

# UNSUPERVISED IMPROVEMENT OF AUDIO-TEXT CROSS-MODAL REPRESENTATIONS

Zhepei Wang<sup>‡</sup>, Cem Subakan<sup>b,‡,x</sup>, Krishna Subramani<sup>‡</sup>, Junkai Wu<sup>‡</sup>, Tiago Tavares<sup>†</sup>,  
Fabio Ayres<sup>†</sup>, Paris Smaragdis<sup>‡,α</sup>

<sup>‡</sup>University of Illinois at Urbana-Champaign, <sup>b</sup>Université Laval, <sup>‡</sup>Concordia University,  
<sup>x</sup>Mila-Quebec AI Institute, <sup>†</sup>Inspire, <sup>α</sup>Amazon Web Services

## ABSTRACT

Recent advances in using language models to obtain cross-modal audio-text representations have overcome the limitations of conventional training approaches that use predefined labels. This has allowed the community to make progress in tasks like zero-shot classification, which would otherwise not be possible. However, learning such representations requires a large amount of human-annotated audio-text pairs. In this paper, we study unsupervised approaches to improve the learning framework of such representations with unpaired text and audio. We explore domain-unspecific and domain-specific curation methods to create audio-text pairs that we use to further improve the model. We also show that when domain-specific curation is used in conjunction with a soft-labeled contrastive loss, we are able to obtain significant improvement in terms of zero-shot classification performance on downstream sound event classification or acoustic scene classification tasks.

**Index Terms**— Audio-text representation learning, data augmentation, contrastive learning, sound event classification, acoustic scene classification

## 1. INTRODUCTION

Representation learning methods such as Self-Supervised Learning (SSL) [1] expand the limited scope of supervised learning by learning representations that can be applied to a large variety of downstream tasks. However, the mainstream SSL methods in the literature typically train an encoder on uni-modal data [2, 3, 4].

Learning cross-modal representations that involve text adds the additional flexibility of incorporating language in downstream tasks. This enables downstream tasks such as zero-shot classification possible, where the model is able to perform classification without being restricted by a pre-defined and explicitly annotated label set. The cross-modal representations can also be used in other tasks, such as audio-to-text and text-to-audio retrieval.

Learning cross-modal representations has been explored in computer vision under prior works including CLIP [5], Florence [6], and ALIGN [7]. AudioCLIP [8] extends the CLIP framework to the audio domain by incorporating an audio encoder to learn a joint embedding space for audio, vision, and language using aligned data across the three domains. CLAP [9], on the other hand, learns audio-text embeddings directly without depending on the image domain. It aims to maximize the similarity of the text-and-audio representations that correspond to the paired audio and caption within a given batch. To train high-quality representations, these methods require a large number of paired audio and text items. While in-the-wild audio-text pairs exist at an extensive scale (e.g., captions from online video, metadata from audio datasets), the text is likely

to be irrelevant to the sound events presented in the audio and is therefore unsuitable for training audio representations. Collecting high quality captioned audio is therefore an expensive task, hence data availability might constitute a bottleneck for scaling the audio-text pretraining. To overcome the limitation, Wav2CLIP [10] uses audio-image pairs from video clips to learn audio-text correspondence by distilling from the pre-trained CLIP model. VIP-ANT [11] extends Wav2CLIP by mining additional audio-text pairs from video and text data using pre-trained CLIP. The LAION-CLAP [12] augments the training data by performing keyword-to-caption generation from tags or labels of audio clips using a pre-trained language model. These approaches, however, assume the existence of either an additional anchor modality with paired annotations or a pre-trained model for generating text.

In this work, we explore the possibility of improving the zero-shot classification performance of audio-text representations using unpaired text and audio. For this purpose, we first train an initial teacher model using paired audio clips and captions; we then use this teacher model to automatically align textual descriptions to in-the-wild audio files, using the pairs aligned with higher confidence to train a new model. We then argue that this performance can be further improved by curating a refined, domain-specific dataset that is more akin to the zero-shot classification domain.

