# Impact of Lossy Compression Errors on Passive Seismic Data Analyses

Abdul Hafiz S. Issah and Eileen R. Martin 20

#### Abstract

New technologies such as low-cost nodes and distributed acoustic sensing (DAS) are making it easier to continuously collect broadband, high-density seismic monitoring data. To reduce the time to move data from the field to computing centers, archival requirements, and speed up interactive data analysis and visualization, we are motivated to investigate the use of lossy compression on passive seismic array data. In particular, there is a need to not only just quantify the errors in the raw data but also the characteristics of the spectra of these errors and the extent to which these errors propagate into results such as detectability and arrival-time picks of events. We compare three types of lossy compression: sparse thresholded wavelet compression, zfp compression, and low-rank singular value decomposition compression. We apply these techniques to compare compression schemes on two publicly available datasets: an urban dark fiber DAS experiment and a surface DAS array above a geothermal field. We find that depending on the level of compression needed and the importance of preserving large versus small seismic events. different compression schemes are preferable.

Cite this article as Issah, A. H. S., and E. R. Martin (2024). Impact of Lossy Compression Errors on Passive Seismic Data Analyses, Seismol. Res. Lett. 95, 1675-1686, doi: 10.1785/0220230314.

Supplemental Material

#### Introduction

are needed to capture key features of the defeorts in lossy Innovations in seismic instrumentation have given rise to a variompression of eismic data have been skewed toward active ety of ways of gathering data. Two major innovations are the sisjemic data and in the pasthave often focused on wavelet of low-costseismic nodes and fiber-optic distributed acoustic decomposition and the discrete cosine transform (Bosman sensing (DAS)both of which enable the recording of high-fre- and Reiter, 1993; Donoho et al., 1999; Averbuch et al. 2001). quency data on many closely spaced sen Forsexamplethe More recently, the advancement of machine learning has given Penn State Fiber-Optic foR Environment SEnsEing (FORESEE) to more efforts toward using autoencoders for the compresproject recorded urban environmental seismic data at 500 sasion of active seismic data (Valentine and Trampatia). ples per second across 4189 m of fiber-optic cable at 2 m spacifigthis article, we provide a quantitative assessment of the (Zhu et al., 2021; Spica et al., 2023) and the PoroTomo Natusalitability of three methods—wavelet decomposition, zfp float-Laboratory at Brady's Hot Springs shared geothermal production point compression, and low-rank singular value decompoand microseismicity data at 1000 samples per second across@i\u00e9m(SVD) approximation—for compressing passive seismic of cable at 1.021 m spacing (Coleman, 2016; Feigl et al., 2016) AThis assessment focuses on the errors in analyses rather sample of just nine DAS experiments (including these two) pulban the error in the raw data. We choose these methodiso, duced more than 750 TB of DAS data from 2015 to 2020d, part, because of the potentialor analytical bounds on these this substantially faster rate of data acquisition (in comparisoerrors' propagation. We give a brief overview of these methods, with long-period seismometers or nodalrays) has inhibited their comparison in different metrics for assessing the integrity geoscientists' ability to quickly access, analyze, and visualized meson structed data and the effect of compression on event detection. In addition, we have released open-source software new data sources (Lindsey and Mart20,21).

One solution to reduce the volume of seismic data is com-to enable easy application of these error analysis workflows to pression with higher compression ratios commonly achieved additional compression schemes and datasets in the future. through the use of lossy compression techniques. However, lossy compression techniques introduce data error we need to

quantitatively compare the various options for reducing data 1. Colorado School of Mines, Golden, Colorado, U.S.A., https://orcid.org/0000storage and data movement during passive seismic processing data 1,0002-4794-2754 (AHSI); https://orcid.org/0000-0002-3420-4971 (ERM); storage and data movement during passive seismic processing griginal Tech, Blacksburg, Virginia, U.S.A. Compression of seismic data can be achieved by the transformorresponding author: aissah@mines.edu

mation of data into a sparse representation so that fewer bytesseismological Society of America

## Background

Here, we provide an introduction to three types of compression that are frequently applied to reduce the size of spatiotemporal scientific data: wavelet compression (1D and 2D), zfp compression, and SVD compression. We use two publicly available DAS These data include one urban dark fiberdataset, called the FORESEE data, and one geothermal microseismicity monitoring

dataset, called the Brady's Hot Springs data. Here, we provid@Dawavelet compression overview of the compression schemes and datasets. Throughout additional redundancies 2D datasets, wavelet

size of the original data to the size of the compressed data.

#### 1D wavelet compression

Given a time-domain signal f(t), we can represent it in terms of sition can be computed as follows: a spanning set of wavelet function shis spanning set is generated using scaled and translated versions of a mother wavelet  $\overset{XJ}{}$   $\overset{X^{j-1}}{}$   $\overset{X^{j-1}}{}$   $\overset{X^{j-1}}{}$  function,  $\psi$ , and a father wavelet function,  $\varphi$ . The discrete  $\overset{j}{}$   $\overset{m0}{}$   $\overset{n0}{}$   $\overset{n0}{}$   $\overset{n0}{}$ wavelet transform of fcan be computed as follows:

in which the inner products  $\eta$  dhf,  $\psi_{in}$  i and  $\psi_{in}$  in  $\psi_{in}$  in in which  $\psi^1$ ,  $\psi^2$ ,  $\psi^3$ , and  $\Phi x$ ,y are the 2D wavelets and are denote the detail and approximation wavelet coefficients, respen puted from the 1D wavelets as  $\psi^1 x$ ,  $\psi x \psi y$ . tively. -L > 0 is the number of nested subspaces, correspond  $\mathbf{g}^2$  to  $\mathbf{v} = \mathbf{v} + \mathbf{$ the number of scaling levels represented in the wavelet basist pathyl is related to the number of discrete points in f, that is, N<sup>⊥</sup>.2 The functions  $\mu$  and  $\phi_n$  are the scaled and translated version nation coefficients Similar to the 1D casewe can threshold of the mother and father wavelets, respectively, and are defirtbe as efficients to achieve compression with approximation follows:

$$\psi_{j,n}x$$
  $\frac{1}{2^j}\psi$   $\frac{x-n}{2^j}$  and  $\phi_{j,n}x$   $\frac{1}{2^j}\phi$   $\frac{x-n}{2^j}$ : 2

high amplitude in time (Mallat, 2009). This property allows the 2D compression considering the differentchannels asa application of the wavelet transform for denoising (Donoho, 1995) or for the preservation of events when used for compression Villasenor et al. (1996). For a predetermined threshold Tzfp floating-point compression we can construct an approximation of our signal,

