# AN OVERVIEW OF THE SKA PROJECT: WHY TAKE ON THIS SIGNAL PROCESSING CHALLENGE?

*Steven J. Tingay*

International Centre for Radio Astronomy Research, Curtin University
Bentley, Western Australia, Australia, 6102

## ABSTRACT

The Square Kilometre Array (SKA) is an international project to build and operate the most powerful radio telescope ever conceived, in order to answer a range of fundamental questions in cosmology and astrophysics. The SKA will be a facility built in two phases, primarily in two host countries (Australia and South Africa), and in four technology elements (telescope systems). In this paper I briefly outline the science goals and technical challenges of the SKA. I will focus most attention on the Big Data signal processing challenges of the SKA, highlighting the need for the astronomy community to reach out into other disciplines such as the broad signal processing community, in order to obtain fresh and inter-disciplinary perspectives on these challenges.

***Index Terms***— Square Kilometre Array, Cosmology, Astrophysics, inter-disciplinary research

## 1. INTRODUCTION

Human-kind has for thousands of years sought to understand the Universe in which we live. In that time, humans have built ever more powerful instruments to aid in this quest, fuelled in recent centuries by exponential improvements in technology. Telescopes now comprise electronic and detector technologies that allow humans to observe objects across almost the entire approximately 14 billion year history of the Universe. We can look back in time and space, virtually to the Big Bang, the origin of all of space and time, matter and energy. Detector technology, advanced information and communications technologies, and our knowledge of physics allow us to ask questions regarding the nature, constituents, and history of the Universe, uncovering ever deeper versions of physics. A grand cycle can be demonstrated. Our understanding of physics leads to more advanced technologies. More advanced technologies lead to an ability to ask and answer deeper questions, leading to a better understanding of physics. This is a cycle measured in decades and centuries. The SKA [1][1] is one of the current set of facilities seeking to participate in this grand cycle.

The SKA is a telescope that will operate at radio wavelengths and, like all other modern astronomy instrumentation, is best considered as a multi-disciplinary program due, largely, to the technologies involved. However, the greatest benefits can be reaped from such projects if the multi-disciplinary approach can be maximally encouraged, not just limited to the utilisation of clever "widgets". Facilities such as the SKA, as will be seen in the next sections, produce so much data, requiring so much analysis, that it would probably be naive to expect that the world's several thousand radio astronomers have all the answers to the many analysis challenges we face. Thus, the motivation behind the special session convened at ICASSP in 2015 is to put the scale and scope of the SKA challenges in front of the international signal processing community, to look for diversity and innovation in finding answers to the questions we face.

The idea of an ICASSP special session emerged from a plenary talk I gave at the 2014 IEEE Workshop on Statistical Signal Processing[2]. I had a number of interesting conversations with Doug Gray, Victor Solo, Ba-Ngu Vo, Louis Scharf, and others that indicated a lot of cross-over interest in the SKA signal processing challenges. While there was a lot of interest, it was also clear that radio astronomers and signal processing engineers/mathematicians quite often had very different vocabularies to describe the same techniques and processes. Doug Gray proved to be a skilled interpreter between myself and others at the workshop. We concluded that we needed to continue the discussion between our communities and ICASSP 2015 was a perfect opportunity.

## 2. WHAT IS THE SKA AND HOW WILL IT WORK?

The SKA will be a radio telescope utilising interferometric techniques. Radio interferometers use arrays of individual antennas to synthesis a radio telescope with a sensitivity proportional to the sum of the individual antenna collecting areas and an angular resolution proportional to the maximum distance between individual antennas. In general it is good to have a large number of antennas in the array and arranged optimally for the scientific goals of interest.

---

[1] http://www.skatelescope.org

[2] http://www.ee.unimelb.edu.au/SSP2014/

The individual antennas can take a number of different forms, mostly dependant on the frequency of the radio waves relevant for particular observations. For example, at low radio frequencies (hundreds of MHz), simple dipole antennas may be used to directly intercept the cosmic radio waves. At higher frequencies (GHz and above), parabolic dish antennas are typically used to concentrate the radio waves at the focal point.

The signals from each antenna (a digital representation of the temporal variation in the band-limited signal in radio frequency, a voltage, at each antenna induced by the electric field caused by the radio waves) are transported to a central location and combined to form the fundamental observables of an interferometer, the co-called visibilities. Visibilities are the complex quantities produced when the voltage representations from each antenna are multiplied together (measuring the mutual coherence of the electric field between different pairs of spatial locations). Typically the band-limited signals are broken into different sub-bands (channelised) before multiplication, a process collectively referred to as correlation. For an array of $N$ individual antennas, $N(N-1)/2$ independent pairs of antennas (baselines) can be formed, with correlations produced for each baseline, for each sub-band. The antennas can measure the two orthogonal components of the polarised electric field, so the correlations can also be formed for each polarisation and between polarisations, to recover all information contained in the radio waves. Ensembles of visibilities represent the spatial frequency components of the structure of the cosmic radio sources and can be inverted, via Fourier transforms, into images of those radio sources. The physics of this image formation step is identical to how a magnifying lens produces an image. One of the standard references for radio interferomery is listed in [2], to which the reader is refereed for a much deeper description of radio interferometers and image formation using such instruments.

