

A Study on Joint Modeling and Data Augmentation of Multi-Modalities for Audio-Visual Scene Classification

Qing Wang¹, Jun Du¹, Siyuan Zheng¹, Yunqing Li¹, Yajian Wang¹, Yuzhong Wu^{2,4}, Hu Hu³,
Chao-Han Huck Yang³, Sabato Marco Siniscalchi^{3,5}, Yannan Wang², Chin-Hui Lee³

¹NELSLIP, University of Science and Technology of China, Hefei, China,

²Tencent Ethereal Audio Lab, Tencent Corporation, China,

³School of Electrical and Computer Engineering, Georgia Institute of Technology, GA, USA,

⁴DSP & Speech Technology Laboratory, The Chinese University of Hong Kong, Hong Kong,

⁵Computer Engineering School, University of Enna Kore, Italy

jundu@ustc.edu.cn

Abstract

In this paper, we propose two techniques, namely joint modeling and data augmentation, to improve system performances for audio-visual scene classification (AVSC). We employ pre-trained networks trained only on image data sets to extract video embedding; whereas for audio embedding models, we decide to train them from scratch. We explore different neural network architectures for joint modeling to effectively combine the video and audio modalities. Moreover, data augmentation strategies are investigated to increase audio-visual training set size. For the video modality the effectiveness of several operations in RandAugment is verified. An audio-video joint mixup scheme is proposed to further improve AVSC performances. Evaluated on the development set of TAU Urban Audio Visual Scenes 2021, our final system can achieve the best accuracy of 94.2% among all single AVSC systems submitted to DCASE 2021 Task 1b.

Index Terms: audio-visual scene classification, acoustic-visual joint modeling, joint data augmentation

1. Introduction

Acoustic scene classification (ASC) task aims to identify the environment classes of audio recordings, which can be used for various demands of contextualization and personalization. The Detection and Classification of Acoustic Scenes and Events (DCASE) Challenge has organized ASC related tasks for years [7, 8, 9] under different scenarios. In this year, DCASE 2021 Challenge proposes a new audio-visual scene classification (AVSC) task [3] by leveraging on additional information of video modality, which makes a big difference compared to previous tasks. Since multi-modal methods can greatly boost the performance compared to single modality [3], AVSC should be more promising in challenging realistic applications.

For an ASC task, the most competitive methods [10, 6] in the previous DCASE Challenges extracted representative audio embedding by inputting log-Mel spectrogram or Mel-frequency cepstral coefficients (MFCC) features into deep neural networks for classification. Data augmentation strategies such as mixup [4] and SpecAugment [5] were also introduced to enhance generalization ability of models. For visual scene classification (VSC), many convolutional neural networks (CNNs) based methods [12, 13, 14] transferred from object recognition task were adopted. Large in-domain datasets [18, 20] were also utilized for pre-training, which could enrich the extracted feature.

Additionally, there exist other group of networks [15, 16] which were specifically developed for scene classification. For the new AVSC task [3], audio recordings and corresponding video clips are both provided, which means new model architectures and new training strategies are needed. Multi-modal fusion approaches can be summarized into three categories: early fusion [34], late fusion [35] and hybrid fusion [36]. Since early fusion strategy can learn the correlation of different modalities at the feature level, it has become a very popular way to tackle multi-modal fusion.

In this paper, we propose a simple yet effective multi-modal approach for AVSC task, which mainly consists of audio module, visual module and modality fusion module. Under the premise of multi-modal input, we explore the selection of unimodal representation to improve system performance. Our observation is that the combination of audio feature extraction with our customized fully convolutional neural networks (FCNN) and video feature extraction with DenseNet [14] can achieve the best results. The paradigm of joint fine-tuning after pre-training is also introduced in AVSC task, which can greatly enhance the performance in Task 1b of DCASE 2021 Challenge. Moreover, we design multi-modal data augmentation strategies that fully enrich the input diversity of two modalities. In particular, a joint mixup strategy is proposed to synchronously generate new audio and video data, thus adding modality correlation at the input level. Evaluated on DCASE 2021 Task 1b, experimental results show that our system achieves the best accuracy of 94.2% on the development set among all single systems.

