

Robot Hearing Through Optical Channel in a Cocktail Party Environment

Xiao Guo¹, Siyi Ding¹, peng ti¹, Kenan Li¹, and Xiaoping Hong¹

¹Southern University of Science and Technology

June 18, 2022

Abstract

The cocktail party problem refers to a challenging process when the human sensory system tries to separate a specific voice from a loud mixture of background sound sources. The problem is much more demanding for machines and has become the holy grail in robotic hearing. Despite the many advances in noise suppression, the intrinsic information from the contaminated acoustic channel remains difficult to recover. Here we show a simple-yet-powerful laser-assisted audio system termed REAL (Robot Ear Accomplished by Laser) to probe the vibrations of sound-carrying surfaces (mask, throat and other nearby surfaces) in optical channel, which is intrinsically immune to acoustic background noises. Our results demonstrate that REAL can directly obtain the audio-frequency content from the laser without acoustic channel interference. The signals can be further transcribed into human-recognizable audio by exploiting the internal time and frequency correlations through memory-enabled neural networks. The REAL system would enable a new way in human-robot interaction.

Xiaoping Hong Email: hongxp@sustech.edu.cn

ToC Figure



The cocktail party problem, hearing a specific speaker in a loudly noisy environment, is extremely difficult for machines and has become the holy grail in audio processing research. A robotic ear termed REAL (Robot Ear Accomplished by Laser) is devised to accurately recognize human voices in cocktail party environments by remotely probing the vibrations of sound-carrying surfaces by optical means.

Introduction

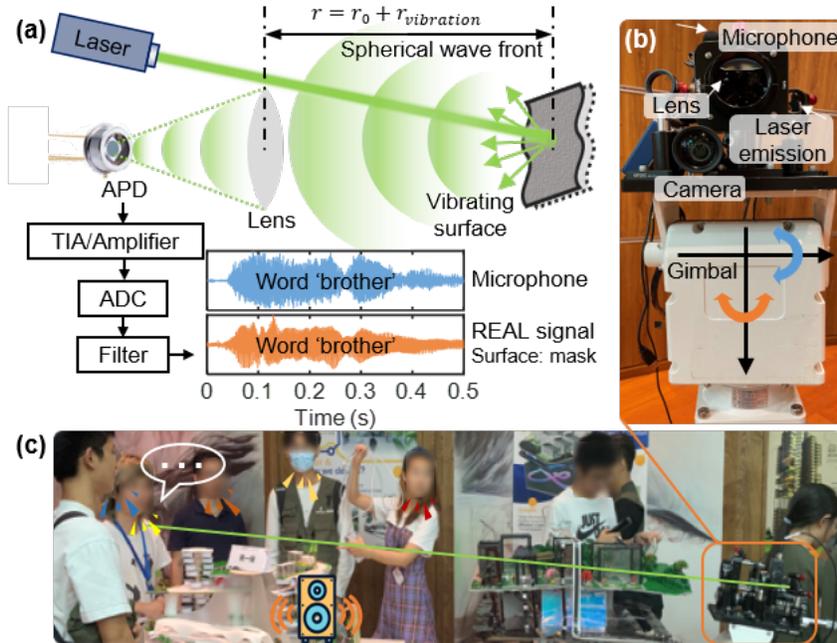
Voice control is the most intuitive and efficient way in human-robot interaction. However, the robot with onboard microphones will be difficult to interact with verbally in a noisy environment, e.g. service robot in a crowded airport. This is an example of cocktail party problem^[1,2] which is the holy grail in robotic hearing

and perception.^[3,4] Various techniques have been applied by mimicking the human sensing mechanism, e.g. speaker localization through microphone arrays^[5-7] similar to human ears^[8], voice feature extraction with algorithms^[9-11] as in human brains^[12] and sensor fusion with audio and visual modalities.^[13,14] These methods alleviated the problem within a limited scope, but the voice from a person at a distance can be much weaker than noises from other nearby sources because sound pressure typically decreases with squared distance. Comparing microphones, optical detection could remotely capture vibration signals from the source while remaining unaffected by the surrounding acoustic noises. As early as 1880, the optical telecommunication pioneer Alexander Graham Bell invented the first apparatus named photophone^[15] that used modulated light to reproduce sound from 231 meters away. Since then, many intriguing applications of remote sound sensing with light were adopted, especially after the invention of laser.^[16] Laser doppler vibrometers (LDV) were developed to remotely probe the surface vibrations with interferences to acquire sounds.^[17-19] Due to the interferometric nature, the sensitivity to detect vibrations from smooth surfaces (spectacular reflections) is satisfactory. However, the detection can be much more complicated for large scattering surfaces where the returned signal is mixed in phases, and in this case, the detection sensitivity is reduced significantly due to speckles (see Supporting Information for analysis on LDV with rough surfaces). Additionally, the sophisticated setup of LDV makes it too costly and bulky to be massively deployed in consumer robots. On the other hand, direct intensity modulation measurements could be more sensitive and cost-effective than interferometric methods. Simpler and lower laser microphones^[20] have also been attempted to accomplish similar functions as LDV by monitoring intensity modulation from specular reflections, however the strict back reflection requires a perpendicular mirror-like surface and is prone to optical misalignment and fluctuations. Recently, Nassi et al. demonstrated a direct passive sound recovery^[21] from a photodetector which is telescoped at a bulb that measures the minute changes of its brightness caused by sound-excited bulb surface vibrations. However, this method involves a necessary illuminating bulb next to the speaker. None of the existing optical techniques are suitable for voice commanding a robot in a cocktail party environment. The REAL system we proposed has better performance on large scattering vocal surfaces, simpler and more affordable construction and greater adaptability. We will illustrate the principle and construction of REAL and demonstrate signals of REAL operating both on the speaker’s facial masks and on their throats respectively. Furthermore, the REAL signal could be transcribed through a memory-enabled neural network to enhanced voice contents in a noisy cocktail party environment. To our best knowledge, REAL is the first optical channel solution with the potential to solve the cocktail party problem.

