

# Design your e-infrastructure!

<https://indico.egi.eu/indico/event/4434/>

Use case:

Break out group coordinator:

Bjorn / Enol

...



Amsterdam, 9. May, 2019.

# Group members

- ASTRON
- NenuFAR
- PROCESS
- EISCAT\_3D
- EGI

# First break-out

**Background and Users**

# Who will be the user? Can the users be characterised? How many are they?

	ASTRON	NenuFAR	PROCESS	EISCAT
Users	<ul style="list-style-type: none"> <li>• Astronomers</li> <li>• Astrophysicists</li> <li>• Particle physicists</li> </ul>	<ul style="list-style-type: none"> <li>• As ASTRON</li> </ul>	<ul style="list-style-type: none"> <li>• Pipeline developers</li> <li>• Astronomers</li> <li>• Medical</li> </ul>	<ul style="list-style-type: none"> <li>• Space physics (near and far)</li> <li>• Plasma physics</li> <li>• Atmospheric physics</li> </ul>
Characteristics	<ul style="list-style-type: none"> <li>• Transition from local to cloud services (zero downloads)</li> <li>• Pipeline developers</li> <li>• Backend providers</li> <li>• Provide own workflows</li> <li>• Production managers</li> <li>• Support for L0 through L4 data is needed</li> </ul>	<ul style="list-style-type: none"> <li>• ±same as ASTRON</li> <li>• Experience using LOFAR → used to working remotely</li> </ul>	<ul style="list-style-type: none"> <li>• Prepare for Exa-scale</li> <li>• 2 sets of users</li> <li>• 1. Pipeline dev</li> <li>• 2. End-end users</li> <li>• User can build code (1)</li> <li>• Unable to scale up beyond local resources (1)</li> <li>• End-end users just use portals (high level of abstraction)</li> </ul>	<ul style="list-style-type: none"> <li>• 3 sets of users</li> <li>• 1. just download and use data locally</li> <li>• 2. Pipeline devs</li> <li>• 3. Operational services... job scheduling++</li> <li>• Events trigger pipeline for storing data and not throwing it away</li> </ul>
Numbers	<ul style="list-style-type: none"> <li>• O(1,000s)</li> </ul>	<ul style="list-style-type: none"> <li>• O(100s)</li> </ul>	<ul style="list-style-type: none"> <li>• O(1,000s)</li> </ul>	<ul style="list-style-type: none"> <li>• 1. O(10,000)</li> <li>• 2. O(100)</li> </ul>

# What value will the envisaged system deliver for them (the whole setup)? What will the system exactly deliver to them?

	ASTRON	NenuFAR	PROCESS	EISCAT
What should EGI do	<ul style="list-style-type: none"> <li>Integrate Scientific Analysis Platform with <b>AAI</b> and Notebooks</li> <li>HTC for batch jobs</li> <li>Payment models?</li> <li>Hosting platform?</li> <li>Training to use EGI services</li> <li>Training for service management</li> <li>Training infrastructure / environment</li> </ul>	<ul style="list-style-type: none"> <li>Data transfer and storage mechanisms. Archiving and online access for analysis</li> <li>AAI and Group Management</li> <li>Scientific Analysis Platform – Cloud Container Compute</li> <li>Interactive analysis platform for pipeline development (Notebooks)</li> </ul>	<ul style="list-style-type: none"> <li>High-throughput Data transfer</li> <li>HPC access: schedule jobs from k8s clusters</li> <li>AAI and SSO to access different systems (HTC, Cloud, HPC)</li> <li>Cloud container compute (docker + k8s)</li> </ul>	<ul style="list-style-type: none"> <li>Analysis platform</li> <li>AAI and SSO to access different systems (HTC, Cloud, HPC)</li> <li>See left</li> </ul>
Value / Gaps			<ul style="list-style-type: none"> <li>A way to</li> </ul>	

# How should they use the system?

ASTRON	NenuFAR	PROCESS	EISCAT
<ul style="list-style-type: none"><li>A user can log on to the SAP, he/she can search for available data, select the data, start up a Notebook / HTC batch processing platform, which can then see the data and he/she can analyse it easily, and publish the data/workflow via Research Object (group of of DOIs for different elements)</li></ul>	<ul style="list-style-type: none"><li>As ASTRON + access to own data and public data. User can choose Notebook / container etc that contains preferred workflow environment, including data.</li><li>Need to deliver / fetch data container to workflow</li></ul>	<ul style="list-style-type: none"><li>Pipeline developer logs on to the system and will have access to a db of containers/notebooks (configs) and resources to spin up a system for a end-end user</li><li>End-end user will log on to a webUI and select a dataset and pipeline launch the processing and retrieve the result from the webUI</li></ul>	<ul style="list-style-type: none"><li>A user has to log on to a portal, search across different levels of data (events), then uses these data in an interactive analysis platform (Notebooks / HTC), publish the data in a Research Object.</li></ul>

# What's the timeline for development, testing and large-scale operation?

(Consecutive releases can/should be considered.)

ASTRON	NenuFAR	PROCESS	EISCAT
Grab from PPTs			

# Second break-out

**Design and implementation plan**



# What should the first version include? - The most basic product prototype imaginable already bringing value to the users

(the so-called Minimal Viable Product - MVP)

ASTRON	NenuFAR	PROCESS	EISCAT
<ul style="list-style-type: none"> <li>• Check-in (AAI)</li> <li>• Notebooks</li> <li>• Online storage for Notebooks (LOFAR) → Onedata/ DataHub?</li> </ul>	<ul style="list-style-type: none"> <li>• Check-in (AAI)</li> <li>• Data (online) storage and archiving (+4PB/year)</li> <li>• Data transfer service between infrastructures (I/O needs?)</li> </ul>	<ul style="list-style-type: none"> <li>• Check-in</li> <li>• Data transfer/staging (16TB) between HTC/HPC sites (want to move data from Poznan to Krakow)</li> <li>• Data storage</li> <li>• Cloud container compute</li> </ul>	<ul style="list-style-type: none"> <li>• Check-in</li> <li>• DIRAC</li> <li>• Data transfer</li> <li>• Notebooks (experimental) + Matlab kernel (NB to sort out license)</li> </ul>

- MVPs should be easily scalable (at least think about scaling, i.e. don't make it difficult to scale).

