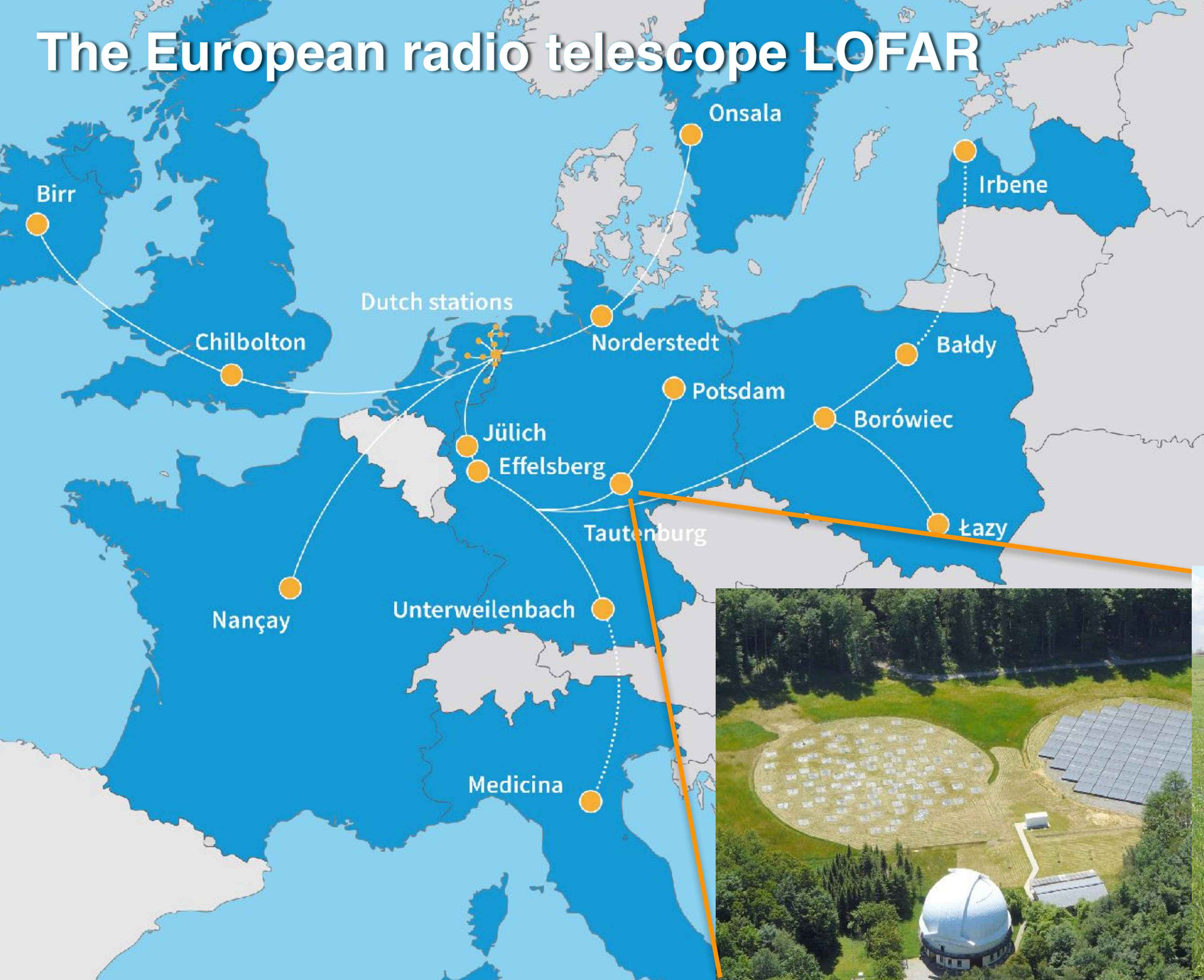


Workflow 1: LoTSS image processing

- LoTSS: Multi-year observing campaign to image the northern sky at very low radio frequencies
- About 1500 ‘observations’ (8hrs pointing)
- Data taken with the Low Frequency Array (LOFAR) imaging requires substantial computing (several mio core-h / year)
- Our aim: Demonstrate that Compute4PUNCH may serve as ‘science data processor’
- Pipeline available as container