Figure 1. Robot Ear Accomplished by Laser (REAL) in a cocktail party environment. a) Working principle of the REAL system. The signal is gathered from the scattered laser photons bouncing off a vibrating surface (throat, mask, etc.). b) A gimbalized REAL system containing a camera on a gimbal to continuously point the laser to the tracked target surface. c) REAL system operates in a typical cocktail party environment.

Figure 1a illustrates the working principle of the REAL system. A low-power collimated laser beam is targeted at the remote vibrating surface. A telescoping lens system is aligned with the remote spot to collect the back-scattered photons and focus them onto the avalanche photodetector (APD), from which the AC signal is amplified and processed as the REAL audio (see Methods for REAL signal processing). This signal is the result of back-scattered light intensity change as the surface vibrates. If we assume the scattering surface is Lambertian,^[22] the collected back-scattered laser power will decay with squared distance $P_c \propto P_0/r^2$, where P_0 is the scattering laser power exiting normal to the surface, and r is the distance between the surface spot and the collecting lens. As the surface vibrates, the relative distance changes ($r = r_0 + \Delta r$) and this Δr results in a change in collected laser power $\Delta P_c \sim P_0(1/r^2 - 1/(r + \Delta r)^2)$. As an example, the waveform of the pronounced word ‘brother’ gathered by the REAL system (detecting from the speaker’s mask) is shown in Figure 1a. A microphone waveform is also shown for comparison. Because of its construction simplicity, the cost of the REAL system is much lower than a conventional LDV system and the miniaturization is readily available. (see Supporting Information for comparison of REAL and LDV)

As a robotic ear, the REAL system needs to continuously capture a specific voice source. In Figure 1b, the system is mounted on a motorized gimbal, and a camera is used to detect and track the throat or the mask of the speaking person. The detected target position is fed into the control loop of the gimbal, which points the



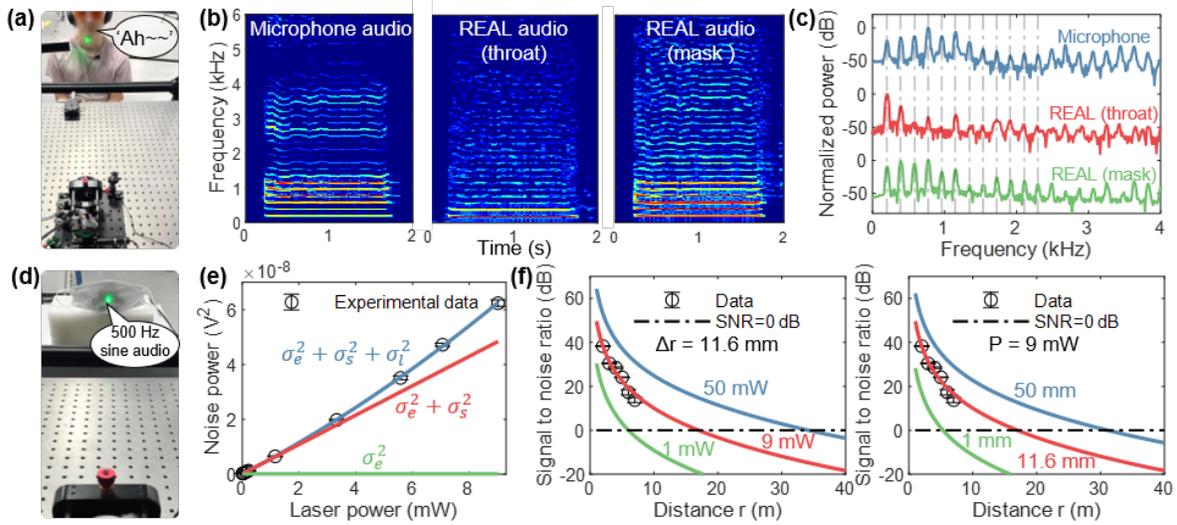
laser to the target continuously as the target is moving. A microphone is attached to REAL system to collect the audio signal for comparison, as well as further augment the REAL signal by fusing the two independent modalities. Figure 1c shows a cocktail party scenario where the gimbaled REAL system operates to ‘hear’ a specific person remotely without acoustic channel interference.

