# Learning neural audio features without supervision

*Sarthak Yadav[1], Neil Zeghidour[2]*

[1]University of Glasgow, UK
[2]Google Research, Paris, France

s.yadav.2@student.gla.ac.uk, neilz@google.com

## Abstract

Deep audio classification, traditionally cast as training a deep neural network on top of mel-filterbanks in a supervised fashion, has recently benefited from two independent lines of work. The first one explores "learnable frontends", i.e., neural modules that produce a learnable time-frequency representation, to overcome limitations of fixed features. The second one uses self-supervised learning to leverage unprecedented scales of pre-training data. In this work, we study the feasibility of combining both approaches, i.e., pre-training learnable frontend jointly with the main architecture for downstream classification. First, we show that pretraining two previously proposed frontends (SincNet and LEAF) on Audioset drastically improves linear-probe performance over fixed mel-filterbanks, suggesting that learnable time-frequency representations can benefit self-supervised pre-training even more than supervised training. Surprisingly, randomly initialized learnable filterbanks outperform mel-scaled initialization in the self-supervised setting, a counter-intuitive result that questions the appropriateness of strong priors when designing learnable filters. Through exploratory analysis of the learned frontend components, we uncover crucial differences in properties of these frontends when used in a supervised and self-supervised setting, especially the affinity of self-supervised filters to diverge significantly from the mel-scale to model a broader range of frequencies.

**Index Terms**: self-supervised learning, audio, sound, learnable audio frontend

## 1. Introduction

Mel-filterbanks have long remained the features of choice for machine learning and audio signal processing. Mel-filterbanks first pass a spectrogram through a bank of triangular bandpass filters that are logarithmically scaled to model human perception of pitch [1], referred to as the mel-scale, providing shift-invariance and robustness to deformations [2]. The dynamic range of the resulting coefficients is then typically compressed by a logarithm, to replicate human sensitivity to loudness. However, mel-filterbanks suffer from inherent limitations of fixed features. Not only the mel-scale has been reconsidered several times over its history [3, 4] but logarithmic compression has also been shown to be suboptimal with respect to $n^{th}$ root non-linearities [5, 6]. These limitations, paired with the development of deep learning methods, have fostered a growing corpus of work on learning an audio frontend from raw waveform signals. Most of these contributions focus on learning a filterbank as an alternative to mel-scale triangular filters [7, 8, 9, 10, 11, 12, 13], while some propose replacing logarithmic compression by a trainable non-linearity [14, 15], or learning all operations (filtering, pooling, compression) in an end-to-end fashion [16].

Learnable alternative to fixed features are typically trained and evaluated in a supervised setting. In that context, the lack of strong priors (e.g., mel-scale, log compression) is compensated by strong labelling of the training data, which allows exploring the space of models to reach a parameterization that can consistently outperform fixed features [16]. In recent years, there has been a rise in learning audio representations with self-supervision [17, 18, 19, 20, 21, 22, 23], capable of learning general-purpose acoustic representations without extensive manual data annotation and have demonstrated excellent few-shot learning performance. However, these methods either adopt large-scale architectures without any structure that would be akin to audio features [17, 18] or rely on mel-filterbanks [19, 24, 23]. Thus, the question of learning strong audio frontends without labelled data still remains unanswered.

In this work, we explore learning audio features in a self-supervised setting, to study whether 1) learnable audio frontends can outperform audio features even without labelled data 2) at convergence, the learned filterbanks and compression nonlinearity differ between the supervised and the self-supervised setting. More specifically, we explore two previously proposed learnable audio frontends: SincNet [12], and LEAF [16], both in a supervised multi-label classification setting, and in a self-supervised framework using contrastive learning, on Audioset [25]. On a linear probe transfer task, learnable frontends significantly outperform fixed mel-filterbanks, confirming experimental results previously limited to the supervised setting [12, 16]. When exploring the role of filter initialization, we observe that while randomly initialized LEAF and SincNet frontends do not improve supervised learning performance, they surprisingly improve performance for self-supervised training over their mel-scale initialized counterpart, a counter-intuitive result that suggests reconsidering the mel-scale. Moreover, frequency analysis of the learned filters shows that they converge to similar configurations in a supervised setting, regardless of their initialization. In contrast, self-supervised learning converges to more diverse configurations, with the randomly initialized kernels modelling a broader range of frequencies. In the upcoming sections, we present the proposed approach, followed by an in-depth analysis of our experimental setting and results on the benchmarked datasets. Finally, we conclude with an inspection of the learned components of the audio frontends.

