



**HAL**  
open science

## A comparison of lossless compression algorithms for altimeter data

Mathieu Thevenin, Stephane Pigoury, Olivier Thomine, Flavien Guillon

► **To cite this version:**

Mathieu Thevenin, Stephane Pigoury, Olivier Thomine, Flavien Guillon. A comparison of lossless compression algorithms for altimeter data. 2022, pp.1094. 10.5194/egusphere-2022-1094 . cea-04945367

**HAL Id: cea-04945367**

**<https://cea.hal.science/cea-04945367v1>**

Submitted on 13 Feb 2025

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License



# A Comparison of Lossless Compression Algorithms for Altimeter Data

Mathieu Thevenin<sup>1</sup>, Stephane Pigoury<sup>2</sup>, Olivier Thomine<sup>3</sup>, and Flavien Guillon<sup>4</sup>

<sup>1</sup>CEA SPEC, Bat 772 F91191 Gif-sur-Yvette - France

<sup>2</sup>Subnet SAS, 21 av. de la belle image, 94440 Marolles-en-Brie - France

<sup>3</sup>Insa Rouen Normandie, 685 Av. de l'Université, 76800 Saint-Étienne-du-Rouvray - France

<sup>4</sup>CNES, Centre Spatial de Toulouse, 18 avenue Edouard Belin, 31401 Toulouse, France

**Correspondence:** Mathieu Thevenin [mathieu.thevenin@cea.fr](mailto:mathieu.thevenin@cea.fr)

**Abstract.** Satellite data transmission is usually limited between hundreds of kilobits-per-second (kb/s) and several megabits-per-second (Mb/s) while the space-to-ground data volume is becoming larger as the resolution of the instruments increases while the bandwidth remains limited, typically. The Surface Water and Ocean Topography (SWOT) altimetry mission is a partnership between the National Aeronautics and Space Administration (NASA) and the Centre National des Études Spatiales (CNES) which uses the innovative KaRin instrument, a  $K_a$  band (35.75 GHz) synthetic aperture radar combined with an interferometer. Its launch is expected for 2022 for oceanographic and hydrological levels measurement and it will generate 7 TeraBytes-per-day, for a lifetime total of 20 PetaBytes. That is why data compression needs to be implemented at both ends of satellite communications. This study compares the compression results obtained with 672 algorithms, mostly based on the Huffman coding approach which constitute the state-of-the-art for scientific data manipulation, including Computational Fluid Dynamics (CFD). We also have incorporated data preprocessing such as shuffle and bitshuffle, and a novel algorithm known as *SL6*.

## 1 Introduction

Satellites for Earth observations is a topic for research and development, which has been historically taken-on by academia and institutions. More recently by private operator are working on this topics. The observation devices are based on different technologies such as visible light collection, radar and hyperspectral imaging, interferometry. These technologies can be combined to be used together. This, in association with an increasing number of satellites in service, has led to an important rise in the amount of data collected that has to be transmitted to the ground, which is then processed and stored in a database for further use (Sudmanns et al. (2020)). Traditionally, satellite data transmission relies on radio-frequencies (Else (1968)). Data may transit on Data Relay Satellites (DRS) like in International Telecommunication Union (2017); Radhakrishnan et al. (2016) or can be directly sent to ground datacenters or to terminals (Fraire et al. (2019)) in unidirectional or bidirectional manners, depending on the application. Despite the most recent advances, the available bandwidth remains limited; indeed, lightweight terminals only have a few hundred of kilobits-per-second ( $kb/s$ ) capacity (WMO (2018)). Of course, more important communication links can reach a few Megabits-per-second ( $Mb/s$ ) (under 3 GHz frequency bands) to a few hundreds of  $Mb/s$  for



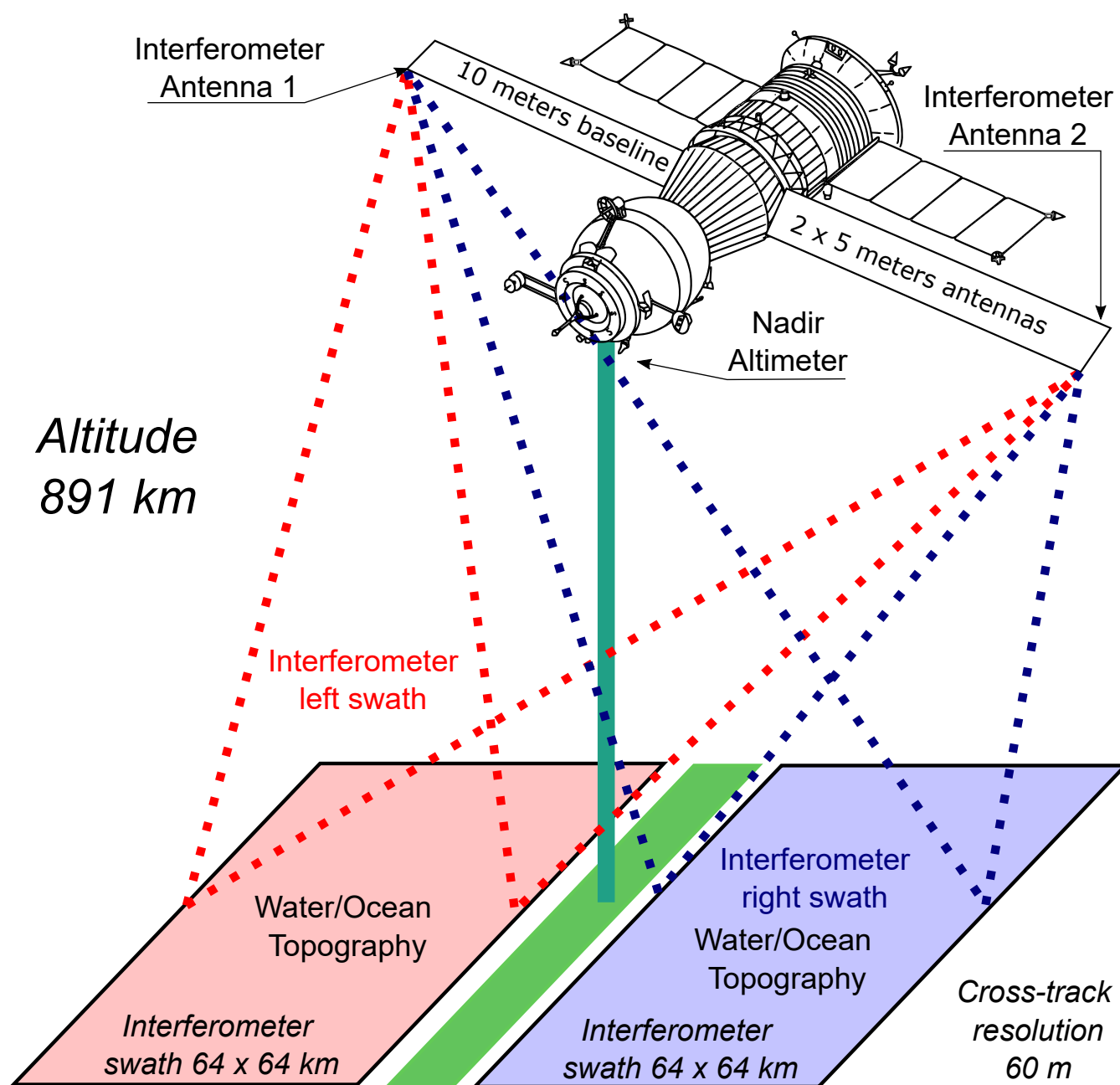
the inter-orbital datalinks (usually over 10 GHz frequency bands). In some cases, the satellite has a limited window of time to  
25 transmit the data to the ground – *i.e.* where it is aligned with the reception antenna. Data compression needs to be used and  
implemented at both ends of the transmission devices (space and ground). The compression scheme is usually application and  
data dependent like in Huang (2011), in observation, the biggest volume of data is transmitted from the space to the ground.

The Surface Water and Ocean Topography (SWOT) (Vaze et al. (2018)) mission is a partnership between the National Aero-  
nautics and Space Administration (NASA) and the Centre National des Études Spatiales (CNES), and continues the long history  
30 of altimetry missions with an innovative instrument known as KaRin (Fjørtoft et al. (2014)), which is a  $K_a$  band (35.75 GHz)  
synthetic aperture radar associated to an interferometer as illustrated by the Figure 1. The SWOT mission launch is foreseen for  
year 2022 and addresses both oceanographic and hydrological communities. It aims at accurately measuring the water level of  
the oceans, the rivers and the lakes. It is expected that the SWOT mission will generate about 20 Petabyte (PB) of data during  
the mission lifetime which corresponds approximately over 7 TB-per-day. Even if the data format is not fully defined yet for  
35 the SWOT mission, the kind of data generated by such missions is usually stored in Hierarchical Data Format 5 (HDF5) (Trott  
et al. (1996)) files. The data volume issue has been addressed by the implementation of different compression schemes (De-  
varajan et al. (2019); Welton et al. (2011)) in the tools for the manipulation of the HDF5 formats of the Computational Fluid  
Dynamics (CFD) General Notation System (CGNS) and benchmarked by previous works like in Delaunay et al. (2019); Di and  
Cappello (2018). The most recent one shows that the combination of shuffle preprocessing and deflate lossless compression  
40 (level 4) provides good results. Moreover, it states that in the case of using a lossy compression, a required precision must be  
defined by the scientists as a Number of Significant Digits (NSD) for each dataset variable. Naturally, when using lossy com-  
pression algorithms, the full precision of the data is not always available, which can be problematic for a scientific use of the  
data. In return of the lost of precision, lossy algorithms generally have a better compression rate than the lossless ones. Since  
scientific data needs to keep the original precision, lossless compression algorithms are preferable when possible. A key point  
45 that need to be evaluated is the consistency of the compression level depending on the nature of the data and the homogeneity  
of compression time and throughput. Indeed, variability in the compression rates leads to non deterministic results when it  
comes to data transfer, especially for a limited time frame.

This paper proposes to investigate different lossless compression algorithms for scientific data, – especially in CFD – that  
are representative of the usual earth observation data. This article extends the work presented in Delaunay et al. (2019) by  
50 exploring the algorithms that are traditionally used in HDF5 and by adding another one to the benchmark.

If the compression level is an important metric, it does not provide any information on the potential benefit of on-the-fly  
compression/decompression during data processing. That is why key points and the use of specific metrics for evaluating the  
data compression algorithms performances will be investigated:

- compression level defined here by  $C_r$ ;
- 55 – compression throughput;
- the homogeneity of the compression regarding the fields of different types (float, integer) and nature;
- the memory usage;



**Figure 1.** Illustration of the SWOT Ka-band Radar interferometer (KaRin) and nadir altimeter measurement concept and the products generated (blue and red shapes on the ground) inspired from Vaze et al. (2018).



- the ability to compress data on-the-fly, in other words, does it needs the storage of the full data in memory;

The major contributions of this paper are: a) the proposal of a thorough bench methodology for lossless compression algo-  
60 rithm for scientific data; b) the proposal of metrics that goes beyond the compression rate; c) a selection of lossless compression  
algorithms suitable for advanced CFD data and d) discussion of different scenarios that can benefit from compression algo-  
rithms in the domain of the earth observation.

This paper is organized in five sections. The first one is this introduction. The second one depicts the methodology followed  
for this research. The third section depicts the results and focuses on the most interesting algorithms. The fourth section  
65 discusses the way how the most efficient algorithms can be used in different environments and infrastructures, including High-  
Performance Computing (HPC) and embedded systems. Finally, the last section provides the conclusions of this research.

## 2 Methodology

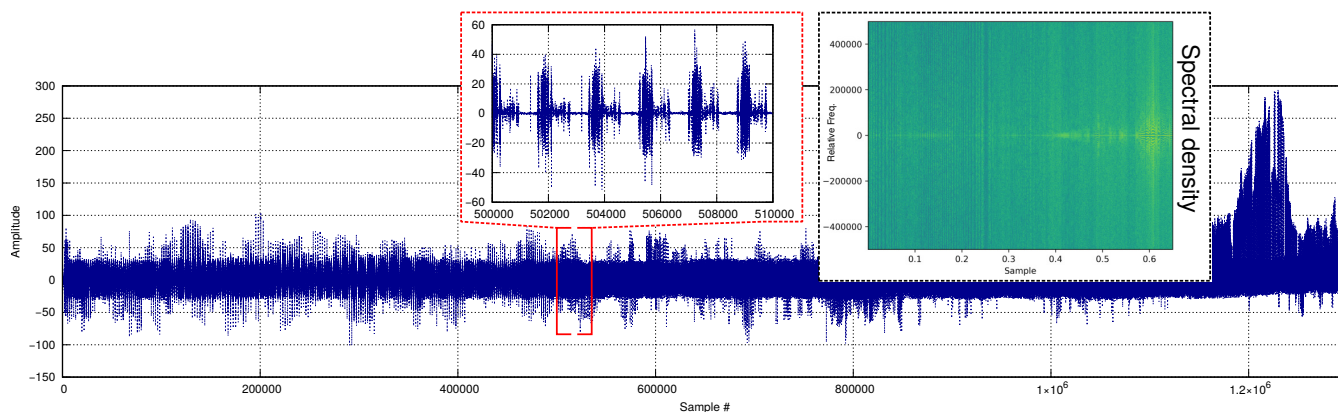
The section describes the methodology developed in this study. It is made of four subsections: first it starts with a description  
of the dataset; followed by a description of the metrics that allow us to choose the compression algorithms that perform the  
70 best. Finally, the testbench and data management approach is explained followed by the results we have obtained.

