SKAO

# SKA SWG Update

Robert Braun, SKAO Science Director 21 June 2022

### **SKA Science Update**

- New SWG co-chairs
- Science Data Challenge progress (Philippa, Anna, Simon)
- SWG Banners (Tyler)
- AOB



### **New SWG Co-Chairs**

- Cradle of Life:
  - Welcome to Cherry Ng (University of Toronto)
  - Thanks to Laurent Lamy!
- HI Galaxies:
  - Welcome to Betsy Adams (ASTRON)
  - Thanks to Paolo Serra!



### **SDC2 results paper**

- Results and analysis from SDC2 now in preparation for submission to MNRAS
- 12 finalist teams from over 40 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

#### SKA Science Data Challenge 2: analysis and results

P. Hartley, A. Bonaldi, R. Braun, [Order TBC:] D. Cornu<sup>1</sup>, B. Semelin<sup>1</sup>, X. Lu<sup>1</sup>, S. Aicardi<sup>2</sup>, P. Salomé<sup>1</sup>, A. Marchal<sup>3</sup>, J. Freundlich<sup>4</sup>, F. Combes<sup>1,5</sup>, C. Tasse<sup>6,7</sup>, C. Heneka<sup>8</sup>, M. Delli Veneri<sup>9</sup>, A. Soroka<sup>10</sup>, F. Gubanov<sup>10</sup>, A. Meshcheryakov<sup>11</sup>, B. Fraga<sup>12</sup>, C.R. Bom<sup>12</sup>, M. Brüggen<sup>8</sup>, A. K. Shaw<sup>13</sup>, N. Patra<sup>14</sup>, A. Chakraborty<sup>15</sup>, R. Mondal<sup>16</sup>, S. Choudhuri<sup>17</sup>, A. Mazumder<sup>15</sup>, M. Jagannath<sup>18</sup>, M. J. Hardcastle<sup>19</sup>, J. Forbrich<sup>19</sup>, L. Smith<sup>20</sup>, V. Stolyarov<sup>20,21</sup>, M. Ashdown<sup>20</sup>, J. Coles<sup>20</sup>, H. Håkansson<sup>22</sup>, A. Sjöberg<sup>22</sup>, M. C. Toribio<sup>23</sup>, M. Önnheim<sup>22</sup>, M. Olberg<sup>23</sup>, E. Gustavsson<sup>22</sup>, M. Lindqvist<sup>23</sup>, M. Jirstrand<sup>22</sup>, J. Conway<sup>23</sup> K. M. Hess<sup>24,25,26</sup>, R. J. Jurek<sup>27</sup>, S. Kitaeff<sup>28</sup>, P. Serra<sup>29</sup>, A. X. Shen<sup>30,31</sup>, J. M. van der Hulst<sup>25</sup>, T. Westmeier<sup>28</sup>, A. Alberdi<sup>33</sup>, J. Cannon<sup>34</sup>, L. Darriba<sup>33</sup>, J. Garrido<sup>33</sup>, J. Gósza<sup>35</sup>, D. Herranz<sup>36</sup>, M. G. Jones<sup>37</sup>, P. Kamphuis<sup>38</sup>, D. Kleiner<sup>29</sup>, I. Márquez<sup>33</sup>, J. Moldón<sup>33</sup>, M. Pandey-Pommier<sup>39</sup>. M. Parra<sup>33</sup>, J. Sabater<sup>40</sup>, S. Sánchez<sup>33</sup>, A. Sorgho<sup>33</sup>, L. Verdes-Montenegro<sup>33</sup>, G. Fourestey<sup>41</sup>, A. Galan<sup>41</sup>, C, Gheller<sup>29</sup>, D, Korber<sup>41</sup>, A, Peel<sup>41</sup>, M, Sargent<sup>41</sup>, E, Tollev<sup>41</sup>, B, Liu<sup>42</sup>, R, Chen<sup>42</sup>, B, Peng<sup>42</sup>, L, Yu<sup>42</sup>, H. Xi<sup>42</sup>, K. Yu<sup>43</sup>, O. Guo<sup>43</sup>, W. Pei<sup>43</sup>, Y. Liu<sup>43</sup>, Y. Wang<sup>43</sup>, X. Chen<sup>43</sup>, X. Zhang<sup>44</sup>, S. Ni<sup>44</sup>, J. Zhang<sup>44</sup>, L. Gao<sup>44</sup>, M. Zhao<sup>44</sup>, L. Zhang<sup>45</sup>, H. Zhang<sup>45</sup>, X. Wang<sup>45</sup>, J. Ding<sup>45</sup>, S. Zuo<sup>46</sup>, Y. Mao<sup>46</sup>, A. Vafaei Sadr<sup>47</sup>, M. Kunz<sup>47</sup>, B. Bassett<sup>48</sup>, D. Crichton<sup>49</sup>, V. Nistane<sup>47</sup>, N. Oozeer<sup>35</sup>, S. Jaiswal<sup>50</sup>, B. Lao<sup>50</sup>, J. N. H. S. Aditya<sup>50</sup>, Y. Zhang<sup>50</sup>, A. Wang<sup>50</sup>, and X. Yang<sup>50</sup> 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 unrivalled 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 neutral hydrogen (HI) sources in a simulated data product representing a 2000 h SKA MID spectral line observation. Through the generous support of eight international supercomputing facilities, participants were able to undertake the Challenge using dedicated computational resources. This model not only supported the accessible provision of a realistically large dataset, but also provided the opportunity to test several aspects of the future SKA Regional Centre network. Sitting alongside the main challenge, 'reproducibility awards' were made in recognition of those pipelines which demonstrated Open Science best practice. The Challenge saw over 100 finalists develop a range of new and existing techniques, in results which highlight the strengths of multidisciplinary and collaborative effort. The winning strategy - combining predictions from two independent machine learning techniques - 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.



