

one observatory
two telescopes
three continents

SKA Science Data Challenges

Philippa Hartley
SKAO Scientist

Swiss SKA Days 2022



The SKA data journey

SKA LOW



SKA MID

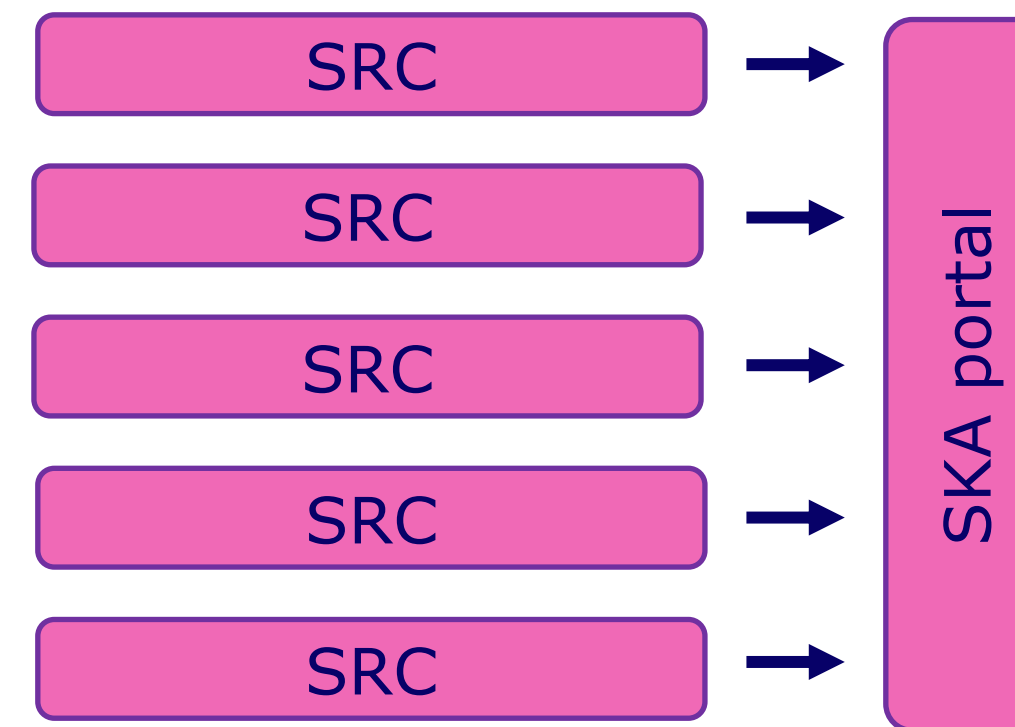
SDP: Science Data Processor



SDP prototype, Cambridge, UK

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

