# Optimizing Scientific Workflows

## Integrating REANA with PUNCH Infrastructure

Dr Arman Khalatyan /Analysis Facilities Workshop / 19. 06 2024

# Research interests



Projects

HPC
BigDATA
ML
MPI
OpenMPI
Visualisation

HESTIA:
High-resolution
Environmental
Simulations of The
Immediate Area
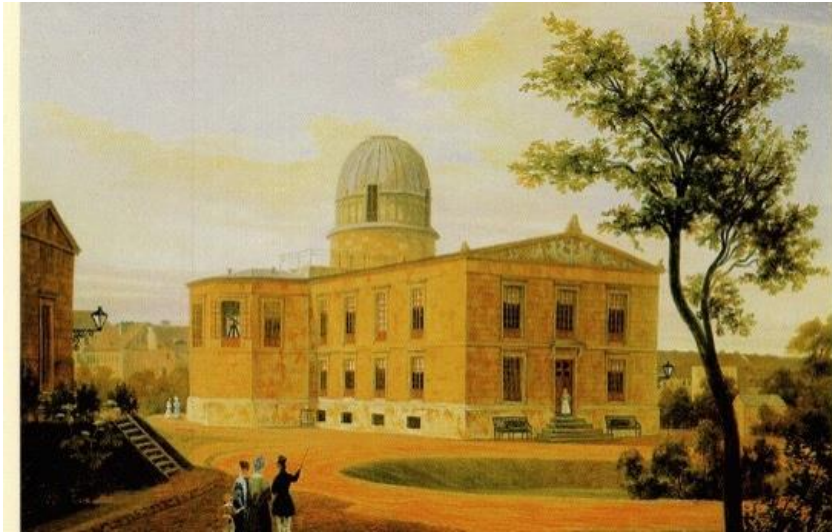
StarHorse:
Photo-astrometric
distances,
extinctions, and
astrophysical
parameters for Gaia
stars brighter than G
= 18

>3000 Cores
GPUs-T4,A100,RTX8000

colab.aip.de cloud.aip.de vr.aip.de

Nationale Forschungs-
Daten Infrastruktur NFDI

PUNCH 4NFDI

# From Berlin to Babelsberg



The Berliner Sternwarte in Berlin-Dorotheenstadt, today Berlin-Mitte.



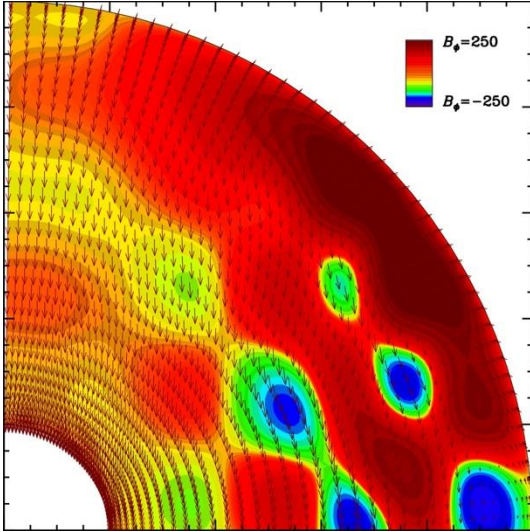The new Sternwartengebäude in Babelsberg, built 1913.

The Berliner Sternwarte, founded in 1700, moved to 1913 Babelsberg because the growing city made scientific observations difficult. Light pollution and vibrations from traffic being the main reasons

# Changing Names



From the Zentralinstitut für Astrophysik (1969) after the wall came down the Astrophysikalische Institut Potsdam (1992) was founded. In 2011 ithe AIP was renamed to **Leibniz-Institut für Astrophysik Potsdam (AIP),** to emphasize the membership with the Leibniz-Gemeinschaft.
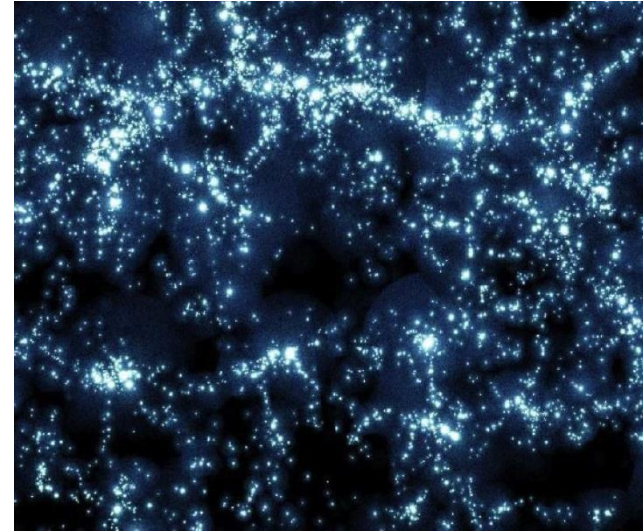
# Research Areas