2. The proposed AVSC approach

2.1. Acoustic-Visual Joint Modeling

For the AVSC task, we design a multi-modal system as shown in Fig.1. Our AVSC model mainly consists of three parts: audio module, visual module and modality fusion module. We will elaborate them in the following subsections.

2.1.1. Audio Embedding

In the audio module, we extract the log-Mel filter bank (LMFB) features of the raw data I_A with delta and delta-delta operations, forming the input I_A^{LMFB} . For high-level feature representations, we employ audio extractor f_A on I_A^{LMFB} , and obtain the audio

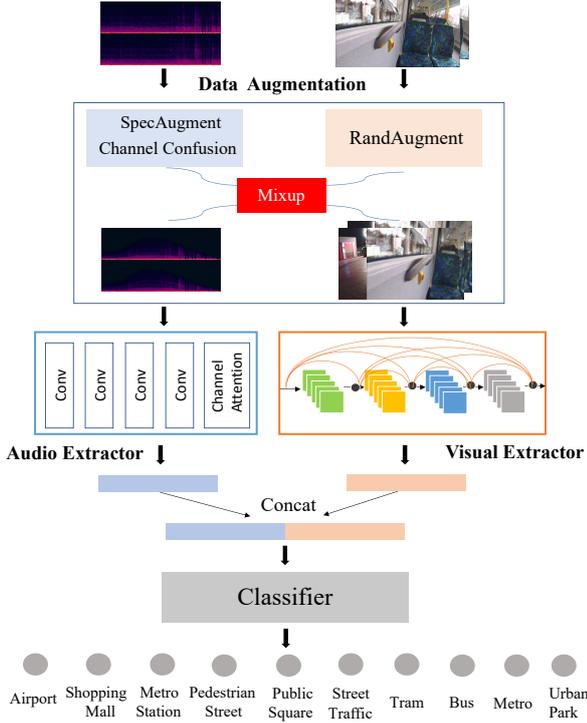


Figure 1: *The proposed AVSC system.*

embedding E_A as:

$$E_A = f_A(I_A^{\text{LMFB}}) \quad (1)$$

FCNN model is applied to extract high-level audio embedding, which achieved promising performance for the ACS task in our previous work [6]. The FCNN model mainly consists of four convolutional blocks with total 9 stacked convolutional layers and a channel attention block. Each convolutional layer is followed by a batch normalization and ReLU activation function. Dropout is used to alleviate over-fitting. Channel attention [32] is applied before the final global average pooling layer. In addition to the FCNN model which is trained from scratch, we also investigate an audio feature extractor VGGish [33] pre-trained on AudioSet [22].

2.1.2. Video Embedding

In the visual module, high-level video embedding E_V is calculated given the input image I_V by using visual extractor f_V as follows,

$$E_V = f_V(I_V). \quad (2)$$

We explore the effect of different networks (discussed in Section 3.2), and select DenseNet [14] as our video embedding extractor. DenseNet proposed dense connection between layers in a feed-forward fashion, which encouraged the maximum information flow in network. The transfer learning results on many downstream tasks, such as object detection and instance segmentation have proved its effectiveness. To enhance the generalization ability of the model, we adopt models pre-trained on in-domain large databases and apply them in this AVSC task by transfer learning.

2.1.3. Modality Fusion

With the high-level feature embedding of audio and video, modality fusion module is needed to integrate information from two sources. To fully exploit the complementarity of two modalities, we concatenate audio embedding E_A and video embedding E_V into fusion embedding E_F , and feed it into a classifier, which is a four-layer multi-layer perceptron (MLP) network with size of 512, 128, 64 and 10, respectively.

$$E_F = [E_A, E_V] \quad (3)$$

$$p = \text{MLP}(E_F) \quad (4)$$

The output probability vector $p = [p_1, p_2, \dots, p_{10}]$ is a 10-dimensional vector, corresponding to the number of classes, i.e. airport, shopping mall, metro station, pedestrian street, public square, street traffic, tram, bus, metro and urban park.

