




AlgoSCR: an algorithm for solar contamination removal from radio interferometric data

Anh Phan ¹★, Santanu Das,^{1,2}★ Albert Stebbins,² Peter Timbie,¹★ Reza Ansari,³ Shifan Zuo,⁴ Jixia Li,⁵ Trevor Oxholm,¹ Fengquan Wu,⁵ Xuelei Chen ^{5,6,7} Shijie Sun,^{5,6} Yougang Wang ⁵ and Jiao Zhang⁸

¹Department of Physics, University of Wisconsin Madison, 1150 University Ave, Madison WI 53703, USA

²Fermi National Accelerator Laboratory, P.O. Box 500, Batavia IL 60510, USA

³Université Paris-Saclay, CNRS/IN2P3, IJCLab, F-91405 Orsay, France

⁴Department of Astronomy and Tsinghua Center for Astrophysics, Tsinghua University, Beijing 100084, P. R. China

⁵National Astronomical Observatory, Chinese Academy of Science, 20A Datun Road, Beijing 100101, P. R. China

⁶University of Chinese Academy of Sciences, Beijing 100049, P. R. China

⁷Center of High Energy Physics, Peking University, Beijing 100871, P. R. China

⁸School of Physics and Electronics Engineering, Shanxi University, Taiyuan 030006, P. R. China

Accepted 2022 February 25. Received 2022 February 22; in original form 2021 June 24

ABSTRACT

Hydrogen intensity mapping is a new field in astronomy that promises to make three-dimensional maps of the matter distribution of the Universe using the redshifted 21 cm line of neutral hydrogen gas (HI). Several ongoing and upcoming radio interferometers, such as Tianlai, CHIME, HERA, HIRAX, etc., are using this technique. These instruments are designed to map large swaths of the sky by drift scanning over periods of many months. One of the challenges of the observations is that the daytime data are contaminated by strong radio signals from the Sun. In the case of Tianlai, this results in almost half of the measured data being unusable. We try to address this issue by developing an algorithm for solar contamination removal (AlgoSCR) from the radio data. The algorithm is based on an eigenvalue analysis of the visibility matrix and hence is applicable only to interferometers. We apply AlgoSCR to simulated visibilities, as well as real daytime data from the Tianlai dish array. The algorithm can reduce strong solar contamination by about 95 per cent without seriously affecting other weaker sky signals and thus makes the data usable for certain applications.

Key words: instrumentation: interferometers – methods: analytical – methods: data analysis – cosmology: observations – radio continuum: general.

1 INTRODUCTION

Cosmologists study the Universe on the largest observable distance scales in order to understand its origin and evolution. In the past few decades, cosmic microwave background instruments have mapped almost the entire sky with high sensitivity and fine angular resolution. These maps measure the intensity and polarization fluctuations at the last scattering surface and remain a primary tool for studying the Universe. However, for understanding the nature of dark matter and dark energy, it is essential to study the evolution of structure as a function of time. Galaxy redshift surveys have been extremely successful in mapping the large-scale structure of the Universe by cataloging the distribution of luminous galaxies in redshift space. These maps can be used, for example, to observe the characteristic baryon-acoustic oscillation signal, which can be used as a standard ruler to extract cosmological parameters. However, as we map larger and more distant volumes of the Universe, the method faces multiple challenges. For example, the galaxies become fainter and spectral lines are redshifted to wavelengths that are difficult to detect from the ground.

Hydrogen intensity mapping, a radically different technique, creates 3D maps using the 21 cm emission of neutral hydrogen (HI) without resolving individual galaxies. This line is unique in cosmology as, for $\lambda > 21$ cm, it is the dominant astronomical line emission for all redshifts. Hence, to a good approximation the wavelength of a spectral feature can be converted to a redshift without having to first identify the atomic transition. In principle, HI intensity mapping could be used to make 3D maps of matter at all redshifts up into the ‘dark ages’ ($z \approx 100$), even before galaxies have formed.

The first HI intensity mapping observations began over a decade ago (Abdalla & Rawlings 2005; Peterson, Bandura & Pen 2006; Chang et al. 2008; Mao et al. 2008; Morales 2008) and interest has continued to grow (Ansari et al. 2018; Slosar et al. 2019; Liu & Shaw 2020). A number of dedicated projects have been launched to detect the signal and turn the technique into a useful cosmological tool. These are mainly interferometers, such as CHIME (Bandura et al. 2014; Newburgh et al. 2014), Tianlai (Chen 2011; Xu, Wang & Chen 2014; Das et al. 2018; Li et al. 2020; Wu et al. 2021), MWA (Tingay et al. 2013), LWA Eastwood et al. (2018), HERA (DeBoer et al. 2017), HIRAX (Newburgh et al. 2016), and PUMA (Slosar et al. 2019), but they also include single dishes with multiple feed antennas, such as BINGO (Battye et al. 2012; Dickinson 2014; Wuensche et al. 2019) and FAST (Hu et al. 2020). Intensity mapping

* E-mail: anh@wisc.edu (AP); sanjone@gmail.com (SD); pttimbie@wisc.edu (PT)



Figure 1. Left: A top view photograph of the Tianlai arrays, which consist of the dish array and the cylinder array. The photo was taken with a drone at a height of 280 m above the ground. The arrays saw first light in 2016. The position of the calibration noise source is indicated by the white arrows on the left. Right: A schematic diagram of the Tianlai dish array. The dishes are arranged in two concentric circles of radius 8.8 m and 17.6 m around a central dish. The dishes have dual-linear polarization feed antennas with one axis oriented parallel to the altitude axis (horizontal, H, parallel to the ground, or E–W in the figure) and the other orthogonal to that axis (vertical, V, N–S in the figure). For example, the red line shows one of the baselines that is the H polarization of dish 4 correlated with the H polarization of dish 9: [4H 9H]. The above image is reproduced from Wu et al. (2021).

instruments can address questions at a variety of redshift ranges. At $z \sim 10$, they probe the Epoch of Reionization, star formation, and galaxy assembly, while at lower redshifts they trace large-scale structure for studies of dark energy, etc. (Battye, Davies & Weller 2004; Abdalla & Rawlings 2005; Peterson et al. 2006; Chang et al. 2008; Mao et al. 2008; Morales 2008; Bull et al. 2015).

So far, the HI signal has not been detected using intensity mapping by itself. Intensity mapping observations, in the post-recombination epoch, have detected HI when cross-correlated with galaxy redshift surveys (Masui, McDonald & Pen 2010; Masui et al. 2013; Anderson et al. 2018). A number of challenging systematic effects must be overcome to allow autocorrelation detections. The foremost of these is separating the HI signal from Galactic and extra-Galactic astronomical foregrounds, which are ~ 4 – 5 orders of magnitude brighter (Liu & Shaw 2020). The Sun represents an astronomical foreground that is even brighter and of a different character.

The daytime data from radio interferometer arrays in general, and the Tianlai dish array in particular, are contaminated by the solar signal, making the data unusable for most astronomical analyses. The lost data have a significant impact on observing efficiency; reaching the required survey sensitivity means observing the sky for almost twice the number of days. This penalty is particularly problematic for HI intensity mapping, where long integration times (months or years) are necessary to detect the HI signal. Furthermore, this data loss prevents obtaining continuous, 24-h data sets, which allow dense coverage of the $u-v$ plane and facilitate detection of periodic signals. Furthermore, not having 24 h of continuous usable observations prevents the application of m-mode map-making techniques (Shaw et al. 2014). The (u,v) coverage can still be quite good with nighttime data, and one can recover full 24-h RA coverage by combining nighttime data from observations about 6 months apart. In this paper, we try to remove the solar contamination from the daytime data from a radio interferometer. We have used the data from the Tianlai dish array as our test sample. However, the problem is not unique to Tianlai; the same algorithm may be used for other radio interferometric observations. While daytime observations with single dish radio telescopes are also plagued by the Sun’s signal, this algorithm is applicable only to interferometer arrays.

The Tianlai Project is led by the National Astronomical Observatory of China (NAOC). It consists of two pathfinder radio interferometers: an array of cylinder antennas and an array of dishes, at a radio-quiet site in Xinjiang, China (Chen 2012; Li et al. 2020; Wu et al. 2021). The objective is to obtain high-fidelity 3D images of the northern sky using HI intensity mapping. The analysis described in this paper concentrates on the dish array data. The Tianlai dish array consists of 16 steerable, 6-m diameter dishes; a schematic is shown in Fig. 1, which also shows the dish numbering scheme. We use these dish numbers for referring to different baselines in the paper. The dish array currently operates between 685 MHz and 810 MHz, corresponding to redshift $0.75 < z < 1.07$, divided into 512 equally spaced frequency bins of width 244 kHz ($\delta z = 0.0002$). The 16 dual-polarization feeds yield 32 autocorrelation visibilities and $32 \times (32 - 1)/2 = 496$ cross-correlation visibilities, which are currently sampled every second. The system noise temperatures for the dish antennas are 80–85 K (Zhang et al. 2016; Das et al. 2018; Li et al. 2020; Wu et al. 2021).

Different methods have been proposed to remove broad-band radio frequency interference (RFI) and the solar contamination (Briggs, Bell & Kesteven 2000; Paciga et al. 2011) based on singular value decomposition (SVD) and other methods in the time and frequency domain. The solar contamination is extremely strong in the Tianlai data and methods using SVD decomposition in the time and frequency domain remove the background signal along with the solar contamination. The objective of this paper is to describe an eigenvalue-based approach that operates in the baseline space for removing solar contamination from radio interferometric data without affecting the background signal. We propose an algorithm, AlgoSCR, that can remove most of the solar contamination, provided the Sun is the strongest source in the sky along with other weaker sources. The paper is organized as follows. In the second section, we discuss the solar contamination problem in the Tianlai dish array in detail. The third and the fourth sections give the detailed algorithm for removing the solar contamination. We also show the results of our analysis on the real Tianlai data. To test what fraction of the Sun signal can be removed by our algorithm, and how much signal from other cosmic sources is removed by it, in Section 5 we perform two tests. First, we apply it to Sun-contaminated Tianlai data and compare

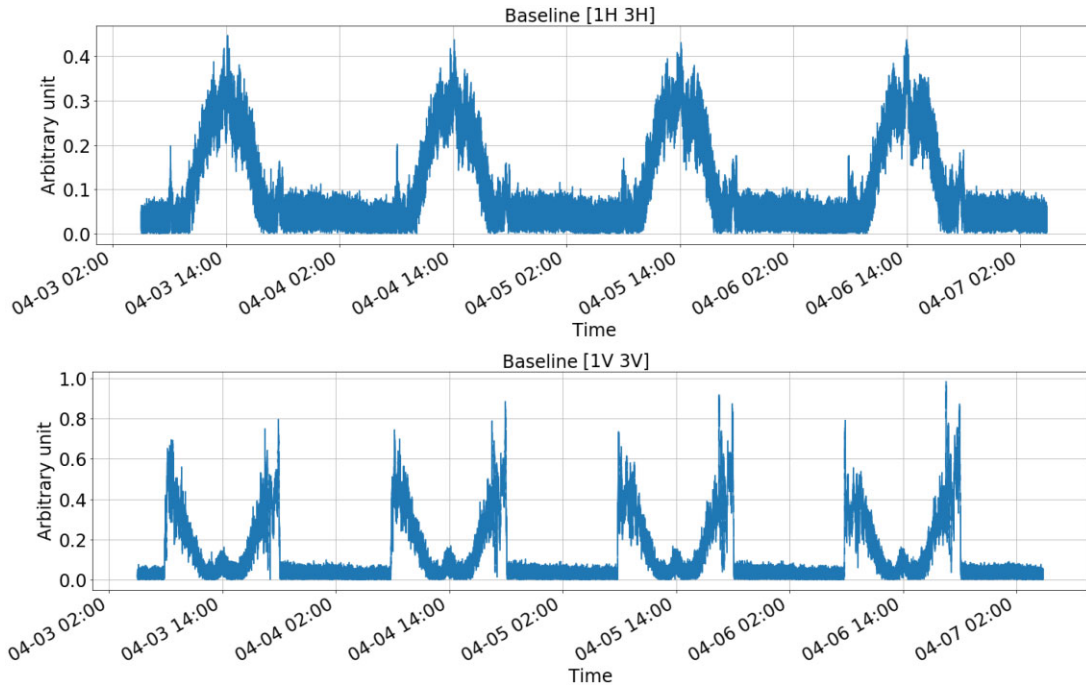


Figure 2. Value of the (uncalibrated) cross-correlation visibility amplitudes averaged over the 10 central frequency bands during 4 d of observations of the NCP in 2019 April. Integration time is 1 s. Each plot corresponds to a different baseline, as indicated. The baseline numbering scheme appears in Fig. 1. Time is given in local time.

the cleaned visibilities to the same sidereal times observed during nighttime. Secondly, we apply it to simulated data where the amount of solar contamination and foreground point sources are known. In the Discussion section, we assess the efficacy of the method and the issues that we face when applying it. We also describe some future directions to pursue with this approach.

2 THE SOLAR CONTAMINATION PROBLEM

The Tianlai data show strong contamination from the solar signal during the daytime. In Fig. 2, we plot the sum of the absolute visibility from 10 frequency channels (out of 512 channels) at the centre of the band for 4 consecutive days (total 96 h) for two different baselines. The horizontal axis shows the time in hours, starting at the beginning of the observations. The two different plots are for two representative visibilities. We can see a roughly smooth visibility amplitude for about 10 h every day and then a sudden increase in the absolute visibility and a noisy pattern for about the next 14 h.

The plots clearly show that the daytime signal is several times stronger than the night. The shape of the contamination pattern also varies with baseline. Some of the baselines show a bumpy feature with the strongest visibility occurring near noon, whereas for other baselines the signal is strongest during Sunrise and Sunset and shows a ‘dip’ feature during the daytime. The top plot is the autopolarization visibility corresponding to two horizontal feeds, whereas the bottom plot shows an autopolarization visibility from two vertical feeds. The autopolarization signals from similar feeds on other baselines show roughly similar types of patterns, except for a couple of baselines. The data are taken during observations of the North Celestial Pole (NCP) in 2019 April. During this observing period, the path of the Sun is located at an angle of approximately 85° from the direction of the main beam. The plot gives an overview of the magnitude of the solar contamination problem in the Tianlai dish array.

