

Remote recovery of audio signals from videos of optical speckle patterns: a comparative study of signal recovery algorithms

CONCETTA BARCELLONA,^{1,3} DONATUS HALPAAP,² PABLO AMIL,²
ARTURO BUSCARINO,¹ LUIGI FORTUNA,¹ JORDI TIANA-ALSINA,²
AND CRISTINA MASOLLER^{2,4} 

¹*Dipartimento di Ingegneria Elettrica Elettronica e Informatica, University of Catania, 95125, Catania, Italy*

²*Departament de Física, Universitat Politècnica de Catalunya, Rambla de Sant Nebridi 22, 08222, Terrassa, Barcelona, Spain*

³*concetta.barcellona@unict.it*

⁴*cristina.masoller@upc.edu*

Abstract: Optical remote sensors are nowadays ubiquitously used, thanks to unprecedented advances in the last decade in photonics, machine learning and signal processing tools. In this work we study experimentally the remote recovery of audio signals from the silent videos of the movement of optical speckle patterns. This technique can be used even when in between the source and the receiver there is a medium that does not allow for the propagation of sound waves. We use a diode laser to generate a speckle pattern on the membrane of a loudspeaker and a low-cost CCD camera to record the video of the movement of the speckle pattern when the loudspeaker plays an audio signal. We perform a comparative analysis of six signal recovery algorithms. In spite of having different complexity and computational requirements, we find that the algorithms have (except for the simplest one) good performance in terms of the quality of the recovered signal. The best trade-off, in terms of computational costs and performance, is obtained with a new method that we propose, which recovers the signal from the weighted sum of the intensities of all the pixels, where the signs of the weights are determined by selecting a reference pixel and calculating the signs of the cross-correlations of the intensity of the reference pixel and the intensities of the other pixels.

© 2020 Optical Society of America under the terms of the [OSA Open Access Publishing Agreement](#)

1. Introduction

Optical remote sensing is a key technology that allows to probe the environment noninvasively and over long distances between probe and probed objects. Popular techniques include lidar, hyperspectral imaging and near-infrared spectroscopy. Applications of optical remote sensing span such diverse topics as atmospheric boundary layer assessments, power grid and infrastructure monitoring, offshore wind profile measurements, agricultural, fishing, forest and biodiversity monitoring, archaeological discovery, biomedicine, lighting and many more [1–10].

Diode lasers are particularly attractive for sensing applications due to the fact that they are compact, energy-efficient, low-cost, and cover a wide range of wavelengths. Because of these advantages, here we implement the setup proposed by Rodriguez-Cobo et al [11] and Robles-Urquijo et al [12] for the remote recovery of audio signals from videos of speckle patterns [13], but using a diode laser. In our setup the output of the laser is projected to the membrane of a loudspeaker and the movement of the speckle pattern when the speaker plays a sound is video-recorded with a CCD camera.

Our goal is to perform a comparative study of different signal processing algorithms for the recovery of the audio signal from the silent video. We propose new recovery algorithms and

compare with those that have been proposed in the literature [12,14]. The technique that was proposed by researchers at MIT [14] allows to extract sound from videos of minute vibrations. It does not use laser-generated speckle, and therefore, the video scene needs to be naturally illuminated (in contrast, here we illuminate the video scene with a laser from a distance, as in [12]). Another technique proposed here uses the algorithm IsoMap [15] for dimensionality reduction, that reconstructs time-varying features in a low-dimensional space. This method is very reliable and precise, but has high a computational cost. A simple alternative is the calculation of the mean intensity of each video frame, and we find that it gives, in most cases, good results. We also propose a new algorithm that is based on the cross-correlation between the time series of the intensity of a reference pixel and the time series of the intensity of all other pixels. This algorithm gives a good trade-off between the quality of the recovery signal and the computational cost. The method can be further improved by considering as reference pixel the one with the largest variability.

