

# 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  
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

## Contents

12	1. Introduction	2
13	1.1. The Paradigms of Data Science	2
14	1.2. The Rise of Big Data in Astronomy	3
15	1.3. An Overview of The Four Surveys	3
16	2. Methods	6
17	2.1. The Optical Big Data Pipeline	6
18	2.1.1. The Sloan Digital Sky Survey (SDSS)	6
19	2.1.2. The Rubin/Large Synoptic Survey Telescope (LSST)	9
20	2.2. The Radio Big Data Pipeline	12
21	2.2.1. The MeerKAT	12
22	2.2.2. The Square Kilometre Array (SKA)	14
23	3. Results	14
24	3.1. The Optical Data Results	14
25	3.1.1. The Sloan Digital Sky Survey (SDSS)	14
26	3.1.2. The Rubin/Large Synoptic Survey Telescope (LSST)	16
27	3.2. The Radio Data Results	16
28	3.2.1. The MeerKAT	16
29	3.2.2. The Square Kilometre Array (SKA)	16
30	4. Discussion	16
31	4.1. Open Source Policies and Transparency	16

32	4.1.1. The SDSS Policy	16
33	4.1.2. The LSST Policy	17
34	4.1.3. The MeerKAT Policy	17
35	4.1.4. The SKA Policy	17
36	4.2. The Rise of AI/ML in Surveys	17
37	A. Python Code for SDSS Data Retrieval (Figure 3)	17

## A. Python Code for SDSS Data Retrieval (Figure 3)

## 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 ([Hey et al. 2009](#)), 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 (Hey et al. 2009), 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 (Verbunt & Van Gent 2010). However, empirical evidence can be compromised by human error, the precision of the instruments, and, most importantly, the relatively slow pace of data acquisition compared to subsequent paradigms.

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

The third paradigm is simulation evidence (Hey et al. 2009), a relatively recent development. Simulation models natural phenomena that are too complex to model analytically or compute by hand. It allows interpolation and extrapolation of data using computational techniques grounded in

known physical laws. For example, in astronomy, N-body simulations are used to study the complex 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 <sup>1</sup>, The Sloan Digital Sky Survey (SDSS) <sup>2</sup>, The Legacy Survey of Space and Time (LSST) <sup>3</sup>, and The Square Kilometre Array (SKA) <sup>4</sup>. 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 <sup>5</sup>. This telescope is stationed in Chile, where it observes the southern hemisphere instead. Similar to the foundational

<sup>1</sup> <https://www.skao.int/en>

<sup>2</sup> <https://www.sdss.org/>

<sup>3</sup> <https://www.lsst.org/>

<sup>4</sup> <https://www.skao.int/en>

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

108 telescope at Apache Point, this telescope has a  $2.1^\circ$  field of view but only uses one corrector lens  
 109 ([Bowen & Vaughan 1973](#)).

110 The third telescope is the NMSU 1-meter Telescope <sup>6</sup>. The NMSU telescope is stationed at the  
 111 Apache Point Observatory alongside the foundational telescope. The NMSU telescope is designed to  
 112 observe bright stars that are too bright for the aforementioned two telescopes to observe ([Majewski  
 et al. 2017](#)).

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

121 The MeerKAT <sup>10</sup> is an important precursor telescope to the SKA ([Jonas & the MeerKAT Team  
 122 2018](#)) MeerKAT became fully operational in 2018 in the Northern Cape Province of South Africa.  
 123 MeerKAT comprises 64 antennas distributed over a radius of approximately 600 miles ([Goedhart  
 124 2025](#)). These antennas operate across frequency bands ranging from 350 MHz to 3500 MHz ([Goedhart  
 125 2025](#)).

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

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

139 Another survey, ThunderKAT <sup>13</sup>, aims to find, identify and understand high-energy radio trans-  
 140 sentists, usually grouped with observations at similar wavelengths. Examples include supernovae,  
 141 microquasars, and similar events ([Woudt et al. 2018](#)).

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

<sup>6</sup> <https://newapo.apo.nmsu.edu/>

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

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

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

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

<sup>11</sup> <https://science.uct.ac.za/laduma>

<sup>12</sup> <https://mals.iucaa.in/>

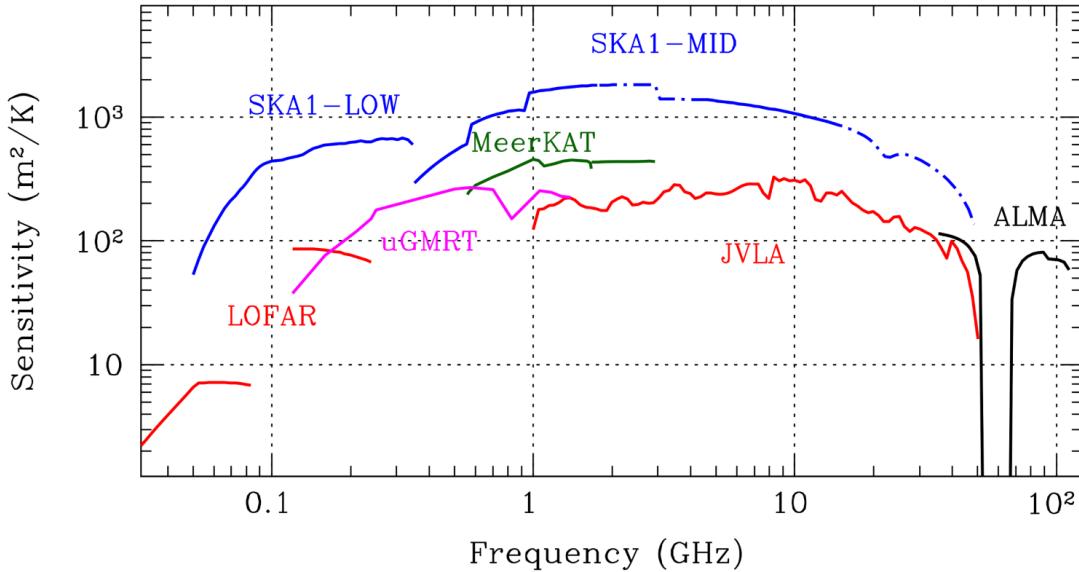
<sup>13</sup> <https://www.physics.ox.ac.uk/research/group/meerkat>

<sup>14</sup> <https://mhongoose.astron.nl/>

<sup>15</sup> <https://www.astron.nl/halogas/>

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 MIGHTEE<sup>17</sup>. 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 3 from the SKA Official website<sup>b</sup>, SKA1 sensitivity compared to existing facilities at similar frequencies

- <sup>a</sup> <https://www.skao.int/en/science-users/118/ska-telescope-specifications>
- <sup>b</sup> <https://www.skao.int/en/science-users/118/ska-telescope-specifications>

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  $9.6 \text{ deg}^2$  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 unprecedent 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

<sup>16</sup> <https://www2.mpia-hd.mpg.de/THINGS/Overview.html>

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

165 MeerKAT to their larger successor telescopes, Rubin/LSST and SKA and consider the evolution of  
 166 data challenges.