## Cosmic Magnetic Fields

Research on solar, stellar and galactic magnetic fields and magnetohydrodynamic (MHD) mechanisms.
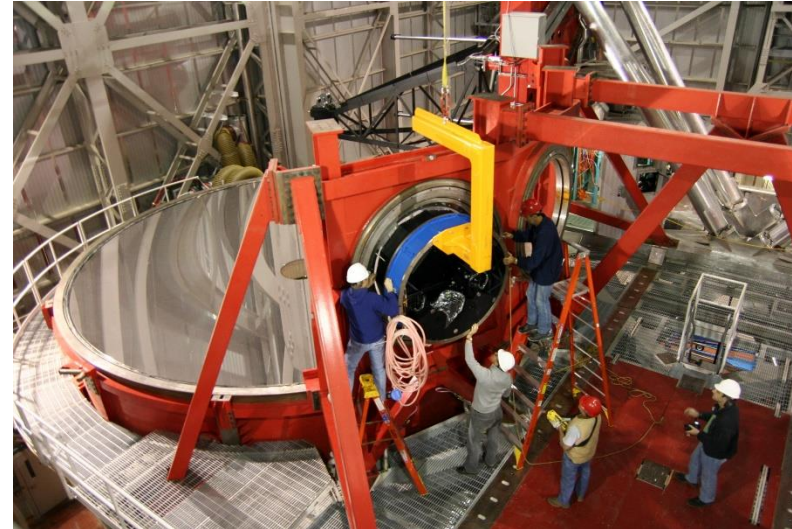**https://www.aip.de**



## Extragalactic Astrophysics

Active galaxies and quasars. Galactic archaeology and extragalactic research based on high resolution simulations.
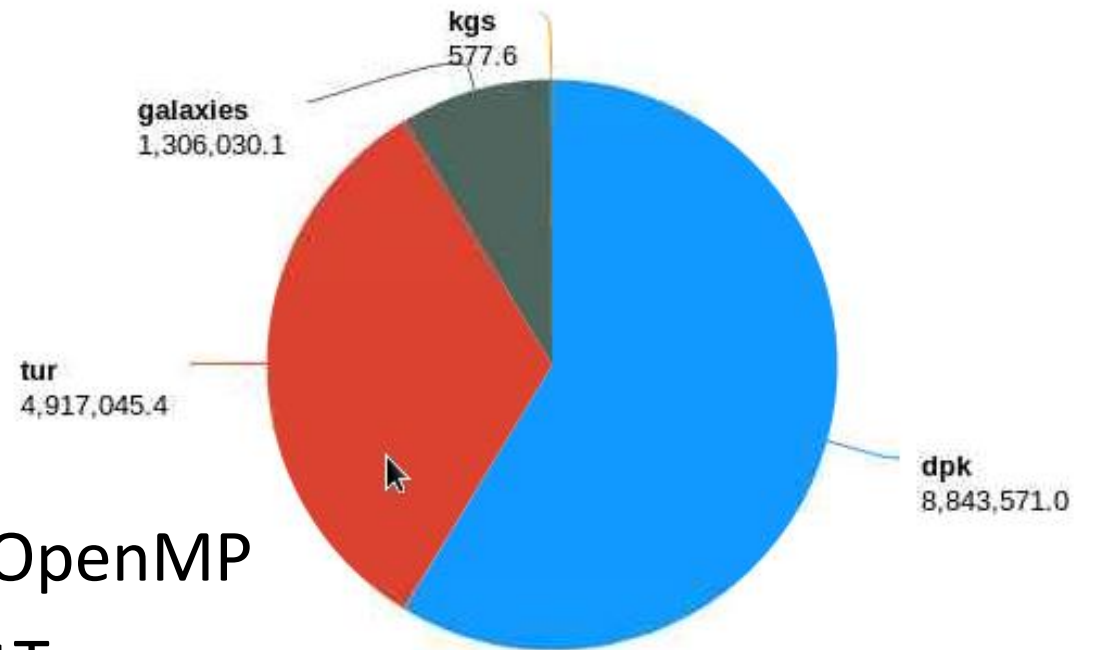
# Research technology and infrastructure



„Development of Research technology and infrastructure" ensures the scientific endeavors of AIP and also its participation in international astronomical projects. AIP has construction workshops and labs for instrumentation, especially with fiber optics, ex: AIP is lead institute of the **4MOST** project a fibre-fed spectroscopic survey facility on the VISTA telescope, op-2021)

# Who is using most of the CPU time?

- Cosmology:
  - MHD+Gravity+Gasdynamics
  - Starformation, Cosmic Rays,BH…
- Magneto-hydrodynamics: MHD
- Data processing from telescopes




- Adaptive unstructured mesh, MPI-OpenMP

- Magneto-hydrodynamics: AMR-OctTree

- Data processing from telescopes: python, c, java, other

# Languages

AIP/Clusters
- Python
- C/C++
- Fortran
- IDL
- Java
- Perl
- R
- ?Cuda?