The paper is organized as follows: in Sec. 2 the experimental setup is described; in Sec. 3 six methods to process the video frames are presented, which are referred to as Differential Processing (DIF) [11], Mean Intensity (MI), 1-pixel method (1 PIX), Cross-Correlation (CC), Machine Learning (ML) using dimensionality reduction and MIT method (MIT) [14]; Sec. 4 presents the measures used to compare the performance of the algorithms; Sec. 5 presents the results of the comparison and Sec. 6 draws the conclusions of the paper.

2. Experimental setup

The setup is shown in Fig. 1(a). We directed the collimated light from a diode laser (Thorlabs HL6750MG, $\lambda = 685$ nm) driven by a Thorlabs ITC4001 controller, through an iris of approximately 1 mm diameter at an angle of 40 deg onto the membrane of a loudspeaker (Logitech X-210 subwoofer). A white paper was glued to the speaker membrane to avoid the absorption of the light by the black surface of the speaker membrane. We imaged the speckle pattern that originates from the roughness of the paper by a lens (focal distance $f = 100$ mm) onto the sensor of an 8-bit CMOS camera (IDS UI-1222LE-M, pixel size $6 \mu\text{m}$). The diameter of the lens was $D \sim 50$ mm and thus the angular resolution was $1.22\lambda/D = 1.590 \times 10^{-5}$ rad. An example of a speckle pattern is shown in Fig. 1(b).

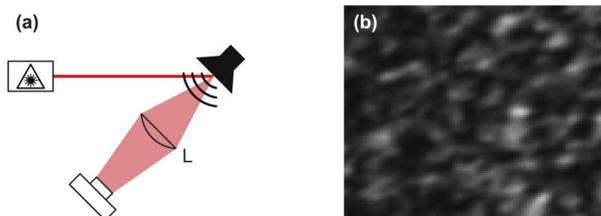


Fig. 1. (a) Experimental setup: A collimated laser beam shines onto a white paper glued to the membrane of a loudspeaker. The paper surface roughness generates a speckle pattern that is imaged by lens L onto the sensor of a CMOS camera. (b) Example a speckle image recorded by the camera. Examples of recorded videos and recovered audio signals can be found [here](#).

The camera is run in video mode in order to observe the dynamical change of the speckle pattern when a sound is played by the speaker. A high video frame rate is needed to reconstruct high-frequency signals because the frame rate corresponds to the sampling frequency, and a sampling frequency of $2f$ is necessary to reconstruct a signal of frequency f . At 800 frames per second (fps), for example, it is possible to recover audio frequencies of up to 400 Hz.

We reconstruct from video two types of audio signals: sinusoidal with frequencies between 60 Hz and 140 Hz and a song with a predominance of bass frequencies (Another One Bites the Dust by Queen). Examples of recorded videos and recovered audio signals can be found [here](#). For the sinusoidal signals we used 360.62 fps and recorded an area of 100×100 pixels; for the song, the recorded area was 560×16 pixels and the frame rate was 886.05 fps.

3. Recovery methods

To describe the signal recovery methods we consider the following notation: $I_i[n]$ is the intensity of the i -th pixel in frame n and $x[n] = x(nT_s)$ is the sampling of the audio signal $x(t)$ emitted by the speaker, with T_s being the sampling time. We use $y[n]$ to denote the recovered signal at time nT_s , that is a function of the pixel intensities, which are functions of the audio signal. Neglecting delays, constant terms, and using a linear approximation we have:

$$I_i[n] = f_i(x[n]) + \sigma_i \epsilon_i[n] \approx \alpha_i x[n] + \sigma_i \epsilon_i[n]. \quad (1)$$

Here α_i is a constant in time that can depend on the pixel and the signals $\epsilon_i[n]$ account for noise and intensity discretization; we assume they are uncorrelated white noise with unity variance. σ_i is the strength of the noise that can also depend on the pixel.

3.1. Differential processing scheme (DIF)

This method, used in [11], assumes a proportionality between speckle pattern variations and the perturbation to be measured. It is based in computing the difference, pixel-to-pixel, between two consecutive frames. The recovered signal is $y[n] = \sum_i |I_i[n] - I_i[n-1]|$. If Eq. (1) holds we see that

$$y[n] = \sum_i |\alpha_i (x[n] - x[n-1]) + \sigma_i (\epsilon_i[n] - \epsilon_i[n-1])|. \quad (2)$$

Therefore, $y[n]$ does not vary linearly with $x[n]$, which leads to harmonic generation in the recovered signal due to the nonlinearity of the reconstruction.