167 **2. METHODS**

168 **2.1. *The Optical Big Data Pipeline***

169 This section carefully examines how each of the optical focused surveys collects and processes its  
 170 data. By describing the nature, scope, and type of the data. Next, we discuss how each survey  
 171 collects and archives its data, followed by an explanation of their general data processing methods.  
 172 Finally, we consider the use of real-time data processing.

173 As noted Above, we can consider the SDSS to consist of three main components: the 2.5 m  
 174 foundational telescope, the Irénée du Pont telescope, and the NMSU 1-meter telescope. The SDSS  
 175 has evolved rapidly overtime. In the early 2000's the SDSS was primarily focused on optical imaging,  
 176 but ever since the mid 2010's it has shifted towards gathering spectral data (??) As discussed  
 177 previously, we will consider how each telescope collects data, how the SDSS processes the data, and  
 178 makes the data accessible through an open data policy.

179 **2.1.1. *The Sloan Digital Sky Survey (SDSS)***

180 *Data Collection and Storage*—**(1) The SDSS 2.5 m Foundational Telescope:** The SDSS camera  
 181 contains 54 2048 x 2048 charge-coupled devices (CCDs) and 24 2048 x 400 CCDs. A CCD is an  
 182 imaging detector that converts incoming light into an electronic signal. When photons strike the  
 183 CCD, they generate electrons through the internal photoelectric effect. The accumulated charge is  
 184 measured per pixel and is stored as a digital value (Lesser 2015).

185 In addition to the CCD imaging data the SDSS collects spectra using a pair of fiber-fed double  
 186 spectrographs, over a wavelength range from 3800 to 9200 Å and at field angles between 0 and 90°  
 187 (Gunn et al. 2006). The spectrograph fibers are positioned in pre-selected objects of the field. The  
 188 optical performance of these spectrographs, which are summarized in Table 5 from Gunn et al. (2006).

189 In Figure 5,  $\lambda$  represents wavelength, and “Angle” refers to the field angle.  $f_b$  denotes the best-  
 190 focus distance.  $h/dh$  represents the lateral color,  $D$  denotes the longitudinal difference from the best  
 191 focus, and  $\epsilon$  is the root mean square (rms) image diameter. Smaller values of the lateral color and  
 192 longitudinal difference indicate sharper images. Both of these quantities remain close to zero for  
 193 most wavelengths and field angles, except between roughly 5300 and 6500 Å(Gunn et al. 2006). This  
 194 demonstrates the high optical accuracy of the SDSS spectrographs. The 2.5 m telescope collects  
 195 imaging and spectroscopic data at a rate of about 20Gb/hr (Lupton et al. 2007).

196 **(2) The Irénée du Pont Telescope:** The du Pont telescope data collection has evolved overtime  
 197 (Bowen & Vaughan 1973). It used to collect imaging data, which I will talk about first. In recent  
 198 times, such as the 16th SDSS data release, the Du Pont telescope is used for collecting spectra data  
 199 (?), which is what I'll talk about later.

200 During its optical era, the telescope is a modified Ritchey-Chrétien design with Gascoigne correcton  
 201 and a 100-inch primary mirror (Bowen & Vaughan 1973). Approximately 40% of the light is reflected  
 202 to the secondary mirror, resulting in only a 16% loss of light at that stage. (Bowen & Vaughan 1973)

203 The du Pont Telescope used 18.9 inch nonvignetted plates in order to minimize vignetting (Bowen  
 204 & Vaughan 1973). Vignetting is the process where light beds through the lense of a telescope. The

$\lambda$ (Å)	Angle (arcmin)	$f_b$ (mm)	$h/dh$ (mm)	$D$ (mm)	$\epsilon$ (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
4600.....	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
6500.....	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), Telescope Optical Performance for the Spectrographic Mode: Average Focus

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 applied was a Gascoigne corrector plate. The plate helped with data collection. The Gasciogne corrector plate was abled to be moved, which could 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 gave 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 minimized astigmatism is described in Bowen's paper as

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

Where  $\Delta L$  is the change in seperation in millimeters and  $\Delta n$  is the difference between a refractive index of 1.47 and the index wanted.

The last technology the du Pont Telescope used was conical baffles. The reason for this was to promote shielding in the telescope (Bowen & Vaughan 1973). As explained in the Bowen paper, shielding was necessary in order to protect the photographic plate from light that escaped from the secondary lense due to long time exposure. the conic baffles were located in the space between the primary and secondary lenses in the plate (Bowen & Vaughan 1973). Theoretically, the conic baffles

had the disadvantage of producing a diffraction pattern. However, as explained by Bowen, this did not majorly affect the images of stars (Bowen & Vaughan 1973).

Towards the mid 2010's the du Pont telescope has shifted focus towards spectra data (?). The du Pont telescope alongside the Foundational telescope used spectrographs during the APOGEE-2 survey and sampled approximately 400,000 stars (?). The du Pont telescope uses a fiber-optic system consisting of 300 short fibers that are used throughout the night (?). These fibers can collect observations up to 10 plates per night which are stored on five cartridges (?). The camera uses a 1024 x 1024 pixel ccd, with the most effective wavelength for the camera being approximately 7600 Å (?)

**(3) The New Mexico State University (NMSU) 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++ (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. The NMSU telescope was used to observe stars too bright to observe with the larger-aperture telescopes. It is no longer used in SDSS-V (?).

*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). The data pipeline of the SDSS can be further sub-divided as 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. This subpipeline uses vxWorks to initialize the processing sequence 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. This pipeline processes the data collected from the Photometric telescope. These data are 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 (Lupton et al. 2001). The SSC reorganizes the star cutouts collected from previous pipelines in order to prepare data for the subsequent processing (Lupton et al. 2001).

The Astrometric pipeline follows and estimates the average position of stars using data collected from the Astroline and SSC pipelines. It then converts the pixel coordinates to celestial coordinates ( $\alpha, \delta$ ) (Lupton et al. 2001).

The next stage is the Postage Stamp Pipeline (PSP). The PSP estimates data quality by calculating factors such as the flat field vectors, bias drift, and sky levels (Lupton et al. 2001).