# Which components/services already exist in this architecture?

ASTRON	NenuFAR	PROCESS	EISCAT
<ul style="list-style-type: none"> <li>• Check-in (RCAuth included)</li> <li>• Notebooks</li> <li>• DataHub (Onedata)</li> </ul>	<ul style="list-style-type: none"> <li>• Check-in (needs to talk to eduTEAMS)</li> <li>• DataHub (Onedata)</li> <li>• FTS</li> </ul>	<ul style="list-style-type: none"> <li>• Check-in</li> <li>• FTS</li> <li>• Cloud container compute (k8s)</li> <li>• DataHub (Onedata)</li> </ul>	<ul style="list-style-type: none"> <li>• Check-in</li> <li>• DIRAC (data catalogue)</li> <li>• FTS</li> <li>• Notebooks (Matlab license dependent)</li> </ul>
<b>Gaps in components / services?</b>			
<ul style="list-style-type: none"> <li>• Access to e.g. LOFAR LTA via Check-in (being worked on in EOSC LOFAR CC)</li> </ul>	<ul style="list-style-type: none"> <li>• Capacity allocation (+4PB/yr)</li> </ul>	<ul style="list-style-type: none"> <li>• HPC token translation in Check-in</li> </ul>	<ul style="list-style-type: none"> <li>• Checking out data from DIRAC to make visible in Notebooks and vice versa (→ Onedata)</li> </ul>

- DIRAC can manage batch jobs to HTC
- Submit Notebook as a job via DIRAC to HTC

# Which components/services are under development (and by who)?

ASTRON	NenuFAR	PROCESS	EISCAT
<ul style="list-style-type: none"> <li>Finding data from different communities (ESCAPE).</li> <li>SAP interface needs to be built (ESCAPE WP5) – python/django framework, OIDC</li> <li>SAP integration with Check-in, Notebooks (EGI/ESCAPE)</li> </ul>	<ul style="list-style-type: none"> <li>Data management (NenuFAR)</li> <li>Reproduce existing workflow configurations in EGI infrastructure (EGI/NenuFAR)</li> </ul>	<ul style="list-style-type: none"> <li>Configure Check-in on pipeline development platform (PROCESS/EGI)</li> <li>Some Pipelines under dev (PROCESS)</li> <li>webUI under dev (PROCESS)</li> </ul>	<ul style="list-style-type: none"> <li>Tool to populate DIRAC with metadata (EISCAT)</li> <li>Use cloud services to run containerised applications (EOSC-hub).</li> <li>Adapt VM environment specific to EISCAT applications (EOSC-hub)</li> </ul>
<b>Beyond the MVP</b>			
<ul style="list-style-type: none"> <li>Connect to HTC platforms (EGI/ESCAPE)</li> <li>Publish Research Objects (data + pipeline + publications, assign everything with a DOI)</li> </ul>	<ul style="list-style-type: none"> <li>Notebooks</li> <li>Publishing Research Objects</li> </ul>	<ul style="list-style-type: none"> <li>Scalability</li> <li>More pipelines</li> </ul>	<ul style="list-style-type: none"> <li>Publish Research Objects</li> <li>Notebooks with Matlab</li> </ul>

# Are there gaps in the EGI service catalogues that should be filled to implement the use case? Which service provider could fill the gap?

## ASTRON

## NenuFAR

## PROCESS

## EISCAT

- PRACE for HPC
  - Seamless integration
- Main work is about combining/integrating existing services – there were no major gaps identified...
- How to deal with system heterogeneity?
- Information portal about service providers.  
Up-to-date documentation about service providers and available resources.  
Machine readable specs.  
Users need to know where to submit jobs
  - Available hardware
  - Available soft / middleware
  - HTC technology support
- Pilot framework to help users define requirements (sandbox / fedcloud)
  - Improve documentation
- CI/CD (Continuous Integration/Delivery) training

# Proposed proof of concept

## Baseline infrastructure:

- 2 HTC sites; 2 cloud sites; 1 HPC site
  - Enable federated AAI with Check-in
- Enable Data Transfer Service / DataHub to work with all the sites

## Use-case specific extension:

- Configure DIRAC and Jupyter Notebooks as user interfaces
- Connect PROCESS environment
- Configure DIRAC and Notebooks to access the 5 sites:
  - User to discover dataset with DIRAC
  - Check-out test dataset from DIRAC and import them to Notebooks
  - Run exploratory analysis in Notebooks
  - Go back to DIRAC, discover full dataset
  - Run analysis with full dataset in DIRAC or Notebooks

# Tasks

- Compute sites who can support us (ideas)
  - EOSC-hub LOFAR CC members
  - PROCESS members: Krakow and Poznan
  - EISCAT\_3D VO sites  
Gergely
- Enable federated AAI on HPC (where exactly?)
- Integrate 'data import/export' feature into Notebooks (from DIRAC)
  - Enol, Michal, Andrei
- Integrate batch computing back-end with Notebooks
  - Enol (new EGI Task force!)
- Run Notebook code in DIRAC?

+ complete based on Bjorn's slide #11 ?

# NenuFAR

- Storage + Archive is the priority
  - France Grilles is evaluating it