Suppose that the scene label is one-hot vector $y = [y_1, y_2, \dots, y_{10}] \in \mathbb{R}^{10}$, the cross entropy loss for classification is calculated as follows:

$$L = - \sum_{i=1}^{10} y_i \log p_i \quad (5)$$

It is worth noting that we train the whole AVSC model without freezing the audio or visual extractor, which means the parameters of both extractors are updated together with the classifier parameters during model training. Unlike feature-based approaches [23, 24], our fine-tuning strategy makes two feature extractors more task-specific on the challenge data set. Moreover, fine-tuning the two modalities together can achieve better modality fusion performances.

2.2. Audio-Video Data Augmentation

In our previous work [6], data augmentation methods are proven to be effective in ASC task of DCASE 2020 Challenge. Thus in this task, we adopt several different data augmentation methods for audio and video modalities. In addition, a joint mixup strategy is proposed, which is effective to improve the generalization abilities of models.

2.2.1. Audio Data Augmentation

When training the audio extractor, four techniques that generate extra data are used as listed below: (i) pitch shifting, where we perform it on audio clips to randomly shift the pitch based on the uniform distribution; (ii) speed changing, where we perform it on audio waveforms to randomly change the speed of audio recordings; (iii) noise adding, where we perform it on audio waveforms to add random Gaussian noise; and (iv) audio mixing, where we mix two audio recordings from the same scene class to generate a new sample with the same label. The other two on-the-fly data perturbation techniques includes: (i) SpecAugment [5]: in this study, we do not perform time warping. For time and frequency masking, we set the masking parameter to 10 % of the dimensions. It is applied on LMFB features and is performed on batch level; (ii) channel confusion: two channels of input audio feature are randomly swapped. The two on-the-fly audio data perturbation strategies are also applied when fine-tuning audio extractor.

2.2.2. Video Data Perturbation

In order to improve the robustness of the visual extractor, we investigate data augmentation techniques in RandAugment [30],

which was first used for image object classification and object detection tasks. RandAugment is a search algorithm to find the best data augmentation strategy. There are two optimal parameters called N and M . For each image in each mini-batch, select N data augmentation methods with M magnitudes. Due to the difference between the object image and scene image, the accuracy drops after directly applying RandAugment to the AVSC task. Therefore, we evaluate each sub-policy for its usefulness in scene classification. Through detailed ablation studies, three sub-policies in the search space, namely Sharpness, Contrast and Identity Mapping are adopted. For each image, two operations are randomly selected to be applied in sequence to image data. Since these video augmentation techniques do not produce any extra offline data, we also call it data perturbation to make a distinction from commonly used data augmentation strategies. It is worth noting that all the video data perturbation techniques mentioned above are adopted for fine-tuning visual extractor.

2.2.3. Audio-Video Joint Mixup

Mixup was firstly proposed for improving the robustness of deep neural networks, which has been successfully applied in various tasks, such as unimodal image classification and adversarial examples generating. However, in multi-modal tasks, mixup has not been investigated yet as far as we know. With the change of input, regular mixup strategy cannot be directly applied in our AVSC task, let alone the verification of its effectiveness.

