# SKA data challenges for SRCs



A. Bonaldi
Project scientist

**SQUARE KILOMETRE ARRAY**

Exploring the Universe with the world's largest radio telescope

# The SKA's data journey

- ## Data flow challenges

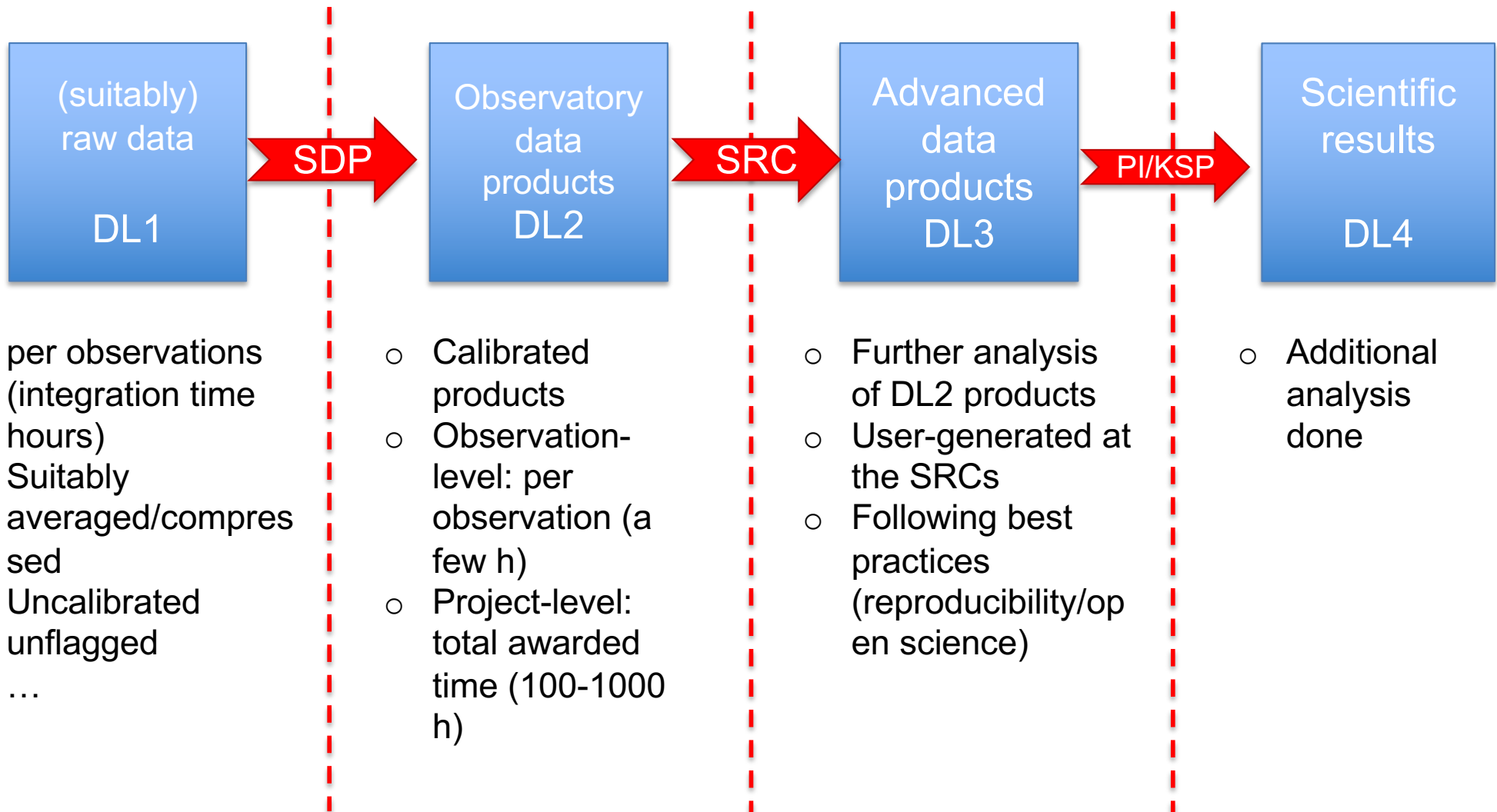# 50 x data rate reduction by Science Data Processors



- But, have a series of buffers
- *Raw voltage data can be stored for about 2 minutes*

