

# SNOWPAC: a multiscale cubic B-spline wavelet compressor for astronomical images

Jesus Pulido,<sup>1★</sup> Caixia Zheng,<sup>2★</sup> Paul Thorman<sup>3</sup> and Bernd Hamann<sup>1</sup>

<sup>1</sup>*Department of Computer Science, University of California, One Shields Avenue, Davis, CA 95616, USA*

<sup>2</sup>*College of Information Sciences and Technology, Northeast Normal University, 2555 Jingyue Street, Changchun 130117, China*

<sup>3</sup>*Departments of Physics and Astronomy, Haverford College, 370 Lancaster Avenue, Haverford, PA 19041, USA*

Accepted 2020 February 7. Received 2020 February 2; in original form 2019 January 9

## ABSTRACT

As more advanced and complex survey telescopes are developed, the size and scale of data being captured grows at increasing rates. Across various domains, data compression through wavelets has enabled the reduction of data size and increase in computation efficiency. In this paper, we provide qualitative and quantitative tests of a new wavelet-based image compression method compared against the current standard for astronomical images. The analysis is improved by making use of state-of-the-art object detection systems to accurately measure the impact of the compression. We find that a combination of lossy wavelet-based methods, efficient quantization, and lossless dictionary compressors can preserve up to 98 per cent of astronomical objects at a 10:1 compression ratio. This significant reduction in file size also preserves astronomical object properties better than existing methods. These methods help further reduce future workloads for image-heavy processing pipelines.

**Key words:** methods: data analysis – techniques: image processing – astrometry.

## 1 INTRODUCTION

The Large Synoptic Survey Telescope (LSST) is a wide-field survey telescope that is currently being constructed and is projected to produce hundreds of gigabytes of data per night. Transferring and processing raw data daily becomes challenging as network and input/output (I/O) systems become bottlenecks for this large data transfer problem. In the past, lossless data compressors provided sufficient data reduction but as data sets become larger and more complex, lossy methods have become a necessity. Both transferring data between on-site systems and sharing data for collaborative research imposes significant bottlenecks making smart but lossy data reduction schemes attractive.

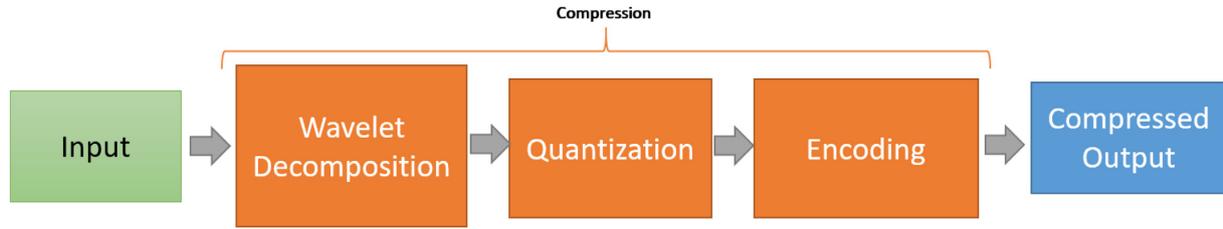
The Flexible Image Transport System (FITS) data format is the standard astronomical data container used in large by the astronomy community. The subject of astronomical image data compression is not entirely new.

There are a range of tools for generating FITS files that employ different compression techniques. CFITSIO is a tool developed by Pence (1999) that provides subroutines for reading and writing FITS data files. A more recent codevelopment has been integrated called FPACK and has been proposed for FITS file compression (Pence, Seaman & White 2011). This tool provides several general purpose compression algorithms: GZIP, RICE, HCOMPRESS, and PLIO. The most popular and simplest method used for general lossless

compression is GZIP (GNU 1997). Although lossless, GZIP provides low compression ratios for large, dynamic data sets. Generally, the RICE compression algorithm (Rice, Yeh & Miller 1993) provides better lossless compression ratios and is faster than GZIP. HCOMPRESS (White, Postman & Lattanzi 1992) can be either a lossy or lossless compressor based on the H-transform (Fritze et al. 1977). HCOMPRESS uses a generalized 2D Haar, quantization, and quad-tree coding to achieve its lossy compression. This method generally compresses slightly better than RICE at the cost of more compute time.

Previous work by Pence (2009) conducted a feasibility study for the use of a more efficient compression algorithm named BZIP2 in astronomical images. That work concluded that BZIP2 was an order of magnitude slower compared to RICE and required significantly larger block sizes to achieve same compression efficiency. Unlike other compressors, JPEG2000 is keyed as a true image compression standard that intends to take advantage of multidimensional data. Similar applications by Peters & Kitaeff (2014), Kitaeff et al. (2015), and Vohl, Fluke & Vernardos (2015) analysed the impact of JPEG2000 lossy compression on radio astronomy imagery. Since JPEG 2000 is a non-standard method for astronomical images, several challenges may appear such as the inability to process floating-point images, requiring the conversion to integers prior to processing. This simplification of the data may lead to issues when handling images with very high dynamic ranges in intensities. Methods such as those in Vohl et al. (2017) overcome this by adding 32-bit support, but do not go in depth on the effects of achieving ‘extreme’ 35:1 compression. Compared to traditional observational

\* E-mail: [jpulido@ucdavis.edu](mailto:jpulido@ucdavis.edu) (JP); [zhengcx789@nenu.edu.cn](mailto:zhengcx789@nenu.edu.cn) (CZ)



**Figure 1.** The image compression pipeline has an image as input, applies cubic B-spline wavelet decomposition, and performs floating-point-to-integer quantization for an initial lossy step. After, the encoding step efficiently stores lossy coefficients in a lossless format.

sky survey data, radio astronomy imagery is significantly noisier and lacks the structures needed for general purpose compression algorithms to excel.

Other methods have suggested the deviation from the FITS file format to a more compression-friendly container such as Hierarchical Data Format 5 (HDF5) in Price, Barsdell & Greenhill (2014) and Masui et al. (2015). The advantage of switching to HDF5 is its availability of internal data filters. One such filter explored in both works is BITSHUFFLE, a lossless compressor that shifts binary data for more efficient encoding. As with other lossless compressors, BITSHUFFLE is hardly able to achieve a compression ratio beyond 2:1, unable to achieve as much compression as other lossy methods. Finally, singular value decomposition methods (SVD) such as those in Kolev, Tsvetkova & Tsvetkov (2012) and Morii et al. (2017) are used to compress astronomical movie data in a lossy representation.

When lossy compression is used, the goal is to preserve as many astronomical objects as possible without significantly altering their science-critical characteristics, such as the total flux, the spatial extent and profile of the flux, and the orientation of the source profile on the sky. In a typical pipeline, object detection and extraction is performed by widely used tools such as Source Extractor (SEXTRACTOR) by Bertin & Arnouts (1996) or the techniques of Zheng et al. (2015), which attempt to extract and characterize these astronomical objects.

The work presented in this paper introduces a new wavelet-based lossy image compression method for astronomical images. This new method is able to achieve extreme levels of compression while preserving features of interest (astronomical objects). To accurately measure the impact of compression, state-of-the-art detection systems are used to accurately measure differences, both qualitatively and quantitatively.

This paper is structured as follows. Section 2 gives an overview on the construction of the method and usage to achieve these results. Section 3 evaluates the performance, cost, and accuracy of our method by comparing the obtained results with those of competing methods. Section 4 provides conclusions and points out possible directions for future work.