### 2.1 Dataset

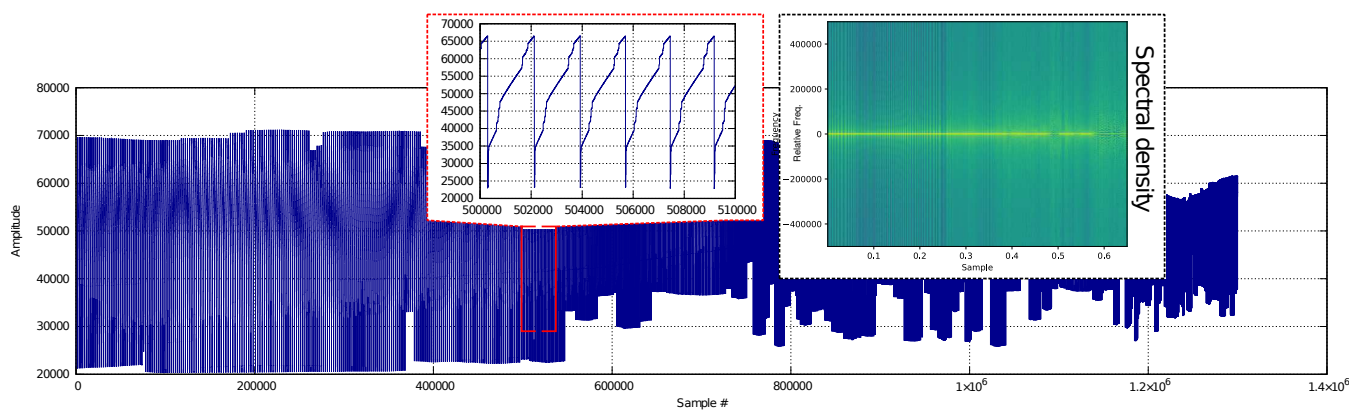
The scientific community relies on CGNS which is a binary unstructured hierarchical format (Diane Poirie and et al. (1998);  
Christopher Rumsey, Bruce Wedan and Poinot (2012)) implementing Advanced Data Format (ADF) (Owen and Daniel (1998))  
and HDF5 (Folk et al. (2011)) to store observational data. It aims at providing a standard for recording and recovering computer  
75 data associated to the numerical solution of the equations of fluid dynamics. It also implements shapes, up to three dimensions:  
0-D point; 1-D line; 2-D triangle and quadrangle; 3-D tetrahedron, pyramid, pentahedron, hexahedron. Due to the data volume,  
not only the transmission and storage are problematic, the access and the read/write time are significant bottlenecks for both  
post-processing and simulations (Soumagne et al. (2010)).

The dataset used for this study corresponds to the SWOT mission products. The SWOT mission will generate two types of  
80 products: the high-resolution products, which are dedicated to the hydrology thematic, and the low-resolution products, which  
are mostly dedicated to the oceanography domain. Basically, L0 data contain raw telemetry; L1 Single Look complex means  
that each pixel encodes its magnitude (I and Q) and therefore contains both amplitude and phase information. Each I and Q  
value is encoded using 16 bits per pixel but stored in 32-bit floating point datawords. The Pixel Cloud product (L2\_HR\_PIXC)  
contains data from the KaRin instrument configured in High Resolution (HR) mode – *i.e.* the HR mask is enabled; it corre-  
85 sponds to the pixels that are detected as being over water. The "Pixel Cloud product" is organized into sub-orbit tiles for each  
swath and each pass, and this is an intermediate product between the L1 Single Look Complex products and the L2 lake/river  
ones. As illustrated in Figure 1, the product granularity is  $64\text{ km} \times 60\text{ km}$  tile.

The dataset used for this study is a simulated golden dataset generated using data obtained from the SWOT mission and  
although, it might not be the most recent, it is based on the one used in Delaunay et al. (2019); for better comparisons. It



**Figure 2.** Illustration of the samples of the Pixel Cloud Height File (64 bits Float), the signal appears to be periodic but with quite chaotic characteristics.

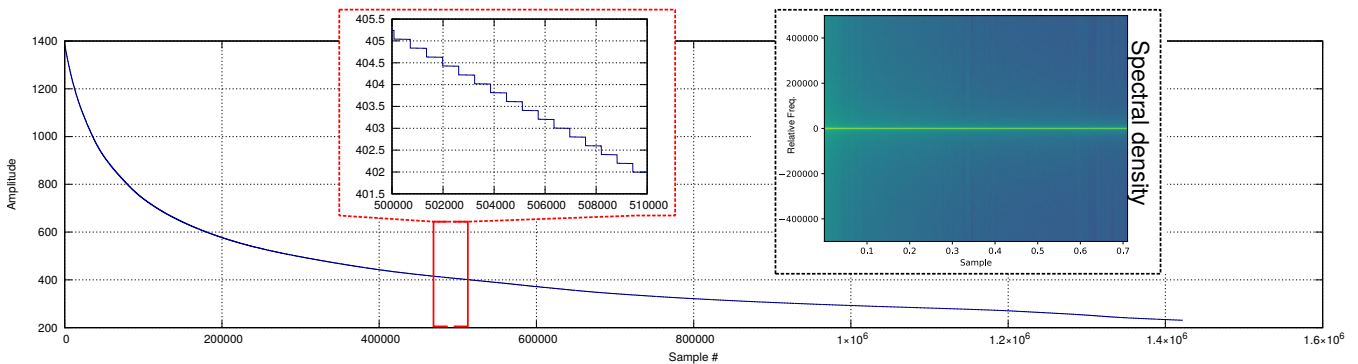


**Figure 3.** Illustration of the samples of the Pixel Cloud L2H CrossTrack File (64 bits Float), the signal seems to be locally continuous and periodic.

90 consists of four HDF5 files. Two of them contain the measurements and the two others contains synthetic generated data divided into fields as follow:

- *signal 1 (s1)*: synthetic data that contains 106 954 752 samples of little-endian (LE) 32-bit floating-point (FP) (427.8 MB), its name is : s1 : signal;
- *signal 3D (s3D)*: synthetic data, that contains 1048576 samples of LE 32-bit FP (4.2 MB), the related field is s3D : signal;

95



**Figure 4.** Illustration of the samples of the `swot pixel_area` file (64 bits Float), the data of this file are continuous, a good compression ratio is expected.

- 12 fields of real world measurement data, that contain 1300111 samples of LE 64-bit FP (10.4 MB), their names are `pixel_cloud:classification,coherent_power,cross_track,dheight_dphase,dlatitude_dphase,dlongitude_dphase,height,illumination_time,incidence_angle,latitude,longitude, pixel_area`;
- 7 fields of experimental data that contain 1421888 samples of LE 64-bit FP (11.4 MB), their names are `SWOT_L2:cross_track,height,illumination_time,latitude,longitude,pixel_area,range_index`;
- 9 fields of experimental data that contains 1300111 samples of LE 32-bit FP (5.2 MB), their name are `pixel_cloud:continuous_classification,num_med_looks,num_rare_looks,phase_noise_std,power_left,power_right,sigma0,x_factor_left,x_factor_right`;
- 1 field of experimental data that contains 2600222 samples in two dimension of LE 32-bit FP (10.4 Mb), its name is `pixel_cloud: ifgram`.

In order to estimate the compression performance on each type of data, every single field had been extracted and compression/decompression was performed on each of them in the testbench. The Figures 2, 3 and 4 illustrates how different the fields are from each other (extracted from the SWOT file). Figure 2 shows a high entropy and low correlation between the samples, while Figure 3 shows a high correlation and quite high entropy in the signal. Finally, Figure 4 shows high correlation and a low entropy signal. Field extraction is performed using the HDFtools and generates a binary file encoded in the original format of the related field.

## 2.2 Lossless Compression Algorithms

Several compression approaches have been proposed through the past decades in the domain of earth observation and scientific data. HDF5 data representation is widely used in the community, thus naturally came the need of compressing this file structure.





115 Within years, CFD tools, that are used to produce the HDF5, have started to implement compression and decompression algorithms. Many different approaches are proposed for data compression in Jayasankar et al. (2021). Today's state-of-the-art works revolve around the use of Huffman, entropy encoders, arithmetic encoders, they are implemented in LZ, LZ4 etc. algorithms. They consist of building a dictionary of redundancies in the signal, then, only the dictionary and the indexes of the words are transmitted. By themselves, these approaches perform well, but they require the analysis of huge portions of the data before  
120 being able to start the compression. Moreover, they are more well-suited for strings than floating-point encoded samples. For this reason, some works focused on using data preconditioners to reorganize the data. Shuffle, for example, analyses the entire file or chunks of data, to reorganize the values to get consecutive similar-valued samples, this way, it permits the encoders compressor to perform well. Among the most recent advances, a bit-level filter was proposed in Masui (2017), initially associated with a lossy compression algorithm (Masui et al. (2015)), it has been demonstrated that it can be efficiently combined  
125 with lossless ones as given in Delaunay et al. (2019). Because of their impact on the compression level when combined with Huffman or entropy encoders; the shuffle and the bitshuffle preconditioners are implemented in the netCFD tools.

Delta coding is a basic approach that consists of encoding the derivative of the signal. This way, all the values are zero-centered, and most of the bits of the encoded samples are set to zero. This approach is lossless on integer encoded, but lossy when used on floating-point encoded data, even if the loss is limited. Recently, a variant inspired by delta-coding was proposed  
130 in the domain of the HPC (Lloyd et al. (2018)) to reduce the required bandwidth for data transmission between nodes. Thanks to a parallel implementation, it sometimes outperforms the other approaches. But again, LZ4fast performs better in many situations. The delta coding approach can also be combined as done in Patauner et al. (2011) where Huffman coding and delta coding are used together. If the traditional delta coding is not really efficient for floating point coded samples, and more generally inferior to the traditional lossless compression schemes, a more recent approach, named here *SL6* (Thomine et al.  
135 (2016, US Patent 20,190,044,532)), that consists in encoding the difference between the signal and an approximated value makes possible to slightly reduce the number of bits used to transmit the signal. A whole compression and decompression framework is available for *SL6* as well in hardware and software under commercial license.

To summarize, the Huffman-based entropy-encoders and arithmetic-encoders based algorithms will be compared to the recent class of algorithms known as *SL6* (Thomine et al. (2016, US Patent 20,190,044,532)). The comparison will be done on  
140 the SWOT dataset which contains integer and floating point number. The shuffle and bitshuffle preconditioners will be used for the LZ family algorithms, as it has been shown in the past that they help them to compress integer and floating-point samples.

### 2.3 The Testbench Process

Since it provides most of today's compression algorithms, the benchmark we developed is derived from the opensource *lzbench* as it was done in Kunkel (2017). Originally, *lzbench* focuses on the LZ77/LZSS/LZMA compression algorithms and performs  
145 in-memory, thus the results are independent from the disk reading and writing times. All the compression algorithms are compiled from sources. This way they all use the same compiler and the same options. Files are compressed and decompressed and compared to the original ones, the time is measured using C primitives. A total of 52 LZ-family algorithms were tested, the LZ family algorithms provided by the *lzbench* are: *blosclz*, *brieflz*, *brotli*, *crush*, *csc*, *density*, *fastlz*, *gipfeli*, *libdeflate*, *lizard*, *LZ*,





150 *lz4, lz4fast, lz4hc, lzf, lzfse, lzg, lzham, lzjb, lzlib, lzma, lzmat, lzo1, lzo1a, lzo1b, lzo1c, lzo1f, lzo1x, lzo1y, lzo1z, lzo2a, lzrw, lzsse2, lzsse4, lzsse8, lzvn, memcpy, pithy, quicklz, lzfse, lzg, lzham, lzjb, lzlib, lzma, lzmat, lzo1, lzo1a, lzo1b, lzo1c, lzo1f, lzo1x, lzo1y, lzo1z, lzo2a, lzrw, lzsse2, lzsse4, lzsse8, lzvn, pithy, quicklz, shrinker, slz\_zlib, snappy, ucl\_nrv2b, ucl\_nrv2d, ucl\_nrv2e, wflz, xpack, xz, yalz77, yappy, zlib, zling and zstd.* The version of these algorithms are the latest available version. For each of them, the most common variants were tested. The variants are usually identified using one or several parameters.