- *Raw visibility data can be stored for about 2 weeks*

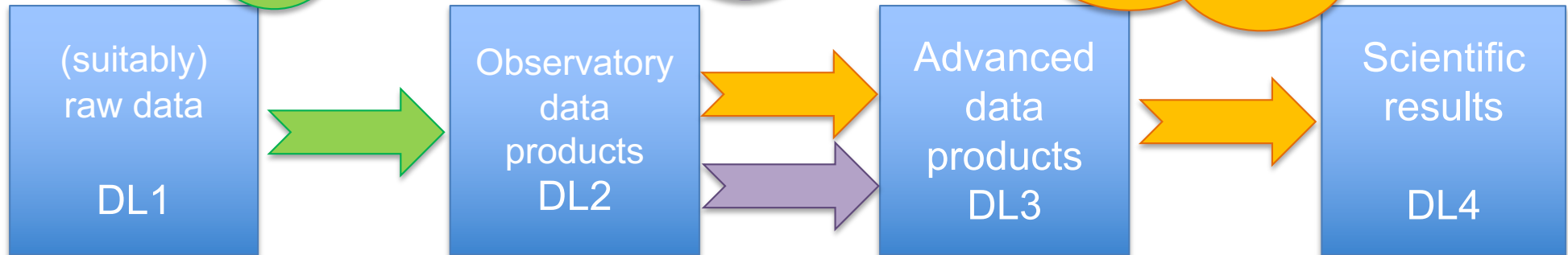- *Final data products will be stored forever*

# SKA data challenges: Data Layers (DL)

| (suitably) raw data  DL1 | → SDP → | Observatory data products DL2 | → SRC → | Advanced data products DL3 | → PI/KSP → | Scientific results  DL4 |
|---|---|---|---|---|---|---|

- o per observations (integration time hours)
- o Suitably averaged/compressed
- o Uncalibrated
- o unflagged
- o …

- o Calibrated products
- o Observation-level: per observation (a few h)
- o Project-level: total awarded time (100-1000 h)

- o Further analysis of DL2 products
- o User-generated at the SRCs
- o Following best practices (reproducibility/open science)

- o Additional analysis done

SDP de-risking "Big data" challenge

SRC data challenges

SWG activities "Science data" challenge

(suitably) raw data

DL1

Observatory data products DL2

Advanced data products DL3

Scientific results

DL4

- Focus on efficiency/scalability
- Calibration strategy and implementation
- Data size 10-100 TBs

- Focus on algorithm development
- Proposal-Specific processing
- Data size few TBs

- Technical: data movement, format, protocols, security, databases
- Algorithms - Best practices

Exploring the Universe with the world's largest radio telescope

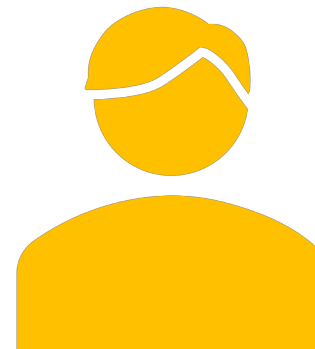Data challenges
Coordination group

SDP challenges    SRC challenges    Science challenges
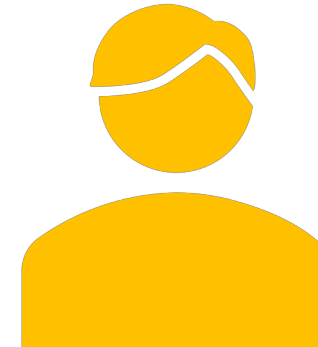
# Data challenges
## Coordination group

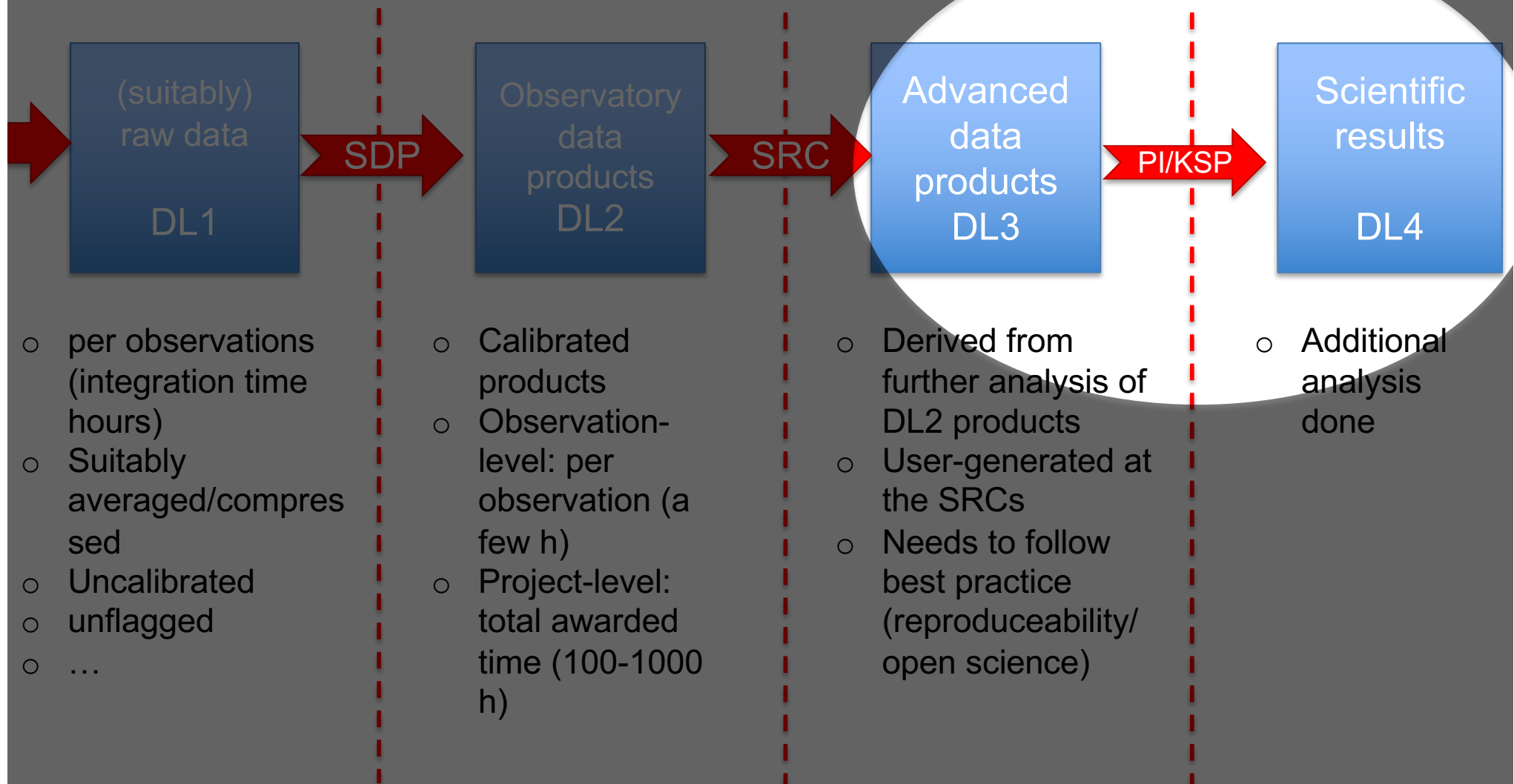SDP challenges        SRC challenges        Science challenges

