

Compression techniques with a focus on astronomical data: An overview

Heinrich Strauss
Department of Computer Science
University of Cape Town
hstrauss@cs.uct.ac.za

ABSTRACT

Radio-Astronomical Observation data in its raw format is comprised of noisy data tiles which are cleaned before being subjected to manual and/or automated analysis. We consider in this synthesis the methods employed which provide good compression rates while maintaining data integrity for the observed data and noise and their applicability in the “General-Purpose on GPU” programming paradigm.

1. INTRODUCTION

The amount of astronomical data captured has increased quite dramatically in the last 30 years: early sensors would capture images at 512×512 with a 16-bit value per observed pixel. More recently, it is not uncommon to find image plates which have been captured at 16384×16384 with 16-bit, or even 32-bit, values per image plate per radio source.

Much of the captured data is noise from various sources: atmospheric noise, noise from the sensing Charged-Coupling Devices (CCDs), radio interference from other transmitters in observed spectra, interference on the transmitting channel to the central storage, among others. While much of this is discarded during the data-cleanup phase prior to processing, data custodians are often uncomfortable with completely removing these data-points, since it is time-dependent observational data which is not obtainable again and may be employed to explain researched phenomena[4].

No compression algorithm can work equally well on every data set[17]. We investigate possible methods of compressing the raw data from the receivers, with a strong focus on lossless compression methods, so that data can be archived as raw data in the event they are needed again. If this can be transformed into a massively parallelisable form, General-Purpose Computing on a GPU (GPGPU) methods may be of value in speeding up the handling of these data. This has already implemented for signal convolution, as shown by Harris, et al.[7] and for detection by Resnick, et al.[16]

1.1 Flexible Image Transport System (FITS)

Flexible Image Transport System (FITS) is an open standard format for interchange of scientific data[19]. Since 1981, it has evolved to accommodate the varying types of data in the scientific fields and is still used in the astronomical data field as the de-facto standard for information interchange.

The data format depends on a number of header “cards” interleaved in the data blocks which describe the format of the data. These are limited to 80 7-bit ASCII characters for each card and there are 36 cards per data block. The individual FITS files are described by the (implementation-based) metadata encoded in these and contain one or more image plates. This allows for arbitrary data to be encoded into the FITS files and compresses well under standard data-file compression implementations. This is possible since the cards can be used to determine the data-type stored and therefore make predictions on the location of lossless compressible data.

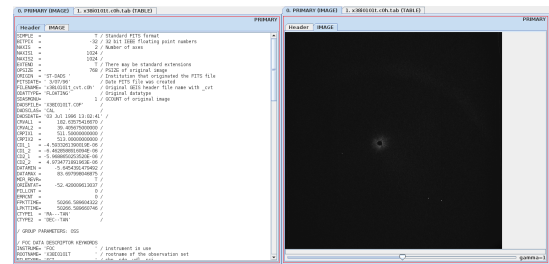


Figure 1: A FITS header file and the output from the PRIMARY table.

By experimentation over the data-set provided by the Sloan Digital Sky Survey (SDSS), the ratio is about 1:2 by experimentation on the SDSS dataset of 2009.

The compressibility stems from the adjacent points, which are often close to each other (and to the background mean) in value, with noise or reference objects (e.g. galaxies, stars) sparsely-interleaved after noise-suppression within the original data. By decreasing the entropy in the data from the raw data we obtain much smaller file sizes after compression.

1.2 Hierarchical Data Format 5 (HDF5)

The HDF5 format is a frequently used data storage format, evolved from the HDF4 format[5]. It supports features such as unlimited number of records, split files and parallel-IO (through MPI-IO) which makes it more applicable to scientific research data-storage.

2. NOISE

The noise sources listed in the introduction produce varying types of Gaussian and Poisson-noise, among others, which are distinguishable by the distribution of the noise. Gaussian Noise is considered random enough to preclude compression[18]. It is therefore prudent to exclude as much Radio-Frequency Interference (RFI) and “stable noise” (such as that from the sensors) to improve the compressibility of the data.

Most commonly, 16-bit integers or single-precision floating point numbers are used to describe pixels on each 2D image plate captured[13]. If floating point numbers are used, the accuracy cannot be ensured, so lossy algorithms (and their implicitly higher compression ratios) become quite attractive.

