

# Modeling the Extragalactic Epoch of Reionization Foreground

Patricia A. Carroll

A dissertation  
submitted in partial fulfillment of the  
requirements for the degree of

Doctor of Philosophy

University of Washington

2016

Reading Committee:

Miguel Morales, Chair

Mario Juric

Matthew McQuinn

Program Authorized to Offer Degree:  
Astronomy

©Copyright 2016

Patricia A. Carroll

University of Washington

**Abstract**

Modeling the Extragalactic  
Epoch of Reionization Foreground

Patricia A. Carroll

Chair of the Supervisory Committee:  
Professor Miguel Morales  
Physics

The Epoch of Reionization represents a largely unexplored yet fundamental chapter of the early universe. During this period, spanning several hundred million years, the first stars and galaxies formed and the Hydrogen-dominated intergalactic medium transitioned from a predominantly neutral to ionized state. Modern efforts to study exactly when and how reionization occurred are largely focused on the distribution of neutral Hydrogen gas and its evolution in response to the increasing abundance of luminous objects and ionizing flux. The Murchison Widefield Array is a low frequency radio interferometer designed as a first generation EoR experiment. The predominant systematic difficulty in making a detection of the primordial HI signal is the overwhelmingly bright emission from the intervening foreground galaxies and quasars. This thesis presents novel survey methods used to create a highly precise and reliable catalog of discrete extragalactic sources for the purposes of both calibration and foreground removal.

# TABLE OF CONTENTS

	Page
List of Figures . . . . .	iii
List of Tables . . . . .	vi
Chapter 1: Introduction . . . . .	1
1.1 The Epoch of Reionization . . . . .	1
1.2 Mapping the EoR . . . . .	3
1.3 The EoR Foregrounds: Radio Source Populations . . . . .	4
1.4 Radio Interferometry . . . . .	11
1.5 Low-Frequency Southern Sky Surveys . . . . .	21
1.6 This Thesis . . . . .	23
Chapter 2: Observations & Data Reduction . . . . .	29
2.1 Pre-processing . . . . .	30
2.2 Fast Holographic Deconvolution . . . . .	31
Chapter 3: Source Finding & Characterization . . . . .	36
3.1 Component Clustering . . . . .	37
3.2 Cross-Snapshot Association . . . . .	38
Chapter 4: Reliability Classification . . . . .	43
4.1 Machine Learning . . . . .	44
4.2 KATALOGSS Sample Selection . . . . .	53
Chapter 5: Radio Survey Cross-Matching . . . . .	57
5.1 Positional Update & Matching Algorithm . . . . .	57
5.2 Flagging & Visual Inspection . . . . .	59
5.3 Results . . . . .	63

Chapter 6:	The EoR0 Foreground Catalog . . . . .	67
6.1	Flux Scale . . . . .	67
6.2	Eddington Bias . . . . .	68
6.3	Completeness . . . . .	71
6.4	Spectral Index Distribution . . . . .	73
6.5	Astrometry . . . . .	73
6.6	The Catalog . . . . .	82
6.7	Caveats . . . . .	94
6.8	Summary . . . . .	94
Chapter 7:	Unidentified & Ultra Steep Spectrum Radio Sources . . . . .	96
7.1	New Radio Source Detections . . . . .	96
7.2	Galaxy Clusters . . . . .	124
7.3	HzRG Candidates . . . . .	127
Chapter 8:	EoR Foreground Modeling and Removal . . . . .	135
8.1	Sky Model Tests . . . . .	135
8.2	A Multi-Survey Master Catalog . . . . .	142
8.3	Extended Source Models . . . . .	144
8.4	Summary . . . . .	149
Chapter 9:	Conclusions . . . . .	155
9.1	EoR Foreground Modeling . . . . .	155
9.2	Future Directions . . . . .	157

## LIST OF FIGURES

Figure Number	Page
1.1 A timeline of the early universe . . . . .	2
1.2 The EoR Power Spectrum . . . . .	5
1.3 The chromatic foreground wedge . . . . .	6
1.4 AGN observational types . . . . .	7
1.5 Example radio spectral energy distributions . . . . .	8
1.6 A radio galaxy: VLA and HST image of Hercules A . . . . .	10
1.7 A starburst galaxy: VLA image of the Sculptor galaxy . . . . .	12
1.8 Example of PSF side lobes . . . . .	14
1.9 An MWA tile . . . . .	15
1.10 MWA tile layout . . . . .	16
1.11 MWA snapshot uv coverage . . . . .	17
1.12 MWA primary beam response . . . . .	19
1.13 Telescope pointing diagram . . . . .	20
1.14 Coverage of southern radio surveys . . . . .	24
2.1 Beam response across the zenith pointing. . . . .	30
2.2 FHD and EoR pipeline flow chart . . . . .	33
2.3 EoR0 field image . . . . .	35
3.1 DBSCAN clustering diagram . . . . .	39
3.2 Source candidate detection counts . . . . .	41
3.3 SNR distributions . . . . .	42
4.1 Side lobe excess near the brightest sources . . . . .	45
4.2 Feature distributions . . . . .	48
4.3 Feature variance of PCA components . . . . .	49
4.4 Principle component density distributions . . . . .	50
4.5 Bayes Information Criterion for model selection . . . . .	51

4.6	GMM classification results . . . . .	52
4.7	Adaboost classification results . . . . .	53
4.8	Side lobe excess near the brightest sources . . . . .	54
4.9	Feature distributions with reliability classification . . . . .	56
5.1	Example of PUMA match process and result. . . . .	60
5.2	Three Sculptor group galaxies. . . . .	65
6.1	PUMA extrapolated 182 MHz flux ratios. . . . .	69
6.2	Eddington bias flux density correction. . . . .	70
6.3	Eddington bias in the spectral index distribution. . . . .	72
6.4	Differential source counts compared to NVSS. . . . .	74
6.5	Differential source counts compared to VLSSr. . . . .	75
6.6	Spectral index distribution from PUMA SED fits. . . . .	76
6.7	Two-point spectral index distributions of PUMA matches. . . . .	77
6.8	KGS position bias and correction. . . . .	79
6.9	Modeled position bias and correction. . . . .	81
6.10	Examples of PUMA match results. . . . .	85
7.1	Member counts of galaxy clusters likely hosting a KGS source. . . . .	125
7.2	Spectral index distributions of sources near and far from galaxy clusters. . .	126
7.3	The affect of galaxy cluster richness on the spectral properties of the radio sources they host. . . . .	128
8.1	MWACS–KGS PS difference . . . . .	138
8.2	MWACS–KGS normalized PS difference . . . . .	139
8.3	MWACS vs. KGS flux scale . . . . .	140
8.4	Schematic of foreground contaminated modes by relative beam position . . .	141
8.5	Power spectrum difference including widefield foregrounds. . . . .	143
8.6	Multi-survey master catalog of foregrounds . . . . .	145
8.7	MWACS–Master PS difference . . . . .	146
8.8	MWACS–Master normalized PS difference . . . . .	147
8.9	Number and flux density maps of NGC 253 components . . . . .	149
8.10	Rendered image of NGC 253 components compared to NVSS. . . . .	150
8.11	Extended source model for NGC 253 . . . . .	151
8.12	NGC 253 point and extended model residual images. . . . .	151

8.13 NGC 253 point and extended model PS difference . . . . . 152  
8.14 NGC 253 point and extended model normalized PS difference. . . . . 153

## LIST OF TABLES

Table Number	Page
1.1 Radio spectral index by source types . . . . .	9
1.2 Overview of existing large sky surveys at low radio frequencies . . . . .	23
4.1 Principle components of features for reliability classification . . . . .	50
4.2 Properties of reliability classes . . . . .	55
5.1 Final composition of KGS PUMA match types. . . . .	64
6.1 Modeled position bias polynomial coefficients. . . . .	80
6.2 KGS EoR0 catalog sample. . . . .	91
7.1 New radio sources. . . . .	97
7.2 HzRG candidates . . . . .	130

## ACKNOWLEDGMENTS

This scientific work makes use of the Murchison Radio-astronomy Observatory, operated by CSIRO. We acknowledge the Wajarri Yamatji people as the traditional owners of the Observatory site. Support for the operation of the MWA is provided by the Australian Government Department of Industry and Science and Department of Education (National Collaborative Research Infrastructure Strategy: NCRIS), under a contract to Curtin University administered by Astronomy Australia Limited. We acknowledge the iVEC Petabyte Data Store and the Initiative in Innovative Computing and the CUDA Center for Excellence sponsored by NVIDIA at Harvard University.

This research has made use of the NASA/IPAC Extragalactic Database (NED), which is operated by the Jet Propulsion Laboratory, California Institute of Technology, under contract with the National Aeronautics and Space Administration. We also acknowledge the use of NASA's SkyView facility (<http://skyview.gsfc.nasa.gov>) located at NASA Goddard Space Flight Center.

This thesis was supported by National Science Foundation grants AST-0847753, AST-1410484, and AST-1506024 as well as the American Australian Association Sir Keith Murdoch fellowship and the University of Washington Graduate Opportunities and Minority Achievement Program dissertation fellowship.

## DEDICATION

To Benjamin Orion, my partner in life and adventure,  
for showing me the wonders of the world under my feet.

## Chapter 1

# INTRODUCTION

### ***1.1 The Epoch of Reionization***

Our universe came into existence with the Big Bang and, in a fraction of a second, rapidly inflated and formed a hot ionized plasma composed of protons, electrons, and other subatomic particles. This plasma expanded and cooled for several hundred thousand years until eventually it became energetically favorable for the oppositely charged protons and electrons to combine. Very quickly, the fire ball of hot plasma transitioned to a neutral atomic Hydrogen gas (HI) and the universe entered the Cosmic Dark Ages.

The Cosmic Dark Ages lasted for hundreds of millions of years during which time gravity worked to form stars, galaxies, and quasars. The sudden influx of energetic photons from these first luminous objects marks the Cosmic Dawn and the beginning of a second major phase transition. Ultraviolet photons were absorbed by Hydrogen atoms, transferring enough energy to break the electron–proton bond. In this way, the Hydrogen began to re-ionize in what is known as the Epoch of Reionization (EoR).

The Epoch of Reionization represents a largely unexplored yet foundational chapter of the early universe. When did it occur? How long did it last? What was the dominant source of ionizing flux? How did ionized regions grow? These seemingly basic questions are the first line of inquiry into a deeper understanding of the fundamental physics governing our universe. Their answers encompass an intricate interplay of radiative processes, cosmological expansion, and galaxy and stellar formation and evolution. Over the past several decades, rich theory has been developed to describe the EoR and guide the design of observational tests. This thesis is a part of ongoing work to analyze the first rounds data from a first generation EoR instrument.

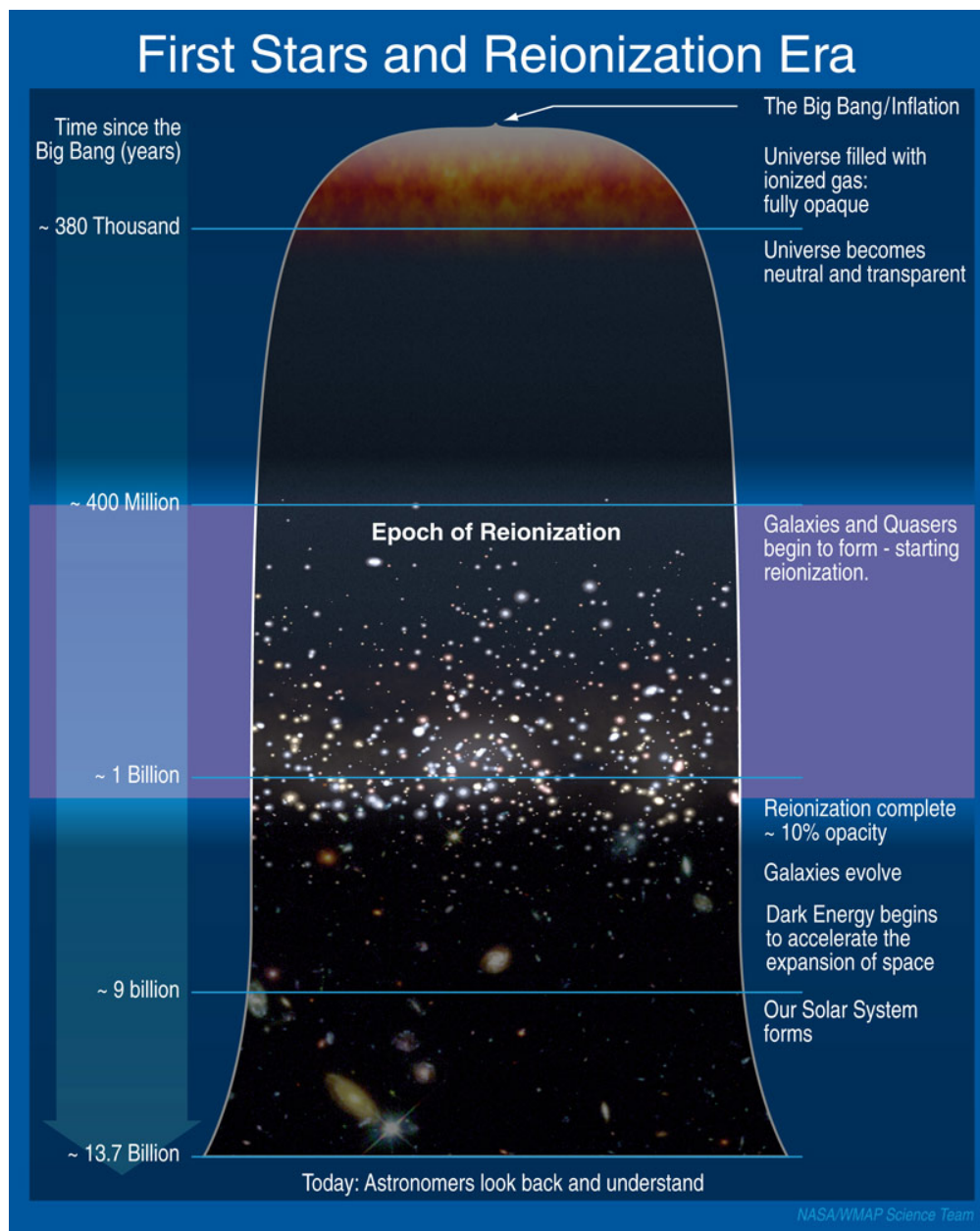


Figure 1.1: About 400 million years after the Big Bang, the first stars and galaxies began to form and re-ionize the Hydrogen gas. This Epoch of Reionization is a fundamental yet little understood chapter of the early universe (credit: NASA/WMAP Science Team).

## 1.2 Mapping the EoR

Modern efforts to detect the EoR are aimed at measuring a faint signal from the neutral Hydrogen gas. The electron in a neutral Hydrogen can occupy one of two spin states at slightly different energy levels. An electron in the higher energy state will eventually decay to the lower energy state “flipping” its spin and emitting a weak photon with a wavelength  $\lambda = 21$  cm (or frequency  $\nu = 1420.4$  MHz). This radiative half-life of this transition is 11 million years, but the sheer abundance of Hydrogen in the universe promises a detectable signal given a sensitive enough instrument. This promise drove the design of the Murchison Widefield Array (MWA). The MWA is predicted to require  $\sim 1000$  hours of carefully processed high quality data to make a statistically significant detection [3].

Due to cosmological redshift of the HI line and the redshift-distance relation, the frequency dimension becomes a spatial dimension. This is the basic concept underlying 21 cm tomography, a form of 3D imaging. By observing some sky area over a frequency range corresponding to the redshifted EoR signal, we probe the 3D spatial distribution of HI.

Early on, at high redshift or lower frequencies, we should detect more power from HI on large scales. Later, as the EoR progresses and ionized hydrogen (HII) regions grow, the neutral abundance decreases and the size scale of HI structure shrinks. This leads us to the EoR power spectrum, an analytic tool we can use to characterize the abundance and structure of HI in time and space.

Before exploring the EoR power spectrum, we must first consider a complicating matter. HI is far from the dominant source of emission within the relevant frequency range,  $\sim 80 - 300$  MHz. Intervening galaxies and quasars overwhelm the HI signal by 4-5 orders of magnitude. We must rely on the different spectral behavior of the foregrounds to separate them from the desired cosmological signal. Foregrounds, for the most part, have a smooth radio continuum, meaning there are no sharp features like emission lines. In contrast, the spatial structure of neutral Hydrogen along the line-of-sight is lumpy and redshift translates this to a lumpy spectrum.

These differing spectral structures ought to separate astrophysical foregrounds from the EoR signal in Fourier space. Smooth spectrum foregrounds will occupy low Fourier modes along the line-of-sight ( $k_{\parallel}$  modes), while the lumpy HI structure will occupy high modes. The other two dimensions are the sky dimensions, Right Ascension (RA) and Declination (Dec). Since cosmological structure has no preferred direction, the 2D sky can be collapsed into a single dimension perpendicular to the line-of-sight.

The resulting 2D Fourier transform of the reduced cosmological volume (parallel and perpendicular to the line-of-sight) gives us the EoR 2D ( $k_{\parallel}$ ,  $k_{\perp}$ ) power spectrum illustrated in Figure 1.2. Due to the spectral smoothness of the radio continuum, foreground power is largely constrained toward the bottom, in the lowest  $k_{\parallel}$  modes, and these can be avoided.

Unfortunately reality isn't quite so ideal. The instrument itself introduces frequency dependent or "chromatic" structure to the otherwise smooth foregrounds. This causes foreground power to bleed into higher  $k_{\parallel}$  modes in a wedge shape pattern, contaminating the EoR signal [12, 25, 43, 29, 41, 17, 39, 32, 13, etc.]. The chromatic foreground wedge is illustrated in Figure 1.3.

A small "window" remains uncontaminated at low  $k_{\perp}$  but given the weakness of the signal, we aim to directly remove foreground power in order to mitigate the amount of contamination and widen the EoR window. Foreground removal requires a high fidelity model of the galaxies and quasars as seen by the MWA. The creation of such a model is the central focus of this thesis.

### ***1.3 The EoR Foregrounds: Radio Source Populations***

It is important to understand, physically, what these foregrounds are. The extragalactic radio continuum is dominated by emission from galaxies hosting active galactic nuclei (AGN), powered by accretion onto supermassive black holes at their core. Several observational subclasses of AGN exist including radio galaxies, blazars, and quasars each with their own various sub-types. Significant evidence suggests that the observational diversity of AGN is in large part due to orientation effects as illustrated in Figure 1.4. AGN Unification

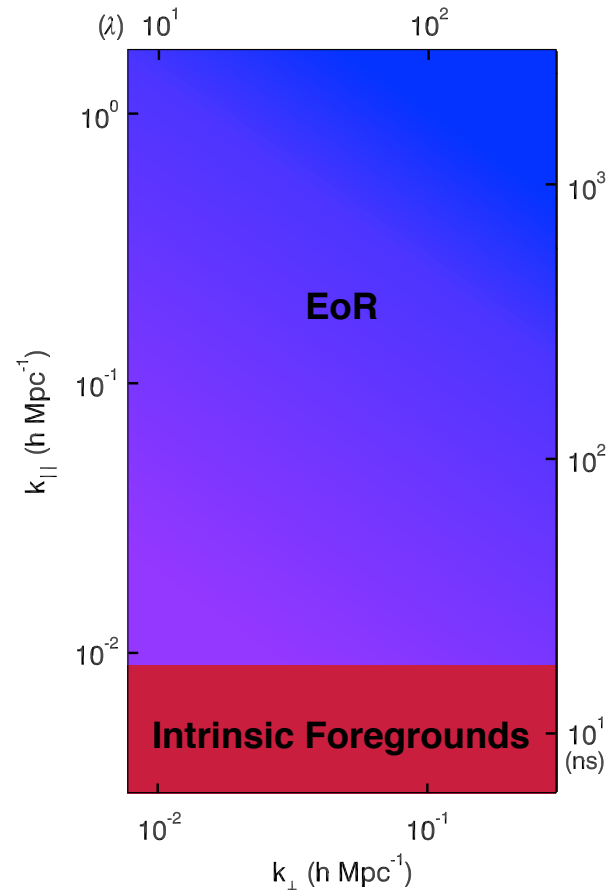


Figure 1.2: The 2D  $k$  power spectrum theoretical distribution of power. Smooth spectrum foregrounds are restricted to the lowest  $k_{\parallel}$  modes. The non-smooth line-of-sight redshifted EoR HI signal occupies high  $k_{\parallel}$  modes.

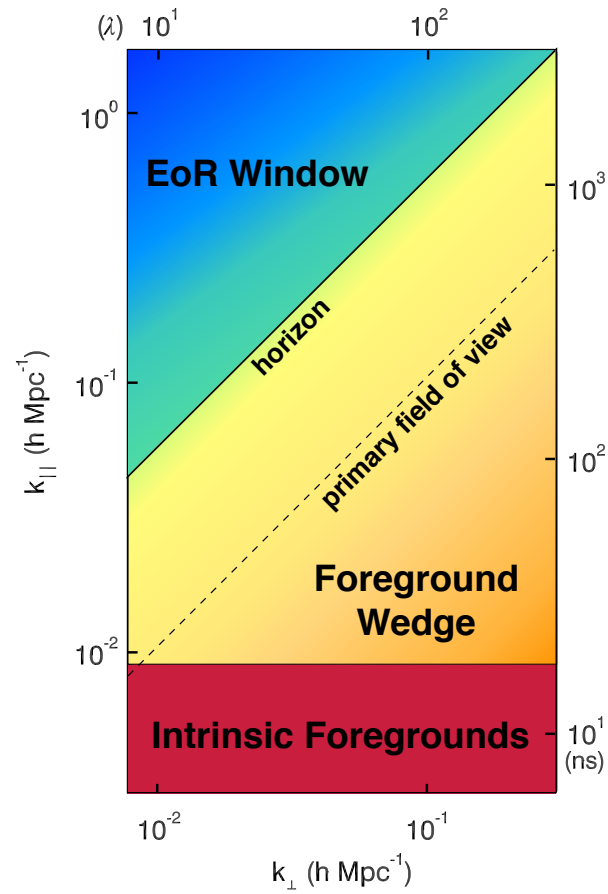


Figure 1.3: Chromatic, frequency dependent, instrument effects introduce structure to the otherwise smooth-spectrum foregrounds, causing power to bleed into higher  $k_{\parallel}$  modes. This shape in the 2D EoR k-power spectrum is referred to as the foreground “wedge”.

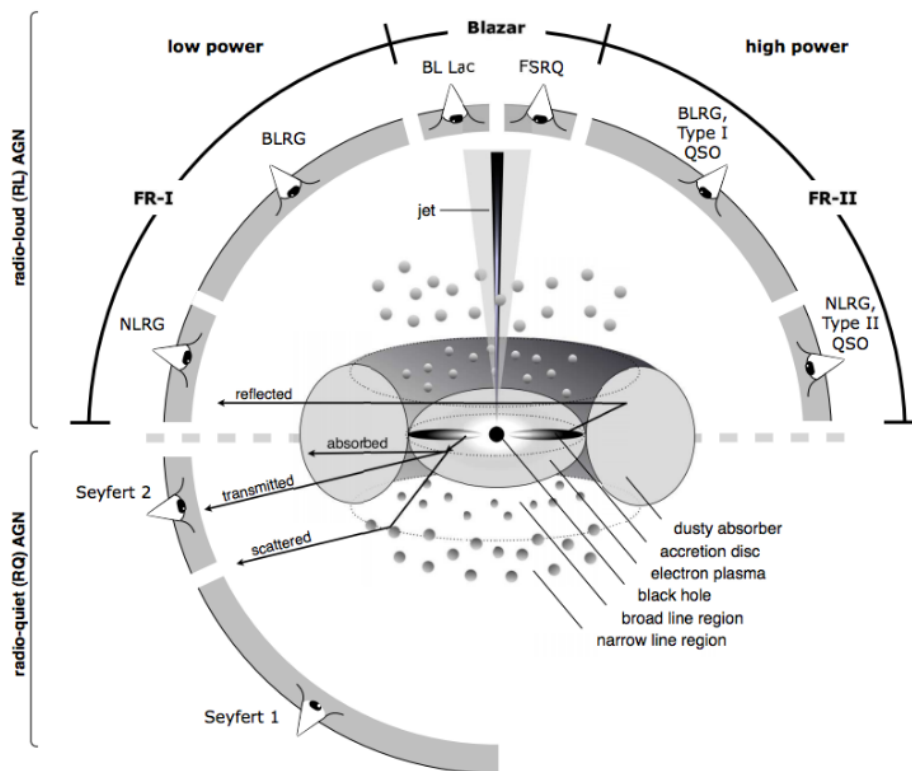


Figure 1.4: Schematic representation of orientation effects on AGN observational classification. Graphic from (Beckmann & Shrader 2013), courtesy of Marie-Luise Menzel (MPE)

models attempt to disentangle this variation from intrinsic physical differences that could be associated with evolutionary state or environment.

These different subclasses have different morphological and spectral characteristics. Our average radio source is a “point” source with a power-law relation between flux density and frequency. A point source is not resolved by the observing instrument, meaning it’s angular scale on the sky is smaller than the instruments ability to focus. A point source is the simplest to model as it is simply a delta function convolved with the instruments point spread function (PSF) often approximated by a 2D Gaussian. The power-law spectrum arises from Synchrotron emission, characterized by the relation  $S \propto \nu^\alpha$ ; where  $S$  is the flux density,  $\nu$  the frequency, and  $\alpha$  the spectral index. Most radio sources exhibit a power-law

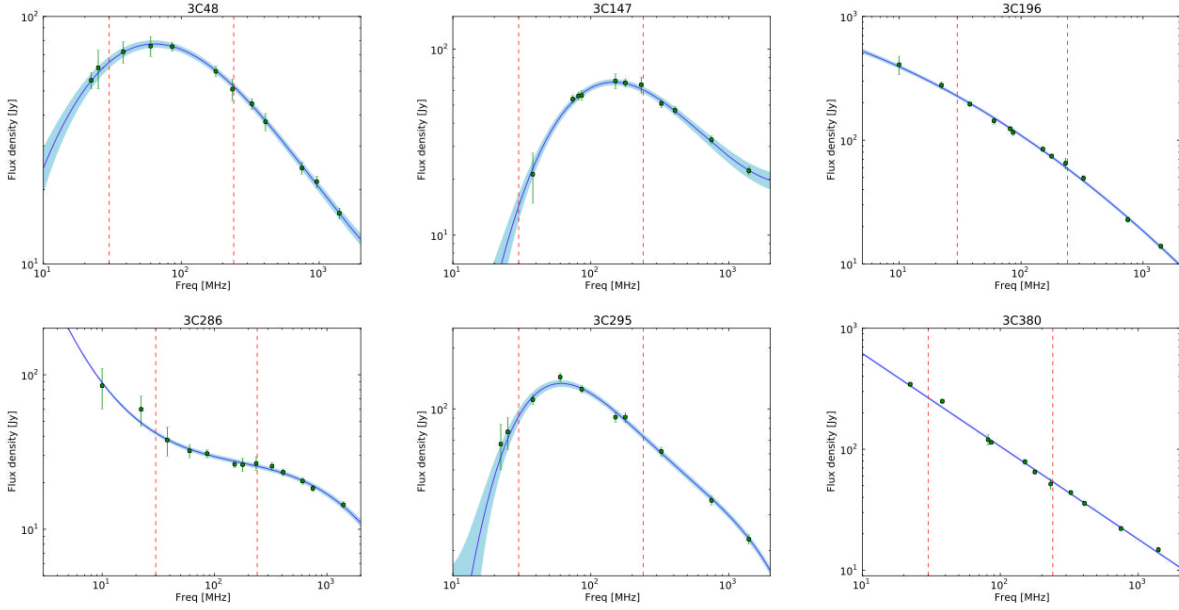


Figure 1.5: Example radio spectra of six bright radio sources from the 3C catalog (figure from [34]). Many radio sources exhibit some degree of curvature over large frequencies ranges. 3C380 illustrates a well fit power-law spectrum across the full frequency range with spectral index  $\alpha = -0.77$ . Red dashed lines mark the Low Frequency Array (LOFAR) coverage between 30-240 MHz. The MWA covers 80-300 MHz.

synchrotron spectrum (also referred to as the spectral energy distribution or SED) with spectral indexes in a tight distribution about a typical median of -0.8.

A small but significant number of sources deviate from this ideal picture in known ways. The fraction of sources with spectral index  $< -1$  or  $> -0.5$  are exceptional but not uncommon in large surveys. Many exhibit curvature – steepening, or flattening of the spectral index – over large frequency ranges. A variety of radio spectral shapes is shown in Figure 1.5. Flat, steep, ultra-steep, and peaked spectra are commonly used to select rare-type AGN sub-populations (see Table 1.1).

Gigahertz Peaked Spectrum (GPS) and Compact Steep Spectrum (CSS) sources are

Source Type	Spectral Index
Gigahertz Peaked Spectrum (GPS)	$\gtrsim 0.5$
Compact Steep Spectrum (CSS)	$\lesssim -0.5$
Flat Spectrum Radio Quasar (FSRQ)	$-0.5 < \alpha < 0.5$
Ultra-Steep Spectrum (USS)	$\lesssim 1.4$

Table 1.1: Sub-populations of radio AGN and approximate spectral index ranges.

young, compact, and powerful radio sources with spectra that peak in the Gigahertz or Megahertz range. The peak is thought to be caused by synchrotron self-absorption, when the density of the emission region is high. Flat Spectrum Radio Quasars (FSRQs) are core dominated, synchrotron self-absorbed, and are commonly associated with gamma ray emission (e.g. BL Lac objects). The sparse population of ultra-steep spectrum (USS) sources is of particular interest because they have been shown to trace high redshift ( $z > 2$ ) galaxies, and are thus a fundamental part of understanding galaxy formation and evolution. Lastly, radio AGN are known to vary considerably in brightness over long periods of time. This results in an erratic broad-band SED when considering data taken years or decades apart, and potentially a loss of fidelity in any foreground model over time.

In terms of morphology, radio galaxies and star forming galaxies often cannot be approximated as point sources. Radio galaxies are characterized by jets and lobes on large scales compared to their optical extent. Figure 1.6 shows radio observations of a classic nearby radio galaxy, Hercules A. At a large enough distance, radio lobes can usually be approximated as point sources, but at the resolution limit they begin to blend into a single elongated source that requires more careful modeling.

Star forming galaxies on the other hand, are not AGN. These sources have relatively low surface brightness with radio emission coming from star forming regions in a non-smooth distribution. Figure 1.7 shows one bright nearby starburst galaxy, the Sculptor Galaxy (NGC 253). Radio galaxies and star-forming galaxies are often partially resolved and require

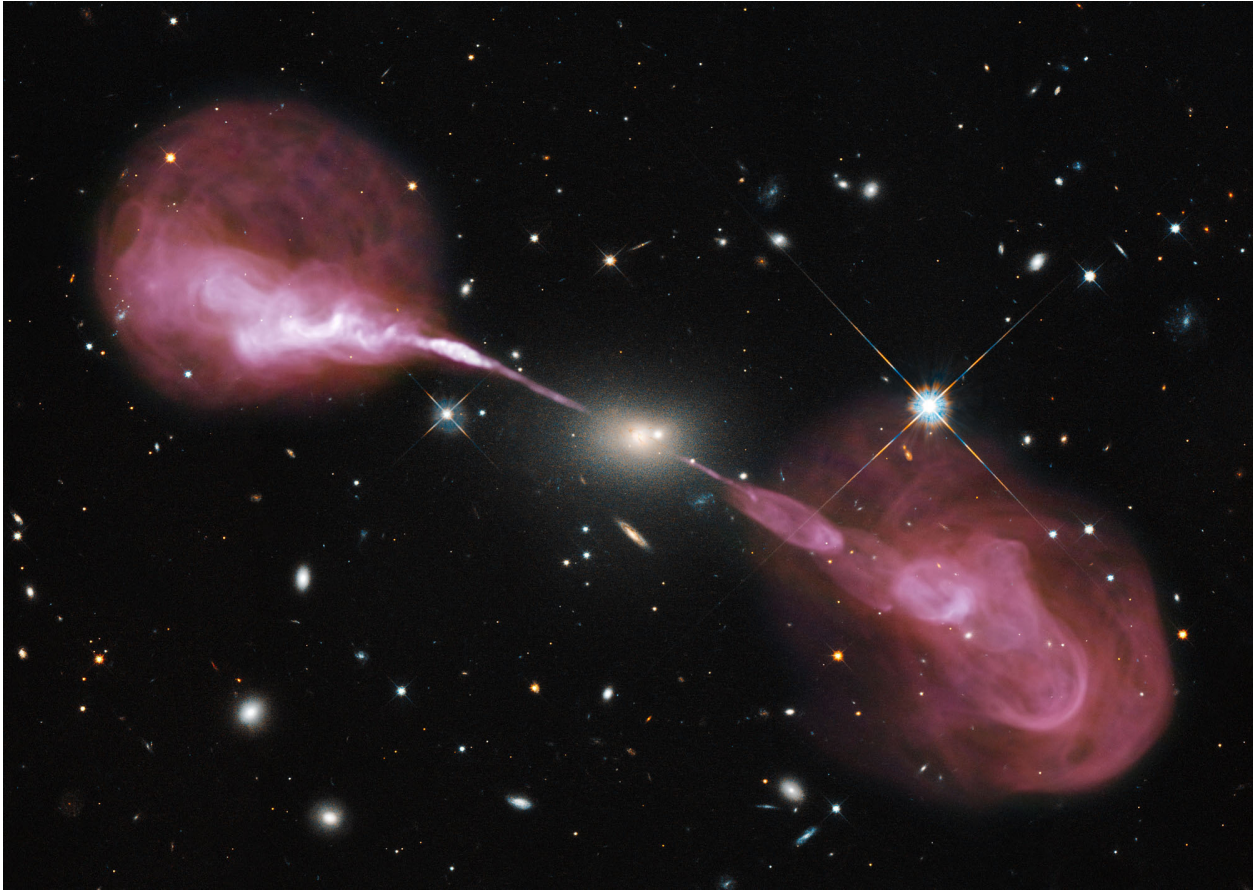


Figure 1.6: This classic image of Hercules A shows the giant elliptical galaxy at optical wavelengths as seen by the Hubble Space Telescope with radio imaging of the jets and lobes from the VLA overlaid in magenta. Credit: NASA, ESA, S. Baum and C. ODea (RIT), R. Perley and W. Cotton (NRAO/AUI/NSF), and the Hubble Heritage Team (STScI/AURA).

more complex models to accurately describe their spatial extent.

#### **1.4 Radio Interferometry**

Radio telescopes come in many shapes and sizes. The angular scale that any telescope can resolve is directly proportional to wavelength and inversely proportional to aperture diameter. In the radio regime, particularly at low frequencies, it is impractical to build a single dish capable of resolving discrete sources on small angular scales. Instead, interferometric antenna arrays are used.

The maximum distance between any two antenna in an array is the effective diameter and determines the resolution of the instrument. Each possible pair of antennas represents a single baseline. A baseline gives a single point in the  $u$ - $v$  plane, where  $u$  is the east-west separation and  $v$  is the north-south separation (directions are as aligned with the celestial coordinates Right Ascension (RA) and Declination (Dec)). The combinations of all baselines determines the  $uv$ -coverage of the array. Each baseline is sensitive to a specific scale and angle on the sky. Long baselines are sensitive to small scale structure and short baselines are sensitive to large scale structure. So, ideally, we want a full and smooth  $uv$  distribution to cover all scales and angles on the sky. The combined baseline response gives the synthesized array beam response, which describes the overall sensitivity of the instrument.

The array and primary beam are referred to throughout this thesis and can be confusing terms. The primary beam refers to the area of sky to which the instrument is most sensitive and is therefore also called its primary field of view. This maximum scale is determined by the size of the antenna. The center of the beam, where the telescope is pointed, has the peak response – that is, the maximum ratio of antenna temperature to sky brightness temperature – and this falls off radially outward. The beam shape is approximately Gaussian but reaches a null before increasing again in the side lobes. The peak side lobe response is a small fraction of the primary beam response and is usually ignored. In fact, the primary beam is usually ignored below the 50% response level, or half-beam although that is not the case in this work.

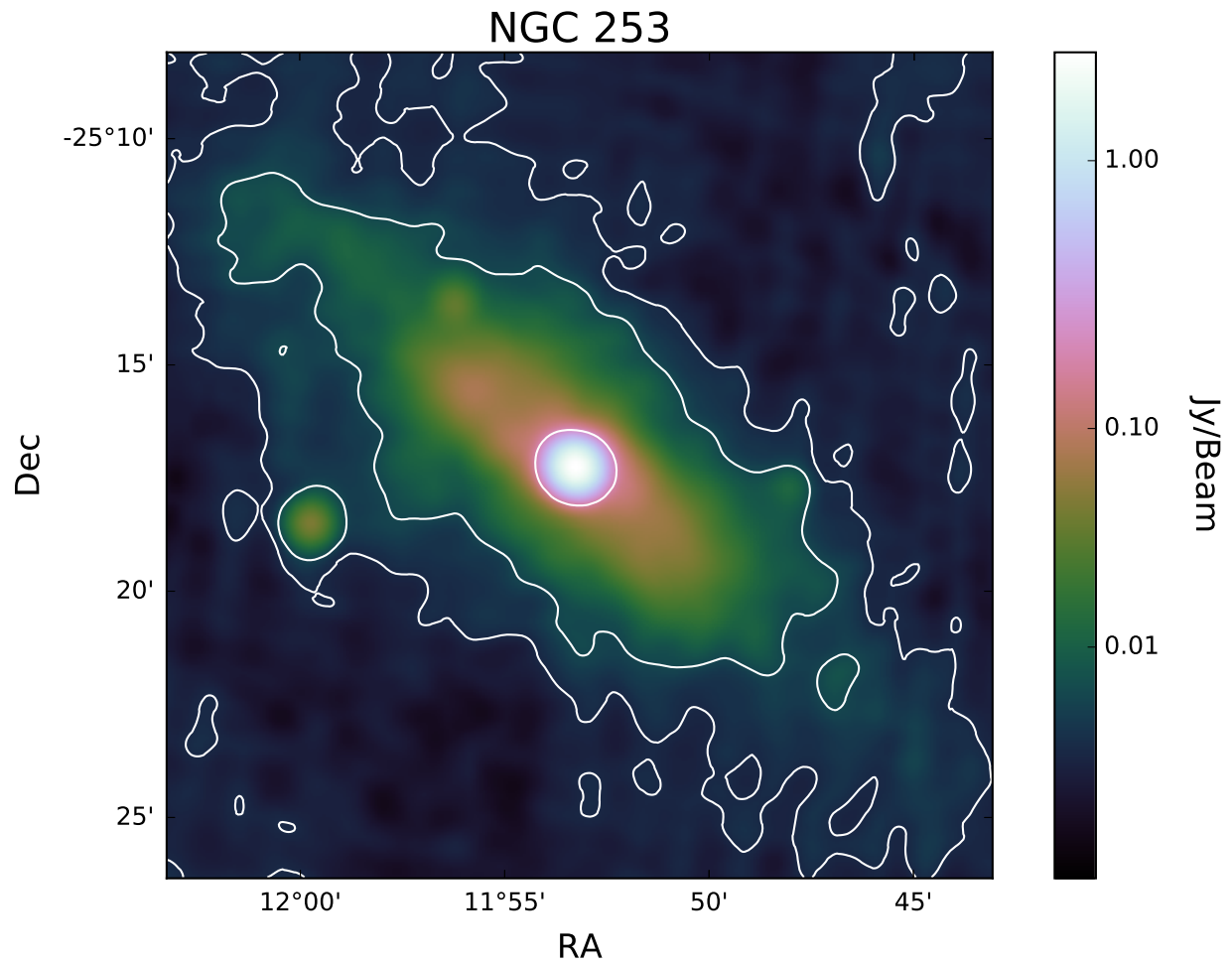


Figure 1.7: Nearby star-forming galaxies are morphologically complex to model, as shown in this 1.4 GHz image of the local starburst galaxy NGC 253

The array beam is the point spread function (PSF), or the shape a point source will take. The size of the PSF is determined by the resolution limit according to the longest baseline. For the brightest sources, the side lobes of the PSF can result in small peaks above the noise level that mimic real sources. An example of PSF side lobes in an MWA image is shown in Figure 1.8. These can be challenging to identify and remove from large radio surveys and represent a significant source of contamination. Properly modeling the array beam is important both to ensure correct scaling to recover the true sky brightness on large scales, and to fit and extract sources from the images on small scales.

Throughout this thesis, “beam” or “beam power” refers to the relative sensitivity within the primary field-of-view and side lobes, and can take any value between 0 and 1. The exception to this rule is any mention of “beam-width”. This refers to the FWHM of the PSF. Likewise, “side lobe” will most often refer to the side lobes of the PSF. Side lobes of the primary beam will be specified as such.

#### *1.4.1 The Murchison Widefield Array*

The MWA is a novel low-frequency radio interferometer located in a radio quiet region of Western Australia. Its design is tuned to the EoR but well-suited for survey science [7]. It boasts a large field of view ( $\sim 30^\circ$ ), wide frequency range (80-300 MHz), high resolution in frequency and time (40kHz and 0.5 sec), good instantaneous uv-coverage, and high surface brightness sensitivity.

The array has 2048 individual antennas arranged in 128 4x4 tiles. Figure 1.9 shows a single MWA tile. The tiles are located within a 1.5 km radius as shown in Figure 1.10. The layout of tiles were chosen to optimize uv-coverage with respect to physical building constraints [4, 3, 40], and is concentrated near the core to enhance sensitivity to the large angular scales of HI regions at high redshift. The instantaneous uv-coverage is shown in Figure 1.4.1 for the full array and the central 112 tiles used for the EoR studies, and the primary beam response at zenith is shown in Figure 1.4.1.

The MWA tiles are stationary and there is no mechanical movement. Instead, to steer

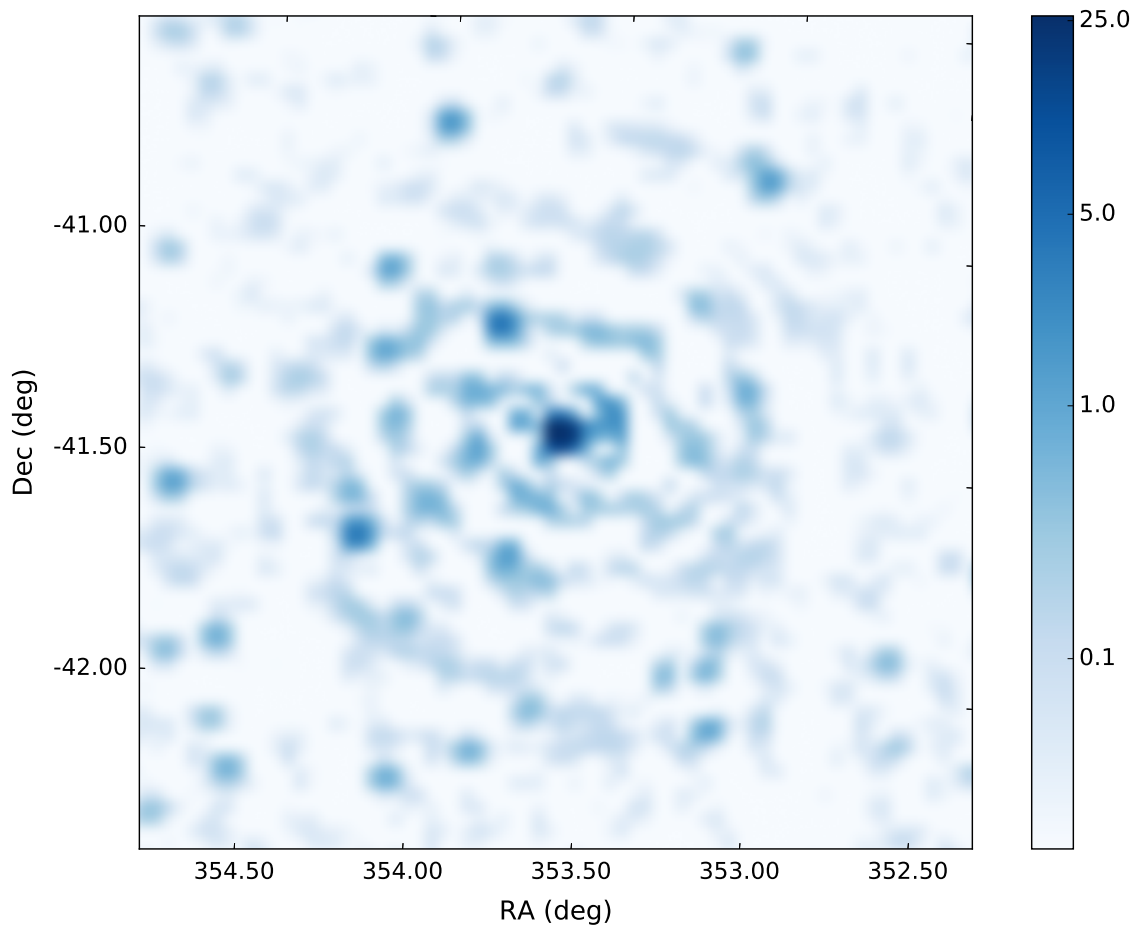


Figure 1.8: This is an MWA image of a 26 Jy source. The image is a weighted mean of 71 snapshot observations. The PSF rotates between observations resulting in arcs of side lobe sources that are confused with real sources.

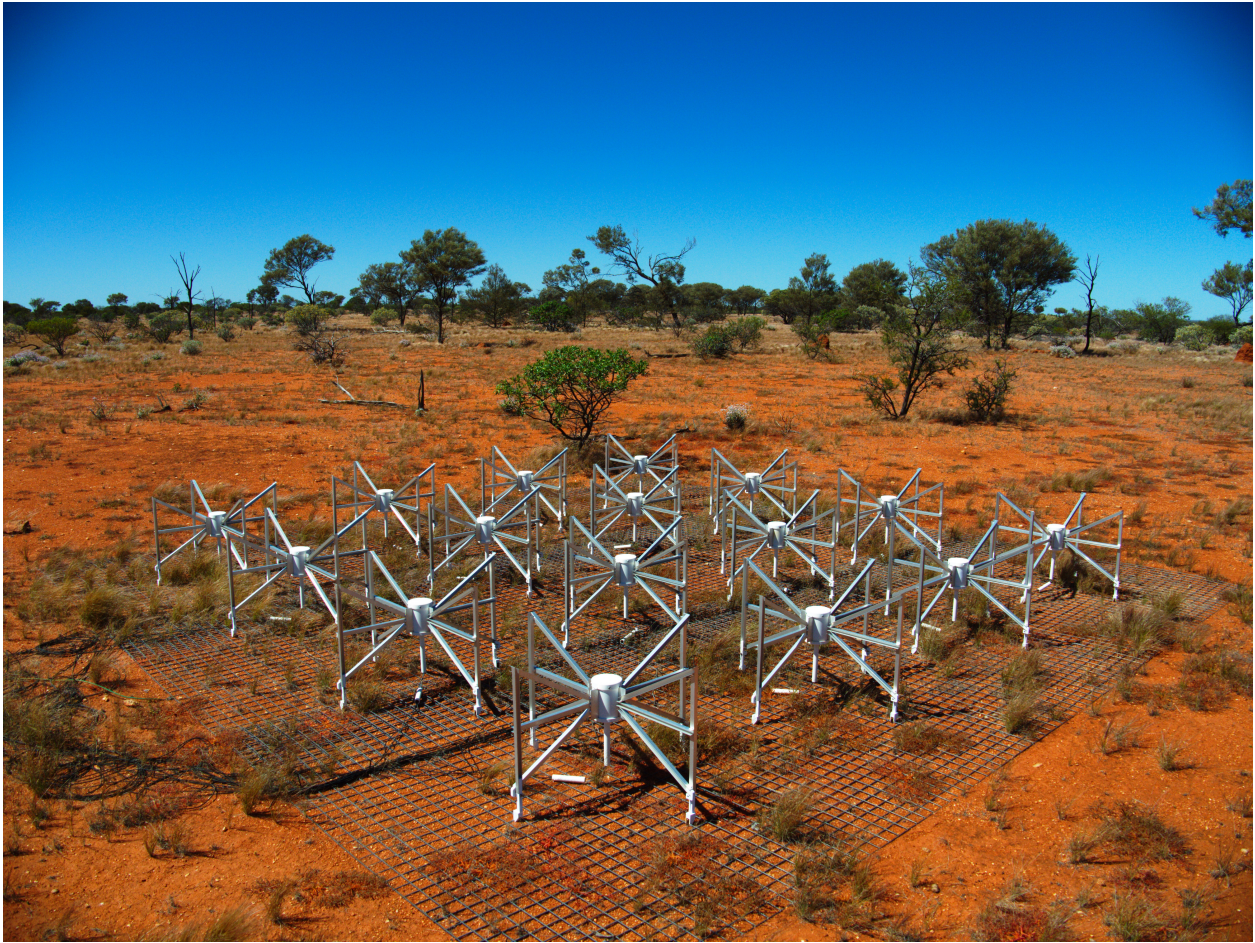


Figure 1.9: A single MWA tile element consists of  $4 \times 4$  dipole antennas on a  $5 \times 5$  m ground screen.

the telescope a time delay is introduced electronically that corresponds to the differential in arrival times of plane waves from an off-axis direction. Electronically steered arrays benefit from reduced system noise, a wider field of view, and the ability to observe in many directions simultaneously. This concept is illustrated in Figure 1.4.1.

Due to its large field of view and good uv-coverage, the MWA is a great survey instrument capable of observing large areas of sky very quickly.

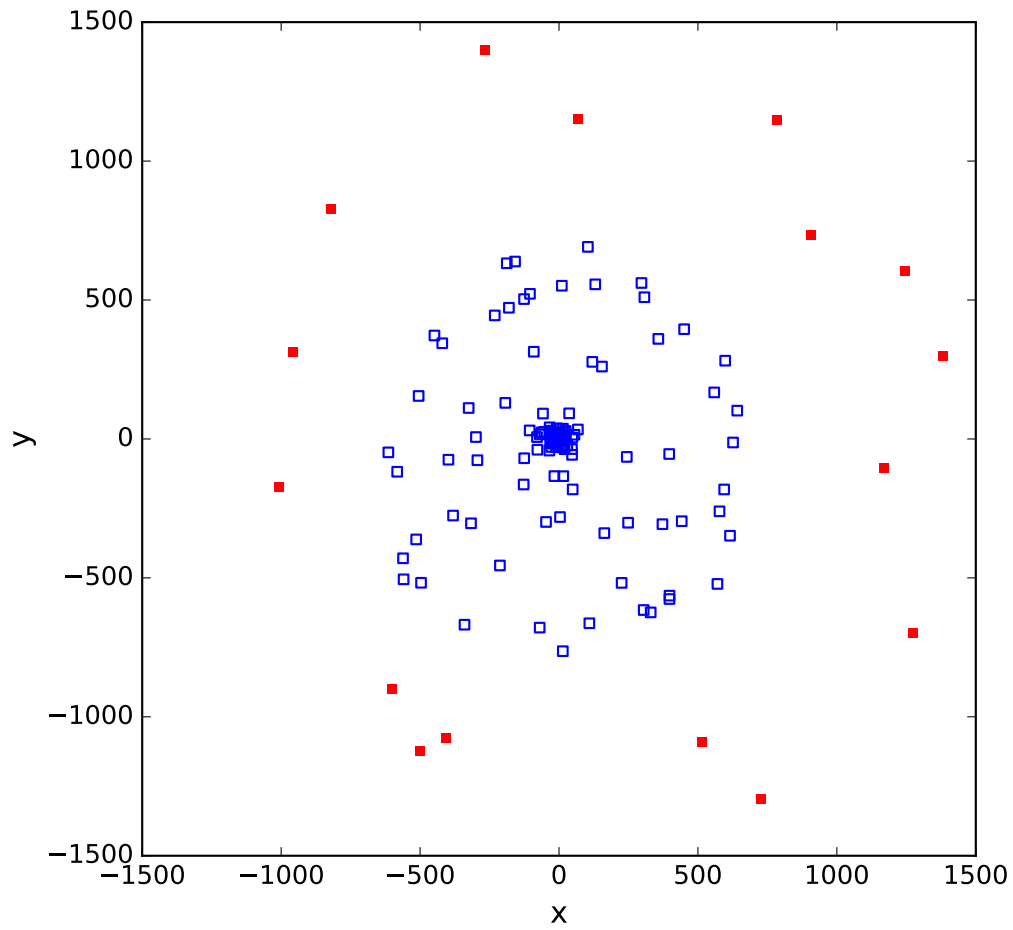
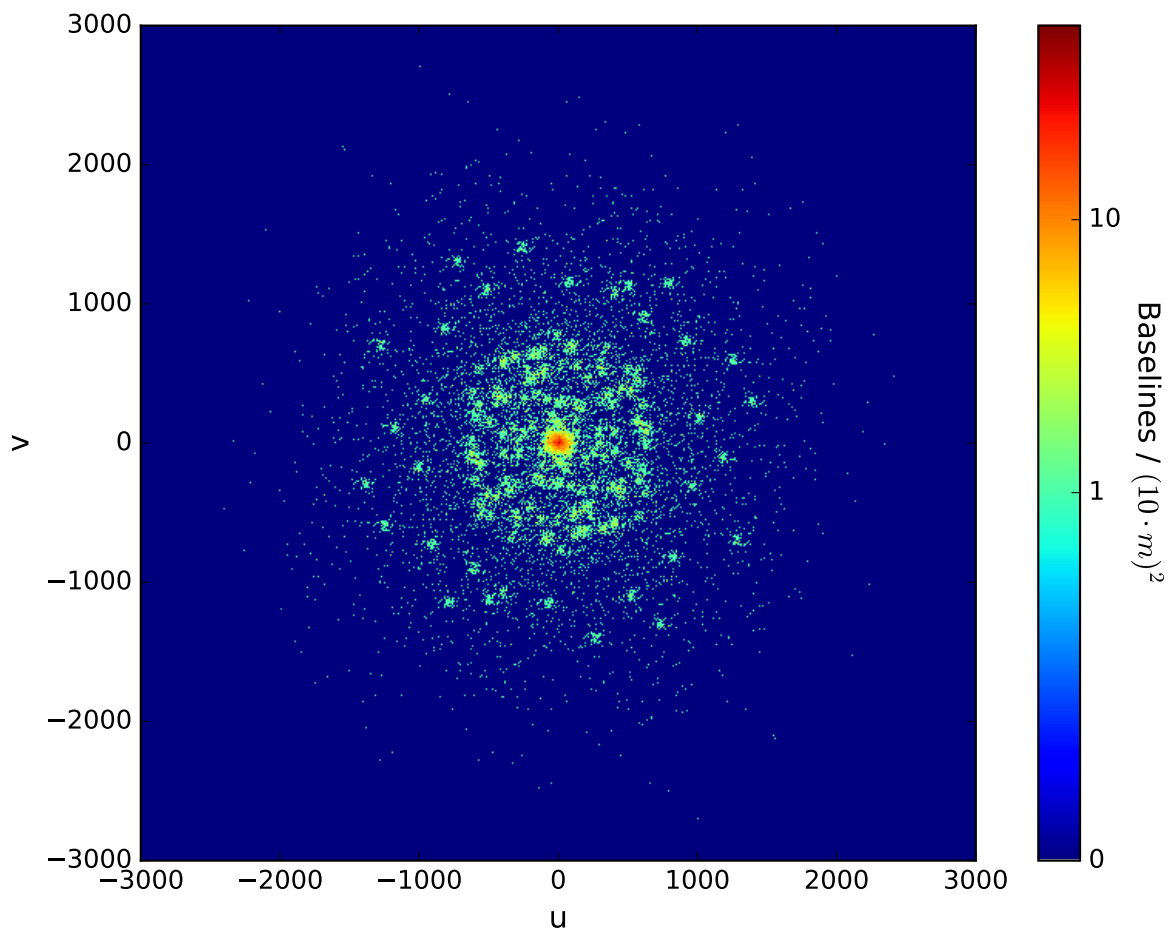
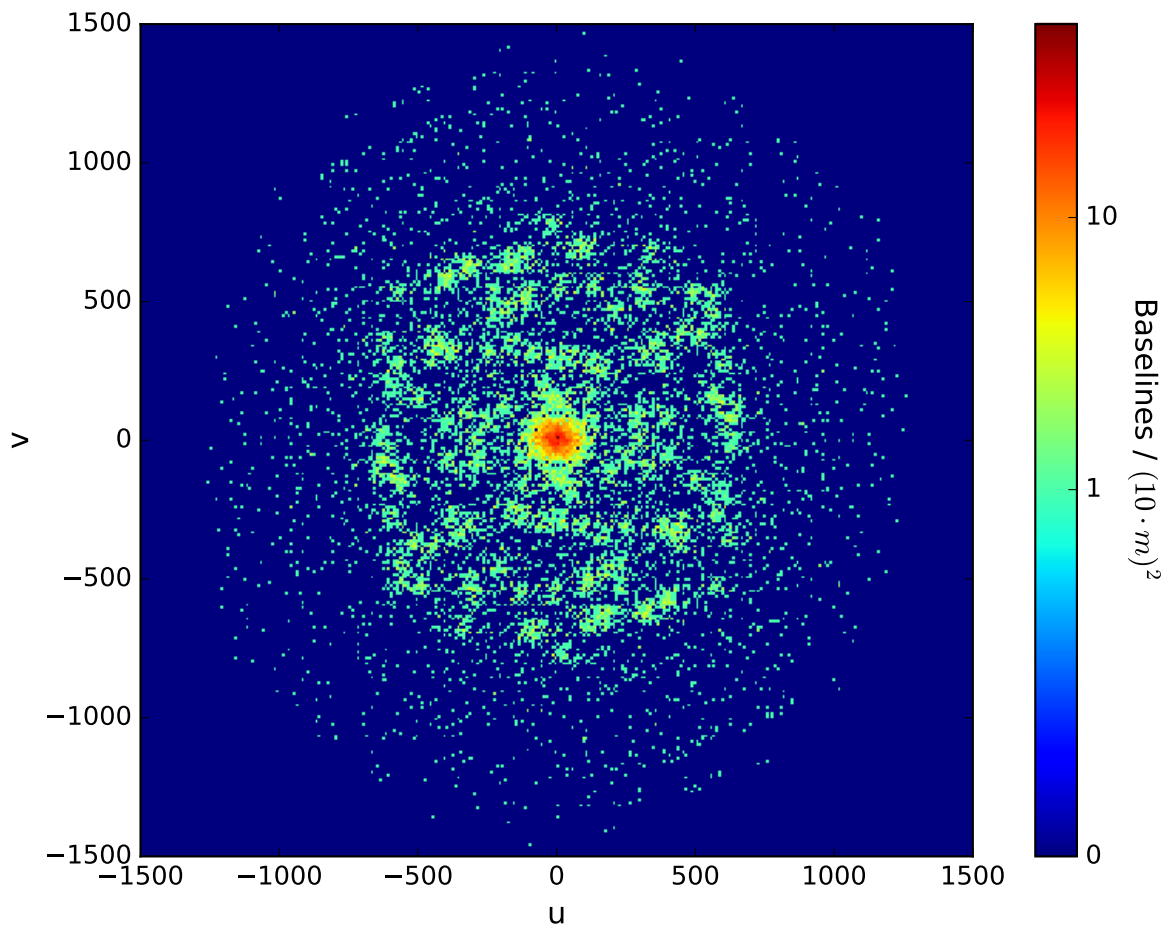


Figure 1.10: The MWA tile layout. Tiles are concentrated near the center but extend to a radius of 1.5 km. The 16 red tiles are not used in the EoR analysis.