155 The *SL6* algorithm is also benched, again using different parameters, originally designed for the compression of messages between nodes of HPC clusters, its lightweight properties and low computing and memory requirements are properties that make it interesting for the SWOT mission. The parameters are the size of the block and the interpolation mode as explained in Thomine et al. (2016, US Patent 20,190,044,532). Moreover, we added the shuffle and bitshuffle preconditioners. This lead to a total of 672 different versions.

160 The results are compared to the in-memory data copy *memcpy* which provides no compression but maximum throughput – 4.2 GB/s measured

165 The workstation used to execute the testbench is an Intel(R) Xeon(R) CPU E5-2620 v3 running at 2.40 GHz processor, 64 GB of memory. Since this study aims at selecting the algorithms that could be embedded on an onboard computer for a mission, the multicore capacities of the architecture were disabled. Other compression algorithms are excluded from this study. For example the case of *Zarr* as it is designed to be directly used in Python code, and thus cannot not provide a throughput comparable to the others. *Blosc* (Howison (2013)) is a compression schemes that relies on other compression codecs. It implements fast data accesses to exploit the processor cache and Single Instruction Multiple Data (SIMD) instruction (SSE, AltiVec *etc.*) and implements shuffle and bitshuffle filters. The compression itself is performed using external codecs, usually *FastLZ*; thus, potentially, any compression algorithm could be used with *Blosc*. This is not evaluated here since it would be a step forward. Another key point is that the algorithms we are considering here are available as C libraries or source code and thus can be  
170 suitable both for HPC or embedded applications.

## 2.4 Definition of the Metrics

To ensure a fair comparison of the different algorithms, several metrics are considered, including metrics we specifically propose in this paper. First, the ability to compress, expressed in percent, is calculated by  $C_r = \frac{size_{orig} - size_{comp}}{size_{orig}} * 100$  where  $size_{comp}$  refers to the size of the file after compression and  $size_{orig}$  refers to the size of the algorithm before compression.  
175 From this measurement, one can calculate the mean ( $M$ ) and the standard deviation ( $\sigma$ ) for a set of fields, which also allows to derivate the coefficient of variation calculated by  $\widehat{C}_v = \sigma/M$ . Secondly, the compression throughput  $C_{throughput}$ , or compression speed, is calculated using the size of the original data divided by time required by the algorithm to fully perform compression:  $C_{throughput} = size_{orig}/time_{comp}$ . This is the same approach for the decompression with  $D_{throughput} = size_{orig}/time_{decomp}$ . Here again, the corresponding  $C_v$  can be calculated from the standard deviation  $\sigma$ .

180 As one of the targeted application is the on-board compression and data transmission, it is important to put emphasis on the algorithms with consistent performances, whatever the type of data. In other words, an algorithm providing extremely different compression  $C_r$  and/or times, depending on the nature of the data, should get a lower score than an algorithm which



has consistent results among the whole dataset. Intuitively, it means the highest throughput and the highest compression  $C_r$  are preferable but with the lowest standard deviation for both of them. Thus, we propose the metric  $H$ -score for a given set of data (fields) calculated on the normalized results (maximum set to 1) using the Equation 1, with,  $\sigma(C_r)$  the standard deviation of the compression rate,  $T$  the throughput and  $\sigma(T)$  its the standard deviation. A graphic intuitive illustration of the variables used in the  $H$ -score is given in Figure 5. Moreover, we defined  $\alpha$  and  $\beta$  coefficients, which are set to one. They can be used to weight either the compression  $C_r$  ( $\alpha$ ) or the throughput ( $\beta$ ). Finally, the  $\rho$  coefficient, set to 1000, is simply used to ease the reading of the results. A low or negative  $H$ -score means the algorithms results are heterogeneous within the selected data, on the opposite side, a higher  $H$ -score means the algorithms results are homogeneous among the dataset while providing good performances.

$$H - score = \underbrace{\rho}_{1000} \times (\overline{C_r} - \underbrace{\alpha}_1 \cdot \sigma(rate)) \cdot (\overline{T} - \underbrace{\beta}_1 \cdot \sigma(T)) \quad (1)$$

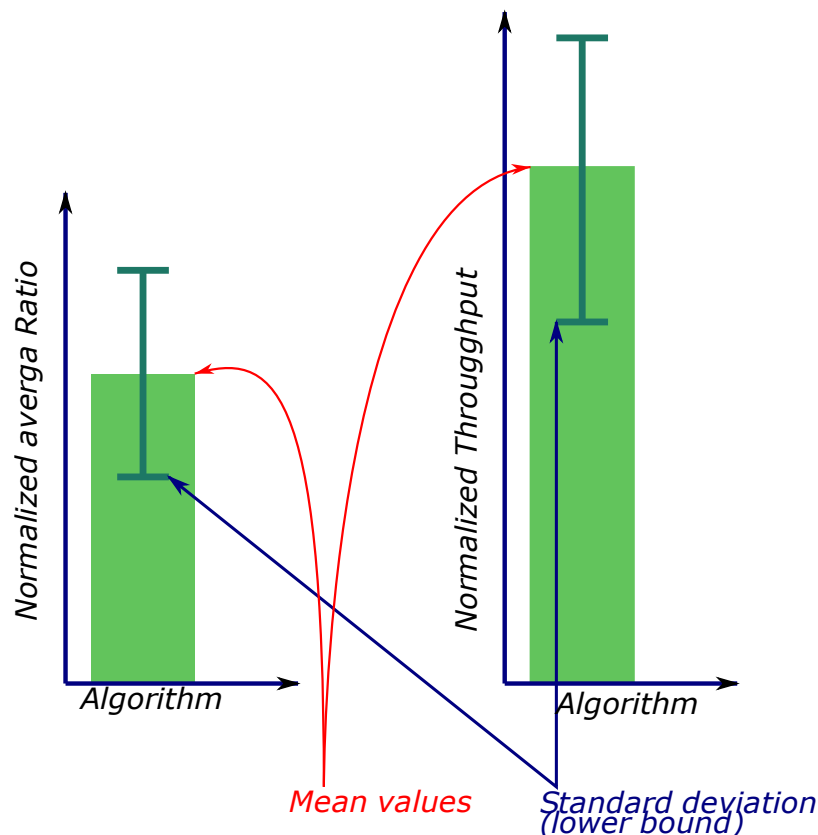
The third metric that was considered in this paper, and proposed by Thomine, is called the  $tt_{tx}$  for transmission-throughput-threshold and is expressed in Bytes-per-second. It is used to determine if it is worth compressing the data when a certain transfer  $C_r$  is available for data transmission. In other words, if  $tt_{tx}$  is higher than the media throughput (or speed), then the time used to compress the data, to transmit and to decompress it, is lower than transmitting the original data; in other words, the higher  $tt_{tx}$  is better. We can derive a similar metric for writing the data to a memory or to a disk. In this case, as decompression does not need to be performed, its time is removed from the calculation; it gives  $tt_{wr}$  which stands for writing-throughput-threshold, also expressed in bits-per-second.  $tt_{tx}$  and  $tt_{wr}$  are given in the Equations 2 and 3,  $Size_{orig}$  refers to the original size of the data,  $Size_{comp}$  the size of the compressed data,  $Time_{comp}$  and  $Time_{decomp}$  the time required to compress and decompress it.

$$tt_{tx} = \frac{Size_{orig} - Size_{comp}}{time_{comp} + Time_{decomp}} \quad (2)$$

$$tt_{wr} = \frac{Size_{orig} - Size_{comp}}{Time_{comp}} \quad (3)$$

## 205 2.5 Data management

The number of algorithms is quite important and for each algorithm, many variation are considered. For example, the impact of the algorithms parameters (more than 17 variants for ZStd), but also, the use of preconditioning filters such as shuffle and bitshuffle. Since the volume of data produced by our extended testbench is quite important we decided to store all the results in a relational database we have specifically designed for the analysis of the results. This way, all the data can be extracted and analyzed using a combination of SQL queries. In the same vein, the derivated metrics such as compression  $C_r$ , compression and decompression throughput, ( $tt_{tx}$  and  $tt_{wr}$  can be calculated on-the-fly directly in the queries. Moreover, to make it easier, we have designed a graphical user interface (Figure 7) to browse through the results, it allows an instant displaying, filtering,



**Figure 5.** Illustration of the variables used to calculate the  $H$ -score; intuitively, the lower-bound of error bars made out of the standard deviation ( $\sigma$ ) is subtracted from the mean compression  $C_r$  and speeds for a given dataset. The result of the subtraction is multiplied. Coefficients ( $\alpha$ ,  $\beta$ ) are added to weight the  $\sigma$  values (set to 1) as given in the Equation 1.

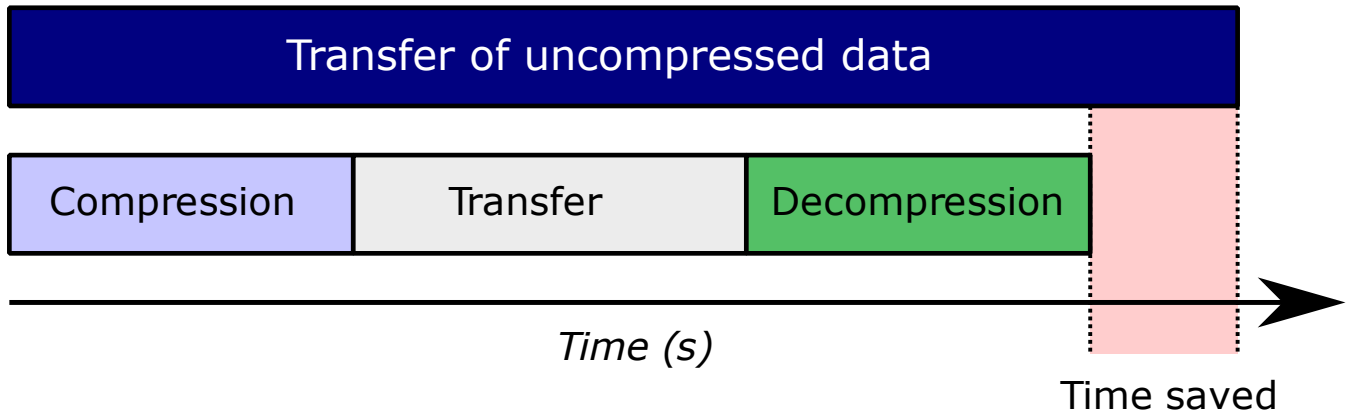
ranking and analysis of the different metrics used. This user interface was designed for internal use only, but can be distributed on demand.

### 215 3 Results obtained using the Testbench

This section presents the results we have obtained with the test-bench. First, a general overview of the results is introduced; then, we focus on the fields that appear to be difficult to compress.

#### 3.1 General Results

As all 672 variants of the algorithms we have benchmarked cannot be presented in a concise way, we choose to focus on the  
220 most interesting of them, based on the results obtained using the metrics defined in the previous section:

