

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

NOTE: I fixed the things you told me to fix, I added a new figure to sdss section, I added lsst section

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 Redo the Du Pont section, it's outdated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	2
1.1. The Paradigms of Data Science	2
1.2. The Rise of Big Data in Astronomy	3
1.3. An Overview of The Four Surveys	3
2. Methods	6
2.1. The Optical Big Data Pipeline	6
2.1.1. The Sloan Digital Sky Survey (SDSS)	6
2.1.2. The Rubin/Large Synoptic Survey Telescope (LSST)	9
2.2. The Radio Big Data Pipeline	10
2.2.1. The MeerKAT	10
2.2.2. The Square Kilometre Array (SKA)	12
3. Results	12
3.1.	12
4. Discussion	12
4.1. Open Source Policies and Transparency	12

32	4.1.1. The SDSS Policy	12
33	4.1.2. The LSST Policy	13
34	4.1.3. The MeerKAT Policy	13
35	4.1.4. The SKA Policy	13
36	A. Python Code for SDSS Data Retrieval (Figure 3)	13

37 1. INTRODUCTION

38 The concept of data has long been central throughout the history of astronomy. Data allows scientists
 39 to discover natural laws in the universe, have control over events, and make reliable predictions.
 40 It has played a critical role in other time-sensitive fields such as medicine and engineering, where
 41 accurate data is essential for decision-making and design. Although the nature of data varies funda-
 42 mentally across different fields, one trend has remained consistent: the continual evolution of data
 43 science. As explained in The Fourth Paradigm (Hey et al. 2009), this evolution can be character-
 44 ized through four successive paradigms. In the following sections, I describe the progression of data
 45 acquisition across these paradigms and illustrate each using examples from astronomy. I will then
 46 explain how SKA and LSST fit into this trajectory and exemplify the emerging era of data-intensive
 47 discovery.

48 1.1. *The Paradigms of Data Science*

49 The first and most primitive paradigm, as described by Hey et al. (2009), is empirical evidence.
 50 Empirical evidence refers to data collected through traditional means, such as direct observation
 51 or experimentation. The primary purpose of empirical evidence is to identify patterns that allow
 52 scientists to develop a fundamental understanding of the natural world. Throughout much of human
 53 history, empirical evidence has dominated knowledge generation. An example of the first paradigm
 54 in astronomy is the career of Tycho Brahe, a Danish astronomer. Throughout his career in the 16th
 55 century, Brahe collected and cataloged data on the position of astronomical bodies using naked-eye
 56 observations. Tycho Brahe’s catalogue was accurate to only around 1’ precision and took decades to
 57 acquire (Verbunt & Van Gent 2010). However, empirical evidence can be compromised by human
 58 error, the precision of the instruments, and, most importantly, the relatively slow pace of data
 59 acquisition compared to subsequent paradigms.

60 The second paradigm is analytical evidence. Analytical evidence is obtained by constructing math-
 61 ematical formulas and theoretical frameworks based on empirical data Hey et al. (2009). Unlike the
 62 empirical evidence, which merely demonstrates that phenomena occur, the second paradigm seeks to
 63 explain why they occur. An example of the second paradigm in astronomy is the work of Johannes
 64 Kepler, who used Brahe’s empirical observations to derive the laws of planetary motion (Hey et al.
 65 2009). By transforming raw observational data into mathematical laws, Kepler exemplified how
 66 analytical evidence advances scientific understanding beyond description to explanation.

67 The third paradigm is simulation evidence Hey et al. (2009), a relatively recent development.
 68 Simulation models natural phenomena that are too complex to model analytically or compute by
 69 hand. It allows interpolation and extrapolation of data using computational techniques grounded in
 70 known physical laws. For example, in astronomy, N-body simulations are used to study the complex
 71 dynamical evolution of planetary systems and galaxies.

The fourth and most recent paradigm is data-intensive science Hey et al. (2009). This paradigm is characterized by the unprecedented scale, velocity, and complexity of data acquisition, driven in part by exponential advances in computational power and detector technologies, often associated with Moore's law Hey et al. (2009). Unlike earlier paradigms, which focused on observation, theory, or simulation, data-intensive science emphasizes the ability to manage, analyze, and interpret vast datasets that exceed the capacity of traditional methods. While this exponential growth in data has enabled transformative discoveries, it also introduces significant challenges related to storage, processing, and accessibility.

1.2. *The Rise of Big Data in Astronomy*

Astronomy has become data intensive. Modern observatories may now generate petabyte-scale data that need new strategies for data management and analysis Hey et al. (2009). The fourth paradigm enables discoveries from interpreting massive data sets. However, these advances also expose alarming issues, including bottlenecks in the data pipeline, storage challenges, increased skills needed to handle the data, and open access concerns. The field of astronomy is both a beneficiary and a victim of this data-intensive transition.

