

# Evaluating existing science platforms for SKA development



**Team Tangerine, SKAO Prototyping Team  
(Total members ~ 40)**

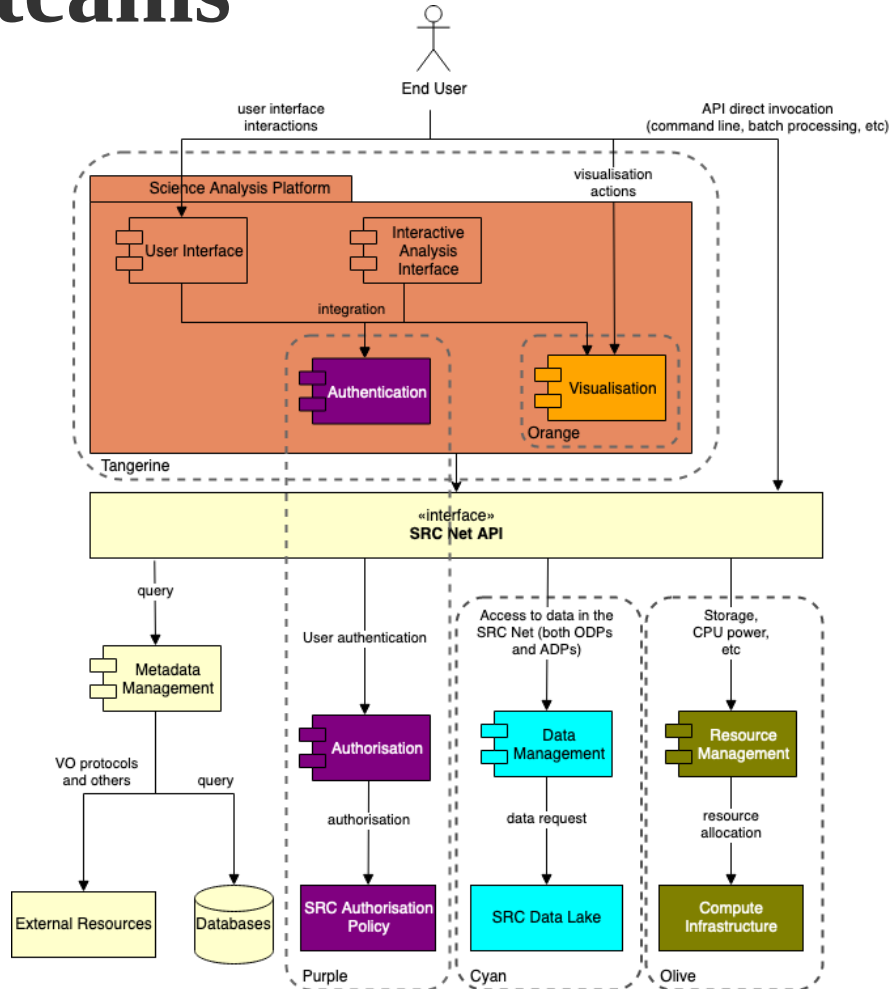
- Presented by: Rohit Sharma, FHNW, Windisch

12<sup>th</sup> Jan 2023

Winter SKACH Meeting, University of Basel, Switzerland



# Prototyping teams



# Contents of the talk



- What is a Science platform?
- Requirements of SRCNet
- Challenges
- List of existing science platforms
- Down listing the science platforms
- Evaluating Criteria for system architecture
- Evaluating selected platforms



# What is a Science platform?

- Analysis platform permitting scientists a collaborative handling of large and diverse datasets with access to large-scale computing resources.
- Linked to a specific science projects
- Must provide features that are relevant to the user base and avoid any pitfalls

# General Criteria for a Science Platform



- Consistency
- Scalable
- Reproducibility
- Usability
- Reliability
- Accessibility

# Some Requirements at high-level SRC system requirement



- Provide accessible and straightforward interfaces / must be available to widest possible user base
- Provide expert-compute users with maximal flexibility
- Provide appropriate interfaces to the lower-level services and tools deployed across the SRCNet
- Platform will only provide users with access to services and/or data products to which those users are entitled and will include limited, public, anonymous user access.

