

# Data discovery and (interactive) analysis work in SRCNet and at ASTRON

Yan Grange

# What is a science platform?

Team Tangerine (SRCNet) wrote a vision document on this!

## Attributes

- consistency
- scalability
- reproducibility
- usability
- Reliability

## Guiding principles:

- Highly collaborative
- end-to-end
- accessible (with focus on inclusion)

### SRC Science Analysis Platform Vision Document

de Boer, J.<sup>1</sup>, Cimpan, I.<sup>2</sup>, Das, A.<sup>3</sup>, Fabbro, S.<sup>4</sup>, Grange, Y. G.<sup>1\*</sup>, Hardcastle, M. J.<sup>5</sup>, Sharma, R.<sup>6</sup>, Skipper, C. J.<sup>2</sup>, Swinbank, J. D.<sup>1</sup>, Webster, B.<sup>5</sup>

<sup>1</sup>ASTRON, the Netherlands Institute for Radio Astronomy, Oude Hoogeveensedijk 4, 7991 PD Dwingeloo, The Netherlands

<sup>2</sup>Jodrell Bank Centre for Astrophysics, Alan Turing Building, The University of Manchester, Manchester, M13 9PL, UK

<sup>3</sup>École polytechnique fédérale de Lausanne, Rte Cantonale, 1015 Lausanne, Switzerland

<sup>4</sup>NRC Herzberg Astronomy and Astrophysics, 5071 West Saanich Road, Victoria, BC V9E 2E7, Canada

<sup>5</sup>Centre for Astrophysics Research, University of Hertfordshire, College Lane, Hatfield, AL10 9AB, UK

<sup>6</sup>fachhochschule Nordwestschweiz, Bahnhofstrasse 6, 5200 Windisch, Switzerland

\* corresponding author

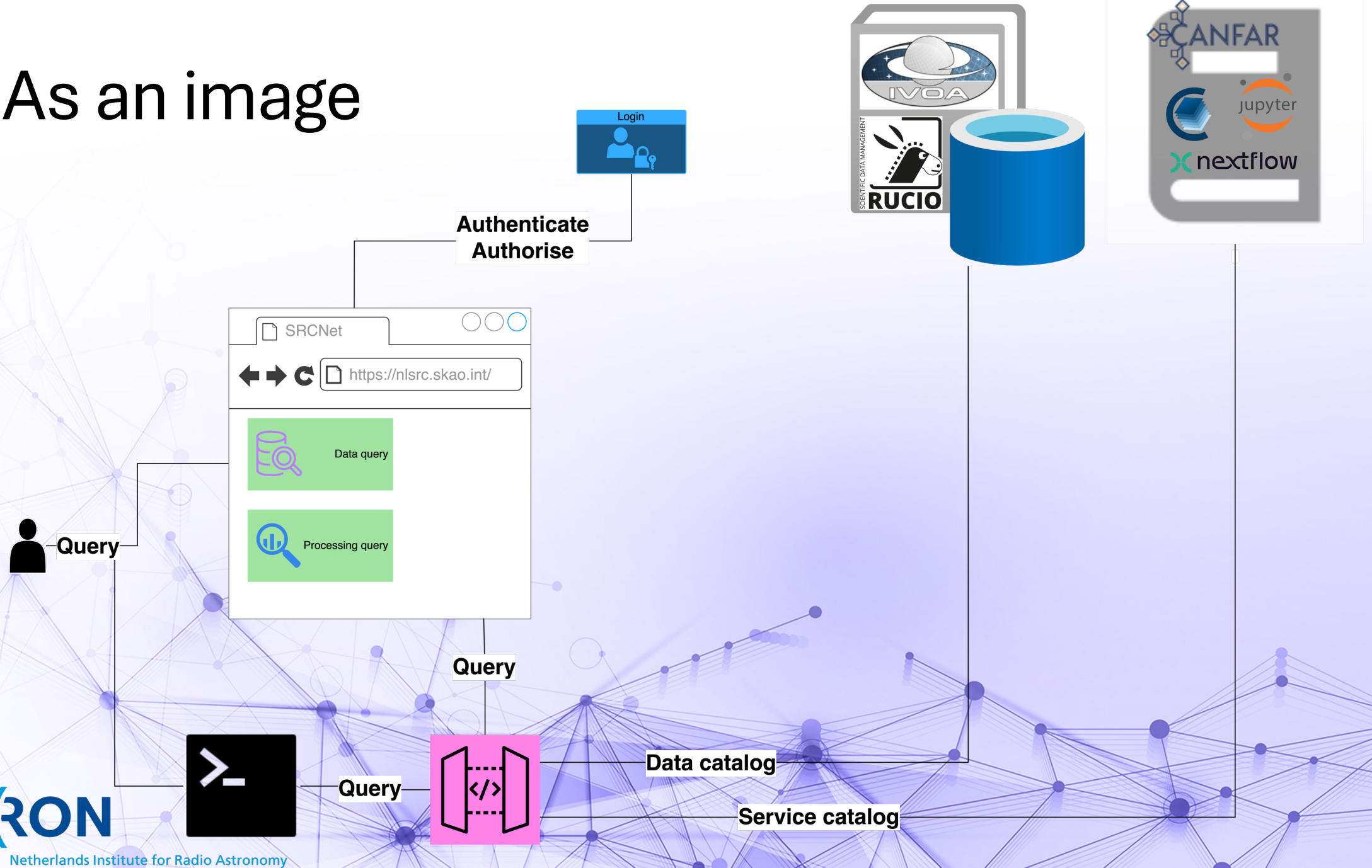
**Abstract.** This document describes the vision for the [Square Kilometer Array \(SKA\)](#) Regional Centres Science Analysis Platform. It is intended to set the broad terms of reference for the platform and to provide guidance for both development teams and other stakeholders. Among the features and services that are expected to be included are data querying and discovery tools, some form of notebook interface, user-managed software environments, workflow management, and a comprehensive set of APIs enabling access to all low-level platform functionality. This document is not a design specification, and the features and services described herein will be further refined, or could be discarded, at a later stage of development.

