

The New Age of Big Data In Astronomy: A Review of the SKA & Rubin

MATHEW ICHO

The University of Illinois at Urbana-Champaign

ABSTRACT

I'm making my abstract my to do list for now

1. Fix the SDSS part based on provided notes, also make real-time processing section
2. Find a paper or part of it that describes how SDSS stores the data. I don't think I've done that sufficiently
3. Do the LSST part
4. Do the Results section, look at data collected. I plan on coding A LOT for this section
5. Redo the Du Pont section, it's outdated

Contents

1. INTRODUCTION

The concept of data has long been central throughout the history of astronomy. Data allows scientists to discover natural laws in the universe, have control over events, and make reliable predictions. It has played a critical role in other time-sensitive fields such as medicine and engineering, where accurate data is essential for decision-making and design. Although the nature of data varies fundamentally across different fields, one trend has remained consistent: the continual evolution of data science. As explained in The Fourth Paradigm (?), this evolution can be characterized through four successive paradigms. In the following sections, I describe the progression of data acquisition across these paradigms and illustrate each using examples from astronomy. I will then explain how SKA and LSST fit into this trajectory and exemplify the emerging era of data-intensive discovery.

1.1. *The Paradigms of Data Science*

The first and most primitive paradigm, as described by (?), is empirical evidence. Empirical evidence refers to data collected through traditional means, such as direct observation or experimentation. The primary purpose of empirical evidence is to identify patterns that allow scientists to develop a fundamental understanding of the natural world. Throughout much of human history, empirical evidence has dominated knowledge generation. An example of the first paradigm in astronomy is the career of Tycho Brahe, a Danish astronomer. Throughout his career in the 16th century, Brahe collected and cataloged data on the position of astronomical bodies using naked-eye observations. Tycho Brahe's catalogue was accurate to only around 1' precision and took decades to acquire (?). However, empirical evidence can be compromised by human error, the precision of

35 the instruments, and, most importantly, the relatively slow pace of data acquisition compared to
 36 subsequent paradigms.

37 The second paradigm is analytical evidence. Analytical evidence is obtained by constructing mathematical
 38 formulas and theoretical frameworks based on empirical data (?). Unlike the empirical
 39 evidence, which merely demonstrates that phenomena occur, the second paradigm seeks to explain
 40 why they occur. An example of the second paradigm in astronomy is the work of Johannes Kepler,
 41 who used Brahe's empirical observations to derive the laws of planetary motion (?). By transforming
 42 raw observational data into mathematical laws, Kepler exemplified how analytical evidence advances
 43 scientific understanding beyond description to explanation.

44 The third paradigm is simulation evidence (?), a relatively recent development. Simulation models
 45 natural phenomena that are too complex to model analytically or compute by hand. It allows
 46 interpolation and extrapolation of data using computational techniques grounded in known physical
 47 laws. For example, in astronomy, N-body simulations are used to study the complex dynamical
 48 evolution of planetary systems and galaxies.

49 The fourth and most recent paradigm is data-intensive science (?). This paradigm is characterized
 50 by the unprecedented scale, velocity, and complexity of data acquisition, driven in part by exponential
 51 advances in computational power and detector technologies, often associated with Moore's law (?).
 52 Unlike earlier paradigms, which focused on observation, theory, or simulation, data-intensive science
 53 emphasizes the ability to manage, analyze, and interpret vast datasets that exceed the capacity of
 54 traditional methods. While this exponential growth in data has enabled transformative discoveries,
 55 it also introduces significant challenges related to storage, processing, and accessibility.

56 1.2. *The Rise of Big Data in Astronomy*

57 Astronomy has become data intensive. Modern observatories may now generate petabyte-scale
 58 data that need new strategies for data management and analysis (?). The fourth paradigm enables
 59 discoveries from interpreting massive data sets. However, these advances also expose alarming issues,
 60 including bottlenecks in the data pipeline, storage challenges, increased skills needed to handle the
 61 data, and open access concerns. The field of astronomy is both a beneficiary and a victim of this
 62 data-intensive transition.

63 As mentioned above, the exponential growth of data acquisition can be attributed to Moore's law
 64 (?). Moore's law predicts that integrated circuit chip density doubles approximately each year at a
 65 fixed price point (?). (?) questioned whether technical development would sustain the growth.

66 Moore's law can be seen in many data-intensive fields, including astronomy. It explains both the
 67 recent development of big data in astronomy, and predicts future challenges.

68 This paper therefore seeks to review the rise of big data in astronomy and the technical and
 69 scientific issues surrounding it by examining four case studies: MeerKAT ¹, The Sloan Digital Sky
 70 Survey (SDSS) ², The Legacy Survey of Space and Time (LSST) ³, and The Square Kilometre Array
 71 (SKA) ⁴. These facilities represent the scope of contemporary astronomical data, the methods of its
 72 acquisition, their relative successes, the ongoing challenges, and the solutions currently in use.

¹ <https://www.skao.int/en>

² <https://www.sdss.org/>

³ <https://www.lsst.org/>

⁴ <https://www.skao.int/en>

73 1.3. *An Overview of The Four Surveys*

74 The SDSS is vital to this paper, as it is one of the earliest large-scale optical surveys that marks
 75 the start of the fourth paradigm. The SDSS is a precursor to LSST. The SDSS consists of three main
 76 telescopes.

77 The first of the three is The Sloan Foundation 2.5m Telescope. The Telescope is stationed at the
 78 Apache Point Observatory in New Mexico, where it observes the sky in the northern hemisphere. It
 79 is able to observe a 3° field of view by use of two corrector lenses (?).