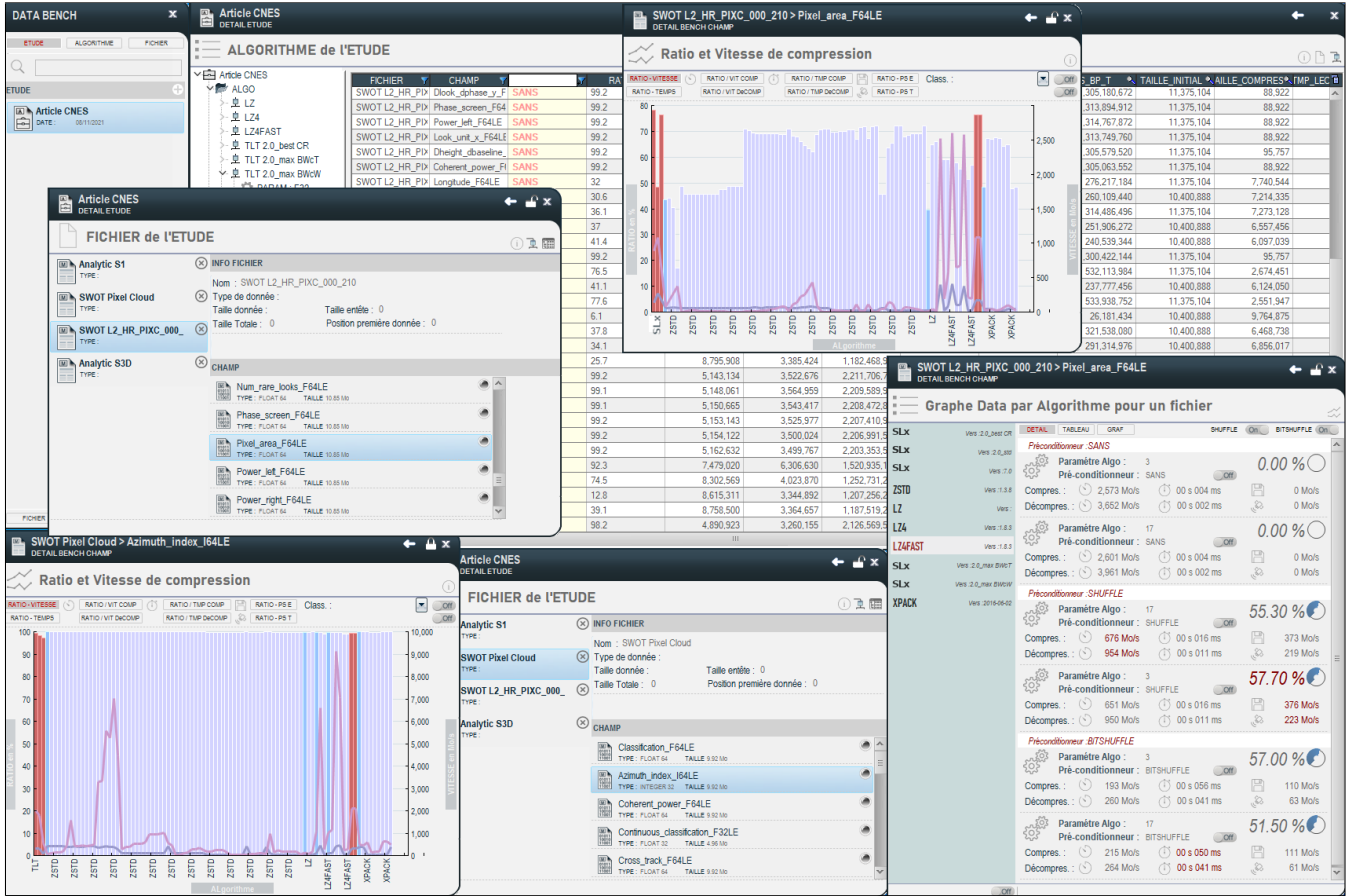


**Figure 6.** illustration of the time saved by compressing, transferring and decompressing; the calculation of the  $tt_{tx}$  is used to emphasize the threshold from which time is saved; when it comes to data writing into a memory or to a disk,  $tt_{wr}$  is used, the decompression time can be omitted since it is supposed to be done at the reading process.

- lz4fast associated to shuffle, parameter 17;
- LZ4fast associated to bitshuffle, parameter 17;
- Zstd parameter 22;
- Zstd associated to shuffle, parameter 1;
- 225 - Zstd associated to bitshuffle, parameter 1;
- LZ associated to shuffle, parameter 9;
- LZ associated to bitshuffle, parameter 9;
- SL6 with a polynomial sample estimator.

The measured time resolution on the workstation devoted to the bench is  $1.7 \mu s \pm 0.4 \mu s$ . Table 1 provides the results  
 230 for the S1 fields of the signal file. It shows a certain variability in the compression  $C_r$  (mean = 63.34 %;  $\sigma = 2.36$ ) and a higher one in the compression and decompression speeds:  $mean = 559263973$ ;  $\sigma = 477136146.9$  for the compression and  $mean = 718261809$ ;  $\sigma = 513374554$  for the decompression. One can see the standard deviation on the compression speed is extremely high and close to the mean value.

The results obtained on the SWOT pixel\_cloud file are quite different in terms of ranking, due to the nature of  
 235 the data which is slightly different. They are given in Table 2 where we can see, again, an homogeneous compression  $C_r$  (mean = 54.21;  $\sigma = 3.65$ ); heterogeneous compression speed (mean = 383153349;  $\sigma = 313197839$ ). This is the same for the decompression speed with (mean = 652996328;  $\sigma = 509295674$ ).



**Figure 7.** screenshot of the graphical user interface designed for the analysis of the data produced by the testbench; a SQL database is used to store the results.

The entropy level associated to the fields is higher and the fields are more complex to compress using the traditional Huffman-based approaches. The LZ4Fast and LZ associated to the shuffle filter obtain good results, but *SL6*, which exploits the correlation between samples scores higher in terms of throughput with a compression  $C_r$  close to the others.

A similar behavior is observed on the SWOT L2HR file, as given in Table 3: the compression  $C_r$  are quite homogeneous  $\sigma = 0.95$  while a great variability appears in the compression and decompression speed with a  $\sigma = 639\,050\,509$ , which is close to the mean value 680465679 for compression speed;  $\sigma = 825\,523\,818$  for a mean value of 849856905 for decompression speed. It means that the throughput highly depends on the algorithms and on the nature of the data.

However, the ranking is different, ZStd with no filter provides the highest throughput, but lowest compression  $C_r$ . Closely followed by *SL6* which has a higher (among the highest) compression  $C_r$  and is faster at decompressing. The other algorithms are slightly slower as the compression  $C_r$  are similar or a few percent lower than the best ones.



**Table 1.** Summary of the results obtained for the S1 signal file, original size 8388608 B.

File	Par.	Preproc.	$C_r$ (%)	comp. time (ns)	Decomp. time (ns)	Comp. speed (Bytes/s)	Decomp. speed (Bytes/s)	Comp. size (Bytes)
ZSTD	1	none	58.50	10477144	5958746	800657889	1407780765	3481550
ZSTD	1	shuffle	65.68	10155354	8260379	826028123	1015523380	2878818
ZSTD	1	bitshuffle	65.34	38472020	30087028	218044387	278811453	2907213
LZ	9	none	58.24	254577532	33740694	32951093	248619901	3503334
LZ	9	shuffle	66.84	132626450	25893111	63249887	323970650	2781908
LZ	9	bitshuffle	65.58	170177748	48532690	49293213	172844489	2887199
LZ4	0	shuffle	62.15	7712785	7685319	1087623731	1091510710	3174678
LZ4	0	bitshuffle	64.91	32678452	29960194	256701511	279991778	2943353
LZ4FAST	3	shuffle	62.10	6445824	7242951	1301401962	1158175445	3178959
LZ4FAST	3	bitshuffle	64.75	32926048	29876013	254771177	280780705	2957242
LZ4FAST	17	shuffle	62.11	6516309	7435122	1287325079	1128240801	3178408
LZ4FAST	17	bitshuffle	64.39	31675455	29930478	264829913	280269764	2987309
SL6	Poly.	none	63.12	10572072	5160822	793468679	1625440288	3093395

**Table 2.** Summary of the results obtained for the SWOT pixel\_cloud signal file, size 208017760 Bytes.

File	Param.	Preproc.	$C_r$ (%)	comp. time (ns)	Decomp. time (ns)	Comp. speed (Bytes/s)	Decomp. speed (Bytes/s)	Comp. size (Bytes)
ZSTD	1	shuffle	49.39	444398317	249599814	468088541	833405108	105280832
ZSTD	1	bitshuffle	47.15	1181907250	756880412	176001763	274835703	109940201
LZ	9	none	39.00	21218526054	845143541	9803591	246133053	126894903
LZ	9	shuffle	50.46	21854440692	654839210	9518329	317662346	103043840
LZ	9	bitshuffle	47.79	14458854352	1299657146	14386877	160055874	108602087
LZ4	0	shuffle	46.00	291760551	205508459	712974250	1012210208	112324430
LZ4	0	bitshuffle	45.97	881074723	730831535	236095480	284631615	112399104
LZ4FAST	3	shuffle	45.77	291303202	201869174	714093627	1030458271	112805305
LZ4FAST	3	bitshuffle	45.76	877821970	729865607	236970328	285008306	112832614
LZ4FAST	17	shuffle	44.83	253447978	202266864	820751310	1028432220	114772156
LZ4FAST	17	bitshuffle	44.92	837072886	729876936	248506150	285003882	114585439
SL6	-	none	42.98	239367624	103952649	869030475	2001081858	118608676

Unsurprisingly, it appears the set of algorithms which shows the best performance are similar to the one presented in Delaunay et al. (2019) even if the method and the test-bench differs from the one presented in this paper. The rest of the article will focus on the SWOT fields as they seem to provide the lowest compression  $C_r$  compared to the others – which compress quite well whatever algorithm is used.

### 3.1.1 Compression $C_r$

First of all, the unweighted average compression  $C_r$  for the fields of the SWOT file is 74.36 %, with a  $\sigma = 37.39$ , the details are given in Figure 8. As LZ4 and LZ4fast are not able to compress 17 fields out of 72, they obtain the worst results with an average compression  $C_r$  of 67.22 and 66.51 ( $\sigma = 45.81$  and  $\sigma = 46.49$ ). The median compression  $C_r$  for the SWOT file fields



**Table 3.** Summary of the results obtained for the SWOT L2HR signal file, original size 459 269 824 Bytes.

File	Param.	Preproc.	$C_r$ (%)	comp. time (ns)	Decomp. time (ns)	Comp. speed (Bytes/s)	Decomp. speed (Bytes/s)	Comp. size (Bytes)
ZSTD	1	none	89.44	245 576 988	187 199 279	1 870 166 369	2 453 373 894	48 500 312
ZSTD	1	shuffle	92.65	504 396 287	475 060 914	910 533 713	966 759 863	33 760 218
ZSTD	1	bitshuffle	92.26	1 833 936 246	1 680 779 517	250 428 457	273 248 109	35 545 956
LZ	9	none	90.51	7 658 339 217	1 170 823 423	59 969 898	392 262 245	43 603 304
LZ	9	shuffle	92.78	6 934 283 461	1 574 703 688	66 231 764	291 654 759	33 159 715
LZ	9	bitshuffle	92.46	8 688 358 240	2 503 710 877	52 860 369	183 435 647	34 640 701
LZ4	0	shuffle	91.88	4 610 883 354	457 298 287	996 056 005	1 004 311 271	37 294 444
LZ4	0	bitshuffle	91.37	1 717 883 548	1 646 306 534	267 346 308	278 969 812	39 618 481
LZ4FAST	3	shuffle	91.78	453 166 592	463 579 997	1 013 467 965	990 702 418	37 753 332
LZ4FAST	3	bitshuffle	91.27	1 715 480 989	1 649 380 612	267 720 731	278 449 874	40 096 510
LZ4FAST	17	shuffle	91.64	451 584 784	457 452 644	1 017 017 934	1 003 972 389	38 406 860
LZ4FAST	17	bitshuffle	90.75	1 690 834 724	1 648 500 436	271 623 132	278 598 546	42 475 508
SL6	-.	none	92.19	254 777 476	173 152 489	1 802 631 187	2 652 400 937	35 870 622

**Table 4.** Average compression  $C_r$ ,  $\sigma$  and the coefficient of variation ( $\widehat{C}_v$ ) obtained on the SWOT file obtained with the most efficient algorithms, ordered by the couple ( $C_r, \sigma$ ) and illustrated in Figure 8

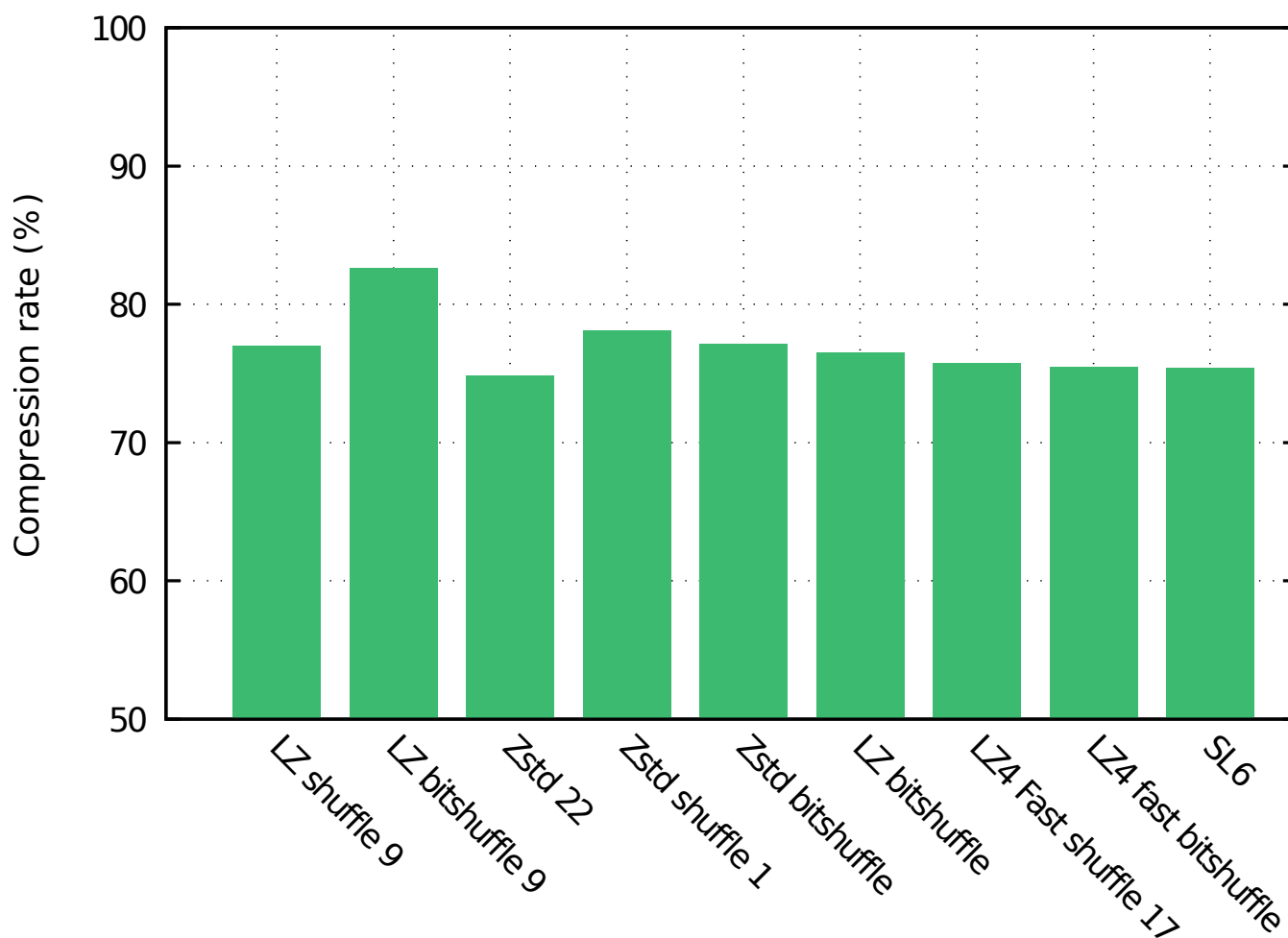
Algorithm	Average $C_r$ (%)	$\sigma$ (%)	$\widehat{C}_v$
LZ shuffle 9	76.98	34.31	0.45
LZ bitshuffle	76.46	35.18	0.46
Zstd 22	74.83	38.05	0.51
Zstd shuffle 1	78.95	32.98	0.42
Zstd bitshuffle 1	77.12	34.06	0.46
LZ4fast shuffle 17	75.71	35.36	0.47
LZ4fast bitshuffle 17	75.42	34.80	0.46
SL6 polynomial	75.33	33.93	0.45