Document number	.....	TBC
Version	.....	1.0
Date	.....	2023-01-13
Status	.....	Released

# In concreto

- The science platform is aimed at making the data and processing available to users with different levels of expertise
  - In the SKA/LOFAR case teams are an important concept.
- Making large data sets accessible
  - Using APIs for programmatic access
  - And a user interface for interactive interaction
- With science-ready data products, the focus is mainly at analysis work rather than large-scale pipeline processing
- **NB:** Within the scope of this talk “science platform” refers to the whole system while “science gateway” refers to the interface to the user (i.e. frontend + API layer)

# As an image



# Galaxy

- Science gateways are being used in other fields. For instance Galaxy, which is focused on life science.

The screenshot displays the Galaxy web interface. At the top, the navigation bar includes 'Galaxy', 'Workflow', 'Visualize', 'Data', 'Help', and 'User'. The left sidebar contains a search bar for tools and a list of tool categories: Inputs, Get Data, Send Data, Collection Operations, Expression Tools, GENERAL TEXT TOOLS (highlighted), Text Manipulation, Filter and Sort, Sub-sample sequences files, Bigwig extremes to bed features, Filter data, Select lines, GFF, Join, Subtract and Group, and Datamash. The main workspace shows a workflow titled 'Very random workflow (unsaved changes)' on a grid. The workflow consists of four steps: 1. 'Download and Extract Reads in BAM' (output: output\_collection (input)), 2. 'Sub-sample sequences files' (output: Sequence file), 4. 'Sub-sample sequences files' (output: Sequence file). The right sidebar contains navigation icons for Upload, Tools, Workflows, Workflow Invocations, Visualization, Histories, History Multiview, Datasets, and Pages.