Result and application

Frequency characterization

A systematic investigation on REAL’s frequency responses is needed to justify the new approach. Short-time Fourier transform (STFT) is a commonly used tool to understand both the temporal and frequency response of audio signals.^[23] In **Figure 2a-b**, the STFT responses of the microphone, the REAL signal on the throat and the REAL signal on the mask are obtained when the speaker pronounces ‘Ah~’. Sharp vocal resonances from all three STFT spectrograms could be identified. In Figure 2c, the frequency cross-sections of the signals at time 1 s were plotted for a better comparison. Although the envelopes of the frequency spectra are slightly different, the peaks are prominent and accurate. From the point of view of speech synthesis, these frequency responses already provide enough information to understand the content of the speech.^[24] Notably, the throat signal is reasonably good in low frequencies (< 1 kHz) but slightly lacking in the higher frequencies, while the mask signal is more uniform across the vocal spectrum. This different response is related to the respective elastic and damping properties of the surface materials, and in the case of the human throat, biological tissues attenuate higher frequency contents while transmitting lower-frequency information.^[25] We discovered that many other objects could be used with the REAL system, such as plastic bags, packaging boxes and papers, given that they are close to and excited mainly by the target sound source. A detailed analysis of the frequency response of various materials is presented in the Supporting Information.

Figure 2. Performance characterization of REAL. a) REAL collects audios from the speaker. b) STFT spectrograms of ‘Ah~’ recorded by microphone, REAL on throat and REAL on mask. c) Frequency cross-



section of the signals in b) at time 1s. d) SNR characterization of REAL system from a mask on a loudspeaker playing 500 Hz sine audio. e) The relationship between laser power and total noises, which includes electronic noise (s_e), shot noise (s_s) and laser power fluctuation noise (s_l). f) The REAL SNR as a function of measurement distance, with different laser power (P , left) and vibration magnitude ($[?]r$, right)

Signal-to-noise ratio characterization

Understanding the noise or signal-to-noise ratio (SNR) is one of the most fundamental considerations in constructing an artificial sensory system. In REAL, the noise of the system comes from several sources, namely target movement noise (s_m), electronic noise (s_e), laser fluctuation noise (s_l) and shot noise (s_s). Since the total noise is the geometric mean of all these independent noises, the largest noise source dominates. The target movement noise is due to the body movement during speech and usually has a smaller magnitude in the audio frequency band compared to s_l and s_s (see Methods for SNR analysis of REAL). Electronic noises including the Johnson–Nyquist noise depend on temperature and fabrication of the chips. They are independent of the laser power. Laser fluctuation can play a role and in our case. Our laser fluctuation is also an important factor to consider. In our case, laser power fluctuation within the audio band is measured to be $\sim 0.005\%$. Since the APD signal is proportional to the collected laser power, this noise power is approximately proportional to square of laser power ($s_l^2 \sim P^2$). Additionally, shot noise should be considered which stems from the particle nature of photon and electron detection; the noise power is proportional to laser power ($s_s^2 \sim P$) with a gain. Because the shot noise is the ultimate noise of a given signal, the gain of the APD should be properly set so that the shot noise dominates the system noises for an optimal SNR performance. In Figure 2e, the noise power (s_{total}^2) in this REAL system is measured as a function of the laser power. A noise model considering the respective noises (see Methods for SNR analysis of REAL) was used to fit the experimentally measured total noise. The close-to-linear relationship between the noise power and the laser power is an indication that the shot noise indeed dominates. With this noise model, a fixed amplitude vibration (11.6 mm) from a mask excited at 500 Hz by a loudspeaker (Figure 2d) is used to examine the signal model, i.e. vibration modulated intensity change, and hence the signal to noise ratio. The experimental SNRs detected at different distances are plotted in Figure 2f and our proposed SNR model could well fit the experimental data based on the theory of spherical wave propagation (see Methods for SNR analysis of REAL). Additionally, based on the proposed model, a few other parameter configurations (laser power and vibration amplitude) are calculated in Figure 2f as a guide to help design similar systems. If we set $\text{SNR} = 0$ dB as the detection limit, our current configuration would allow a remote detection range of 17 meters.

Application I: REAL signal from speakers' masks

Wearing masks in public places such as in airports or hospitals (**Figure 3a**) is a new norm since coronavirus disease COVID-19. The mask will cause a noticeable reduction in loudness and clarity of speech^[26] which is more susceptible to background noises. Additionally, masks prevent facial expressions and lipreading,^[27,28] which makes speech understanding more difficult without audio-visual understanding. We demonstrate REAL as an ideal tool to probe the audios in these cases. In Fig 3b, a speaker wearing a mask and face shield in an environment with loud background noise from a loudspeaker. One microphone is placed near the speaker's mouth to collect clear audio as ground truth, and another microphone is placed on the REAL platform to represent the ear of the listener. Figure 3c shows that the microphone ear of the listener completely failed while the audio measured by REAL is similar to the ground truth both in waveform and STFT spectrogram. Audios obtained from REAL on masks could often be understood by humans directly without additional processing (Supplementary Video 1). In many cases, the audios (Supplementary Audio 1 and 2) can be accurately transcribed using speech-to-text services such as the Google Cloud platform (<https://cloud.google.com/speech-to-text>) (see **Figure S3**).