Figure 1.11: The snapshot uv-coverage of the MWA for the full array (a) and central 112 tiles (b).



(a) Full Array



(b) Central 112 tiles

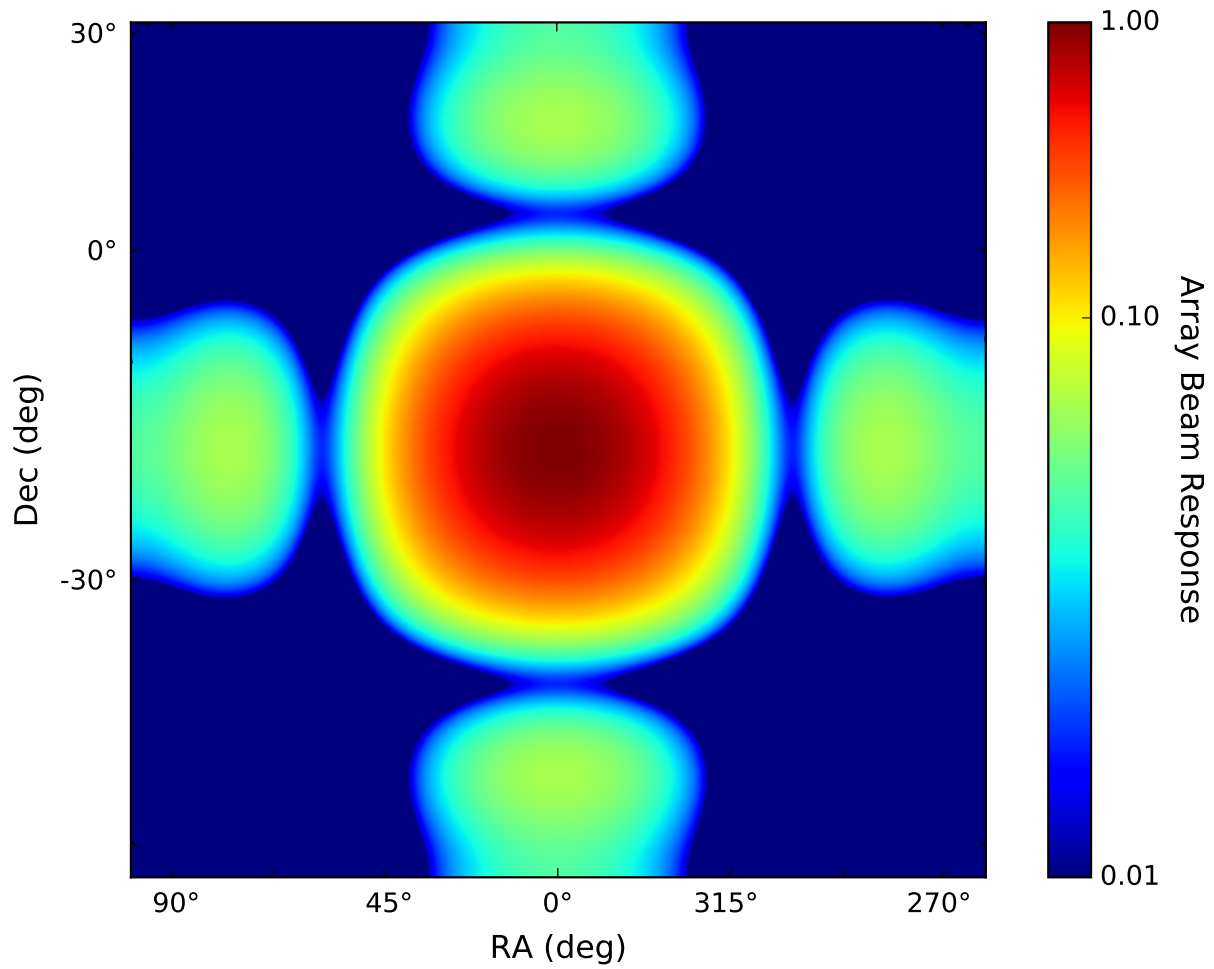


Figure 1.12: The MWA array primary beam response and first side lobes at zenith. Coordinates reflect the MWA EoR field at RA=0hr.

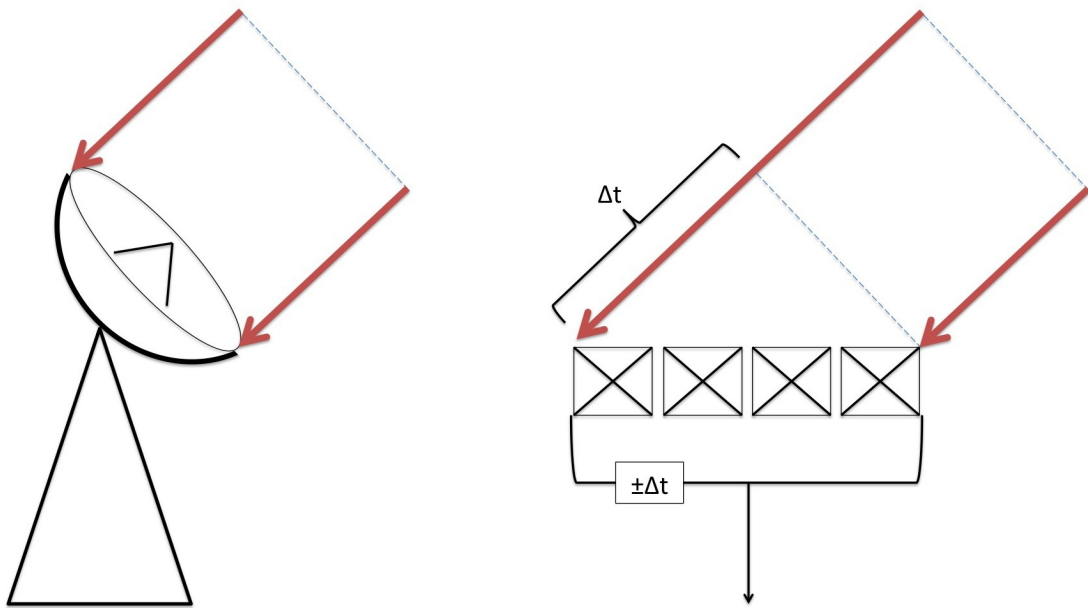


Figure 1.13: The MWA steers by adding a time delay to the signal that corresponds to the differential in arrival times of radio waves from an off-axis direction (right). This is in contrast to a dish antenna that is mechanically steered (left). Electronically steered arrays benefit from reduced system noise, a wide field of view, and the ability to observe in many directions simultaneously.

### 1.5 *Low-Frequency Southern Sky Surveys*

Established large sky surveys in the southern hemisphere at low radio frequencies are few and no one survey provides a sufficient model of foregrounds within the MWA band. It is worthwhile to discuss comparable surveys however, to highlight their strengths and weaknesses and to use as a point of reference as we dive into a new and somewhat nontraditional survey of the EoR foregrounds with the MWA. Ultimately these surveys are complementary to an MWA EoR foreground survey by providing SED information and enhancing reliability through cross-matching. Here I summarize each survey by frequency coverage, resolution, depth, and other noteworthy characteristics.

- **NRAO Very Large Array Sky Survey**

The NRAO Very Large Array (VLA) Sky Survey (NVSS [11]) at 1.4 GHz is our benchmark survey in terms of accuracy and completeness depth. The NVSS covers the full sky north of  $-40^\circ$  declination. The survey boasts a resolution of  $45''$  and rms noise levels 0.45 mJy/beam. Positional uncertainties in RA and Dec are at the sub-arcsecond scale for sources above 15 mJy and the completeness limit is 2.5 mJy.

- **VLA Low-frequency Sky Survey (redux)**

The VLA Low-frequency Sky Survey (VLSS [10]) at 74 MHz was re-processed using a corrected primary beam model among other improvements and released in 2014 as the VLSS redux (VLSSr [22]). The VLSSr complements the NVSS, covering nearly the entire sky north of  $-30^\circ$ , but it is not considered complete below  $-10^\circ$ . The resolution is  $75''$  and average rms noise level is 100 mJy/beam. The completeness limit varies but can be estimated at approximately 1 Jy. Nearly all VLSSr sources are also detected in the NVSS.

- **Sydney University Molonglo Sky Survey**

The Sydney University Molonglo Sky Survey (SUMSS [24]) at 843 MHz covered the full sky south of  $-30^\circ$  using the Molonglo Observatory Synthesis Telescope (MOST).

The resolution is variable with declination, but closely matched to the NVSS at  $45''$ . The catalog is considered to be complete to  $18 \text{ mJy/beam}$  above  $-50^\circ$  degrees. The catalog is also considered to be very reliable with an estimated 96% of spurious sources rejected using a decision tree classifier. Positional uncertainties are  $< 10''$  and only  $1\text{--}2''$  for sources brighter than  $20 \text{ mJy/beam}$ , while the uncertainty on the flux density scale is quoted at 3%.

- **The Molonglo Reference Catalog**

The Molonglo Reference Catalog (MRC [23]) at 408 MHz is unique in that it covers the full southern sky and was the best matched survey to the sky coverage, resolution, and frequency of the MWA EoR observations prior to the MWACS (see below). It was therefore the first choice for sky model subtraction in the initial stages of the EoR pipeline development. The MRC is relatively shallow, complete to 1 Jy with a detection limit of 0.7 Jy, but the reliability is believed to be better than 99.9%. The resolution is  $\sim 3'$ , but the positions are tied to the NVSS with a standard error between  $3\text{--}10''$ .

- **MWA Commissioning Survey**

The MWA Commissioning Survey (MWACS [20]) was the first large sky survey undertaken by the MWA science team. The relatively low  $2\text{--}3'$  angular resolution of the MWA increases sensitivity to low surface brightness objects that may be resolved out of other surveys. As a result, an MWA specific survey will most accurately capture the full extent of foreground emission from non-point sources. The MWACS covers  $\sim 6100 \text{ deg}^2$  below  $-14^\circ$  declination. The resolution is  $3'$  and the rms noise is about  $40 \text{ mJy/beam}$ . The catalog contains 180 MHz flux densities and spectral index measurements across three frequency bands centered at 119, 150, and 180 MHz. The MWACS quickly proved itself a superior EoR foreground model to the MRC. As its title suggests, however, the survey was completed during the MWA commissioning phase and, as a product of a partially built instrument, the catalog suffers from large uncertainties and biases.

Survey	$\nu$ (MHz)	$N_{\text{sources}}$	Dec	FWHM	$S_{\text{complete}}$
VLSSr	74	211,050	$\delta > -30^\circ$	75''	$\sim 1$ Jy/beam
MRC	408	12,141	$\delta_{1950} > -85^\circ$	180''	$\sim 1$ Jy/beam
SUMSS	843	95,491	$\delta < -30^\circ$	45''	18 mJy/beam
NVSS	1400	1,773,484	$\delta > -40^\circ$	45''	2.5 mJy/beam
MWACS	180	14,110	$\delta < -14^\circ$	180''	$\sim 500$ mJy/beam

Table 1.2: Overview of established large sky surveys at low radio frequencies in the southern hemisphere

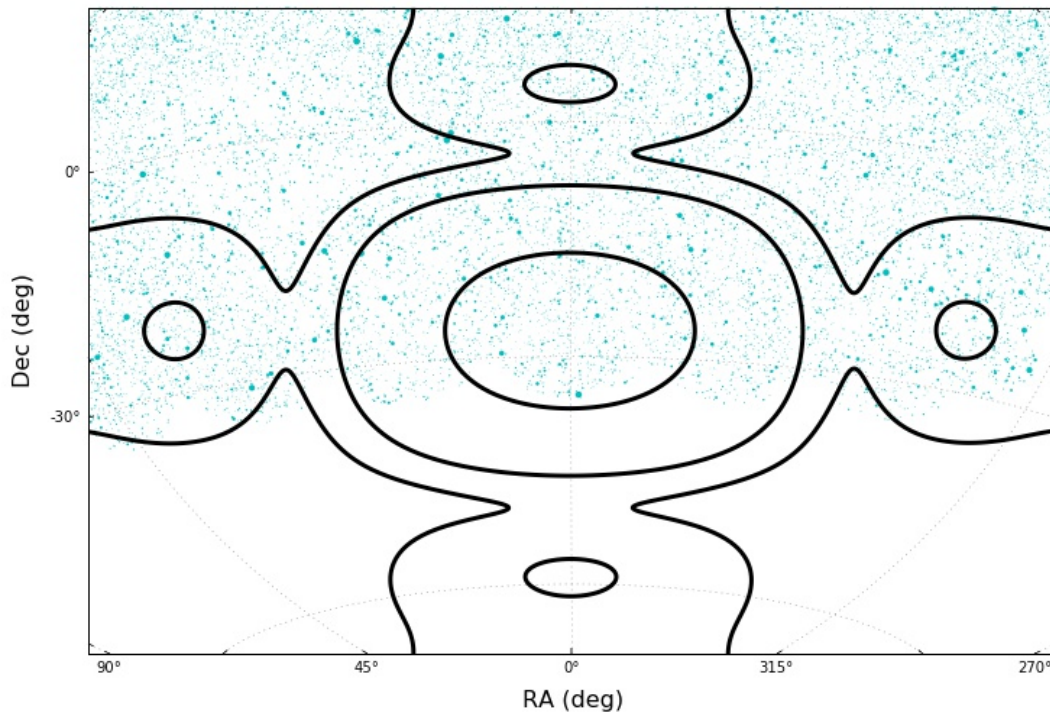
The deepest surveys are the NVSS and SUMSS with complementary sky coverage but their frequency and resolution are poorly matched to the MWA. The MRC has the most complete coverage of the southern sky but is very shallow, complete to only  $\sim 2$  Jy at 182 MHz. The VLSSr provides an excellent reference at even lower frequencies than the MWA and gives broad spectral coverage in combination with the NVSS, but it is also shallow and limited in its southern extent. The MWACS is the most true in frequency and resolution to MWA EoR observations but has limited sky coverage and large measurement uncertainties.

No single available survey provides a sufficient foreground model for MWA EoR analysis. Observations for the GaLactic and Extragalactic MWA All-sky survey (GLEAM, [44]) had commenced at the start of this work but a catalog would not be available for years. A limited survey tuned to the needs of the EoR0 field analysis was therefore initiated.

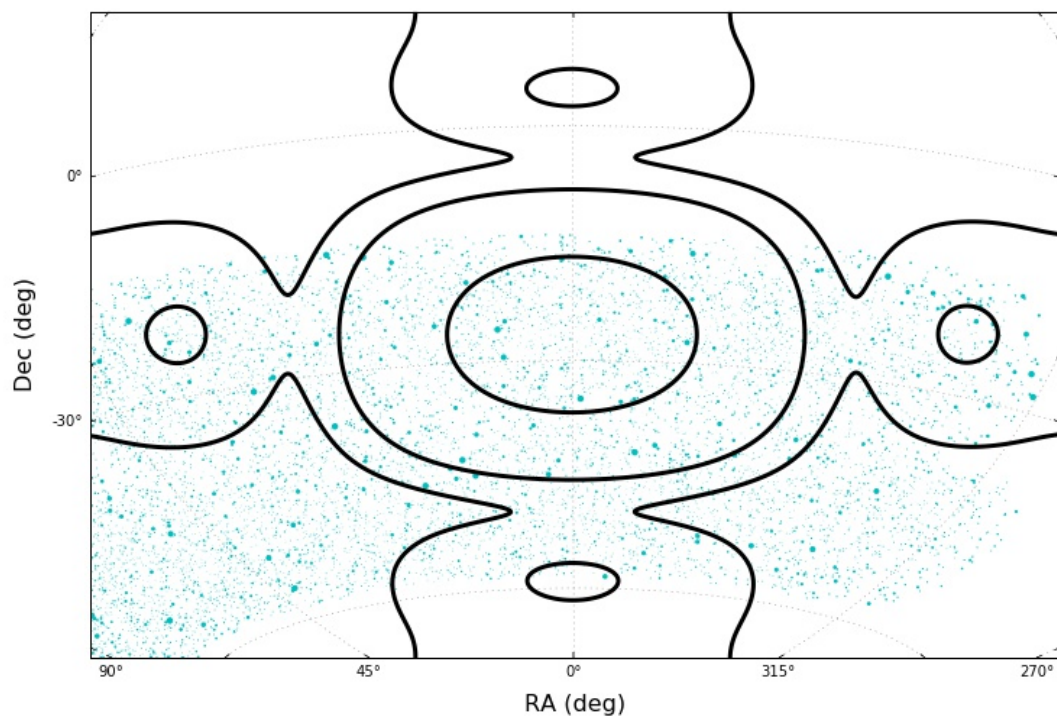
## 1.6 This Thesis

This thesis focuses on a survey of the MWA EoR0 field, one of three MWA EoR fields. The fields were specifically chosen to be relatively sparsely populated by bright or complex foregrounds, both discrete and diffuse. While bright point sources can be helpful for calibration, the first side lobes may easily exceed a signal-to-noise ratio (SNR) of 5, mimicking confident detections of true sources and increasing the level of side lobe confusion and false

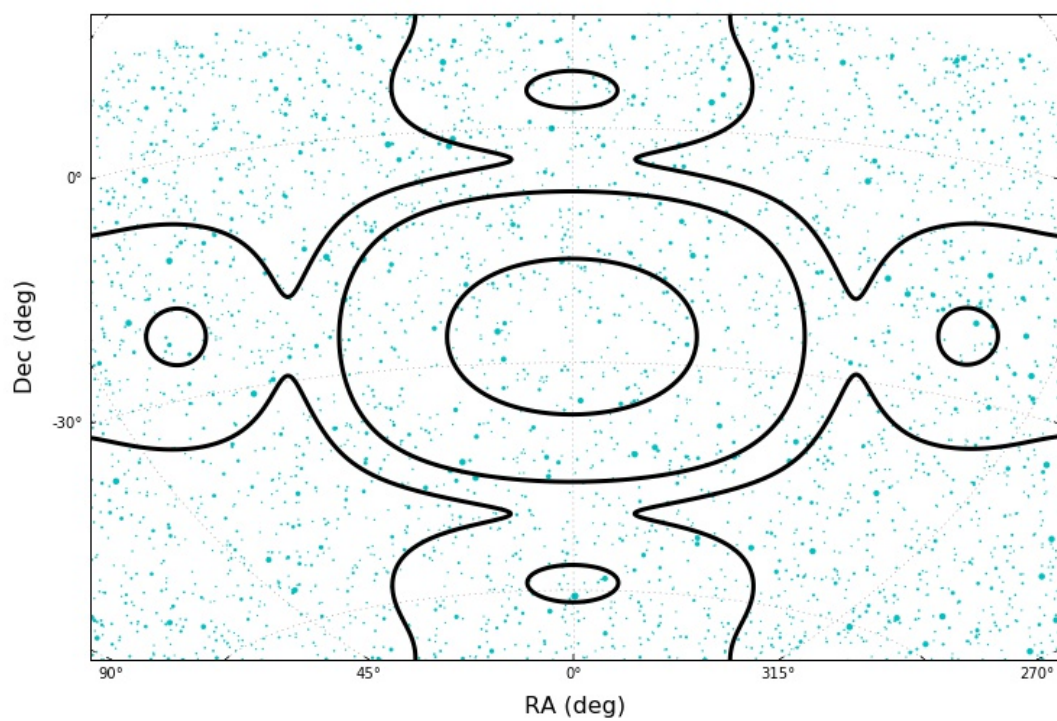
Figure 1.14: Sky coverage of low-frequency radio large sky surveys in the southern hemisphere. Surveys include MWACS (b), MRC (c), SUMSS (d), NVSS (e), and VLSSr (a). To illustrate survey depth, point diameters are scaled to source flux density projected to 182 MHz with an assumed average spectral index  $\alpha = -0.8$  (sizes are clipped at 20 Jy). The MWA EoR0 field zenith beam contours are overlaid at 1%, 5% and 50% beam response for reference.



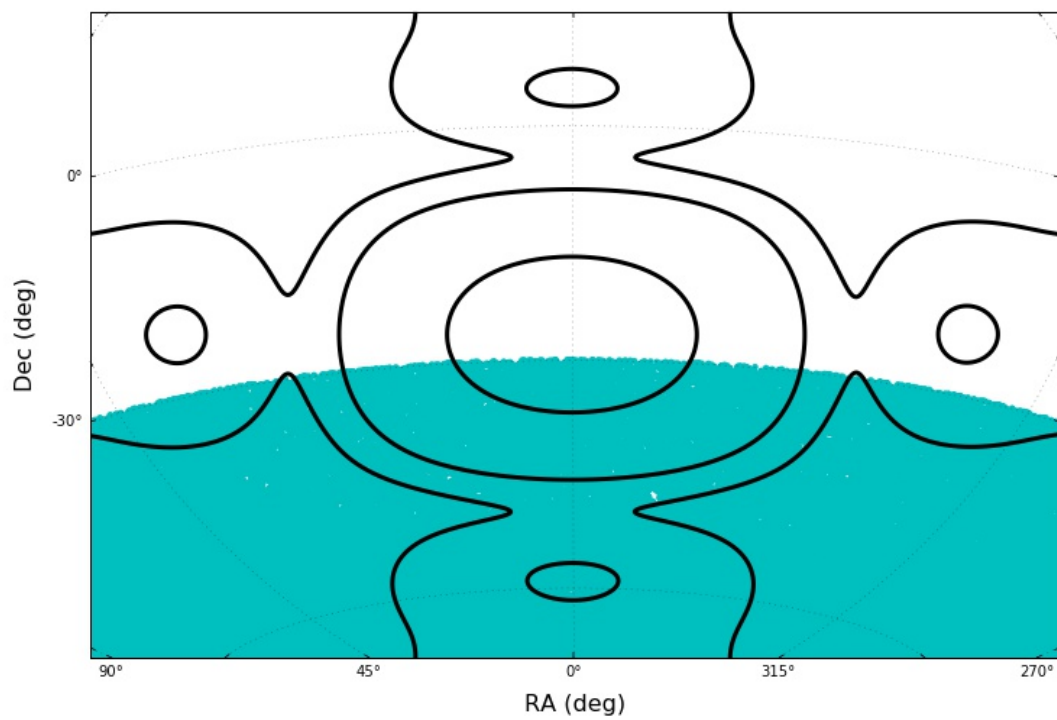
(a) VLSSr



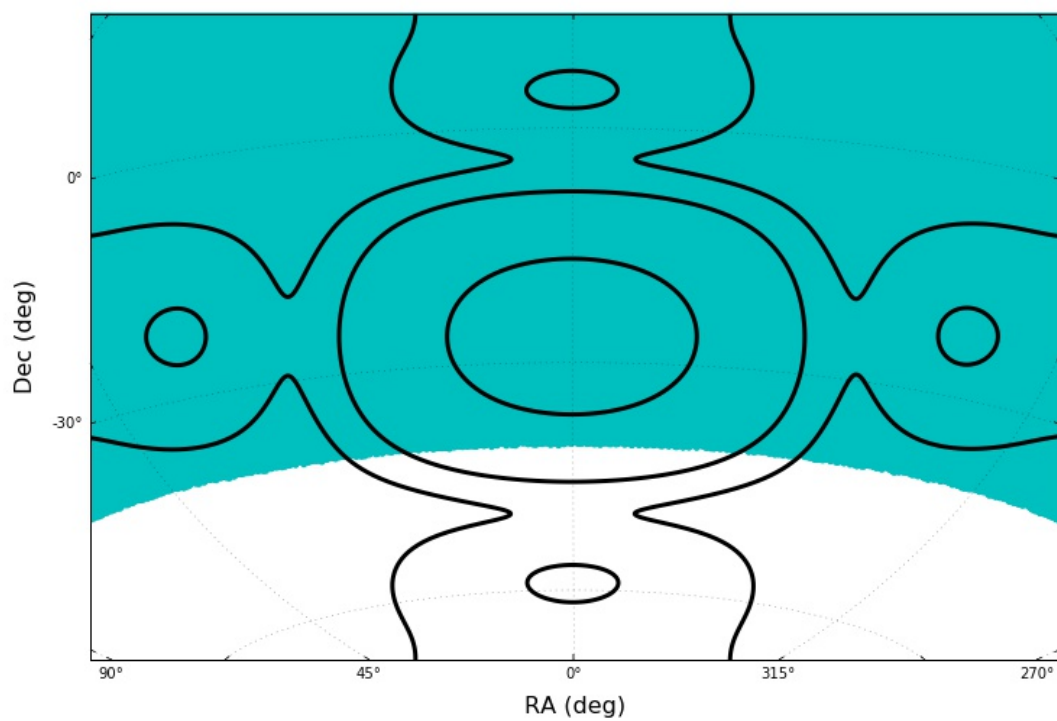
(b) MWACS



(c) MRC



(d) SUMSS



(e) NVSS

source contamination. Further, the brightest sources are often relatively nearby radio galaxies or starburst galaxies, with partially resolved complex emission regions that are difficult to model and subtract from the image data. The MWA EoR0 field pointing was chosen to minimize the number of bright and extended sources.

The MWA EoR0 field is centered at  $RA = 0$  hr and  $Dec = -27^\circ$ . It has no extremely bright or morphologically complex sources in the primary field of view. The brightest source is 28 Jy, and the most complex is the nearby Sculptor galaxy (NGC 253). The FWHM of the antenna beam is approximately  $20^\circ$ , but sources in the edges of the beam and first few side lobes are clearly visible and must be subtracted [30]. For this catalogue we concentrate on identifying sources in the primary beam but push out to the 5% power point, nearly to the first beam null. This spans  $\sim 40^\circ$  covering approximately  $1400 \text{ deg}^2$  of sky.

The survey is based on MWA EoR0 snapshot observations. Chapter 2 describes these observations, pre-processing, and initial data reduction. Due to the nature of the EoR power spectrum and the difficulty of making a detection, it is important that any foreground model be as accurate and reliable as possible. New methods are devised for source finding that attempt to circumvent errors introduced through imaging and source fitting in the image plane. Source finding, measurement, and association across snapshot observations are described in Chapter 3.

To maximize reliability of the final catalog, Chapter 4 describes a machine learning classifier designed to self-consistently assign all source candidates to a reliability class. This classification is used for source selection in combination with cross-matching to comparable radio surveys. Cross-matching is discussed in Chapter 5. This process allows for the robust identification and removal of contamination as well as the discovery of new radio detections. The final catalog is presented in Chapter 6 along with an in-depth analysis of the astrometric and flux scale accuracy of the catalog.

We dive into the EoR analysis pipeline in Chapter 8. The EoR0 survey catalog is input for both calibration and foreground removal and the improvement is measured against the MWA commissioning survey catalog within the primary beam. We then motivate the need

to account for widefield foregrounds, particularly in the side lobes of the primary beam, and build a “master” foreground catalog by cross-matching the multiple available radio surveys. Lastly, we demonstrate a simple approach to build extended source models for bright complex sources, and include a model for NGC 253 in the master foreground catalog. The 182 MHz master catalog based on the EoR0 field survey is the current standard in use for calibration and foreground modeling and removal in the U.S. MWA EoR analysis.

In Chapter 7, we diverge from the EoR focus to explore some of the most interesting sources uncovered through the catalog creation and analysis. We consider the spectral properties and potential identifications of our ultra-steep spectrum population, including 25 newly discovered radio sources. This thesis is summarized in Chapter 9, where we review the novelty of the methods employed, our results and lessons learned, and pose open questions such as the impact of AGN variability on the fidelity of the EoR foreground model over time.

## Chapter 2

### OBSERVATIONS & DATA REDUCTION

The data for this survey were taken on August 23, 2013, early in the MWA science operations and analysis pipeline development. From this night, approximately 3.5 hours of observations, roughly centered on field transit at zenith, were selected and designated the “golden set”. The purpose of the golden set was to establish a high quality and consistent reference data set for EoR pipeline development and performance analysis while holding observational conditions constant. The observational conditions were nearly ideal. The field was at high elevation, minimizing wide field effects, and the ionosphere was stable.

The disk of the Milky Way, like other bright sources in the side lobes, “throws” power into the the primary field of view if it is not correctly modeled and subtracted. Lacking a sufficient model, the golden set observations are restricted to a window where the galactic plane had set. Given the extended nature of the disk, residual galactic contamination was still found to afflict the first pointing. The first and last pointings, lowest in elevation, were later found to also suffer from wide field effects that were ultimately detrimental to the survey quality. For these reasons, the first and last pointings ( $\sim 60$  minutes) of the golden were dropped for this thesis.

A total of 2.5 hours of observations were selected, comprised of seventy-five consecutive 2 min snapshots. The observations are centered at 182 MHz with 31 MHz bandwidth (138.9-197.7 MHz). Snapshots are consecutive 112 seconds of integration with 8 second gaps in-between. The MWA EoR observation strategy uses a ‘drift-and-shift’ survey pattern. The MWA antenna can only point in distinct locations in elevation and azimuth [40], so the sky is allowed to drift through the antenna beam for approximately 30 minutes (15 snapshots) before the antenna is re-pointed. Because the sky drifts  $15^\circ$  through a pointing, a source

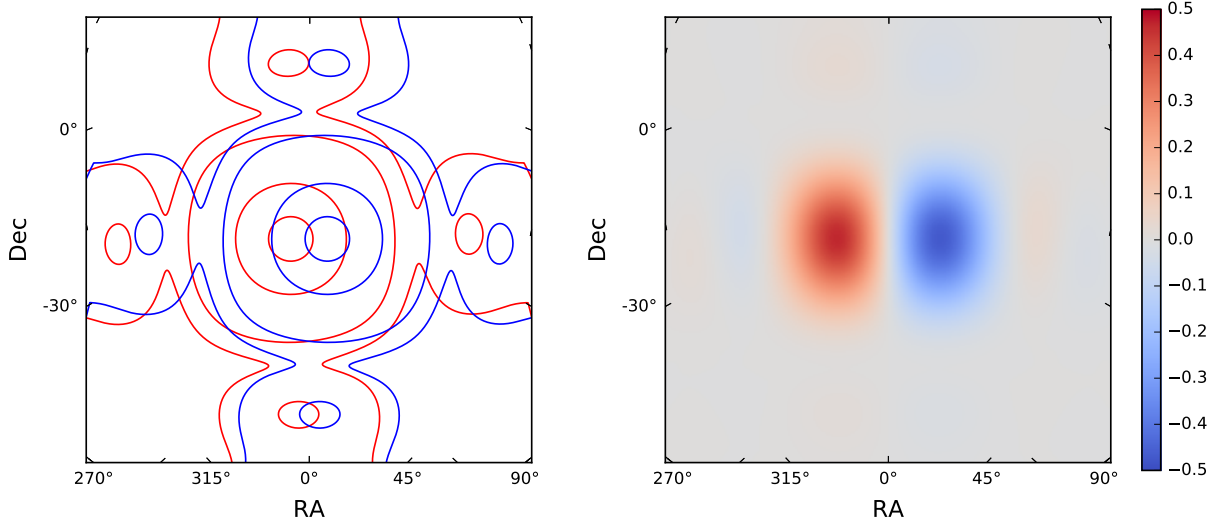


Figure 2.1: Left: Beam response contours for the first and last snapshot of the zenith pointing. Lines are drawn at 1, 5, 50 and 90% beam response. Right: The difference in beam response between the first and last snapshot of the zenith pointing. The sensitivity of the instrument to a source can change up to 50% across a pointing and even more so as the sky rotates over many hours.

can nearly traverse the  $20^\circ$  FWHM of the primary beam. The beam sensitivity to a source changes across a pointing by 25% on average and can exceed 50%. Many sources move into or out of the field of view entirely. Figure 2.1 shows the difference in the beam response between a snapshot at the beginning of a pointing and a snapshot at the end of a pointing. These 75 snapshot observations covered 5 distinct antenna pointings.

## 2.1 Pre-processing

The raw MWA visibility data must first be flagged for radio frequency interference (RFI). RFI flagging is performed by AOFlagger within the COTTER software package [28]. AOFlagger was initially developed for the Low Frequency Array (LOFAR). LOFAR, located in the

Netherlands with remote stations throughout Europe, relies heavily on RFI flagging and removal of contaminated data. The MWA on the other hand, is located at the Murchison Radio Observatory (MRO) in Western Australia. The MRO site, remote and sparsely populated, is one of the most radio-quiet sites on the planet with the infrastructure required to host such a facility. This opposing strategy of RFI avoidance means that only  $\sim 3\%$  of data were contaminated.

After flagging out RFI, the data are averaged in time and frequency. The initial 0.5 sec cadence and 40 kHz frequency resolution, are reduced to 2 sec and 80 kHz in exchange for a more manageable data volume. The snapshot data are then written to UVFITS file format for transfer and further processing.

## **2.2 Fast Holographic Deconvolution**

Fast Holographic Deconvolution (FHD [36]) is a software package developed at the University of Washington for calibration, imaging, and source removal within the EoR analysis pipeline<sup>1</sup>. FHD is an efficient implementation of A-projection [26, 5, 27]. A-projection effectively performs direction dependent calibration by integrating model visibilities with the known antenna response or holographic beam pattern. By introducing a Holographic Mapping function, FHD pre-computes the time limiting step of forming and gridding the model visibilities resulting in a significant increase in speed.

FHD has two branches for source removal that are used for different purposes. The first is full deconvolution, in which flux density peaks are identified and iteratively “cleaned” from the dirty visibilities. The second is “Firstpass” in which an input catalog of sources is simply modeled and subtracted from the visibility data. These branches in relation to the EoR pipeline and this survey are illustrated in Figure 2.2. Firstpass is not used for the EoR foreground survey or cataloging process, but it is introduced here as the second branch of FHD and is detailed in §8 where it is used for testing catalog performance within the EoR

---

<sup>1</sup>This thesis rests on an understated but significant amount of time spent testing FHD functionality, developing quality control checks, and tracing bugs.

analysis pipeline.

### *Full Deconvolution*

Full deconvolution is used to build the EoR0 field survey. It is the first step to source finding, measurement, and selection. An input foreground catalog is still required for the purpose of calibration. The MWA Commissioning Survey [20] was the most reliable choice at the time and was therefore used to establish the flux scale and the mapping of sources onto the sky despite large measurement uncertainties and missing coverage north of  $14^\circ$  declination. Biases that propagated into the EoR0 field survey were later revealed and corrected for.

During each deconvolution iteration, the brightest pixels are identified and fit with a Gaussian-approximated array beam shape to determine the flux density. An amplitude gain of 0.2 is applied to create a source component. The source components are forward modeled through the direction dependent instrument model to update the residual sky image. The deconvolution loop is stopped if the rms of the residual image increases or if the source fitting fails and no valid components can be extracted. Otherwise, deconvolution halts after a maximum of 500 iterations or a maximum of  $3 \times 10^5$  components are subtracted.

Standard CLEAN components are located at pixel centers which are unlikely to be the true peak location of the source. Negative compensatory components are needed to correct for the offset. FHD operates differently. The position of the Gaussian fit to the source is centroided to a floating point location. All components are therefore strictly positive-valued and can be expected to describe the true flux density distribution of the source. This difference is important to the source finding process described in §3.

In four of the 75 snapshot observations, the deconvolution loop halted prematurely due to lack of convergence on a single bright and extended source (later identified as the local starburst galaxy, NGC 253). These deconvolved less than  $10^5$  components and were excluded from the remaining analysis. The other 71 snapshots deconvolved at least  $2 \cdot 10^5$  components with an average of  $2.7 \cdot 10^5$ . The deconvolution threshold is an average of 57 mJ at beam-center and increases as  $1/beam$ . A gain factor of 0.2 was used for component extraction

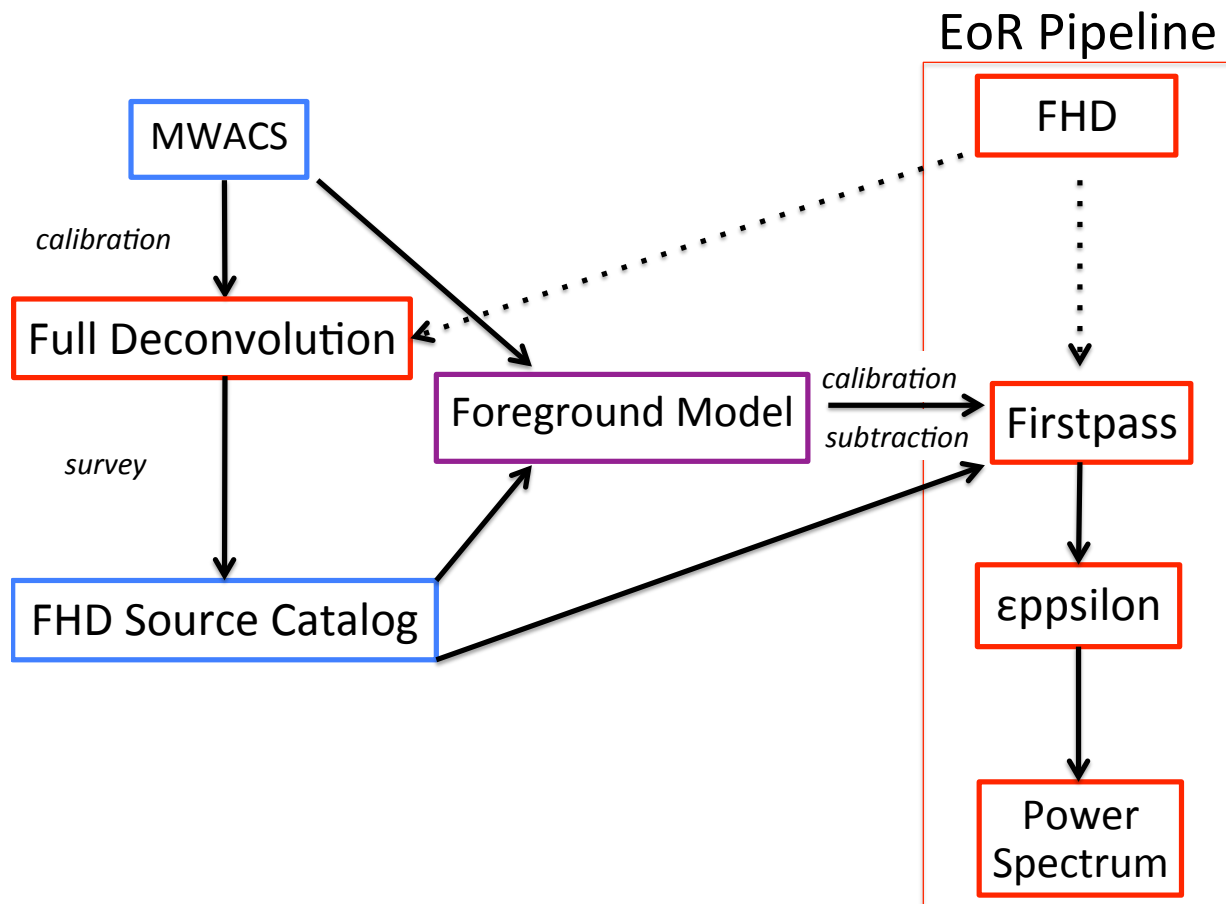


Figure 2.2: A flow chart illustrating the two source removal branches of FHD: Full Deconvolution, and Firstpass. I use full deconvolution to build a new foreground catalog of sources in the EoR0 field. This catalog is then fed into Firstpass for source removal before the power spectrum analysis with  $\epsilon$  in the EoR analysis pipeline.

but this was allowed to vary slightly under certain conditions and was recorded for each component along with its flux density and RA/Dec position.

The output of FHD for each snapshot includes: the deconvolved component arrays, residual and restored images, beam maps, and the observation meta data. The weighted mean restored image of the EoR0 field is shown in Figure 2.3, illustrating point sources in the primary beam and first side lobes. The average beam power across all snapshots is contoured in red for reference.

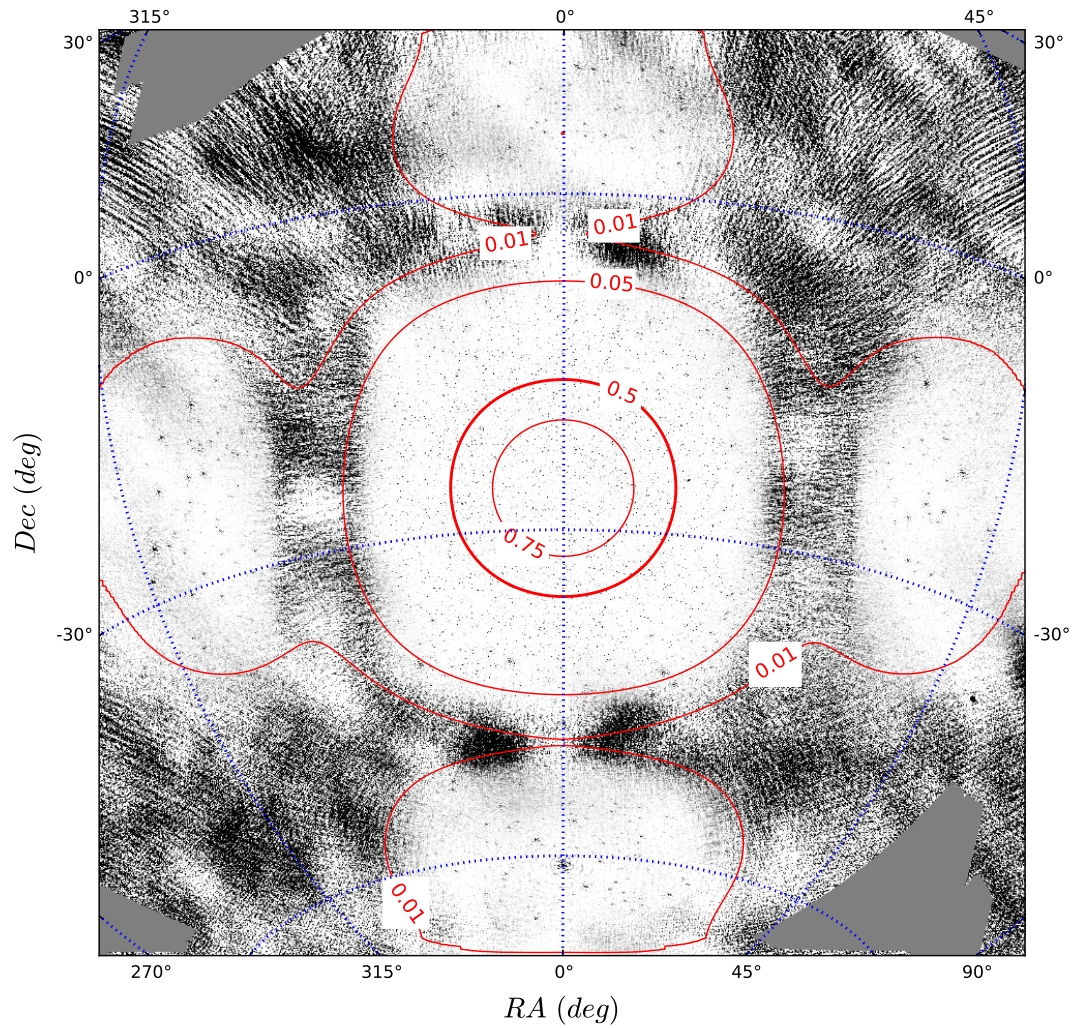


Figure 2.3: The weighted mean Stokes I image of the EoR0 field showing the primary beam and first side lobes. The average beam power is contoured in red for reference. Traditionally source finding would be limited to  $beam > 0.5$ , but sources further afield must also be modeled and subtracted for the EoR analysis. The color scale is linear in units of Jy/pixel and clipped to the range 0–1 Jy.

## Chapter 3

### SOURCE FINDING & CHARACTERIZATION

Finding and extracting sources from radio survey data is typically performed in the image plane. After deconvolution, the clean components are usually convolved with a Gaussian approximated PSF and restored to the residual image. Given multiple observations, images may be combined through stacking or mosaicing to drive down the rms noise and increase the signal-to-noise ratio. One of numerous available source extraction software packages can then be used to identify peaks in the image and fit sources for shape, flux density, and the local background level. Advanced source extraction algorithms work very well for the purpose of measuring sources in restored images, but there are certain assumptions and inherent uncertainties in this process.

A point source can, in general, be accurately modeled as a delta function convolved with a Gaussian approximated PSF. This works best when a majority of the source's flux has been deconvolved and restored in a single snapshot. When a significant fraction of flux is below the deconvolution limit and remains in the residual image, any error in the modeled PSF compared with the true PSF will propagate into the flux density estimate. When many observations are stacked or mosaiced together before source fitting, the PSF becomes difficult to predict and model as it can change with direction and time. Often, sources are simply fit with a two-dimensional Gaussian and the PSF shape itself is ignored in source extraction. This still assumes a Gaussian shape and is still subject to an errors in the PSF model introduced in the restored image. A significant fraction of sources at radio frequencies are extended, partially resolved, or otherwise morphologically complex; e.g. radio galaxy jets and lobes, star-forming galaxies, and radio relics and halos in galaxy clusters. A Gaussian model won't capture this complexity and cannot accurately measure the peak and integrated

flux density. We therefore take a different approach.

FHD deconvolution automatically clusters deconvolved components into sources in order to efficiently update the model with each iteration. Multiple components of a single point source can be combined into only a single component that needs to be modeled and subtracted. During early pipeline development, the resulting crude source list was analyzed statistically for quality control or to create restored images as a qualitative check. This simple grouping of clean components was not meant to accurately perform source finding and measurement and the resulting source lists were contaminated with spurious noise peaks and side lobes.

When a survey of the EoR0 field was initiated, the existing functionality of grouping clean components into sources was the basis for the development of a more robust clustering approach to source finding. In contrast to typical image based source finding, clustering is insensitive to source shape or extent, or the assumed PSF model beyond the deconvolution process. Further, because source finding is performed separately on each observation, we retain the detection frequency, a strong indicator of source reliability and a driving factor in identifying true sources from contamination. True sources will appear in most observations, while false sources should appear in few. This chapter describes the process of source finding and measurement by clustering components into sources and associating sources across snapshots.

### ***3.1 Component Clustering***

DBSCAN [15], a density-based hierarchical clustering algorithm, was used to identify spatially isolated clusters from the array of source components produced by FHD. No restored image was used for source extraction. DBSCAN hierarchical clustering works by identifying local maxima in the density distribution of components and builds clusters hierarchically from these cores.

The input parameters for DBSCAN are the neighborhood radius within which a point is considered to be a part of the same cluster, and the minimum number of points required

within that radius for a new core point to be formed. Distances were calculated assuming euclidian geometry, as this approximation is valid on arcminute scales. Figure 3.1 illustrates how DBSCAN forms a cluster.

For this application the minimum number of components was set to one so as not to exclude sources with only a single component extracted, and the neighborhood radius was set to half of the 2.3' array beam width. This radius was found to maximize the number of sources detected in all observations. A radius of one-quarter the beam width maximized the number of sources detected uniquely in all observations, meaning that some close pairs (e.g. radio lobes) will be blended. These are addressed in §5.

Approximately 5000 clusters identified in each snapshot, for a total of  $3.55 \times 10^5$  across all observations. For each cluster, or source candidate, the flux-weighted mean position, standard deviation, and rms width were calculated and the detection was flagged as extended if the rms width exceeded the beam width. The component fluxes were summed and, if the source was not extended, a correction was added to account for the estimated residual source flux. Given the gain factor  $g_i$  for each component and the number of components deconvolved  $N_{comp}$ , the true flux  $S_{true}$  of a point source can be estimated according to Equation 3.1.

$$S_{true} = \left( \sum_{i=1}^{N_{comp}} S_i \right) / \left( 1 - \prod_{i=1}^{N_{comp}} (1 - g_i) \right) \quad (3.1)$$

### 3.2 Cross-Snapshot Association

A second round of spatial clustering was then performed to match associated detections across snapshots and create the final sample set of source candidates. Various clustering radii were tested to optimize source association but minimize contamination. A radius equal to half of the beam width was found to maximize the number of sources detected in all observations, but a radius equal to one quarter of the beam width maximized the number of sources detected *uniquely* in all observations. This is an important distinction.

Consider a radio galaxy with two lobes separated by a small angular distance close to the

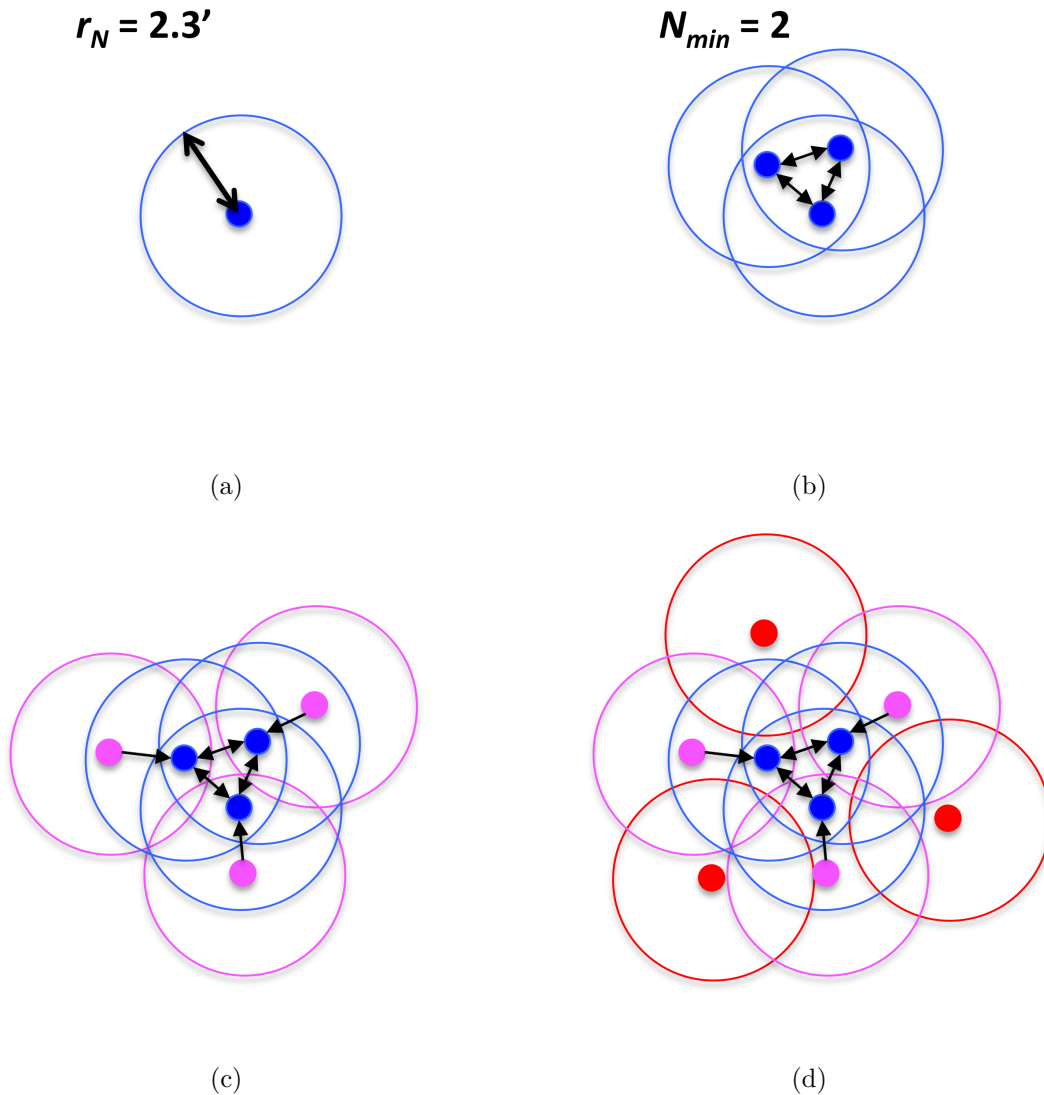


Figure 3.1: DBSCAN builds a cluster by first identifying a “core” point near a local maximum in the density distribution and then searching for other points within a specified radius  $r_N$  (a). A core point must have a minimum number ( $N_{min}$ ) of other points within  $r_N$  to form, and clustered points that also have  $r_N$  neighbors become core points themselves (blue points in (b)). A clustered point that falls within  $r_N$  of a core point but does not itself have  $N_{min}$  neighbors is “density-connected” to the cluster (purple points in (c)). If a point is not density-connected to a cluster it is labeled as background noise except in the case of  $N_{min} = 1$  (red points in (d)). For component clustering, we set  $N_{min} = 1$  because a true source may have only a single component deconvolved. On association across many snapshots, we set  $N_{min} = 2$  to discard many spurious detections.

initial clustering radius. Due to variations in the beam response and positional uncertainty between snapshots, the lobes may be clustered into one single source in some snapshots but not others. Then on cross-association, using a smaller clustering radius can result in three apparent sources (the lobe and their centroid combination), with each detected in only a fraction of observations. I therefore used the larger radius, as it is easier to combine multiple components than to split one. A total of 9490 spatially isolated source candidates were identified with detections in at least two snapshots. A roughly equivalent number were detected only once and discarded as noise. If multiple detections of a source were found within a snapshot the detections were combined into a single source per observation by summing their flux and centroiding their position.