SRC: SKA Regional Centre



Distributed facilities

600 PB/yr

User data
products up to
TBs in size

Key Science Projects



Scientific analysis



SKA 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**

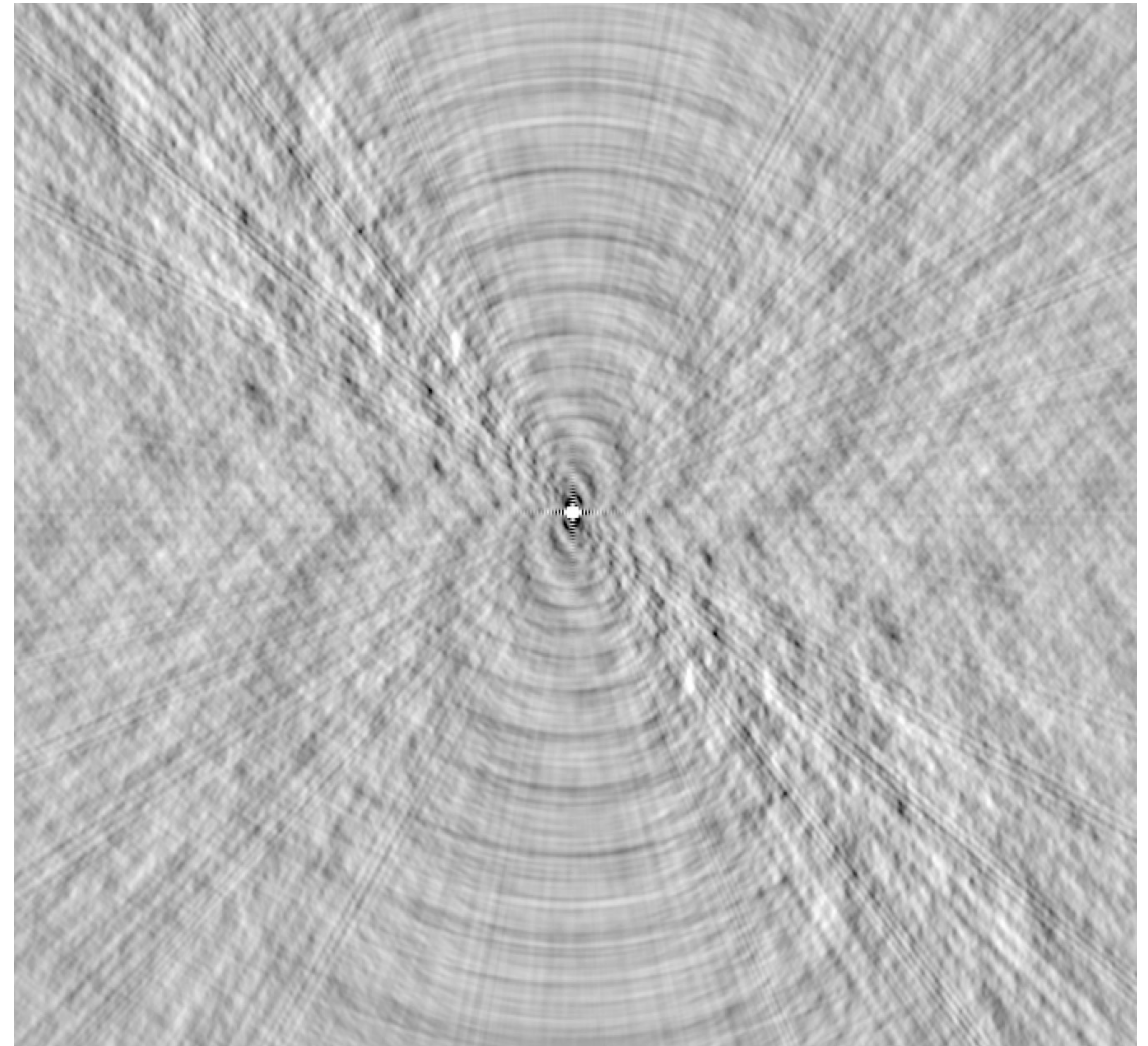
SDC data products are made publicly available for the long term