The data is fed into the frames Pipeline. The frames pipeline does a majority of the work, processing the data from all the previous pipelines and producing the complete image datasets and cataloging the images (Lupton et al. 2001).

The calibration pipeline takes data from the MT and Frames pipeline and converts the counts into calibrated flux densities (Lupton et al. 2001).

**(2) Spectroscopy Data Pipeline:**

263 *Real-Time Processing—*

264       2.1.2. *The Rubin/Large Synoptic Survey Telescope (LSST)*

265 *Data Collection and Storage*—The SDSS collected around 16 TB of data over a decade in their data  
 266 release 7 (Juric et al. 2017). Yet the LSST is expected to collect 20 TB of data per night (NSF-DOE  
 267 Vera C. Rubin Observatory 2025).

268 The LSST pipeline consists of approximately 750000 in Python and uses relevant libraries such as  
 269 SciPy<sup>18</sup> and AstroPy<sup>19</sup> (NSF-DOE Vera C. Rubin Observatory 2025). The pipeline also contains  
 270 approximately 220000 lines of C++ to ensure efficient performance (NSF-DOE Vera C. Rubin Ob-  
 271 servatory 2025). The tool pybind11<sup>20</sup> is needed to parse from Python to C++, and ndarray objects  
 272 are able to be converted from C++ arrays (NSF-DOE Vera C. Rubin Observatory 2025).

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

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

284 *Data Processing*—There are multiple tasks which define how the LSST processed data to find objects.  
 285 These are all defined in the `meas_algorithms` package (NSF-DOE Vera C. Rubin Observatory 2025).  
 286 The task that first handles processing the catalogued images is the `SourceDetectionTask` (NSF-  
 287 DOE Vera C. Rubin Observatory 2025). This task uses Gaussian smoothing in the point spread  
 288 function. It then convolves the collected image with the point spread function in order to suppress  
 289 potential noise (NSF-DOE Vera C. Rubin Observatory 2025).

290 Another task that the LSST uses to process data is `MaskStreaksTask` (NSF-DOE Vera C. Rubin  
 291 Observatory 2025). The task serves to mask pixels from streaks from other satellites. It identifies  
 292 streaks using a Canny Filter and the Kernel-Based Hough Transform (NSF-DOE Vera C. Rubin  
 293 Observatory 2025; Fernandes & Oliveira 2008). This task is combined with the deblending of collected  
 294 images allows for the LSST to accurately identify objects.

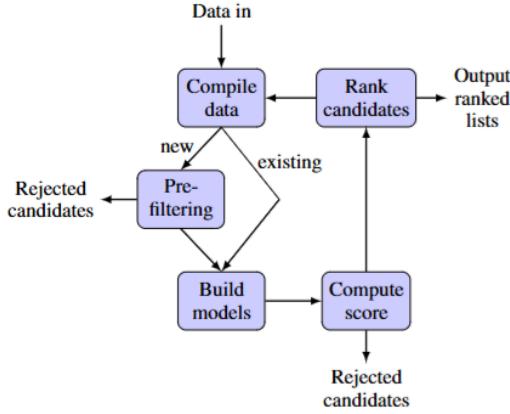
295       **New text starts here**

296 *Real-Time Processing—Should I leave this here or put it in discussions?* Due to the sheer  
 297 volume of data processed by the LSST, processing data in real time is important. Multiple programs  
 298 have been made to circumvent this issue, one I will highlight is AAS2TRO. The goal of the AAS2RT0  
 299 program is to use the Danish-1.54 m telescope to quality control LSST data (?).

<sup>18</sup> <https://scipy.org/>

<sup>19</sup> <https://www.astropy.org/>

<sup>20</sup> <https://pybind11.readthedocs.io/en/stable/index.html>



**Figure 3.** Figure 2 from Sedgewick et al. (2025), A graph showing the AAS2RTO process

The first step of AAS2TRO is compiling an unordered set of candidates based on user-given parameters, and computing an interest value for each candidate (?). The AAS2RTO receives data about new observations from alert brokers such as FINK <sup>21</sup> and Lasair <sup>22</sup> (?).

The LSST sends raw data to these alert brokers in packets. These packets consist of information about a detected astrophysical event which underwent a change of brightness or position (?). The brokers then filter the data for false positive alerts and send the filtered data to AAS2RTO (?).

The next step of AAS2TRO is to prefilter the data (?). This is done using a rough scoring function in order to prefilter bad data effectively before using more costly models (?). Then, AAS2RTO uses fit models in order to further identify good observations. After that, AAS2RTO computes the scores of observations in order to rank observations in descending order and removes negative observations (?).

An application of AAS2TRO, as described in (?), is the use of classifying type Ia supernovae (SNe Ia). The prefilter step uses a specialized model. The model works by combining four functions, which I will describe, into one and determining whether the score is positive (good data), or negative (bad data) (?). The following equations have been generated with ChatGPT and have been checked for errors compared to the original equations in ?. The first equation is the magnitude of the observation:

$$x_{\text{mag}} = 10^{0.5 \times (18.5 - m)} \quad (2)$$

Where where  $m$  is the magnitude as detected in the data. This equation promotes brighter supernovae over fainter ones. Supernovae fainter than 18.5 magnitude are then removed from the scoring list, but not completely deleted, as they may become brighter overtime (?).

The next equation is the peak brightness:

$$x_{\text{peak}} = A \times \exp \left[ -\frac{(t_{\text{obs}} - t_0)^2}{2\sigma^2} \right] \quad (3)$$

<sup>21</sup> <https://fink-broker.org/about/>

<sup>22</sup> <https://lasair-ztf.lsst.ac.uk/>

322 Described by a gaussian distribution, this equation serves to emphasize objects that are near their  
 323 predicted peak using light-curve fit model such as the Spectral Adaptive Lightcurve Template (SALT)  
 324 model (?).  $t_{\text{obs}}$  describes the observation time of the supernovae (?).  $A$  is the amplitude parameter,  
 325 set to  $A = 30$ .  $\sigma$  determines how far away the observation time can be from the predicted peak  
 326 time, this is usually set as  $\sigma = 1\text{day}$  (?). The equation quantifies the importance of how close to its  
 327 predicted peak the supernovae is by assigning  $x_{\text{peak}}$  to its expected value, but also other values for  
 328 different cases (?). If  $|t_{\text{obs}} - t_0| > 4$ , then  $x_{\text{peak}}$  is set to  $10^{-2}$  (?). Lastly, if the SALT model fails,  
 329  $x_{\text{peak}}$  is set to 1.0 (?).