3.2. Mean intensity (MI)

The recovered signal is the mean intensity of frame n , $y[n] = \frac{1}{N} \sum_i I_i[n]$, with N being the number of pixels. If Eq. (1) holds:

$$y[n] = \left(\sum_i \frac{\alpha_i}{N} \right) x[n] + \frac{1}{N} \sum_i \sigma_i \epsilon_i[n], \quad (3)$$

which shows that $y[n]$ varies linearly with $x[n]$ but, considering that the α_i values are likely positive or negative, the sum $(\sum_i \frac{\alpha_i}{N})$ will likely be small and the signal, $x[n]$, will likely be hidden by the noise.

3.3. 1-pixel method (1 PIX)

This method consists in choosing a single pixel r and using only the intensity of that pixel to reconstruct the signal as $y[n] = I_r[n]$. The selected pixel is chosen to maximize the variance of the signal $I_r[n]$, i.e., is the one with the largest variance. If Eq. (1) holds, $y[n] = \alpha_r x[n] + \sigma_r \epsilon_r[n]$. Assuming that the noise strengths are all the same ($\sigma_i = \sigma$), then choosing the pixel with the largest intensity variance corresponds to maximizing the signal to noise ratio (SNR).

3.4. Cross-Correlation (CC)

In this approach we select a reference pixel, r , and compute the Pearson correlation, ρ , between the intensity of each pixel, i , and the intensity of the reference pixel r : $C_i = \rho(I_i[n], I_r[n])$. If Eq. (1) holds and we neglect the noise, then

$$C_i = \text{sign}(\alpha_i \alpha_r). \quad (4)$$

We use this property to find the signs of the coefficients β_i in a weighted mean intensity that is used to reconstruct the signal,

$$y[n] = \sum_i \beta_i I_i[n]. \quad (5)$$

We select the weights such that they maximize the SNR. Using Eq. (1) and assuming that the noise strength is the same for all the pixels ($\sigma_i = \sigma$) we have

$$y[n] = \left(\sum_i \beta_i \alpha_i \right) x[n] + \sigma \sum_i \beta_i \epsilon_i[n]. \quad (6)$$

If the noise signals $\epsilon_i[n]$ are uncorrelated, then $\text{SNR} \propto (\sum_i \beta_i \alpha_i) / \sqrt{\sum_i \beta_i^2}$. Therefore, SNR will be maximized if we chose the weights such that $\beta_i \propto \alpha_i$.

To estimate α_i we use the Root Mean Square (RMS) of the signal $I_i[n]$, $\text{RMS}_i = \sqrt{\frac{1}{L} \sum_{n=1}^L I_i^2[n]}$ (with L being the number of frames in the video), that with the linear approximation is $\text{RMS}_i = \sqrt{P^2 \alpha_i^2 + \sigma^2}$, where P is the RMS power of the signal $x[n]$. From here we see that $|\alpha_i| \propto \sqrt{\text{RMS}_i^2 - \sigma^2}$.

Assuming that there are values of α_i that are small enough to be neglected, $\min_i (\text{RMS}_i)$ is an estimation of the noise strength, $\sigma \sim \min_i (\text{RMS}_i)$, and therefore, $|\alpha_i| \propto \sqrt{\text{RMS}_i^2 - \min_i^2 (\text{RMS}_i)}$.