- Chair: R. Bolton
- Priorities and strategies from the SRC community
- Coordination through SRCSC

# SKA data challenges: Data Layers (DL)

DL0=raw data out of CSP

| (suitably) raw data DL1 | SDP → | Observatory data products DL2 | SRC → | Advanced data products DL3 | PI/KSP → | Scientific results DL4 |
|---|---|---|---|---|---|---|

- per observations (integration time hours)
- Suitably averaged/compressed
- Uncalibrated
- unflagged
- …

- Calibrated products
- Observation-level: per observation (a few h)
- Project-level: total awarded time (100-1000 h)

- Derived from further analysis of DL2 products
- User-generated at the SRCs
- Needs to follow best practice (reproduceability/open science)
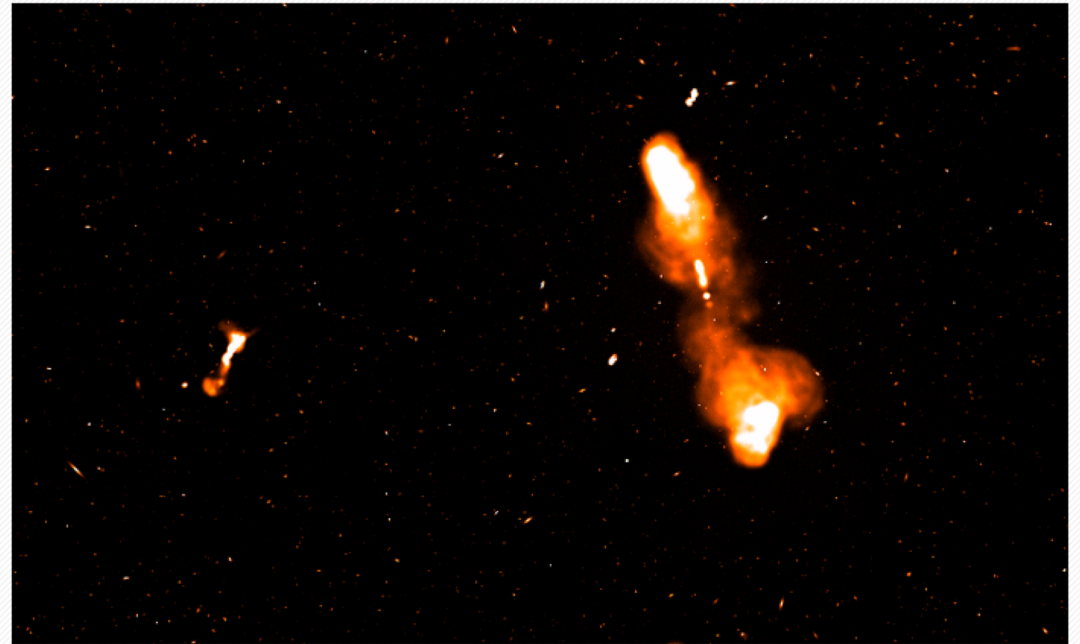
- Additional analysis done

# Science data challenge 1 (SDC1)

- Science-ready (SRC) imaging product

- Radio continuum, SKA Mid

- Not too challenging data sizes

- 1 pointing, 3 freqs, 3 depths

- Source finding

- Source identification, classification & characterization

Home » Latest News » SKA launches first Science Data Challenge for astronomy community

Print this page

## SKA Launches First Science Data Challenge For Astronomy Community

A snapshot from the SKA Science Data Challenge image, showing a large Active Galactic Nucleus (AGN) as if observed by SKA-mid at 1.4 GHz. (Credit: SKA Organisation)

**SKA Global Headquarters, 26 November 2018** – The Square Kilometre Array Organisation (SKAO) is today releasing its first ever Science Data Challenge, giving astronomers a taste of the highly detailed images the SKA will produce.