Energy Efficiency across Programming Languages How Do Energy, Time, and Memory Relate?
Rui Pereira et al 2017
https://doi.org/10.1145/3136014.3136031

**Table 4.** Normalized global results for Energy, Time, and Memory

| Total | | | | | | | |
|---|---|---|---|---|---|---|---|
| | **Energy** | | | **Time** | | | **Mb** |
| (c) C | 1.00 | | (c) C | 1.00 | | (c) Pascal | 1.00 |
| (c) Rust | 1.03 | | (c) Rust | 1.04 | | (c) Go | 1.05 |
| (c) C++ | 1.34 | | (c) C++ | 1.56 | | (c) C | 1.17 |
| (c) Ada | 1.70 | | (c) Ada | 1.85 | | (c) Fortran | 1.24 |
| (v) Java | 1.98 | | (v) Java | 1.89 | | (c) C++ | 1.34 |
| (c) Pascal | 2.14 | | (c) Chapel | 2.14 | | (c) Ada | 1.47 |
| (c) Chapel | 2.18 | | (c) Go | 2.83 | | (c) Rust | 1.54 |
| (v) Lisp | 2.27 | | (c) Pascal | 3.02 | | (v) Lisp | 1.92 |
| (c) Ocaml | 2.40 | | (c) Ocaml | 3.09 | | (c) Haskell | 2.45 |
| (c) Fortran | 2.52 | | (v) C# | 3.14 | | (i) PHP | 2.57 |
| (c) Swift | 2.79 | | (v) Lisp | 3.40 | | (c) Swift | 2.71 |
| (c) Haskell | 3.10 | | (c) Haskell | 3.55 | | (i) Python | 2.80 |
| (v) C# | 3.14 | | (c) Swift | 4.20 | | (c) Ocaml | 2.82 |
| (c) Go | 3.23 | | (c) Fortran | 4.20 | | (v) C# | 2.85 |
| (i) Dart | 3.83 | | (v) F# | 6.30 | | (i) Hack | 3.34 |
| (v) F# | 4.13 | | (i) JavaScript | 6.52 | | (v) Racket | 3.52 |
| (i) JavaScript | 4.45 | | (i) Dart | 6.67 | | (i) Ruby | 3.97 |
| (v) Racket | 7.91 | | (v) Racket | 11.27 | | (c) Chapel | 4.00 |
| (i) TypeScript | 21.50 | | (i) Hack | 26.99 | | (v) F# | 4.25 |
| (i) Hack | 24.02 | | (i) PHP | 27.64 | | (i) JavaScript | 4.59 |
| (i) PHP | 29.30 | | (v) Erlang | 36.71 | | (i) TypeScript | 4.69 |
| (v) Erlang | 42.23 | | (i) Jruby | 43.44 | | (v) Java | 6.01 |
| (i) Lua | 45.98 | | (i) TypeScript | 46.20 | | (i) Perl | 6.62 |
| (i) Jruby | 46.54 | | (i) Ruby | 59.34 | | (i) Lua | 6.72 |
| (i) Ruby | 69.91 | | (i) Perl | 65.79 | | (v) Erlang | 7.20 |
| (i) Python | 75.88 | | (i) Python | 71.90 | | (i) Dart | 8.64 |
| (i) Perl | 79.58 | | (i) Lua | 82.91 | | (i) Jruby | 19.84 |

# Data Scales

- Cosmological simulations: >5PB(active)

