

Improving Domain Generalization for Sound Classification with Sparse Frequency-Regularized Transformer

Honglin Mu, Wentian Xia, Wanxiang Che

Research Center for Social Computing and Information Retrieval

Harbin Institute of Technology, China

{hlmu, wtxia, car}@ir.hit.edu.cn

Abstract—Sound classification models’ performance suffers from generalizing on out-of-distribution (OOD) data. Numerous methods have been proposed to help the model generalize. However, most either introduce inference overheads or focus on long-lasting CNN-variants, while Transformers has been proven to outperform CNNs on numerous natural language processing and computer vision tasks. We propose FRITO, an effective regularization technique on Transformer’s self-attention, to improve the model’s generalization ability by limiting each sequence position’s attention receptive field along the frequency dimension on the spectrogram. Experiments show that our method helps Transformer models achieve SOTA generalization performance on TAU 2020 and Nsynth datasets while saving 20% inference time.

Index Terms—Sound classification, acoustic scene classification, transformer, attention, domain generalization

I. INTRODUCTION

Machine learning has made significant strides in the realm of sound classification. However, factors such as variable ambience and device-specific sound distortion [1] often impede models that excel on the training set from generalizing effectively upon deployment. This challenge has garnered substantial attention, as evidenced by the DCASE Challenges [2]–[5]. These competitions introduce a cross-device acoustic scene classification (ASC) task, where models are trained to discriminate the scene of the input fragment, e.g., metro or square, while recording devices differ between the training set and the test set. Performance on the unseen devices (i.e., domains) within the test set serves as an indicator of the model’s generalization capabilities.

We investigated domain generalization methods in DCASE 2019-2022, finding that most sound classification generalization methods, despite improve models’ generalization ability, either introduce overheads [6]–[9] or based on durable model structures such of Resnets [10]–[13]. We focus on the newer model structure optimized for domain generalization.

CNN-variants, such as VGG [14], Resnet [15], and FCNN [8], are de facto domination backbones in previous work. Koutini et al. [16] found that CNNs suffer from overfitting on biased data in the ASC task, which could be alleviated by restricting CNNs’ receptive field. Similarly, McDonnell et al. [11] bolstered CNN’s performance on OOD devices by limiting the model’s perception of frequency. The

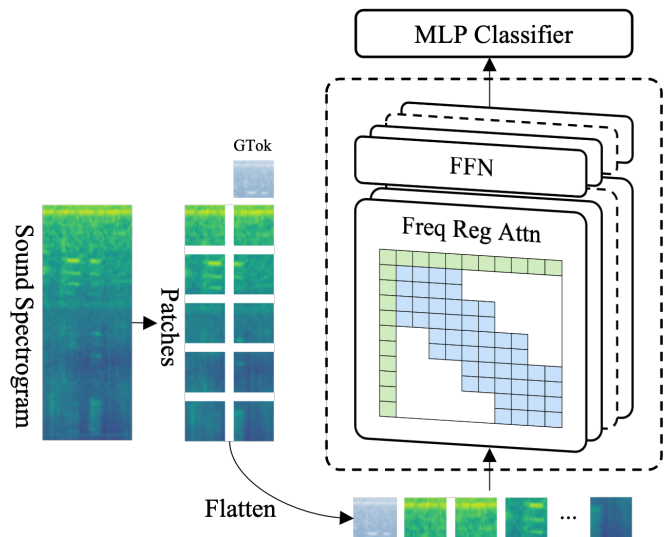


Fig. 1: The Frequency Regularized Transformer. *GTok* is abbreviation for global token.

proposed CNN variant includes two separate inference paths, each probing half range on the spectrogram’s frequency and fusing with the other before the classification head. With Transformer [17] applied to the text [18]–[20], image [21] and speech [22], [23], its power in the sound classification has also been discovered [24], [25]; the Patchout faSt Spectrogram Transformer (PaSST) [24] surpassed CNN variants in performance and efficiency through dropping out patches from the spectrogram. While Transformer-based models challenge sound classification tasks, their generalization degradation issue remains to be investigated.