is between 99.1 % (SL6) and 99.9 % (Zstd with shuffle). The other approaches are consistent in term of compression  $C_r$  and standard deviation. The best result is obtained using Zstd associated to shuffle (78.08;  $\sigma = 32.98$ ) and SL6 (75.33;  $\sigma = 33.93$ ). The Figure 9 illustrates the results obtained field by field. It is clear for the reader to see that they perform almost equally in terms of compression  $C_r$ , even if SL6 tends to have a bit more homogeneous results than the others.

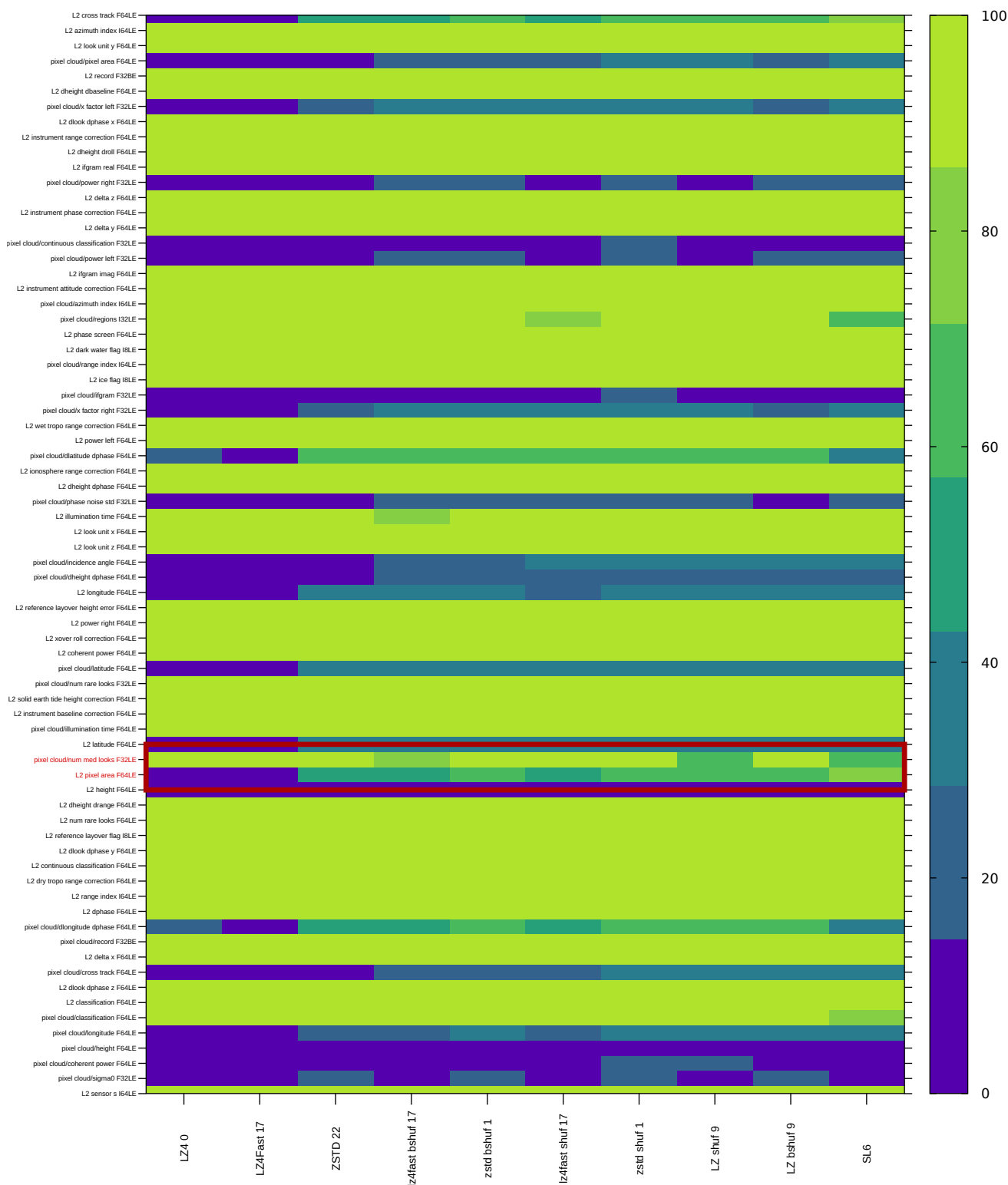
### 260 3.1.2 Compression speed

The average throughput results obtained for the SWOT file are quite heterogeneous from an algorithm to one another. They vary from 1.7 GB/s to only a few dozen MB/s for the slowest, the median equals to 238 MB/s. They are summarized in Table 5 and displayed in Figure 10. It can be seen that the Huffman algorithms associated with shuffle and bitshuffle filters perform quite well. For example, the LZ4Fast (parameter 17) associated with bitshuffle provides the lowest  $\sigma$  but is relatively slow





**Figure 8.** Average compression  $C_r$  (see Table 4 for the values) for the fields comprised in the SWOT file. It is clearly visible that all the algorithms perform well on this dataset, the lowest standard deviations are observed for *SL6* and *Zstd* associated to the shuffle with. The Pithy algorithm provides good results on some fields too, unfortunately it is not able to compress some fields of the SWOT file.



**Figure 9.** Matrix of the compression  $C_r$  of the fields of the SWOT file, it can be seen that the compression  $C_r$  varies in a similar manner whichever algorithm is used, even if some are slightly better; the most interesting fields to look at are the one marked in red as they emphasize the differences between the algorithms.



**Table 5.** Average throughput (TP) ( $MB/s$ ) and associated standard deviation  $\sigma$  obtained on the SWOT file obtained with the most efficient algorithms, ordered by the couple  $(C_r, \sigma)$ , results are plotted in the Figure 10.

Algorithm	Average TP ( $MB/s$ )	$\sigma$	$\widehat{C}_v$
LZ shuffle 9	212.89	249.491.17	
LZ bitshuffle 9	120.24	203.091.69	
Zstd 22	80.75	67.120.83	
Zstd shuffle 1	833.98	344.460.41	
Zstd bitshuffle 1	225.78	54.330.24	
LZ4Fast shuffle 17	965.33	231.470.24	
LZ4Fast bitshuffle 17	254.47	21.590.08	
SL6 polynomial	1431.24	753.440.53	

265 (254.46 MB/s). Indeed, it is outperformed by the same with the shuffle preconditioner. The *SL6* algorithm is 48 % faster than the fastest Huffman-based, even when it is associated to shuffle and bitshuffle filters. The LZ4Fast associated to the bitshuffle has the most homogeneous results for the SWOT dataset.

A closer look at the results obtained for the SWOT `pixel_area`, in Figure 11 field shows more variability in the compression throughput, depending on the algorithms. We choose to display this field in particular since it emphasizes the disparity of the results. Here, the *SL6* algorithms ranks first and second best performance (more than 2.5 GB/s for 48.3 % compression  $C_r$  for the first), the second has a higher compression  $C_r$  (76.5 %) for 1.75 GB/s.

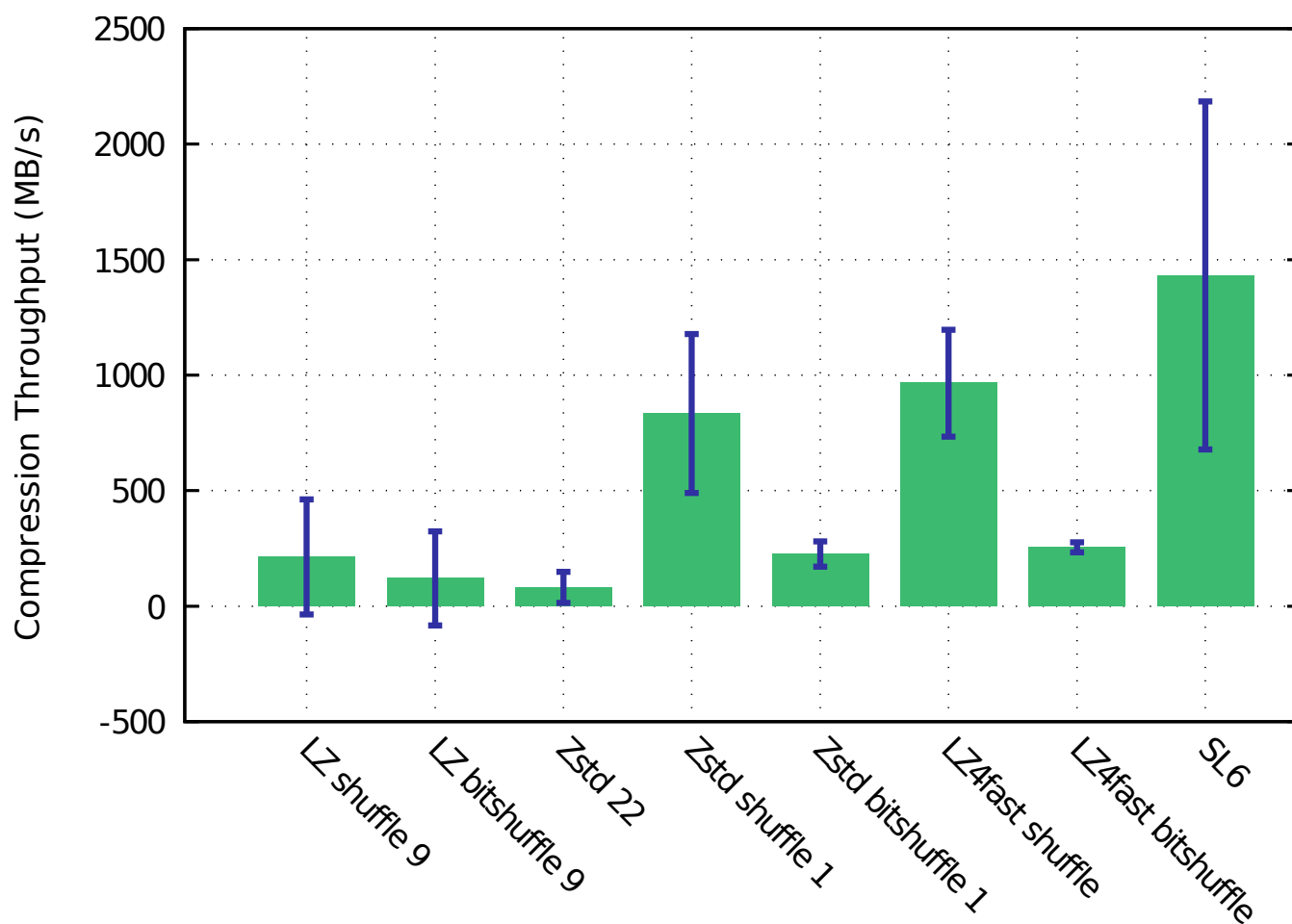
Since it is a stream algorithm, the compression throughput of *SL6* is constant for a given word size, moreless a few percent. The throughput for 8-bit is the lowest 329 MB/s, while the figures on 64-bit data is 1817 MB/s, 32-bit data words lead to 1008 MB/s of throughput. This is easily understandable by the fully deterministic nature of the algorithm and the fact that the algorithm locally compress the data using blocks of predetermined size.

