

one observatory
two telescopes
three continents

SKAO Science Data Challenges

Philippa Hartley
SKAO Scientist

SKA-China Workshop on SKA Science and Operations
Wednesday 28th September 2022



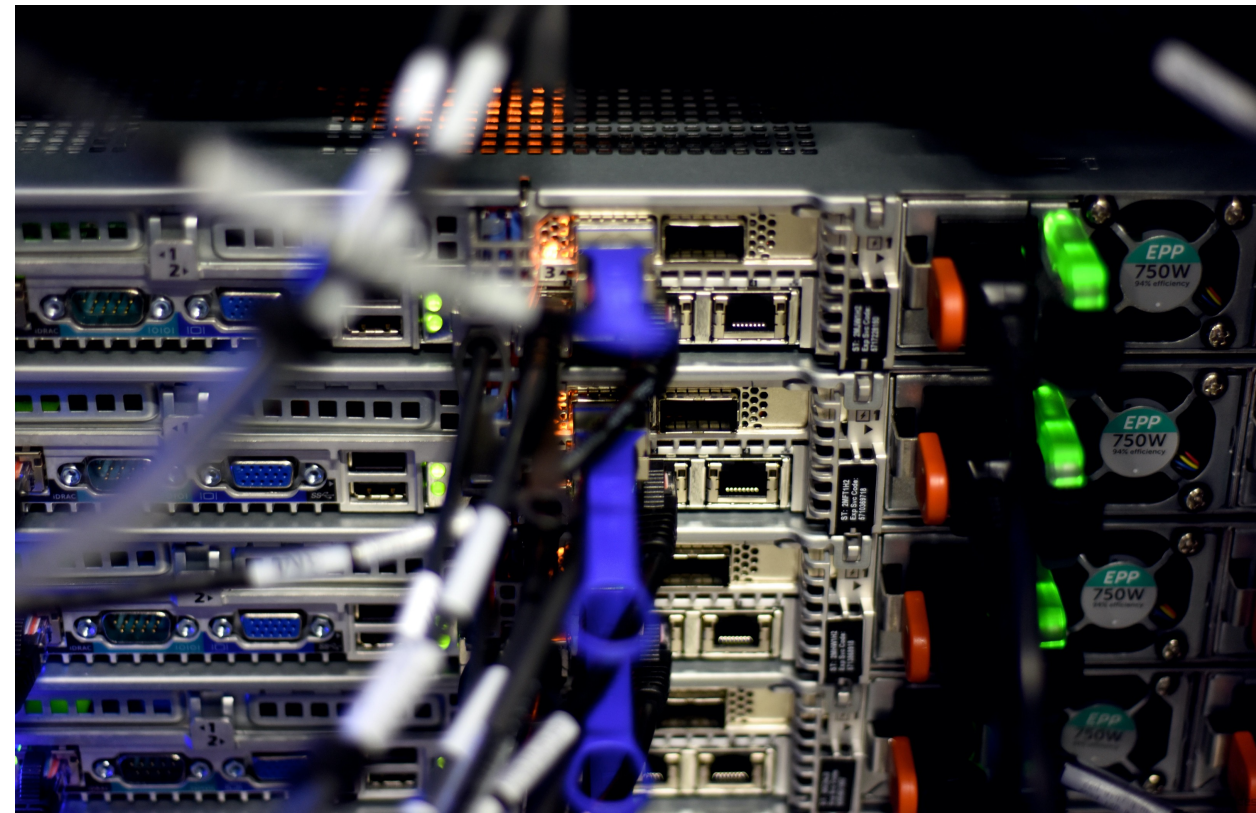
The SKAO data journey

SKA LOW



SKA MID

SDP: Science Data Processor



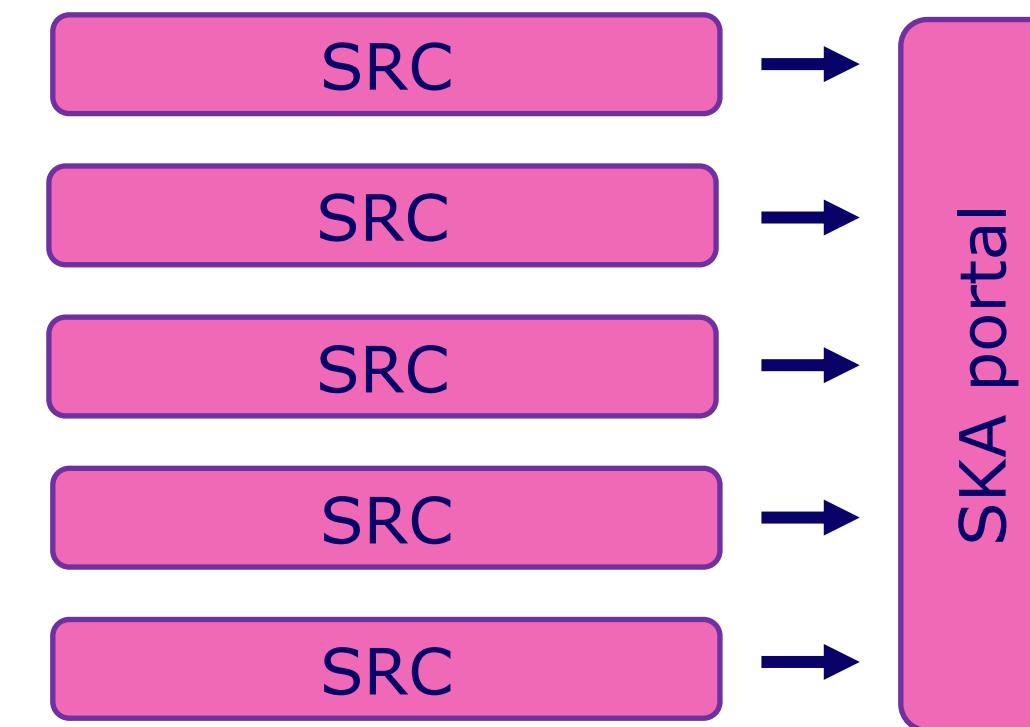
SDP prototype, Cambridge, UK



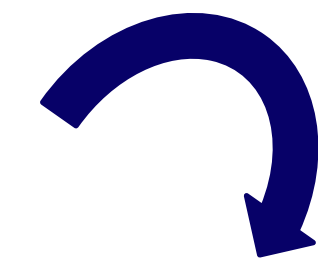
5 + 9 Tb/s
Approx 300
high definition
movies per
second!

600 PB/yr

SRC: SKA Regional Centre



Distributed facilities



User data
products up to
TBs in size

Key Science Projects



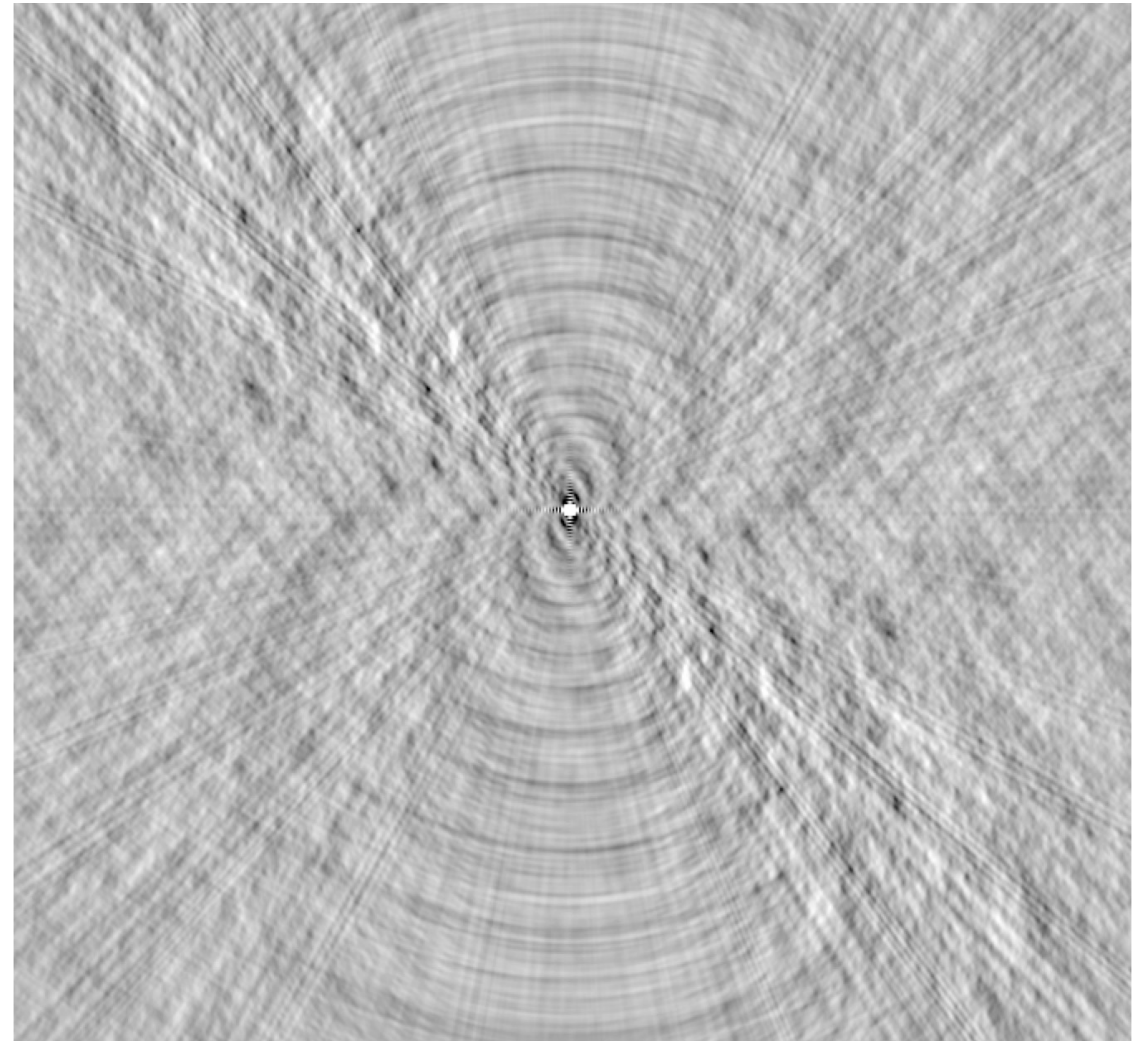
Scientific analysis



SKA telescope characteristics

SKA-unique features of the data products:

- In the **image plane**, not visibilities
- “**Benign**” dirty beam
- Deconvolved down to **8h exposures**
- Very deep -> **towards confusion limit**
- Very **large number** of sources to detect and classify



SKA MID 1.4 GHz beam



SKAO Science Data Challenges

Primary goals:

- Familiarise the science community with **size and complexity of SKA data**
- Support the **design** of future SKA observations
- Drive the development of **data analysis techniques**

Additional benefits:

- Familiarise the science community with **data access models**
- Test SKA Regional Centre **prototyping**
- Encourage best practices for **Open Science** and **reproducibility**

SDC data products are made publicly available for the long term



Science Data Challenge 1 (SDC1)

Continuum emission

Continuum science with the SKAO:

- explore our Universe from the nearest star-forming galaxies to high-redshift clusters of galaxies
- examine structure formation and evolution on all scales.

Extragalactic continuum Science Working Group:
<https://www.skao.int/en/science-users/science-working-groups-focus-groups/107/extragalactic-continuum>

Credits: ESO, NRAO

Centaurus A



SKAO

64 MeerKAT dishes

The MID telescope

15 m

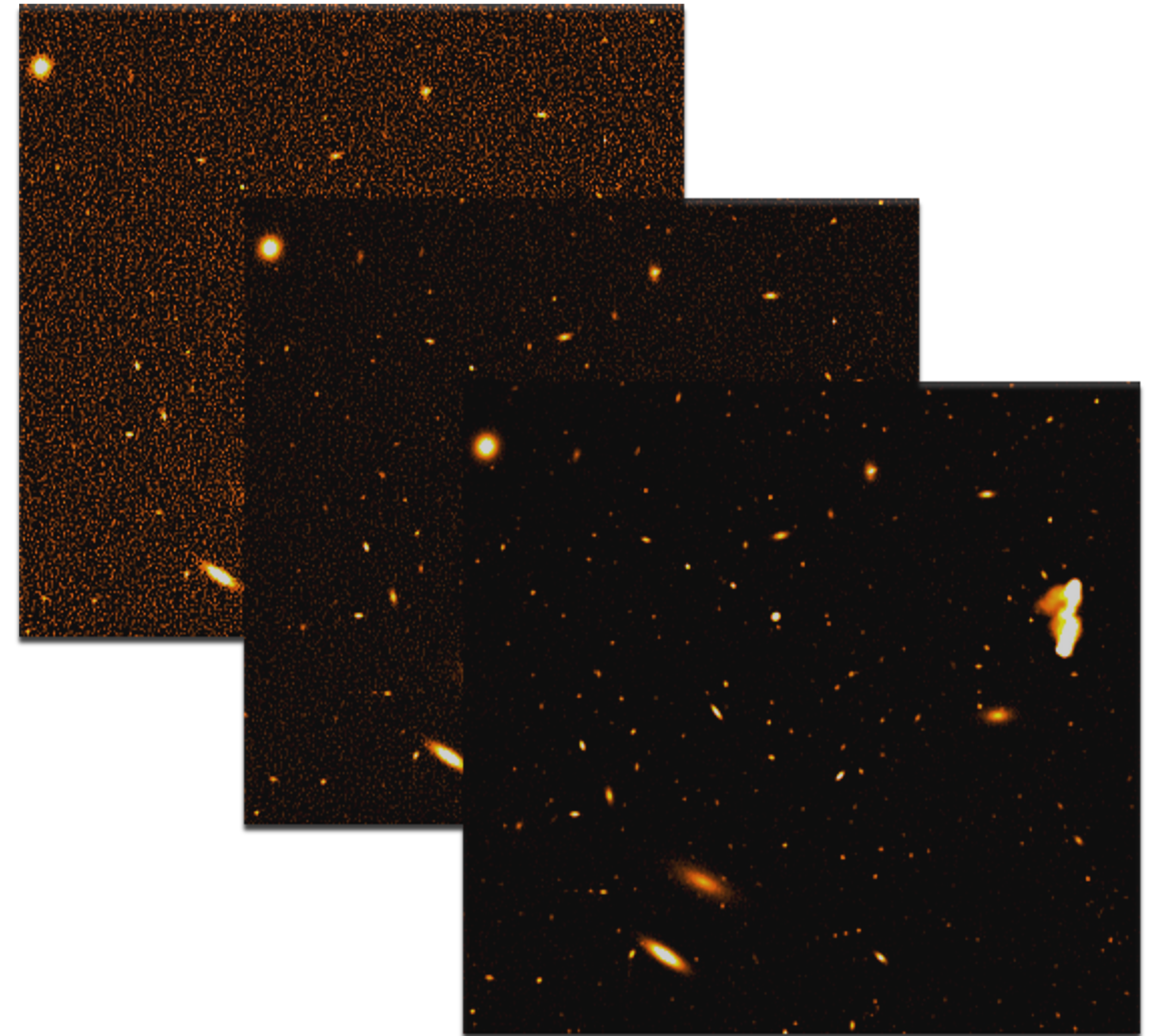
133 SKAO dishes



Science Data Challenge 1 (SDC1)

Continuum emission

- **Continuum emission** images, simulating observations for SKA MID Bands 1, 2 and 5
- Images populated by star forming galaxies (**SFGs**) and active galactic nuclei (**AGN**)
- 3 telescope **integrations** each: 8, 100 and 1000h
- High telescope sensitivity → highly **crowded** images
- **The challenge:** to **find and characterise sources**
- **[SDC1 website](#)**



Zoom-in of the 1.4 GHz maps, showing the same region of the sky with different telescope integrations: 8, 100, 1000 h from left.