For each source candidate, a  $3\sigma$  clip was first applied to the flux distribution of all detections to exclude coincident noise or side lobes. The RA and Dec positions and uncertainties were taken as the mean and standard deviation. The flux error for each detection was estimated from the residual image by fitting the background rms as a function of antenna beam power. This error was used to find the weighted mean flux, standard deviation, and standard deviation of the mean. The standard deviation  $\sigma_S$  is the intrinsic scatter of the measured flux, but is poorly constrained for small  $N_{det}$ . The standard deviation of the mean  $\sigma_{\bar{S}}$  therefore gives a complementary, if indirect, measure of flux uncertainty and is related to the average beam power of the detections.

The resulting set of 9490 source candidates is still highly contaminated. Figure 3.2 shows the distribution of  $N_{det}$ . While the maximum number of sources are detected in all observations, the second peak is at only 2 snapshot detections with no clear dividing point. Real sources will move out of the field or fall below the detection threshold as the beam drifts, but the excess at small  $N_{det}$  is a result of spatially coincident noise fluctuations and side lobe contamination. The two measures of SNR show a very clear bimodality as well and a strong correlation with  $N_{det}$ . This is shown in Figure 3.3. Together, these are the most important indicators of source reliability as we will explore in the next chapter.

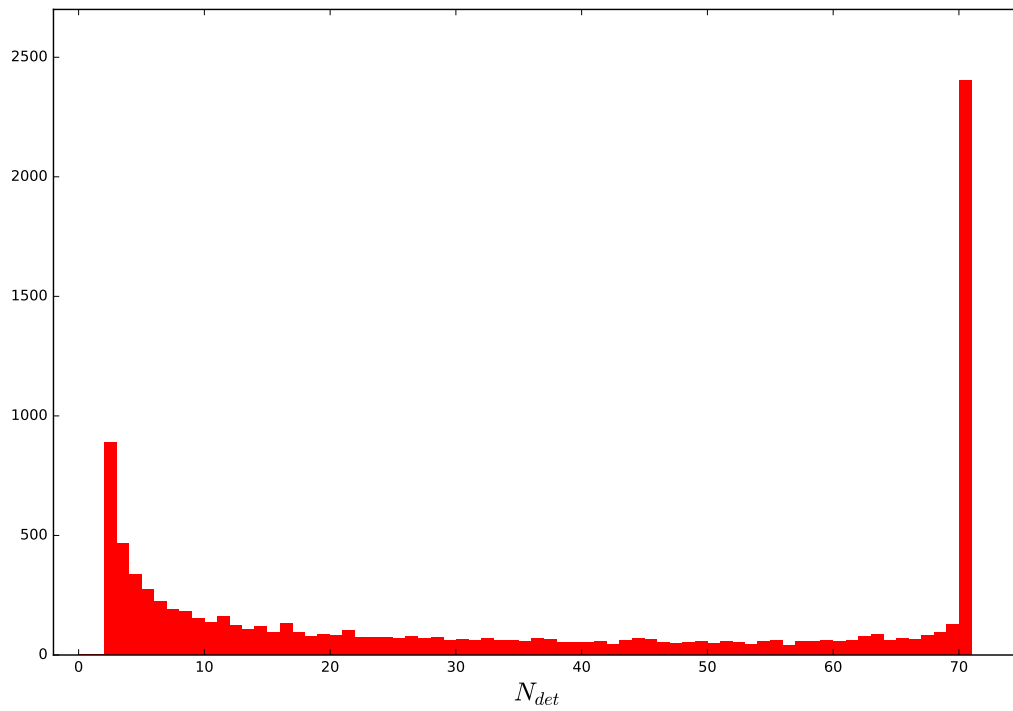


Figure 3.2: The distribution of  $N_{det}$  for the 9490 source candidates. The maximum number of sources are detected in all observations, but the second maximum is at only 2 snapshot detections with no clear dividing point. Real sources will move out of the field of view or fall below the detection threshold as the beam drifts, but the excess at small  $N_{det}$  is a result of spatially coincident noise fluctuations and side lobe contamination.

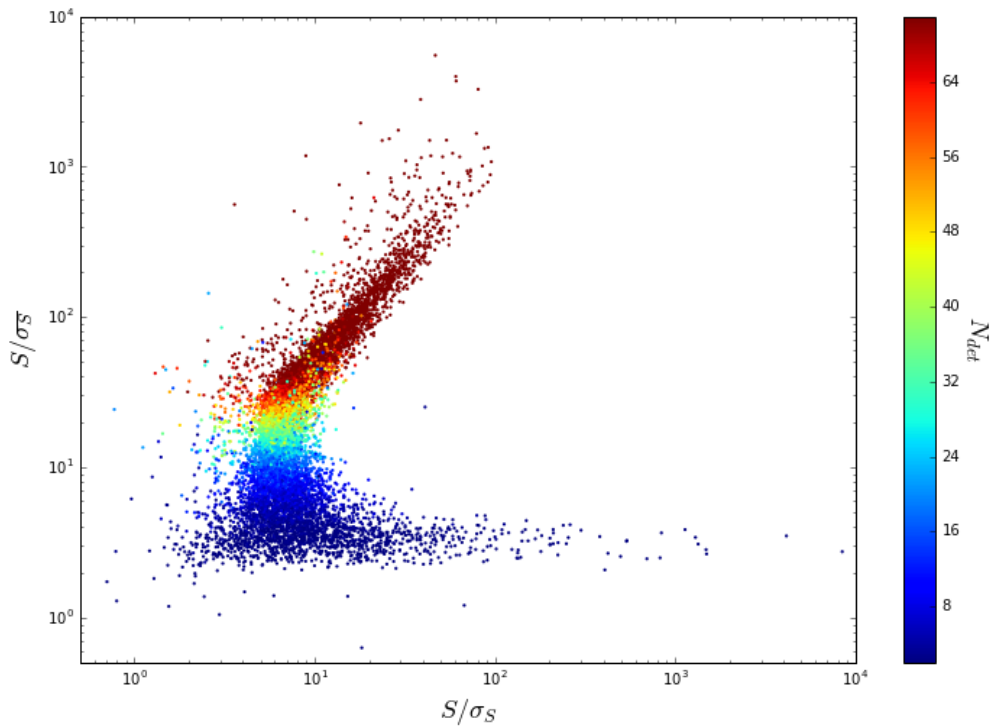


Figure 3.3: The two measures of SNR are correlated for sources detected in most observations. Sources with few detections ( $< 10$ ) are poorly constrained and may have artificially high or low  $\sigma_S$ , whereas  $\sigma_{\bar{S}}$  is estimated from the residual background. As an indirect measure of the error however,  $\sigma_{\bar{S}}$  may not be accurate for a given source particularly in the presence of side lobes.

## Chapter 4

### RELIABILITY CLASSIFICATION

The quality of a blind sky survey can be described in terms of completeness and reliability. Completeness refers to the probability that a true source is detected. If all sources considered to be detectable are detected, the catalog is 100% complete. Reliability describes the probability that a detected source is in fact true. If no false positives contaminate the catalog, it is 100% reliable. A survey with both high completeness and high reliability is ideal, but maximizing completeness often comes at the expense of reliability and vice versa.

The EoR0 foreground survey emphasizes reliability. The reason for this is that reliability impacts both calibration and foreground power removal. The EoR power spectrum analysis relies on the assumption that foregrounds are constant and spectrally smooth, having no sharp features in the frequency dimension. False sources are not constant and the spectral behavior of noise, including side lobes and other artifacts, is complex not well understood. False sources in the sky model would result in poor calibration and improper flux scaling. It is unclear if and how the subtraction of false sources would impact the EoR window.

Reliability is typically estimated with respect to a comparison survey. Cross-matching to overlapping surveys can help to identify contamination, but this is limited by the reliability of the comparison survey, how "detectability" is determined relative that survey, and the incidence of false matches. Although most true sources seen by the MWA should be detected in deeper, higher resolution surveys, we can expect to find ultra steep spectrum and diffuse sources that are not. Requiring a match will result in a loss of completeness as well as the loss of particularly interesting detections of rare-type sources.

While helpful, we use cross-matching only in conjunction with a self-consistent reliability determination. By convention in the physical sciences,  $\text{SNR} > 5$  is the standard statistical

confidence level required for "discovery". This  $5\sigma$  rule eliminates statistical anomalies to one part in 3.5 million but operates on the assumption of a Gaussian noise distribution that can be sufficiently measured or predicted. In reality, the noise behavior is complex and non-thermal sources of noise can lead to false confidence in spurious detections.

Particularly problematic in radio images are the side lobes of bright sources that result from imperfect calibration and deconvolution. These side lobes can be significantly brighter than the background rms noise level. They may also be spatially coincident between consecutive observations due to the short integrations (small  $uv$  rotation) and low resolution. This false confidence and seemingly non-spurious behavior make it challenging to identify contamination in an automated way.

A significant number of noise and side lobe sources clearly contaminate the sample of source candidates. In the EoR0 field image shown in Figure 2.3, you can see small clusters of sources near the edge of the field and in the side lobes. These are side lobes surrounding bright sources. An excess in the number density of source candidates is found out to a  $1^\circ$  radius of sources brighter than 10 Jy (Figure 4.1, black) and a  $5\sigma$  selection appears to be insufficient on its own to exclude the brightest and most problematic contaminants.

Fortunately we have much more information available to us to assess reliability in a self-consistent way. Below I outline an unsupervised machine learning classification scheme developed to assess source reliability given multiple observable parameters. In addition to SNR, reliability determination relies heavily on the detection rate given the many observations. Specific parameters are also engineered to target side lobes. The resulting classifications are later used in conjunction with cross-matching to inform the final catalog selection.

#### **4.1 Machine Learning**

Machine learning based classification algorithms come in two forms; supervised and unsupervised. Supervised methods rely on a data set with known classifications to train the model. Unsupervised methods have no training data on which to learn, and instead search for clusters or patterns in parameter space to differentiate subsets of the sample. I begin with

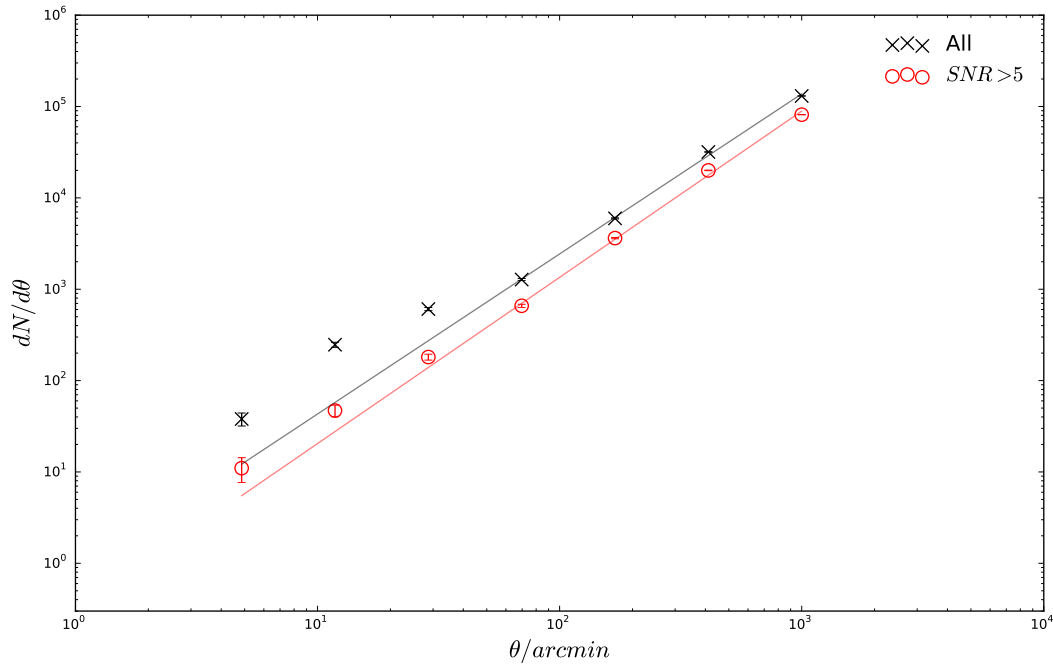


Figure 4.1: The combined number density of source candidates as a function of distance relative to sources brighter than 10Jy. An excess is observed in the distribution of all source candidates (black) out to  $\sim 1^\circ$  due to the presence of side lobes. A  $5\sigma$  cut (red) greatly reduces these contaminants but a significant excess is still present and residual contamination cannot be ruled out.

unsupervised cluster finding and use the results to train a more robust ensemble classifier.

The classification steps are broken down as follows; feature engineering and standardization, dimensionality reduction, initial cluster finding (unsupervised classification), and training an ensemble classifier.

#### 4.1.1 Feature Engineering

Machine learning algorithms require a set of input features that describe the population. Feature engineering is the process of scaling, combining, or otherwise manipulating fundamental parameters to increase the predictive potential of the model. The process of feature engineering and selection is somewhat of a black art driven by domain insight and the question at hand as much as the data available.

Features input to complex models are typically developed and honed through many iterations of trial and error. I ultimately define 9 features based on observable measurements that result in a well-modeled distribution and interpretable results. These are:

1. **Log Flux Density** The log of the weighted mean flux density of all source detections,  $\log_{10}(S/\text{Jy})$ .
2. **Log Signal to Noise** The log of the ratio of the mean flux density to the standard deviation,  $\log_{10}(S/\sigma_S)$ .
3. **Log Signal to Noise of the mean** The log of the ratio of the mean flux density to the standard deviation of the mean,  $\log_{10}(S/\sigma_{\bar{S}})$ .
4. **Number of Detections** The number of snapshots in which a source was detected,  $N_{\text{det}}$ .
5. **Expected Number of Detections** The estimated cumulative probability that the mean source flux density lies above the deconvolution limit in each snapshot,  $N_{\text{exp}}$ .

6. **Detection Rate** We define a weighted measure of the detection rate  $r_{\text{det}} = N_{\text{det}}/\sqrt{N_{\text{exp}}}$ . The square root in the denominator down-weights sources that drift out of the field (i.e. 71 of 71 expected is more reliable than 2 of 2 expected).
7. **Local Density.** The number density of sources within a  $1^\circ$  radius of the source candidate  $\rho_N$  ( $\pi \text{ deg}^2$ ).
8. **Distance to Brightest Neighbour.** Distance to the brightest source within a  $1^\circ$  radius of the source candidate,  $d_{\text{bright}}$  (deg).
9. **Flux Density Ratio to Brightest Neighbour.** The flux density ratio between the source candidate and the brightest neighbour within a  $1^\circ$  radius,  $S/S_{\text{bright}}$ .

The final three features are engineered to differentiate likely side lobe sources that typically occupy regions of high number density ( $\rho_N$ ) in close proximity to a much brighter source ( $d_{\text{bright}}$  and  $S/S_{\text{bright}}$ ). The distributions of all pairs of input features are shown in Figure 4.2.

#### 4.1.2 Dimensionality Reduction

The features were standardized by subtracting the mean and dividing by the standard deviation to put them onto the same scale. Principle component analysis was then used to reduce the parameter space to three dimensions prior to fitting a model to the distribution. Figure 4.3 shows the fraction of the total variance explained by each PCA component. The first component accounts for nearly 50% of the total variance, while the first three account for 83%. We select the first three components, beyond which the variance ratio begins to level out. Reducing the parameter space to three dimensions allows for much simpler model fitting and visualization compared to higher dimensions and the 17% information loss is recovered at a later stage.

Each principle component (denoted  $C_0$ ,  $C_1$ , and  $C_2$ ) is a linear combination of the input features weighted by the set coefficients given in table 4.1. The number density distribution

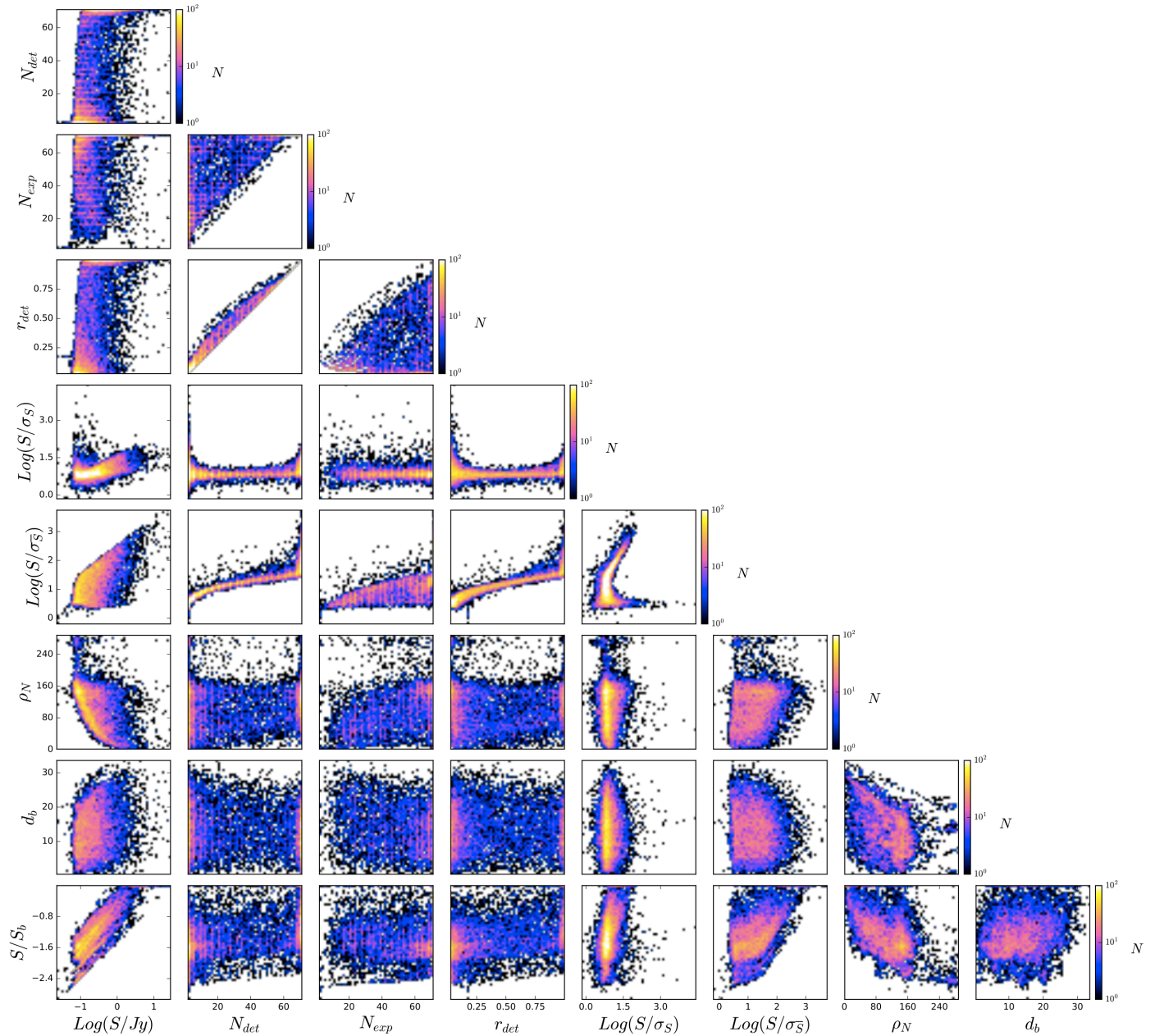


Figure 4.2: The 2D distributions of all 9 input features. The color is log-scaled to highlight structure.

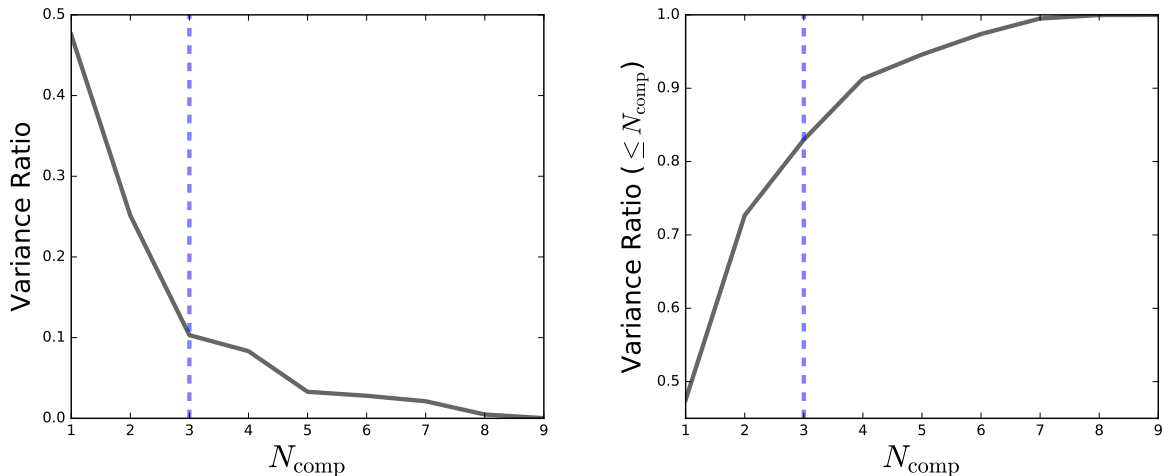


Figure 4.3: The fraction (left) and cumulative fraction (right) of total feature variance explained by each PCA component. The first component explains nearly half of the variance, while the first three explain 83%.

of the principal components is shown in Figure 4.4. The first component  $C_0$  is the most divisive and is most strongly weighted by the detection rate and signal to noise of the mean.

#### 4.1.3 Initial Unsupervised GMM Classification

Next, I fit a Gaussian mixture model (GMM) to the data in the reduced parameter space. Various other methods may be used to make this initial classification (e.g. k-means, nearest neighbors, quadratic discriminant analysis) however I find that the distribution is sufficiently represented by a ten component GMM. At ten components a sharp minimum is observed in the Bayes information Criterion (BIC; Figure 4.5), an indicator of goodness of fit that is sensitive to over-fitting. Every source candidate was labelled according to the GMM class it most probably belonged to. The distribution of classifications in the reduced parameter space is shown in Figure 4.1.3.

While the BIC suggests the GMM is a sufficient approximation of the data distribution, it

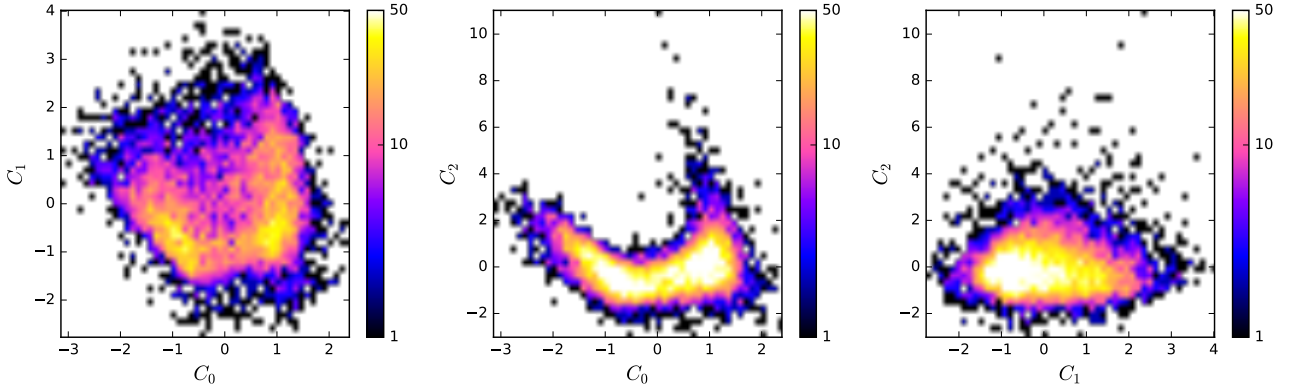


Figure 4.4: The PCA component number density distributions of the source candidates. Each component is a linear combination of the 9 input features with coefficients listed in Table 4.1.

	$S(Jy)$	$S/\sigma_S$	$S/\sigma_{\bar{S}}$	$r_{det}$	$N_{detected}$	$N_{detectable}$	$\rho_N$	$d_{bright}$	$S/S_{bright}$
Scale	$\log_{10}$	$\log_{10}$	$\log_{10}$	$linear$	$linear$	$linear$	$linear$	$linear$	$\log_{10}$
$C_0$	-0.391	0.034	-0.462	-0.458	-0.081	-0.245	-0.479	0.01	-0.355
$C_1$	0.341	-0.544	-0.181	-0.202	-0.482	0.051	-0.15	0.36	0.356
$C_2$	0.058	0.111	-0.235	-0.232	-0.049	0.835	-0.064	-0.408	0.062

Table 4.1: The principle component coefficients of all input 9 features described in §4.1.1.

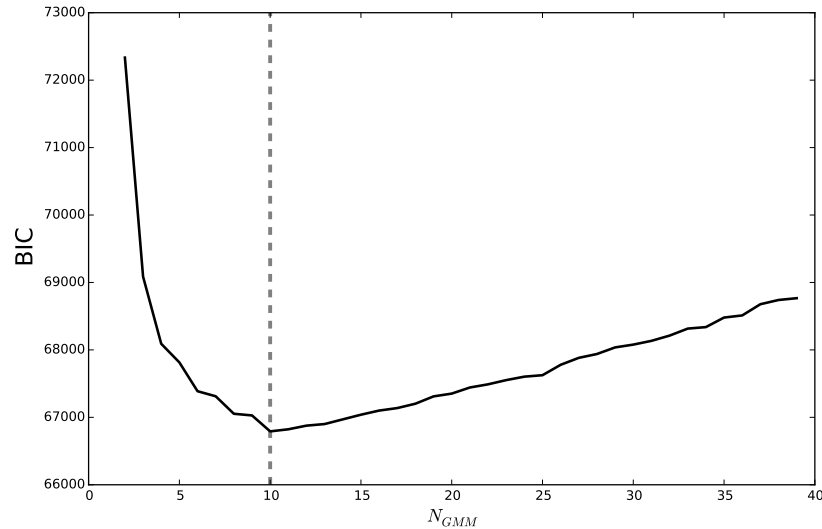


Figure 4.5: The optimum number of Gaussian components was chosen to minimize the Bayes Information Criterion. A strong minimum is found at 10 components.

is naive to assume Gaussianity with full knowledge that input feature distributions are non-Gaussian. The boundaries between classes appear to be forced by the Gaussian assumption rather than true to the underlying distribution. I therefore use the GMM classes as input to train a Decision Tree ensemble classifier.

#### 4.1.4 A Decision Tree Ensemble Classifier

Both Random Forest (RF) and AdaBoost (adaptive boosting) classifiers were tested. RF averages over many decision trees created on sub-samples of the data to build a more accurate classifier [8]. RF is robust against label noise (mis-classifications on the training set). AdaBoost builds a strong classifier from a decision tree by iteratively adjusting feature weights to focus on the outliers [16, 45]. AdaBoost is therefore sensitive to mis-classifications. Because the decision tree favors the most distinguishing features, we used all nine original inputs to minimize information loss.

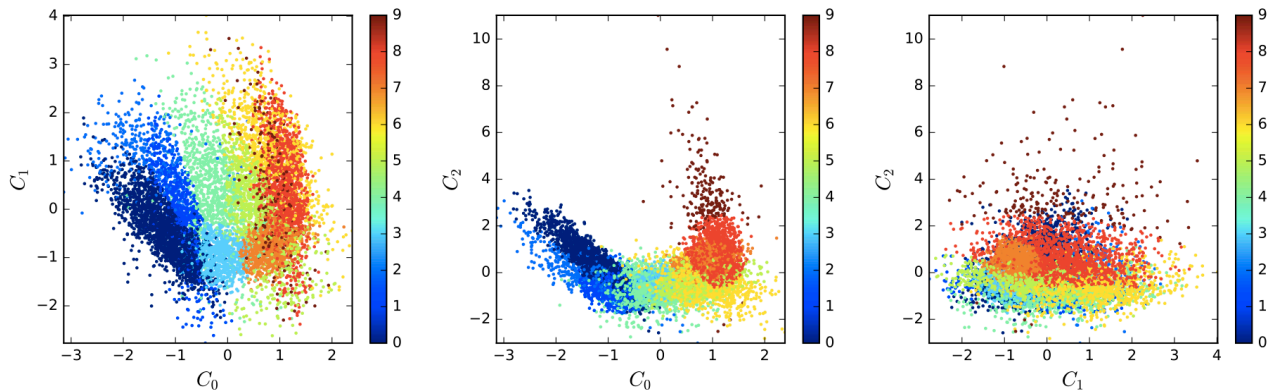


Figure 4.6: The PCA component distributions of the source candidates with color representing the 10 component GMM classification result.

The variable input parameters to the RF and AdaBoost classifiers were not fully optimized, but variations on the number of trees or iterations were explored. The RF classifier tended to change the GMM negligibly or result in odd boundaries. The AdaBoost classifier was subject to over-fitting if too many iterations were allowed. An AdaBoost classifier with 50 iterations reduced the artificial footprint of the Gaussian model assumption while maintaining reasonable agreement.

The source candidates were split randomly 9:1 into training and testing sets. The classifier learns on the training set and is used to predict the test set. This was repeated for 5000 iterations so that each source was independently classified an average of 500 times and assigned to the cluster with the highest mean probability. The resulting classifications relative to the principle components are shown in Figure 4.1.4.

#### 4.1.5 Interpretation

To interpret the reliability classes we look at their average properties, feature distribution, and spatial distribution. A selection of the input features are shown in Figure 4.9 and the

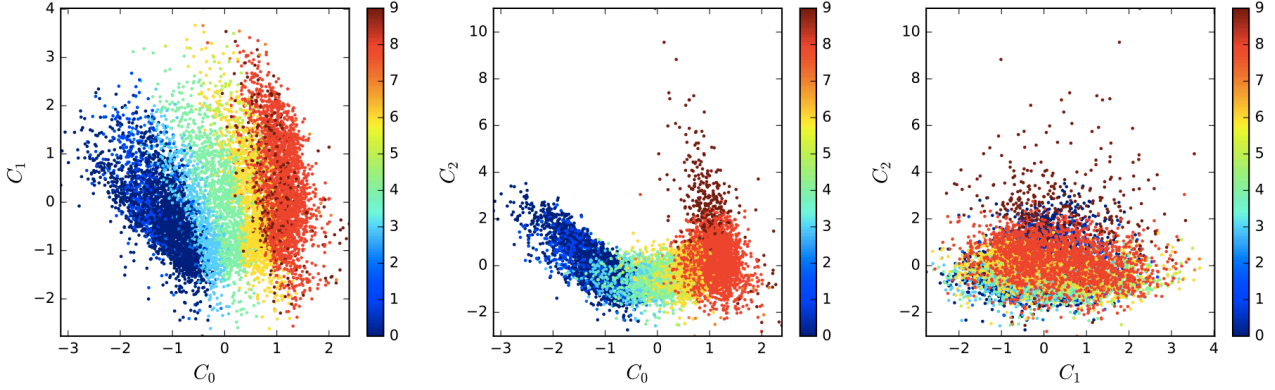


Figure 4.7: The PCA component distributions of the source candidates with color representing the Adaboost revised classification result.

median values for each class are listed in table 4.2. Lower numbered classes tend to be the most reliable overall in terms of detection rate and SNR. Classes  $R0 - R2$  are highly reliable sources observed and detected in nearly all snapshots. Classes  $R3 - R6$  capture fainter sources with reliability decreasing further afield. Classes  $R5$  and  $R6$  appear to identify real sources in high density regions but may be subject to contamination. Sources with  $R_{class} > 6$  may still be real but observed too few times near the detection threshold. False side lobe sources are mostly restricted to classes  $R8$  and  $R9$ . Figure 4.8. Class  $R9$  sources in particular stand out, having artificially high SNR due to a very small number of detections.

## 4.2 KATALOGSS Sample Selection

The described clustering and machine learning methods used to extract and identify reliable source candidates in snapshot data has been termed KATALOGSS; for KDD Astrometry, Trueness, and Apparent Luminosity Of Galaxies in Snapshot Surveys. We abbreviate this to KGS when referring to the KATALOGSS based catalog of sources.

We make an initial cut on the initial set of 9490 source candidates to include only those

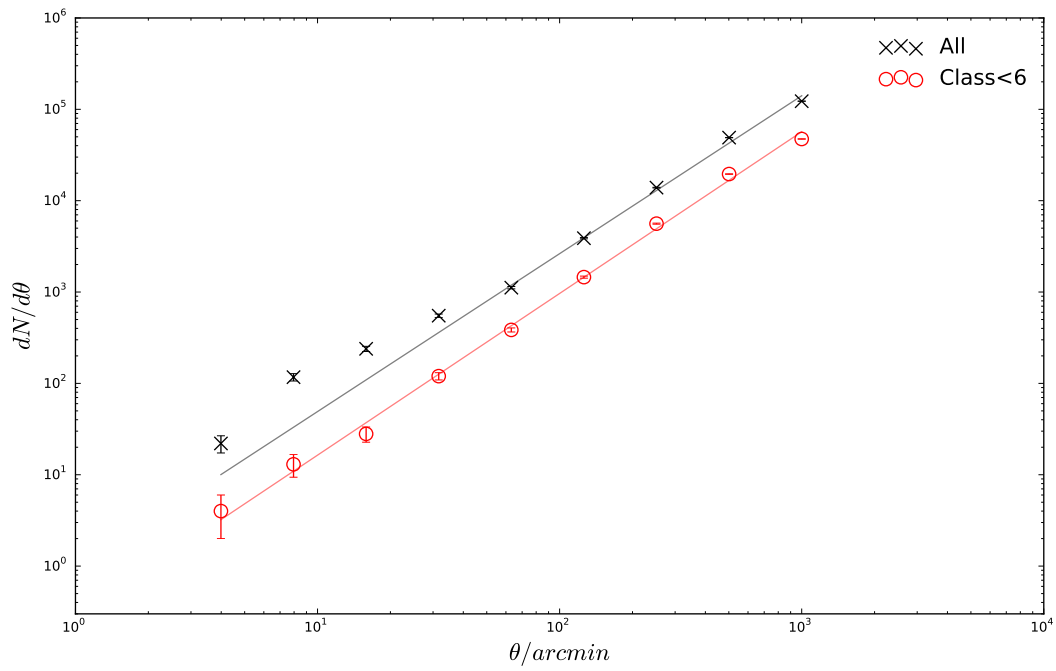


Figure 4.8: The combined number density of source candidates as a function of distance relative to sources brighter than 10Jy. The excess observed in the distribution of all source candidates (black) out to  $\sim 1^\circ$  due to the presence of side lobes disappears by selecting reliability class  $R < 6$ .

$R_{class}$	$N$	$N_{cum}$	$r_{det}$	$N_{det}$	$N_{exp}$	$S(Jy)$	$S/\sigma_S$	$S/\sigma_{\bar{S}}$
0	2353	2353	0.99	70	71	0.35	10.7	57.7
1	416	2769	0.99	70	71	0.93	19.5	131.7
2	172	2941	0.97	69	71	0.27	7.9	39.1
3	794	3735	0.80	56	71	0.20	6.7	26.4
4	1162	4897	0.60	38	60	0.20	6.5	18.4
5	204	5101	0.38	21	38	0.31	6.6	11.2
6	1435	6536	0.32	18	50	0.12	6.3	10.8
7	21	6557	0.11	3	6	0.21	4.5	3.8
8	2630	9187	0.09	4	33	0.14	6.9	4.6
9	303	9490	0.05	2	32	0.11	36.7	3.3

Table 4.2: Source counts ( $N$ ) and median properties for the ten classifications. Lower classes are more reliable while higher classes tend to be fainter or far from field center. Classes  $R8$  and  $R9$  appear to capture sporadic noise and side lobe contaminants, but also faint sources near the detection threshold. The reliability classification is used to inform the final catalog selection.

detected with high confidence ( $S > 5\sigma_S$  and  $S > 5\sigma_{\bar{S}}$ ) or reliability ( $R_{class} < 7$ ; all of which meet the  $S > 5\sigma_{\bar{S}}$  criteria). This reduces the KGS sample to 7466 source candidates. The reliability classification is informative but not infallible, and we can still expect side lobes to potentially contaminate  $R_{class} > 4$ . The plan is to validate source reliability and weed out stubborn contamination through cross-matching in order to assess performance of the classifier as a standalone tool. In the case of a source with no cross-match to another radio catalog, or large uncertainty in the match, we can use the reliability class to definitively include or exclude it from the final catalog.

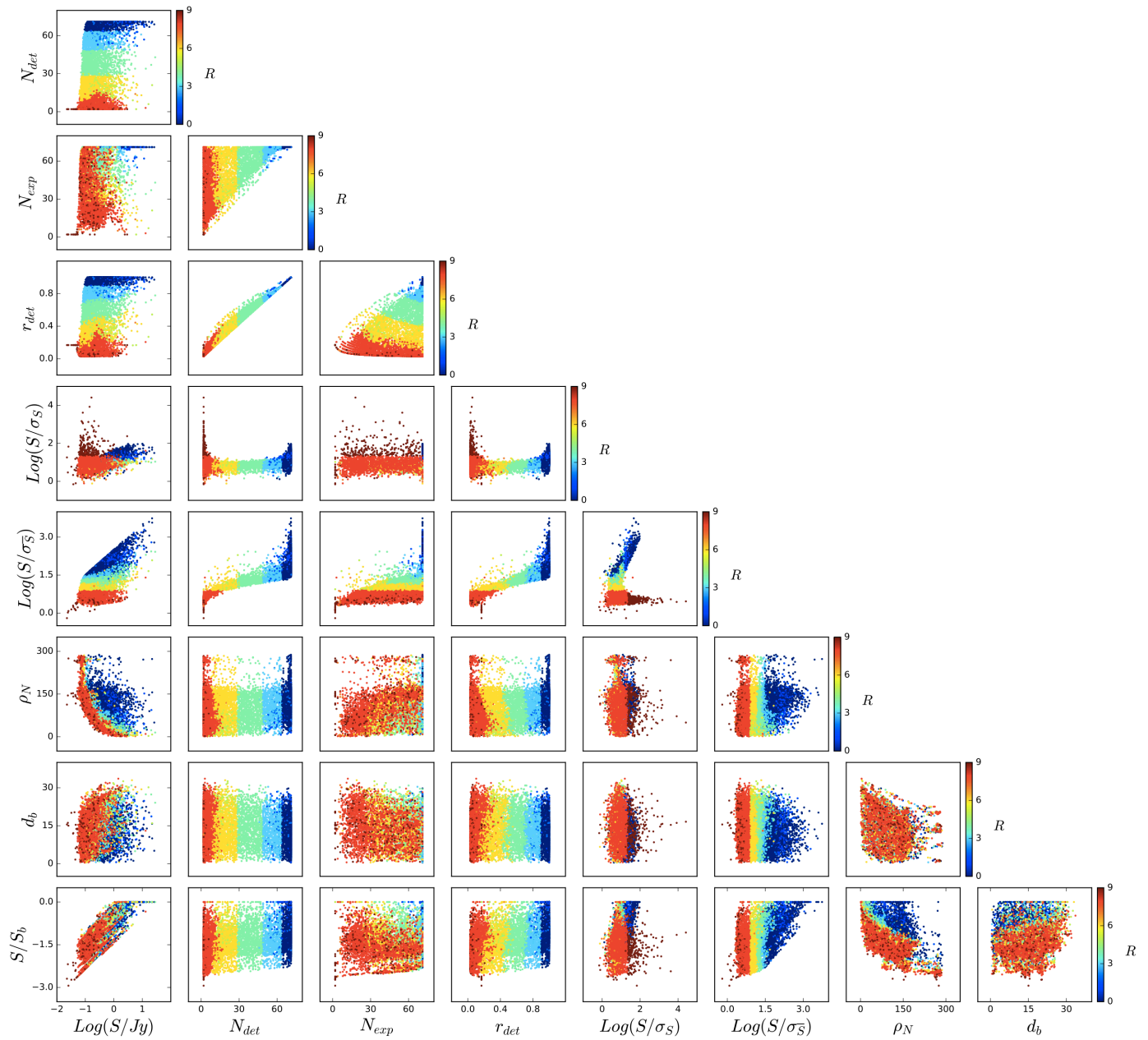


Figure 4.9: The 9 input features cross-plotted and colored by their final classifications ( $R_0 - R_9$ ). Blue are the most reliable and red the least.

## Chapter 5

### RADIO SURVEY CROSS-MATCHING

The 7466 selected KGS source candidates in the EoR0 field cover an area of  $\sim 1400 \text{ deg}^2$ , centered at RA = 0 hr and Dec =  $-27^\circ$ , reaching a depth of  $\sim 60 \text{ mJy}$  at 182 MHz in the center of the field. As discussed in §1, this coverage is not well matched by any other radio survey. The VLSSr and NVSS catalogs cover the northern half of the field ( $\gtrsim -30^\circ$ ) while SUMSS covers the southern half ( $\lesssim -30^\circ$ ). The MRC catalog covers the full sky area but is only complete to 1 Jy at 408 MHz or  $\sim 2 \text{ Jy}$  at 182 MHz. When taken together, however, these surveys can be helpful in assessing the performance of the reliability classifier and identifying contamination. In this chapter we cross match the KGS source candidates to all four comparison surveys using both positional probabilities and goodness of fit to a power law spectral model.

#### **5.1 Positional Update & Matching Algorithm**

Cross matching was performed using the Positional Update and Matching Algorithm (Line et al., in prep.). PUMA uses a combination of positional and spectral information to statistically test whether sources from multiple surveys in close proximity to one-another are true matches.

Initially, PUMA attempts to match sources purely by position. A positional cross match is performed using STILTS [37] by selecting all sources within a radius of 2.3' from the base KGS source. The choice of this radius is somewhat arbitrary (equal to the PSF FWHM) and intentionally liberal. For each initial cross match result, the probability  $P$  that all catalogs are describing the same source is calculated following Budavari et al. (2008, [9]), taking account of the positional errors. When matching  $N_{\text{cat}}$  catalogs, it can be shown that the

Bayes factor is given by

$$P(H|D) = \frac{B P(H)}{1 + B P(H)}. \quad (5.1)$$

The prior  $P(H)$  is defined in terms of the scaled full sky number of sources in each catalog,  $n_i = 4\pi N_i/\Omega_i$ , where  $N_i$  is the number of sources in the catalog and  $\Omega_i$  the catalog survey area. The prior is given by,

$$P(H) = \frac{n_0}{\prod_{i=1}^{N_{\text{cat}}-1} n_i}, \quad (5.2)$$

where  $n_0$  is the scaled source count of the base catalog.

At this point, if a KGS source is matched to only one source from any catalog and  $P > 0.95$  it is accepted without further investigation. This is labelled as an **isolated** match. If  $P < 0.95$ , the SED is fit with a power law model of the form  $\log S \propto \alpha \log \nu$  using linear least squares. The fit is considered good if the reduced Chi-square statistic  $\chi_{\text{red}}^2$  is less than 2. Due to uncertainty on  $\chi_{\text{red}}^2$  given the small number of data points and uncertainty on the errors, the fit is additionally considered good if the residuals  $\epsilon$  are less than 0.1.

$$\epsilon = \frac{1}{N_{\text{cat}}} \sum_{i=1}^{N_{\text{cat}}} \left( \frac{|f_i - \langle f_i \rangle|}{f_i} \right) \quad (5.3)$$

The quoted uncertainties on position cannot always be trusted to fairly represent the true uncertainty, so if  $0.8 < P < 0.95$  we investigate the spectral energy distribution (SED) by fitting a power law spectral model,  $S \propto \nu^\alpha$ , using weighted least squares. If the fit is good, the source is accepted as an **isolated** match. Note that if a match is only found in one other catalog, this fit always passes as there can be no residuals. Steps are taken later in §5.2 to account for any issues that could arise here.

Multiple matches to a single catalog may occur due to confusion at the lower resolution or coincidental false source contamination. In the case where multiple sources from a comparison catalog are matched to a single KGS source, PUMA first attempts to remove any false matches by fitting the spectral model to each possible combination of sources. If one match combination has smaller residuals than all others, as well as having  $P > 0.95$ , it is accepted as the **dominant** match.

If no **dominant** match is found, it is possible that a source is resolved into multiple components in the higher resolution catalogs. This is the common case for radio galaxies and star-forming galaxies with structure that is unresolved by the MWA. To test this, the spectral model is fit to the cumulative flux of the matches at each frequency. If the fit is good, the source is accepted as a **multiple** match.

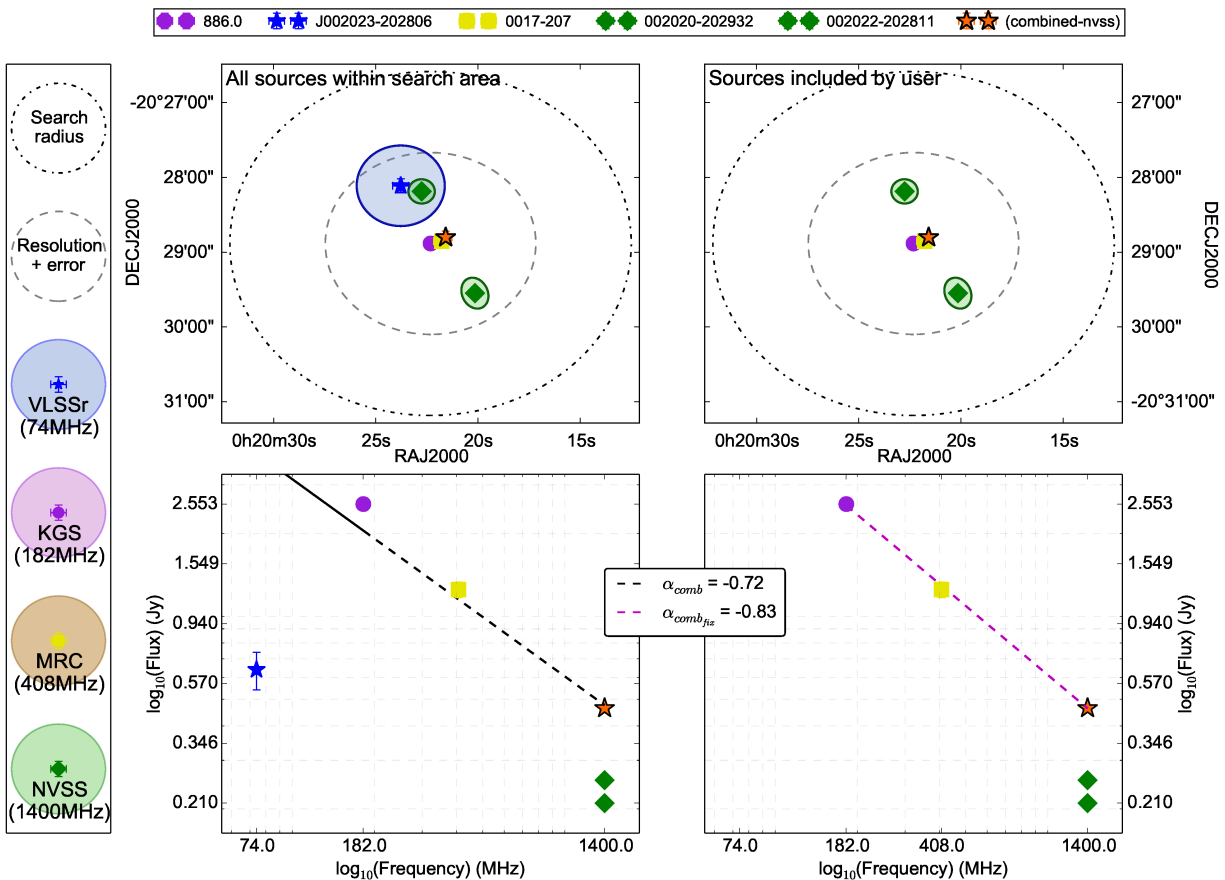
## 5.2 *Flagging & Visual Inspection*

To check the robustness of the PUMA decisions we visually inspected any potentially suspect matches or atypical sources. These include 66 sources automatically flagged by PUMA when a confident match decision could not be made and 45 sources with a STILTS match that was automatically rejected by PUMA. For the sake of reliability, we also double check PUMA accepted matches that we manually flagged as outliers. These include all 900 sources accepted by PUMA as a **multiple** match and 205 defined as having either *a*) spectral index in the 1% tails of the distribution,  $\alpha < -1.46$  or  $\alpha > -0.17$ ; *b*) positional offset from NVSS or SUMSS  $> 3 \sqrt{\sigma_{\text{RA}}^2 + \sigma_{\text{Dec}}^2}$ .

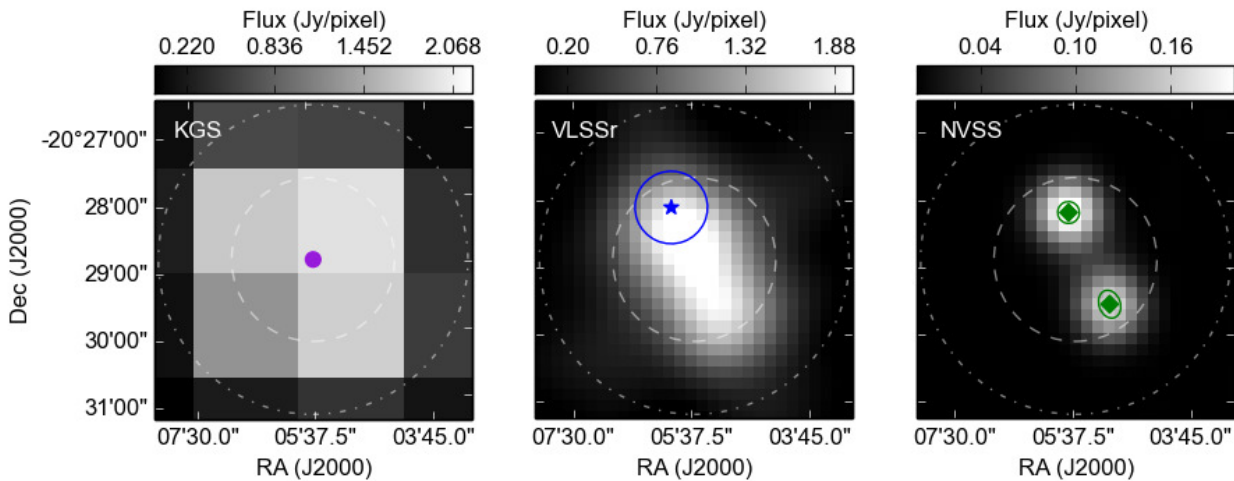
In addition to visualizing the catalog information and PUMA results, we looked at postage stamps of VLSSr, SUMMS, NVSS, and MWA images. Where appropriate, the PUMA decision or match information was modified. This could include removing matched sources that appeared spurious or ignoring a catalog in a **multiple** match that appeared to be missing a source visible in its image. Figure 5.1 shows an example of a source where both the catalog information and images were inspected and used to modify the catalog match.

We removed 24 false **isolated** matches on visual inspection. These were typically sources

Figure 5.1: To investigate flagged matches, catalog information on source position, shape, and flux were plotted (a) and, for complicated sources, postage stamp images were obtained and compared (b). In this example of a **multiple** match, two unresolved sources appear in both the VLSSr and NVSS images, but one is missing from the VLSSr catalog. The centroid position of the NVSS sources agrees well with the KGS and MRC positions. The flux of the NVSS sources are combined and the VLSSr source is excluded.



(a) Example of visualized position (top) and SED (bottom) information for a complicated match before (left) and after (right) manual modification. Ellipses indicate the reported major/minor axis and position angle.



(b) Example postage stamp images inspected for complicated matches. The white dash/dotted circles correspond to the search radius and resolution+error as indicated in (a)

with poor reliability, visually identified as bright side lobes, and coincidentally matched to either a true source or to apparent side lobe contamination in a comparison survey image. We erred on the side of reliability in these decisions. Approximately 10% of multiple matches were able to be deconstructed into two or more components. In §3, we found that a radius of one-quarter the beam width maximized the number of sources detected uniquely in all observations. If multiple source candidates were found by using the tighter clustering radius, these were similarly cross-matched and substituted in manually if the overall match result improved. A total of 90 sources were replaced with 192 counterparts. The catalog includes the column `R_cluster` to indicate the clustering radius used: `h` for a radius of one-half beam-width, and `q` for one-quarter beam-width. The reliability class of replacement sources was predicted independently from the replaced source using the original classifier.

### 5.2.1 Extended Sources

We visually inspected all `multiple` matches and the unusual morphology of NGC 7793 resulted in its identification as a near face-on spiral galaxy in the Sculptor group. A subse-

quent search on the positions of other group members revealed detections of NGC 253 and NGC 55. KGS postage stamp images are shown in Figure 5.2 along with the locations of NVSS detections.

The Sculptor Group is a nearby group of star forming galaxies and the Sculptor Galaxy (NGC 253) is one of the brightest and morphologically complex sources in the EoR0 field. Its extended morphology resulted in failure of the deconvolution loop discussed in §2 for four snapshot observations, and require a more complex treatment in the foreground model as we will explore in §8. The low frequency emission from the Sculptor galaxies will be further investigated in Kapinska et al. (in prep.).

### 5.2.2 *Unmatched Sources*

There are 167 source candidates not matched to another catalog within the initial STILTS search radius, and another 12 STILTS matches that were automatically rejected by PUMA. The majority of these are classified  $R7 - R9$  and visual inspection supports their exclusion as noise or side lobe contaminants. Twenty-five are determined to be real and are included in the final catalog. Of these, 20 are reliably classified  $R0 - R6$ . Five appear to be extended emission blended with a brighter source, but are reliably detected independently of that source.

Among the new detections, we chose to include five faint sources classified  $R8$  after careful consideration. These are interesting and illustrative. Detected in few observations, the mean and standard deviation are poorly constrained. The flux is also likely to be over-estimated due to Eddington bias near the detection threshold. In combination, these effects seem to have resulted in artificially high SNR measures, allowing their inclusion in the source candidate sample. We chose to keep these simply because they appear to be real in the images and were deemed deserving of follow up. Many similar sources did not make the initial candidate selection. In terms of a foreground model, the level of contamination we risk through their inclusion is negligible.

The new source detections are likely be diffuse and/or ultra steep spectrum (USS) sources

associated with galaxy clusters or high redshift radio galaxies (HzRGs). As these are likely to be some of the most scientifically interesting sources in the field, their properties and potential associations in other wavebands are explored in §7.