As SNR will be maximized if the weights are $\beta_i \propto \alpha_i$, we chose $|\beta_i| = \sqrt{\text{RMS}_i^2 - \min_i^2 (\text{RMS}_i)}$. The last step is to use the sign of the correlation C_i to set the sign of β_i such that

$$\beta_i = \text{sign}(C_i) |\beta_i| = \text{sign}(C_i) \sqrt{\text{RMS}_i^2 - \min_i^2 (\text{RMS}_i)}. \quad (7)$$

Using Eq. (4), $\beta_i = \text{sign}(\alpha_i \alpha_r) |\beta_i|$ and substituting in Eq. (6) we obtain

$$y[n] = \text{sign}(\alpha_r) \left(\sum_i |\alpha_i \beta_i| \right) x[n] + \sum_i \text{sign}(\alpha_i \alpha_r) |\beta_i| \sigma_i \epsilon_i[n], \quad (8)$$

We note that (in contrast with the unweighted mean intensity), by using the weights defined in Eq. (7) all the terms in the sum that multiplies $x[n]$ are positive. To obtain a good performance with this method we need to choose a reference pixel with high RMS, otherwise, the signs of the cross-correlation coefficients (used to set the signs of the weights) will be mainly random and not informative.

3.5. Machine learning method (ML)

In this approach each frame n is described by a point in a N -dimensional space whose coordinates are the values of the N pixels, $\{I_1[n] \dots I_N[n]\}$. The recorded video constituted by L frames is then represented by a sequence of L points in this high dimensional space (i.e., a path or trajectory). The points are in fact located on a low dimensional manifold because the values of the pixels depend on a low number of variables, such as the position ($x(t)$) and the velocity of

the membrane of the loudspeaker (expressed as the derivative of the position ($x'(t)$). To recover the coordinates of each frame in the low dimensional manifold we used a well-known algorithm for manifold learning known as IsoMap [15]. We fed the algorithm with the coordinates of each frame, $\{I_1[n] \dots I_N[n]\}$, and used the Euclidean pixel-to-pixel distance to calculate the distance between any two frames. For each frame the algorithm returns a set of m “features”, $\{g_1[n] \dots g_m[n]\}$ ($m \ll N$), which are the coordinates of the frame in the reduced space. The audio signal can then be recovered from the trajectory in the reduced space. Here we use $m = 2$ as we expect that the intensities $I_i[n]$ depend on both $x(t)$ and $x'(t)$ due to the time integration made by the camera and the rolling shutter capturing method. We use the first feature to recover the signal, $y[n] = g_1[n]$, because the features are ordered in decreasing order of power, and we expect that the dependence to the position is the one that generates the highest power.

This method has the advantage that it can deal with nonlinearities (as the IsoMap method is capable of learning nonlinear manifolds), but has the disadvantage of having a high computational cost as the computing time increases quadratically with the number of frames being analyzed.

3.6. MIT method

This method, described in detail in [14], is based on the complex steerable pyramid [16,17] which analyzes the motion of each pixel of the video in a certain direction and scale. Then a weighted average is performed for each direction and scale, and all directions and scales are synchronized (with a time-shift) to an arbitrary reference and added together to generate the recovered signal. In [14] the authors also proposed post-processing the recovered signal by using denoising algorithms. In order to perform a fair comparison with the other methods the denoising algorithms have not been used.

4. Performance measures

To compare the performance of the different methods, for the sinusoidal audio signals we use the Signal-to-Noise and Distortion ratio (SINAD) [18], while for the song, we use the Pearson cross-correlation coefficient.

The SINAD measures (in dB) the ratio between the power of the signal and the rest of the power that corresponds to noise and distortion. In the Fast Fourier Transform (FFT), the power of