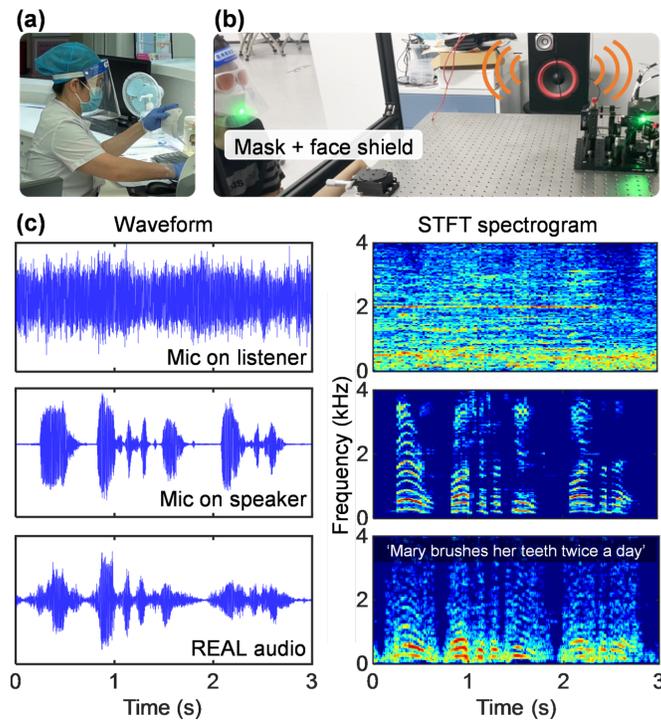


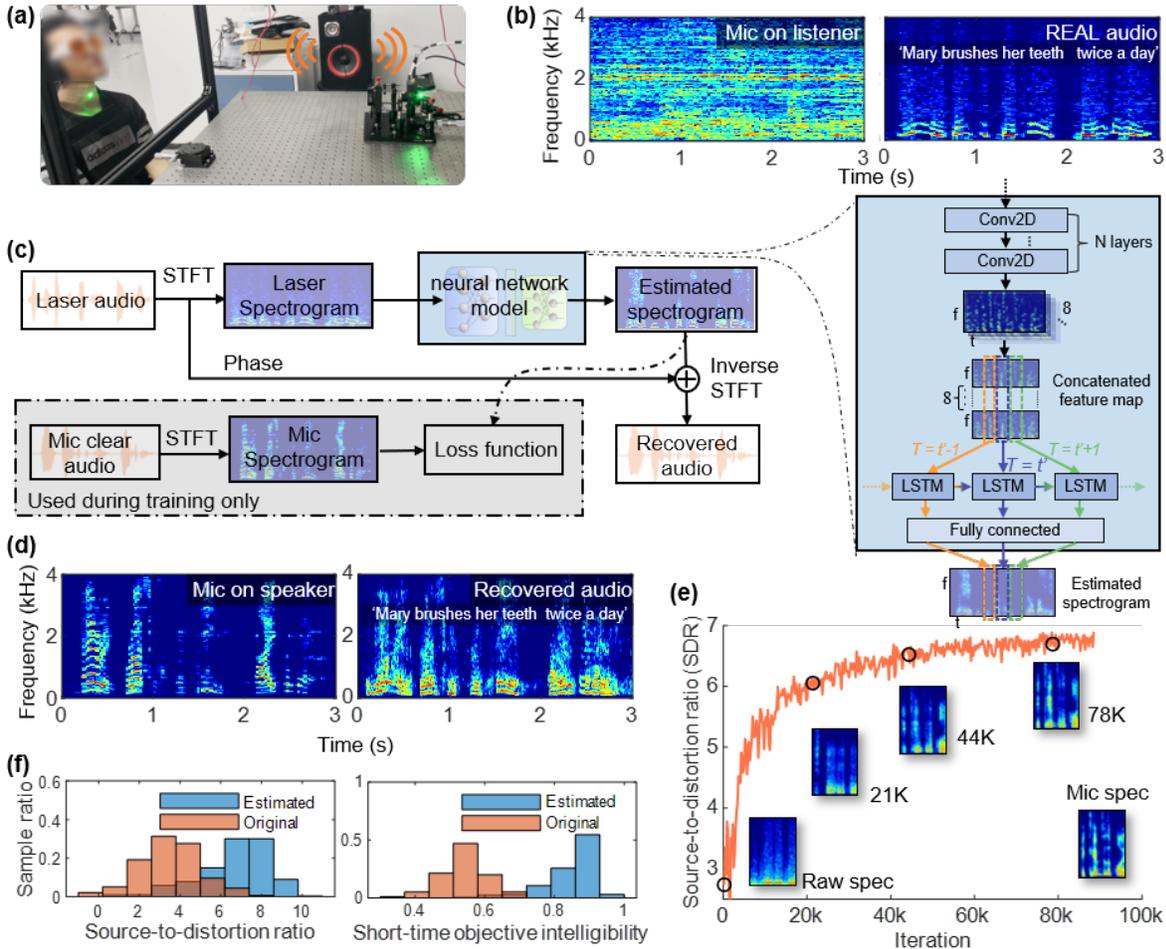
Figure 3. REAL captures audio from masks. a) A typical scenario. b), Speaker wearing the mask (KF94) and face shield speaks in an extremely noisy environment (from the loudspeaker) while REAL measures the mask vibration. c), Waveforms and STFT spectrograms of the audios recorded with a listener microphone, a speaker microphone (ground truth) and REAL signal respectively.

