

# Joint Time-Frequency Scattering

Joakim Andén, Vincent Lostanlen, and Stéphane Mallat

**Abstract**—In time series classification and regression, signals are typically mapped into some intermediate representation used for constructing models. Since the underlying task is often insensitive to time shifts, these representations are required to be time-shift invariant. We introduce the joint time-frequency scattering transform, a time-shift invariant representation which characterizes the multiscale energy distribution of a signal in time and frequency. It is computed through wavelet convolutions and modulus non-linearities and may therefore be implemented as a deep convolutional neural network whose filters are not learned but calculated from wavelets. We consider the progression from mel-spectrograms to time scattering and joint time-frequency scattering transforms, illustrating the relationship between increased discriminability and refinements of convolutional network architectures. The suitability of the joint time-frequency scattering transform for time-shift invariant characterization of time series is demonstrated through applications to chirp signals and audio synthesis experiments. The proposed transform also obtains state-of-the-art results on several audio classification tasks, outperforming time scattering transforms and achieving accuracies comparable to those of fully learned networks.

**Index Terms**—Acoustic signal processing, continuous wavelet transform, convolutional neural networks, supervised learning.

## I. INTRODUCTION

To extract information from signals, we typically map them into a lower-dimensional representation space where we construct model. The suitability of these representations depends on their ability to capture signal structure relevant to the task in question, such as classification or regression. For time series, this often includes the signal's time-frequency geometry. Figure 1 shows a time-frequency decomposition, the wavelet transform, applied to two audio recordings. Both are recordings of a person laughing, so their time-frequency structure is similar, but they also exhibit significant variability. We would like to construct representations invariant to this type of variability but which adequately capture the time-frequency structure of the signals.

An especially important form of variability is time-shifting (and time-warping deformations). Indeed, many time series classification and regression tasks are invariant to these transformations. This work will therefore study representations that are time-shift invariant.

Initial work on audio classification computed representations from time-frequency decompositions, such as windowed Fourier transforms. These include mel-spectrograms, mel-frequency

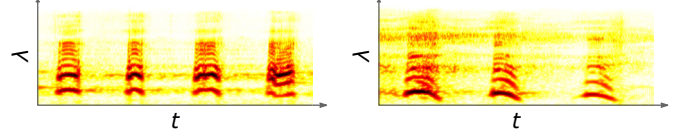


Fig. 1: The wavelet transform amplitudes, or scalograms, of two recordings as a function of time  $t$  and log-frequency  $\lambda$ . Both recordings are of one person laughing.

cepstral coefficients (MFCCs) [1], modulation spectrograms [2], [3] and correlograms [4], [5]. More recent work employs deep convolutional neural networks—cascades of filter banks alternated with nonlinearities [6], [7], [8]. Filters are learned from data, so each network is adapted to then task, often resulting in excellent performance [9]. However, learning typically requires large training sets and extensive computational resources.

This work provides a bridge between traditional time-frequency representations and deep convolutional neural networks. In particular, we implement the mel-spectrogram as a convolutional network and extend it by adding certain filters to that network which increase its discriminative power while maintaining the amount of time-shift invariance. These filters are not learned but fixed according to the invariance and discriminability needs of the task. This simplifies analysis and interpretation of the network. Fixed filters also reduces the associated computational burden since no training is necessary.

A convolutional network cascades convolutions, subsampling operators, and pointwise nonlinearities (such as rectifiers) [10], [11]. Its convolution kernels, or filters, are optimized over a training set. Section II-A describes how the wavelet transform is computed by a similar cascade of convolutions, but with fixed filters. A wavelet transform is thus a convolutional network with filters specified by certain time-frequency topology.

To impose time-shift invariance, we compute the modulus of the wavelet transform, known as the scalogram, and average in time. As shown in Section II-B, this yields a variant of the popular mel-spectrogram.

Although powerful, mel-spectrograms do not capture large-scale temporal structure, such as amplitude modulation. In Section II-C, the time scattering transform extends the mel-spectrogram through multiscale modulation coefficients [12], [13]. Instead of averaging the scalogram, it applies a second wavelet transform in time, takes the modulus, and averages. This representation is more discriminative and performs well for several classification tasks [13], [14], [15], [16]. Extending the wavelet transform network now lets us implement both mel-spectrograms and time scattering as convolutional networks.

A significant limitation of the time scattering transform is its restriction to convolutions along the time axis. In other

This work is supported by the ERC InvariantClass 320959.

J. Andén is with the Flatiron Institute, New York, NY, USA (e-mail: janden@flatironinstitute.org).

V. Lostanlen is with the Cornell Lab of Ornithology, Cornell University, Ithaca, NY, USA and the Music and Audio Research Laboratory, New York University, New York, NY, USA (e-mail: vincent.lostanlen@nyu.edu).

S. Mallat is with the Département d'Informatique, Ecole Normale Supérieure, Paris, France, the Collège de France, Paris, France, and the Flatiron Institute, New York, NY, USA (e-mail: mallat@di.ens.fr).

words, its convolutional network is actually a tree, with each node having only a single parent. A consequence is that time scattering cannot separate signals subjected to time shifts which vary in frequency, which is shown in Section III-A. To remedy this, we must capture time and frequency structure jointly.

With this goal in mind, we introduce the joint time-frequency scattering transform. As described in Section III-B, it replaces the one-dimensional, channel-by-channel wavelet decomposition of the scalogram by a two-dimensional wavelet transform. Its construction is inspired by the cortical transform of Shamma et al. [17], [18], which provides neurophysiological models of auditory processing in the mammalian brain. The corresponding joint scattering network introduces additional filters into the time scattering network, breaking its tree structure and increasing its discriminative power. To illustrate this, Section III-C shows how the joint scattering transform captures the chirp rate of frequency-modulated excitations.

The representational power of the proposed transform is further demonstrated in Section IV through synthesis experiments. Here, a signal is synthesized from a target scattering transform by minimizing the distance of its transform to that target. The resulting synthesized signals show how certain structures which are not captured by the mel-spectrogram and time scattering are better characterized by the joint scattering transform.

Section V concludes by evaluating the joint time-frequency scattering transform on several audio classification tasks. These include classification of phone segments, musical instruments, and acoustic scenes. The joint transform outperforms the mel-spectrogram and time scattering while achieving results comparable to, or better than, state-of-the-art convolutional networks. All figures and tables may be reproduced using software available at <http://www.di.ens.fr/data/software/>.

## II. TIME-SHIFT INVARIANT REPRESENTATIONS

Section II-A defines the wavelet transform, a representation well suited for time series with multiscale structure. The modulus of the wavelet transform, known as the scalogram, is averaged in time to yield the time-shift invariant mel-spectrogram, as described in Section II-B. Section II-C reviews the time scattering transform, introduced in Andén and Mallat [13], which extends the invariant mel-spectrogram. Instead of just averaging the scalogram, it also applies a second wavelet transform, demodulates, and averages the result in time. These representations are cascades of convolutions and non-linearities and may thus be implemented as deep convolutional networks with fixed filters.

### A. Wavelet Transform Filter Bank

The wavelet transform of a signal is obtained by convolving it with a set of dilated bandpass filters known as wavelets. It captures both short, transient structures and long-range oscillations in a localized manner. In the frequency domain, the ratio between center frequency and bandwidth, the  $Q$  factor, is the same for all filters. These transforms are therefore constant- $Q$  transforms [19]. Wavelet filter banks provide good models for cochlear function in mammals [20], [17], [18], [21] and form the basis for many audio representations [22]. The transform

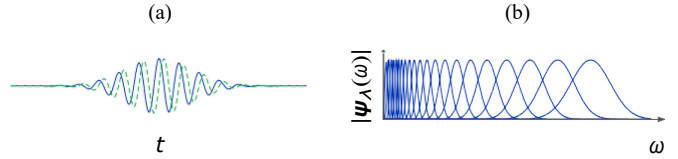


Fig. 2: (a) Real and imaginary parts of a Morlet mother wavelet with  $Q = 4$ . (b) The wavelet filters in the frequency domain.

may be computed using a multirate filter bank, as has been described in several works [22], [23].

Let  $\mathbf{x}(t)$  be a continuous signal for  $t \in \mathbb{R}$ . Its Fourier transform is given by  $\mathbf{\hat{x}}(\omega) = \int_{\mathbb{R}} \mathbf{x}(t) e^{-2\pi i \omega t} dt$  for  $\omega \in \mathbb{R}$ . Following Andén and Mallat [13], we consider a complex analytic wavelet  $\psi(t)$  whose Fourier transform  $\hat{\psi}(\omega)$  is concentrated in the interval  $[2^{-1/Q}, 1]$  for some  $Q \geq 1$ . Dilating  $\psi(t)$  by factors  $2^{-\lambda}$  now yields the wavelet filter bank

$$\psi_{\lambda}(t) = 2^{\lambda} \psi(2^{\lambda} t) \iff \hat{\psi}_{\lambda}(\omega) = \hat{\psi}(2^{-\lambda} \omega), \quad (1)$$

for  $\lambda \in \mathbb{R}$ . Consequently,  $\hat{\psi}_{\lambda}(\omega)$  is concentrated in  $[2^{\lambda-1/Q}, 2^{\lambda}]$ . This interval has approximate center  $2^{\lambda}$  and bandwidth  $2^{\lambda}/Q$ . We therefore need  $Q$  filters to cover an octave, independent of frequency. Since  $\hat{\psi}_{\lambda}(\omega)$  is concentrated around  $2^{\lambda}$ , we refer to  $\lambda$  as the wavelet's log-frequency index.

We are typically interested only in structures shorter than some fixed time scale  $T$ . In time,  $\psi_{\lambda}(t)$  has approximate duration  $2^{-\lambda}Q$ . We therefore require  $\lambda$  to satisfy  $2^{-\lambda}Q \leq T$ . Unfortunately, certain low frequencies are then not covered by any wavelet. For audio signals, these frequencies typically contain a small amount of energy and may be safely ignored. In the following, we instead add a set of constant-bandwidth filters covering these frequencies (see Andén and Mallat [13]).

In numerical experiments, we use the Morlet wavelet due to its near-optimal time-frequency localization [22], [13]. Figure 2 shows a sample Morlet wavelet and its wavelet filter bank.

We now define the continuous wavelet transform of  $\mathbf{x}(t)$  as

$$\mathbf{x} * \psi_{\lambda}(t) \quad (2)$$

for  $\lambda$  such that  $2^{-\lambda}Q \leq T$ . It captures the local oscillations of  $\mathbf{x}(t)$  at time  $t$  and frequency  $2^{\lambda}$  with resolution  $2^{-\lambda}Q$  and  $2^{\lambda}/Q$  in time and frequency, respectively. In audio applications, we typically set  $Q \approx 8$  to better resolve oscillatory components.

Now let  $\mathbf{x}[n]$  be a discrete signal for  $n \in \mathbb{Z}$ . Its discrete-time Fourier transform is  $\mathbf{\hat{x}}(\omega) = \sum_{n \in \mathbb{Z}} \mathbf{x}[n] e^{-2\pi i n \omega}$  for  $\omega \in [-1/2, 1/2]$ .

We now define a discrete analog of the continuous wavelet transform (2), implemented as a multirate filter bank.

To achieve this, we consider the multiresolution pyramid obtained by averaging  $\mathbf{x}[n]$  at different scales  $2^j$ . We initialize the finest scale to  $\mathbf{a}_0[n] = \mathbf{x}[n]$ . For  $j > 0$ ,  $\mathbf{a}_j[n]$  is obtained from  $\mathbf{a}_{j-1}[n]$  through convolution by a lowpass filter  $\mathbf{h}[n]$  whose transfer function  $\hat{\mathbf{h}}(\omega)$  is concentrated in  $[-1/4, 1/4]$ . We then subsample by 2 to obtain

$$\mathbf{a}_j[n] = \mathbf{a}_{j-1} * \mathbf{h}[2n]. \quad (3)$$

Note that  $\mathbf{a}_j[n] = \mathbf{x} * \mathbf{h}_j[2^j n]$  for some filter  $\mathbf{h}_j[n]$  defined by

$$\hat{\mathbf{h}}_j(\omega) = \hat{\mathbf{h}}(2^j \omega).$$

$$p=0$$

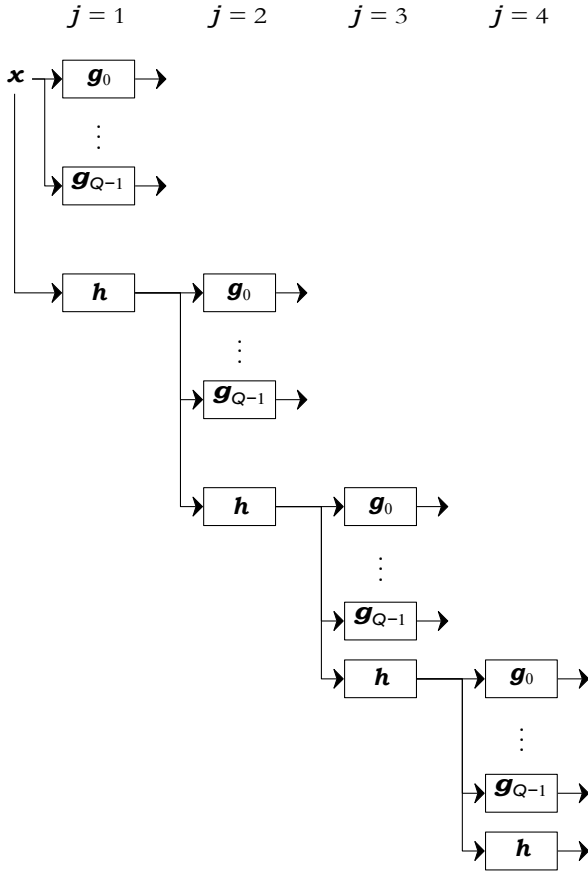


Fig. 3: Multirate filter bank computing wavelet coefficients for  $J = 4$ . Each block corresponds to a filter convolution subsampled by 2 where a boxed  $h$  is a low-pass filter and a boxed  $g_k$  is a band-pass filter. The depth corresponds to the octave index  $j$  while  $k = 0, \dots, Q - 1$  is the suboctave index.