## 2. Method

The proposed model has two core components: i) a common *neural backbone* which consists of an audio frontend followed by an encoder, and ii) a task-dependent shallow MLP head. Our model is trained either in a supervised or a self-supervised setting, and is described further in this section.

### 2.1. Audio frontends

A frontend $\mathcal{F}_\psi$ maps an input waveform $x \in \mathbb{R}^T$ sampled at a frequency $F_s$ (Hz) to a $(M \times N)$-dimensional feature space,
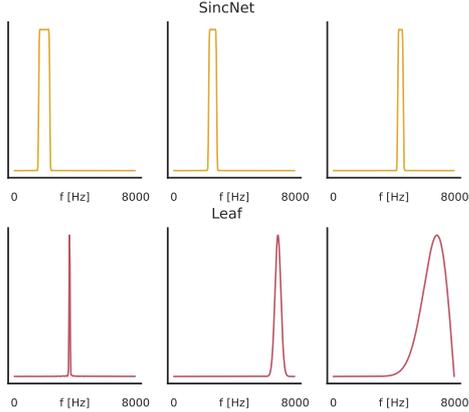
Figure 1: *Frequency response of randomly selected SincNet filters and LEAF filters, with truncated normal and uniform initialization schemes, respectively.*

where $M$ denotes the number of temporal frames and $N$ represents the number of frequency bins. Along with log-compressed mel-filterbanks as our baseline, we explore two recent learnable audio frontends, SincNet and LEAF.

**SincNet** [12] utilizes learnable sinc function based band-pass filters with a rectangular frequency response (see Figure 1) that are parameterized by cut-off frequencies $[f_1, f_2]$:

$$g[n, f_1, f_2] = 2f_2 sinc(2\pi f_2 n) - 2f_1 sinc(2\pi f_1 n), \quad (1)$$

where $sinc(x) = sin(x)/x$, and the cut-off frequencies are initialized on the mel-scale. This is followed by a leaky ReLU non-linearity and a max-pooling layer for down-sampling.

**LEAF** [16] proposes a learnable frontend that utilizes convolution with $N$ learnable complex-valued Gabor filters with a gaussian frequency response (Fig 1), each of length $W$, which are parameterized by their center frequencies $\eta_n$ and inverse bandwidths $\sigma_n$ as follows:

$$\varphi_n(t) = e^{i2\pi\eta_n t} \frac{1}{\sqrt{2\pi}\sigma_n} e^{-\frac{t^2}{2\sigma_n^2}}, \quad (2)$$

where $n = 1, \ldots, N$ and $t = -W/2, \ldots, W/2$, and the center frequencies are initialized in the $[0, 1/2]$ range (normalized units) on the mel-scale. This is followed by a squared modulus operator, which gets the input back to the real-valued domain, followed by a learnable Gaussian low-pass pooling layer. Finally, a Per-Channel Energy Normalization layer [14], applies learned compression and normalization to the coefficients. Hence, while SincNet only learns the filterbank, LEAF jointly learns filtering, pooling, compression and normalization.

**Random initialization schemes**: Previous work [8, 10] has examined how random initialization affects performance for learnable audio frontends, albeit for other frontends and only in the supervised setting. In the same light, we explore using filters randomly initialized from uniform and truncated-normal distributions for both the SincNet and LEAF audio frontends across both the supervised and the self-supervised training regimes.

### 2.2. Backbone architecture

The complete *neural backbone* $f(.)$ consists of the audio frontend $\mathcal{F}_\psi$ followed by an encoder $enc(.)$, and maps the input waveform $x \in \mathbb{R}^T$ into a latent representation $h = f(x) =$

$enc(\mathcal{F}_\psi(x)) \in \mathbb{R}^D$ and is common across all the experiments. A linear classifier head is added on top of the neural backbone in the supervised setting. We use the lightweight EfficientNetB0 [26] network with around $4M$ parameters as the convolutional encoder in this work.