As mentioned above, the exponential growth of data acquisition can be attributed to Moore's law Hey et al. (2009). Moore's law predicts that integrated circuit chip density doubles approximately each year at a fixed price point Moore (2006). Moore (2006) questioned whether technical development would sustain the growth.

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

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

1.3. *An Overview of The Four Surveys*

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

The first of the three is The Sloan Foundation 2.5m Telescope. The Telescope is stationed at the Apache Point Observatory in New Mexico, where it observes the sky in the northern hemisphere. It is able to observe a 3° field of view by use of two corrector lenses (Gunn et al. 2006).

The SDSS also uses the Irénée du Pont telescope at Las Campanas Observatory ⁵. This telescope is stationed in Chile, where it observes the southern hemisphere instead. Similar to the foundational telescope at Apache Point, this telescope has a 2.1° field of view but only uses one corrector lens Bowen & Vaughan (1973).

¹ <https://www.skao.int/en>

² <https://www.sdss.org/>

³ <https://www.lsst.org/>

⁴ <https://www.skao.int/en>

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

The third telescope is the NMSU 1-meter Telescope ⁶. The NMSU telescope is stationed at the Apache Point Observatory alongside the foundational telescope. The NMSU telescope is designed to observe bright stars that are too bright for the aforementioned two telescopes to observe (Majewski et al. 2017).

the SDSS is made up of multiple subsurveys. The eBoss survey ⁷, a continuation of BOSS, uses spectrographs to observe light in a wavelength range of 3600-10,400 Å (Dawson et al. 2016). An additional subsurvey is APOGEE-2, a continuation of APOGEE. It uses spectrographs similar to eBOSS, but APOGEE-2 collected near-infrared spectra (Majewski et al. 2017). MaNGA ⁸ is a subsurvey that collects integral field unit spectra of 10,000 nearby galaxies (Bundy et al. 2014a). MARVELS ⁹ is another SDSS subsurvey, it was built specifically to obtain radial velocity measurements of stars with high-precision in hopes of finding exoplanets (Bundy et al. 2014b).

The MeerKAT ¹⁰ is an important precursor telescope to the SKA (Jonas & the MeerKAT Team 2018) MeerKAT became fully operational in 2018 in the Northern Cape Province of South Africa. MeerKAT comprises 64 antennas distributed over a radius of approximately 600 miles (Goedhart 2025). These antennas operate across frequency bands ranging from 350 MHz to 3500 MHz (Goedhart 2025).

MeerKAT has conducted and continues to conduct ten major survey projects (Jonas & the MeerKAT Team 2018). For conciseness, this discussion will focus on five of these surveys. One is the LADUMA ¹¹ survey. The objective of the LADUMA survey is to use HI obversations to research galaxy evolution over approximately 9.8 billion years (Blyth et al. 2018). LADUMA has used MeerKAT's Phase 1 receivers, which cover 0.9-1.75 GHz. It later transitioned to longer observations in Phase 4, which cover the 0.58-2.5 GHz band (Blyth et al. 2018). Although the LADUMA survey is still ongoing, a portion of the data has already been released and will be discussed in the Methods section.

The MeerKAT absorbtion line survey ¹² (MALS) is a survey of HI and OH absorbers at a redshift of $z < 0.4$ and $z < 0.7$. HI is a descriptive tracer of the cold neutral medium in a galaxy (Gupta et al. 2021). The cold neutral medium contains the physical conditions of the interstellar medium of each galaxy. This, in turn, allows scientists to estimate star formation rate in the galaxy (Gupta et al. 2021).

Another survey, ThunderKAT ¹³, aims to find, identify and understand high-energy radio transients, usually grouped with observations at similar wavelengths. Examples include supernovae, microquasars, and similar events (Woudt et al. 2018).

Another notable MeerKAT survey is MHONGOOSE ¹⁴. This survey aims to catalogue the properties of HI gas using 30 nearby star-forming spiral and dwarf galaxies. MHONGOOSE is remarkable for its higher sensitivity compared to previous surveys such as HALOGAS ¹⁵ and THINGS ¹⁶ (De Blok

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/>

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>

et al. 2024). This sensitivity is crucial for investigating how low-column-density gas influences the cosmic web and galactic accretion processes (De Blok et al. 2024).

The final MeerKAT survey considered here is MIGHTEE¹⁷. MIGHTEE spans 900-1670 MHz, achieving a resolution of approximately 6 arcseconds. MIGHTEE seeks to study the evolution of active galactic nuclei, neutral hydrogen, and the properties of cosmic magnetic fields (MIG ????).