Science Data Challenge 1 (SDC1)

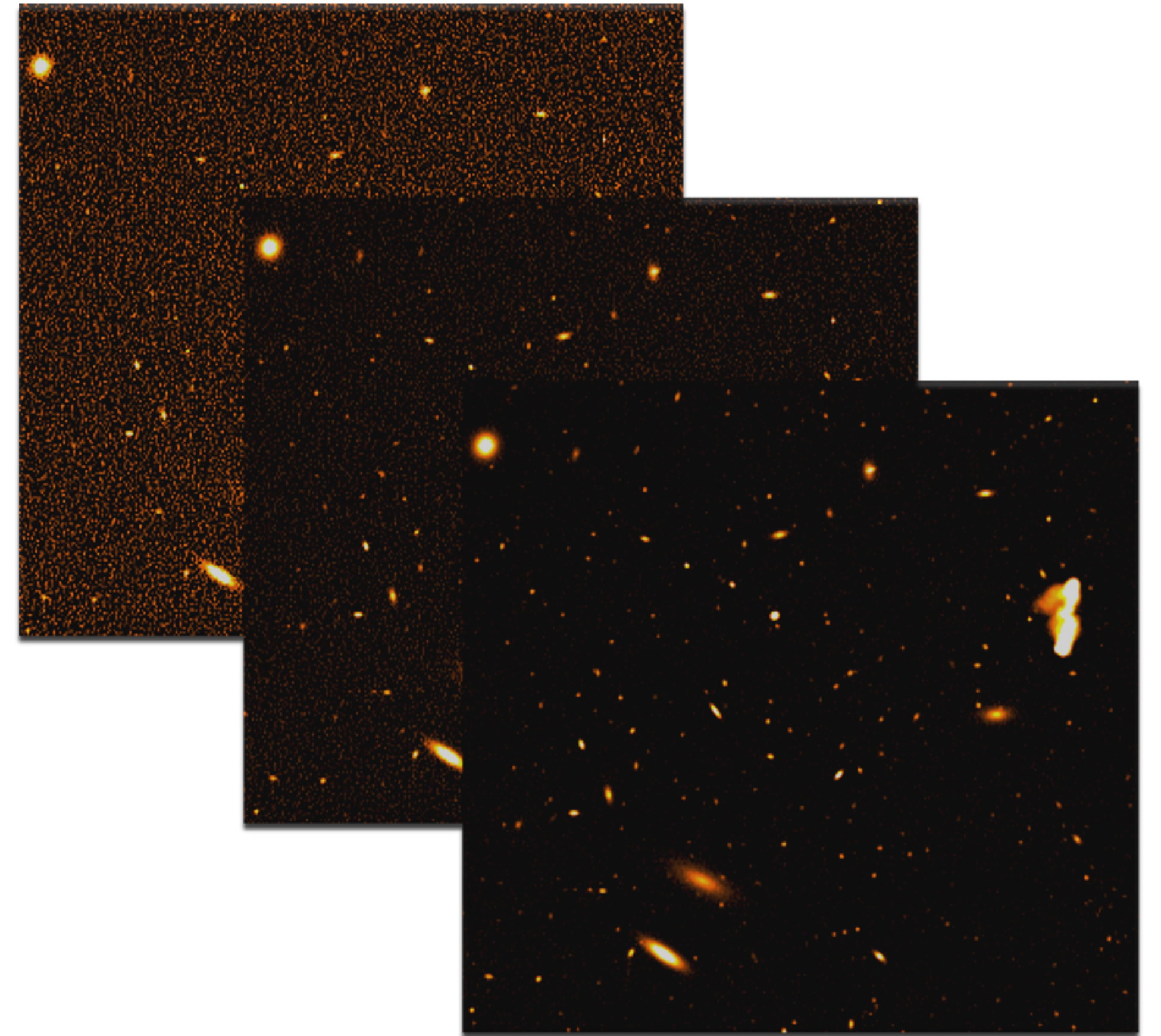
Continuum emission

- **25** square degrees
- **1.5, 0.60 and 0.0913** arcsec FWHM
- Dirty beam sidelobes **4×10^{-4}**

- No simulated calibration errors
- No simulated pointing errors

	560 MHz	1400 MHz	9200 MHz
8 h	2880	710	430
100 h	810	200	120
1000 h	255	73	38
confusion	15	0.36	0.0002

Noise RMS, nJy/beam



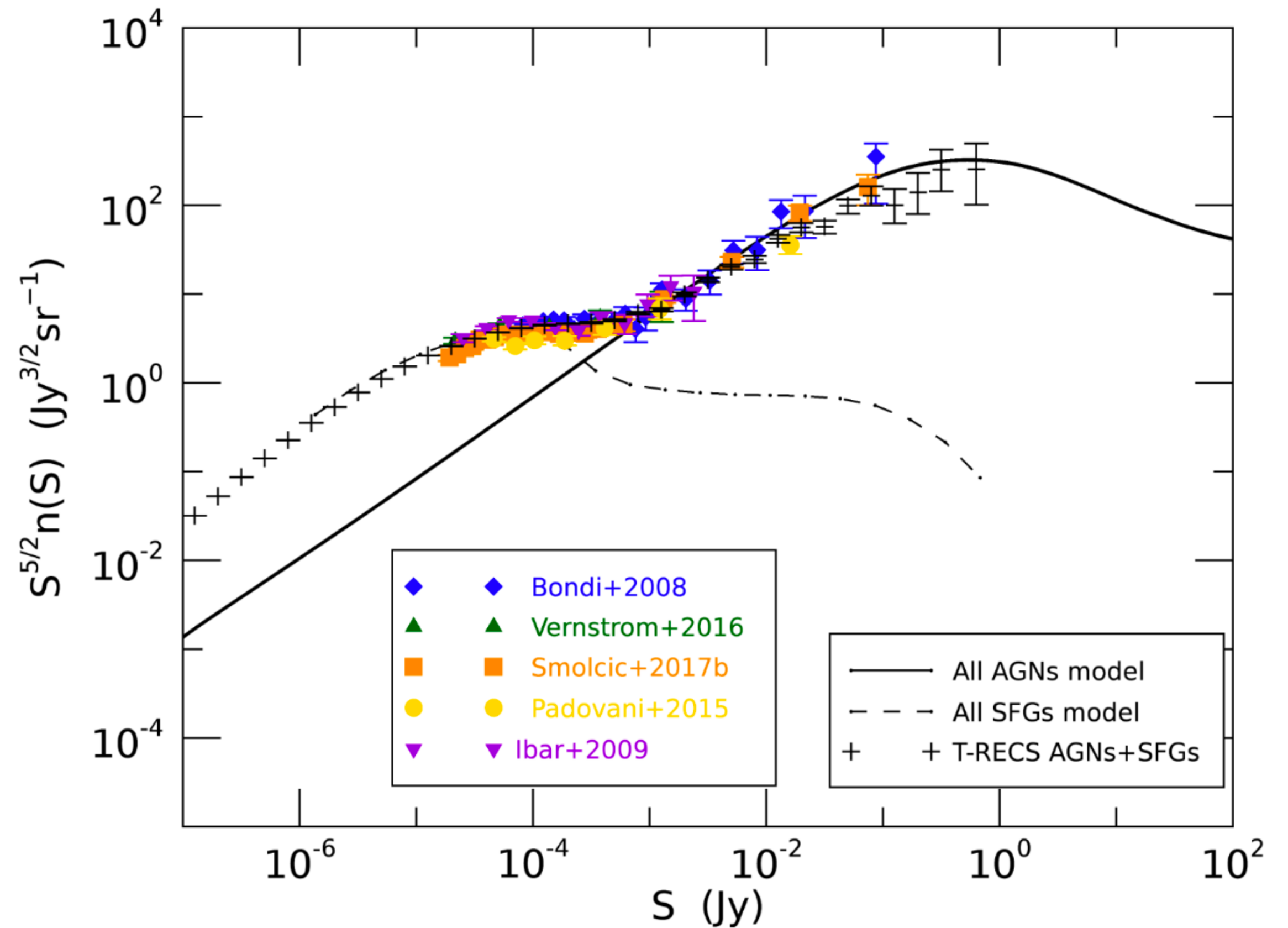
Zoom-in of the 1.4 GHz maps, showing the same region of the sky with different telescope integrations: 8, 100, 1000 h from left.



SDC1 simulations: continuum catalogue

The Tiered Radio Extragalactic Continuum Simulation (T-RECS) *Bonaldi+ 2019*

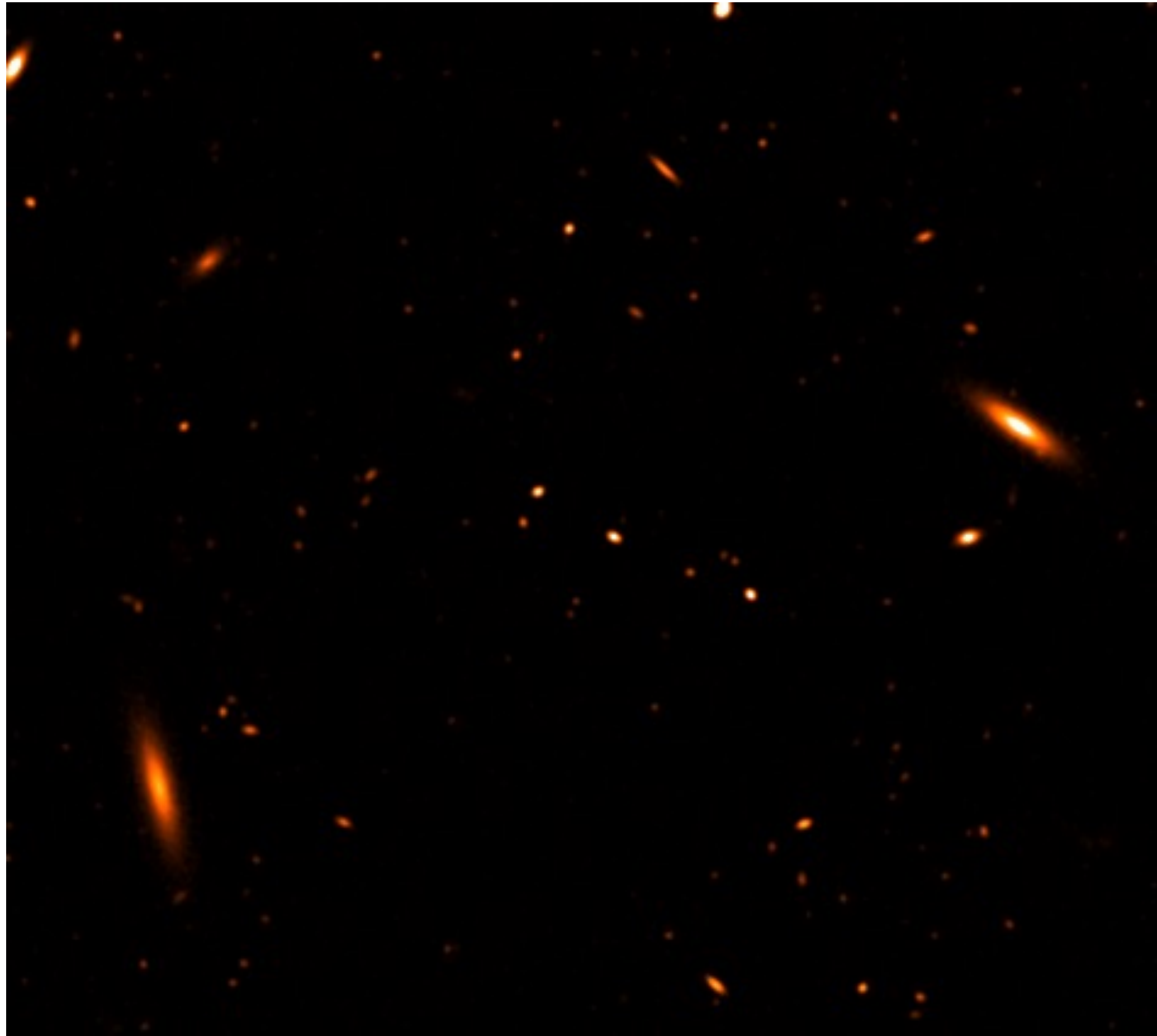
- Good agreement with source counts at the **sub-mJy level**
- **Polarisation**
- Realistic **clustering** via P-millennium simulation (Baugh et al. 2018)