### 5.3 Results

Table 5.1 details the number of each type of match found by PUMA. The majority of sources (87%) are matched to a single counterpart in other catalogs (isolated or **dominant**) with a 98.6% automatic success rate. When confusion occurred (**multiple** matches), PUMA chose the proper match combination in 84% of cases. Most modifications were required for complex and extended sources.

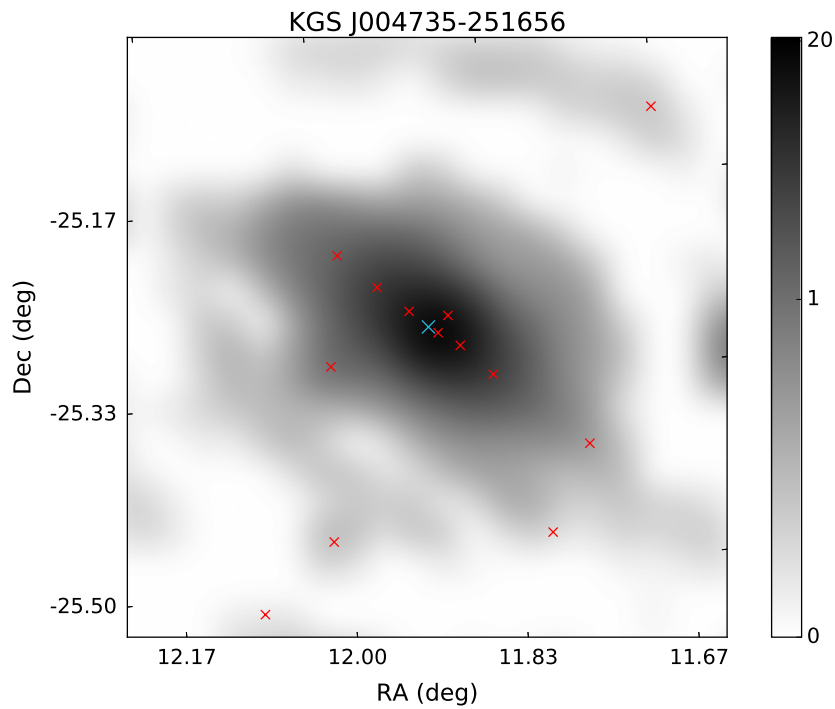
The flux and uncertainty of matched sources are included in the KGS catalog along with the measured broad band spectral index (SI). The position of matches to the NVSS or SUMSS catalogs are also included for reference (a flux weighted mean position is reported for **multiple** matches) and used to assess astrometric precision. The catalog includes an **Inspected** flag signifying the level of inspection: 0 if none; 1 if the catalog data and images were inspected; or 2 if match was modified.

By specifically targeting and inspecting outliers we were able to correct match mistakes and further clean the sample of contamination, but we also made note of several very real and particularly interesting sources that have been marked for follow-up.

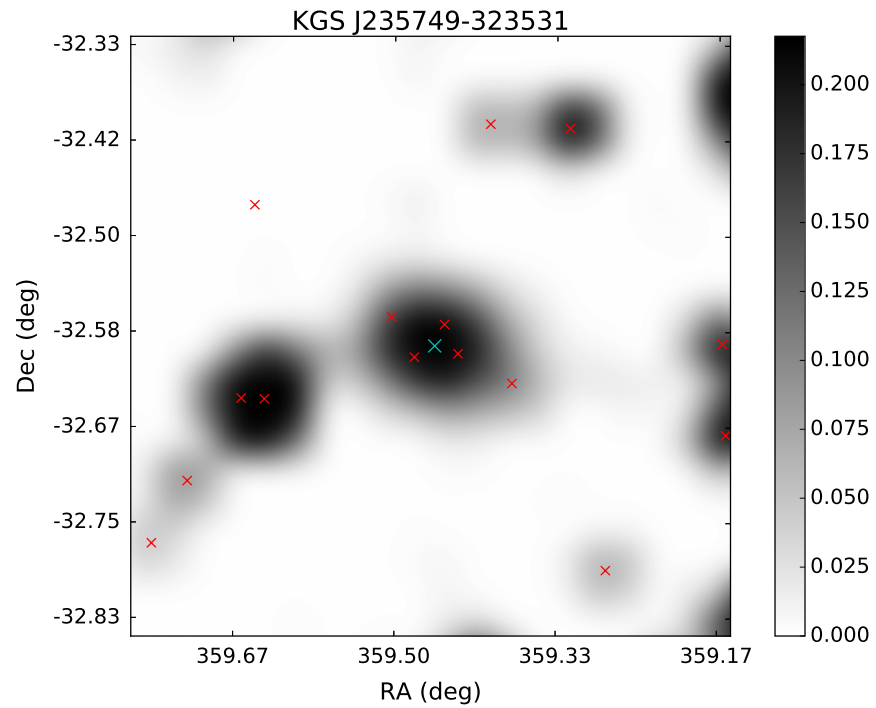
<b>Match result</b>	<b>Count</b>	<b>(%)</b>	<b>Modified</b>	<b>(%)</b>
isolated	6119	(82.8)	75	(1.2)
dominant	310	(4.2)	11	(3.5)
multiple	940	(12.7)	153	(16.3)
none	25	(0.34)		
<b>Total</b>	<b>7394</b>	<b>(100)</b>	<b>239</b>	<b>(3.2)</b>

Table 5.1: The total number and percent of catalog sources in each PUMA decision category, and the number and percent of each category for which the automatic decision was modified. The majority of sources (87%) are matched to a single counterpart in other catalogs (isolated or dominant) with a 98.6% automatic success rate. When confusion occurred (multiple matches), PUMA chose the proper match combination in 84% of cases. Most modifications were required for complex and extended sources.

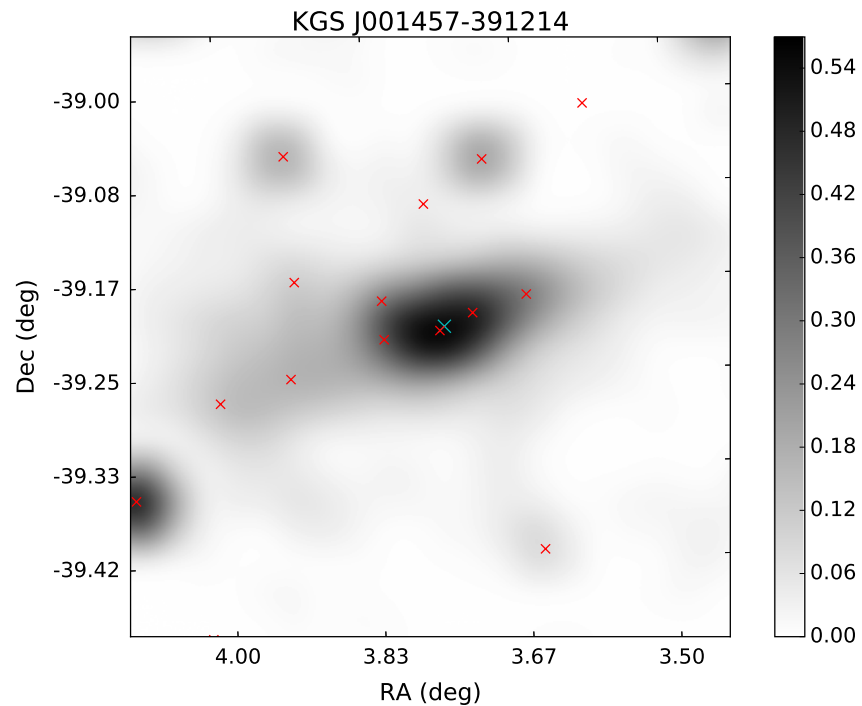
Figure 5.2: Images of the three identified Sculptor Group galaxies. Cyan markers indicate the KGS mean position. Red markers indicate the locations of NVSS sources. NGC 253 (a) is the sixth brightest source in the catalog and is shown on a log scale to emphasize the extended emission. Images are 20x20 pixels, smoothed with a cubic interpolation.



(a) NGC 253 is the sixth brightest source in the catalog and is shown on a log scale to emphasize the extended emission. A more complex model will be required for this source.



(b) NGC 7793 is one of the brightest galaxies within the Sculptor Group and hosts a micro-quasar in its outer disk.



(c) Star-forming galaxy NGC 55 is commonly associated with the Sculptor group though distance measurements suggest it is not gravitationally bound.

## Chapter 6

### THE EOR0 FOREGROUND CATALOG

The final catalog consists of 7394 sources, of which 7369 have confident matches to one or more of the comparison catalogs between 74 MHz and 1400MHz. Since nearly all sources are matched, we can explore the relative flux scale reliability of the catalog, completeness, spectral index distribution, and astrometric accuracy. In this way, we are able to robustly identify and correct for any systematic biases and gain a more thorough understanding of the intrinsic low frequency behavior of our source sample.

#### 6.1 Flux Scale

To investigate how each matched catalog contributed to the spectral index fit, the flux density at each frequency was extrapolated using the fitted parameters. To be sure of a true matched SED, only `isolated` sources were used. A ratio between the reported catalog flux density and extrapolated flux density was then calculated, as shown in Figure 6.1.

On average there is no significant bias in the distributions, however it is interesting to note the width and skewness in these distributions. The VLSSr and KGS skew somewhat low and the MRC skews somewhat high. For sources detected in more than two catalogs, the spectral index fit used the quoted flux densities and uncertainties of the comparison catalogs. NVSS and SUMSS have lower uncertainties than the lower frequency catalogs and VLSSr has the largest. This is evident in the spread of the flux ratio distributions; NVSS has a tight distribution centered at one, whereas VLSSr has a broader distribution. Clearly NVSS is being fit preferentially over the other catalogs. Although the median values are all consistent with unity, systematic effects are clearly present.

It is difficult to distinguish catalog flux biases from intrinsic spectral curvature effects,

but it appears that only a sub-population of sources are affected rather than there being an overall shift in the distribution. A systematic under or over estimation due to the original flux scaling or calibration could account for this. No attempt has been made to match flux scales across catalogs since all are tied to or derived from the Baars scale [2]. Late in the analysis however, it was realized that unlike the original VLSS, the VLSSr is tied to the RCB scale [33]. The re-scaling negligibly impacts the resulting SED fit and overall SI distribution due to the low weighting of the VLSSr data points and the fact that all but three VLSSr matches are also matched to the NVSS. Nonetheless, in the following sections we divide all VLSSr flux densities by a factor of 1.1 to place them on the Baars scale and increase the flux density uncertainty by 5% following [22].

Alternatively, if a significant portion of SEDs display some degree of intrinsic positive curvature, the ratio distributions would be impacted predictably. The lowest frequency flux densities (VLSSr and KGS) would typically be overestimated by a power-law fit, while the central frequencies (MRC and SUMSS) would be underestimated. This is consistent with the flux ratio distributions observed, however a more careful treatment is needed to conclude intrinsic curvature over systematic effects. We do not find evidence for flux scale bias relative to the comparison catalogs.

## **6.2 Eddington Bias**

Source flux densities near the detection threshold will be systematically high due to Eddington bias. At low apparent flux densities, statistical fluctuations below the sensitivity limit go undetected resulting in an overestimation of the true mean [14]. To estimate a correction for this effect, we numerically solve for the true flux density that would most likely result in the observed average over all detections.

The probability density function of a source detection is assumed to be Gaussian, centered on the true flux density and with standard deviation equal to the background rms. The expectation value of the measurement is then the mean of the probability density function above the detection threshold. For each observation, the detection threshold is estimated as

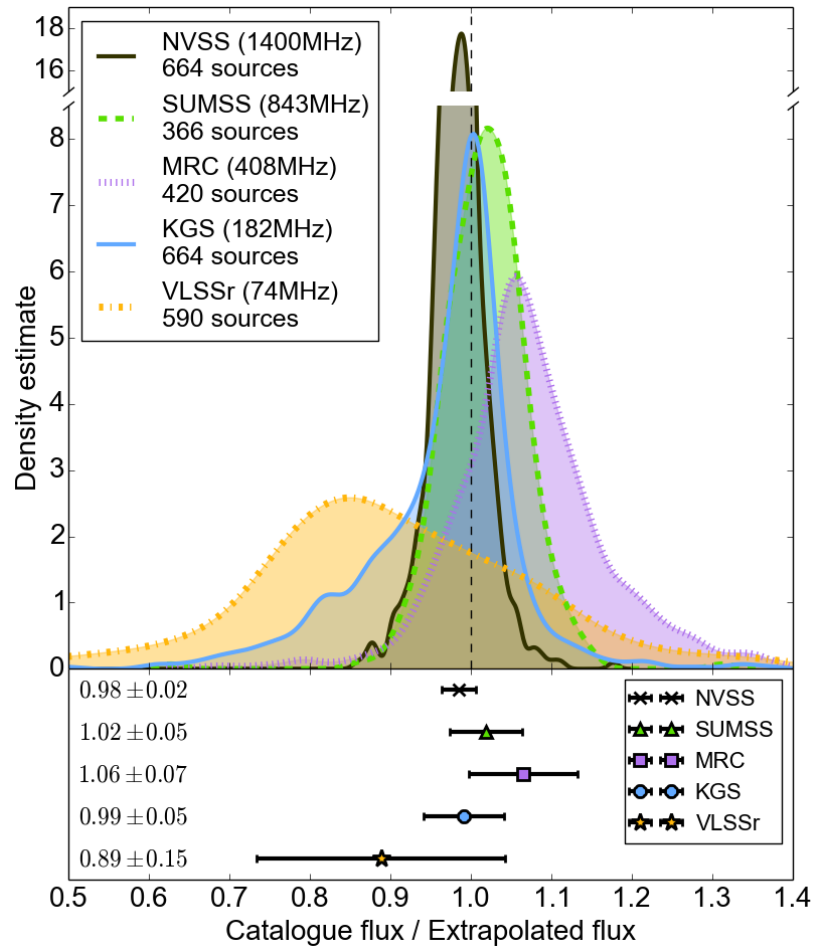


Figure 6.1: The ratio between observed flux density and extrapolated flux density from a fit to the SED is shown for every time a catalog appeared in a match with at least two other catalogs for *isolated* sources. The upper panel shows a univariate kernel density estimation of each distribution (note broken y axis due to the sharp peak in the NVSS ratio distribution), while the lower panel shows the median and median absolute deviation of each distribution. The KGS spectral index agrees very well with no indication of flux bias on average.

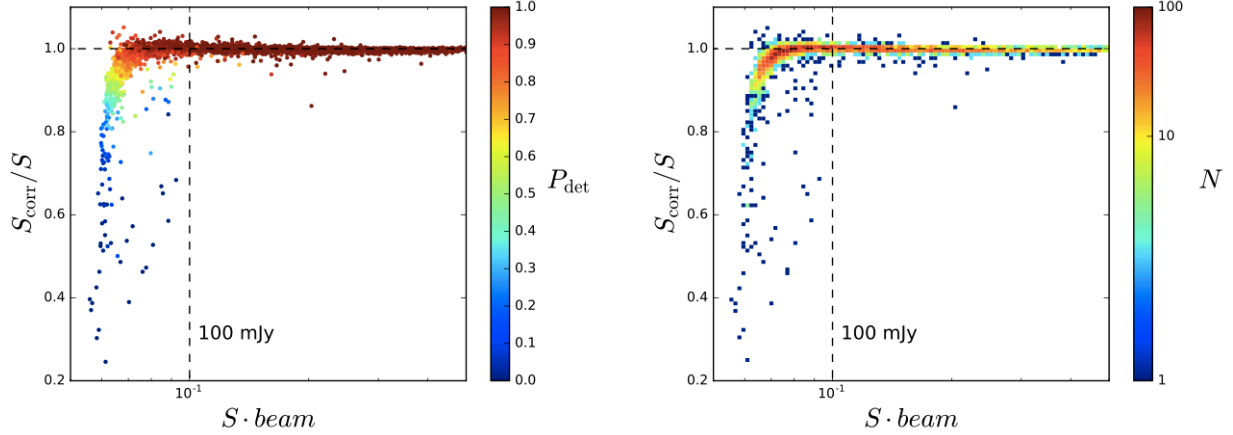


Figure 6.2: The ratio of the corrected flux density to the measured flux density ( $S_{\text{corr}}/S$ ) versus the apparent flux density ( $S \cdot \text{beam}$ ). In the left plot, points are colored by the overall probability of detection of a source assuming the true flux density  $S_{\text{corr}}$ . The right plot shows the number density.

a function of position in the beam and the background rms is estimated from the residual images within a 20x20 pixel box around the source position.

The reported source flux density is the weighted mean of the snapshot detections. Using this as the initial guess for the true flux density of the source, we find the weighted mean of the expectation values for all detections as described above, and numerically solve for the true flux density that minimizes the difference between this expected mean and that observed. Figure 6.2 shows the estimated correction versus the apparent flux density ( $S \cdot \text{beam}$ ) for all `isolated` sources. We find the bias affects sources with apparent flux density below about 100 mJy, but the vast majority of sources are changed by less than 10%.

To gauge the accuracy of the true flux density estimates, we looked at the 182-1400 MHz spectral index distribution for `isolated` sources before and after correction. Above  $S \cdot \text{beam} = 100$  mJy (Figure 6.3, right), there is no significant difference in the distributions before and after correction. Below this threshold, Eddington bias is evident in the shift of the spectral

index distribution toward more negative values (Figure 6.3, left). After correction, the median value agrees well with that of the unbiased distribution. The remaining discrepancy may be explained by selection effects and larger uncertainties at low apparent flux density, as well as imperfect modeling of the correction.

Of the 2548 sources that would be corrected, the difference in flux density exceeds the standard error  $\sigma_S$  for only 177 sources (50 exceed  $3\sigma_S$ ). For this reason, Eddington bias is not a major concern. The overall median SI change is small, from  $-0.850$  to  $-0.843$ . However, the bias is significant for individual sources within the catalog, particularly at low apparent flux density. It is important to note that the validity of the correction factor for any source is contingent on there being a large enough number of detections that the mean is sufficiently constrained. It is therefore the least reliable for the most affected sources.

The catalog contains a column `EB_corr` that may be multiplied by the flux density to approximate a correction for Eddington Bias. The correction factor is 1 by default for a source if  $(S \cdot beam) > 150$  mJy or if it is a `multiple` match (i.e. not a point source). We recommended adding the reported flux density uncertainty in quadrature with the absolute difference between the original and corrected flux density values. In the following sections we use the corrected flux

### 6.3 Completeness

In order to assess completeness, we compare source counts to the predicted source counts of the NVSS and VLSSr surveys projected to 182 MHz (Figures 6.3–6.5). Counts are considered only within the overlapping survey areas and each catalog is projected using the median two-point spectral indexes  $\alpha_{74}^{182} = -0.60$  and  $\alpha_{182}^{1400} = -0.85$  for all matched `isolated` sources. Again, we have used flux densities corrected for Eddington Bias so as not to bias the estimated limit.

Because the sensitivity and thus detection threshold goes as  $1/beam$ , the completeness falls off steadily below  $\sim 1$  Jy compared to the NVSS. Within the half-power point, we find the catalog is complete to approximately 80 mJy. Source counts appear to be comparable to

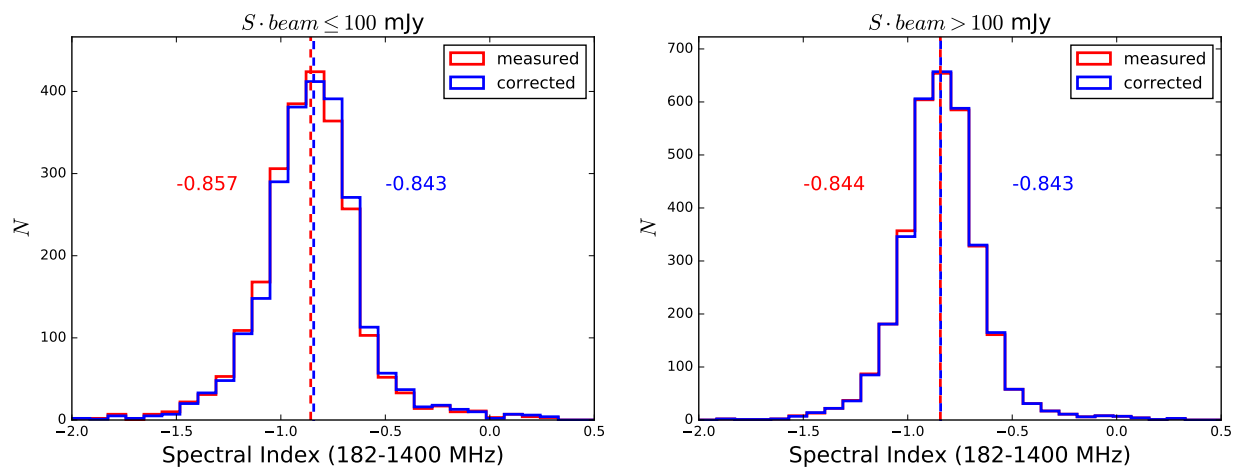


Figure 6.3: The two-point spectral index distributions of *isolated* sources with matches to NVSS at 1400 MHz before (blue) and after (red) correction for Eddington bias. Sources with apparent flux density  $S \cdot beam < 100$  mJy (left) tend to be the most affected. After correction, the median spectral index agrees well with sources above this threshold (right) for which the estimated correction is negligible.

the VLSSr above 200 mJy, below which KGS sources are likely to go undetected at 74 MHz within the overlapping footprint.

#### 6.4 Spectral Index Distribution

The spectral index (SI) distribution found by the match performed in §5 is shown in Figure 6.6. We find an overall median of -0.85, however it is difficult to compare spectral index measurements across different frequency ranges in surveys with differing sensitivity limits and flux scales.

In Figure 6.7, we show the SI distribution for all two-point SI measurements among matches to `isolated` KGS sources. The VLSSr catalog has been corrected to the Baars scale and the KGS flux densities have been corrected for Eddington bias. Further, we select only sources detected at an average beam power greater than 0.5 and flux density  $S > 200$  mJy. The two-point median spectral index is seen to range considerably, from -0.6 to -0.95, with a trend toward steeper spectra at higher frequencies. The lowest frequency measurements  $\alpha_{74}^{182}$  and  $\alpha_{182}^{408}$  give an average of -0.70 at 182 MHz.

By fitting a second order polynomial to the subset of 883 `isolated` sources in 3 or more catalogs, we predict a median 182 MHz spectral index of -0.71 with an interquartile range between -0.88 and -0.53 (the mean and standard deviation are  $-0.71 \pm 0.32$ ).

These results are consistent with Offringa et al. (in prep), who directly measure the sub-band (132–198 MHz) spectral index for a highly comparable set of sources in the center of the MWA EoR0 field. They find a median of -0.70 at 168 MHz (mean of  $-0.687 \pm 0.275$ ). For comparison, the Low Frequency Array (LOFAR) MSSS MVF survey [18] finds median values of  $\alpha_{30}^{158} = -0.66$  and  $\alpha_{119}^{158} = -0.77$  (mean values of -0.60 and -0.70 respectively) among 628 sources with  $S_{150} > 200$  mJy.

#### 6.5 Astrometry

Catalog matching allows us to approximate the best astrometric position of a source based on higher frequency, higher resolution counterparts. Of the 7369 matched sources, all but

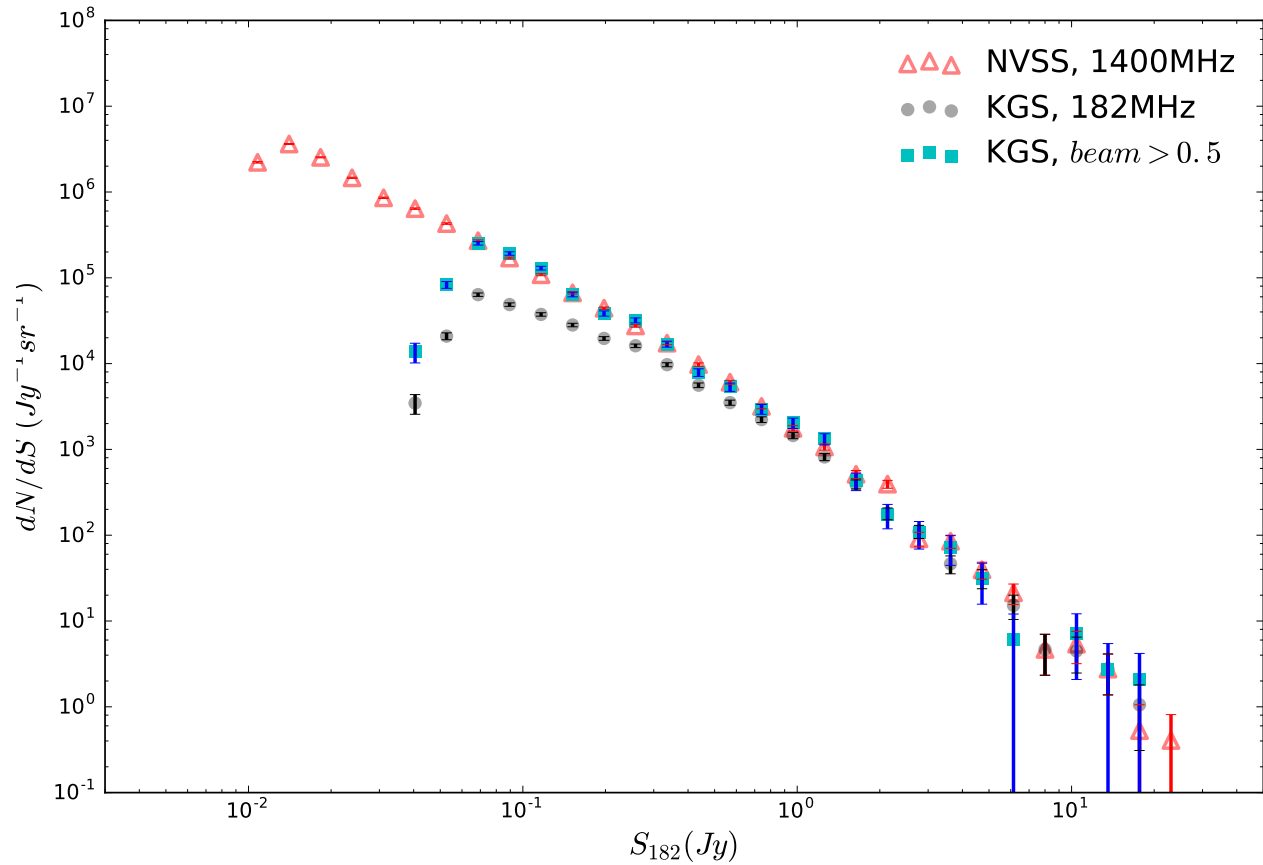


Figure 6.4: Differential source counts of the catalog (grey points) compared to NVSS (red open triangles). Counts are made within the overlapping footprint and the flux density values are projected to 182 MHz using the median 2-point spectral index of all **isolated** matches. NVSS is complete far below the KGS detection limit and is here used as a basis for comparison. KGS completeness falls off as  $1/\text{beam}$  below 1 Jy, however within half-beam (blue squares) the source counts are comparable to NVSS to approximately 80 mJy.

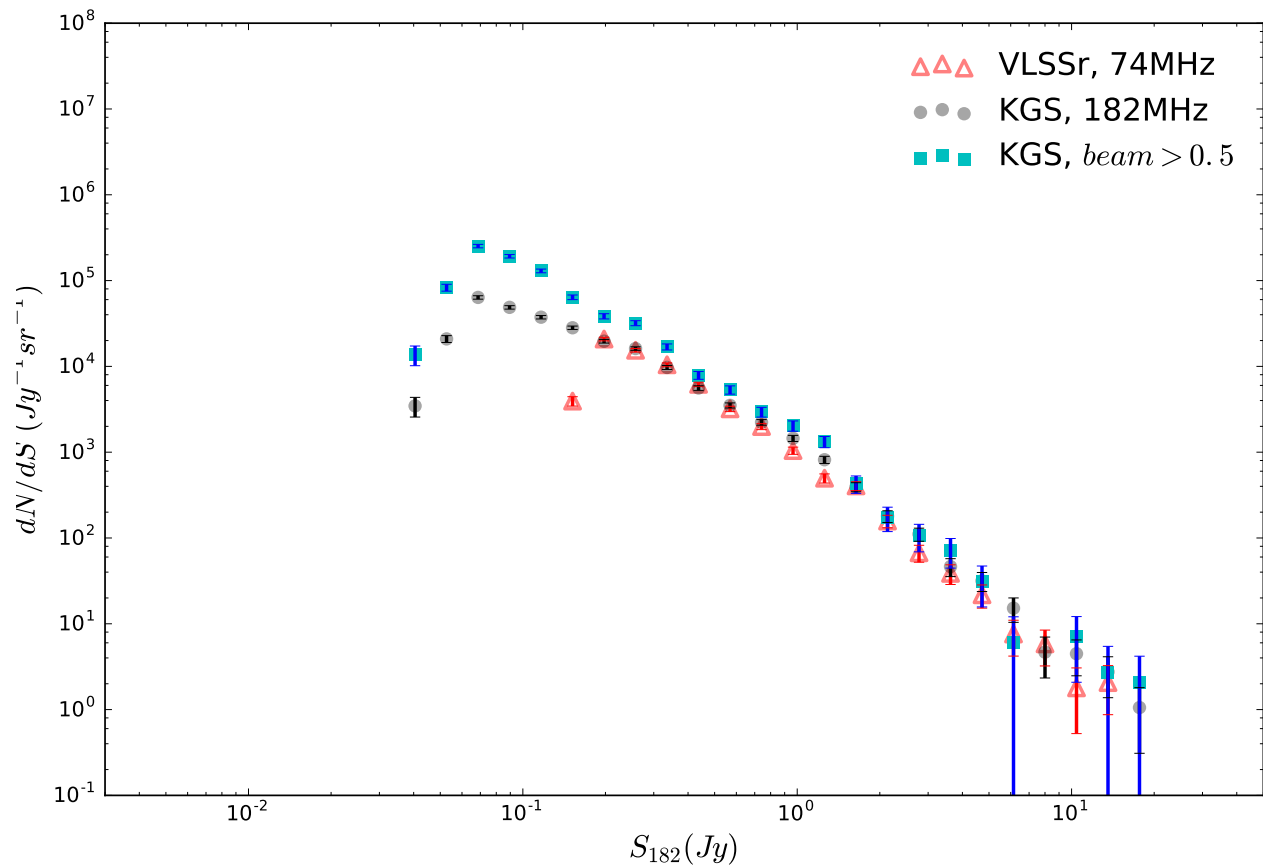


Figure 6.5: Differential source counts of the catalog (grey points) compared to VLSSr (red open triangles). Counts are made within the overlapping footprint and the flux density values are projected to 182 MHz using the median 2-point spectral index of all *isolated* matches. The two surveys are comparable to  $S_{182} = 200 \text{ mJy}$  below which KGS sources are more likely to go undetected in the VLSSr.

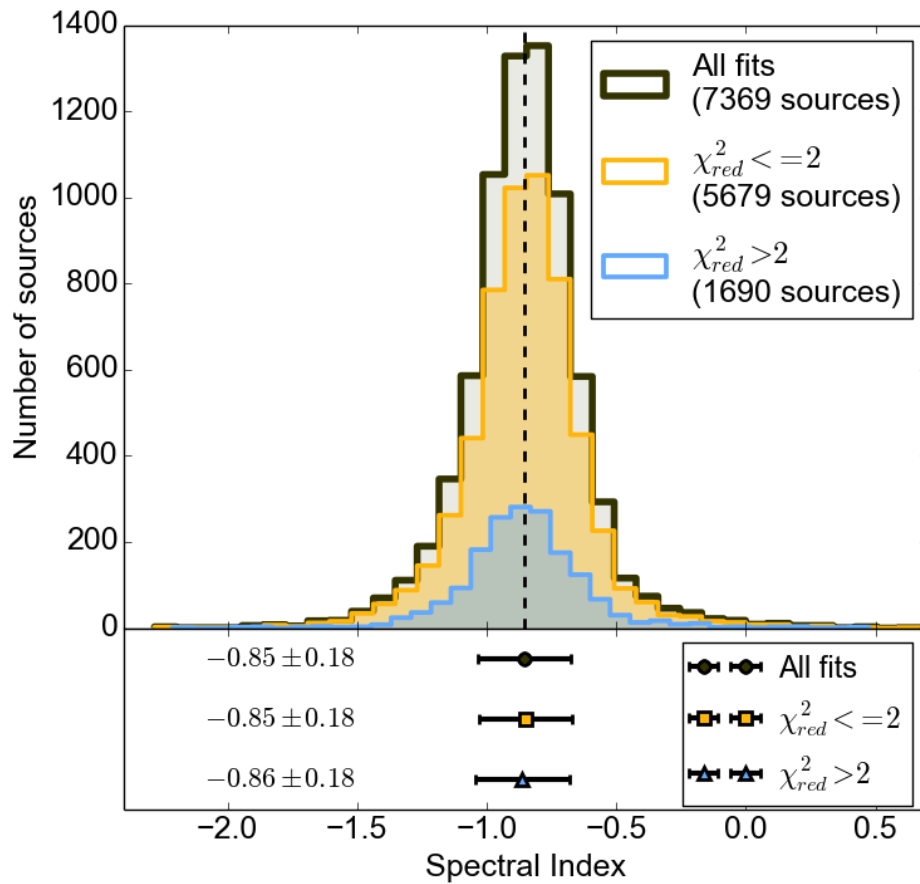


Figure 6.6: The SI distribution derived by matching to VLSSr, MRC, SUMSS and NVSS. The full sample is shown in black, with good spectral fits ( $\chi_{red}^2 \leq 2$ ) shown in gold and poor spectral fits ( $\chi_{red}^2 > 2$ ) shown in blue. The mean and standard deviation of SI distributions are shown in the lower panel. There is no evidence for biases in the SI distribution based on this or other cuts explored.

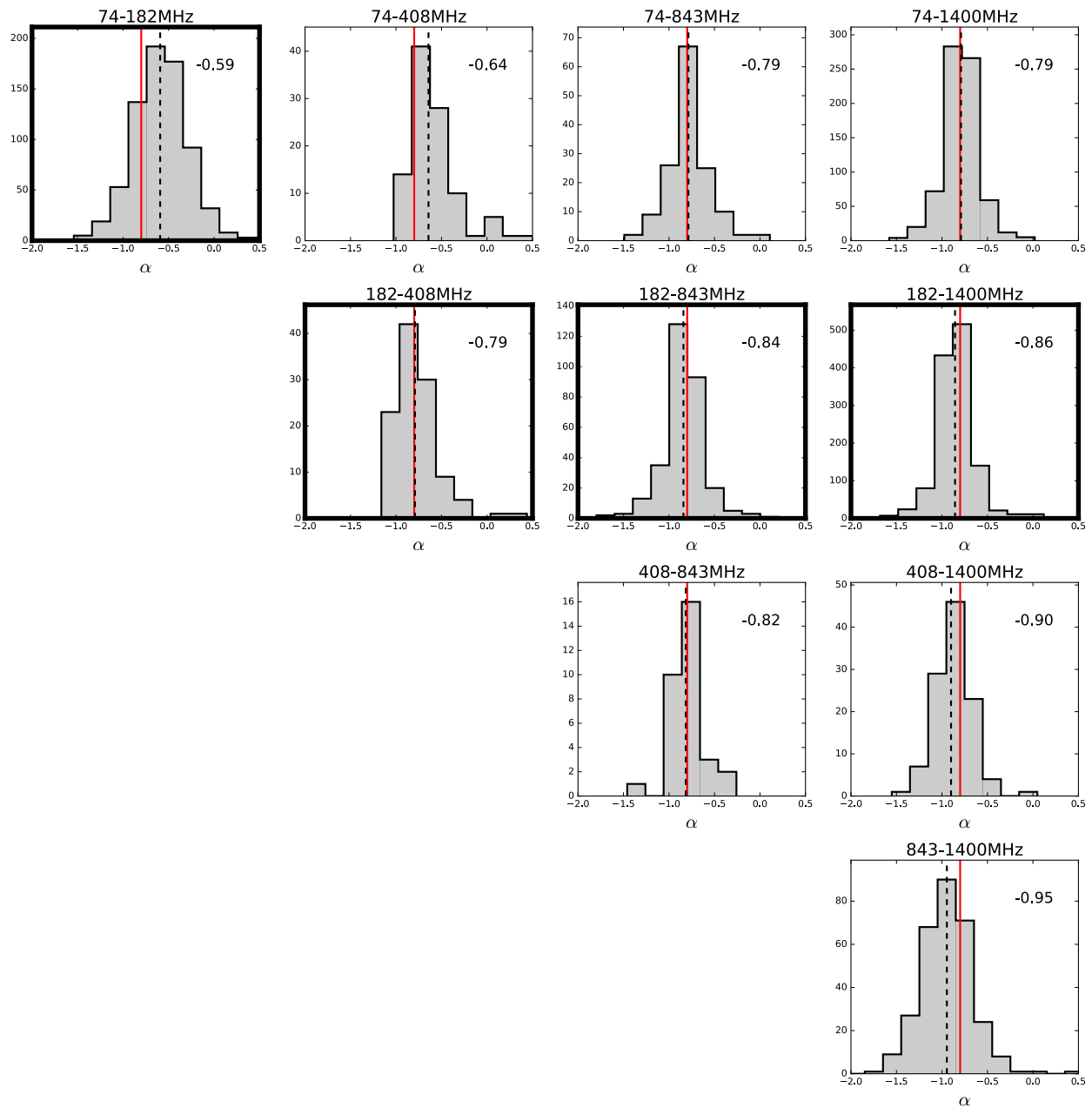


Figure 6.7: The two-point SI distributions for all catalog matches to isolated KGS sources with  $S \cdot \text{beam} > 200$  mJy and  $\text{beam} > 0.5$ . Bold axes indicate distributions using the KGS 182 MHz flux density. The median values are marked by dashed black lines and the red lines mark -0.8 for reference. The median becomes increasingly negative values toward higher frequencies.

three include a match to either the NVSS or SUMSS catalogs.

Figure 6.8(a) shows the distribution of offset distances in RA and Dec from the NVSS or SUMSS match to isolated sources. The median offset is  $\sim 10''$  in either dimension. While this is less than the median errors  $\sigma_{\text{RA}} = 19''$  and  $\sigma_{\text{Dec}} = 15''$ , a north-eastward systematic bias is clearly apparent. This is illustrated by a vector field in Figure 6.8(b).

The source of this offset can be traced to errors in the MWACS catalog used for calibration. Considering MWACS isolated matches on position, the median offsets for MWACS within half-beam are  $\Delta\text{RA} = -8''$  and  $\Delta\text{Dec} = -10''$ . KGS offsets outside of half-beam are much larger in RA. We expect the root of the problem is with the calibration catalog, compounded by wide-field mapping projection errors far from beam-center.

The offsets  $\Delta\text{RA}$  and  $\Delta\text{Dec}$  are found to be well modeled by a second order polynomial as a function of (RA, Dec) position. The models are shown in Figure 6.9 and the polynomial coefficients are given in Table 6.1. The models are fit to `isolated` sources and are used to approximate a positional correction. The distribution of offsets after correction are shown in the bottom panels of Figure 6.8. The bias is dramatically reduced and the the median offset is  $< 1''$  compared to  $\sim 10''$  in either dimension. The median absolute offsets are  $|\Delta\text{RA}| = 4''$  and  $|\Delta\text{Dec}| = 3''$ . In the catalog, we report the bias-corrected KGS position as well as the matching catalog (NVSS or SUMSS) position for comparison.

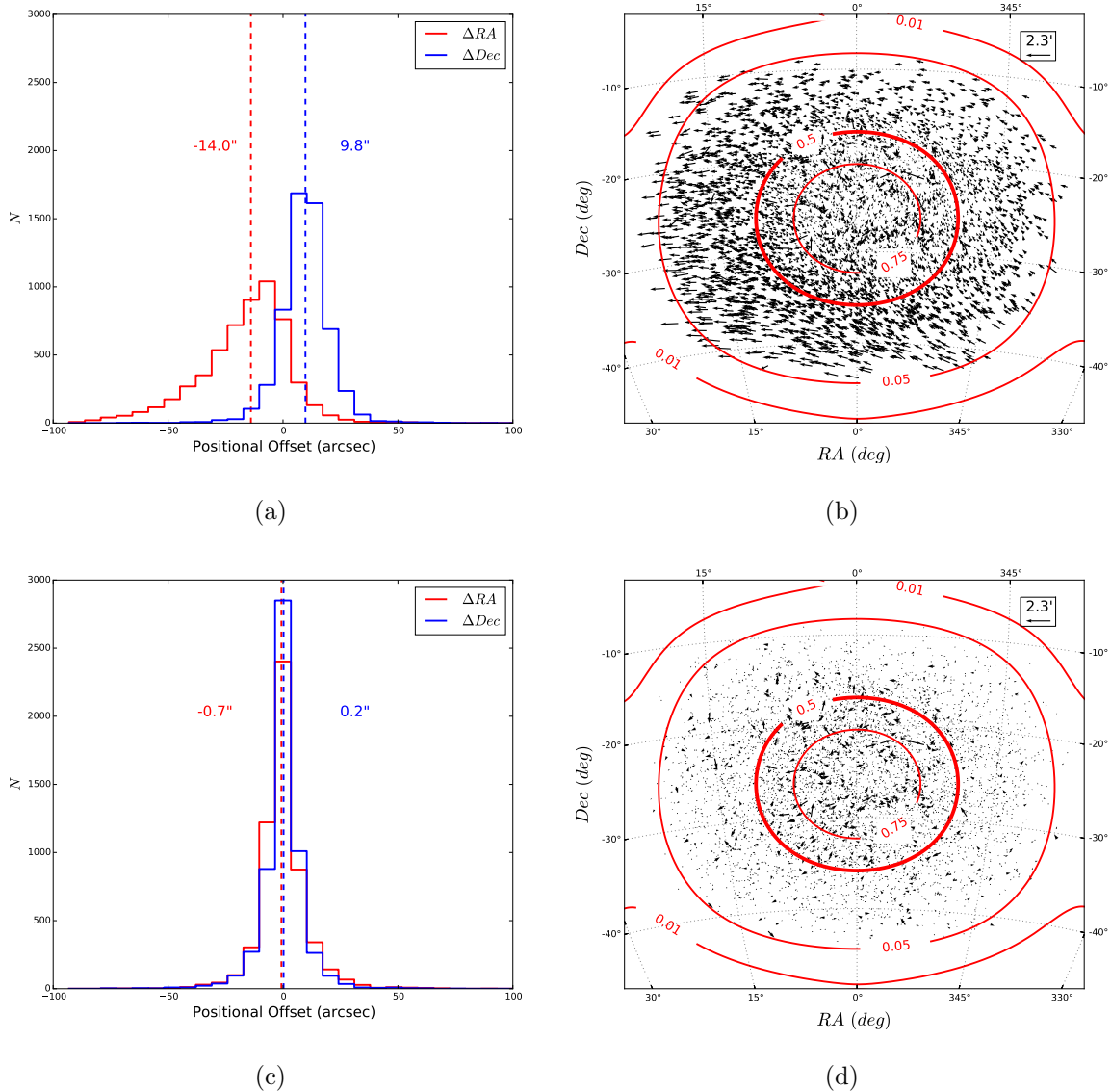


Figure 6.8: The distribution of positional offsets compared to either NVSS or SUMSS counterparts for `isolated` sources before (top ) and after (bottom) bias correction. The vector field arrows are scaled relative to the beam width of 2.3', indicated in the top right. The average analytic MWA beam power across all snapshots is contoured for reference. A systematic north-eastward directional bias is apparent in (a) and (b). After correction, the bias is dramatically reduced. The median offset in (c) is  $< 1''$  either dimension and no apparent bias is evident in (d).

	$\alpha^0$	$\alpha^1$	$\alpha^2$	$\delta^0$	$\delta^1$	$\delta^2$
$\Delta \alpha$	-1.416e4	78.52	-0.1098	-14.12	-0.1961	1.205e-2
$\Delta \delta$	-512.5	3.196	-4.919 e-3	-1.708	4.872e-3	4.126e-3

Table 6.1: The polynomial coefficients of the fit to the position bias. The input coordinates,  $\alpha$  and  $\delta$ , are the right ascension and declination in degrees, where  $\alpha$  is adjusted to the range  $-30^\circ \lesssim \alpha \lesssim 30^\circ$  for continuity. The offset magnitudes,  $\Delta\alpha$  and  $\Delta\delta$  are in arcseconds, such that the corrected positions are  $\alpha + \Delta\alpha/3600$  and  $\delta + \Delta\delta/3600$ . The modeled and corrected positions are shown in Figures 6.8–6.9.

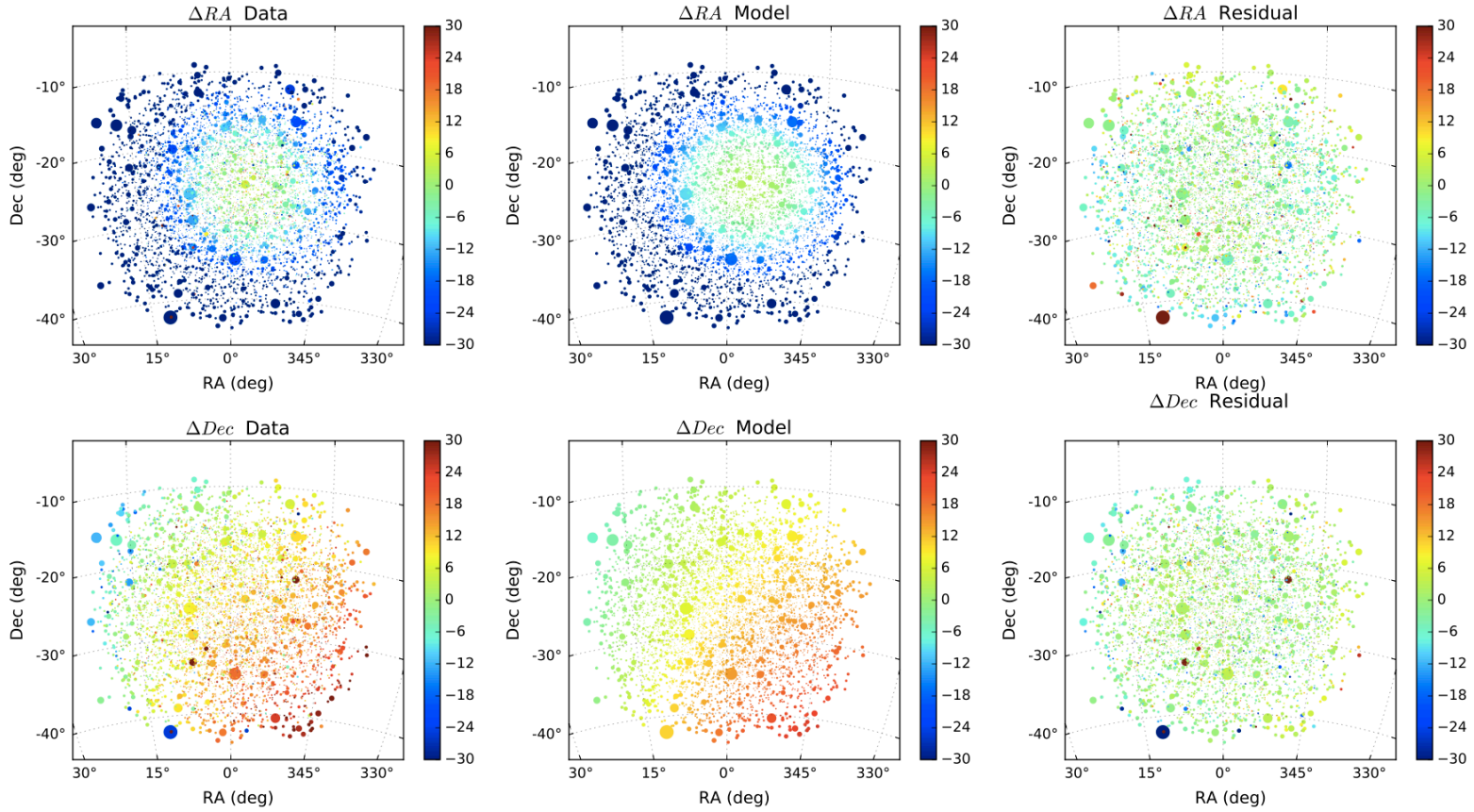


Figure 6.9: Sources scattered in RA and Dec from the NVSS and SUMSS counterparts (left) are well modeled by a 2D polynomial (middle). The offsets after correction are near-zero with no visible bias (right). Outliers with large residual offsets  $> 3\sqrt{\sigma_{\text{RA}}^2 + \sigma_{\text{Dec}}^2}$  were flagged and inspected in §5. The polynomial coefficients are given in Table 6.1.

## 6.6 *The Catalog*

Table 6.2 lists a subset of the catalog selected to represent a diverse sample. The PUMA position and spectral cross-match information for these are illustrated in Figure 6.10. The complete catalog of 7394 sources is included in the electronic supplement. The columns are:

1. **Name** Source name.
2. **RAJ2000** Corrected mean J2000 Right Ascension in degrees of the snapshot detections.
3. **DECJ2000** Corrected mean J2000 Declination in degrees of the snapshot detections.
4. **e\_RAJ2000** Standard deviation of the measured Right Ascension for all snapshot detections in arcseconds.
5. **e\_DECJ2000** Standard deviation of the measured Declination for all snapshot detections in arcseconds.
6. **S\_182** Weighted mean integrated 182 MHz flux density measured in Jy.
7. **e\_S\_182** Standard deviation of the measured flux density for all snapshot detections in Jy.
8. **EB\_corr** Estimated flux density correction factor for Eddington Bias.
9. **R\_class** The reliability classification (0-9).
10. **Beam** The mean relative beam response (0–1) at the source location.
11. **N\_det** Number of snapshots the source was detected in.
12. **Match\_Type** Type of match: `isolated`, `dominant`, `multiple`, `combine`, or `none`.

13. **Inspected** 0 if not visually inspected; 1 if the catalogue data and images were inspected; 2 if the match was modified by the authors.
14. **Match\_RAJ2000** J2000 Right Ascension in degrees of the catalogue match to NVSS or SUMSS.
15. **Match\_DECJ2000** J2000 Declination in degrees of the catalogue match to NVSS or SUMSS.
16. **e\_Match\_RAJ2000** Uncertainty in Right Ascension of the NVSS or SUMSS catalogue match in arcseconds.
17. **e\_Match\_DECJ2000** Uncertainty in Declination of the NVSS or SUMSS catalogue match in arcseconds.
18. **SI** Spectral index  $\alpha$  from a power law spectral index fit  $S \propto \nu^\alpha$  to all catalogue matches.
19. **e\_SI** Error on the spectral index parameter.
20. **S\_74** Flux density in Jy of the VLSSr catalogue match.
21. **e\_S\_74** VLSSr flux density error in Jy.
22. **S\_408** Flux density in Jy of the MRC catalogue match.
23. **e\_S\_408** MRC flux density error in Jy.
24. **S\_843** Flux density in Jy of the SUMSS catalogue match.
25. **e\_S\_843** SUMSS flux density error in Jy.
26. **S\_1400** Flux density in Jy of the NVSS catalogue match.

27. **e\_S\_1400** NVSS flux density error in Jy.

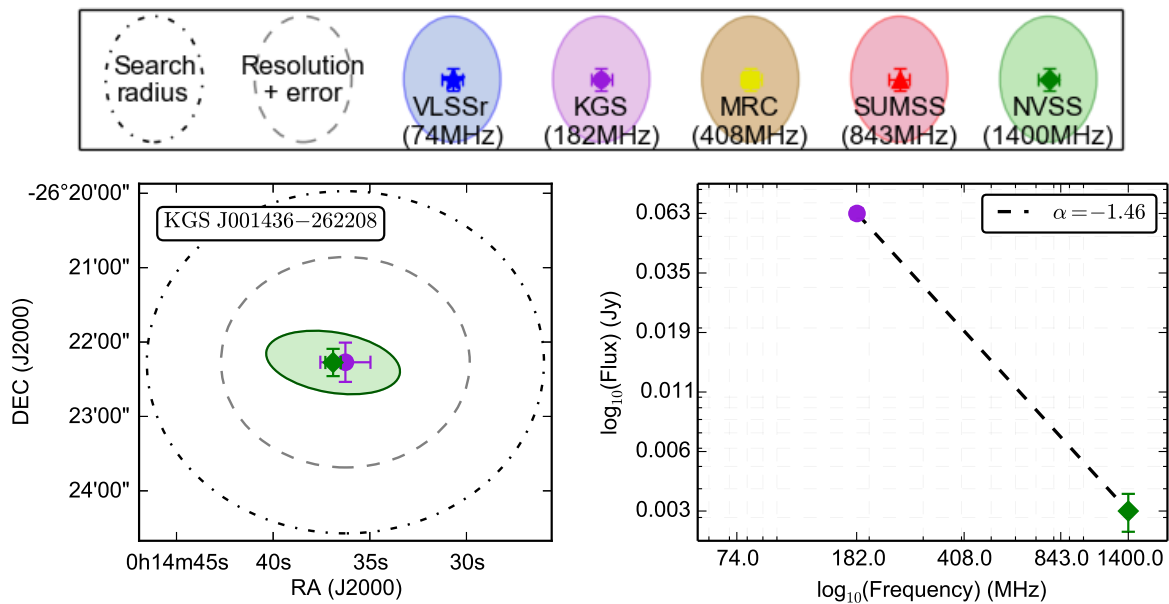
28. **VLSSr** VLSSr source name.

29. **MRC** MRC source name.

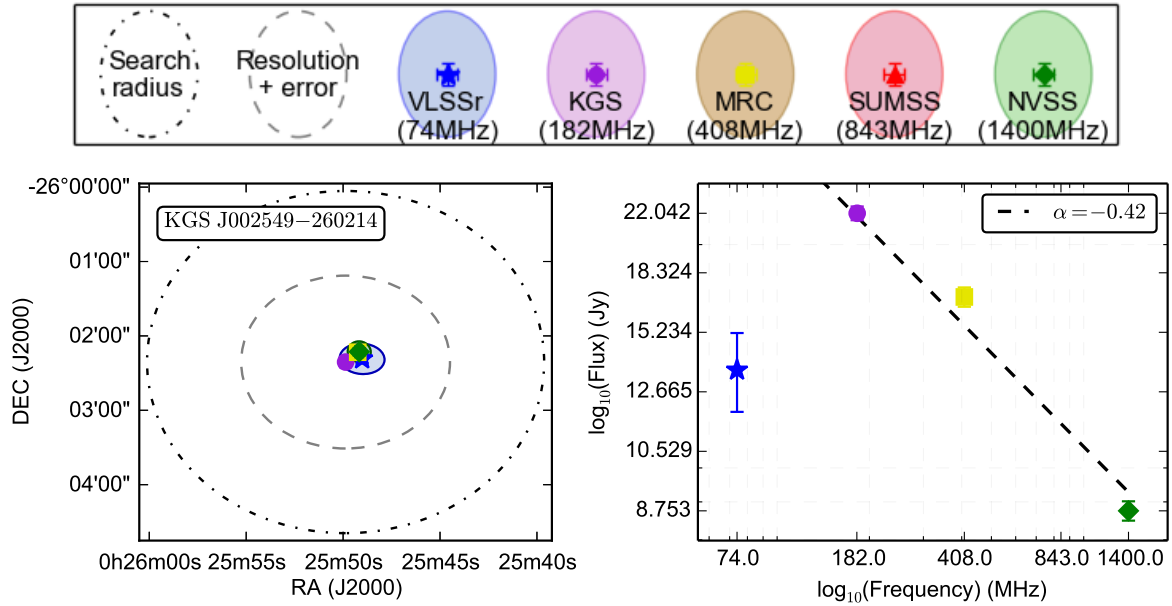
30. **SUMSS** SUMSS source name.