Developed by the SKAO's Project Science team, the challenge requires the analysis of a series of high resolution images created through data simulations. Researchers are invited to download the images and use their own software to find, identify and classify the sources.

The key aim of the series of Data Challenges is to prepare the science community for the kind of data products they will receive from SKA observations, and to gather valuable feedback which will inform the development of data reduction procedures.
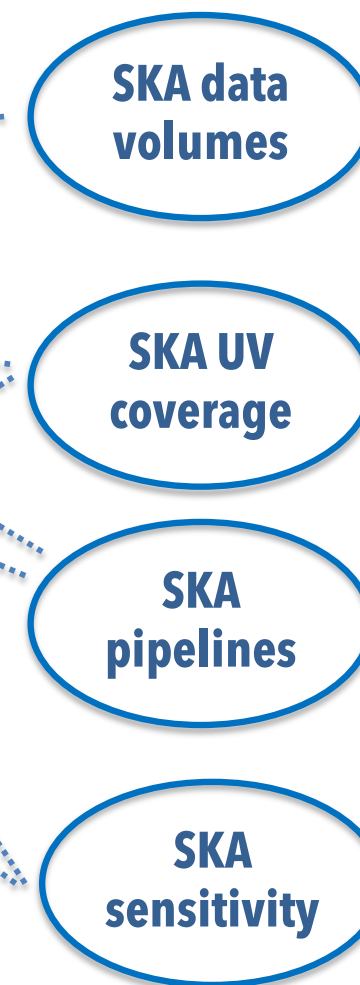
# Science data challenge 1 (SDC1)

❖ SKA unique features map into the data products:

  ✓ In the image plane, not visibilities

  ✓ "Benign" dirty beam

  ✓ Deconvolved down to 8h exposures

  ✓ Very deep -> confusion limited

  ✓ Very big number of sources to detect and classify

❖ SDC1 goals:

  ✓ Get the community familiar with this data product

  ✓ Develop efficient methods for source finding and source characterization

  ✓ Begin talking about best practices for deploying pipelines @ SRCs

**SKA data volumes**

**SKA UV coverage**

**SKA pipelines**

**SKA sensitivity**

## Data

| | | |
|---|---|---|
| 560 MHz, 8 hours | 4 Gb | DOWNLOAD |
| 560 MHz, 100 hours | 4 Gb | DOWNLOAD |
| 560 MHz, 1000 hours | 4 Gb | DOWNLOAD |

B1

| | | |
|---|---|---|
| 1400 MHz, 8 hours | 4 Gb | DOWNLOAD |
| 1400 MHz, 100 hours | 4 Gb | DOWNLOAD |
| 1400 MHz, 1000 hours | 4 Gb | DOWNLOAD |

B2

| | | |
|---|---|---|
| 9200 MHz, 8 hours | 4 Gb | DOWNLOAD |
| 9200 MHz, 100 hours | 4 Gb | DOWNLOAD |
| 9200 MHz, 1000 hours | 4 Gb | DOWNLOAD |

B5

- Short
- Medium
- Long

## Ancillary data

| | | |
|---|---|---|
| 560 MHz, primary beam | 300 Kb | DOWNLOAD |
| 560 MHz, synthesized | 4 Gb | DOWNLOAD |
| 1400 MHz, primary beam | 300 Kb | DOWNLOAD |
| 1400 MHz, synthesized | 4 Gb | DOWNLOAD |
| 9200 MHz, primary beam | 300 Kb | DOWNLOAD |
| 9200 MHz, synthesized | 4 Gb | DOWNLOAD |

Data access: from

**https://astronomers.skatelescope.org/**

Data reside on the Italian Center for Astronomical Archive (IA2) operated by INAF