330 The third function goes as follows:

$$331 \quad x_{\text{rise}}^k = \frac{\sum_{i=1}^{N_k-1} [m_{i+1}^k < m_i^k]}{N_k - 1} \quad (4)$$

332 Where  $k$  is the photometric band.  $m_i^k$  is the magnitude of the  $i$ th detection in the  $k$  band.  $N^k$  is  
 333 the number of detections in the  $k$  band (?). The numerator of the equation is a boolean expression.  
 334 If the statement is true, it equals 1, and 0 otherwise (?). This equation is used to quantify if the  
 335 supernovae is becoming brighter or dimmer based on its current and earlier observations (?).

336 The last equation is described as:

$$337 \quad x_{\text{span}} = \begin{cases} 1, & T < 20 \text{ days} \\ L(T; r, x_m), & \text{otherwise} \end{cases} \quad (5)$$

338 Where  $L$  is defined as:

$$339 \quad L(x; r, x_m) = \frac{1}{1 + \exp(-r(x - x_m))} \quad (6)$$

340  $x_m$  is the day in which the brightness peaks. Similar to  $x_{\text{peak}}$ ,  $x_{\text{span}}$  describes how far along a  
 341 supernovae is only with respect to the first time it is observed, and serves to quantify the goodness  
 342 of the supernovae observation when the SALT model fails (?). Observations with time greater than  
 343 30 days are discarded (?)

344 These four functions combine into one function, the SNe Ia score:

$$345 \quad S_{\text{Ia}} = S_{\text{base}} x_{\text{mag}} x_{\text{peak}} x_{\text{rise}} x_{\text{span}} \quad (7)$$

346 Where  $S_{\text{base}} = 1$  (?). Next, the final score function is calculated from  $S_{\text{Ia}}$  and the visibility function  
 347 (?). The first function is the visibility function described as follows:

$$348 \quad x_{\text{vis}} = \left( \frac{A_{\text{vis}}}{(a_{\text{ref}} - a_{\min})(t_{\text{SR}} - t_{\text{obs}})} \right)^{-1} \quad (8)$$

349  $A_{\text{vis}}$  is equal to:

$$350 \quad A_{\text{vis}} = \int_{t_{\text{obs}}}^{t_{\text{SR}}} [a(t) - a_{\min}] dt \quad (9)$$

$a_{ref} = 90^\circ$  is the reference altitude used for normalization.  $a(t)$  is the altitude at a given time.  $a_{min}$  is the minimum altitude that the observation can be detected at.  $t_{SR}$  is the expected time of sunrise.  $t_{obs}$  is the time of the observation (?).

Finally, the final score,  $S_{DK154}$ , is calculated using the following equation:

$$S_{DK154} = S_{\text{Ia}} x_{vis} \quad (10)$$

After all that, the AAS2RTO ranks the observations based on descending order and removes negative scores.

The AAS2RTO is only one of many programs used to deal with real-time processing for the LSST. The rise of prevalence with these programs demonstrate the challenges of big data in astronomy, which will be discussed further in the discussion section.

## 2.2. *The Radio Big Data Pipeline*

### 2.2.1. *The MeerKAT*

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

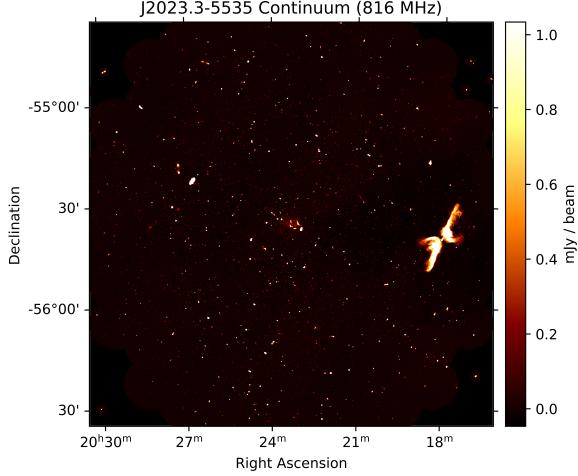
**(1) The Calibration Pipeline:** The pipeline starts by choosing an antenna as a reference (Ratcliffe 2021). The first scans for a given antenna are then averaged and the fourier transform is applied. The peak-to-noise ratio is the ratio of the data maximum to the peak rms noise (Ratcliffe 2021). The antenna with the highest median peak-to-noise ratio over every baseline phase is chosen as the reference antenna (Ratcliffe 2021). This antenna will have all of its phase calibration solutions set to zero (Ratcliffe 2021).

The purpose of the calibration pipeline is to ensure that instrumental antenna errors are calibrated. At the start of every observation, the reference antenna is evaluated based on certain flags, such as data loss (Ratcliffe 2021). If 80 percent or more of the data is flagged, a new reference antenna is determined using the aforementioned process (Ratcliffe 2021).

**(2) The Continuum Pipeline:** The second pipeline is the continuum pipeline. This pipeline produces continuum images using the OBIT software package (Ratcliffe 2021; Cotton 2008). The OBIT package reads the UV data. The OBIT `MFIImage` task is used to perform wide-band, wide-field imaging (Ratcliffe 2021; Cotton & Schwab 2010). Once the data are read, they are split into scans, which are then averaged and combined into a dataset.

This dataset is then split into eight 107 MHz intermediate frequencies (Ratcliffe 2021). The UV data's headers contain the number of antennas, the number of baselines, the number of channels, and the number of polarization (Ratcliffe 2021). In the 48-antenna Meerkat's array, short baselines dominate (Ratcliffe 2021). This allows baseline-dependent averaging to reduce data volume (Ratcliffe 2021). This step reduces the data volume by about 3-4 times (Ratcliffe 2021).

The `MFIImage` task uses joint frequency deconvolution to handle wide-band effects in wide-band images (Ratcliffe 2021). Normally, Meerkat uses approximately 140 circular image facets with a size of 6 arcminutes to cover 1 degree from the phase center (Ratcliffe 2021). Additionally, there are facets of 1 degree that use the SUMSS or NVSS catalogue to cover a radius of 2.5 degrees from the phase center (Ratcliffe 2021). These 1-degree facets are used for phenomena that are anticipated to



**Figure 4.** An example UHF continuum thumbnail image produced automatically by the continuum imaging pipeline is shown above, (Ratcliffe 2021)

have a flux density greater than 5 mJy (Ratcliffe 2021). The facets integrate wide-band imaging effects. The frequency band is split into 10 components (Ratcliffe 2021). During joint-frequency deconvolution, the brightest sources in the dataset are found and subtracted from the slices of data individually using the CLEAN algorithm (Ratcliffe 2021).

