



REGIONAL
CENTRE
NETWORK

SKA Regional Centres

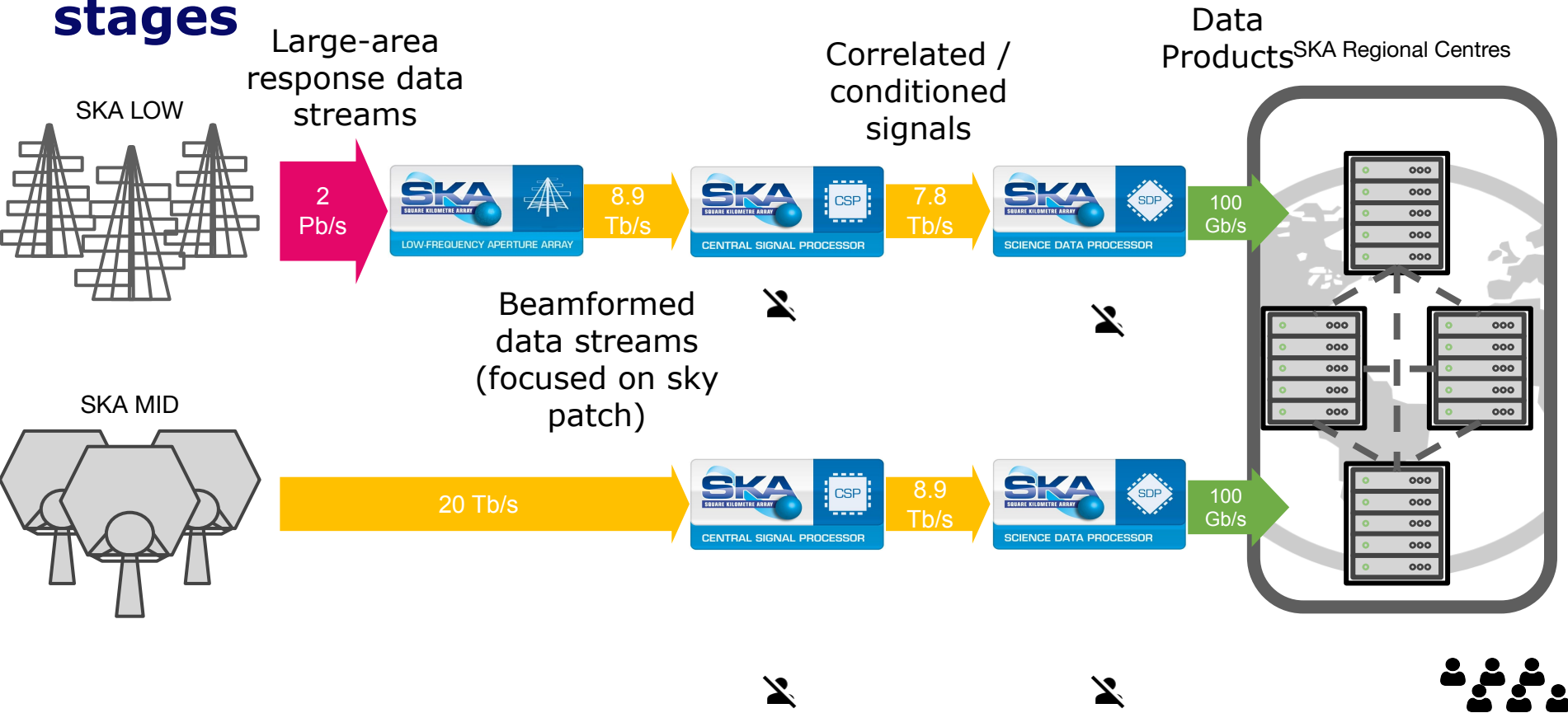
Dr Rosie Bolton

Head of Data Operations group, SKAO

27/09/2022



SKA Regional Centres: SKAO data processing stages



The SRC Vision

"To ensure that scientists can access SKA data products and use them to make discoveries"



What are we trying to achieve?

Globally distributed science teams
Diverse skill base - not all expert radio astronomers

Know who user is
Understand group membership
Respect proprietary periods
Support public access

"To ensure that scientists can access SKA data products and use them to make discoveries"

get data products from Observatory
allow inclusion of new data products
protect data products - ensure long term reliability
know where copies are
users and data products will be globally distributed
access depends on credentials not geography
meet SKA Construction timeline

data products analysed on SRC compute resources
allocation and tracking of resources to user groups
interactive sessions
user-defined software
batch workflows
user-defined data collections
create and save new data products and software
...and many more uses we need to keep talking with community about



What functions do SRCs need to provide?



The Role of SRCs: Data Intensity vs. User Flexibility

SRCs will bridge the gap between the highly data intensive **pre-defined workflows** generating **SKA data products** in the SDP, and the **iterative flexible, user-led data analysis** required to produce scientific results

"Size" or Data intensity

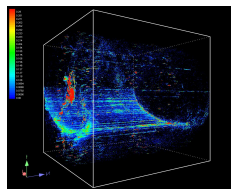


Image cubes

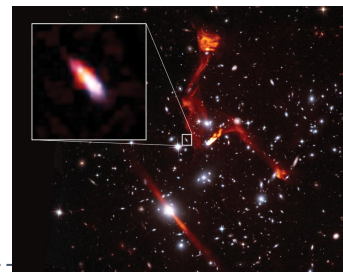
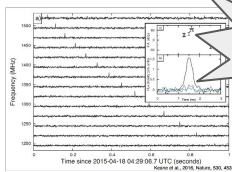


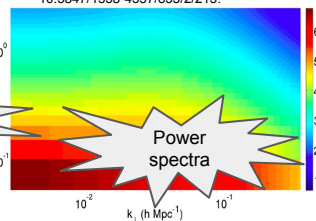
Image cut-outs



Time series data

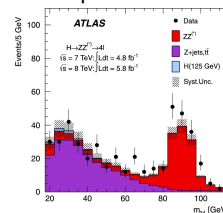
Catalogues / Source List

Credit: Heywood et al.; Sophia Dagnello, NRAO/AUI/NSF; STScI.
Paul, Sourabh et al. (2016). ApJ. 833.
10.3847/1538-4357/833/2/213.



Power spectra

Plots and scripts for publication



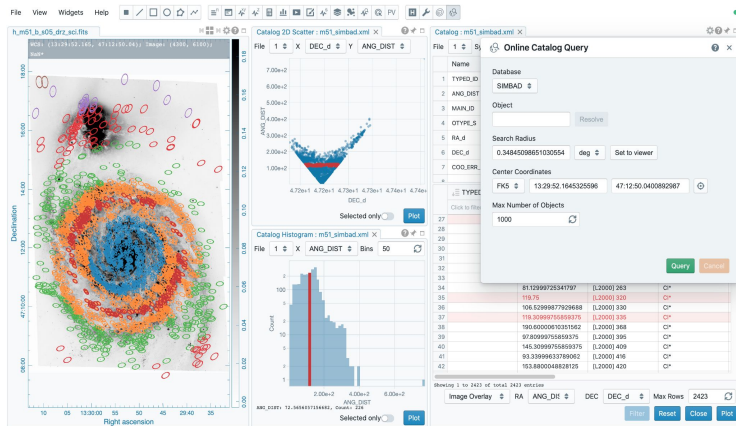
"Pre-defined" (Observatory Data Products)

"User-defined" (Advanced Data Products)



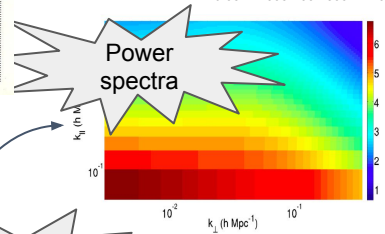
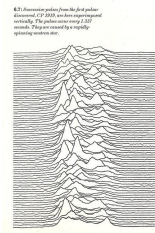
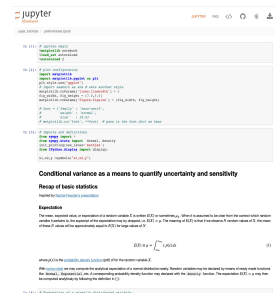
The Role of SRCs: Collaboration platform

SRCs will provide collaborative tools backed up by powerful compute and data management supporting **wide range of use cases**



Users will not have access to the SDP or to Raw SKA data!