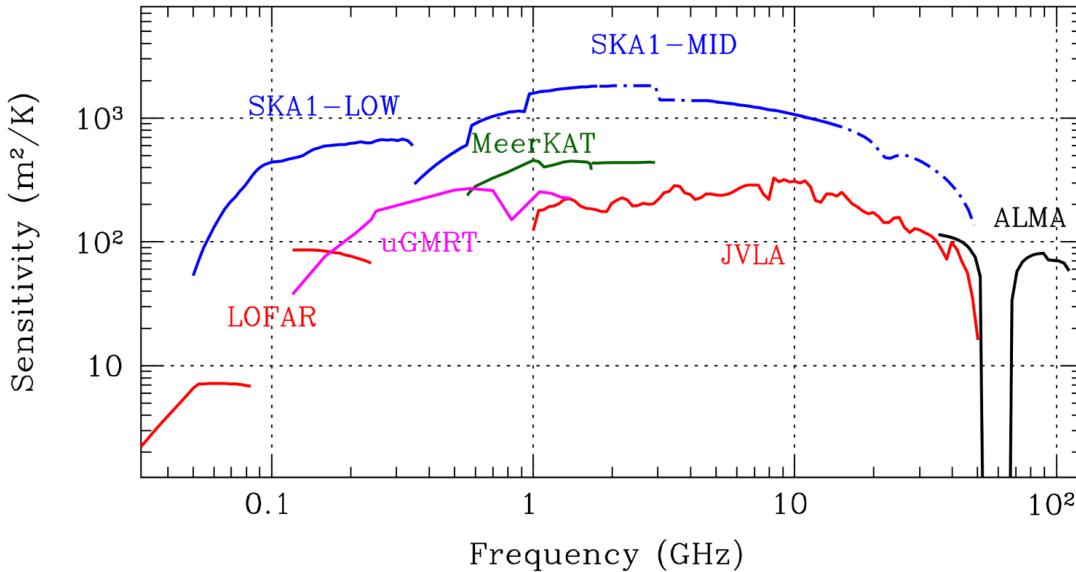


Figure 1. Figure from the SKA Official website, demonstrating the sensitivity compared to similar observatories

The SKA has built on technical and scientific achievements paved by MeerKAT and other radio interferometers. The SKA covers an area of approximately 131,205 antennas (SKA ????). The SKA represents the start of a new frontier for big data in astronomy. As an interferometer it uses aperture synthesis, which allows for the signals from antennas to be phased, this allows to reduce noise (Dewdney et al. 2009). The SKA will be discussed in further detail in the Methods section.

Alongside the SKA is its optical big data counterpart, the LSST. As noted above, the LSST is a successor to the SDSS. Rubin/LSST, however, has much more sophisticated goals. The LSST plans to address four key scientific issues: investigating dark energy and dark matter, cataloguing the solar system, collecting data for sky surveys, and mapping the Milky Way. To achieve this, the LSST uses a 3.2-gigapixel camera with a sampling of 9.6 deg² field of view (Ivezić et al. 2019). These cameras are equipped with highly resistant sensors reinforced with silicon (Ivezić et al. 2019). Rubin/LSST has an unpreceded data rate for an optical telescope.

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

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

166

2. METHODS

167

2.1. *The Optical Big Data Pipeline*

168

169

170

171

172

173

174

175

176

177

178

179

180

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

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

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

181

182

183

184

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

185

186

187

188

189

190

191

(1) The SDSS 2.5 m Telescope: According to Gunn et al. (2006), the SDSS camera contains “30 2048 x 2048 Scientific Imaging Technologies Charge-Coupled Devices (CCDs) and 24 2048 x 400 CCDs” Gunn et al. (2006). A CCD is a detector that converts incoming light into an electronic signal. When photons strike the CCD, they generate electrons through the photoelectric effect. Using applied voltages, the resulting charge is measured based on the number of electrons produced. That measurement is then converted into a digital value and stored as a pixel, forming an image Lesser (2015).

192

193

194

195

Another major innovation that enables the SDSS to collect data is the pair of fiber-fed double spectrographs, which record imaging data across wavelengths from 3800 to 9200 Å and at field angles between 0 and 90° Gunn et al. (2006). present measurements of the optical performance of these instruments, which are summarized in Figure 2.

196

197

198

199

200

201

202

203

In Figure 5, λ represents the wavelength of light, and “Angle” refers to the field angle. f_b is the best-focus distance, which is the position that provides the sharpest image. h/dh represents the lateral color, and D indicates the longitudinal difference from the best focus. Finally, ϵ is the root mean square (rms) image diameter. Among these parameters, the lateral color and longitudinal difference are the most important for image quality because smaller values indicate sharper images. Based on the data from Gunn et al. (2006), 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.

204

205

The combination of these two innovations alongside others help the 2.5 m telescope collect data at a rate of about 20Gb/hr Lupton et al. (2007).

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.

(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 Bowen & Vaughan (1973). 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. Bowen & Vaughan (1973) 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 Bowen & Vaughan (1973). 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 Richards (2020). Because of the nonvignetted plates, the du Pont Telescope experiences an exceptionally low 3% percent loss of light Bowen & Vaughan (1973).