80 The SDSS also uses the Irénée du Pont telescope at Las Campanas Observatory ⁵. This telescope
 81 is stationed in Chile, where it observes the southern hemisphere instead. Similar to the foundational
 82 telescope at Apache Point, this telescope has a 2.1° field of view but only uses one corrector lens (?).

83 The third telescope is the NMSU 1-meter Telescope ⁶. The NMSU telescope is stationed at the
 84 Apache Point Observatory alongside the foundational telescope. The NMSU telescope is designed to
 85 observe bright stars that are too bright for the aforementioned two telescopes to observe (?).

86 the SDSS is made up of multiple subsurveys. The eBoss survey ⁷, a continuation of BOSS, uses
 87 spectrographs to observe light in a wavelength range of 3600-10,400 Å (?). An additional subsurvey
 88 is APOGEE-2, a continuation of APOGEE. It uses spectrographs similar to eBOSS, but APOGEE-2
 89 collected near-infrared spectra (?). MaNGA ⁸ is a subsurvey that collects integral field unit spectra
 90 of 10,000 nearby galaxies (?). MARVELS ⁹ is another SDSS subsurvey, it was built specifically to
 91 obtain radial velocity measurements of stars with high-precision in hopes of finding exoplanets (?).

92 The MeerKAT ¹⁰ is an important precursor telescope to the SKA (?) MeerKAT became fully
 93 operational in 2018 in the Northern Cape Province of South Africa. MeerKAT comprises 64 antennas
 94 distributed over a radius of approximately 600 miles (?). These antennas operate across frequency
 95 bands ranging from 350 MHz to 3500 MHz (?).

96 MeerKAT has conducted and continues to conduct ten major survey projects (?). For conciseness,
 97 this discussion will focus on five of these surveys. One is the LADUMA ¹¹ survey. The objective of
 98 the LADUMA survey is to use HI obversations to research galaxy evolution over approximately 9.8
 99 billion years (?). LADUMA has used MeerKAT's Phase 1 receivers, which cover 0.9-1.75 GHz. It
 100 later transitioned to longer observations in Phase 4, which cover the 0.58-2.5 GHz band (?). Although
 101 the LADUMA survey is still ongoing, a portion of the data has already been released and will be
 102 discussed in the Methods section.

103 The MeerKAT absorbtion line survey ¹² (MALS) is a survey of HI and OH absorbers at a redshift
 104 of $z < 0.4$ and $z < 0.7$. HI is a descriptive tracer of the cold neutral medium in a galaxy (?). The
 105 cold neutral medium contains the physical conditions of the interstellar medium of each galaxy. This,
 106 in turn, allows scientists to estimate star formation rate in the galaxy (?).

⁵ <https://www.lco.cl/irenee-du-pont-telescope/>

⁶ <https://newapo.apo.nmsu.edu/>

⁷ <https://www.sdss4.org/surveys/eboss/>

⁸ <https://www.sdss4.org/surveys/manga/>

⁹ <https://www.sdss4.org/surveys/marvels/>

¹⁰ <https://www.sarao.ac.za/science/meerkat/>

¹¹ <https://science.uct.ac.za/laduma>

¹² <https://mals.iucaa.in/>

107 Another survey, ThunderKAT ¹³, aims to find, identify and understand high-energy radio transients,
 108 usually grouped with observations at similar wavelengths. Examples include supernovae,
 109 microquasars, and similar events (?).

110 Another notable MeerKAT survey is MHONGOOSE ¹⁴. This survey aims to catalogue the properties
 111 of HI gas using 30 nearby star-forming spiral and dwarf galaxies. MHONGOOSE is remarkable for
 112 its higher sensitivity compared to previous surveys such as HALOGAS ¹⁵ and THINGS ¹⁶ (?). This
 113 sensitivity is crucial for investigating how low-column-density gas influences the cosmic web and
 114 galactic accretion processes (?).

115 The final MeerKAT survey considered here is MIGHTEE ¹⁷. MIGHTEE spans 900-1670 MHz,
 116 achieving a resolution of approximately 6 arcseconds. MIGHTEE seeks to study the evolution of
 117 active galactic nuclei, neutral hydrogen, and the properties of cosmic magnetic fields (?).

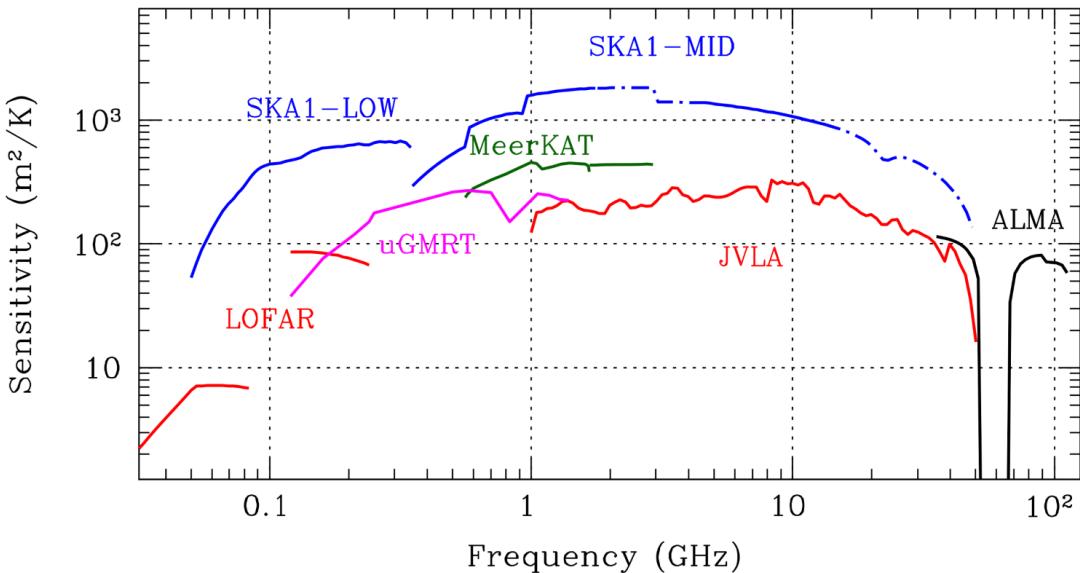


Figure 1. Figure from the SKA Official website, demonstrating the sensitivity compared to similar observatories