Workflows notebooks



Paul, Sourabh et al. (2016), ApJ 833, 10.3847/1538-4357/833/2/13.

Catalogues / Source List

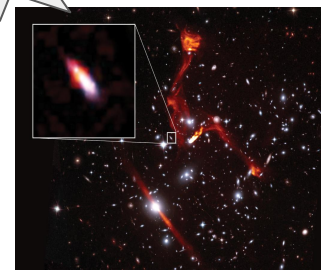
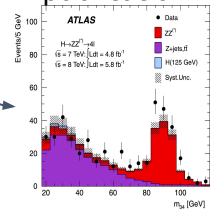


Image cut-outs

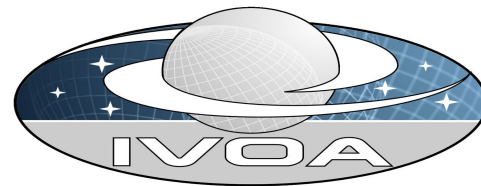
Credit: Heywood et al., Sophia Dagnello, NRAO/AUI/NSF; STScI.

Plots for publication

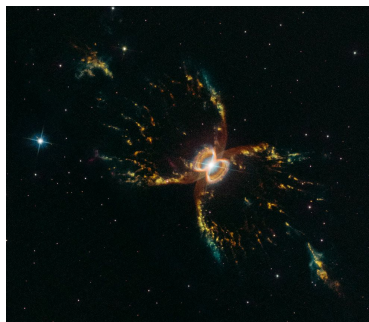


The Role of SRCs: Support data product (re-)use

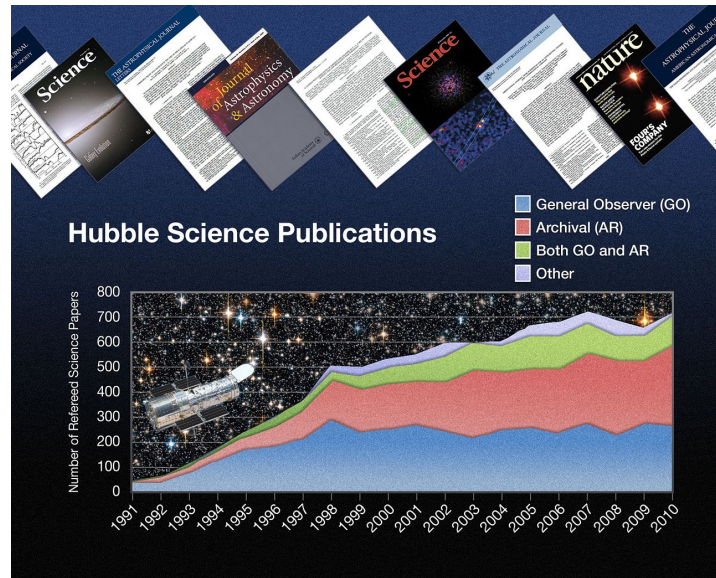
- All SKA Data Products will (in time) become public - this is likely to be the biggest science generator long term
- Build SKA science archive around International Virtual Observatory Alliance standards
- Ensure interoperability with other archives and other experiments
- Example from Hubble: Science archive users produce more scientific publications annually than users with dedicated observations



*FAIR principles
(Find, Access, Interoperate, Re-use)*



Southern Crab Nebula imaged by WFC3 - Hubblesite.org



SRC Users

- SKA Programme Users
 - With planned SKA projects
 - User has account in SRCs
 - Estimate SRC workflows and resources in SKA proposal
 - SRC workflows include generating Project-Level Data Products (SKAO responsible for SW for these)
 - Also analysis of OLDPs to produce Advanced Data Products
 - Resources allocated to the collaboration (SKA Project number)
- SRC Users: Public SKA data
 - Requested resources to run big analysis project on SRCNet
- Public users
 - "low" resource requirements, possibly anonymous access to SRC data collections



SRC Capabilities



SKA Regional Centers: Data management

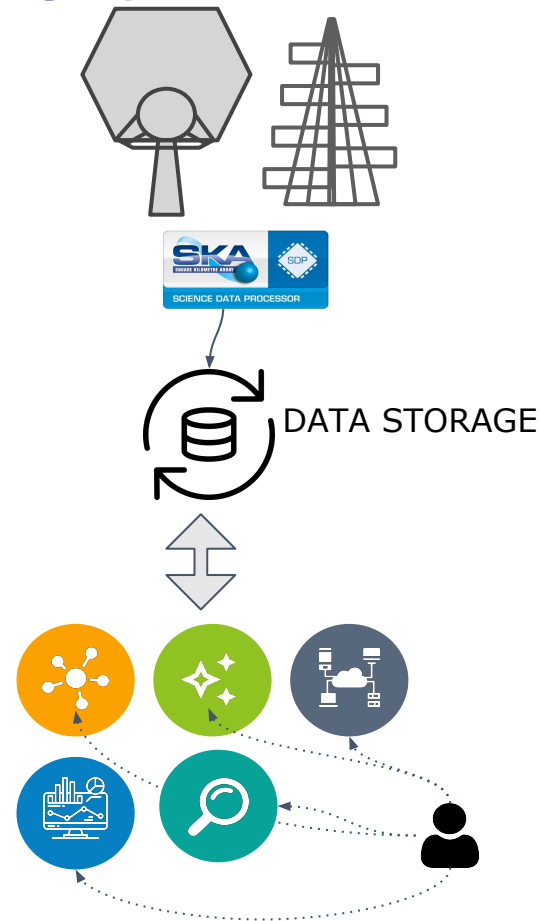
Storing SKAO data growing at up to 700 PBytes each year will be a challenge (plus user-generated data too).

Several million dollars per year in new data, for one copy

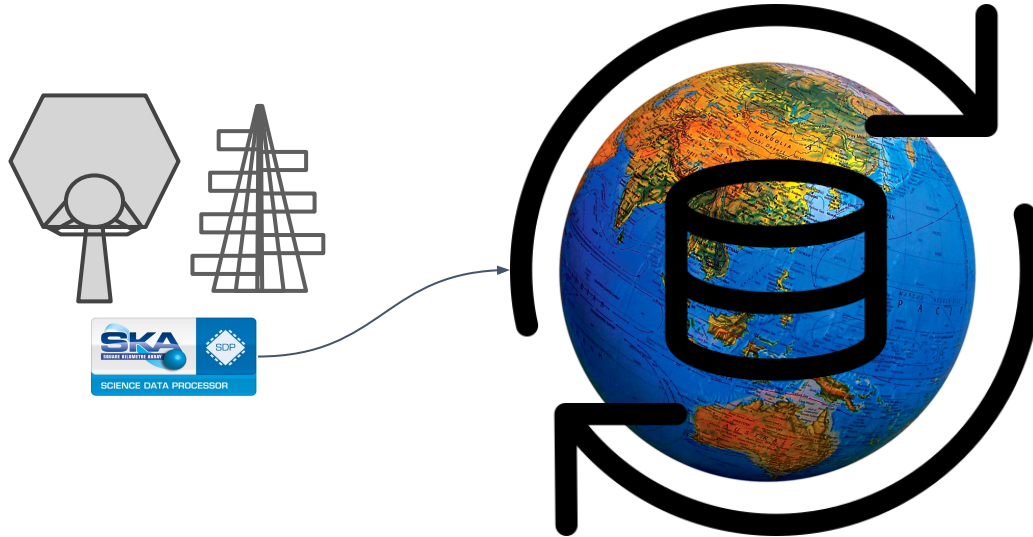
Global data management within SRCNet should enable best possible use to be made of available storage resources