### 2.3. Self-supervised pre-training

For self-supervised pretraining, we use COLA [23], a contrastive approach for learning general-purpose audio representations. Given the similarity function

$$s(x, x') = g(f(x))^\top W g(f(x')), \quad (3)$$

where $f(.)$ is the *neural backbone*, $g(.)$ is a shallow neural network that maps $h$ onto a space $z = g(h) \in \mathbb{R}^G$, and $W \in \mathbb{R}^G$ are bilinear similarity parameters, COLA learns a latent space such that the similarity $s(x, x^+)$ between an anchor-positive pair is greater than the similarity $s(x, x^-)$ between the same anchor and negative samples (unrelated examples) by optimizing the following objective function:

$$\mathcal{L} = -\log \frac{\exp\left(s(x, x^+)\right)}{\sum\limits_{x^- \in \mathcal{X}^-(x) \cup \{x^+\}} \exp\left(s(x, x^-)\right)}, \quad (4)$$

where $\mathcal{X}^-(x)$ refers to the set of negative samples. As opposed to existing methods that use perturbations on the anchor/positive samples [27, 24, 22], in COLA, the anchor and its corresponding positive samples are segments from the same audio clip without any perturbations applied, while segments from all other clips in the batch act as negatives, removing the need for benchmarking multiple augmentation strategies.

## 3. Experiments

To evaluate the relative improvements of learnable frontends over fixed features based on the level of supervision, we compare SincNet, LEAF and mel-filterbanks in two settings: 1) Supervised training on the AudioSet dataset [25], and 2) self-supervised COLA pre-training on AudioSet followed by linear-probing (only training a linear classifier on top of a frozen backbone) on the SpeechCommandsV2 dataset [28]. We choose SpeechCommandsV2 given that it has a moderate number of fixed-sized data samples (over 105000, 1-second audio clips across 35 classes) which simplifies experimentation, and as shown previously in [23] it poses a challenging transfer task.

### 3.1. Implementation details

The input signal sampled at $F_s = 16$ kHz is used for all experiments. Log-compressed mel-filterbanks with 64 filters with a window of 25 ms and a stride of 10 ms are used as a baseline. To facilitate comparison, LEAF and SincNet frontends also have $N = 64$ learnable filters with kernels of width $W = 401$ ($\approx 25$ ms at 16 kHz), and their corresponding pooling layers are initialized with a stride of 160 samples (10 ms at 16 kHz). For per-channel energy normalization, the "sPCEN" variant as proposed in [16] is used, and is here on referred to simply as "PCEN". To study how the per-channel energy normalization in the LEAF frontend affects performance and model characteristics, we also run an extra set of experiments with a fixed PCEN layer for each initialization scheme.

The neural *backbone* $f(.)$ is common for all the experiments, maps into an embedding of size 1280 and uses stochastic depth [29] for regularization. All supervised AudioSet models

Table 1: *Test mAP, mAUC and d-prime (± std over 3 runs) for supervised classification on AudioSet. LEAF\* denotes LEAF frontend with a fixed PCEN layer.*

| frontend | init | mAP | mAUC | d' |
|----------|------|-----|------|-----|
| mfbanks | N/A | **0.380±0.001** | **0.971±0.001** | **2.660±0.010** |
| SincNet | mel | 0.358±0.008 | 0.967±0.001 | 2.593±0.022 |
| SincNet | uniform | 0.351±0.007 | 0.965±0.000 | 2.567±0.006 |
| SincNet | truncn | 0.353±0.007 | 0.966±0.001 | 2.578±0.019 |
| LEAF* | mel | 0.373±0.002 | 0.969±0.001 | 2.639±0.019 |
| LEAF* | uniform | 0.358±0.003 | 0.967±0.001 | 2.595±0.019 |
| LEAF* | truncn | 0.358±0.008 | 0.967±0.001 | 2.593±0.022 |
| LEAF | mel | **0.380±0.002** | **0.970±0.001** | **2.653±0.019** |
| LEAF | uniform | 0.365±0.003 | 0.968±0.001 | 2.616±0.016 |
| LEAF | truncn | 0.368±0.002 | 0.968±0.000 | 2.622±0.011 |