SKA 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



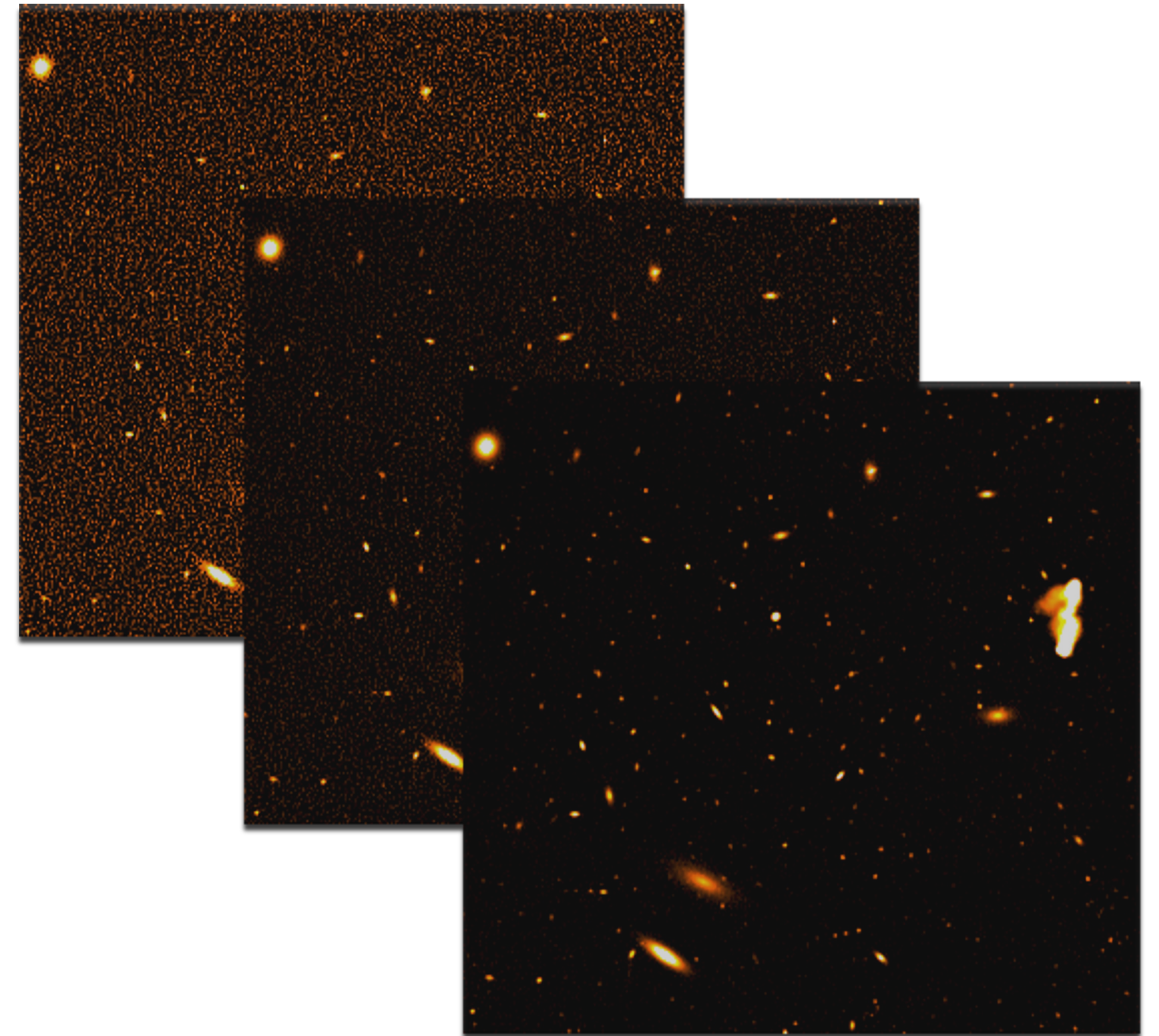
SKA MID 1.4 GHz beam



Science Data Challenge 1

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