### 3.1.3 H-score

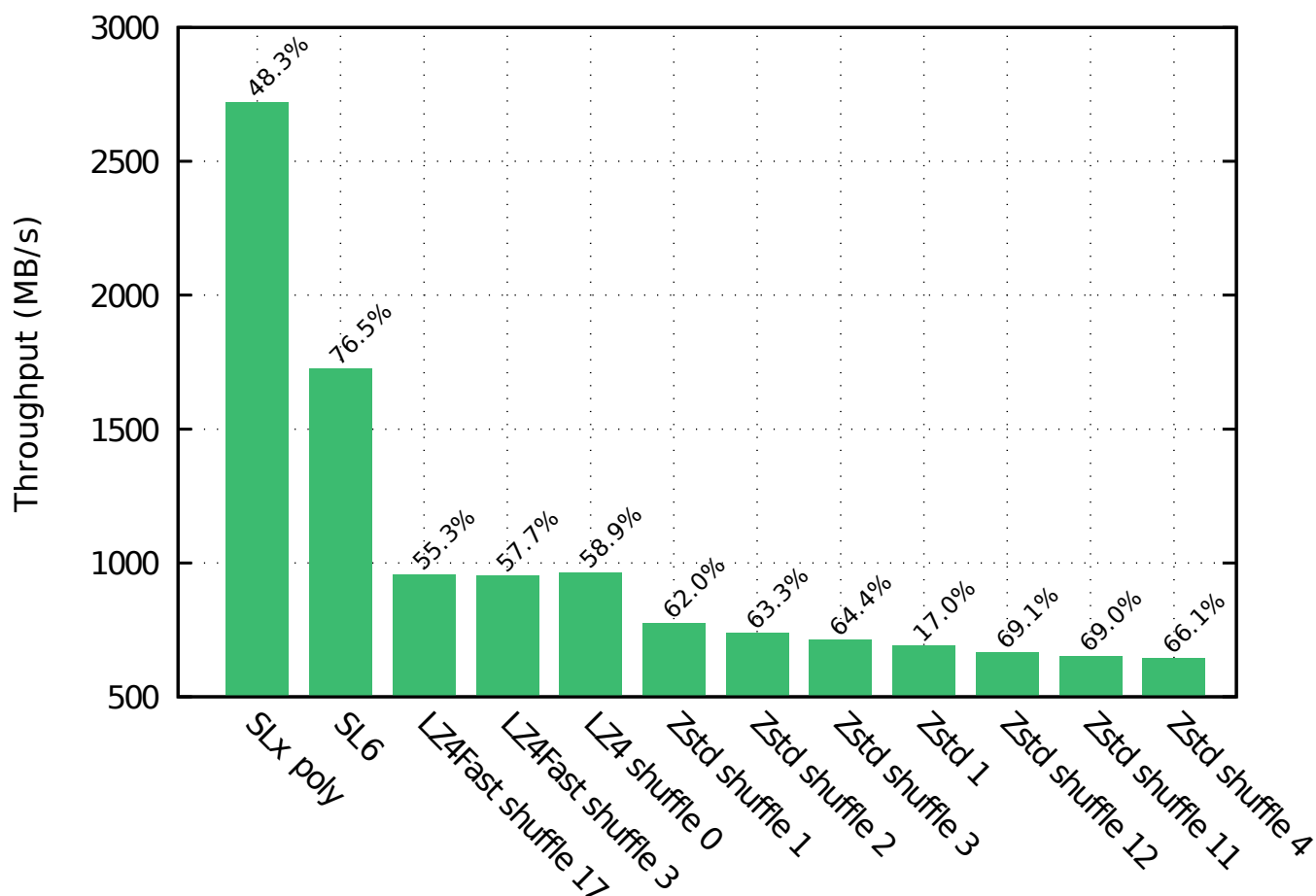
The *H*-score defined in the previous section emphasize the algorithms that provide consistent compression  $C_r$  and consistent throughput. This is extremely important when the dataset can be heterogeneous and when the targeted application is real-time compression on a single-core. Since *SL6* provides quite high throughput with a compression  $C_r$  comparable to the other algorithms – and for both a reasonable  $\sigma$  – it scores higher than the others. It is followed by LZ4Fast associated to the shuffle filter (parameter 17), which scores 2 % lower than the *SL6* algorithm. The results are visible in Table 6 and in Figure 12.

## 3.2 Threshold Throughput

The results for the transmission-threshold-throughput ( $tt_{tx}$ ), that was depicted in subsection 2.4 is considered in this subsection. When it is higher than the media throughput, the time used to compress the data, to transmit and to decompress it is lower than the time required for the transmission of the original data. This metric is highly correlated with the speed of the hardware used



**Figure 10.** Plot of the average compression throughputs for the *swot* fields, the standard deviation is given as error bars, on *SL6* the errors bars are quite important since we mixed words of data of different size (8, 32 and 64 bits). The errors bars for a given datasize (for example 64 bits) is constant with minimal standard deviation.

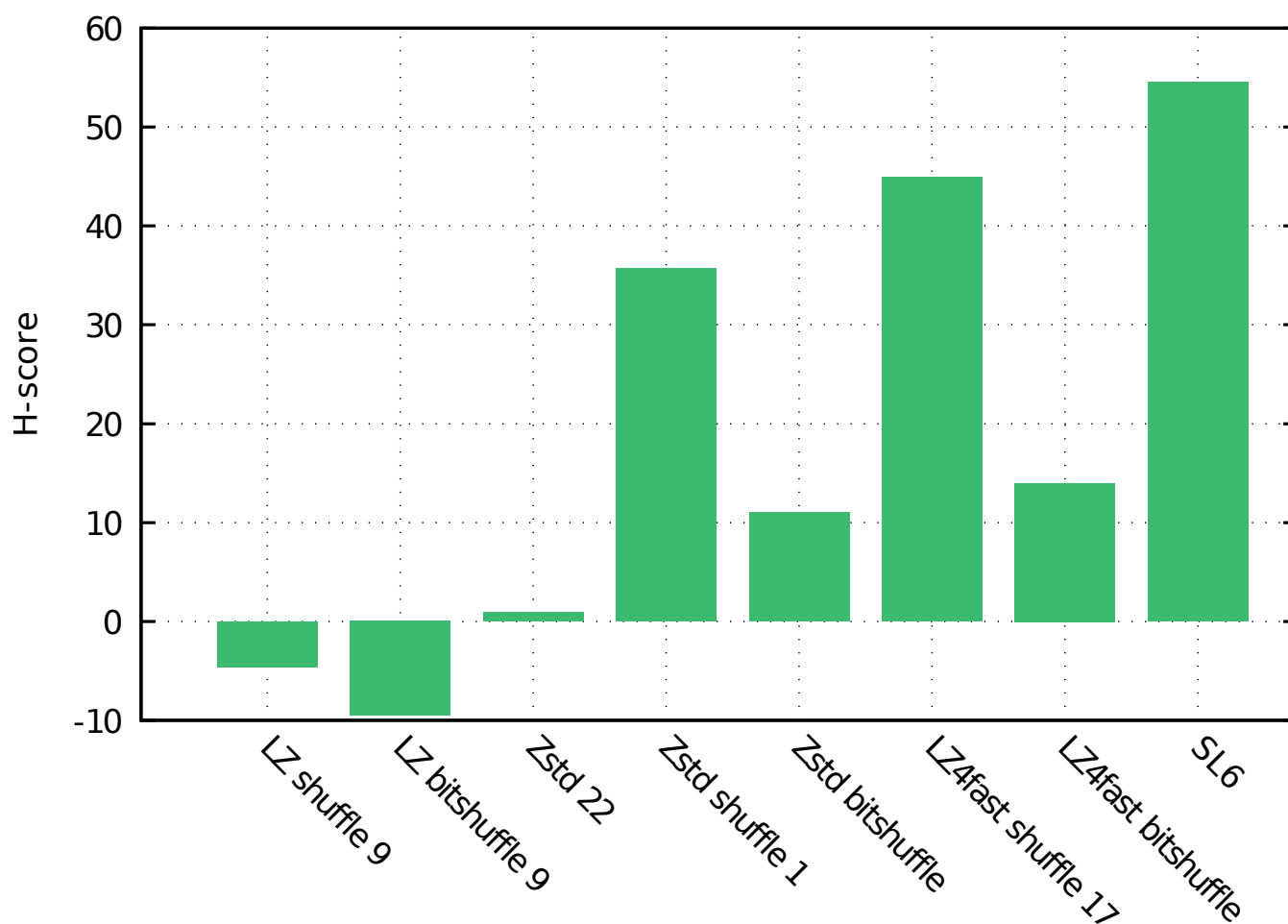


**Figure 11.** Plot of the decompression speeds for the most efficient algorithms for the SWOT `pixel_area` field, the compression  $C_r$  are given as labels on the top of each histogram boxes, some of the algorithms provides good results with different parameters and preprocessing (ZStd and LZ4 Fast); the Pithy algorithm has been excluded since it is not able to compress some other fields of the file. SLx was tested with different blocks sizes.

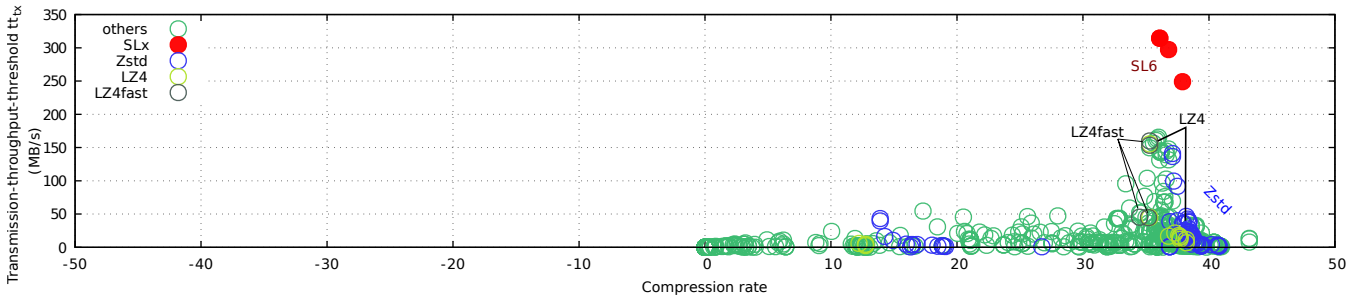


**Table 6.** *H*-score calculated on all the fields of the SWOT file. Since *SL6* has the highest throughput and is among the most homogeneous results, it outperforms the others approaches.

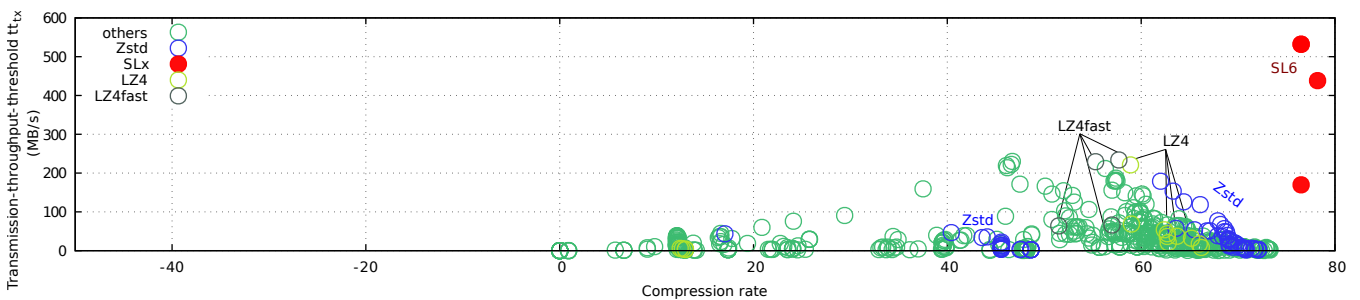
Algorithm	<i>H</i> -score
LZ shuffle	-4.62
LZ bitshuffle	-9.53
Zstd 22	0.92
Zstd shuffle 1	35.65
Zstd bitshuffle 1	11.02
LZ4Fast shuffle 17	44.92
LZ4Fast bitshuffle 17	13.96
<i>SL6</i> polynomial	54.52



**Figure 12.** Plot of the *H*-scores calculated on the results for the SWOT file.



**Figure 13.** Ranking of the algorithms,  $tt_{tw}$  (y) vs. the compression  $C_r$  (x) for the SWOT latitude field.



**Figure 14.** Ranking of the algorithms,  $tt_{tw}$  (y) vs. the compression  $C_r$  (x) for the SWOT L2 pixel\_area.

to execute the algorithms. Here, since the hardware is the same for all of them, these calculated metrics are comparable. Of course, the higher is better, it means that even with a extremely fast transmission channel, compressing and uncompressing the data on-the-fly is worth it. For example, on the Figure 13, it can be seen the  $tt_{tx}$  of the *SL6* algorithm (in red) is higher than 300 MB/s which corresponds to more than 2.4 Gb/s. Transmitting on a media that has a throughput lower than 2.4 Gb/s would benefit from compression using this algorithm. This is the combination of a good compression  $C_r$  (around 40 %) (X axis) and a low computation time. If now the media throughput is 120 GB/s (1 Gb/s), then using all the algorithm above this threshold would be more efficient than transmitting the original data, it means roughly the best LZ4fast LZ4 and Zstd variants, a few others one and *SL6*. As the compression time depends on the data, the, we have plotted the  $tt_{tx}$  for the pixel\_area field in Figure 14. As reminder, the pixel\_area field is quite smooth which makes it easier for the algorithms like the *SL6* which performs mathematical signal processing. Here, for the *SL6* algorithm, the  $tt_{tx}$  is higher than 500 MB/s which corresponds to more than 4 Gb/s for a compression  $C_r$  of more than 75 %.