<sup>1</sup> Our contributions in this paper are as follows,

- We propose domain-unspecific and domain-specific curation methods and show the improvement on three downstream zero-shot audio classification tasks.
- We propose a soft-labeled training objective that can avoid learning with hard labels in cases where the batch contain similar data items. We show that this training objective significantly improves the performance under domain-specific curation strategies.
- We show that the proposed curation methods significantly improve the model even in the case where the teacher model is trained on a small amount of paired data (%10 of the data).

## 2. METHODOLOGY

### 2.1. Contrastive Language-Audio Pretraining and Zero-Shot Classification

CLAP, “Contrastive Language-Audio Pretraining”, learns a joint latent space between text and audio by maximizing the similarity between the text and audio latent representation for the text caption and its corresponding audio signal. Let,  $X_t, X_a$  respec-

<sup>1</sup>Implementation available at [https://github.com/zhepei/wav2clip\\_curator](https://github.com/zhepei/wav2clip_curator)

tively denote a batch of text and audio. In the CLAP model, the latent representation is obtained by passing the text and audio through the text and audio encoders  $f_t(\cdot)$ , and  $f_a(\cdot)$  such that  $L_t = f_t(X_t)$ ,  $L_a = f_a(X_a)$ , where  $L_t \in \mathbb{R}^{N \times T}$ ,  $L_a \in \mathbb{R}^{N \times A}$ , such that  $T$  is the latent dimensionality of text,  $A$  is the latent dimensionality of audio, and  $N$  is the batch size. CLAP trains a joint latent space by passing  $L_t$  and  $L_a$  through fully-connected layers such that  $t = \text{MLP}_t(L_t)$ ,  $a = \text{MLP}_a(L_a)$ , where  $\text{MLP}(\cdot)$  denotes the multi-layer perceptron transformation layers,  $t \in \mathbb{R}^{N \times d}$ , and  $a \in \mathbb{R}^{N \times d}$  respectively denote the latent variables with same latent dimensionality  $d$ . The model then tries to maximize the diagonal entries on the matrix  $C = ta^\top$ . This translates into the following training loss function,

$$\mathcal{L}(C) = \frac{1}{2} \sum_{i=1}^N \left( \log(\text{Softmax}_t(C/\tau)_{i,i}) + \log(\text{Softmax}_a(C/\tau)_{i,i}) \right), \quad (1)$$

where  $\text{Softmax}_t(\cdot)$  and  $\text{Softmax}_a(\cdot)$  respectively denote Softmax functions along text and audio dimensions,  $\tau$  is a learnable temperature scaling parameter, and the  $C_{i,i}$  denotes the diagonal elements of the  $C$  matrix. We show the training forward pass pipeline in the left panel of Figure 1.

The CLAP model is able to perform zero-shot classification by simply calculating the similarity of a given audio to a fixed set of text prompts constructed from class labels. That is, the classification decision is simply taken to be  $\hat{c} = \arg \max_j t_j^\top a_{\text{test}}$ , where  $\hat{c}$  is the zero-shot classification decision,  $a_{\text{test}}$  is the embedding for the test audio, and  $t_j$  is the text embedding corresponding to the label of class  $j$ . We show the pipeline of zero-shot classification in the right panel of Figure 1.

## 2.2. Unsupervised Improvement of Cross-Modal Audio-Text Representations

In this section, we describe the different strategies we employ to curate a paired dataset from unpaired text and audio. This curated dataset is used to train a student model in order to improve upon the zero-shot performance of the teacher model. We call this curated dataset the *Improvement-Set*. We follow two different strategies to curate this dataset. We call the first strategy the ‘‘Domain-Unspecific’’ (DU) dataset curation, where we form pairs using in-the-wild audio and text. The second strategy is ‘‘Domain-Specific’’ (DS) where we explicitly try to find audio recordings that are more relevant to the zero-shot classification task at hand.

### 2.2.1. Domain-Unspecific Improvement-Set Curation

For this strategy, we use a teacher CLAP model to curate an Improvement-Set. We match the captions of the training data with audio from a large dataset such as AudioSet [13]. We compute the cosine similarity between each pair of audio and text embeddings, and keep pairs with similarity above a threshold  $\sigma \in [0, 1]$ . We show this strategy in Figure 2. As we showcase in the experiments, this strategy provides an improvement over the base model; however, with a domain-specific refinement of the Improvement-Set, we can further enhance the zero-shot performance.