Self-calibration is then performed in two rounds. The first round uses CLEAN components with a flux density above 1 mJy. This yields approximately 1000 CLEAN components (Ratcliffe 2021). The data are then self calibrated in a second round, using a CLEAN threshold of  $100 \mu\text{Jy}$ . The second round yields approximately 10000 CLEAN components (Ratcliffe 2021). After this process, the field continuum images are produced. Lastly, the self calibrated data is converted into AIPS format (Ratcliffe 2021). The sky model, made up of the CLEAN components, is stored in AIPS CC format. The flux density of all the 10 frequency components are summed (Ratcliffe 2021). The merged flux density is then fitted with a second degree polynomial over frequency and subtracted from the image (Ratcliffe 2021). Finally, the fully processed images are converted into FITS and PNG files and archived (Ratcliffe 2021).

(3) **The Spectral Pipeline:** The purpose of the spectral imaging pipeline is to produce high-quality spectral line images effectively (Ratcliffe 2021). Spectral channels are independent of each other, and can be processed in parallel (Ratcliffe 2021). However, the raw resulting data are received in time-major order, and this data structure must be transposed in channel-major order (Ratcliffe 2021). In order to solve this issue, a visibility writer is used. The writer stores the visibilities on a Ceph cluster<sup>23</sup> with chunks each across 64 spectral channels (Ratcliffe 2021). This choice in the chunk size is not optimized, but it does avoid issues with RAM and memory allocation (Ratcliffe 2021).

<sup>23</sup> <https://ceph.io/en/>

412 Using CUDA <sup>24</sup>, every channel then is imaged separately. The use of CUDA allows for the processing  
 413 to speed up due to the usage of NVIDIA GPUs (Ratcliffe 2021). The iteration over W slice in each  
 414 chanel is shown in psuedo-code in Fig 4 (Fig. 1 of (Ratcliffe 2021))

1. For each W slice<sup>1</sup>
  - 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.

415 **Figure 5.** The inner cyclic algorithm used by the MeerKAT spectral line pipeline, (Ratcliffe 2021)

416  
 417 The chunk system only partially solves the problem of transposing the data and further steps are  
 418 necessary (Ratcliffe 2021). The visibilities in each chunk are ordered by channel, w slice (wide-field  
 419 imaging plane), and baseline. The data chunks inherently store measurements such as UVW coordi-  
 420 nates and parallactic angles (Ratcliffe 2021). At this point, conservative baseline-depedent averaging  
 421 is also applied to the visibilities (Ratcliffe 2021). The coordinates of every visibility are solved for,  
 422 then visibilities with matching coordinates are merged, which is the last of the preprocessing of the  
 423 visiblities (Ratcliffe 2021).

#### 424 2.2.2. *The Square Kilometre Array (SKA)*

425 New text starts here, I dont know why my figure has 2 footnotes, I have compiled  
 426 multiple times and I have only put one footnote:

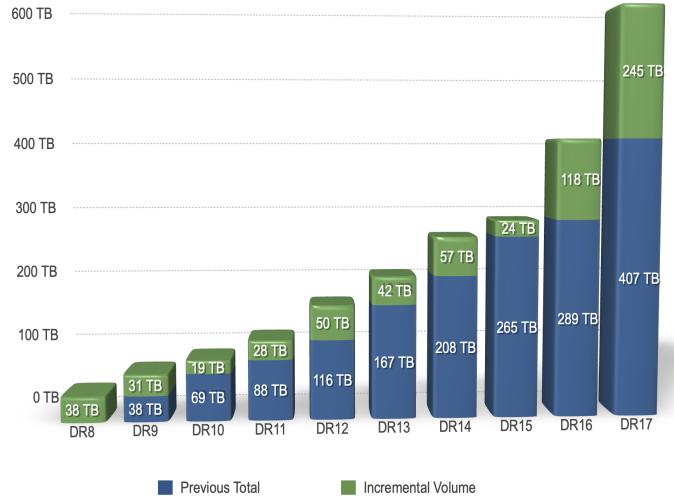
### 427 3. RESULTS

#### 428 3.1. *The Optical Data Results*

##### 429 3.1.1. *The Sloan Digital Sky Survey (SDSS)*

430 Since 1998, the SDSS has published nineteen data releases and five generations, each with their own  
 431 goals. Most recently, SDSS-IV planned to use spectroscopic surveys to detect cosmological objects  
 432 (?). The sixteenth data release is the first one that utilizes the SDSS-IV (?). While the seventeenth  
 433 data release is the last one that utilizes the SDSS-IV. I plan to compare these two data releases.  
 434 The reason I'm not comparing the eighteenth data release is because it was mainly used to set the  
 435 foundation for future SDSS-V data releases by introducing new models, functions, strategies, and  
 436 more (?). The nineteenth release is also not applicable to this comparison because it is only a preview

<sup>24</sup> <https://developer.nvidia.com/cuda-zone>



**Figure 6.** A graph from the SDSS-IV<sup>b</sup> website, showing the increase of data volume throughout multiple data releases

- a [https://www.sdss4.org/dr17/data\\_access/volume/](https://www.sdss4.org/dr17/data_access/volume/)
- b [https://www.sdss4.org/dr17/data\\_access/volume/](https://www.sdss4.org/dr17/data_access/volume/)

of what is achievable by the SDSS-V (?) I believe it is important to compare these two data releases because they quantify the increase of data acquisition within the same generation.

According to Figure 6, it is clear that every data release has grown in terms of volume data. The first few generations increased only slightly, and the later generations increased exponentially. This is a clear example of Moore's Law. The sixteenth data release accounts for approximately 18.1% of the SDSS-IV's total data, and the seventeenth data release accounts for approximately 37.7% of the SDSS-IV's total data. The seventeenth data release consists of over 46 million new files, which is the reason for the massive increase in data volume (?). This increase in data can be summarized in four reasons (?).

The first reason is because the seventeenth data release accounts for the entirety of the APOGEE-2 survey, which had an additional 879,437 infrared spectral measurements (?). The second reason is because the MaNGA survey also completed (?). The MaNGA survey released had 11273 cubes compared to the sixteenth data release's 4824 cubes (?). These cubes hold 3D data, which two spacial dimensions and one wavelength dimension (?). The third reason is because of the eBOSS survey (?). Although the number of data is the same between the sixteenth and seventeenth data release, the seventeenth data release contains 25 value-added catalogues (VAC) that were either updates, or added (?). These VACs contain numerous amounts of processed data that astronomers then use, so each VAC accounts for a substantial portion of the seventeenth data release's total data (?). The last reason is the seventeenth data release includes all the previous SDSS data releases' data (?). The data has been reprocessed with, what was at the time, the newest pipelines dedicated towards data processing (?). The updated pipelines accounts for a substantial amount of new data, as more data can be derived from more effective processes (?)