3. COMPRESSION METHODS OVERVIEW

3.1 Wavelet Transforms

For at least the decade past, Wavelet Transforms have been commonly used to encode the astronomical data in such a way that compression is feasible[18]. The most commonly observed transform in the literature is the 2-D Haar-Transform (H-Transform). This standard processing logic for compression (Hcompress) as suggested by White and Percival, which could be used in either lossless or lossy mode depending on whether or not integer arithmetic is used[20], is as follows:

1. Manipulate the image to obtain roughly equal noise per pixel
2. Apply the H-Transform
3. Encode the output using quadtrees and optionally compress the result (which should be quite amenable to standard compression techniques).

This results in simple arithmetic calculations, provided the size of the image being viewed is square with size an integral power of 2. White and Percival estimate that the H-Transform will require about $16N^2/3$ operations for an $N \times N$ square[20]. In the case where the image is not a square, it can be trivially extended along the non-conforming axes to yield a square image for the purposes of applying the transform. All stored data can, therefore, considered to be stored in a square matrix format, without loss of generality.

The output of the H-Transform is strongly biased towards zeroes (wherever the background is almost consistent, such as noise-free astronomy images), so RLE should efficiently compress the data.

3.2 Haar Transform

The Haar Transform, or H-Transform[20], is given by the wavelet function:

$$\psi(t) = \begin{cases} 1 & 0 \leq t < 1/2, \\ -1 & 1/2 \leq t < 1, \\ 0 & \text{otherwise.} \end{cases}$$

This is run recursively on a dataset and yields a form of the data which is completely invertible, given integer arithmetic. An easy interpretation of this transform-space is whether or not a given square contains values which are brighter than their background. This is used to automatically find points of interest (such as stars,

galaxies and clouds) to further investigate. The transformed data can then be run through a quad-tree encoding so that the image may be retrieved progressively. This allows a quick overview of the individual image plates to ascertain whether interesting points are present. In the event that nothing of consequence is identified on the image plate, the data retrieval can be aborted in favour of another image where less transmission bandwidth will be used. [20].

It should not be inferred that the H-transform is the best transform for any particular image. It is chosen only to exemplify a Wavelet Transform and the importance of selecting a good transform should not be underestimated for any particular implementation.

3.3 Quad-Tree Coding

Quad-Tree-based implementations are often suggested as the de facto way of encoding astronomical data[13][10].

For a $2^n \times 2^n$ square, we can view and encode the internal data recursively as 4 equally-sized squares of size $2^{n-1} \times 2^{n-1}$. The data can then be traversed as follows:

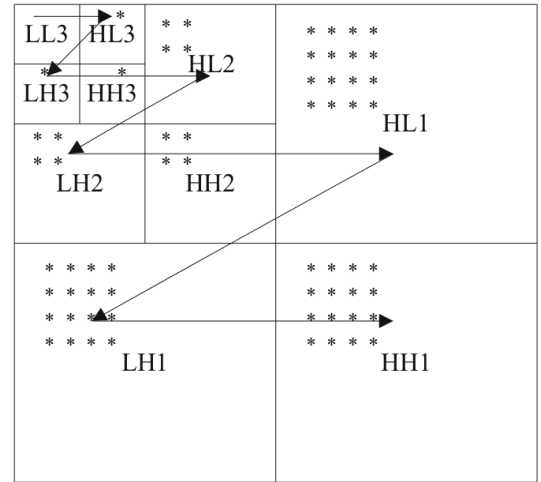


Figure 2: Traversing a Quad-Tree[4]

If we have a region that is uniform in colour (or value), we see that the Quad-Tree traversal will encode them as successive bits of equal value. This allows for easier Run-Length Encoding (RLE). The base-case for a single pixel is stored as the value of the relevant pixel.

In analogue to progressive JPEG encoding[3], this allows individual detail levels to be displayed and the image will become clearer at the scale of each successively smaller Quad-Tree.

4. DISCUSSION

4.1 Du and Ye’s Approach

Du and Ye[4] showed that with two transforms, Integer Wavelet Transform (IWT) and Embedded Zero-tree Wavelet (EZW), and a Quad-Tree representation, we can achieve a more compressible data-set with their algorithm IWT+EZW+TS+A. This yields an increased compression factor of 1.46 over GZIP and 1.20 over the 5/3 transform[3], which was ratified in the JPEG2000 lossless standard, over the LAMOST dataset.