Avoid (reduce) unnecessary duplication

Support mirroring of popular data products to enhance user experience



SKA Regional Centers: Important for Operation of SKA



Pushing data out of SDP staging areas into SRC storage is essential to keep the SKA telescopes running

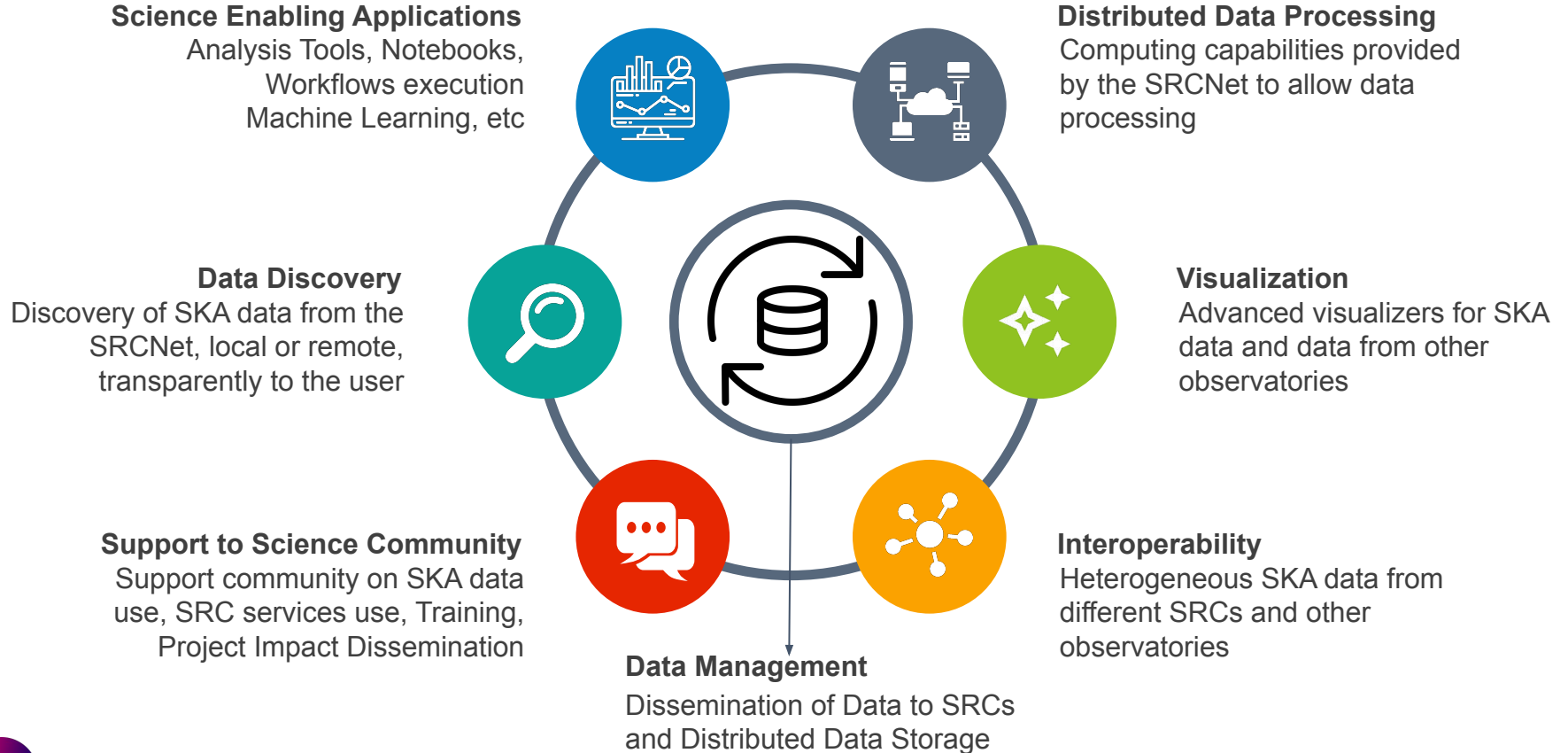
- If staging area gets full, SKA operations will be impacted

We need to be able to predict link usage

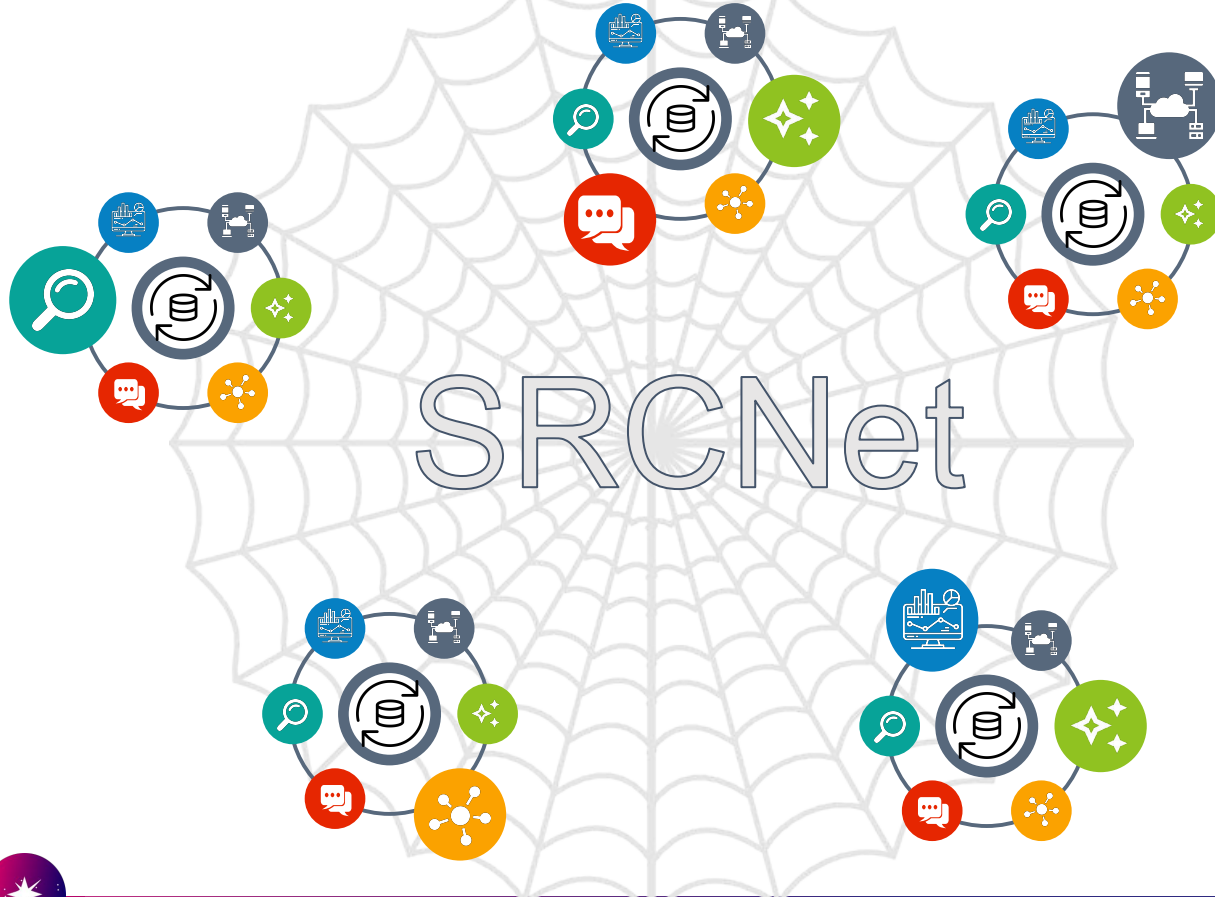
- Which data products are sent where?
- How reliable are links?
- How many copies? From SDP to each SRC or from SDP to one SRC then less time-critical copying



SKA Regional Centre Capabilities



SRC Network global capabilities



Collectively meet the needs of the global community of SKA users

Anticipate heterogeneous SRCs, with different strengths



How might SRC resources be managed?