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 collection of light in a wanted wavelength Bowen & Vaughan (1973). Given a seperation of 1000 mm from the end of the corrector plate to the focus gives an image with a minimized astigmatism for a refractive index of $n = 1.47$ Bowen & Vaughan (1973). At a given wavelength, the change of length which minimizes astigmatism is described in Bowen's paper as

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

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

The last technology the du Pont Telescope uses is conical baffles. The reason for this is to promote shielding in the telescope [Bowen & Vaughan \(1973\)](#). As explained in the Bowen paper, shielding is necessary in order to protect the photographic plate from light that escapes from the secondary lens due to long time exposure. As explained in the Bowen paper, the conic baffles are "located in the space between the incoming beam as it approaches the primary and the return beam from the secondary to the plate" [Bowen & Vaughan \(1973\)](#). Theoretically, the conic baffles have the disadvantage of producing a diffraction pattern. However, as explained by Bowen, this should not majorly affect the images of stars [Bowen & Vaughan \(1973\)](#).

(3) The New Mexico State University (NMSU) Telescope: The NMSU telescope takes the most technologically advanced approach to collecting data compared to the 2.5 m telescope and the du Pont Telescope. The NMSU telescope uses a camera that has a 2048 x 2048 CCD. The camera is controlled by a linux computer, which is connected by fiber optic cables [Holtzman et al. \(2010\)](#).

The data collection of the NMSU telescope is almost fully automated using C++, the only time it is not is when engineering is being done on the telescope [Holtzman et al. \(2010\)](#). The NMSU telescope has a camera which analyzes the brightness level of the sky to see if it is dark enough to start collecting data. When the sky becomes dark enough, the telescope initiates its program. The telescope goes through the list and observes said objects [Holtzman et al. \(2010\)](#). Objects can, however, be observed as many times as requested.

Data Processing—The SDSS processes its data through an innovative acquisition system that records and organizes observations in real time while maintaining strict quality control [Gunn et al. \(2006\)](#). This system ensures that data are processed and stored efficiently without any loss of precision. [Gunn et al. \(2006\)](#). The data pipeline of the SDSS can be described by two different fields of data, the imaging pipeline and the spectroscopy pipeline

(1) Imaging Data Pipeline: The Imaging data Pipeline itself consists of multiple subpipelines. The first subpipeline is the Astroline. The astroline uses vxWorks in order to initialize the processing sequence. This happens by composing star cutouts and column quartiles collected from the CCD's mentioned before ([Lupton et al. 2001](#))

The second subpipeline is the MT pipeline. the MT Pipeline processes the data collected from the Photometric Telescope. This data is used to calculate important parameters for the 2.5 m Telescope scans, such as extinction and zero-points ([Lupton et al. 2001](#)).

The third pipeline is the Serial Stamp Collecting (SSC) Pipeline. the SSC reorganizes the star cutouts collected from previous pipelines. The SSC does this in order to prepare data for the upcoming pipelines ([Lupton et al. 2001](#)).

Next is the Astrometric Pipeline. The Astrometric pipeline processes the average location of stars using the data collected from the astroline and SSC pipelines. It then converts the pixel data from the images into celestial coordinates (α, δ) ([Lupton et al. 2001](#)).

263 After that is the Postage Stamp Pipeline (PSP). The PSP estimates the quality of the data collected
 264 by calculating factors such as the flat field vectors, bias drift, and sky levels ([Lupton et al. 2001](#)).

265 After all that is done, the data is fed into the Frames Pipeline. The Frames pipeline does a majority
 266 of the work, processing the data from all the previous pipelines and producing the complete datasets
 267 of images. It does this by correcting the frames based on the data before and cataloging the images
 268 ([Lupton et al. 2001](#)).

269 Lastly, the processed data is then ran through the Calibration pipeline. The Calibration pipeline
 270 takes data from the MT and Frames pipeline. The Calibration pipeline converts the counts into more
 271 measurable quantities such as flux ([Lupton et al. 2001](#)).

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

275 (2) Spectroscopy Data Pipeline:

276 *Real-Time Processing*—

277 2.1.2. *The Rubin/Large Synoptic Survey Telescope (LSST)*

278 *Data Collection and Storage*—Despite the recent times in creation. The SDSS collected around 16 TB
 279 of data over a decade in their data release 7 ([Juric et al. 2017](#)). Yet the LSST is expected to collect
 280 20 TB of data per night ([NSF-DOE Vera C. Rubin Observatory 2025](#)).

281 The LSST pipeline is written with around 750000 in Python in order to use relevant libraries such
 282 as SciPy and AstroPy ([NSF-DOE Vera C. Rubin Observatory 2025](#)). The pipeline is also written
 283 with around 220000 lines in C++ in order to ensure efficient performance ([NSF-DOE Vera C. Rubin
 Observatory 2025](#)). In order to combine these two languages, pybind11 allows for the transition from
 285 Python to C++, and ndarray objects are able to be converted from C++ arrays ([NSF-DOE Vera C.
 Observatory 2025](#)).