Table 1 The compression ratios of LAMOST simulation data

| 250 fibers | Magnitudes ≈ 17 | | | | | | Magnitudes ≈ 20 | | | | | |
|--------------|-------------------------|-------|-------|---------------|-------|-------|-------------------------|-------|-------|---------------|-------|-------|
| | 360 nm–580 nm | | | 570 nm–890 nm | | | 360 nm–580 nm | | | 570 nm–890 nm | | |
| | E1 | E2 | E3 | E1 | E2 | E3 | E1 | E2 | E3 | E1 | E2 | E3 |
| IWT+EZW+TS+A | 2.760 | 2.760 | 2.777 | 2.677 | 2.650 | 2.649 | 2.851 | 2.850 | 2.850 | 2.691 | 2.684 | 2.684 |
| GZIP | 1.886 | 1.886 | 1.920 | 1.737 | 1.679 | 1.679 | 2.175 | 2.076 | 2.076 | 1.800 | 1.768 | 1.768 |
| JPEG2000 | 2.367 | 2.367 | 2.415 | 2.228 | 2.130 | 2.130 | 2.329 | 2.326 | 2.326 | 2.290 | 2.267 | 2.267 |

Table 2 The compression ratios of SDSS observation data

| Exposures of the first 320 fibers | 370 nm–590 nm | | | | | 580 nm–920 nm | | | | |
|-----------------------------------|---------------|-------|-------|-------|-------|---------------|-------|-------|-------|-------|
| | E1 | E2 | E3 | E4 | E5 | E1 | E2 | E3 | E4 | E5 |
| IWT+EZW+TS+A | 2.764 | 2.737 | 2.769 | 2.750 | 2.730 | 2.592 | 2.566 | 2.618 | 2.597 | 2.559 |
| GZIP | 1.718 | 1.642 | 1.738 | 1.680 | 1.619 | 1.469 | 1.420 | 1.499 | 1.459 | 1.409 |
| JPEG2000 | 2.124 | 2.021 | 2.160 | 2.077 | 1.989 | 1.726 | 1.716 | 1.740 | 1.735 | 1.715 |

Figure 3: IWT+EZW+TS+A Performance[4]

Since two Wavelet Transforms are performed, it is reasonable to expect that the processing time complexity will increase over the H-transform.

4.2 Pence, Seaman and White

An analysis of the compression methods commonly employed is undertaken by Pence, Seaman and White in [13]. They show that the Rice compression algorithm, implemented from the CFITSIO library[14], yields better compression than GZIP in all tested cases against real-world data gathered from the NOAO Mozaic CCD (16-bit) and the NEWFIRM camera (32-bit). Hcompress yields smaller file sizes than Rice at the expense of much greater processing time.

TABLE 1
COMPRESSION STATISTICS FOR 16-BIT INTEGER IMAGES

| | Rice | Hcompress | Tiled-GZIP | Host-GZIP |
|-------------------------|------|-----------|------------|-----------|
| Compression ratio | 2.11 | 2.18 | 1.53 | 1.64 |
| Relative compression | 1.0 | 2.8 | 5.6 | 2.6 |
| CPU time | | | | |
| Relative uncompression | 1.0 | 3.1 | 1.9 | 0.85 |
| CPU time | | | | |

Figure 4: Compression of 16-bit integer images[13]

TABLE 2
COMPRESSION STATISTICS FOR 32-BIT INTEGER IMAGES

| | Rice | Hcompress | Tiled-GZIP | Host-GZIP |
|-------------------------|------|-----------|------------|-----------|
| Compression ratio | 3.76 | 3.83 | 2.30 | 2.32 |
| Relative compression | 1.0 | 5.2 | 7.8 | 4.7 |
| CPU time | | | | |
| Relative uncompression | 1.0 | 3.4 | 2.2 | 1.3 |
| CPU time | | | | |

Figure 5: Compression of 32-bit integer images[13]

If 32-bit values become more common in image plates, the compression ratios presented become higher, as there are more redundant bits in the data captured near zero, which make up a large proportion of the captured points. The Rice algorithm is the most efficient algorithm by processing time and average compression ratio. The increase in value representation size implicitly doubles the required transmission and storage bandwidth. At sufficiently large-scale, network bandwidth and storage speed will become an issue. Other than aligning the data for improved processing or compressibility, the changes required to current implementations are minimal.

4.3 General Discussion

The methods described above hinge strongly on Wavelet transforms, Quad-Trees and classic compression techniques. The need for lossless compression is also a recurring theme in the literature. We have seen a natural progression of compression algorithms with Richard L. White contributing to two of the three common implementations [20][2][13]. As one would expect, the newer algorithms are customised for the data they are compressing, and so yield better results.

The suggested pattern for compression of the astronomical data is therefore

- Transform the data using one or more Wavelet Transforms
- Encode the data into a Quad-Tree to decrease entropy and allow progressive decomposition
- Apply standard compression techniques to reduce the amount of bandwidth being consumed.

This has the obvious drawback of being computationally intensive to store and retrieve image plates.