As a result,  $\hat{h}_j(\omega)$  is concentrated in  $[-2^{-j-1}, 2^{-j-1}]$  and  $h_j[n]$  has approximate duration  $2^{j+1}$ .

The high frequencies of  $a_{j-1}[n]$  lost when convolving with  $h[n]$  are captured by  $Q$  bandpass filters  $g_0[n], \dots, g_{Q-1}[n]$ . Each has a transfer function  $\hat{g}_k(\omega)$  concentrated in  $[2^{-(k+1)/Q-1}, 2^{-k/Q-1}]$ . After convolving  $a_{j-1}[n]$  with  $g_k[n]$ , the result is subsampled by 2, yielding

$$d_{j,k}[n] = a_{j-1} * g_k[2n], \quad (4)$$

for  $j > 0$  and  $0 \leq k < Q$ . One may verify that

$$d_{j,k}[n] = x * g_{j,k}[2^j n], \quad (5)$$

where  $\hat{g}_{j,k}(\omega) = \hat{h}_{j-1}(\omega) \hat{g}_k(2^{j-1} \omega)$ . These filters are concentrated in intervals  $[2^{-j-(k+1)/Q}, 2^{-j-k/Q}]$ . In time, they have approximate duration  $2^j Q$ . Since we are only concerned with local variability below time scale  $T$ , we require  $2^j Q \leq T$ . This specifies the maximum depth  $J = \log_2(T/Q)$  of the cascade.

Figure 3 illustrates this multirate filterbank cascade. Each box corresponds to a convolution and subsampling by 2 according to (3) or (4). First,  $x[n]$  is convolved with  $g_0[n], \dots, g_{Q-1}[n]$  and subsampled to yield the highest octave of bandpass coefficients  $d_{1,0}[n], \dots, d_{1,Q-1}[n]$ . Convolution with  $h[n]$  and subsampling provides the remaining low frequencies, and

the process is repeated. As we progress through this cascade, the depth corresponds to the octave index  $j$ .

Combining the bandpass outputs yields the discrete wavelet transform in (5) for  $1 \leq j \leq J$  and  $0 \leq k < Q$ . This is similar to the output of the continuous wavelet transform. Indeed, if we sample a continuous band-limited signal  $x(t)$  at unit intervals, its discrete wavelet transform (5) approximates the continuous transform (2) for  $\lambda = -j - k/Q \leq -1$  provided

possible to construct filters  $h[n]$  and  $g_0[n], \dots, g_{Q-1}[n]$  such that this correspondence holds for large  $j$  [22]. The result is an approximation of the continuous wavelet transform using the convolutional network illustrated in Figure 3.

### B. Mel-Spectrogram

The lack of time-shift invariance of the wavelet transform hinders its generalization power for classification. For most classification tasks, shifting a signal in time does not modify its class. To reduce variability when constructing models, the signal representation must therefore be made time-shift invariant. In Andén and Mallat [13], this is achieved by computing the modulus and applying a lowpass filter. Let us review this construction and study how this may be implemented in a

The amplitude of the wavelet transform is the scalogram:

$$X(t, \lambda) = |x * \psi_\lambda(t)|. \quad (6)$$

Figure 1 shows two sample scalograms. Since the wavelets are analytic, applying the complex modulus performs a Hilbert demodulation, capturing the temporal envelope of each subband. The scalogram  $X(t, \lambda)$  therefore describes the time-frequency intensity of  $x(t)$  at time  $t$  and log-frequency  $\lambda$ .

Unfortunately, the scalogram is not time-shift invariant. Indeed, shifting a signal  $x(t) \mapsto x_c(t) = x(t - c)$  also shifts its scalogram  $X(t, \lambda) \mapsto X_c(t, \lambda) = X(t - c, \lambda)$ . To ensure invariance, we average in time to obtain

$$Mx(t, \lambda) = X(\cdot, \lambda) * \varphi_T(t) = |x * \psi_\lambda| * \varphi_T(t), \quad (7)$$

where  $\varphi_T(t) = T^{-1} \varphi(T^{-1}t)$  for some lowpass filter  $\varphi(t)$  of duration 1, so  $\varphi_T(t)$  has duration  $T$ . This is the mel-spectrogram  $Mx(t, \lambda)$  of  $x(t)$ . For  $|c| \ll T$ , it satisfies  $Mx_c(t, \lambda) \approx Mx(t, \lambda)$ , so it is locally invariant to time-shifts. The underlying wavelet structure of the mel-spectrogram also ensures stability to time-warping deformations [13].

The mel-spectrogram was originally introduced for speech classification [1] and was motivated by psychoacoustic studies. It has since found widespread use in various audio classification tasks [24], [25], [26]. Traditionally, the mel-spectrogram is computed through frequency averaging of the windowed Fourier transform amplitude, also known as the spectrogram. However, it has recently been shown that they may be approximated by the time-averaged scalogram coefficients (7) [13], [27], [28]. Note that this formulation makes the time-shift invariance of the mel-spectrogram explicit. Indeed, the amount of invariance is directly controlled by the duration  $T$  of the lowpass filter  $\varphi_T(t)$ . We shall use this wavelet-based variant of the mel-spectrogram in the following.

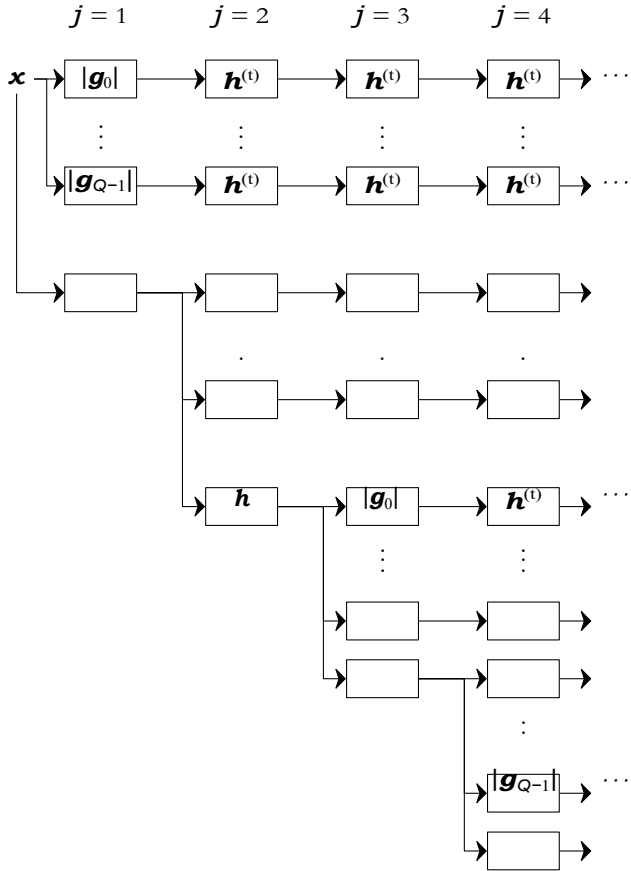


Fig. 4: Mel-spectrogram implemented as a convolutional network. Each  $|g_k|$  block convolves by a band-pass filter  $g_k[n]$ , computes the modulus, and subsamples by 2. Blocks containing  $h$  or  $h^{(t)}$  convolve by a low-pass filter and subsample by 2.

We now define the discrete mel-spectrogram using the discrete wavelet transform. The resulting convolutional network is shown in Figure 4. Instead of just convolving by  $g_k[n]$ , this network also applies a modulus and subsamples by 2. The whole operation is denoted by a boxed  $|g_k|$ . The result is then passed through a sequence of lowpass filters  $h^{(t)}[n]$  alternated with subsampling operators, approximating the convolution by  $\varphi_T(t)$ . The output is  $JQ + 1$  signals of form  $|x * g_{j,k} * h_{j-j}[2^{-j}n]|$ , where  $j$  is the depth at which the modulus was applied. If the filters are chosen as in Section II-A, this approximates  $\mathbf{M}x(t, \lambda)$  for a bandpass  $x(t)$ .

For real  $g_k[n]$ , we may replace the modulus with a rectified linear unit. Indeed, averaging a rectified bandpass signal approximates its Hilbert envelope, so the result is similar [29].

### C. Time Scattering

The mel-spectrogram discards a large amount of potentially useful information when averaging  $\mathbf{X}(t, \lambda)$  along  $t$  in (7), removing any high-frequency structure. The time scattering transform extends the mel-spectrogram and partially recovers this lost structure while maintaining invariance and stability [12], [13]. This is achieved in Andén and Mallat [13] by convolving the scalogram with a second set of wavelets, taking the modulus, and averaging to create second-order time

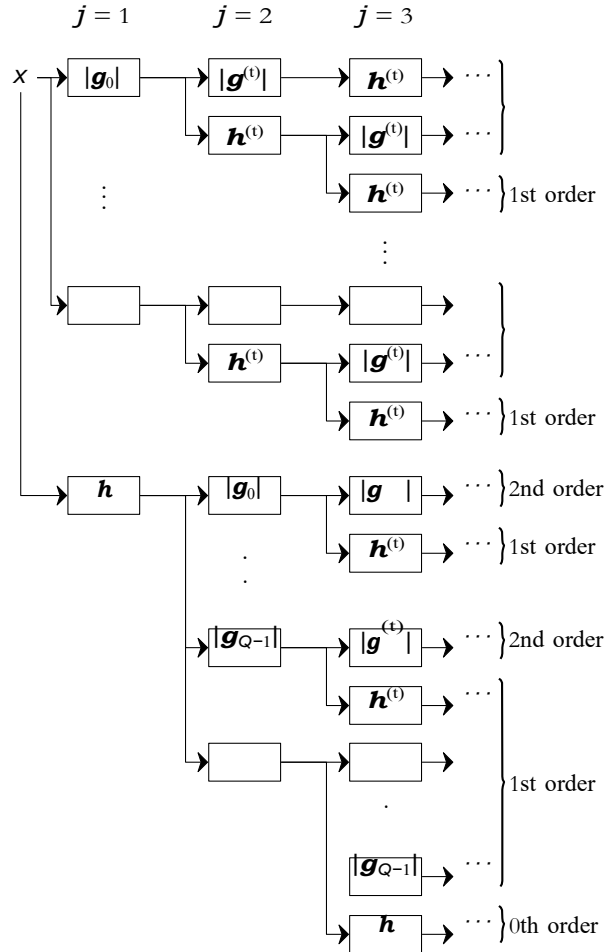


Fig. 5: A time scattering network. Each block with  $|g_k|$  or  $|g^{(t)}|$  outputs the modulus of the input convolved with a band-pass filter, subsampled by 2. Blocks with  $h$  and  $h^{(t)}$  convolves the input with a low-pass filter and subsample by 2.

scattering coefficients. Let us rederive this representation and implement it as a convolutional network extending that of the mel-spectrogram (see Figure 4).

The first-order time scattering coefficients coincide with the mel-spectrogram  $\mathbf{M}x(t, \lambda)$  and are given by

$$\mathbf{S}_1 x(t, \lambda) = \mathbf{X}(\cdot, \lambda) *_{\mathcal{T}} \varphi_T(t).$$

The lost high frequencies of  $\mathbf{X}(t, \lambda)$  are recovered by convolving with a new set of wavelets, defined from a Morlet mother wavelet  $\psi^{(1)}(t)$  by  $\psi_\mu^{(1)}(t) = 2^\mu \psi^{(1)}(2^\mu t)$  for  $\mu \in \mathbb{R}$ . Each  $\psi_\mu^{(1)}(t)$  has a center frequency of approximately  $2^\mu$ , so we refer to  $\mu$  as their log-frequency. Unlike their first-order counterparts  $\psi_\lambda(t)$ , the second-order wavelets  $\psi_\mu^{(1)}(t)$  have  $Q = 1$ . As a result, they are better adapted to structures in  $\mathbf{X}(t, \lambda)$ , which are less oscillatory and more localized in time compared to those in  $x(t)$ .