## 2 METHOD

### 2.1 Implementation

A combination of several techniques can be used in a pipeline; see Fig. 1, where the data compression pipeline is depicted with its main components marked in the centre. Once a FITS image is read, the largest data reduction occurs in the initial wavelet decomposition and wavelet coefficient floating-point-to-integer quantization. Using the native hierarchical structure generated by a wavelet decomposition, the highest magnitude wavelet coefficients are preserved with high precision and the lowest magnitude at

a lower precision via quantization. Once completed, a lossless encoding step is performed where, optionally, further compression can be achieved, at the cost of compute time. Like compression, decompression only requires us to decode the stored coefficients and perform a wavelet reconstruction, returning the coefficients to real space.

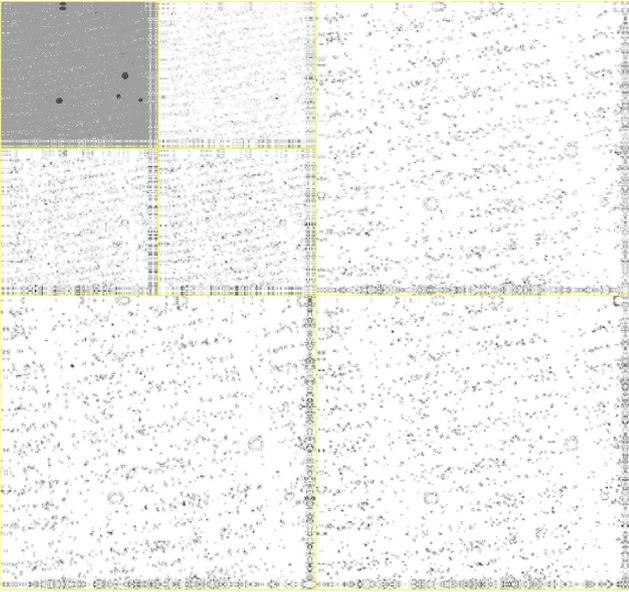
We present a data reduction pipeline with flexible components, trading-off lower file sizes for higher computation time. As a result of the compression, a significantly larger subset of astronomical imagery can be stored, processed, and transferred. The complete method named SNOWPAC (SpliNe Wavelet Packing And Compression; Pulido 2019) is made available in both C++ and MATLAB variants. The end of a typical processing pipeline may contain additional astronomical object detection components that quantify the number of stars and galaxies in a certain region of the sky. The usage of wavelet methods in our pipeline leads to advantageous properties such as data streaming, selective decompression of data subsets, and denoising of astronomical images.

### 2.2 Wavelet compression

The wavelet transform is a generalization of the Fourier transform, using bases that represent both location and spatial frequency (Daubechies 1992). Typical wavelet signals contain several vanishing moments that allow for a sparse but accurate representation of an input data set, with only a small number of coefficients. A signal is decomposed through multiple steps involving ‘folds’ at the largest scale until data are reduced to the smallest scale. A 2D example for an astronomical image is shown in Fig. 2. Here, a fold is performed twice for each dimension until a set of low- and high-pass coefficients is obtained. The upper-left quadrant demonstrates the capability of wavelets to preserve coherent features (astronomical objects) and, at the same time, to filter out small-scale signal behaviour usually correlated to Gaussian noise. Additionally, this behaviour creates a multiscale structure where each fold represents features at a certain resolution. This capability allows a compressor to reduce data to extreme levels without losing important objects, even at medium to high levels of dynamic range. Additionally, a 2D transform can be expanded to 3D data arising in multiple applications, such as stacked imagery, supporting the signal acquisition, and preservation of many images of the same sky region.

Unlike the classical and simple Haar wavelet, biorthogonal B-spline wavelets are an extension with similar usage and implementation, providing the capability to capture more complex behaviour via their basis functions (Cohen, Daubechies & Feauveau 1992).

In a recent study by Pulido et al. (2016), several multiresolution representation methods, including higher order B-spline wavelets, were tested for their ability to capture a broad range of quantities pertaining to turbulent data representation using a reduced set of

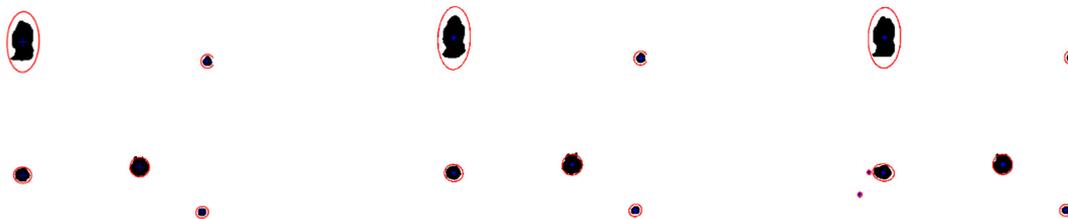


**Figure 2.** 2D data decomposition. Decomposition of an astronomical image using two levels produces a series of low- and high-pass coefficients using cubic B-spline wavelets. Astronomical objects are preserved, while small-scale noisy artefacts are filtered out.

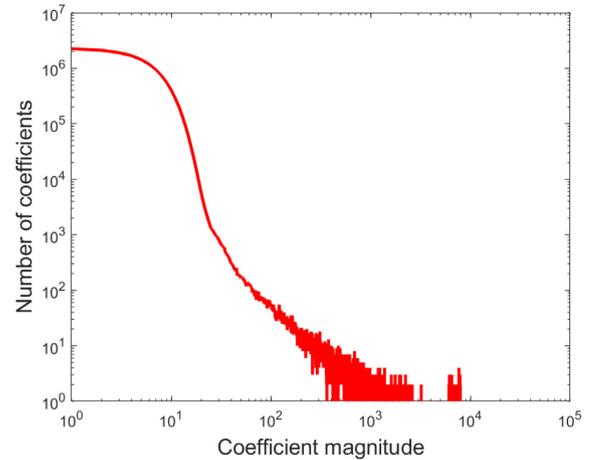
coefficients. The higher order B-spline families consistently scored at or near the top based on the metrics considered. One of the most important metrics was the preservation of coherent structures using the least number of coefficients through hard thresholding. Here, the fourth family of B-spline wavelets (cubic) is considered as the method of choice for compression and analysis. A cubic B-spline wavelet transform can capture and preserve directionality well. This characteristic is particularly important in the effort to preserve the structures of astronomical objects with high accuracy. We employ the discrete version as described by Shensa (1992) in our processing pipeline.

### 2.3 Hard thresholding and quantization

Hard thresholding selects only distinct coefficients up to a certain cut-off threshold, while setting to zero all other coefficients, after a wavelet decomposition has been performed. The primary criterion for selecting a cut-off threshold is a coefficient's (absolute) magnitude. Coefficient magnitude can be understood as signal energy for a frequency (or basis function). Thresholding is done by sorting the coefficients by magnitude and selecting the largest ones, either via a cut-off percentage or a threshold value.



**Figure 3.** Coefficient quantization versus thresholding. The shown objects result from performing lossy coefficient operations, illustrating their preservation properties. Quantization of all coefficients using a 10:1 ratio (left), hard thresholding using 10 per cent coefficients (centre), and original data (right) show that pronounced stars and galaxies indicated by red ellipses are preserved after a binary image operation during object detection.



**Figure 4.** Coefficient magnitude sorting. This plot shows the magnitude-sorted cubic B-spline wavelet coefficients for an LSST data set. Many low-magnitude redundant coefficients exist, making possible a lower precision representation using quantization without significant data information loss.

This lossy compression approach is not sufficient for achieving a large enough reduction for astronomical images. Controlling the precision of a quantized data representation of coefficients is preferred.

Fig. 3 shows a comparison between hard thresholding for a fixed percentage of coefficients (10 per cent relative to original size) and a floating-point-to-integer quantization for a compression ratio near 10:1 and the effect on object detection. Both reduction techniques produce similar results in terms of detected objects, but there are differences. Most significantly, an object's shape is better preserved when using all wavelet coefficients at a reduced precision rather than using a subset of large magnitude coefficients via hard thresholding.