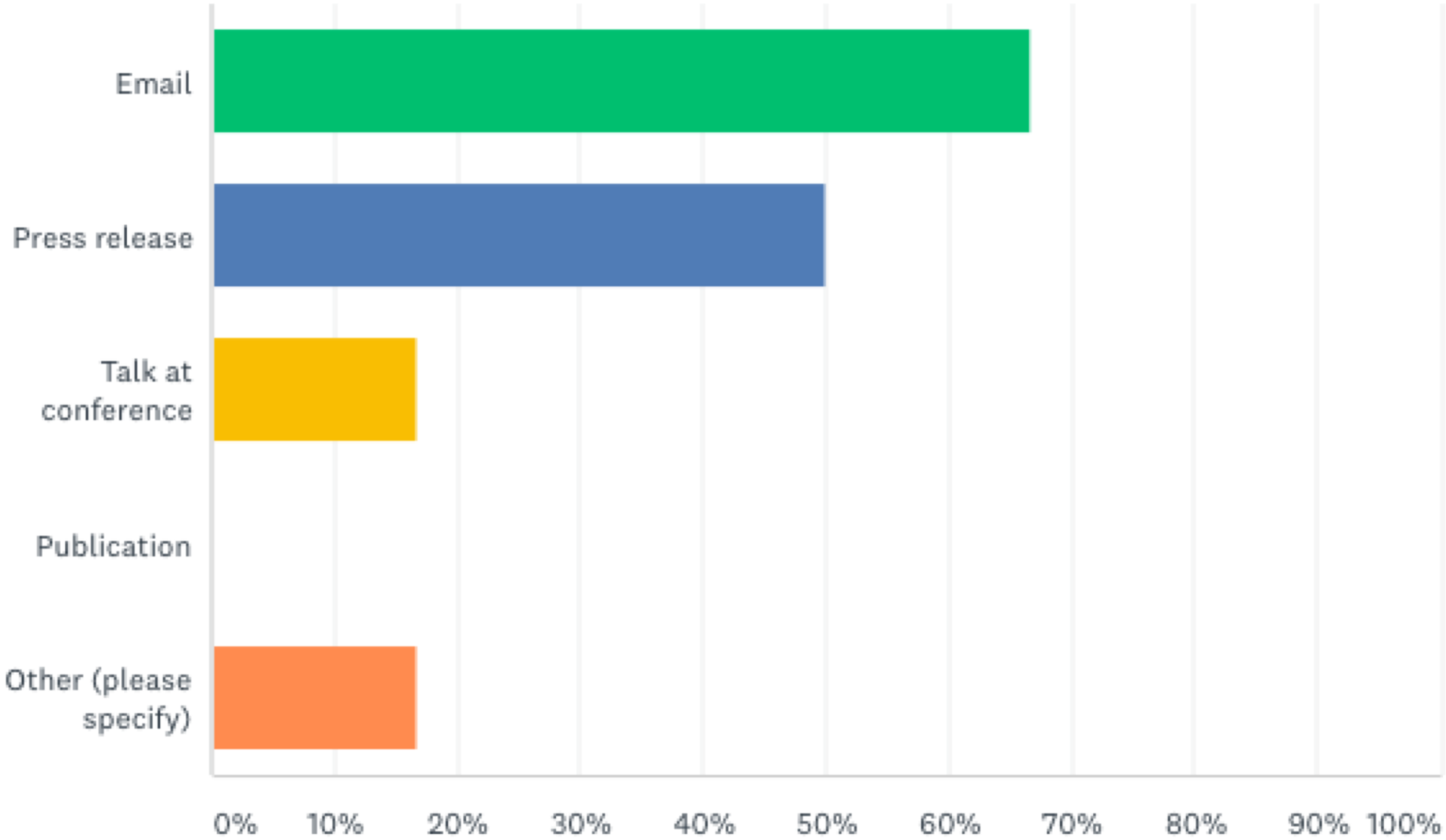## Training set

| | | |
|---|---|---|
| 560 MHz, truth catalogue | 54 Mb | DOWNLOAD |
| 1400 MHz, truth catalogue | 14 Mb | DOWNLOAD |
| 9200 MHz, truth catalogue | 340 Kb | DOWNLOAD |

Truth table for a 5% sky area: training set

# SDC1 communication strategy

ANNA BONALDI & ROBERT BRAUN, FOR THE SKAO SCIENCE TEAM *

SKA Organization, Jodrell Bank, Lower Withington, Macclesfield, Cheshire, SK11 9DL, United Kingdom