Convoluting  $\mathbf{X}(t, \lambda)$  with these wavelets along  $t$ , we obtain  $\mathbf{X}(\cdot, \lambda) * \psi_\mu^{(1)}(t)$ . To ensure local invariance to translation, we take another modulus and average using  $\varphi_T(t)$ , which yields

$$\mathbf{S}_2 x(t, \lambda, \mu) = ||x * \psi_\lambda| * \psi_\mu^{(1)}| * \varphi_T(t). \quad (8)$$

These are the second-order time scattering coefficients. They describe the variability of  $\mathbf{X}(t, \lambda)$  along  $t$  at frequency  $2^\mu$ , where  $\lambda$  is the first-order, or acoustic, log-frequency, while  $\mu$  is the second-order, or modulation, log-frequency. As before, we limit ourselves to scales shorter than  $T$  by enforcing  $2^{-\mu} \leq T$ .

Concatenating all first- and second-order scattering coefficients  $\mathbf{S}_1 \mathbf{x}$  and  $\mathbf{S}_2 \mathbf{x}$  of  $\mathbf{x}(t)$  yields the time scattering transform  $\mathbf{S} \mathbf{x}$ . Higher-order scattering coefficients may be defined [12], but these are of negligible energy [30] and do not greatly affect classification results [13]. The scattering transform exhibits the same amount of time-shift invariance and time-warping stability as the mel-spectrogram described previously [12], [13]. It is more discriminative than the mel-spectrogram, however, since it captures amplitude modulations in  $\mathbf{X}(t, \lambda)$  along  $t$ . As a result, the time scattering transform enjoys better performance for classification of audio [13], biomedical [14], and other types of time series [15], [16].

Other approaches capture temporal structure in the scalogram using Fourier transforms [2], [3] or second-order moments [4], [5], [31]. However, these lack the time-warping stability or noise robustness of the scattering transform [13], [12].

Extending the mel-spectrogram convolutional network of Figure 4, we define the network of a discrete time scattering transform. The result is shown in Figure 5. To implement the second-order wavelets  $\psi_\mu^{(0)}(t)$ , we use the network of Figure 3, but with a single bandpass filter  $\mathbf{g}^{(0)}[n]$  and a lowpass filter  $\mathbf{h}^{(0)}[n]$ . These are constructed to approximate convolutions with  $\psi_\mu^{(0)}(t)$  for  $\mu = -j \leq -1$  as described in Section II-A.

As before,  $\mathbf{x}[n]$  is first decomposed in the  $|\mathbf{g}_k|$  boxes by convolution with  $\mathbf{g}_0[n], \dots, \mathbf{g}_{Q-1}[n]$  followed by modulus and subsampling by 2. However, instead of averaging their outputs, they are further convolved with  $\mathbf{g}^{(0)}[n]$  followed by modulus and subsampling, denoted by  $|\mathbf{g}^{(0)}|$  boxes. These coefficients are then averaged using lowpass filters  $\mathbf{h}^{(0)}[n]$  which alternate with subsampling operators. This yields the second-order scattering coefficients of  $\mathbf{x}[n]$  for the highest octave in  $\lambda$  and the highest octave in  $\mu$ . We obtain lower octaves in  $\mu$  by applying a sequence of convolutions with  $\mathbf{h}^{(0)}[n]$  alternated with subsampling operators before convolving with  $\mathbf{g}^{(0)}[n]$ . Similarly, lower octaves in  $\lambda$  are obtained by applying a sequence of convolutions by  $\mathbf{h}[n]$  and subsampling operators before the decomposition by  $\mathbf{g}_0[n], \dots, \mathbf{g}_{Q-1}[n]$ . The outputs of this convolutional network approximate the continuous time scattering transform  $\mathbf{S} \mathbf{x}$  of  $\mathbf{x}(t)$ .

### III. JOINT REPRESENTATIONS IN TIME AND FREQUENCY

While successfully describing temporal modulation, the time scattering transform fails to capture more sophisticated time-frequency structure, as shown in Section III-A. It fails because it decomposes the scalogram as a set of one-dimensional time series. Section III-B introduces the joint time-frequency scattering transform, which instead decomposes the scalogram in both time and log-frequency. Its convolutional network representation introduces connections between nodes in each layer, maintaining the amount of time-shift invariance but increasing its discriminability. This property is demonstrated in Section III-C, where we show how the proposed transform accurately captures frequency-modulated excitations.

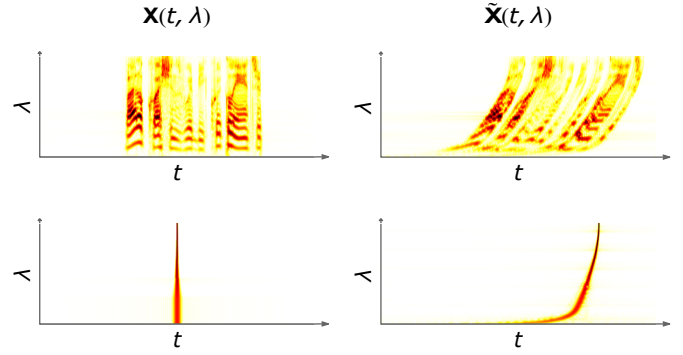


Fig. 6: Effect of frequency-dependent time-shifts  $\tau(\lambda)$  on scalograms of a speech recording (top) and a Dirac delta function (bottom). The two columns correspond to the original signal  $\mathbf{x}(t)$  and the transformed signal  $\tilde{\mathbf{x}}(t)$ , respectively.

#### A. Loss of Time-Frequency Structure

The time scattering convolutional network in Figure 5 has a tree structure; that is, each node only has one parent. In contrast, a general convolutional network sums contributions from multiple nodes in a layer to produce a node in the next layer. Due to this tree structure, the time scattering transform is not sensitive to certain time-frequency deformations.

To see this, we suppose that  $\mathbf{x}(t)$  is transformed into  $\tilde{\mathbf{x}}(t)$  whose scalogram  $\tilde{\mathbf{X}}(t, \lambda)$  is an approximate translation of  $\mathbf{X}(t, \lambda)$  by  $\tau(\lambda)$  in each frequency band. In other words,  $\tilde{\mathbf{X}}(t, \lambda) \approx \mathbf{X}(t - \tau(\lambda), \lambda)$ . Such transformations are illustrated in Figure 6 for a speech signal and a Dirac delta function. This time-frequency warping misaligns the speech harmonics and transforms the delta function into a chirp. Although  $\mathbf{x}(t)$  differs markedly from  $\tilde{\mathbf{x}}(t)$ , this is not detected by time scattering if  $|\tau(\lambda)| \ll T$ . Indeed, the effect of the frequency-varying time shift disappears when averaging by  $\phi_T(t)$ . Computing the scattering transforms  $\mathbf{S} \mathbf{x}$  and  $\mathbf{S} \tilde{\mathbf{x}}$  for  $T$  equal to the signal length yields relative differences  $|\mathbf{S} \tilde{\mathbf{x}} - \mathbf{S} \mathbf{x}| / |\mathbf{S} \mathbf{x}|$  of 0.07 and 0.09 for the speech signal and the delta function, respectively.

Detection of time-frequency warping requires measurement of scalogram variability across frequency. In particular, the second-order wavelet convolution (8) in time must be replaced by a convolution in time and log-frequency.

#### B. Joint Time-Frequency Scattering

Existing methods for capturing a signal's time-frequency geometry are not always suitable for classification. For example, McDermott and Simoncelli [31] compute higher-order moments of the scalogram across frequencies. Through synthesis experiments, this representation is shown to provide a good model for audio textures. However, higher-order moments are not robust to noise, reducing the descriptor's usefulness for classification.

An alternative approach, motivated by neurophysiological studies in the audio cortex of ferrets, is the cortical transform of Shamma et al. [17]. It decomposes the scalogram in both time and log-frequency using two-dimensional Gabor wavelets. The cortical transform and related representations have brought significant improvements over mel-spectrograms in tasks from



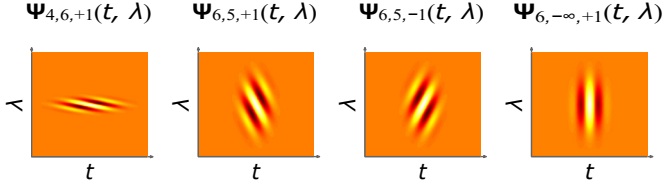


Fig. 7: Real parts of four time-frequency wavelets  $\Psi_{\mu,\mathcal{E},s}(t, \lambda)$ . White-yellow is negative, orange is zero, and red-black is positive.

speech classification [32], [33] to timbre analysis [18], [34], [35]. Unfortunately, the lack of time-shift invariance and time-warping stability limits the performance of this approach.

In the following, we adapt the cortical transform within the scattering framework, allowing us to address its invariance and stability. This also lets us analyze its discriminative power.

We first decompose the scalogram  $\mathbf{X}(t, \lambda)$  using a two-

wavelets. Two-dimensional Morlet wavelets are also used in the

success in natural image classification [36], [37]. In this case, however, the wavelets are obtained by rotating and uniformly

scalogram. Indeed, rotation does not preserve the relationship

not the scalogram of some other signal.

We instead define our wavelets separably, with independent

mother wavelet  $\Psi(t, \lambda) = \psi^{(t)}(t) \psi^{(f)}(\lambda)$  is the product of two one-dimensional functions in time and log-frequency. Both the time  $\psi^{(t)}(t)$  and the frequency  $\psi^{(f)}(\lambda)$  wavelets are Morlet wavelets with  $Q = 1$ . Dilating by  $2^{-\mu}$  along  $t$ , dilating by  $2^{-\mathcal{E}}$  along  $\lambda$ , and reflecting according to  $S$  yields the wavelet

$$\Psi_{\mu,\mathcal{E},s}(t, \lambda) = 2^{\mu+\mathcal{E}} \psi^{(t)}(2^\mu t) \psi^{(f)}(s 2^\mathcal{E} \lambda), \quad (9)$$

where the spin  $s = \pm 1$  specifies the oscillation direction (up or down). The frequency of the wavelet along  $t$  is  $2^\mu$ , so  $\mu$  is the log-frequency of  $\Psi_{\mu,\mathcal{E},s}(t, \lambda)$ . Its frequency along  $\lambda$  is  $2^\mathcal{E}$ , so we refer to it as a “quefrequency.” Consequently,  $f$  is the “log-quefrequency” of  $\Psi_{\mu,\mathcal{E},s}(t, \lambda)$ .

As before,  $\mu$  satisfies  $2^{-\mu} \leq T$ . Along  $\lambda$ , we fix some maximum log-frequency scale  $F$ , measured in octaves, and let  $2^{-\mathcal{E}} \leq F$ . At this maximum scale, we include a lowpass filter to capture average structure along  $\lambda$ . Specifically, we set

$$\Psi_{\mu,-\infty,+1}(t, \lambda) = 2^\mu \psi^{(t)}(2^\mu t) \varphi_F(\lambda). \quad (10)$$

Note that these are only defined for  $s = +1$ . Figure 7 shows a few sample two-dimensional wavelets  $\Psi_{\mu,\mathcal{E},s}(t, \lambda)$ .

The two-dimensional wavelet transform of the scalogram  $\mathbf{X}(t, \lambda)$  computes convolutions  $\mathbf{X} * \Psi_{\mu,\mathcal{E},s}(t, \lambda)$ . It captures the joint variability of  $\mathbf{X}(t, \lambda)$  at log-frequency  $\mu$  and log-quefrequency  $f$  with spin  $S$ . To ensure time-shift invariance and time-warping stability, we take the complex modulus and average, obtaining the second-order joint time-frequency scattering coefficients

$$\mathbf{S}_2 \mathbf{x}(t, \lambda, \mu, f, s) = |\mathbf{X} * \Psi_{\mu,\mathcal{E},s}(\cdot, \lambda)| * \varphi_T(t). \quad (11)$$

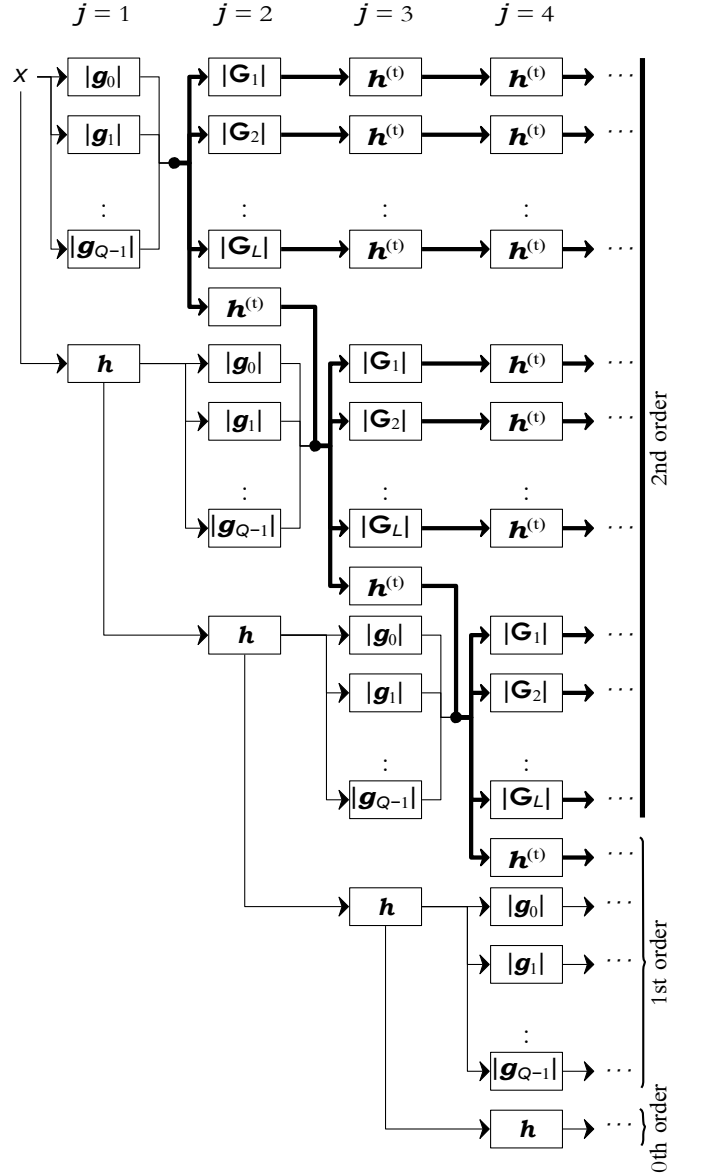


Fig. 8: A joint time-frequency scattering network. Each  $|g_k|$  block convolves a one-dimensional signal with the band-pass filter  $g_k$  and outputs its modulus. The outputs are aggregated into two-dimensional arrays shown by thick lines. A  $|G_{\mathcal{E}}|$  block convolves a two-dimensional array with the band-pass filter  $G_{\mathcal{E}}$  and outputs its modulus. The  $h$  and  $h^{(t)}$  blocks convolve only in time. All blocks subsample their output in time by 2.