Table 2: *Test accuracy (± std over 3 runs) of COLA pre-training + linear-probe results on SpeechCommandsV2. LEAF\* denotes a LEAF frontend with a fixed PCEN layer.*

| frontend | init | Accuracy |
|----------|------|----------|
| mfbanks | N/A | 61.2±2.6 |
| SincNet | mel | 61.6±2.8 |
| SincNet | uniform | 64.5±2.2 |
| SincNet | truncn | 64.2±3.6 |
| LEAF* | mel | 57.4±3.2 |
| LEAF* | uniform | 62.5±2.3 |
| LEAF* | truncn | 63.6±2.8 |
| LEAF | mel | 75.0±1.2 |
| LEAF | uniform | **76.1±3.3** |
| LEAF | truncn | 75.8±0.4 |



Figure 2: *Learned filter center frequencies for LEAF and Sinc-Net frontends, for COLA(top) and Supervised (bottom) settings. Mel center frequencies added for reference. default denotes default mel-scale based initialization. Error bands show standard deviations across three runs (best viewed in electronic format)*

are trained on 5-second random crops with a batch size of 1024 for a total of 50 epochs, with SpecAugment [30] and MixUp [31]. For COLA pre-training, randomly cropped 960 ms segments from an input clip are used as the anchor-positive pairs and passed through the backbone $f(.)$. We use a projection head $g(.)$ with $G = 512$ units, followed by layer normalization [32] and tanh activation. All COLA models are trained with a batch size of 2048 to provide a wider set of negative samples, for 50 epochs. For linear-probe evaluation on the SpeechCommandsV2 dataset, a linear classifier is trained directly on top of the COLA pre-trained neural backbone.

We train all models with an AdamW optimizer [33], with a linear warmup to a base learning rate of 5e-4 followed by a cosine decay [34] to 0 on a single TPU-v3 machine and repeated at least thrice.

### 3.2. Results

Table 1 reports results on the supervised multi-label classification task on Audioset for all the frontends. The LEAF frontend (row 2) matches the performance of log-compressed mel-filterbanks. For both the learnable frontends, random initialization schemes, although not as good as the default mel-scaled initialization, work remarkably well. It is worth noting that LEAF outperforms SincNet across the board, with even the random initialization schemes performing better than the default mel-scaled SincNet frontend. Overall, our observations on using strong priors (mel-scale) for initialization confirms the findings of previous work [8, 12, 11].
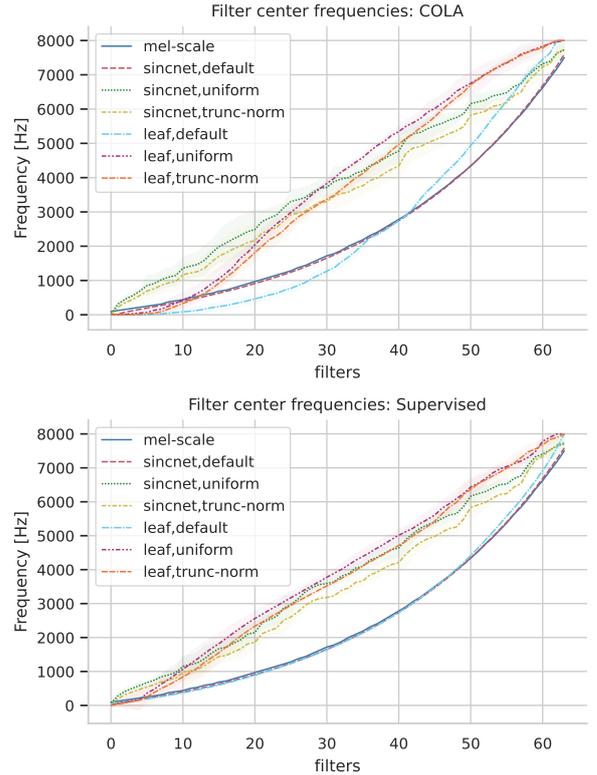
Table 2 reports linear-probe results for speech command classification on pre-trained COLA embeddings. We observe:

- *Learnable frontends significantly improve downstream accuracy over fixed features*: The learnable frontends provide significant improvements in linear-probe performance over mel-filterbanks, with the mel-scale initialized LEAF frontend offering a 14% absolute increment in classification accuracy.

- *Random initialization outperforms mel-scaled initialization*: Surprisingly, for both LEAF and SincNet, the proposed random initialization schemes outperform initializing filter frequencies on the mel-scale, with a more prominent absolute improvement of $\approx$ 3% observed for SincNet. Not only is this trend inversely correlated with the supervised results, but it is also counter-intuitive as initializing a filterbank on an auditory scale should provide a stronger starting point for the optimization process than initializing randomly.