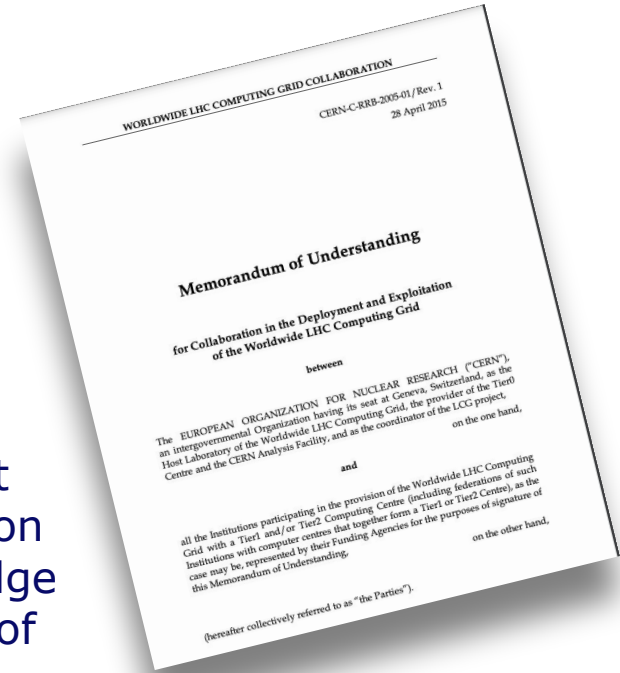
How might resources be obtained?

We hope that each SRC project will pledge resources into the SRCNet pool

MoU to cover deployment and use of SRCs

Public pledging system to gather resource offers, rolling (e.g. look ahead 2-5 years, revise each year)

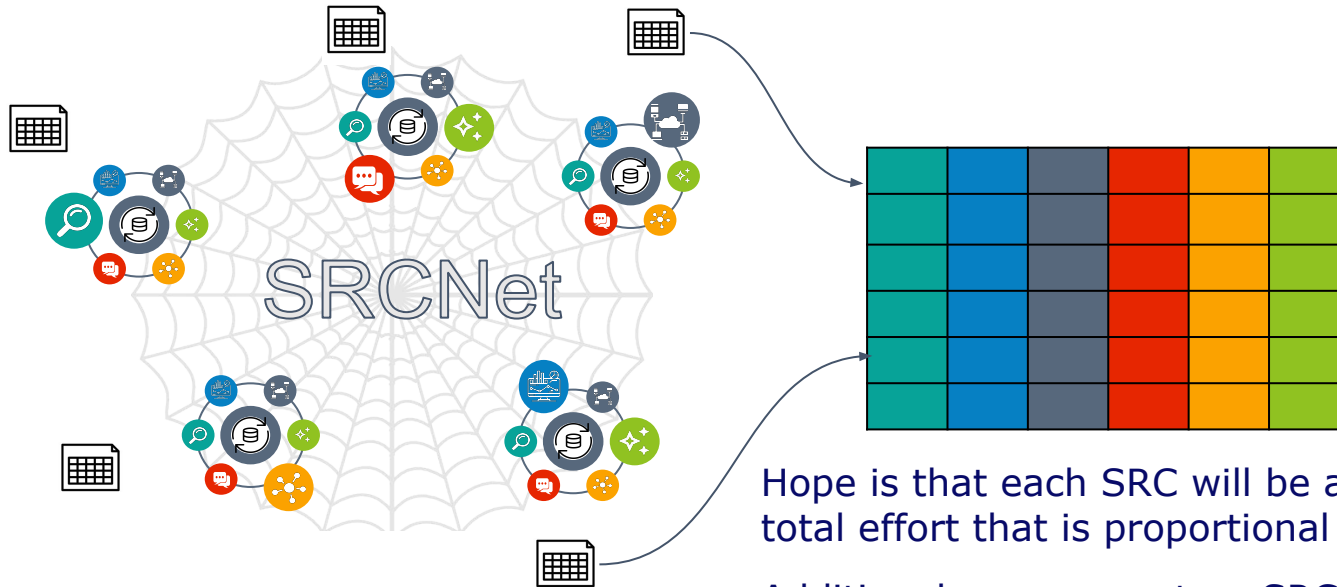
- Programme users should scale with SKA commitment
 - SRC pledges hope to match or exceed this fraction
 - e.g. 5% stakeholder in SKA, we hope would pledge resources at least (or exceeding) estimated 5% of the total required SRCNet resources to support users described previously.



Pledging

Each SRC to pledge resources into global pool to support SRCNet activities

Users can access resources across SRCNet according to their research needs and permissions



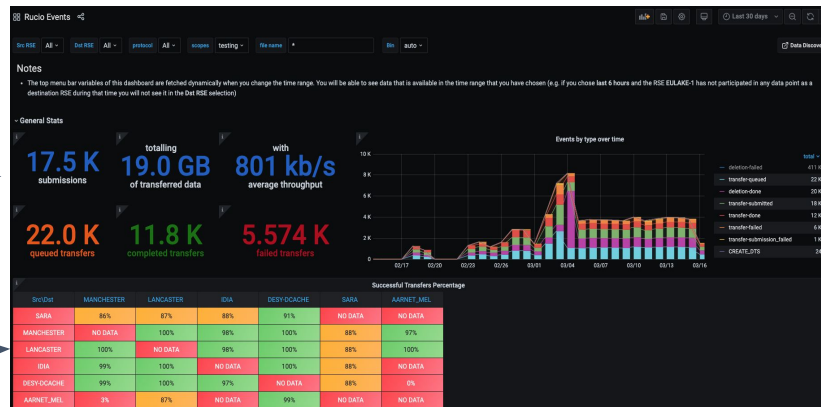
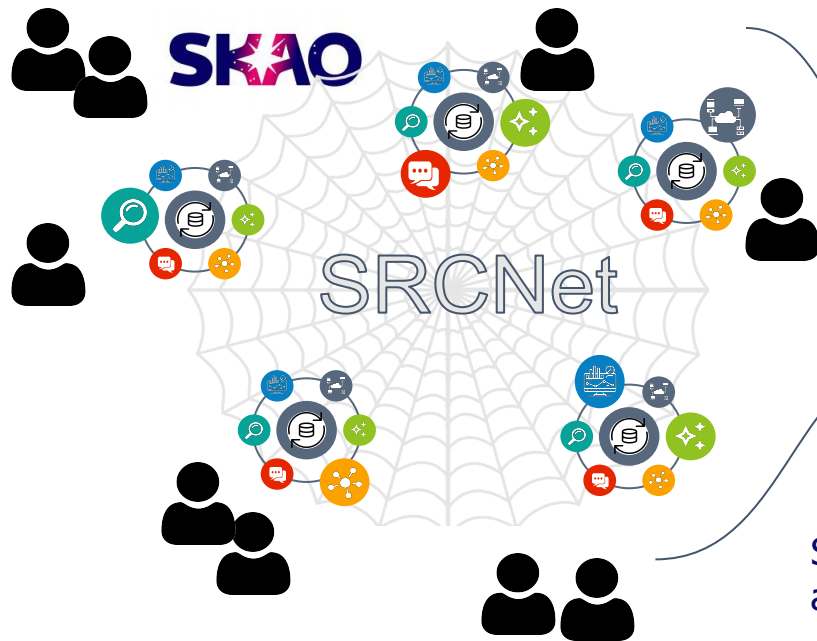
Hope is that each SRC will be able to contribute a total effort that is proportional to their SKA fraction

Additional resources at an SRC could be given to the pool or prioritised to support national interests



Operations

Personnel within each SRC project will be identified to be part of the SRC Operations Group (SOG) - meeting regularly to discuss issues, share tasks, see and test global system health



(an example dashboard from our data management prototype, details not important, but nice to see that we are using UK grid storage endpoints in our Rucio prototype which is itself run off **IRIS resources at STFC cloud**)