The European radio telescope LOFAR



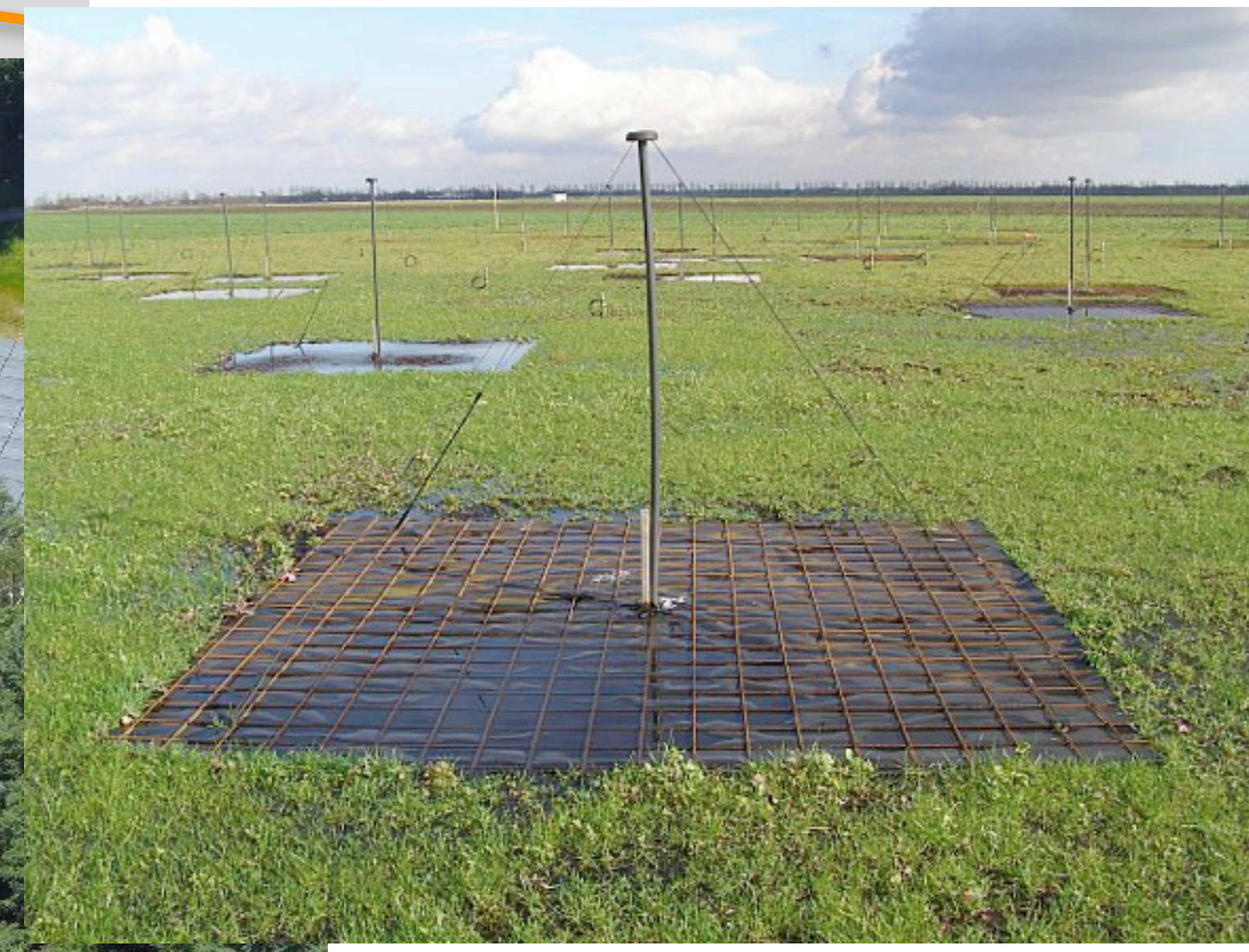
~ 50 stations in Europe

Correlator in Groningen (NL)