### 2.2.2. Domain-Specific Improvement-Set Curation

For this strategy, the main idea is to narrow down the captions used for the Improvement-Set by calculating text-to-text similarities be-

tween the captions from the training set and in-domain labels for a zero-shot classification task. Once these similarities are obtained, the first domain-specific (DS) strategy is to narrow down the training set by picking the subset of audio-text pairs that are most related to the downstream task through the text modality. We illustrate this procedure in Figure 3.

Another option is to augment the domain-specific Improvement-Set by involving in-the-wild audio data. The process to create the Augmented Domain-Specific (ADS) Improvement-Set, as demonstrated in Figure 4, is as follows:

1. We calculate the similarities between the text embeddings of class labels from the in-domain dataset and the captions from the teacher’s training set. We take the most similar captions from this set by thresholding the similarity values.
2. We find the correspondences between the domain-specific captions and the in-the-wild audio using the teacher CLAP.

## 2.3. Using Soft Labels in the CLAP Loss

An underlying issue with the original CLAP learning objective in Equation (1) is that it neglects local similarities for data within the same batch and treats all negative samples equally. It is possible to sample audio-text pairs with similar content from the same batch (i.e. multiple clips of ‘‘dog-barking’’), while the objective in Equation (1) penalizes the model for not discriminating these similar instances. This issue is exacerbated when the training data is less diverse. For instance, with the DS or ADS Improvement-Set, it is more likely to sample similar inputs within a batch.

The original learning framework of CLAP is equivalent to solving an  $N$ -class classification problem for a batch size of  $N$ , where the label for the  $i$ -th sample of each batch is the one-hot vector  $\mathbf{y}_i \in \mathbb{R}^N$  with only the  $i$ -th entry equal to 1. Similar to prior studies in image-text pretraining [14, 15], we propose a label-softening technique when training the student model. The soft labels are obtained from the similarity between input samples within each batch, as illustrated in Figure 5. For each input batch, we obtain the soft targets by estimating the intra-modal similarity using the teacher audio ( $\tilde{a}$ ) and text ( $\tilde{t}$ ) embeddings. We define,

$$\tilde{C}_a = \tilde{a}\tilde{a}^\top, \quad \tilde{C}_t = \tilde{t}\tilde{t}^\top, \quad (2)$$

where  $\tilde{C}_a, \tilde{C}_t \in \mathbb{R}^{N \times N}$  represent the intra-batch, intra-modal similarity matrices for audio and text, respectively. For the  $i$ -th sample in this batch, we define the audio-to-text soft label as,

$$\tilde{\mathbf{y}}_{a \rightarrow t}^{(i)} = (1 - \beta)\mathbf{y}_i + \beta\tilde{C}_a^{(i)}, \quad (3)$$

where  $\tilde{C}_a^{(i)}$  is the  $i$ -th row of the matrix  $\tilde{C}_a$  and the coefficient  $\beta \in [0, 1]$  adjusts the weights of the hard and soft labels; in symmetry, the text-to-audio soft label can be formulated as,

$$\tilde{\mathbf{y}}_{t \rightarrow a}^{(i)} = (1 - \beta)\mathbf{y}_i + \beta\tilde{C}_t^{(i)}. \quad (4)$$

The overall learning objective is therefore,

$\tilde{\mathcal{L}}(C) = \frac{1}{2}(\tilde{\mathcal{L}}_t(C) + \tilde{\mathcal{L}}_a(C))$ , with

$$\begin{aligned} \tilde{\mathcal{L}}_t(C) &= \frac{1}{N} \sum_{i=1}^N \mathcal{D}_{KL}(\tilde{\mathbf{y}}_{a \rightarrow t}^{(i)} \parallel \text{Softmax}_t(C/\tau)^{(i)}), \\ \tilde{\mathcal{L}}_a(C) &= \frac{1}{N} \sum_{i=1}^N \mathcal{D}_{KL}(\tilde{\mathbf{y}}_{t \rightarrow a}^{(i)} \parallel \text{Softmax}_a(C/\tau)^{(i)}), \end{aligned} \quad (5)$$

where  $\mathcal{D}_{KL}$  denotes KL-Divergence, and  $C$  denotes the similarity matrix obtained from the student model.