# The SDC1 teams!

17 teams registered to SDC1

# The SDC1 teams!

9 teams submitted results by the deadline of 30th April
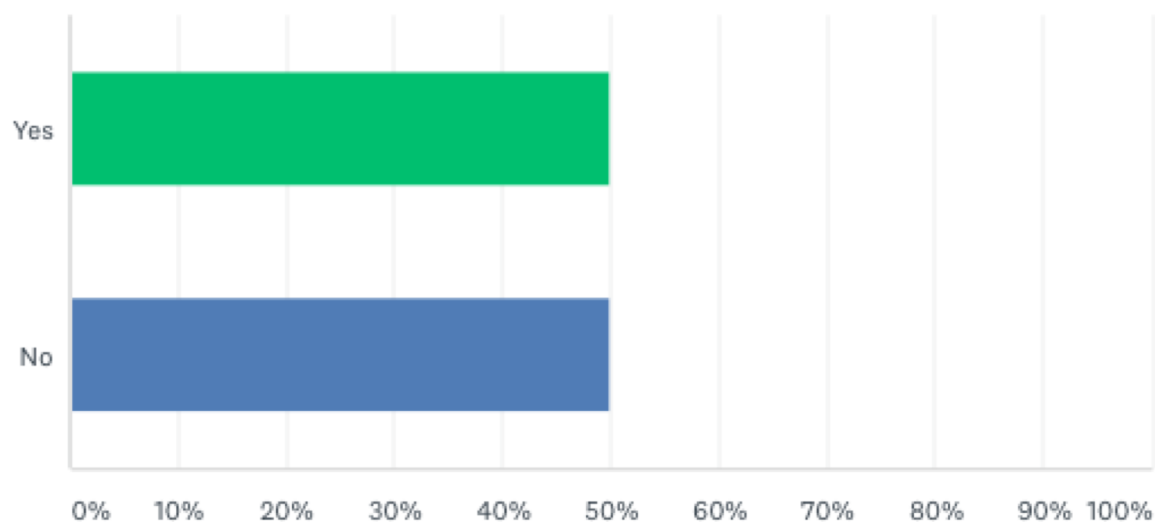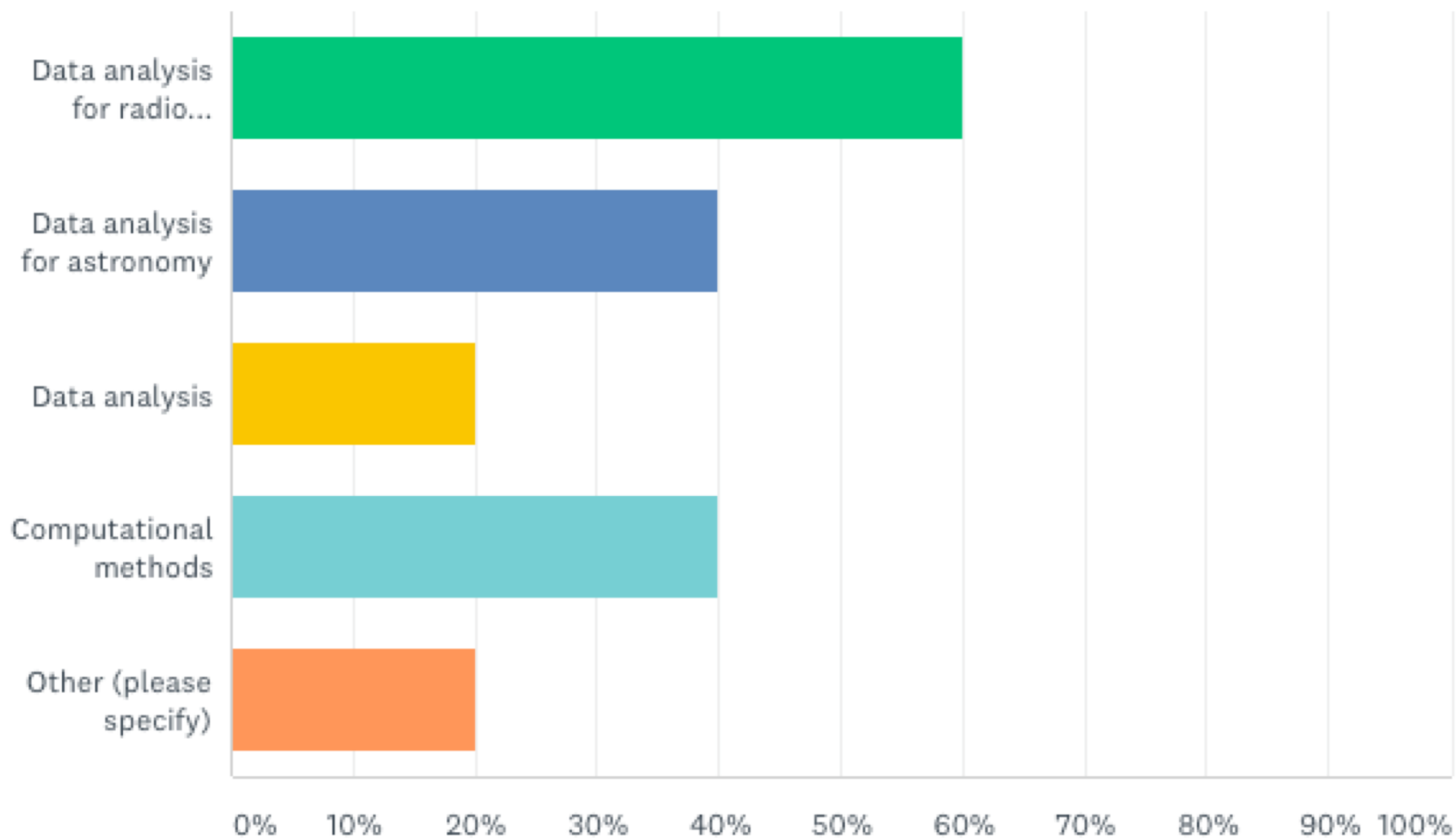
# Team's provenance

Are there people in your team who are formally affiliated with an SKA Science Working Group?