118 The SKA has built on technical and scientific achievements paved by MeerKAT and other radio
 119 interferometers. The SKA covers an area of approximately 131,205 antennas (?). The SKA represents
 120 the start of a new frontier for big data in astronomy. As an interferometer it uses aperture synthesis,
 121 which allows for the signals from antennas to be phased, this allows to reduce noise (?). The SKA
 122 will be discussed in further detail in the Methods section.

123 Alongside the SKA is its optical big data counterpart, the LSST. As noted above, the LSST is a
 124 successor to the SDSS. Rubin/LSST, however, has much more sophisticated goals. The LSST plans
 125 to address four key scientific issues: Investigating dark energy and dark matter, cataloguing the solar
 126 system, collecting data for sky surveys, and mapping the Milky Way. To achieve this, the LSST uses
 127 a 3.2-gigapixel camera with a sampling of 9.6 deg² field of view (?). These cameras are equipped

¹³ <https://www.physics.ox.ac.uk/research/group/meerkat>

¹⁴ <https://mhongoose.astron.nl/>

¹⁵ <https://www.astron.nl/halogas/>

¹⁶ <https://www2.mpi-a-hd.mpg.de/THINGS/Overview.html>

¹⁷ <https://www.mighteesurvey.org/home>

128 with highly resistant sensors reinforced with silicon (?). Rubin/LSST has an unprecedeted data rate
 129 for an optical telescope.

130 The SDSS, MeerKAT, SKA, and LSST generate unprecedented data rates and allow the exper-
 131 imentation of complex astrophysical events and phenomena. In the following Methods section, I
 132 describe how the data are collected, processed, analyzed, and stored. I then compare SDSS and
 133 MeerKAT to their larger successor telescopes, Rubin/LSST and SKA and consider the evolution of
 134 data challenges.

135 2. METHODS

136 2.1. *The Optical Big Data Pipeline*

137 This section carefully examines how each of the optical focused surveys collects and processes its
 138 data. It begins by describing the nature, scope, and type of the data. Next, it discusses how each
 139 survey collects and archives its data, followed by an explanation of their general data processing
 140 methods. Finally, it considers the use of real-time data processing. By comparing these surveys, this
 141 study highlights the rapid growth of big data in astronomy, a trend that has created challenges for
 142 data storage, processing, and analysis. These challenges will be discussed further in the discussion
 143 section.

144 The first survey examined is the Sloan Digital Sky Survey (SDSS). The SDSS I plan on splitting
 145 the SDSS into its three main components, the 2.5 m Telescope, the Irénée du Pont Telescope, and the
 146 NMSU 1-meter telescope. As discussed previously, I will explain how each telescope collects data,
 147 then I will explain how the SDSS processes the data both generally and in real-time. Lastly, I will
 148 talk about the open data policy the SDSS has employed.

149 2.1.1. *The Sloan Digital Sky Survey (SDSS)*

150 *Data Collection and Storage*—Although the SDSS is a complex survey, it can be divided into several
 151 major components, each of which contributes to the collection of astrophysical data. The 2.5-meter
 152 telescope plays a central role in the operations of the SDSS. It was designed to conduct precise optical
 153 observations of the sky over many years.

154 **(1) The SDSS 2.5 m Telescope:** According to Gunn et al. (2006), the SDSS camera contains
 155 “30 2048 x 2048 Scientific Imaging Technologies Charge-Coupled Devices (CCDs) and 24 2048 x 400
 156 CCDs” (?). A CCD is a detector that converts incoming light into an electronic signal. When photons
 157 strike the CCD, they generate electrons through the photoelectric effect. Using applied voltages, the
 158 resulting charge is measured based on the number of electrons produced. That measurement is then
 159 converted into a digital value and stored as a pixel, forming an image (?).

160 Another major innovation that enables the SDSS to collect data is the pair of fiber-fed double
 161 spectrographs, which record imaging data across wavelengths from 3800 to 9200 Å and at field angles
 162 between 0 and 90° (?). present measurements of the optical performance of these instruments, which
 163 are summarized in Figure 2.

164 In Figure 5, λ represents the wavelength of light, and “Angle” refers to the field angle. f_b is the
 165 best-focus distance, which is the position that provides the sharpest image. h/dh represents the
 166 lateral color, and D indicates the longitudinal difference from the best focus. Finally, ϵ is the root
 167 mean square (rms) image diameter. Among these parameters, the lateral color and longitudinal