‘Online’ processing

connected via fast data links
data rate at station: 4 Gb/s

100,000 very simple antennas

currently *upgrade to LOFAR2.0*

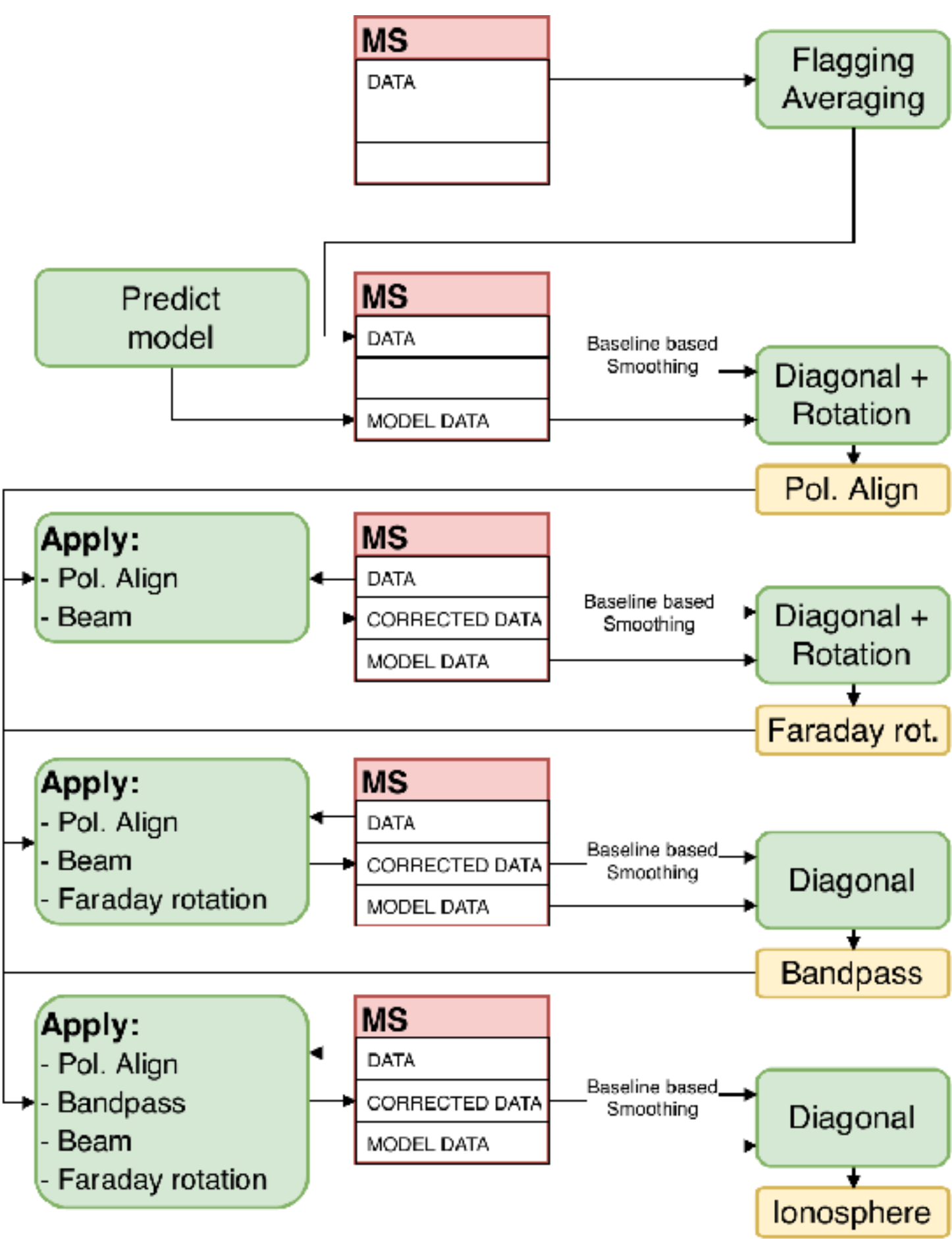
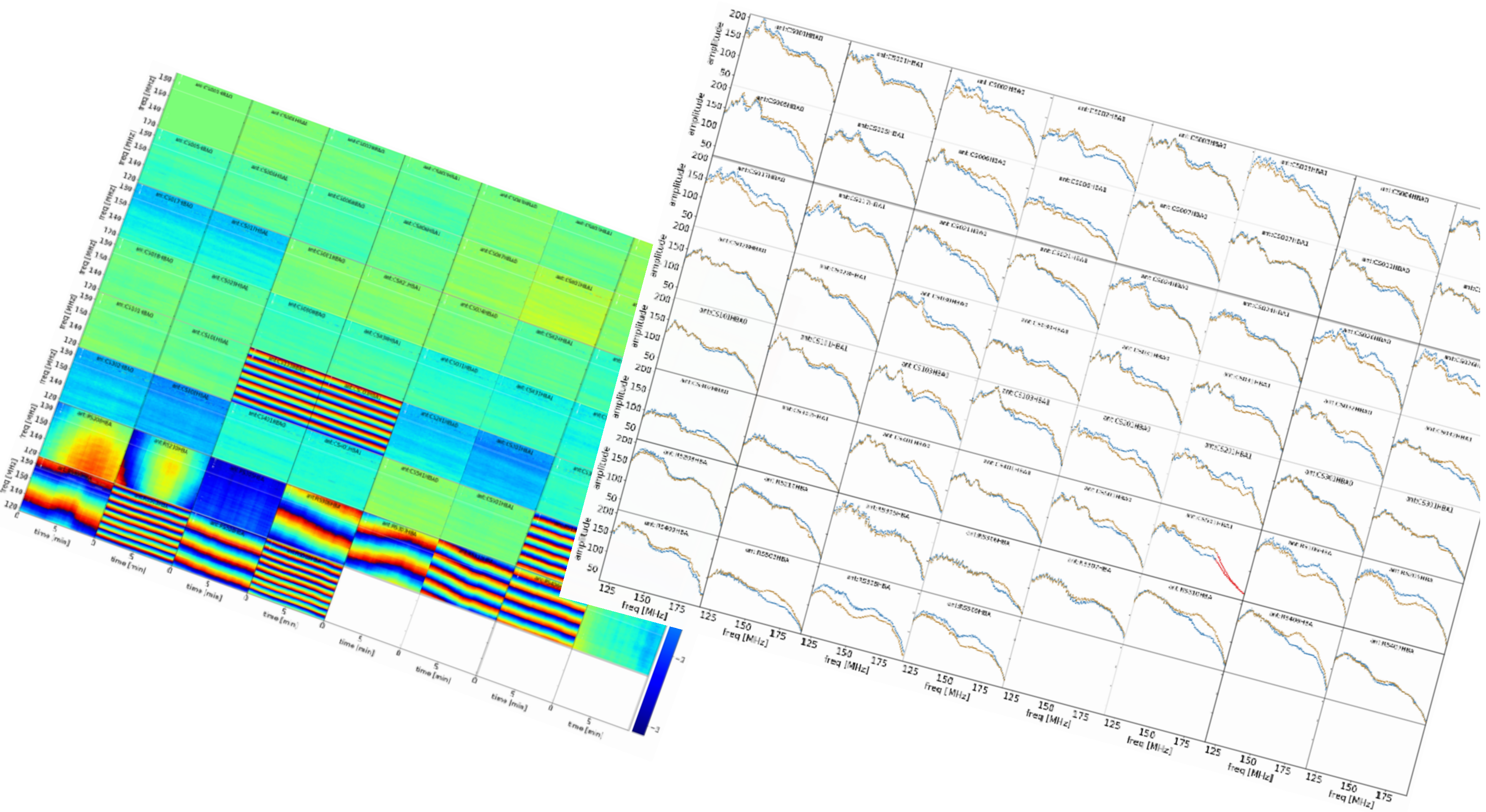