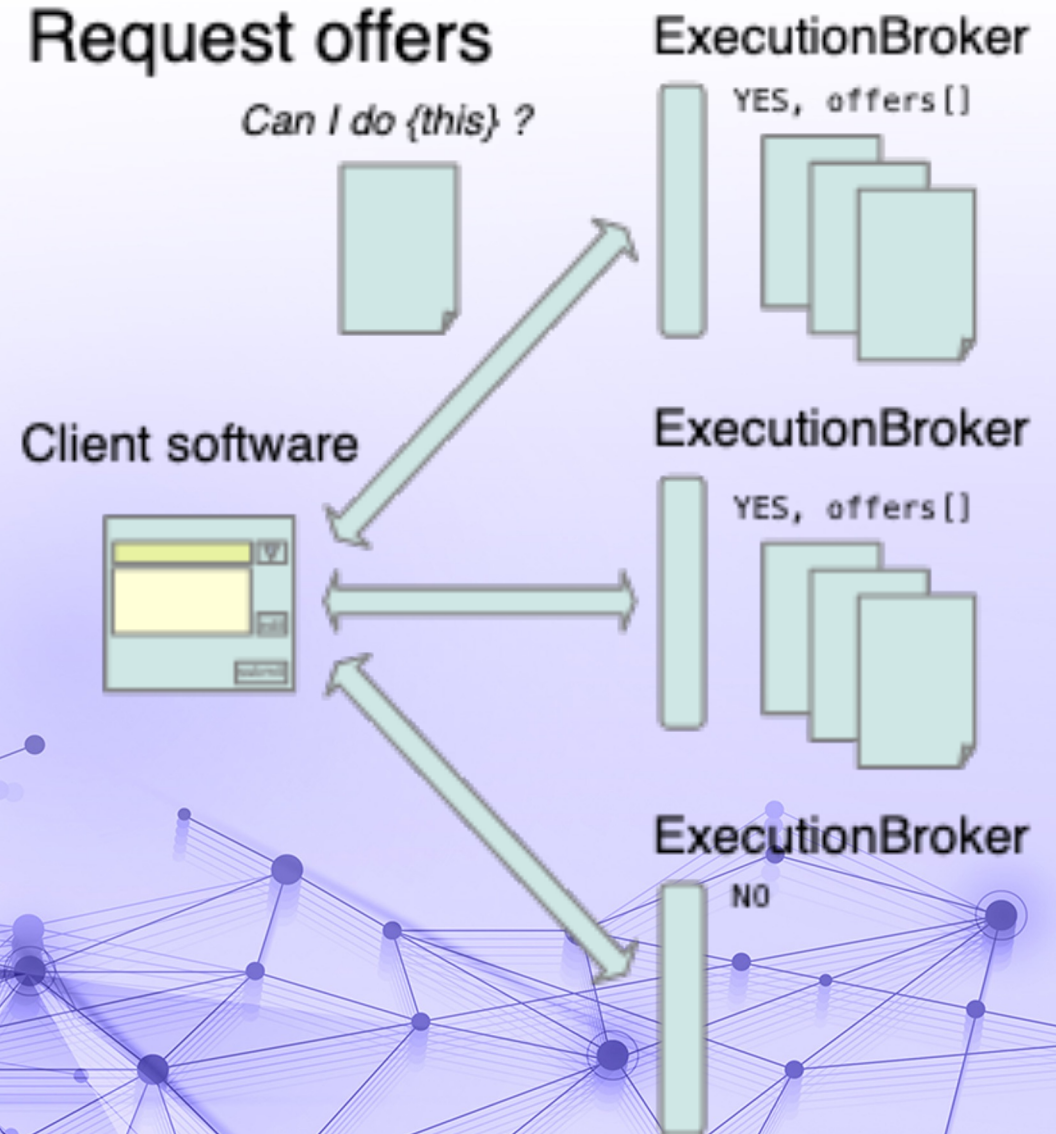
# VO standards (small intermezzo)

*The Virtual Observatory (VO) standards have been designed to “astronomers to interrogate multiple data centers in a seamless and transparent way (...) and gives data centers a standard framework for publishing and delivering services using their data.” (from [ivoa.net](http://ivoa.net))*

- Apart from standards for defining and sharing tables, and observational (meta) data some relevant standards that interact with the lower-level parts of an archive are:
  - Server-side Operations for Data Access (SODA): a low-level data access capability or server side data processing that can act upon the data files, performing various kinds of operations: filtering/subsection, transformations, pixel operations, and applying functions to the data.
  - The Universal Worker Service pattern (UWS) defines how to manage asynchronous execution of jobs on a service.