The Wavelet transforms and Quad-Tree implementations seem to lend themselves well to parallelising. Compression requires a dictionary to map expanded data to compressed data, which is complicated if the non-entropic data are divided between computational nodes. The simplest strategy is to define large enough blocks of data and have these manipulated independently. This would introduce delay into the data stream, which precludes real-time observation and manipulation of the data source if interesting points are observed during the capture phase, which is an intended use-case in future.

The common input format for astronomy data compression is the FITS format, which allows the data to be stored in a patterned data structure. With the optimisation of the header fields, the compressibility factor can be tweaked, since the generally 16-bit integer or floating-point fields can be transformed for compressibility. This can be extended to any container format, including HDF5, where header information can consume a sizeable amount of storage relative to the file size.

The tendency of the data to approximate the background levels plays a large role in the predictability and repeatability of the data and the improved performance of Rice over other, more commonly implemented algorithms shows that this can be exploited more aggressively by choosing an evolved compression algorithm[13]. Since the observed values are often close to zero, we should take care to ensure that denormalisation of any IEEE-754 floating point values does not affect the accuracy of the data stored. It is also important that denormalisation be consistently used throughout any compression implementation to preserve the lossless nature of the transform.

Finding an optimal Wavelet function is also of critical importance, as can be seen from the work of Du and Ye[4]. By choosing the IWT and EZW transforms, a significant enhancement in the overall compression rate over the H-Transform was noted. Even the JPEG lossless 5/3 transform improved over GZIP by a factor of 1.22. As long as the wavelet is invertible, we are able to perform lossless compression, however the Quad-Tree conversion and compressibility factor would remain as potential hazards to ability

to parallelise. These problems have already been implemented in a parallelisable method, as we now show.

The encoding and alignment of the data allow parallelisable decomposition of the Quad-Tree implementation. This is of importance, as Single-Instruction, Multiple Data methods can be used for the construction of the Quad-Tree. This in turn opens the possibility of a GPGPU implementation for this, a concept which was explored in [9]. Although no concrete examples are shown by Ibarra and Kim, Kelly and Breslow[10] describe this in some detail.

The importance of classic data compression techniques should not be underestimated, as even a modest improvement over uncompressed data allows more data to be transmitted and stored within the same hardware confines. GZIP does not perform well compressing this type of image plates compared to newer methods, such as Rice. This is quite understandable, given that the dictionary would have to account for pixel differences of even 1-bit and GZIP has no means of compressing the differences between values of nearby points[13]. The requirement by the data custodians that the compressed data is accurate precludes lossy algorithms. As long as integer arithmetic is used, the chance that the Wavelet transform is invertible is excellent[1]. Alternately, specifying clear accuracy requirements for the values captured would allow for lossy compression methods to be explored. Adaptive Filtering, which identifies points of interest and weights the distribution of the compressed data towards these points, would be a natural path to investigate.

Franaszek, Robinson and Thomas tested parallel data compression with a shared dictionary in [6]. Execution times were quite comparable to the common serial implementations. If this can be parallelised using GPGPU methods, we hope to see improved performance, though large memory buffer copies may make the task difficult.

Percival and White showed that the data convergence period over a transmission medium is greatly reduced when sending the transformed data instead of the raw image data[15]. If this can be used to reconstruct the data where they are to be manipulated, it would extend the lifetime of the infrastructure, both for storage and transmission.

With regard to the stored data, the HDF5 format allows flexibility of data at the cost of needing to customise implementations to support the storage format. Since the data format is likely to be static for a particular implementation, this can be stored in a more efficient transformed format and extracted into usable data at request, or processed directly from the HDF5 data in strides. This requires in-depth analysis of the data formats, which will be dealt with during further analysis. Hardware implementations already exist for Huffman-coded data, as described in [11]. Naturally, disk-based and network-based I/O are in the critical path for the data processing; these are orders of magnitude slower than the transfer of the data within a computing node. Even presently, transmission at a rate of over 40 gigabits per second is often prohibitively expensive and storage systems are similarly impeded, even in the advent of Solid-State Disks.

Many of the compression techniques discussed do not apply to streaming data, as they require data lookahead in order to build the compression algorithm's dictionary[13]. The work of Du and Ye[4], however, alludes to an approach to gain better compression:

Hunt and Rodríguez[8] suggest fast piecewise linear predictors to compress the streaming data by heuristically selecting from a number of independent compression predictors. Oseret and Timsit[12] suggest Object-Based Compression. Hunt and Rodríguez claim a *lossless* compression ratio of between 1:1.3 and 1:3.2 averaged over five data sets. Oseret and Timsit's method is *selectively lossless*, but claims compression ratios of about 1:200.