287 The python enviroment of the LSST pipeline uses a package named `rubin-env`. This package gives
 288 the user all the code needed to run LSST’s data. In order to execute the code, the pipeline consists
 289 of multiple packages that each serve their own purpose. The LSST has defined a class labeled `Task`,
 290 which is used to define algorithms ([NSF-DOE Vera C. Rubin Observatory 2025](#)).

291 One instance of a task is the `PipelineTask`, which serves to organize subtasks. These subtasks each
 292 have their own purpose ([NSF-DOE Vera C. Rubin Observatory 2025](#)). The most important subtask,
 293 labeled `daf_butler`, handles the data storage. This subtask is titled by the LSST as The Data Butler.
 294 The Butler serves as a database to store collected data. It stores objects with data IDs similar to
 295 SQL, with headers that hold useful information ([NSF-DOE Vera C. Rubin Observatory 2025](#)). An
 296 example of this would be a data coordinate labeled `instrument="LSSTCam"`, `exposure=299792458`,
 297 `detector=42`, `band=z`, `day_obs=20251011`.

298 *Data Processing*—This task first Before the LSST analyzes data, it removes noise. An example of
 299 noise would be **LOOK AT 5.1 TO FINISH THIS**.

300 There are multiple tasks which define how the LSST processed data to find objects. These are all
 301 defined in the `meas_algorithms` package ([NSF-DOE Vera C. Rubin Observatory 2025](#)). The task
 302 that first handles processing the catalogued images is the `SourceDetectionTask` ([NSF-DOE Vera](#)

C. Rubin Observatory 2025). This task uses Gaussian smoothing in the point spread function. It then convolves the collected image with the point spread function in order to suppress potential noise (NSF-DOE Vera C. Rubin Observatory 2025).

Another task that the LSST uses to process data is `MaskStreaksTask` (NSF-DOE Vera C. Rubin Observatory 2025). The task serves to mask pixels from streaks from other satellites. It identifies streaks using a Canny Filter and the Kernel-Based Hough Transform **Note: Fix this citation** (NSF-DOE Vera C. Rubin Observatory 2025). This task is combined with the deblending of collected images allows for the LSST to accurately identify objects.

311 *Real-Time Processing—*

312 2.2. *The Radio Big Data Pipeline*

313 2.2.1. *The MeerKAT*

314 *Data Processing*—The MeerKAT data processing pipeline is split into three parts, the calibration
 315 pipeline, the continuum pipeline, and the spectral pipeline (Ratcliffe 2021).

316 **(2) The Continuum Pipeline:** The second pipeline is the continuum pipeline. This pipeline
 317 aims to produce continuum images using the orbit software package (Ratcliffe 2021; Cotton 2008).
 318 The orbit package which interfaces UV data. It does this by using the `MFIImage` task, which does
 319 wide-band, wide-field imaging (Ratcliffe 2021; Cotton & Schwab 2010). Once the data is read, it is
 320 split into scans, which are then averaged and combined into a dataset.

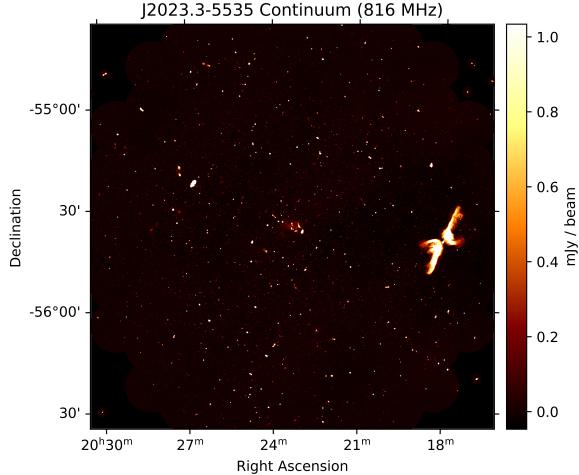


Figure 3. An example UHF continuum thumbnail image produced automatically by the pipeline is shown above, (Ratcliffe 2021)

321 This dataset is then split into 8 107 MHz intermediate frequencies (Ratcliffe 2021). The UV data
 322 then stores headers of `ntimes`, `nbaselines`, `nchannels`, `npolarisation` (Ratcliffe 2021). Given
 323 the Meerkat's design of 48 antennas, short baselines dominate. This causes an inflation in data
 324 volume (Ratcliffe 2021). In order to circumvent this, a baseline dependent averaging step is applied
 325 (Ratcliffe 2021). Time averaging the data with respect to the baseline before merging them into a
 326 dataset. This step reduces the data volume by about 3-4 times with a minimal loss of 1