An obvious conclusion of this strong daytime visibility is that the telescopes are responding to the Sun’s illumination of their far sidelobes. For the baselines measuring correlations of the H polarization, the antenna sidelobes are aligned with the direction of Sun near noon, providing a strong visibility at midday, while for the V polarization the Sun falls between two side lobes at noon, producing stronger signal during Sunrise and the Sunset. These effects are consistent with the expected responses of the feed antennas, which are essentially orthogonally oriented crossed dipoles.

In Fig. 3, we show a cut through the simulated beam pattern for a single dish measured using an electromagnetic (EM) simulation package (CST¹), corresponding to the Sun’s track during the daytime. We can see that for one of the polarizations, we are getting a low amplitude during the midday whereas for the other polarization (bottom plot), the amplitude is comparatively high at noon. The simulated patterns shown in Fig. 3 do not exactly replicate the observed pattern of Fig. 4, because the sidelobes from these EM simulations do not exactly match those of the real beam. The sidelobes at this particular angle are also highly cluttered. A couple of degrees change in the path gives rise to a very different shape in the sidelobes, making it difficult to reconstruct the exact pattern through such an EM simulation.

This daily response to the Sun signal is relatively constant over a period of a year. Fig. 4 shows the directive gain of the dish antennas as computed by an EM simulation. The Sun enters the sidelobes of the antennas over a range of polar angles for which the beam patterns are relatively flat. The simulation is consistent with measurements of the daytime visibilities at different times of the year. Using the eigenvalue analysis described below, Fig. 5 shows the contribution by the Sun to the visibility for a typical baseline during 2018 January and then again in 2019 April. As the paths of the Sun through the

¹<https://www.3ds.com/products-services/simulia/products/cst-studio-suite/>

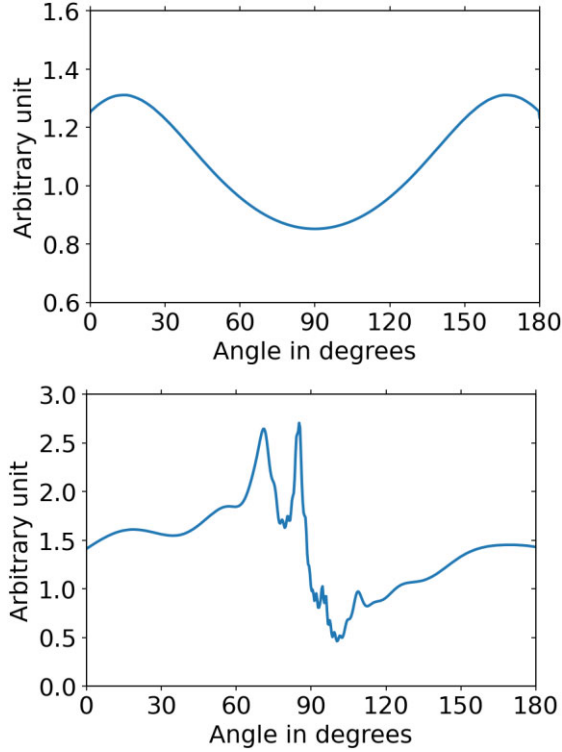


Figure 3. Simulated antenna directivity along the Sun’s track across the beam during daytime (180° corresponding to 12 h) when the dish array is observing the NCP in April. The scale is in arbitrary linear units (not in dB). The plot shows that for one of the polarizations, the signal from the Sun has a ‘dip’ in power during the daytime, whereas for the other polarization, there is high power during certain parts of the daytime. Even though the plots do not show the exact features that we observe during daytime (Fig. 2), we must remember that the sidelobes from the electromagnetic simulations are not exactly the same as those of the real antenna.

sidelobes of the antennas are different at different times of the year, the visibilities are also slightly different, but the overall patterns of the signals are similar. We can see that the amplitudes are within ~ 30 per cent of each other. Part of this amplitude variation was induced by the variation of system gain, which is caused by the different air temperature in January and April. The fast oscillating fringes from the Sun are present in both plots.

The complex visibility for a typical baseline is shown as a ‘waterfall plot’ in Fig. 6 for a 24-h period. We can see that the daytime data are dominated by bright fringes caused by the Sun. On the other hand, the pattern in the nighttime data comes from the much dimmer radio sky and has a very different character. The dominant fringes in the nighttime data come from a combination of weak sources near the NCP and bright sources far from the NCP, particularly Cassiopeia A (Cas A) and Cygnus A (Cyg A).

3 REMOVING SUN CONTAMINATION USING EIGENVALUE ANALYSIS

We start by defining the notation used in this paper. The visibility matrix is given by

$$\mathbf{V} = [\mathbf{D}^s \mathbf{G}]^\dagger [\mathbf{D}^s \mathbf{G}] + \langle [\mathbf{N}]^\dagger [\mathbf{N}] \rangle, \quad (1)$$

where \mathbf{D}^s is the voltage signal from the antenna in matrix form, \mathbf{G} is a direction-independent complex gain matrix, \mathbf{N} describes the noise from the receivers, and † represents the conjugate transpose. The

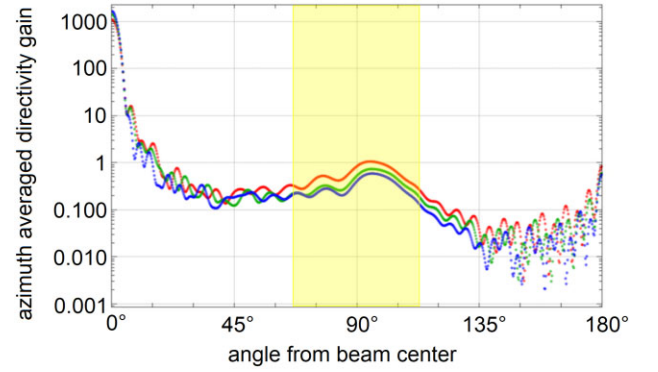


Figure 4. Simulated beam directivity as a function of beam angle θ from the beam centre of the antennas for three different frequencies, 700 (red), 750 (green), and 800 MHz (blue). Each plot shows the absolute co-polar directive gain averaged over the azimuthal angle. The angle is the polar angle calculated from the centre of the beam. The yellow shaded region shows the range of polar angles for which the Sun appears in the sidelobes of the beam, ranging from 66.55° at the Summer Solstice to 113.45° at the Winter Solstice, when the dish array observes the NCP. The gain is relatively flat over this range of angles and causes the Sun signal to vary by only a factor of about 6 over the year.

individual components are given by

$$\mathbf{V}_{(i,j)} = \langle E_i^* E_j \rangle. \quad (2)$$

E_i represents the complex voltage from receiver i , with E_i^* being its complex conjugate. E_i is given by

$$E_i = \left(\sum_s D_i(\vec{\omega}_s) e^{i\mathbf{k}\cdot\mathbf{r}_i} F_s \right) G_i + N_i, \quad (3)$$

where F_s is the electric field of the radio wave coming from a source on the celestial sphere, $D_i(\vec{\omega}_s)$ is the primary beam of antenna i , and this is a function of the direction vector $\vec{\omega}_s$, \mathbf{k} is the three-dimensional wavenumber, the Fourier dual to the position vector \mathbf{r}_i of feed i .

The intensity of the source at any frequency ν , is given by

$$I_s(\nu) = |F_s(\nu)|^2 = F_s^*(\nu) F_s(\nu). \quad (4)$$

For extended sources, we need to integrate over different directions for calculating E_i :

$$E_i = \left(\int D_i(\vec{\omega}_s) e^{i\mathbf{k}\cdot\mathbf{r}_i} F_s d\omega_s \right) G_i + N_i. \quad (5)$$

The visibility is an ensemble average of the $E_i^* E_j$, i.e.

$$\begin{aligned} \mathbf{V}_{(i,j)} &= \langle E_i^* E_j \rangle_{\tau_{\text{int}}} \\ &= \left[\frac{1}{\tau_{\text{int}}} \int_0^{\tau_{\text{int}}} E_i^* E_j dt \right], \end{aligned} \quad (6)$$

where τ_{int} is the integration time, which is constant for any time and frequency bin (t, ν) . For the current Tianlai setup, the integration time is 1 s. The asterisk (*) represents the complex conjugate and the bracket $\langle \rangle$ represents the ensemble average.

Here, we should note that the visibilities from different astrophysical sources are additive. Provided there is only one point source on the sky, the visibility matrix, i.e. $\mathbf{V}_{(i,j)}$, at any time can be written as an outer product of the electric field from the source measured at different feed antennas. Therefore, if the visibility matrix is decomposed into its corresponding eigenvalues and eigenvectors, there should be only one non-zero eigenvalue. In the presence of other weaker sources, the largest eigenvalue should correspond to

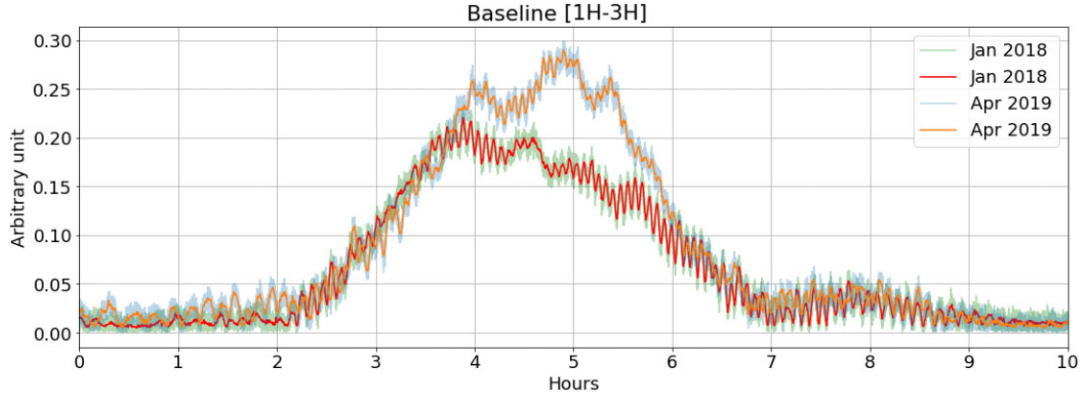


Figure 5. The amplitude of the daytime visibilities in 2018 January and 2019 April. Due to the difference in the time of Sunrise in January and April, the 0 h of each curve is adjusted so that the Sun signals from both data sets peak at about the same time. The green and blue curves are the amplitude of the visibility obtained from the telescope for 2018 January and 2019 April, respectively. The ‘fast oscillation’ fringes are seen in both observations. The red and orange curves are corresponding spline fits to better highlight the fast oscillation fringes.

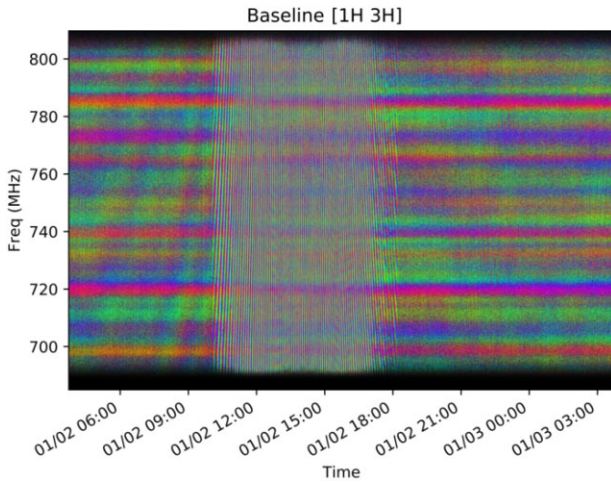


Figure 6. The complex visibility for a typical baseline plotted over a 24-h period in 2019 January. We represent the phase of the complex visibility by hue (colour) and the amplitude by value (brightness) in an HSV (hue, saturation, value) display of the colour model (see Fig. B2 for details). The local time proceeds linearly from left to right, with a sampling interval of 1 s. The frequency increases linearly from bottom (685 MHz) to top (810 MHz) in 512 equally spaced frequency bins. The time interval from about 10:00–18:00 is dominated by the Sun.

the Sun signal and the eigenvector corresponding to the largest eigenvalue will roughly point towards the direction of that source in the eigenspace. The contributions from additional, weaker sources and noise may alter the direction slightly.

Comparing the visibility amplitudes between the daytime and the nighttime data in Fig. 2, we can infer that the largest contribution to the daytime signal is from the Sun, entering through the antenna side-lobes. Therefore, in the eigendecomposition of the visibility matrix, the largest eigenvalue should represent the solar contamination.

3.1 Issues with the autocorrelation signal

The voltage from the feeds contains a contribution from the receiver noise. Therefore, the measured signal or voltage E_i for a given feed i is the sum of the sky signal, $E_{\text{Sky } i}$ and the instrument noise, N_i , i.e. $E_i = E_{\text{Sky } i} + N_i$.

Under the assumption that the noise terms from separate feeds are uncorrelated, we can say that the ensemble average of the noise from feed i and feed j is zero, i.e. $\langle N_i^* N_j \rangle \approx 0$. Therefore, the visibility for cross-correlated feed i and j , where $i \neq j$, is $\mathbf{V}_{(i,j)} \approx \langle E_{\text{Sky } i}^* E_{\text{Sky } j} \rangle$.

However, for the autocorrelations, the visibilities, $\mathbf{V}_{(i,i)}$ are dominated by the positive noise term $\langle N_i^* N_i \rangle$. The amplitudes of the autocorrelation signals are much higher than those of the cross-correlation signals. Therefore, in an eigendecomposition of the visibility matrix, the eigenvectors are dominated by the noise signals from the autocorrelation, as the sky signals are typically much smaller than the noise.