Application II: REAL signal from speakers' throats

More generally, REAL can directly capture the vibrations of the throat surfaces (**Figure 4a**). It is reminded that our laser power is far below the maximum permissible exposure of skin from IEC 60825-1:2014 and can be considered safe. In the simulated cocktail party environment, Figure 4b shows the respective STFT spectrogram from the microphone on the listener and the REAL system (as the listener). The speaker's

speech cannot be distinguished from the listener microphone, while the REAL signal resembles the speaker’s speech only in the lower frequency (Figure 4d) due to the throat filtering. In addition, the REAL audio lacks the unvoiced signal (such as /sh/) because the unvoiced audio was generated in the mouth but not the throat. Direct understanding of the REAL audio without priors is difficult. However, REAL captured information should be adequate to understand the context because intrinsically the position of the vocal resonances, their timing characteristics and the speech pattern are correlated in this multi-modality space. To recover human-understandable speech, we propose a data-driven model (Figure 4c) on the STFT spectrogram to learn the mapping relationship from the REAL signal to the ground truth audio. A convolutional neural network (CNN) is used to capture both the frequency and short-time temporal patterns of human voice commands, with additional long short-term memory (LSTM) to correlate information at longer intervals in speeches. The model is expected to learn how to enhance the high-frequency texture and supplement unvoiced information. The details of the proposed model are described in Methods. The result is presented in Figure 4d, where the recovered audio from REAL is well augmented compared to Figure 4b. A video demonstrating this experiment is provided in Supplementary Video 2. We evaluate the performance using the source-to-distortion ratio (SDR)^[29] and short-time objective intelligibility (STOI),^[30] which are two commonly used evaluation metrics in speech enhancement tasks. As shown in Figure 4e, with iterations a final SDR score of around 6.8 is obtained for the recovered REAL audios in the testing set, representing a significant increase in the content clarity.^[31] Figure 4f shows SDR and STOI histograms of the 293 testing samples (original and recovered), demonstrating the competence and accuracy of the audio recovery model (see Methods for evaluation metrics). Finally, we note that this model is capable of real-time inference with reasonable hardware to ensure REAL could operate with the robot’s onboard computing platform (see Supporting Information for real-time analysis).

Figure 4. REAL captures audio from the human throat. a) REAL captures audio from the speaker’s throat in a noisy environment. b) STFT spectrogram of two audios recorded with microphone (on REAL) and REAL respectively. c) Block diagram of training and inference pipeline of the neural network model to recover REAL audio. d) Spectrogram of the speaker microphone signal (ground truth) and the recovered audio from the model. e) The audio quality metric SDR of the testing dataset and the sentence (‘He came to the point’) with corresponding spectrograms at different training iterations. f) Metrics histograms demonstrating improved audio quality from the original audios to recovered audios.



Conclusion

Our results demonstrate REAL can recover the speech signal by exploiting the back-scattered intensities from vibrating surfaces. With strong resistance to acoustic noise and the ability to collect specific audio signals over long distances, REAL provides a feasible solution to tackle the cocktail party problem in the optical channel. It is demonstrated that REAL could direct ‘hear’ the voices from masks and throats in a noisy environment, where the noise characteristics are fully considered in the hardware and the neural networks could help in signal recovery. Further work could include utilizing additional sensing modalities to enhance the overall detection accuracy such as the audio-visual cues and microphone array. With the high signal quality, simple construction, affordability and miniaturization readiness, we anticipate the REAL system will foster a new way in human-robot interaction, benefiting applications in speaker identification, speech understanding and accelerating the development of voice-guided home and field robots.

Materials and methods

Components and construction of REAL

A low-power collimated laser (such as a laser pointer) can be used as a REAL laser. The detection field of view of the telescoping lens system must align well to cover the surface laser spot. A color filter should be added before the APD to isolate the background light noise, and the laser power and spot size need to be adjusted to ensure laser safety according to IEC 60825-1:2014. Choose an appropriate gain of APD to ensure that shot noise dominates other noises such as electronic noise, while not introducing much excess noise. A low-noise transimpedance amplifier (TIA) and a secondary amplifier might be necessary to amplify the signal. The entire system needs to have a bandwidth beyond doubling the desired highest audio frequency.