These coefficients describe the time-frequency geometry of  $\mathbf{x}(t)$  at time  $t$  and log-frequency  $\lambda$ . They retain the time-shift invariance and time-warping stability of the second-order time scattering coefficients, but with increased discriminative power.

Concatenating the first-order time scattering coefficients  $\mathbf{S}_1 \mathbf{x}$  and the second-order time-frequency scattering coefficients  $\mathbf{S}_2 \mathbf{x}$  yields the complete joint time-frequency scattering transform  $\mathbf{S} \mathbf{x}$  of  $\mathbf{x}(t)$ . As for time scattering, we may define higher-order coefficients, but these are often of limited use for classification. For each  $t$ , there are  $O(Q \log_2 T)$  first-order coefficients and  $O(Q(\log_2 T)^2 \log_2 F)$  second-order coefficients.

We now define a convolutional network to provide a discrete implementation of the joint scattering transform. In the time scattering network (see Figure 5), we approximate the convolution of  $\mathbf{X}(t, \lambda)$  with  $\boldsymbol{\psi}_\mu^{(i)}(t)$  along  $t$  by cascading discrete filters  $\mathbf{h}^{(i)}[n]$  and  $\mathbf{g}^{(i)}[n]$ , alternated with subsampling operators. The joint transform network incorporates additional filters along the discrete log-frequency  $m = Q\lambda = -jQ - k \in \mathbb{Z}$ , where, as before,  $j$  and  $k$  are the octave and subband indices of  $\lambda$ .

In a given layer, the modulus bandpass outputs of the previous layer are arranged along time  $n$  and log-frequency  $m$  into a two-dimensional array. This array is then filtered along  $m$  by different filters  $2^\varepsilon \boldsymbol{\psi}^{(f)}(s2^\varepsilon m/Q)$ . It is also filtered by  $\boldsymbol{\varphi}_F(m/Q)$  to account for  $f = -\infty$ . The sampling interval of the filters is  $1/Q$ , since this is the spacing of the discretized log-frequencies  $\lambda = m/Q$ . Each frequency-filtered array is then filtered by  $\mathbf{g}^{(i)}[n]$  along  $n$ .

Combining these into two-dimensional filters, we get

$$\begin{aligned} \mathbf{G}_{\varepsilon, s}[n, m] &= \mathbf{g}^{(i)}[n] 2^\varepsilon \boldsymbol{\psi}^{(f)}(s2^\varepsilon m/Q) \\ \mathbf{G}_{-\infty, +1}[n, m] &= \mathbf{g}^{(i)}[n] \boldsymbol{\varphi}_F(m/Q), \end{aligned}$$

where  $f \in \mathbb{Z}$  such that  $-\log_2 F \leq f \leq \log_2 Q$  (to ensure that  $1/Q \leq 2^{-\varepsilon} \leq F$ ) and  $s = \pm 1$ . Abusing notation slightly, we renumber this set of discrete filters as  $\mathbf{G}_1[n, m], \dots, \mathbf{G}_L[n, m]$ . These filters capture all log-quefrequencies along  $\lambda$ , but only high frequencies along  $n$ . The missing low frequencies are absorbed by  $\mathbf{h}^{(i)}[n]$ , which averages along  $n$ , leaving  $m$  intact.

Using these filters, we construct the convolutional network shown in Figure 8, extending the time scattering network of Figure 5. Small circles denote aggregation of multiple time series into a two-dimensional array, while the arrays themselves

are thick lines. We denote by a boxed  $|\mathbf{G}_\varepsilon|$  convolution with  $\mathbf{G}_\varepsilon[n, m]$  for  $f = 1, \dots, L$ , followed by a complex modulus and subsampling by 2 along  $n$ . Similarly, a boxed  $\mathbf{h}^{(i)}$  denotes lowpass filtering along  $n$  by  $\mathbf{h}^{(i)}[n]$  followed by subsampling.

Starting with a signal  $\mathbf{x}[n]$ , we first compute its decomposition using the first-order blocks  $|\mathbf{g}_0|, \dots, |\mathbf{g}_{Q-1}|$ , extracting the highest octave of the signal. We then combine these into a two-dimensional array which is decomposed by  $|\mathbf{G}_1|, \dots, |\mathbf{G}_L|$ . The outputs of  $|\mathbf{G}_1|, \dots, |\mathbf{G}_L|$  are then forwarded to a succession of  $\mathbf{h}^{(i)}$  blocks which implement the averaging by  $\boldsymbol{\varphi}_T[n]$ . The original array is also decomposed by  $\mathbf{h}^{(i)}$ , and the result is concatenated to the first-order outputs of the second layer (that is, the second octave of the original signal). We then repeat the process on this array. As before, an appropriate choice of  $\mathbf{g}^{(i)}[n]$  and  $\mathbf{h}^{(i)}[n]$  ensures that the network accurately approximates the continuous joint scattering transform.

The important difference between this network and the time scattering network is the presence of within-layer connections. These break the tree structure, increasing discriminative power through better characterization of time-frequency geometry. Returning to the frequency-warped signals of Figure 6, the joint network separates the original and transformed signals, with  $|\mathbf{S}\tilde{\mathbf{x}} - \mathbf{S}\mathbf{x}|/|\mathbf{S}\mathbf{x}|$  of 0.41 and 0.90, compared to 0.07 and 0.09 for time scattering. This network therefore has same time-shift invariance as time scattering, but with better discriminability. Note that this increased discriminative power may not always be desirable. For example, frequency-dependent time-shifts (as

shown in Figure 6) or similar transformations may not be relevant for the classification task. In this case, replacing the time scattering transform with a joint time-frequency scattering transform would needlessly increase the number of model parameters, potentially requiring more training data to train an accurate classifier. On the other hand, the high-quefreny joint coefficients approximate the standard second-order time scattering coefficients. As a result, the types of structures captured by the time scattering transform are equally well characterized by the joint transform, so little discriminative power is lost by replacing the former by the latter.

The invariance and discriminability properties of the transform are controlled by three parameters:  $Q$ ,  $T$ , and  $F$ . The number of wavelets per octave,  $Q$ , depends on the time-frequency localization of the input signal. For example, if the signal is highly oscillatory (that is, well-localized in frequency, but not necessarily in time), a higher value for  $Q$  is appropriate. This is the case in audio, but not necessarily for biomedical or geophysical time series, which are more localized in time.

The averaging scale  $T$  controls the maximum length of the signal structure captured by the transform. In other words, if the relevant structure in a classification problem occurs at very small scales,  $T$  should be kept small. This is the case in phone segment classification (see Section V-B), where the object of interest, the phone, is of short duration. For other signals, such as musical instrument recordings (see Section V-C), there are relevant structures at larger scales. The  $T$  parameter also controls the length of the lowpass filter  $\boldsymbol{\varphi}_T(t)$  and therefore determines the amount of desired time-shift invariance.

The frequency scale  $F$  has a similar role, controlling the maximum frequency extent of the signal structure captured by the transform. If we expect relevant structures to spread out over several octaves, a large value for  $F$  is needed. This is the case for speech signals, where plosive phones occupy a large part of the frequency domain. For other signals, such as environmental sounds, relevant structures may be confined within an octave, so a small  $F$  is more appropriate.

The output of a scattering network may be used as input to another convolutional network whose filters are subsequently optimized for some classification task. This yields a large convolutional network taking raw waveforms as input and whose first few layers are fixed. By fixing certain layers, the network has fewer parameters to optimize and could then be trained using less data. Previous work training convolutional networks on raw waveforms have yielded mel-like filters in the first few layers [38], [39], providing some support for this idea. Other attempts at explicitly incorporating wavelets into convolutional network architectures have also demonstrated the viability of the approach [40], [41], [42], [43]. In addition, the success of transfer learning [44], [45], [46], [47] suggests that there exist certain universal representations which perform well for a wide range of tasks. The joint scattering network provides a way to construct such a representation while enforcing certain time-shift invariance and time-frequency discriminability conditions.



### C. Frequency Modulation

The above construction is similar to that of traditional convolutional networks except that filters are not learned from data. These filters provide the time-shift invariance and time-warping stability of the time scattering transform, but the joint transform is also more discriminative. To illustrate this, we show how the joint time-frequency scattering transform captures frequency modulation structure ignored by time scattering.

Let  $\mathbf{x}(t) = \exp(2\pi i \xi(t))$  be a frequency-modulated excitation with instantaneous phase  $\xi(t)$ . At time  $t$ , its instantaneous frequency is given by  $\xi^1(t)$ , while the relative change in this frequency, the (relative) chirp rate, is  $\xi^1(t)/\xi^1(t)$ . Frequency modulation occurs in a variety of signals, such as speech, animal calls, music and radar signals [48].

We now consider a particular case of frequency modulation: the exponential chirp. Here  $\xi(t) = 2^{at}$ , so it has instantaneous frequency  $\xi^1(t) = a \log(2) 2^{at}$  and constant chirp rate  $\xi^1(t)/\xi^1(t) = a \log(2)$ . We note that an arbitrary frequency-modulated excitation may be locally approximated by an exponential chirp by setting  $a = (\log 2)^{-1} \xi^1(t)/\xi^1(t)$ .

For exponential chirps, we have the following result.

**Theorem 1.** Let  $\psi_\lambda(t)$  and  $\Psi_{\mu, \varepsilon, s}(t, \lambda)$  be defined as in (1) and (9). We require that  $\psi(t)$  have compact support, that  $\|\psi\|_\infty, \|\psi^1\|_1, \|\psi^{(1)}\|_1, \|\psi^{(1)}\|_\infty$  are bounded, and that  $\text{supp } \psi^{(1)}(\lambda) \subset [-A, A]$  for some  $A > 0$ . Further, we assume that  $\psi^{(1)}(t)$  is the product of a positive envelope  $|\psi^{(1)}(t)|$  and  $\exp(2\pi i t)$ . Let  $\mathbf{x}(t) = \exp(2\pi i 2^{at})$  for some  $a \in \mathbb{R}$ . The joint scattering transform (11) then satisfies

$$\mathbf{S}_2 \mathbf{x}(t, \lambda, \mu, f, s) = \frac{c_0 E(t, \lambda)}{a} \left( \frac{s 2^{\mu - \varepsilon}}{a} + \varepsilon(t, \lambda, \mu, f, s) \right),$$

where

$$E(t, \lambda) = |\psi^{(1)}_\mu| * \varphi_T \left( t - \frac{\lambda}{a} + \frac{\log \log 2^a}{\log 2^a} \right),$$

$$\|\varepsilon\|_\infty < C \left( |a| 2^{-\lambda + 2^{-A}} + 2^{2\mu} |a|^{-2} + 2^{2\mu - \varepsilon} |a|^{-2} \right),$$

for  $C > 0$  depending only on  $\psi(t)$ ,  $\psi^{(1)}(t)$ , and  $\psi^{(1)}(\lambda)$ , and  $c_0 = \int_{\mathbb{R}} |\psi(2^u)| du$ .

The proof is given in Appendix A. The result relies on approximating  $\mathbf{X}(t, \lambda)$  by  $|\psi(\log(2^a) 2^{-\lambda + at})|$ . Since  $|\psi(\omega)|$  is maximized at  $\omega = 1$ , this forms a ridge  $\lambda = at$  with slope  $a$ , as illustrated in Figure 9(a,b). In the joint transform, this ridge only activates certain second-order wavelets  $\Psi_{\mu, \varepsilon, s}(t, \lambda)$ . Indeed, only wavelets whose slope  $-s 2^{\mu - \varepsilon}$  aligns with  $\lambda = at$  yield large coefficients. Taking the complex modulus and averaging in time preserves this slope information.