It is not possible to ignore these autocorrelation signals or simply set them to 0 during the eigenvalue decomposition. To overcome this difficulty, we replace the corresponding terms in the visibility matrix by the following quantity as a proxy for the autocorrelation visibilities:

$$\mathbf{V}_{(i,i)} = \frac{1}{n} \sum_{k,j} \text{abs} \left[\frac{\mathbf{V}_{(i,k)} \mathbf{V}_{(j,i)}}{\mathbf{V}_{(j,k)}} \right], \quad \forall i \neq j \neq k. \quad (7)$$

The receiver noise component in the correlation matrix is bypassed by using equation (7) to replace the autocorrelations. But its long-term effects remain uninvestigated. Here, n is the number of values over which we are doing the sum, i.e. the number of (j, k) pairs. This brings the level of the amplitude of the autocorrelation to the order of the cross-correlation amplitude and we can do a meaningful eigenvalue decomposition.

3.2 DC offset in the visibility

If there is no strong source in the sky, then the real and the imaginary parts of the visibility are expected to randomly fluctuate around 0. However, often in radio interferometers, there are some DC offsets in the real and imaginary components of the visibility. The offsets may originate from a variety of systematic effects, and cross-coupling of signals between the antennas is one of them. In the Tianlai data, we see it in multiple baselines as coloured horizontal stripes in the waterfall plots of the complex visibility (see Fig. 6).

In the top plot of Fig. 7, we show the real and the imaginary parts of the visibility from a transit of Cas A observed by baseline [5H 7V]. The amplitude of the visibility during the transit is expected to form a Gaussian profile. However, as there is some DC offset, we can expect the plot to show some wavy feature modulating the Gaussian.

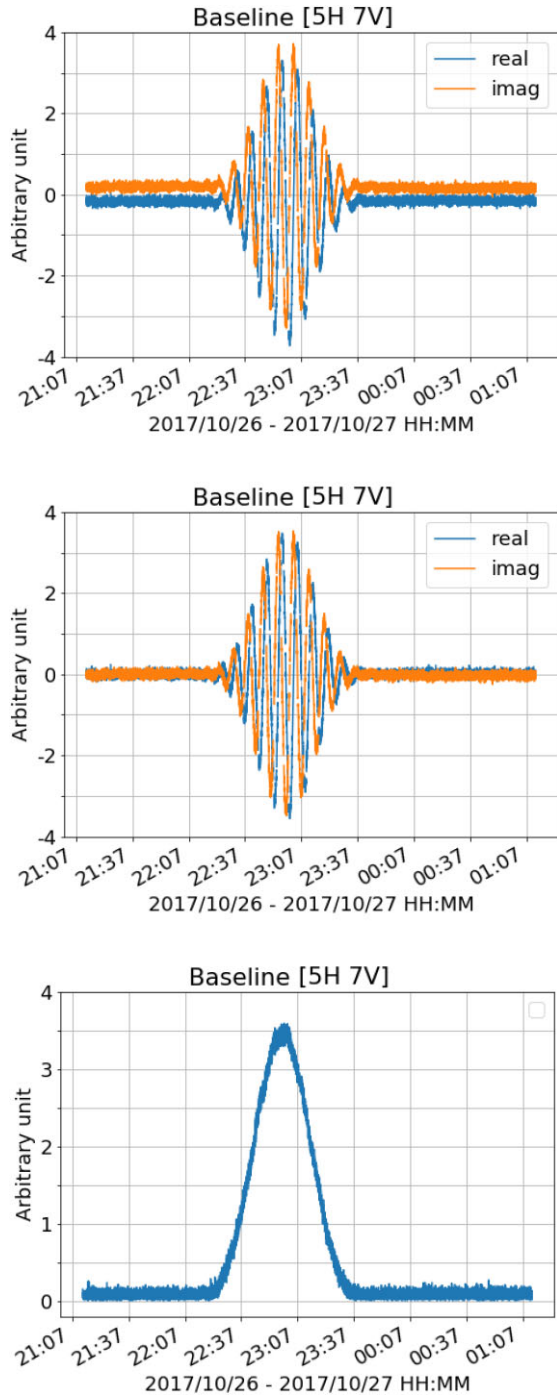


Figure 7. Top: The real and imaginary components of the raw visibility of baseline [5H 7V] during transit of Cas A in 2017 October. We can see that there is a small DC offset in both the real and imaginary components. Middle: The real and imaginary components of the raw visibility after removing the offset from each of the components. Bottom: Amplitude of the visibility of baseline [5H 7V] after removing the mean. We can see a perfectly Gaussian transit peak.

To prevent this, we need to remove the DC offset. In the middle panel of Fig. 7, we show the real and the imaginary components of the visibility, after subtracting the mean of the nighttime data from both the real and imaginary components of the visibility. We remove the nighttime mean from each frequency channel and each baseline.

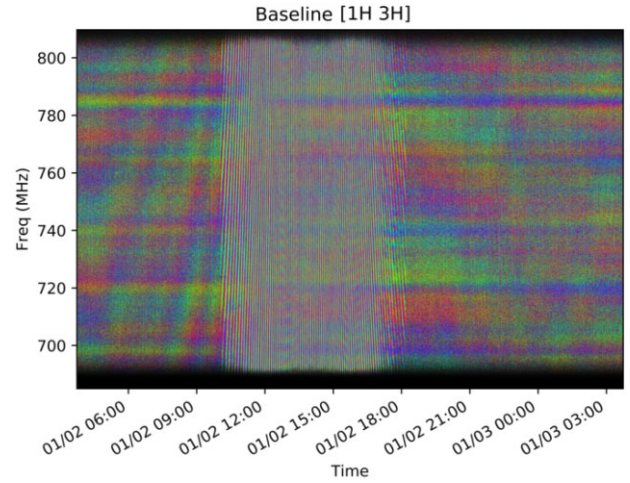


Figure 8. The waterfall plot of the complex visibility (same as Fig. 6) after the nighttime mean subtraction. We can see that most of the horizontal stripes, which probably are caused by crosstalk, are now gone from the waterfall plot. The structures from the sky are more prominent.

The amplitude of the visibility, after DC offset removal, shows a Gaussian peak during the transit of Cyg A, as expected and is shown in the bottom panel of the same figure.

We are investigating the source of the DC offset. However, the nighttime mean subtraction substantially reduces the night-to-night variation in absolute terms and as a fraction of the remaining signal, as discussed in Wu et al. (2021). This nightly mean subtraction removes much of the correlated noise as well as a significant fraction of the signal (gain times sky). Because the sky signal should be the same at the same local sidereal time, it does not contribute to the nightly variation that can be caused by variations in gain or correlated noise. If the variations were due only to gain fluctuations, we would not see a decrease in fractional variation. Thus, much of the subtracted signal is correlated noise.

The presence of this DC offset may also introduce an error in the eigendecomposition and it must be removed before running the Sun removal algorithm described below. We subtract the mean value of the real and imaginary parts of the visibility for each night of data. We do not include the daytime data when computing the mean, because it is contaminated by the Sun. However, the DC offset is very stable over each night and from night to night. So, we remove the nightly mean from the entire 24 h of data, including the daytime data.

In Fig 8, we show a waterfall plot of the complex visibility from one baseline after the nighttime mean removal. We can see the nighttime structures more prominently after the mean subtraction.

Note that, for simplicity, we have considered only the autopolarization signals. If we use both the autopolarization and cross-polarization signals, we expect to get two large eigenvalues, each corresponding to one of the polarizations. However, at present, we are in the process of understanding different systematic effects involved in measuring the cross-polarization signals in the Tianlai data. Different systematic effects, e.g. mutual coupling between two feeds in the same dish, which are in close proximity to each other, are more complicated for cross-polarization data than the same polarization and require detail investigation both in data analysis and the instrumentation level. These are beyond the scope of this paper. Therefore, in this paper, we set the cross-polarization signal

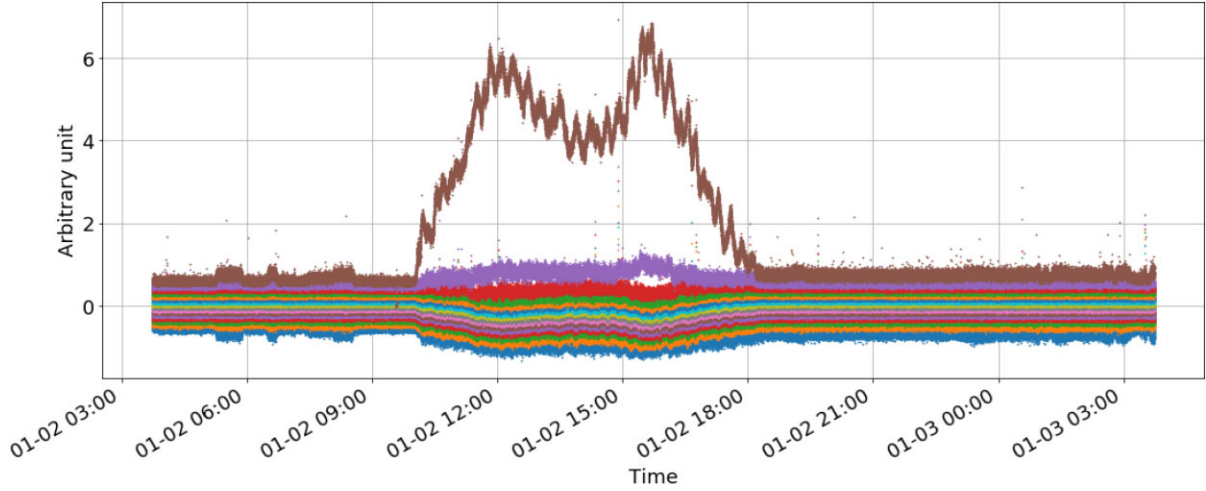


Figure 9. Plot of all the 16 eigenvalues in the eigendecomposition of the horizontal polarization visibility $\mathbf{V}^{(H)}$. We can see that one of the eigenvalues is much larger than the other eigenvalues during daytime. This particular eigenvalue is coming from the solar contamination of the daytime data. We can see that the other eigenvalues are also affected during daytime. This happens due to the change in the eigenvectors, one of which (the eigenvector corresponding to the largest eigenvalue) is oriented towards the Sun during the daytime. The plot here is shown at the central frequency (747.5 MHz) of the observed Tianlai Dish Array band. All other one-dimensional plots also use this frequency.

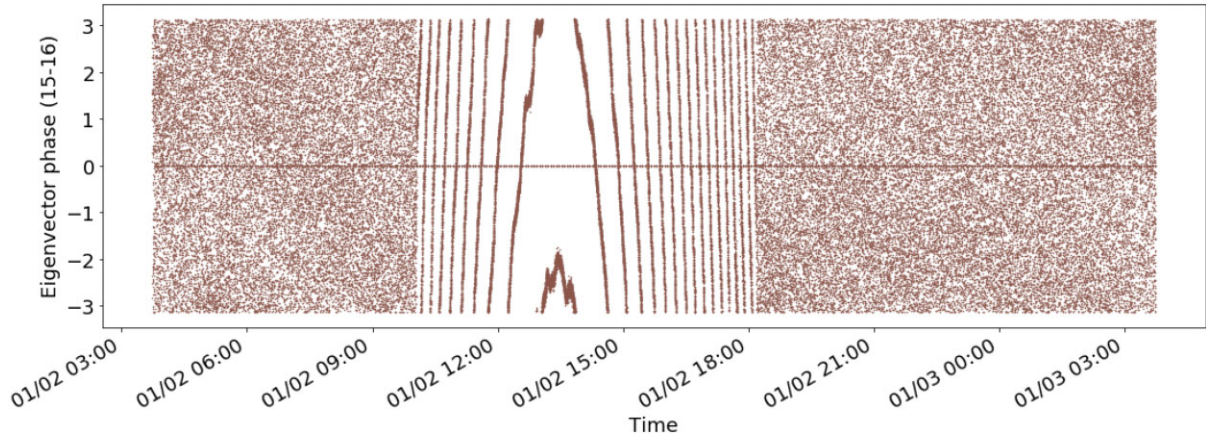


Figure 10. Phase plot of the 15th component of the eigenvector corresponding to the largest eigenvalue of the horizontal polarization visibility, $\mathbf{V}^{(H)}$. We can see the strong fringes during daytime, which confirms that the eigenvalue is coming from a single strong source, the Sun. During night, as there is no single strong source, the phase is varying randomly. The horizontal line in the centre is caused by the calibration noise source, which is turned on and off periodically.

to 0, making the visibility matrix look like

$$\mathbf{V} = \begin{bmatrix} \mathbf{V}^{(H)} & \mathbf{0} \\ \mathbf{0} & \mathbf{V}^{(V)} \end{bmatrix}, \quad (8)$$

where the superscripts H and V refer to the horizontal or vertical polarization, respectively. \mathbf{V} is the visibility matrix at each time and frequency (t, ν).

3.3 Understanding the eigenvalue decomposition

For removing the solar contamination from the real data, we first remove the DC offset from all the cross-correlation channels. We write the visibility matrix for each frequency and at every time bin in the form shown in equation (8). We replace the autocorrelation signals using the formula given in equation (7) and decompose the matrix into eigenvalues and eigenvectors as $\mathbf{V} = \mathcal{E}\Lambda\mathcal{E}^{-1}$. At this point, it should be noted that the eigenvalue decomposition is invariant under a $U(1)$ transformation, i.e. if we multiply the full eigenvector matrix, \mathcal{E} , by a factor of $e^{i\psi}$ for any real ψ , then the corresponding eigenvalue

matrix Λ will remain invariant. Therefore, without loss of generality, we choose the first component of the eigenvector for each time and frequency component to be real and positive.

Also, in our case the visibility, $\mathbf{V}^{(X)}$ (where $X = \{H, V\}$) is a block diagonal matrix. Therefore, the eigenvalues and the eigenvectors of the matrix will be the eigenvalues and eigenvectors from each of the blocks, i.e. $\mathbf{V}^{(X)} = \mathcal{E}^{(X)}\Lambda^{(X)}\mathcal{E}^{(X)-1}$. For each t and ν , \mathcal{E} is a $n \times n$ matrix whose i -th column is the complex normalized eigenvector, $\mathcal{E}_i^{(X)}$ of \mathbf{V} , and Λ is the diagonal matrix whose diagonal elements, $\Lambda_{ii} = \lambda_i$, are the corresponding eigenvalues.