Here in this section, we propose a novel joint mixup algorithm, successfully applying the mixup to the audio and visual inputs at the same time. Our joint mixup method mixes the input of two modalities synchronously, further increasing the diversity of input for better enhancing our joint AVSC model. In more detail, a joint audio-video example can be constructed by using the following formula:

$$(x_{ij}^{\text{ma}}, x_{ij}^{\text{mv}}) = \alpha \times (x_i^{\text{a}}, x_i^{\text{v}}) + (1 - \alpha) \times (x_j^{\text{a}}, x_j^{\text{v}}) \quad (6)$$

$$t_{ij}^{\text{m}} = \alpha \times t_i + (1 - \alpha) \times t_j \quad (7)$$

where $(x_i^{\text{a}}, x_i^{\text{v}})$ are audio and video embedding of the i -th sample and t_i is the corresponding label. $x^{\text{ma}}, x^{\text{mv}}$ and t^{m} denote the audio, video embedding and scene label after the joint mixup, respectively. α is the mixed ratio and set to 0.4. To build the multi-modal systems, 20% of the training data are employed with audio-video joint mixup.

3. Experimental Results and Analysis

3.1. Experimental Setup

The data set used for the AVSC task is TAU Urban Audio-Visual Scenes 2021 [3], which consists of 34 hours of data with time-synchronized audio and video content. There are about 8k 10-second audio clips for training and 3k test audio clips recorded in a binaural way using a 48kHz sampling rate. For the video clips, every second contains 30 frames. We split all the data into 1-second samples without overlap to match the development set of Task 1b in DCASE 2021 Challenge.

For each input audio clip, we use the Librosa [31] library to extract the LMFB features and compute log-Mel delta and delta-delta operations without padding, which generates a feature tensor shape of $39 \times 128 \times 6$. For each input visual clip, the first frame and the fifteenth frame images are extracted and resized to 224×224 patches to calculate the video embedding. Then two video embedding are added together to serve as the

Table 1: An accuracy (‘Acc.’ in %) comparison of acoustic-visual joint modeling. The first two columns correspond to audio and video models, where ‘VGGish’ is pre-trained on AudioSet while ‘FCNN’ is trained from scratch. ‘Pre_V’ denotes pre-trained video models. ‘Fine_A’ and ‘Fine_V’ denote fine-tuning the audio and video extractors, respectively.

Audio	Video	Pre_V	Fine_A	Fine_V	Acc.
VGGish [33]	-	-	✓	-	59.3
FCNN	-	-	-	-	75.2
-	VGG	-	-	-	76.2
-	VGG	✓	-	-	80.3
FCNN	VGG	✓	-	-	87.4
FCNN	ResNet	✓	-	-	91.6
FCNN	ResNeSt	✓	-	-	92.2
FCNN	DenseNet	✓	-	-	92.5
FCNN	DenseNet	✓	✓	-	93.1
FCNN	DenseNet	✓	-	✓	92.9
FCNN	DenseNet	✓	✓	✓	93.2

final visual feature of the input video data. All our models are trained using PyTorch toolkit. And Adam optimizer is used during training. For our audio-visual joint model, we set a small learning rate of 1e-5 to fine-tune both audio and video extractors, with a weight decay of 1e-5 and batch size of 32.

3.2. Results on Acoustic-Visual Joint Modeling

We consider three aspects for acoustic-visual joint modeling: the selection of audio extractor, the selection of visual extractor, and the joint training strategy. Here in this section, we conduct a series of experiments to show the effectiveness of our joint modeling. Table 1 shows the experimental results on the development set.

For audio embedding, we compare VGGish [33] which is pre-trained on AudioSet and FCNN which is trained from scratch. All these preliminary experiments are conducted without data augmentation methods. From the top two rows of Table 1, we can see that the FCNN model trained from scratch performs better than pre-trained VGGish. Therefore, in the following experiments, FCNN trained from scratch with the official data is used to extract audio embedding.

For visual embedding, large in-domain data sets, such as ImageNet [19] are always adopted for pre-training. Firstly, we make a comparison of video models with and without pre-training. From the third and fourth rows of Table 1, VGG model pre-trained on ImageNet achieves higher accuracy than that trained from scratch, which demonstrates that pre-training is helpful for extracting better visual embedding. In this study, we adopt another scene-image data set Places365 [20] for pre-training and compare different visual embedding. ResNet, ResNeSt and DenseNet are all pre-trained on the Places365 data set. We compare the performance of four outstanding pre-trained networks for video feature extracting as shown between the fifth to eighth rows of Table 1. The results show that, in the AVSC model, FCNN for audio extractor and DenseNet for visual extractor achieve the best performance with an accuracy of 92.5%.