Fig. 4 shows the coefficients sorted by magnitude. For a typical sky survey data set, many low-magnitude coefficients exist that can be reduced in precision via quantization, with little impact as long as high-magnitude coefficients are preserved with higher precision. This principle is used to achieve lossy compression.

A wavelet transform produces a floating-point data set. Typically, FITS images are integer data. First, we convert them to floating-point representation. For efficient data reduction we apply floating-point-to-integer quantization across all wavelet coefficients (see Section 2.3). By performing requantization, we efficiently encode integer coefficients.

We use a highly efficient floating-point-to-integer quantization method that shifts the decimal of the floating-point representation to the right, to preserve the integer portion and truncate the fractional half. This method is available in HDF 5's (Folk et al. 2011)

**Table 1.** Floating-point-to-integer quantization. Using a weight as input, a quantization operation shifts the decimal during conversion to control the amount of bits needed for storage.

Weight	Original value: 116618.821101		
	Shifted	Q integer	Bits
0.01	1166.188211	1166	12
0.1	11661.88211	11661	15
1	116618.8211	116618	18
10	1166188.211	1166188	22
100	11661882.11	11661882	25

‘scale\_offset’ filter and ensures minimal additional loss in precision when combined with wavelets, as most of a coefficient’s entropy is captured in its integer component. For example, a coefficient value of 3027.567812346 is converted to the integer 3027567, using a three-decimal scale off-set shift of 1000. Table 1 lists the number of bits used to represent a floating-point value as an integer.

This technique allows us to control the precision of coefficients. By using the natural wavelet hierarchy, coefficients at the lowest bands have the largest magnitudes and importance. Therefore, one should preserve them at higher precision compared to those at the highest bands, having lowest magnitudes. Example weights for an extraction of eight levels are [100, 10, 10, 10, 10, 10, 1, 1], where more bits are used at the lowest and less bits at the highest levels.

## 2.4 Encoding

The simplest and most common encoding method available is run-length coding (RLE), which efficiently packs wavelet coefficients after quantization. RLE counts the number of repeated byte spaces in a data set and compacts the representation with a single value.

As a more valuable option, we use off-the-shelf, dictionary-based lossless compression methods to expand the overall capabilities of our method. The LZ4 method (Collet 2011) is characterized by its extremely high speed, coming at the expense of low compression ratios; the BZIP2 method (Steward 1996) leads to high compression ratios, coming at the expense of low speed. The GZIP technique, traditionally applied to sky survey data sets, exhibits performance that lies somewhere in the middle, producing average compression ratios and computation times. When combined with the floating-point-to-integer quantization method used in the previous section, applying lossless compression on the lossy represented coefficients demonstrates to be an efficient way to pack coefficients for storage.

## 2.5 Decompression

Decompression of astronomical images using the pipeline shown in Fig. 5 takes much less effort than compression. Decompression merely decodes the coefficients for reconstruction and reapplies inverse weights for requantization.

## 2.6 Object detection

Analysis on astronomical imagery is typically performed after data acquisition, including object detection of stars and galaxies. By using these methods, we can quantify the impact on compression and decide which method best preserves additional properties in images.

In Zheng et al. (2015), an improved method for detecting objects of interest (galaxies and stars) in astronomical images was recently

presented. This new method combined various global and local subroutines to improve detection over existing methods. In this work, after a global detection scheme is applied, refinement is done by dividing the entire image into several irregularly sized subregions using the watershed segmentation method. A more refined detection procedure is performed in each subregion by applying adaptive noise reduction and a layered strategy to detect bright objects and faint objects, respectively. Finally, a multithreshold technique is used to separate objects that initially blended together. On both simulated and observational data, this method detected more real objects than SEXTRACTOR at comparable object counts. The method also had an increased chance of successfully detecting very faint objects, up to 2 mag fainter than SEXTRACTOR on similar data. Because of the improved detection properties of this method and its measurable ability to extract faint objects, it is best suitable for testing compression.

## 3 TESTS AND RESULTS

### 3.1 Data sets

The following data sets are used to evaluate the accuracy of these methods.

#### 3.1.1 DLS data set

The Deep Lens Survey (DLS) data set (Wittman et al. 2002) used for testing is taken from a deep *BVRz* imaging survey that covers 20 deg<sup>2</sup> to a depth of approximately 28th mag (AB) in *BVR* and 24.5 in *z* (AB). Our subsample includes four FITS files of size 4096 × 4096. The survey was taken by the National Optical Astronomy Observatory’s (NOAO) Blanco and Mayall 4-m telescopes. The images correspond to four waveband filters of the same sky area: *R*, *V*, *B*, and *z*.

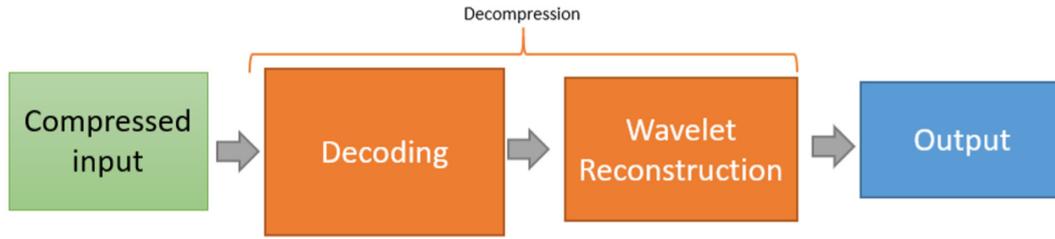
Because of the observational nature, there is no ground truth catalogue. To generate one, we compare the objects detected in different bands to gather a coinciding subset of detected objects that are likely to be real. Since the *R* band is used as a detection band by the DLS survey due to its excellent depth and seeing, detection results of other bands are often verified against it. To compare all types of compression methods, we focus on the *R*-band FITS file<sup>1</sup> due to its prominence in sky survey analysis. In addition, for several levels of lossy compression we detect objects and their properties to compare to the original.

#### 3.1.2 LSST data set

The LSST data set (LSST Science Collaborations 2009) is a simulated set of images that use target parameters for the future LSST project. Using GALSIM, the lens properties and other characteristics of the telescope can be simulated to produce expected outputs for the upcoming observational data sets. Our sample set of FITS files<sup>2</sup> is of size 4000 × 4000 containing several bands. Unlike the DLS data set, which represents a series of stacked images, this simulated data has noise added to replicate that of the future telescope. Each pixel in an image represents about 0.2 arcmin of real space, so we can expect a single fits file to represent 800 × 800 arcmin<sup>2</sup> or 0.05 deg<sup>2</sup> depth covered by these images.

<sup>1</sup><http://dls.physics.ucdavis.edu/imdownload.html>

<sup>2</sup><https://www.lsst.org/scientists/simulations/phosim>



**Figure 5.** Decompression pipeline. The decompression pipeline requires much less effort than the compression pipeline. Decompression can be performed efficiently and, in some cases, can omit reconstruction altogether for fast data analysis.

**Table 2.** Comparison between compression methods. Traditional lossless compression methods result in larger file sizes, while our lossy method achieves a high compression ratio while preserving 98 per cent correct objects.

Compression type	Size (MB)	CRatio	Time (s)
None	64	1.0	–
LZ4	17.92	3.57	11.87
GZIP	15.04	4.26	18.89
BZIP2	9.54	6.71	14.42
SNOWPAC	6.55	9.77	10.32

The LSST data contain the simulated raw output of a single survey image therefore it is difficult to run object detection for validation on it. Without a form of pre-processing, object detection methods will introduce a significant amount of false positives (FP) without stacking. The object detection method used for this analysis is intended to be used on stacked sets of images. As a pre-processing step, we have emulated the effects of image stacking on one of the regions for the LSST data, which in-turn removes the most obvious pixel-sized artefacts. Thereafter, we can run compression and object detection schemes on this data.