SNR model of REAL

The signal of REAL is caused by the varying detection power when the vibrating surface displaces. The signal model is expressed as

$$\Delta U = U - U' = \frac{e\eta}{hv} P_0 MR \left[\frac{1}{r^2} - \frac{1}{(r + \Delta r)^2} \right] = A \left[\frac{1}{r^2} - \frac{1}{(r + \Delta r)^2} \right]$$

Where U is APD output voltage, e is elementary charge, h is the quantum efficiency of APD, M is the gain of the APD, h is the Planck constant, v is the photon frequency, $\frac{P_0}{r^2}$ is the received laser power when the surface is static, R is the TIA load resistance, and r is the distance between REAL and the speaker, Δr is surface displacement. We could use constant $A = \frac{e}{hv} P_0 MR$ to simplify the equation. A is 2.48 Vm^2 when the emitted laser power is 9 mW in our setup.

We use the root mean square s to quantify noise amplitude. REAL signal DU is affected by laser power fluctuation s_l , optical shot noise s_s , electronic noise s_e , target movement noise s_m and laser pointing noise s_p . Total noise is expressed as

$$s = \sqrt{s_l^2 + s_s^2 + s_e^2 + s_m^2 + s_p^2}$$

where we assumed these noises are independent. In noise analysis, the emitted laser power P is proportional to the collected laser power P_c given that r is fixed. The collected laser power is $P_c = \frac{(1+n)P_0}{r^2}$, where n represents the power fluctuation ratio and is much smaller than one. Therefore, the laser fluctuation is proportional to laser power: $s_l^2[?]P^2$. The shot noise s_s is determined by the Poisson process and goes as squared root of laser power: $s_s^2[?]P$. The electronic noise s_e is mainly determined by the electronic thermal noise and can be assumed as a constant. The target movement noise s_m is caused by the target movement. Our simulated s_m^2 is $3.1265 \times 10^{-12} \text{ V}^2$ in audio frequency range, which is caused by a movement (Δr) with a moving speed of 0.1 m/s, when laser power is 9 mW and distance (r) is 10 m. Since this noise is relatively small and different each time, we do not introduce s_m in the total noise model. The laser pointing noise s_p is usually caused by the unstable lasing mode change within the laser cavity and can be ignored from a stable laser. Therefore, the relationship between noise and laser power is as follows

$$s^2 = s_l^2 + s_s^2 + s_e^2 = aP^2 + bP + c$$

For our setup, we can determine the coefficient a, b, c by measuring the total noise under different laser power (Figure 2e). Here only the audio frequency range of 0.1~10 kHz is considered. From the fitting, it is clear that

the shot noise is dominant in our setup (Figure 2e). It means that REAL has obtained its best SNR because shot noise cannot be reduced by signal amplification. SNR model (in power) of REAL can be obtained by combining Equation 1 and Equation 3.

REAL signal process

The AC signal from the APD is amplified, passed to analog-to-digital convertor for digitalization and processed. As **Figure 5** shown, the first step of processing is to detrend by subtracting the low-frequency (below 30 Hz) envelope from the raw signal. The second step is to use a bandpass filter (0.1~7 kHz) to extract the signal in the audio frequency range. The final step uses spectral subtraction to reduce noise.^[32] After processing, the signal is normalized into the interval $[-1, 1]$ and converted into audio.

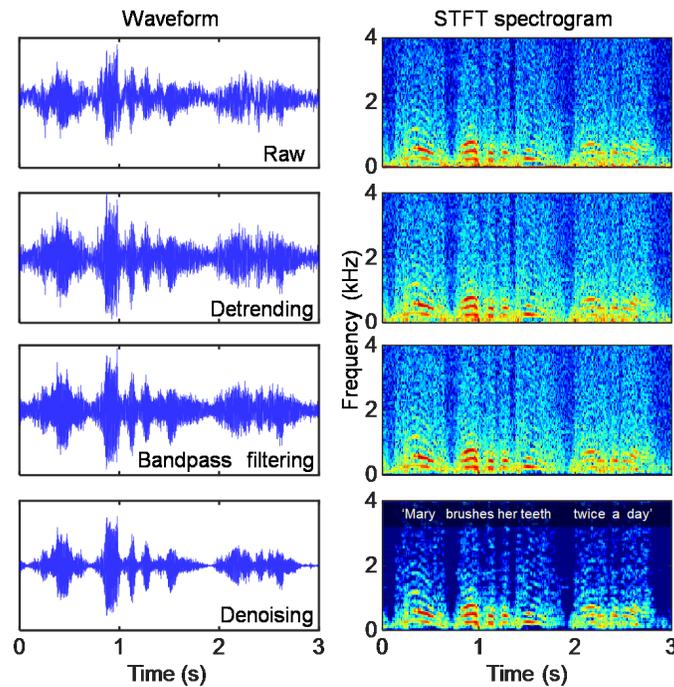


Figure 5. The signal processing of REAL audio. The raw signal from APD is processed sequentially by detrending, bandpass filtering, denoising and normalization.