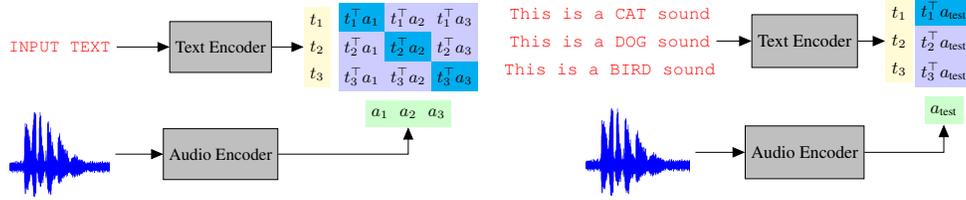


Figure 1: **(left)** The training of the CLAP model for learning cross-modal representations. **(right)** Zero shot classification with the CLAP model

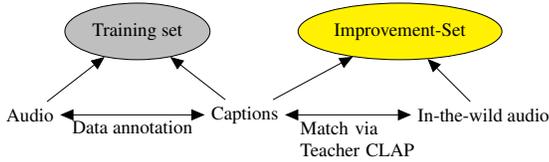


Figure 2: Domain-Unspecific (DU) Curation of Improvement-Set. With the pre-trained teacher CLAP model, the embedding of each in-the-wild audio clip is matched against the embeddings of all training captions. The new audio-text pairs with similarity above a certain threshold  $\sigma$  are included in the Improvement-Set.

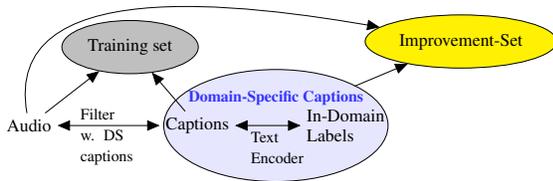


Figure 3: Domain-Specific (DS) Curation of Improvement-Set with Domain-Specific Audio. The Improvement-Set consists of audio-caption pairs from the training set that are most relevant to the downstream task by measuring the similarity between the embeddings of the training captions and in-domain labels of the downstream task.

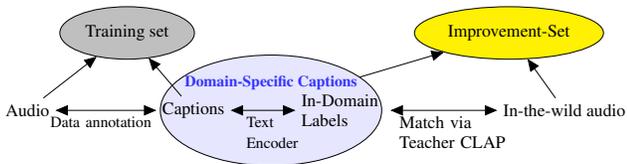


Figure 4: Augmented Domain-Specific (ADS) Curation of Improvement-Set. First, DS Curation is performed to obtain the subset of captions from the training set that are most related to the downstream task. Then, in-the-wild audio clips are aligned with this subset of captions with the pre-trained teacher CLAP model.

### 3. EXPERIMENTAL CONFIGURATIONS

#### 3.1. Training Data

For training the teacher CLAP, we follow the original paper [9] by using 128k paired audio and text captions from FSD50k [16], ClothoV2 [17], AudioCaps [18], and MACS [19]. For the unsupervised curation in Section 2.2, we use the captions from the teacher training set and recordings from the unbalanced training set

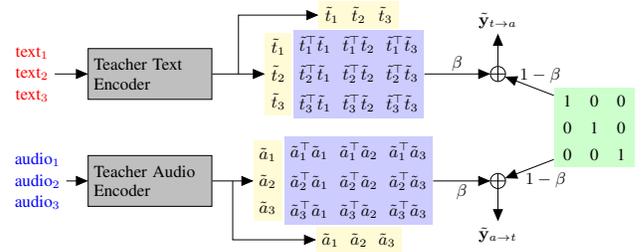


Figure 5: Computation pipeline for soft-labeled loss function

Table 1: Number of audio-caption pairs used in the curated Improvement-Set. DU represents the domain-unspecific curation; DS refers to the domain-specific curation containing the subset of the teacher training set relevant to each downstream task; ADS is the domain-specific curation from the in-the-wild audio.