31. **NVSS** NVSS source name.

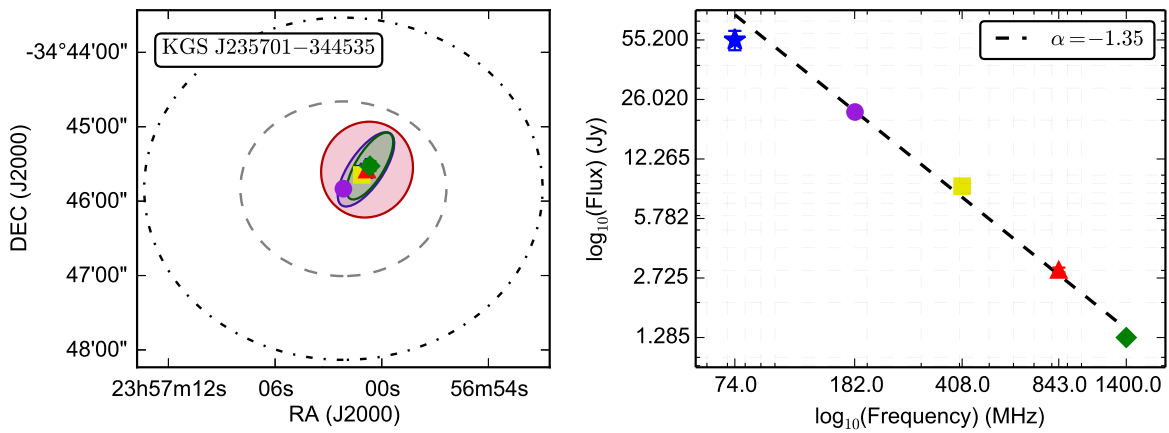
Figure 6.10: The PUMA match results for the 10 sources listed in Table 6.2, selected to demonstrate a variety of possible match and source types. The left column shows the uncorrected catalog positions, errors, and reported shape when available. The gray dashed line indicates the approximate PSF FWHM about the KGS source position. The dot-dashed black line marks the 2.3' initial search radius.



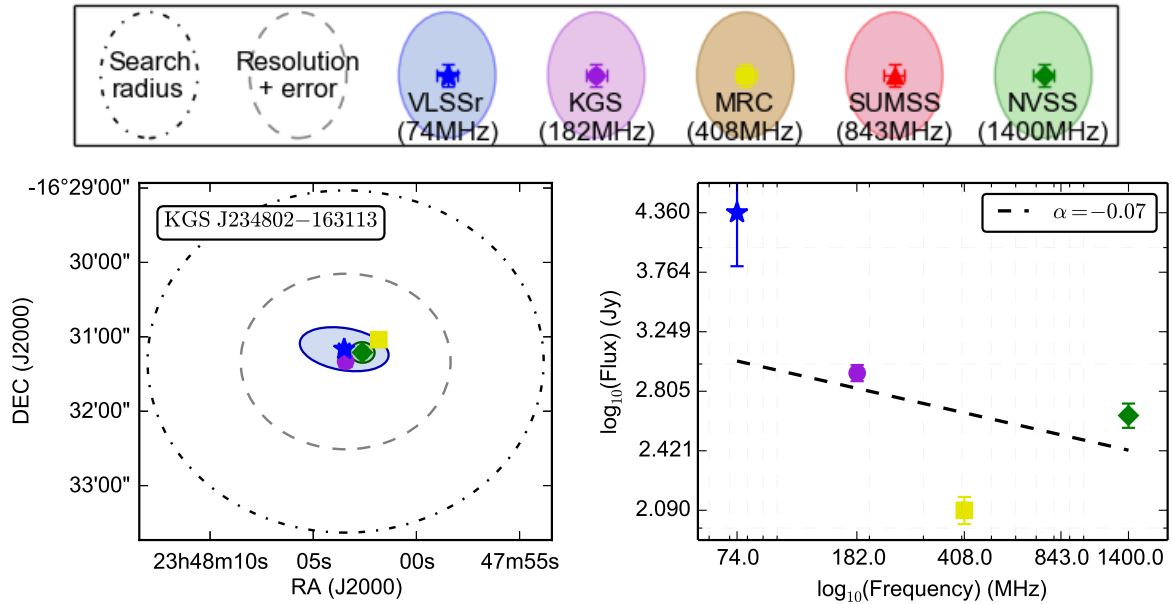
(a) **KGS J001436-262208** (NVSS J001436-262216): The faintest source in the catalog at 63 mJy with excellent positional agreement to an NVSS point source. The spectral index is steep at  $\alpha=-1.46$ , but this does not account for the Eddington bias correction listed in Table 6.2.



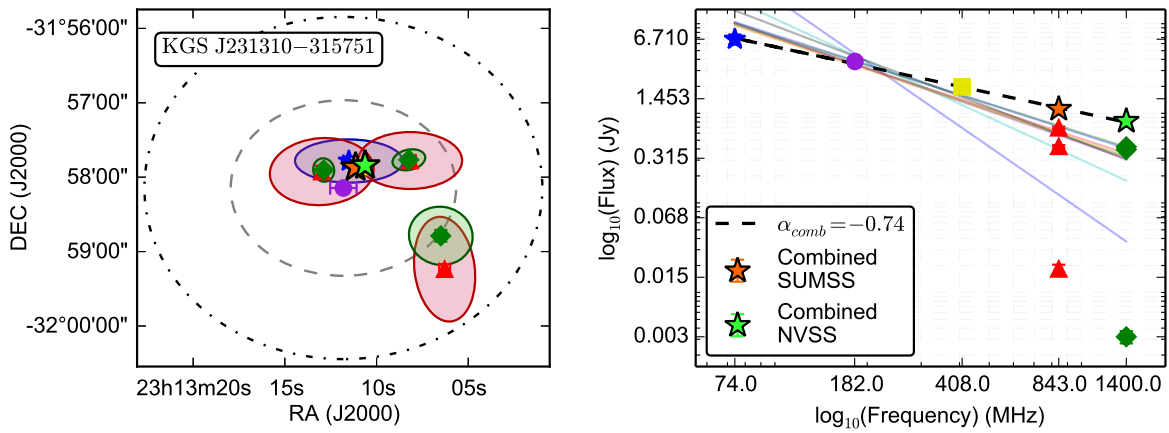
(b) **KGS J002549-260214** (PKS B0023-263): A strongly peaked spectrum source.



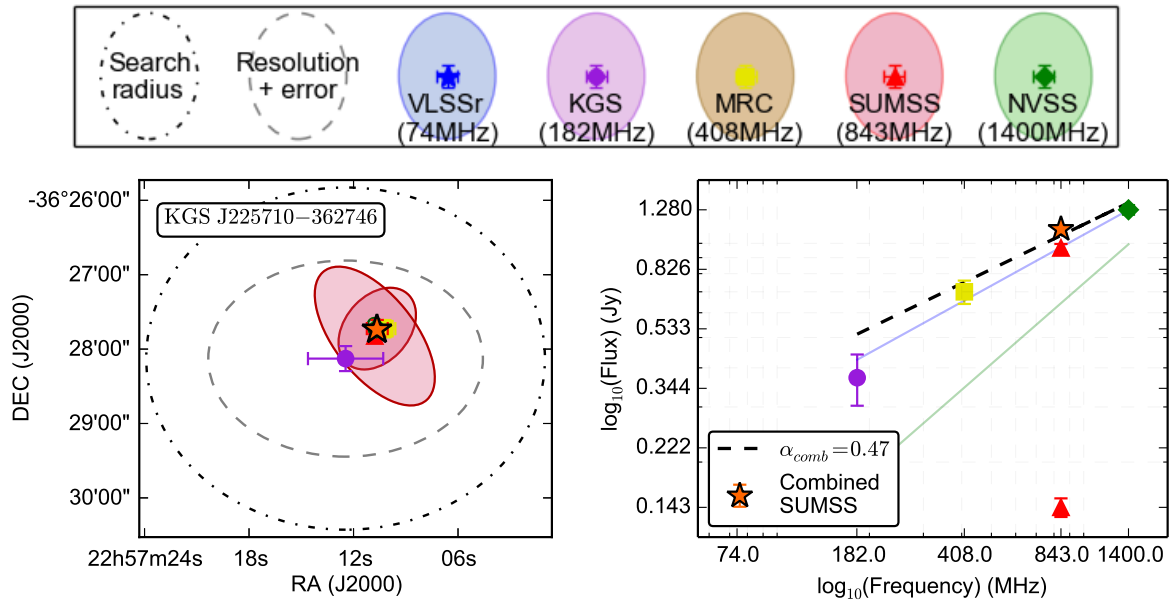
(c) **KGS J235701-344535** (PKS B2354-350): This source demonstrates excellent positional agreement with a counterpart observed in all five comparison catalogs.



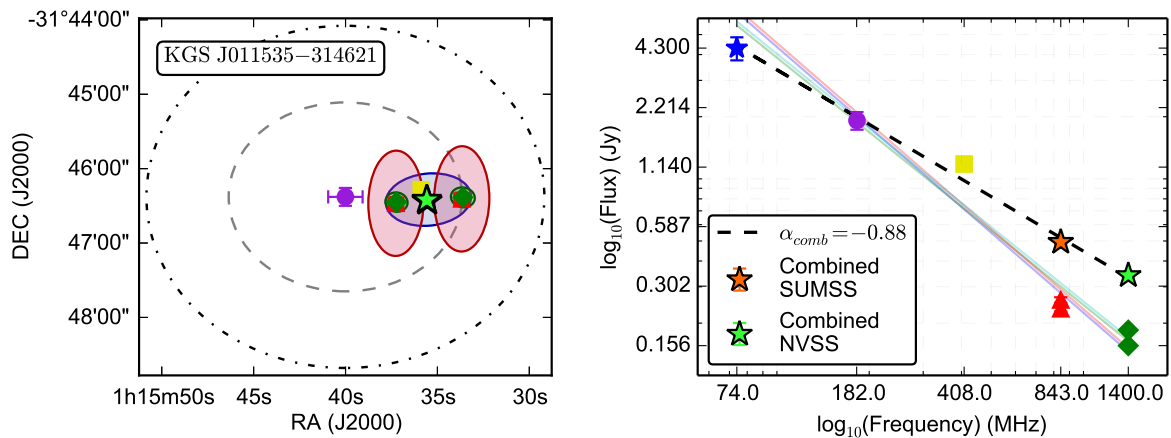
(d) **KGS J234802-163113** (PKS B2345-167): A good positional match with a poorly fit spectrum. This was subsequently identified as a variable QSO [42].



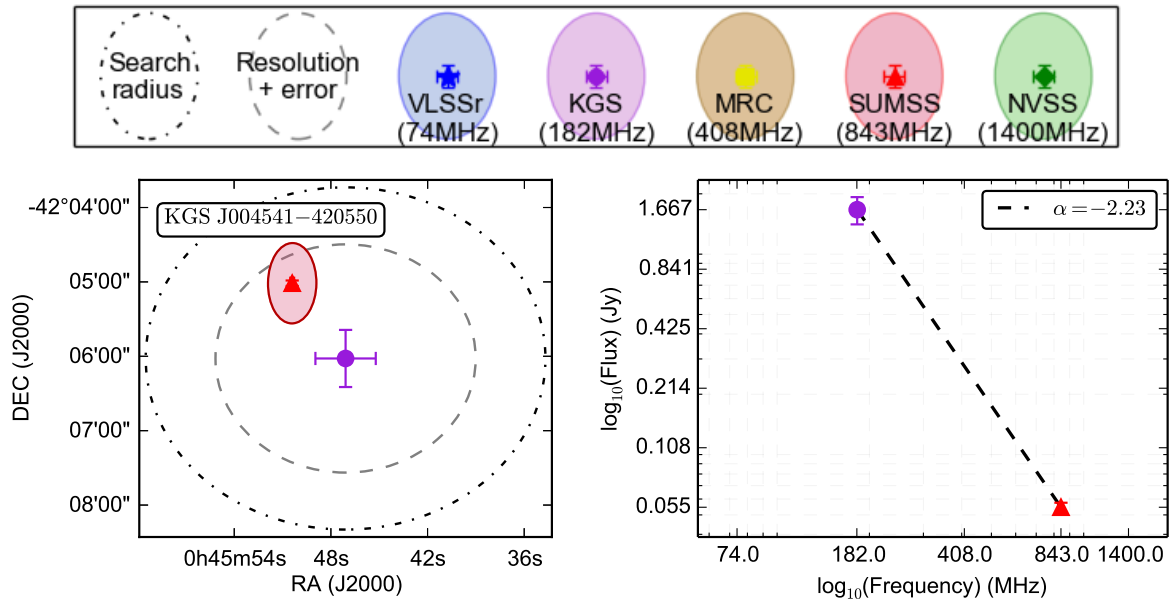
(e) **KGS J231310-315751** (PKS B2310-322): An extended source well-matched to a double at higher frequencies. The combined spectrum is well-fit by a power law. A confusing source (lower-right) within the initial search radius is ignored.



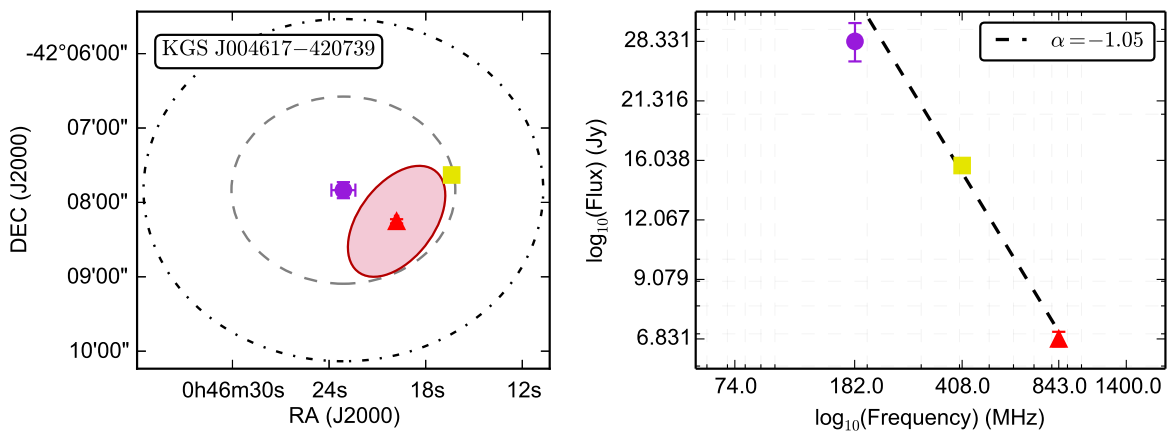
(f) **KGS J225710-362746** (PKS B2254-367): A source with a positive spectral index  $\alpha \sim 0.5$ .



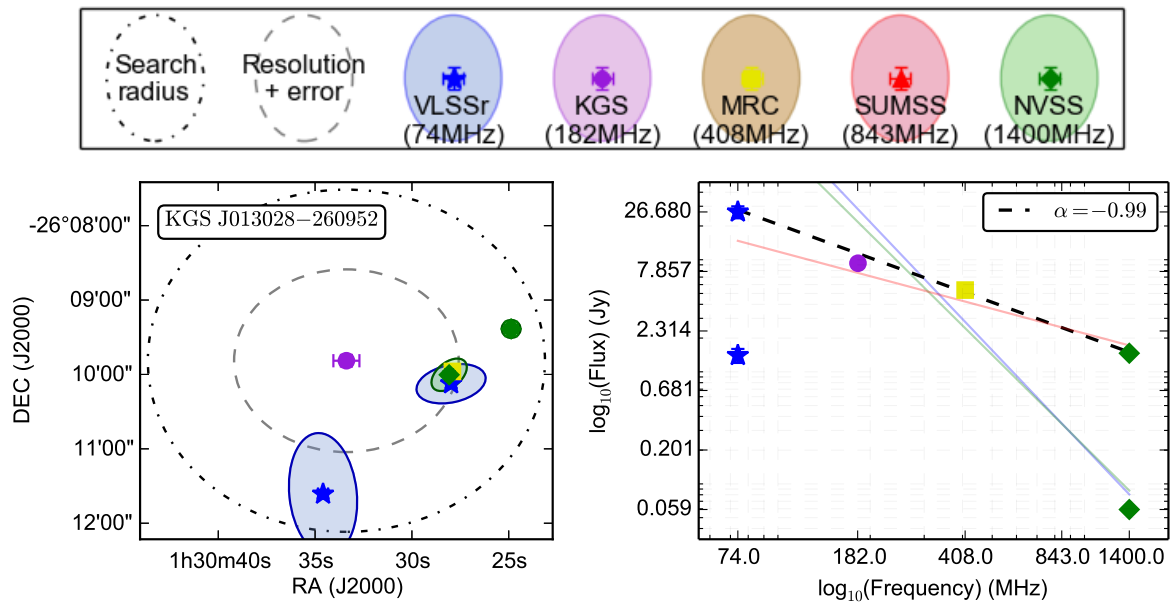
(g) **KGS J011535-314621** (PKS B0113-320): An example of a multiple match detected in all catalogs. When the components are combined, there is good spectral agreement with a power law fit, confirming the match despite the positional offset.



(h) **KGS J004541-420550** (SUMSS J004550-420501): This is an extremely steep spectrum source detected near the edge of the field with  $\alpha=-2.23$ . The flux density may be subject to error in the beam model.



(i) **KGS J004617-420739** (PKS B0043-424): This is the brightest source in the catalog at 28.3 Jy. It is near the edge of the field and exhibits a large positional bias. Two components are resolved in SUMSS and combined.



(j) **KGS J013028-260952** (PKS B0128-264): A good demonstration of the Bayesian positional match selection for a source with strong positional bias. The selected match has a posterior probability of 0.99 and is confirmed by the spectral fit.

Table 6.2: A sample subset of the catalog data. Ten sources were chosen to represent a diversity of characteristics. Ordering is by increasing distance from field centre ( $0^h, -27^\circ$ ). The corresponding cross-match results are shown in Figure 6.10.

Name	KGS_RAJ2000	KGS_DECJ2000	e_RAJ2000	e_DECJ2000	S_182	e_S_182	EB_corr	R_class	Beam
KGS J001436-262208	3.65004	-26.36889	19.4	15.8	0.063	4e-3	0.45	8	0.94
KGS J002549-260214	6.45528	-26.03734	4.0	0.7	22.04	0.47	1.00	0	0.77
KGS J235701-344535	359.25422	-34.75985	2.9	1.4	22.21	0.36	1.00	0	0.57
KGS J234802-163113	357.01100	-16.52027	4.3	1.8	2.93	0.06	1.00	0	0.45
KGS J231310-315751	348.29558	-31.96420	10.8	1.8	3.79	0.22	1.00	1	0.38
KGS J225710-362746	344.29165	-36.46284	32.4	10.1	0.37	0.07	1.00	8	0.20
KGS J011535-314621	18.89879	-31.77273	14.0	7.2	1.92	0.19	1.00	4	0.15
KGS J004541-420550	11.42311	-42.09745	28.1	23.0	1.67	0.26	0.98	4	0.08
KGS J004617-420739	11.57300	-42.12757	11.2	6.5	28.33	2.59	1.00	4	0.08
KGS J013028-260952	22.61847	-26.16470	10.1	4.7	9.30	0.84	1.00	6	0.09

N_det	Match_Type	Inspected	Match_RAJ2000	Match_DECJ2000	e_Match_RAJ2000	e_Match_DECJ2000	SI	e_SI
5	isolated	1	3.65380	-26.37120	6.1	11.2	-1.46	
71	isolated	0	6.45490	-26.03690	0.4	0.7	-0.42	0.10
71	isolated	0	359.25280	-34.75880	0.7	0.7	-1.35	0.05
71	isolated	1	357.01090	-16.52020	0.4	0.7	-0.07	0.10
71	multiple	1	348.29430	-31.96400	0.7	0.7	-0.74	0.03
7	multiple	1	344.29450	-36.46190	0.7	0.7	4.300	0.550
41	multiple	1	18.89830	-31.77370	0.7	0.7	-2.23	
31	isolated	1	11.46000	-42.08360	1.8	2.2	-0.88	0.07
43	isolated	2	11.58250	-42.13760	1.4	1.8	-1.05	0.14
13	isolated	0	22.61700	-26.16680	0.4	0.7	-0.99	0.06

**Table 6.2 continued.**

S_74	e_S_74	S_408	e_S_408	S_843	e_S_843	S_1400	e_S_1400	VLSSr	MRC	SUMSS	NVSS
						3.2E-3	6E-4			001436...	
13.55	1.65	17.00	0.51			8.75	0.26	J002549...	0023-263		002549...
55.20	6.72	8.70	0.35	3.02	0.09	1.28	0.04	J235700...	2354-350	J235700...	235700...
4.36	0.54	2.09	0.07			2.64	0.08	J234803...	2345-167		234802...
6.71	0.86	1.97	0.10	1.12	0.03	0.82	0.02	J231311...	2310-322	J231308	231306...
		0.70	0.06	1.11	0.03	1.28	0.05		2254-367	J225710...	225710...
4.30	0.55	1.18	0.07	0.49	0.01	0.34	7E-3	J011535...	0113-320	J011533...	011533...
				5.5E-2	3E-3					J004550...	
		15.65	0.39	6.83	0.237				0043-424	J004613...	
26.68	3.26	5.36	0.17			1.46	0.05	J013027...	0128-264		013028...

**Table 6.2 continued.**

## 6.7 Caveats

There are two important caveats to emphasize for potential users of the catalog.

### 6.7.1 Primary Beam Model

The purpose of this survey was primarily to build a foreground model for the EoR analysis. As such, we've elected to include sources covering the full field, out to 5% of the peak beam response. The accuracy of the source flux density measurements relies on the accuracy of the model of the primary beam response. In-situ measurements for beam sensitivity characterization are in progress but at the time of this analysis an analytic MWA beam shape was assumed. As sources move through the beam, trends in the light curves suggests a 10–20% error near the edge of the field ( $beam < 0.2$ ). This error has not been factored into the flux density error reported in the catalog.

### 6.7.2 Extended Sources

Among the 13% of the sources flagged as `multiple` and visually inspected, many exhibit extended morphologies that are not well represented by the sub-components indicated in the higher resolution catalogs. The MWA has many short baselines and much higher surface-brightness sensitivity than most radio telescopes. This has already led to the discovery of a number of large sources that were resolved out in previous surveys (e.g. a dying giant radio galaxy presented in [19]). Diffuse emission picked up by the MWA will make interpreting the flux densities between surveys problematic for extended sources and care should be taken. For this reason, we have limited the above analysis to `isolated` sources.

## 6.8 Summary

In this chapter, I've presented the KGS MWA EoR0 catalog of radio sources at 182 MHz. Careful analysis revealed a small number of sources affected by Eddington bias in excess of their standard errors among sources with apparent flux density  $S \cdot beam \lesssim 100$  mJy, as well

as a significant positional bias that was traced to the MWACS sources used for calibration. Both biases were able to be modeled and corrected. There is no evidence for flux scale bias relative to the comparison catalogs, although there is strong evidence for spectral flattening toward lower frequencies. This could be indicative of intrinsic spectral curvature but a deeper analysis is required to sync the flux scales and rule out systematic effects. The overall broad-band spectral index distribution is dominated by sources matched to the NVSS and SUMSS, with a median of -0.85. The median spectral index at 182 MHz is predicted to be -0.71, which is consistent with independently processed sub-band spectral index measurements.

## Chapter 7

# UNIDENTIFIED & ULTRA STEEP SPECTRUM RADIO SOURCES

The astrophysical classifications of sources have little relevance to the EoR foreground model. However, in producing the KGS catalog, we have uncovered several new radio detections and can easily select for interesting and rare type sources based on their broad-band spectral properties. We temporarily diverge from the EoR focus to explore the ultra-steep spectrum (USS;  $\alpha \lesssim -1.4$  with  $S \propto \nu^\alpha$ ) population.

We expect the majority of point sources with power-law spectra above  $-30^\circ$  declination to be detected in the NVSS or VLSSr surveys. Steep spectrum sources that are missed in the NVSS should be detectable in the VLSSr and vice versa. However, the VLSSr completeness varies considerably and below  $-30^\circ$ , we can expect to detect ultra-steep spectrum sources that fall below the SUMSS completeness limit (18 mJy at 843 MHz). There are 25 previously undetected radio sources in the KGS and another 203 sources with an estimated 182 MHz spectral index  $< -1.4$ .

USS sources are known to trace galaxy clusters as well as high redshift ( $z$ ) radio galaxies (HzRGs). In this chapter we explore the USS population and discuss their spectral properties and potential identifications. We highlight observed trends and the predictive potential of broad-band spectral index measurements for USS source classification.

### **7.1 New Radio Source Detections**

The properties of the 25 sources with no previous radio detection are listed in Table 7.1.

Name	$RA(deg)$	$Dec(deg)$	$S(mJy)$	$\sigma_S(mJy)$	$f_{EB}$	$N_{det}$	$Beam_{mean}$	$R_{class}$	GClstr/Grp
KGS J233620-313606	354.0857	-31.6017	424	57	1.00	71	0.67	0	Abell S1136
KGS J232803-145208	352.0139	-14.8689	270	59	0.93	30	0.25	4	*
KGS J000958-353932	2.4944	-35.6590	164	27	1.00	47	0.50	4	Abell 2730
KGS J231311-230716	348.2981	-23.1212	147	47	1.00	22	0.54	6	Abell S1099
KGS J235156-165850	357.9836	-16.9808	136	14	0.98	14	0.51	6	*
KGS J235021-194846	357.5881	-19.8128	123	19	1.00	57	0.69	2	'
KGS J233116-192443	352.8168	-19.4120	118	20	0.95	11	0.59	6	
KGS J231928-302751	349.8685	-30.4644	117	20	0.98	19	0.61	6	
KGS J001054-341312	2.7254	-34.2201	116	18	0.97	14	0.60	6	2PIGG SGP 5843
KGS J233617-244958	354.0732	-24.8328	98	16	1.00	47	0.82	4	
KGS J232926-255814	352.3617	-25.9707	93	8	0.99	17	0.80	6	
KGS J234703-305612	356.7634	-30.9368	90	16	0.98	32	0.81	4	
KGS J001640-215455	4.1674	-21.9155	82	12	0.97	7	0.84	8	
KGS J000215-275242	0.5636	-27.8786	79	14	0.99	44	0.94	4	2PIGG SGP 2684
KGS J002837-261426	7.1576	-26.2406	76	13	0.87	5	0.85	8	
KGS J234344-263049	355.9365	-26.5138	73	17	0.94	14	0.87	6	
KGS J235556-224242	358.9862	-22.7119	70	9	0.85	5	0.88	8	
KGS J234709-281746	356.7876	-28.2963	69	7	0.84	6	0.94	8	Abell 4037/4038
KGS J234851-232934	357.2135	-23.4929	68	6	0.66	6	0.91	8	
KGS J002451-204048	6.2162	-20.6801	240	23	1.00	49	0.61	4	Abell 0027
KGS J000821-193833	2.0904	-19.6427	217	65	0.99	34	0.68	4	Abell 0002
KGS J000412-151811	1.0506	-15.3031	211	40	1.00	13	0.39	6	
KGS J001702-312239	4.2596	-31.3775	155	31	0.99	63	0.74	3	2PIGG SGP 7135
KGS J001007-282942	2.5317	-28.4951	132	28	1.00	37	0.91	4	2PIGG SGP 8480
KGS J235139-255937	357.9129	-25.9938	81	20	0.98	11	0.92	6	Abell 2667

Table 7.1: There are 19 isolated sources (top) and 6 sources of apparently extended or diffuse emission (bottom). Postage stamp images are shown in the discussion. For 11 sources, an Abell galaxy cluster or 2PIGG galaxy group member is located within 2' of the source position. These are indicated in the GClstr/Grp column.

\* A 2MASX extended IR source and several possible galaxy counterparts are found within the search radius.

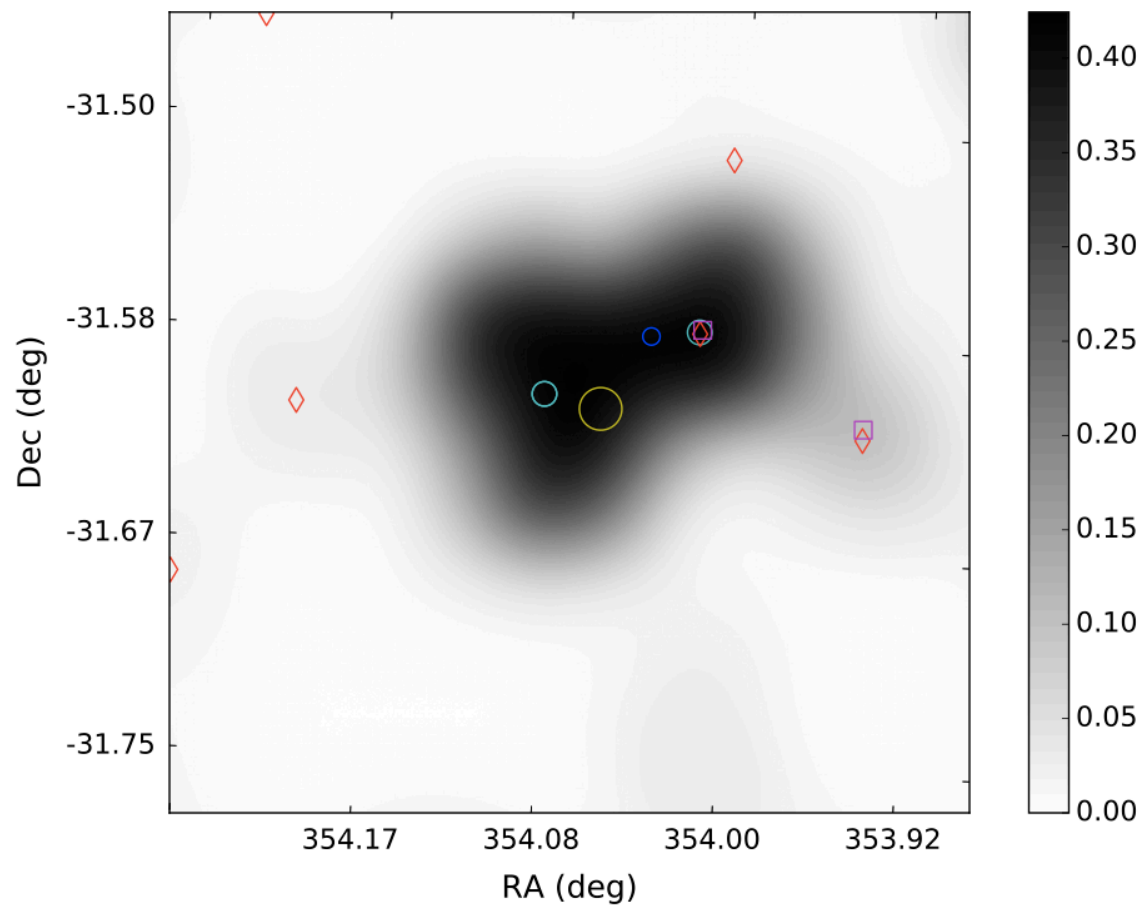
A search was made around their positions in the NASA Extragalactic Database<sup>1</sup> to find potential associations within a 2.3' search radius. The match radius corresponds to the beam width, but we note that this does not imply a true counterpart. In most cases, the positional error is much smaller than the beam width.

Below, we present postage stamp images and discuss possible cross-match identifications for each source. Images are smoothed with a cubic interpolation for easier visualization. All KGS sources are marked with a cyan circle. Nearby sources in the MWACS (blue circle), NVSS (red diamond), SUMSS (magenta square), VLSSr (orange triangle), and Abell cluster catalog (yellow circle) are also indicated.

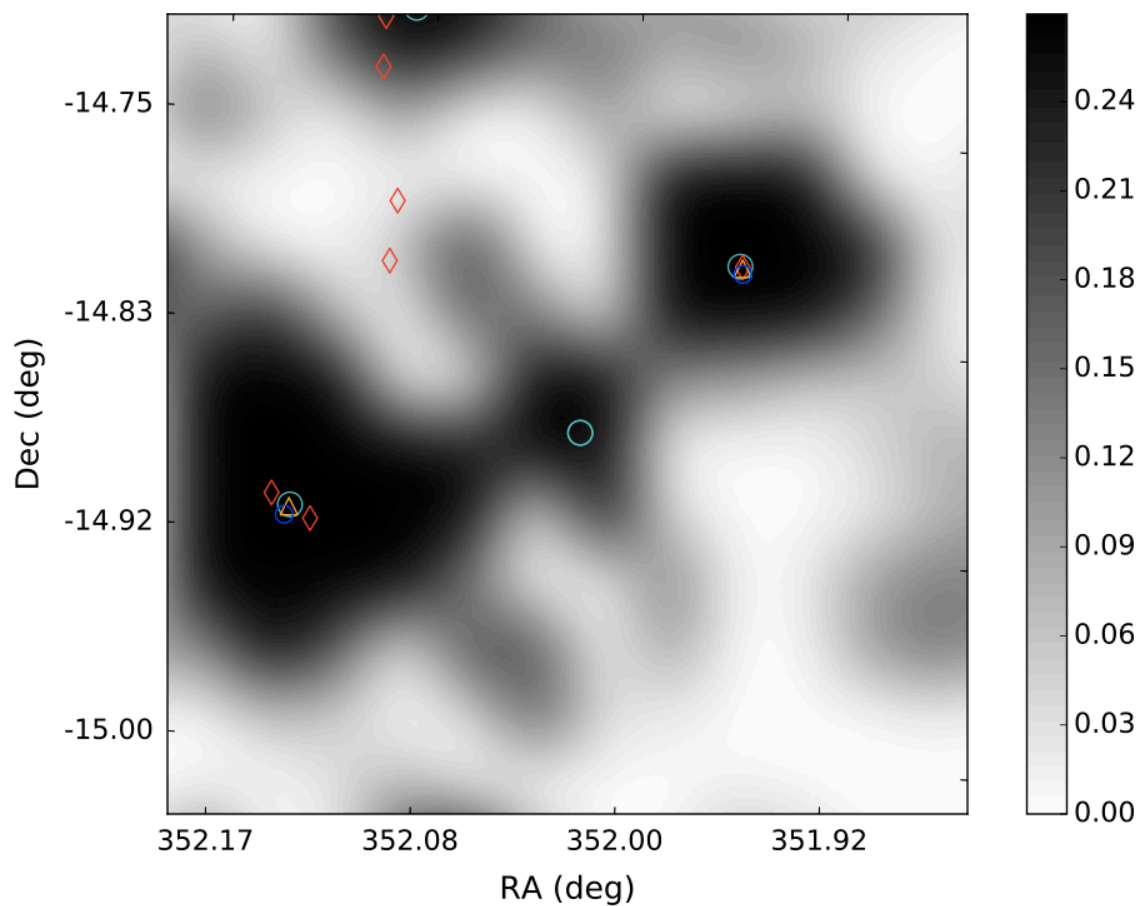
---

<sup>1</sup>The NASA/IPAC Extragalactic Database (NED) is operated by the Jet Propulsion Laboratory, California Institute of Technology, under contract with the National Aeronautics and Space Administration.

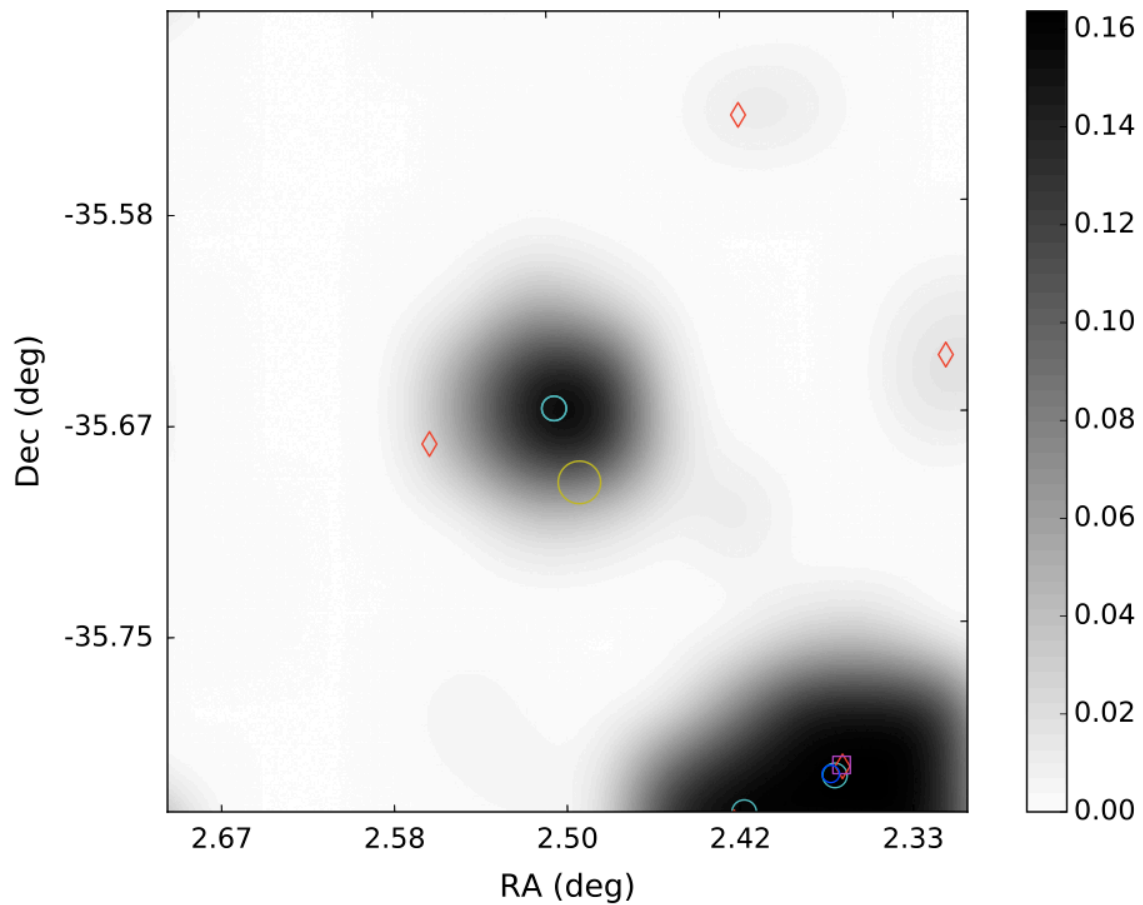
**KGS J233620-313606:** A 421 mJy source most likely associated with the galaxy cluster Abell S1136. The elongated shape suggests it may be a blended double. The cluster center is located at a distance of  $0.92'$  with a radius  $R_c \equiv 1.72/z = 27.5'$ . The source is most closely matched to GALEXASC J233618.66-313604.6 at  $17''$  and the x-ray source SW J233617-313626 at  $0.7'$ . Several other cluster members and sources at all wavelengths are also found within the  $2.3'$  search radius.



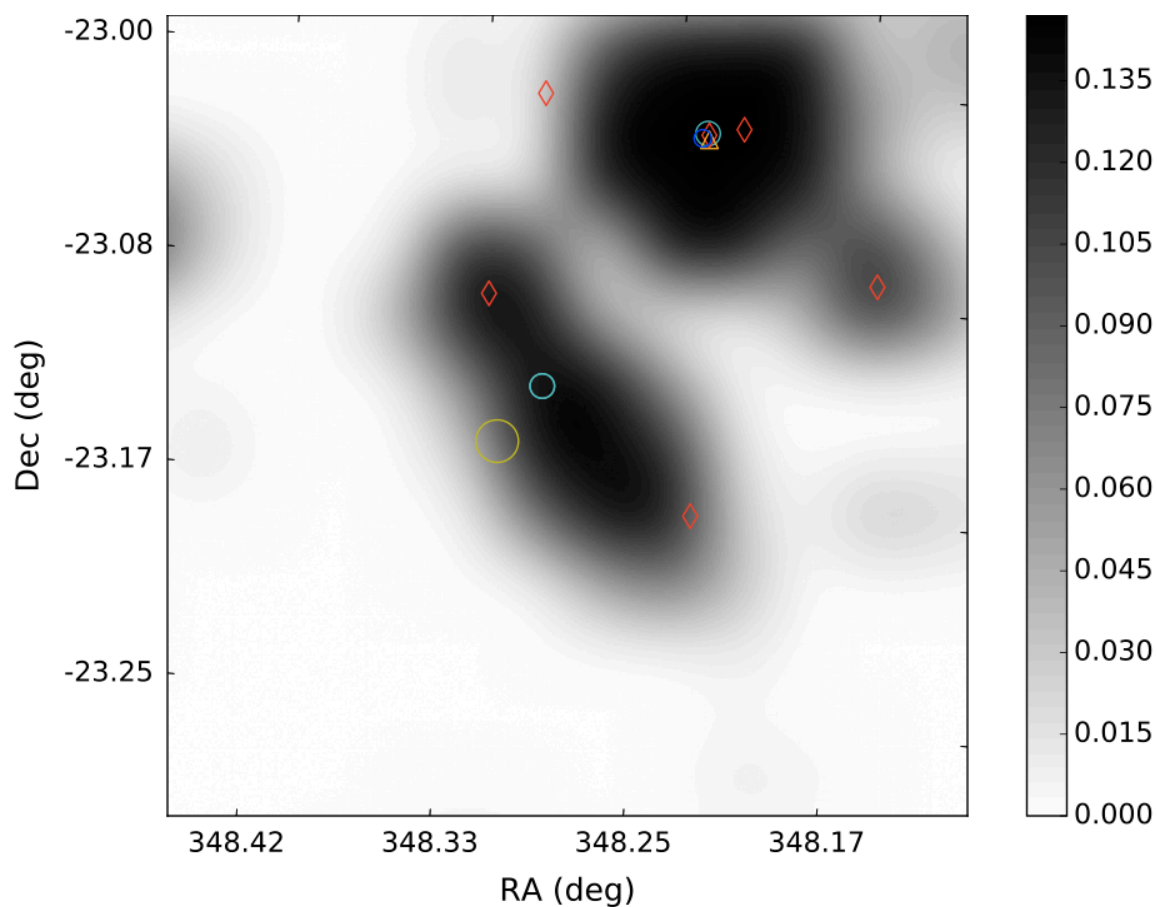
**KGS J232803-145208:** A 249 mJy source most closely matched to the galaxy APMUKS(BJ) B232526.08-150914.0 at 0.6' from the source position. Five APMUKS(BJ) galaxies are found within 2.3' including one extended IR source 2MASX J23280750-1452221 at 1' separation.



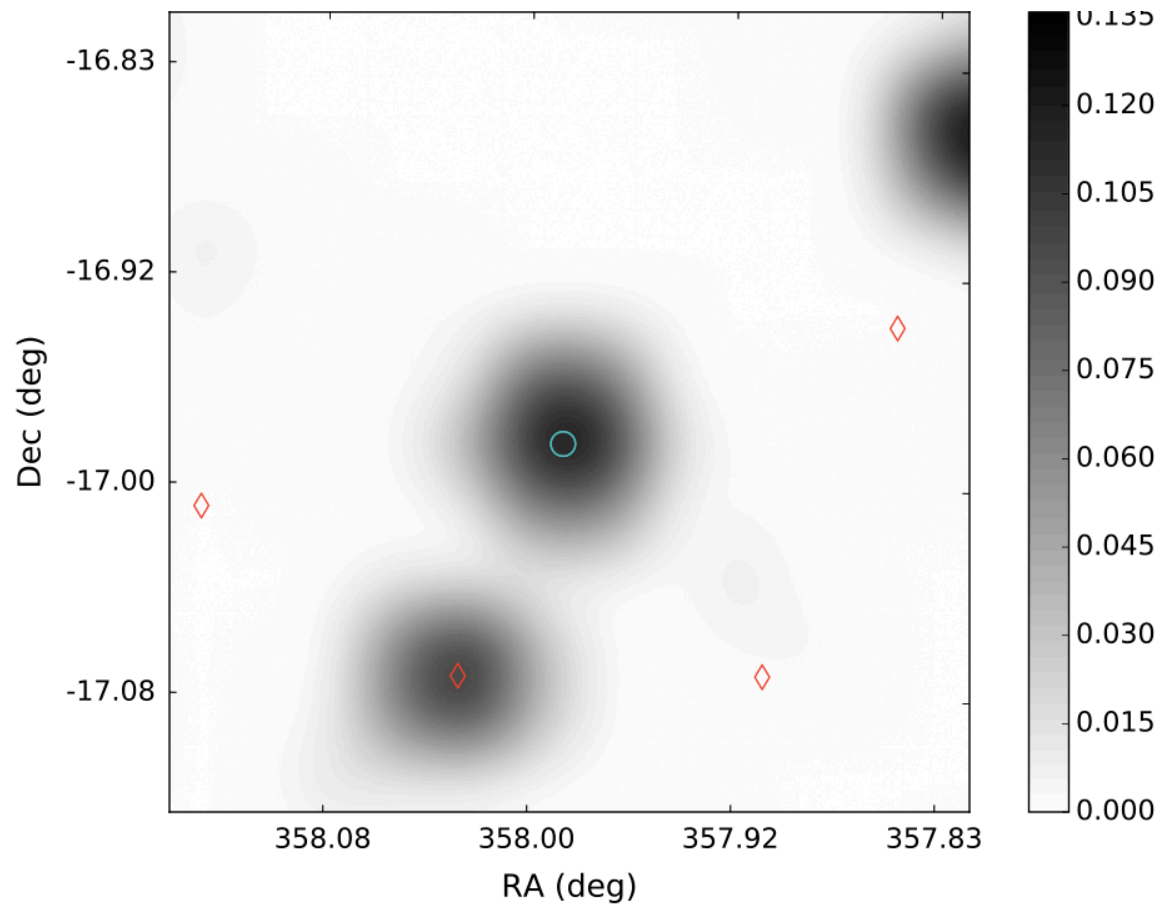
**KGS J000958-353932:** A 164 mJy source 19" from the galaxy cluster member EDCC 408:[CGN95] 000726.8-3556 and extended IR source 2MASX J00095865-3539515. EDCC 408 is cross-identified with Abell 2730, centered 1.8' from the source position with a 14' estimated cluster radius. Many other sources are found within 2.3' including 12 other cluster members.



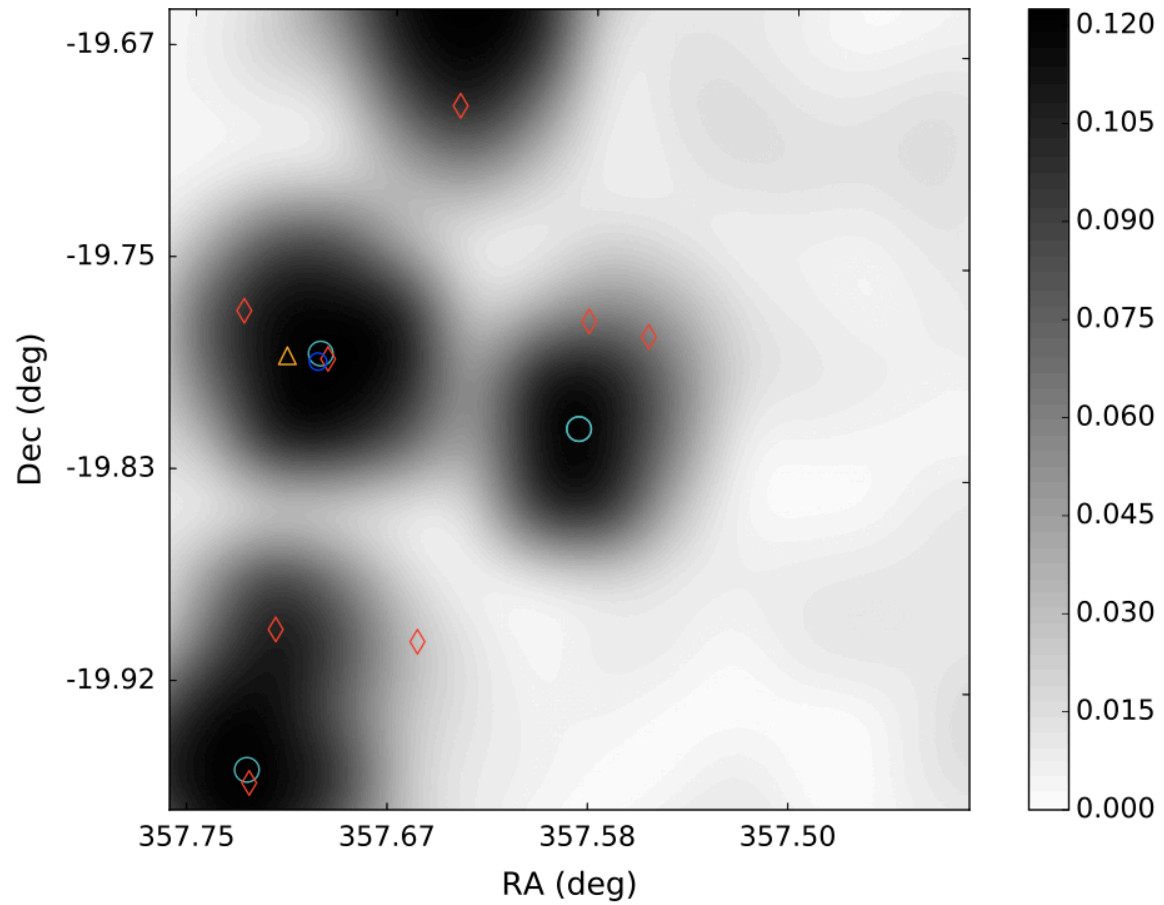
**KGS J231311-230716:** A 147 mJy blended double most closely matched to the UV source GALEXASC J231312.42-230715.1 with 20" separation. The center of the galaxy cluster ABELL S1099 is found at a distance of 1.8' with a cluster radius  $R_c = 15.6'$ . Two cluster members CMW2004 388 and 339 are cross-identified with 2MASX extended IR sources and located at 1.35' and 1.82' respectively from the source position.



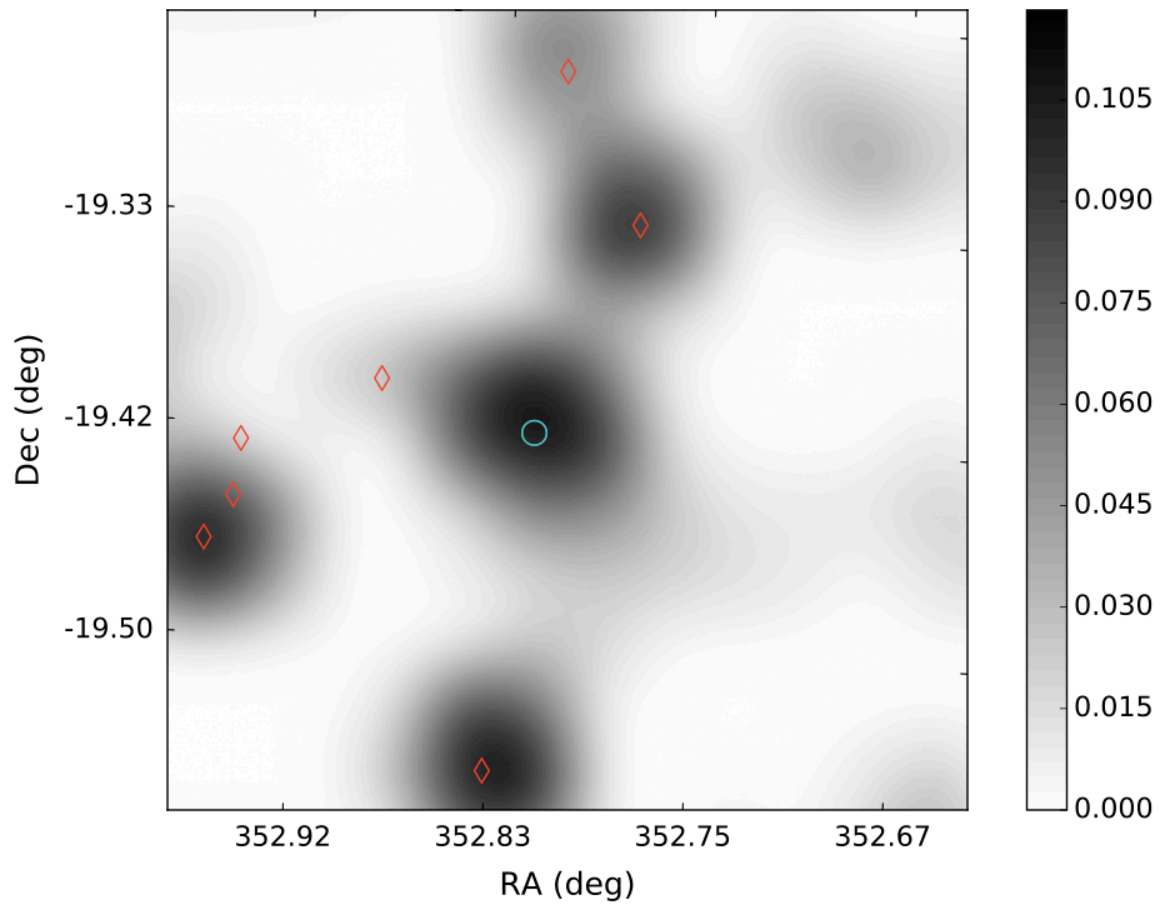
**KGS J235156-165850:** A 133 mJy source unmatched within  $30''$ . The extended IR source 2MASX J23515772-1657374 is  $1.3'$  from the source position and many MRSS galaxy and GALEXASC UV sources are located with the search radius.



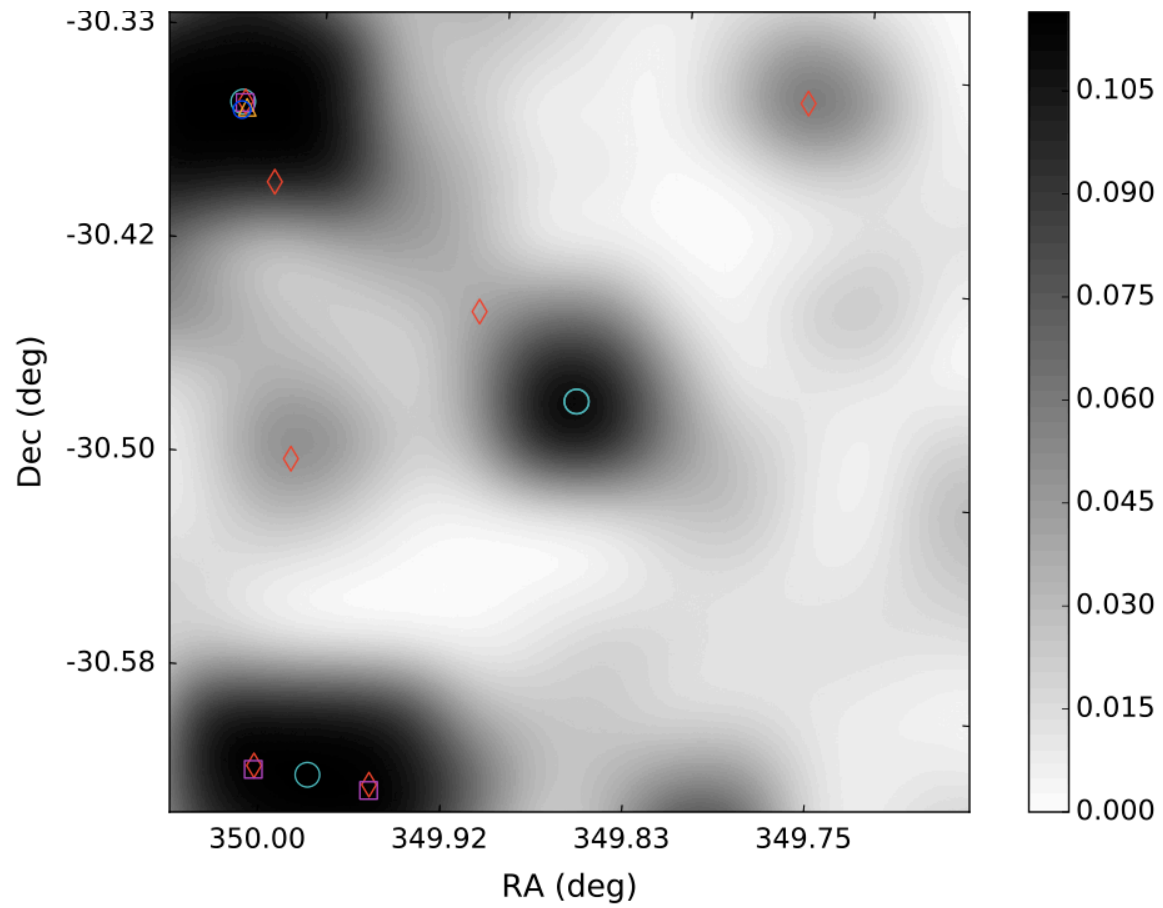
**KGS J235021-194846:** A 123 mJy source unmatched within 30". Eleven MRSS and AP-MUKS identified galaxies are found within 2.3'.