Answered: 6    Skipped: 0

# Team's expertise

# SDC1 Leaderboard

**Ox-ICRAR**
**UK / Australia**

Catherine Hale, Aaron Robotham, Matt Jarvis

*The team acknowledges STFC and Oxford Hintze Centre for Astrophysical Surveys (funded through the Hintze Family Charitable Foundation), ICRAR (for ASVR)*

**ICRAR**
**Australia**

Chen Wu, Ivy Wong

Richard Dodson, Ellie Gholami

*The team acknowledges CSIRO/Bracewell, ICRAR/Pleiades*

**EngageSKA**
**Portugal**

Bruno Coelho

Domingos Barbosa, Sonia Antón, Jorge Bruno Morgado, Valério A. R. M. Ribeiro, Dzianis Bartashevych, João Paulo Barraca, Miguel Bergano, Dalmiro Maia

*R.C. acknowledges support from the Advanced EU Network of E-infrastructures for Astronomy with the SKA (AENEAS), funded by the European Commission Framework Horizon 2020. RIA under grant agreement n. 731016 and from the ENGAGE SKA grant POCI-01-0145-FEDER-022217 funded by COMPETE 2020 and FCT, Portugal. S.A. & B.R. acknowledge support from the FCT Investigator through the exploratory project of reference IF/01051/2014. The team acknowledges support from project UID/04434/2019 funded by FCT. The team has made use of the ENGAGE SKA Compute Cluster. We acknowledge support from FEDER grant POCI-01-0145-FEDER-022217 funded by COMPETE 2020 and FCT, Portugal.*

**ARCIt-CACAO**
**Italy**

Sandra Burkutean, Marcella Massardi, Jan Brand, Elisabetta Liuzzo, Matteo Bonato

Eugenio Schisano, Andrea Giannetti, Kazi Rygl, Nicola Marchili

| **hs** **Germany** | **IITK** **India** | **IPM** **Iran** | **JLRAT** **China** | **Shanghai** **China** |
|---|---|---|---|---|
| Vesna Lukic, Marcus Brüggen | Pankaj Jain, Pritpal Kaur, Mohit Panwar Prabhakar | Zahra Bagheri, Hadis Goodarzi, Elham Saremi, Somayeh Sheikhnezami, Vajiheh Sabz Ali | Lei Yu, Bin Liu | Tao An, Sumit Jaiswal, Yang Lu, Prashanth Mohan, Baoqiang Lao |
| *the team acknowledges the use of the LSC2 cluster at University of Hamburg* | | *The team acknowledges funding to Institute for Research in Fundamental Sciences (IPM)* | | *National Key R&D Programme of China under grant number 2018YFA0404603* |