Let us consider another chirp  $\tilde{\mathbf{x}}(t) = \exp(2\pi i 2^{\tilde{a}t})$ . We may obtain  $\tilde{\mathbf{x}}(t)$  from  $\mathbf{x}(t)$  using a frequency-dependent time-shift of its scalogram  $\mathbf{X}(t, \lambda)$  as in Section III-A. Here, we take

$$\tau(\lambda) = \lambda \left( \frac{1}{a} - \frac{1}{\tilde{a}} \right) - \frac{\log \log 2^a}{\log 2^a} + \frac{\log \log 2^{\tilde{a}}}{\log 2^{\tilde{a}}}.$$

As we saw in Section III-A, the time scattering transform is not sensitive to such changes. In other words, the scattering transform discards information on slope, rendering it unsuitable for

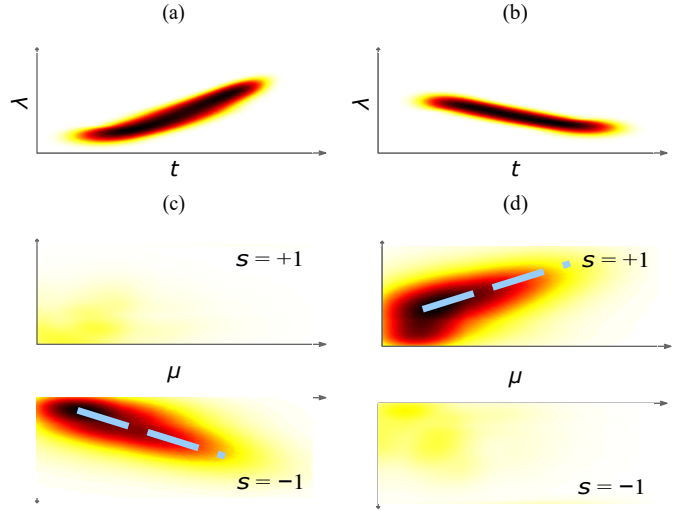


Fig. 9: Scalograms of two exponential chirps with chirp rates (a)  $a = 4$  and (b)  $a = -2$ . (c, d) Corresponding second-order joint time-frequency scattering coefficients  $\mathbf{S}_2 \mathbf{x}(t, \lambda, \mu, f, s)$  for fixed  $t$  and  $\lambda$ . The dotted lines satisfy  $s 2^{\mu - \varepsilon} = -a$ .

describing frequency modulation. The same applies to related representations which also decompose each subband of  $\mathbf{X}(t, \lambda)$  separately, such as mel-spectrograms, MFCCs, and modulation spectrograms. This information loss is fundamentally due to the tree structure of their convolutional networks.

Theorem 1 states that, for fixed  $t$  and  $\lambda$ ,  $\mathbf{S}_2 \mathbf{x}(t, \lambda, \mu, f, s)$  is approximately proportional to  $|\psi^{(1)}(-s 2^{\mu - \varepsilon} a^{-1})|$ . Since  $\psi^{(1)}$  is concentrated around frequency 1, this is maximized for  $-s 2^{\mu - \varepsilon} a^{-1} = 1$ . In other words, a ridge is present along  $s 2^{\mu - \varepsilon} = -a$ . Frequency modulation structure in the form of the chirp rate  $a$ , is thus encoded in the second-order joint time-frequency scattering coefficients. Consequently, they are sensitive to frequency-dependent time-shifts  $\mathbf{X}(t, \lambda) \mapsto \mathbf{X}(t - \tau(\lambda), \lambda)$  even when  $|\tau(\lambda)| \ll T$ , since these change  $a$ .

Figure 9(c,d) displays a subset of the second-order joint scattering coefficients for the chirps whose scalograms are shown in Figure 9(a,b). These coefficients do indeed show a maximum along the predicted ridge. At low  $f$  and high  $\mu$ , the approximation does not hold, but for most of the frequency range, it is accurate. We thus see how the chirp rate  $a$  is captured by the joint scattering coefficients in a natural way.

### IV. AUDIO TEXTURE SYNTHESIS

Section III-A showed how mel-spectrograms and time scattering transforms do not adequately capture time-frequency structure. As  $T$  increases, this problem becomes more serious, necessitating the introduction of the joint time-frequency scattering transform. In this section, we illustrate the representational power of this transform using texture synthesis experiments.

With the aim of generating realistic soundtracks of arbitrary duration, audio texture synthesis has many applications in virtual reality and multimedia design [49]. In computational neuroscience, it also offers a testbed for the comparative evaluation of biologically plausible models for auditory perception

[31]. Given a signal  $\mathbf{x}(t)$  and a time-shift invariant representation  $\Phi\mathbf{x}$  of  $\mathbf{x}$ , the texture synthesis problem may be formulated as the minimization of the error  $E(\mathbf{y}) = \|\Phi\mathbf{y} - \Phi\mathbf{x}\|$  between  $\Phi\mathbf{x}$  and the representation  $\Phi\mathbf{y}$  of the synthesized signal  $\mathbf{y}(t)$ . Here,  $\Phi$  can be a scattering transform  $\mathbf{S}$ , a mel-spectrogram  $\mathbf{M}$ , or some other representation. Note that minimizing  $E(\mathbf{y})$  does not imply that  $\mathbf{y}(t)$  approximates  $\mathbf{x}(t)$  in any way; since  $\Phi$  is a time-shift invariant representation, this is not possible. Instead, we expect  $\mathbf{y}(t)$  to contain examples of the time-frequency structures captured in  $\Phi(\mathbf{x})$ .

The state of the art in the domain is held by McDermott and Simoncelli [31], who define  $\Phi$  using a set of summary statistics. These statistics are similar to the time scattering transform as they are calculated using cascades of constant- $Q$  filterbanks and pointwise nonlinearities. However, unlike the scattering transform, which simply averages in time, McDermott and Simoncelli also compute higher-order statistical moments: variance, skewness, kurtosis, and correlation coefficients across frequency bands. These coefficients are very sensitive to outliers in the data, which reduces their applicability to classification.

To synthesize  $\mathbf{y}(t)$ , we first initialize using random Gaussian noise with power spectral density matching the first-order scattering coefficients  $\mathbf{S}_1\mathbf{x}(t, \lambda)$  of the target waveform  $\mathbf{x}(t)$ , since these coefficients are present in all the considered representations. We then iteratively refine the signal by gradient descent [50]. Because the modulus nonlinearity is not convex, the error  $E(\mathbf{y})$  is not convex; consequently, gradient descent only converges to a local minimum of  $E(\mathbf{y})$ . However, this local minimum is typically of low error, with  $E(\mathbf{y})$  equal to around  $0.02 \times \|\Phi\mathbf{x}\|$  for typical audio recordings. We found empirically that the convergence rate is increased using a fixed momentum term and a “bold driver” learning rate policy [51].

Gradient descent in a scattering network can be implemented by backpropagation from deeper to shallower layers. Like in a deep convolutional network, the gradient backpropagation of the convolution with each wavelet  $\mathbf{g}_k(t)$  corresponds to a convolution with the adjoint filter  $\mathbf{g}_k^\dagger(t) = \bar{\mathbf{g}}_k(-t)$ , obtained by time reversal and complex conjugation of  $\mathbf{g}_k(t)$ .

Figure 10 shows the synthesized scalograms of three sounds for various values of  $T$ . Here, time-frequency scattering outperforms time scattering for  $T$  greater than 1 s. Again, we do not expect these synthesized signals to reproduce the originals in the top row due to the imposed time-shift invariance. In particular, speech is more intelligible due to better reconstruction of articulations, individual notes in a musical scale are more salient, and broadband impulses such as dog barks keep their typical amplitude envelopes and inter-onset intervals. Compared to the representation of McDermott and Simoncelli [31], time-frequency scattering achieves similar quality, but does not have the same sensitivity to outliers. Indeed, the contractivity of the wavelet transform and the modulus ensures the scattering transform’s robustness to additive noise [12], [13].

## V. SUPERVISED CLASSIFICATION

We evaluate the performance of the joint time-frequency scattering transform on various classification tasks. It is shown to enjoy significantly greater accuracy compared to baseline

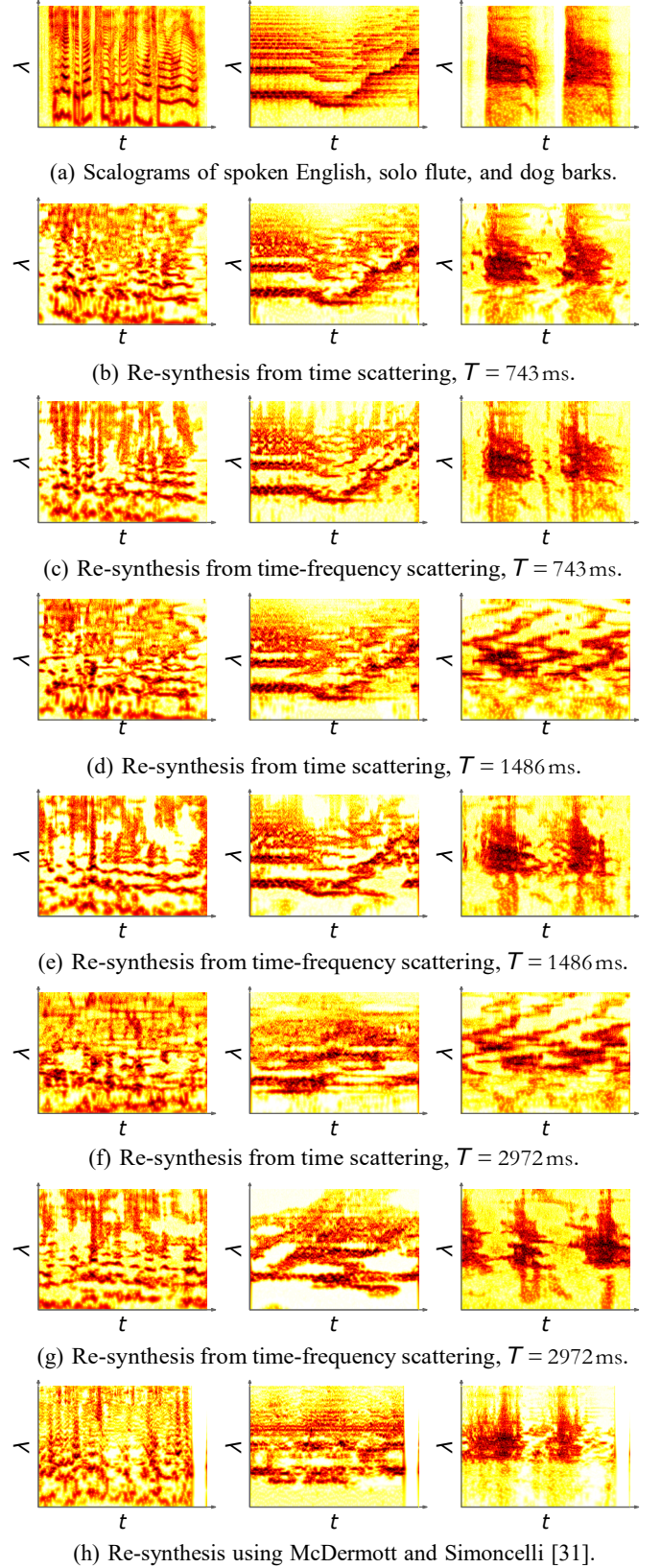


Fig. 10: Scalograms of audio re-synthesis using time scattering, time-frequency scattering, and McDermott and Simoncelli [31]. Synthesis is performed at various time scales  $T$  and inputs: spoken English (left), solo flute (middle), and dog barks (right).

MFCC and time scattering approaches. In fact, the proposed transform performs comparably to state-of-the-art learned convolutional networks whose training requires significant computational resources and large training sets. As a result, the joint scattering transform provides a good alternative when such an expensive training step is infeasible or undesirable.

#### A. Frequency Transposition Invariance

In addition to time-shifting and time-warping, signals are also transformed by frequency-shifting and frequency-warping. Frequency-shifting, also known as frequency transposition, changes the pitch but leaves subband envelopes intact. This shifts the scalogram  $\mathbf{X}(t, \lambda)$  by a fixed amount  $\eta$  in log-frequency, giving  $\mathbf{X}(t, \lambda - \eta)$ . While certain tasks are sensitive to pitch, like speaker identification, others, like speech recognition in non-tonal languages, require invariance to transposition.

The time scattering transform is rendered transposition invariant and stable to frequency-warping by applying a second scattering transform along log-frequency  $\lambda$ . The result is the separable time and frequency scattering transform, introduced in Andén and Mallat [13]. Note that we may skip the averaging step of this second scattering transform. Indeed, the averaging step is a linear map that can be learned by the classifier given enough training data [13].