What remains to be shown is how much of the noise in the data can be separated from the data and if the static noise, such as sensor noise, can be eliminated before compression with no loss of information. Given the size of the datasets, we also need to ascertain whether GPGPU-based methods have enough memory bandwidth to be able to compress the data at realtime rates to enable interactive manipulation in future.

5. REFERENCES

- [1] ADAMS, M. D., AND KOSSSENTNI, F. Reversible integer-to-integer wavelet transforms for image compression: performance evaluation and analysis. *Image Processing, IEEE Transactions on* 9, 6 (2000), 1010–1024.
- [2] DEUTSCH, L. P. GZIP file format specification version 4.3.
- [3] DILLEN, G., GEORIS, B., LEGAT, J.-D., AND CANTINEAU, O. Combined line-based architecture for the 5-3 and 9-7 wavelet transform of JPEG2000. *Circuits and Systems for Video Technology, IEEE Transactions on* 13, 9 (2003), 944–950.
- [4] DU, B., AND YE, Z. A novel method of lossless compression for 2-D astronomical spectra images. *Experimental Astronomy* 27, 1-2 (2009), 19–26.
- [5] FOLK, M. Introduction to HDF5. *The National Center for Supercomputing Applications* (1998).
- [6] FRANASZEK, P., ROBINSON, J., AND THOMAS, J. Parallel compression with cooperative dictionary construction. In *Data Compression Conference, 1996. DCC'96. Proceedings* (1996), IEEE, pp. 200–209.
- [7] HARRIS, C., HAINES, K., AND STAVELEY-SMITH, L. GPU accelerated radio astronomy signal convolution. *Experimental Astronomy* 22, 1-2 (2008), 129–141.
- [8] HUNT, S., AND RODRÍGUEZ, L. S. Fast piecewise linear predictors for lossless compression of hyperspectral imagery. In *Geoscience and Remote Sensing Symposium, 2004. IGARSS'04. Proceedings. 2004 IEEE International* (2004), vol. 1, IEEE.
- [9] IBARRA, O. H., AND KIM, M. H. Quadtree building algorithms on an SIMD hypercube. In *Parallel Processing Symposium, 1992. Proceedings., Sixth International* (1992), IEEE, pp. 22–27.
- [10] KELLY, M., AND BRESLOW, A. Quadtree Construction on the GPU: A Hybrid CPU-GPU Approach. *Retrieved June 13* (2011).
- [11] LEFURGY, C., BIRD, P., CHEN, I.-C., AND MUDGE, T. Improving code density using compression techniques. In *Microarchitecture, 1997. Proceedings., Thirtieth Annual IEEE/ACM International Symposium on* (1997), IEEE, pp. 194–203.
- [12] OSERET, E., AND TIMSIT, C. Optimization of a lossless object-based compression embedded on GAIA, a next-generation space telescope. In *Optical Engineering+ Applications* (2007), International Society for Optics and Photonics, pp. 670003–670003.
- [13] PENCE, W., SEAMAN, R., AND WHITE, R. Lossless astronomical image compression and the effects of noise.

Publications of the Astronomical Society of the Pacific 121, 878 (2009), 414–427.

- [14] PENCE, W. D. New image compression capabilities in CFITSIO. In *Astronomical Telescopes and Instrumentation* (2002), International Society for Optics and Photonics, pp. 444–447.
- [15] PERCIVAL, J., AND WHITE, R. Efficient transfer of images over networks. In *Astronomical Data Analysis Software and Systems II* (1993), vol. 52, p. 321.
- [16] RESNICK, G., KUTTEL, M., AND MARAIS, P. GPU Accelerated Source Extraction in Radio Astronomy: A CUDA Implementation. *University of Cape Town, South Africa* (2010).
- [17] SHANNON, C. E., AND WEAVER, W. A mathematical theory of communication, 1948.
- [18] STARCK, J.-L., AND BOBIN, J. Astronomical data analysis and sparsity: from wavelets to compressed sensing. *Proceedings of the IEEE* 98, 6 (2010), 1021–1030.
- [19] WELLS, D., GREISEN, E., AND HARTEN, R. FITS-a flexible image transport system. *Astronomy and Astrophysics Supplement Series* 44 (1981), 363.
- [20] WHITE, R. L., AND PERCIVAL, J. W. Compression and progressive transmission of astronomical images. In *1994 Symposium on Astronomical Telescopes & Instrumentation for the 21st Century* (1994), International Society for Optics and Photonics, pp. 703–713.