### 3.2 Compression types

In Section 2.4, we discussed the usage of various alternative lossless compression schemes to efficiently encode our coefficients. One might ask, why not use these directly on the original data? Table 2 explores that comparison between these traditional compression techniques alongside the wavelet encoding pipeline we have introduced in this paper.

Compared to the alternative lossless methods, we can use our lossy method combined with the high efficiency of BZIP2 encoding to achieve compression ratios as high as 9.77:1 preserving over 98 per cent correct objects, compared to stand-alone applications of LZ4 (3.57:1), GZIP (4.26:1), and BZIP2 (6.71:1). Additionally, because the quantization of our coefficients reduces the complexity of the representation of the data through bit reduction, BZIP2 operates much faster compared to the original data. As a result, our method overall, i.e. computation of the B-spline wavelet transform, quantization, and encoding, operates faster when paired with BZIP2 (10.32 s) compared to lossless BZIP2 stand-alone on the original data (14.42 s). When comparing a MATLAB implementation of SNOWPAC against the lossy methods implemented in C inside of FPACK, we found that both HCOMPRESS and RICE to be faster in compute time. A faster, more optimized C++ version of SNOWPAC is planned to be released in the future to achieve faster performance compared to the current MATLAB version.

An additional advantage our method is the ability to achieve further compression by reducing the precision of coefficients further, and is explored in the follow-up sections. When the wavelet coefficients are quantized to integers, the complexity of the data is reduced therefore running an off-the-shelf compressor such as BZIP2 or LZ4 will reduce the overall compression time compared to running them on the original data set. The outcome is that our method has a lower total compute time and a higher compression ratio at the cost of losing numerical precision.

The choice of the quantization precision (weights) for the smallest coefficients has the most impact on compression ratios on data sets, since that is the step of our method that poses the most risk of losing data fidelity. These changes are evaluated in the next sections.

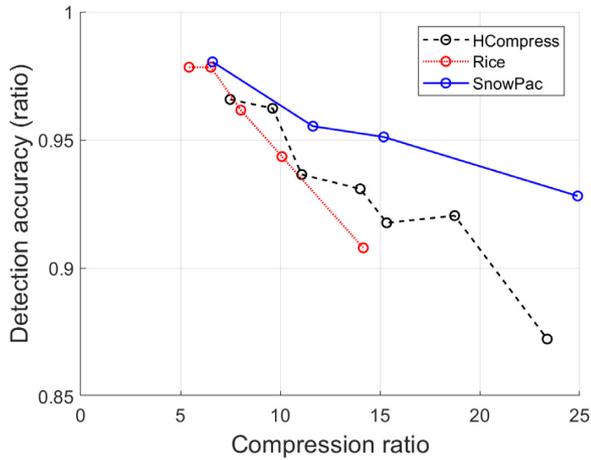
Besides lossless methods, we also compare against lossy methods available in FPACK (Pence et al. 2011), which include HCOMPRESS and RICE. Both methods were tested using recommended parameters from the FPACK user manual, which include RICE quantization levels of 8, 4, 2, 1, 0.5, and HCOMPRESS scales of 0.5, 1, 1.5, 2, 2.5, 3, 4 and a texture size of  $256 \times 256$ . The lossy compressor PLIO is available in FPACK, but was unable to process our both data sets reporting a ‘data out of range for PLIO compression (0–2\*\*24)’ error. In addition, the JPEG2000 implementation available in MATLAB was unable to process floating point FITS imagery. Because of these issues, both of these compressors were omitted.

### 3.3 Evaluation – data set 1

We evaluate the impact of our compression scheme by controlling the quantization amounts for wavelet coefficients on the DLS data set. We then compare several quantities related to the detection of objects in the imagery, such as the number of objects detected and preserved compared to the original, uncompressed data set.

Fig. 6 compares the ratio of correctly detected objects per compression ratio relative to the original data set. The accuracy ratio is a function of objects that are verified against the ground truth, over the total number of objects in the ground truth. The compression ratio is a function of the original file size over the compressed file size. When comparing our method against HCOMPRESS and RICE, we are able to achieving a higher compression ratio while preserving a significant amount of correct objects. Additionally, our method is able to achieve extreme levels of compression (25:1 ratio) and stay above an accuracy ratio of 0.925. We observe a staircase effect with HCOMPRESS, where accuracy is greatly affected for decimal-valued scales such as 0.5, 1.5, and 2.5.

To compare all methods as equitably as possible, we have selected the compressed data sets closest to a 10:1 compression ratio as seen in Table 3. From this data, the original ground truth contained 1433 objects. The general behaviour across all lossy compression methods shows that the total number of detected objects increases as



**Figure 6.** Detection accuracy versus compression ratio on DLS. Detection rates for our method are better preserved as we achieve higher levels of compression.

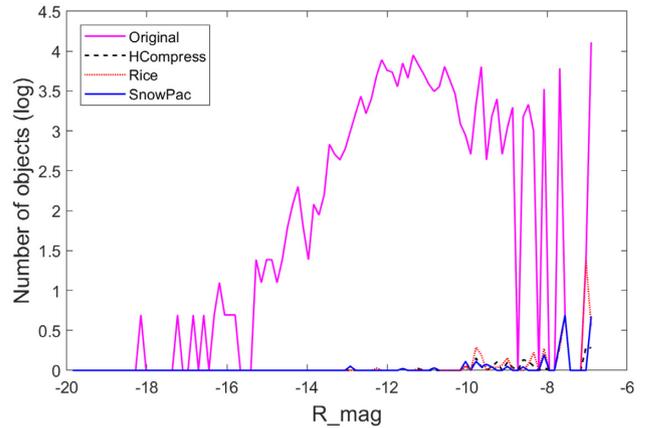
**Table 3.** DLS 10:1 compression ratio comparison. At near a 10:1 compression ratio, all methods begin to introduce noise resulting in more objects being detected incorrectly. From these compression types, our method resulted in the highest compression ratio while nearly matching HCOMPRESS in accuracy, and surpassing RICE in all metrics.

Compressor	CRatio	Total	Correct	Accuracy	FP
None	1.0	1433	–	–	–
HCOMPRESS	9.61	1487	1379	0.962	108
RICE	10.08	1471	1352	0.943	119
SNOWPAC	11.64	1453	1369	0.955	84

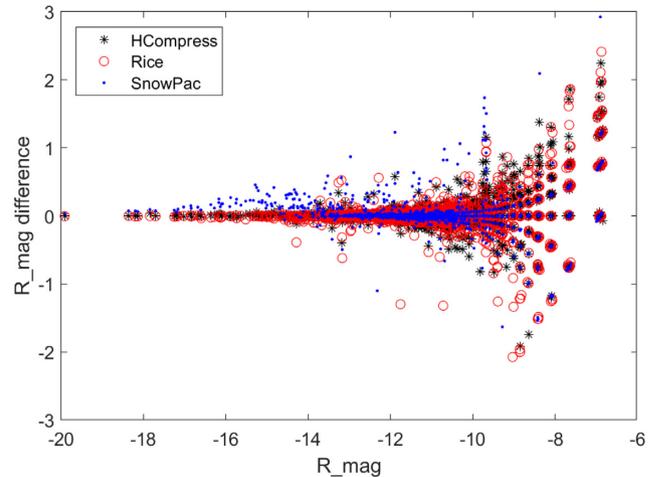
more noise is introduced through data precision loss. As expected, the number of correctly detected objects drops and increases the likelihood of FP. Our method shows that we are able to achieve the highest compression ratio, reduce the number of FP in detected objects, and achieve near the highest accuracy (true positives) for correct objects while compressing an additional 2.00 ratio over HCOMPRESS.