This work focuses on improving the Transformer’s domain generalization ability in sound classification tasks. Inspired by Koutini et al. and McDonnell et al., we propose the Frequency-Regularized Transformer which fuses frequency regularization into self-attention to provide both robustness and efficiency. We limit the Transformer’s receptive field along the frequency on the sound’s spectrogram by applying two types of masks to the Transformer’s self-attention, one of which perceives neighbor frequencies while the other considers the entire

range. Such a strategy alleviate the over-fitting problem that the model suffers facing OOD data. Our contributions can be summarized as follows:

- We propose the **F**requency-**R**egularIzed **T**ransfOrmer (FRITO), a Transformer variant enhanced on the generalization ability for sound classification, which achieves SOTA on unseen domains of TAU 2020 [26] and Nsynth [27] datasets.
- We implement the sparse form of this method’s common use case, saving 20% of inference time and 21% memory compared to its full-attention version.

Our code for this paper will be publicly available at <https://github.com/hlmu/FRITO>.

II. FREQUENCY-REGULARIZED TRANSFORMER (FRITO)

A. Backbone

Our proposed method utilizes the Vision Transformer (ViT) [21] as its backbone. ViT is a state-of-the-art deep learning model that has adapted the Transformer architecture for computer vision tasks. It achieves this by dividing the image into a fixed grid of patches, which are then added with positional encoding and fed into a standard Transformer block. This enables ViT to capture long-range dependencies between image patches and achieve impressive results on various vision tasks.

Building on this foundation, the PaSST model [24] leverages ViT in the audio domain by splitting the audio spectrogram into patches and employing disentangled time and frequency positional encoding. In our approach, we take inspiration from both ViT and PaSST and modify the attention block to restrict its receptive field on the frequency dimension. This approach enables our model to mitigate overfitting issues and generalize better to out-of-domain data.

Our model architecture is illustrated in Figure 1, and it is composed of modified ViT blocks that process the audio spectrogram patches.

B. Frequency Attention Regularization

In this section, we introduce our model’s regularization scheme, which restricts the receptive field of the Transformer’s each input position by adding masks to its self-attention. Our scheme includes local and global attention, similar to the efficient transformers [28], [29]. This method differs from previous work in that it restricts the receptive field along the frequency dimension, aiming to improve the model’s generalization ability instead of dealing with long sequences.

Given a piece of audio’s Mel-spectrogram, we split it into patches and add positional encoding following PaSST. Let such patch matrix be $P \in \mathbb{R}^{h \times w \times d}$ where h, w , and d represent rows of patches, columns of patches, and Transformer’s hidden size, respectively. We then flatten the patch matrix into $x_{P(1,1)}, \dots, x_{P(h,w)} \in \mathbb{R}^{hw \times d}$, and prefix k global tokens x_{g_1}, \dots, x_{g_k} which will be explained shortly. Then the Transformer’s input sequence can be represented as:

$$X = (x_{g_1}^1, \dots, x_{g_k}^k, x_{P(1,1)}^{k+1}, \dots, x_{P(h,w)}^t) \in \mathbb{R}^{t \times d}$$

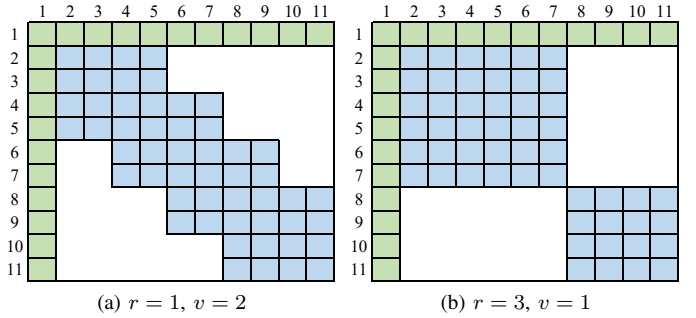


Fig. 2: Example attention masks M for the input sequence, made from spectrogram patches with 5 rows in frequency and 2 columns in time, one global token prefixed. Blank and colored blocks represent $-\infty$ and 0, respectively. Green blocks represent global attention, while the blue ones represent local attention. (a) Under the conditions of $r = 1$ and $v = 2$, directly adjacent rows are interiorly accessible. (b) When $r = 3$ and $v = 1$, the patches in the first three rows are visible to each other, and the patches in the last two rows are mutually sensible.

where t is the Transformer’s input sequence length.