# Challenges Involved



- Data Volume: 700 PB/year
- Data Distribution of ODPs – image cubes, uv-grids, calibrated vis, pulsar timing solutions, calibration data and metadata
- Typical design and working of the platform is driven by a scientific projects, but SKA has many science projects

# Existing Science Platforms



List of various suggested science platforms

		Science Platform	Acronym (if any)	Links, references, etc	Notes	
1	SP1	Service for Web-based Analysis	SWAN	<ul style="list-style-type: none"><li>ScienceDirect Article: SWAN</li><li><a href="https://www.sciencedirect.com/science/article/abs/pii/S0167739X16307105?via=ih30hub">https://www.sciencedirect.com/science/article/abs/pii/S0167739X16307105?via=ih30hub</a></li><li><a href="https://swan.web.cern.ch/swan/">https://swan.web.cern.ch/swan/</a></li><li><a href="https://swan.docs.cern.ch/intro/what_is/">https://swan.docs.cern.ch/intro/what_is/</a></li></ul>		
2	SP2	SciServer		<ul style="list-style-type: none"><li>ScienceDirect Article: SciServer - paywall-free access via arXiv at SciServer (Taghizadeh-Popp et al, 2020)</li><li><a href="https://www.sciserver.org/">https://www.sciserver.org/</a></li></ul>		
3	SP3	Rubin Science Platform		<ul style="list-style-type: none"><li>LSST Science Platform Vision Document</li><li>LSST L&amp; Plans (LineA Workshop): <a href="https://youtu.be/nwH4RWGyU7n-2302">https://youtu.be/nwH4RWGyU7n-2302</a></li><li>Rubin Science Platform: <a href="https://data.lsst.cloud/">https://data.lsst.cloud/</a></li><li>Rubin (LineA Workshop 2021 link): <a href="https://youtu.be/nwH4RWGyU7n-1050">https://youtu.be/nwH4RWGyU7n-1050</a></li></ul>		
4	SP4	European Strategy Forum on Research Infrastructure - Science Analysis Platform	ESCAPE ESAP	<ul style="list-style-type: none"><li><a href="https://projectescape.eu/services/esfri-science-analysis-platform">https://projectescape.eu/services/esfri-science-analysis-platform</a></li><li>Backlog + JupyterHub + OpenData: <a href="https://indico.in2p3.fr/event/23333/sessions/14433/attachments/62070/84839/lapp-escape-wp2-arturoes-220121.pdf">https://indico.in2p3.fr/event/23333/sessions/14433/attachments/62070/84839/lapp-escape-wp2-arturoes-220121.pdf</a></li><li>ADASS proceedings – need to track down a link.</li><li>LineA Workshop 2021 link: <a href="https://youtu.be/YHJw93BwZ0E7n-2935">https://youtu.be/YHJw93BwZ0E7n-2935</a></li></ul>		
5	SP5	CyVerse		<ul style="list-style-type: none"><li>Original paper: PLOS</li><li><a href="https://cyverse.org/">https://cyverse.org/</a></li><li>Science APIs: <a href="https://cyverse.org/Science-APIs">https://cyverse.org/Science-APIs</a></li><li>Data Commons Service: <a href="https://datacommons.cyverse.org/">https://datacommons.cyverse.org/</a></li><li>LineA Workshop 2021 link: <a href="https://youtu.be/YHJw93BwZ0E7n-4944">https://youtu.be/YHJw93BwZ0E7n-4944</a></li></ul>		
6	SP6	MeerKat Toolbelt		<ul style="list-style-type: none"><li>Introduction to MeerKat Toolbelt</li><li>LineA Workshop 2021 link: <a href="https://youtu.be/ITDjYSUuK7n-123">https://youtu.be/ITDjYSUuK7n-123</a></li></ul>		
7	SP7	China Virtual Observatory	China VO	<ul style="list-style-type: none"><li>ScienceDirect Article: China-VO (arXiv)</li><li><a href="http://www.china-vo.org/doc.html">http://www.china-vo.org/doc.html</a></li></ul>		
8	SP8	Canadian Advanced Network for Astronomical Research	CANFAR	<div>16 SP16 WholeTale</div> <div>17 SP17 ESO Science Archive ESO</div> <div>18 SP18 Pangeo</div> <div>19 SP19 NEWT NERSC</div>	<ul style="list-style-type: none"><li><a href="https://wholetale.org">https://wholetale.org</a></li><li><a href="http://archive.eso.org/cms.html">http://archive.eso.org/cms.html</a></li><li><a href="https://pangeo.io">https://pangeo.io</a></li></ul> <ul style="list-style-type: none"><li>The ESO Archive doesn't (as far as I know...) have notebooks or workflows, but it does provide data discovery, exploration, access, etc.</li><li>Geophysics, not astronomy, but many of the technologies and considerations are likely similar.</li><li>Project looks quite abandoned (last commit is from 2019 across all forks) so may be interesting architecture-wise.</li></ul>	