Fig. 9 shows 16 eigenvalues calculated from the horizontal polarization matrix ($\mathbf{V}^{(H)}$) as a function of time. The plot clearly shows that one eigenvalue is much higher than the other values during daytime. We can undoubtedly infer that the major contribution to the power in that particular eigenvalue comes from the solar contamination, as the Sun is by far the strongest source in the sky during daytime.

Fig. 10 shows the phase of one of the components of the eigenvector that corresponds to the largest eigenvalue: $\mathcal{E}_{(15,16)}^{(X)}$. Eigenvalues are

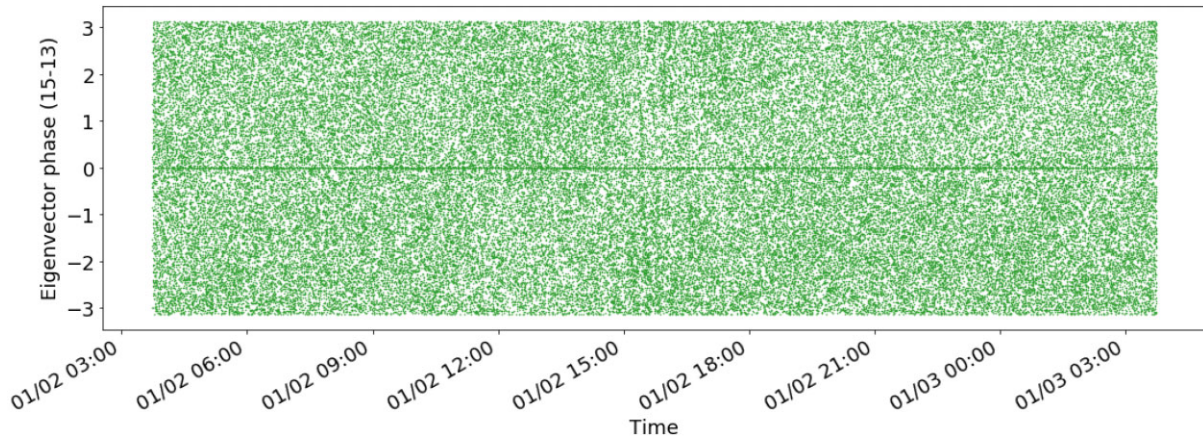


Figure 11. Phase plot of the 15th component of the eigenvector corresponding to the fourth largest eigenvalue in the horizontal polarization visibility $V^{(H)}$. Here, we do not see any fringes, showing that no individual strong source is contributing to this particular eigenvalue.

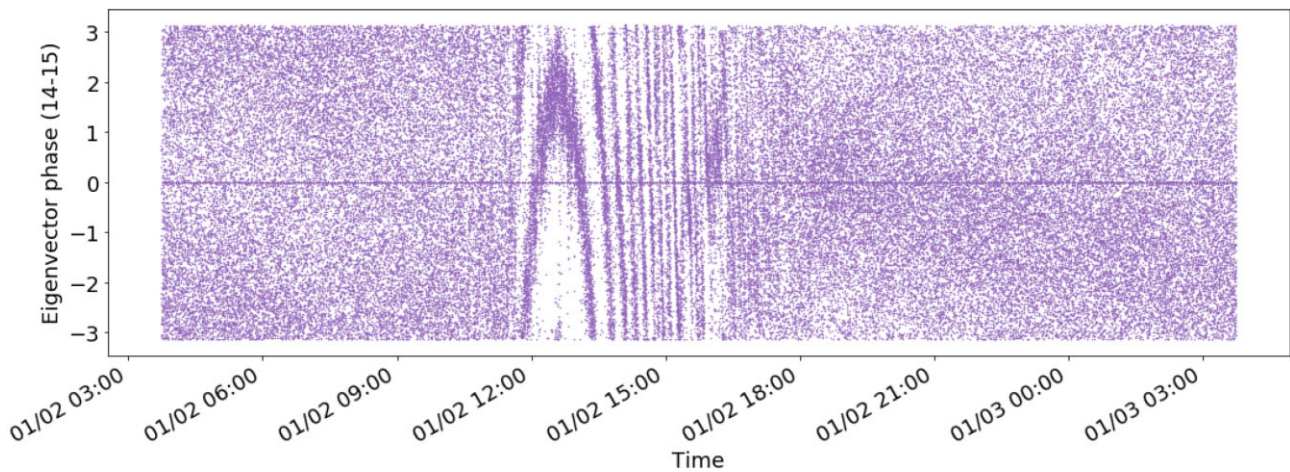


Figure 12. The phase of one component of the eigenvector corresponding to the second largest eigenvalue. We can see weak fringes, indicating that some of the solar contribution is present in the second largest eigenvalue.

sorted according to the daytime amplitude. The 16th eigenvalue is the largest and we have shown the 15th component of the corresponding eigenvector.

Clear fringes are visible in the daytime data, showing that the daytime signal in the eigenvector is coming from a single strong source. In the nighttime data, we can see that the phase varies completely randomly, proving the absence of any single strong source at the nighttime data.

In Fig. 11, we show the phase from one of the components of the fourth largest eigenvalue, $\mathcal{E}_{(15,13)}^{(X)}$. Unlike Fig. 10, no fringes are visible, indicating that there is no single strong source being detected by that particular baseline and the signal is coming from the background sky. The same thing is true for any of the other smaller eigenvalues. In Fig. 12, we have plotted the phase from one component of the eigenvector corresponding to the second largest eigenvalue. We can find weak fringes during the daytime, indicating that some of the Sun signal has ‘leaked’ into this eigenvector. Ideally, the second eigenvalue represents the second strongest sources in the sky, and this leakage may be due to the presence of other sources and the background noise, which includes diffuse sources from the sky and thermal noise. In addition, the re-normalization of the autocorrelation signal using equation (7) is another possible cause of this leakage. Finally, it may also be that some of the solar radiation

is being reflected from the ground and illuminating the feeds from a different direction from the main Sun signal.

If we plot the phase from any component of the eigenvector corresponding to the third largest eigenvalue, we can still see some fringe pattern in the daytime data. However, these fringes are much weaker in comparison to $\mathcal{E}_{(15,x)}^{(X)}$ showing that the leakage of solar power is mostly restricted to the second largest eigenvalue.

3.4 A first attempt to subtract the solar contamination signal

Because the solar signal is contributing mostly to the largest eigenvalue, as a first step in removing the Sun signal, we can set the largest eigenvalue during the daytime data to 0 and then reconstruct the visibility. In Fig. 13, we show the largest eigenvalue as a function of time (in blue during the daytime). The red curve shows the value after setting the largest eigenvalue during daytime to be 0. All the other eigenvalues are kept fixed. In Fig. 14, we show the waterfall plot of the complex visibility of one baseline that we recover after this step. The plot shows that most of the contamination is removed. However, some solar contamination signal is still discernible in the visibility plot. We can see clear, faint fringes for the baseline plotted in Fig. 14.

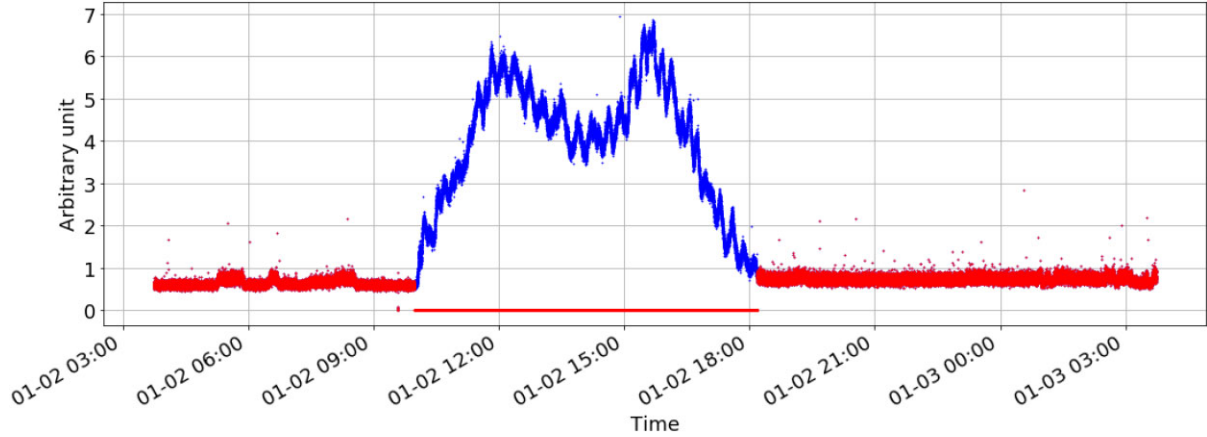


Figure 13. Blue: The original largest eigenvalue from the eigendecomposition during daytime. Red: The same largest eigenvalue but with zero value during the daytime. This step simulates removing the Sun signal in equation (11), since the largest eigenvector during the daytime points in the direction of the Sun in the eigenspace.

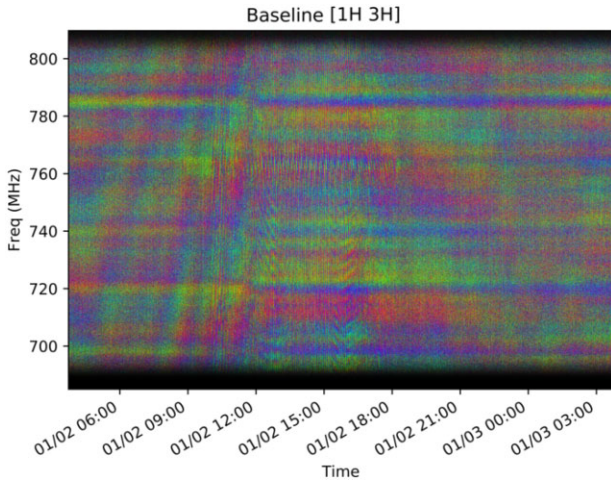


Figure 14. The waterfall plot of the complex visibility after zeroing the largest eigenvalue during daytime for a typical baseline. Nightly mean subtraction has been applied. There is still some residual solar contamination signal, which is causing the weak fringes during the daytime.

4 IMPROVING THE SUN SUBTRACTION

As we can see from Fig. 14, some of the solar signal still remains in the visibility matrix, mainly the signal that leaked to the second largest or even to the third largest eigenvalue. We attempt to remove this residual signal through the following two steps.

4.1 Smoothing the eigenvalues and eigenvectors

The problem with the direct eigenvalue removal method described above is that it is based on the assumption that there is a single source on the sky. This is not true in this case, as the visibility matrix contains signals from other sources as well as instrument noise. In Fig. 15, one component of the eigenvector corresponding to the largest eigenvalue is shown in light green over a period of 4 h. The data, sampled every second, are noisy. However, as the Sun moves smoothly and the beam is not expected to be structured on small scales, we expect the eigenvector to vary smoothly with time. These fluctuations in the eigenvalue probably come from noise. The long-term (minute level and longer) fluctuations originate in the structure of the sidelobes of the telescopes.

As the noise in the visibility matrix may cause the Sun signal to leak from the largest eigenvalue to other eigenvalues, in this section we try to reduce the effect of noise. For doing that, we fit a smooth curve (black line) through the eigenvectors corresponding to the Sun signal. This smoothed signal from the largest eigenvector is then subtracted from the original visibility to construct the Sun-removed visibility.

The cleaning routine can be summarized as follows. The visibility matrix $\mathbf{V}^{(X)}$ is first decomposed into the eigenvalue and the eigenvectors for each time and each frequency bin,

$$\mathbf{V}^{(X)} = \mathcal{E}^{(X)} \Lambda^{(X)} \mathcal{E}^{(X)-1}, \quad X = \{H, V\}. \quad (9)$$

Suppose \mathcal{E}_S is the eigenvector corresponding to the largest eigenvalue, λ_S . As shown in Fig. 15, the direction of \mathcal{E}_S will vary in every second. The n -dimensional complex eigenvectors have only $2n - 1$ degrees of freedom as we have already set the first component to be real and positive. As the eigenvectors are unit vectors, the total number of independent components becomes $2n - 2$. If we fit a smooth line through each of the $2n - 1$ components, then we will overfit and the amplitude of the eigenvectors will not be 1. To keep the eigenvector normalized while doing the fitting, we express each (complex) component of the eigenvector in n -dimensional spherical coordinates and then fit a smooth line through the tangents of the angles in spherical coordinates and convert back to Cartesian space. This gives the black line, shown in Fig. 15. Smoothing in spherical coordinates ensures that the normalization of the eigenvector is preserved during the smoothing procedure. (See Appendix A for details.)

Let the smoothed components of the largest eigenvectors be $\tilde{\mathcal{E}}_S$. If $\tilde{\lambda}_S$ is the contribution to the visibility from the direction of the eigenvector $\tilde{\mathcal{E}}_S$, then we can write, $\tilde{\lambda}_S = \tilde{\mathcal{E}}_S^T \mathbf{V} \tilde{\mathcal{E}}_S$. If we assume that this smoothed component comes from the Sun signal, then the contribution to the visibility from the Sun is given by

$$\tilde{\mathbf{V}}_S = \tilde{\lambda}_S [\tilde{\mathcal{E}}_S \otimes \tilde{\mathcal{E}}_S]. \quad (10)$$

After subtracting the Sun signal, the contribution to the visibility from the rest of the radio sky and noise is given by

$$\mathbf{V}_{\text{sky}} = \mathbf{V} - \tilde{\mathbf{V}}_S. \quad (11)$$

In Fig. 16, we show the complex visibility after removing the Sun signal using this particular algorithm. In comparison to the simplest algorithm, of just removing the largest eigenvalue, this new algorithm works better. However, we can see that some of the Sun signal is still present in the visibility.