DU	ESC-50		UrbanSound8K		TUT17	
	DS	ADS	DS	ADS	DS	ADS
$5.0 \times 10^6$	$15.6 \times 10^3$	$1.30 \times 10^6$	$9.7 \times 10^3$	$1.1 \times 10^6$	$8.7 \times 10^3$	$0.3 \times 10^6$

from AudioSet [13] by excluding recordings that are corrupted or contained in AudioCaps for the teacher training, which results in more than 1.9 million clips; notice that no label information from AudioSet is used during curation. We use a similarity threshold  $\sigma = 0.7$  for the domain-unspecific curation and a threshold between 0.6 and 0.75 for the domain-specific experiments. The numbers of audio-caption pairs of the curated datasets are outlined in Table 1.

During training, each input audio recording is resampled to 44.1 kHz. If longer than 5 seconds, a random 5-second segment of the recording is chosen; if shorter, the recording is zero-padded.

#### 3.2. Training Details

Following CLAP [9], we use the CNN14 [20] as the audio encoder and the BERT [3] as the text encoder, each followed by a two-layer MLP as the projection layer. All modules are initialized following the setups in [9]. The embedding vectors have a dimension of 1024. The temperature parameter  $\tau$  is initialized to 0.007, and the coefficient  $\beta$  for soft labels is set to 0.3. The models are trained with 2 Quadro RTX 6000 GPUs using a batch size of 64 per GPU. We perform three runs for each setup and obtain the average results.

For student training with the DU Improvement-Set, we randomly replace the batches with samples replayed from the teacher training set to combat catastrophic forgetting [21]. For student training with ADS Improvement-Set, the replay dataset is selected as the subset of audio-text pairs that are relevant to the downstream task, namely, the DS Improvement-Set. Note that with the addition of

Table 2: Accuracy (percentage) on zero-shot evaluation for teacher and student models based on the teacher trained using full dataset. DU, DS, ADS, and SL denote Domain-Unspecific, Domain-Specific, Augmented Domain-Specific, and soft-labeled loss. Best performances are highlighted in bold.

Model	Zero-Shot Evaluation Set		
	ESC-50	UrbanSound8K	TUT17
CLAP teacher	81.9 ± 0.9	74.8 ± 1.2	29.8 ± 1.3
SL	83.1 ± 1.2	73.9 ± 2.6	30.1 ± 2.1
DU	82.4 ± 1.4	73.9 ± 0.2	31.5 ± 1.0
DU+SL	83.0 ± 0.5	74.9 ± 1.4	29.9 ± 1.9
DS	78.8 ± 0.5	73.2 ± 1.1	29.8 ± 2.6
DS+SL	83.5 ± 0.6	75.5 ± 1.4	31.8 ± 2.6
ADS	84.2 ± 0.5	74.2 ± 2.1	32.5 ± 1.0
ADS+SL	<b>85.1 ± 0.7</b>	<b>77.4 ± 0.6</b>	<b>36.0 ± 1.8</b>

replay, we observed a performance gain of 0.6%, 0.9%, and 2.1% under the ADS curation strategy, on ESC-50, Urbansound8K, and TUT17, respectively compared with the models trained without using replay. We noticed similar trends with other curation strategies as well, so we used replay in all of our results.

### 3.3. Downstream Evaluation

We consider the following datasets for downstream evaluation with sound event classification and acoustic scene classification tasks:

- ESC-50 [22] with 2000, 5-second audio clips from 50 environmental sound classes.
- UrbanSound8K [23] with 8732, 4-second recordings from 10 possible urban sound classes.
- TUT Acoustic Scenes 2017 (TUT17) [24] with 6300, 10-second clips from 15 possible scene classes.

We measure the performance of the models using the accuracy from zero-shot evaluation described in 2.1. Following [9], we add a prefix prompt “*this is a sound of [class label]*” to each label in the downstream dataset before computing text embeddings.

## 4. RESULTS AND DISCUSSIONS

### 4.1. Student Training from Full-dataset Teacher

In Table 2, we showcase the improvements obtained when the teacher model is trained with the full training set.<sup>2</sup> When further training the teacher CLAP model with the soft-labeled loss (SL), the accuracy increases by 1.2% on ESC and by 0.3% on TUT17 while decreasing by 0.9% on UrbanSound8K compared to the teacher model, indicating that soft-labeled loss alone cannot yield significant and consistent performance gains.

We next observe that the domain-unspecific (DU) improvement of the CLAP model improves the zero-shot performance in certain cases. On ESC-50, especially with the addition of soft labels (SL), we are able to increase the zero-shot accuracy from 81.9% to 83.0%. However, we do not observe an improvement for all three downstream datasets with DU.