327 The aforementioned MFImage task uses joint frequency deconvolution to handle wide-band effects
 328 in wide-band images (Ratcliffe 2021). Normally, Meerkat uses around 140 circular facets with a size
 329 of around 6 arcminutes which cover 1 degree from the phase center (Ratcliffe 2021). Additionally,
 330 there are facets with around 1 degree that use the SUMSS or NVSS catalogue for a radius of 2.5
 331 from the center of the phase Ratcliffe (2021). These 1 degree facets are used for phenomena that
 332 are anticipated to have a flux density greater than 5 mJy Ratcliffe (2021). Due to the fact that these
 333 facets deal with wide-band images, they also deal with wide-band imaging effects. These effects
 334 are bypassed by splitting the frequency band used into 10 components (Ratcliffe 2021). When the
 335 joint frequency deconvolution is initiated, the brightest light sources in the dataset are found and
 336 subtracted from the slices of data individually. This is called the CLEAN algorithm (Ratcliffe 2021).

337 Self-calibration is then performed in two rounds in order to transform data. The first round consists
 338 of using the CLEAN algorithm with a spectral density of 1 mJy. This yields approximately 1000
 339 CLEAN components (Ratcliffe 2021). The data is then self calibrated once more, using a CLEAN
 340 algorithm with a spectral density of 100 μ Jy. The second round yields approximately 10000 CLEAN
 341 components (Ratcliffe 2021). After this process, the data is imaged and used to create the continuum
 342 images. Lastly, the self calibrated data is converted into AIPS format, which consists of tables with
 343 headers such as UV data, polarisation, etc (Ratcliffe 2021). The data is then converted into numpy
 344 arrays for each timestamp with shape (nif, npol, nantennas) (Ratcliffe 2021). The sky model, made up
 345 of the CLEAN components, is stored in AIPS CC format. The flux density of all the 10 components
 346 of the frequency are summated (Ratcliffe 2021). The merged flux density is then fitted with a second
 347 degree polynomial to measure frequency and is used to subtract from the image (Ratcliffe 2021).
 348 Finally, the fully processed images are converted into FITS and PNG files and archived (Ratcliffe
 349 2021).

350 **(3) The Spectral Pipeline:** The purpose of the spectral pipeline is to speed up image creation
 351 while keeping control of quality (Ratcliffe 2021). The issue with the imaging data, at this point is
 352 that they're independent of each other, which causes issues with accessing data when trying to run
 353 said images (Ratcliffe 2021). The raw data, also known as visibilities. are received in time-major
 354 order, but this data structure must be transposed in order for the data to be in channel-major order
 355 (Ratcliffe 2021). In order to solve this issue, a visibility writer is used. The writer stores the visibilities
 356 on a Ceph cluster with chunks each across 64 channels (Ratcliffe 2021). This choice in the chunk size
 357 is not optimized, but it does avoid issues with RAM and memory allocation (Ratcliffe 2021).

358 The chunk system only partially implements the solution of transposing the data. The solution
 359 is complete after another series of events are done (Ratcliffe 2021). The visibilities in each chunk
 360 are ordered by channel, W slice, and baseline. This allows for both the accessibility of the data
 361 to improve, and it inherently stores measurements such as UVW coordinateas and parallactic angles
 362 (Ratcliffe 2021). While all that is occurring, conservative baseline-depedent averaging is also applied to
 363 the visibilities (Ratcliffe 2021). The coordinates of every visibility are solved for, then visibilities with
 364 matching coordinates are merged, which is the last of the preprocessing of the visiblities (Ratcliffe
 365 2021).

366 After all of that is done, By using CUDA, a parallel processing software, every channel is imaged
 367 separately. The use of CUDA allows for the processing to speed up due to the usage of NVIDIA
 368 GPUs (Ratcliffe 2021). Lastly, the data follows a cycle as described in the figure below in order to
 369 be updated.

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, ([Ratcliffe 2021](#))

370

2.2.2. *The Square Kilometre Array (SKA)*

371

3. RESULTS

372

3.1.

373

4. DISCUSSION

374

4.1. *Open Source Policies and Transparency*

375

4.1.1. *The SDSS Policy*

376

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

377

378

379

380

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

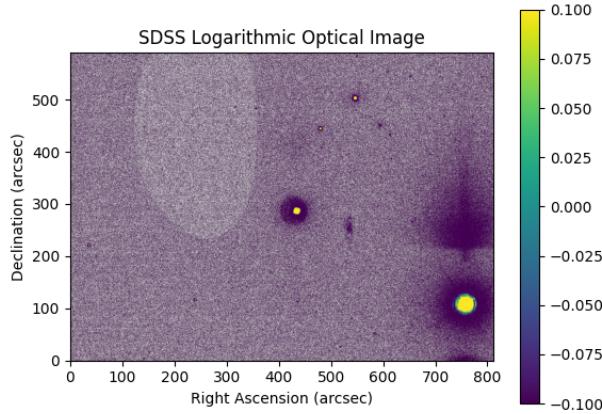


Figure 5. A spectrographical image obtained using data collected from SDSS

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

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

389 4.1.2. *The LSST Policy*

390 Similar to the SDSS, the LSST has declared a commitment to having their software open source
 391 and their pipeline transparent.