We constrain Transformer’s receptive field by adding a mask matrix M on its self-attention weights. Let $Q, K, V \in \mathbb{R}^{t \times d}$ be the Transformer’s queries, keys, and values, respectively. Then the attention of our Frequency-Regularized Transformer can be written as:

$$ATTN_{FRITO}(X) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}} + M\right)V \quad (1)$$

where $M \in \{0, -\infty\}^{n \times n}$ is the attention mask matrix shown as 2. Sequence token X_*^j is visible to X_*^i iff $M(i, j)$ equals zero. The following sections provide a detailed design of the mask matrix M .

C. Local Frequency Attention

This subsection first presents how we select the receptive field of each token, and then introduces the M matrix representation for this scheme. Specifically, we restrict the receptive field S of a patch $P_{(a,b)}$ on the patch matrix P to its neighbors close in frequency. The receptive field S can be expressed using the following formula, where a and b represent arbitrary row and column numbers:

$$S(a, b) = \left\{ (p, q) \mid p \in \left[\left\lfloor \frac{a-1}{r} \right\rfloor r + 1, \left\lfloor \frac{a+r-1}{r} \right\rfloor rv \right], \right. \\ \left. q \in [1, w] \right\}$$

where r is the size of a row cluster with r rows of internally visible patches; v is the overlap factor, indicating the number of row clusters that $P_{(a,b)}$ can observe; the optimal value of r and v are chosen from experiments.

These receptive field rules are then applied to the attention mask M for the flattened input sequence X , represented as:

$$\begin{aligned} M[iw + k + 1 : (i + rv)w + k + 1, \\ iw + k + 1 : (i + rv)w + k + 1] &= 0, \\ i &\in \{0, r, 2r, \dots, \lfloor \frac{h-1}{r} \rfloor r\} \end{aligned}$$

where k is the number of global tokens, which will be explained shortly. A visualized example is shown in Figure 2.

D. Global Frequency Attention

Despite local frequency insights, the model needs overall perception to perform the classification task. Similar to BIGBIRD-ETC’s approach [29], we add global tokens $x_{g_1}, x_{g_2}, \dots, x_{g_k}$ at the beginning of the flattened patch tokens to perceive the entire sequence. Global tokens share a mutual sight with all other tokens, and its attention mask M can be expressed as:

$$M[i, :] = 0, M[:, i] = 0, i \in \{1, 2, \dots, k\}$$

Values on other positions in M not defined by local and global frequency attention are all $-\infty$. Figure 2 visualize the M matrix with an example.

E. Sparse Frequency-Regularized Attention

The above approach restricts the model’s receptive field by adding a mask to the attention weight matrix. However, its naïve implementation does not improve the speed of reasoning. While the M matrix has a relatively large number of $-\infty$ value, its sparsity can be utilized. This paper adopts a sparse attention operation for the case where $v = 1$ to improve the model training and reasoning speed.

When $v = 1$, the input sequence X excluding the global tokens can be segmented into $l = \lfloor \frac{h}{r} \rfloor$ mutually imperceptible blocks denoted as X_1, X_2, \dots, X_l . The Q, K, V in Equation 1 can be split along the time dimension into $Q_i, V_i, K_i \in \mathbb{R}^{r \times d}, i \in 1, 2, \dots, l$, respectively, corresponding to the local attention of the i -th block. In this case, the attention formula in Equation 1 can be rewritten as:

$$ATTN_{SLocal}(X) = \left\| \left\| \text{softmax} \left(\frac{Q_i K_i^T}{\sqrt{d}} \right) V_i \right\| \right\|$$

where $\|$ represents concatenation along the time dimension.

The above formula only considers local attention, but FRITO also requires global attention. We aggregate the global tokens into a block represented as X_g , its query and key matrix denoted as Q_g and K_g , respectively. The attention formula can then be rewritten as the final sparse version:

$$\begin{aligned} ATTN_{SFRITO}(X) = & \text{softmax} \left(\frac{Q_g K_g^T}{\sqrt{d}} \right) V_{\oplus} \\ & \left\| \left\| \text{softmax} \left(\frac{Q_i (K_g \parallel K_i)^T}{\sqrt{d}} \right) V_i \right\| \right\| \end{aligned}$$

where both \oplus and $\|$ are concatenation operation.