λ (Å)	Angle (arcmin)	f_b (mm)	h/dh (mm)	D (mm)	ϵ (mm)
4000.....	0.00	-0.007	0.000	0.135	0.036
	30.00	-0.143	0.004	0.081	0.030
	45.00	-0.424	0.005	0.015	0.025
	60.00	-0.978	0.005	0.076	0.028
	70.00	-1.536	0.004	0.148	0.036
	80.00	-2.265	0.002	0.231	0.049
	90.00	-3.203	-0.004	0.325	0.065
4600.....	0.00	-0.007	-0.000	-0.058	0.031
	30.00	-0.143	0.002	-0.035	0.027
	45.00	-0.424	0.002	-0.006	0.024
	60.00	-0.978	0.002	0.033	0.025
	70.00	-1.536	0.002	0.065	0.027
	80.00	-2.265	0.001	0.101	0.030
	90.00	-3.203	-0.001	0.141	0.035
5300.....	0.00	-0.007	0.000	0.000	0.029
	30.00	-0.143	-108.818	0.000	0.026
	45.00	-0.424	-163.322	0.000	0.024
	60.00	-0.978	-217.855	0.000	0.025
	70.00	-1.536	-254.241	0.000	0.027
	80.00	-2.265	-290.713	0.000	0.026
	90.00	-3.203	-327.372	0.000	0.025
6500.....	0.00	-0.007	-0.000	0.062	0.031
	30.00	-0.143	-0.002	0.037	0.027
	45.00	-0.424	-0.002	0.007	0.024
	60.00	-0.978	-0.002	-0.035	0.029
	70.00	-1.536	-0.002	-0.068	0.034
	80.00	-2.265	-0.001	-0.106	0.036
	90.00	-3.203	0.002	-0.149	0.040
9000.....	0.00	-0.007	0.000	0.131	0.036
	30.00	-0.143	-0.004	0.078	0.029
	45.00	-0.424	-0.004	0.014	0.026
	60.00	-0.978	-0.004	-0.074	0.036
	70.00	-1.536	-0.004	-0.145	0.046
	80.00	-2.265	-0.002	-0.226	0.056
	90.00	-3.203	0.003	-0.317	0.068

Figure 2. Figure 5 from Gunn et al. (2006), showing results of the SDSS spectrographs given Wavelength and Angle.

difference are the most important for image quality because smaller values indicate sharper images. Based on the data from (?), both of these quantities remain close to zero for most wavelengths and field angles, except between roughly 5300 and 6500 Å, which demonstrates the high optical accuracy of the SDSS spectrographs.

The combination of these two innovations alongside others help the 2.5 m telescope collect data at a rate of about 20Gb/hr (?).

(2) The Irénée du Pont Telescope: Unlike the 2.5 m telescope, the Du Pont Telescope does not rely on CCDs to collect data. According to Bowen's 1973 paper, the telescope is described as a modified Ritchey-Chrétien design with Gascoigne correctors (?). The Du pont telescope uses a 100-inch primary mirror. Approximately 40% of the light is reflected to the secondary mirror, obtaining only a 16% loss of light at that stage. (?) The combination of light from the two aforementioned mirrors are then sent to a 20 inch x 20 inch plate, where monocromatic images are formed.

The du Pont Telescope uses 18.9 inch nonvignetted plates in order to minimize vignetting (?). Vignetting is the process where light beds through the lense of a telescope. The bending form a cone of light, which causes images to be darker near the edges and brighter in the center of the image (?). Because of the nonvignetted plates, the du Pont Telescope experiences an exceptionally low 3% percent loss of light (?).

Another technology the du Pont Telescope applies is a Gascoigne corrector plate. The plate helps with data collection. The Gasciogne corrector plate is able to be moved, which can help optimize the

187 collection of light in a wanted wavelength (?). Given a seperation of 1000 mm from the end of the
 188 corrector plate to the focus gives an image with a minimized astigmatism for a refractive index of n
 189 = 1.47 (?). At a given wavelength, the change of length which minimizes astigmatism is described
 190 in Bowen's paper as

$$191 \quad \Delta L = 590\Delta n/(n - 1) = -1250\Delta n \quad (1)$$

192 Where ΔL is the change in seperation in millimeters and Δn is the difference between a refractive
 193 index of 1.47 and the index wanted.

194 The last technology the du Pont Telescope uses is conical baffles. The reason for this is to promote
 195 shielding in the telescope (?). As explained in the Bowen paper, shielding is necessary in order to
 196 protect the photographic plate from light that escapes from the secondary lense due to long time
 197 exposure. As explained in the Bowen paper, the conic baffles are "located in the space between
 198 the incoming beam as it appraoches the primary and the return beam from the secondary to the
 199 plate" (?). Theoretically, the conic baffles have the disavantage of producing a diffraction pattern.
 200 However, as explained by Bowen, this should not majorly affect the images of stars (?).

201 **(3) The New Mexico State University (NMSU) Telescope:** The NMSU telescope takes the
 202 most technologically advanced approach to collecting data compared to the 2.5 m telescope and the
 203 du Pont Telescope. The NMSU telescope uses a camera that has a 2048 x 2048 CCD. The camera is
 204 controled by a linux computer, which is connected by fiber optic cables (?).

205 The data collection of the NMSU telescope is almost fully automated using C++, the only time
 206 it is not is when engineering is being done on the telescope (?). The NMSU telescope has a camera
 207 which analyzes the brightness level of the sky to see if it is dark enough to start collecting data.
 208 When the sky becomes dark enough, the telescope initiates its program. The telescope goes through
 209 the list and observes said objects (?). Objects can, however, be observed as many times as requested.

210 *Data Processing*—The SDSS processes its data through an innovative acquisition system that records
 211 and organizes observations in real time while maintaining strict quality control (?). This system
 212 ensures that data are processed and stored efficiently without any loss of precision. (?). The data
 213 pipeline of the SDSS can be described by two different fields of data, the imaging pipeline and the
 214 spectroscopy pipeline

215 **(1) Imaging Data Pipeline:** The Imaging data Pipeline itself consists of multiple subpipelines.
 216 The first subpipeline is the Astroline. The astroline uses vxWorks in order to initalize the processing
 217 sequence. This happens by composing star cutouts and column quartiles collected from the CCD's
 218 mentioned before (?).

219 The second subipline is the MT pipeline. the MT Pipeline processes the data collected from the
 220 Photometric Telescope. This data is used to calculate important parameters for the 2.5 m Telescope
 221 scans, such as extinction and zero-points (?).