392 4.1.3. *The MeerKAT Policy*

393 4.1.4. *The SKA Policy*

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

```

395 1 #Import relevant libraries/functions
396 2 from astroquery.sdss import SDSS
397 3 from astropy import coordinates as coords
398 4 import astropy.units as u
399 5 import matplotlib.pyplot as plt
400 6 import numpy as np
401 7
402 8 #Initialize Right Ascension and Declination
403 9 ra = 20
404 10 dec = -10
405 11
406 12 #Convert ra and dec into a SkyCoord Object

```

```

408 13 coord = coords.SkyCoord(ra, dec, unit='deg', frame = 'icrs')
409 14
410 15 #Query the SDSS System to find object given coordinates in a radius of
411      0.01 degrees
412 16 result = SDSS.query_region(coord, radius=0.01*u.deg, spectro=True)
413 17 print(result)
414 18 #Retrieve the Image from found object, make into a FITS File
415 19 image = SDSS.get_images(matches=result, band=['u', 'g', 'r', 'i', 'z'])
416 20
417 21
418 22 #Retrieve hdulist from FITS file
419 23 hdulist = image[0]
420 24
421 25 #Retrieve the image data from the hdulist
422 26 imageData = hdulist[0].data
423 27
424 28 #Log the image data in order to get rid of background
425 29 imageDataLog = np.log10(imageData) + 1e-8
426 30
427 31 #Save the header
428 32 header = hdulist[0].header
429 33
430 34
431 35 #Obtain the relevant headers
432 36
433 37 #Retrieve pixel scale numbers, divided amongst two parts for ra and dec
434 38 CD1_1 = header['CD1_1']
435 39 CD1_2 = header['CD1_2']
436 40 CD2_1 = header['CD2_1']
437 41 CD2_2 = header['CD2_2']
438 42
439 43 #Width of image in pixels
440 44 keyWordNAXIS1 = header['NAXIS1'] #[pixels]
441 45
442 46 #Height of image in pixels [pixels]
443 47 keyWordNAXIS2 = header['NAXIS2'] #[pixels]
444 48
445 49 #Normalize the pixel scale, then multiply by 3600 to convert units
446 50 CDELT1ArcSec = np.linalg.norm([CD1_1,CD1_2]) * 3600 #[arcsec/pixels]
447 51 CDELT2ArcSec = np.linalg.norm([CD2_1,CD2_2]) * 3600 #[arcsec/pixels]
448 52
449 53 #Set up the image
450 54 plt.xlabel("Right Ascension (arcsec)")
451 55 plt.ylabel("Declination (arcsec)")
452 56 plt.title('SDSS Logarithmic Optical Image')
453 57 vmin2 = np.percentile(imageDataLog, 85)
454 58 vmax2 = np.percentile(imageDataLog, 98)

```

```

455 59 plt.imshow(imageDataLog, cmap='viridis', extent = (0, CDELT1ArcSec *
456      keyWordNAXIS1, 0, CDELT2ArcSec * keyWordNAXIS2), vmin = vmin2, vmax =
457      vmax2)
458 60 plt.colorbar()
459 61
460 62 plt.show()