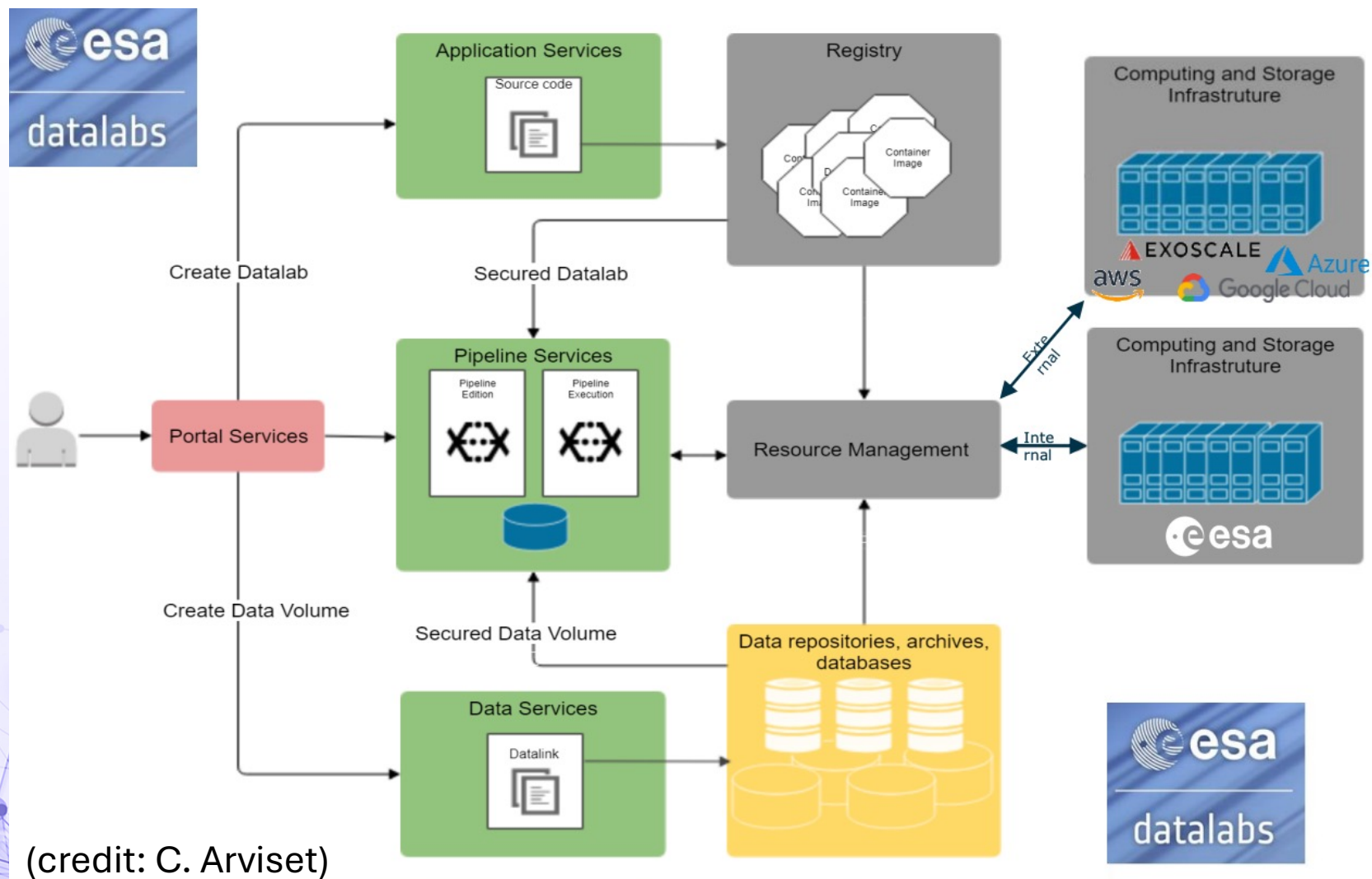
# The Execution Broker standard

- Currently in DRAFT (i.e. not endorsed as an official standard yet), and being developed in the SRCNet as a way too provision resources.
- The goal is to abstract away as many details away from the user
- This feels like quite a bit of scope...



# ESA Datalabs

- Platform for access to ESA data





# ESA Datalabs

- Notebooks
- VNC
- Pipelines (using CWL)
- Data collections
  - Directly linked to applications as “virtual directory”
- Group/team management

The image displays two screenshots of the ESA Datalabs web interface. The top screenshot shows the 'Pipeline launch' page for a JWST pipeline. It includes a 'Change version' button and a graph of the pipeline workflow. The graph shows a central 'Detector1' node connected to several input nodes: 'Keyword in input file', 'Pipeline input FITS file', 'Detector1 CRDS override list', 'CRDS cache overrides', 'CRDS cache', and 'CRDS Context'. 'Detector1' outputs to 'Detector1 intermediate output', which then feeds into an 'Image2' node. The 'Image2' node has two outputs: 'FITS file output' and 'Keyword in output file'.

The bottom screenshot shows a data visualization window titled 'SAOImage ds9'. The window displays a table of metadata for a file named 'hst\_9771\_a4\_acs\_wfc\_f606w\_j8mqa4\_drc.fits.gz[SCI]'. The table includes fields for Object, Value, FK5 coordinates (alpha, delta), Physical coordinates (x, y), Image coordinates (x, y), and Frame 1 coordinates (x, y). Below the table is a grid of toolbars for file, edit, view, frame, bin, zoom, scale, color, region, wcs, illustrate, analysis, and help. The main area shows a tilted astronomical image of a star field with a prominent white diagonal line.