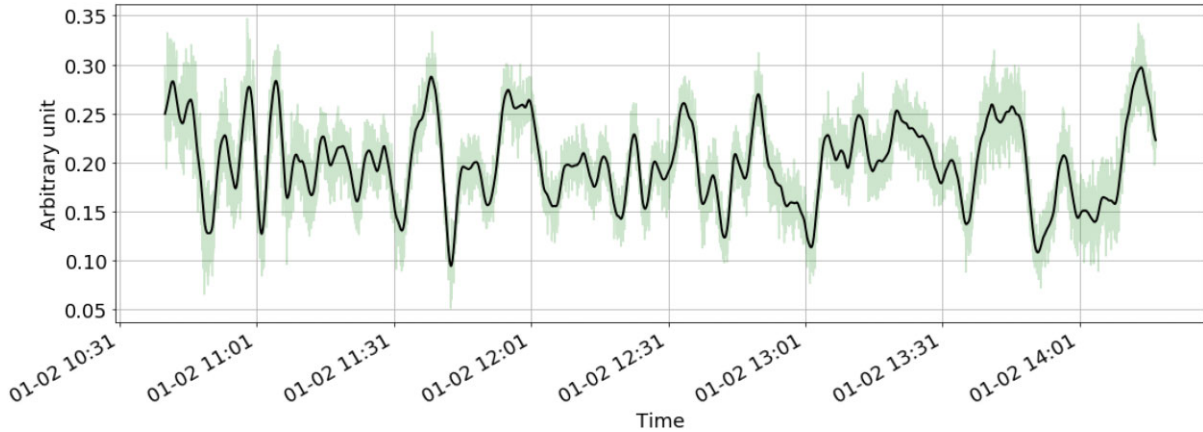


Figure 15. A 4-h segment of the amplitude of one component of the eigenvector corresponding to the largest eigenvalue is shown in the light green curve. The sampling interval is 1 s. The random fluctuations in the data come from the noise. The black curve shows this component after smoothing the data, as described in Section 4.1.

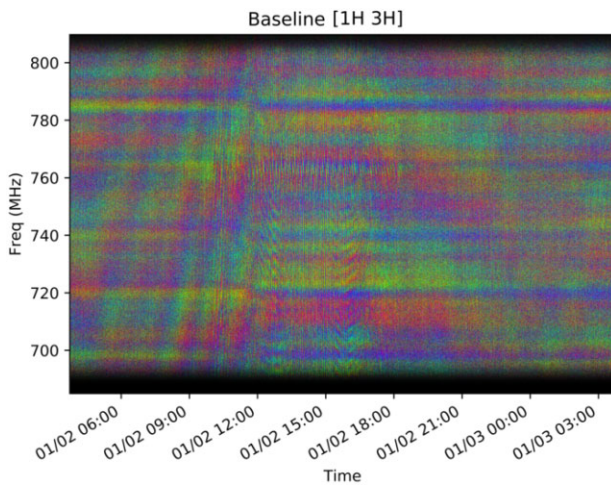


Figure 16. The complex visibility after smoothing the components of the largest eigenvector in spherical coordinates (described in Section 4.1, before scaling by a gain factor, described in Section 4.2) for a typical baseline. Nightly mean subtraction has been applied. There is still some signal from the Sun that has not been removed.

4.2 Scaling the Sun signal from eigenvalue analysis

The above eigenvalue analysis is based on the assumption that the signal coming from Sun is contained in the largest eigenvalue. As discussed before, this assumption is not correct because of the leakage of power into other eigenvalues.

To overcome this issue, we consider that during the daytime the signals from the sky are much smaller than the solar signal. Therefore, the Sun signal, $\mathbf{V}_S(t, \nu)$, calculated from our analysis should roughly match with the visibility $\mathbf{V}(t, \nu)$ during the daytime as the other signals are negligible in comparison to the $\mathbf{V}_S(t, \nu)$, provided that there are no other strong sources during the day. To do that, we introduce a scaling (gain) factor, $g = Ae^{i\phi}$, for each 1000 s (about 15 min) of daytime data and minimize

$$\chi^2 = \sum_{t, \nu} [\Re(\mathbf{V} - g\mathbf{V}_S(t, \nu))]^2 + \sum_{t, \nu} [\Im(\mathbf{V} - g\mathbf{V}_S(t, \nu))]^2. \quad (12)$$

Here, $\Re()$ and $\Im()$ are the real and imaginary parts of the quantity inside the bracket. We get 36 gain factors (g), calculated from 10 h

(36 000 s) of daytime data. In Fig. 17, we show the plot of g over 10 h of daytime, with circular dots. The smooth lines show the interpolated data. We can see that g varies smoothly throughout the day. The expectation is that the $|g|$ should be very close to 1 and very smooth, and the phase variation should be very small. This is because the Sun and the sky move smoothly through the beams over the day. As long as the Sun signal is strong enough in comparison to the background sky, we can expect that the power leakage will vary smoothly and the gain variation should also be smooth. As the signal in the largest eigenvector and leaked power both are coming from Sun, we can expect the phase variation to be minimum. Fig. 17 shows that the assumption is a good one in this case. However, near sunrise and sunset the amplitude and the phase change rapidly, possibly because the Sun signal is weaker at those times.

The interpolated g is used as a multiplication factor to determine the solar contribution $g \times \tilde{\mathbf{V}}_S(t, \nu)$, which is finally subtracted from $\mathbf{V}(t, \nu)$. This gives our final Sun-removed signal from the daytime data, i.e.

$$\mathbf{V}_{\text{sky}}(t, \nu) = \mathbf{V} - g_{\text{int}}(t, \nu)\tilde{\mathbf{V}}_S(t, \nu). \quad (13)$$

In Fig. 18, we show the complex visibility after the solar contamination removal using equation (13). We can see by visual inspection that most of the contamination signal is removed and the fringes from the weaker sources in the background sky are visible. This is the best that we get from AlgoSCR. However, on closer inspection we can see that a small amount of the Sun signal is still present in the data in the form of weak fringes. In the next section, we will make a first estimate of the performance of our solar signal subtraction and its effect on the signals from the fainter sources.

5 TESTING THE EFFICIENCY OF THE ALGORITHM

5.1 Comparison with uncontaminated data

Due to the orbital motion of Earth, the solar signal contaminates observations made at different sidereal times, or sky orientations. Therefore, to quantify the fraction of the solar contamination removed by AlgoSCR, in this section we compare the Sun-removed visibility to the uncontaminated visibility observed during the same sidereal time at an interval of four sidereal months. In the left plot

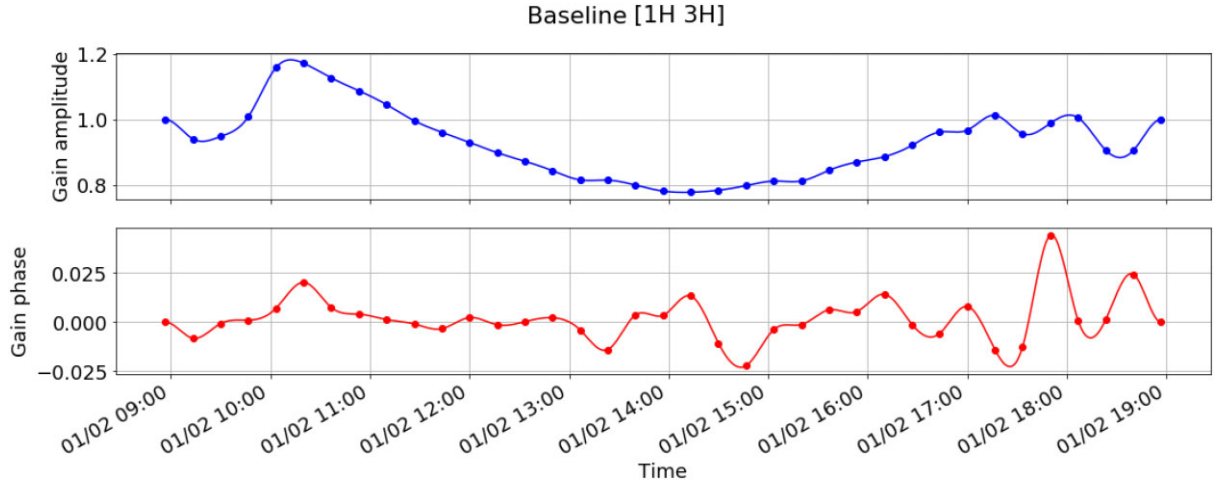


Figure 17. Plot of the gain $g = Ae^{i\phi}$ after χ^2 minimization for 10 h during the daytime. Blue: Plot of the gain amplitude A . Red: Plot of the gain phase ϕ in radians. The dots represent the points of χ^2 minimization that occur every 1000 s. These gain values are extrapolated to the intervening points for a total of 36 000 s.

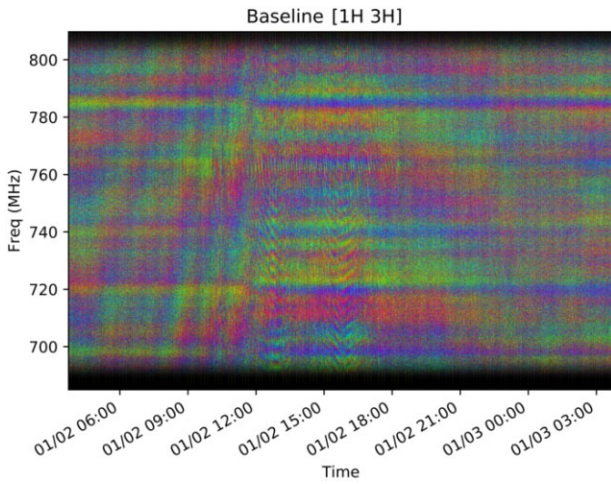


Figure 18. The waterfall plot of the complex visibility after χ^2 optimization scaling the Sun signal by a gain factor (Section 4.2) for a typical baseline.

of Fig. 19, we show the raw complex visibility from 2018 January (same as Fig. 6). The middle plot is the visibility after the solar contamination removal using AlgoSCR (same as Fig. 18). Finally, the plot on the right shows the complex visibility of the sky observed in

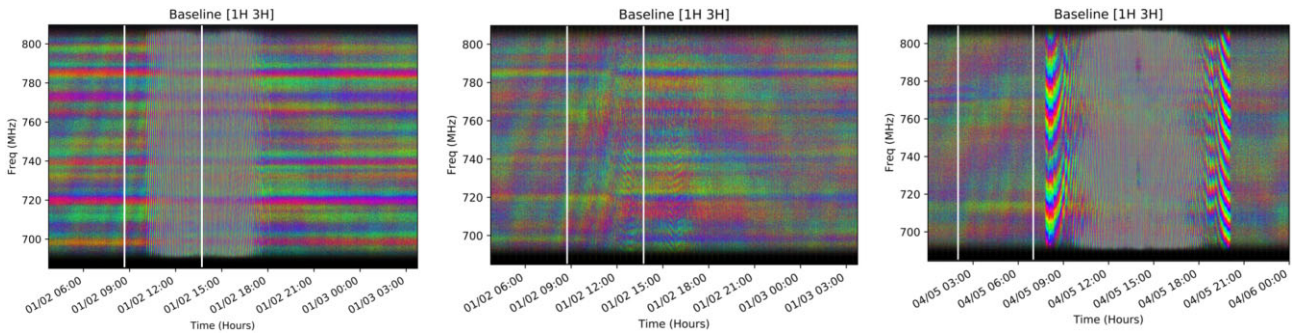


Figure 19. Left: Waterfall plot of original (Sun-contaminated) 2018 January complex visibility over 24 h (same as Fig. 6). Middle: waterfall plot after applying removing solar contamination with AlgoSCR (same as Fig. 18). Right: waterfall plot of visibility data taken in 2019 April. The areas between the white lines show periods of sidereal time when observations occur during daytime in January, and nighttime in April.

2019 April. All plots show the complex visibilities over one sidereal day. About 5 h of 2019 April, visibility plot is not contaminated by the solar signal, and we mark that period with two white lines in all the plots.

The plots show that visibilities of the same sky orientations are similar. The large fringes in the Sun-cleaned plot coincide with those from nighttime of 2019 April. Here one should note that, in the plot in the middle, we can also see some additional weak, rapidly oscillating fringes, which are not present in the 2019 April data. However, upon careful inspection, we can see that these ‘fast fringes’ are originally present in the left plot (before Sun removal). Therefore, the algorithm does not introduce any obvious additional signal.

For understanding these results quantitatively, we take 1-min time averages and 384-bin frequency averages and then calculate the ratio of the residual Sun signal to the background signal over the same 5-h range of sidereal time:

$$\frac{\sum_{t,v} |\mathbf{V}_{\text{Sky}}(t, \nu) - \mathbf{V}_{\text{org-2019}}(t, \nu)|}{\sum_{t,v} |\mathbf{V}_{\text{org-2019}}(t, \nu)|} = 0.37. \quad (14)$$

Here, $\mathbf{V}_{\text{org-2019}}$ is the original (i.e. uncleaned) visibility from 2019 April. The ratio of the original daytime signal to the background

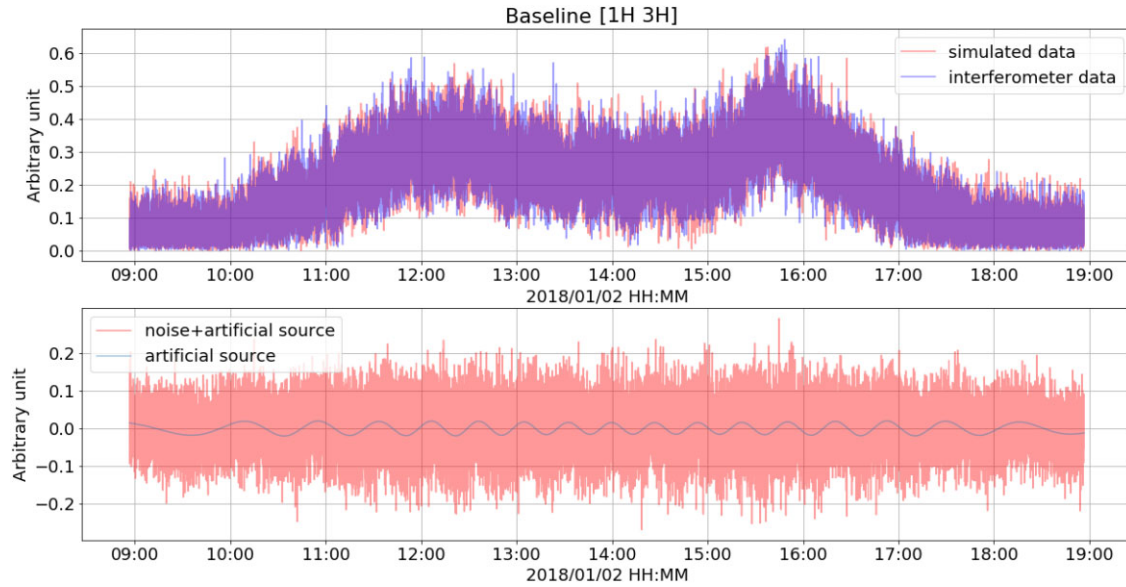


Figure 20. Top: The amplitude of the visibility for baseline [1H 3H]. The actual data from the Tianlai dish array are coloured blue, while the simulated data are shown in red. The actual data are very similar to the simulation; regions of overlap appear purple. Bottom: In red is the real part of the noise plus the artificial sources. The blue line shows the real part of the signal from the artificial sources that is added to the data. The imaginary part (not shown) is similar to the real part.

signal over the same sidereal time is:

$$\frac{\sum_{t,v} |\mathbf{V}_{\text{org-2018}}(t, v) - \mathbf{V}_{\text{org-2019}}(t, v)|}{\sum_{t,v} |\mathbf{V}_{\text{org-2019}}(t, v)|} = 7.48. \quad (15)$$

Assuming that the difference measured by equation (14) is dominated by residual Sun contamination, we see that AlgoSCR reduces the solar contamination by about 95 per cent.