Sparse FRITO is computationally equivalent to its full-attention version. In particular, since Q, K , and V are only split along the time dimension, the weights obtained from the full-attention model can be directly used for sparse attention inference.

III. EXPERIMENT SETUP

We evaluate our method on two publicly available datasets, TAU 2020 [26] and Nsynth [27].

A. Nsynth

The Nsynth dataset is a collection of three different sound sources, *acoustic*, *electronic*, and *synthetic*, consisting of sounds from various instruments played at 16kHz, including *bass*, *brass*, *flute*, *guitar*, *keyboard*, *mallet*, *organ*, *reed*, *string*, *synth_lead*, and *vocal*. In addition, each instrument except *synth_lead* includes segments from *acoustic* sound source. We require models to infer on unseen *electronic* and *synthetic* source after training them to classify instruments on *acoustic* source. Original data split mixes *synthetic*, *acoustic*, and *electronic* sources in train, development, and test set. While we aim to measure models’ generalization ability on unseen domains, i.e. sound sources, we re-divide the train, development, and test set in the following way:

We retain all sounds from the *acoustic* source, incorporating them into the training set, while utilizing *synthetic* and *electronic* sounds, with markedly different sound characteristics, as out-of-domain samples for the development set. The train-development proportion adheres to the original train-test ratio established in the original split. Detailed dataset split information can be found in Table I.

TABLE I: Statistics of the re-split Nsynth dataset

Instrument	Train Set	Dev Set		
	Acoustic	Synthetic	Acoustic	Electronic
Bass	180	792	20	110
Brass	13740	0	20	70
Flute	6552	104	20	35
Guitar	13323	108	20	345
Keyboard	8488	59	20	657
Mallet	27702	107	20	341
Organ	156	0	20	469
Reed	14242	156	20	76
String	20490	0	20	84
Vocal	3905	121	20	140

B. TAU 2020

The TAU 2020 dataset is directly used by DCASE 2020 Task 1A [3] and DCASE 2021 Task 1A [4] to evaluate the generalization properties of systems across a number of different devices. It is an environmental sound classification dataset that requires the model to classify the sound environment, e.g., *metro*, *shopping mall*, and *public square*, of the input audio. The dataset is recorded by multiple recording devices. The development set contains devices with unique acoustic characteristics not presented in the training set, so the performance on it can reflect the model’s

generalization ability. We use this dataset following DCASE 2020’s official instructions [30].

C. Model Initialization

We use the pre-trained weight `passt_s_swa_p16_128_ap476` from the PaSST model, which has 8 attention heads and a 768-dimensional encoder, derived from DeiT-B \uparrow 384 [31], finetuned on Audioset [32]. For the PaSST experiment, we set $s_patchout_t=10$, $s_patchout_f=5$ as in their work. We finetune models on 4 Tesla-A100-80GBs, on a batch size containing 800 seconds of sound, using the Adam optimizer with a maximum learning rate of $1e-3$.

IV. RESULTS

A. Nsynth

Our FRITO method demonstrates optimal overall performance on the Nsynth dataset, as shown in Table II. The dataset comprises three domains: *acoustic*, *electronic*, and *synthetic*. We train models on the *acoustic* domain, and evaluate their generalization performance on the remaining two domains. Consistent with previous work’s experimental conclusions [24], Transformer-based methods outperform the Resnet overall. Among these Transformer-based methods, PaSST exhibits a degradation in the *Synthetic* domain compared to ViT, while our method remains less affected.

TABLE II: Accuracy on the Nsynth dataset. *Electronic* and *Synthetic* do not appear in the training set.