Object	Value
FK5	$\alpha$ 14:07:01.7110 $\delta$ +35:03:56.288
Physical	x 4964.9 y 2824.6
Image	x 4964.9 y 2824.6
Frame 1	x 0.134588 y 0

# ESAP

- Developed in the ESCAPE project
- The development inspired the work on ADEX and SRCNet, but also other ESCAPE partners (e.g. CTA) are planning to reuse the concept.

For acronym soup lovers:

*“ESCAPE ESAP” is the European Science Cluster of Astronomy and Particle physics  
European Strategy Forum on Research Infrastructures research infrastructures  
European Strategy Forum on Research Infrastructures Science Analysis Platform.*

*(thanks, John, for working this out!)*

The screenshot shows the ESCAPE ESAP website header with navigation links: Archives, Multi Query, Interactive Analysis, Batch Analysis, Asynchronous Jobs, and IVOA-SAMP. There are also social media icons and a user profile for 'Logout Klaas Kliffen'. Below the header are four project cards:

- WSRT-Apertif**: Includes an image of a radio telescope and text about Apertif Surveys, mentioning data from imaging and time-domain surveys, high-time resolution filterbank data, and raw observations in MS standard format.
- ASTRON VO**: Includes the ASTRON Virtual Observatory logo and text explaining that the VO defines standards for downloading astronomical data, including image surveys in FITS format.
- Zooniverse**: Includes a target icon and text about the Zooniverse Classification Database, describing it as the world's largest and most popular platform for people-powered research.
- Virtual Observatory (VO)**: Includes the IVOA logo and text defining the VO as a set of standards for downloading astronomical data. A button below reads 'Visit Virtual Observatory (VO) Archives'.

(credit: K. Kliffen)

# The SRCNet

- The numbers for SKA may be bigger, but the image looks very much like the LOFAR setup.



(credit: J. Salgado)

# SRCNet components

- AAI through IAM
- Data management through Rucio
- Visualisation (CARTA, VISIVO, etc.)
- Several internally developed APIs (e.g. a service catalogue)
- CANFAR/Azimuth (see next slide) for provisioning of resources
  - Or providing direct access to e.g. a Jupyter lab instance
- Software repository for processing tools
- ...

# CANFAR, Azimuth

SRC | Net

## Science Portal

### Active Sessions

No interactive sessions found

### New Session

type  notebook  
desktop  
carta  
contributed  
notebook

container image

name

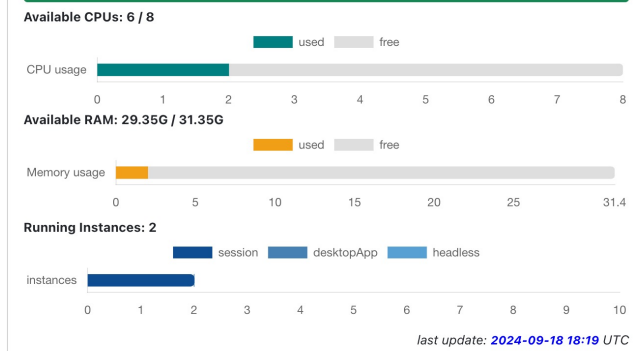
memory 8

# cores 2

Launch Reset

The dashboard shows a grid of platform cards for 'ska-src-cambridge'. The cards are organized into two rows. The first row contains two cards: 'mjh-carta' (Linux Workstation) and 'mjh-test' (Linux Workstation). The second row contains five cards: 'bb23-lapcat-4mm' (Slurm), 'bb23-pybsdf-8gm' (Slurm), 'bo307-k8s' (Kubernetes), 'bo307-k8s-test' (Kubernetes), and 'bo307-manila-benchmark' (Linux Workstation). Each card displays the platform logo, name, description, and status (READY or ERROR). A 'Details' button is present on each card.

### Platform Load



(credit: M. Hardcastle)

# ADEX

- <https://sdc.astron.nl/adex-next/> (better to live-demo than to talk too much about it)

# SRCNet science gateway

- <https://gateway.srcdev.skao.int/> (better to live-demo than to talk too much about it)

# What does this mean for the people here?

- Many of the systems, tools, software used expect the data to be accessible like it is on disk, but generally it is not.
  - Remote/network mounts can get (very!) laggy.
- Data products themselves may be very large. In the VO world, this is generally solved by offering a cutout service (only deliver the pixels the user wants rather than downloading a full image).
  - If the data is remotely stored, are there smart ways of doing this?
- Using a platform for interfacing offers flexibility of backend implementations. At what level is it needed to be homogeneous?
- How much information can/should we hide from the user?
  - Staging, data locality have effects on this.
- How is this mapped to software access, and users writing their own code?
  - In SRCNet we are looking into EESSI for this
- ... and probably much more I did not think of