1.4 GHz differential source counts



SDC1 simulations: source morphology



SFGs: Exponential Sersic profiles
Flat spectrum AGN: Gaussians and points

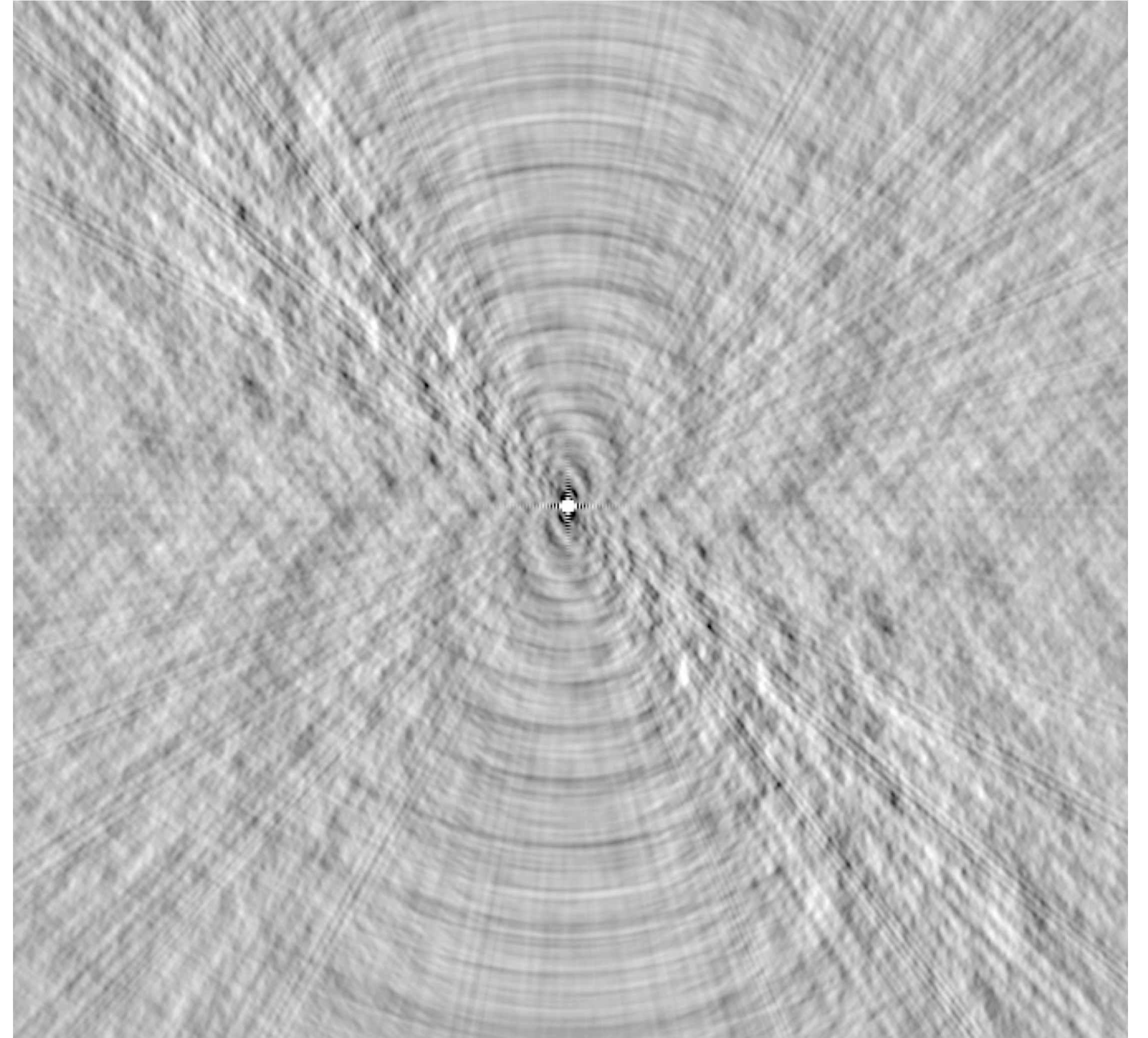


Steep spectrum AGN: DRAGNs
database, University of Manchester



SDC1 simulations: observed continuum

- Addition of **noise**
- '**dirty**' beam simulated
- Imperfect **deconvolution** signatures



SKA MID 1.4 GHz beam



SDC1 results

Continuum observations

Main findings:

- Very **crowded** skies demand new approaches
- Variety of methods including **latest machine learning** techniques
- **Complementarity** of methods: tendency to score well either on finding galaxies *or* measuring them

Square Kilometre Array Science Data Challenge 1: analysis and results

A. Bonaldi,^{1,2*} T. An³, M. Brüggen⁴, S. Burkutean⁵, B. Coelho⁶, H. Goodarzi⁷, P. Hartley¹, P. K. Sandhu⁸, C. Wu⁹, L. Yu¹⁰, M. H. Zhoolideh Haghighi⁷, S. Antón^{11,6}, Z. Bagheri^{7,12}, D. Barbosa⁶, J. P. Barraca^{6,13}, D. Bartashevich⁶, M. Bergano⁶, M. Bonato⁵, J. Brand⁵, F. de Gasperin⁴, A. Giannetti⁵, R. Dodson⁹, P. Jain⁸, S. Jaiswal³, B. Lao³, B. Liu¹⁰, E. Liuzzo⁵, Y. Lu³, V. Lukic⁴, D. Maia¹⁴, N. Marchili⁵, M. Massardi⁵, P. Mohan³, J. B. Morgado¹⁴, M. Panwar⁸, Prabhakar⁸, V. A. R. M. Ribeiro^{6,15}, K. L. J. Rygl⁵, V. Sabz Ali⁷, E. Saremi⁷, E. Schisano¹⁶, S. Sheikhezami^{17,7}, A. Vafaei Sadr¹⁸, A. Wong¹⁹, O. I. Wong^{9,21,20}

Affiliations are at the end of the paper

Monthly Notices of the Royal Astronomical Society, Volume 500, Issue 3, January 2021, Pages 3821–3837



Science Data Challenge 2 (SDC2)

Neutral hydrogen (HI)

HI Science with the SKAO:

- Study the formation and evolution of galaxies by mapping the 21cm spectral line of neutral atomic hydrogen in emission and absorption, over cosmic time
- Structure and gas of cold gas in, around and between galaxies

HI Galaxy Science Science Working group:

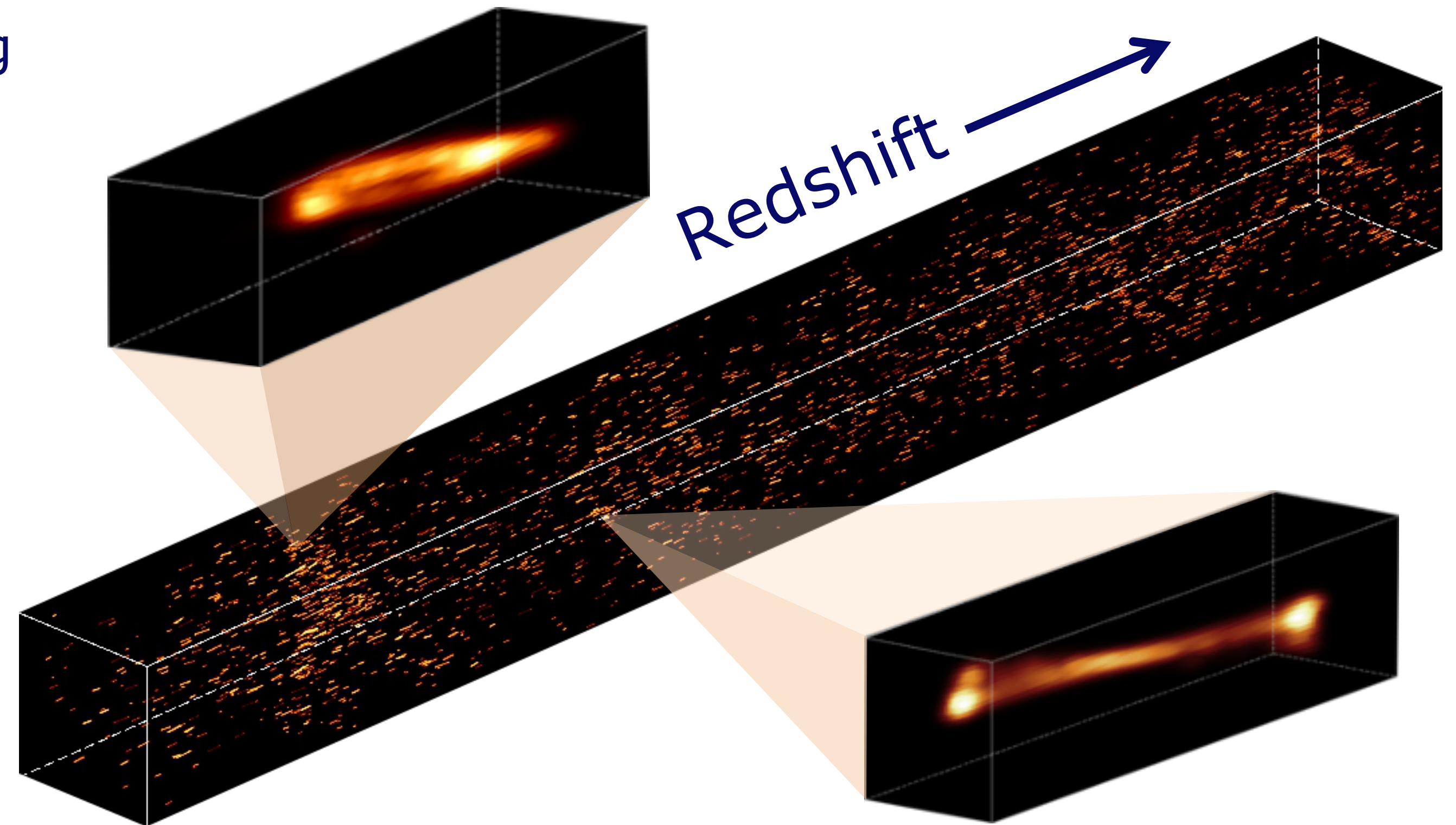
<https://www.skao.int/en/science-users/science-working-groups-focus-groups/111/hi-galaxy-science>

The M81 group

Science Data Challenge 2

Neutral hydrogen (HI)

- **21cm spectral line** image cube, simulating deep **SKA MID** observations (redshift 0.25 to 0.5)
- Image cube populated by **HI** content of **galaxies**
- **2000 h** integration time across **20 sq deg** field of view
- The challenge: to **find and characterise HI sources**
- **Data volume = 1 TB**
- [SDC2 website](#)