In his presentation at the 2014 IEEE Workshop on Statistical Signal Processing, Louis Scharf gave a couple of illuminating slides that presented exactly this process in relation to a completely different problem, describing the process as multi-channel spectral correlation, where the multi-channel aspect maps to the multiple antenna aspect of a radio interferometer. Thus, the process of correlating the data for a radio interferometer is the first relevant signal processing challenge we encounter that may be of interest to the signal processing community. The inversion of spatial frequency measurements into the image plane is also likely to be of interest to the signal processing community. I discuss these aspects in slightly more depth in later sections of this paper and other participants in the ICASSP special session will deal with these issues in more detail, being the experts in these specialised fields.

Interferometry has been used in radio astronomy for more than forty years. Examples of large interferometric arrays are the Giant Metre-wave Radio Telescope (GMRT) in India [3] (Figure 1) and the Very Large Array (VLA) in the US [4]

(Figure 2). The GMRT and VLA consist of $N = 30$ and $N = 27$ antennas, respectively. These are relatively large-N arrays. The SKA will have $N \sim 900$ for its low frequency component and $N \sim 100$ and $N \sim 250$ for its two higher frequency components. As the amount of data produced by an interferometer scales as $N^2$, it is easy to see that the SKA will produce vast amounts of data. This is one of the biggest challenges for the SKA, producing, processing, transporting, storing, analysing, and interpreting these data. In the next section, I'll briefly quantify this challenge.



**Fig. 1**. Antennas of the Giant Metre-wave Radio Telescope (GMRT: http://gmrt.ncra.tifr.res.in/).



**Fig. 2**. Antennas of the Very Large Array (VLA). Image courtesy of NRAO/AUI.

## 3. THE SKA DATA CHALLENGE

Equation 1, below, approximately illustrates the rate of numbers produced by a radio interferometer, $R$, per second (a complex visibility is made up of two numbers, representing the real and imaginary components of the visibility):

$$R \propto N(N-1)n_s \frac{B}{\Delta B}\frac{1}{\Delta t} \quad \mathrm{s}^{-1},$$

where: $N$ is the number of antennas in the array; $n_s$ is the number of polarisation parameters (products of the polarisation components); $B$ is the frequency width of the band-limited signal (in MHz); $\Delta B$ is the frequency width of the sub-bands generated by the channelisation process (in MHz), the radio frequency sampling period of the visibilities; and $\Delta t$ is the temporal sampling period of the visibilities. The choices of these various parameters depend on the exact science goals and can be varied per observation. Thus, the calculations presented here should only be considered illustrative.

For the VLA, with $N = 27$, $n_s = 4$, $B = 1000$, $\Delta B = 0.1$, and $\Delta t = 10$, $R \sim 3 \times 10^6$ numbers per second. The precision at which these numbers are stored and processed is somewhat dependent of the requirements of the experiment. For example, one can convert $R$ to a data rate by choosing a number of bytes to represent each number. With a 16 bit floating point representation, the data rate becomes approximately 50 Mbps, or approximately 20 GB per hour.

Contrast this with the first operational precursor telescope for the SKA, the Murchison Widefield Array (MWA)[3] [5] (Figure 3). The MWA has $N = 128$, $n_s = 4$, $B = 30$, $\Delta B = 0.01$, and $\Delta t = 2$, giving $R \sim 100 \times 10^6$ numbers per second, 30 times larger than for the VLA.



**Fig. 3**. Antennas of the Murchison Widefield Array (MWA).

Finally, consider one of the four parts of the SKA, the low frequency array (Figure 4), in its Phase 1 form. This segment of the SKA will operate at low frequencies (the MWA is a precursor for this segment of the SKA) and will have $N = 866$, $n_s = 4$, $B = 300$, $\Delta B \sim 0.01$, and $\Delta t \sim 1$, giving $R \sim 100 \times 10^9$ numbers per second, 1000 times larger than for the MWA. At 16 bit precision, SKA-low will produce of order a terabit per second of data, or $\sim$400 TB per hour. That is $\sim$10 PB per day or $\sim$4 exabytes per year!

[3]http://www.mwatelescope.org



**Fig. 4**. Antennas of the future SKA-low. Image credit: ICRAR/Swinburne University of Technology/ASTRON/SKA.

Two other SKA arrays (dish-based arrays operating at higher frequencies) will be operational at the same time as SKA-low and will produce comparable amounts of data. And all this only represents the Phase 1 construction of the SKA, the first 10%. The final Phase 2 construction will produce arrays $\sim$10 times larger.