REAL audio neural network model

The overall network design relies on the use of CNNs and LSTM to learn meaningful time-frequency information from laser spectrograms to estimate recovered audios. CNNs and LSTM have been extensively used in speech processing, especially on the task of speech enhancement and source separation.^[33-35] This section describes the full details of our laser audio processing network architecture.

The recorded dataset contained laser audios (as training and testing sets) and corresponding microphone audios as ground truth for supervised learning. Time synchronization should be ensured for each pair of recorded trials. In preprocessing, all audios are converted to single-channel with a sampling rate of 16 kHz, and the temporal sequence of each audio pair (REAL and microphone) is aligned with MATLAB,^[36] which estimates time delay that maximizes the cross-correlation between a pair of signals. The synchronized audios are then divided into segments of 3-second pieces, each pair is regarded as a sample instance of the final dataset. We discard all segments with utterances shorter than 1 second. After filtering 924 valid sample

instances are obtained, followed by a splitting procedure that randomly resamples and divides the dataset into training and testing subsets with a ratio of 7:3. To be clear, the REAL signal in the audio dataset used for the neural network model does not require spectral subtraction since the network is capable of denoising.

We used a modified version of the Voicefilter architecture^[34] as the base of our network (Figure 4c). The speaker voice embedding was removed from the original architecture since it was aimed at separating the voice of a speaker of interest from several concurrent speakers, while in our situation the aim is to recover the laser signal to the original audio. The soft mask prediction from the earlier work is also removed due to the same reason. The network is trained to minimize the difference between the estimated magnitude spectrogram and the target magnitude spectrogram. The phase of the estimated magnitude spectrogram is directly assigned from the laser audio. The complete model is composed of 8 convolutional layers, 1 LSTM layer, and 1 fully connected layer, each with ReLU activations except the last layer, which has a sigmoid activation. The feature map outputs from the convolutional layers are concatenated along the frequency axis, and the concatenated full feature map is then fed as the input to the following LSTM layers in time frames orders. To reduce overfitting, a dropout layer with a 20% drop rate is added between the convolutional layer blocks. We apply an initial learning rate of 10^{-3} , and decrease by a factor of 10 every 100 epochs. This model is implemented with PyTorch (<https://pytorch.org/>) deep learning framework^[37]. Adam solver in PyTorch is used to train our model and minimize the cross-entropy loss.

To evaluate the performance of different audio enhancement models, we use two metrics: source-to-distortion ratio (SDR) and short-time objective intelligibility (STOI). SDR is a common metric to evaluate speech enhancement performance^[10,34,38] and is typically expressed as an energy ratio (in dB) between the target signal and the total error from interference, noise and artifacts^[34]. Intelligibility is another indicator describing the effectiveness in sound perception, and STOI is the state-of-the-art intelligibility metric^[5]. STOI uses a discrete Fourier transform-based time-frequency decomposition to measure the correlation between the short-time temporal envelopes of a clean utterance and a separated utterance^[30]. The value can be interpreted as percent correct between 0 and 1. It is clear that in both metrics the recovered REAL signal improved significantly after enhancement and allowed the understanding of REAL signals from throats.

Acknowledgements

This work was supported by SUSTech startup Fund Y01966105 and DJI-joint Lab Fund K2096Z028. X. Guo thanks DJI for DJI-scholarship and Jiawei Wang for assistance and discussions. X. Guo and S. Ding contributed equally to this work.

Conflict of interest

The authors declare no conflicts of interest.

Supporting Information

[10.22541/au.165468550.08934241/v1](https://doi.org/10.22541/au.165468550.08934241/v1)