To render the joint time-frequency scattering transform transposition invariant, we similarly apply a second scattering transform along  $\lambda$  for the first-order coefficients  $\mathbf{S}_1\mathbf{x}$ . For the second order  $\mathbf{S}_2\mathbf{x}$ , however, we simply average along  $\lambda$ , since the two-dimensional wavelet decomposition already captures the relevant frequency structure. The resulting representation then has the necessary transposition invariance and frequency-warping stability properties. Again, if the training set is large enough, the final averaging steps can be learned by the classifier.

#### B. Phone Segment Classification

An individual phone in speech is short, on average 40 ms in duration. For phone identification, we therefore require the invariance scale  $T$  to be of this order. Since  $T$  is small, there is less room for the type of misalignment seen in Section III-A. We therefore expect the joint time-frequency scattering to provide only limited improvement over time scattering.

To evaluate, we use the TIMIT dataset, which contains recordings of spoken phrases, each labeled with its constituent phones and their locations [52]. Given a phone segment, we wish to classify the phone according to the standard protocol [53], [54]. This task is simpler than continuous speech recognition, but provides a good framework for evaluating representations. The training and evaluation sets consist of 3696 and 192 phrases, respectively. We use a 400–phrase validation set to optimize hyperparameters (see Andén and Mallat [13]).

Instead of the raw scattering transform, we use their logarithm, known as the log-scattering transform, as input to the classifier [13]. We compute these coefficients over 192 ms intervals centered on each segment with  $T = 32$  ms. All coefficients are concatenated into a single vector together with the logarithm of the segment duration. This vector is then used for classification. The same processing is also performed for

Representation	Error (%)
Delta-MFCCs	18.3
State of the art [58]	11.9
Time scattering	17.3
Separable time and freq. scattering	16.1
Joint time-freq. scattering	15.7

TABLE I: Error rates for phone segment classification. All representations are computed with  $T = 32$  ms and  $Q = 8$ .

separable time and frequency scattering as well as joint time-frequency scattering. We set the maximum frequency scale  $F$  to 4 octaves. As a baseline, we compute Delta-MFCCs, which supplement standard MFCCs with first and second time derivatives [55]. These are computed with the same windows and concatenation as the log-scattering coefficients.

For each representation, we train a support vector machine (SVM) [56] with a Gaussian kernel. Here and in the following, we use a modified implementation of the LIBSVM library [57].

Results are shown in Table I. Delta-MFCCs have an error rate of 18.3%, while the state-of-the-art representation, a convolutional neural network, achieves 11.9% [58]. The time scattering transform obtains an error rate of 17.3%, which is improved by scattering along log-frequency to give 16.1%. Finally, we obtain an error of 15.7% for the joint time-frequency scattering transform. As mentioned earlier, the amount of time-frequency structure present in an individual phone is small, but there is enough to give a small improvement to the joint transform. This is partly due to the fact that certain phones (such as plosives) are characterized by their onset, which exhibits sophisticated time-frequency structure.

For this task, the joint time-frequency scattering transform does not outperform the state-of-the-art learned convolutional network. Note, however, that the only learning involved for the scattering transform is training the SVM. The scattering network weights are fixed, providing a simpler representation with acceptable performance. Another important difference is that the state-of-the-art result was obtained by simultaneously estimating the labels for all phone segments in an utterance. As a result, this network has access to more context about each segment that it can use to improve classification performance. Combining these two approaches—a scattering transform as input to a more adaptive deep neural network—could yield even better performance as fewer parameters need to be estimated. Indeed, replacing mel-spectrograms by scattering transforms in deep neural networks have improved performance for several tasks [59], [60], [61].

#### C. Musical instrument classification

The timbre of a musical instrument is essentially determined by its shape and materials. Both remain constant during a musical performance. Therefore, musical instruments may be modeled as dynamical systems with constant parameters. The task of musical instrument classification is to retrieve these parameters while remaining invariant to changes in pitch, intensity, and expressive technique induced by the performer.

In a musical instrument, the response of the vibrating body to an excitation is typically nonlinear. As a result, sharp onsets

Representation	Error (%)
Delta-MFCCs	39.3
Time convolutional networks	38.2
Time-frequency convolutional networks	28.3
Spiral convolutional networks [65]	26.0
Time scattering	38.0
Time-frequency scattering	22.0

TABLE II: Error rates for musical instrument classification. All representations are computed with  $T = 3$  s,  $Q = 12$ , and  $F = 4$  octaves.

produce distinctive time-frequency patterns which are not adequately captured by short-term audio descriptors operating on scales  $T \approx 20$  ms, a typical window size for MFCCs. Joint time-frequency scattering, on the other hand, captures such patterns up to the scale  $T \approx 3$  s of a short musical phrase.

To illustrate this, we apply it to automatic instrument classification in solo phrases with a taxonomy of eight instruments. In line with the cross-collection methodology of Bogdanov et al. [62], we train and validate all models on the MedleyDB v1.1 dataset [63] and test them on the solosDb dataset [64]. This is the evaluation setting of Lostanlen and Cella [65].

Results are shown in Table II. It appears that all models which do not explicitly decompose in both time and log-frequency (Delta-MFCCs, time scattering, and a convolutional network of temporal convolutions on the scalogram) perform comparably, with errors around 38%. Introducing decompositions along the log-frequency axis through time-frequency convolutional networks and spiral convolutional networks, we obtain error rates of 28.3% and 26.0, respectively [65]. The improvement likely stems from the fact that musical instruments carry important discriminative information in the temporal evolution of their spectral envelopes as well as frequency modulation structures, both of which are captured by joint decompositions in time and log-frequency. The joint time-frequency scattering transform further reduces the error to 22.0%. The small size of the training set makes optimizing a convolutional network difficult, which may partially explain the improved accuracy of the joint scattering transform compared to the fully learned convolutional networks.

#### D. Acoustic Scene Classification

Environmental sounds and acoustic scenes are characterized by larger-scale time-frequency structures. These recordings typically stretch over several seconds, each composed of shorter sound events which characterize the scene. This could be birdsong in a park, car horns in a street, or the scraping of chairs in a café. To differentiate between different sequences of such events, we must characterize longer-range structures. As discussed above, this is not possible using standard representations, such as MFCCs or time scattering, which do not adequately capture time-frequency structure.

We evaluate the joint scattering transform on three acoustic scene datasets: UrbanSound8K (US8K) [66], ESC-50 [26], and DCASE2013 [67]. US8K and DCASE2013 have 10 classes each, while ESC-50 contains 50 classes, ranging from gun shots and subway stations to crying babies and supermarkets.

Representation	US8K	ESC-50	DCASE2013
Delta-MFCCs [66], [26]	46.0	56.0	42
Salamon and Bello [68]	21.0	—	—
SoundNet [46]	—	25.8	12
L <sup>3</sup> network [47]	—	20.7	7
Time scatt.	26.9 ± 4.1	39.3 ± 2.2	12
Separable time and freq. scatt.	22.8 ± 3.0	26.0 ± 2.7	6
Joint time-freq. scatt.	19.6 ± 2.9	21.8 ± 2.0	5

TABLE III: Average and standard deviation of error rates for scene classification on US8K, ESC-50, and DCASE2013.

Both US8K and ESC-50 contain several thousand recordings of approximate duration 4 s. DCASE2013, on the other hand, contains 100 (public) training samples and 100 (private) evaluation samples, each of duration 30 s. All recordings being relatively long, they may exhibit sophisticated time-frequency structures that are discriminative for classification.

For US8K and ESC-50, we compute scattering transforms with  $Q = 8$  and  $T = 4$  s. We choose a large value for  $T$  because there are long, texture-like structures in this dataset that we would like to characterize. To ensure some transposition invariance, we explicitly average the separable and joint transforms over  $F = 1$  octave (US8K) or  $F = 2$  octaves (ESC-50). Here, we do not want to choose a large frequency scale  $F$  since some pitch information is necessary to distinguish certain sounds. Since  $T$  equals the clip duration, each clip yields a single scattering vector, which is fed into the classifier.

For DCASE2013, we compute scattering transforms with  $Q = 4$ ,  $T = 1.5$  s, and frequency averaging over  $F = 8$  octaves where applicable. We must select parameters different from those of US8K and ESC-50 due to the much smaller size of DCASE2013. Choosing smaller values for  $Q$  and  $T$  limits the complexity of the time-frequency structure captured by the transform, while choosing a large  $F$  and averaging along frequency creates additional invariance to transposition. Since  $T$  is much smaller than the recording duration (30 s), this yields multiple scattering vectors which are classified separately. The overall class is then obtained by majority voting.

Delta-MFCCs are computed for all datasets as a baseline. For each representation, we train a linear SVM with hyperparameters optimized by cross-validation on the training set.

The error for US8K and ESC-50 is calculated through cross-validation on pre-specified folds. For these datasets, we use the data augmentation scheme of Salamon and Bello [68], but without pitch-shifting, since transposition invariance is already enforced. We calculate the DCASE2013 error on the evaluation subset in accordance with previous work [46], [47].

Results are shown in Table III. The Delta-MFCCs have error rates of 46.0%, 56.0%, and 42% for US8K, ESC-50, and DCASE2013, respectively. State-of-the-art convolutional networks, on the other hand, obtain 21.0%, 20.7% and 7%.

The standard time scattering transform yields accuracies of 26.8% (US8K), 39.3% (ESC-50), and 12% (DCASE2013), improving on Delta-MFCCs by better capturing the temporal structure of each subband. Adding a scattering transform along the log-frequency axis improves results to 22.8% (US8K), 26.0% (ESC-50), and 6% (DCASE2013). This improvement is expected since these sounds exhibit significant pitch variability

which is not discriminative to each class.

The joint time-frequency scattering transform performs even better, giving errors of 19.6% (US8K), 21.8% (ESC-50), and 5% (DCASE2013). This is partly because environmental sounds are often characterized by dynamic filters which evolve in time, creating a spectrotemporal filter. The mechanical and biological nature of these sounds also results in frequency modulation. Both phenomena are examples of time-frequency geometry which are well-characterized by the joint scattering transform. From a different perspective, the recordings in these datasets are sensitive to frequency-dependent time-shifts (see Section III-A). Indeed, taking a signal with many transients, such as a jackhammer in a street scene, and misaligning its subbands yields a completely different sound. A representation sensitive to such transformations is therefore expected to perform better.

Again, the joint scattering transform performs comparably to learned convolutional neural networks. However, learned networks require significant computational resources to train and certain expertise in designing the network. Both SoundNet and the  $L^3$  network are pretrained on large external datasets,

requiring several days of computation on graphics processing

units. In contrast, the joint scattering transform has a fixed network structure, so the only training needed is for the SVM, requiring at most a few hours. By considering the invariances of the problem (time-shifting, frequency transposition) and the structures we would like to capture (joint time-frequency geometry), we obtain good performance without costly pretraining.

## VI. CONCLUSION

We introduced a joint time-frequency scattering transform, a time-shift invariant descriptor with state-of-the-art classification performance for a wide range of audio datasets. Important improvements are obtained for classification tasks involving large-scale signal structures. Time-frequency scattering descriptors also recover complex signals including audio textures.

A joint time-frequency scattering has a computational structure similar to deep convolutional networks [69], but is calculated with fixed wavelet filters. It thus requires less training data to obtain accurate classification results. However, when more training examples are available, learned convolutional networks provide state-of-the-art results. Indeed, these networks

adapt the representation to each classification problem. Taking into account prior information on time-frequency geometry could help improve their performance.

## APPENDIX A

**Lemma 1.** Let  $\psi_\lambda(t)$  be as defined in Theorem 1. For  $\mathbf{x}(t) = \exp(2\pi i 2^{at})$ , we then have

$$|\mathbf{x} * \psi_\lambda|(t) = |\phi(\log(2^a) 2^{at-\lambda})| + \varepsilon(t, \lambda), \quad (12)$$

where

$$|\varepsilon(t, \lambda)| \leq C|a|2^{-\lambda} \quad (13)$$

for some constant  $C > 0$  which only depends on  $\psi(t)$ .

*Proof.* If  $a = 0$ ,  $\mathbf{x} * \psi_\lambda(t) = \phi(0) \exp(2\pi i)$ . We therefore assume that  $a \neq 0$ . If  $\text{supp } \psi \subset [-\Delta, \Delta]$ , we have

$$\mathbf{x} * \psi_\lambda(t) = \int_{|u| \leq 2^{-\lambda}\Delta} \exp(2\pi i 2^{a(t-u)}) \psi_\lambda(u) du. \quad (14)$$

For  $u$  close to zero, the derivative of  $2^{a(t-u)}$  is approximately  $-\log(2^a) 2^{at}$ . We exploit this to integrate (14) by parts. Let  $\mathbf{g}(u) = \exp(2\pi i 2^{at}(2^{-au} + u \log(2^a)))$ . We then have

$$\begin{aligned} \mathbf{x} * \psi_\lambda(t) &= \int_{|u| \leq 2^{-\lambda}\Delta} \mathbf{g}(u) \exp(-2\pi i u \log(2^a) 2^{at}) \psi_\lambda(u) du \\ &= \mathbf{g}(2^{-\lambda}\Delta) \phi_\lambda(\log(2^a) 2^{at}) - \int_{|u| \leq 2^{-\lambda}\Delta} \mathbf{g}'(u) \mathbf{l}(u) du, \end{aligned} \quad (15)$$

where  $\mathbf{l}(u) = \int_{-2^{-\lambda}\Delta}^u \exp(-2\pi i v \log(2^a) 2^{at}) \psi_\lambda(v) dv$ .