Special requirements for this workflow

- Data are stored on tapes at FZ Jülich (long-term archive, also other sites)
- Data volume of one observation $\sim 4 - 16$ TB
- processing of one observations typically once or a few times
- current software design requires powerful (single) nodes
and sufficient scratch space ($\sim 5 \dots 20$ TB)
processing a few 10000 core-h

Workflow 1: Radio interferometry analysis workflow

- ‘Historically’ grown software suite
- CWL pipeline



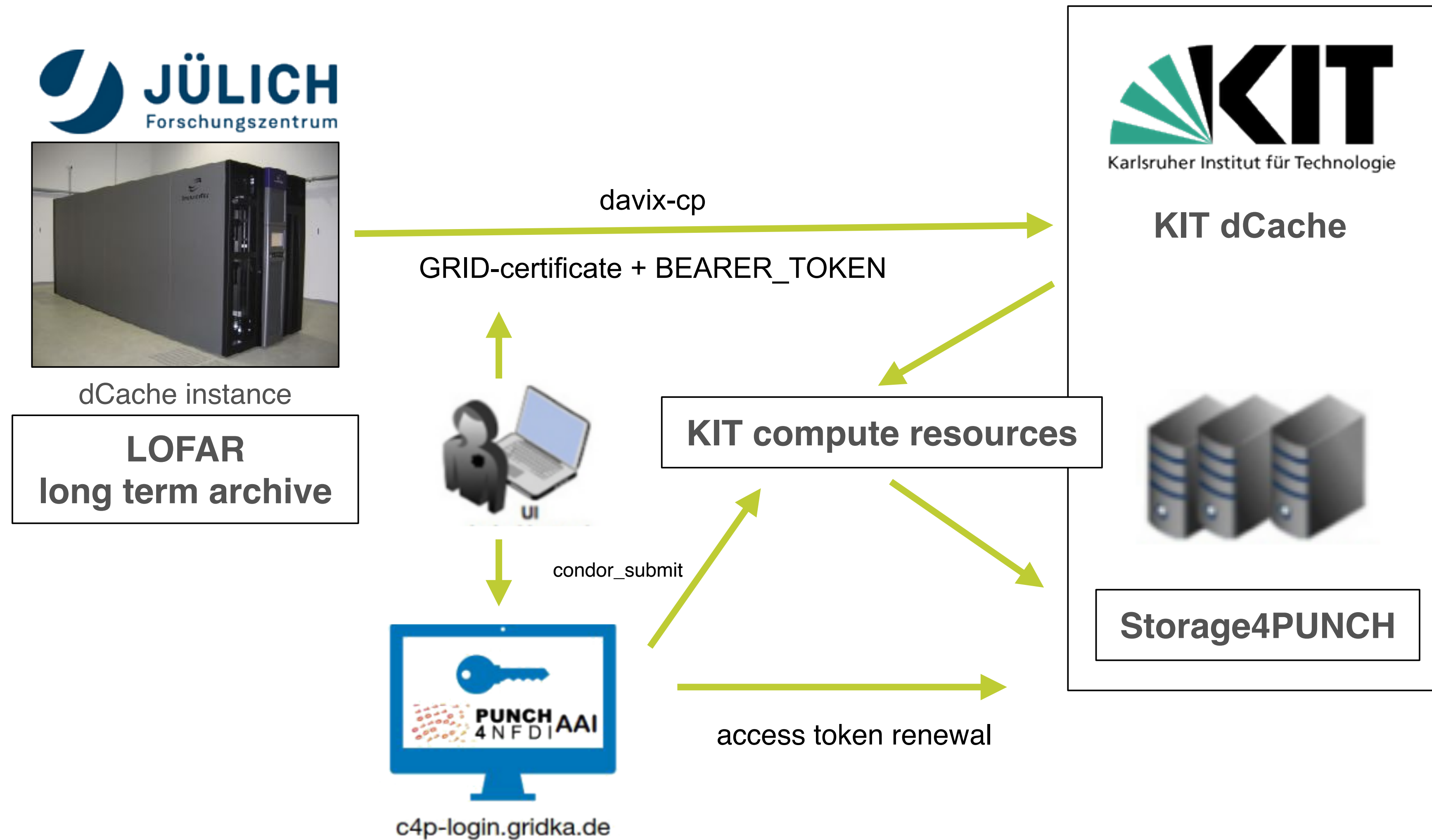
Slide: A. Drabent

<https://git.astron.nl/RD/LINC>



Workflow 1: LoTSS radio interferometry processing

Data transfer challenge :



Slide: A. Drabent

Workflow 1: LoTSS radio interferometry processing

Documentation :

C

C4P_LOFAR_Processing_Documentation

Star

0

Fork

0

main