459      The comparison of only two consecutive generations has shown a massive increase in data volume.  
 460      This trend is likely to explode with the release of SDSS-V and its new components, which we will  
 461      discuss later in the discussion section.

462            3.1.2. *The Rubin/Large Synoptic Survey Telescope (LSST)*

463            3.2. *The Radio Data Results*

464            3.2.1. *The MeerKAT*

465            3.2.2. *The Square Kilometre Array (SKA)*

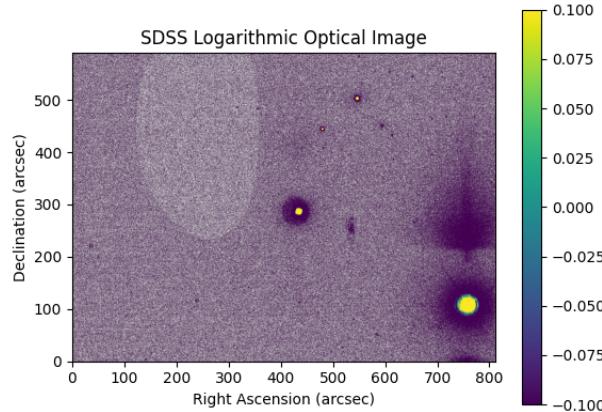
466      New text ends here:

467            4. DISCUSSION

468            4.1. *Open Source Policies and Transparency*

469            4.1.1. *The SDSS Policy*

470      The SDSS-IV collaboration states that all of its software is open source under the BSD-3 license.  
 471      However, the SDSS outlines practices for users who wish to reuse or extend the SDSS software. Most  
 472      importantly, proper citation of software and websites is required.



**Figure 7.** A spectrographical image obtained using data collected from SDSS

473      The SDSS has also implemented digital object identifiers (DOI) in all software code. These DOIs  
 474      allow software and data to be easily identified, which is important for ownership. The SDSS team  
 475      also promotes transparency in coding by Git and SVN <sup>25</sup> to version code repositories.

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

<sup>25</sup> <https://subversion.apache.org/>

479

#### 4.1.2. The LSST Policy

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

#### 4.1.3. The MeerKAT Policy

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

The South African Radio Astronomy Observatory (SARAO) has expressed that MeerKAT data products must be open source with adequate acknowledgement through two steps (Camilo 2024). First and foremost, all MeerKAT data used must come from the MeerKAT ADS Library<sup>26</sup> (Camilo 2024). Secondly, any publication that contains MeerKAT data, the author must state the following statement: "The MeerKAT telescope is operated by the South African Radio Astronomy Observatory, which is a facility of the National Research Foundation, an agency of the Department of Science and Innovation." (Camilo 2024).

Data products can be released into the MeerKAT archive by contacting the SARAO. (Camilo 2024). The SARAO has developed a help desk <sup>27</sup> from which users can ask general questions (Camilo 2024). If users have questions on policies of the MeerKAT, they are advised to communicate them to the SARAO chief scientist (Camilo 2024). There is also a SARAO Users Committee, reporting to the SARAO MD (Camilo 2024).

#### 4.1.4. The SKA Policy

### 4.2. The Rise of AI/ML in Surveys

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

```

501 1 #Import relevant libraries/functions
502 2 from astroquery.sdss import SDSS
503 3 from astropy import coordinates as coords
504 4 import astropy.units as u
505 5 import matplotlib.pyplot as plt
506 6 import numpy as np
507 7
508 8 #Initialize Right Ascension and Declination
509 9 ra = 20
510 10 dec = -10
511 11
512 12 #Convert ra and dec into a SkyCoord Object
513 13 coord = coords.SkyCoord(ra, dec, unit='deg', frame = 'icrs')
514 14
515 15 #Query the SDSS System to find object given coordinates in a radius of
516 16      0.01 degrees
517 result = SDSS.query_region(coord, radius=0.01*u.deg, spectro=True)

```

<sup>26</sup> <https://ui.adsabs.harvard.edu/public-libraries/wmc9yO6IQ3mUZCPx7MQRxg>

<sup>27</sup> <https://commons.datacite.org/doi.org?query=client.uid%3Awhno.ljncxe&resource-type=dataset>

```

519 17 print(result)
520 18 #Retrieve the Image from found object, make into a FITS File
521 19 image = SDSS.get_images(matches=result, band=['u', 'g', 'r', 'i', 'z'])
522 20
523 21
524 22 #Retrieve hdulist from FITS file
525 23 hdulist = image[0]
526 24
527 25 #Retrieve the image data from the hdulist
528 26 imageData = hdulist[0].data
529 27
530 28 #Log the image data in order to get rid of background
531 29 imageDataLog = np.log10(imageData) + 1e-8
532 30
533 31 #Save the header
534 32 header = hdulist[0].header
535 33
536 34
537 35 #Obtain the relevant headers
538 36
539 37 #Retrieve pixel scale numbers, divided amongst two parts for ra and dec
540 38 CD1_1 = header['CD1_1']
541 39 CD1_2 = header['CD1_2']
542 40 CD2_1 = header['CD2_1']
543 41 CD2_2 = header['CD2_2']
544 42
545 43 #Width of image in pixels
546 44 keyWordNAXIS1 = header['NAXIS1'] #[pixels]
547 45
548 46 #Height of image in pixels [pixels]
549 47 keyWordNAXIS2 = header['NAXIS2'] #[pixels]
550 48
551 49 #Normalize the pixel scale, then multiply by 3600 to convert units
552 50 CDELT1ArcSec = np.linalg.norm([CD1_1,CD1_2]) * 3600 #[arcsec/pixels]
553 51 CDELT2ArcSec = np.linalg.norm([CD2_1,CD2_2]) * 3600 #[arcsec/pixels]
554 52
555 53 #Set up the image
556 54 plt.xlabel("Right Ascension (arcsec)")
557 55 plt.ylabel("Declination (arcsec)")
558 56 plt.title('SDSS Logarithmic Optical Image')
559 57 vmin2 = np.percentile(imageDataLog, 85)
560 58 vmax2 = np.percentile(imageDataLog, 98)
561 59 plt.imshow(imageDataLog, cmap='viridis', extent = (0, CDELT1ArcSec *
562      keyWordNAXIS1, 0, CDELT2ArcSec * keyWordNAXIS2), vmin = vmin2, vmax =
563      vmax2)
564 60 plt.colorbar()
565 61
566 62 plt.show()