References

- [1] E. C. Cherry, *J. Acoust. Soc. Am.* **1953**, *25*, 975-979.
- [2] S. Haykin, & Z. Chen, *Neural Comput.* **2005**, *17*, 1875-1902.
- [3] H. DuRant, J. You, *Science* **2014**, *346*, 184-185.
- [4] J. C. Middlebrooks, J. Z. Simon, A. N. Popper, R. R. Fay, *The auditory system at the cocktail party*, Springer, New York, USA **2017**.
- [5] D. Wang, J. Chen, *IEEE/ACM Trans. Audio Speech Lang. Process.* **2018**, *26*, 1702-1726.
- [6] J. Huang, X. Zhang, F. Guo, Q. Zhou, H. Liu, B. Li, *IEEE Trans. Instrum. Meas.* **2014**, *64*, 2035-2043.
- [7] I. Dokmanić, R. Scheibler, M. Vetterli, *IEEE J. Sel. Top. Signal Process.* **2015**, *9*, 825-836.
- [8] J. Blauert, *Spatial hearing: the psychophysics of human sound localization*, MIT press, Cambridge, USA **1997**.
- [9] J. Xu, J. Shi, G. Liu, X. Chen, B. Xu, in AAAI Conf. Artif. Intell. (AAAI) New Orleans, **2018**.
- [10] Y. Luo, N. Mesgarani, *IEEE/ACM Trans. Audio Speech Lang. Process.* **2019**, *27*, 1256-1266.
- [11] C. Han, O. James, Y. Luo, H. Jose, D. M. Ashesh, M. H. Mesgarani, Han, *Sci. Adv.* **2019**, *5*, eaav6134.
- [12] J. K. Bizley, Y. E. Cohen, *Nat. Rev. Neurosci.* **2013**, *14*, 693-707.
- [13] A. Ephrat, M. Inbar, L. Oran, D. Tali, K. Wilson, A. Hassidim, T. F. William, R. Michael, *ACM Trans. Graph.* **2018**, *37*, 1-11.
- [14] R. Gu, S. Zhang, Y. Xu, L. Chen, Y. Zou, D. Yu, *IEEE J. Sel. Top. Signal Process.* **2020**, *14*, 530-541.
- [15] A. G. Bell, *Science* **1880**, *1*, 130-134.
- [16] A. L. Schawlow, C. H. Townes, *Phys. Rev.* **1958**, *112*, 1940-1949.
- [17] W. Li, M. Liu, Z. Zhu, T. S. Huang, in 18th Int. Conf. Pattern Recognit. (ICPR) IEEE, Hong Kong, **2006**.
- [18] Y. Qu, T. Wang, Z. Zhu, *IEEE/ASME Trans. Mechatron.* **2010**, *16*, 1110-1119.
- [19] L. Sun, J. Du, Z. Xie, Y. Xu, *J. Signal Process. Syst.* **2018**, *90*, 975-983.
- [20] J. M. Moses and K. Trout, *The Phys. Teach.* **2006**, *44*, 600-603.
- [21] B. Nassi, Y. Pirutin, A. Shamir, Y. Elovici, B. Zadov, *IACR Cryptol. ePrint Arch.* **2020**, 708.
- [22] J. H. Lambert, *Photometria, sive de Mensura et gradibus luminis, colorum et umbrae*, Sumptibus viduae E. Klett, Augsburg, GER 1760.
- [23] J. B. Allen, L. R. Rabiner, *Proceedings of the IEEE* **1997**, *65*, 1558-1564.
- [24] M. Morise, F. Yokomori, K. Ozawa, *IEICE Trans. Inf. Syst.* **2016**, *99*, 1877-1884
- [25] D. Ivo, H. S. Jae, S. Stefan, I. Sebastian, G. Rahel, P. Flurin, E. Albrecht, M. H. Alexander, R. Christof, *Hear. Res.* **2017**, *355*, 1-13.
- [26] R. M. Corey, U. Jones, A. C. Singer, *J. Acoust. Soc. Am.* **2020**, *148*, 2371-2375.
- [27] L. L Mendel, J. A. Gardino, S. R. Atcherson, *J. Am. Acad. Audiol.* **2008**, *19*, 686-695.
- [28] J. Chodosh, B. E. Weinstein, J. Blustein, *BMJ-Brit. Med. J.* **2020**, *370*, m3326.

- [29] E. Vincent, R. Gribonval, C. Févotte, *IEEE/ACM Trans. Audio Speech Lang. Process.* **2006**, *14*, 1462-1469.
- [30] C. H. Taal, R. C. Hendriks, R. Heusdens, J. Jensen, in Int. Conf. Acoust. Speech Signal Process. (ICASSP), Texas, **2010**.
- [31] S. R. Park, J. Lee, *Arxiv*, abs/1609.07132, **2016**.
- [32] S. Boll, *IEEE/ACM Trans. Audio Speech Lang. Process.* **1979**, *27*, 113-120.
- [33] A. V. D. Oord, D. Sander, Z. Heiga, S. Karen, V. Oriol, G. Alex, K. Nal, S. Andrew, K. Koray, *Arxiv*, abs/1609.03499, **2016**.
- [34] Q. Wang, M. Hannah, W. Kevin, S. Prashant, Z. Wu, H. John, S. Rif, W. Ron, Y. Jia, L. M. Ignacio, *Arxiv*, abs/1810.04826, **2018**.
- [35] J. Shen, R. Pang, R. J. Weiss, M. Schuster, N. Jaitly, Z. Yang, Z. Chen, Y. Zhang, Y. Wang, R. RJ, S. Rif, A. Yannis, Y. Wu, in Int. Conf. Acoust. Speech Signal Process. (ICASSP), Alberta, **2018**.
- [36] D. J. Higham, N. J. Higham, MATLAB guide. SIAM, Philadelphia, 2016).
- [37] P. Adam, G. Sam, C. Soumith, C. Gregory, Y. Edward, D. Zachary, Z. Lin, D. Alban, A. Luca, L. Adam, Automatic differentiation in pytorch. in 31st Conf.on Neural Inf. Process. Syst. (NIPS), California, **2017**.
- [38] D. Yu, M. Kolbæk, Z. H. Tan, J. Jensen, in Int. Conf. Acoust. Speech Signal Process. (ICASSP), New Orleans, **2017**.