The SKA will push astronomy into exascale signal processing and computing in the next decade. We will need to be clever with both new signal processing algorithms and practical implementations of those algorithms, in order to deal sensibly with the amount of data processed.

## 4. OPPORTUNITIES FOR SIGNAL PROCESSING

In what areas of signal processing will we have to change our ways in order to meet the challenges posed by the SKA? Where could radio astronomers do with advice and help from the signal processing community?

Above I briefly covered the process of correlation, whereby the digital representation of the electric field variations at each antenna are combined to form visibilities. This is a very highly compute intensive task, usually implemented in very high speed Application Specific Integrated Circuits or Field Programmable Gate Arrays, but increasingly in software using General Purpose Graphical Processing Units or even Central Processing Units. The input data rates to the correlator are massive (60 Gbps for the MWA, for example) and, as has been seen above, the output visibility data rates are very large in the SKA case. While the correlation process is well known, can the signal processing community suggest improvements? Further, can we incorporate additional functionality into the correlator? The correlator may be a natural place in the system to detect and remove human-made radio frequency interference. Separating human-made signals from cosmic signals can result in higher quality data for analysis [6]. For example, there is a chance that the radio signals from hydrogen gas in the early Universe reach us in the FM broadcast band (87.5 - 108 MHz). Further, the signal processing analysis of data at the correlator could be used to detect

sources of transient radio emission, usually explosive events involving compact objects such as neutron stars and black holes [7]. These applications are generally highly sophisticated signal processing tasks in both the time and frequency domains.

I also briefly touched above on the image formation process, whereby the visibility data are transformed into the image plane, to reveal the structure of the sky (and therefore the Universe) at radio wavelengths. This appears to me to be a ripe area for a fresh look at the way radio astronomers work, with many aspects to re-consider. In a conventional sense, the act of forming images from visibilities requires the calibration of those visibilities; the visibilities encode information about the sky, but also encode deformations of the wavefront caused by the Earth's atmosphere and ionosphere, as well as instrumental imperfections. These non-astronomical effects require correction before the visibilities can be used to form an image.

Up to the current point in time radio interferometers have recorded visibility data for post-observation transformation into images. However this will not be possible for the SKA, given the data rates out of the correlator. Thus, visibility calibration and image formation will have to occur in real-time. This is a massive signal processing challenge.

More radically, why do we make images? It is, of course, the image plane that appeals to human intuition. However, the visibility (conjugate) plane contains the same amount of information as the image plane. The visibility plane is the actual measurement domain for interferometers and is where we best understand the noise properties of our data. And when you look closely at how astronomers implement the transform from visibilities to images, you realise that we start with pristine visibility data and take pretty nasty shortcuts that actually degrade the quality of our data. A good example is the quantisation of the visibility plane in order to utilise Fast Fourier Transforms in the image formation step. Perhaps a fresh look at this process may lead to techniques to recover the information about the sky without leaving the visibility plane, utilising the data in their purest form and maybe saving money (or maybe costing more money!). Signal processing techniques in radio astronomy are now being looked at from the Information Theory point of view [8], but there is a long way to go.

Finally, once the astrophysical information regarding the sky has been recovered from the measurements, this information will be held in vast catalogs containing probably many hundreds of millions of events, as functions of time, frequency, and morphology. Sifting through these data, matching the radio objects and events with catalogs from other telescopes at other wavelengths (i.e. from optical telescopes, X-ray telescopes etc etc), recognising trends in populations, recognising patterns in the data, and performing exception detection to discover new phenomena, will present massive challenges in the domain of Big Data analytics. Big Data an-

alytics is now big business in the retail, finance, economics, scientific, and government sectors, with signal processing, statistics, and Information Theory at the heart of the matter. Radio astronomy must be able to learn from these other sectors and, as radio astronomy pushes the envelope with Big Data, perhaps these sectors can learn something from radio astronomy?

## 5. CONCLUDING REMARKS

In this brief contribution I've tried to outline the basic nature of the SKA and the challenges it faces in terms of generating and analysing the data it will produce. The amount of data involved is directly driven by the astrophysical and cosmological questions being asked, pushing to the limits of physics in the same way the the Large Hadron Collider does. These days, the limit of physics is Big Science, requiring Big Data.

The SKA will be a data intensive facility, with very large data rates and very large computational requirements. Most of the data processing for the SKA will be familiar to the signal processing community, including time series analysis, multi-channel spectral analysis, image formation and analysis, separation of signals in time/frequency/Fourier spaces, statistical signal processing and Information Theory, and Big Data analytics. Thus, it appears to me that radio astronomy and the SKA are areas in which the signal processing community may find some interesting problems to work on, with mutual benefit.