- *Importance of a trainable PCEN layer*: It is expected that a fixed PCEN (all experiments marked LEAF*) will adversely affect performance, as observed for the supervised setting in Table 1. However, as evident from Table 2, the adverse effect on performance is even more significant in the self-supervised setting, with a drastic reduction in linear-probe accuracy when a fixed PCEN layer is used. To the best of our knowledge, this work is the first to show the benefits of learning PCEN compression and normalization in a self-supervised setting.
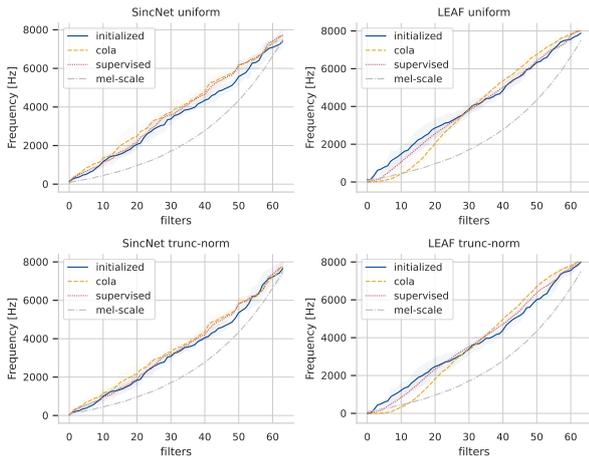
Figure 3: *Comparing random initialization schemes at initialization and convergence for the SincNet (left) and LEAF frontends (rights)*

### 3.3. Analysis of the learned frontend components

*i. Visualizing filter center frequencies*

We visualize center frequencies learned by the first convolution layer in Leaf and SincNet frontends for both COLA and supervised settings across all initialization schemes allowing us to analyze and compare frequencies modelled by the frontends under different settings (see Figure 2). When training for supervised classification, both frontends converge to center frequencies that follow a scale similar to their initialization (e.g., mel or uniform). On the other hand, self-supervised learning yields much more diverse filterbanks at convergence. It is also worth noting that for the LEAF frontend, even for mel-scale initialization, the center frequencies diverge away much further from the mel-scale for COLA in comparison to supervised training. The fact that contrastive loss draws learnable filters much further from their initialization than cross-entropy suggests that learnable frontends could benefit self-supervised learning even more than supervised classification.

*On random initialization:* Figure 3 shows randomly initialized LEAF and SincNet frontends, respectively, at initialization and at convergence. SincNet center frequencies at convergence overlap with their initialization throughout the frequency spectrum for both uniform and truncated normal initializations. LEAF filters are much more distinct in their center frequencies modelled across the COLA and supervised training regimes. In the supervised setting, randomly initialized LEAF filters converge closer to their initialization, in contrast to COLA pre-training where they diverge to model a broader set of center frequencies, especially in the 0-2000 Hz and the 3500-7000 Hz range, improving performance over mel-scale initialization.

These observations support empirical observations made in Section 3.2: across frontend configurations, there is a much larger variation in linear-probe performance as compared to supervised training. Moreover, while we could have expected learnable filters to converge to an auditory, logarithmic scale akin to the mel-scale, they instead converge to an almost linear scale which yet provides higher downstream accuracy.

*ii. Self-supervised learning of PCEN*

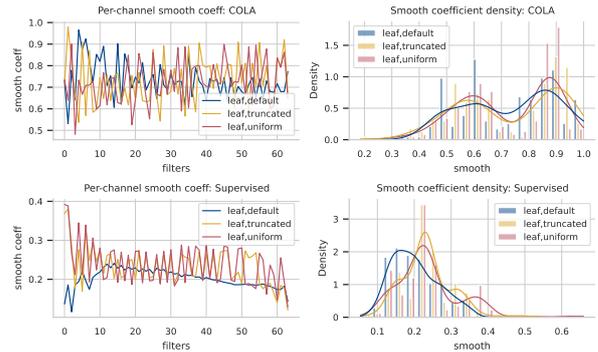While a PCEN layer learns several parameters including