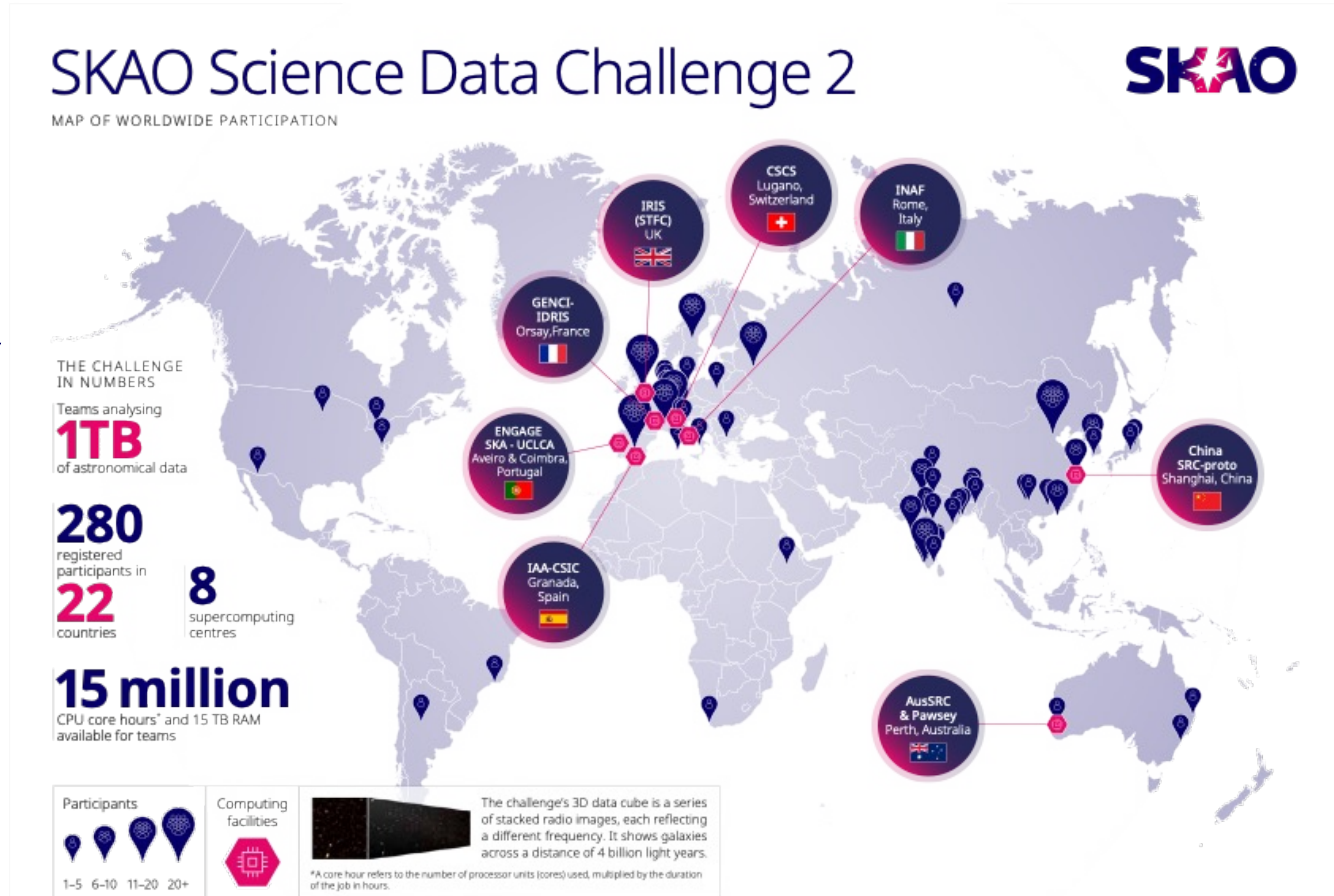


Sample noise-free simulated HI image cube



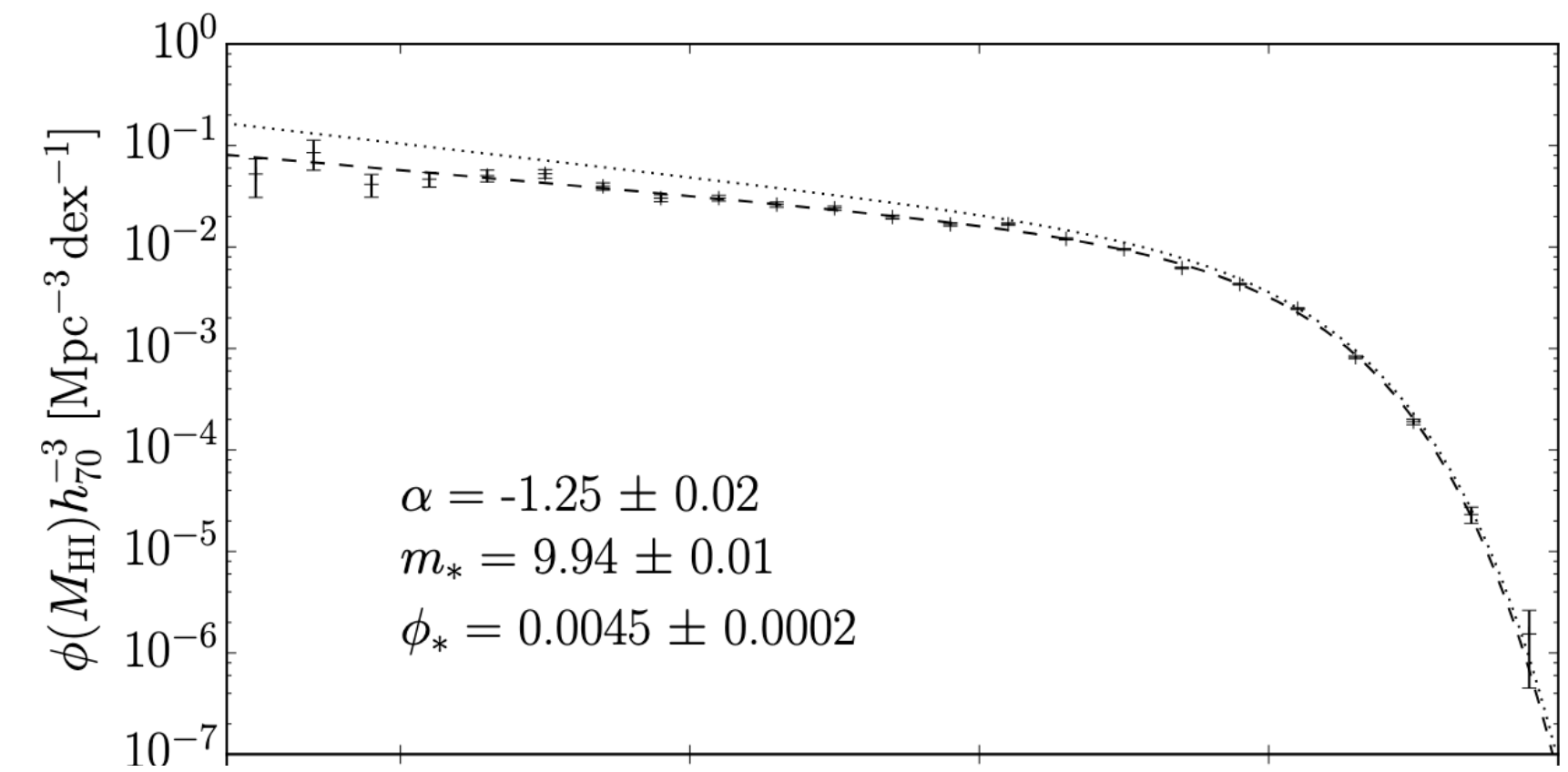
SDC computational facility partners

- **Support from eight** international computing facilities essential to success of SDC2
- Enabled accessible provision of **realistically large dataset**
- Test aspects of the future **SKA Regional Centre** model, e.g.:
 - Community data access
 - New technologies for distributed platform
 - SKA is committed to Open Science best practice

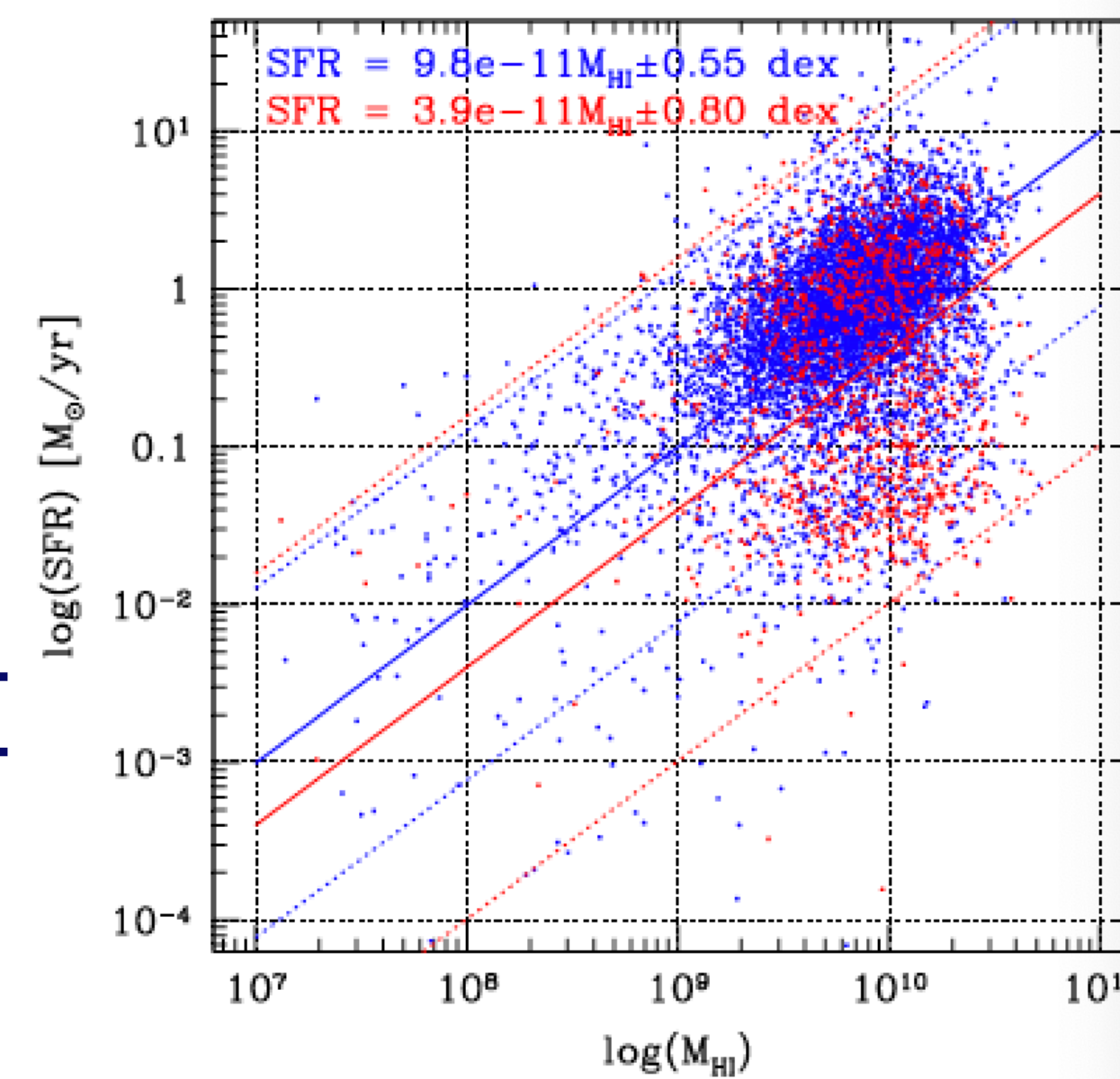


SDC2 simulations: HI and continuum catalogues

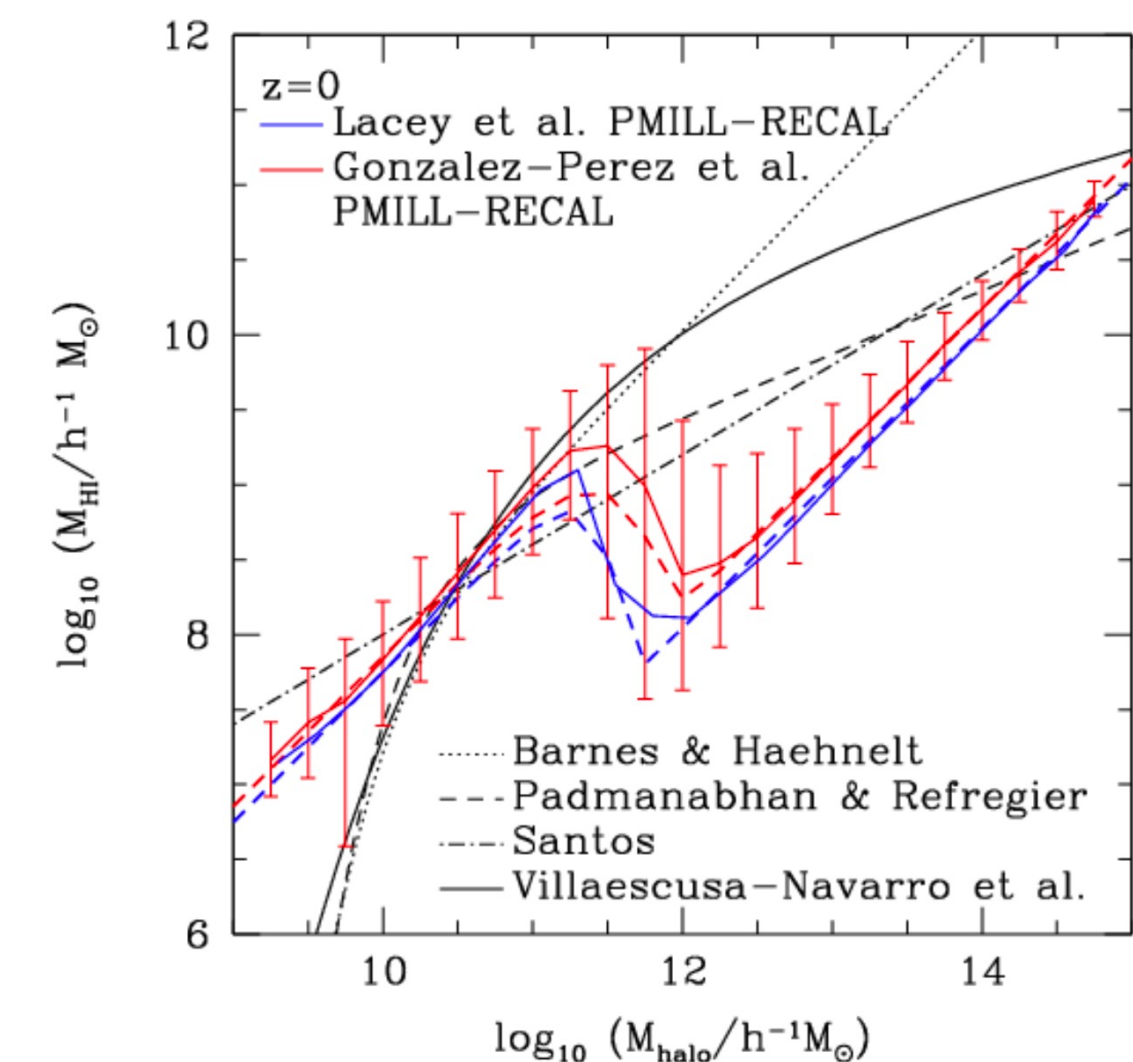
- HI and continuum **catalogues** produced using T-RECS modules
- **HI catalogue:** HI masses drawn from ALFALFA survey HI mass function
- **Continuum catalogue:** HI masses *predicted* from ALFAFLA $M_{\text{HI}}\text{-SFR}$ (for SFGs) and P-Millennium $M_{\text{HI}}\text{-}M_{\text{halo}}$ (for AGN)
- Large scale structure via P-millennium simulation
- Catalogues **cross-matched** by HI mass