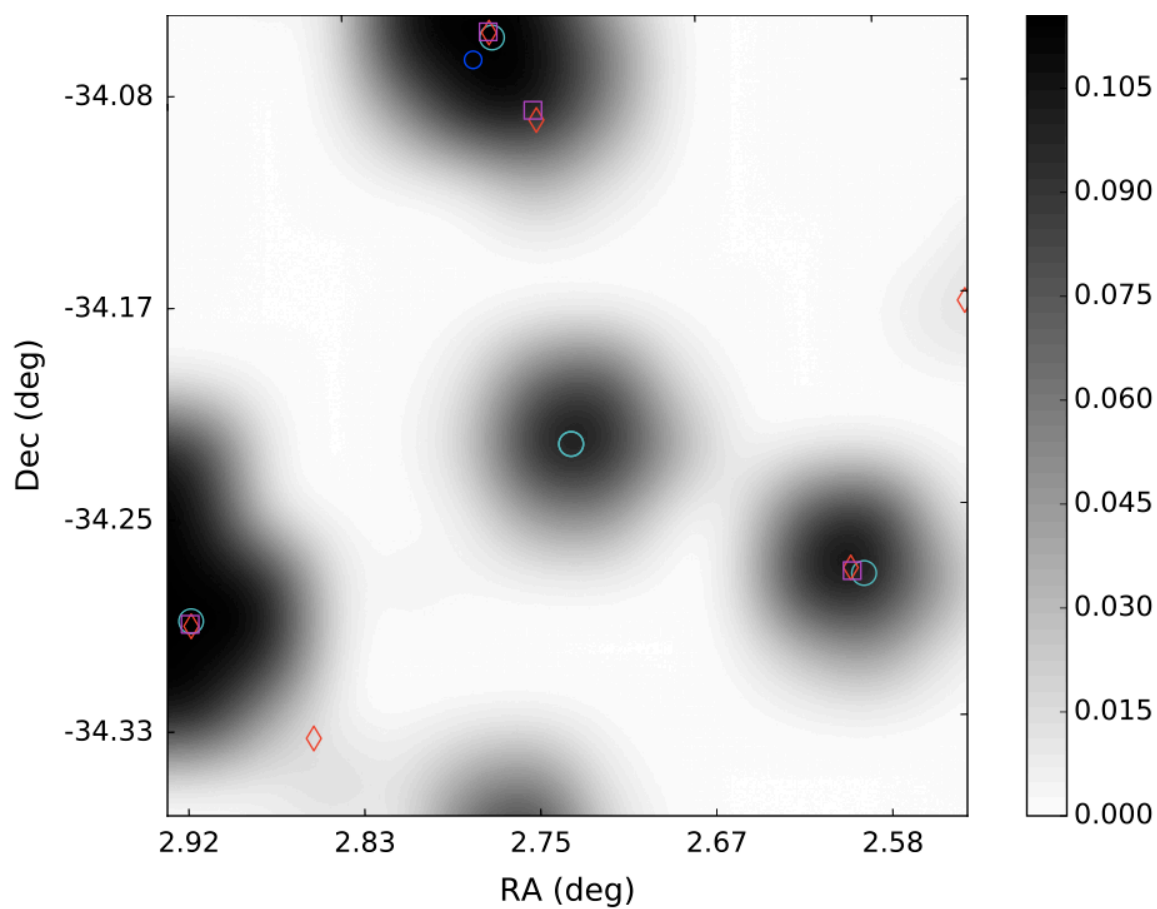
**KGS J233116-192443:** A 113 mJy source matched most closely to the UV source GALEX-ASC J233115.57-192441.2 at 7". Several other UV sources and three MRSS galaxies are located within 2.3'.



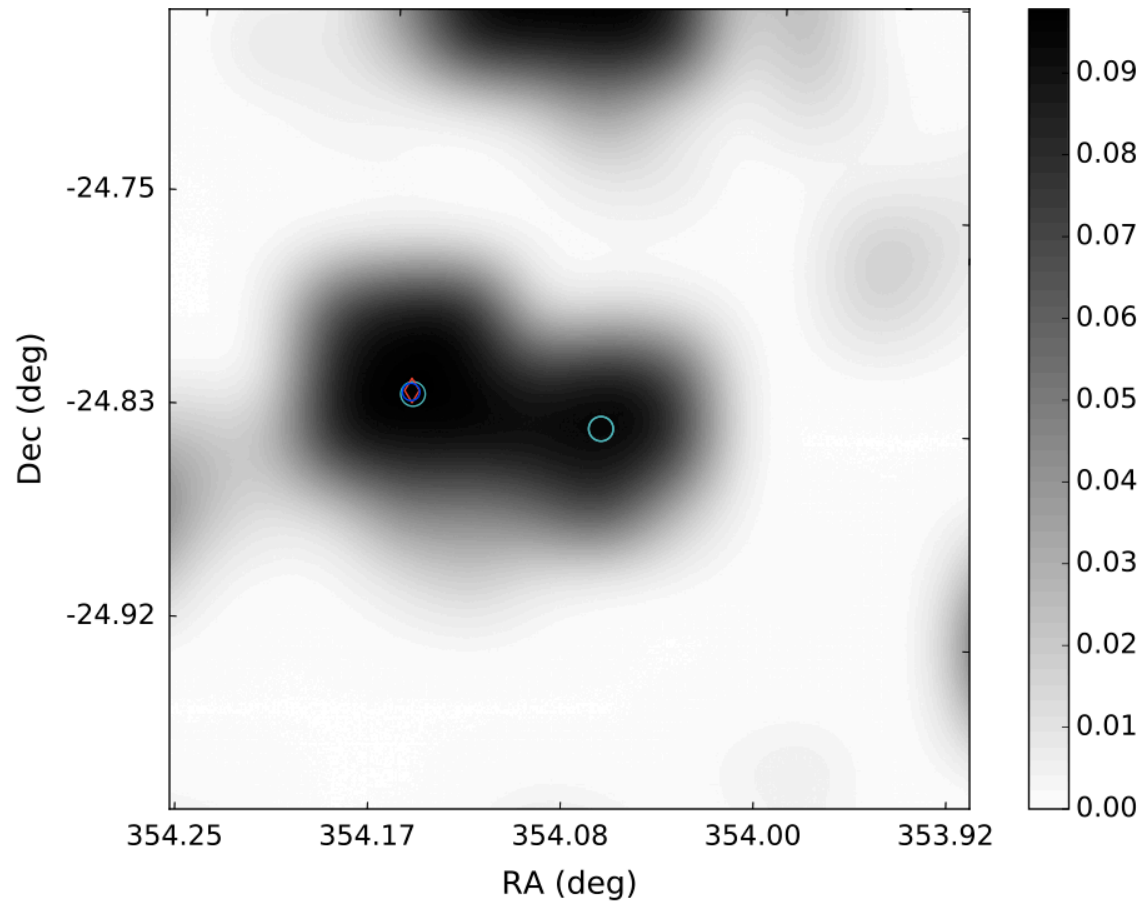
**KGS J231928-302751:** A 114 mJy source most closely matched to the quasar 2QZ J231927.7-302845 ( $z = 1.059$ ) at a  $0.9'$ . Several other UV sources and one MRSS galaxy are found within  $2.3'$ .



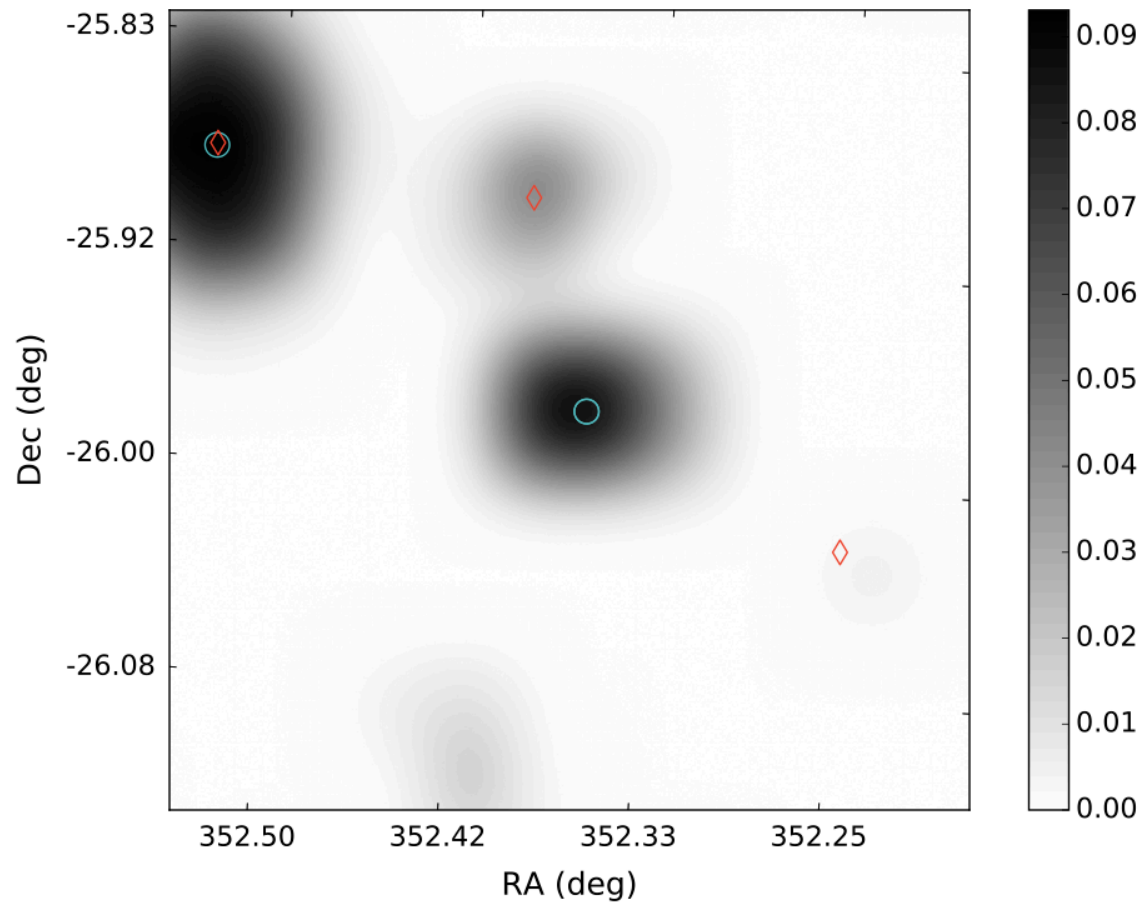
**KGS J001054-341312:** A 112 mJy source most closely matched to 2MASX J00105255-3413132 at 19" and 2dFGRS S495Z294 at 28". The latter is one of 4 members of the 14-member galaxy group 2PIGG SGP 5843 within 1.8'. The galaxy cluster and x-ray source APMCC 014 is 1.8' from the source position. A GALEXASC UV source and 1WGA X-ray source are also found at 20" and 26" respectively.



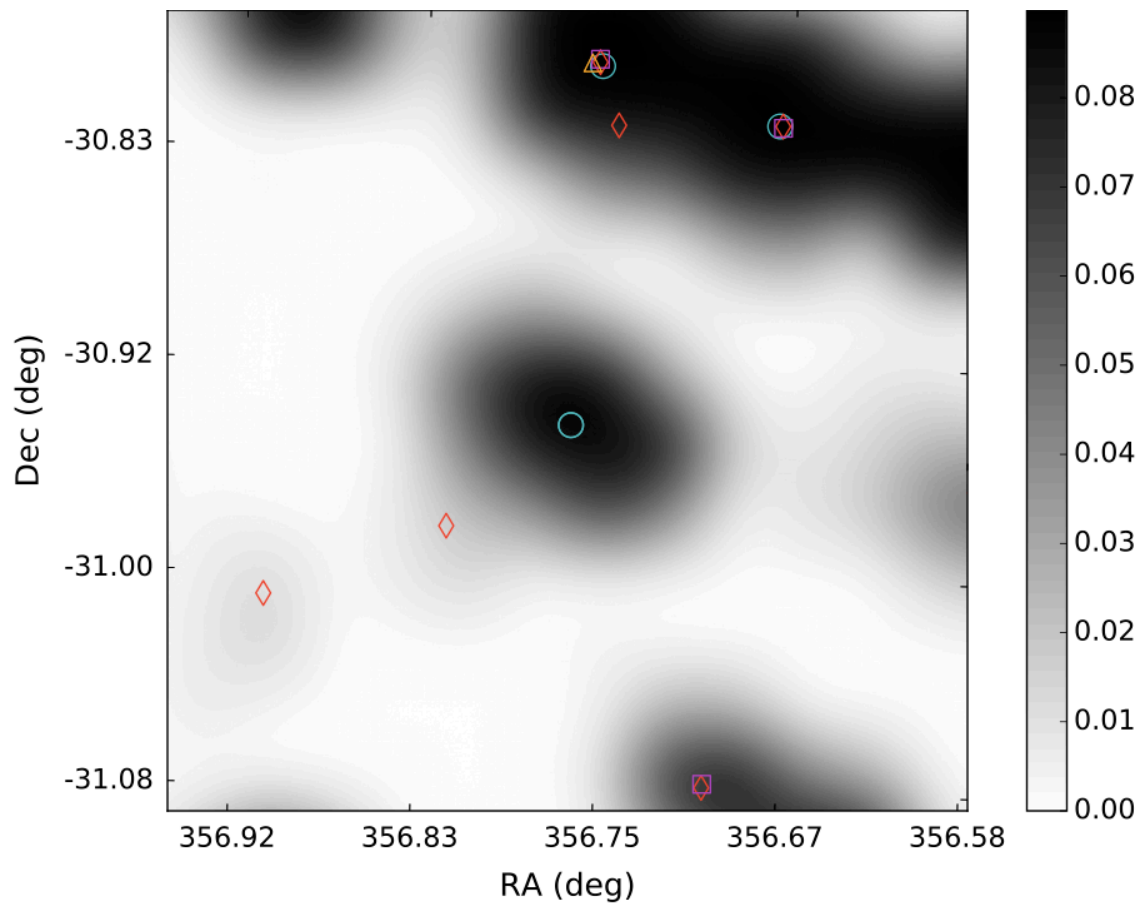
**KGS J233617-244958:** A 98 mJy source. A total of 24 galaxies and UV sources are found between 11" and 2.3'.



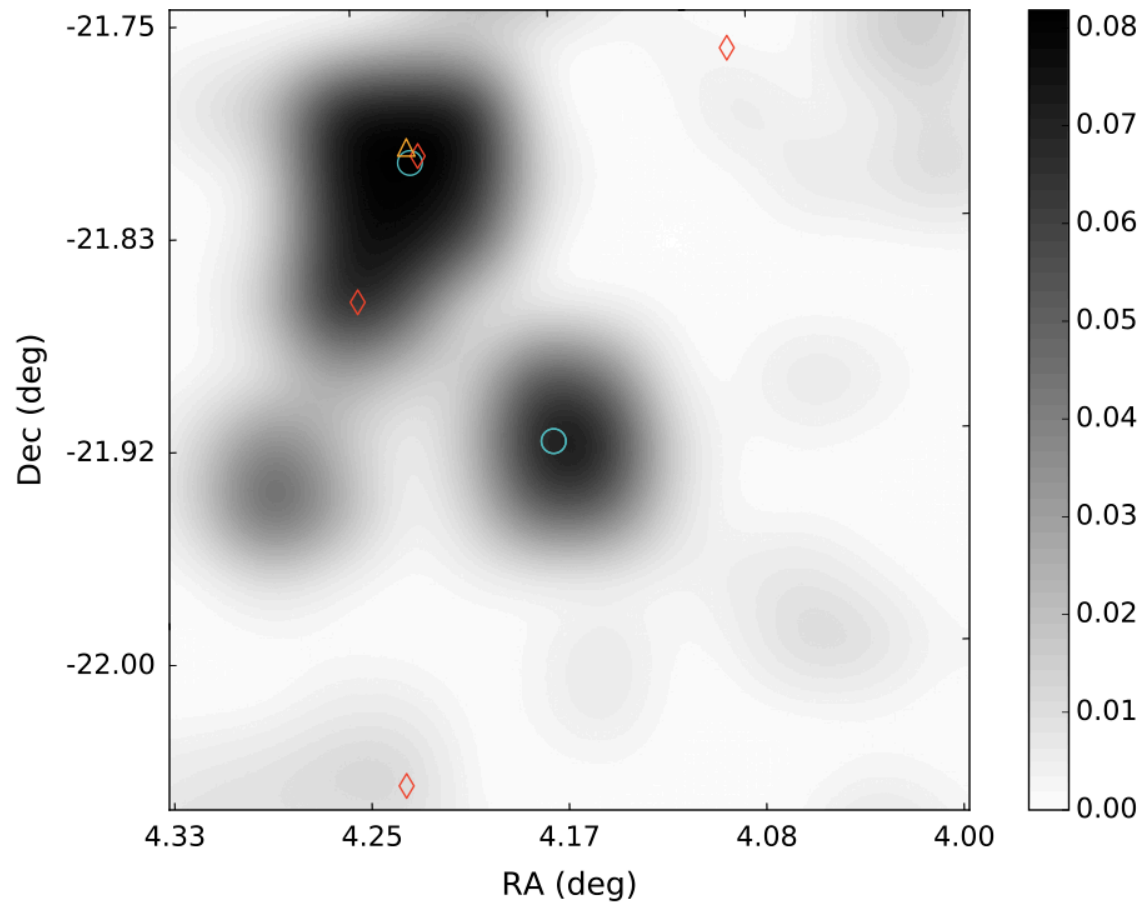
**KGS J232926-255814:** A 92 mJy source most closely matched to the galaxy 2dFGRS S128Z262 at a distance of  $40''$ . Five other galaxies and 4 UV sources are also located within the search radius.



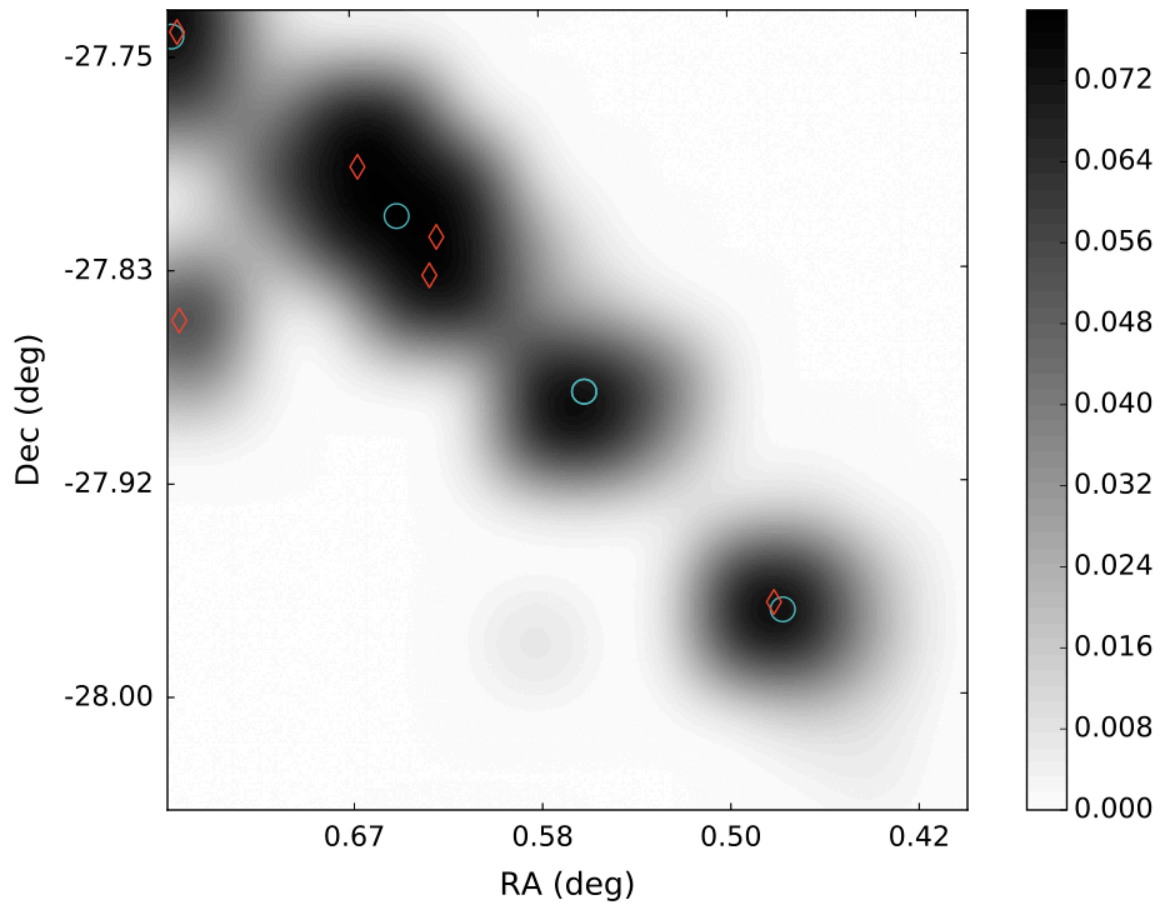
**KGS J234703-305612:** A 88 mJy source unmatched within  $30''$ . Seven galaxies are found within the search radius including two un-grouped 2dFGRS identifications at  $1'$  and  $2'$  from the source position. Two additional UV sources are also located within the search radius.



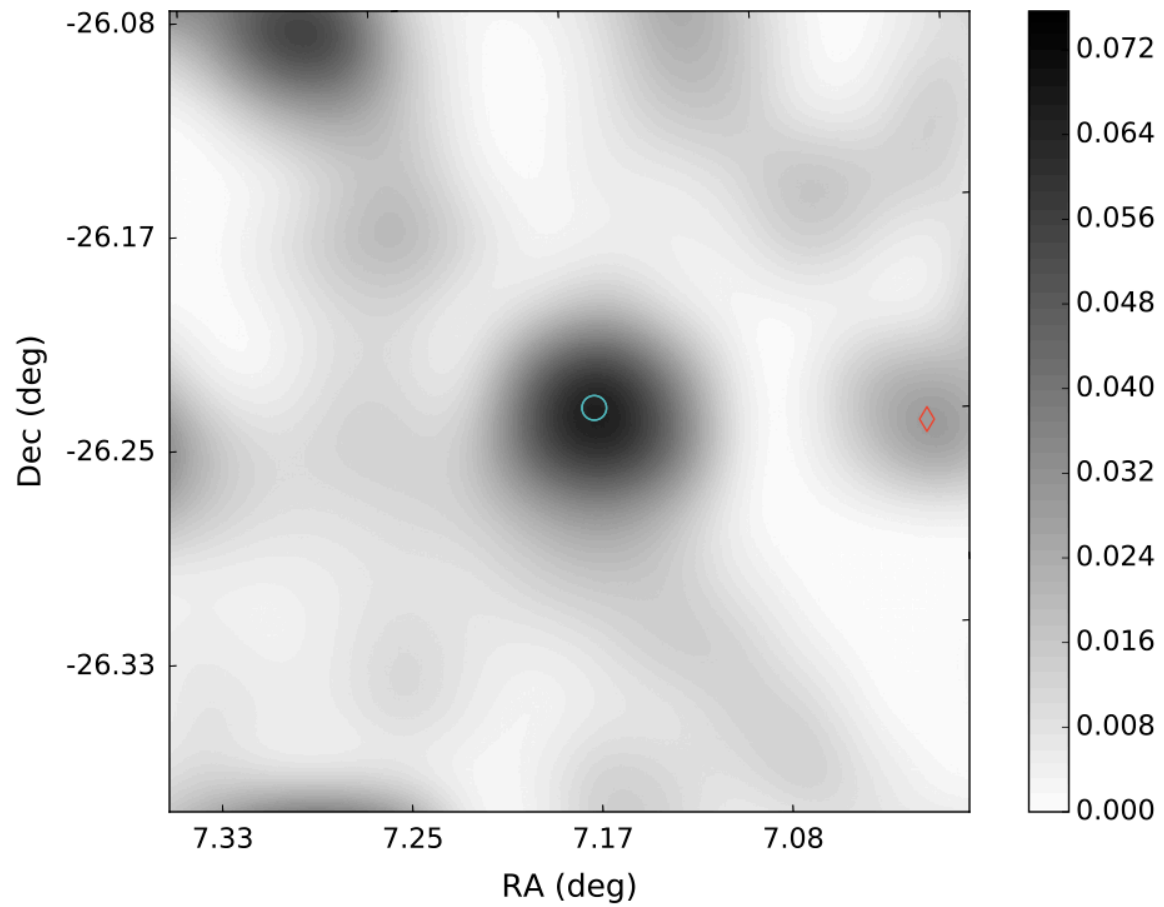
**KGS J241640-215455:** An 80 mJy source most closely matched to the UV source GALEX-ASC J001640.19-215516.7 at 21". Two optical galaxies and 26 UV sources are located within the search radius.



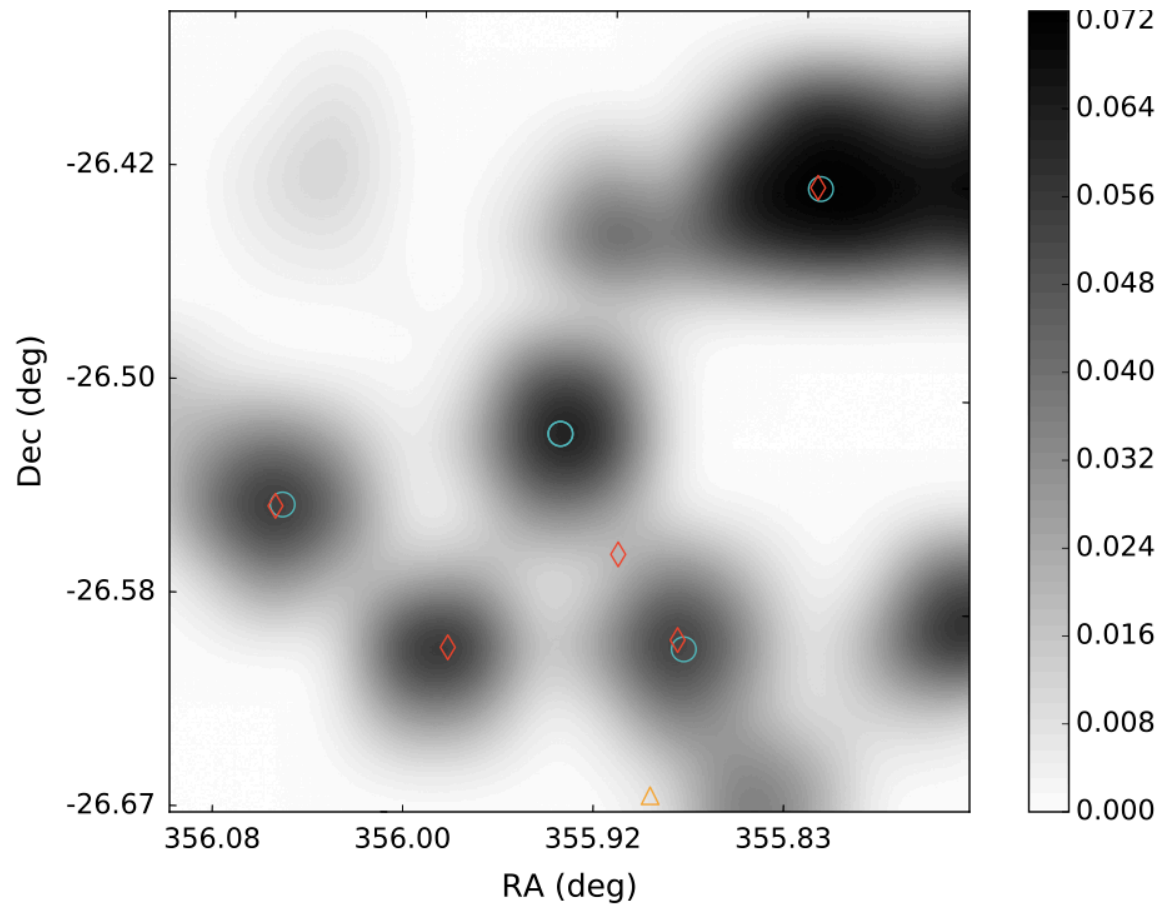
**KGS J000215-275242:** A 79 mJy source unmatched within 30". The galaxy 2dFGRS S198Z035 2' from the source position is part of the 32 member galaxy group 2PIGG SGP 2684. Five optical galaxies and 4 UV sources are found within the search radius.



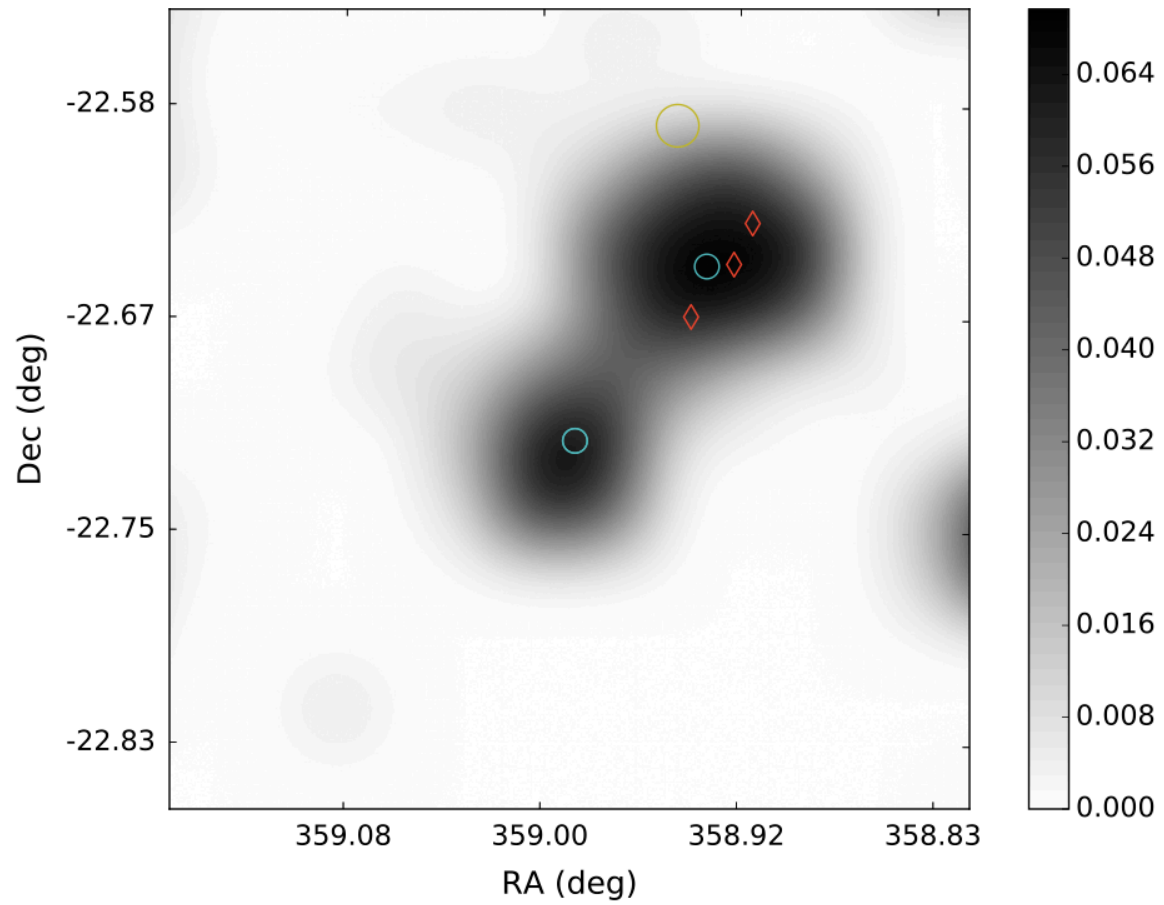
**KGS J002837-261426:** A 66 mJy source unmatched within  $30''$ . Four optical galaxies and many UV sources are found within the search radius.



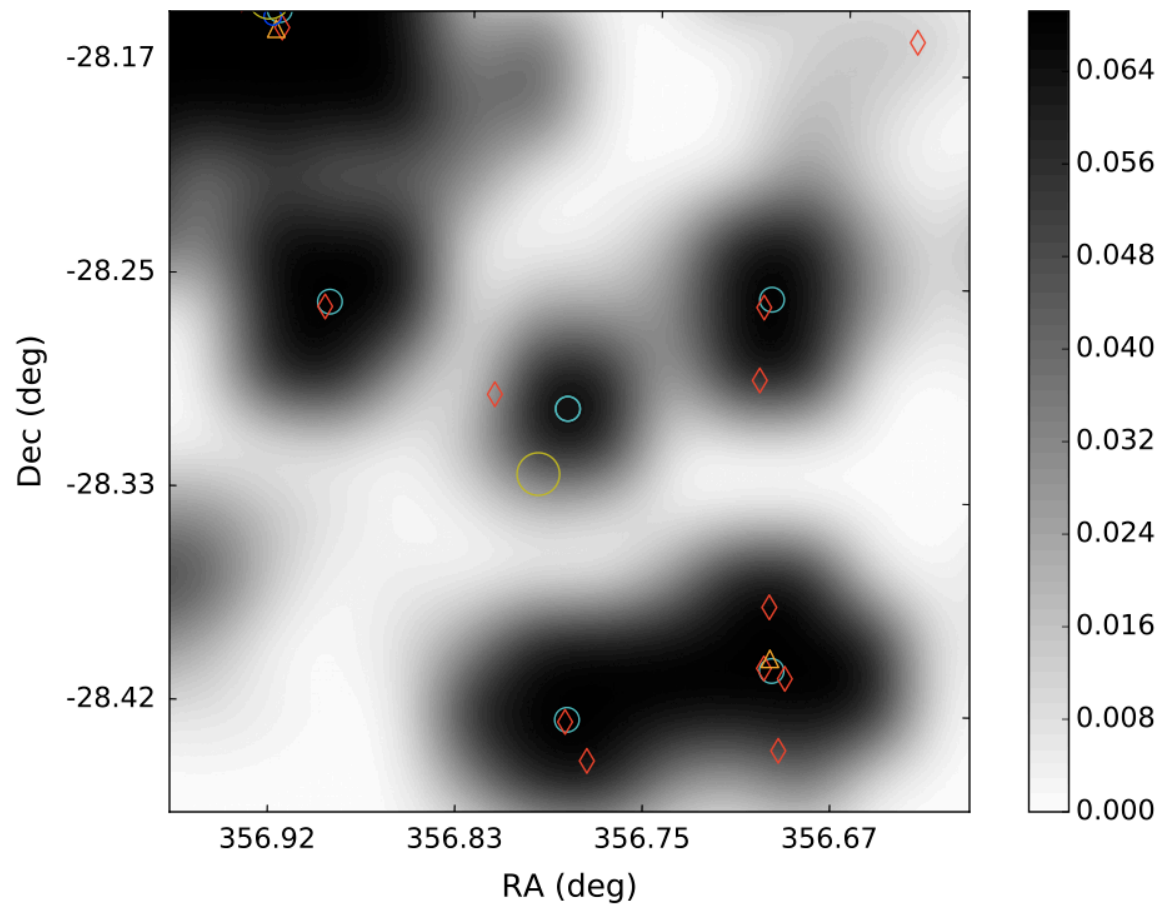
**KGS J234344-263049:** A 68 mJy source most closely matched to the galaxy 2dFGRS S194Z277 at 0.84'. There are 5 galaxies and 4 UV sources within the search radius.



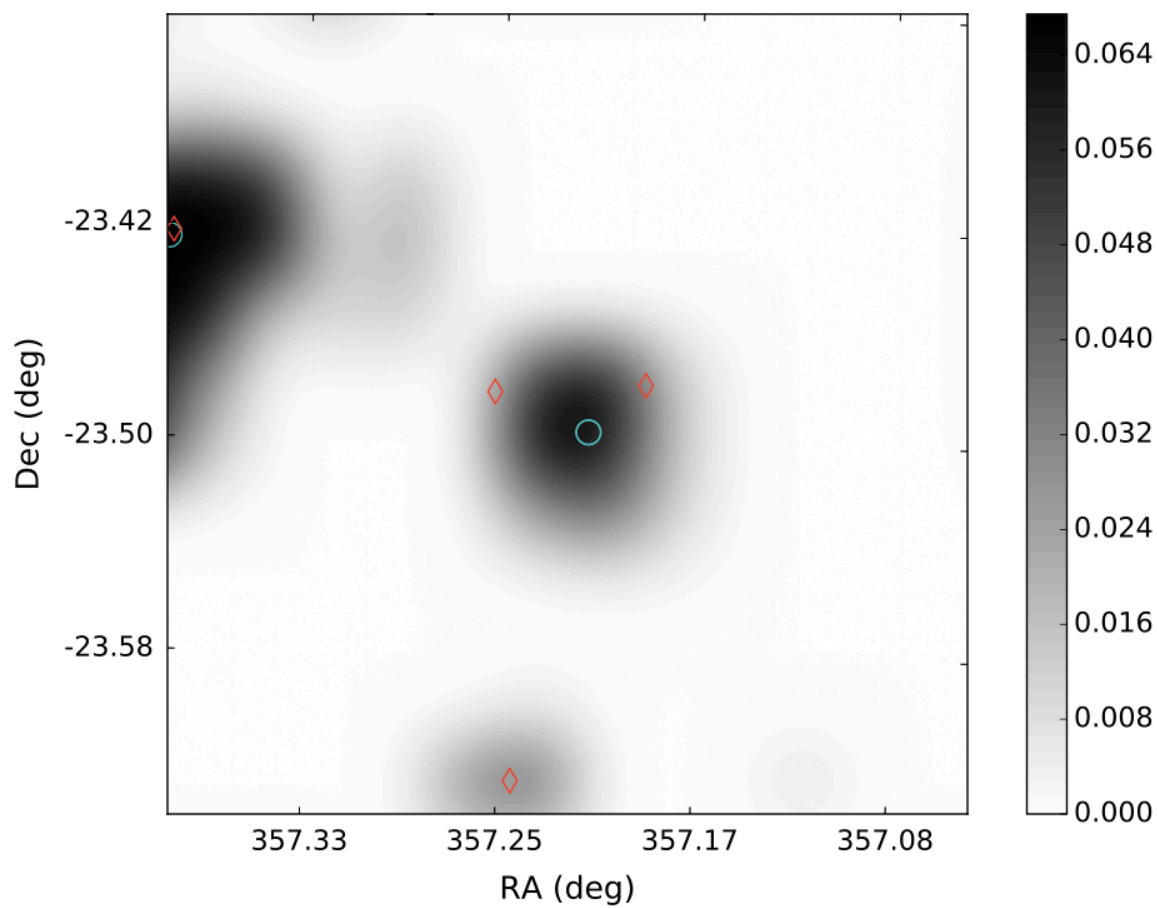
**KGS J235556-224242:** A 60 mJy source with 4 galaxies and 8 UV sources within the search radius.



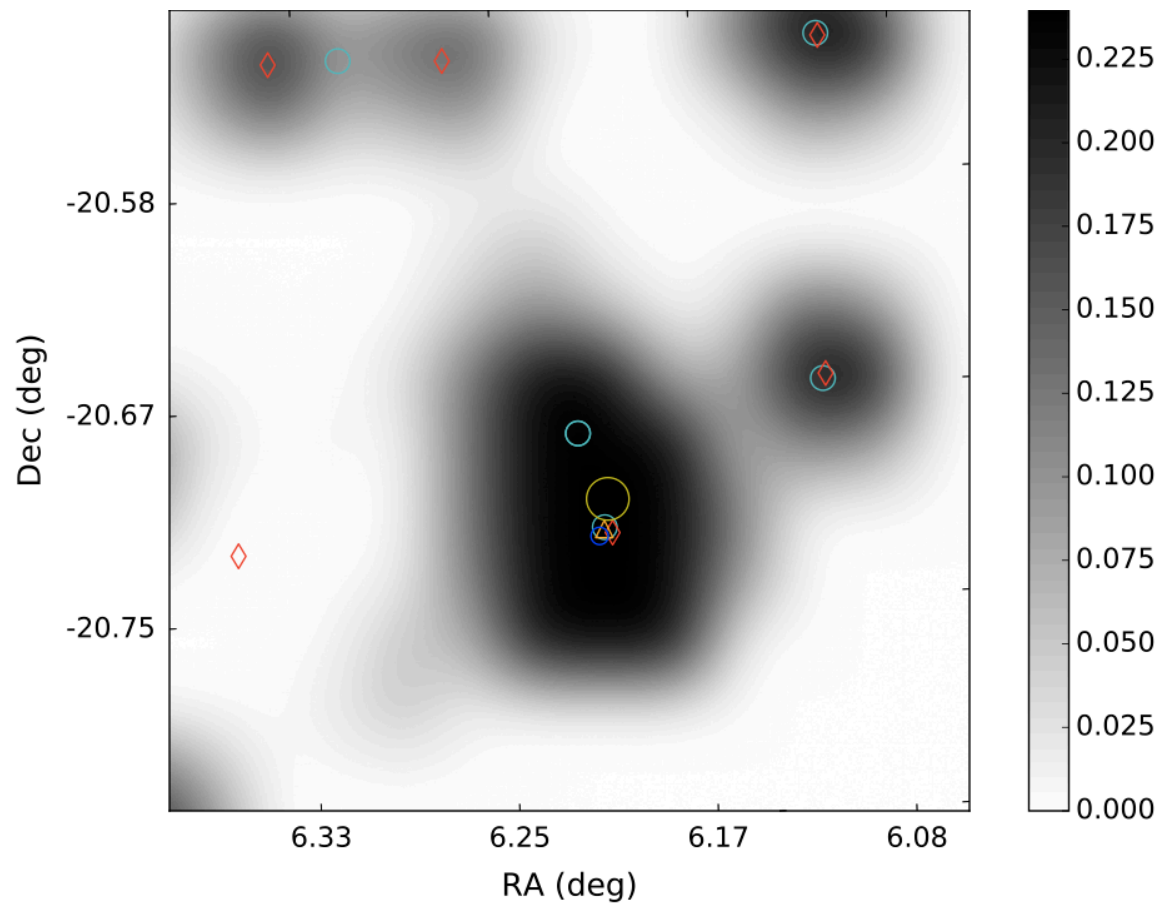
**KGS J234709-281746:** A 69 mJy source most closely matched to the X-ray source 2XMM J234709.2-281814 at a distance of  $27''$ . The radio source ABELL 4038:[SPS89] 06iii is 860 mJy at 1.5 GHz located  $1.735'$  from the source position. Abell 4037 is centered  $1.67'$  from the source position with a radius  $R_c = 59'$ . Numerous other sources including several cluster members and are also found.



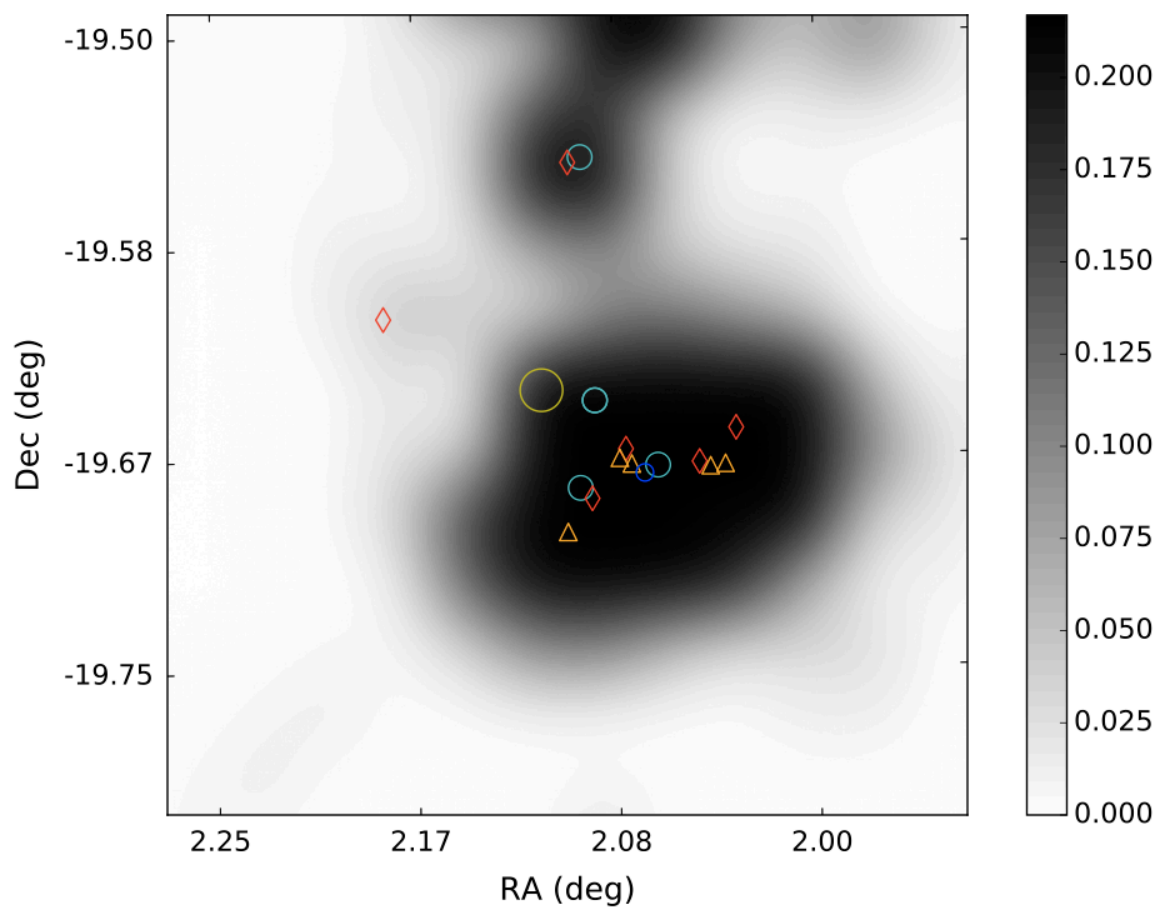
**KGS J234851-232934:** A 67 mJy source matched to two GALEXASC UV sources within  $30''$ . The radio source NVSS J234845-232827 at  $1.69'$  has a low probability of association. Two galaxies and 4 other UV sources are found within  $2.3'$ .



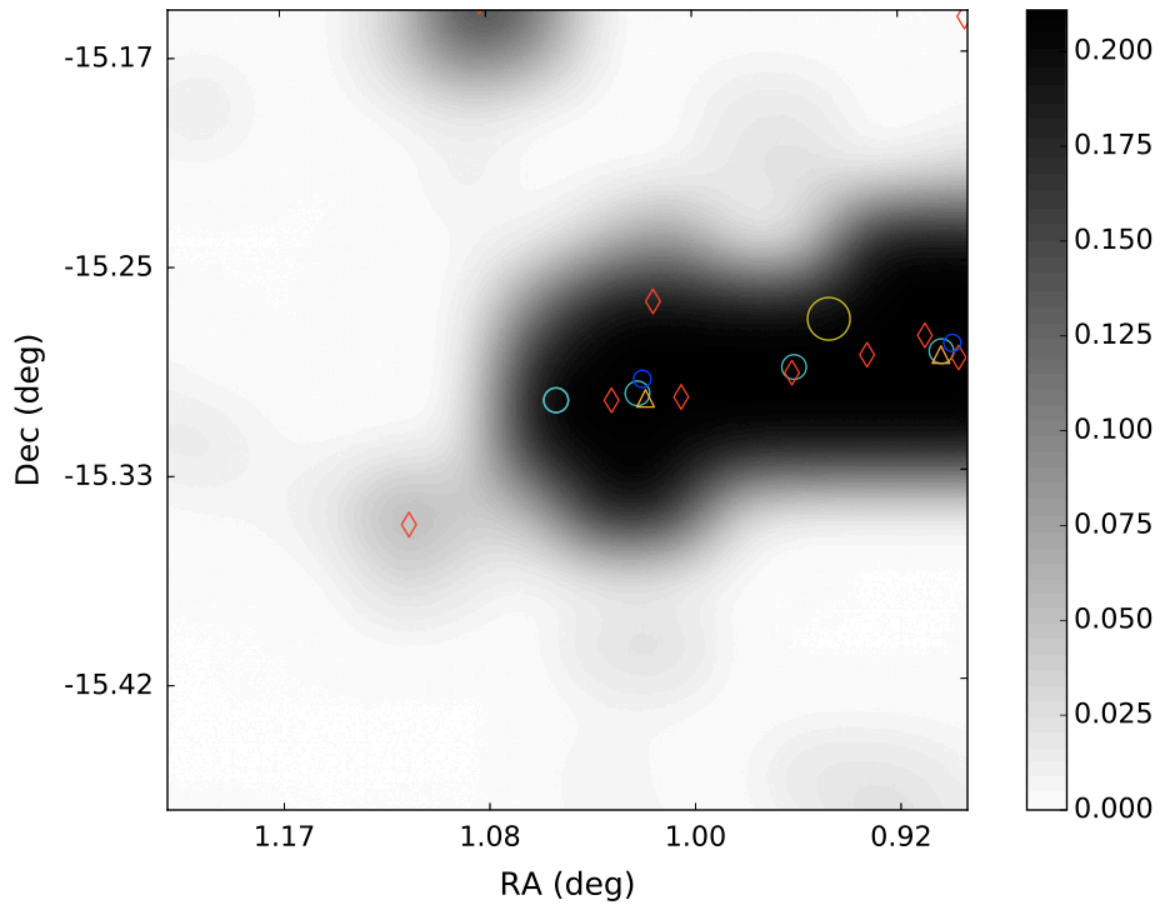
**KGS J002451-204048:** A 240 mJy confused source blended with a 2 Jy source just beyond the 2.3' search radius. It is coincident with the center of the galaxy cluster Abell 0027 at only 14" separation. Several cluster members and numerous other sources are found within the search radius.



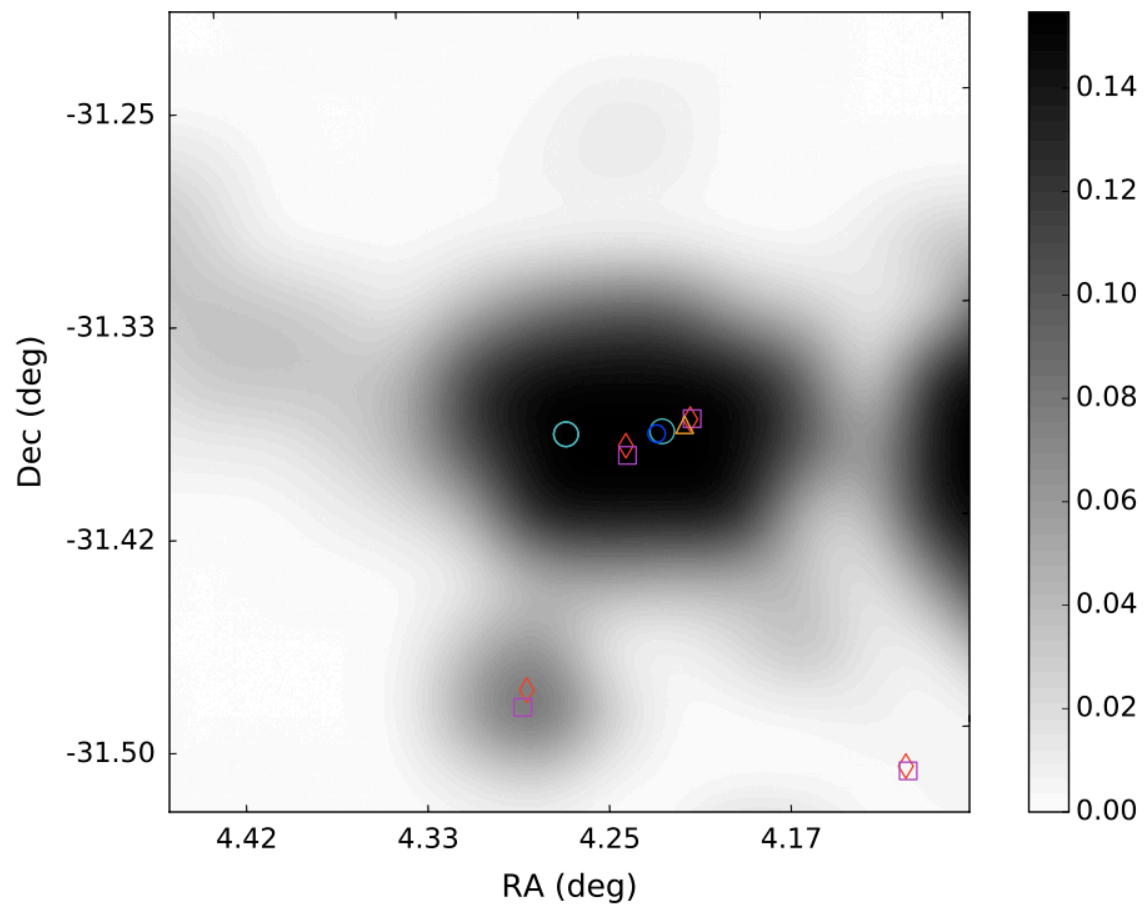
**KGS J000821-193833:** A 217 mJy confused source possibly associated with the galaxy cluster Abell 0002. It is most closely matched to GALEXASC and 2MASX identified galaxies at 7.4" and 8" respectively from the source position. Abell 0002 has a cluster radius  $R_c = 14'$  centered 1.2' from the source position. This source may be associated with the cluster member and radio source PKS 0005-199 but is not matched at 1.67' separation. Numerous other sources at all wavelengths are found within the search radius.