It is always interesting to reach out to another field and discover that the main barrier to engagement are the different nomenclatures that evolve in different disciplines to describe exactly the same underlying physics and mathematics. I have found that this is the case between radio astronomy and signal processing engineers. That this difference exists is interesting from a sociological point of view and a little depressing as it again illustrates the silo nature of science. However, more interesting is finding ways to address the differences and cross the boundaries between these disciplines. Of course, step one is to start talking and keep talking. Thus, I greatly appreciate the chance to run the special session at ICASSP 2015 and interact with the signal processing community. I hope that a few from this community can come along to some radio astronomy workshops and conferences in the future.

Step two, I think, is actually starting to work together on some problems of common interest. I've outlined a few in this paper and we have the means to start making this happen. For example, the Murchison Widefield Array (MWA) is the first operational precursor for the SKA, already starting to seriously move into the Big Data domain, where we have to begin to face the challenges I've outlined. The MWA has been operational since mid-2013 and in its first two years of observations will have collected 3 PB of data, quite enough to be going on with to explore some of the signal processing challenges discussed here.

## 6. REFERENCES

[1] R. T. Schilizzi, P. E. F. Dewdney, and T. J. W. Lazio, "The square kilometre array," in *Society of Photo-Optical Instrumentation Engineers (SPIE) Conference Series*, July 2010, vol. 7733 of *Society of Photo-Optical Instrumentation Engineers (SPIE) Conference Series*.

[2] A Richard Thompson, James M Moran, and George W Swenson, *Interferometry and Synthesis in Radio Astronomy; 2nd ed.*, Wiley-VCH, Weinheim, 2001.

[3] R. Nityananda, "The GMRT: Background, Status, and Upgrades," in *The Low-Frequency Radio Universe*, D. J. Saikia, D. A. Green, Y. Gupta, and T. Venturi, Eds., Sept. 2009, vol. 407 of *Astronomical Society of the Pacific Conference Series*, p. 389.

[4] R. Perley, P. Napier, J. Jackson, B. Butler, B. Carlson, D. Fort, P. Dewdney, B. Clark, R. Hayward, S. Durand, M. Revnell, and M. McKinnon, "The Expanded Very Large Array," *IEEE Proceedings*, vol. 97, pp. 1448–1462, Aug. 2009.

[5] S. J. Tingay, R. Goeke, J. D. Bowman, D. Emrich, S. M. Ord, D. A. Mitchell, M. F. Morales, T. Booler, B. Crosse, R. B. Wayth, C. J. Lonsdale, S. Tremblay, D. Pallot, T. Colegate, A. Wicenec, N. Kudryavtseva, W. Arcus, D. Barnes, G. Bernardi, F. Briggs, S. Burns, J. D. Bunton, R. J. Cappallo, B. E. Corey, A. Deshpande, L. Desouza, B. M. Gaensler, L. J. Greenhill, P. J. Hall, B. J. Hazelton, D. Herne, J. N. Hewitt, M. Johnston-Hollitt, D. L. Kaplan, J. C. Kasper, B. B. Kincaid, R. Koenig, E. Kratzenberg, M. J. Lynch, B. Mckinley, S. R. Mcwhirter, E. Morgan, D. Oberoi, J. Pathikulangara, T. Prabu, R. A. Remillard, A. E. E. Rogers, A. Roshi, J. E. Salah, R. J. Sault, N. Udaya-Shankar, F. Schlagenhaufer, K. S. Srivani, J. Stevens, R. Subrahmanyan, M. Waterson, R. L. Webster, A. R. Whitney, A. Williams, C. L. Williams, and J. S. B. Wyithe, "The Murchison Widefield Array: The Square Kilometre Array Precursor at Low Radio Frequencies," *PASA*, vol. 30, pp. 7, Jan. 2013.

[6] J. F. Bell, P. J. Hall, W. E. Wilson, R. J. Sault, R. J. Smegal, M. R. Smith, W. van Straten, M. J. Kesteven, R. H. Ferris, F. H. Briggs, G. J. Carrad, M. W. Sinclair, R. G. Gough, J. M. Sarkissian, J. D. Bunton, and M. Bailes, "Base Band Data for Testing Interference Mitigation Algorithms," *PASA*, vol. 18, pp. 105–113, 2001.

[7] C. J. Law and G. C. Bower, "All Transients, All the Time: Real-time Radio Transient Detection with Interferometric Closure Quantities," *ApJ*, vol. 749, pp. 143, Apr. 2012.

[8] C. M. Trott, R. B. Wayth, J.-P. R. Macquart, and S. J. Tingay, "Source Detection in Interferometric Visibility Data. I. Fundamental Estimation Limits," *ApJ*, vol. 731, pp. 81, Apr. 2011.