For joint training, the extractor for each modality has two choices: freezing or fine-tuning the parameters. Based on the best extractors (FCNN and DenseNet), we compare three different joint training strategies shown in the bottom three rows of Table 1. We can conclude that fine-tuning can greatly im-

Table 2: An accuracy (in %) comparison on the DCASE development set of different video data perturbation strategies to modify each individual video feature.

Index	Transformation	Accuracy
1	A-FCNN+V-DenseNet	93.2
2	TranslateX	92.6
3	TranslateY	92.6
4	Solarize	90.4
5	ShearX	91.5
6	ShearY	91.0
7	Sharpness	93.3
8	Rotate	92.0
9	Posterize	88.5
10	Invert	91.3
11	Equalize	92.7
12	Cutout	91.5
13	Contrast	93.5
14	Color	92.3
15	Brightness	92.0
16	AutoContrast	93.2

prove performances in audio-visual models. Fine-tuning only the audio or visual extractor achieves better results than no fine-tuning. Fine-tuning both audio and visual extractors can achieve the best performance of 93.2%. That is what we adopt for our follow-up experiments.

3.3. Results on Audio-Video Data Augmentation

We have shown that data augmentation is effective for the ASC task in [6]. On the other hand, not all data augmentation methods are valid for the VSC task. Accordingly, we now investigate the effectiveness of the 15 data perturbation strategies in RandAugment. Table 2 shows the accuracy for 15 video data perturbation policies of the joint audio-visual model. The audio modality in the joint model is trained without any data augmentation. There are two parameters, the number of data augmentation methods N and the magnitude M . The parameters N and M are set to 2 and 14, respectively. The top row in Table 2 provides the accuracy of the joint system without any data augmentation. By comparing the results in Table 2, we can observe that ‘Sharpness’ and ‘Contrast’ are effective to improve system performances for scene classification, while other methods in RandAugment are not suitable for this task. Eventually, we adopt ‘Sharpness’, ‘Contrast’ and ‘Identity Mapping’ as the sub-policies of RandAugment.

The experimental results of audio-visual joint model with different data augmentation methods are presented in Table 3. Clearly, both audio and video data augmentation methods are effective, while the use of both can further improve the performance. Specifically, the accuracy when adopting both audio and video data augmentation is improved to 93.9%. We apply the mixup method for the audio-visual joint model. It has been proved effective for the audio system. Nevertheless, there is no idea about what proportion of mixing can bring the greatest improvement for video clips. Accordingly, the performance comparisons of video model when using different mixup percentage are listed in Table 3. The best accuracy of 94.2% can be achieved when doing mixup on 20% of the data, while the accuracy of audio-visual joint model baseline is 93.2%. Consequently, we set the joint mixup percentage of 20% for the final

Table 3: Experimental results of accuracy (in %) for joint audio-visual data augmentation. The middle three columns denote whether audio augmentation, video perturbation and joint mixup (in term of percentage of training data) are applied. ‘✓’ and ‘-’ denote with and without each operation, respectively.

System	Audio	Video	Mixup	Accuracy
	-	-	-	93.2
	✓	-	-	93.4
A-FCNN	-	✓	-	93.7
+	✓	✓	-	93.9
V-DenseNet	✓	✓	20%	94.2
	✓	✓	50%	93.8
	✓	✓	100%	94.0

Table 4: An accuracy (in %) comparison of the state-of-the-art techniques.

System	Model	Accuracy
1	Baseline [3]	77.0
2	Zhang et al. [23]	94.1
3	Yang et al. [26]	93.9
4	Pham et al. [25]	93.9
5	Proposed	94.2

audio-visual model based on the results of Table 3.