### **SDC2 signal-to-noise analysis**

- Expressing SDC2 outcomes in terms of source signalto-noise values
- Challenging to define a meaningful measure of integrated signal-to-noise, since it is intimately tied to the quality of a "matched filter"
- SKA noise properties unlike current telescopes
  - RMS noise remains ~constant when spatially smoothing between range of about 0.4 to 100 arcsec FWHM (at 1.2 GHz)
- Possible implications for source finding approaches
  - Need to include accurate noise models to optimize detection strategy
- See SDC2 paper for discussion



## Science Data Challenge 3 (SDC3)

- SDC3 will consist of two tiers:
  - SDC3 "foregrounds" (SDC3a)
    - Simulations for the challenge are maturing, and we are aiming for a challenge **start date of October 2022,** for a duration of six months
    - Computational costs are likely to vary significantly according to approaches (see later slides)
  - SDC3 "inference" (SDC3b)
    - SDC3 inference will be run after a short break following the conclusion of "foregrounds", so it will take place **during 2023.**
    - Again, computational costs are likely to vary between teams
- Working with members of EoR SWG to develop both tiers
- **SKAO** is developing a website dedicated to SDC3, with Challenge registration due to open soon



### **Status of SDC3 Foregrounds data products**

- The datasets that will be provided to participants are reaching maturity
- Currently under review by EoR specialists
- Positive feedback from review will allow us to confirm target start date of October 2022
- Data products that will be available to participants:
  - 3 GB image cube
  - 1 TB visibility set



### SDC3 foregrounds dataset feedback

• Preliminary dataset provided to a few EoR SWG experts for feedback



### **Foreground All-scale Radio Modeller (FARM)**

"FARM adapts SKAObserve script into a python-based, modular, reusable, and extensible CLI- or GUI-based foreground-simulation tool"

- Can use different telescopes
- Executions to simulate the foreground use a configuration file to specify parameters
- Can 'switch on/off' different effects/foreground components
- Specify calibration errors/effects
- Can generate some components from scratch, or uses models

#### [directories]

root\_name = "example\_1" # root name for all output files (not directory)
telescope\_model = "/Users/SDC3/Data/TelescopeModels/V512/telescope.tm"
output\_dcy = "/Users/SDC3/Desktop/test\_output\_farm"

```
# observation not yet incorporated in to FARM
[observation]
time = 2027-01-01T10:00:01 # Start time (UTC) YYYY-MM-DDTHH:MM:SS
t_total = 600 # Total on source time [s]
n_scan = 1 # Number of observational 'scans' of target field
min_gap_scan = 1200 # Minimum time-gap between consecutive scans [s]
min_elevation = 30.0 # Minimum elevation [deg]
```