Alongside detection rates, we examine the effect on objects'  $R_{\text{mag}}$  ranges and their quantities in Fig. 7. We map the distribution of objects throughout different  $R_{\text{mag}}$  ranges and observe the residual behaviour. All methods perform similarly, showing that the total flux of each object is being preserved extremely well, as represented by the very small residuals. This behaviour holds true up to a certain point; an  $R_{\text{mag}}$  of  $-10$ . After that point, much elevated amounts of residuals are observed for RICE and HCOMPRESS that represent larger deviations from the original  $R_{\text{mag}}$  distribution. Our method, SNOWPAC, is able to achieve some of the lowest residuals across all levels in latter  $R_{\text{mag}}$  ranges while achieving the highest compression ratio overall.

A per object  $R_{\text{mag}}$  analysis is shown in Fig. 8, where each object's original  $R_{\text{mag}}$  is compared as a function of the  $R_{\text{mag}}$  difference. All methods perform within a  $(+/-)$  1  $R_{\text{mag}}$  difference up to about an  $R_{\text{mag}}$  of  $-12$ , then objects and their quantities begin to deviate for some objects. Beyond this range, our method has more objects clustered near zero signifying a better preservation of  $R_{\text{mag}}$  compared to the other methods. One point to note is the range between  $-17$  and  $-13$ , where SNOWPAC has slightly more error than HCOMPRESS and RICE. Despite this small fluctuation in



**Figure 7.**  $R_{\text{mag}}$  difference of compressed DLS data. The original  $R_{\text{mag}}$  quantity (magenta) is used to compare against the residual of lossy compression methods. SNOWPAC (blue) has the least amount of missing  $R_{\text{mag}}$  quantity for its detected objects compared to HCOMPRESS (black) and RICE (red). Bright and large objects are preserved well in general, but fainter objects classified by their  $R_{\text{mag}}$  are better preserved.

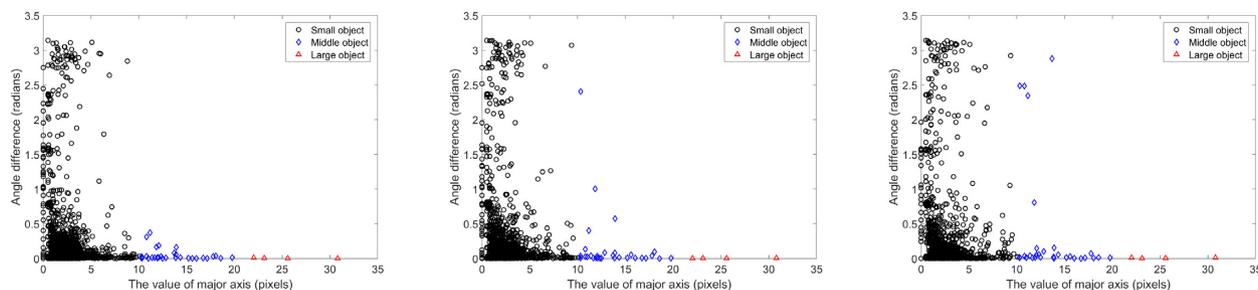


**Figure 8.**  $R_{\text{mag}}$  difference scatter plot for compressed DLS data. The original  $R_{\text{mag}}$  quantity is compared for each detected object and the difference is plotted. In general, all methods perform well up to  $R_{\text{mag}}$  of  $-12$ , where then quantities tend to deviate more than  $(+/-)$  1. Notably, SNOWPAC has more objects with less error at larger ranges, but introduces a small amount of error in ranges between  $-17$  and  $-13$  compared to the other methods.

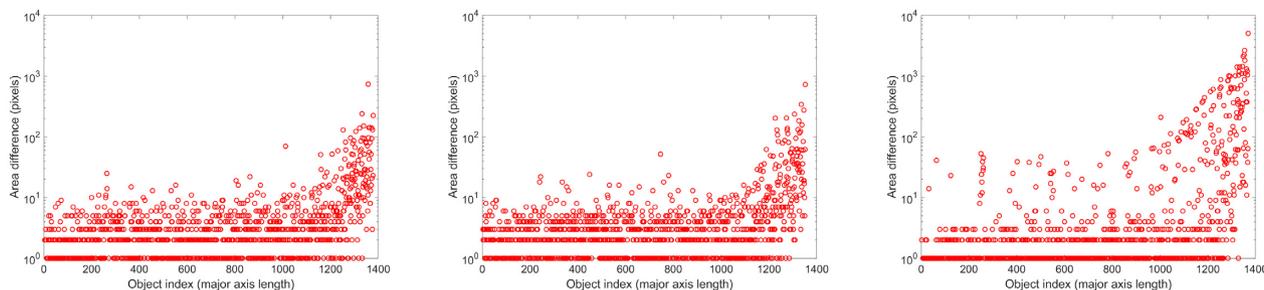
error, the detection results were not affected as observed in the rest of the analysis in this section.

To further evaluate the quality of lossy compression, three compression methods near the same compression ratio (10:1) are explored in derived object properties in Figs 9–11.

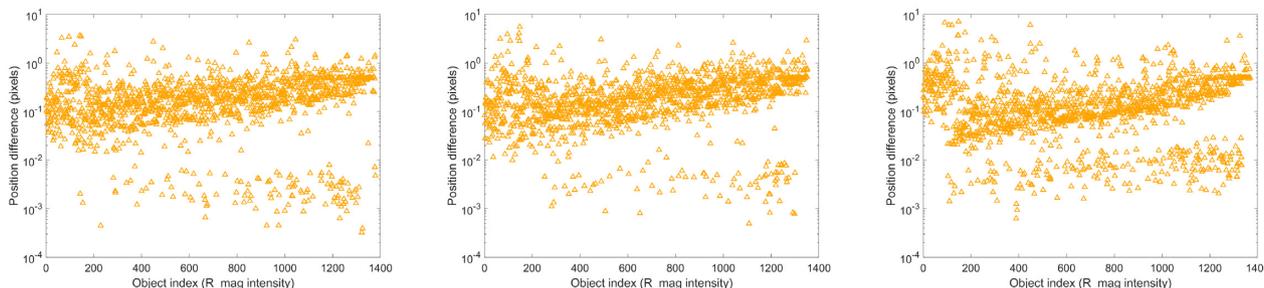
One important property of detected and classified objects is their ellipticity, which is used to distinguish stars from galaxies and must be preserved, along with the angle of the major axis, in order to properly measure weak gravitational lensing. Fig. 9 compares the difference of the major axis in ellipticity and how it degrades. Lossy compression of medium and large objects shows to have minimal impact on the direction of the major axis. The clustering of objects towards the bottom of the y-axis shows that this property is preserved very well across different compression types. When evaluating the smallest objects, high levels of lossy compression may cause these



**Figure 9.** Object angle difference. From left to right: HCOMPRESS, RICE, and SNOWPAC on DLS data. Angle distributions exhibit clustering towards the bottom and left, indicating that many objects are detected with low error. Larger amounts of clustering for smaller objects can be observed for SNOWPAC, indicating a lower error compared to the rest. Large and middle-sized objects have similar behaviour across all compression methods.



**Figure 10.** Object area difference. From left to right: HCOMPRESS, RICE, and SNOWPAC on DLS data. Sorted by the length of the major axis of objects, larger objects have their area most affected by lossy compression. While both HCOMPRESS and RICE have tighter bounds in area difference, there is more variance in preservation quality compared to SNOWPAC. While SNOWPAC has a few objects with higher residual area, a larger portion of objects are between 1 and 10 compared to the other methods.