```

## REFERENCES

- 568     ????, The MIGHTEE Survey,  
 569      <https://www.mighteesurvey.org/home>
- 570     ????, SKA Telescope Specifications,  
 571      <https://www.skao.int/en/science-users/118/ska-telescope-specifications>
- 572     Blyth, S., Baker, A. J., Holwerda, B., et al. 2018,  
 573      in Proceedings of MeerKAT Science: On the  
 574      Pathway to the SKA — PoS(MeerKAT2016)  
 575      (Stellenbosch, South Africa: Sissa Medialab),  
 576      004, doi: [10.22323/1.277.0004](https://doi.org/10.22323/1.277.0004)
- 577     Bowen, I. S., & Vaughan, A. H. 1973, Applied  
 578      Optics, 12, 1430, doi: [10.1364/AO.12.001430](https://doi.org/10.1364/AO.12.001430)
- 579     Bundy, K., Bershadsky, M. A., Law, D. R., et al.  
 580      2014a, The Astrophysical Journal, 798, 7,  
 581      doi: [10.1088/0004-637X/798/1/7](https://doi.org/10.1088/0004-637X/798/1/7)
- 582      —. 2014b, The Astrophysical Journal, 798, 7,  
 583      doi: [10.1088/0004-637X/798/1/7](https://doi.org/10.1088/0004-637X/798/1/7)
- 584      Camilo, F. 2024
- 585      Cotton, W. D. 2008, Publications of the  
 586      Astronomical Society of the Pacific, 120, 439,  
 587      doi: [10.1086/586754](https://doi.org/10.1086/586754)
- 588      Cotton, W. D., & Schwab, F. R. 2010
- 589      Dawson, K. S., Kneib, J.-P., Percival, W. J., et al.  
 590      2016, The Astronomical Journal, 151, 44,  
 591      doi: [10.3847/0004-6256/151/2/44](https://doi.org/10.3847/0004-6256/151/2/44)
- 592      De Blok, W. J. G., Healy, J., Maccagni, F. M.,  
 593      et al. 2024, Astronomy & Astrophysics, 688,  
 594      A109, doi: [10.1051/0004-6361/202348297](https://doi.org/10.1051/0004-6361/202348297)
- 595      Dewdney, P., Hall, P., Schilizzi, R., & Lazio, T.  
 596      2009, Proceedings of the IEEE, 97, 1482,  
 597      doi: [10.1109/JPROC.2009.2021005](https://doi.org/10.1109/JPROC.2009.2021005)
- 598      Fernandes, L. A., & Oliveira, M. M. 2008, Pattern  
 599      Recognition, 41, 299,  
 600      doi: [10.1016/j.patcog.2007.04.003](https://doi.org/10.1016/j.patcog.2007.04.003)
- 601      Goedhart, S. 2025, MeerKAT Specifications
- 602      Gunn, J. E., Siegmund, W. A., Mannery, E. J.,  
 603      et al. 2006, The Astronomical Journal, 131,  
 604      2332, doi: [10.1086/500975](https://doi.org/10.1086/500975)
- 605      Gupta, N., Jagannathan, P., Srianand, R., et al.  
 606      2021, The Astrophysical Journal, 907, 11,  
 607      doi: [10.3847/1538-4357/abcb85](https://doi.org/10.3847/1538-4357/abcb85)
- 608      Hey, T., Tansley, S., & Tolle, K. 2009, Microsoft  
 609      Research
- 610      Holtzman, J. A., Harrison, T. E., & Coughlin,  
 611      J. L. 2010, Advances in Astronomy, 2010,  
 612      193086, doi: [10.1155/2010/193086](https://doi.org/10.1155/2010/193086)
- 613      Ivezic, Z., Kahn, S. M., Tyson, J. A., et al. 2019,  
 614      The Astrophysical Journal, 873, 111,  
 615      doi: [10.3847/1538-4357/ab042c](https://doi.org/10.3847/1538-4357/ab042c)
- 616      Jonas, J., & the MeerKAT Team. 2018, in  
 617      Proceedings of MeerKAT Science: On the  
 618      Pathway to the SKA — PoS(MeerKAT2016)  
 619      (Stellenbosch, South Africa: Sissa Medialab),  
 620      001, doi: [10.22323/1.277.0001](https://doi.org/10.22323/1.277.0001)
- 621      Juric, M., Kantor, J., Lim, K.-T., et al. 2017
- 622      Lesser, M. 2015, Publications of the Astronomical  
 623      Society of the Pacific, 127, 1097,  
 624      doi: [10.1086/684054](https://doi.org/10.1086/684054)
- 625      Lupton, R., Gunn, J. E., Ivezić, Z., et al. 2001,  
 626      The SDSS Imaging Pipelines, arXiv,  
 627      doi: [10.48550/arXiv.astro-ph/0101420](https://arxiv.org/abs/astro-ph/0101420)
- 628      Lupton, R. H., Ivezić, Z., Gunn, J., et al. 2007
- 629      Majewski, S. R., Schiavon, R. P., Frinchaboy,  
 630      P. M., et al. 2017, The Astronomical Journal,  
 631      154, 94, doi: [10.3847/1538-3881/aa784d](https://doi.org/10.3847/1538-3881/aa784d)
- 632      Moore, G. E. 2006, IEEE Solid-State Circuits  
 633      Society Newsletter, 11, 33,  
 634      doi: [10.1109/N-SSC.2006.4785860](https://doi.org/10.1109/N-SSC.2006.4785860)
- 635      NSF-DOE Vera C. Rubin Observatory. 2025,  
 636      PSTN-019: The LSST Science Pipelines  
 637      Software: Optical Survey Pipeline Reduction  
 638      and Analysis Environment, NSF-DOE Vera C.  
 639      Rubin Observatory,  
 640      doi: [10.71929/RUBIN/2570545](https://doi.org/10.71929/RUBIN/2570545)
- 641      Ratcliffe, S. 2021, SDP Pipelines Overview,  
 642      <https://skaafrica.atlassian.net/wiki/spaces/ESDKB/pages/73333333/SDP+Pipelines+Overview>
- 643      Richards, S. 2020, What Is Vignetting?
- 644      Verbunt, F., & Van Gent, R. H. 2010, Astronomy  
 645      and Astrophysics, 516, A28,  
 646      doi: [10.1051/0004-6361/201014002](https://doi.org/10.1051/0004-6361/201014002)
- 647      Woudt, P. A., Fender, R., Corbel, S., et al. 2018,  
 648      in Proceedings of MeerKAT Science: On the  
 649      Pathway to the SKA — PoS(MeerKAT2016)  
 650      (Stellenbosch, South Africa: Sissa Medialab),  
 651      013, doi: [10.22323/1.277.0013](https://doi.org/10.22323/1.277.0013)