#### [observation.field]

```
ra0 = "00:00:00.0" # HH:MM:SS.S
dec0 = "-30:00:00.00" # DD:MM:SS.SS
frame = "fk5" # Astropy-compatible frame string e.g. fk4, fk5, galactic etc
nxpix = 1024 # Number of pixels in x/right ascension
nypix = 1024 # Number of pixels in y/declination
cdelt = 28.125 # arcsec
```

```
[observation.correlator]
freq_min = 115e6 # Hz
freq_max = 169e6 # Hz
nchan = 217 # Number of evenly spaced channels across whole bandwidth, int
chanwidth = 250e3 # channel bandwidth for smearing calculation, Hz
t_int = 60.0 # visibility integration time, s
```





Slide / 10

### **SDC Computational support model**

- SDC2 received invaluable support from international computing facilities
  - Enabled SKAO to provide a 1 TB data product in an accessible way
  - Test cases for SKA Regional Centre (SRC) prototyping
- SDC3 will also receive computational support for participating teams
  - Computational cost estimates of Foregrounds analysis are currently being finalized
  - A new registration and time allocation model
    - Streamlining the computational resource allocation process and help teams to get the most from their resources



### **SDC3 Computational costs**



### **SDC3 Computational costs: current estimates**

with thanks to EoR SWG members for estimates

- We expect ~10 teams to take part in SDC3a
- We expect that some teams will deal only with the smaller (~3 GB) image dataset
  - Estimated compute cost: ~ <100 CPUh per full pipeline run; <64 GB RAM
  - i.e. ~3 hours on 30 core VM
- Some teams will rely on the larger (~1 TB) visibilities dataset (or both)
  - Estimated compute cost: cost dominated by one step of the processing <10k CPUh; <128 GB RAM.
  - i.e. 20 days in 30 core VM
  - This step is not likely to be run many times



### **SDC3 registration**

As part of SDC3 registration participants will be invited to submit a short proposal, which:

- Allows us to assess teams' software and hardware requirements
- Enables us to match teams to HPC facilities (who themselves have provided us with the month-by-month availability of their facilities)
- Facilitates more effective use of HPC facilities



### **SDC3 Registration**

- Registration will open soon
- Via SDC3 website (under construction)





### **SDC3 computing support allocation**

All participants assigned to HPC-facility. Asked to conduct a 'mock' run on server in month 1, then given permanent account on that facility for remainder of the challenge



### **SRC prototyping activities**

- The data challenges provide the opportunity to test some of techniques being prototyped for future SKA Regional Centre (SRC) Network
- The SRC development work is now being conducted within the 'SAFe' methodology
  - SDC development work could interface via raising features or by identifying prototype products that would be useful to test during a science data challenge
- We would be cautious to ensure that any prototype testing would happen at a suitable level of maturity of a prototyped product, and within a timescale that suits all supporting facilities.

### • No major prototyping planned for SDC3a

• Possibly testing by participants of e.g. image viewing software (CARTA)



### **Reproducibility, Open Science and best practice**

- The data challenges provide the opportunity to familiarise the science community with some of the best practices concerning software and data processing
- Option to again include '<u>reproducibility awards</u>' alongside main challenge:
  - Possibly revise format of award
  - For example: can SKAO team re-run a team's pipeline?
- Idea of an environmental award to run alongside the main challenge





Slide / 19

## **Any Other Business**

- Upcoming meetings
  - EAS2022 "S7: Building bridges: The lifecycle of dust and gas in the Milky Way with ALMA and SKA", 27 June – 1 July (<u>https://eas.unige.ch/EAS\_meeting/session.jsp?id=S7</u>)
  - EAS2022 "SS23: Towards the SKA Observatory: Artificial Intelligence in radio astronomy", 27 June 1 July,

(<u>https://eas.unige.ch/EAS\_meeting/session.jsp?id=SS23)</u>

- EAS2022 "SS5: Neutral hydrogen: the next generation of science and simulations", 27 June – 1 July, (<u>https://eas.unige.ch/EAS2022/session.jsp?id=SS5</u>)
- "Coordinated Surveys of the Southern Sky", 10 14 October, (https://www.cadc-ccda.hia-iha.nrc-cnrc.gc.ca/en/meetings/getMeetings.html?number=6792)

• ...

• Other Points?



We recognise and acknowledge the Indigenous peoples and cultures that have traditionally lived on the lands on which our facilities are located.  $\bullet$ 



• •

•

 $\bullet$ 

•

 $\bullet$