The error in approximation E is then

$$\begin{array}{cccc} X^J & \cancel{X}^{J-1} & \cancel{X}^{J-1} \\ & \tilde{d}_j n \psi_{j,n} & \tilde{a}_J n \varphi_{J,n} & 4 \\ jL1 & n0 & n0 & \end{array}$$

datasets for testing and comparing these compression tech- If a large enough T is selected, this approximation has few nonniques' effects on the characteristics of ompressed data and zero coefficients that can be encoded to achieve compression. the errors incurred in the results of event detection workflows. Hence, the threshold determines the amount of compression and the associated error in the approximation.

this articlethe compression factor is defined as the ratio of theompression can be performed using 2D wavelet decomposition (Villasenor et al., 1996). In this approach, separable 2D waveletsderived from 1D wavelets are used to decompose the data in both dimensions. For 2D data, f, the wavelet decom-

let 
$$X^{j}$$
  $\mathscr{R}^{-1}$   $\mathscr{R}^{-1}$   $d_{j}^{-1}$   $m, n\psi_{j,m,n}^{1}$   $d_{j}^{2}$   $m, n\psi_{j,m,n}^{2}$   $j$   $L^{1}$   $m^{0}$   $n^{0}$   $\mathcal{R}^{j-1}$   $\mathcal{R}^{j-1}$   $d_{j}^{3}$   $m, n\psi_{j,m,n}^{3}$  ...  $a_{j}$   $m, n\Phi_{j,m,n}$   $a_{j}$ 

 $\Psi^3$ x,y  $\Psi$ x $\Psi$ y, The coefficients \$\daggamma m,n, d^2m,n, and d^3m,n are referred to as detail coefficients, and,a are called approxierror.

High-dimensional wavelet compression has been studied to achieve higher compression rates for active source seismic data organized by streamer number, shot number, sensor, and time Typically, this representation has a small number of large-dimensions (Villasenor et al., 1996). For passive DAS data, we magnitude wavelet coefficients around and leading to areas of investigate how 1D wavelet compression does in comparison dimension in space in addition to the recording in time.

The zfp compression technique uses processes such as block transforms and embedded codingommonly used in image compression, to perform compression that is suitable for a variety of floating-point scientific data as detailed by Lindstrom (2014). We briefly outline the process he Feor d-dimensional data, the data array is sectioned into blocks of affalues, which are assumed to be approximately continuous within any block. Each block is compressed separately; the following are the steps during compression that may introduce some errors:

1. Convert floating point values in the 4 d block to scaled integers with a common exponent.

- Perform a block transform to introduce some sparsity into The FORESEE data the integer representation.
- 3. Encode the numbers in the sparse representation only. The sploration, we used data from the FORESEE urban DAS many bits as required to save nonzero entri**Es**e process also allows for specified bits to be allocated for each block~46 TB of data (Zhu et al., 2021). The data examplesin (fixed rate), the number of binary exponents to encode (fixed precision), and the maximum error allowed for each one with a few events due to thunder earthquakesand the floating point value (fixed accuracy).

Although Wade (2020) released a zfp formatind implementation designed to handle triggered active-source seismisequentdata exploration, analysis, and interpretation. This data formats, here we use the general p implementation to align with a wider variety of zfp error analysis studies. Although this article focuses on analyzing the errors introduced in passive seismic data workflows, other teams have pure an environment (Zhu and Stensrud, 2019; Zhu et al., 2021). viously studied the numericaerror due to zfp (Diffenderfer et al., 2019) and its effect on workflows for severalfluid dynamics and plasma physics problems as well as climate mod-study the use of lossy compressed data for microseismicity

#### SVD compression

used in various scientific fields (e.gsignal processing image compressiondata mining, and machine learning) because of seismic activities that have been previously studied and cataits ability to identify and capture the highest possible amount loged (Li and Zhan, 2018). We take advantage of the existing of variability in the data. Seismic data collected through DAS workflow and resulting catalog to compare the detectability matrices. The application of SVD to these matrices offers a unique opportunity to decouple the channelsand samples and enable efficient processing (Mart20,19).

expensive; and omized SVD is a powerful and efficient algorandomized projection methods to quickly approximate the dominant singular vectors and singular values. One of the pristeps such as imaging and inversidationately, this can lead mary benefits of randomized SVD is its ability to efficiently decomposelarge matrices storing high-dimensional data range of scientific applications. Given data stored a matrix  $D \in R^{N_c \times N_t}$ , in which  $N_c$  is the number of channels, is the rank; and N<sub>t</sub> is the number of time sample so achieve compression, we construct a low-rank (rank, r) approximation, using randomized SVD, in which  $U \in R^{N_c \times r}$ ,  $\Sigma \in R^{r \times r}, \ V \in R^{r \times N_t}$ . Then we combine U and  $\Sigma$  into an  $(N_c \times r)$  matrix. Storing  $(N_c \times r)$  and  $(r \times N_t)$  matrices provides a compression factor of  $\frac{N_c N_t}{N_t N_c}$ . "Compression factor" data to the size of the compressed data.