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. Zhooldideh 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

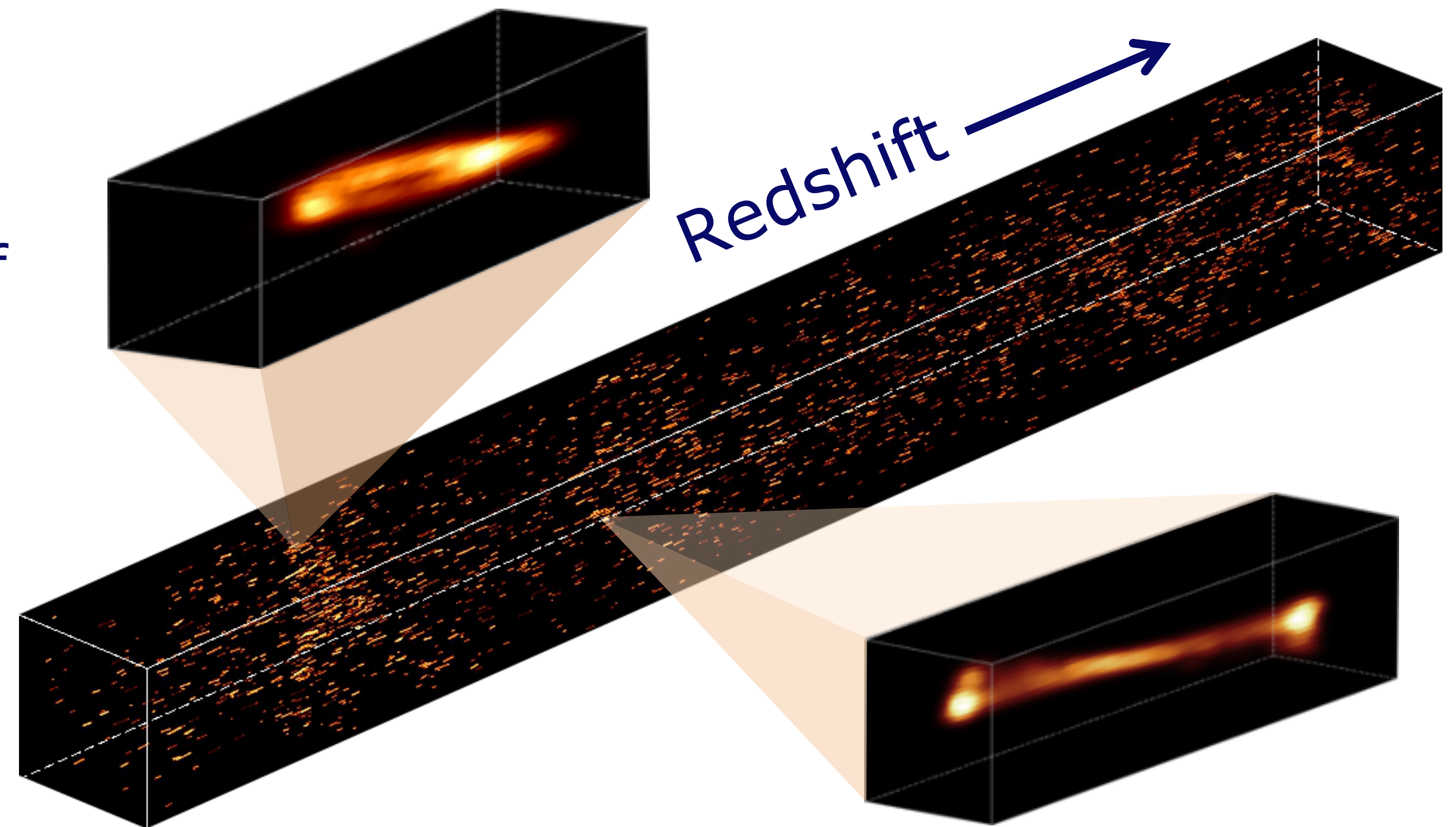
Monthly Notices of the Royal Astronomical Society, Volume 500, Issue 3, January 2021, Pages 3821–3837