Hi mass function, Jones+ 2018



Data from ALFALFA, Giovanelli+ 2005

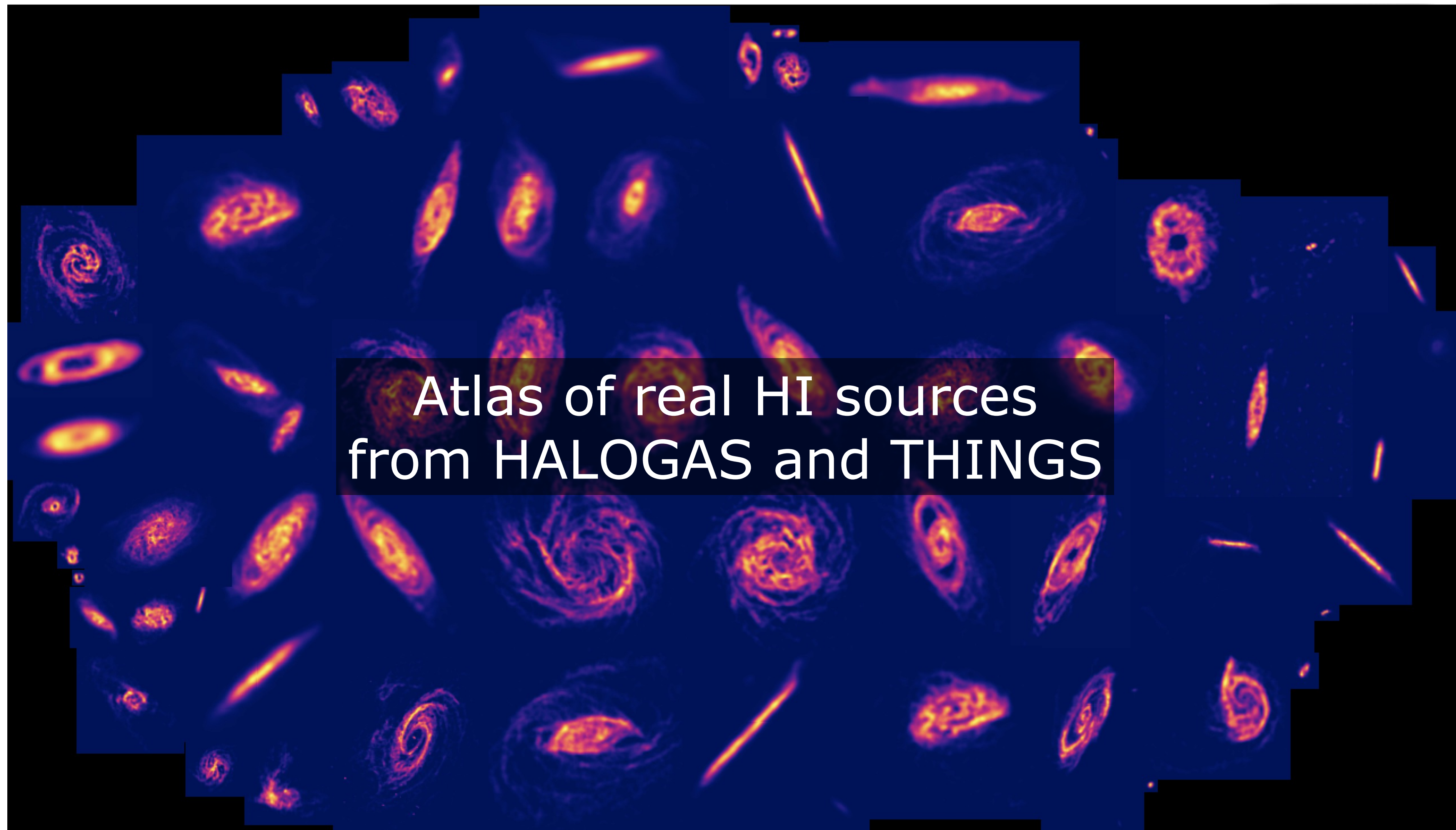


P-Millennium, Baugh+2018



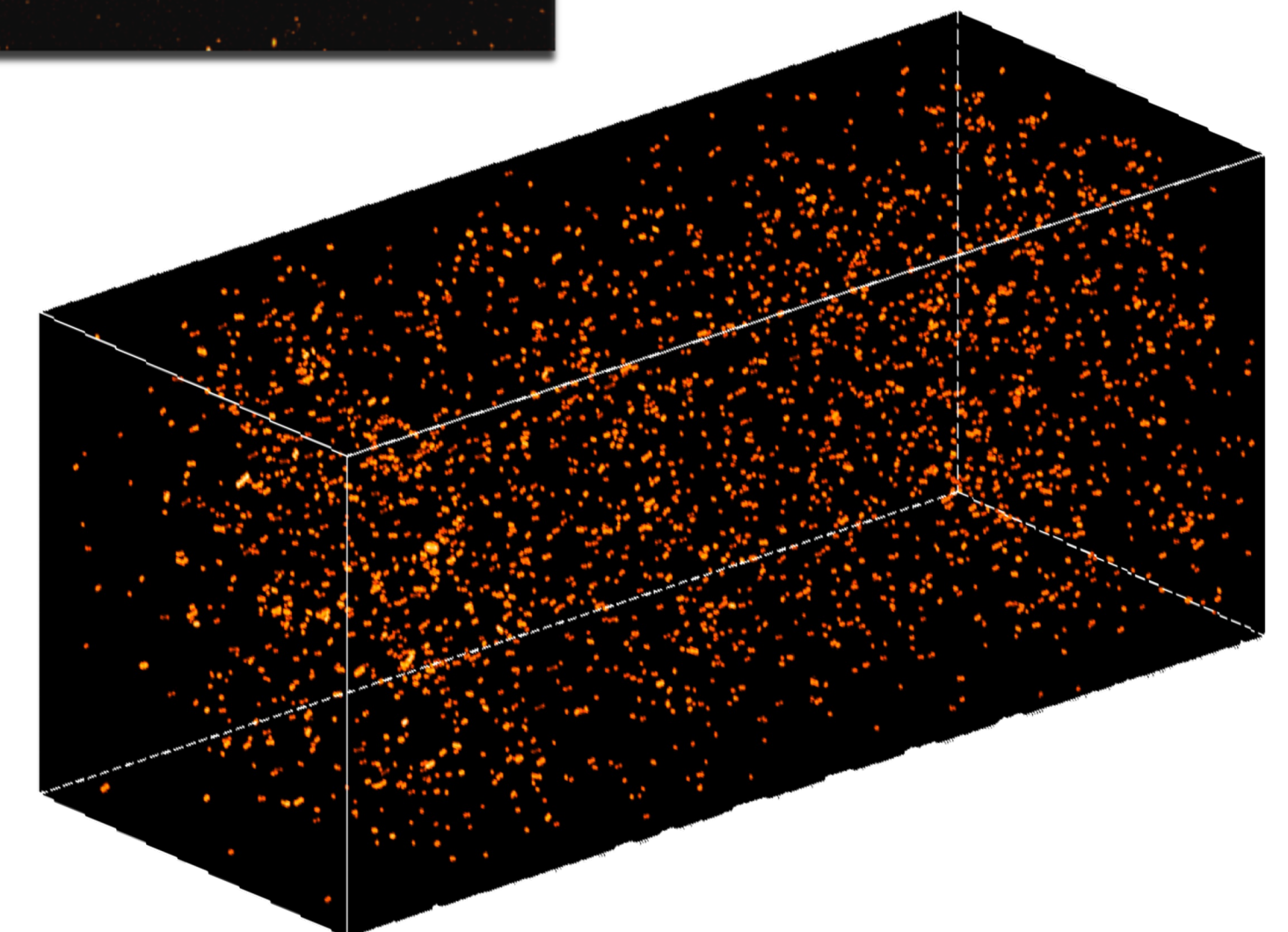
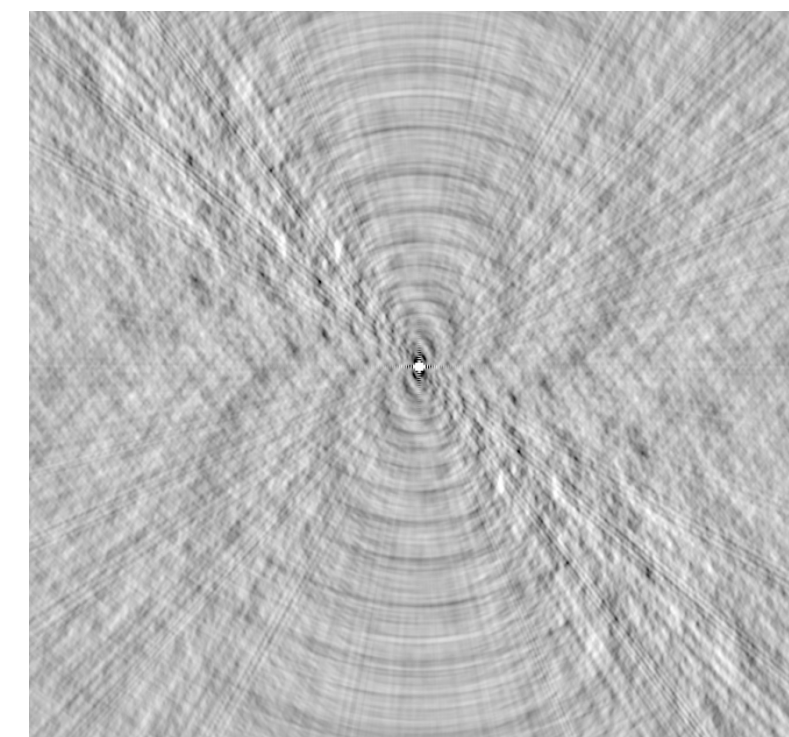
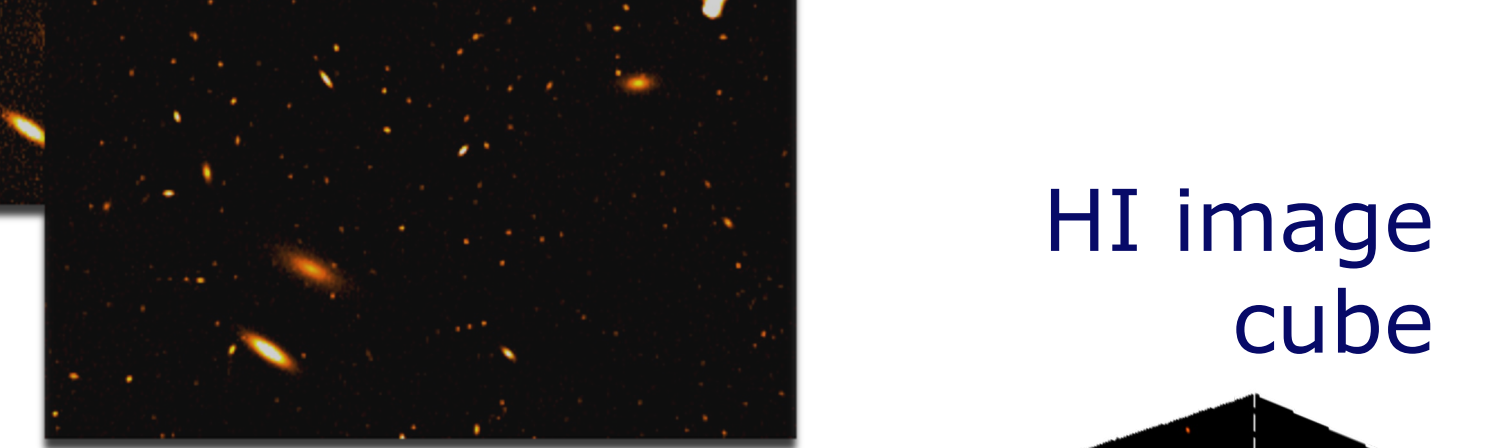
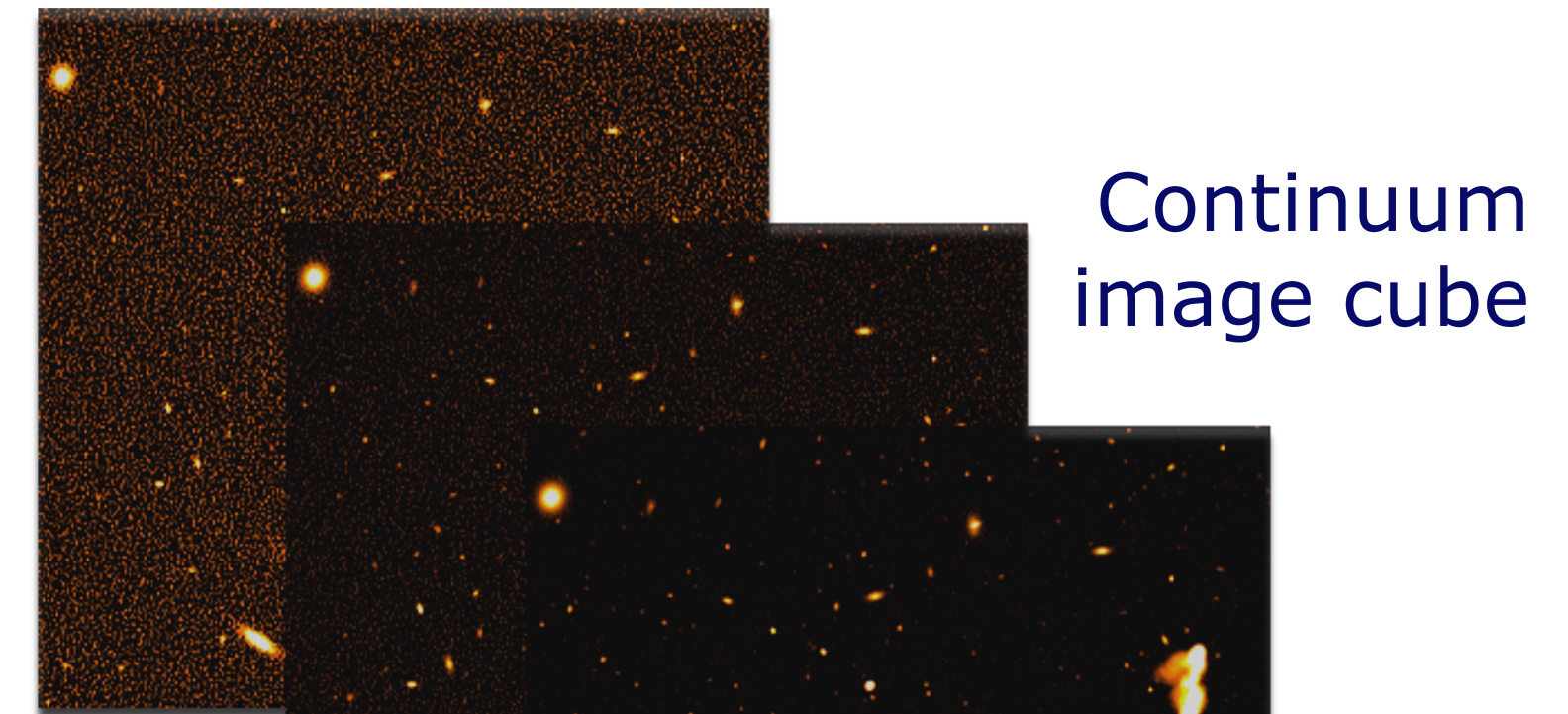
SDC2 simulations: source morphology

Heald et al. 2011;
Walter et al. 2008



SDC2 simulations: observed HI

- HI **absorption** signatures calculated and net absorption/emission cube produced
- Imperfect **continuum subtraction**
- Residual **deconvolution** signatures using a simulated 'dirty' beam
- **Noise** added, with **RFI** flagging simulated



Reproducibility awards

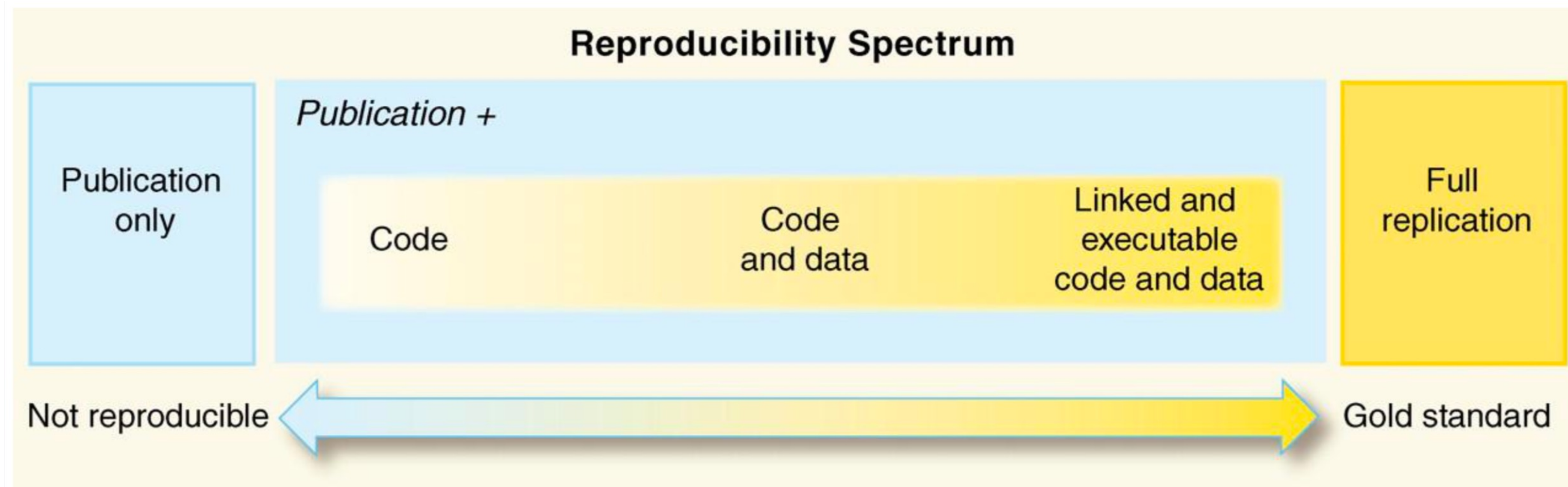


In partnership
with the Software
Sustainability
Institute



www.software.ac.uk

- An essential part of the scientific method, **reproducibility** leads to better, more efficient science.
- **Reusability** generalises this principle to create software that can be adapted by others, allowing previous work to be built upon for the future: a key feature of **Open Science**
- SKA is committed to delivering on the **FAIR** principles for scientific data management



Credit: Rachael Ainsworth



Reproducibility awards



In partnership
with the Software
Sustainability
Institute



www.software.ac.uk

Reproducibility:

Is the software:

- Well-documented
- Easy to install
- Easy to use

Reusability:

Does the software:

- Use an open licence
- Have findable code
- Use code standards
- Use built-in tests