5.2 Comparison with simulated data

We test the efficiency of AlgoSCR when applied to simulated data sets. This provides us with another way to check the fraction of the solar contaminant signal that is removed and to determine how much of the sky signal we are erroneously removing by the analysis.

5.2.1 Construction of simulated data

For constructing a simulated visibility signal \mathbf{V}_{sim} , we assume that the electric field at each feed antenna contains contributions from Sun, the sky, and noise. We have assumed that the noise variance is the same throughout the analysis.

The receiver noise is modelled as Gaussian noise in the electric field, $E_{\text{noise } i}$, at the feed antenna. We consider the noise contribution to the electric field to be Gaussian in each sample.

In the Tianlai dish array, the integration time in the correlator is 1 s. The correlator takes in the data that are collected every few microseconds and averages them in an interval of 1 s. To simulate this process, we add Gaussian random noise in the electric field with a sampling interval of 10 ms. We then calculate the noise contribution to the visibility as $V_{\text{noise}(i,j)} \equiv \langle E_{\text{noise } i}^* E_{\text{noise } j} \rangle_{\tau_{\text{int}}}$, where $\langle \rangle_{\tau_{\text{int}}}$ represents the ensemble average over integration period, $\tau_{\text{int}} = 1$ s. Here, we have 100 data points for every second on which the average is carried out. This method also ensures that the autocorrelation visibilities follow a χ^2 distribution and the cross-correlation visibilities follow a product normal distribution. The mean and variance of $E_{\text{noise } i}$ are chosen empirically so that the simulated visibility, \mathbf{V}_{sim} ,

matches the observed visibility. The mathematical details on how to calculate the visibilities from artificial point sources in the sky are shown in Appendix B.

To create the simulated Sun signal, we have taken the largest eigenvalue and corresponding eigenvector from equation (A6) from the Tianlai dish array data and treated it as the solar signal. The electric field for the Sun, thus calculated, is added to the simulated noise.

For the simulated artificial sources, we assume that the telescope array is pointed at the NCP. The artificial sources are three made-up sources near the NCP. All the artificial sources are visible within the main beam, which is assumed to be Gaussian. Their brightness is chosen so that the amplitude of their combined visibility is about 10 times smaller than the noise. (An analysis with different source strengths is presented in the next section.) The artificial source visibilities are frequency- and baseline-dependent, just as visibilities from real sources on the sky. We also assume that the visibilities for the Sun and artificial sources are uncorrelated, i.e. there is no cross-term between the Sun and the artificial sources. This makes the visibilities for the Sun and artificial sources additive, as shown in equation (16). Please check Appendix B for details.

$$\mathbf{V}_{\text{sim}} = \mathbf{V}_{\text{noise}} + \mathbf{V}_{\text{S}} + \mathbf{V}_{\text{artificial sources}} \quad (16)$$

5.2.2 Results from the simulated data

We generated simulated data as shown in Fig. 20. The top panel of Fig. 20 shows the amplitude of the visibility for baseline [1H 3H] for both simulated and actual Tianlai dish array data: the plot in blue shows the actual complex visibilities from the Tianlai dish array, and the red plot shows the simulated data in our simulation (see equation 16). The bottom panel shows the real part of the signal from the artificial sources (in blue). The real part of the combined signal (simulated noise and the visibility of the artificial sources that is added to the Sun) is shown in red.

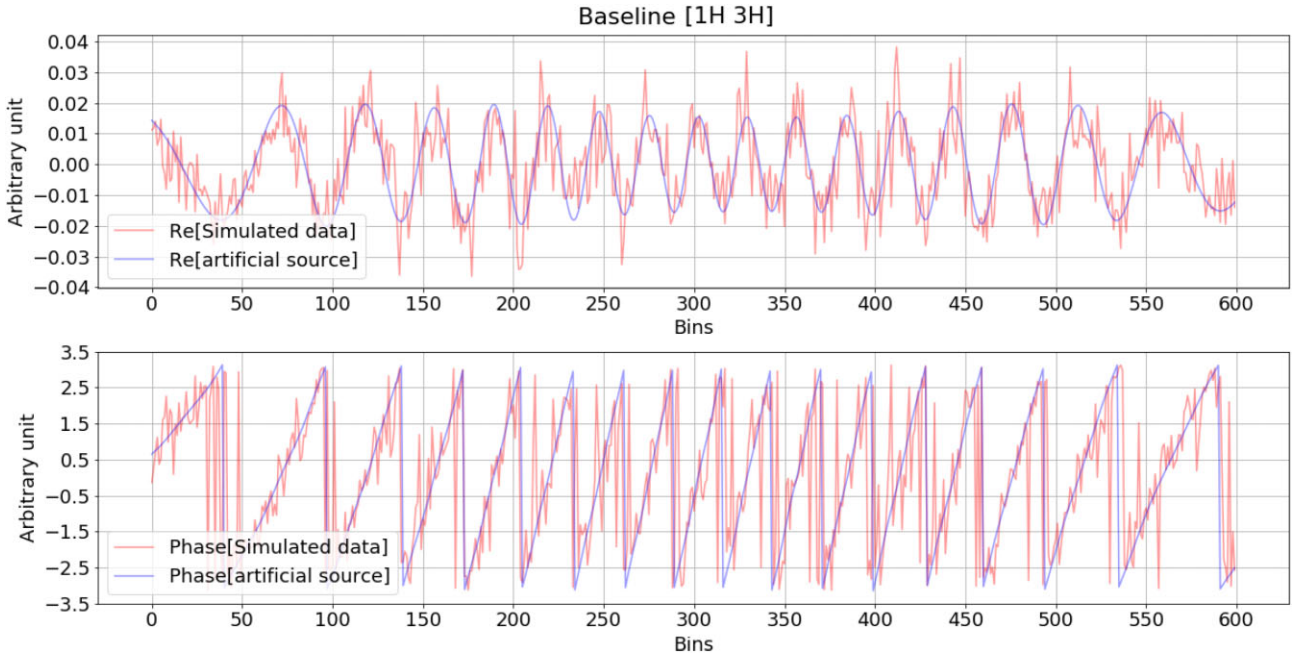


Figure 21. The real part (top) and the phase (bottom) of the Sun-removed visibility from simulated data for baseline [1H 3H] after 60-s averaging. The signal from the artificial sources is shown in blue.

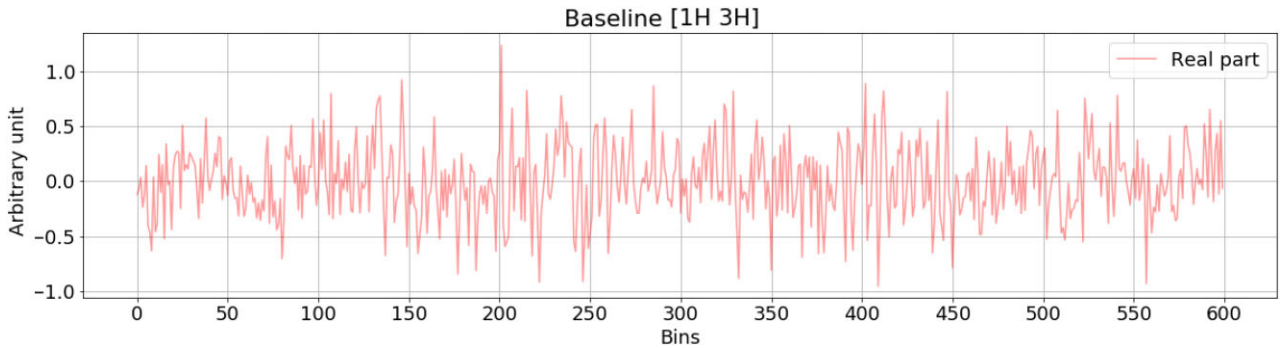


Figure 22. Plot of the real part of $(V_{\text{sim}} - V_{\text{org}})/\sigma$, where σ^2 is the variance of the added noise. We can see that the ratios for the real part are roughly within 3. As the noise is Gaussian, we can expect that the noise signal should be with 3σ . We conclude that the Sun removal algorithm does not introduce any significant additional noise in this analysis.

We apply AlgoSCR to the simulated data. Top panel in Fig. 21 shows the real part of the visibility (in red) after applying the Sun removal algorithm, along with the real part of the visibility of the artificial sources (in blue) for baseline [1H 3H]. As the nature of the imaginary part will be similar, we have not explicitly shown it in the plot. The phase is shown in the bottom panel of the same figure.

The ratio of the difference between the simulated and the original visibility $V_{\text{sim}} - V_{\text{org}}$, and the noise standard deviation, σ , is plotted in Fig. 22. We can see that the ratio is within 3. As the injected noise is Gaussian, we can expect that most of the visibility should also fall within 3σ . Therefore, Fig. 22 ensures that the recovery of the signal using AlgoSCR does not introduce additional noise.

The red plot in Fig. 23 shows the difference between the amplitude of the simulated visibility (including the Sun) and the visibility that we are getting after applying AlgoSCR. This gives the contribution from the Sun in our simulated data. The blue curve shows the Sun signal that we introduced for generating the simulated data. We can see that the plots match very well. Top plot is constructed using the data from each second and the bottom plot is after averaging the data over a minute.

In the next set of plots, Fig. 24, we show the complex visibilities for one baseline, before and after the solar contamination removal by AlgoSCR. The visibility data show that the artificial sources that we had introduced are clearly visible after the solar signal removal, even though the source strength was much smaller than the Sun signal and the noise. This plot shows qualitatively the potential of the Sun removal algorithm. In the next section, we quantify the amount of the signal from artificial sources that is removed.

5.2.3 Comparing efficiency of the method for different external source strengths

Here, we compare the efficiency of AlgoSCR in recovering the artificial sources for different source strengths. For this analysis, we use the real daytime visibility data taken by the Tianlai dish array as the base visibility. To these data, we add the artificial visibility signal with different source strengths. We assume that sources are not correlated with the visibility data and the visibilities are additive.

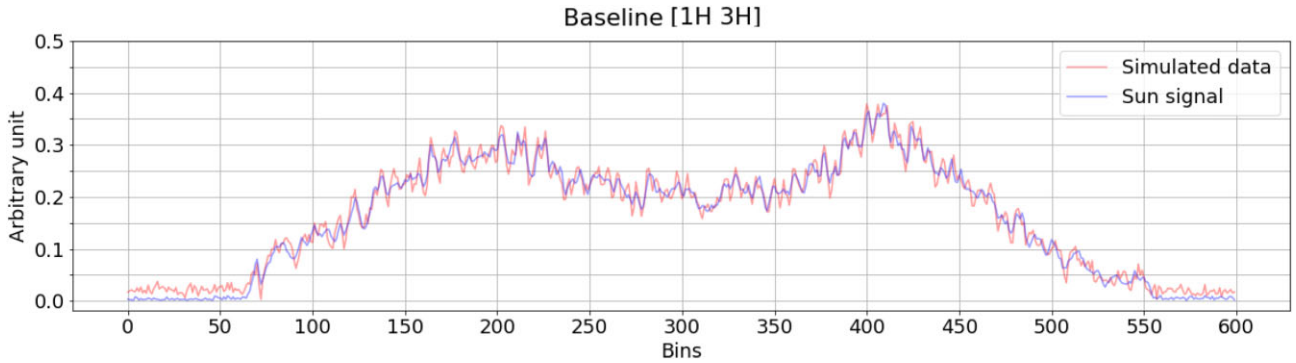


Figure 23. The difference between the simulated visibility (shown in Fig. 20) and the Sun-removed visibility (shown in Fig. 21) is shown in red with 1-s averaging. The signal from the largest eigenvector, which is used as the Sun signal during the daytime, is shown in blue. The data are averaged in 60-s time bins to reduce the noise.

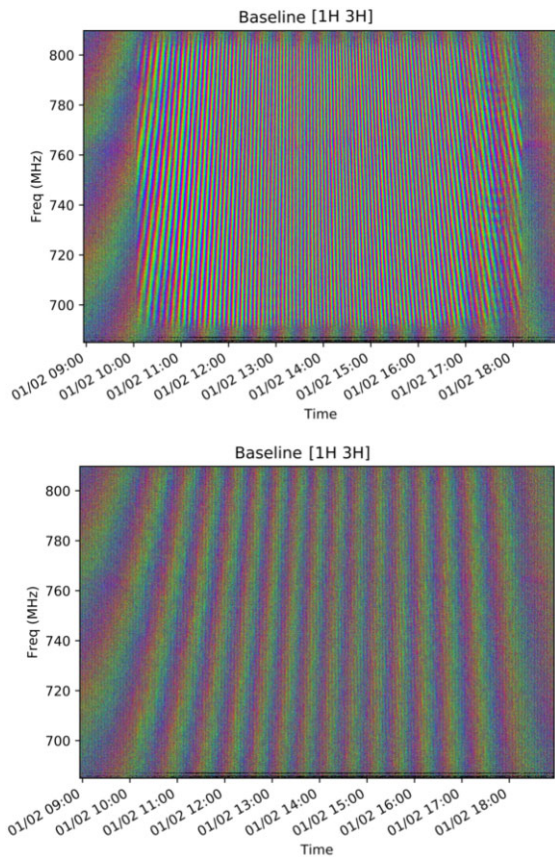


Figure 24. Top: Simulated daytime visibilities for a typical baseline. The bright fringes are from the Sun and the weaker fringes are from the artificial sources. Bottom: Sun-removed visibility using AlgoSCR. The fringes from the artificial sources are clearly visible.