```

REFERENCES

- 462 ???, The MIGHTEE Survey,
463 <https://www.mighteesurvey.org/home>
- 464 ???, SKA Telescope Specifications,
465 <https://www.skao.int/en/science-users/118/ska-telescope-specifications>
- 466 Blyth, S., Baker, A. J., Holwerda, B., et al. 2018,
467 in Proceedings of MeerKAT Science: On the
468 Pathway to the SKA — PoS(MeerKAT2016)
469 (Stellenbosch, South Africa: Sissa Medialab),
470 004, doi: [10.22323/1.277.0004](https://doi.org/10.22323/1.277.0004)
- 471 Bowen, I. S., & Vaughan, A. H. 1973, Applied
472 Optics, 12, 1430, doi: [10.1364/AO.12.001430](https://doi.org/10.1364/AO.12.001430)
- 473 Bundy, K., Bershadsky, M. A., Law, D. R., et al.
474 2014a, The Astrophysical Journal, 798, 7,
475 doi: [10.1088/0004-637X/798/1/7](https://doi.org/10.1088/0004-637X/798/1/7)
- 476 —. 2014b, The Astrophysical Journal, 798, 7,
477 doi: [10.1088/0004-637X/798/1/7](https://doi.org/10.1088/0004-637X/798/1/7)
- 478 Cotton, W. D. 2008, Publications of the
479 Astronomical Society of the Pacific, 120, 439,
480 doi: [10.1086/586754](https://doi.org/10.1086/586754)
- 481 Cotton, W. D., & Schwab, F. R. 2010
- 482 Dawson, K. S., Kneib, J.-P., Percival, W. J., et al.
483 2016, The Astronomical Journal, 151, 44,
484 doi: [10.3847/0004-6256/151/2/44](https://doi.org/10.3847/0004-6256/151/2/44)
- 485 De Blok, W. J. G., Healy, J., Maccagni, F. M.,
486 et al. 2024, Astronomy & Astrophysics, 688,
487 A109, doi: [10.1051/0004-6361/202348297](https://doi.org/10.1051/0004-6361/202348297)
- 488 Dewdney, P., Hall, P., Schilizzi, R., & Lazio, T.
489 2009, Proceedings of the IEEE, 97, 1482,
490 doi: [10.1109/JPROC.2009.2021005](https://doi.org/10.1109/JPROC.2009.2021005)
- 491 Goedhart, S. 2025, MeerKAT Specifications
- 492 Gunn, J. E., Siegmund, W. A., Mannery, E. J.,
493 et al. 2006, The Astronomical Journal, 131,
494 2332, doi: [10.1086/500975](https://doi.org/10.1086/500975)
- 495 Gupta, N., Jagannathan, P., Srianand, R., et al.
496 2021, The Astrophysical Journal, 907, 11,
497 doi: [10.3847/1538-4357/abcb85](https://doi.org/10.3847/1538-4357/abcb85)
- 498 Hey, T., Tansley, S., & Tolle, K. 2009, Microsoft
499 Research
- 500 Holtzman, J. A., Harrison, T. E., & Coughlin,
501 J. L. 2010, Advances in Astronomy, 2010,
502 193086, doi: [10.1155/2010/193086](https://doi.org/10.1155/2010/193086)
- 503 Ivezic, Ž., Kahn, S. M., Tyson, J. A., et al. 2019,
504 The Astrophysical Journal, 873, 111,
505 doi: [10.3847/1538-4357/ab042c](https://doi.org/10.3847/1538-4357/ab042c)
- 506 Jonas, J., & the MeerKAT Team. 2018, in
507 Proceedings of MeerKAT Science: On the
508 Pathway to the SKA — PoS(MeerKAT2016)
509 (Stellenbosch, South Africa: Sissa Medialab),
510 001, doi: [10.22323/1.277.0001](https://doi.org/10.22323/1.277.0001)
- 511 Juric, M., Kantor, J., Lim, K.-T., et al. 2017
- 512 Lesser, M. 2015, Publications of the Astronomical
513 Society of the Pacific, 127, 1097,
514 doi: [10.1086/684054](https://doi.org/10.1086/684054)
- 515 Lupton, R., Gunn, J. E., Ivezić, Z., et al. 2001,
516 The SDSS Imaging Pipelines, arXiv,
517 doi: [10.48550/arXiv.astro-ph/0101420](https://arxiv.org/abs/astro-ph/0101420)
- 518 Lupton, R. H., Ivezić, Z., Gunn, J., et al. 2007
- 519 Majewski, S. R., Schiavon, R. P., Frinchaboy,
520 P. M., et al. 2017, The Astronomical Journal,
521 154, 94, doi: [10.3847/1538-3881/aa784d](https://doi.org/10.3847/1538-3881/aa784d)
- 522 Moore, G. E. 2006, IEEE Solid-State Circuits
523 Society Newsletter, 11, 33,
524 doi: [10.1109/N-SSC.2006.4785860](https://doi.org/10.1109/N-SSC.2006.4785860)
- 525 NSF-DOE Vera C. Rubin Observatory. 2025,
526 PSTN-019: The LSST Science Pipelines
527 Software: Optical Survey Pipeline Reduction
528 and Analysis Environment, NSF-DOE Vera C.
529 Rubin Observatory,
530 doi: [10.71929/RUBIN/2570545](https://doi.org/10.71929/RUBIN/2570545)
- 531 Ratcliffe, S. 2021, SDP Pipelines Overview,
532 <https://skaafrica.atlassian.net/wiki/spaces/ESDKB/pages/>
- 533 Richards, S. 2020, What Is Vignetting?
534 Verbunt, F., & Van Gent, R. H. 2010, Astronomy
535 and Astrophysics, 516, A28,
536 doi: [10.1051/0004-6361/201014002](https://doi.org/10.1051/0004-6361/201014002)

- 538 Woudt, P. A., Fender, R., Corbel, S., et al. 2018,
539 in Proceedings of MeerKAT Science: On the
540 Pathway to the SKA — PoS(MeerKAT2016)
541 (Stellenbosch, South Africa: Sissa Medialab),
542 013, doi: [10.22323/1.277.0013](https://doi.org/10.22323/1.277.0013)