c4p_lofar_processing_documentation

+

History

Find file

Edit

Code

update documentation

Alexander Drabent authored 1 minute ago

d50bb017

Name	Last commit	Last update
LoTSS_processing	Add script for untarring data	4 minutes ago
README.md	update documentation	1 minute ago

README.md

Compute4PUNCH Documentation -- Usecase radioastronomical data reduction for the LOFAR Two-Metre Sky Survey (LoTSS)

Login to Compute4PUNCH Login Node

Set-up oidc-agent

Project information

2 Commits

1 Branch

0 Tags

3 KiB Project Storage

README

Auto DevOps enabled

Add LICENSE

Add CHANGELOG

Add CONTRIBUTING

Add Kubernetes cluster

Configure Integrations

Created on

September 03, 2024

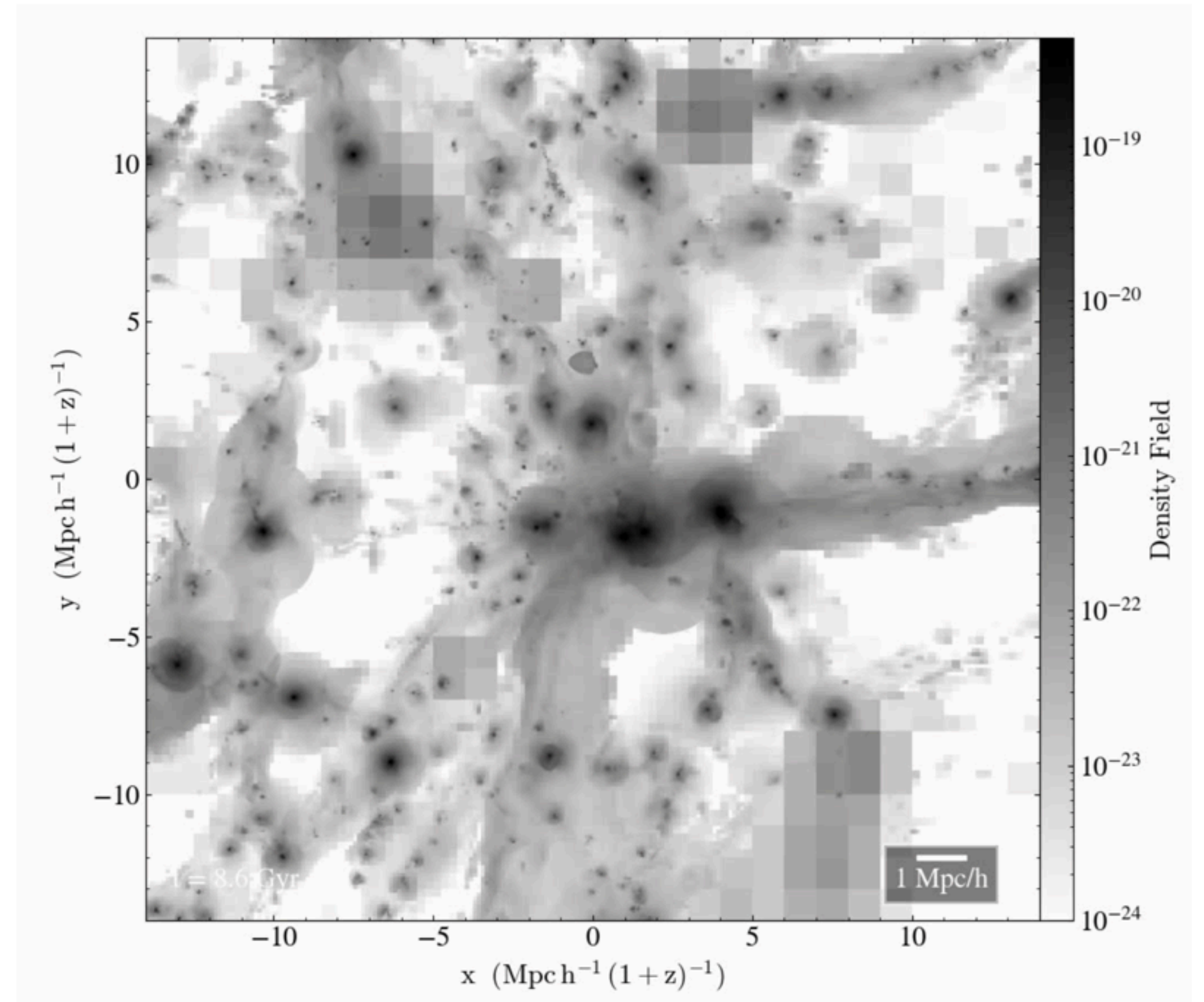
https://gitlab-p4n.aip.de/compute4punch/tutorials/c4p_lofar_processing_documentation