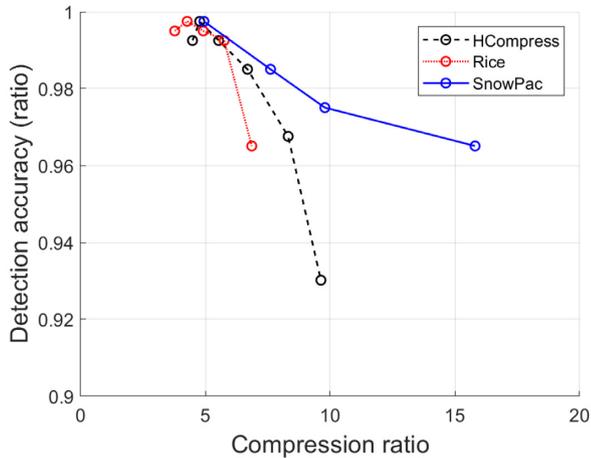
**Figure 11.** Object position difference. From left to right: HCOMPRESS, RICE, and SNOWPAC on DLS data. Sorted by R\_mag intensity, most objects tend to have their centre positions preserved pretty well. As subpixel accuracy is used, SNOWPAC is observed to have consistently lower difference amounts compared to the rest.

objects that are oftentimes circular in nature to swap their major and minor axes. Although we have corrected this by comparing only the shortest difference between these two axes, the small and circular nature still prompts for larger shifts in direction for a small subset of the faintest objects. The clustering for the small objects towards the bottom-left shows a greater preservation of ellipticity for our method compared to HCOMPRESS and RICE.

Fig. 10 derives the area for each detected object relative to the original data set. Objects in the figure are sorted by the length of the major axis, meaning larger objects will be to the right of the plot. At a 10:1 compression ratio, areas of the largest objects start to have noticeable differences. While SNOWPAC has more objects with lower error, the maximums of the larger objects are impacted more than HCOMPRESS and RICE. This can be attributed to the underlying B-spline-based method that tends to erode high-frequency regions

around objects. Larger objects seem to be the most susceptible to this but overall more objects with lower errors are achieved.

Fig. 11 uses a flux-weighted, mean position for the area pixels to derive the centre position of objects allowing for subpixel accuracy. At a 10:1 compression ratio, all objects are within 2 pixels accuracy to the original with the majority being within a single pixel accuracy and below. The distribution of object positions shows two visible clusters of data points forming, those between a  $10^{-2}$  and  $10^0$  error and another between  $10^{-3}$  and  $10^{-2}$ . All three methods work within similar ranges in accuracy, but SNOWPAC achieves more objects with lower error (bottom) and has a tighter set of clustered points for slightly less accurate objects (top) compared to the rest. The tighter clustering of object errors showing low variance is preferred as it makes error more predictable. Finally, SNOWPAC is observed to have more objects with a smaller difference clearly visible by



**Figure 12.** Detection accuracy versus compression ratio. The constant, sparse object nature of the LSST data allows our method to achieve higher detection rates and compression ratios. Other methods exhibit significant drop-offs in accuracy. The RICE compressor is unable to reduce the data any further.

**Table 4.** LSST 10:1 compression ratio comparison. The best compression possible by RICE is 6.86 therefore that is used. At near a 10:1 compression ratio, our method allows for higher number of correct objects, a higher accuracy ratio and better compression ratio.

Compressor	CRatio	Total	Correct	Accuracy	FP
None	1.0	401	–	–	–
HCOMPRESS	9.63	374	373	0.930	1
RICE	6.86	389	387	0.965	2
SNOWPAC	9.79	396	391	0.975	4

the clustering seen between a difference of 0.001 and 0.01 while achieving the highest compression ratio overall.

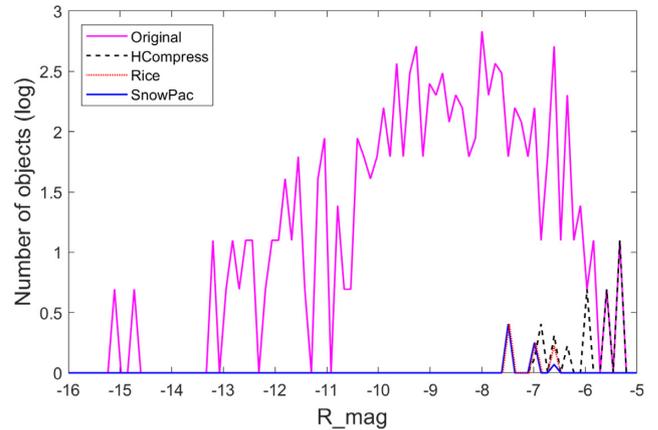
### 3.4 Evaluation – data set 2

The LSST data set poses several more challenges due to its noisy, non-stacked nature but contains a lesser amount of astronomical objects per observable area. We compress the data set at various compression ratios and show the results in Fig. 12.

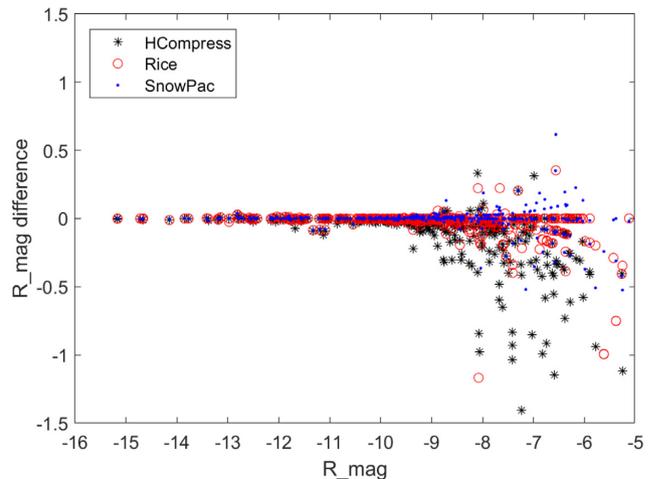
The LSST data are less object dense when compared to the DLS due to its higher resolution and smaller subsection of the sky. Objects also tend to be more elliptical. We extract a total of 401 objects and were able to preserve nearly 97.5 per cent of the originally detected objects with a 10:1 compression ratio. Because of the data being less object dense, we achieve higher compression results than HCOMPRESS and RICE; however, lossy compression in general may compromise the future ability to stack multiple images to detect objects too faint to be detected on a single image. The data sets with input parameters closest to a 10:1 compression ratio were selected and analysed starting in Table 4.

The highest compression ratio achievable by RICE was 6.86 and therefore used in this comparison. Both our method and HCOMPRESS achieved a near 10:1 ratio, without method compressing slightly better. Despite having the best compression for this comparison, we achieve a 0.975 accuracy ratio when versus HCOMPRESS’ 0.930 and RICE’s 0.965.