The magnitude of the second term in (15) is bounded by

$$2\pi 2^{at} |\log(2^a)| (1 - 2^{-|a|2^{-\lambda}\Delta}) 2^{1-\lambda}\Delta \max_{|u| \leq 2^{-\lambda}\Delta} |\mathbf{l}(u)|. \quad (16)$$

Integrating  $\mathbf{l}(u)$  by parts and taking the modulus gives

$$|\mathbf{l}(u)| \leq \frac{|\psi_\lambda|_\infty + |\psi_\lambda^1|_{l_1}}{2\pi |\log(2^a)| 2^{at}} = 2^\lambda \frac{|\psi|_\infty + |\psi^1|_{l_1}}{2\pi |\log(2^a)| 2^{at}},$$

since  $a \neq 0$ ,  $|\psi_\lambda|_\infty = |\psi|_\infty$  and  $|\psi_\lambda^1|_{l_1} = 2^{-\lambda} |\psi^1|_{l_1}$ .

Plugging this into (16), we obtain

$$2\Delta(1 - 2^{-|a|2^{-\lambda}\Delta}) (|\psi|_\infty + |\psi^1|_{l_1}). \quad (17)$$

Given that  $1 - 2^{-u} < \log(2)u$  for all  $u > 0$ , this simplifies to

$$2 \log(2) \Delta^2 (|\psi|_\infty + |\psi^1|_{l_1}) |a| 2^{-\lambda}. \quad (18)$$

Since  $|\mathbf{g}(u)| = 1$ , the modulus of the first term in (15) is  $|\phi_\lambda(\log(2^a) 2^{at})| = |\phi(\log(2^a) 2^{at-\lambda})|$ . The triangle inequality then establishes (12) with  $|\varepsilon(t, \lambda)|$  bounded by (18).  $\square$

**Lemma 2.** Define  $\psi_\lambda(t)$ ,  $\Psi_{\mu, \varepsilon, s}(t, \lambda)$ , and  $C_0$  as in Theorem 1 and let

$$t_0(\lambda) = \frac{\lambda}{a} - \frac{\log \log 2^a}{\log 2^a}.$$

Given

$$\mathbf{Y}(t, \lambda) = |\phi(\log(2^a) 2^{at-\lambda})|,$$

its two-dimensional wavelet modulus decomposition satisfies

$$\begin{aligned} |\mathbf{Y} * \Psi_{\mu, \varepsilon, s}|(t, \lambda) &= \frac{C_0 |\psi_\mu^{(t)}|(t - t_0(\lambda))}{a} \phi^{(t)}\left(\frac{s 2^{\mu-\varepsilon}}{a}\right) + \varepsilon(t, \lambda, \mu, f, s), \end{aligned} \quad (19)$$

where  $|\psi_\mu^{(t)}|(t) = 2^\mu |\psi^{(t)}|(2^\mu t)$ , and

$$|\varepsilon(t, \lambda, \mu, f, s)| \leq C(2^{2\mu} |a|^{-2} + 2^{2\mu-\varepsilon} |a|^{-2}),$$

for some  $C > 0$  depending only on  $\psi(t)$ ,  $\psi^{(t)}(t)$ ,  $\psi^{(t)}(\lambda)$ .

*Proof.* Since  $|\phi(\omega)|$  is maximized at  $\omega = 1$ , fixing  $\lambda$ , the maximum of  $\mathbf{Y}(t, \lambda)$  is at  $t_0(\lambda)$ . For small enough  $\mu$ ,  $\mathbf{Y}(t, \lambda)$  approximates a Dirac delta function centered at  $t_0(\lambda)$ . We exploit this when convolving  $\mathbf{Y}(t, \lambda)$  by  $\psi_\mu^{(t)}(t)$ .

Approximating  $\psi^{(t)}(u)$  with its value at  $u = t - t_0(\lambda)$  gives

$$\begin{aligned} \mathbf{Y}(\cdot, \lambda) * \psi^{(t)}(t) &= \int_{\mathbb{R}} \mathbf{Y}(t - u, \lambda) \psi_\mu^{(t)}(u) du \\ &= \int_{\mathbb{R}} |\phi(\log(2^a) 2^{a(t-u)-\lambda})| \times \\ &\quad (\psi_\mu^{(t)}(t - t_0(\lambda)) + \varepsilon_1(t, \lambda, \mu, u)) du, \end{aligned}$$



where  $|\varepsilon_1(t, \lambda, \mu, u)| \leq |t - t_0(\lambda) - u| |\psi_\mu^{(1)}|_\infty$ . Setting  $\varepsilon_2(t, \lambda, \mu) = \int_{\mathbb{R}} |\phi(\log(2^a) 2^{a(t-u)-\lambda})| |\varepsilon_1(t, \lambda, \mu, u)| du$  gives  $Y(\cdot, \lambda) * \psi_\mu^{(1)}(t) = c_0 \sigma^{-1} \psi_\mu^{(1)}(t - t_0(\lambda)) + \varepsilon_2(t, \lambda, \mu)$ , (20)

using a change of variables, where  $c_0 = \int_{\mathbb{R}} |\phi(2^u)| du < \infty$ . We bound  $\varepsilon_2(t, \lambda, \mu)$  through

$$|\varepsilon_2(t, \lambda, \mu)| \leq \int_{\mathbb{R}} |\phi(\log(2^a) 2^{a(t-u)-\lambda})| |\varepsilon_1(t, \lambda, \mu, u)| du \leq |\psi_\mu^{(1)}|_\infty \int_{\mathbb{R}} |\phi(\log(2^a) 2^{a(t-u)-\lambda})| |t - t_0(\lambda) - u| du.$$

The change of variables  $t - t_0(\lambda) - u \mapsto \sigma^{-1}u$  now gives

$$|\varepsilon_2(t, \lambda, \mu)| \leq |\sigma|^{-2} 2^{2\mu} |\psi^{(1)}|_\infty \int_{\mathbb{R}} |\phi(2^u)| |u| du, \quad (21)$$

where we have used  $|\psi_\mu^{(1)}|_\infty = 2^{2\mu} |\psi^{(1)}|_\infty$ .

We now convolve (20) by  $\psi_{\varepsilon, s}^{(f)}(\lambda) = 2^\varepsilon \psi^{(f)}(s 2^\varepsilon \lambda)$ . At high  $f$ , this wavelet will mostly capture phase variation. To see this, we factorize  $\psi_\mu^{(1)}(t)$  into an envelope and a phase, yielding  $|\psi_\mu^{(1)}(t) \exp(2\pi i 2^\mu t)|$ . The convolution then becomes

$$\begin{aligned} c_0 \sigma^{-1} \psi_\mu^{(1)}(t - t_0(\cdot)) * \psi_{\varepsilon, s}^{(f)}(\lambda) \\ = c_0 \sigma^{-1} \int_{\mathbb{R}} |\psi_\mu^{(1)}(t - t_0(\lambda - \gamma))| \times \\ \exp(2\pi i 2^\mu (t - t_0(\lambda - \gamma))) \psi_{\varepsilon, s}^{(f)}(\gamma) d\gamma. \end{aligned} \quad (22)$$

We now make the approximation

$|\psi_\mu^{(1)}(t - t_0(\lambda - \gamma))| = |\psi_\mu^{(1)}(t - t_0(\lambda)) + \varepsilon_3(t, \lambda, \mu, \gamma)|$ , where  $|\varepsilon_3(t, \lambda, \mu, \gamma)| \leq |\psi_\mu^{(1)}|_\infty |t_0(\lambda - \gamma) - t_0(\lambda)|$ . Plugging this into (22), we obtain

$$\begin{aligned} c_0 \sigma^{-1} |\psi_\mu^{(1)}(t - t_0(\lambda))| \times \\ \int_{\mathbb{R}} \exp(2\pi i 2^\mu (t - t_0(\lambda - \gamma))) \psi_{\varepsilon, s}^{(f)}(\gamma) d\gamma + \varepsilon_4(t, \lambda, \mu, f, s), \end{aligned} \quad (23)$$

where  $\varepsilon_4(t, \lambda, \mu, f, s)$  equals

$$c_0 \sigma^{-1} \int_{\mathbb{R}} \varepsilon_3(t, \lambda, \mu, \gamma) \exp(2\pi i 2^\mu (t - t_0(\lambda - \gamma))) \psi_{\varepsilon, s}^{(f)}(\gamma) d\gamma.$$

Since  $t_0(\lambda - \gamma) = t_0(\lambda) - \sigma^{-1}\gamma$ , the first term in (23) is

$$c_0 \sigma^{-1} |\psi_\mu^{(1)}(t - t_0(\lambda))| e^{i 2^\mu (t - t_0(\lambda))} \psi_{\varepsilon, s}^{(f)}(-2 \sigma^{-1} \gamma). \quad (24)$$

The same property of  $t_0(\lambda)$  lets us bound  $\varepsilon_4(t, \lambda, \mu, f, s)$  by

$$\begin{aligned} |\varepsilon_4(t, \lambda, \mu, f, s)| \leq c_0 |\sigma|^{-2} |\psi_\mu^{(1)}|_\infty \int_{\mathbb{R}} |\gamma| |\psi_{\varepsilon, s}^{(f)}(\gamma)| d\gamma \\ = c_0 |\sigma|^{-2} 2^{2\mu - \varepsilon} |\psi^{(1)}|_\infty \int_{\mathbb{R}} |\gamma| |\psi^{(f)}(\gamma)| d\gamma, \end{aligned} \quad (25)$$

which follows from change of variables and from  $|\psi_\mu^{(1)}|_\infty = 2^{2\mu} |\psi^{(1)}|_\infty$ . We must also convolve  $\varepsilon_2(t, \lambda, \mu)$  with  $\psi_{\varepsilon, s}^{(f)}(\lambda)$ .

Since  $|\psi_{\varepsilon, s}^{(f)}|_1 = |\psi^{(f)}|_1$  for all  $f, s$ , we have

$$|\varepsilon_2(t, \cdot, \mu) * \psi_{\varepsilon, s}^{(f)}(\lambda)| \leq |\varepsilon_2(t, \cdot, \mu)|_\infty |\psi^{(f)}|_1. \quad (26)$$

Combining (20) with (24) and taking the modulus yields (19) since  $\phi_{\varepsilon, s}(\omega) = \phi(s 2^{-\varepsilon} \omega)$ , where the bound on

$\varepsilon(t, \lambda, \mu, f, s)$  follows from (21), (25), (26), and the triangle inequality.  $\square$

*Proof of Theorem 1.* Lemma 1 gives

$$\mathbf{X}(t, \lambda) = |x * \psi_\lambda|(t) = |\phi(\log(2^a) 2^{a(t-\lambda)})| + \varepsilon_1(t, \lambda),$$

where  $|\varepsilon_1(t, \lambda)| \leq C_1 |a| 2^{-\lambda}$  for some  $C_1 > 0$ . We now convolve  $\mathbf{X}(t, \lambda)$  with  $\psi_\mu^{(1)}(t)$  in time  $\psi_{\varepsilon, s}^{(f)}(\lambda)$  in log-frequency and take the modulus. Lemma 2 approximates the convolution of the first term. For the second term, we observe that

$$|\varepsilon_1(\cdot, \lambda) * \psi_\mu^{(1)}| \leq |\varepsilon_1(\cdot, \lambda)|_\infty |\psi_\mu^{(1)}|_1 \leq C_2 |a| 2^{-\lambda},$$

for some  $C_2 > 0$ , since  $|\psi_\mu^{(1)}|_1 = |\psi^{(1)}|_1$  for all  $\mu$ . Now,

$$\begin{aligned} |\varepsilon_1 * \psi_{\mu, \varepsilon, s}(t, \lambda)| &\leq C_2 |a| \int_{\mathbb{R}} 2^{-(\lambda - \mu)} |\psi_{\varepsilon, s}^{(f)}(\mu)| d\mu \\ &= C_2 |a| 2^{-\lambda} \int_{\mathbb{R}} 2^\mu |\psi_{\varepsilon, s}^{(f)}(\mu)| d\mu \\ &\leq C_2 |a| 2^{-\lambda} 2^{2^{-A}} |\psi_{\varepsilon, s}^{(f)}|_1 = C_3 |a| 2^{-\lambda + 2^{-A}}, \end{aligned}$$

for some  $C_3 > 0$ , since  $\psi_{\varepsilon, s}^{(f)}$  is supported on  $[-A, A]$ .

As a result,

$$\begin{aligned} |\mathbf{X} * \psi_{\mu, \varepsilon, s}(t, \lambda)| \\ = \frac{c_0}{a} |\psi_\mu^{(1)}(t - t_0(\lambda))| \phi^{(f)}\left(-\frac{s 2^{\mu - \varepsilon}}{a}\right) + \varepsilon_2(t, \lambda, \mu, f, s), \end{aligned} \quad (27)$$

where

$$|\varepsilon_2(t, \lambda, \mu, f, s)| \leq C (|a| 2^{-\lambda + 2^{-A}} + |a|^{-2} 2^{2\mu} + |a|^{-2} 2^{2\mu - \varepsilon}).$$

Since this bound is constant in  $t$  and  $|\psi_\mu^{(1)}|_1 = |\psi^{(1)}|_1$  for all  $\mu$ , it still holds after convolving (27) with  $\psi_T^{(f)}(t)$ .  $\square$

## ACKNOWLEDGMENTS

The authors would like to thank J. Salamon for sharing his data augmentation code and A. Barnett for helpful discussions on oscillatory integrals.

