

An infrastructure for the reproducible scientific workflows.

Dr. Arman Khalatyan

Researcher at eScience/Supercomputing/IT

Leibniz-Institut für Astrophysik Potsdam (AIP), Germany

17.09.2021

Scientific life



Scientific life (top to down)



Scientific products



Research technology and infrastructure



Compute Cluster Simulations in Cosmology & Magneto-Hydro-Dynamic Leibniz: 1700 Cores Newton: >3000 Cores GPUs: >8(NEW)





Storage / Archives / Database

Simulation and observation data

ca. >5 Petabyte Storage,

Virtual Research Environment for CLUES, MUSE

Databases: CosmoSim, RAVE, Gaia, Plates Archives

Down to top problem



Can we obscure the HW and Infra form the users?



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.





Local workflow: "..but it works on my laptop"





Global reproducible workflow: "...works everywhere"



Global Workflow of StarHorse team

A Bayesian code to estimate the photoastrometric 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

An interactive StarHorse data access example (github+gaia.aip.de+binder or googlecloud)



What is **docker** as an astronomer? A streamlit example app.



Point your browser to localhost:8581 You have an interactive plots

Docker

DOCKER COMPONENTS



What about kubernetes?

• in Astrophysics infrastructur still in the same stage as "Do Inc." was in 2014.

Why?

- It was complex
- Rapid development in the Ir
- No LTS

Situation is matured in 20

Because of <u>https://www.cnc</u>

CLOUD NATIVE

 We are ready to adopt some from industry into to scienti





Kubernetes



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: <u>distill.pub</u> 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

Upcoming colab.aip.de vision

> DEPLOYMENT cocalc-kubernetes-server, #1	1.1 < 0.01 0.1 Gib Memory Cores CPU Kib/s Network	1 pod	:
 POD project-02cf92a0-d4c6-4735-850d- 5f94d564324c 	140 < 0.01 0.02 Mib Memory Cores CPU Kib/s Network	1 pod	:
 project-5279d32f-ffdf-4843-859c- fbbb901c3dec 	140 < 0.01 0.02 Mib Memory Cores CPU Kib/s Network	1 pod	:
 project-642e581e-0242-482c- af3d-55e29e59340d 	140 < 0.01 0.04 Mib Memory Cores CPU Kib/s Network	1 pod	:
Quota:CPU/SPACE/RAM	Project2 Quota:CPU/SPACE/RAM		
pod1	pod2	•	

What's next?

- Large Kubernetes clusters deployments on demand via PUNCH4NFDI
- Parallel disk IO to deal with TBs of data
 - Connection with Lustrefs, S3
- Easy to maintain
- Clone service on demand
- Reuse/reshare unused resources
- Snapshots for versioning of the data-pipelines
- 24/7 HA
- Special hardware share: GPU, SSD, NVME
- LTS Container Registry for each publication
- Share docker containers with publications