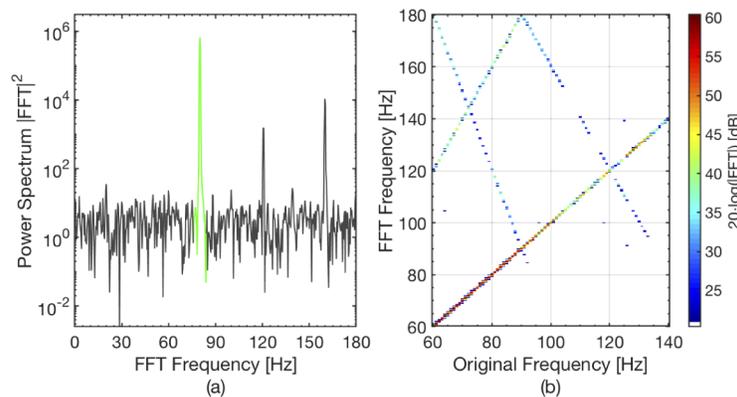


Fig. 2. Evaluation of the Signal-to-Noise and Distortion ratio (SINAD) from the Fast Fourier Transform (FFT) of the recovered signal from the MI method: (a) the green line indicates the signal, i.e., the section of FFT centered at the applied frequency. In this example the applied frequency is $f_0 = 80$ Hz. (b) Spectrogram obtained with the Mean Intensity method that illustrates the appearance of additional peaks due to aliasing and harmonics.

the signal is within a small frequency range centered at the frequency f_0 of the sinusoidal signal that is applied to the speaker, see Fig. 2(a), while all the other power is the noise and distortion:

$$SINAD = \frac{P_{signal}}{P_{nd}} \quad (9)$$

In the FFT, in addition to a peak at f_0 , peaks at other frequencies occur for all the methods due to harmonics and aliasing. An example, that corresponds to the MI method, of how the peaks vary with f_0 is shown in Fig. 2(b). The power of these peaks is included in the noise and distortion term, P_{nd} .

5. Results

We first compare the performance of the recovery methods when a sinusoidal signal is applied to the speaker. Figure 3(a) shows that the performance of each method (except for the Differential Processing scheme) increases with the volume of the audio signal, but if the volume is too high, the performance decreases. This decrease is likely due to a nonlinear response of the speaker that produces sudden variations on the amplitude of the recovered signal with respect to the increase/decrease of the frequency of the input signal [19]. The performances of the Machine learning (ML) and Cross-correlation (CC) methods are comparable and represent the best choices. The lowest performance is obtained with the Differential Processing scheme (DIF). This can be understood because, as it was explained in Sec. 3.1, the signal recovered with this method is not a linear function of the applied signal, which leads to harmonic generation.

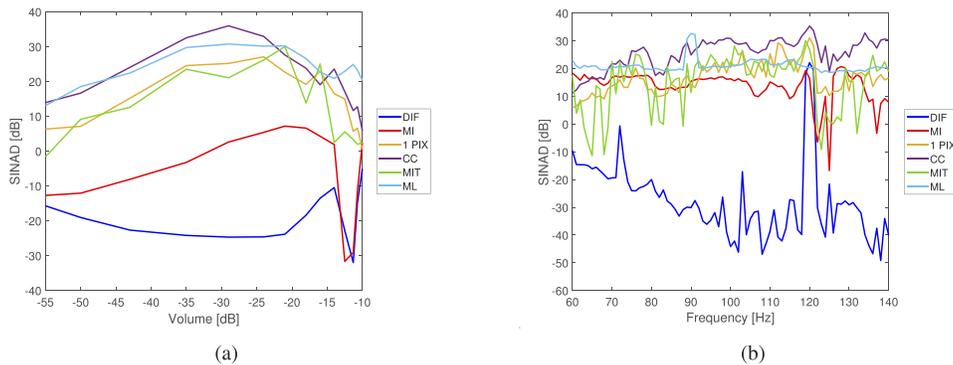


Fig. 3. (a) Performance of the recovery algorithms quantified with the Signal-to-Noise and Distortion ratio (SINAD) when a sinusoidal signal is played by the loudspeaker and the volume of the signal is increased, while the signal frequency is kept constant (100 Hz). (b) SINAD for increasing signal frequency when the signal volume is kept constant (-30dB, dB refers to sound pressure levels).

The results obtained when increasing the frequency of the signal (while keeping constant the volume), shown in Fig. 3(b), are consistent and confirm that ML and CC give the best results. Very similar performances are obtained when the frequency is decreased (not shown).

In Table 1 we summarise the computational costs of the methods. While DIF is the method with the lowest computational cost, considering both, performance and computational cost, we can conclude that the best trade-off is obtained with the Cross-Correlation method (CC) because it has the highest SINAD while the computational time is reasonably low.