### 3.3 Memory Usage

Memory usage of all tested algorithms is lower than 2 MB as well as for compression and decompression. Because of the necessity of building a dictionary, the Huffman-based ones often reach 2 MB. Their memory usage is not deterministic and depends on the data. The lowest is *SL6* with less than 128 kB (for the 64-bit floating point) version. Most of the memory usage





**Table 7.** Average results for the fields from the SWOT file that are the most difficult to compress, the compression  $C_r$  and the standard deviation, the compression and decompression throughput and the H-score are displayed.

Algorithm	Rate %	$\sigma$ %	Comp. tp MB/s	$\sigma$ MB/s	Decomp. tp MB/s	$\sigma$ MB/s	H-score
LZ bitshuffle 9	28.46	23.27	180.41	5.68	24.35	71.17	0.19
Zstd 22	21.35	27.03	4.51	0.95	768.78	580.69	-0.18
Zstd shuffle 1	31.96	23.77	371.82	109.61	741.43	117.69	19.46
Zstd bitshuffle 1	28.63	23.51	155.22	14.21	264.38	19.56	6.54
LZ4 shuffle	25.97	24.85	604.12	158.14	1009.12	64.52	4.51
LZ4Fast shuffle 17	23.82	25.03	741.24	150.69	1028.61	50.26	-6.49
LZ4 bitshuffle 0	27.26	22.50	217.77	7.02	278.56	4.56	9.10
SL6 polynomial	24.35	15.74	540.86	127.92	1320.52	230.79	32.21

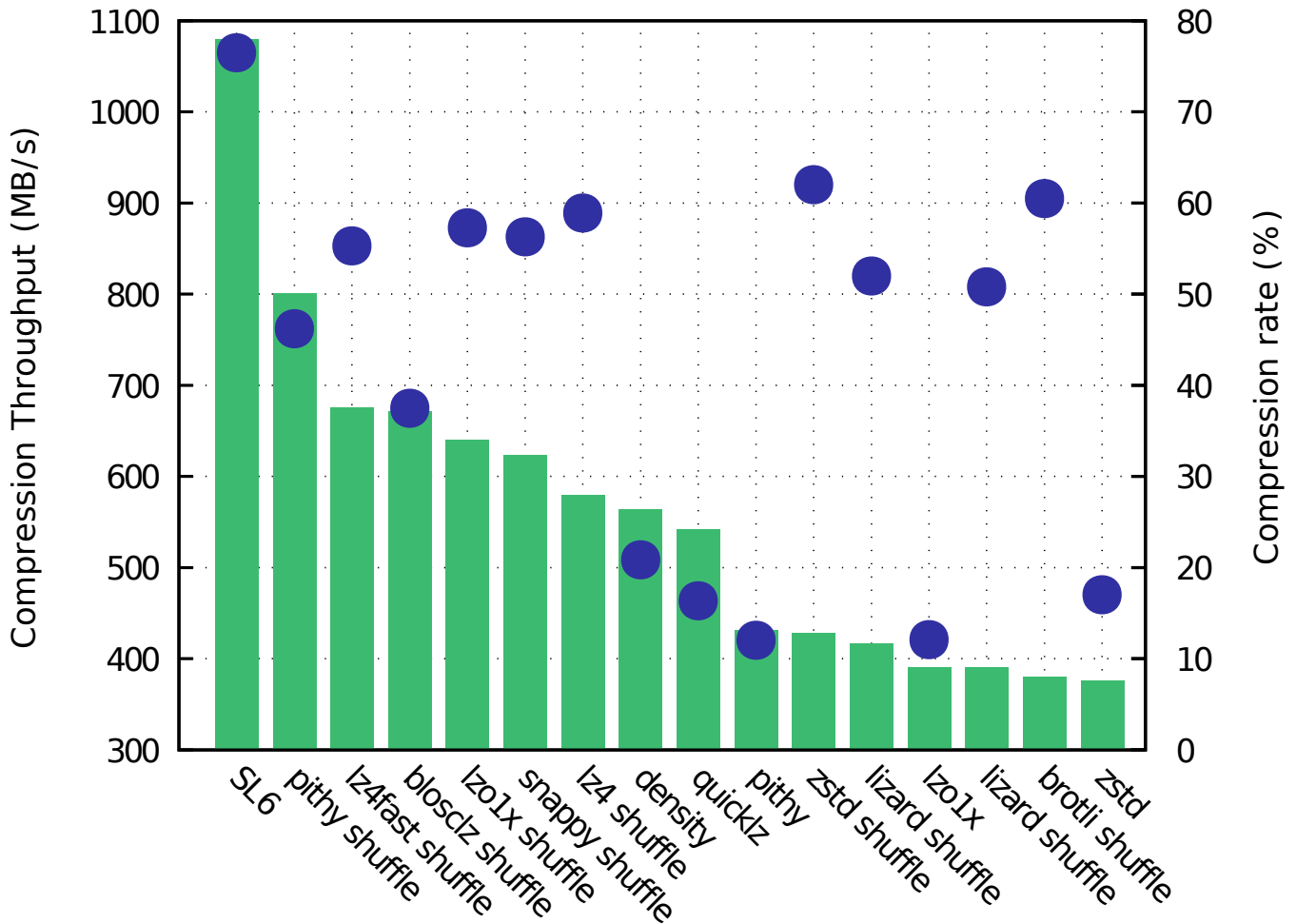
is used for the input and output buffers. Another key point is that the *SL6* memory usage is constant and does not depend on the data but only of the size of the words (8, 16, 24, 32, 64 bits), while the choice of the parameters has a limited impact on the memory usage of less than 5 %. Alternatively, the fact the Huffman based methods have to construct a dictionary requires memory, and makes embedded implementation bit more complex, especially on hardware targets. This is usually cope by working on chunks of data, in that cas, the compression level can be lower. Consequently, this feature makes the *SL6* approach ideal for embedded systems.

### 3.4 Results on relatively complex fields

The entropy of most of the fields of the SWOT file is quite low, which explains the high compression  $C_r$ . We decided in this subsection to focus on fields for which algorithms have more difficulties to compress. The average results are displayed in Table 7. The ranking differs from Tables 5 and 6. Indeed, for example, the LZ4Fast shuffle has a low  $H$ -score because of the important variability in the compression  $C_r$ : some of the fields are almost not compressed while the LZ4 bitshuffle now appears in the ranking at the third position. Figures 15 and 16 show that the compression  $C_r$  of the *SL6* algorithm is comparable to the other Huffman-based approaches associated to shuffle and bitshuffle filter. However, it shows superior performances in terms of compression speed. The high score of *SL6* is mainly due to the consistent compression  $C_r$  among the dataset associated to one of the best throughput ( $3^{rd}$ ). Some of the fields are (almost) not compressed by the LZ4 shuffle and LZ4Fast shuffle. Consequently, these algorithms tend to copy portions of the fields without any modification, leading to a quite high throughput.

## 4 Discussion

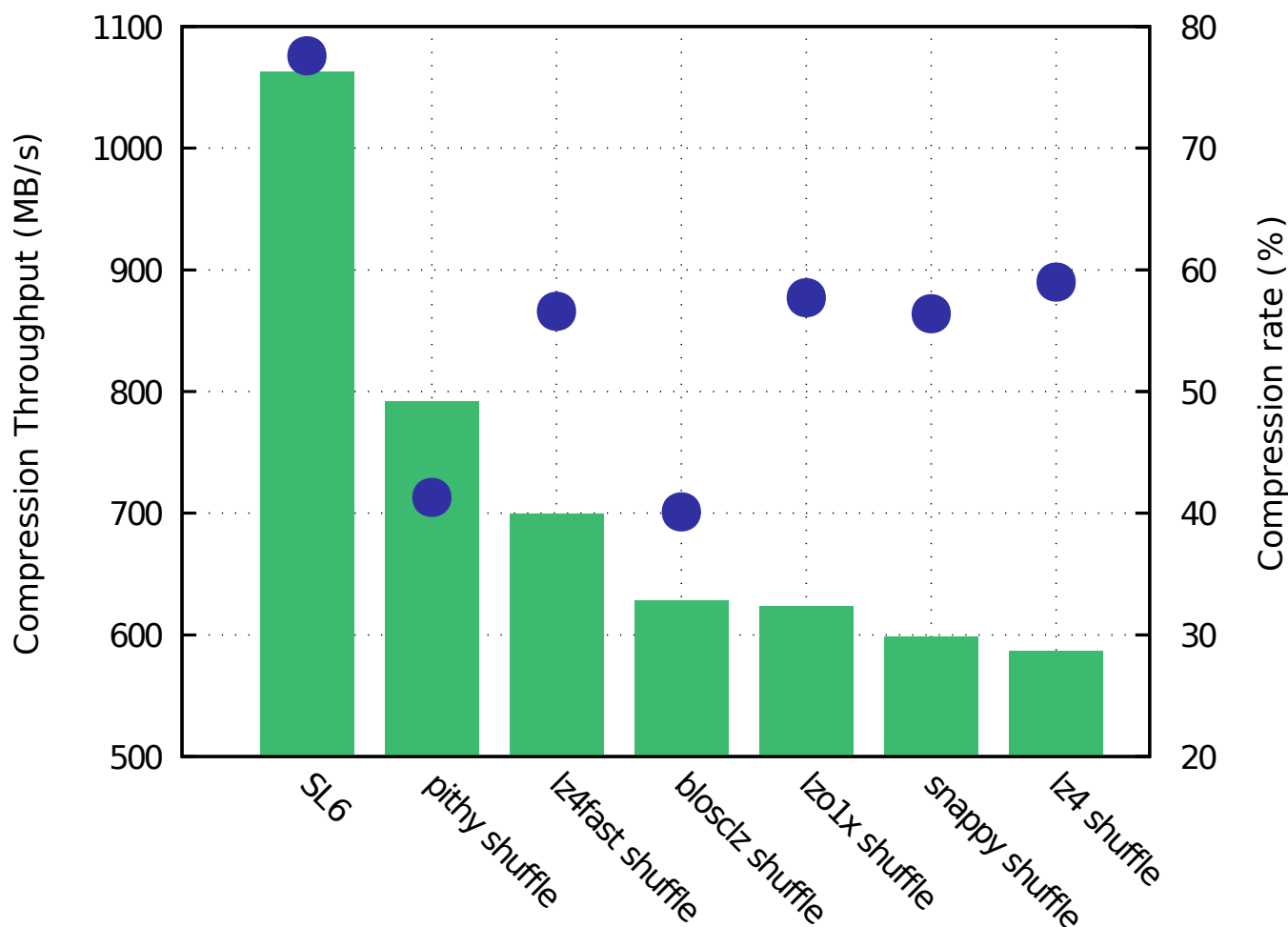
The results obtained with the extended testbench performed in this work show that the state-of-the-art algorithms performs similarly in term of compression rate. Indeed, if there is a certain level of variability among the different fields of the data



**Figure 15.** Plot of the compression speeds over 300 GB/s, as histogram boxes, and the associated  $C_r$  as blue points, for the SWOT Pixel\_Area field.

acquired by the swot mission, the average rate for the best algorithms is comprised between 75 % and 82 %. Considering that metric only, the LZ algorithm associated to bitshuffle (parameter 9) and Zstd associated to a shuffle preprocessing (parameter 1) provide the best results. However, the worst one does less than 10 % lower. Thus, the compression rate only is not a metric to be considered provided that we pick the algorithm among the top ten best compression rate. It is worth to note that some algorithms are not able compress some fields, this is the case of Pithy, some variants of the LZ4 etc. They have been excluded from the final results, even if for some fields they provide good performances typically the Pithy algorithm associated to the shuffle filter as visible on the Figures 15 and 16.