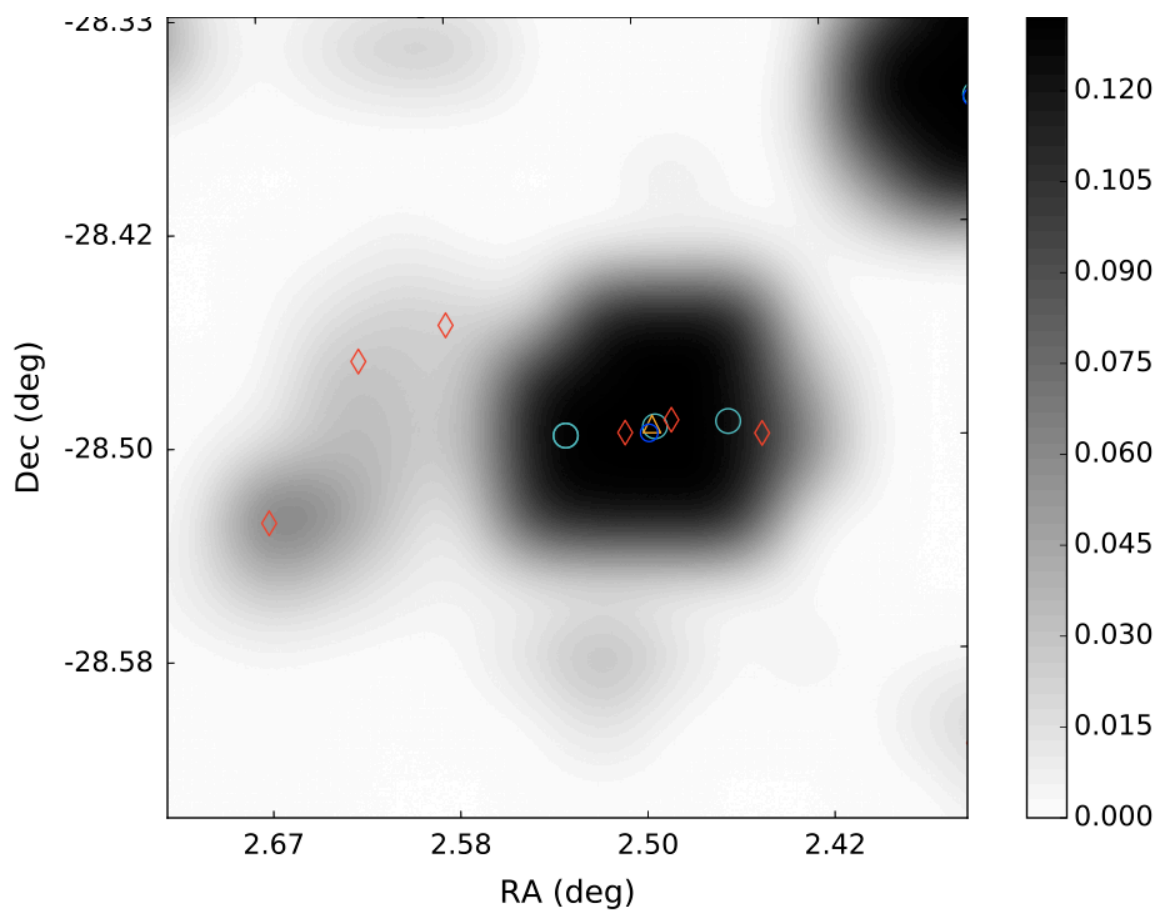
**KGS J000412-151811:** A 211 mJy confused source possibly associated with PMN J0004-1518 (cross identified in the NVSS and VLSS) but unmatched at 1.5'. Many galaxies and UV sources are found within the search radius.



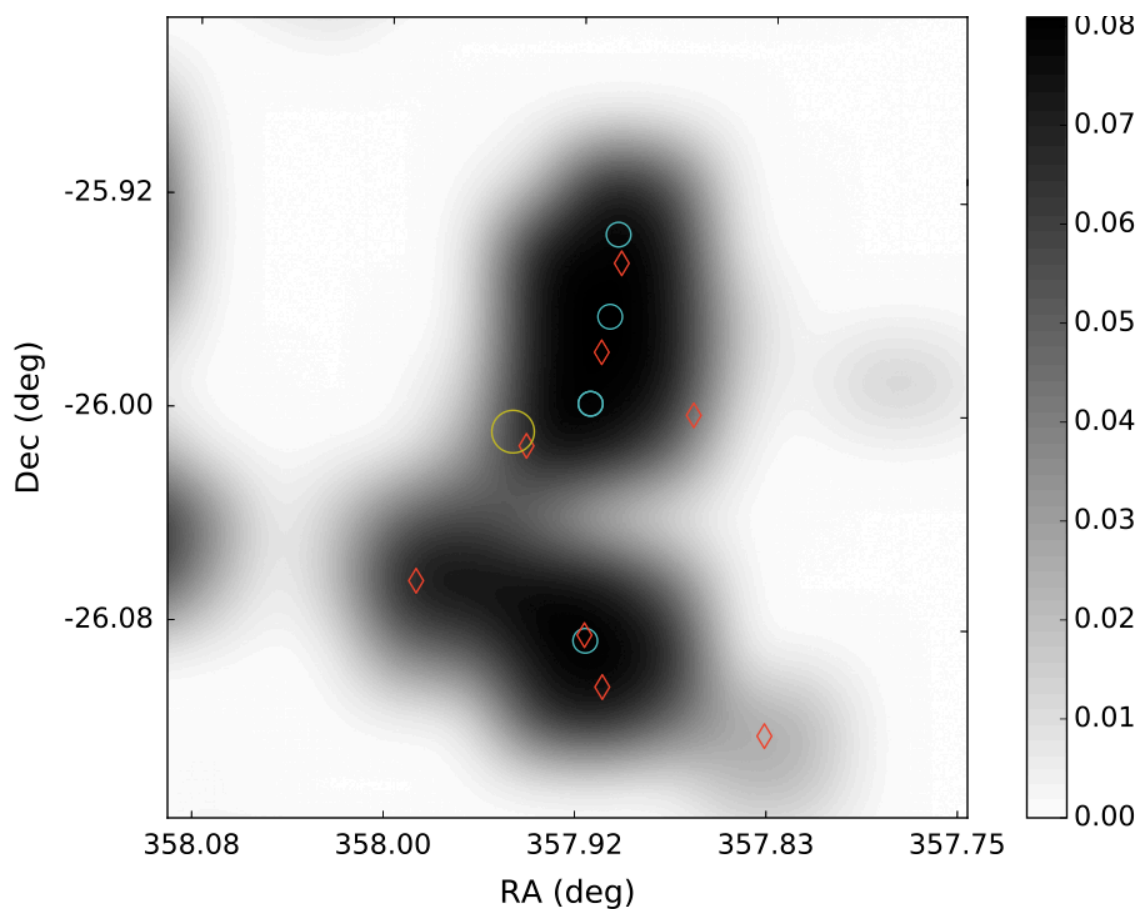
**KGS J001702-312239:** A 155 mJy confused source. It is most closely matched to 2dFGRS S436Z221, part of the 116-member galaxy group 2PIGG SGP 7135, at 1.2' from the source position. Several other sources are found within the search radius, including the 5 other 2PIGG member galaxies within 1.7'.



**KGS J001007-282942:** A 132 mJy confused source most closely matched to a UV source at  $27''$  separation. The galaxy 2dFGRS S279Z095 is found at  $0.79'$  and is part of the 2-member galaxy group 2PIGG SGP 8480. Two NVSS radio sources lie within the search radius but are not matched. The x-ray source 1WGA J0009.9-2829 and numerous galaxies and UV sources are also found.



**KGS J235139-255937:** An 81 mJy confused source closely matched to 2dFGRS S132Z149 at 22" separation. This is cross identified as an X-ray and extended IR source, and is a member of the galaxy cluster Abell 2667. The nearest radio source is NVSS J235145-260038 at 1.8' from the source position. Many other sources at all wavelengths are found within the search radius.



## 7.2 Galaxy Clusters

Among the 25 new radio detections, seven are cross-matched to galaxy clusters, and another four to galaxy groups. We flag two others with several possible galaxy counterparts that are also cross-matched to a 2MASS extended IR source (2MASX), commonly cross-identified with clusters and groups. These represent approximately half of all new detections and are among the brightest. Only one source of confused or extended emission is not cross-matched to a known cluster or group, but appears to occupy a high density environment 6.6' from the center of Abell 2699. Seven sources find an Abell galaxy cluster [1] within 2'. This represents 28% compared to 1.1% (82) of all matched catalog sources.

The distribution of cluster member counts for sources within 2' of a cluster center is shown in Figure 7.1. The mean cluster count is 73 for the 7 new detections compared to 60 for matched sources. To test the significance of this difference, we make  $10^4$  random draws of seven from the 82 source sample, and calculate the mean cluster member count for each draw. The probability that a mean of 73 can be explained by random chance is only 18%.

These findings suggest that we are preferentially finding sources associated with galaxy clusters, and that these clusters are richer on average. This is consistent with previous findings that USS sources at low redshift are found in high density environments and that richer clusters tend to host steeper spectrum radio sources [6]. Figure 7.2 shows the spectral index distribution for all `isolated` sources within 5' of an Abell cluster, compared to sources more than  $1^\circ$  from a cluster center. The distribution is steeper on average but also broader near clusters.

We further find evidence of a correlation between spectral curvature and cluster richness. Figure 7.3 shows the estimated curvature versus spectral index and cluster member count for sources within 5' of a cluster center. Positive curvature is almost strictly limited to sources with steeper than average spectra. The pattern is constant regardless of which catalogs a source was matched. Positive curvature can be explained if the source has multiple components with different spectral index. The steeper spectrum component will dominate

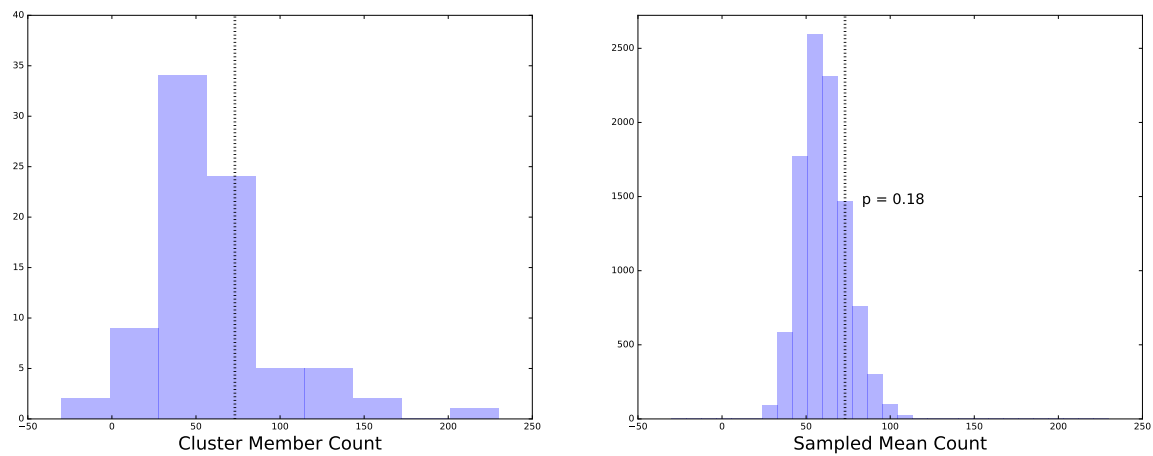


Figure 7.1: Left: The distribution of cluster member counts for the 82 matched KGS sources within  $2'$  of an Abell galaxy cluster center. The 7 unmatched KGS radio sources have a mean cluster member count of 73, indicated by the black dotted line. Right: The distribution of mean cluster member counts for  $10^4$  random draws of 7 from the 82 source sample. The probability that the mean of the new detections can be explained by random change is 18%.

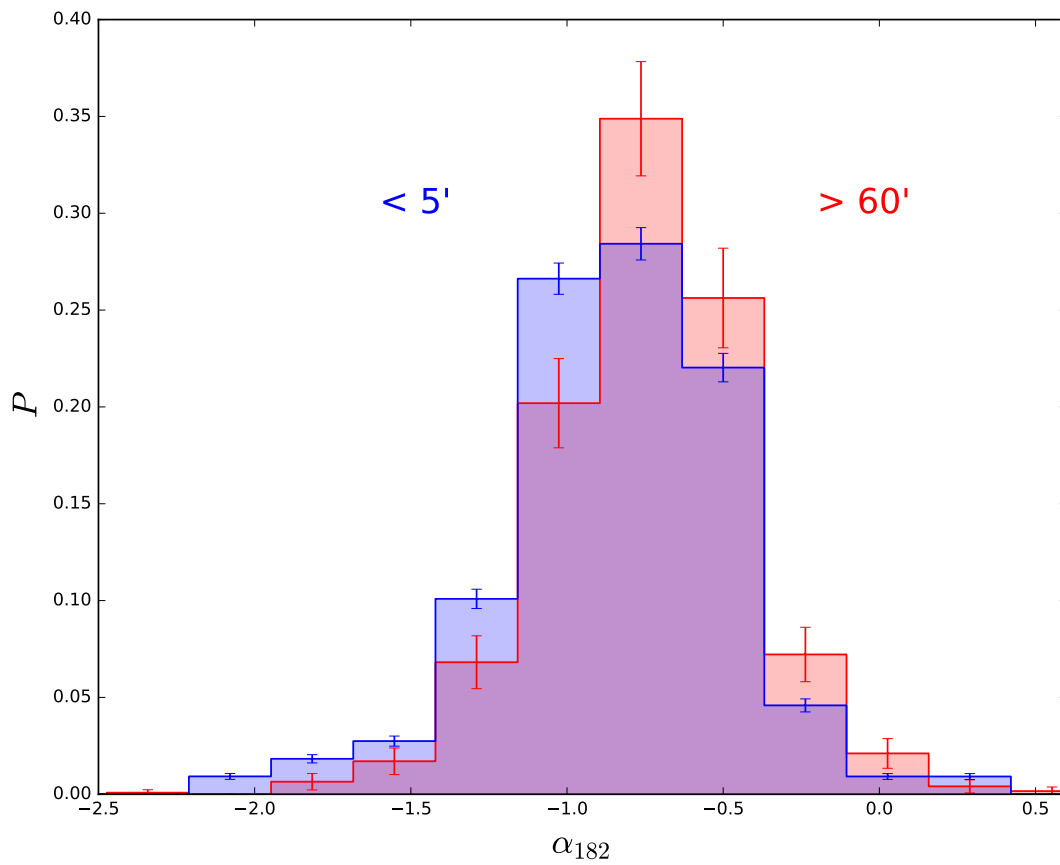


Figure 7.2: The spectral index distribution for all **isolated** sources within 5' of an Abell cluster, compared to sources more than  $1^\circ$  from a cluster center. The distribution is steeper on average but also broader near clusters.

at low frequencies while the shallower spectrum will dominate at higher frequencies. This is consistent with a source in a high density environment such as a galaxy cluster.

### 7.3 *HzRG Candidates*

New detections that are not associated with a galaxy cluster or group can still be presumed to be ultra-steep spectrum, and are candidate HzRGs. HzRGs are faint, subtend a small angular scale, and have been shown to exhibit single power-law spectra across to  $<100$  MHz with no indication of curvature. Broad band spectral curvature measurements allow for a more robust selection of HzRG candidates. There are 96 candidates, including 18 new detections, that meet the following criteria:

- The 182 MHz estimated spectral index or limit is less than -1.4.
- The PUMA match is not classified as `multiple`.
- The angular separation from the nearest Abell cluster center is  $> 2'$ .
- The estimated curvature is small, within the interquartile range of sources in 3 or more catalogs.

Of the 96 HzRG candidates, 91 were only detected in one other survey and therefore have no measure of curvature. Among sources detected in three or more catalogs, 92 met the first three criteria but only 5 met the curvature selection. This is not surprising given that a source must be relatively bright to be detected in more than two surveys (with the exception of the few sources with both SUMSS and NVSS cross matches), precluding their candidacy as high redshift radio galaxies. This underscores the usefulness of curvature measurements in selecting specifically for relatively rare HzRGs among USS sources.

The candidates are searched for counterparts in NED within a  $30''$  radius. We exclude sources matched to a galaxy cluster or group member, or matched to more than one galaxy or UV source. Known correlations between redshift and k-band or r-band magnitude are

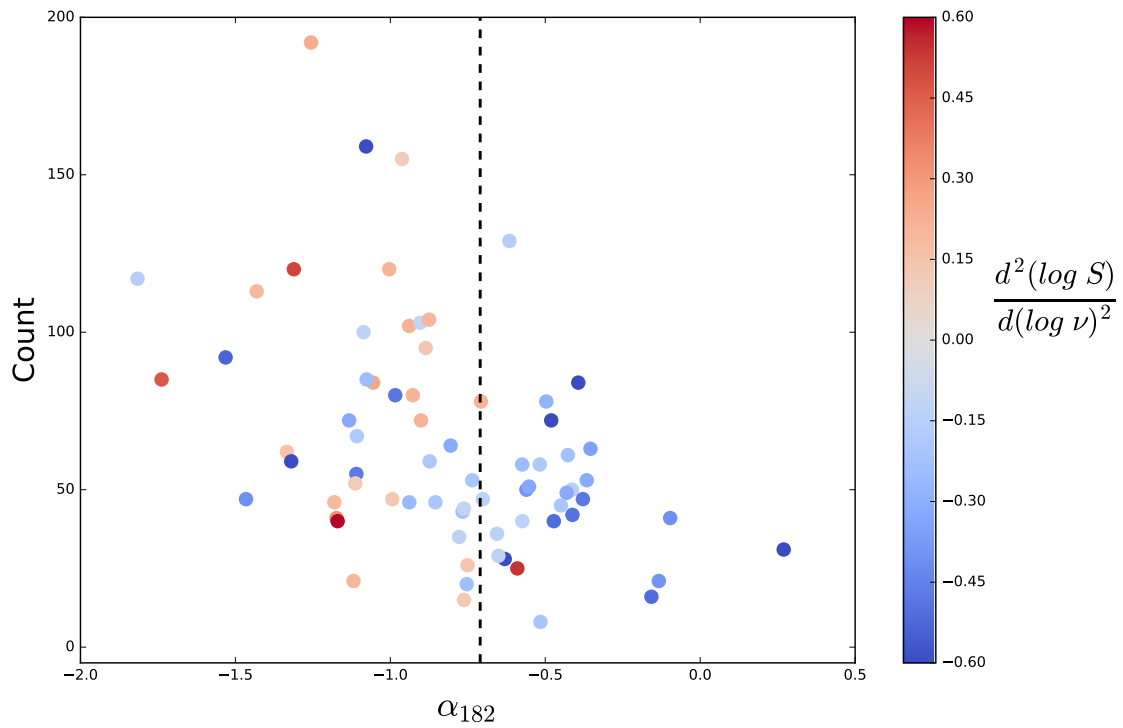


Figure 7.3: The estimated curvature versus spectral index and cluster member count for sources within 5' of an Abell cluster center. Positive curvature is almost strictly limited to sources with steeper than average spectra. The pattern is constant regardless of which catalogs a source was matched, and does not appear to be caused by a selection effect.

used to further exclude sources matched to a 2MASS extended IR sources or a galaxy with r-band magnitude  $< 19$  mag.

The remaining 73 HzRG candidates are listed in Table 7.2. We note the source KGS J000419-310303 is  $25''$  from the QSO 2QZ J000417.8-310317 at  $z = 1.155$ . It is detected in the NVSS at 2.7 mJy with spectral index  $\alpha = -1.6$ . This list may serve as a starting point for further investigation and eventual follow up spectroscopy.

Table 7.2: A conservative selection of high redshift radio galaxy candidates. Columns are KGS name, flux density and error, reliability class, number of detections, estimated 182 MHz spectral index, distance to the nearest Abell galaxy cluster, average beam response of all detections, Eddington bias correction applied to the flux density, and NED cross-match catalog, type, and photometry where available.

Name	$S$ (mJy)	$\sigma_S$ (mJy)	$R_{\text{class}}$	$N_{\text{det}}$	$\alpha_{182}$	$d_{\text{gc}}$ (arcmin)	Beam	$f_{\text{EB}}$	Match ( $< 30''$ )	Type	phot
KGS J001702-312239	153.0	5.1	3	63		10.5	0.74	0.989			
KGS J231928-302751	114.2	20.8	6	19		73.5	0.61	0.98			
KGS J000215-275242	77.8	14.9	4	44		12.0	0.94	0.985			
KGS J232803-145208	249.3	32.7	4	30		124.4	-0.75	0.925			
KGS J001007-282942	132.2	4.6	4	37		27.6	0.91	1.0	GALEXMSC	UvS	
KGS J001640-215455	79.6	12.9	8	7		59.0	0.84	0.973	GALEXASC	UvS	
KGS J000412-151811	210.9	40.2	6	13		6.7	0.39	1.0	GALEXMSC	G	20.39 b_j
KGS J002837-261426	65.9	16.4	8	5		78.7	0.85	0.871			
KGS J235156-165850	132.6	15.2	6	14		91.9	0.51	0.975			
KGS J233116-192443	112.7	21.1	6	11		46.5	0.59	0.953	GALEXASC	UvS	
KGS J234703-305612	87.6	16.3	4	32		31.7	0.81	0.975			
KGS J235021-194846	122.5	19.8	2	57		37.5	0.69	1.0			
KGS J232926-255814	91.8	8.3	6	17		16.4	0.8	0.986			
KGS J235556-224242	59.2	13.9	8	5		7.8	0.88	0.849			
KGS J234344-263049	68.4	8.9	6	14		38.6	0.87	0.938			

Name	$S$ (mJy)	$\sigma_S$ (mJy)	$R_{\text{class}}$	$N_{\text{det}}$	$\alpha_{182}$	$d_{\text{gc}}$ (arcmin)	Beam	$f_{\text{EB}}$	Match ( $< 30''$ )	Type	phot
KGS J002001-255255	80.4	16.4	4	29	-2.211	14.1	0.89	0.969			
KGS J235021-352635	131.1	22.8	6	29	-1.963	32.4	0.52	0.955	GALEXASC	UvS	
KGS J232100-163405	221.8	23.5	5	26	-1.842	26.2	0.29	1.0			
KGS J000619-174112	100.9	19.3	8	7	-1.835	14.8	0.56	0.889	NVSS	RadioS	
KGS J235124-354628	167.1	26.2	3	50	-1.819	45.3	0.49	0.999	GALEXASC	UvS	
KGS J002612-222223	138.5	18.9	3	61	-1.752	45.2	0.7	1.0	GALEXASC	G	19.09 b_J
KGS J234115-350444	111.5	39.0	6	21	-1.698	85.8	0.44	0.759	NVSS	RadioS	
KGS J234602-183855	124.0	18.6	3	28	-1.673	68.4	0.6	0.99	NVSS	RadioS	
KGS J002055-241233	73.8	11.4	8	6	-1.661	80.4	0.92	0.963			
KGS J002145-222145	90.1	14.8	6	16	-1.623	33.5	0.79	0.985	NVSS	RadioS	
KGS J000419-310303	70.9	11.5	6	10	-1.603	34.2	0.85	0.922	2QZ ( $z = 1.155$ )	QSO	20.66
KGS J232830-185047	130.8	25.3	6	19	-1.592	36.4	0.52	0.958	GALEXASC	UvS	
KGS J234412-221718	73.2	11.0	8	5	-1.584	25.1	0.82	0.913	GALEXASC	UvS	
KGS J000756-191013	89.9	15.0	6	14	-1.579	29.0	0.66	0.907	NVSS	RadioS	
KGS J232014-290409	112.1	8.7	6	22	-1.567	76.4	0.67	0.988	GALEXMSC	UvS	

(Table 7.2 continued.)

Name	$S$ (mJy)	$\sigma_S$ (mJy)	$R_{\text{class}}$	$N_{\text{det}}$	$\alpha_{182}$	$d_{\text{gc}}$ (arcmin)	Beam	$f_{\text{EB}}$	Match ( $< 30''$ )	Type	phot
KGS J000428-294036	83.4	7.8	8	5	-1.556	26.8	0.89	0.995	NVSS	RadioS	
KGS J233153-250511	68.6	17.4	8	4	-1.552	25.3	0.85	0.888	GALEXASC	UvS	
KGS J000022-280034	142.1	18.9	0	71	-1.541	5.7	0.94	1.0	NVSS	RadioS	
KGS J002347-293925	96.8	22.3	3	37	-1.539	71.4	0.77	0.934	MRSS	G	19.3 r
KGS J001219-253108	67.7	9.0	6	21	-1.529	49.5	0.95	0.976	NVSS	RadioS	
KGS J000817-192907	204.6	27.3	0	71	-1.506	9.5	0.67	0.996	NVSS	RadioS	
KGS J002615-255216	265.2	5.2	0	69	-1.499	75.4	0.77	1.0			
KGS J230854-245128	152.9	11.8	5	23	-1.493	27.0	0.52	0.995			
KGS J001835-351610	132.1	17.3	8	7	-1.483	14.2	0.51	0.967	GALEXASC	UvS	
KGS J002642-252424	79.3	8.2	6	21	-1.478	83.0	0.85	0.943	NVSS	RadioS	
KGS J234723-264739	85.2	13.1	3	52	-1.477	29.3	0.92	0.996	NVSS	RadioS	
KGS J002347-291530	75.1	21.2	6	15	-1.477	72.3	0.77	0.851	NVSS	RadioS	
KGS J230459-235458	143.8	24.7	8	9	-1.476	14.1	0.49	0.972	NVSS	RadioS	
KGS J234555-183824	115.2	9.0	8	3	-1.475	69.0	0.59	0.971			
KGS J235638-384553	219.1	29.0	6	20	-1.474	55.2	0.3	0.971	GALEXASC	UvS	

(Table 7.2 continued.)

Name	$S$ (mJy)	$\sigma_S$ (mJy)	$R_{\text{class}}$	$N_{\text{det}}$	$\alpha_{182}$	$d_{\text{gc}}$ (arcmin)	Beam	$f_{\text{EB}}$	Match ( $< 30''$ )	Type	phot
KGS J010439-254046	251.5	48.6	5	16	-1.473	66.1	0.34	0.987	NVSS	RadioS	
KGS J000354-171710	148.5	29.5	4	42	-1.472	35.7	0.53	0.991	NVSS	RadioS	
KGS J000912-172635	116.2	22.9	7	2	-1.471	31.4	0.56	0.891	GALEXASC	G	19.57 b_J
KGS J002451-254816	73.6	18.4	6	9	-1.467	56.3	0.88	0.892	GALEXASC	UvS	
KGS J000159-210427	85.3	10.9	7	2	-1.466	18.6	0.79	0.961			
KGS J002054-253943	64.4	11.5	8	5	-1.458	3.7	0.95	0.903	MRSS	G	18.9 r
									GALEXASC	UvS	
KGS J002802-262846	89.2	15.3	6	28	-1.465	62.3	0.79	0.974	MRSS	G	19.5 r
KGS J232309-260325	111.1	18.2	4	41	-1.457	11.8	0.69	1.0	GALEXASC	UvS	
KGS J233055-222411	86.7	16.2	8	7	-1.452	33.9	0.71	0.943	NVSS	RadioS	
KGS J231255-293331	128.5	9.1	8	6	-1.45	32.0	0.56	0.985	NVSS	RadioS	
KGS J234011-180428	135.5	24.1	6	26	-1.44	133.6	0.53	0.978	NVSS	RadioS	
KGS J233425-251804	84.4	13.8	6	19	-1.439	29.3	0.82	0.975	GALEXMSC	UvS	
									MRSS	G	19.7 r
KGS J232321-275300	78.4	19.5	8	6	-1.436	98.5	0.77	0.892	NVSS	RadioS	

(Table 7.2 continued.)

Name	$S$ (mJy)	$\sigma_S$ (mJy)	$R_{\text{class}}$	$N_{\text{det}}$	$\alpha_{182}$	$d_{\text{gc}}$ (arcmin)	Beam	$f_{\text{EB}}$	Match ( $< 30''$ )	Type	phot
KGS J232437-243816	124.5	16.3	4	41	-1.434	31.7	0.7	0.99	MRSS	G	19.2 r
									GALEXMSC	UvS	
KGS J003140-273444	89.9	17.8	6	22	-1.428	8.0	0.76	0.975	NVSS	RadioS	
KGS J001558-260853	82.4	9.3	6	8	-1.427	12.6	0.95	0.996			
KGS J000950-282926	121.9	21.9	6	18	-1.424	30.2	0.87	1.0	GALEXMSC	UvS	
KGS J002607-263520	54.5	21.3	8	5	-1.422	48.2	0.84	0.725	MRSS	G	19.2 r
KGS J000845-300733	299.0	19.4	0	71	-1.418	56.2	0.85	1.0	NVSS	RadioS	
KGS J002306-271628	67.9	12.9	6	12	-1.415	36.5	0.9	0.892	NVSS	RadioS	
KGS J001938-193851	104.8	16.5	6	15	-1.412	32.8	0.64	0.955	MRSS	G	19.1 r
KGS J000103-293209	85.1	13.5	4	41	-1.411	32.1	0.9	0.998			
KGS J235807-263706	141.5	16.3	0	71	-1.41	64.3	0.96	1.0	NVSS	RadioS	
KGS J000922-245728	59.3	16.8	6	13	-1.403	100.8	0.95	0.848	NVSS	RadioS	
KGS J232725-163129	174.0	35.3	8	11	-1.402	50.0	0.38	0.97	NVSS	RadioS	

(Table 7.2 continued.)

## Chapter 8

### EOR FOREGROUND MODELING AND REMOVAL

The purpose of the KGS EoR field survey was to create a more accurate and reliable sky model for improved calibration and foreground power removal. Before this catalog was complete, the sky model was built from the MWA Commissioning Survey. As discussed in §1 and revealed in §6.5, there are known problems with this catalog including large measurement errors, biases, and lack of coverage at the northern extent of the field. In this chapter, we compare the KGS catalog with the MWACS catalog in terms of foreground power removal within the primary beam. We then motivate the need for a combined master catalog in order to subtract power from widefield sources in the side lobes of the primary beam. Finally, we address the fact that not all sources can be accurately modeled as point sources, and demonstrate a simple approach to build an extended model for NGC 253.

#### **8.1 Sky Model Tests**

We compare the two catalogs as standalone sky models input to FHD firstpass for both calibration and foreground subtraction. FHD firstpass takes an input source list containing coordinates, flux densities, frequency, and spectral index. The spectral index is used to extrapolate the flux density to 182 MHz. Each source is modeled as a point source, or delta function, to generate model visibilities as seen by the instrument. This visibility model is used first for direction-dependent calibration, establishing the overall flux scale and mapping of the observed snapshot (dirty) visibilities, and then subtraction to remove the foreground power.

Visibilities are gridded in  $u$ ,  $v$ , and frequency ( $uvf$ ) cubes before generating the power spectrum. By differencing the dirty, model, and residual cubes, we can compare two pro-

cessing runs. As a rule, we subtract the “test” cube from the “standard” cube such that a positive difference in the residual indicates improved foreground power subtraction. We use the 2D  $k$ -power spectrum introduced in §1.2 to visualize these differences.

In some cases, the calibration and overall flux scale will change significantly between two processing runs. If the flux scale is not comparable between the two runs, the residual difference becomes complicated and cannot be interpreted strictly in terms of foreground power removal. It is therefore helpful to normalize by taking the ratio of the residual to the dirty cubes before differencing. Both the direct difference and ratio difference power spectra are useful figures of merit (FoM) when assessing any change made within the FHD- $\epsilon$  pipeline.

### 8.1.1 Primary Beam

The first test is a direct comparison between the new KGS catalog and the original MWACS catalog, limiting calibration and source subtraction to the primary beam. The KGS catalog covers the full EoR0 field to 5% average beam power. It is complete to  $\sim 80$  mJy within 50% beam power, but this limit falls off to  $\sim 1$  Jy at the edge of the field. The MWACS covers nearly the full field, but lacks coverage north of  $-14^\circ$ . The completeness limit is variable and several bright sources are absent. Overall, MWACS is shallower than the KGS but it is deeper near the edge of the field, where the beam sensitivity and KGS completeness fall off.

The difference and ratio difference power spectra are shown in Figures 8.1–8.2. A very definitive improvement in foreground power subtraction is seen in both FOMs. This is observed in the residual wedge below the dashed line, marking the extent of the primary beam. All modes except the lowest  $k_\perp$  modes are positive, indicating improvement. The negative residual modes at low  $k_\perp$  can be explained by the differences at the edge of the primary beam, where MWACS is the deeper and more complete catalog.

The direct difference of the model visibilities shows more power in the KGS catalog. There are 7303 sources with a total 182 MHz flux density of 3631 Jy in the KGS sky model, compared to 2525 sources and 2426 Jy in the MWACS sky model. Although more power is

subtracted with the KGS, the difference of the dirty visibilities shows that there is also an overall flux scale change. Among 1881 isolated matched sources between the two catalogs, the MWACS is seen to systematically underestimate source flux density compared to the KGS (Figure 8.3). This bias is even more evident when considering sources only within half-beam. The overall flux scale is therefore lower using MWACS compared to KGS, resulting in a negative difference of the dirty cubes. The fact that the residual foreground wedge difference is positive tells us that the foreground model improvement outweighs the flux scale discrepancy. The ratio difference normalizes the flux scales and confirms the increase in relative foreground power removal using the KGS.

### 8.1.2 *Widefield*

We have found that the KGS is an unequivocal improvement for foreground power removal in the MWA EoR0 field. However, the KGS is still limited in the spatial extent and depth of its coverage. Recently, it has become clear that widefield foreground sources – particularly those occupying the side lobes of the primary beam – are also important to model and subtract. Sources occupy higher  $k_{\parallel}$  modes in the foreground wedge the further they are from beam center. Those furthest afield are especially problematic as they encroach on the EoR window [38]. This is illustrated in Figure 8.4.

To gain widefield coverage in our sky model in the absence of targeted snapshot observations of the side lobes, we can fill in the blanks as best we can with other catalogs. Before the KGS catalog was finalized, a precursor set of source candidates was selected and combined with the MWACS and MRC catalogs to create the first version of a “master” foreground catalog. The KGS precursor catalog was composed of 2993 sources detected in at least 90% of snapshots in which the flux density was expected to lie above the detection threshold. The three catalogs were associated by clustering, as in §3.2, using a 3’ neighborhood radius corresponding to the approximate MWACS PSF FWHM.

For each cluster the KGS, MWACS, or MRC source was kept in that order of preference. The Culgoora Circular Array (CCA [35]) survey is a targeted survey on select MRC sources

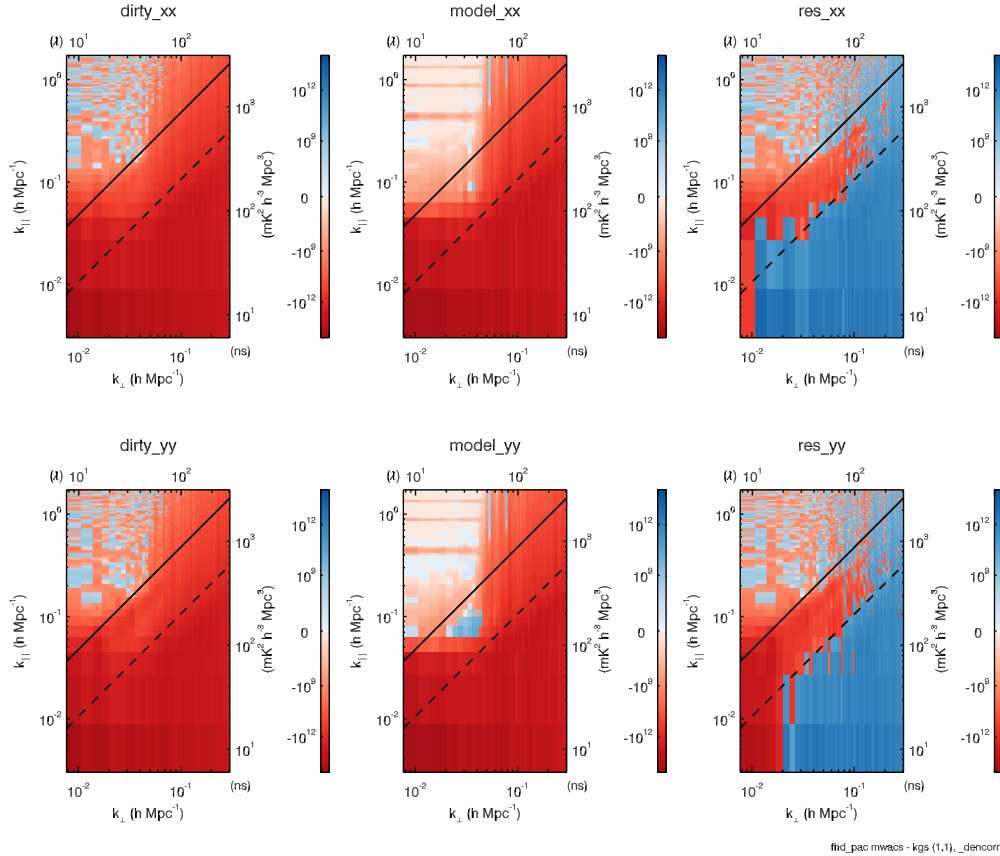


Figure 8.1: The MWACS–KGS power spectrum difference. The top and bottom rows are the XX and YY polarizations respectively. Left: The power spectrum of the difference of the dirty visibility cubes. Middle: The power spectrum of the difference of model visibility cubes (middle). Right: The power spectrum of the difference of the residual (dirty - model) visibility cubes. In each difference the test KGS catalog run is subtracted from the standard MWACS catalog run.

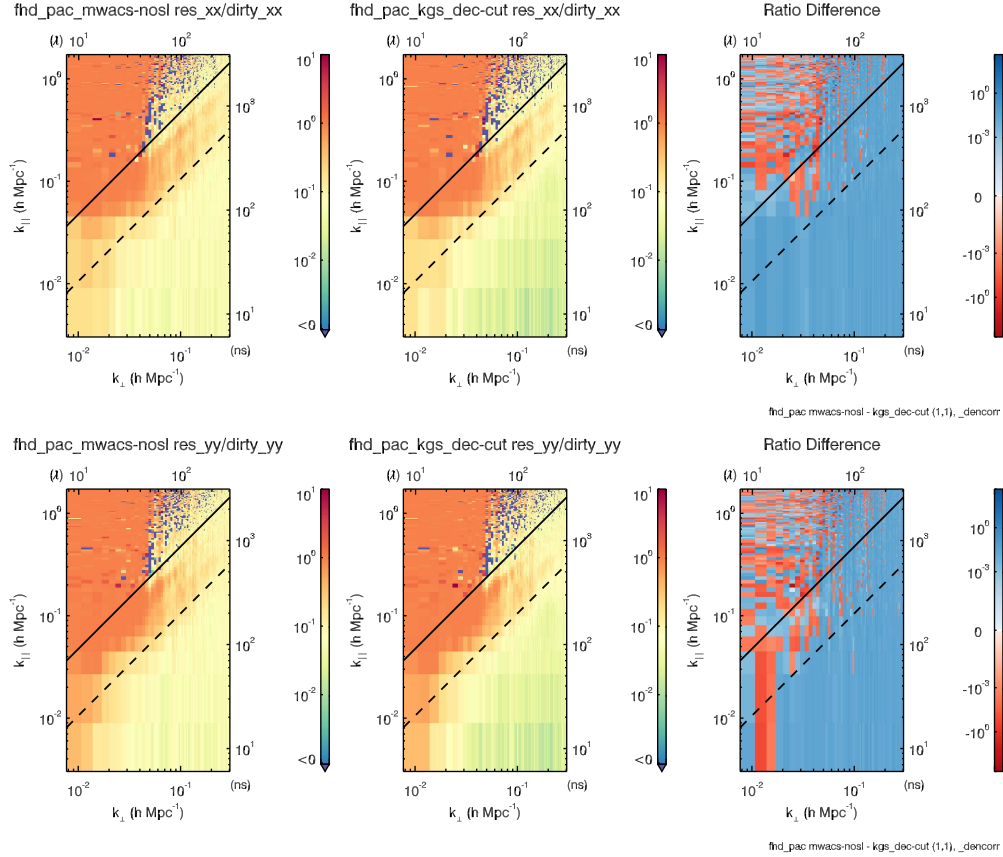


Figure 8.2: The MWACS–KGS normalized residual power spectrum difference. The top and bottom rows are the XX and YY polarizations respectively. Left: The power spectrum of the ratio of the residual to dirty visibility cubes for the standard MWACS catalog run. Middle: The power spectrum of the ratio of the residual to dirty visibility cubes for the test KGS catalog run. Right: The power spectrum of the difference between the residual-to-dirty ratio cubes (KGS run subtracted from MWACS run). Taking the ratio effectively normalizes for flux scale calibration differences. The ratio difference is positive in the wedge and foreground dominated modes, indicating better foreground power removal using the KGS catalog.

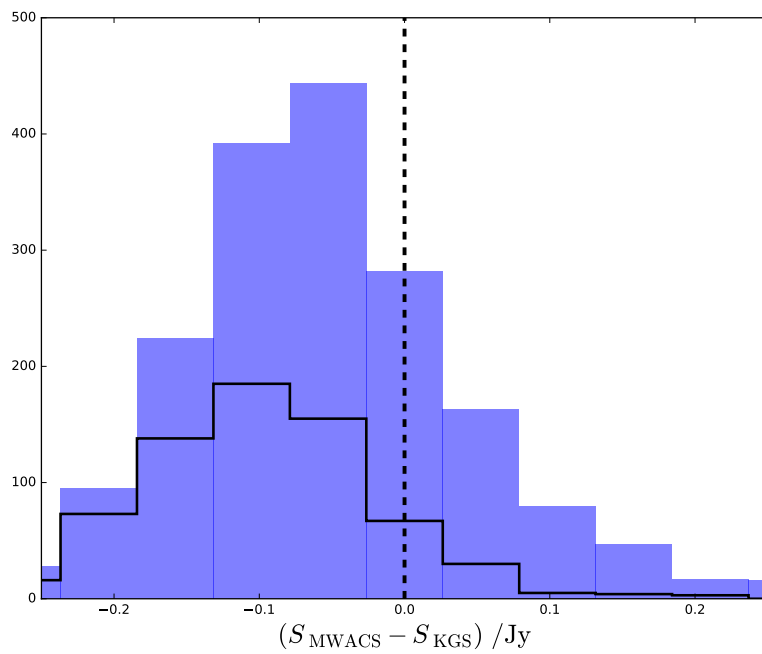


Figure 8.3: The difference between MWACS and KGS flux density for 1881 isolated matched sources. The MWACS 180 MHz flux density was extrapolated to 182 MHz using the reported catalog spectral index. The MWACS source flux densities are systematically low compared to KGS. The effect is more pronounced for sources within half-beam (black line).

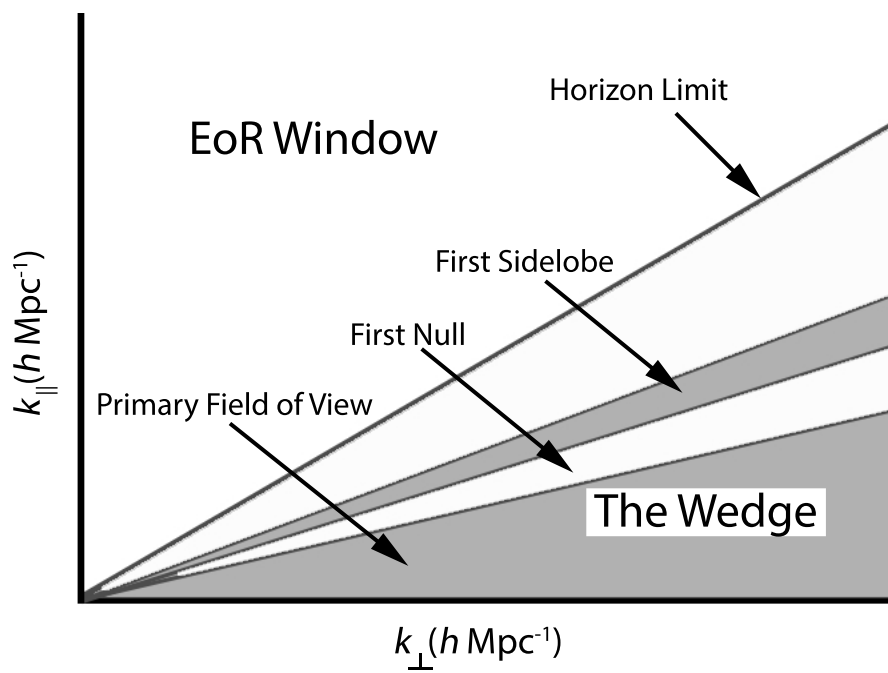


Figure 8.4: A schematic of the 2D  $k$ -power spectrum showing the modes occupied by foreground sources depending on where they are located in the beam. Widefield sources contaminate higher  $k_{\parallel}$  than sources closer to beam center. Figure credit: J. C.Pober.

at 160 MHz and 80 MHz. Because of the better frequency agreement with the MWA and two-point spectral index measurements, these were included in place of their MRC counterparts. MWACS and CCA cataloged spectral index measurements were used to extrapolate the flux density to 182 MHz. MRC sources were cross-matched to SUMSS and the estimated two-point spectral index was used or, if no clear match was found, a spectral index of -0.8 was assumed.

The master catalog was used to demonstrate the impact of including widefield foreground sources in the sky model by Pober et al. (2016 [30]). A sky model consisting of the 4600 sources in the primary beam was compared to a sky model including a total of  $\sim 8500$  sources through the first side lobe. The power spectrum difference is shown in Figure 8.5. A positive residual difference is seen in  $k_{\parallel}$  modes above the primary beam line, corresponding to decreased residual foreground power in first side lobes.

## 8.2 A Multi-Survey Master Catalog

With the development of the PUMA cross matching software and the finalization of the KGS catalog, an updated version of the master catalog was built. By cross-matching with PUMA and taking advantage of the higher positional accuracy of NVSS and SUMSS, we gain more robust spectral index measurements and consistency in the overall flux scale and mapping.

The KGS foreground catalog serves as the base, covering  $\sim 1400 \text{ deg}^2$  centered on the EoR0 field, and flux densities are corrected for Eddington bias. We then cross-match the MWACS to the VLSSr, MRC, SUMSS, and NVSS catalogs using PUMA as described in §5 and a  $3'$  initial search radius. Lacking enough information to assess their reliability, MWACS sources unmatched to any of the comparison surveys are discarded. Sources within the primary beam are included only if there is no KGS match and  $S \cdot \text{beam} < 100 \text{ mJy}$ . MWACS positions are known to be biased and are therefore corrected to the NVSS, SUMSS, VLSSr, or MRC match position in that order of preference. The updated position error is added in quadrature with the offset distance. Sources that are matched by STILTS but that are flagged by PUMA are assumed to be real and complex. These are included using

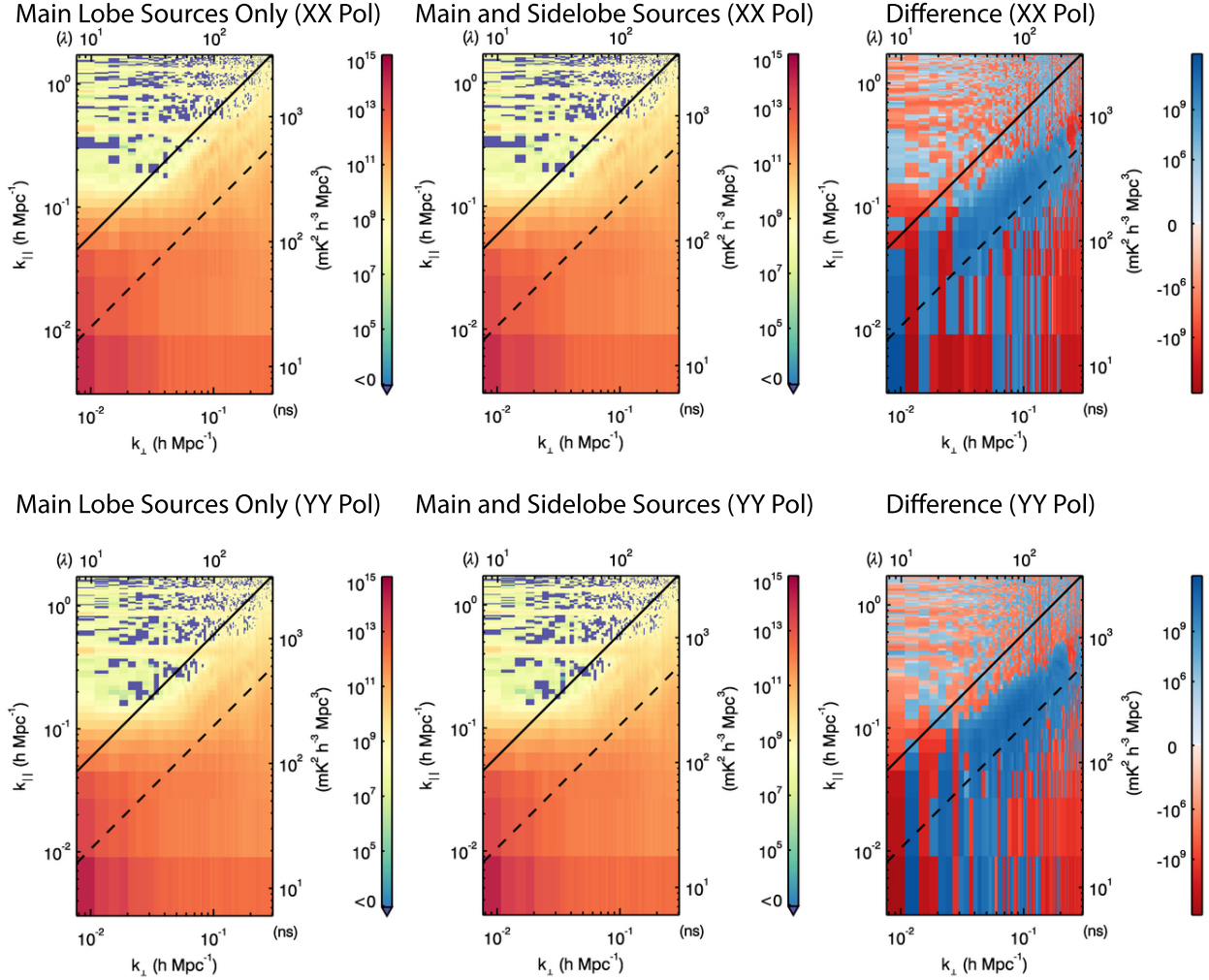


Figure 8.5: The difference in the 2D  $k$ -power spectrum when including sources in the side lobes of the primary beam in the sky model. The dashed black line marks the extent of the primary beam. The solid line marks the approximate horizon limit. A positive (blue) residual difference is seen in  $k_{\parallel}$  modes above the primary beam line, corresponding to the first side lobes.

their original catalog position. To adjust the flux density to 182 MHz, we use the reported MWACS spectral index with the exception of outliers. The 1% tails of the spectral index distribution are considered unphysical. For these we use the PUMA measured spectral index where available or -0.71 is assumed.

Next, we turn to the MRC. At 408 MHz, the MRC was initially problematic as a foreground model due to flux density extrapolation errors. By cross-matching to multiple surveys, we gain more robust spectral index measurements, substantially improving the 182 MHz flux density estimate. The MRC is cross-matched to the combined KGS+MWACS catalog, NVSS, SUMSS, and VLSSr. Due to poor match results of a few bright sources, any source matched to the KGS/MWACS within a radius of 5' was discarded. The PUMA estimated broad-band spectral index was used to extrapolate the flux density to 182 MHz where available or -0.8 was assumed. Although the MRC positions show no bias relative to NVSS or SUMSS, isolated match positions are corrected for consistency.

Lastly, any MRC source in the combined catalog was replaced by its CCA counterpart. If a CCA spectral index was not reported for a source, -0.71 was assumed to extrapolate the flux density to 182 MHz. There are 28962 sources in the updated master catalog. The contribution and coverage of each survey are illustrated in Figure 8.6 relative to the EoR0 field.

As a final step, the master catalog was updated with an extended source model for NGC 253 described in §8.3. We tested the master catalog against MWACS by running first-pass on the 15 snapshots from the zenith pointing while allowing calibration and subtraction in the side lobes of the primary beam. The resulting power spectrum difference and ratio plots are shown in Figures 8.7–8.8.

### **8.3 Extended Source Models**

We have so far modeled all sources as point sources. A minor but significant 13% of KGS sources are `multiple` matches, with complex and extended structure that is not adequately represented as a point source. There are few bright extended sources in the MWA EoR fields

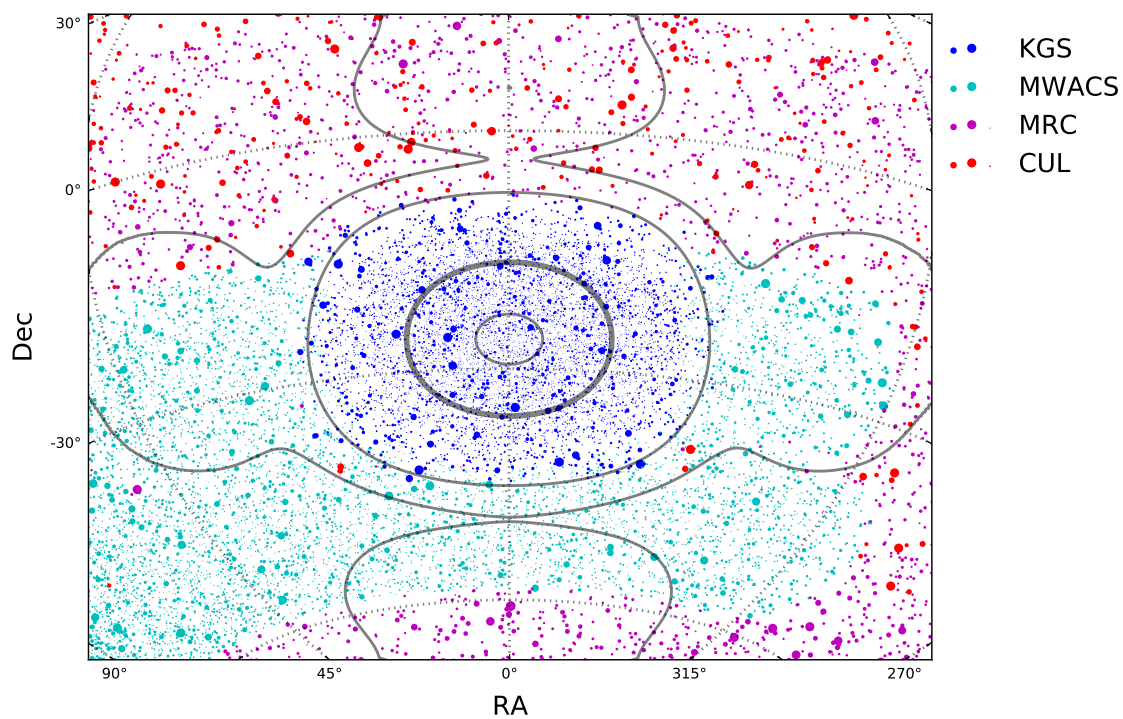


Figure 8.6: A multi-survey master catalog of foreground sources. The coverage of the contributing catalogs are shown by color. Point size is proportional to source flux density clipped at 20 Jy.