To study the use of lossy compressed data for passive data encoding takes sparsity into consideration and only uses astudy. This study's publicly available data were continuously recorded between April2019 and October 2021 resulting in Figure 1 illustrate two instances of recordings in the dataset, other with 30 min of passive data with vehicles and noise due to infrastructure. These examples provide visual sights into the characteristics of the recorded datacilitating subdatasetcontains recordings of signals like these and others from both natural and anthropogenic sourcesproviding a comprehensive representation the seismic activity in the

#### The Brady's Hot springs data

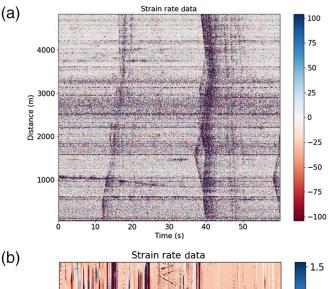
els (Laney et al., 2013; Baker et al., 2016; Poppick et al., 2020) tection and extraction of arrival times for tomographic imaging, we used data from the Brady's Hot Springs geothermal field. This was recorded in 2016 in Nevada as part of an investigation The SVD can be used to decompose data represented in mainto the feasibility of using DAS for cost-effective monitoring of form into the product of three matrices: the left singular vectogs othermal reservoirs. The data consist of ~8.7 km of fiber-optic a diagonal singular value matrix, and the right singular vectorsable deployed horizontally in a shallow backfilled trench, with in which the singular values represent the amount of variation in in the channel spacing. The data, which were shared publicly, the data explained by each singular vector. SVD has been wighted sampled at 1000 samples per second and recorded ~40 TB of data in 15 days (Colema@016). The dataset recorded microcan be organized as "channels by samples" to obtain seismicodetents and arrival times in various types of compressed data.

#### Computational Experiments

Errors are expected in lossy compressed seismic data, and these Although constructing the full SVD can be computationally errors can affect the quality of the data as wells the results obtained from seismic processing workflows. These errors can rithm for computing partial SVD of large-scale matrices usingcause distortions in the seismic waveforms and may propagate through the seismic processing workflow, affecting subsequent to interpretation errors based on seismic propertiesaging, and inversion results. Therefore, it is crucial to carefully evalu-(Halko et al., 2011). This makes it an attractive tool for a wideate the effects of lossy compression on seismic data and the extent to which these propagate through processing workflows. Here, we present several experiments designed to explore how errors from varying compression types and ratios propagate into (1) the level of noise in the data.(2) the distribution of error across frequency rangeand (3) errors in key metrics for microseismic event detection.

#### Norm of error in data

as used here is defined as the ratio of the size of the original To investigate the level of noise introduced by different compression schemeat various compression rateswe studied



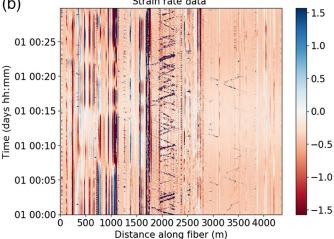


Figure 1. Data from the Penn State Fiber-Optic foR Environment SEnsEing (FORESEE)-urban project (a) example of 60 s of continuous data recorded at 03:33:35 on 15 April 2019. This shows recordings of multiple thunder earthquakes starting at 10 and 40 s. It also shows some traffic noise around 1000 m from 0 to 25 s. (b) Twenty min data used for frequency preservation experiment recorded from UTC 00:00:04 to 00:20:04 on 1 August 2019. The coherent signals here are traffic noises between 1500 and 2500 m along the fiber. These visual resentations provide valuable insights into the characteristics of the recorded data, offering us a better understanding of what to anticipate during the analysis and interpretation of the compressed data. The color version of this figure is available only in the electronic edition.

10 days of data from the FORESEE urban projectecorded from UTC 23:51:35 on 04/09/2019 to UTC 00:07:35 on 04/21/2019. We selected this particular time range because it contains a diverse range of ecorded signals that we wish to preserve when compressing passive seismic of at a each compression schemethe comparison workflow is as follows:

- 1. Begin with data  $D \in R^{N_c \times N_t}$ .
- 2. Compress D.

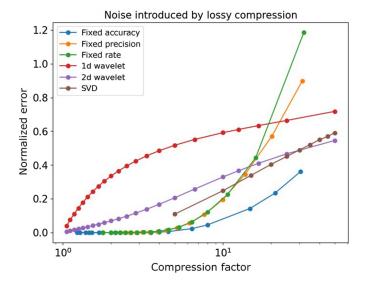
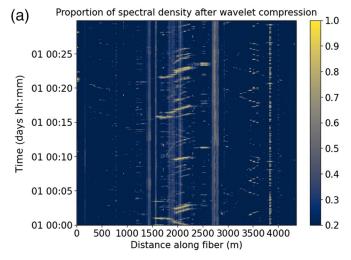


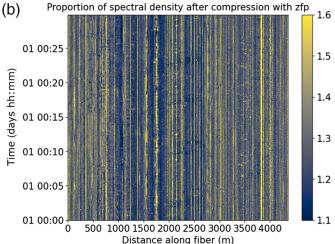
Figure 2. Frobenius norm of noise introduced in data at various levels of compression for compression with wavelet decomposition, singular value decomposition, and zfp floating point compression (three modes—fixed accuracy, fixed precision, and fixed rate). The color version of this figure is available only in the electronic edition.

- 3. Reconstruct compressed date R<sup>N<sub>c</sub>×N<sub>t</sub></sup>.
- Compute the normalized Frobenius norm error defined as ||D − D̃||<sub>F</sub> = ||D||<sub>F</sub>.

The results of this test are presented in Figure 2. Among the three modes of zfp compression, fixed accuracy maintains lower errors than fixed precision and fixed rate at similar compression factors. This can be attributed to the fixed accuracy mode encoding as many bit planes as necessary to achieve a specific absolute error; the other modes encode a fixed amount of bit planes irrespective of error level. The drawback to the fixed accuracy mode is that it requires complete knowledge of the range of values in the data to provide an absolute error tolerance, making it unsuitable for some on-the-fly compression of streaming data. When comparing 2D and 1D wavelet compression, it is observed that 2D wavelet compression results in a lower error the same compression rate improved performance could be attributed to the additional dimension available forcompression. which allows for more effective use of space and time redundancy (Villasenor et al.1996).

