Computing Requirements for SKA Science Data Challenges
SKAO Science Team
15 August 2019

Science Data Challenges (SDC) represent a subset of the broader range of SKA Data Challenges that are outlined in SKA-TEL-SKO-00001016. The primary motivation here is to provide increasingly realistic data products to the scientific community. The first SDC (hereafter SDC1) was released in November 2018 and is described in SKA-TEL-SKO-00001001 while the analysis of submissions to SDC1 is described in SKA-TEL-SKO-00001081.

SDC1 makes use of a state-of-the-art model (Bonaldi et al, 2019 MNRAS 482, 2) of extragalactic continuum sources that is defined over the full frequency range of the SKA. Simulated "full-field" images (about 2x the primary beam FWHM) at moderately high angular resolution (0.6 FWHM arcsec at 1400 MHz) were generated at three central frequencies within the coverage of SKA1-Mid (560, 1400 and 9200 MHz) each sampled with 30% fractional bandwidth. This was done with image sizes of 32768 pixels on a side, so 4.3 GByte (a factor of 2 to 4 less than would be ideal from the perspective of beam sampling and sky coverage). Populating these simulated sky images to a depth suitable to represent what could be seen in very deep ($1000^h$) integrations with the SKA sensitivity required introducing some $10^7$ discrete sources. A full-size SKA1-Mid database (an eight-hour track with 133+64 antennas, 4 polarisations, together with time and frequency sampling that keep smearing effects suitably low at the edge of the primary beam $\Delta t = 0.14s$ and $(\Delta \nu/\nu) = 10^{-5}$) has $10^{14}$ visibilities and occupies some 2.4 PByte for a single precision representation. (In the case of SKA1-Low, a 4-hour track and $\Delta t = 0.8s$ and $(\Delta \nu/\nu) = 6 \; 10^{-5}$ yield database sizes about 10 times smaller.) This is much larger than could be simulated directly with current resources. The (u,v) coverage of such a track was instead approximated with 1 polarisation and very crude time and frequency sampling of $\Delta t = 30s$ and $(\Delta \nu/\nu) = 0.0033$, yielding a factor of 283,000 reduction in data volume to only 8.4 GByte. This very crude sampling implied that the source model for the sky could not be added to the visibilities prior to imaging but instead that the source model could only be introduced into the simulation in the image plane. The $10^7$ discrete sources in the sky model catalogue were "painted" on to an image grid for each observing frequency using a Gaussian smoothing kernel to ensure adequate representation of the most compact structures. The "dirty" synthesized beam from the simulated (u,v) coverage was used to convolve the "residual" sky brightness (that part fainter than 3 times the expected RMS noise level of an 8 hour observation) and this was added to the "restored" image consisting of the brightest parts of the gridded sky brightness convolved with a Gaussian restoring beam. Both types of convolution were preceded by a linear deconvolution to undo the smearing of the Gaussian smoothing kernel that was used for sky model gridding. For SDC1, no random or systematic calibration errors were introduced, but only Gaussian thermal noise consistent with the (u,v) coverage and of the expected amplitude.

SDC1 was undertaken with the following resources:
1) SKA Science1 server: 2 x 18 core Xeon E5-2695 CPUs, 256 GB RAM, 22 TB disk RAID
2) RB Laptop: 6 core Intel i9, 32GB RAM, 4TB SSD
3) UnivMan server: 32 core Xeon E5-2640 CPUs, 126 GB RAM, various storage servers

Several weeks run time on each of these platforms (to complete multiple debugging and tuning iterations) was used to complete SDC1. What is vital is rapid turn-around (minutes to hours) of debugging iterations to allow efficient tuning and testing. Final run-times should not exceed days in order to stay practical. Debugging requires high bandwidth communications to allow both graphical and deep dive numerical investigation of each iteration.

For the next SDCs, we aim to increase the level of realism in various ways. The time and frequency sampling of the (u,v) coverage need to improve significantly, so that it becomes possible to begin embedding the sky model directly into the visibilities prior to imaging. This will also permit introduction of some residual calibration errors into the simulations to provide a more representative data product for user analysis. We wish to probe the regime of database sizes of 100s of GB to about 1 TB, together with image sizes of typically 64k – 128k pixels per side (16 – 64 GB per image). This will be best-matched by machine memories that are 100's of GB coupled with extremely fast I/O to a staging storage system of some 10 TB (something like NVMe SSDs) and a longer term (slower) storage of some 100 TB (something like disk RAID). The network connection between the storage system and the desktop environment where debugging occurs needs to be very fast. Data volumes being regularly transferred between these environments in both directions are 100's of GB. We also need to begin more realistic simulations of SKA1-Low with its time-varying station beam response on the sky.

This would benefit greatly from the use of a system with integrated GPUs to allow the OSKAR package to run efficiently.

Summarising the desired specifications to enable effective future SDCs:
  1) 10s of CPU cores
  2) several GPUs
  3) 100s GB RAM
  4) 10 TB ultra-fast storage (NVMe SSD or equivalent)
  5) 100 TB longer term storage (disk RAID or equivalent)
  6) 1 – 10 Gb/s network connectivity between compute and analysis environments