## REFERENCES

- 653     ????, The MIGHTEE Survey,  
 654      <https://www.mighteesurvey.org/home>  
 655     ????, SKA Telescope Specifications,  
 656      <https://www.skao.int/en/science-users/118/ska-telescope-specifications>  
 657  
 658      Blyth, S., Baker, A. J., Holwerda, B., et al. 2018,  
 659      in Proceedings of MeerKAT Science: On the  
 660      Pathway to the SKA — PoS(MeerKAT2016)  
 661      (Stellenbosch, South Africa: Sissa Medialab),  
 662      004, doi: [10.22323/1.277.0004](https://doi.org/10.22323/1.277.0004)  
 663      Bowen, I. S., & Vaughan, A. H. 1973, Applied  
 664      Optics, 12, 1430, doi: [10.1364/AO.12.001430](https://doi.org/10.1364/AO.12.001430)  
 665      Bundy, K., Bershadsky, M. A., Law, D. R., et al.  
 666      2014a, The Astrophysical Journal, 798, 7,  
 667      doi: [10.1088/0004-637X/798/1/7](https://doi.org/10.1088/0004-637X/798/1/7)  
 668      —. 2014b, The Astrophysical Journal, 798, 7,  
 669      doi: [10.1088/0004-637X/798/1/7](https://doi.org/10.1088/0004-637X/798/1/7)  
 670      Camilo, F. 2024  
 671      Cotton, W. D. 2008, Publications of the  
 672      Astronomical Society of the Pacific, 120, 439,  
 673      doi: [10.1086/586754](https://doi.org/10.1086/586754)  
 674      Cotton, W. D., & Schwab, F. R. 2010  
 675      Dawson, K. S., Kneib, J.-P., Percival, W. J., et al.  
 676      2016, The Astronomical Journal, 151, 44,  
 677      doi: [10.3847/0004-6256/151/2/44](https://doi.org/10.3847/0004-6256/151/2/44)  
 678      De Blok, W. J. G., Healy, J., Maccagni, F. M.,  
 679      et al. 2024, Astronomy & Astrophysics, 688,  
 680      A109, doi: [10.1051/0004-6361/202348297](https://doi.org/10.1051/0004-6361/202348297)  
 681      Dewdney, P., Hall, P., Schilizzi, R., & Lazio, T.  
 682      2009, Proceedings of the IEEE, 97, 1482,  
 683      doi: [10.1109/JPROC.2009.2021005](https://doi.org/10.1109/JPROC.2009.2021005)  
 684      Fernandes, L. A., & Oliveira, M. M. 2008, Pattern  
 685      Recognition, 41, 299,  
 686      doi: [10.1016/j.patcog.2007.04.003](https://doi.org/10.1016/j.patcog.2007.04.003)  
 687      Goedhart, S. 2025, MeerKAT Specifications  
 688      Gunn, J. E., Siegmund, W. A., Mannery, E. J.,  
 689      et al. 2006, The Astronomical Journal, 131,  
 690      2332, doi: [10.1086/500975](https://doi.org/10.1086/500975)  
 691      Gupta, N., Jagannathan, P., Srikanth, R., et al.  
 692      2021, The Astrophysical Journal, 907, 11,  
 693      doi: [10.3847/1538-4357/abcb85](https://doi.org/10.3847/1538-4357/abcb85)  
 694      Hey, T., Tansley, S., & Tolle, K. 2009, Microsoft  
 695      Research  
 696      Holtzman, J. A., Harrison, T. E., & Coughlin,  
 697      J. L. 2010, Advances in Astronomy, 2010,  
 698      193086, doi: [10.1155/2010/193086](https://doi.org/10.1155/2010/193086)  
 699      Ivezić, Ž., Kahn, S. M., Tyson, J. A., et al. 2019,  
 700      The Astrophysical Journal, 873, 111,  
 701      doi: [10.3847/1538-4357/ab042c](https://doi.org/10.3847/1538-4357/ab042c)  
 702      Jonas, J., & the MeerKAT Team. 2018, in  
 703      Proceedings of MeerKAT Science: On the  
 704      Pathway to the SKA — PoS(MeerKAT2016)  
 705      (Stellenbosch, South Africa: Sissa Medialab),  
 706      001, doi: [10.22323/1.277.0001](https://doi.org/10.22323/1.277.0001)  
 707      Juric, M., Kantor, J., Lim, K.-T., et al. 2017  
 708      Lesser, M. 2015, Publications of the Astronomical  
 709      Society of the Pacific, 127, 1097,  
 710      doi: [10.1086/684054](https://doi.org/10.1086/684054)  
 711      Lupton, R., Gunn, J. E., Ivezić, Z., et al. 2001,  
 712      The SDSS Imaging Pipelines, arXiv,  
 713      doi: [10.48550/arXiv.astro-ph/0101420](https://arxiv.org/abs/astro-ph/0101420)  
 714      Lupton, R. H., Ivezić, Ž., Gunn, J., et al. 2007  
 715      Majewski, S. R., Schiavon, R. P., Frinchaboy,  
 716      P. M., et al. 2017, The Astronomical Journal,  
 717      154, 94, doi: [10.3847/1538-3881/aa784d](https://doi.org/10.3847/1538-3881/aa784d)  
 718      Moore, G. E. 2006, IEEE Solid-State Circuits  
 719      Society Newsletter, 11, 33,  
 720      doi: [10.1109/N-SSC.2006.4785860](https://doi.org/10.1109/N-SSC.2006.4785860)  
 721      NSF-DOE Vera C. Rubin Observatory. 2025,  
 722      PSTN-019: The LSST Science Pipelines  
 723      Software: Optical Survey Pipeline Reduction  
 724      and Analysis Environment, NSF-DOE Vera C.  
 725      Rubin Observatory,  
 726      doi: [10.71929/RUBIN/2570545](https://doi.org/10.71929/RUBIN/2570545)  
 727      Ratcliffe, S. 2021, SDP Pipelines Overview,  
 728      <https://skaafrica.atlassian.net/wiki/spaces/ESDKB/pages/>  
 729      Richards, S. 2020, What Is Vignetting?  
 730      Verbunt, F., & Van Gent, R. H. 2010, Astronomy  
 731      and Astrophysics, 516, A28,  
 732      doi: [10.1051/0004-6361/201014002](https://doi.org/10.1051/0004-6361/201014002)  
 733      Woudt, P. A., Fender, R., Corbel, S., et al. 2018,  
 734      in Proceedings of MeerKAT Science: On the  
 735      Pathway to the SKA — PoS(MeerKAT2016)  
 736      (Stellenbosch, South Africa: Sissa Medialab),  
 737      013, doi: [10.22323/1.277.0013](https://doi.org/10.22323/1.277.0013)