Workflow 2: Postprocessing of cosmological simulations

- Large outputs of numerically expensive simulations
Postprocessing: requires large analysis facilities
- Data: in Storage4PUNCH or external
- Analysis software in container
- Entry point: reana
can access storage with token
(‘secret’ in reana)

reana



Slide: P. Gupta

Workflow 2: Postprocessing of cosmological simulations

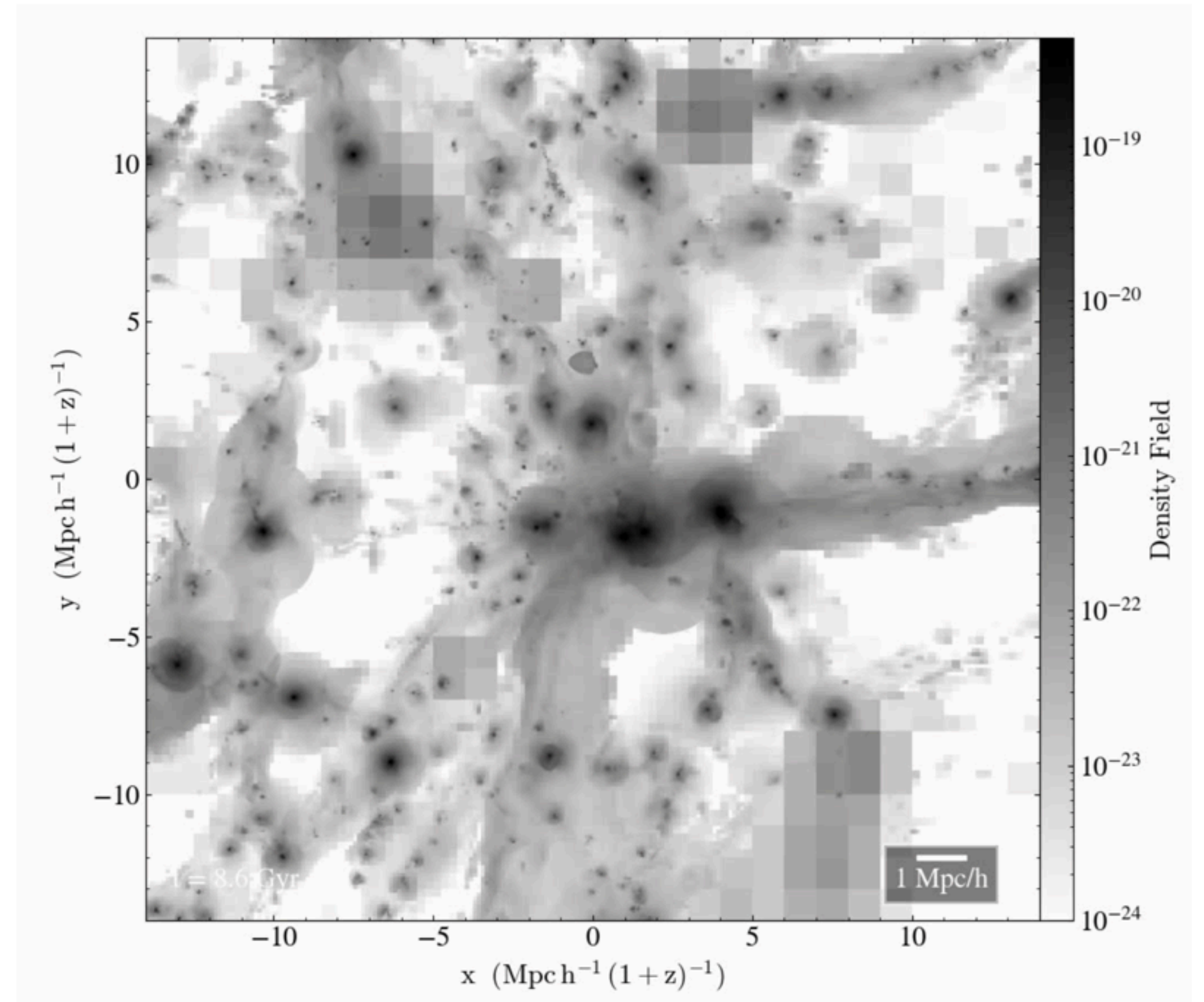
reana.yaml

```
inputs:
  files:
    - multiple_plot_test.py
    - download_data.sh
    - upload_data.sh
  parameters:
    maincode: multiple_plot_test.py
    getdata: download_data.sh
    uploadResult: upload_data.sh

workflow:
  type: serial
  specification:
    steps:
      - name: Make Directories and Download Simulat
        environment: 'docker.io/tlsprateek/davixdoc
        commands:
          - mkdir Simulated_Data
          - mkdir Output
          - bash "${getdata}"

      - name: Prepare Plots
        environment: 'docker.io/tlsprateek/ytprojec
        commands:
          - python3 "${maincode}"

      - name: Upload the Output
        environment: 'docker.io/tlsprateek/davixdoc
        commands:
```



Example 3: Jupyter Environment for ML-PPA

Machine Learning-based Pipeline for Pulsar Analysis