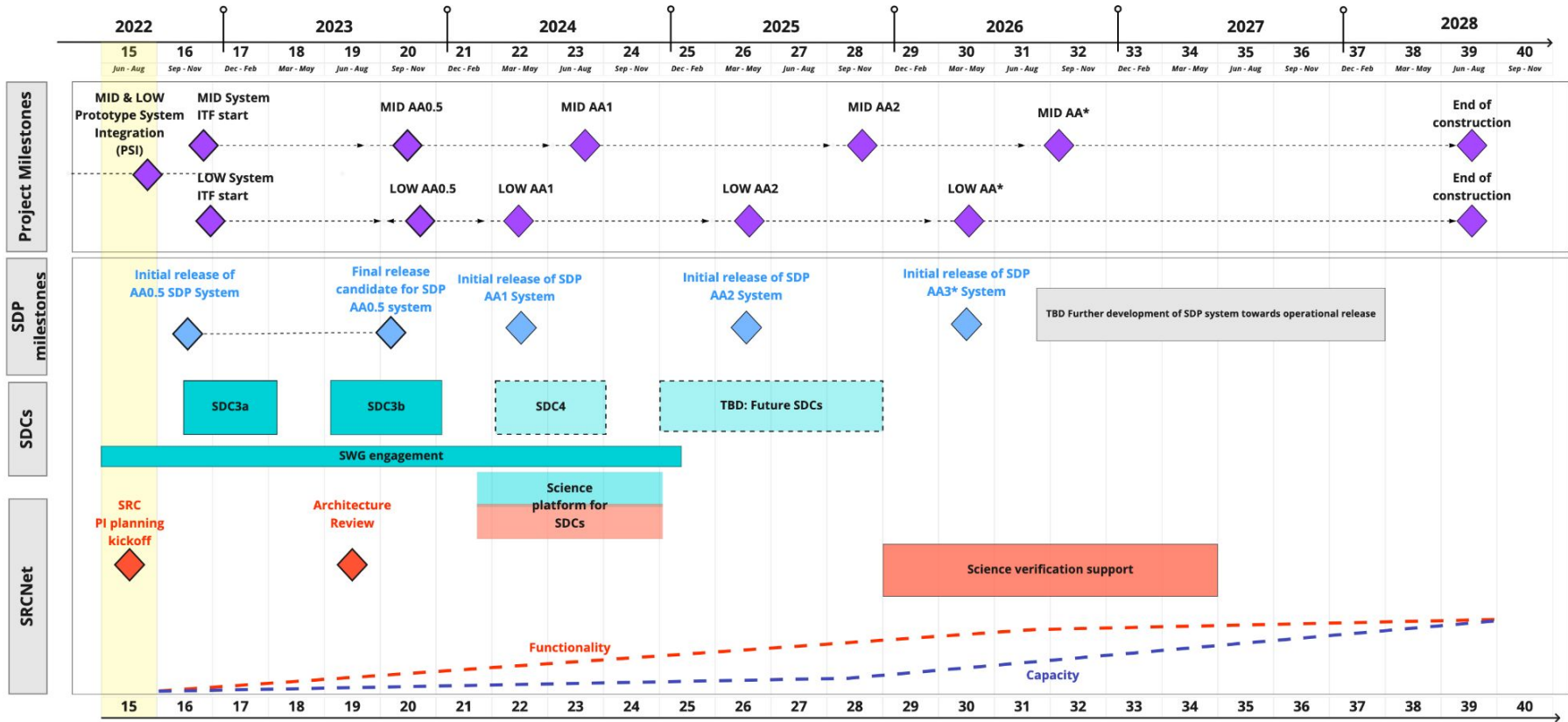
SOG will be led from SKAO Ops, with a team from across each SRC project and SKAO.



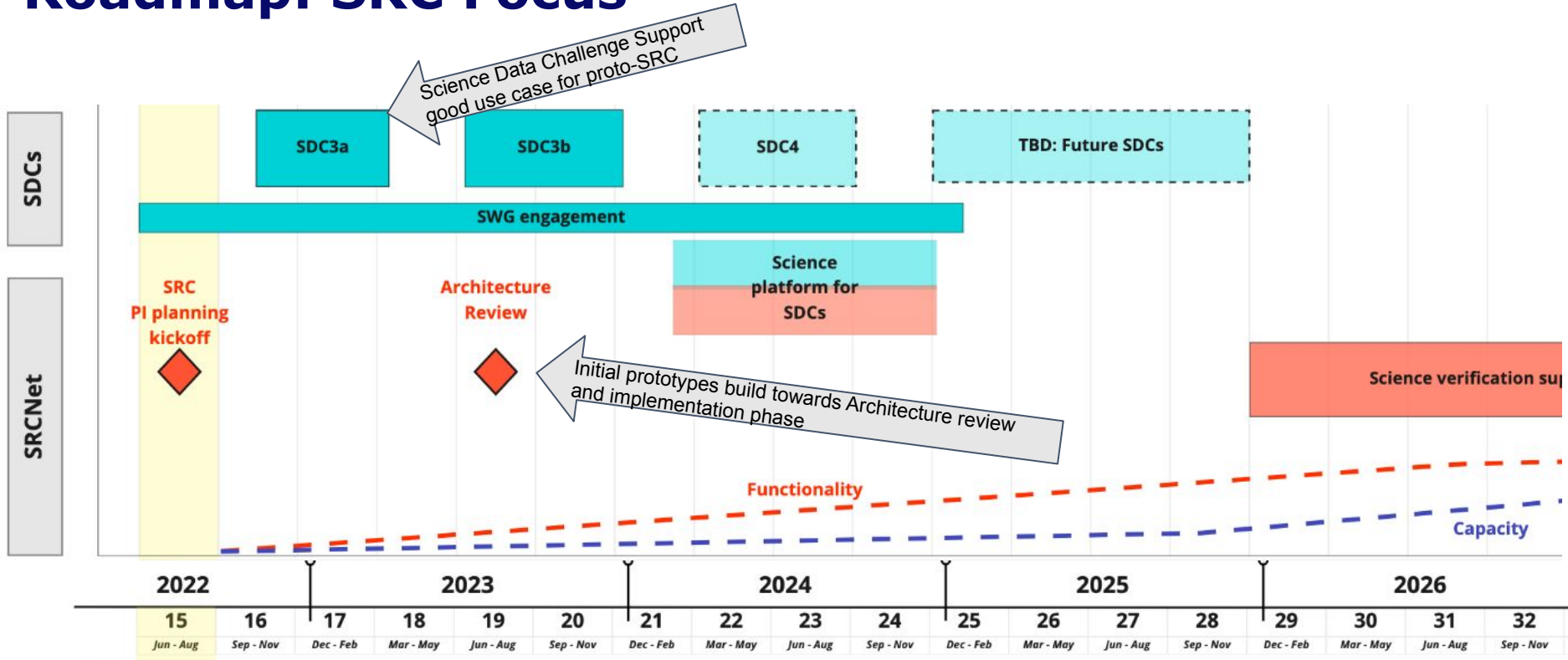
Development plans and progress



Roadmap



Roadmap: SRC Focus



First Prototyping Phase: 2022-2023

Work *now being undertaken* by development teams to prototype key technologies that will enable selection as SRC functionality and scale grows.