Figure 4: *Learned per-channel smoothing coefficients (left) and overall smoothing coefficient density plot (right) of the learned PCEN layer for COLA (top) and Supervised (bottom). Notice the alternating pattern on the left, similar to [14]*

gain normalization and dynamic range compression offset and exponent (often denoted $\alpha$, $\delta$ and $r$, respectively), we focus on the smoothing coefficient $s$, which is the exponential moving average parameter that governs the extent of smoothing of the input time-frequency representation based on its past values. Figure 4 shows per-channel smoothing coefficients as well as their density histogram at convergence. Smoothing parameters tend to oscillate between channels, a pattern similar to that observed by [14], which they posit as an effort by the model to focus on different features to obtain more discriminative information. It's worth noting that this alternating behaviour is more pronounced for the COLA setting and for the randomly initialized frontends. The smoothing coefficient density plot in Figure 4 shows a cleaner picture: COLA learns coefficients that are significantly more spread out over the $[0., 1.]$ range, even for the default mel-scaled initialization, while supervised learning converges to values centered around $0.2$.

These observations, coupled with the fact that a fixed PCEN layer drastically reduces self-supervised performance (see Sec 3.2), suggest that PCEN possibly enables self-supervised learning to effectively capture discriminative information, and in the examined case, more so than supervised learning.

## 4. Conclusion

In this work, we study the feasibility of training learnable audio frontends in a self-supervised fashion. When pre-training a convolutional encoder jointly with SincNet or LEAF on Audioset, the downstream linear-probe performance on a speech command classification task improves significantly. We conduct an exploratory analysis of the learned filters and normalization coefficients, highlighting major differences in learning dynamics between the supervised and the self-supervised setting. Interestingly, our experiments suggest that while a mel-scale initialization of learnable filters expectedly improves the final performance in the supervised setting, self-supervised learning rather benefits from random initialization.

## 5. Acknowledgements

# 6. References

[1] S. S. Stevens and J. Volkmann, "The relation of pitch to frequency: A revised scale," *The American Journal of Psychology*, vol. 53, no. 3, pp. 329–353, 1940.

[2] J. Andén and S. Mallat, "Deep scattering spectrum," *IEEE Transactions on Signal Processing*, vol. 62, no. 16, pp. 4114–4128, 2014.

[3] D. O'shaughnessy, *Speech communications: Human and machine (IEEE)*. Universities press, 1987.

[4] S. Umesh, L. Cohen, and D. Nelson, "Fitting the mel scale," in *1999 IEEE International Conference on Acoustics, Speech, and Signal Processing. Proceedings. ICASSP99 (Cat. No. 99CH36258)*, vol. 1. IEEE, 1999, pp. 217–220.

[5] R. Schluter, I. Bezrukov, H. Wagner, and H. Ney, "Gammatone features and feature combination for large vocabulary speech recognition," in *2007 IEEE International Conference on Acoustics, Speech and Signal Processing-ICASSP'07*, vol. 4. IEEE, 2007, pp. IV–649.

[6] J. G. Lyons and K. K. Paliwal, "Effect of compressing the dynamic range of the power spectrum in modulation filtering based speech enhancement," in *Ninth Annual Conference of the International Speech Communication Association*, 2008.

[7] D. Palaz, M. M. Doss, and R. Collobert, "Convolutional neural networks-based continuous speech recognition using raw speech signal," in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2015, pp. 4295–4299.

[8] T. Sainath, R. J. Weiss, K. Wilson, A. W. Senior, and O. Vinyals, "Learning the speech front-end with raw waveform cldnns," 2015.

[9] R. Balestriero, R. Cosentino, H. Glotin, and R. Baraniuk, "Spline filters for end-to-end deep learning," in *International conference on machine learning*. PMLR, 2018, pp. 364–373.

[10] N. Zeghidour, N. Usunier, G. Synnaeve, R. Collobert, and E. Dupoux, "End-to-end speech recognition from the raw waveform," in *Interspeech 2018*, 2018.

[11] N. Zeghidour, N. Usunier, I. Kokkinos, T. Schaiz, G. Synnaeve, and E. Dupoux, "Learning filterbanks from raw speech for phone recognition," in *2018 IEEE international conference on acoustics, speech and signal Processing (ICASSP)*. IEEE, 2018, pp. 5509–5513.

[12] M. Ravanelli and Y. Bengio, "Speaker recognition from raw waveform with sincnet," in *2018 IEEE Spoken Language Technology Workshop (SLT)*. IEEE, 2018, pp. 1021–1028.