3.4. Overall Comparison

The top row in Table 4 lists the AVSC accuracy of the official baseline system. Based on the original OpenL3 publication, either audio and video embedding is extracted using a single modality. Then it connects the sub-networks using two fully-connected feed-forward layers. Compared with the baseline system, our proposed model achieves promising performance by a large margin. Furthermore, we compare the system results of other teams with the top rankings for DCASE 2021 Task 1b. They are shown in the 2nd, 3rd and 4th rows in Table 4. The accuracy of our final AVSC system is listed in the bottom row. The experimental results clearly demonstrate that the proposed audio-visual joint model with audio-video data augmentation achieves quite a competitive performance when compared with all the state-of-the-art benchmark systems.

4. Conclusion

In this paper, we propose a novel approach to the AVSC task with acoustic-visual joint modeling and data augmentation strategies. Based on multi-modal inputs, we compare different audio and video embedding and achieve the best matching result: FCNN for the audio modality and DenseNet for the video modality. Joint modeling with fine-tuning on both modalities works best with good accuracy improvements. Moreover, we successfully apply RandAugment with ‘Sharpness’ and ‘Contrast’ policies in training AVSC models. A joint mixup strategy for multi-modal AVSC is then proposed for better modality interactions at the input level. Finally, our system obtains the state-of-the-art performance with an accuracy of 94.2% on the development set of DCASE 2021 Task 1b.