9	SP9	French Space Agency - CNES Platforms		<div>20 SP20 Astronomy Commons Platform UW DIRAC</div>	<a href="https://arxiv.org/abs/2206.14392">https://arxiv.org/abs/2206.14392</a>	Small scale but automated deployment on AWS
10	SP10	CDS Astronomy Data Centre	CDS	<div>21 SP21 REANA</div> <div>22 SP22 Galaxy</div>	<ul style="list-style-type: none"><li><a href="https://www.reanahub.io">https://www.reanahub.io</a></li></ul> <a href="https://galaxyproject.org/">https://galaxyproject.org/</a>	<ul style="list-style-type: none"><li>Another CEBS effort. Maybe complementary to SWAN – this provides workflows, while SWAN provides notebooks (! think...).</li></ul> Galaxy is a science platform aimed primarily at the life sciences.
11	SP11	NASA HEASARC			<a href="https://sciencelibrary.org/index.php?option=com_content&amp;view=article&amp;id=405/galaxy-project-enabling-an-active-global-research-community&amp;catid=32w">https://sciencelibrary.org/index.php?option=com_content&amp;view=article&amp;id=405/galaxy-project-enabling-an-active-global-research-community&amp;catid=32w</a>	
12	SP12	NOIRLab Astro Data Lab		<div>23 SP23 WLOG: a distributed infrastructure for large-scale scientific computing WLOG</div>	<ul style="list-style-type: none"><li><a href="https://wlog.web.cern.ch/">https://wlog.web.cern.ch/</a></li><li>Technical Documents: <a href="https://wlog-docs.web.cern.ch/7d6-technical_documents">https://wlog-docs.web.cern.ch/7d6-technical_documents</a></li><li>LineA Workshop 2021 link: <a href="https://youtu.be/YHJw93BwZ0E7n-1015">https://youtu.be/YHJw93BwZ0E7n-1015</a></li><li><a href="https://twiki.cern.ch/HEPTape">https://twiki.cern.ch/HEPTape</a></li><li><a href="https://wlog.web.cern.ch/activities/working-groups">https://wlog.web.cern.ch/activities/working-groups</a></li></ul>	<ul style="list-style-type: none"><li>Deployed and used in CEBS</li></ul>
13	SP13	European Space Astronomy Centre - Science Data Centre	ESAC - SDCC	<div>24 SP24 nanoHUB</div>	<a href="https://nanohub.org/simulate/">https://nanohub.org/simulate/</a>	Modelling and simulation gateway for nanotechnology scientists and engineers. Jupyter notebook-type interface running transparently in a scientific computing cloud.
14	SP14	Spanish SDC	SPSRC	<div>25 SP25 The Virtual Imaging Platform VIP</div>	Paper at: <a href="https://docs.google.com/a/ind.edu/viewer?as=v&amp;pid=pdf&amp;srcid=bnQuZWR1fG3c3cyMDE4fG404NmOGRMjDz0T0EjNTRnZjI">https://docs.google.com/a/ind.edu/viewer?as=v&amp;pid=pdf&amp;srcid=bnQuZWR1fG3c3cyMDE4fG404NmOGRMjDz0T0EjNTRnZjI</a>	A medical imaging platform allowing new applications to be imported through Docker.
15	SP15	European Space Agency	ESA		Platform at: <a href="https://vip.creatis.insa-lyon.fr/documentation/">https://vip.creatis.insa-lyon.fr/documentation/</a>	
16	SP16	WholeTale		<div>26 SP26 InterMine Science Gateway</div>	<a href="http://intermine.org/documentation/">http://intermine.org/documentation/</a>	A science platform for biology data, providing web-enabled access to large databases via SQL, with tools for building queries. No notebook element, but provides extensive API access from Perl, Python, Ruby, and Java.
17	SP17	ESO Science Archive	ESO	<div>27 SP27 HubZero</div>	<a href="https://hubzero.org">https://hubzero.org</a>	The generalised version of the nanohub platform. It is somewhat unclear to me what HubZero exactly is but I think it is a framework that can be used to build a Science Analysis Platform/Gateway.
18	SP18	Pangeo		<a href="https://eexplore.jeeo.org/document/5432299">https://eexplore.jeeo.org/document/5432299</a>		