Fig. 13 compares the flux property of objects that were detected and their  $R_{\text{mag}}$  ranges. As shown, all of the methods are able to



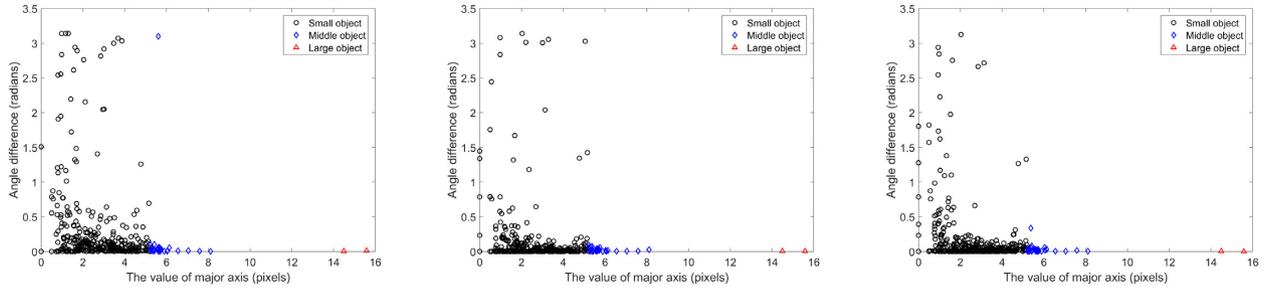
**Figure 13.**  $R_{\text{mag}}$  of compressed LSST data. The original  $R_{\text{mag}}$  quantity (magenta) is compared against the residual of various lossy compression methods. All methods are able to reproduce the original data well, but only SNOWPAC (blue) has the least amount of missing  $R_{\text{mag}}$  compared to HCOMPRESS (black) and RICE (red). Specifically, objects with the smallest flux are better preserved.



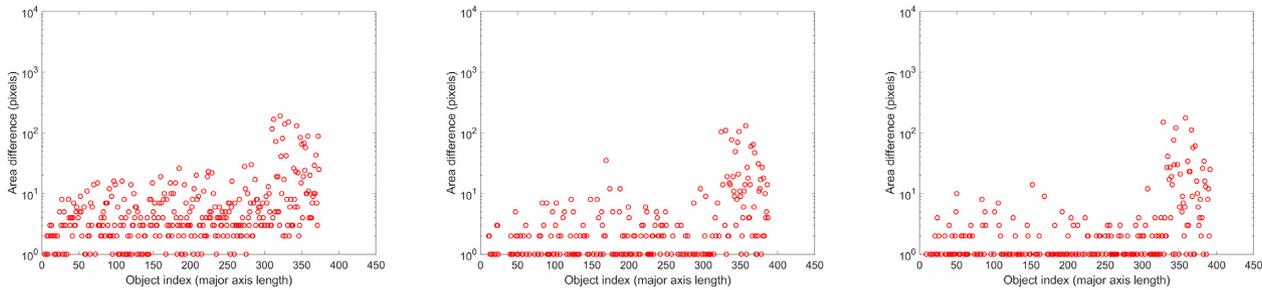
**Figure 14.**  $R_{\text{mag}}$  difference scatter plot for compressed LSST data. The original  $R_{\text{mag}}$  quantity is compared for each detected object and the difference is plotted. The distribution of delta  $R_{\text{mag}}$  shows that SNOWPAC preserves this quantity more accurately than HCOMPRESS and RICE.

preserve flux well up to a certain range. For ranges that correspond to some of the faintest objects, SNOWPAC is able to preserve the  $R_{\text{mag}}$  slightly better despite being able to achieve a higher compression ratio (9.79) compared against both RICE (6.86) and HCOMPRESS (9.63). The largest amounts of residual can be found for both of the other methods when they failed to detect objects in the  $-5.0$  to  $-7.0$  range. Effectively, objects in those flux ranges are completely lost while our method is able to preserve even at that scale.

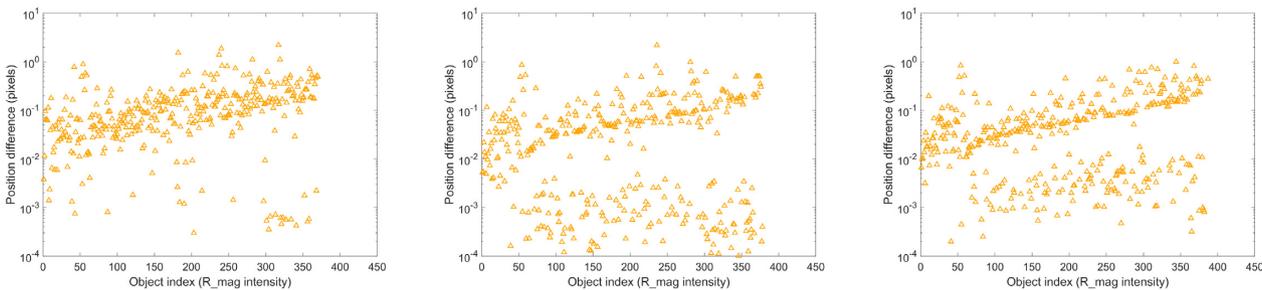
For objects that are detected across all methods, Fig. 14 plots each object’s original  $R_{\text{mag}}$  as a function of the  $R_{\text{mag}}$  difference. Despite having the highest compression ratio overall, SNOWPAC shows to have the least amount of delta  $R_{\text{mag}}$  compared to HCOMPRESS and RICE. There are a notable quantity objects with less error and the overall spread is less. This holds true despite RICE unable to compress beyond a ratio of 6.86, while our method reaches 9.79.



**Figure 15.** Object angle difference. From left to right: HCOMPRESS, RICE, and SNOWPAC on LSST data. Our method is able to achieve higher amounts of clustering towards the bottom portion of the angle-difference plot. SNOWPAC and RICE have similar amounts of error, despite having varying greatly in compression ratios.



**Figure 16.** Object area difference. From left to right: HCOMPRESS, RICE, and SNOWPAC on LSST data. Sorted by the length of the major axis of objects, larger objects' areas are affected the most by lossy compression. Even while having the highest compression ratio, our method is able to preserve objects with higher precision having less impact on the object area computation.



**Figure 17.** Object position difference. From left to right: HCOMPRESS, RICE, and SNOWPAC on LSST data. Sorted by the value of R\_mag, most objects for all methods have a centre computation error of at most 1.5 pixels. The majority of centres are within subpixel accuracy, with both SNOWPAC and RICE having the lowest. RICE has lower observable error only because it is unable to compress above a ratio of 6.86, compared to SNOWPAC's 9.79.

Additional insight can be gathered by analysing the objects and their ellipticity. Fig. 15 compares the quality of the major axis when computing the ellipticity of detected objects. In general, lossy compression of all types has little impact on the direction of the major axis signified by the clustering on the lower left. Nevertheless, the direction of the major axis is preserved very well for our method and RICE. SNOWPAC achieves similar results as RICE, despite having a larger compression ratio (9.79 versus 6.86). Additionally, when comparing SNOWPAC and HCOMPRESS with their similar compression ratios, SNOWPAC has tighter clustering for small object errors signifying a greater preservation of object properties.

Likewise, Fig. 16 compares the area of each detected object relative to the original data set. The index of objects is sorted by the length of the major axis, objects to the right of the figure are the largest in size. While achieving the highest compression ratio, our

method is also able to preserve the physical properties of objects the best by having the lowest error in area computations.

Fig. 17 analyses the effects of lossy compression on deriving the centre position of objects. The index of objects is sorted by their respective R\_mag value, meaning left-to-right denotes fainter to brighter objects. In general, medium- to high-R\_mag objects are affected the most by compression and the centre-pixel computation of objects. At most, objects are 1.5 pixels off-centre and the majority lie within subpixel accuracy. Similar to the DSL data, the LSST data show an observable clustering behaviour of two clusters of objects with error. Comparing HCOMPRESS and SNOWPAC, though both achieve similar compression ratios, SNOWPAC has more objects tightly clustered towards the upper regions and more numerous amounts of objects towards the bottom region. Although RICE is able to achieve more objects with lower error as seen by the clustering of object errors towards the lower regions, our method is able to

nearly match RICE's detection performance with over an additional +3.0 compression ratio in savings.