Login to Compute4PUNCH

Access via `c4p-login.gridka.de` and forward a port (e.g., `9998:localhost:9998`).



Submit Job Script

Use `condor_submit` for 1-hour runtime, minimal resources.



Install & Start Jupyter

Download Miniconda, install Jupyter, and start the server on forwarded port.



Secure SSH Login

Authenticate via `oidc-ssh` to access Compute4PUNCH.



Submit and Access Job

Allocate resources with `condor_submit` and SSH into the job.



Port Forwarding & Notebook Access

Forward port 9998 and connect via local browser to Jupyter Notebook.

https://gitlab-p4n.aip.de/punch/intra-docs-content/-/blob/master/files/TA5/Compute4PUNCH_Access_Usability_JupyterNotebooks.pdf

Slide: A. Redelbach
G. Ravindra Dange

Conclusions

- Storage4PUNCH & Compute4PUNCH prototypes are available
- Prototypes needs to be tested and commissioned by community specific workflows (Demonstrator workflows)
- Several use case examples are implemented, for instance
 - LoTSS observation processing → data transfer
 - Postprocessing of cosmological simulation output → reana
 - ML-PPA : → enable jupyter notebooks
- Lessons learned:
 - Requires extra effort to implement workflows
 - allows to use federated infrastructure