Data Management service:

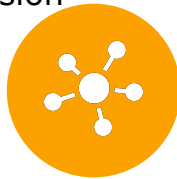
Replication, distribution, synchronisation of data products and location index



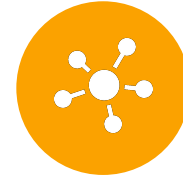
Teams working together with alignment driven by the common roadmap

Central Services and Software Distribution:

SW infrastructure, compute provision



Federated **Authentication and Authorization**: identity management, compatible with SKAO



Data Visualisation and discovery - performance at SKA scale



Data Analysis: Science Extraction, Processing in Notebooks



SKA Regional Centers: Data management

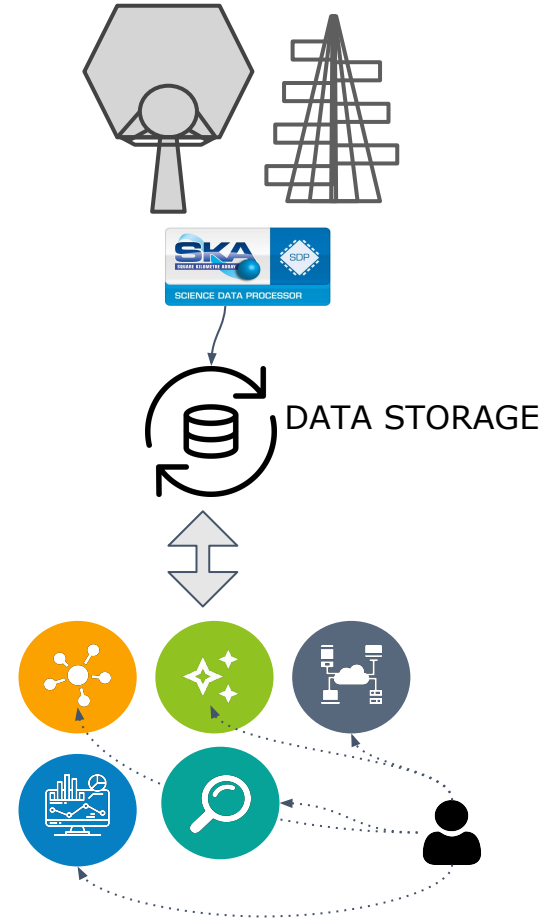
Storing SKAO data growing at up to 700 PBytes each year will be a challenge (plus user-generated data too).

Several million dollars per year in new data, for one copy

Global data management within SRCNet should enable best possible use to be made of available storage resources

Avoid (reduce) unnecessary duplication

Support mirroring of popular data products to enhance user experience



SRC Rucio prototype

Hope for Chinese site by December 2022!



Well suited to centralised Operations model for data management

Performed long-haul transfers, Rucio stress tests, subscriptions (via our automated test framework)

Integrating storage from national SRC efforts to increase understanding and inform assessment



Progress - Identity management (Authentication and Authorisation)

- "IAM"* service deployed in UK (Rutherford Appleton Lab)
<https://ska-iam.stfc.ac.uk/login>
- Landscape report created ready for wider sharing and feedback

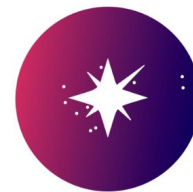
*<https://indigo-iam.github.io/docs/v/current/about.html>

Executive Summary

Science has large scale research activities that require pooling global computing and data storage resources, such as the [Large Hadron Collider](#) (LHC), [climate research](#), and [NGS](#). Naturally no one wants to register with every remote service with a new password every time: a method is needed to connect the user's existing identities to resources across the world, possibly in different trust domains.

Authentication and authorisation infrastructure federations provide the means to connect users to resources. NRENs have operated (usually national) identity federations for decades and these are interconnected through activities like eduGAIN, but (usually) only cover academic member organisations.

...



Welcome to **SKA IAM Prototype**

Sign in with your SKA IAM Prototype credentials

[Forgot your password?](#)

Or sign in with

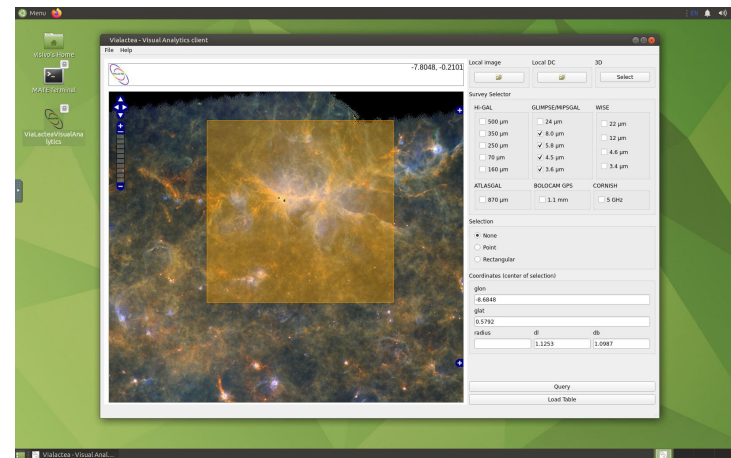
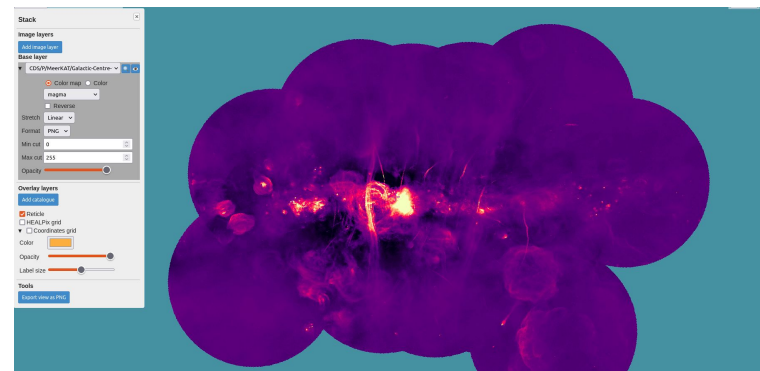
Not a member?

[Privacy Policy](#)



Progress - Visualisation

- Comparative review of three tools where teams have expertise: CARTA, ALADIN and visIVO
- Collection of data products (with links) for use in demos
- Demo of these tools to the SRCNet team of teams
- Containerised CARTA deployment documented
- All three tools deployed at CHINA SRC, and access given to external users (user pass, IP address whitelist)



Progress - Science Platform

- Evaluation of (30) existing platforms considering many different criteria

Shortlist of 8 (or so) for further work

- Science Platform vision document ready for review

Next steps are to consider architectural aspects of science platform vision, update vision to incorporate this.