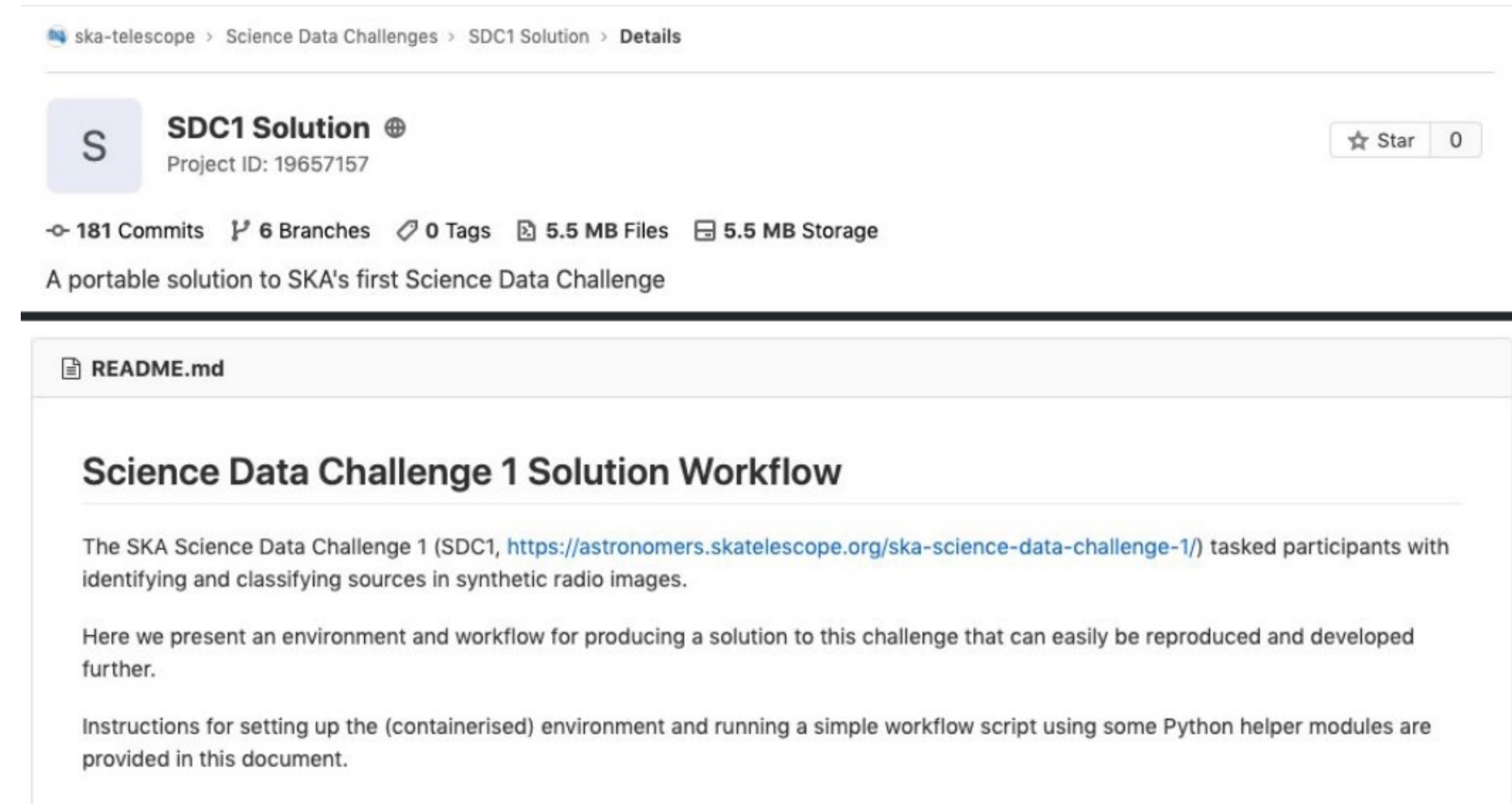
	Reproducibility of the solution		Can the software pipeline be re-run easily to produce the same results? Is it:
			<ul style="list-style-type: none"> • Well-documented Research software documentation best practice • Easy to install Top tips for packaging software • Easy to use Top tips for documentation
Well-documented	High-level description of what/who the software is for is available		er people to develop new projects? Does it:
	High-level description of what the software does is available		Using an open source licence
	High-level description of how the software works is available		code Choosing a repository for your project
	Documentation consists of clear, step-by-step instructions		Writing readable source code
	Documentation gives examples of what the user can see at each step e.g. screenshots or command-line excerpt		Software
	Documentation uses monospace fonts for command-line inputs and outputs, source code fragments, function names, class names etc), BSD 3-Clause
	Documentation is held under version control alongside the code		pository
Easy to install	Full instructions provided for building and installing any software		reader
	All dependencies are listed, along with web addresses, suitable versions, licences and whether they are mandatory or optional		available online
	All dependencies are available		ustainable third-party repository
	Tests are provided to verify that the installation has succeeded		re: Introduction to GitHub
	A containerised package is available, containing the code together with all of the related configuration files, libraries, and dependencies required. Using .e.g. Docker/Singularity		opers
Easy to use	A getting started guide is provided outlining a basic example of using the software e.g. a README file		well
	Instructions are provided for many basic use cases		
	Reference guides are provided for all command-line, GUI and configuration options		les or packages
			age and variable names
			y to the architecture or design
		Testing	Source code has unit tests
			Software recommends tools to check conformance to coding standards e.g. A 'linter' such as PyLint for Python



SDC1 solution


- **Fully reproducible** solution by Alex Clarke (SKAO)
- Utilises **containerisation** in order to package up all software and environment
- Demonstrates **best practices** in research and software development

<https://gitlab.com/ska-telescope/sdc/sdc1-solution>



The screenshot shows the GitLab repository page for 'SDC1 Solution' under the 'ska-telescope' organization. The breadcrumb trail is 'ska-telescope > Science Data Challenges > SDC1 Solution > Details'. The repository name is 'SDC1 Solution' with a Project ID of 19657157. It has 181 commits, 6 branches, 0 tags, 5.5 MB of files, and 5.5 MB of storage. A 'Star' button shows 0 stars. The description is 'A portable solution to SKA's first Science Data Challenge'. The 'README.md' file is open, showing the title 'Science Data Challenge 1 Solution Workflow'. The text in the README describes the SKA Science Data Challenge 1 (SDC1) and provides instructions for setting up the containerised environment and running a workflow script.

ska-telescope > Science Data Challenges > SDC1 Solution > Details

SDC1 Solution 
Project ID: 19657157 ☆ Star 0

181 Commits 6 Branches 0 Tags 5.5 MB Files 5.5 MB Storage

A portable solution to SKA's first Science Data Challenge

README.md

Science Data Challenge 1 Solution Workflow

The SKA Science Data Challenge 1 (SDC1, <https://astronomers.skatelescope.org/ska-science-data-challenge-1/>) tasked participants with identifying and classifying sources in synthetic radio images.

Here we present an environment and workflow for producing a solution to this challenge that can easily be reproduced and developed further.

Instructions for setting up the (containerised) environment and running a simple workflow script using some Python helper modules are provided in this document.



SDC2 results

Neutral hydrogen (HI)

- **12** finalist teams from over **50** institutions
- High level findings:
 - **Complementary** methods
 - Mix of **new and existing** techniques; **machine learning and non-machine learning**
 - **SoFiA package** very popular thanks to excellent documentation and ease of use
 - Analysis of **biases** and **HI mass** recovery with redshift
- Results and analysis from SDC2 prepared for **submission to MNRAS**

SKA Science Data Challenge 2: analysis and results

P. Hartley⁰, A. Bonaldi⁰, R. Braun⁰, J. N. H. S. Aditya⁵⁰, S. Aicardi², L. Alegre⁴⁰, A. Chakraborty¹⁵, X. Chen⁴³, S. Choudhuri¹⁷, A. O. Clarke⁰, J. S. Collinson⁰, D. Cornu¹, L. Darriba³³, M. Delli Veneri⁹, J. Forbrich¹⁹, G. Fourestey⁴¹, B. Fraga¹², A. Galan⁴¹, J. Garrido³³, C. Gheller²⁹, F. Gubanov¹⁰, H. Håkansson²², M. J. Hardcastle¹⁹, C. Heneka⁸, D. Herranz³⁶, K. M. Hess^{24,25,26}, M. Jagannath¹⁸, S. Jaiswal⁵⁰, R. J. Jurek²⁷, D. Kober⁴¹, S. Kitaef²⁸, D. Kleiner²⁹, B. Lao⁵⁰, X. Lu¹, A. Mazumder¹⁵, J. Moldón³³, R. Mondal¹⁶, S. Ni⁴⁴, M. Önnheim²², M. Parra³³, N. Patra¹⁴, A. Peel⁴¹, P. Salomé¹, S. Sánchez-Expósito³⁶, M. Sargent^{41,51,52}, B. Semelin¹, P. Serra²⁹, A. K. Shaw¹³, A. X. Shen^{30,31}, A. Sjöberg²², L. Smith²⁰, A. Soroka¹⁰, V. Stolyarov^{20,21}, E. Tolley⁴¹, M. C. Toribio²³, J. M. van der Hulst²⁵, A. Vafaei Sadr⁴⁷, L. Verdes-Montenegro³³, T. Westmeier²⁸, K. Yu⁴³, L. Yu⁴², L. Zhang⁴⁵, X. Zhang⁴⁴, Y. Zhang⁴⁰, A. Alberdi³³, M. Ashdown²⁰, C.R. Bom¹², M. Brügger⁸, J. Cannon³⁴, R. Chen⁴², J. Coles²⁰, F. Combes^{1,5}, J. Conway²³, J. Ding⁴⁵, J. Freundlich⁴, L. Gao⁴⁴, Q. Guo⁴³, E. Gustavsson²², M. Jirstrand²², M. G. Jones³⁷, G. Józsa³⁵, P. Kamphuis³⁸, M. Lindqvist²³, B. Liu⁴², Y. Liu⁴³, Y. Mao⁴⁶, A. Marchal³, I. Márquez³³, A. Meshcheryakov¹¹, M. Olberg²³, N. Oozeer³⁵, M. Pandey-Pommier³⁹, W. Pei⁴³, B. Peng⁴², J. Sabater⁴⁰, A. Sorgho³³, C. Tasse^{6,7}, A. Wang⁵⁰, Y. Wang⁴³, H. Xi⁴², X. Yang⁵⁰, H. Zhang⁴⁵, J. Zhang⁴⁴, M. Zhao⁴⁴, S. Zuo⁴⁶

Affiliations can be found after the references

Accepted XXX. Received YYY; in original form ZZZ

ABSTRACT

The Square Kilometre Array Observatory (SKAO) will explore the radio sky to new depths in order to conduct transformational science. SKAO data products made available to astronomers will be correspondingly large and complex, requiring the application of advanced analysis techniques in order to extract key science findings. To this end, SKAO is conducting a series of Science Data Challenges, each designed to familiarise the scientific community with SKAO data and to drive the development of new analysis techniques. We present the results from Science Data Challenge 2 (SDC2), which invited participants to find and characterise 233245 neutral hydrogen (HI) sources in a simulated data product representing a 2000 h SKA MID spectral line observation from redshifts 0.25 to 0.5. Through the generous support of eight international supercomputing facilities, participants were able to undertake the Challenge using dedicated computational resources. Alongside the main challenge, ‘reproducibility awards’ were made in recognition of those pipelines which demonstrated Open Science best practice. The Challenge saw over 100 participants develop a range of new and existing techniques, in results which highlight the strengths of multidisciplinary and collaborative effort. The winning strategy – which combined predictions from two independent machine learning techniques to yield a 20 percent improvement in overall performance – underscores one of the main Challenge outcomes: that of method complementarity. It is likely that the combination of methods in a so-called ensemble approach will be key to exploiting very large astronomical datasets.



Science Data Challenge 3

Epoch of Reionisation (EoR)

- Studies of the first generations of galaxies
- Footprints in diffuse neutral hydrogen gas, tracing the evolution of cosmic structure in the $6 < z < 30$ range

EoR Science Working Group:

<https://www.skao.int/en/science-users/science-working-groups-focus-groups/106/epoch-reionization>

SKAO

The LOW telescope

131,000 x 2 m
antennas



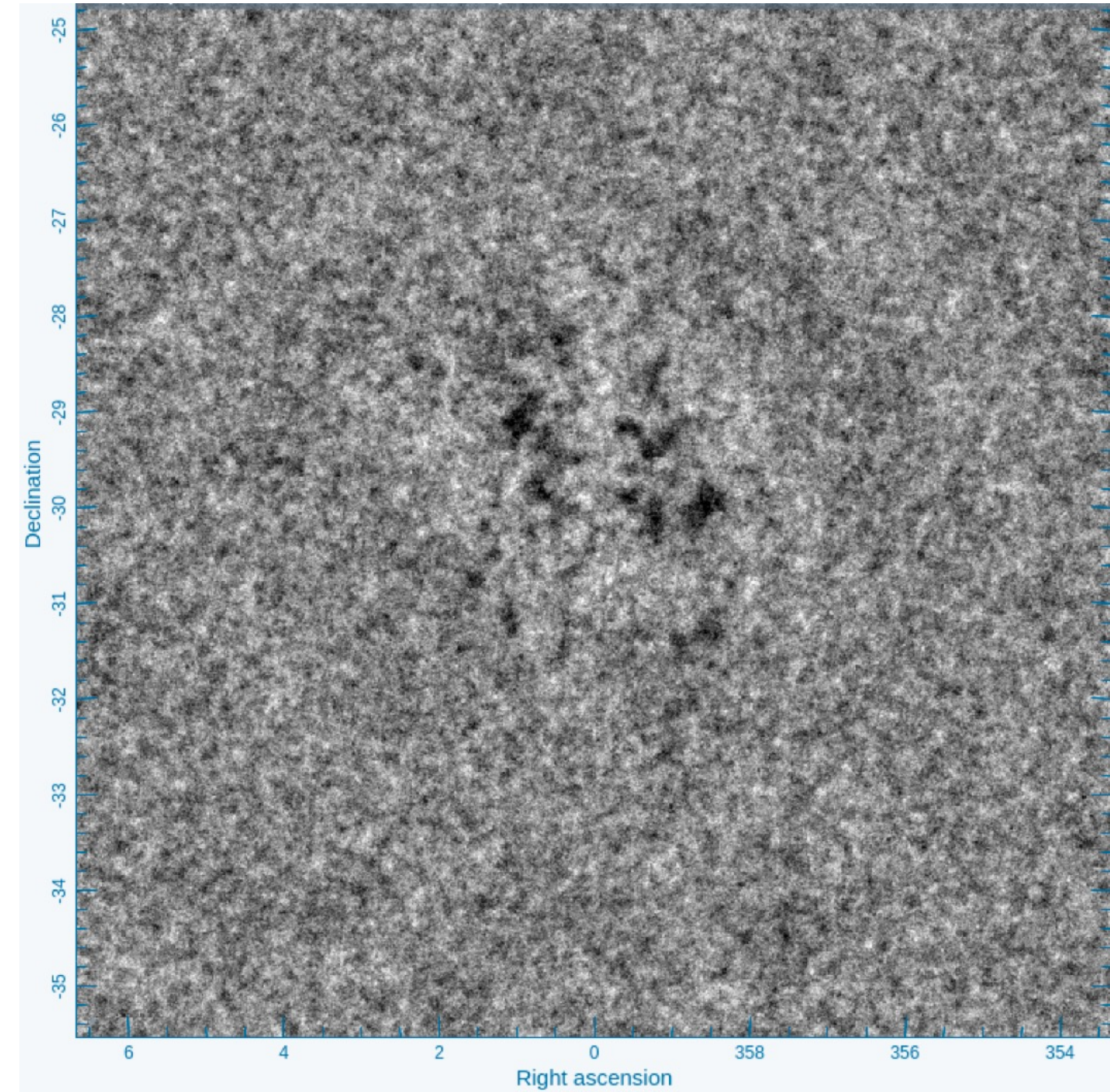
Science Data Challenge 3 (SDC3)