Finally the performance of the recovery methods when a song is played by the speaker is quantified by the cross-correlation and the results are presented in Table 2. We note that for low volume (-39.9 dB) the methods MIT, ML, CC and 1 PIX perform reasonably well since the

Table 1. Computational costs of the recovery algorithms in terms of mean time of execution and standard deviation. All the methods were run using MatLab in a portable computer with an Intel i7-7700HQ processor and 16 GB of RAM.

	1 PIX	DIF	MI	MIT	ML	CC
Mean Time [s]	0.156	0.024	0.207	5.789	54.23	0.206
σ [s]	0.009	0.004	0.006	0.173	1.104	0.011

correlation coefficient is $\rho \sim 0.42$. For moderate volume, the MIT, CC and 1 PIX correlation coefficients decrease slightly (to $\rho = 0.35$) while the ML method keeps the same performance as for low volume. For the highest volume, the ML method outperforms the other ones. The DIF and MI are not suitable methods to recover the original song, regardless of the signal volume.

Table 2. Performance of the recovery algorithms quantified with the cross-correlation when a song is played by the loudspeaker, for three volume levels.

	1 PIX	DIF	MI	MIT	ML	CC
-39.9 dB	0.4198	0.0415	0.2079	0.4321	0.4282	0.4308
-26.6 dB	0.3566	0.0591	0.1691	0.3032	0.4313	0.3540
-19.7 dB	0.1392	0.0601	0.0982	0.3084	0.7942	0.0972

We note that the performance of the 1 PIX and CC methods decays with the volume of the speaker. Because the methods are aimed at recovering weak signals, their performance is compared for small amplitudes.

Like in other diffraction experiments, the size of the speckle in the pattern depends on several factors: the wavelength of the light, the roughness of the speckle-generating surface as well as the imaging optics. For different imaging geometries, the size of the speckle might change, but as long as we make sure that the imaged speckle grains have at least the size of several pixels, we do not expect changes in the relative performance of the described methods.

6. Conclusions

We have performed an experimental study of the remote recovery of audio signals from the silent videos of the movement of optical speckle patterns. We have compared six signal recovery algorithms in terms of the quality of the recovered signal and the computational costs. We have considered the differential processing method used in [11,12] which is simple to implement, but introduces undesired nonlinearities. We have analyzed how the mean intensity of the speckle pattern can recover the audio signal and have shown that a single observed pixel, if it is appropriately selected (the pixel with largest variance) can be sufficient for obtaining a reasonably good recovery of the audio signal. We have proposed two new recovery methods, one is based on a weighted sum of all the pixels, where the signs of the weights are determined by the cross-correlation (CC) of the intensity of the different pixels; the other method is based on a machine learning (ML) algorithm for dimensionality reduction. We have compared these methods with that proposed in [14]. We have found that the best trade-off, in terms of computational costs and performance, is the Cross-Correlation method (CC).

As shown in Fig. 3, some methods show a significant lower SINAD for certain frequencies. It would be interesting, for future work, to investigate if this effect can be linked to the occurrence of nonlinear resonance effects, such as the jump resonance [19] displayed by a nonlinear, non-autonomous system whose frequency response is characterised by one or more hysteresis regions, depending on the increasing/decreasing trend of the frequency at the input signal.

Our optical technique for remote sound recover can allow to make non-contact vibration measurements with the advantage of being an inexpensive setup that uses a low-cost diode laser and a low-cost CCD camera.

Funding

Ministerio de Ciencia, Innovación y Universidades (PGC2018-099443-B-I00); Institució Catalana de Recerca i Estudis Avançats (Academia); European Commission (675512).

Acknowledgments

C.B., D.H. and P.A. contributed equally to this work.

Disclosures

The authors declare no conflicts of interest.