Method	Overall	Electronic	Synthetic
Resnet [8]	0.276	0.251	0.317
ViT [24]	0.301	0.351	0.221
PaSST [24]	0.285	0.355	0.173
FRITO- <i>r6-v1</i>	0.307	0.379	0.192

B. TAU 2020

Performance on the TAU 2020 dataset is presented as Table III. Our regularized method achieved the optimal performance on the unseen domains $S4 - S6$, with 3.5% improvement compared to Resnet and 0.5% improvement in generalization performance compared to PaSST. Transformer-based models in our experiments have significantly better generalization performance than Resnet, which is in accord with results on the Nsynth dataset.

C. Sparse Attention Efficiency

Our Sparse FRITO-*r6-v1* demonstrates a 20% improvement in inference speed and a 21% reduction in memory usage compared to its full-attention version, as outlined in Table IV. We employ the same parameters as FRITO-*r6-v1* on the Nsynth dataset for Sparse FRITO. However, it should be noted that the full-attention FRITO-*r6-v1* experiences a 6% degradation in inference speed compared to the vanilla ViT, which can be attributed to the overhead of computing the attention mask M .

TABLE III: Accuracy on TAU 2020 development set. A, B, C are real recording devices, while $S1-S6$ are virtual devices generated from A, B, C . $S4-S6$ only appears in the development set. The results marked with a \star symbol correspond to our own run.

Method	Overall	S4-S6	S1-S3	A & B & C
Resnet [8]	0.746	0.710	0.736	0.784
ViT [24]	0.763	-	-	-
ViT \star	0.763	0.739	0.770	0.779
PaSST [24]	0.756	-	-	-
PaSST \star	0.756	0.740	0.737	0.791
FRITO- <i>r1-v8</i> \star	0.761	0.745	0.735	0.804

TABLE IV: Comparison of inference speed and memory

Method	Speed	Mem
ViT [24]	0%	0%
FRITO- <i>r6-v1</i>	-6%	+0%
Sparse FRITO- <i>r6-v1</i>	+14%	-21%

V. CONCLUSION

Our research is centered on enhancing the domain generalization capabilities of Transformer models in the sound classification task. By constraining the receptive field of the Transformer’s self-attention, our model attains superior generalization performance on two publicly available datasets. Furthermore, we investigate sparse attention operations to optimize the model’s inference speed and memory consumption.

Despite the empirical success, a theoretical foundation for the principle of limiting the receptive field remains to be established. Future work could endeavor to derive the relationship between the receptive field and generalization performance, providing a more rigorous understanding of the underlying mechanisms at play.

REFERENCES

- [1] (2019) Dcase2019 challenge introduction. [Online]. Available: <https://dcase.community/challenge2019/index>
- [2] M. Mandel, J. Salamon, and D. P. W. Ellis, *Proceedings of the Detection and Classification of Acoustic Scenes and Events 2019 Workshop (DCASE2019)*. NY, USA: New York University, October 2019.
- [3] N. Ono, N. Harada, Y. Kawaguchi, A. Mesaros, K. Imoto, Y. Koizumi, and T. Komatsu, *Proceedings of the Fifth Workshop on Detection and Classification of Acoustic Scenes and Events (DCASE 2020)*, Tokyo, Japan, November 2020.
- [4] F. Font, A. Mesaros, D. P. Ellis, E. Fonseca, M. Fuentes, and B. Elizalde, *Proceedings of the 6th Workshop on Detection and Classification of Acoustic Scenes and Events (DCASE 2021)*, Barcelona, Spain, November 2021.
- [5] M. Lagrange, A. Mesaros, T. Pellegrini, R. S. Gaël Richard, and D. Stowell, *Proceedings of the 7th Workshop on Detection and Classification of Acoustic Scenes and Events (DCASE 2022)*. Nancy, France: Tampere University, November 2022.
- [6] J. Mukhoti, V. Kulharia, A. Sanyal, S. Golodetz, P. Torr, and P. Dokania, “Calibrating deep neural networks using focal loss,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 15 288–15 299, 2020.
- [7] S. Suh, S. Park, Y. Jeong, and T. Lee, “Designing acoustic scene classification models with cnn variants,” *Tech. Rep., DCASE2020 Challenge*, 2020.