Regarding how the different compression schemes compare, we found that the three modes of zfp compression introduce the least noise up to a compression rate of ~15×, followed by SVD,2D wavelet and then 1D wavelet compression in the same rangeHowever, at higher compression rates fp compression begins to incur errors at a higher rate, with only fixed accuracy mode maintaining lower errors than SVD and 2D wavelet compression. The errors incurred by wavelet compression may be attributed to its denoising qualityHowever, at





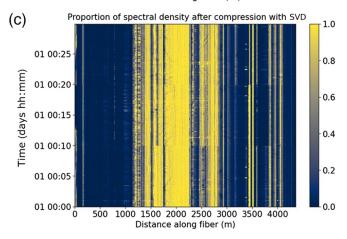


Figure 3. Frequency preservation with respect to compression for (a) wavelet compression, (b) zfp compression, and (s)ingular value decomposition (SVD) compression. These images are  $N_t = T_w \times N_c$  data in time and channel, respectively, obtained by averaging  $P_{\!\!\!\!D=D}$  along the frequency axis. The color version of this figure is available only in the electronic edition.

higher compression rates; ompression may be depleting the real signals to some extentSVD produces errors that grow at a nearly linear rate relative to the log of the compression

factor. The factors governing the error bound is an area that may be explored in future studies.

To facilitate meaningful comparisons, we selected a specific mode for wavelet compression and zfp compression, assuming that the various modes within each compression method yield comparable error characteristics. For wavelet compression, we opted for the 2D mode because ofts observed tendency to introduce a lower norm of noise, as evidenced by the previous analysis. On the other hand, for zfp compression, we chose the fixed precision mode instead of the fixed accuracy mode, despite the latter's exhibiting the least norm of noise. The fixed precision mode allows for relative error controland can be applied to unexplored datasets without requiring specific fine-tuning for optimal tolerable error and compression ratio. This consideration becomes crucial if the selected compression method is intended for use on data as they are being gathered.

#### Errors in frequency content

To assess the extent to which the frequency content of the data is preserved following lossy compressianmethodology was used whereby the proportion of the power spectrum retained at each frequency was calculated and then averaged across frequencies. The step-by-step process is as follows for each file containing a window of data in time:

- 1. Obtain  $\tilde{D}$  previously defined by compression and subsequent decompression of D.
- 2. Identify the power spectrum P  $\tilde{D}$  of  $\tilde{D}$  in predetermined time windows,  $T_w$  (5 s for this experiment) such that P  $\tilde{D} \in N_c \times N_t = 2 \times N_t = T_w$ .
- 3. Divide the power spectrum at each frequency by the original power spectrum that is,  $P_{\tilde{D}=D} = \frac{P \tilde{D}}{PD}$ .

The outcome was a 3D data cube for  $\mathcal{P}_{=D}$  that captured information across time windows of size  $\mathcal{N}_{=D}$ , frequencies of size  $\mathcal{N}_{=D}$ , and channels  $\mathcal{N}_{=D}$ . Weighted averaging was then performed along each dimension to investigate the patterns of frequency preservation across the various dimensions.

The data obtained by averaging the ratios of the compressed to original energy across all frequencies can be represented in two dimensions (time by channels), revealing the variations in frequency preservation as depicted in Figure 3 the patterns observed in this 2D matrix closely resemble the trends in the data shown in Figure 1. In the case of wavelet and low-rank compression, a general reduction in energy can be attributed to the thresholding operation used in both the compression methods. In addition, whereas wavelet decomposition exhibits better preservation of small events, SVD compression tends to preserve mostly high-amplitude events. On the other hand, zfp compression leads to a general increase in energy but preserves the energy around the identified events to levels close to precompression levels his may be due to the spurious frequencies introduced at multiples of a quarter of the sampling rate as

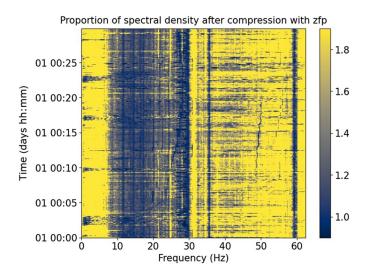


Figure 4. Frequency preservation with respect to compression for zfp compression. This image is constructed from  $N=T_w \times N_t=2$ data in time and frequency, respectively, obtained by averaging  $P_{\tilde{n}=D}$  along the channel axis. The color version of this figure is available only in the electronic edition.

tion analysiseven weak events that may be buried within the background noise can be identified and accurately located (Gibbons and Ringdal2006). To evaluate the impact of noise introduced by different lossy compression schemes on event detectivenconducted a series of experiments using a template matching workflow. In particular, we investigate two questions:

coefficient of the template waveform and recorded seismic data at different time intervals. This correlation coefficient is nor-

malized to account for variations in signal amplitude and noise

of template matching lies in its ability to enhance the detection capability for microseismic eventBy applying cross-correla-

levels and provides a quantitative measure of similarity

between the template and the recorded datae significance

- 1. To what extent is the array-wide detection significance of varying-size events impacted?
- 2. Are there biases or increases in the variability of event times picked across the array as higher compression ratios are used for any types of compression?

To test microseismic event detection via template matching,

depicted in Figure 4This effect can be attenuated when there we use the Brady's Hot springs data. This dataset has been preare other strong frequencies corresponding to events but tendiscusly studied and various microseismic activities have been to be exaggerated in noisy parts when there are not many otheraloged providing a baseline catalog for comparison (Li and strong frequencies on trast, such false frequencies are not Zhan, 2018). Our goals in this study are to apply a similar observed in wavelet or low-rank compression. workflow as outlined by Li and Zhan (2018) and to compare

(i.e., columns of the data matrix). This discrepancy may ter-preserved frequency characteristics compared with othersion performance for the compression schemes considered. when those channels have high-amplitude events (e.g., vehicles). Notably, the channels that were better preserved through compression as evaluated in the time domain also showed improved representation in the frequency domain. Despite compression being applied in smaller time windows, at 08:39:13The resulting array-wide normalized cross correthe channels thatdemonstrated good preservation remained lation is shown in Figure 5. consistent throughout the entire time period under consideration (i.e., well-preserved channels continued to be well-pre- collections of point sensors with large amplitude normalized served in most other time windows). This consistency cant portion of the overall variation of the data. However, when many spatially distributed ray paths are needed (e.g., we use a simple metricalculating the array-wide average of many source-receiver pairs for imaging) ormalization prior compressed data are high quality.