References and links

1. S. Emeis, K. Schäfer, and C. Münkel, "Surface-based remote sensing of the mixing-layer height a review," *Meteorol. Zeitschrift* **17**(5), 621–630 (2008).
2. Z. Zalevsky, Y. Beiderman, I. Margalit, S. Gingold, M. Teicher, V. Mico, and J. Garcia, "Simultaneous remote extraction of multiple speech sources and heart beats from secondary speckles pattern," *Opt. Express* **17**(24), 21566–21580 (2009).
3. Y. L. Pichugina, R. M. Banta, W. A. Brewer, S. P. Sandberg, and R. M. Hardesty, "Doppler lidar-based wind-profile measurement system for offshore wind-energy and other marine boundary layer applications," *J. Appl. Meteor. Climatol.* **51**(2), 327–349 (2012).
4. V. Klemas, "Fisheries applications of remote sensing: an overview," *Fish. Res.* **148**, 124–136 (2013).
5. C. Kuenzer, M. Ottinger, M. Wegmann, H. Guo, C. Wang, J. Zhang, S. Dech, and M. Wikelski, "Earth observation satellite sensors for biodiversity monitoring: potentials and bottlenecks," *Int. J. Remote. Sens.* **35**(18), 6599–6647 (2014).
6. W. Emery and A. Camps, *Introduction to Satellite Remote Sensing* (Elsevier, 2017), first edition ed.
7. M. Eady, B. Park, and S. Choi, "Rapid and early detection of salmonella serotypes with hyperspectral microscopy and multivariate data analysis," *J. Food Protection* **78**(4), 668–674 (2015).
8. S. C. Murray, "Optical sensors advancing precision in agricultural production," *Photon. Spectra* **51**, 48 (2018).
9. C. A. Hostetler, M. J. Behrenfeld, Y. Hu, J. W. Hair, and J. A. Schullien, "Spaceborne lidar in the study of marine systems," *Annu. Rev. Mar. Sci.* **10**(1), 121–147 (2018).
10. N. Masini, F. T. Gizzi, M. Biscione, V. Fundone, M. Sedile, M. Sileo, A. Pecci, B. Lacovara, and R. Lasaponara, "Medieval archaeology under the canopy with lidar. the (re)discovery of a medieval fortified settlement in southern Italy," *Remote Sens.* **10**(10), 1598 (2018).
11. L. Rodriguez-Cobo, M. Lomer, and J. M. Lopez-Higuera, "Fiber specklegram-multiplexed sensor," *J. Lightwave Technol.* **33**(12), 2591–2597 (2015).
12. I. Robles-Urquijo, M. Lomer, L. Rodriguez-Cobo, and J. M. Lopez-Higuera, "Non-contact vibration analysis using speckle-based techniques," in *2017 25th Optical Fiber Sensors Conference (OFS)*, (2017), pp. 1–4.
13. J. W. Goodman, *Speckle Phenomena in Optics: Theory and Application* (Roberts and Company Publishers, 2007).
14. A. Davis, M. Rubinstein, N. Wadhwa, G. Mysore, F. Durand, and W. T. Freeman, "The visual microphone: Passive recovery of sound from video," *ACM Trans. Graph.* **33**(4), 1–10 (2014).
15. J. B. Tenenbaum, V. De Silva, and J. C. Langford, "A global geometric framework for nonlinear dimensionality reduction," *Science* **290**(5500), 2319–2323 (2000).
16. J. Portilla and E. P. Simoncelli, "A parametric texture model based on joint statistics of complex wavelet coefficients," *Intl. J. Comput. Vision* **40**(1), 49–70 (2000).
17. E. P. Simoncelli, W. T. Freeman, E. H. Adelson, and D. J. Heeger, "Shiftable multiscale transforms," *IEEE Trans. Inf. Theory* **38**(2), 587–607 (1992).
18. W. Kester, "Understand sinad, enob, snr, thd, thd + n, and sfdr so you don't get lost in the noise floor," <https://www.analog.com/media/en/training-seminars/tutorials/MT-003.pdf> (2008). [online; accessed 2020-01-15].
19. A. Buscarino, C. Famoso, L. Fortuna, and M. Frasca, "Multi-jump resonance systems," *Int. J. Control* **93**(2), 282–292 (2020).