- Observations(preserve/provide service):
  - GAIA: until now about 500TB, serving 200TB (DB) , soon >1PB (+2025)
  - Applaus: 200TB (photo plates archive)
  - 4MOST: GAIA+?
  - Pepsi: 250TB+
  - Stella- robotic telescope: 150TB+
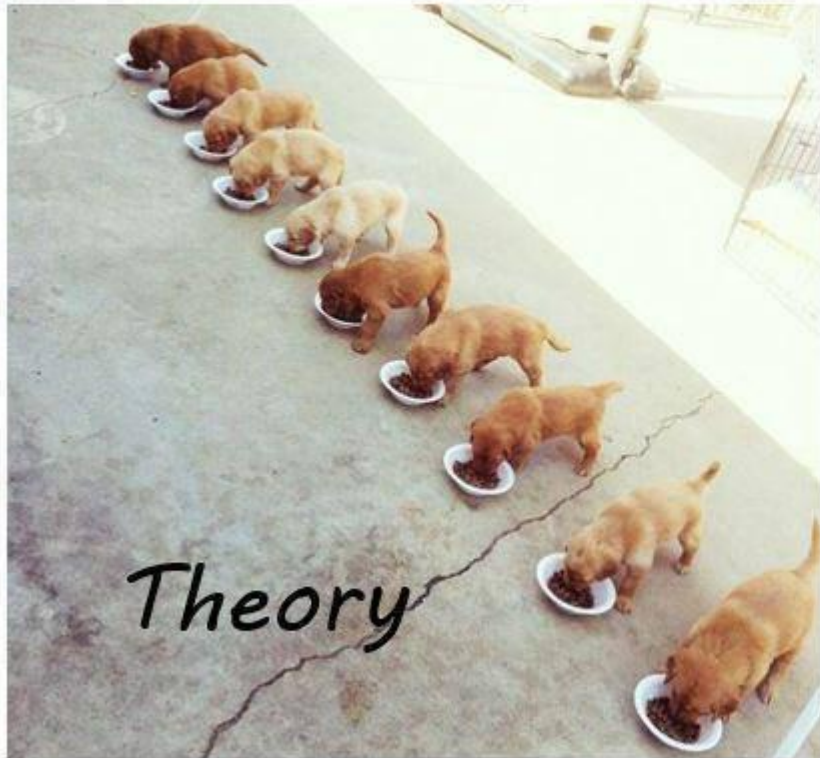- Data sizes:
  - 1k up to +32GB single files

**Note:**
for **5PB** one need to get also backup:
- 5PB-**300000**€ storage system
- **50000**€ For backup system
- power consumption: 500TB - 0.4kWh, **5y** to keep 5PB up and running is about: **160000**€ for power consumption

# HPC admins: Users software

# Scientific life (top to down)



| Idea | → | Collaborate | → | Publish |

| Reading | → | Prototyping | → | Reproduce results of others | → | Share data | → | Plots |

| Data access | → | Develop software | → | Funding | → | Hardware | → | Infrastructure |

# Infrastructure (down to top)



Funding: **EFRE**

Hardware

Software

Users

Compute nodes

NVIDIA

Storage

HPC Cloud friendly

Cloud friendly? Recycle your hardware

Rocky Linux

oVirt

docker

RANCHER

lustre™

MinIO Architecture

SERVER 1    SERVER 2    SERVER 3    SERVER 32

# What we did at AIP before 2023?



**User with data and algorithms**

colab.aip.de

Hide complexity

GaiaData

**WEB BROWSER**
- Python
- Interactive plots
- Share
- Latex
- HPC
- GPU

**Libraries**

LaTeX

Panel · hvPlot · HoloViews · GeoViews · Datashader · Param · Colorcet

DASK

K Keras · TensorFlow

eano · Astropy · scikit learn

PYTORCH

**Hardware with accelerators**

docker

PROXMOX

NVIDIA

HPC clusters, LustreFS, IB,10Gbit

Since 2016, single docker

COCALC
Collaborative Calculation and Data Science
by Sagemath, Inc.

# CoCalc Integrated Tool: LLM as an assistant



- ai.aip.de hosting local LLM at AIP

- Based on **ollama**

- **And more...**

- colab.aip.de: 250 users, over 900+ projects over 6 years
- Ai.aip.de: over 75 users in 3 weeks

# The whole complexity is obscured from the users

Users want the all tools in one place
- Data+LaTex+Code
- Collaborators to share
- Article versions
- Cluster access
- Easy publishing for the demo notebooks

Possible solution:
- dask+Kubernetes
- CoCalc project

CoCalc is a web-based cloud computing and course management platform for computational mathematics. Part of the Sage project, it supports editing of Sage worksheets, LaTeX documents and Jupyter notebooks.

# Global Workflow of StarHorse team

**A Bayesian code to estimate the photo-astrometric distances, extinctions, and astrophysical parameters for Gaia DR2 stars**
F.Anders et al. (2019)



## Getting the data

Get the list of the files: `wget --no-check-certificate http://data.aip.de/data/starhorse/fits/list-fits.txt`

Download the data: wget --no-check-certificate -i list-fits.txt

- Access examples: starhorse_db
- cmd_from_db: launch binder | Launch on Google Colab
- cmd_from_db_chunking: launch binder | Launch on Google Colab

https://data.aip.de/projects/starhorse2019.html

6 weeks  3000 cores get
**400 000 000** Stellar parameters

# What we learn from notebooks+jupyter?

- No versioning ( even py codes are not versioned)
- No git (it is somehow possible but no one does this)
- No share
- No modularity
- 2-3 years cant run, for got the parameters in the cell.
- Astronomers during prototyping are writing terrible codes.

# What about kubernetes?

- in Astrophysics infrastructure we are still in the same stage as **"Docker Inc." was in 2014**.

    Why?