After running AlgoSCR to remove the Sun, we use a χ^2 statistic to compare the signal with the source visibility that was originally inserted. The plot of the reduced χ^2 for baseline [1H 3H] is shown in Fig. 25 against the source strength. Here, χ^2 is defined as $\chi^2 = \frac{1}{n_t \sigma^2} \sum_{t,v} |\mathbf{V}_{\text{org}}(t, v) - \mathbf{V}_{\text{sim}}(t, v)|^2$ during 10 h of daytime. Here, σ^2 is the noise variance and n_t is the number of time-steps, which is the number of degrees of freedom in this case. As we are sampling each second for a total of 10 h, the number of degrees of freedom $n_t = 36000$.

We can see that the χ^2 value is small for the cases in which the artificial source amplitudes are small compared to the Sun signal amplitude. As the artificial source strengths increase, the fit gets worse. This is because our analysis is based on the assumption that the solar signal is the only dominant signal. As the strength of the artificial sources increases, the assumption slowly breaks down. In such cases, the largest eigenvalue starts to capture signal from the artificial sources. When the artificial sources are larger than the Sun, the largest eigenvalue provides the contribution from the artificial sources and not the Sun. In such cases, we are essentially removing the artificial sources and thus the χ^2 grows quadratically.

In Fig. 26, we have plotted the same χ^2 , where instead of dividing by σ^2 we have divided by $|\mathbf{V}_{\text{org}}|$. Here, we can see that the χ^2 is lowest when the strength of the artificial source is about 40 per cent of the Sun signal. When the source strength is small, the χ^2 is dominated by the noise and the χ^2 is high. On the other hand, when the strength of the artificial sources is high compared to the Sun contamination as described before, the recovery gets worse.

6 DISCUSSION

While developing AlgoSCR, we explored multiple techniques and came across various issues. Here, we discuss some of the points that are important in the context of optimizing AlgoSCR.

In Section 4, we address the issue of removing the residual Sun signal after subtraction of the largest eigenvalue. Here, we introduce the concept of the multiplication factor g . This procedure raises the question of what will happen if instead of filtering out just the largest eigenvalue, we filter out a smooth component from the two largest eigenvalues. We find that removing the two largest components after smoothing, then the signal from the second largest component, which includes some radio sources, also gets removed, i.e. we will be removing the components from other radio sources and hence the method will not work.

In Section 4.2, while choosing the gain values, g , we calculate the gains at intervals of 15 min and then interpolate. We find that the results are fairly insensitive to the choice of time interval (say, 10 min or 30 min), as is expected because the gain varies smoothly throughout the day. However, if we choose a long time interval (several hours) for setting the gains, then we expect the results to worsen as the gain may change significantly in that time. However, we have not simulated these cases.

Another important fact that came up during our analysis is that the Sun removal algorithm works better with more baselines, i.e.

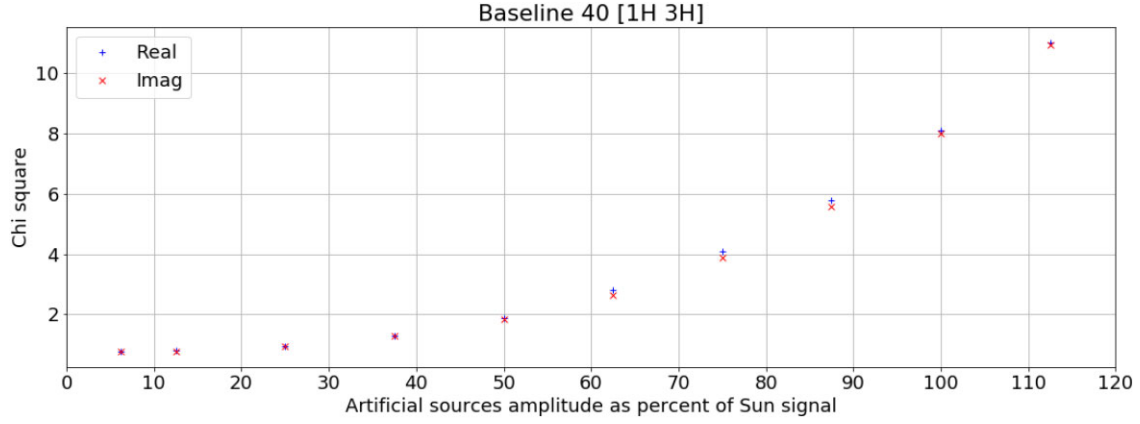


Figure 25. Plots of $\chi^2 = \frac{1}{n_t \sigma^2} \sum_{t,v} (V_{\text{org}} - V_{\text{sim}})^2$ from the real and imaginary parts of the visibility for different artificial source amplitudes. n_t is the number of sample points in the time direction and σ^2 is the noise variance. The amplitude of the original visibility V_{org} and the simulated visibility V_{sim} are shown in Fig. 21. We can see that the χ^2 is increasing as we increase the amplitude of the artificial sources.

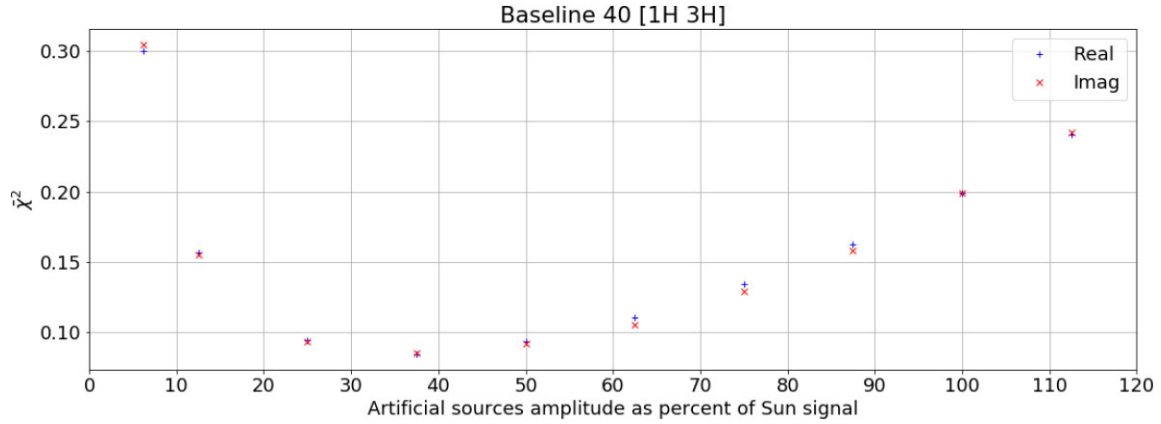


Figure 26. Plots of $\bar{\chi}^2 = \frac{1}{n_t |V_{\text{org}}|} \sum_{t,v} (V_{\text{org}} - V_{\text{sim}})^2$ from the real and imaginary parts of the visibility for different artificial source amplitudes. Here, the plots are normalized by $|V_{\text{org}}|$ instead of σ^2 . We can see that the $\bar{\chi}^2$ is lowest at about 40, indicating that the recovery is best when the amplitude of the source is about 40 per cent of the Sun signal.

if we use all 16 dishes from Tianlai instead of, say, 10 dishes, then the effectiveness of the algorithm increases. Analysis with fewer dishes increases the power leakage to other eigenvalues. The exact reason behind this is not known, but it may happen as more baselines reduce the effect of the noise in the eigendecomposition.

In addition, if we increase the integration time from 1 s to a larger value, the results get worse, which may be due to the fact that the Sun is not a point source. This is different from co-adding the signal from multiple days, which is eventually what the Tianlai array is designed to do. However, we have not tested the algorithm on co-added signals yet.

Our analysis shows that AlgoSCR removes most of the solar contamination during the day. However, it is just a first step. We have not yet tested its effect on map-making and power spectrum estimation. A critical next step is to make sure that AlgoSCR does not affect the statistics of the maps. This can be checked by comparing the HI power spectra and other statistical quantities from the maps produced using only nighttime data and the maps produced using the full day data after solar contamination removal. Such an analysis requires foreground subtraction and map making and is outside the scope of the present work.

7 CONCLUSION

In this paper, we present a way to separate out the solar contamination from the daytime data observed by an interferometric radio array using eigendecomposition techniques. The technique is primarily based on the assumption that if the Sun signal is the dominant signal in the sky, along with other weaker sources, and if the signals from the different sources are not correlated, then in the eigendecomposition of the visibility matrix, the largest eigenvalue is from the strongest source, i.e. the Sun. The eigenvector corresponding to the largest eigenvalue points in the direction of that source in the eigenspace. The technique should filter out this largest eigenvalue while retaining the signals from other sources in the sky.

However, antenna gain fluctuations, noise, sidelobe gain patterns, ground reflection, thermal effects on the instruments and cables, and crosstalk between antennas introduce mixing between the largest eigenvalue and other smaller eigenvalues. For these reasons, singling out and removing the Sun signal is not straightforward, and there is some residual contamination from the Sun. Therefore, we apply some novel techniques to remove the leftover Sun signal.

We have tested AlgoSCR in two ways. First, we compared the visibilities obtained by cleaning observations made during sidereal times when the Sun was up with observations made during those

same sidereal times at night. We show that AlgoSCR can reduce the solar contamination by a factor of 95 per cent. Secondly, we used simulations to show that our algorithm is able to remove the solar contamination without removing other, weaker sources in the sky. We showed that the efficacy of the algorithm is maximum when the amplitude of the external source is about the 40 per cent of the solar contamination signal. The fraction of the removed background signal remains significantly small when the external source strength is within a range of 20 per cent to 65 per cent solar contamination.

To the best of our knowledge, this is the first published method for removing solar contamination from radio interferometer data. AlgoSCR can contribute to other ongoing and upcoming radio interferometers for solar contamination removal.

ACKNOWLEDGEMENTS

The authors wish to thank Richard Shaw, Jeff Peterson, and John Marriner for several fruitful discussions about the project during the Tianlai meeting in Guizhou. We wish to thank Kevin Gayley for allowing us to use his electromagnetic simulations of the Tianlai dish array beam pattern. We thank Juyong Zhang and his team at Hangzhou Dianzi University for letting us use the beam mapping data from their drone survey.

Work at UW-Madison and Fermilab is partially supported by NSF Award AST-1616554.

This research was performed using the compute resources and assistance of the UW-Madison Center For High Throughput Computing (CHTC) in the Department of Computer Sciences. The CHTC is supported by UW-Madison, the Advanced Computing Initiative, the Wisconsin Alumni Research Foundation, the Wisconsin Institutes for Discovery, and the National Science Foundation, and is an active member of the Open Science Grid, which is supported by the National Science Foundation and the U.S. Department of Energy's Office of Science.

The Tianlai array is operated with the support of NAOC Astronomical Technology Center. The work at NAOC is supported by the Ministry of Science and Technology of China under grants 2018YFE0120800, 2016YFE0100300, and 2012AA121701, the National Natural Science Foundation of China under grants 11633004, 11473044, 11761141012, 11653003, 11773031, and Chinese Academy of Science grants QYZDJ-SSW-SLH017, XDA15020200, ZDKYYQ20200008.

This research includes personnel and uses resources of the Fermi National Accelerator Laboratory (Fermilab), a U.S. Department of Energy, Office of Science, HEP User Facility. Fermilab is managed by Fermi Research Alliance, LLC (FRA), acting under contract no. DE-AC02-07CH11359.

DATA AVAILABILITY

The data underlying this article will be shared on reasonable request to the corresponding author.

REFERENCES

- Abdalla F. B., Rawlings S., 2005, *MNRAS*, 360, 27
 Anderson C. et al., 2018, *MNRAS*, 476, 3382
 Ansari R. et al., 2018, preprint (arXiv:1810.09572)
 Bandura K. et al., 2014, in Stepp L. M., Gilmozzi R., Hall H. J., eds, Proc. SPIE Conf. Ser. Vol. 9145, Ground-based and Airborne Telescopes V. SPIE, Bellingham, p. 914522
 Battye R. A., Davies R. D., Weller J., 2004, *MNRAS*, 355, 1339

- Battye R. et al., 2012, preprint (arXiv:1209.1041)
 Briggs F. H., Bell J. F., Kesteven M. J., 2000, *AJ*, 120, 3351
 Bull P., Ferreira P. G., Patel P., Santos M. G., 2015, *ApJ*, 803, 21
 Chang T.-C., Pen U.-L., Peterson J. B., McDonald P., 2008, *Phys. Rev. Lett.*, 100, 091303
 Chen X., 2011, *Scientia Sinica Physica Mechanica & Astronomica*, 41, 1358
 Chen X., 2012, *Internat. J. Modern Phys.: Conf. Ser.*, 12, 256
 Das S. et al., 2018, in Zmuidzinas J., Gao J.-R., eds, Proc. SPIE Conf. Ser. Vol. 10708, Millimeter, Submillimeter, and Far-Infrared Detectors and Instrumentation for Astronomy IX. SPIE, Bellingham, p. 1070836
 DeBoer D. R. et al., 2017, *PASP*, 129, 045001
 Dickinson C., 2014, preprint (arXiv:1405.7936)
 Eastwood M. W. et al., 2018, *AJ*, 156, 32
 Hu W., Wang X., Wu F., Wang Y., Zhang P., Chen X., 2020, *MNRAS*, 493, 5854
 Li J. et al., 2020, *Sci. China Phys. Mech. Astron.*, 63, 129862
 Liu A., Shaw J. R., 2020, *PASP*, 132, 062001
 Mao Y., Tegmark M., McQuinn M., Zaldarriaga M., Zahn O., 2008, *PRD*, 78, 023529
 Masui K. W., McDonald P., Pen U.-L., 2010, *Phys. Rev. D*, 81, 103527
 Masui K. W. et al., 2013, *Astrophys. J. Lett.*, 763, L20
 Morales M. F., 2008, *AIP Conf. Proc.*, 1035, 82
 Newburgh L. B. et al., 2014, in Stepp L. M., Gilmozzi R., Hall H. J., eds, Proc. SPIE Conf. Ser. Vol. 9145, Ground-based and Airborne Telescopes V. SPIE, Bellingham, p. 91454V
 Newburgh L. et al., 2016, in Hall H. J., Gilmozzi R., Marshall H. K., eds, Proc. SPIE Conf. Ser. Vol. 9906, Ground-based and Airborne Telescopes VI. SPIE, Bellingham, p. 99065X
 Paciga G. et al., 2011, *MNRAS*, 413, 1174
 Peterson J. B., Bandura K., Pen U. L., 2006, preprint (arXiv:astro-ph/0606104)
 Shaw J. R., Sigurdson K., Pen U.-L., Stebbins A., Sitwell M., 2014, *Astrophys. J.*, 781, 57
 Slosar A. et al., 2019, *BAAS*, 51, 53
 Tingay S. J. et al., 2013, *PASA*, 30, e007
 Wu F. et al., 2021, *MNRAS*, 506, 3455
 Wuensche C., BINGO-Collaboration, 2019, *J. Phys.: Conf. Ser.*, 1269, 012002
 Xu Y., Wang X., Chen X., 2014, *ApJ*, 798, 40
 Zhang J., Ansari R., Chen X., Campagne J.-E., Magneville C., Wu F., 2016, *MNRAS*, 461, 1950

APPENDIX A: SUMMARY OF ALGOSCR

Here, we review the step-by-step procedure for Sun removal using the algorithm described above.