222 The third pipeline is the Serial Stamp Collecing (SSC) Pipeline. the SSC reorganizes the star
 223 cutouts collected from previous pipelines. The SSC does this in order to prepare data for the upcoming
 224 pipelines (?).

225 Next is the Astrometric Pipeline. The Astrometric pipeline processes the average location of stars
 226 using the data collected from the astroline and SSC pipelines. It then converts the pixel data from
 227 the images into celestial coordinates (α, δ) (?).

228 After that is the Postage Stamp Pipeline (PSP). The PSP estimates the quality of the data collected
 229 by calculating factors such as the flat field vectors, bias drift, and sky levels (?).

230 After all that is done, the data is fed into the Frames Pipeline. The Frames pipeline does a majority
 231 of the work, processing the data from all the previous pipelines and producing the complete datasets
 232 of images. It does this by correcting the frames based on the data before and cataloging the images
 233 (?).

234 Lastly, the processed data is then ran through the Calibration pipeline. The Calibration pipeline
 235 takes data from the MT and Frames pipeline. The Calibration pipeline converts the counts into more
 236 measurable quantities such as flux (?).

237 The SDSS imaging pipeline is composed of multiple connected pipelines that operate collaboratively
 238 to transform raw imaging data into structured datasets from which measurable physical quantities
 239 such as flux can be derived.

240 (2) Spectroscopy Data Pipeline:

241 *Real-Time Processing*—

242 2.1.2. *The Rubin/Large Synoptic Survey Telescope (LSST)*

243 *Data Collection and Storage*—Despite the recent times in creation. The SDSS collected around 16
 244 TB of data over a decade in their data release 7 (?). Yet the LSST is expected to collect 20 TB of
 245 data per night (?).

246 The LSST pipeline is written with around 750000 in Python in order to use relevant libraries such
 247 as SciPy and AstroPy (?). The pipeline is also written with around 220000 lines in C++ in order to
 248 ensure efficient performance (?). In order to combine these two languages, pybind11 allows for the
 249 transition from Python to C++, and ndarray objects are able to be converted from C++ arrays (?).

250 The python enviroment of the LSST pipeline uses a package named **rubin-env**. This package gives
 251 the user all the code needed to run LSST's data. In order to execute the code, the pipeline consists
 252 multiple packages that each serve their own purpose. The LSST has defined a class labeled **Task**,
 253 which is used to define algorithms (?).

254 One instance of a task is the **PipelineTask**, which serves to organize subtasks. These subtasks
 255 each have their own purpose (?). The most important subtask, labeled **daf_butler**, handles the data
 256 storage. This subtask is titled by the LSST as The Data Butler. The Butler serves as a database to
 257 store collected data. It stores objects with data IDs similar to SQL, with headers that hold useful
 258 information (?). An example of this would be a data coordinate labeled **instrument="LSSTCam"**,
 259 **exposure=299792458**, **detector=42**, **band=z**, **day_obs=20251011**.

260 *Data Processing*—There are multiple tasks which define how the LSST processed data to find objects.
 261 These are all defined in the **meas_algorithms** package (?). The task that first handles processing
 262 the catalogued images is the **SourceDetectionTask** (?). This task uses Gaussian smoothing in the
 263 point spread function. It then convolves the collected image with the point spread function in order
 264 to suppress potential noise (?).

265 Another task that the LSST uses to process data is `MaskStreaksTask` (?). The task serves to
 266 mask pixels from streaks from other satellites. It identifies streaks using a Canny Filter and the
 267 Kernel-Based Hough Transform (??). This task is combined with the deblending of collected images
 268 allows for the LSST to accurately identify objects.

269 *Real-Time Processing—*

270 2.2. *The Radio Big Data Pipeline*

271 2.2.1. *The MeerKAT*

272 **NEW TEXT STARTS HERE:**

273 *Data Processing*—The MeerKAT data processing pipeline is split into three parts, the calibration
 274 pipeline, the continuum pipeline, and the spectral pipeline (?).

275 **(1) The Calibration Pipeline:** The pipeline starts by choosing an antenna as a reference point
 276 for the other antennas (?). This process is done by collecting the first scan of an antenna. The data
 277 is then averaged and fourier transformed. After that, it is the ratio called the peak-to-noise ratio is
 278 measured. The peak-to-noise ratio is the ratio of the maximum of the fourier transformed data and
 279 the rms noise at a distance far from the peak (?). After tha the process repeats for every antenna and
 280 the one with the highest median with respect to every baseline for the peak-to-noise ratio is chosen
 281 as the reference antenna (?). This antenna will have all of its calibration solutions set to zero (?).

282 The purpose of the calibration pipeline is to ensure that the reference antenna and the other
 283 antennas by extension are callibrated. At the start of every observation, the reference antenna is
 284 evaluated based on certain flags, such as data loss (?). If 80% or more flags are pinged, a new
 285 reference antenna is determined based on the aforementioned process (?).

286 **(2) The Continuum Pipeline:** The second pipeline is the continuum pipeline. This pipeline
 287 aims to produce continuum images using the orbit software package (??). The orbit package which
 288 interfaces UV data. It does this by using the `MFIImage` task, which does wide-band, wide-field imaging
 289 (??). Once the data is read, it is split into scans, which are then averaged and combined into a dataset.