Epoch of Reionisation (EoR)

Developing in collaboration with SKA EoR SWG members

Two parts:

- SDC3 "**Foregrounds**" (SDC3a; SWG Coordinators: C. Trott, V. Jelic)
 - **Foreground removal** exercise
 - SDC3a registration **will open soon**: [SDC3 website](#)
- SDC3 "**Inference**" (SDC3b; SWG Coordinators: A. Mesinger, G. Melema)
 - Extraction of **cosmological parameters**
 - SDC3b launching 2023



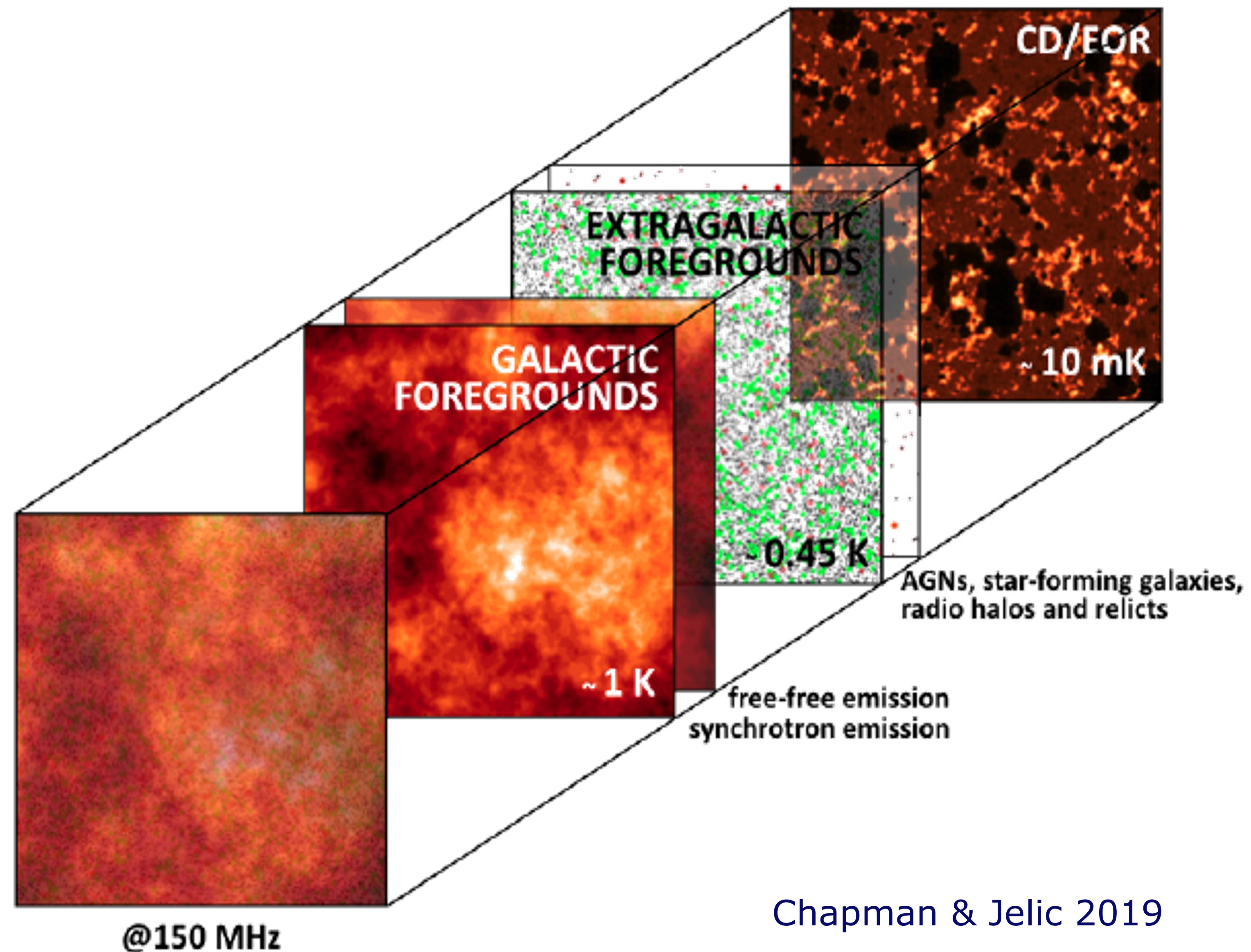
Sample EoR signal with noise added



Science Data Challenge 3a (SDC3a)

“Foregrounds”

- Target participants: SWGs including EoR, Cosmology, Continuum
- Input data: calibrated visibilities and high fidelity image cube
- Challenge will be based on:
 - a) Ability to remove the point source + diffuse foregrounds from the data-set
 - b) Ability to extract the spherical and cylindrical power spectrum



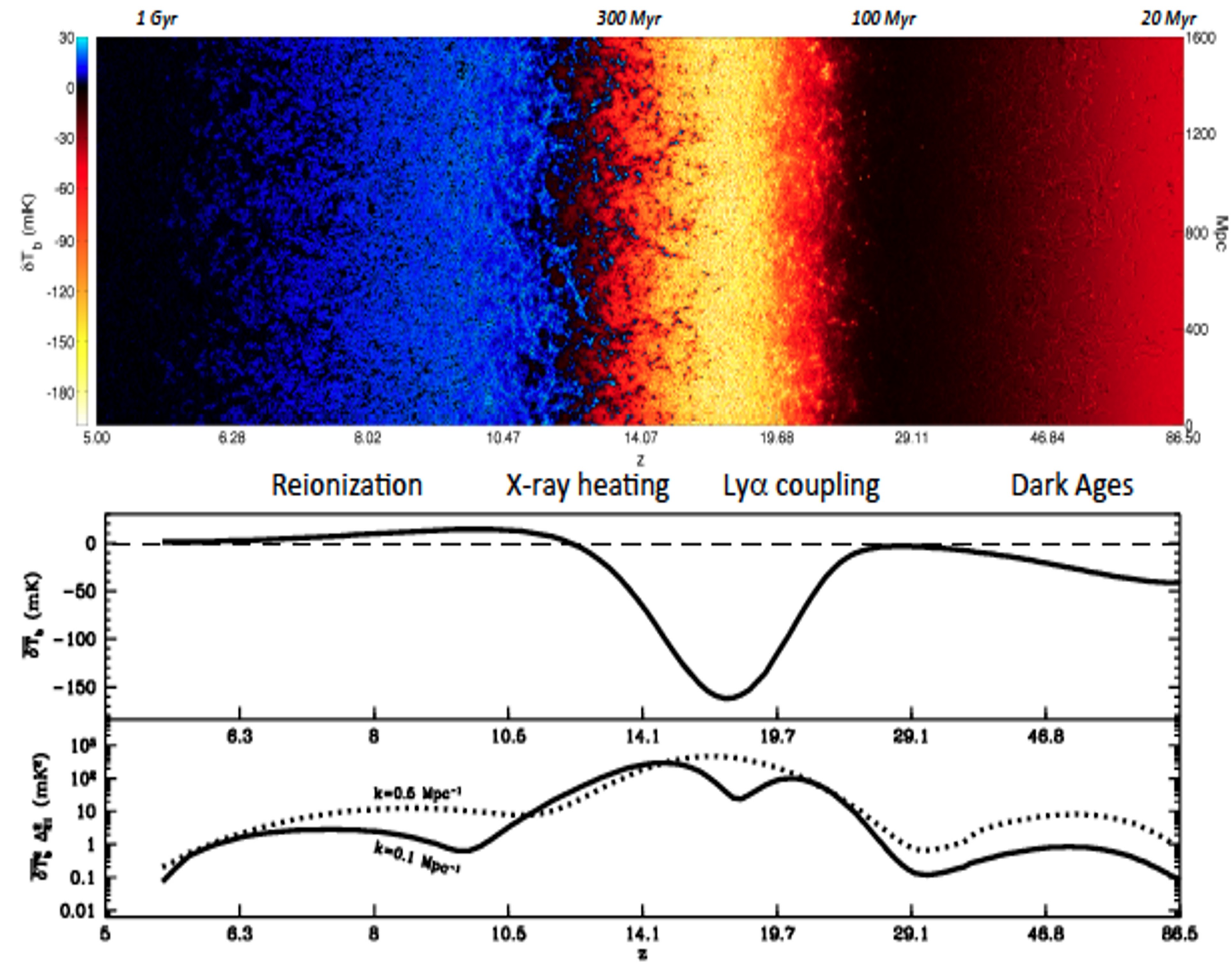
Chapman & Jelic 2019



Science Data Challenge 3a (SDC3b)

“Inference”

- Target Participants: SWGs including EoR
- Input data: EoR power spectrum + noise and residual foreground contamination
- Challenge will be based on:
 - a) ability to extract the input EoR history (ionisation fraction)



Mesinger+ (2016)



Science Data Challenge 3 (SDC3)

Website [SDC3 website](#)



SKA SDC3

Overview

Challenges ▾

Computational resources ▾

Challenge registration ▾

Discussion forum

FAQs

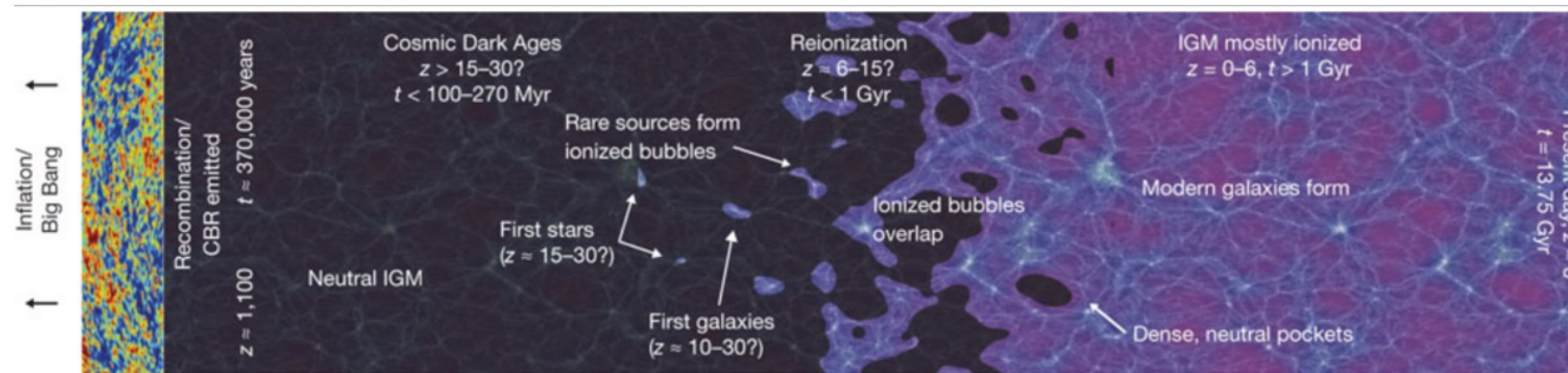


Epoch of Reionisation

Given our current understanding, cosmological history prior to the current state of the universe can be divided into several distinct epochs:

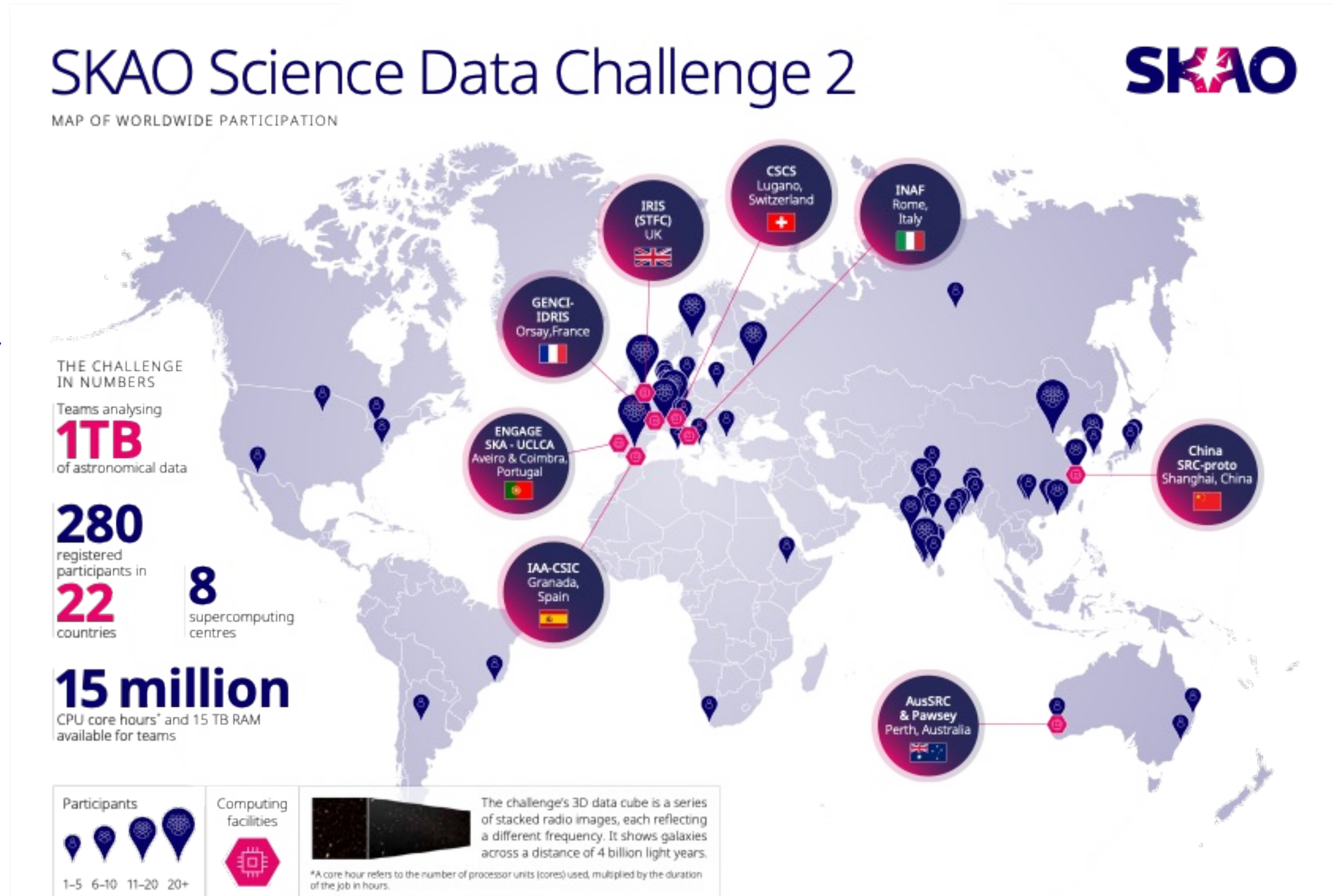
1. Cosmic Dark Ages ($1100 < z < 30$)
2. Cosmic Dawn ($30 < z < 15$)
3. Epoch of Reionisation ($15 < z < 6$)

After the Big Bang, our universe expanded and cooled to the point at which ionised hydrogen recombined to its atomic state, with 21cm emission subsequently being observed in absorption against the background. With the appearance of a significant population of galaxies, the first stars formed and began to heat their surroundings. This heating shifted the absorption profile of the 21 cm Hydrogen line into emission. However, as this heating continued, more of the surrounding gas became ionised to the point at which this emission ceased. Having only roughly constrained the time periods at which these important events occurred, a primary science goal of the SKA is to constrain that timeline, which is where Science Data Challenge 3 comes in.