Platform ID	Science Platform	Evaluation by	Implementation (language)	Major Frameworks and technologies	Open Source	Production status	Multitenant status	Data discovery	Notebooks	Workflows
SP1	QGIS SWAN	@Bibeki, CHN; @SabinasPabon	Python, Javascript	jQuery, Bootstrap, Leaflet, OpenLayers	<ul style="list-style-type: none"> Open Source - The platform is fully developed in the open, with public source repositories, issue tracking, and/or other development systems. Open Source - A version of the platform is available under an OSI approved license. Open Source - The source code or other artifacts are made available, but the project is development is not open without substantial permission to contributors. Open Source - The source code is not available. 	<ul style="list-style-type: none"> Production - The platform is providing services to a substantial user community. Production - The platform is available, which is broadly feature complete, but does not represent a substantial user community. Production - A development version of the platform is available, but it is not yet feature complete and/or insufficiently stable to be taken to production. Production - The platform is not yet reached a minimum viable form. 	<ul style="list-style-type: none"> Production - The platform developers are active in providing support to users and answering questions, being blogs, or other means of communication. Production - The platform are regularly made available. Production - The platform is available in a public or private code repository. They report to their code repository. Production - The platform is available in a public or private code repository. They report to their code repository. Production - The platform is available in a public or private code repository. They report to their code repository. 	<ul style="list-style-type: none"> Production - The platform provides a mechanism for searching and returning lists of (images, source code, etc.) Production - The platform provides a mechanism for searching and returning lists of (images, source code, etc.) Production - The platform provides a mechanism for searching and returning lists of (images, source code, etc.) Production - The platform provides a mechanism for searching and returning lists of (images, source code, etc.) 	<ul style="list-style-type: none"> Production - The platform provides access to Jupyter notebooks or equivalent functionality. Production - The platform provides access to Jupyter notebooks or equivalent functionality. Production - The platform provides access to Jupyter notebooks or equivalent functionality. Production - The platform provides access to Jupyter notebooks or equivalent functionality. 	<ul style="list-style-type: none"> Production - The platform provides functionality for defining and executing processing pipelines on workflow. Production - The platform provides functionality for defining and executing processing pipelines on workflow. Production - The platform provides functionality for defining and executing processing pipelines on workflow. Production - The platform provides functionality for defining and executing processing pipelines on workflow.
SP2	SciServer	@Bibeki, CHN; @KADNET	Java/JSP		<ul style="list-style-type: none"> Open Source - Source code available on GitHub at https://github.com/SciServer 	<ul style="list-style-type: none"> Production - The platform is providing services to a substantial user community. Production - The platform is available, which is broadly feature complete, but does not represent a substantial user community. Production - A development version of the platform is available, but it is not yet feature complete and/or insufficiently stable to be taken to production. Production - The platform is not yet reached a minimum viable form. 	<ul style="list-style-type: none"> Production - The platform developers are active in providing support to users and answering questions, being blogs, or other means of communication. Production - The platform are regularly made available. Production - The platform is available in a public or private code repository. They report to their code repository. Production - The platform is available in a public or private code repository. They report to their code repository. 	<ul style="list-style-type: none"> Production - The platform provides a mechanism for searching and returning lists of (images, source code, etc.) Production - The platform provides a mechanism for searching and returning lists of (images, source code, etc.) Production - The platform provides a mechanism for searching and returning lists of (images, source code, etc.) Production - The platform provides a mechanism for searching and returning lists of (images, source code, etc.) 	<ul style="list-style-type: none"> Production - The platform provides access to Jupyter notebooks or equivalent functionality. Production - The platform provides access to Jupyter notebooks or equivalent functionality. Production - The platform provides access to Jupyter notebooks or equivalent functionality. Production - The platform provides access to Jupyter notebooks or equivalent functionality. 	<ul style="list-style-type: none"> Production - The platform provides functionality for defining and executing processing pipelines on workflow. Production - The platform provides functionality for defining and executing processing pipelines on workflow. Production - The platform provides functionality for defining and executing processing pipelines on workflow. Production - The platform provides functionality for defining and executing processing pipelines on workflow.
SP3	Robin Science Platform	@Bibeki, CHN; @SabinasPabon	Python, C++	libFEMTO, JupyterLab, Dash, FastAPI	<ul style="list-style-type: none"> Open Source - Source code available on GitHub at https://github.com/robin-science 	<ul style="list-style-type: none"> Production - The platform is providing services to a substantial user community. Production - The platform is available, which is broadly feature complete, but does not represent a substantial user community. Production - A development version of the platform is available, but it is not yet feature complete and/or insufficiently stable to be taken to production. Production - The platform is not yet reached a minimum viable form. 	<ul style="list-style-type: none"> Production - The platform developers are active in providing support to users and answering questions, being blogs, or other means of communication. Production - The platform are regularly made available. Production - The platform is available in a public or private code repository. They report to their code repository. Production - The platform is available in a public or private code repository. They report to their code repository. 	<ul style="list-style-type: none"> Production - The platform provides a mechanism for searching and returning lists of (images, source code, etc.) Production - The platform provides a mechanism for searching and returning lists of (images, source code, etc.) Production - The platform provides a mechanism for searching and returning lists of (images, source code, etc.) Production - The platform provides a mechanism for searching and returning lists of (images, source code, etc.) 	<ul style="list-style-type: none"> Production - The platform provides access to JupyterLab notebooks. Production - The platform provides access to JupyterLab notebooks. Production - The platform provides access to JupyterLab notebooks. Production - The platform provides access to JupyterLab notebooks. 	<ul style="list-style-type: none"> Production - The platform provides functionality for defining and executing processing pipelines on workflow. Production - The platform provides functionality for defining and executing processing pipelines on workflow. Production - The platform provides functionality for defining and executing processing pipelines on workflow. Production - The platform provides functionality for defining and executing processing pipelines on workflow.
SP4	SCAPE CASP	@SabinasPabon	Python, Javascript	Django, ReactJS	<ul style="list-style-type: none"> Open Source - Source code available on GitHub at https://github.com/SCAPE-CASP 	<ul style="list-style-type: none"> Production - The platform is providing services to a substantial user community. Production - The platform is available, which is broadly feature complete, but does not represent a substantial user community. Production - A development version of the platform is available, but it is not yet feature complete and/or insufficiently stable to be taken to production. Production - The platform is not yet reached a minimum viable form. 	<ul style="list-style-type: none"> Production - The platform developers are active in providing support to users and answering questions, being blogs, or other means of communication. Production - The platform are regularly made available. Production - The platform is available in a public or private code repository. They report to their code repository. Production - The platform is available in a public or private code repository. They report to their code repository. 	<ul style="list-style-type: none"> Production - The platform provides a mechanism for searching and returning lists of (images, source code, etc.) Production - The platform provides a mechanism for searching and returning lists of (images, source code, etc.) Production - The platform provides a mechanism for searching and returning lists of (images, source code, etc.) Production - The platform provides a mechanism for searching and returning lists of (images, source code, etc.) 	<ul style="list-style-type: none"> Production - The platform provides access to JupyterLab notebooks. Production - The platform provides access to JupyterLab notebooks. Production - The platform provides access to JupyterLab notebooks. Production - The platform provides access to JupyterLab notebooks. 	<ul style="list-style-type: none"> Production - The platform provides functionality for defining and executing processing pipelines on workflow. Production - The platform provides functionality for defining and executing processing pipelines on workflow. Production - The platform provides functionality for defining and executing processing pipelines on workflow. Production - The platform provides functionality for defining and executing processing pipelines on workflow.
SP5	Cytosca	@SabinasPabon; @SabinasPabon	Go, C++, Javascript, Python	gRPC, gRPC-Web, Apache JAX, ReactJS, Material UI	<ul style="list-style-type: none"> Open Source - Source code available on GitHub at https://github.com/cytosca 	<ul style="list-style-type: none"> Production - The platform is providing services to a substantial user community. Production - The platform is available, which is broadly feature complete, but does not represent a substantial user community. Production - A development version of the platform is available, but it is not yet feature complete and/or insufficiently stable to be taken to production. Production - The platform is not yet reached a minimum viable form. 	<ul style="list-style-type: none"> Production - The platform developers are active in providing support to users and answering questions, being blogs, or other means of communication. Production - The platform are regularly made available. Production - The platform is available in a public or private code repository. They report to their code repository. Production - The platform is available in a public or private code repository. They report to their code repository. 	<ul style="list-style-type: none"> Production - The platform provides a mechanism for searching and returning lists of (images, source code, etc.) Production - The platform provides a mechanism for searching and returning lists of (images, source code, etc.) Production - The platform provides a mechanism for searching and returning lists of (images, source code, etc.) Production - The platform provides a mechanism for searching and returning lists of (images, source code, etc.) 	<ul style="list-style-type: none"> Production - The platform provides access to JupyterLab notebooks. Production - The platform provides access to JupyterLab notebooks. Production - The platform provides access to JupyterLab notebooks. Production - The platform provides access to JupyterLab notebooks. 	<ul style="list-style-type: none"> Production - The platform provides functionality for defining and executing processing pipelines on workflow. Production - The platform provides functionality for defining and executing processing pipelines on workflow. Production - The platform provides functionality for defining and executing processing pipelines on workflow. Production - The platform provides functionality for defining and executing processing pipelines on workflow.
SP6	Material Frontend								<ul style="list-style-type: none"> Production - The platform provides access to JupyterLab notebooks. Production - The platform provides access to JupyterLab notebooks. Production - The platform provides access to JupyterLab notebooks. Production - The platform provides access to JupyterLab notebooks. 	<ul style="list-style-type: none"> Production - The platform provides functionality for defining and executing processing pipelines on workflow. Production - The platform provides functionality for defining and executing processing pipelines on workflow. Production - The platform provides functionality for defining and executing processing pipelines on workflow. Production - The platform provides functionality for defining and executing processing pipelines on workflow.