Changes in template matching event detection comparing known seismic events, referred to as templates, wtithe interference. continuously recorded seismic data underlying principle of template matching involves calculating the cross-correlationle compression levels is shown in Figure 6. Three events were

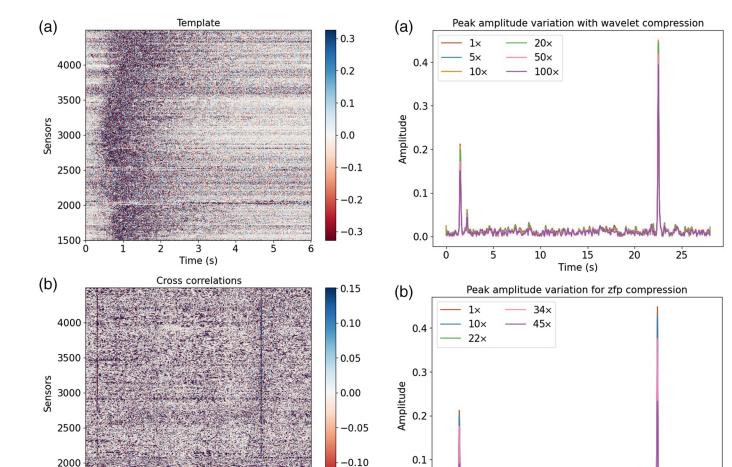
Finding a low-rank approximation of the data can result in the performance of event detection at various compression varying levels of representation fidelity for different channels rates for 2D wavelet, zfp, and SVD compression. By analyzing the results of these experimentee can gain valuable insights explain the observed trend of certain channels exhibiting bet-into the trade-off between compression rates and event detec-

> To provide a detailed view of the template matching process, we calculated the normalized cross correlation of each channel's recording ofthe template eventin Figure 5 with its recording of the noise from UTC 14 March 2016 beginning

Typically, template matching would be carried out on small cross correlations being considered a possible event and potensuggests that these specific columns likely account for a signifially a voting system among sensors to ensure multiple sensors detected the same possible evel that thousands of sensors, the normalized cross correlations at each time, then calculating to compression may be required to ensure that many channells envelope of the resulting time series. Despite time delays in the original event as it moves across the array, if a similar event occurs in the same locatiothe large normalized cross correlations across the array are expected to occur at the same time Template matching enables the detection of similar events byag. Thus, averaging across all sensors does not cause destruc-

The average envelope for each type of compression at multi-

1680



-0.15

0.0

Figure 5. (a) An example template event recorded at Brady Hot

Springs at 08:39:05.24 on 14 March 2016. (b) The array-wide template matching results show the normalized cross correlation of each channel's continuous recording with its template event recording. Verticallines indicate many channels with high similarities at a particular time. The color version of this figure is available only in the electronic edition.

15

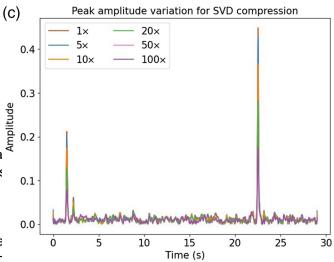
Time (s)

25

1500

detected, even though these events were barely distinguisha in the raw data. We refer to the first event (between 1 and 2 s as event 1, which is the midsize event, the second event (between 2 and 3 s) as event 2, which is the small event, and the third event (between 22 and 23 s) as eventwhich is the largest event. We see that for all three types of compre sion, the three events appear to be largely distinguishable frc... the background noise level. 2D wavelet compression maintains a more constant peak amplitude across compression ratios (Figure 6. The envelope of the average normalized cross correla-

original, 5×, 10×, 20×, 50×, and 100×) than do zfp and SVD compressionalthough there is a smallamount of amplitude loss at higher compression ratios(noticeable at 20×, 50×, and 100×). The zfp compression shows some amplitude loss figure is available only in the electronic edition. in the events, which particularly makes it difficult to



10

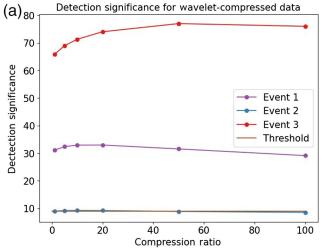
15

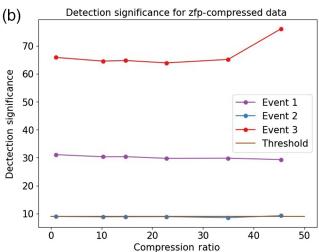
Time (s)

20

25

tions during three events in Figure 5 shows that the event picks are largely similar for various levels of compression I his was tested with (a) 2D wavelet-compressed data, (b) zfp-compressed data, and (c) SVD-compressed data. The color version of this





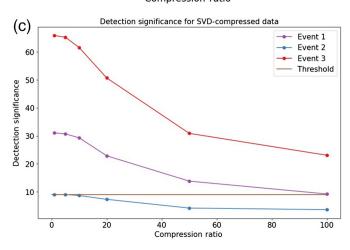


Figure 7. Change in detection significance with compression rate for (a) 2D wavelet-compressed data, (b) zfp-compressed data, (c) low-rank SVD-compressed data. The color version of this figure is available only in the electronic edition.

distinguish event 2but events 1 and 3 are clearly above the significance for the mid- and small-size events The image noise levelat all compression levels (1× original,10×, 22×, 34×, and 45×). The SVD compression leadsto significant energy loss, particularly at the higher compression ratios

(e.g., the peak of events 1 and 3 for 100× compressed data are less than half their original amplitudes).

As in Li and Zhan (2018), we use the detection significance of each event pick to provide a single array-wide value to quantitatively compare the preservation of each event when using various compression schemes and ratios. The detection significance is defined a  $\frac{CC-M}{MAD}$  for a value CCon the array-wide average of the normalized cross correlations at the ith time sample. M is the median of CC and MAD is the median absolute deviation defined as medianjCC - Mj. We set a detection significance minimum threshold of 9, which was set to provide a reasonable bound on false detections following laind Zhan (2018). The detection significance for each of the three events across all compression schemes and compression ratios is shown in Figure 7. All three compression schemes preserve events 1 and 3 (the midsize and large events) as picked events above the threshold, even at high-compression ratios (e45.x and 50x for all and 100× for wavelet and SVD). The SVD-compressed data have the largestdrop in detection significance relative to 2D waveletcompressed and zfp-compressed data.00× SVD compression, midsize event 1 is barely above the detection significance threshold. Even in the original uncompressed data he small event 2 starts out very close to the detection threshol/dith SVD compression at 20× and higher compression and with 2D wavelet compression at the 100× levelent 2 drops substantially below the threshold for detection significande.is not surprising thata barely detectable eveint the raw data could be lost in some of the highly compressed dataough all events compressed with SVD show a drop in detection significance the 2D wavelet compression and zfp compression of the large event 3 data actually show an uptick in detection significance at high-compression ratioslikely indicating some denoising occurring to further emphasize this largest event during the compression process.