- It was complex

- Rapid development in the Industry

- No LTS

    Situation is matured in 2021:

- Because of https://www.cncf.io/



- We are ready to adopt some Infa from industry into to scientific life

# Kubernetes



Kubernetes Master

API Server
Controller Manager
Scheduler
etcd

Developer / Operator

Users

Berliner Philharmoniker

Kubernetes Node

Kubelet   cAdvisor   Kube-Proxy

Pod   Pod   ...   Pod

Kubernetes Node

Kubelet   cAdvisor   Kube-Proxy

Pod   Pod   ...   Pod

Plugin Network (eg Flannel, Weavenet, etc )

Are we special?
Users images are so huge they are filling local host disks

# Microservices: Reproducible science

## Use Cases at AIP

- **Colab.aip.de**
  - Quotas
  - Project isolation
- **Data analysis pipelines with versioning**
  - Reproducible science
  - Pipeline versioning
    - GaiaDR1,2,3
    - RAVEDR1-6
    - StarHorse-18,19,20
- **Publish papers with interactive plots**
  - like binder
  - Example: distill.pub by google
- **Dynamically Scalable webpages**
- **gitlabs @ aip:** CI integration

**Pros:**
- ***Direct GIT integration***
- ***Scalability***
- Modularity
- Distributed development
- Integration
- Save resources/power

**Concerns:**
- Complexity
- Design
- ***Testing, debugging***
- ***Inter-service call latency***

# REANA: Infra

```
helm install reanadev24  reanahub/reana  --create-namespace -n reanadev24  -f values.yaml
```

# SAAS,IAAS and PAAS

| Hardware | Software |
|---|---|
| Compute Nodes GPUs | Linux Ovirt containers kubernetes slurm |
| Storage | MinIO-S3 Lustrefs-IB Glusterfs/nfs |
| Network: Intern/Public | Infiniband 10G 1G |

## User Portals



COCALC    reana NEW

MINIO    GitLab

## Job queues, resource management

slurm workload manager    reana

# REANA and AIP discussion rounds

# REANA: in Action



https://reana-p4n.aip.de/

Use **reana-client** from terminal

Connect to gitlab

https://gitlab-p4n.aip.de/arm2arm/reanatest



Launch from URL

# REANA hosting arbitrary webpage



reana-jailbreak?

# MLFlow: as a ML experiments server



- Can we deploy this within the reana?

- Answer: yes

- Security concerns...

- I loved:
  - Streamlit
  - Mlflow
  - Panel
  - React native page

# Coming Soon:
# Global Workflow of StarHorse team

**Transferring spectroscopic stellar labels to 220 million Gaia DR3 XP stars with XGBoost**
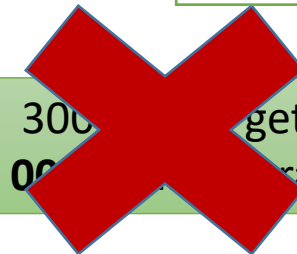


## Getting the data

Get the list of the files: `wget --no-check-certificate http://data.aip.de/data/starhorse/fits/list-fits.txt`