Science Data Challenge 2

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



SDC2 results paper

- 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. Korber⁴¹, S. Kitaeff²⁸, D. Kleiner²⁹, B. Lao⁵⁰, X. Li¹¹, A. Mazumder¹⁵, J. Moldón³³, R. Mondal³³, S. Ni⁴⁴, M. Önnheim²², M. Parra³³, N. Patra¹⁴, A. Peierl⁴¹, P. Salomé¹, S. Sánchez-Expósito³³, M. Sargent^{41,51,52}, B. Semelin¹, P. Serra²⁹, A. K. Sengupta¹³, A. X. Shen^{30,31}, A. Sjöberg²², J. 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. Wu⁴³, L. Yu⁴², L. Zhang⁴⁵, X. Zhang⁴⁴, Y. Zhang⁵⁰, A. Alberdi³³, M. Ashdown²⁰, C.R. Bom¹², M. Brüggen⁸, 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

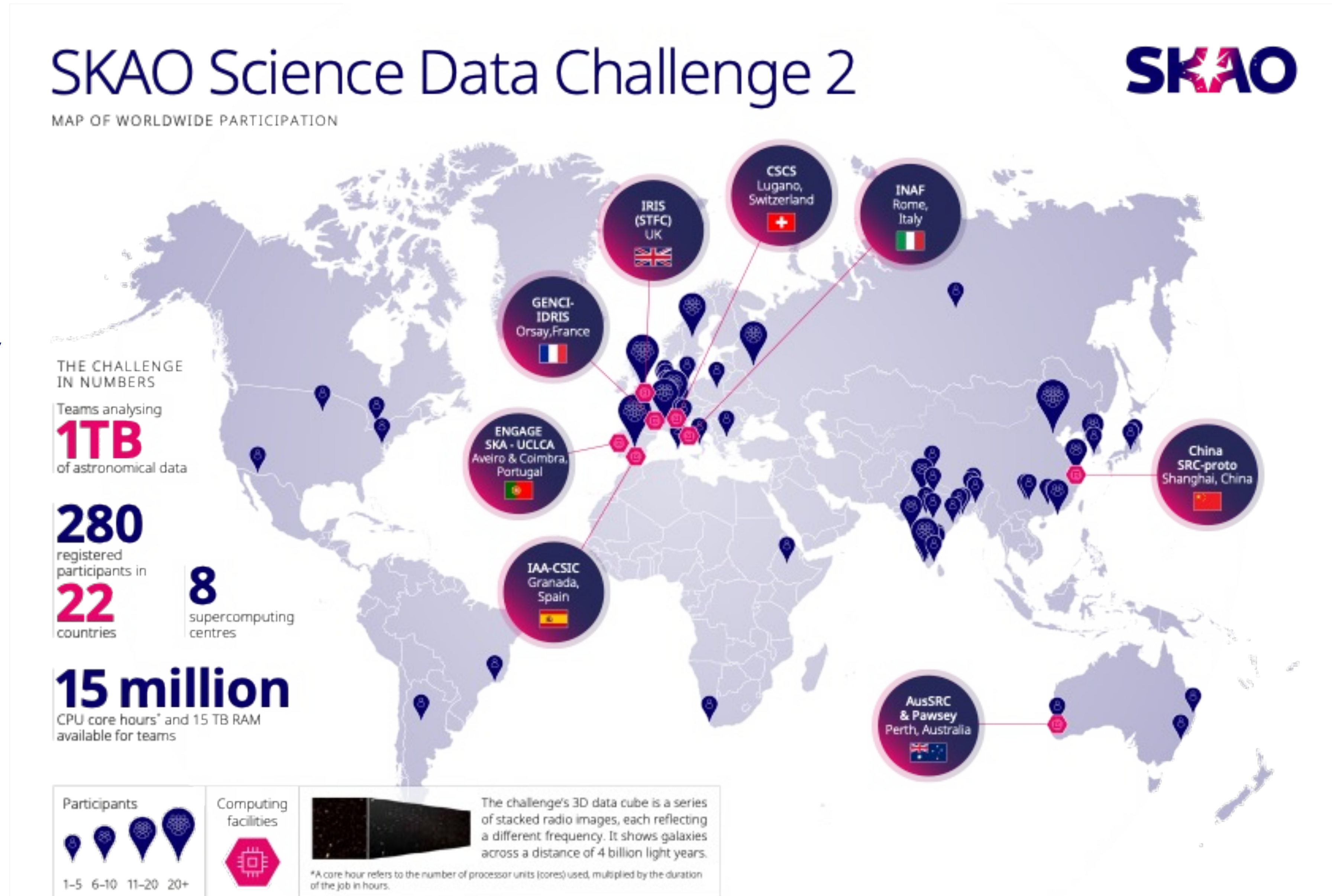
ABSTRACT

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.