SR Science Analysis Platform Vision

Contents

1 Introduction and Context	3
1.1 Purpose and Status of this Document	3
1.2 What is a Science Analysis Platform?	3
1.3 Platform Aims and Objectives	3
1.4 The design of the platform	5
1.5 Accessibility and Inclusion	5
2 Back-end Features and Services	6
2.1 Compute services	6
2.2 Archive and distributed data	6
2.3 User file and database services	6
2.4 Authentication and Authorization	7
3 User-facing Features and Services	7
3.1 Main User Interface	7
3.2 Data Querying and Discovery	8
3.3 Notebook Interface	8
3.4 Software environments	9
3.5 Web APIs	9
3.6 Workflow Management	10
3.7 Resource Management	11
3.8 Groups	11
3.9 Software repository	11



Many layers, shared vision

Users



Science Platform



Interface (API layer)



Metadata query -
Science Data
Discovery



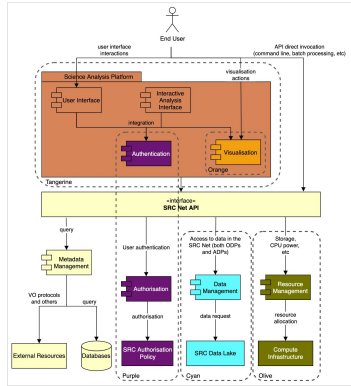
Authentication
Who?
Permissions?



Data Logistics
Globally
distributed
storage sites

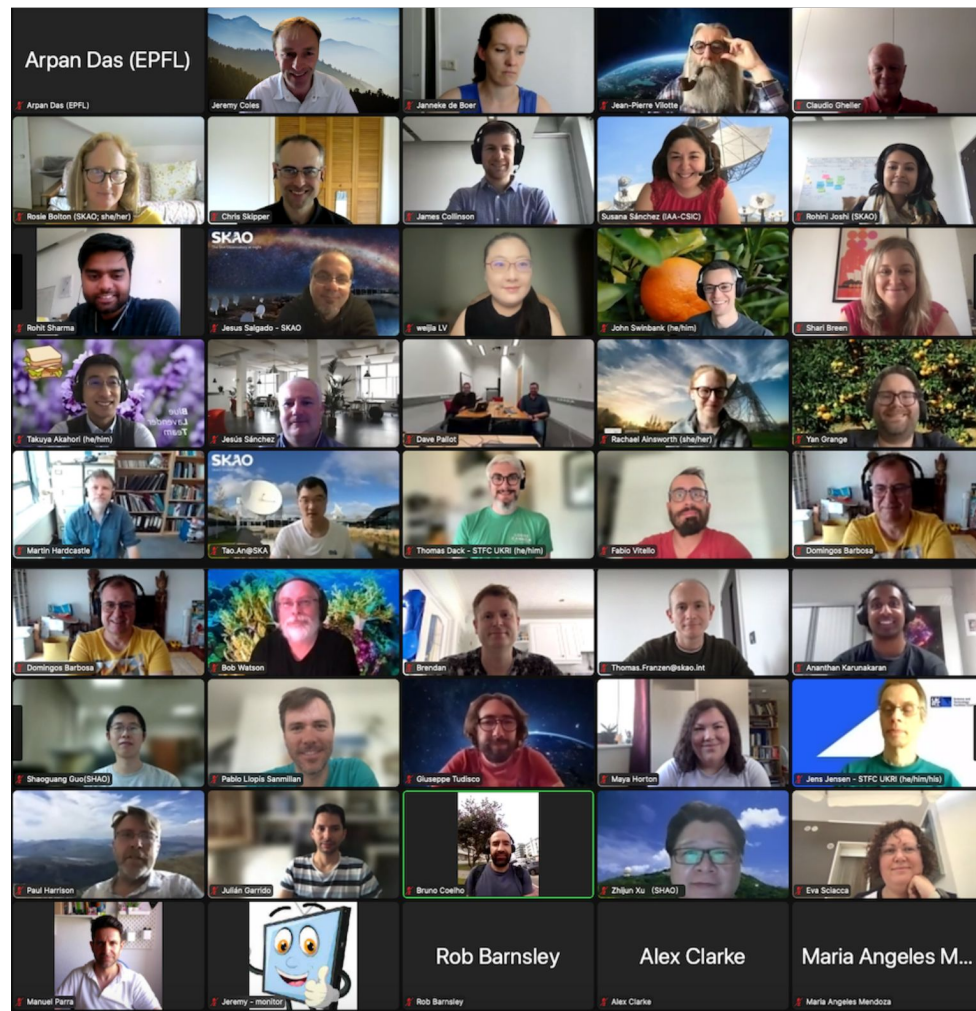


Compute
Resource
Management
Work sharing



SRC Prototyping

- Prototyping started June 2022
- Team members from 12 countries plus SKAO
- about 1000 developer-days per 3-months
- Anticipate growing as national funding to develop SRC nodes is available



Introducing the Blue-Lavender team

Excellent collaboration between **China, Japan, South Korea and Australia** team members

Already achieved:

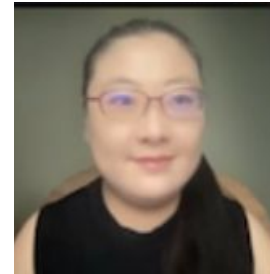
Deploying Visualisation tools at China SRC

Adding AUS-SRC storage into data management prototype

Next: Supporting other teams with visualisation workshop; Adding Chinese and Japanese SRC storage to prototype; defining needs for SDC3 support at China SRC



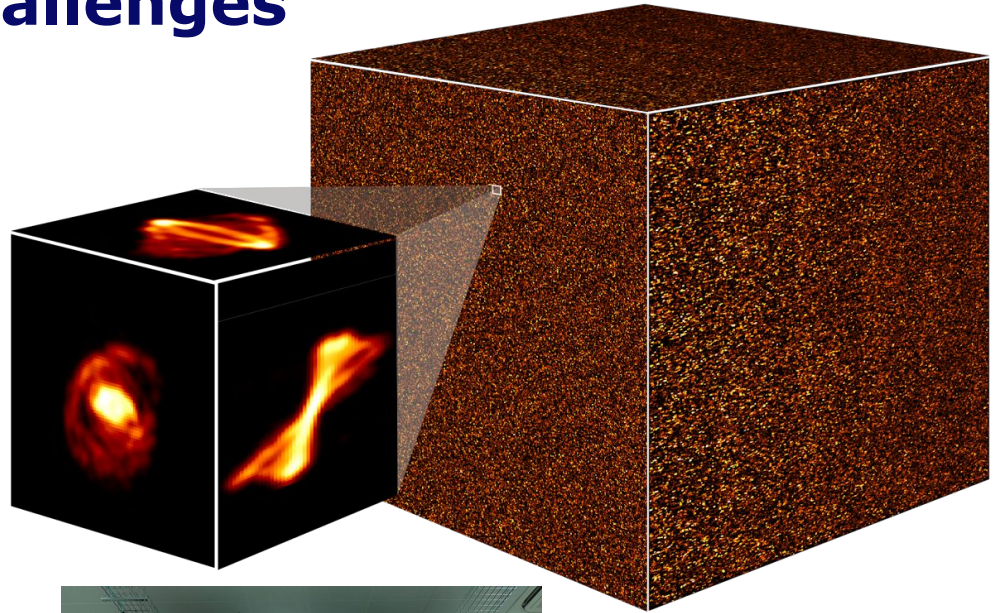
Selected team member images:
An, Tao
Takuya Akahori
Dave Pallot
Austin Shen
Gordon German
Kang Hyunwoo
Lv, Weijia
Guo, Shaoguang



Link to Science Data Challenges

Science Data Challenges are a great way to involve science community (e.g. research teams - hopefully some of you!) with simulated SKA data products

SRC projects offer a route to supporting challenge participants (including China SRC)



Credit: An, Wu, Hong, Nature Astronomy 3, 1030 (2019)



Summary

- SRCs are essential for connecting users with SKA data
- Collaborative model with Observatory - via MoU
- Pledged resources to provide global access
- Prototyping ideas for components now
- Implementation from end 2023, to keep pace with SKA Telescope development



End