# Down Selecting the List



- Galaxy:

- Cyverse:

- ESA Datalabs

- WholeTale

- CERN REANA

- CERN SWAN

- CANFAR

- ESCAPE ESAP

- SciServer

- Open Source

- Maintenance status

- Data discovery

- Notebooks

- Workflows

- Software Distribution

- API access

- Deployment Platform

- Usability and Accessibility

- Production status

A	B	C	D		
Scoring	Platform ID	Science Platform	Implementation Language(s)	Major frameworks and technology	Open Source
18,00		22 Galaxy	Python, javascript, others	RStudio, Jupyter	2
16,33		14 Spanish SRC		Openstack, On-Demand Elastic Clusters, On-Demand	1
16,00		5 CyVerse	Go, C++, Javascript, Python	iRODS, Django, Agave API, React.js, Material-UI	0
14,67		13 ESA Datalabs	most is in Python in the JWST section		2
14,00		16 WholeTale		The Whole Tale platform leverages and extends a v	2
14,00		21 CERN REANA	Python	Kubernetes, HTCondor, Slurm, Docker, Singularity	2
13,67		4 ESCAPE ESAP	Python, Javascript	Django, React.js	2
13,00		24 nanoHUB		Linux, Apache web server, a MySQL database, PH	1
12,00		2 SciServer	C# ASP.NET Javascript		1
11,67		8 CANFAR	Java	kubernetes on OpenStack, Harbor, 0	2
11,67		17 ESO Science Archive	JSky - Java Dataset - web based access		1
11,00		1 CERN SWAN	Python, Javascript	JupyterHub, Spark, EOS, CernVM-FS	1
10,33		12 NOIRLab Astro Data Lab	Python	JupyterHub, Postgres+Q3C	1
10,00		25 The Virtual Imaging Platform	Java, Javascript, Python, HTML		2
10,00		28 AMDA http://amda.irap.omp.eu	PHP, Javascript, HTML	C	1
9,67		3 Rubin Science Platform	Python, C++	kubernetes, JupyterHub, Dask, Firefly	1
9,67		9 French Space Agency -CNES Platforms There are a bu	Java, R	Vue.js, Leaflet.js	1
9,00		29 VirES	Python, javascript HTML	WPS	1
8,00		26 InterMine Science Gateway	Java, SQL	PostgreSQL	2
8,00		27 HubZero	PHP, Javascript, HTML	OpenVZ (www.openvz.org) Apache Httpd Web Se	2
7,00		18 Pangeo Cloud	Python	xarray, iris, Dask, Jupyter	1
5,33		6 MeerKat Toolbelt	A Python 3.8+ (for CASA 6.5+)	CASA Singularity container for parallel processing	1

## Want/must have:

- Fully maintained & well documented (easy to find install, details etc)

- Source code available / open-source

- Rely on the summary of vision doc for criteria taken from vision doc

- System to be able to interact with remote content

## Desired platforms that have:

Modularity workflows and notebooks

## Cross off list the platforms that have:

- Features/modules that can't be developed together OR not useful for us(eg. the platform is intended to be running on laptop)

## Matrix Evaluation

# Assessing Architecture of SP



- **General**

- E.g. Does the platform use virtual machines? Does the platform use containerisation to provide functionality to the user, if so what tool is used (Docker/Singularity)?

- **Federation**

- Is the platform distributed between several sites? Are users registered at individual sites or are they registered to use resources at all sites?

- **Computational Layer**

- How does the computational layer connect to the server layer & presentation layer? What is the underlying infrastructure (HPC, public cloud, private cloud, hybrid)?

- **Presentation Layer**

- Does the presentation layer have direct access to the data layer? Does the presentation layer make use of some kind of pattern (e.g. Gateway pattern) to allow extensibility?

- **User Customization**

- What support is provided for user-supplied software and notebooks? If notebooks are used, can the users add their own dependencies?

- **Server Layer**

- How can the platform be scaled to allow for an increase in the number of users and/or data? Does the platform use orchestration? If so, what tool does it use (e.g. Kubernetes)?

# Evaluation Criteria (under Progress)



- Interfacing with data archive
- Data Quality Control
- Management of User Accounts
- Management of Compute Resources
- Visualization
- Tabular data
- Software access and data processing tools
- User support

Q. What do we want to get out of the assessment?

- To find the perfect platform
- To find components to build upon/from
- To gather ideas
- ... other

Q. What headers for the matrix (columns)

- High level summary of tasks (bullet points = +1)
- Maybe extra points for fancy stuff

Q. Do we need to score all the use cases for the existing platforms, for example the use case 9a: Creating a user account. Is this useful for our assessment?

Marks to be given out of 10, for example:

0 = No, not possible

5 = Yes, exactly what was wanted but with moderate difficulty /  
Can do something similar with ease

10 = Yes, exactly what was wanted and can be done easily

# Summary



- Defined high level essential criteria of SKA science platform
- Downlisted initial science platforms list
- Accessed some science desirable science platforms
- Architecture assessment under progress
- Future: Deployment and testing