- [8] H. Hu, C.-H. H. Yang, X. Xia, X. Bai, X. Tang, Y. Wang, S. Niu, L. Chai, J. Li, H. Zhu *et al.*, “A two-stage approach to device-robust acoustic scene classification,” in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 845–849.
- [9] H.-j. Shim, J.-w. Jung, J.-h. Kim, and H.-J. Yu, “Attentive max feature map and joint training for acoustic scene classification,” in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 1036–1040.
- [10] H. Eghbal-zadeh, K. Koutini, and G. Widmer, “Acoustic scene classification and audio tagging with receptive-field-regularized cnns,” *Tech. Rep., DCASE 2019 Challenge*, 2019.
- [11] M. D. McDonnell and W. Gao, “Acoustic scene classification using deep residual networks with late fusion of separated high and low frequency paths,” in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 141–145.
- [12] B. Kim, S. Yang, J. Kim, and S. Chang, “Qti submission to dcase 2021: Residual normalization for device-imbalanced acoustic scene classification with efficient design,” *arXiv preprint arXiv:2206.13909*, 2022.
- [13] J.-H. Lee, J.-H. Choi, P. M. Byun, and J.-H. Chang, “Hyu submission for the dcase 2022: fine-tuning method using device-aware data-random-drop for device-imbalanced acoustic scene classification,” 2022.
- [14] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *arXiv preprint arXiv:1409.1556*, 2014.
- [15] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [16] K. Koutini, H. Eghbal-zadeh, and G. Widmer, “Receptive-field-regularized cnn variants for acoustic scene classification,” *arXiv preprint arXiv:1909.02859*, 2019.
- [17] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, E. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.
- [18] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” *arXiv preprint arXiv:1810.04805*, 2018.
- [19] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, “Roberta: A robustly optimized bert pretraining approach,” *arXiv preprint arXiv:1907.11692*, 2019.
- [20] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, “Language models are few-shot learners,” *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.
- [21] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, “An image is worth 16x16 words: Transformers for image recognition at scale,” *arXiv preprint arXiv:2010.11929*, 2020.
- [22] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, “wav2vec 2.0: A framework for self-supervised learning of speech representations,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 12 449–12 460, 2020.
- [23] S. Chen, C. Wang, Z. Chen, Y. Wu, S. Liu, Z. Chen, J. Li, N. Kanda, T. Yoshioka, X. Xiao *et al.*, “Wavlm: Large-scale self-supervised pre-training for full stack speech processing,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 16, no. 6, pp. 1505–1518, 2022.
- [24] K. Koutini, J. Schlüter, H. Eghbal-zadeh, and G. Widmer, “Efficient training of audio transformers with patchout,” *arXiv preprint arXiv:2110.05069*, 2021.
- [25] Y. Gong, Y.-A. Chung, and J. Glass, “Ast: Audio spectrogram transformer,” *arXiv preprint arXiv:2104.01778*, 2021.
- [26] A. Mesaros, T. Heittola, and T. Virtanen, “A multi-device dataset for urban acoustic scene classification,” in *Proceedings of the Detection and Classification of Acoustic Scenes and Events 2018 Workshop (DCASE2018)*, November 2018, pp. 9–13.
- [27] J. Engel, C. Resnick, A. Roberts, S. Dieleman, M. Norouzi, D. Eck, and K. Simonyan, “Neural audio synthesis of musical notes with wavenet autoencoders,” in *International Conference on Machine Learning*. PMLR, 2017, pp. 1068–1077.
- [28] R. Child, S. Gray, A. Radford, and I. Sutskever, “Generating long sequences with sparse transformers,” *arXiv preprint arXiv:1904.10509*, 2019.
- [29] M. Zaheer, G. Guruganesh, K. A. Dubey, J. Ainslie, C. Alberti, S. Ontanon, P. Pham, A. Ravula, Q. Wang, L. Yang *et al.*, “Big bird: Transformers for longer sequences,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 17 283–17 297, 2020.
- [30] (2020) Acoustic scene classification task description. [Online]. Available: <https://dcase.community/challenge2020/task-acoustic-scene-classification#subtask-a>
- [31] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jégou, “Training data-efficient image transformers & distillation through attention,” in *International Conference on Machine Learning*. PMLR, 2021, pp. 10 347–10 357.
- [32] J. F. Gemmeke, D. P. 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)*. IEEE, 2017, pp. 776–780.