Figure 8 shows the event detected using the template matching workflow discussed in the original datahen compressed data are used for althe events in this catalogue see similar trends to the one explained by the small-scale classeresults of this are summarized in Figure 9 showing the variation in detection significance with compression. In these plots, we expect a trend with a slope of 1 in the situation in which there are not any changes in detection significance. For wavelet compression (Fig. 9a), we have a lot of points close to this trend even at high-compression rates though points with smaller detection significance show a slight decreased points with highdetection significanceshow some increaseindicating some denoising effectThis leads to smaller events eventually being missed at higher compression rates. SVD compression (Fig. 9c) shows a similar trend but shows more reduction in detection for zfp (Fig. 9b) shows less predictabilitysome smallevents increasein detection significanceeven though most events are observed to be reducing in detection significance.

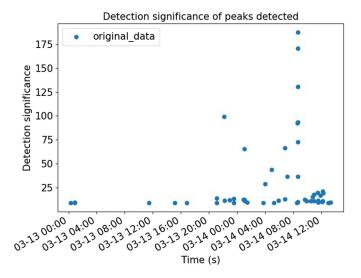
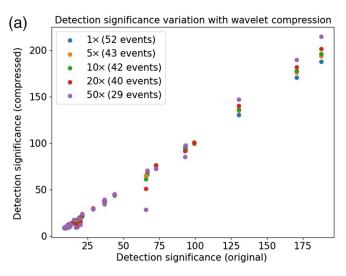
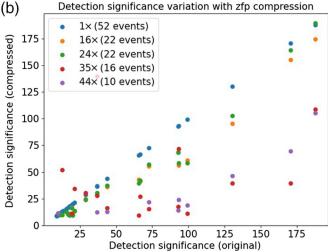


Figure 8. Catalog of events detected by template matching using the uncompressed data. The color version of this figure is available only in the electronic edition.

To address the question of whether event pick times are reable, we need to quantify the distribution of how each event's pick times change throughouthe array of compressed data compared with the event pick times on the originatlata. In particular, we need to know if the median pick time shift is staying very close to 0 (indicating there is no substantial array-wide bias) and compare the rate at which the minimum Q1, Q3, and maximum values are spreading apart for higher compression ratios. For each type of compression and for ea compression ratiowe created a box and whiskers plot for the distribution of each event's changes in pick times across all channels, which is shown in Figure 10F. or each type of compression, these boxplots are overlaid for all three events (colcoded) so that the spread for a small, midsized, and large ev can be compared easily.

There is not an apparent median bias for any of the events or compression typeswith a compression ratio <100×. At 100×, there appear to be a slightly (<0.05 s) late median of picks in SVD compression and a positive skew distribution (based on investigation of mimax, Q1, and Q3) to the picks in 2D wavelet compression at 100 hough the median of the wavelet distribution appears to be unbiased. We see that for types of compression the extremes and quartiles spread out more with higher compression ratios. In the range of 5x-15× compression, all three compression schemes perform reasonably well, with 5× 2D wavelet, 14× zfp, 5× and 10× SVD compression all yielding distributions for all three events that compression rate is increased for (aPD wavelet, (b) zfp, and are completely contained within a ±0.1 second shiftmong these low-ratio compression schemed 10× waveletcompression does perform the worstwith the largest spread for all three events and the small event 2 distribution only has its inner quartile range (Q1-Q3) contained within ±0.1 s,





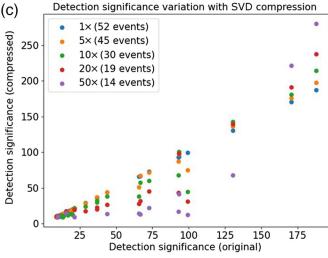
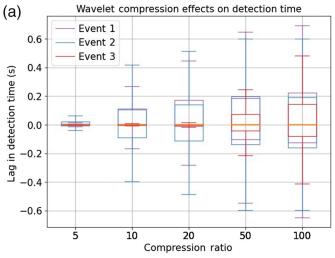
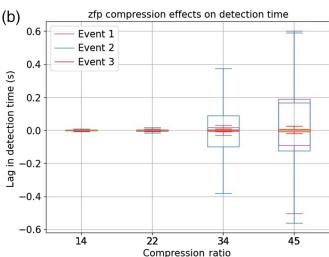


Figure 9. Trends in detection significance and events detected as (c) SVD compressions. The color version of this figure is available only in the electronic edition.





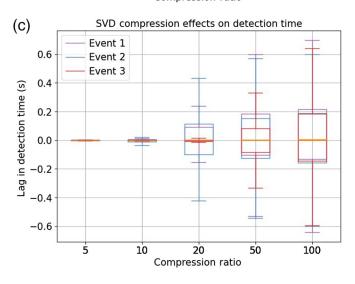


Figure 10. Boxplots show the distribution across allchannels in the array of picked event times from template matching applied to the (a) 2D wavelet-compressed data, (b) zfp-compressed data, and (c) low-rank SVD-compressed data. These are shown at various compression ratios for three events. The color version of this figure is available only in the electronic edition.

whereas its extremes have errors bounded by 0.3 s. In the range 20×-35×,we see that zfp at 22× has the only distributions for all three events that are completely contained within ±0.1 s. although both 2D wavelet and SVD perform similarly well at 20× compression for the large event 3n the range 45×-50×, we see that 2D wavelet at 50× and zfp at 45× have very compact distributions for the large event 3whereas SVD at 50× has a substantial spread in event pick shifts. All three compression schemes in this compression ratio range perform similarly for events 1 and 2, with extremes that exceed ±0.4 s and interguartile ranges that are bounded within ±0.2 s. The zfp did not allow higher compression ratios, but 2D wavelet and SVD compression were tested at the 100× ratiaD wavelet compression outperformed SVD compression on the large event 3 in the sense of providing a more compact distribution of time shifts although 2D wavelet's distribution of pick time shifts is skewed to have larger positive shifts. Both schemes performed similarly at the 100× levelon events 1 and 2, with extreme shifts around ~±0.6 s.