Download the data: wget --no-check-certificate -i list-fits.txt

- Access examples: starhorse_db
- cmd_from_db: launch binder  Launch on Google Colab
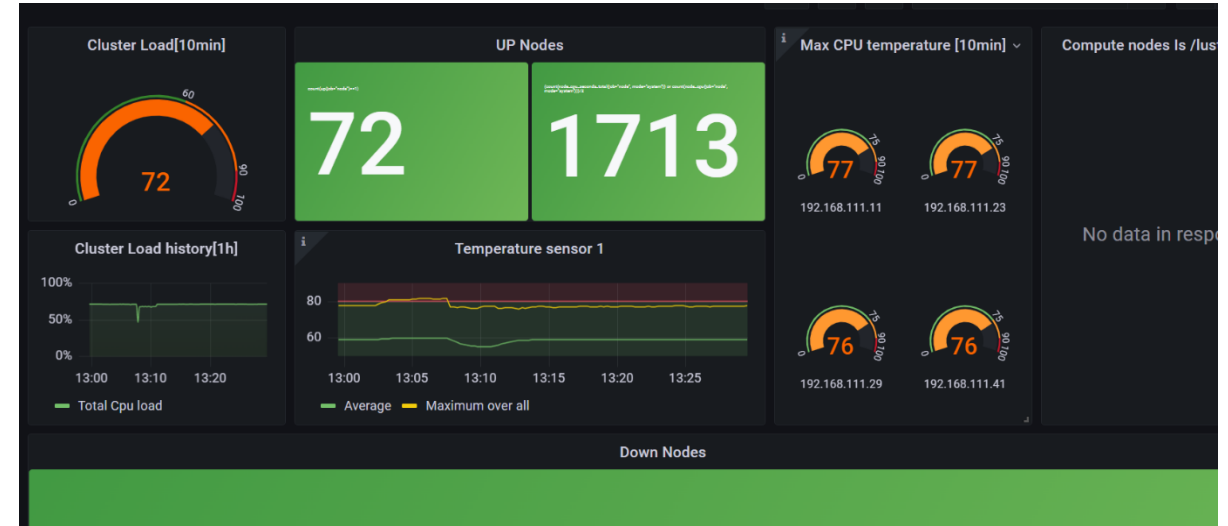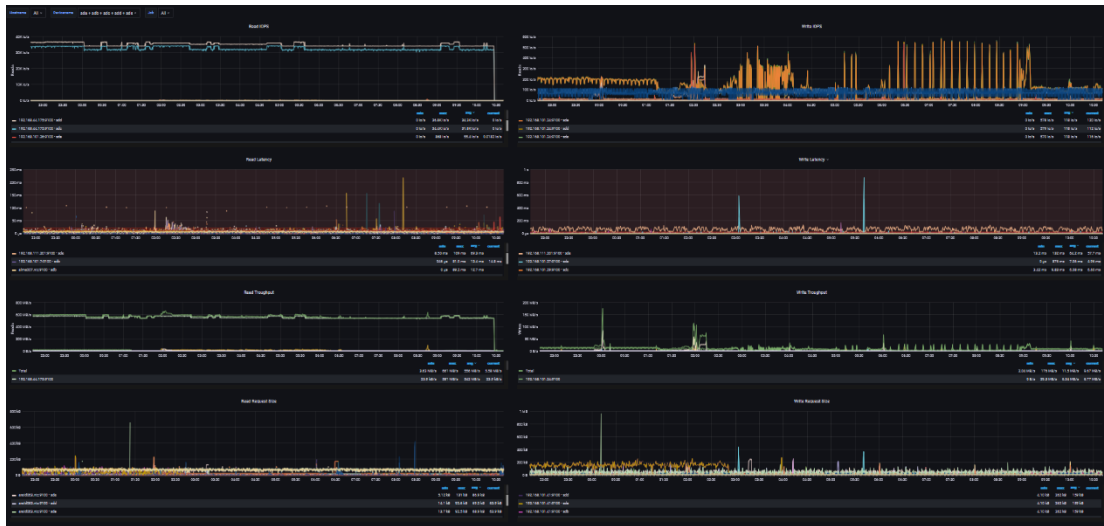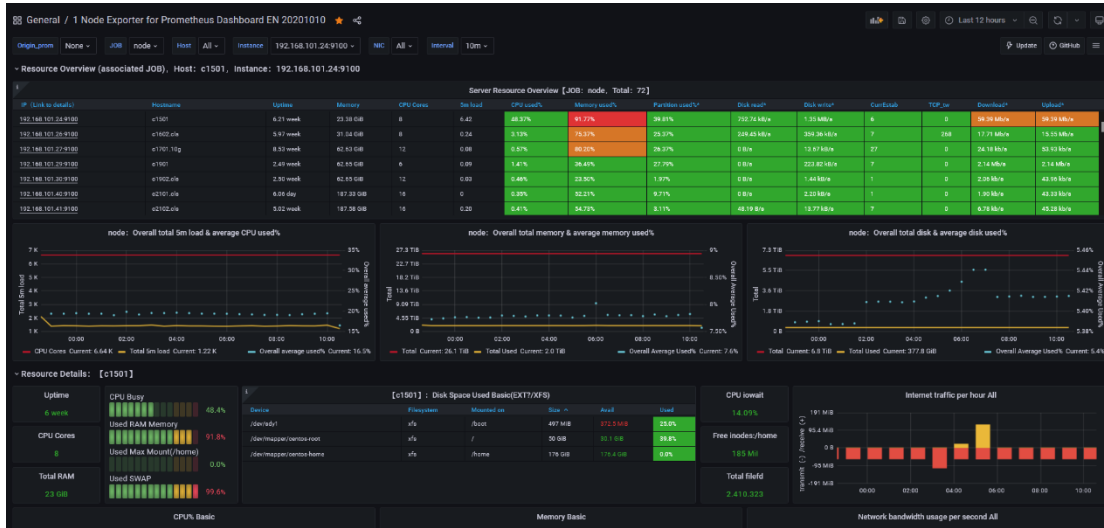- cmd_from_db_chunking: launch binder  Launch on Google Colab