5. References

- [1] A. Smith, C. Jones, and E. Roberts, "Article title," *Journal*, vol. 62, pp. 291–294, January 1920.
- [2] T. Baltrušaitis, C. Ahuja, and L.-P. Morency, "Multimodal machine learning: A survey and taxonomy," *IEEE transactions on pattern analysis and machine intelligence*, vol. 41, no. 2, pp. 423–443, 2018.
- [3] S. Wang, A. Mesaros, T. Heittola, and T. Virtanen, "A curated dataset of urban scenes for audio-visual scene analysis," in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 626–630.
- [4] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz, "mixup: Beyond empirical risk minimization," *arXiv preprint arXiv:1710.09412*, 2017.
- [5] 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," *arXiv preprint arXiv:1904.08779*, 2019.
- [6] 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.
- [7] A. Mesaros, T. Heittola, A. Diment, B. Elizalde, A. Shah, E. Vincent, B. Raj, and T. Virtanen, "Dcase 2017 challenge setup: Tasks, datasets and baseline system," in *DCASE 2017-Workshop on Detection and Classification of Acoustic Scenes and Events*, 2017.
- [8] A. Mesaros, T. Heittola, and T. Virtanen, "A multi-device dataset for urban acoustic scene classification," *arXiv preprint arXiv:1807.09840*, 2018.
- [9] T. Heittola, A. Mesaros, and T. Virtanen, "Acoustic scene classification in dcase 2020 challenge: generalization across devices and low complexity solutions," *arXiv preprint arXiv:2005.14623*, 2020.
- [10] D. Battaglino, L. Lepauloux, N. Evans, F. Mougins, and F. Biot, "Acoustic scene classification using convolutional neural networks," *IEEE AASP Challenge on Detec*, 2016.
- [11] H. Zhang, C. Wu, Z. Zhang, Y. Zhu, H. Lin, Z. Zhang, Y. Sun, T. He, J. Mueller, R. Manmatha *et al.*, "Resnest: Split-attention networks," *arXiv preprint arXiv:2004.08955*, 2020.
- [12] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [13] 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.
- [14] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 4700–4708.
- [15] M. Hayat, S. H. Khan, M. Bennamoun, and S. An, "A spatial layout and scale invariant feature representation for indoor scene classification," *IEEE Transactions on Image Processing*, vol. 25, no. 10, pp. 4829–4841, 2016.
- [16] Y. Liu, Q. Chen, W. Chen, and I. Wassell, "Dictionary learning inspired deep network for scene recognition," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018.
- [17] N. Sun, W. Li, J. Liu, G. Han, and C. Wu, "Fusing object semantics and deep appearance features for scene recognition," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 29, no. 6, pp. 1715–1728, 2018.
- [18] J. Xiao, J. Hays, K. A. Ehinger, A. Oliva, and A. Torralba, "Sun database: Large-scale scene recognition from abbey to zoo," in *2010 IEEE computer society conference on computer vision and pattern recognition*. IEEE, 2010, pp. 3485–3492.
- [19] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
- [20] B. Zhou, A. Lapedriza, A. Khosla, A. Oliva, and A. Torralba, "Places: A 10 million image database for scene recognition," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 6, pp. 1452–1464, 2017.
- [21] Q. Kong, Y. Cao, T. Iqbal, Y. Wang, W. Wang, and M. D. Plumbley, "Panns: Large-scale pretrained audio neural networks for audio pattern recognition," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 28, pp. 2880–2894, 2020.
- [22] 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.
- [23] M. Wang, C. Chen, Y. Xie, H. Chen, Y. Liu, and P. Zhang, "Audio-visual scene classification using transfer learning and hybrid fusion strategy," *DCASE2021 Challenge*, Tech. Rep, Tech. Rep., 2021.
- [24] S. Okazaki, Q. Kong, and T. Yoshinaga, "Ldsvision submissions to dcase21: A multi-modal fusion approach for audio-visual scene classification enhanced by clip variants," *DCASE2021 Challenge*, Tech. Rep, Tech. Rep., 2021.
- [25] L. Pham, A. Schindler, M. Schutz, J. Lampert, and R. King, "DCASE 2021 task 1B: Technique report," *DCASE2021 Challenge*, Tech. Rep., June 2021.
- [26] Y. Yang and Y. Luo, "Scene classification using acoustic and visual feature," *DCASE2021 Challenge*, Tech. Rep., June 2021.
- [27] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *International conference on machine learning*. PMLR, 2015, pp. 448–456.
- [28] G. E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, and R. R. Salakhutdinov, "Improving neural networks by preventing co-adaptation of feature detectors," *arXiv preprint arXiv:1207.0580*, 2012.
- [29] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7132–7141.
- [30] E. D. Cubuk, B. Zoph, J. Shlens, and Q. V. Le, "RandAugment: Practical automated data augmentation with a reduced search space," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2020, pp. 3008–3017.
- [31] B. Mcfee, C. Raffel, D. Liang, D. Ellis, and O. Nieto, "librosa: Audio and music signal analysis in python," in *Python in Science Conference*, 2015.
- [32] Q. Bi, K. Qin, Z. Li, H. Zhang, K. Xu, and G.-S. Xia, "A multiple-instance densely-connected convnet for aerial scene classification," *IEEE Transactions on Image Processing*, vol. 29, pp. 4911–4926, 2020.
- [33] S. Hershey, S. Chaudhuri, D. P. Ellis, J. F. Gemmeke, A. Jansen, R. C. Moore, M. Plakal, D. Platt, R. A. Saurous, B. Seybold *et al.*, "CNN architectures for large-scale audio classification," in *2017 IEEE international conference on acoustics, speech and signal processing (icassp)*. IEEE, 2017, pp. 131–135.
- [34] G. Castellano, L. Kessous, and G. Caridakis, "Emotion recognition through multiple modalities: face, body gesture, speech," in *Affect and emotion in human-computer interaction*. Springer, 2008, pp. 92–103.
- [35] G. A. Ramirez, T. Baltrušaitis, and L.-P. Morency, "Modeling latent discriminative dynamic of multi-dimensional affective signals," in *International Conference on Affective Computing and Intelligent Interaction*. Springer, 2011, pp. 396–406.
- [36] Z.-z. Lan, L. Bao, S.-I. Yu, W. Liu, and A. G. Hauptmann, "Multimedia classification and event detection using double fusion," *Multimedia tools and applications*, vol. 71, no. 1, pp. 333–347, 2014.