#### Discussion

Overall, we see that although zfp has the lowest data errors at lower compression ratioswaveletcompression (especially in 2D) has lower errors at higher compression ratios, and low-rank SVD has an error growth that sits in between zfp and wavelets. Wavelets strongly improve broadband representation of strong events over quiet noise, and SVD tends to have better broadband representation of louder signalsut zfp tends to more evenly distribute errors in the frequency content across loud and quiet eventsUsing a microseismicity dataset could see that after feeding compressed data through a template matching workflow, all types of compression could preserve the eventet smaller compression ratio S.VD compression tended to have the largest drop in detection significance at high-compression ratios, although it still preserved the detection significance of a midsize and large event even at 100× compression. The unbiased picks with increasing variability caused by higher compression ratios from Figure 10 suggest the opportunity to design a postprocessing scheme that romotes spatial coherency in eventpicks across the arrayn this way, more reliable picks can be used from any highly compressed datarticularly as an input to the event location or for tomographic imaging using microseismic event\$his analysis workflow was extended to a 36-hour period recording with 52 events of varying detection significanceand we found that wavelet compression preserved the detection significance better than zfp and SVD compression at similar compression ratios and tended to increase detection significance for larger events and higher compression ratios (i.e., emphasizing and denoising large events). Zfp compression typically led to a reduction in detection significance across all event sizesSVD compression tends to reduce detection significance for smaller to midsize events and tends to increase detection significance for larger eventwith more decrease or increase

in significance for higher compression levelstimately, SVD and wavelet representations have been integrated into a largeanies of the Center for Wave Phenomena (CWP), their support number of analysis workflows that can operate on data directlyade this research possible he authors thank the Virginia Tech in their compressed representation, which may lead us to prefer and Research Computing(ARC) and Colorado School of these in some contexts, but zfp is being increasingly used in Scientific computing so algorithms incorporating zfp may be on the horizon in the coming years. Beyond the scope of this study discussions throughout this work. further analysis should be done on the effects of lossy compression in ambient noise interferometry for imaging

#### Conclusions

New technologies to continuously collect high-resolution seismic data for long periods of time are pushing us to consider lossy compression as a means of reducing data movement tigger, A. H., D. M. Hammerling, S. A. Mickelson, H. Xu, M. B. Stolpe, archival requirements and processing or visualization time (particularly when interactive workflows are desirable) this report, we compare the benefits and drawbacks of wavelet commate simulation data within a large ensemble, Geosci Model pressionzfp compression and SVD compression at compression ratios ranging between 5× and 1007 hese are tested on two public DAS datasets: an urban dark fiber experiment as wellet transforms,in SEG TechnicalProgram Expanded Abstracts, as a geothermalield microseismicity monitoring experiment. errors at low-compression rates versus high-compression rates. Coleman, T. (2016). Brady's geothermal field—metadata for DTS and and we compare these errors because they propagate through an 1261983/ (last accessed February 2023). entire template matching microseismicity detection workflow. Diffenderfer, J., A. L. Fox, J. A. Hittinger, G. Sandersand P. G.

#### Data and Resources

Spica et al. (2023). The Brady's Hot springs data are publicly available PA, P. L., R. A. Ergas, and R. S. Polzer (1999). Development the U.S. Department of Energy's Geothermal Data Repository, which is seismic data compression methodsfor reliable, low-noise, provided alongside example code notebodics accessing the data through Amazon Web Services (AWS) discussed by Coleman (2016) Society of Exploration Geophysicists 1903-1906, doi: 10.1190/ and Feigl et al. (2016). The Python software to reproduce tests and fig<sub>1.1820919</sub>. ures from this article is publicly available at https://github.com/aissapeigl, K., N. Taverna, and M. Rossol (2016). Porotomo natural Issah-SRL-compression-2023.git and is archived at https://zenodo.org/aboratory horizontal and vertical distributed acoustic sensing badge/latestdoi/505936568 in its prestent. This software is built using the following open-sourcepackagesPyWavelets(Lee et al., 2019), ZFPy (Lindstrom, 2014), Numpy, Matplotlib, and h5py. The sumbons, S. J., and F. Ringdal (2006). The detection of low magnitude pressed data for the various compression types considered in the articlet. 165, no. 1, 149-166, These include the time domain data from the Errors in frequency conalko, N., P. G. Martinsson, and J. A. Tropp (2011). Finding structure tent section reconstructed after lossy compression and the channelwis@ith randomness: probabilistic algorithms for constructing normalized cross correlation deach channel's continuous recording with its template event recording used in the event detection analyses.288, doi: 10.1137/090771806. These are pictured for all three compression types consider Allhere Laney, D., S. Langer, C. Weber, P. Lindstrom, and A. Wegener (2013). websites were last accessed in September 2023.

### **Declaration of Competing Interests**

The authors acknowledge that there are no conflicts of interest Analysis, Denver, Colorado, 17-22 November 2013j-12. recorded.

### Acknowledgments

The authors thank the NationaScience Foundation Grant Number Li, Z., and Z.Zhan (2018)Pushing the limit of earthquake detection 2046387,the Air Force Research Laboratory Subaward Number

62681767-227888 through Stanford University, and the sponsor com-(ClARC) for computing resources. The authors also thank the students and faculty of CWP and the Martin Group for helpful

#### References

Averbuch, A. Z., F. Meyer, J.-O. Stromberg, R. Coifman, and A. Vassiliou (2001).Low bit-rate efficient compression for seismic data, IEEE Trans. Image Process10, no. 12, 1801-1814,doi: 10.1109/83.974565

P. Naveau, B. Sanderson, Ebert-Uphoff, S. Samarasinghe, De Simone, et al. (2016). Evaluating lossy data compression on cli-Dev. 9, no. 12, 4381-4403.

Bosman, C., and E. Reiter (1993). Seismic data compression using wave-Washington, DC, 26-30 September 1995, ociety of Exploration Geophysicists, 261-1264 doi: 10.1190/1.1822354.