290 This dataset is then split into 8 107 MHz intermediate frequencies (?). The UV data then stores
 291 headers of `ntimes`, `nbaselines`, `nchannels`, `npolarisation` (?). Given the Meerkat's design of
 292 48 antennas, short baselines dominate. This causes an inflation in data volume (?). In order to
 293 circumvent this, a baseline dependent averaging step is applied (?). Time averaging the data with
 294 respect to the baseline before merging them into a dataset. This step reduces the data volume by
 295 about 3-4 times with a minimal loss of 1

296 The aforementioned `MFIImage` task uses joint frequency deconvolution to handle wide-band effects
 297 in wide-band images (?). Normally, Meerkat uses around 140 circular facets with a size of around 6
 298 arcminutes which cover 1 degree from the phase center (?). Additionally, there are facets with around
 299 1 degree that use the SUMSS or NVSS catalogue for a radius of 2.5 from the center of the phase
 300 (?). These 1 degree facets are used for phenomena that are anciipated to have a flux density greater
 301 than 5 mJy (?). Due to the fact that these facets deal with wide-band images, they also deal with
 302 wide-band imaging effects. These effects are bypassed by splitting the frequency band used into 10
 303 components (?). When the joint frequency deconvolution is initiated, the brightest light sources in

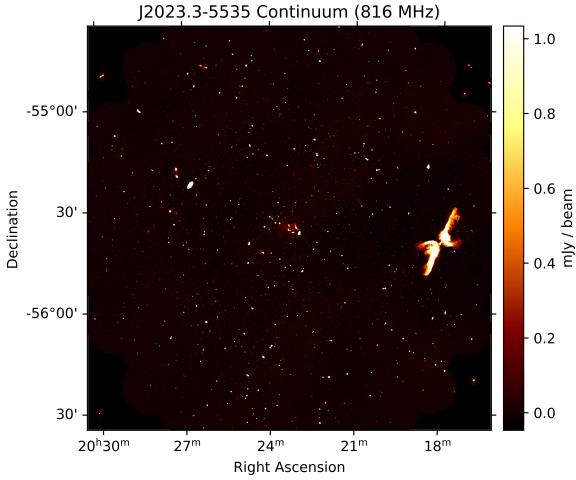


Figure 3. An example UHF continuum thumbnail image produced automatically by the pipeline is shown above, (?)

304 the dataset are found and subtracted from the slices of data individually. This is called the CLEAN
 305 algorithm (?).

306 Self-calibration is then performed in two rounds in order to transform data. The first round consists
 307 of using the CLEAN algorithm with a spectral density of 1 mJy. This yields approximately 1000
 308 CLEAN components (?). The data is then self calibrated once more, using a CLEAN algorithm with
 309 a spectral density of 100 μ Jy. The second round yields approximately 10000 CLEAN components
 310 (?). After this process, the data is imaged and used to create the continuum images. Lastly, the
 311 self calibrated data is converted into AIPS format, which consists of tables with headers such as UV
 312 data, polarisation, etc (?). The data is then converted into numpy arrays for each timestamp with
 313 shape (nif, npol, nantennas) (?). The sky model, made up of the CLEAN components, is stored in
 314 AIPS CC format. The flux density of all the 10 components of the frequency are summed (?).
 315 The merged flux density is then fitted with a second degree polynomial to measure frequency and
 316 is used to subtract from the image (?). Finally, the fully processed images are converted into FITS
 317 and PNG files and archived (?).

318 **(3) The Spectral Pipeline:** The purpose of the spectral pipeline is to speed up image creation
 319 while keeping control of quality (?). The issue with the imaging data, at this point is that they're
 320 independent of each other, which causes issues with accessing data when trying to run said images (?).
 321 The raw data, also known as visibilities, are received in time-major order, but this data structure
 322 must be transposed in order for the data to be in channel-major order (?). In order to solve this
 323 issue, a visibility writer is used. The writer stores the visibilities on a Ceph cluster with chunks each
 324 across 64 channels (?). This choice in the chunk size is not optimized, but it does avoid issues with
 325 RAM and memory allocation (?).

326 The chunk system only partially implements the solution of transposing the data. The solution is
 327 complete after another series of events are done (?). The visibilities in each chunk are ordered by
 328 channel, W slice, and baseline. This allows for both the accessibility of the data to improve, and
 329 it inherently stores measurements such as UVW coordinates and parallactic angles (?). While all
 330 that is occurring, conservative baseline-dependent averaging is also applied to the visibilities (?). The

331 coordinates of every visibility are solved for, then visibilities with matching coordinates are merged,
 332 which is the last of the preprocessing of the visblities (?).

333 After all of that is done, By using CUDA, a parallel processing software, every channel is imaged
 334 separately. The use of CUDA allows for the processing to speed up due to the usage of NVIDIA GPUs
 335 (?). Lastly, the data follows a cycle in order to be updated as described in the figure below.

1. For each W slice¹
 - a. Apply image-plane W term and taper correction to the model image
 - b. FFT the result to get a UV grid.
 - c. For each batch of visibilities
 - i. Predict visibilities by degridding, and subtract them from the measured visibilities in place.
 - ii. Grid the resulting batch of visibilities.
 - d. Inverse FFT the grid.
 - e. Apply image-plane W term and taper correction in the image plane.
 - f. Add the result to the dirty image.
2. Apply CLEAN, adding new components to the model image.

Figure 4. The cyclic algorithm for the Meerkat's spectrographic pipeline, (?)

336 NEW TEXT ENDS HERE:

337 2.2.2. *The Square Kilometre Array (SKA)*

338 3. RESULTS

339 3.1.

340 4. DISCUSSION

341 4.1. *Open Source Policies and Transparency*

342 4.1.1. *The SDSS Policy*

343 The SDSS collaboration states on its official SDSS-IV website ¹⁸ that all of its software be open
 344 source using the open source liscence BSD 3-Clause. However, the SDSS outlines practices for users
 345 who plan to use the SDSS software must abide by. One of the most important ones is the proper
 346 citation of software and websites that were used. The SDSS4 emphasizes the importance of citing
 347 properly as it serves to acknowledge the hard work of the teams behind said projects.