Exp

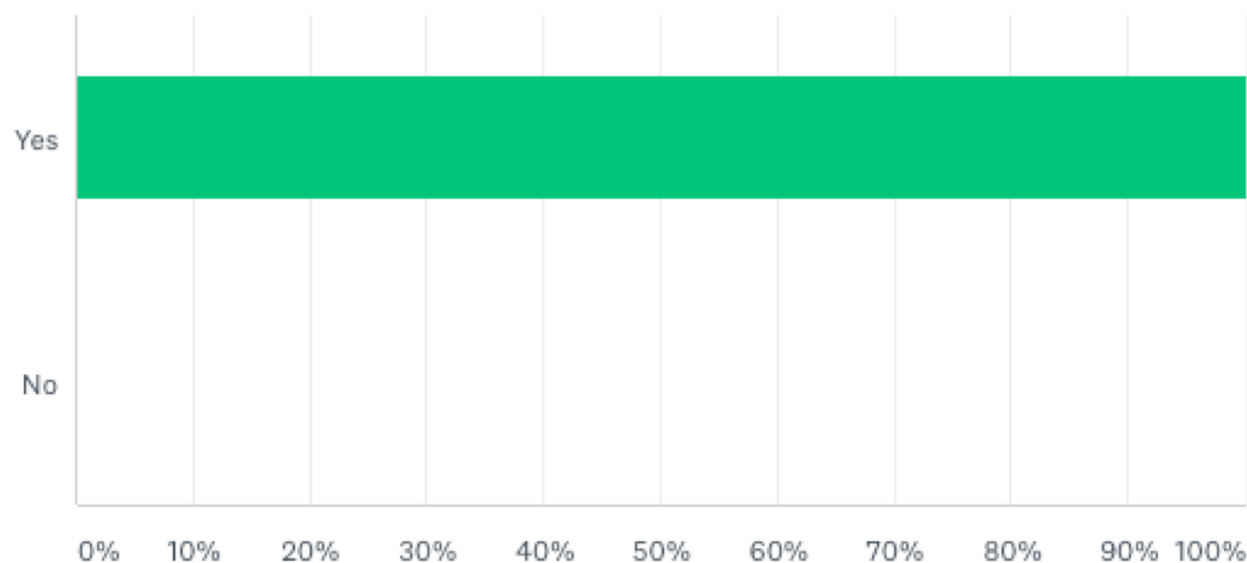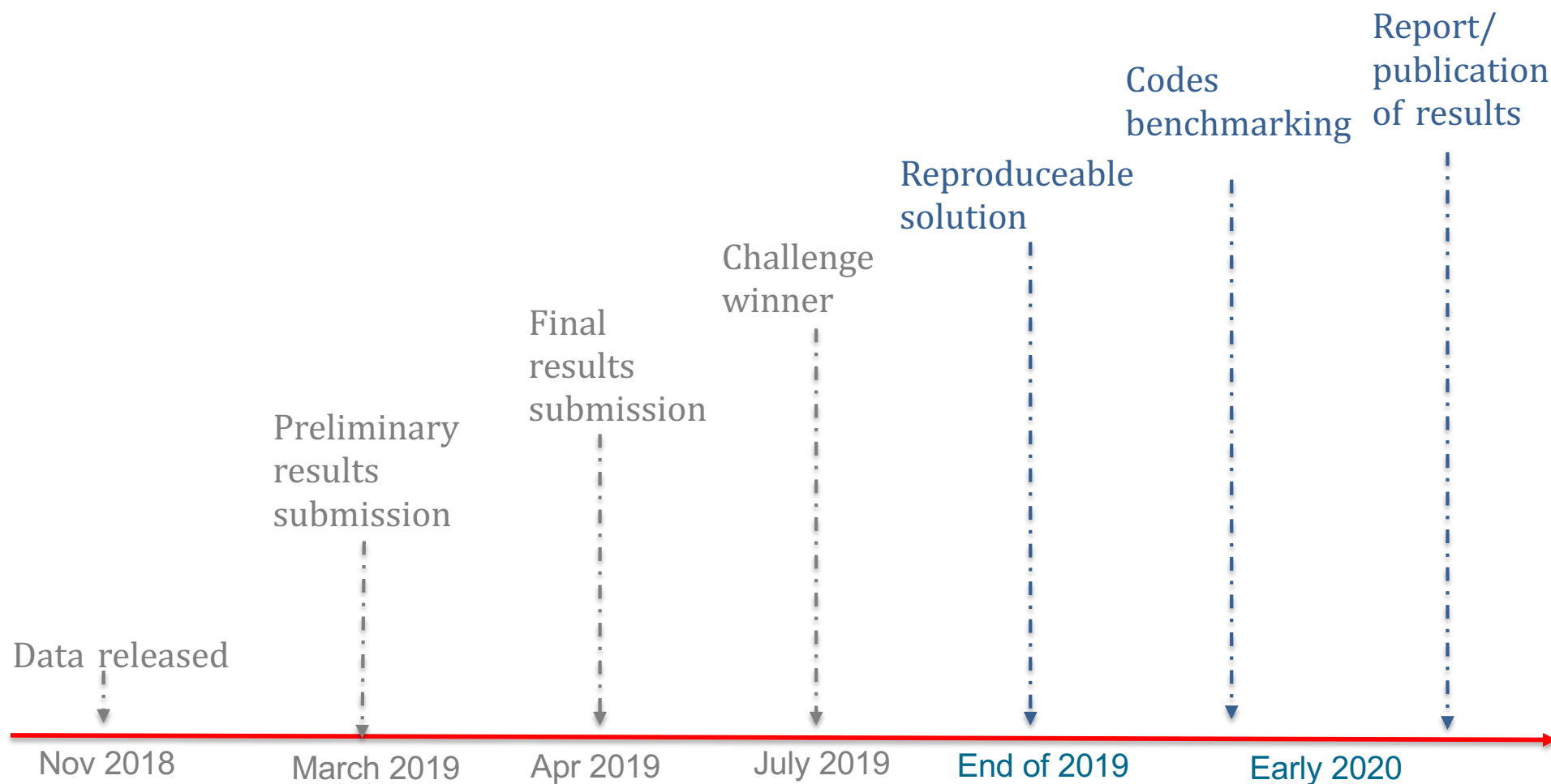# Positive feedback!

Based on your experience with SDC1 would you consider participating in future SKA science data challenges?

Answered: 5    Skipped: 0

# SDC1 timeline and progress



Report/
publication
of results

Codes
benchmarking

Reproduceable
solution

Challenge
winner

Final
results
submission

Preliminary
results
submission

Data released

Nov 2018     March 2019     Apr 2019     July 2019     End of 2019     Early 2020

# Collaborative phase

1  {
2    "cells": [
3      {
4        "cell_type": "code",
5        "execution_count": null,

# Best practices for SDC1

- Open

  - Is your pipeline publicly available (e.g. Github?)

  - How is it licensed?

  - Which licenses/dependencies does it need?

  - Is there a documentation?

- Reproduceable

  - Can you containerize your pipeline?

    - Info on containers circulated

    - Support offered to move towards reproducibility

# Thanks!

# SDP list of observation-level products

- Science Alert Catalogue

- Transient Source Catalogue

- Science Product Catalogue Data Product

- Image Products 1: Image Cubes

- Image Products 2: Uvgrids

- Calibrated Visibilities

- Sieved Pulsar and Transient Candidates

- Pulsar Timing Solutions

- Dynamic spectrum data

- Transient Buffer Data

- Science Data Model

Observation-level products will be combined into project-level products.

Added-value products will be derived from the Observatory data products