Lindstrom (2019). Error analysis of zfp compression for floating-point data, SIAM J. Sci. Comput. 41, no. 3, A1867-A1898. The Penn State Fiber-Optic foR Environment SEnsEing (FORESEE) of the Penn State Fiber-Optic for Environment SEnsEing (FORESEE) of the Penn State Fiber-Optic for Environment SEnsEing (FORESEE) of the Penn State Fiber-Optic for Environment SEnsEing (FORESEE) of the Penn State Fiber-Optic for Environment SEnsEing (FORESEE) of the Penn State Fiber-Optic for Environment SEnsEing (FORESEE) of the Penn State Fiber-Optic for Environment SEnsEing (FORESEE) of the Penn State Fiber-Optic for Environment SEnsEing (FORESEE) of the Penn State Fiber-Optic for Environment SEnsEing (FORESEE) of the Penn State Fiber-Optic for Environment SEnsEing (FORESEE) of the Penn State Fiber-Optic for Environment SEnsEing (FORESEE) of the Penn State Fiber-Optic for Environment SEnsEing (FORESEE) of the Penn State Fiber-Optic for Environment SEnsEing (FORESEE) of the Penn State Fiber-Optic for Environment SEnsEing (FORESEE) of the Penn State Fiber-Optic for Environment SEnsEing (FORESEE) of the Penn State Fiber-Optic for Environment SEnsEing (FORESEE) of the Penn State Fiber-Optic for Environment SEnsEing (FORESEE) of the Penn State Fiber-Optic for Environment SEnsEing (FORESEE) of the Penn State Fiber-Optic for Environment SEnsEing (FORESEE) of the Penn State Fiber-Optic for Environment SEnsEing (FORESEE) of the Penn State Fiber-Optic for Environment SEnsEing (FORESEE) of the Penn State Fiber-Optic for Environment SEnsEing (FORESEE) of the Penn State Fiber-Optic for Environment SEnsEing (FORESEE) of the Penn State Fiber-Optic for Environment SEnsEing (FORESEE) of the Penn State Fiber-Optic for Environment SEnsEing (FORESEE) of the Penn State Fiber-Optic for Environment SEnsEing (FORESEE) of the Penn State Fiber-Optic for Environment SEnsEing (FORESEE) of the Penn SEn

are publicly available on PubDAS, which can be accessed through Glophesm. Theory 41,no. 3, 613-627,doi: 10.1109/18.382009. performance, in SEG Technical Program Expanded Abstracts,

> data, available at https://gdr.openei.org/submissions/98/ast accessed February 2023).

plemental material contains more detailed representations of the com-seismic events using array-based waveform correlation, Geophys. J.

approximate matrix decomposition SIAM Rev.53, no. 2, 217-

Assessing the effects oflata compression in simulations using physically motivated metricsin Proc. of the International Conf. on High PerformanceComputing, Networking, Storage and

Lee, G. R., R. Gommers, F. Waselewski, K. Wohlfahrt, and A. O'Leary (2019). Pywavelets: A python package for wavelet analysis, J. Open Source Software 40. 36, 1237, doi: 10.21105/joss.01237.

with distributed acoustic sensing and template matchingcase

- study at the Brady geothermafield, GeophysJ. Int. 215, no. 3, 1583-1593doi: 10.1093/gji/ggy359.
- Rev.Earth Planet.Sci.49, no. 1, 309-336,doi: 10.1146/annurevearth-072420-065213.
- Lindstrom, P. (2014). Fixed-rate compressed floating-point arrays, IEEE impression using high-dimensional wavelet transform proc. Trans. Vis. Comput. Gr. 20, no. 12, 2674–2683doi: 10.1109/ TVCG.2014.2346458.
- Mallat, S. G. (2009). A Wavelet Tour of Signal ProcessingThe SparseWay, Third Ed., Elsevier/Academi@ress,Amsterdam, Boston.
- Martin, E. R. (2019). A scalable algorithm for cross-correlations of compressed ambient seismic noise, in SEG Technical Program Expanded Abstracts, San Antonio, Texas, 15–20 September 20129u, T., and D. J. Stensrud (2019)Characterizing thunder-induced Society of Exploration Geophysicists doi: 10.1190/segam2019-3216637.1.
- Hammerling (2020). A statistical analysis of lossily compressed cli-ment dynamics by telecommunication fiber-optic sensors: an mate modeldata, Comput. Geosci 145, 104, 599.
- Spica Z. J., J. B. Ajo-Franklin, G. C. Beroza B. Biondi, F. Cheng, B. Gaite, B. Luo, E. R. Martin, J. Shen, C. Thurber, et al. (2023). PubdasA public distributed acoustic sensing datasets repository for geoscience SeismolRes Lett. 94, no. 2A, 983-998.

- Valentine, A. P., and J. Trampert (2012). Data space reduction, quality assessment and searching of seismograms: Autoencoder networks Lindsey, N. J., and E. R. Martin (2021). Fiber-optic seismology, Annu. for waveform data, Geophys J. Int. 189, no. 2, 1183–1202, doi: 10.1111/j.1365-246X.2012.05429.x.
  - Villasenor, J. D., R. A. Ergas and P.L. Donoho (1996) Seismic data of the Data Compression Conference—DCCS96wbird, Utah, U.S.A.,31 March-3 April 1996, IEEE Computer Society Press, 396-405.doi: 10.1109/DCC.1996.488345.
  - Wade, D. (2020). Seismic-ZFPF ast and efficient compression and decompression of seismic data, in First EAGE Digitalization Conferenceand Exhibition, Vienna, Austria, November 30-December 3 2020, 5, doi: 10.3997/2214-4609.202032080.
    - ground motions using fiber-optic distributed acousticsensing array, J. GeophysRes.124, no. 23, 12,810-12,823.
- Poppick, A., J. Nardi, N. Feldman, A. H. Baker, A. Pinard, and D. MZhu, T., J. Shen, and E. R. Martin (2021). Sensing earth and environurban experimentin Pennsylvania, USA, Solid Earth 12, no. 1, 219-235,doi: 10.5194/se-12-219-2021.

Manuscript received 18 September 2023 Published online 2 January 2024