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



Reproducibility awards

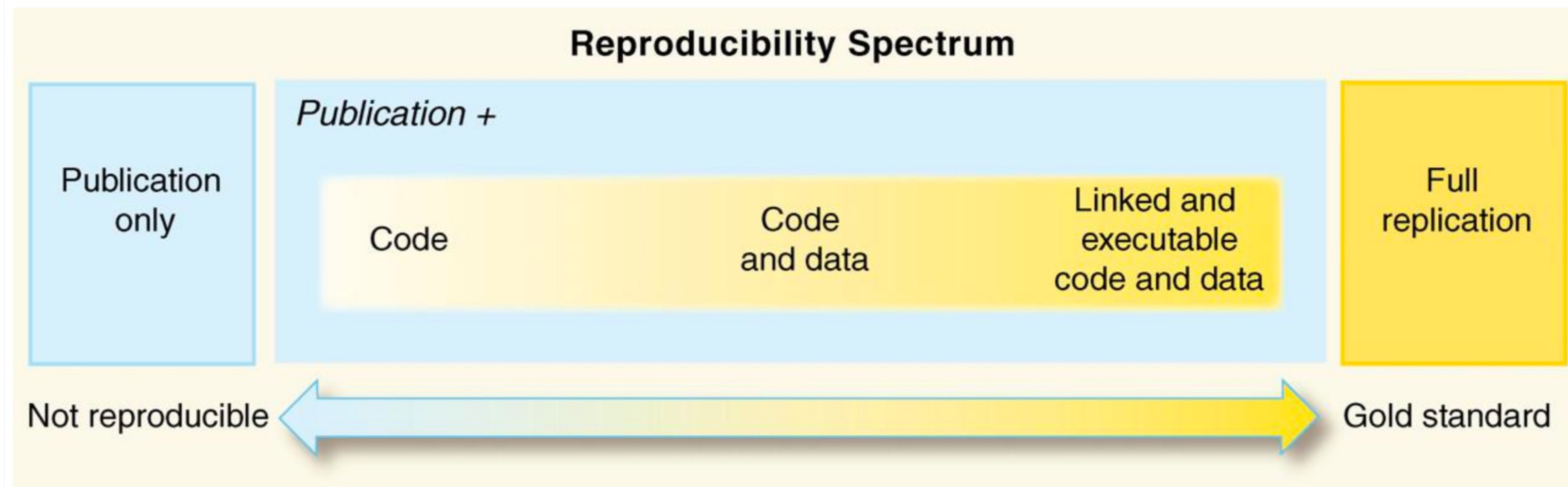


In partnership
with the Software
Sustainability
Institute



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