## REFERENCES

- [1] S. Davis and P. Mermelstein, "Comparison of parametric representations for monosyllabic word recognition in continuously spoken sentences," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 28, no. 4, pp. 357–366, 1980.
- [2] H. Hermansky, "The modulation spectrum in the automatic recognition of speech," in *Proc. ASRU*, IEEE, Dec 1997, pp. 140–147.
- [3] J. Thompson and L. Atlas, "A non-uniform modulation transform for audio coding with increased time resolution," in *Proc. ICASSP*, vol. 5, IEEE, 2003, pp. 397–400.
- [4] M. Slaney and R. Lyon, "On the importance of time–A temporal representation of sound," in *Visual representations of speech signals*, S. B. M. Cooke and M. Crawford, Eds. Wiley, 1993, pp. 95–116.
- [5] R. D. Patterson, "Auditory images: How complex sounds are represented in the auditory system," *Acoustical Science and Technology*, vol. 21, no. 4, pp. 183–190, 2000.
- [6] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proc. IEEE*, vol. 86, no. 11, pp. 2278–2324, Nov 1998.
- [7] J. Schmidhuber, "Deep learning in neural networks: An overview," *Neural Networks*, vol. 61, pp. 85–117, 2015.
- [8] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, pp. 436–444, May 2015.
- [9] A. Graves, A.-R. Mohamed, and G. Hinton, "Speech recognition with deep recurrent neural networks," in *Proc. ICASSP*, IEEE, 2013, pp. 6645–6649.



- [10] Y. LeCun, K. Kavukcuoglu, and C. Farabet, "Convolutional networks and applications in vision," in *Proc. ISCS*. IEEE, 2010.
- [11] V. Nair and G. E. Hinton, "Rectified linear units improve restricted Boltzmann machines," in *Proc. ICML*, 2010, pp. 807–814.
- [12] S. Mallat, "Group invariant scattering," *Comm. Pure Appl. Math.*, vol. 65, no. 10, pp. 1331–1398, 2012.
- [13] J. Andén and S. Mallat, "Deep scattering spectrum," *IEEE Trans. Signal Process.*, vol. 62, pp. 4114–4128, 2014.
- [14] V. Chudáček, J. Andén, S. Mallat, P. Abry, and M. Doret, "Scattering transform for intrapartum fetal heart rate variability fractal analysis: A case-control study," *IEEE Trans. Biomed. Eng.*, vol. 61, no. 4, pp. 1100–1108, 2014.
- [15] R. Talmon, S. Mallat, H. Zaveri, and R. R. Coifman, "Manifold learning for latent variable inference in dynamical systems," *IEEE Trans. Signal Process.*, vol. 63, no. 15, pp. 3843–3856, 2015.
- [16] J. Sulam, Y. Romano, and R. Talmon, "Dynamical system classification with diffusion embedding for ECG-based person identification," *Signal Processing*, vol. 130, pp. 403–411, 2017.
- [17] T. Chi, P. Ru, and S. Shamma, "Multiresolution spectrotemporal analysis of complex sounds," *J. Acoust. Soc. Am.*, vol. 118, no. 2, pp. 887–906, 2005.
- [18] N. Mesgarani, M. Slaney, and S. Shamma, "Discrimination of speech from nonspeech based on multiscale spectro-temporal modulations," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 14, no. 3, pp. 920–930, 2006.
- [19] J. C. Brown, "An efficient algorithm for the calculation of a constant Q transform," *J. Acoust. Soc. Am.*, vol. 92, no. 5, p. 2698, 1992.
- [20] T. Dau, B. Kollmeier, and A. Kohlrausch, "Modeling auditory processing of amplitude modulation. I. Detection and masking with narrow-band carriers," *J. Acoust. Soc. Am.*, vol. 102, no. 5, pp. 2892–2905, 1997.
- [21] E. C. Smith and M. S. Lewicki, "Efficient auditory coding," *Nature*, vol. 439, no. 7079, p. 978, 2006.
- [22] S. Mallat, *A wavelet tour of signal processing*. Academic Press, 1999.
- [23] I. Daubechies, *Ten Lectures on Wavelets*. SIAM, 1992.
- [24] B. Logan, "Mel frequency cepstral coefficients for music modeling," in *Proc. ISMIR*, 2000.
- [25] P. Somervuo, A. Härmä, and S. Fagerlund, "Parametric representations of bird sounds for automatic species recognition," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 14, no. 6, pp. 2252–2263, Nov 2006.
- [26] K. J. Piczak, "ESC: Dataset for environmental sound classification," in *Proc. MM*. ACM, 2015, pp. 1015–1018.
- [27] M. Dörfler, R. Bammer, and T. Grill, "Inside the spectrogram: Convolutional neural networks in audio processing," in *Proc. SampTA*, 2017, pp. 152–155.
- [28] M. Dörfler, T. Grill, R. Bammer, and A. Flexer, "Basic filters for convolutional neural networks applied to music: Training or design?" 2017, submitted. arXiv:1709.02291.
- [29] A. Papoulis, *Signal Analysis*. McGraw-Hill Education, 1977.
- [30] I. Waldspurger, "Exponential decay of scattering coefficients," in *Proc. SampTA*, 2017, pp. 143–146.
- [31] J. McDermott and E. Simoncelli, "Sound texture perception via statistics of the auditory periphery: Evidence from sound synthesis," *Neuron*, vol. 71, no. 5, pp. 926–940, 2011.
- [32] M. Schädler, B. Meyer, and B. Kollmeier, "Spectro-temporal modulation subspace-spanning filter bank features for robust automatic speech recognition," *J. Acoust. Soc. Am.*, vol. 131, no. 5, pp. 4134–4151, 2012.
- [33] M. Schädler and B. Kollmeier, "Separable spectro-temporal Gabor filter bank features: Reducing the complexity of robust features for automatic speech recognition," *J. Acoust. Soc. Am.*, vol. 137, no. 4, pp. 2047–2059, 2015.
- [34] M. Kleinschmidt and D. Gelbart, "Improving word accuracy with Gabor feature extraction," in *Proc. Interspeech*, 2002, pp. 25–28.
- [35] K. Siedenburg, I. Fujinaga, and S. McAdams, "A comparison of approaches to timbre descriptors in music information retrieval and music psychology," *J. New Music Res.*, vol. 45, no. 1, pp. 27–41, 2016.
- [36] J. Bruna and S. Mallat, "Invariant scattering convolution networks," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 35, no. 8, pp. 1872–1886, 2013.
- [37] L. Sifre and S. Mallat, "Rotation, scaling and deformation invariant scattering for texture discrimination," in *Proc. CVPR*. IEEE, 2013, pp. 1233–1240.
- [38] Z. Tüske, P. Golik, R. Schlüter, and H. Ney, "Acoustic modeling with deep neural networks using raw time signal for LVCSR," in *Proc. Interspeech*, 2014, pp. 890–894.
- [39] T. N. Sainath, R. J. Weiss, A. Senior, K. W. Wilson, and O. Vinyals, "Learning the speech front-end with raw waveform CLDNNs," in *Proc. Interspeech*, 2015, pp. 1–5.
- [40] S. Fujieda, K. Takayama, and T. Hachisuka, "Wavelet convolutional neural networks," *arXiv preprint arXiv:1805.08620*, 2018.
- [41] E. Oyallon, E. Belilovsky, S. Zagoruyko, and M. Valko, "Compressing the input for CNNs with the first-order scattering transform," in *Proc. ECCV*. Springer, 2018, pp. 305–320.
- [42] C. Shi and C.-M. Pun, "3D multi-resolution wavelet convolutional neural networks for hyperspectral image classification," *Information Sciences*, vol. 420, pp. 49–65, 2017.
- [43] Y. Liao, X. Zeng, and W. Li, "Wavelet transform based convolutional neural network for gearbox fault classification," in *Prognostics and System Health Management Conference*. IEEE, 2017, pp. 1–6.
- [44] P. Hamel, M. E. P. Davies, K. Yoshii, and M. Goto, "Transfer learning in MIR: Sharing learned latent representations for music audio classification and similarity," in *Proc. ISMIR*, 2013, pp. 9–15.
- [45] A. van den Oord, S. Dieleman, and B. Schrauwen, "Transfer learning by supervised pre-training for audio-based music classification," in *Proc. ISMIR*, 2014, pp. 29–34.
- [46] Y. Aytar, C. Vondrick, and A. Torralba, "SoundNet: Learning sound representations from unlabeled video," in *Proc. NIPS*, 2016, pp. 892–900.
- [47] R. Arandjelović and A. Zisserman, "Look, listen and learn," in *Proc. ICCV*. IEEE, 2017, pp. 609–617.
- [48] M. A. Richards, *Fundamentals of radar signal processing*. McGraw-Hill Education, 2005.
- [49] D. Schwarz, "State of the art in sound texture synthesis," in *Proc. DAFx*, 2011, pp. 221–232.
- [50] J. Bruna and S. Mallat, "Audio texture synthesis with scattering moments," 2013, unpublished. arXiv:1311.0407.
- [51] I. Sutskever, J. Martens, G. Dahl, and G. Hinton, "On the importance of initialization and momentum in deep learning," in *Proc. ICML*, 2013.
- [52] W. Fisher, G. Doddington, and K. Goudie-Marshall, "The DARPA speech recognition research database: specifications and status," in *Proc. DARPA Workshop on Speech Recognition*, 1986, pp. 93–99.
- [53] K.-F. Lee and H.-W. Hon, "Speaker-independent phone recognition using hidden Markov models," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 37, no. 11, pp. 1641–1648, 1989.
- [54] P. Clarkson and P. Moreno, "On the use of support vector machines for phonetic classification," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, vol. 2. IEEE, 1999, pp. 585–588.
- [55] S. Furui, "Speaker-independent isolated word recognition using dynamic features of speech spectrum," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 34, no. 1, pp. 52–59, 1986.
- [56] C. Cortes and V. Vapnik, "Support-vector networks," *Mach. Learn.*, vol. 20, no. 3, pp. 273–297, 1995.
- [57] C. Chang and C. Lin, "LIBSVM: A library for support vector machines," *ACM Trans. on Intell. Syst. and Technol.*, vol. 2, pp. 27:1–27:27, 2011.
- [58] M. Ratajczak, S. Tschieschek, and F. Pernkopf, "Frame and segment level recurrent neural networks for phone classification," in *Proc. Interspeech*, 2017, pp. 1318–1322.
- [59] V. Peddinti, T. N. Sainath, S. Maymon, B. Ramabhadran, D. Nahamoo, and V. Goel, "Deep scattering spectrum with deep neural networks," in *Proc. ICASSP*. IEEE, 2014, pp. 210–214.
- [60] N. Zeghidour, G. Synnaeve, M. Versteegh, and E. Dupoux, "A deep scattering spectrum-deep siamese network pipeline for unsupervised acoustic modeling," in *Proc. ICASSP*. IEEE, 2016, pp. 4965–4969.
- [61] E. Oyallon, E. Belilovsky, and S. Zagoruyko, "Scaling the scattering transform: Deep hybrid networks," in *Proc. ICCV*, 2017.
- [62] D. Bogdanov, A. Porter, P. Herrera, and X. Serra, "Cross-collection evaluation for music classification tasks," in *Proc. ISMIR*, 2016, pp. 379–385.
- [63] R. Bittner, J. Salamon, M. Tierney, M. Mauch, C. Cannam, and J. Bello, "MedleyDB: A multitrack dataset for annotation-intensive MIR research," in *Proc. ISMIR*, 2014, pp. 155–160.
- [64] C. Joder, S. Essid, and G. Richard, "Temporal integration for audio classification with application to musical instrument classification," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 17, no. 1, pp. 174–186, 2009.
- [65] V. Lostanlen and C.-E. Cella, "Deep convolutional networks on the pitch spiral for music instrument classification," in *Proc. ISMIR*, 2016, pp. 612–618.
- [66] J. Salamon, C. Jacoby, and J. P. Bello, "A dataset and taxonomy for urban sound research," in *Proc. MM*. ACM, 2014.
- [67] D. Stowell, D. Giannoulis, E. Benetos, M. Lagrange, and M. D. Plumbley, "Detection and classification of acoustic scenes and events," *IEEE Trans. Multimed.*, vol. 17, no. 10, pp. 1733–1746, 2015.
- [68] J. Salamon and J. P. Bello, "Deep convolutional neural networks and data augmentation for environmental sound classification," *IEEE Signal Process. Lett.*, vol. 24, no. 3, pp. 279–283, 2017.
- [69] S. Mallat, "Understanding deep convolutional networks," *Phil. Trans. R. Soc. A*, vol. 374, no. 2065, p. 20150203, 2016.