348 The SDSS also has implemented Digital Object Identifiers, commonly known as DOIs, into all
 349 software code. DOIs allow for software and data to be easily identified, which is important for
 350 ownership. The SDSS team also promotes transparency in coding by implementing Git and SVN in
 351 order to maintain a record of the development of the software. This not only makes the development
 352 transparent, but also helps users see the evolution of the software.

¹⁸ <https://www.sdss4.org/dr17/software/>

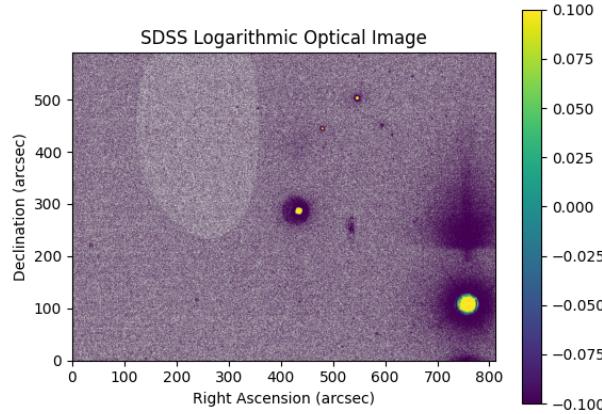


Figure 5. A spectrographical image obtained using data collected from SDSS

353 Overall, the SDSS has demonstrated a strong commitment to making their data and software open
 354 source and transparent. This in turn helps the development of science, by ensuring that knowledge
 355 is accessible to all regardless of resources.

356 4.1.2. *The LSST Policy*

357 **New text starts here:**

358 4.1.3. *The MeerKAT Policy*

359 **New text ends here:**

360 4.1.4. *The SKA Policy*

361 A. PYTHON CODE FOR SDSS DATA RETRIEVAL (FIGURE 3)

```

362 1 #Import relevant libraries/functions
363 2 from astroquery.sdss import SDSS
364 3 from astropy import coordinates as coords
365 4 import astropy.units as u
366 5 import matplotlib.pyplot as plt
367 6 import numpy as np
368 7
369 8 #Initialize Right Ascension and Declination
370 9 ra = 20
371 10 dec = -10
372 11
373 12 #Convert ra and dec into a SkyCoord Object
374 13 coord = coords.SkyCoord(ra, dec, unit='deg', frame = 'icrs')
375 14
376 15 #Query the SDSS System to find object given coordinates in a radius of
377 16      0.01 degrees
378
379 16 result = SDSS.query_region(coord, radius=0.01*u.deg, spectro=True)

```

```

380 17 print(result)
381 18 #Retrieve the Image from found object, make into a FITS File
382 19 image = SDSS.get_images(matches=result, band=['u', 'g', 'r', 'i', 'z'])
383 20
384 21
385 22 #Retrieve hdulist from FITS file
386 23 hdulist = image[0]
387 24
388 25 #Retrieve the image data from the hdulist
389 26 imageData = hdulist[0].data
390 27
391 28 #Log the image data in order to get rid of background
392 29 imageDataLog = np.log10(imageData) + 1e-8
393 30
394 31 #Save the header
395 32 header = hdulist[0].header
396 33
397 34
398 35 #Obtain the relevant headers
399 36
400 37 #Retrieve pixel scale numbers, divided amongst two parts for ra and dec
401 38 CD1_1 = header['CD1_1']
402 39 CD1_2 = header['CD1_2']
403 40 CD2_1 = header['CD2_1']
404 41 CD2_2 = header['CD2_2']
405 42
406 43 #Width of image in pixels
407 44 keyWordNAXIS1 = header['NAXIS1'] #[pixels]
408 45
409 46 #Height of image in pixels [pixels]
410 47 keyWordNAXIS2 = header['NAXIS2'] #[pixels]
411 48
412 49 #Normalize the pixel scale, then multiply by 3600 to convert units
413 50 CDELT1ArcSec = np.linalg.norm([CD1_1,CD1_2]) * 3600 #[arcsec/pixels]
414 51 CDELT2ArcSec = np.linalg.norm([CD2_1,CD2_2]) * 3600 #[arcsec/pixels]
415 52
416 53 #Set up the image
417 54 plt.xlabel("Right Ascension (arcsec)")
418 55 plt.ylabel("Declination (arcsec)")
419 56 plt.title('SDSS Logarithmic Optical Image')
420 57 vmin2 = np.percentile(imageDataLog, 85)
421 58 vmax2 = np.percentile(imageDataLog, 98)
422 59 plt.imshow(imageDataLog, cmap='viridis', extent = (0, CDELT1ArcSec *
423           keyWordNAXIS1, 0, CDELT2ArcSec * keyWordNAXIS2), vmin = vmin2, vmax =
424           vmax2)
425 60 plt.colorbar()
426 61
427 62 plt.show()