(i) For this procedure to work, first we separate the visibility \mathbf{V} into the horizontal and vertical polarizations, $\mathbf{V}^{(H)}$ and $\mathbf{V}^{(V)}$, respectively. If we do not separate the polarizations, the noise and crosstalk in the same dish will give an additional large eigenvalue. The dimension of \mathbf{V} is 32×32 , since we have 16 dual-polarization feeds. The dimension of $\mathbf{V}^{(H)}$ and $\mathbf{V}^{(V)}$ will be 16×16 .

(ii) Remove the nighttime mean from the visibility matrix $\mathbf{V}^{(X)}$: $\mathbf{V}^{(X)} = \mathbf{V}^{(X)} - \langle \mathbf{V}^{(X)} \rangle_{\text{night}}$. Here, the average is over the time direction for different frequency channels. This will remove the crosstalk between the antennas.

(iii) Replace the autocorrelations by equation (7). In practice, if the denominator, $\mathbf{V}_{(i,j)}^{(X)}$, is zero, we replace the term inside the sum by a small number, such as 0.0001.

(iv) Perform an eigen-decomposition of $\mathbf{V}^{(X)}$:

$$\mathbf{V}^{(X)} = \mathcal{E}^{(X)} \Lambda^{(X)} (\mathcal{E}^{(X)})^{-1}. \quad (\text{A1})$$

(v) For each second of integration time, let the largest (normalized) eigenvector corresponding to the largest eigenvalue, $\lambda_{S(t,v)}^{(X)}$, be $\mathcal{E}_{S(t,v)}^{(X)}$. Now $\mathcal{E}_{S(t,v)}^{(X)}$ is a vector containing $n = 16$ complex numbers. For

fitting the smooth line through these vectors, we calculate the tangents, $T_{S(t,v)}^{(X)}$ as:

$$T_{S(t,v)}^{(X)}(i) = \frac{\|\mathcal{E}_{S(t,v)}^{(X)}(i)\|}{\sqrt{\sum_{j=i+1}^n (\|\mathcal{E}_{S(t,v)}^{(X)}(j)\|)^2}} \quad \forall i \in [1, n-1]. \quad (\text{A2})$$

(vi) For smoothing the tangents, $T_{S(t,v)}^{(X)}$, along the time direction, we apply a Butterworth low-pass filter to remove the high-frequency signal. For our data set, the filter order is 2 and the -3 dB cutoff frequency is 0.01 Hz. A 0 phase filtering is done by `scipy's` `filtfilt` function. Let the filtered (smoothed) tangents be $\tilde{T}_{S(t,v)}^{(X)}$.

(vii) Convert the eigenvectors back to Cartesian coordinates. For each second of integration time,

$$\|\tilde{\mathcal{E}}_{S(t,v)}^{(X)}(i)\| = \sin(\tan^{-1}(\tilde{T}_{S(t,v)}^{(X)}(i))) \times \prod_{j=1}^i \cos(\tan^{-1}(\tilde{T}_{S(t,v)}^{(X)}(j))), \quad \forall i \in [1, n-1] \quad (\text{A3})$$

and

$$\|\tilde{\mathcal{E}}_{S(t,v)}^{(X)}(n)\| = \prod_{j=1}^n \cos(\tan^{-1}(\tilde{T}_{S(t,v)}^{(X)}(j))). \quad (\text{A4})$$

We then calculate the real and the imaginary parts of the eigenvectors

$$\tilde{\mathcal{E}}_{S(t,v)}^{(X)}(i) = \|\tilde{\mathcal{E}}_{S(t,v)}^{(X)}(i)\| [\cos(\theta_{S(t,v)}^{(X)}(i)) + i \sin(\theta_{S(t,v)}^{(X)}(i))], \quad (\text{A5})$$

where $\theta_{S(t,v)}^{(X)}(i) = \tan^{-1}(\Im(\tilde{\mathcal{E}}_{S(t,v)}^{(X)}(i)) / \Re(\tilde{\mathcal{E}}_{S(t,v)}^{(X)}(i)))$.

(viii) The contribution to the visibility from the Sun is given by

$$\mathbf{V}_{S(t,v)}^{(X)} = \lambda_{S(t,v)}^{(X)} (\tilde{\mathcal{E}}_{S(t,v)}^{(X)} \otimes \tilde{\mathcal{E}}_{S(t,v)}^{(X)}), \quad (\text{A6})$$

where \otimes denotes the outer product between eigenvector $\mathcal{E}_S^{(X)}$ and itself.

(ix) After removing the Sun contribution, the sky contribution to the visibility is

$$\mathbf{V}_{\text{Sky}}^{(X)} = \mathbf{V}^{(X)} - \mathbf{V}_S^{(X)} = \mathcal{E}^{(X)} \Lambda^{(X)} (\mathcal{E}^{(X)})^{-1} - \mathbf{V}_S^{(X)}. \quad (\text{A7})$$

However, the above steps still leave some Sun signal contamination, as shown in Fig. 14. To better remove this leftover Sun contamination, we multiply the Sun signal in equation (A6) by a complex factor of $g = Ae^{i\phi}$:

$$\mathbf{V}_{\text{Sky}}^{(X)} = \mathbf{V}^{(X)} - Ae^{i\phi} \mathbf{V}_S^{(X)}. \quad (\text{A8})$$

(x) To find A and ϕ , we divide the 10 h of daytime data into 36 intervals (1000 s each) and minimize the χ^2 for each interval. Here, we define the χ^2 as

$$\chi^2 = \sum_{t,v} (\Re[\mathbf{V}^{(H)} - Ae^{i\phi} \mathbf{V}_{S(t,v)}^{(H)}])^2 + \sum_{t,v} (\Im[\mathbf{V}^{(H)} - Ae^{i\phi} \mathbf{V}_{S(t,v)}^{(H)}])^2. \quad (\text{A9})$$

The sum is done over all seconds in the chosen interval and frequency 700.625 MHz to 794.375 MHz. We do not sum the frequency channels before 700 MHz and after 800 MHz, because we do not want to include the edge of the band-pass filter. At the end of this process, we have 36 $[A, \phi]$ pairs.

(xi) We have 36 $[A, \phi]$ pairs corresponding to thirty-six 1000-s intervals in 10 h of daytime data. We use a cubic spline to interpolate a $[A, \phi]$ pair for each second in 10 h of daytime data.

(xii) Subtract the corrected Sun signal:

$$\mathbf{V}_{\text{Sky}(t,v)}^{(H)} = \mathbf{V}^{(H)} - A_{\text{int}(t,v)} e^{i\phi_{\text{int}(t,v)}} \mathbf{V}_{S(t,v)}^{(H)}. \quad (\text{A10})$$

This gives us the final solar contamination removed result, and the results are shown in Fig. 17.

APPENDIX B: CALCULATING THE ARTIFICIAL SOURCE VISIBILITIES

The visibilities of the artificial sources come from three made up sources near the NCP. The three simulated sources are randomly chosen to be at (RA, DEC) = (75.75, 81.25), (79.5, 80.5), and (245.0, 79.75) with constant brightness temperatures across all observed frequencies. We calculated the visibility for each frequency in Tianlai's 512 frequency bins (equally spaced between 685 MHz and 810 MHz).

Each astronomical source exhibits a linearly varying phase with time and frequency, since the visibility is proportional to $e^{-i\varphi}$, where φ is the fringe phase and is defined as

$$\varphi = 2\pi \nu \tau_g(\nu, t) = \frac{2\pi \nu \mathbf{b} \cdot \mathbf{s}}{c}. \quad (\text{B1})$$

$\tau_g(\nu, t)$ is the frequency-independent geometric delay and is equal to

$$\tau_g(\nu, t) \equiv \frac{\mathbf{b} \cdot \mathbf{s}}{c} = \frac{b_x}{c} \cos \delta \cos H(t) - \frac{b_y}{c} \cos \delta \sin H(t) + \frac{b_z}{c} \sin \delta, \quad (\text{B2})$$

where δ is the source declination, $H(t)$ is the source hour angle as a function of sidereal time, $\mathbf{b} = (b_x, b_y, b_z)$ are the baseline

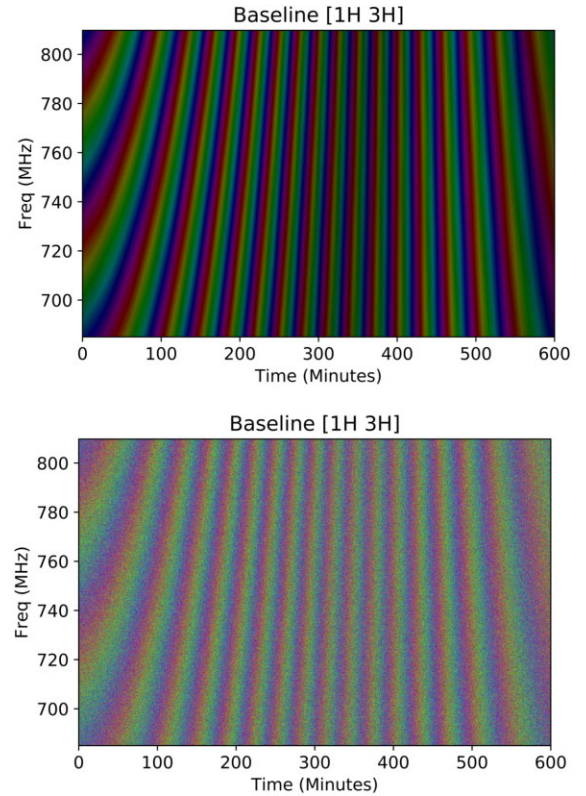


Figure B1. Top: Combined visibilities of three artificial sources for baseline [1H 3H], shown here for 10 h. Bottom: Mock observed simulated visibility with those three artificial sources.

components with units of length in the radial, eastern, and northern polar directions, and \mathbf{s} is the source vector. c is the speed of light. We can also calculate the fringe rates as follows:

$$\frac{\partial \varphi}{\partial t} = \frac{2\pi v}{c} [-b_x \sin H(t) + b_y \cos H(t)] \cos \delta. \quad (\text{B3})$$

For each second of integration time and each frequency, the visibility for dish i and j is calculated as follows

$$\mathbf{V}_{(i,j)} = \sum_{k=1}^n \mathcal{F}_k A(\mathbf{s}) e^{-i\varphi}, \quad (\text{B4})$$

where \mathcal{F}_k is the flux of source k . In our simulations, we used three

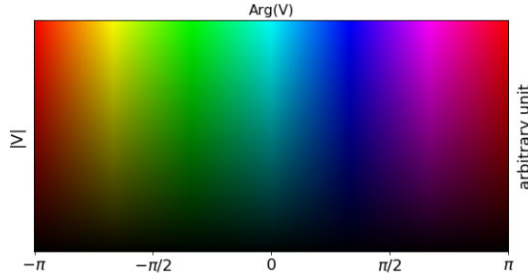


Figure B2. The colour palette used to represent complex visibilities in this paper is shown in this plot. The phase of the complex visibility is represented by the hue and the magnitude is represented by the brightness. For more details, see appendix A of (Wu et al. 2021).

sources with $\mathcal{F} = (546, 170, 128)$ Jansky. $A(\mathbf{s})$ is the gain of the antenna in the direction of the source vector \mathbf{s} . For simplicity, $A(\mathbf{s})$, the simulated main beam gain, is taken as a Gaussian distribution with a standard deviation of 3° ($\text{FWHM} = 7^\circ$), and we assume that all three artificial sources fall within the main beam. Therefore, we did not model the beam sidelobe gain. $A(\mathbf{s})$ is also assumed to be independent of frequency. In the real Tianlai Dish beam pattern, the mainbeam (excluding the sidelobe) FWHM is about 5° (Wu et al. 2021). The simulated baselines are identical to the real Tianlai dish array, and the procedure for calculating the visibility is repeated for every baseline. The waterfall plots of the simulated visibilities for a few typical baselines are shown in Fig. B1 using the same representation of waterfall plots that is used throughout this paper and is described in Fig. B2. As expected, longer baselines give higher fringe rates, and for a given baseline, we see a faster fringe rate at lower frequency.

This paper has been typeset from a $\text{T}_{\text{E}}\text{X}/\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ file prepared by the author.