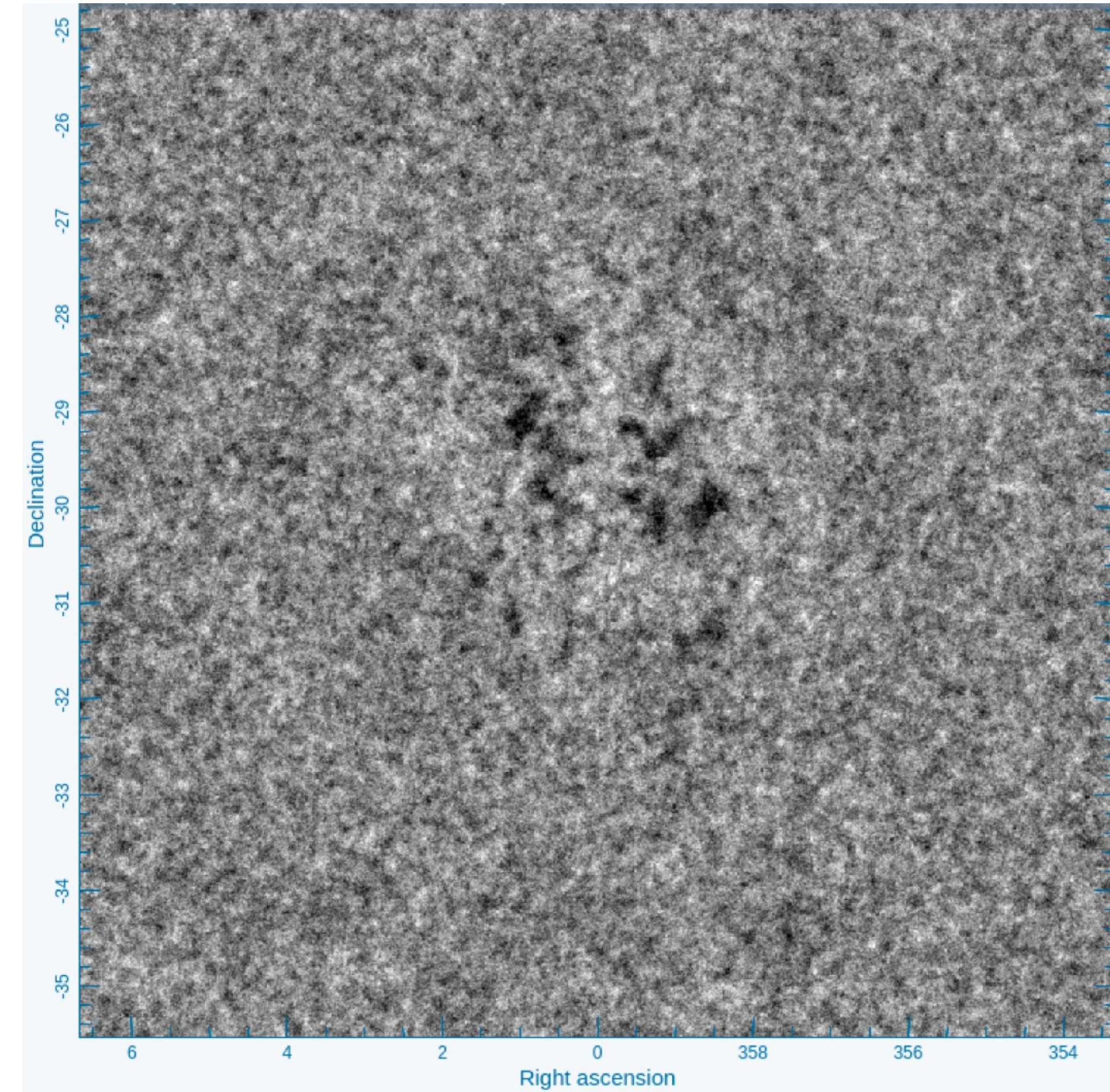
Science Data Challenge 3

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

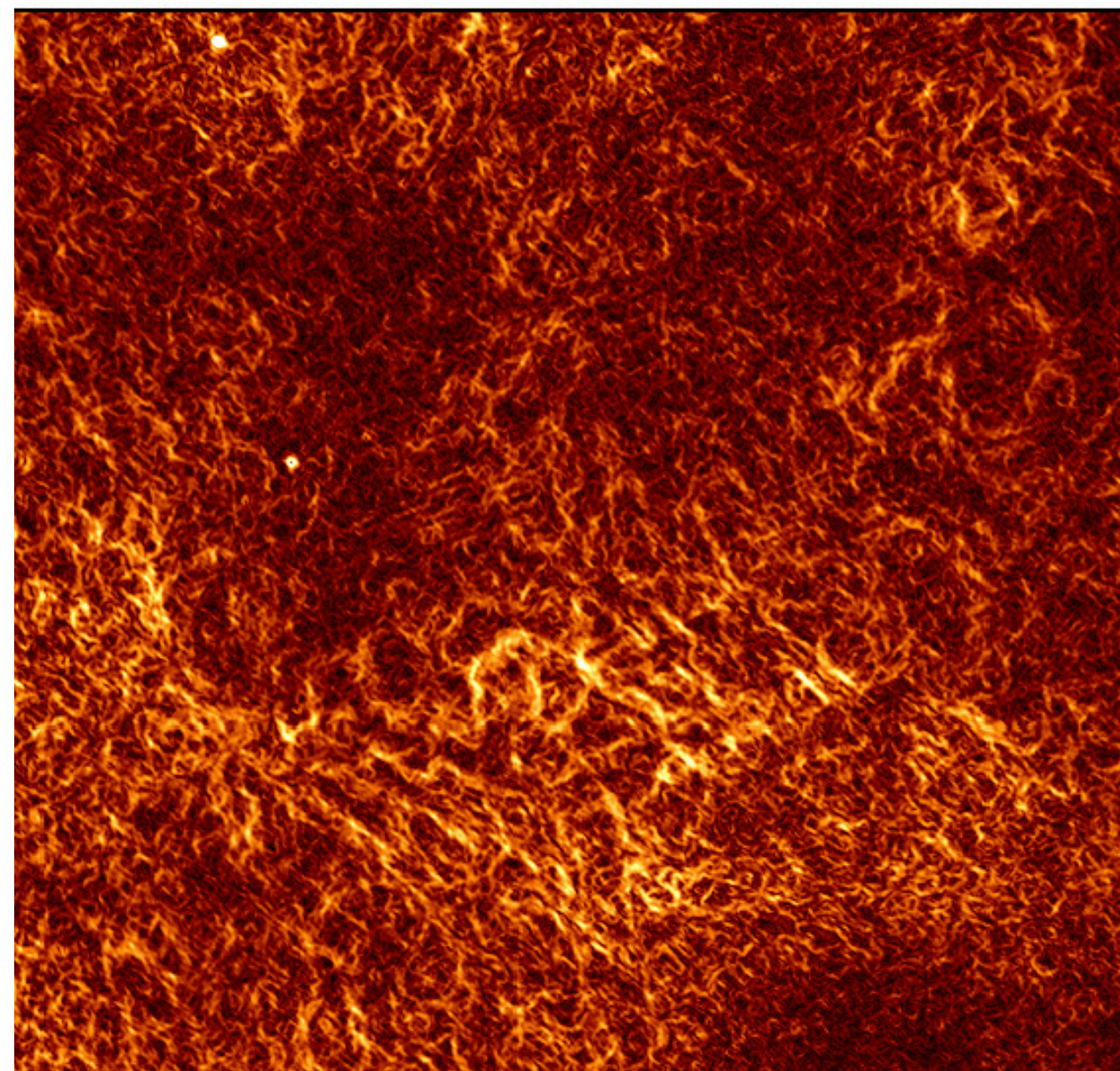


Sample EoR signal with noise added



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