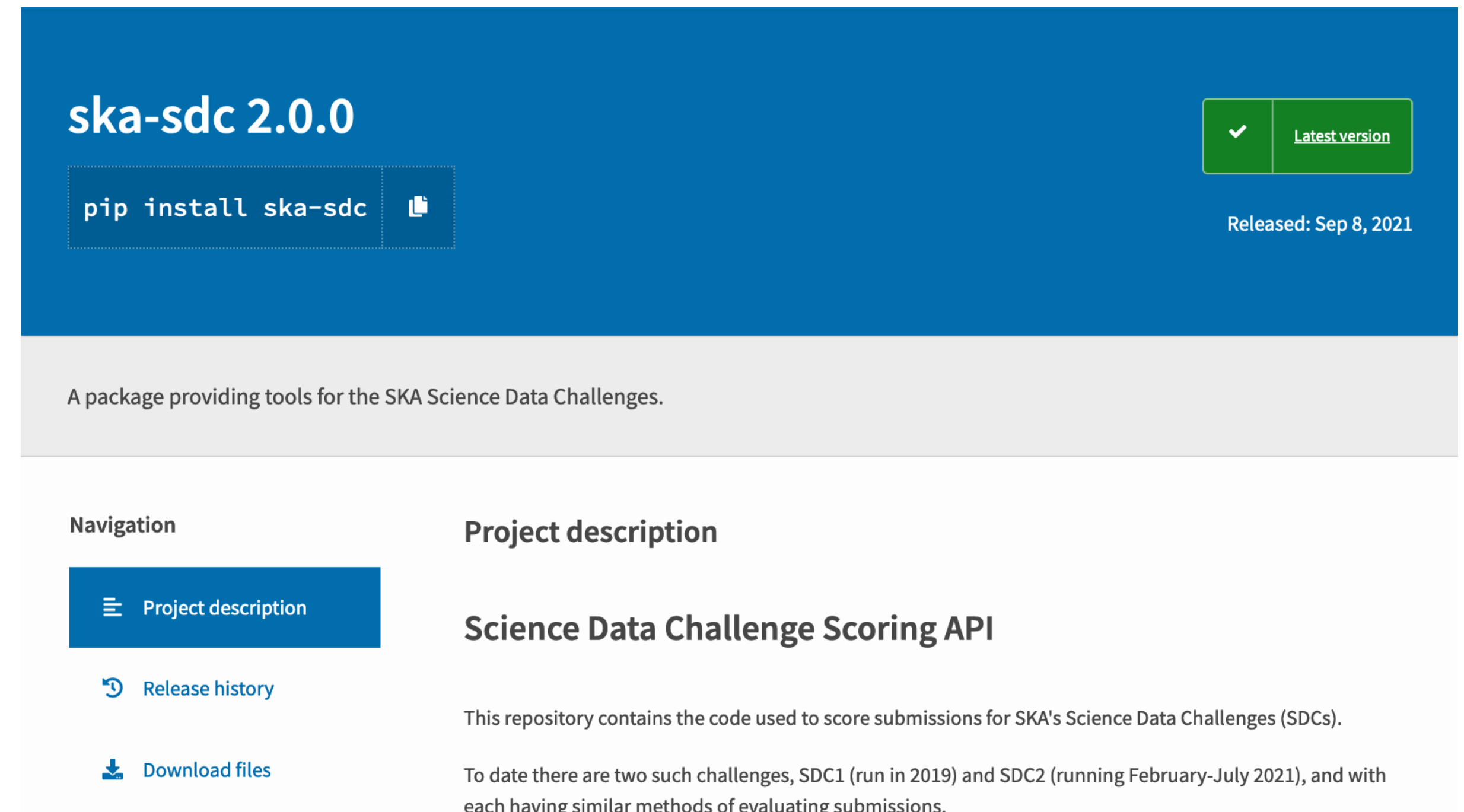
SDC computational facility partners

- **Support from eight** international computing facilities essential to success of SDC2
- Enabled accessible provision of **realistically large dataset**
- Test aspects of the future **SKA Regional Centre** model, e.g.:
 - Community data access
 - New technologies for distributed platform
 - SKA is committed to Open Science best practice



Prototyping for the SKA Regional Centres (SRCs)

- SDC2 **scoring service** made use of **web technologies** for remote scoring by participants
- SDC3 **data transfer** as a test case for **Rucio** solution: “Store, manage, and process data in a heterogeneous distributed environment”

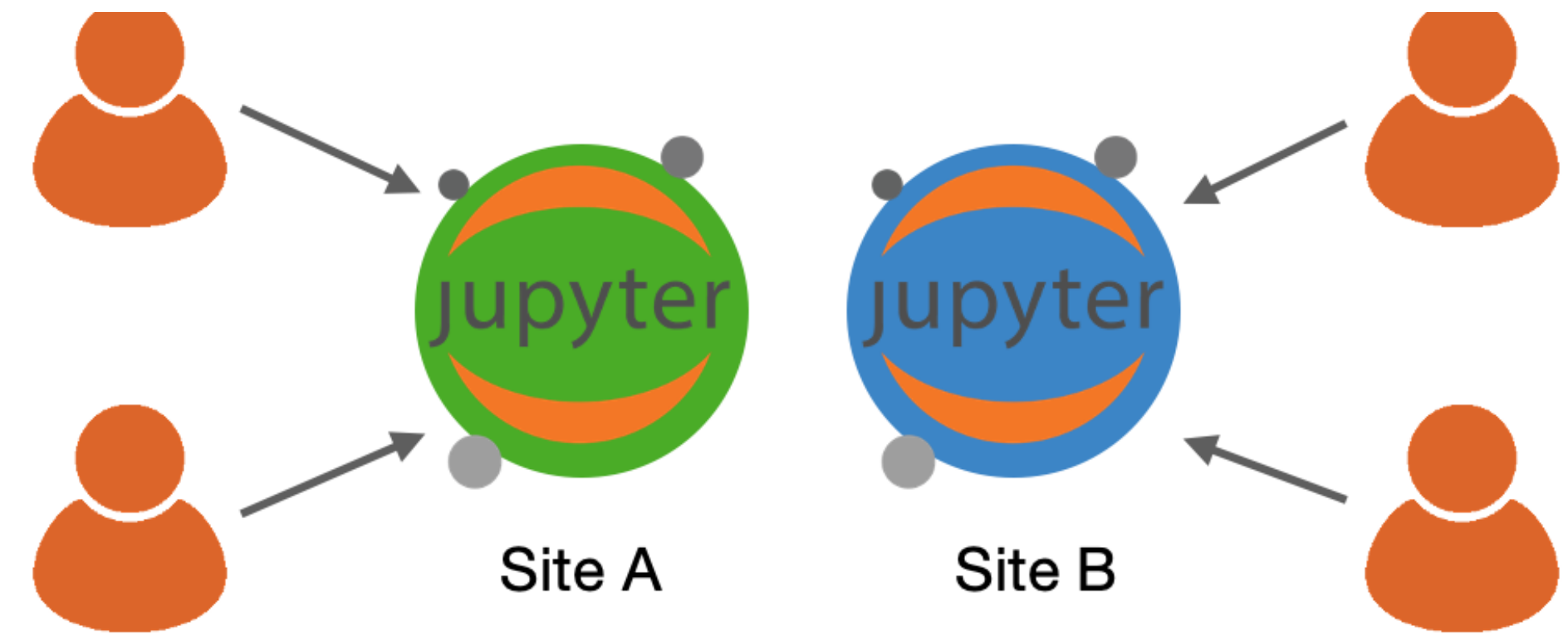


The screenshot shows the PyPI page for the package 'ska-sdc 2.0.0'. The package is marked as the 'Latest version' and was released on Sep 8, 2021. The installation command is 'pip install ska-sdc'. The description states: 'A package providing tools for the SKA Science Data Challenges.' The navigation menu includes 'Project description', 'Release history', and 'Download files'. The project description is titled 'Science Data Challenge Scoring API' and states: 'This repository contains the code used to score submissions for SKA's Science Data Challenges (SDCs). To date there are two such challenges, SDC1 (run in 2019) and SDC2 (running February-July 2021), and with each having similar methods of evaluating submissions.'



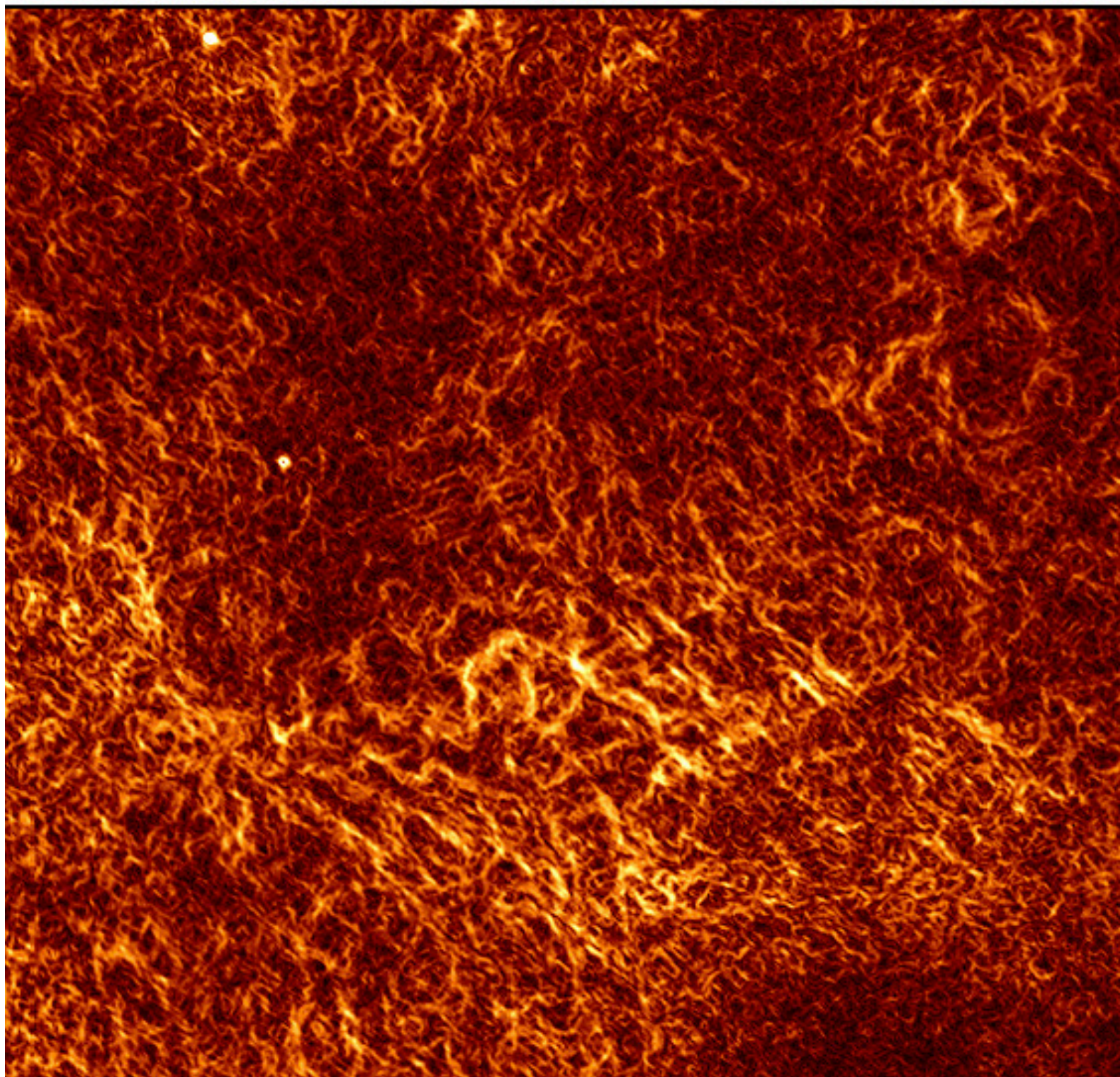
Prototyping for the SKA Regional Centres (SRCs)

- Potential to support JupyterHub environment (to work with **containers**)
 - Further enable reproducibility
 - Would support science community to be able to deploy pipelines
- Deploy **image viewers**, e.g. CARTA



Future Science Data Challenges

- **Cosmic magnetism** SDC (T. Akahori+), **Transients** SDC, and more
- Stay tuned for news and updates!



'Snakes' of cosmic magnetism. Credit: B. Gaensler et al. Quasar schematic. Credit: NASA



Summary

- **Science Data Challenges** goal to prepare the community for the **size and complexity** of SKAO data
- Two data challenges completed: **collaboration and method complementarity** a highlight of both
- **Generous support** from international computing facilities
 - Enables accessible provision of **realistically large dataset**
 - Test aspects of the future **SKA Regional Centre** model
- **EoR Challenge** launching soon
- **Thank you!**