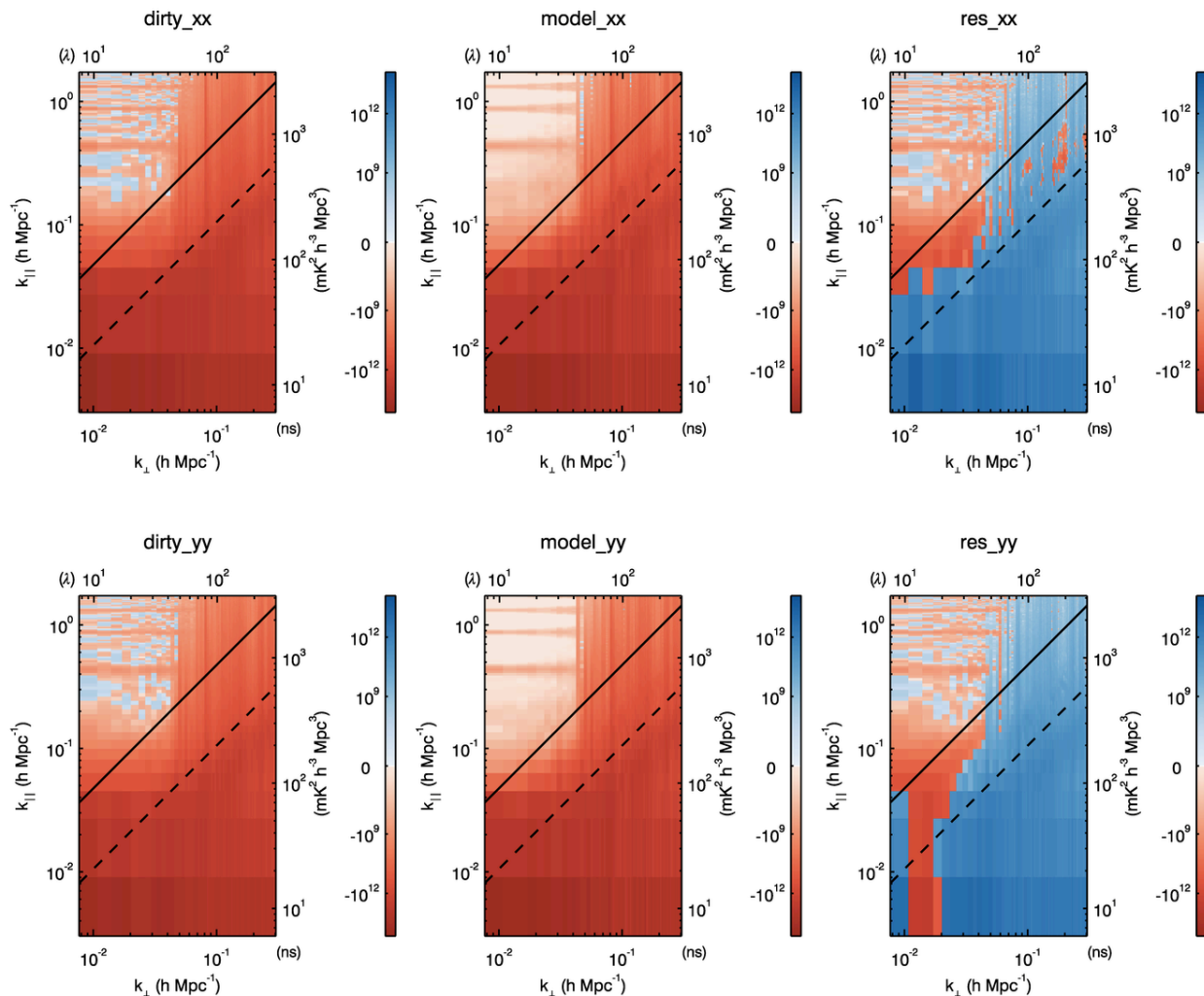


Figure 8.7: The MWACS–Master power spectrum difference. The top and bottom rows are the XX and YY polarizations respectively. Left: The power spectrum of the difference of the dirty visibility cubes. Middle: The power spectrum of the difference of model visibility cubes (middle). Right: The power spectrum of the difference of the residual (dirty - model) visibility cubes. In each difference the test Master catalog run is subtracted from the standard MWACS catalog run.

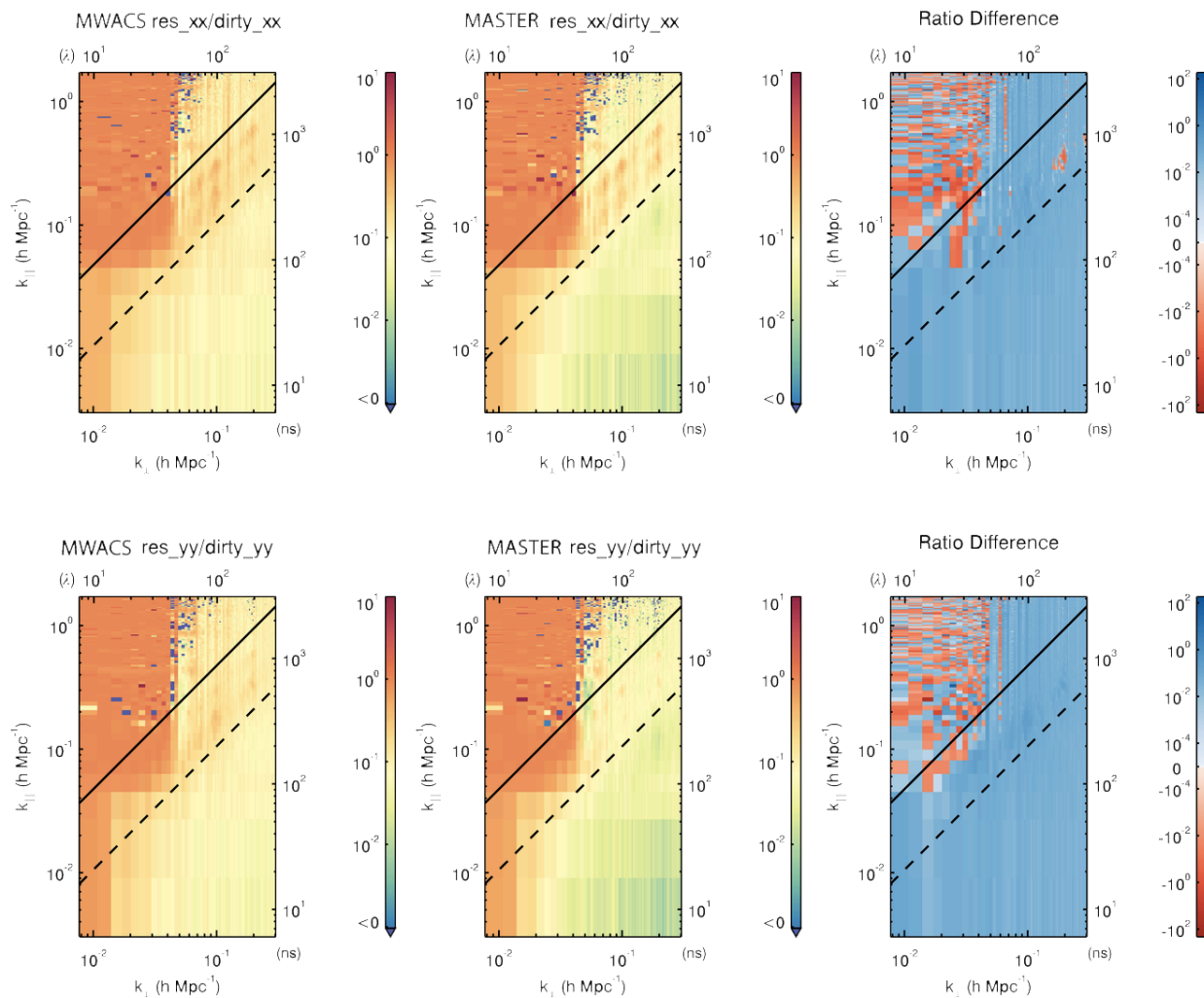


Figure 8.8: The MWACS–Master normalized residual power spectrum difference. The top and bottom rows are the XX and YY polarizations respectively. Left: The power spectrum of the ratio of the residual to dirty visibility cubes for the standard MWACS catalog run. Middle: The power spectrum of the ratio of the residual to dirty visibility cubes for the test Master catalog run. Right: The power spectrum of the difference between the residual-to-dirty ratio cubes (Master run subtracted from MWACS run). Taking the ratio effectively normalizes for flux scale calibration differences. The ratio difference is positive in the wedge and side lobes, indicating better foreground power removal using the KGS catalog.

by design, but where they do occur, a more complex model is required.

In the EoR0 field, there is one problematic source that stands out. The nearby starburst galaxy NGC 253 was introduced in Figure 1.7 and discussed in §5.2.1. Its bright and complex morphology resulted in failed deconvolution for 4 snapshot observations. Several counterparts are identified in the NVSS at higher resolution.

We attempt to model NGC 253 beginning with the original set of deconvolved components from §2.2. Recall that components are centroided at floating point pixel locations and are strictly positive valued. There are 9416 individual components from the 71 snapshots within 15' of the source position (RA = 11.8969°, Dec = -25.2823°). These are mapped to a grid at 10× the original pixel resolution. Figure 8.9 shows the number density and cumulative flux density of components on this map. A remarkable amount of extended structure is captured in the component distribution. A bright core, thin disk, and diffuse halo are apparent. This structure is all but lost after convolving with the instrument PSF to produce a restored image as shown in Figure 5.3(a).

For comparison to the NVSS, the components were separately gridded to a map matching a high resolution NVSS postage stamp image, and the map was smoothed with a Gaussian kernel approximating the NVSS beam (PSF FWHM = 45"). Figure 8.10 shows the resulting rendered image alongside the NVSS postage stamp image. The structure is clearly correlated and deserving of a detailed analysis, but that is beyond the scope of this thesis.

The component flux density map is convolved with a small Gaussian kernel ( $\sigma = 2$  pix) to smooth out Poisson noise and then down sampled by a factor of 4 (Figure 8.11). The map is normalized to a flux density of 18.8 Jy, the beam weighted average of the summed component flux densities from each snapshot. We note that the reported KGS catalog flux density ( $20.8 \pm 1.3$  Jy) for this source was weighted by the estimated background rms from each snapshot, and is higher by  $\sim 2$  Jy.

There are 263 resulting pixels  $\geq 4$  mJy, accounting for 98.95% of the total source flux density. Each pixel becomes a model component and the set of 263 model components forms the extended source model. The model components are included in the catalog as

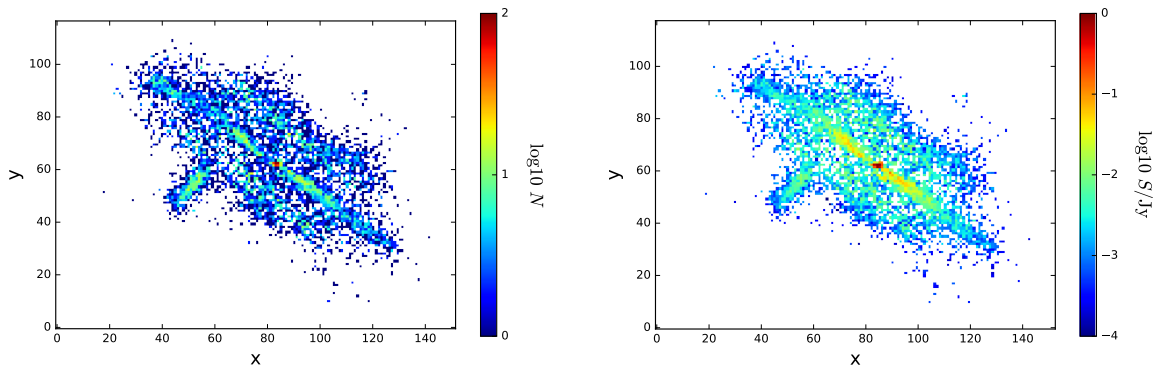


Figure 8.9: The 9416 components deconvolved from all 71 snapshots are gridded at  $10\times$  the original pixel resolution. Number density (left) and flux density (right) maps of NGC 253 components. A remarkable amount of extended structure is retained in the component distribution. This is all but lost after convolving with the instrument PSF to produce a restored image as shown in Figure 5.3(a)

point sources and the original source is removed. The result is a significant improvement in source subtraction. Residual images after first pass subtraction from the zenith snapshot are shown in Figure 8.12 and the power spectrum differences in Figures 8.13–8.14. A decrease in residual power is seen in high  $k_{\perp}$  modes and a change in the flux scale of the dirty difference suggests that calibration is also improved.

#### 8.4 Summary

The KGS EoR0 foreground catalog greatly improves the removal of foreground power from within the primary beam. This was demonstrated by differences in the 2D  $k$ -power spectra compared to the MWACS.

Sources in the side lobes, far from beam center, contaminate higher  $k_{\parallel}$  modes and encroach on the EoR window. For this reason, a combined Master catalog was created. PUMA cross-matching was used to merge the MWACS, MRC, and CCA catalogs with the KGS.

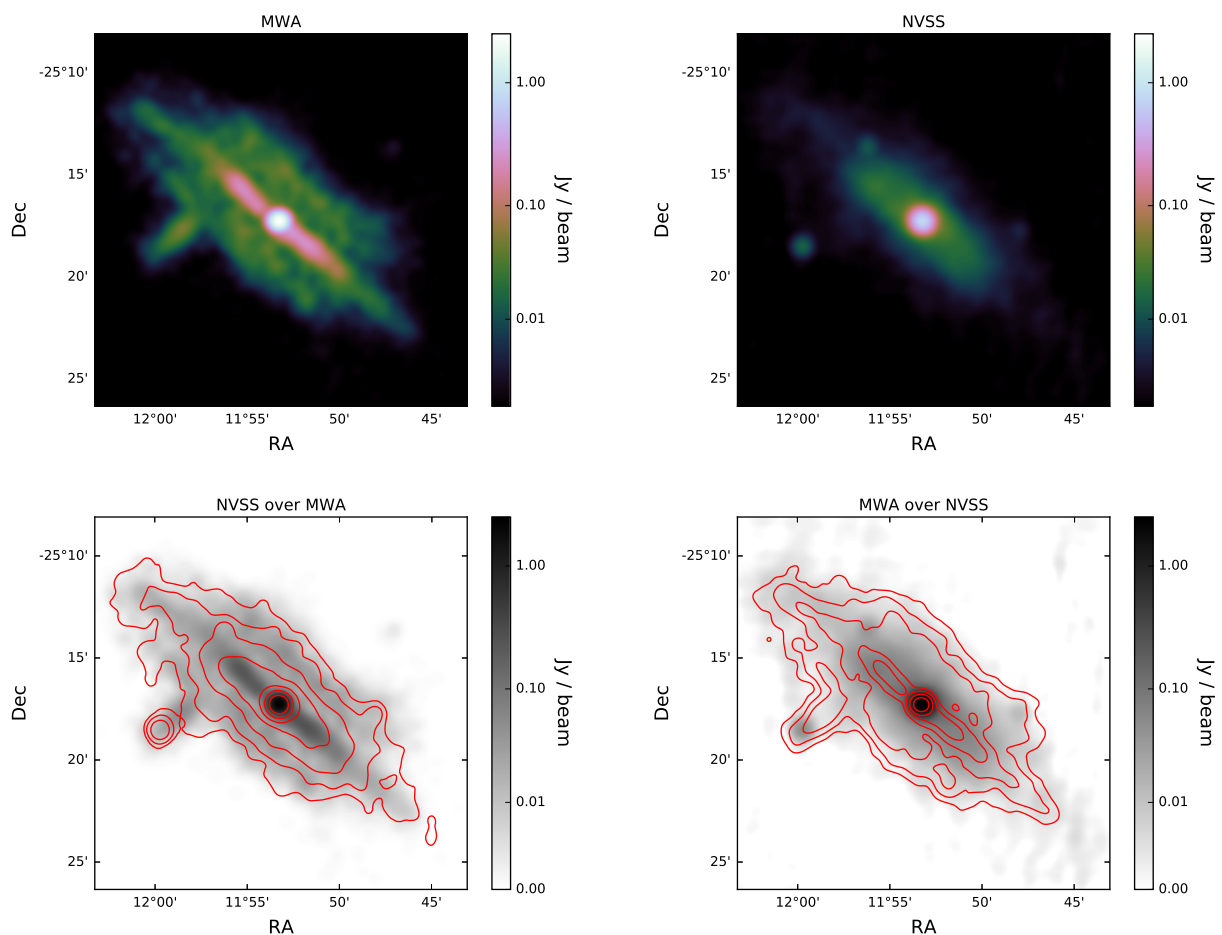


Figure 8.10: For comparison, the components were gridded to a map matching an NVSS postage stamp image and convolved with a Gaussian kernel approximating the NVSS beam (PSF FWHM =  $45''$ ). The MWA rendered image of KGS components are shown on the left with NVSS contours overlaid on the bottom panel. The NVSS image is shown on the right with MWA contours overlaid. Contours are at  $e^{n=1,2,\dots,9}$  mJy.

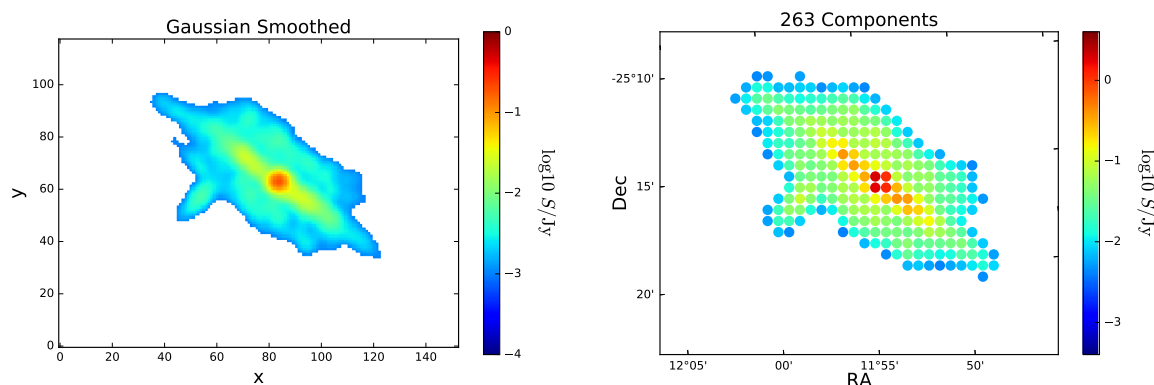


Figure 8.11: Left: The gridded components are smoothed with a small Gaussian kernel ( $\sigma = 2$  pix). Right: The grid is down sampled by a factor of 4 and normalized to conserve flux. The 263 pixels 1 mJy form the new extended source model.

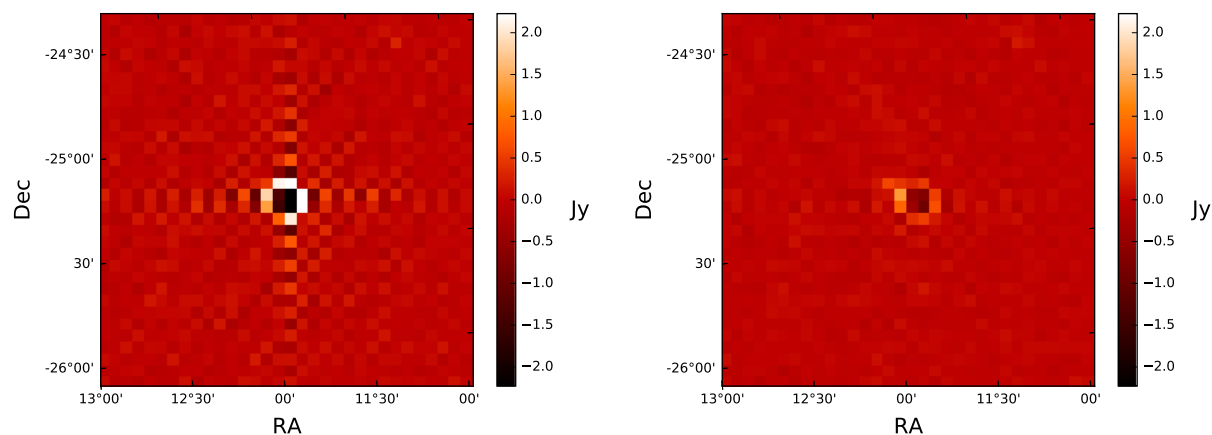


Figure 8.12: Firstpass residual images of the zenith snapshot centered on NGC 253 using a point source model (left) and extended source model (right). The color scale is clipped to show the image artifacts resulting from the point source subtraction. The point source model over-subtracts by 8 Jy/pixel at it's worst.

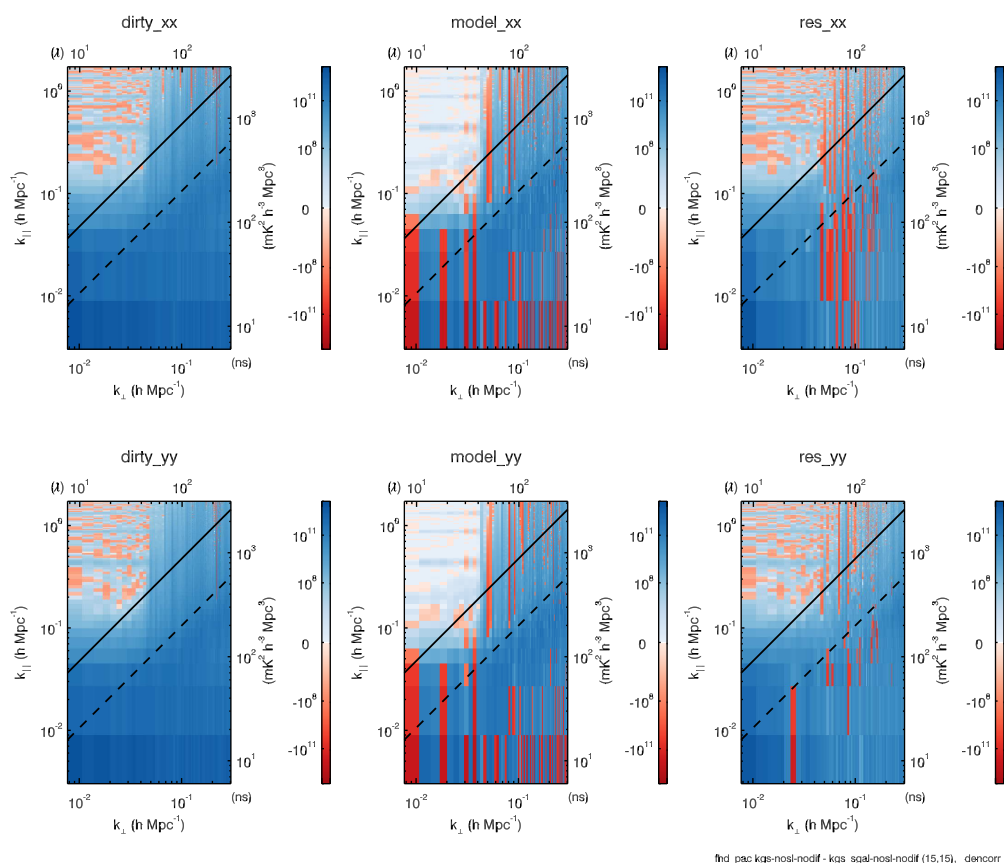


Figure 8.13: The power spectrum difference using an extended source model for NGS 253 compared to a point source model. The KGS catalog is otherwise identical between the two runs. The top and bottom rows are the XX and YY polarizations respectively. Left: The power spectrum of the difference of the dirty visibility cubes. Middle: The power spectrum of the difference of model visibility cubes. Right: The power spectrum of the difference of the residual (dirty - model) visibility cubes. In each difference the extended model run is subtracted from the point source model run. A change in the flux scale of the dirty visibilities suggests that calibration is improved but this makes it difficult to interpret the residual difference.

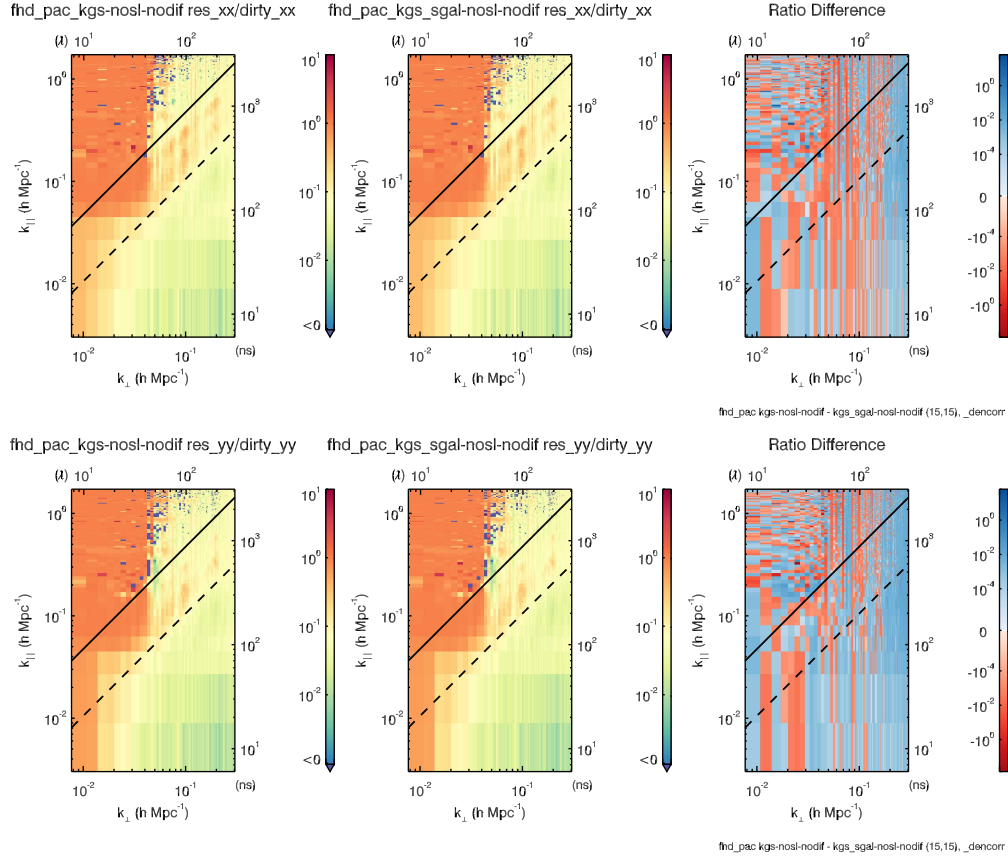


Figure 8.14: The normalized residual power spectrum and difference using an extended source model for NGS 253 compared to a point source model. The top and bottom rows are the XX and YY polarizations respectively. Left: The power spectrum of the ratio of the residual to dirty visibility cubes for the standard point source run. Middle: The power spectrum of the ratio of the residual to dirty visibility cubes for the test extended model run. Right: The power spectrum of the difference between the residual-to-dirty ratio cubes (point source run minus extended model run). Taking the ratio normalizes for the change in flux scale. Little difference is seen overall, but a decrease in residual power is evident at high  $k_{\perp}$  modes, indicating better foreground power removal using the extended model.

These were additionally cross-matched to the NVSS, SUMSS, and VLSSr to improve positional accuracy and spectral index estimates when cataloged values were unavailable or untrustworthy.

Bright and extended sources in the EoR fields are problematic, and cannot be modeled as point sources. An extended model was built for NGC 253. The point source was replaced by a set of 263 extended model components in the KGS and input to Firstpass. A major improvement is seen in the residual image. Residual power is decreased at high  $k_{\perp}$  modes and a change in the flux scale suggests improved calibration.

## Chapter 9

# CONCLUSIONS

### ***9.1 EoR Foreground Modeling***

The MWA EoR effort hinges on the ability to uncover the faint cosmological HI signal in spite of bright complicating foreground sources. It requires a precise and reliable sky model to nail down the calibration and remove foreground power contaminating the Fourier modes most sensitive to the EoR. Existing radio surveys and survey methodology are not tuned to the needs of such an analysis. They do not offer complete coverage or sufficient depth and are subject to false source contamination. This thesis has presented novel methods and a rigorous selection process to survey of the MWA EoR0 field and maximize reliability, completeness, and precision of source measurements.

We processed 2.5 hours of data consisting of 75 2 min snapshots. Each snapshot was individually deconvolved using FHD. Source finding was performed by spatially clustering deconvolved source components into isolated sources, and then similarly clustering to associate detections across all snapshots. This process avoided uncertainties inherent to fitting a simplified source model in a stacked restored radio image, and allowed for a measure of source detection rate valuable to reliability determination. 9490 source candidates were identified in two or more snapshots.

We used machine learning methods to self-consistently assign all source candidates to 10 possible classes that are interpreted in terms of reliability. The reliability classification was used in combination with traditional signal-to-noise measurements to select 7466 source candidates. All candidates were cross-matched to the VLSSr, MRC, SUMSS, and NVSS surveys using Bayesian statistics to assess the match probability. If a positional match was questionable, it was accepted if the the broad-band SED was well fit by a power law. All

outliers were flagged and manually confirmed or modified. The cross-match results supported the reliability classifications and allowed for the inclusion of faint low-reliability sources near the detection threshold.

The final catalog is composed of 7394 sources detected at 182 MHz. A positional bias was uncovered and traced to a bias in the catalog used for calibration. A bias in flux density near the detection threshold, inherent to any survey, was found to significantly affect a small number of sources. Both the positional bias and flux density bias were modeled and corrected. The median broad-band spectral index is -0.85 but this is seen to vary considerably between 74 MHz and 1400 MHz. The median spectral index at 182 MHz is predicted to be -0.71. The KGS catalog is complete to 80 mJy within half-beam power, significantly deeper than the MRC or MWACS previously used for the EoR analysis.

The measurement and modeling of foreground sources have little dependence on their astrophysical classifications. However, in producing the KGS catalog, several new and rare type radio sources were identified. We explored the ultra-steep spectrum (USS) sub-population and identified potential associations with galaxy clusters and high redshift radio galaxies. Observed trends in spectral index, spectral curvature, and source type suggest strong predictive potential for USS source classification.

The final catalog was tested against the MWACS for calibration and foreground power removal in the primary beam. The power spectrum difference figures of merit reveal an overall improvement in the calibration as well as a significant decrease in the amount of residual foreground power. In order to address the need to remove sources in the side lobes of the primary beam an all-sky master catalog was built. The master catalog is based on the KGS in the EoR0 field and built up through cross-matching to the MWACS and MRC. The VLSSr, SUMSS, and NVSS cross-match information is also used to improve the consistency and accuracy of source positions and predicted 182 MHz flux density. By expanding our sky model with the master catalog, foreground power removal is markedly improved within modes nearest the EoR window.

As a final improvement to the foreground sky model, we addressed the fact that not

all sources can be sufficiently modeled as point sources. The nearby star forming galaxy NGC 253 is the brightest extended source in the field. A point source model grossly over-subtracts flux and leaves complex artifacts in the residual image. The FHD deconvolved source components are seen to trace sub-resolution structure remarkably well. We grid and average the source components over all snapshots to produce a multi-component extended source model for NGC 253. Using a more accurate model for one of the brightest sources in the field significantly improves calibration.

## **9.2 Future Directions**

This thesis produced a high fidelity, high reliability catalog of discrete sources for calibration and foreground power removal within the MWA EoR analysis pipeline. The combined master catalog is currently the best catalog available for EoR foreground modeling in the southern sky. It is therefore useful to other southern EoR projects and will be used during pipeline development and early analysis runs with the next generation Hydrogen Epoch of Reionization Array (HERA [31]).

In addition, a new set of tools was developed for source finding, selection, and modeling to address the unique needs of the EoR analysis. These tools will be further developed and applied to an all-sky MWA snapshot survey currently under-way to improve wide-field coverage. A follow-up KGS 2.0 analysis will commence shortly, taking advantage of advances in FHD and better calibration to go even deeper in the EoR0 field.

Despite our best efforts to avoid bright and extended sources, problematic sources like NGC 253 will crop up. For example, the source PKS 2356-61 is a very bright galaxy (131.4 Jy at 145 MHz [21]) occupying the southern side lobe. Subtracting this as a point source results in artifacts that dominate the residual image of the primary beam. Extended source models can now be easily created for this and other troublesome sources from a few snapshots processed through FHD.

Lastly, the fidelity of any foreground model can be expected to degrade over time, as AGN are known to exhibit 5-10% variability in flux density on the order of 1-2 years. Future

investigations will attempt to quantify this impact and suggest a timescale for periodic updates to the foreground catalog.

## BIBLIOGRAPHY

- [1] G. O. Abell, H. G. Corwin, Jr., and R. P. Olowin. A catalog of rich clusters of galaxies. , 70:1–138, May 1989.
- [2] J. W. M. Baars, R. Genzel, I. I. K. Pauliny-Toth, and A. Witzel. The absolute spectrum of CAS A - an accurate flux density scale and a set of secondary calibrators. , 61:99–106, October 1977.
- [3] A. P. Beardsley, B. J. Hazelton, M. F. Morales, W. Arcus, D. Barnes, G. Bernardi, J. D. Bowman, F. H. Briggs, J. D. Bunton, R. J. Cappallo, B. E. Corey, A. Deshpande, L. deSouza, D. Emrich, B. M. Gaensler, R. Goeke, L. J. Greenhill, D. Herne, J. N. Hewitt, M. Johnston-Hollitt, D. L. Kaplan, J. C. Kasper, B. B. Kincaid, R. Koenig, E. Kratzenberg, C. J. Lonsdale, M. J. Lynch, S. R. McWhirter, D. A. Mitchell, E. Morgan, D. Oberoi, S. M. Ord, J. Pathikulangara, T. Prabu, R. A. Remillard, A. E. E. Rogers, A. Roshni, J. E. Salah, R. J. Sault, S. N. Udaya, K. S. Srivani, J. Stevens, R. Subrahmanyam, S. J. Tingay, R. B. Wayth, M. Waterson, R. L. Webster, A. R. Whitney, A. Williams, C. L. Williams, and J. S. B. Wyithe. The EoR sensitivity of the Murchison Widefield Array. , 429:L5–L9, February 2013.
- [4] A. P. Beardsley, B. J. Hazelton, M. F. Morales, R. J. Capallo, R. Goeke, D. Emrich, C. J. Lonsdale, W. Arcus, D. Barnes, G. Bernardi, J. D. Bowman, J. D. Bunton, B. E. Corey, A. Deshpande, L. deSouza, B. M. Gaensler, L. J. Greenhill, D. Herne, J. N. Hewitt, D. L. Kaplan, J. C. Kasper, B. B. Kincaid, R. Koenig, E. Kratzenberg, M. J. Lynch, S. R. McWhirter, D. A. Mitchell, E. Morgan, D. Oberoi, S. M. Ord, J. Pathikulangara, T. Prabu, R. A. Remillard, A. E. E. Rogers, A. Roshni, J. E. Salah, R. J. Sault, N. U. Shankar, K. S. Srivani, J. Stevens, R. Subrahmanyam, S. J. Tingay, R. B. Wayth,

- M. Waterson, R. L. Webster, A. R. Whitney, A. Williams, C. L. Williams, and J. S. B. Wyithe. A new layout optimization technique for interferometric arrays, applied to the Murchison Widefield Array. , 425:1781–1788, September 2012.
- [5] S Bhatnagar, T J Cornwell, K Golap, and J M Uson. Correcting direction-dependent gains in the deconvolution of radio interferometric images. *Astronomy & Astrophysics*, 487:419, August 2008.
- [6] C. G. Bornancini, A. L. O’Mill, S. Gurovich, and D. G. Lambas. Radio galaxies in the Sloan Digital Sky Survey: spectral index-environment correlations. , 406:197–207, July 2010.
- [7] J. D. Bowman, I. Cairns, D. L. Kaplan, T. Murphy, D. Oberoi, L. Staveley-Smith, W. Arcus, D. G. Barnes, G. Bernardi, F. H. Briggs, S. Brown, J. D. Bunton, A. J. Burgasser, R. J. Cappallo, S. Chatterjee, B. E. Corey, A. Coster, A. Deshpande, L. deSouza, D. Emrich, P. Erickson, R. F. Goeke, B. M. Gaensler, L. J. Greenhill, L. Harvey-Smith, B. J. Hazelton, D. Herne, J. N. Hewitt, M. Johnston-Hollitt, J. C. Kasper, B. B. Kincaid, R. Koenig, E. Kratzenberg, C. J. Lonsdale, M. J. Lynch, L. D. Matthews, S. R. McWhirter, D. A. Mitchell, M. F. Morales, E. H. Morgan, S. M. Ord, J. Pathikulangara, T. Prabu, R. A. Remillard, T. Robishaw, A. E. E. Rogers, A. A. Roshi, J. E. Salah, R. J. Sault, N. U. Shankar, K. S. Srivani, J. B. Stevens, R. Subrahmanyam, S. J. Tingay, R. B. Wayth, M. Waterson, R. L. Webster, A. R. Whitney, A. J. Williams, C. L. Williams, and J. S. B. Wyithe. Science with the Murchison Widefield Array. , 30:e031, April 2013.
- [8] Leo Breiman. Random forests. *Machine Learning*, 45(1):5–32, 2001.
- [9] Tamás Budavári and Alexander S. Szalay. Probabilistic CrossIdentification of Astronomical Sources. *The Astrophysical Journal*, 679(1):301–309, May 2008.
- [10] a. S. Cohen, W. M. Lane, W. D. Cotton, N. E. Kassim, T. J. W. Lazio, R. a. Perley, J. J.

- Condon, and W. C. Erickson. The VLA Low-Frequency Sky Survey. *The Astronomical Journal*, 134(3):1245–1262, September 2007.
- [11] J. J. Condon, W. D. Cotton, E. W. Greisen, Q. F. Yin, R. A. Perley, G. B. Taylor, and J. J. Broderick. The NRAO VLA Sky Survey. , 115:1693–1716, May 1998.
- [12] A Datta, J Bowman, and C Carilli. Bright Source Subtraction Requirements for Redshifted 21 cm Measurements. *The Astrophysical Journal*, 724:526–538, November 2010.
- [13] Joshua S Dillon, Adrian Liu, and Max Tegmark. A fast method for power spectrum and foreground analysis for 21 cm cosmology. *Physical Review D*, 87(4):043005, February 2013.
- [14] A. S. Eddington. On a formula for correcting statistics for the effects of a known error of observation. , 73:359–360, March 1913.
- [15] Martin Ester, Hans-Peter Kriegel, Jrg Sander, and Xiaowei Xu. A density-based algorithm for discovering clusters in large spatial databases with noise. pages 226–231. AAAI Press, 1996.
- [16] Yoav Freund and Robert E Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. In *Computational learning theory*, pages 23–37. Springer, 1995.
- [17] Bryna J Hazelton, Miguel F Morales, and Ian S Sullivan. The fundamental multi-baseline mode-mixing foreground in 21 cm Epoch of Reionization observations. *The Astrophysical Journal*, 770(2):156, June 2013.
- [18] G. H. Heald, R. F. Pizzo, E. Orrú, R. P. Breton, D. Carbone, C. Ferrari, M. J. Hardcastle, W. Jurusik, G. Macario, D. Mulcahy, D. Rafferty, A. Asgekar, M. Brentjens, R. A. Fallows, W. Frieswijk, M. C. Toribio, B. Adebahr, M. Arts, M. R. Bell, A. Bonafede, J. Bray, J. Broderick, T. Cantwell, P. Carroll, Y. Cendes, A. O. Clarke, J. Croston,

S. Daiboo, F. de Gasperin, J. Gregson, J. Harwood, T. Hassall, V. Heesen, A. Horneffer, A. J. van der Horst, M. Iacobelli, V. Jelić, D. Jones, D. Kant, G. Kokotanekov, P. Martin, J. P. McKean, L. K. Morabito, B. Nikiel-Wroczyński, A. Offringa, V. N. Pandey, M. Pandey-Pommier, M. Pietka, L. Pratley, C. Riseley, A. Rowlinson, J. Sabater, A. M. M. Scaife, L. H. A. Scheers, K. Sendlinger, A. Shulevski, M. Sipior, C. Sobey, A. J. Stewart, A. Stroe, J. Swinbank, C. Tasse, J. Trüstedt, E. Varenius, S. van Velzen, N. Vilchez, R. J. van Weeren, S. Wijnholds, W. L. Williams, A. G. de Bruyn, R. Nijboer, M. Wise, A. Alexov, J. Anderson, I. M. Avruch, R. Beck, M. E. Bell, I. van Bemmell, M. J. Bentum, G. Bernardi, P. Best, F. Breitling, W. N. Brouw, M. Brüggem, H. R. Butcher, B. Ciardi, J. E. Conway, E. de Geus, A. de Jong, M. de Vos, A. Deller, R.-J. Dettmar, S. Duscha, J. Eislöffel, D. Engels, H. Falcke, R. Fender, M. A. Garrett, J. Grießmeier, A. W. Gunst, J. P. Hamaker, J. W. T. Hessels, M. Hoeft, J. Hörandel, H. A. Holties, H. Intema, N. J. Jackson, E. Jütte, A. Karastergiou, W. F. A. Klijn, V. I. Kondratiev, L. V. E. Koopmans, M. Kuniyoshi, G. Kuper, C. Law, J. van Leeuwen, M. Loose, P. Maat, S. Markoff, R. McFadden, D. McKay-Bukowski, M. Mevius, J. C. A. Miller-Jones, R. Morganti, H. Munk, A. Nelles, J. E. Noordam, M. J. Norden, H. Paas, A. G. Polatidis, W. Reich, A. Renting, H. Röttgering, A. Schoenmakers, D. Schwarz, J. Sluman, O. Smirnov, B. W. Stappers, M. Steinmetz, M. Tagger, Y. Tang, S. ter Veen, S. Thoudam, R. Vermeulen, C. Vocks, C. Vogt, R. A. M. J. Wijers, O. Wucknitz, S. Yatawatta, and P. Zarka. The LOFAR Multifrequency Snapshot Sky Survey (MSSS). I. Survey description and first results. , 582:A123, October 2015.

- [19] N. Hurley-Walker, M. Johnston-Hollitt, R. Ekers, R. Hunstead, E. M. Sadler, L. Hindson, P. Hancock, G. Bernardi, J. D. Bowman, F. Briggs, R. Cappallo, B. Corey, A. A. Deshpande, D. Emrich, B. M. Gaensler, R. Goeke, L. Greenhill, B. J. Hazelton, J. Hewitt, D. L. Kaplan, J. Kasper, E. Kratzenberg, C. Lonsdale, M. Lynch, D. Mitchell, R. McWhirter, M. Morales, E. Morgan, D. Oberoi, A. Offringa, S. Ord, T. Prabu, A. Rogers, A. Roshi, U. Shankar, K. Srivani, R. Subrahmanyan, S. Tingay, M. Water-

- son, R. B. Wayth, R. Webster, A. Whitney, A. Williams, and C. Williams. Serendipitous discovery of a dying Giant Radio Galaxy associated with NGC 1534, using the Murchison Widefield Array. , 447:2468–2478, March 2015.
- [20] N. Hurley-Walker, J. Morgan, R. B. Wayth, P. J. Hancock, M. E. Bell, G. Bernardi, R. Bhat, F. Briggs, A. A. Deshpande, A. Ewall-Wice, L. Feng, B. J. Hazelton, L. Hindson, D. C. Jacobs, D. L. Kaplan, N. Kudryavtseva, E. Lenc, B. McKinley, D. Mitchell, B. Pindor, P. Procopio, D. Oberoi, A. Offringa, S. Ord, J. Riding, J. D. Bowman, R. Cappallo, B. Corey, D. Emrich, B. M. Gaensler, R. Goeke, L. Greenhill, J. Hewitt, M. Johnston-Hollitt, J. Kasper, E. Kratzenberg, C. Lonsdale, M. Lynch, R. McWhirter, M. F. Morales, E. Morgan, T. Prabu, A. Rogers, A. Roshi, U. Shankar, K. Srivani, R. Subrahmanyan, S. Tingay, M. Waterson, R. Webster, A. Whitney, A. Williams, and C. Williams. The Murchison Widefield Array Commissioning Survey: A Low-Frequency Catalogue of 14 110 Compact Radio Sources over 6 100 Square Degrees. , 31:e045, November 2014.
- [21] D. C. Jacobs, J. E. Aguirre, A. R. Parsons, J. C. Pober, R. F. Bradley, C. L. Carilli, N. E. Gugliucci, J. R. Manley, C. van der Merwe, D. F. Moore, and C. R. Parashare. New 145 MHz Source Measurements by PAPER in the Southern Sky. , 734:L34, June 2011.
- [22] W. M. Lane, W. D. Cotton, S. van Velzen, T. E. Clarke, N. E. Kassim, J. F. Helmboldt, T. J. W. Lazio, and A. S. Cohen. The Very Large Array Low-frequency Sky Survey Redux (VLSSr). , 440:327–338, May 2014.
- [23] M. I. Large, L. E. Cram, and A. M. Burgess. A machine-readable release of the Molonglo Reference Catalogue of Radio Sources. *The Observatory*, 111:72–75, April 1991.
- [24] T. Mauch, T. Murphy, H. J. Buttery, J. Curran, R. W. Hunstead, B. Piestrzynski, J. G. Robertson, and E. M. Sadler. SUMSS: a wide-field radio imaging survey of the southern sky - II. The source catalogue. , 342:1117–1130, July 2003.

- [25] Miguel F Morales, Bryna Hazelton, Ian Sullivan, and Adam Beardsley. Four Fundamental Foreground Power Spectrum Shapes for 21 cm Cosmology Observations. *The Astrophysical Journal*, 752(2):137, 2012.
- [26] Miguel F Morales and Michael Matejek. Software holography: interferometric data analysis for the challenges of next generation observatories. *Monthly Notices of the Royal Astronomical Society*, 400(4):1814–1820, December 2009.
- [27] S Myers, C Contaldi, J Bond, U Pen, D Pogosyan, S Prunet, J Sievers, B Mason, T Pearson, A Readhead, and M Shepherd. A Fast Gridded Method for the Estimation of the Power Spectrum of the Cosmic Microwave Background from Interferometer Data with Application to the Cosmic Background Imager. *The Astrophysical Journal*, 591:575–598, July 2003.
- [28] A. R. Offringa, A. G. de Bruyn, M. Biehl, S. Zaroubi, G. Bernardi, and V. N. Pandey. Post-correlation radio frequency interference classification methods. , 405:155–167, June 2010.
- [29] Aaron R Parsons, Jonathan C Pober, James E Aguirre, Christopher L Carilli, Daniel C Jacobs, and David F Moore. A per-baseline, delay-spectrum technique for accessing the 21 cm cosmic reionization signature. *The Astrophysical Journal*, 756(2):165, August 2012.
- [30] J. Pober, B. Hazelton, A. Beardsley, N. Barry, I. Sullivan, and M. Morales. The Importance of Wide-field Foreground Removal for 21 cm Cosmology. accepted to , 2016.
- [31] J. C. Pober, A. Liu, J. S. Dillon, J. E. Aguirre, J. D. Bowman, R. F. Bradley, C. L. Carilli, D. R. DeBoer, J. N. Hewitt, D. C. Jacobs, M. McQuinn, M. F. Morales, A. R. Parsons, M. Tegmark, and D. J. Werthimer. What Next-generation 21 cm Power Spectrum Measurements can Teach us About the Epoch of Reionization. , 782:66, February 2014.

- [32] Jonathan C Pober, Aaron R Parsons, James E Aguirre, Zaki Ali, Richard F Bradley, Chris L Carilli, Dave DeBoer, Matthew Dexter, Nicole E Gugliucci, Daniel C Jacobs, Patricia J Klima, Dave MacMahon, Jason Manley, David F Moore, Irina I Stefan, and William P Walbrugh. Opening the 21 cm Epoch of Reionization window: measurement of foreground isolation with paper. *The Astrophysical Journal*, 768(2):L36, April 2013.
- [33] R. S. Roger, C. H. Costain, and A. H. Bridle. The low-frequency spectra of nonthermal radio sources. , 78:1030, December 1973.
- [34] A. M. M. Scaife and G. H. Heald. A broad-band flux scale for low-frequency radio telescopes. , 423:L30–L34, June 2012.
- [35] O. B. Slee. Radio sources observed with the Culgoora circular array. *Australian Journal of Physics*, 48:143–186, 1995.
- [36] I. S. Sullivan, M. F. Morales, B. J. Hazelton, W. Arcus, D. Barnes, G. Bernardi, F. H. Briggs, J. D. Bowman, J. D. Bunton, R. J. Cappallo, B. E. Corey, A. Deshpande, L. deSouza, D. Emrich, B. M. Gaensler, R. Goeke, L. J. Greenhill, D. Herne, J. N. Hewitt, M. Johnston-Hollitt, D. L. Kaplan, J. C. Kasper, B. B. Kincaid, R. Koenig, E. Kratzenberg, C. J. Lonsdale, M. J. Lynch, S. R. McWhirter, D. A. Mitchell, E. Morgan, D. Oberoi, S. M. Ord, J. Pathikulangara, T. Prabu, R. A. Remillard, A. E. E. Rogers, A. Roshi, J. E. Salah, R. J. Sault, N. Udaya Shankar, K. S. Srivani, J. Stevens, R. Subrahmanyam, S. J. Tingay, R. B. Wayth, M. Waterson, R. L. Webster, A. R. Whitney, A. Williams, C. L. Williams, and J. S. B. Wyithe. Fast Holographic Deconvolution: A New Technique for Precision Radio Interferometry. , 759:17, November 2012.
- [37] M. B. Taylor. STILTS - A Package for Command-Line Processing of Tabular Data. In C. Gabriel, C. Arviset, D. Ponz, and S. Enrique, editors, *Astronomical Data Analysis Software and Systems XV*, volume 351 of *Astronomical Society of the Pacific Conference Series*, page 666, July 2006.

- [38] Nithyanandan Thyagarajan, Daniel C Jacobs, Judd D Bowman, N Barry, A P Beardsley, G Bernardi, F Briggs, R J Cappallo, P Carroll, A A Deshpande, A de Oliveira-Costa, Joshua S Dillon, A Ewall-Wice, L Feng, L J Greenhill, B J Hazelton, L Hernquist, J N Hewitt, N Hurley-Walker, M Johnston-Hollitt, D L Kaplan, Han-Seek Kim, P Kittiwisit, E Lenc, J Line, A Loeb, C J Lonsdale, B McKinley, S R McWhirter, D A Mitchell, M F Morales, E Morgan, A R Neben, D Oberoi, A R Offringa, S M Ord, Sourabh Paul, B Pindor, J C Pober, T Prabu, P Procopio, J Riding, N Udaya Shankar, Shiv K Sethi, K S Srivani, R Subrahmanyam, I S Sullivan, M Tegmark, S J Tingay, C M Trott, R B Wayth, R L Webster, A Williams, C L Williams, and J S B Wyithe. Confirmation of wide-field signatures in redshifted 21 cm power spectra. *The Astrophysical Journal*, 807(2):L28, July 2015.
- [39] Nithyanandan Thyagarajan, N Udaya Shankar, Ravi Subrahmanyam, Wayne Arcus, Gianni Bernardi, Judd D Bowman, Frank Briggs, John D Bunton, Roger J Cappallo, Brian E Corey, Ludi deSouza, David Emrich, Bryan M Gaensler, Robert F Goeke, Lincoln J Greenhill, Bryna J Hazelton, David Herne, Jacqueline N Hewitt, Melanie Johnston-Hollitt, David L Kaplan, Justin C Kasper, Barton B Kincaid, Ronald Koenig, Eric Kratzenberg, Colin J Lonsdale, Mervyn J Lynch, S Russell McWhirter, Daniel A Mitchell, Miguel F Morales, Edward H Morgan, Divya Oberoi, Stephen M Ord, Joseph Pathikulangara, Ronald A Remillard, Alan E E Rogers, D Anish Roshni, Joseph E Salah, Robert J Sault, K S Srivani, Jamie B Stevens, Prabu Thiagaraj, Steven J Tingay, Randall B Wayth, Mark Waterson, Rachel L Webster, Alan R Whitney, Andrew J Williams, Christopher L Williams, and J Stuart B Wyithe. A study of fundamental limits to statistical detection of redshifted H I from the Epoch of Reionization. *The Astrophysical Journal*, 776(1):6, September 2013.
- [40] S. J. Tingay, R. Goeke, J. D. Bowman, D. Emrich, S. M. Ord, D. A. Mitchell, M. F. Morales, T. Boller, B. Crosse, R. B. Wayth, C. J. Lonsdale, S. Tremblay, D. Palot, 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. Rosh, J. E. Salah, R. J. Sault, N. Udaya-Shankar, F. Schlagenhauer, 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. , 30:e007, January 2013.
- [41] Cathryn M Trott, Randall B Wayth, and Steven J Tingay. The impact of point-source subtraction residuals on 21 cm Epoch of Reionization estimation. *The Astrophysical Journal*, 757(1):101, September 2012.
- [42] E. Valtaoja and M. Valtonen. *Variability of Blazars*. Cambridge University Press, 1992.
- [43] Harish Vedantham, N Udaya Shankar, and Ravi Subrahmanyan. Imaging the Epoch of Reionization: limitations from foreground confusion and imaging algorithms. *The Astrophysical Journal*, 745(2):176, 2012.
- [44] R. B. Wayth, E. Lenc, M. E. Bell, J. R. Callingham, K. S. Dwarakanath, T. M. O. Franzen, B.-Q. For, B. Gaensler, P. Hancock, L. Hindson, N. Hurley-Walker, C. A. Jackson, M. Johnston-Hollitt, A. D. Kapińska, B. McKinley, J. Morgan, A. R. Offringa, P. Procopio, L. Staveley-Smith, C. Wu, Q. Zheng, C. M. Trott, G. Bernardi, J. D. Bowman, F. Briggs, R. J. Cappallo, B. E. Corey, A. A. Deshpande, D. Emrich, R. Goeke, L. J. Greenhill, B. J. Hazelton, D. L. Kaplan, J. C. Kasper, E. Kratzenberg, C. J. Lonsdale, M. J. Lynch, S. R. McWhirter, D. A. Mitchell, M. F. Morales, E. Morgan, D. Oberoi, S. M. Ord, T. Prabu, A. E. E. Rogers, A. Rosh, N. U. Shankar, K. S. Srivani, R. Subrahmanyan, S. J. Tingay, M. Waterson, R. L. Webster, A. R. Whitney,

- A. Williams, and C. L. Williams. GLEAM: The GaLactic and Extragalactic All-Sky MWA Survey. , 32:e025, June 2015.
- [45] Ji Zhu, Hui Zou, Saharon Rosset, and Trevor Hastie. Multi-class adaboost. *Statistics and its Interface*, 2(3):349–360, 2009.

## VITA

Patricia Ann Carroll was born to Marjorie and James Carroll in 1987. She was raised with her two brothers Timothy and Robert in a small town in the foothills of the Adirondack mountains of upstate New York. She attended Siena College where she was a standout student athlete, four time recipient of the Presidential merit scholarship, and recipient of the Clare Booth Luce Scholarship for women in STEM. Her love of astronomy began after her first semester of Physics, when she accompanied Professor Rose Finn to Kitt Peak Observatory for an observing run to study star formation in distant galaxies. She later participated in a radio astronomy workshop at Arecibo Radio observatory, and research internships at the American Museum of Natural History and NASA Jet Propulsion Laboratory. She graduated in 2009 with Summa Cum Laude honors, majoring in Physics and minoring in Mathematics and Chemistry. She is a first generation college graduate.

Beginning her graduate career at the University of Washington, Patti was awarded the Graduate Opportunities and Minority Achievement Program dissertation fellowship and joined the Radio Cosmology research group lead by Professor Miguel Morales. In 2012, she spent the summer at ASTRON, the Netherlands Institute for Radio Astronomy, working on the Low Frequency Array (LOFAR) Multi-frequency Snapshot Sky Survey (MSSS) with Dr. George Heald. In 2015, she was awarded the American-Australian Association Sir Keith Murdock fellowship and spent three months working with collaborators in Melbourne and Perth, Australia. Later that year, she devoted another three months to software development for Astropy, an open source Python software package for astronomers, as part of the Google Summer of Code program.

Patti has been active in science outreach over the years; tutoring physics, giving UW Planetarium shows to groups of all ages, acting as staff and student mentor for the Pre-major

in Astronomy program, presenting at the Pacific Science Center planetarium, giving public talks on science policy at UW and her own research at Seattle Town Hall, and advocating for science on behalf of the American Astronomical Society in Washington DC.

After graduating from the University of Washington, Patti will enter the field of data science, but hopes to one day return to academia. She currently lives in Seattle, WA with her wonderful partner, Ben, and dog, Sheba. In her free time, she can be found wielding a camera and exploring the mountains and wilderness of the great Pacific Northwest.