1 A100
1 week
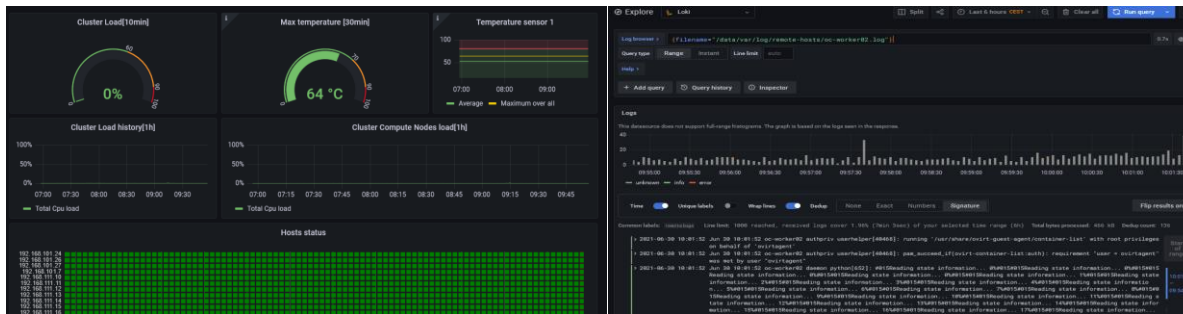Preparation 2 years

6 weeks  300 ... get
**400 000 00** ...rameters

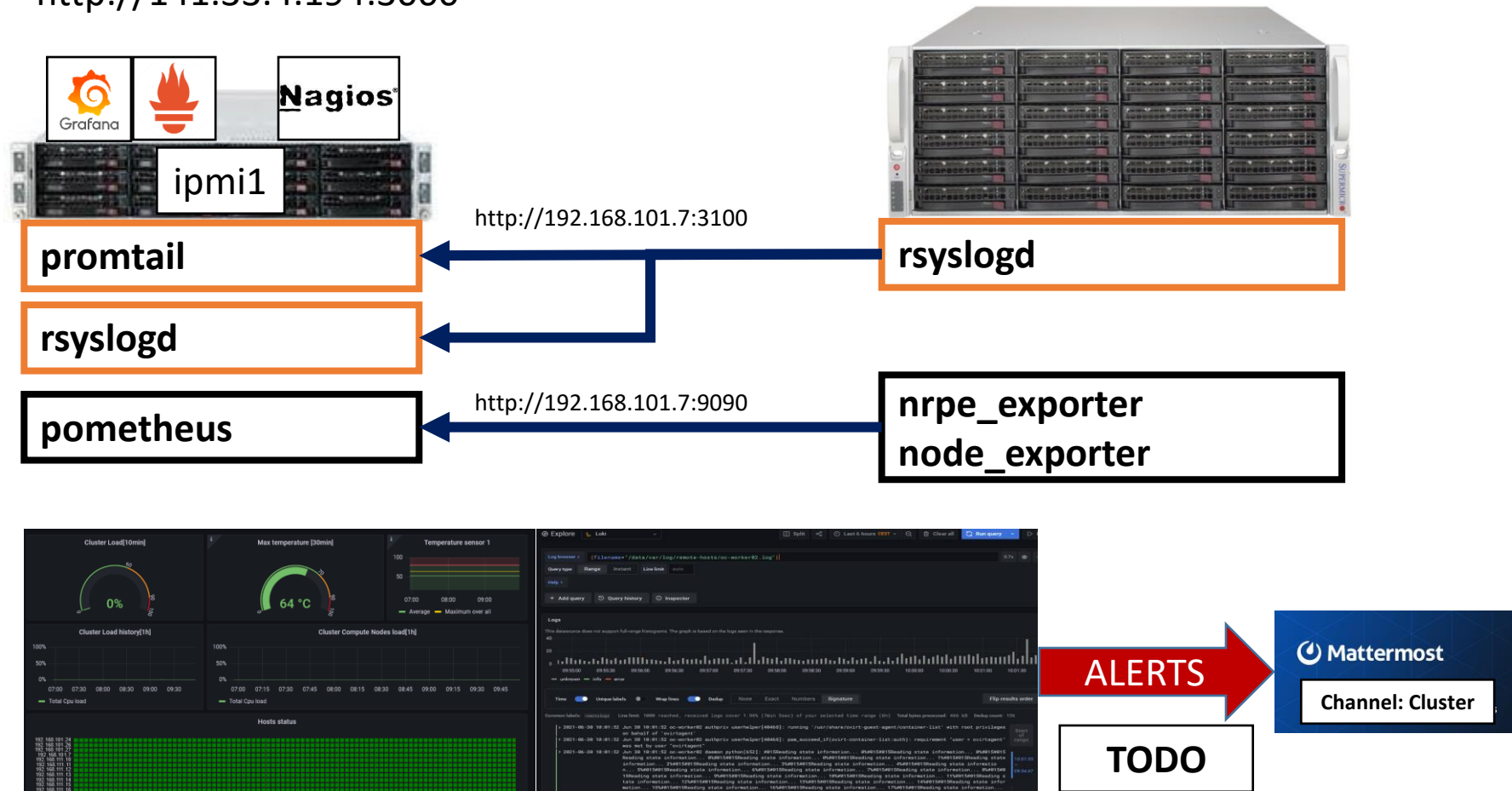# Monitoring is important!