### 3.5 Discussion

The selection of both types of data was made to test the effects of high- and low-density objects with lossy compression. While one data set has a wider view of the sky with many more astronomical objects, the other has a narrower, better well-defined series of objects at a lower count. The selection of these different ranges of data shows the flexibility of our compressor and its ability to preserve features at different density scenarios. Less object-dense data sets such as the LSST are able to achieve significantly higher ratios and accuracy compared to more dense data sets, highlighting one of the strengths of our compressor. Across several bands, we found that general compression performance was similar, which is expected, since it will adapt its top wavelet components to the point spread function (PSF) of the specific image being compressed.

We have selected input parameters for all methods to achieve extreme, medium, and low compression scenarios, shown in Figs 6 and 12. As shown in these sections, conservative compression ratios of about 5.0 on sparse-object data sets can achieve nearly an accuracy ratio of 0.99, and a compression ratio of about 10:1 an accuracy ratio of 0.975. The scaling of accuracy versus compression ratio makes our method the ideal choice for the general purpose compression of astronomical images. File sizes can be significantly reduced, much lower than standard GZIP by simply using higher quantization and the native wavelet hierarchy structure. As observed with DLS, a massive reduction can be made up to 25:1 compression ratio and still preserve 0.93 accuracy ratio of detected objects.

The compute trade-off for using BZIP2 versus LZ4 for encoding quantized coefficients may be significant in lower performing systems. BZIP2, while extremely efficient in compression capabilities, is several magnitudes slower than LZ4. It is in this case that we would recommend using our method with LZ4 for high-performance applications. The general case, when file size and compression ratio are the priority then BZIP2 would suffice for general lossy compression of astronomical images.

### 4 CONCLUSIONS

Future sky survey telescopes will generate truly massive data sets, creating challenges for data representation, storage, transfer, and analysis. Lossy data reduction methods, exceeding the compression capabilities of lossless compression methods, are necessary to reduce data sizes significantly. Lossy data compression has become more acceptable in many domains, and it is therefore crucial to understand what is lost when utilizing lossy methods. We have introduced and characterized a cubic B-spline wavelet-based lossy compression method that achieves high levels of data size reduction without significant loss of astronomically relevant objects.

Using qualitative and quantitative analysis, we have shown that SNOWPAC can conservatively preserve up to 98 per cent of astronomical objects while achieving a 10:1 compression ratio, achieving better performance than both RICE and HCOMPRESS. This level of reduction is critical to support image-intensive data processing pipelines and minimize bottlenecks in integrated systems and network-based collaborative research. The method exhibits high fidelity with respect to original magnitudes and object shapes and orientations, which are critical for weak-lensing applications. The

implementation of our method is freely available (Pulido 2019), in C++ and MATLAB.

The current implementation uses fixed configurations of weights needed for quantization. This current implementation has made little effort to optimize the weights used for quantization, increasing the possibility of higher compression ratios in future versions. It is possible to make the software more flexible by having it explore image data for signs of structures of different scales to automatically adjust the weights, thereby permitting the adaptation of threshold values to the specific signals in an image.

Wavelet methods have also been studied for the purposes of cataloguing and object measuring performed directly on compressed data, since relevant object features are often captured in encoded form in specific scales of a wavelet representation. Taking advantage of this fact should lead to significant data processing acceleration, compared to data processing methods that can be applied only to decompressed images.

### ACKNOWLEDGEMENTS

JP was supported in part by Los Alamos National Laboratory (LANL) and would like to thank them. LANL is operated by Triad National Security, LLC for the US Department of Energy NNSA under contract no. 89233218NCA000001. CZ was supported in part by the Fund of the Jilin Provincial Science and Technology Department under Grant Nos 20190201305JC, 20180520215JH, 20180201089GX, and 20170204018GX.

### REFERENCES

- Bertin E., Arnouts S., 1996, *A&AS*, 117, 393  
 Cohen A., Daubechies I., Feauveau J. C., 1992, *Commun. Pure Appl. Math.*, 45, 485  
 Collet Y., 2011, LZ4 – Extremely Fast Compression, Available at: <http://lz4.github.io/lz4/>  
 Daubechies I., 1992, *Ten Lectures on Wavelets*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA  
 Folk M., Heber G., Koziol Q., Pourmal E., Robinson D., 2011, in AD '11: Proceedings of the EDBT/ICDT 2011 Workshop on Array Databases. ACM, New York, NY, p. 36  
 Fritze K., Lange M., Mostl G., Oleak H., Richter G. M., 1977, *Astron. Nachr.*, 298, 189  
 GNU, 1997, GZIP, Available at: <https://www.gnu.org/software/gzip/>  
 Kitaeff V. V., Cannon A., Wicenc A., Taubman D., 2015, *Astron. Comput.*, 12, 229  
 Kolev V., Tsvetkova K., Tsvetkov M., 2012, *Publ. Astron. Soc. 'Rudjer Boskovic'*, 11, 187  
 LSST Science Collaborations., 2009, preprint ([arXiv:0912.0201](https://arxiv.org/abs/0912.0201))  
 Masui K. et al., 2015, *Astron. Comput.*, 12, 181  
 Morii M., Ikeda S., Sako S., Ohsawa R., 2017, *ApJ*, 835, 1  
 Pence W., 1999, in Mehringer D. M., Plante R. L., Roberts D. A., eds, ASP Conf. Ser. Vol. 172, *Astronomical Data Analysis Software and Systems VIII*. Astron. Soc. Pac., San Francisco, p. 487  
 Pence W., 2009, Feasibility Study of using BZIP2 within the FITS Tiled Image Compression Convention, Available at: <https://heasarc.gsfc.nasa.gov/fitsio/fpack/bzip2report.pdf>,  
 Pence W., Seaman R., White R., 2011, preprint ([arXiv:1112.2671](https://arxiv.org/abs/1112.2671))  
 Peters S. M., Kitaeff V. V., 2014, *Astron. Comput.*, 6, 41  
 Price D. C., Barsdell B. R., Greenhill L. J., 2014, in Taylor A. R., Rosolowsky E., eds, ASP Conf. Ser. Vol. 495, *Astronomical Data Analysis Software and Systems XXIV (ADASS XXIV)*. Astron. Soc. Pac., San Francisco, p. 531  
 Pulido J., 2019, SNWPAC: SPliNe Wavelet Packing and Compression, Available at: <https://github.com/lanl/VizAly-SNWPAC/>

- Pulido J., Livescu D., Woodring J., Ahrens J., Hamann B., 2016, *Comput. Fluids*, 125, 39
- Rice R., Yeh P., Miller W. H., 1993, in *Proceedings of the 9th AIAA Computing in Aerospace*. AIAA, Reston, VA (American Institute of Aeronautics and Astronautics (AIAA-93-4541-CP))
- Shensa M. J., 1992, *IEEE Trans. Signal Processing*, 40, 2464
- Steward J., 1996, Bzip2 File Compression, <http://www.bzip.org>
- Vohl D., Fluke C. J., Vernardos G., 2015, *Astron. Comput.*, 12, 200
- Vohl D., Pritchard T., Andreoni I., Cooke J., Meade B., 2017, *Publ. Astron. Soc. Aust.*, 34, e038
- White R. L., Postman M., Lattanzi M. G., 1992, in *MacGillivray H. T., Thompson E. B., eds, Digitised Optical Sky Surveys*. Kluwer, Dordrecht, p. 167
- Wittman D. M. et al., 2002, in *Tyson J. A., Wolff S., eds, Proc. SPIE Vol. 4836, Survey and Other Telescope Technologies and Discoveries*. SPIE, Bellingham, p. 73
- Zheng C., Pulido J., Thorman P., Hamann B., 2015, *MNRAS*, 451, 4445

This paper has been typeset from a  $\text{\TeX/L\AA\TeX}$  file prepared by the author.