[13] P.-G. Noé, T. Parcollet, and M. Morchid, "Cgcnn: Complex gabor convolutional neural network on raw speech," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 7724–7728.

[14] Y. Wang, P. Getreuer, T. Hughes, R. F. Lyon, and R. A. Saurous, "Trainable frontend for robust and far-field keyword spotting," in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2017, pp. 5670–5674.

[15] V. Lostanlen, J. Salamon, M. Cartwright, B. McFee, A. Farnsworth, S. Kelling, and J. P. Bello, "Per-channel energy normalization: Why and how," *IEEE Signal Processing Letters*, vol. 26, no. 1, pp. 39–43, 2019.

[16] N. Zeghidour, O. Teboul, F. de Chaumont Quitry, and M. Tagliasacchi, "{LEAF}: A learnable frontend for audio classification," in *International Conference on Learning Representations*, 2021. [Online]. Available: https://openreview.net/forum?id=jM76BCb6F9m

[17] A. v. d. Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," *arXiv preprint arXiv:1807.03748*, 2018.

[18] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," in *Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 12 449–12 460.

[19] M. Tagliasacchi, B. Gfeller, F. de Chaumont Quitry, and D. Roblek, "Pre-training audio representations with self-supervision," *IEEE Signal Processing Letters*, vol. 27, pp. 600–604, 2020.

[20] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhotia, R. Salakhutdinov, and A. Mohamed, "Hubert: Self-supervised speech representation learning by masked prediction of hidden units." *IEEE Transactions on Audio, Speech, and Language Processing*, pp. 1–1, 2021.

[21] D. Niizumi, D. Takeuchi, Y. Ohishi, N. Harada, and K. Kashino, "Byol for audio: Self-supervised learning for general-purpose audio representation," in *2021 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2021, pp. 1–8.

[22] L. Wang, P. Luc, Y. Wu, A. Recasens, L. Smaira, A. Brock, A. Jaegle, J.-B. Alayrac, S. Dieleman, J. Carreira *et al.*, "Towards learning universal audio representations," *arXiv preprint arXiv:2111.12124*, 2021.

[23] A. Saeed, D. Grangier, and N. Zeghidour, "Contrastive learning of general-purpose audio representations," in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 3875–3879.

[24] E. Fonseca, D. Ortego, K. McGuinness, N. E. O'Connor, and X. Serra, "Unsupervised contrastive learning of sound event representations," in *2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 371–375.

[25] J. F. Gemmeke, D. P. W. Ellis, D. Freedman, A. Jansen, W. Lawrence, R. C. Moore, M. Plakal, and M. Ritter, "Audio set: An ontology and human-labeled dataset for audio events," in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2017, pp. 776–780.

[26] M. Tan and Q. Le, "Efficientnet: Rethinking model scaling for convolutional neural networks," in *International conference on machine learning*. PMLR, 2019, pp. 6105–6114.

[27] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *ICML 2020: 37th International Conference on Machine Learning*, vol. 1, 2020, pp. 1597–1607.

[28] P. Warden, "Speech commands: A dataset for limited-vocabulary speech recognition," *arXiv preprint arXiv:1804.03209*, 2018.

[29] G. Huang, Y. Sun, Z. Liu, D. Sedra, and K. Q. Weinberger, "Deep networks with stochastic depth," in *European conference on computer vision*. Springer, 2016, pp. 646–661.

[30] D. S. Park, W. Chan, Y. Zhang, C.-C. Chiu, B. Zoph, E. D. Cubuk, and Q. V. Le, "Specaugment: A simple data augmentation method for automatic speech recognition," in *Interspeech 2019*, 2019, pp. 2613–2617.

[31] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz, "mixup: Beyond empirical risk minimization," in *ICLR 2017 : International Conference on Learning Representations 2017*, 2017.

[32] J. L. Ba, J. R. Kiros, and G. E. Hinton, "Layer normalization," *arXiv preprint arXiv:1607.06450*, 2016.

[33] I. Loshchilov and F. Hutter, "Decoupled weight decay regularization," in *International Conference on Learning Representations*, 2019. [Online]. Available: https://openreview.net/forum?id=Bkg6RiCqY7

[34] I. Loshchilov and F. Hutter, "Sgdr: Stochastic gradient descent with warm restarts," in *ICLR 2016 : International Conference on Learning Representations 2016*, 2016.