An other metric that needs to be considered are the memory usage of the algorithms performing compression and decompression. This information is not really relevant if the algorithms are executed by a workstation that can holds several GB of



**Figure 16.** Plot of the compression speeds over 500 MB/s, as histogram boxes, and the associated  $C_r$ , as blue points, for the SWOT CrossTrack field.

330 Random Access Memory (RAM), but it is extremely relevant if they are executed by an embedded system which can be quite  
 limited. Moreover, for space application, memory are prone to faults since cosmic rays can deposit their energy into one or several  
 cells. Thus choosing algorithms that have a low memory usage for such implementations make sense. The *SL6* algorithm  
 has a memory footprint lower than 128 kB while the others vary between 1 MB and 2 MB, which makes it a good candi-  
 date for embedded systems. This is also a good indicator of the capacity of the algorithm to be implemented in an embedded  
 335 programmable microcontroller or hardwired for a Field-Programmable Gate Array (FPGA) implementation.

The couple compression rate and throughput is a good way to rank the algorithms. However, since the analysis of the results  
 we have obtained have shown a quite important variability in the results, we decided to propose the *H*-score to take it into  
 account into a combined metric. It shows that the top-3 algorithms are the *SL6* ( $> 50$ ), the *LZ4Fast* associated of shuffle ( $> 40$ )  
 and the *ZStd* associated to shuffle which scores higher than 35.



340 Another interesting point is that the compression throughput of *SL6* is constant for a given datatype, which is quite under-standable as it operates in stream processing.

Beyond the choice of a compression algorithm for space-to-ground communication, where constraints require approaches that can be easily embedded and have consistent compression and throughput performances, the ground applications related to the SWOT mission is also considered. Compression for ground applications is extremely important too as the SWOT mission  
345 may generate around 5 TB of data per day. Indeed, ultimately, the data will be used in numerical models and computed on HPCs, thus, it is important to be able to compress and decompress the messages on-the-fly without any loss of performance. Compression of HPC messages is a good way to reduce computing time, provided that the compression/decompression is fast enough, that is the reason why we have introduced the transmission-throughput-threshold. In the same vein, to enable smooth storing an manipulating of these data requires, fast compression algorithms are required, ideally that can perform on-the-fly,  
350 with no or almost no perceptible delay for the operators. The *SL6* approach is the only one which has constant compression time for a given datawidth (64, 32, 16 or 8-bit words). It is also able to compress a signal with fixed-size buffers, which participates to its deterministic characteristics in terms of compression time and hardware resources that are required for the *SL6* algorithm execution.

Another key point is the determinism of the compression rate, indeed, the storage and database systems will be sized for an  
355 expected volume of data. The H-score metric that is proposed this paper clearly emphasizes the properties that we are seeking for both space and ground application.

## 5 Conclusion

Since space missions are generating more and more data, for example 7 PB is expected for the swot mission, efficient data compression approaches needs to be explored. Huffman-based algorithms are traditionally used for scientific application, es-  
360 pecially for netCFD data. It has been shown in the past they are quite efficient especially when associated to filters such as shuffle and bitshuffle. This paper compared Huffman based algorithms and the *SL6* approach. The Huffman-based algorithms require memory (shuffle needs to analyze the entire file) and their performances are heterogenous, strongly depending on the data. When associated to shuffle filter, stream processing is not possible. Alternatively, the *SL6* algorithm provides homogeneous results, low memory usage, constant compression time for a given datatype, especially in terms of compression  $C_r$ . On  
365 the dataset that is representative of the SWOT mission, the compression throughput is 35 % faster than the best algorithm of the bench for a similare compression level. It is also often almost twice faster than the other algorithms of this extensive bench.

*Code availability.* The code of the testbench was developped in Python and C. It uses the LZBench available on <https://github.com/inikep/lzbench>. The code for **SLX!** (**SLX!**) can be provided by Subnet, Pigoury@Subnet.fr. The codes for the other compression algorithms are included in the LZBench as opensource.



370 *Data availability.* The SWOT dataset is provided by CNES and is available upon demand.

*Author contributions.* OT is the inventor of the *SL6* compression algorithm, OT performed the benchmark, SP and MT have done the data analysis, FG provided the swot data and the comparison methodology, MT and FG wrote the paper

*Competing interests.* The *SL6* algorithm is under license at Subnet SAS – <http://www.subnet.fr>.

375 *Acknowledgements.* The authors would like to thank the Centre National des Etudes Spatiales (CNES) for funding this research and Subnet for providing the sources of the *SL6* algorithm. They also would like to thanks Nicola MARTIN for the proofreading work.



## References

- Christopher Rumsey, Bruce Wedan, T. H. and Pointot, M.: Recent Updates to the CFD General Notation System (CGNS), in: 50th AIAA Aerospace Sciences Meeting including the New Horizons Forum and Aerospace Exposition, 09-12 January 2012, Aerospace Research Central (ARC), Nashville, Tennessee (USA), <https://doi.org/10.2514/6.2012-1264>, 2012.
- 380 Delaunay, X., Courtois, A., and Gouillon, F.: Evaluation of lossless and lossy algorithms for the compression of scientific datasets in netCDF-4 or HDF5 files, *Geoscientific Model Development*, 12, 4099–4113, <https://doi.org/10.5194/gmd-12-4099-2019>, 2019.
- Devarajan, H., Kougkas, A., and Sun, X.: An Intelligent, Adaptive, and Flexible Data Compression Framework, in: 2019 19th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGRID), pp. 82–91, <https://doi.org/10.1109/CCGRID.2019.00019>, 2019.
- 385 Di, S. and Cappello, F.: Optimization of Error-Bounded Lossy Compression for Hard-to-Compress HPC Data, *IEEE Transactions on Parallel and Distributed Systems*, 29, 129–143, <https://doi.org/10.1109/TPDS.2017.2749300>, 2018.
- Diane Poirie and, S. A., McCarth, D., Smith, M., and Enomoto, F.: The CGNS system, in: 29th AIAA, Fluid Dynamics Conference 15-18 June 1998, Aerospace Research Central (ARC), Albuquerque, NM, U.S.A., <https://doi.org/10.2514/6.1998-3007>, 1998.
- Else, J. A. D. o. N. A. R. C. .: Data Compression for Space Missions, in: *The Space Congress Proceedings, 1968 (5th), The Challenges of the 1970's*, pp. 9.3.1–9.3.15, <https://commons.erau.edu/space-congress-proceedings/proceedings-1968-5th/session-9/3/>, 1968.
- 390 Fjørtoft, R., Gaudin, J.-M., Pourthié, N., Lalaurie, J.-C., Mallet, A., Nouvel, J.-F., Martinot-Lagarde, J., Oriot, H., Borderies, P., Ruiz, C., and Daniel, S.: KaRIn on SWOT: Characteristics of Near-Nadir Ka-Band Interferometric SAR Imagery, *IEEE Transactions on Geoscience and Remote Sensing*, 52, 2172–2185, <https://doi.org/10.1109/TGRS.2013.2258402>, 2014.
- Folk, M., Heber, G., Koziol, Q., Pourmal, E., and Robinson, D.: An Overview of the HDF5 Technology Suite and Its Applications, in: *Proceedings of the EDBT/ICDT 2011 Workshop on Array Databases, AD '11*, p. 36–47, Association for Computing Machinery, New York, NY, USA, <https://doi.org/10.1145/1966895.1966900>, 2011.
- 395 Fraire, J. A., Céspedes, S., and Accettura, N.: Direct-To-Satellite IoT - A Survey of the State of the Art and Future Research Perspectives: Backhauling the IoT Through LEO Satellites, in: *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 11803 LNCS, pp. 241–258, [https://doi.org/10.1007/978-3-030-31831-4\\_17](https://doi.org/10.1007/978-3-030-31831-4_17), 2019.
- 400 Howison, M.: High-Throughput Compression of FASTQ Data with SeqDB, *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 10, 213–218, <https://doi.org/10.1109/TCBB.2012.160>, 2013.
- Huang, B., ed.: *Satellite Data Compression*, Springer, <https://doi.org/10.1007/978-1-4614-1183-3>, 2011.
- International Telecommunication Union: Frequency bands and transmission directions for data relay satellite networks/systems SA Series Space applications and meteorology, Tech. rep., International Telecommunication Union, <http://www.itu.int/ITU-R/go/patents/en>, 2017.
- 405 Jayasankar, U., Thirumal, V., and Ponnurangam, D.: A survey on data compression techniques: From the perspective of data quality, coding schemes, data type and applications, *Journal of King Saud University - Computer and Information Sciences*, 33, 119–140, <https://doi.org/https://doi.org/10.1016/j.jksuci.2018.05.006>, 2021.
- Kunkel, J.: SFS: A Tool for Large Scale Analysis of Compression Characteristics, Tech. Rep. 4, Deutsches Klimarechenzentrum GmbH, Deutsches Klimarechenzentrum GmbH, Bundesstraße 45a, D-20146 Hamburg, 2017.
- 410 Lloyd, T., Barton, K., Tiotto, E., and Amaral, J. N.: Run-Length Base-Delta Encoding for High-Speed Compression, in: *Proceedings of the 47th International Conference on Parallel Processing Companion, ICCP '18*, Association for Computing Machinery, New York, NY, USA, <https://doi.org/10.1145/3229710.3229715>, 2018.



- Masui, K.: Bitshuffle: Filter for improving compression of typed binary data, 2017.
- Masui, K., Amiri, M., Connor, L., Deng, M., Fandino, M., Höfer, C., Halpern, M., Hanna, D., Hincks, A. D., Hinshaw, G., Parra, J. M.,  
415 Newburgh, L. B., Shaw, J. R., and Vanderlinde, K.: A compression scheme for radio data in high performance computing, *Astronomy and  
Computing*, 12, 181–190, <https://doi.org/https://doi.org/10.1016/j.ascom.2015.07.002>, 2015.
- Owen, T. D. and Daniel: Advanced Data Format (ADF) - A portable hierarchical database, in: *AIAA and SAE, 1998 World Aviation  
Conference, Aerospace Research Central (ARC), Anaheim, CA, U.S.A.*, <https://doi.org/10.2514/6.1998-5565>, 1998.
- Patauner, C., Marchioro, A., Bonacini, S., Rehman, A. U., and Pribyl, W.: A Lossless Data Compression System for a Real-Time Application  
420 in HEP Data Acquisition, *IEEE Transactions on Nuclear Science*, 58, 1738–1744, <https://doi.org/10.1109/TNS.2011.2142193>, 2011.
- Radhakrishnan, R., Edmonson, W. W., Afghah, F., Rodriguez-Osorio, R. M., Pinto, F., and Burleigh, S. C.: Survey of Inter-Satellite Commu-  
nication for Small Satellite Systems: Physical Layer to Network Layer View, *IEEE Communications Surveys Tutorials*, 18, 2442–2473,  
<https://doi.org/10.1109/COMST.2016.2564990>, 2016.
- Soumagne, J., Biddiscombe, J., and Clarke, J.: An HDF5 MPI virtual file driver for parallel in-situ post-processing, in: *Lecture Notes in  
425 Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 6305 LNCS,  
pp. 62–71, [https://doi.org/10.1007/978-3-642-15646-5\\_7](https://doi.org/10.1007/978-3-642-15646-5_7), 2010.
- Sudmanns, M., Tiede, D., Lang, S., Bergstedt, H., Trost, G., Augustin, H., Baraldi, A., and Blaschke, T.: Big Earth data: dis-  
ruptive changes in Earth observation data management and analysis?, *International Journal of Digital Earth*, 13, 832–850,  
<https://doi.org/10.1080/17538947.2019.1585976>, © 2019 The Author(s). Published by Informa UK Limited, trading as Taylor & Francis  
430 Group., 2020.
- Thomine, O., Latu, G., and Thevenin, M.: Data Compression, 2016, US Patent 20,190,044,532.
- Trott, A., Moorhead, R., and McGinley, J.: The application of wavelets to lossless compression and progressive transmission  
of floating point data in 3-D curvilinear grids, in: *Proceedings of Data Compression Conference - DCC '96*, pp. 458–,  
<https://doi.org/10.1109/DCC.1996.488390>, 1996.
- 435 Vaze, P., Kaki, S., Limonadi, D., Esteban-Fernandez, D., and Zohar, G.: The surface water and ocean topography mission, in: *2018 IEEE  
Aerospace Conference*, pp. 1–9, <https://doi.org/10.1109/AERO.2018.8396504>, 2018.
- Welton, B., Kimpe, D., Cope, J., Patrick, C. M., Iskra, K., and Ross, R.: Improving I/O Forwarding Throughput with Data Compression, in:  
2011 IEEE International Conference on Cluster Computing, pp. 438–445, <https://doi.org/10.1109/CLUSTER.2011.80>, 2011.
- WMO: Satellite Data Telecommunication Handbook, World Meteorological Organization (WMO), 2018.