```

REFERENCES

- 429 ????, The MIGHTEE Survey,
 430 <https://www.mighteesurvey.org/home>
- 431 ????, SKA Telescope Specifications,
 432 <https://www.skao.int/en/science-users/118/ska-telescope-specifications>
- 433 Blyth, S., Baker, A. J., Holwerda, B., et al. 2018,
 434 in Proceedings of MeerKAT Science: On the
 435 Pathway to the SKA — PoS(MeerKAT2016)
 436 (Stellenbosch, South Africa: Sissa Medialab),
 437 004, doi: [10.22323/1.277.0004](https://doi.org/10.22323/1.277.0004)
- 438 Bowen, I. S., & Vaughan, A. H. 1973, Applied
 439 Optics, 12, 1430, doi: [10.1364/AO.12.001430](https://doi.org/10.1364/AO.12.001430)
- 440 Bundy, K., Bershadsky, M. A., Law, D. R., et al.
 441 2014a, The Astrophysical Journal, 798, 7,
 442 doi: [10.1088/0004-637X/798/1/7](https://doi.org/10.1088/0004-637X/798/1/7)
- 443 —. 2014b, The Astrophysical Journal, 798, 7,
 444 doi: [10.1088/0004-637X/798/1/7](https://doi.org/10.1088/0004-637X/798/1/7)
- 445 Camilo, F. 2024
- 446 Cotton, W. D. 2008, Publications of the
 447 Astronomical Society of the Pacific, 120, 439,
 448 doi: [10.1086/586754](https://doi.org/10.1086/586754)
- 449 Cotton, W. D., & Schwab, F. R. 2010
- 450 Dawson, K. S., Kneib, J.-P., Percival, W. J., et al.
 451 2016, The Astronomical Journal, 151, 44,
 452 doi: [10.3847/0004-6256/151/2/44](https://doi.org/10.3847/0004-6256/151/2/44)
- 453 De Blok, W. J. G., Healy, J., Maccagni, F. M.,
 454 et al. 2024, Astronomy & Astrophysics, 688,
 455 A109, doi: [10.1051/0004-6361/202348297](https://doi.org/10.1051/0004-6361/202348297)
- 456 Dewdney, P., Hall, P., Schilizzi, R., & Lazio, T.
 457 2009, Proceedings of the IEEE, 97, 1482,
 458 doi: [10.1109/JPROC.2009.2021005](https://doi.org/10.1109/JPROC.2009.2021005)
- 459 Fernandes, L. A., & Oliveira, M. M. 2008, Pattern
 460 Recognition, 41, 299,
 461 doi: [10.1016/j.patcog.2007.04.003](https://doi.org/10.1016/j.patcog.2007.04.003)
- 462 Goedhart, S. 2025, MeerKAT Specifications
- 463 Gunn, J. E., Siegmund, W. A., Mannery, E. J.,
 464 et al. 2006, The Astronomical Journal, 131,
 465 2332, doi: [10.1086/500975](https://doi.org/10.1086/500975)
- 466 Gupta, N., Jagannathan, P., Srianand, R., et al.
 467 2021, The Astrophysical Journal, 907, 11,
 468 doi: [10.3847/1538-4357/abcb85](https://doi.org/10.3847/1538-4357/abcb85)
- 469 Hey, T., Tansley, S., & Tolle, K. 2009, Microsoft
 470 Research
- 471 Holtzman, J. A., Harrison, T. E., & Coughlin,
 472 J. L. 2010, Advances in Astronomy, 2010,
 473 193086, doi: [10.1155/2010/193086](https://doi.org/10.1155/2010/193086)
- 474 Ivezic, Z., Kahn, S. M., Tyson, J. A., et al. 2019,
 475 The Astrophysical Journal, 873, 111,
 476 doi: [10.3847/1538-4357/ab042c](https://doi.org/10.3847/1538-4357/ab042c)
- 477 Jonas, J., & the MeerKAT Team. 2018, in
 478 Proceedings of MeerKAT Science: On the
 479 Pathway to the SKA — PoS(MeerKAT2016)
 480 (Stellenbosch, South Africa: Sissa Medialab),
 481 001, doi: [10.22323/1.277.0001](https://doi.org/10.22323/1.277.0001)
- 482 Juric, M., Kantor, J., Lim, K.-T., et al. 2017
- 483 Lesser, M. 2015, Publications of the Astronomical
 484 Society of the Pacific, 127, 1097,
 485 doi: [10.1086/684054](https://doi.org/10.1086/684054)
- 486 Lupton, R., Gunn, J. E., Ivezić, Z., et al. 2001,
 487 The SDSS Imaging Pipelines, arXiv,
 488 doi: [10.48550/arXiv.astro-ph/0101420](https://arxiv.org/abs/astro-ph/0101420)
- 489 Lupton, R. H., Ivezić, Z., Gunn, J., et al. 2007
- 490 Majewski, S. R., Schiavon, R. P., Frinchaboy,
 491 P. M., et al. 2017, The Astronomical Journal,
 492 154, 94, doi: [10.3847/1538-3881/aa784d](https://doi.org/10.3847/1538-3881/aa784d)
- 493 Moore, G. E. 2006, IEEE Solid-State Circuits
 494 Society Newsletter, 11, 33,
 495 doi: [10.1109/N-SSC.2006.4785860](https://doi.org/10.1109/N-SSC.2006.4785860)
- 496 NSF-DOE Vera C. Rubin Observatory. 2025,
 497 PSTN-019: The LSST Science Pipelines
 498 Software: Optical Survey Pipeline Reduction
 499 and Analysis Environment, NSF-DOE Vera C.
 500 Rubin Observatory,
 501 doi: [10.71929/RUBIN/2570545](https://doi.org/10.71929/RUBIN/2570545)
- 502 Ratcliffe, S. 2021, SDP Pipelines Overview,
 503 <https://skaafrica.atlassian.net/wiki/spaces/ESDKB/pages/77071377/SDP+Pipelines+Overview>
- 504 Richards, S. 2020, What Is Vignetting?
- 505 Verbunt, F., & Van Gent, R. H. 2010, Astronomy
 506 and Astrophysics, 516, A28,
 507 doi: [10.1051/0004-6361/201014002](https://doi.org/10.1051/0004-6361/201014002)
- 508 Woudt, P. A., Fender, R., Corbel, S., et al. 2018,
 509 in Proceedings of MeerKAT Science: On the
 510 Pathway to the SKA — PoS(MeerKAT2016)
 511 (Stellenbosch, South Africa: Sissa Medialab),
 512 013, doi: [10.22323/1.277.0013](https://doi.org/10.22323/1.277.0013)
- 513

REFERENCES