# Cluster Monitoring (Grafana)

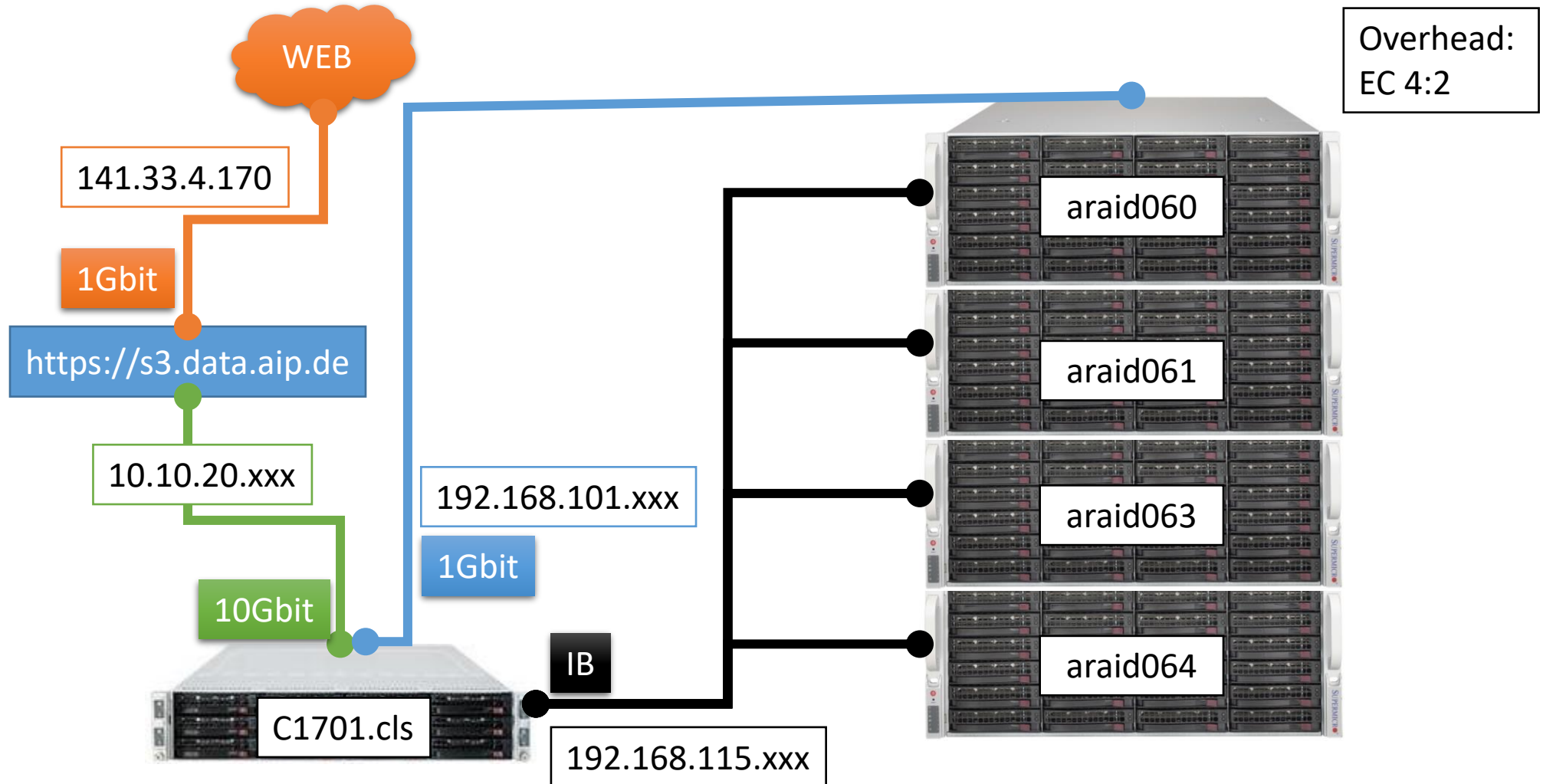# Monitoring stack

# Data management

# S3 storage at AIP: MinIO network

# Moving data from clusters?

# REANA as a main ingredient for NFDI4PUNCH

- Helmholz-AAI is integrated and working w/o problems
- Registered users: >50 users within 3 months
- We are looking for stability tests to announce at AIP
- Gitlab container registry as a main source for containers
- Dev steps:
  - Actively developing HT_Condor integration(Manuel)
  - Any SLURM backend(Arman,Elena)
  - Merging to REANA basecode(Tibor,Marco and team…)
- What is still missing:
  - Workflow shares
  - Data to(from) Storage workflows
  - Easy token management
  - Not implemented the LustreFS integration

Use CLOUD everywhere!!!

Free Software Foundation Europe: fsfe.org

Elena will demonstrate how to use REANA with S3 storage and more

# Questions?