When curating the Improvement-Set with the downstream domain knowledge, we observe that the Domain-Specific (DS) curation alone does not improve the performance: The results validate our hypothesis that learning with hard labels with similar input data

<sup>2</sup>The CLAP model released by Microsoft obtains 82.6%, 73.4%, and 29.6% zero-shot accuracy values as reported in [9]. Notice that our own CLAP teacher model can closely replicate this performance.

Table 3: Zero-shot accuracy for experiments where the teacher model is trained with 10% of the original training data. As a reference, we also include the CLAP teacher model trained with the full dataset from Table 2.

Model	Zero-Shot Evaluation Set		
	ESC-50	UrbanSound8K	TUT17
CLAP teacher (full-dataset)	81.9 ± 0.9	74.8 ± 1.2	29.8 ± 1.3
CLAP teacher (subset)	74.2 ± 1.3	73.5 ± 2.0	30.9 ± 1.6
DU + SL	78.9 ± 0.3	73.7 ± 1.3	28.8 ± 1.1
ADS + SL	<b>81.3 ± 1.0</b>	<b>74.5 ± 0.7</b>	<b>31.3 ± 0.5</b>

from the same batch would adversely impact the quality of learned representations. However, we observe that adding the soft-labeled loss significantly boosts the zero-shot accuracy by exceeding the performance of the teacher model on all three datasets. We finally make the observation that the augmented domain-specific curation along with soft labels (ADS+SL) substantially improves the performance across the board. The results suggest that additional data with domain-specific curation and soft-labeled loss become more effective when used in conjunction, each of which is crucial to performance improvement.

### 4.2. Student Training from Subset Teacher

In prior experiments, we assume the availability of adequate paired audio and text data for training the teacher model. We would also like to investigate how data curation, label softening, and subsequent student training can be impacted if the amount of paired audio and text is limited during teacher training. To this end, we perform the teacher training by randomly sampling 10% of the original  $128 \times 10^3$  paired audio and caption data. We follow the identical setup for the unsupervised curation and label softening as in the experiments involving the complete dataset. We outline the zero-shot accuracy of the teacher and student models in Table 3.

Student models trained with the DU Improvement-Set enhance the performance on ESC-50 (from 74.2% to 78.9%) and UrbanSound8K (from 73.5% to 73.7%) while slightly degrading on TUT17 (from 30.9% to 28.8%); the results can be attributed to the domain-unspecific curation, which yields data that are more relevant to environmental sound classes and distinct from acoustic scene recordings. By incorporating domain-specific knowledge, the ADS strategy further enhances the performance of ESC-50 (to 81.3%) and UrbanSound8K (to 74.5%) beyond that of the teacher model and also exhibits a slight improvement on TUT17 (to 31.3%). Quite notably we observe that with the proposed ADS+SL strategy, the performance on ESC-50 and UrbanSound8K is similar to that of the teacher model trained with the full dataset. These improvements obtained in the student training demonstrate the effectiveness of the data curation and label softening techniques proposed, even in the absence of ample paired multi-modal data for training the teacher model.

## 5. CONCLUSIONS

In this paper, we propose domain-unspecific and domain-specific data curation methods that can effectively improve the zero-shot classification performance of cross-modal representations. We have identified that using domain-specific dataset curation combined with a soft-labeled loss significantly improves the performance over a baseline teacher model. This observation holds even in the case where the baseline teacher is trained with 10% of the original training data. We plan to further improve our approach by incorporating general text corpora into the curation pipeline and exploring more advanced algorithms for audio-text matching in future work.

## 6. REFERENCES

- [1] J. Gui, T. Chen, Q. Cao, Z. Sun, H. Luo, and D. Tao, "A survey of self-supervised learning from multiple perspectives: Algorithms, theory, applications and future trends," 2023.
- [2] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *Proceedings of the International Conference on Machine Learning (ICML)*, 2020.
- [3] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics*, 2019.
- [4] Y. Gong, C.-I. Lai, Y.-A. Chung, and J. Glass, "Ssast: Self-supervised audio spectrogram transformer," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 10, 2022, pp. 10 699–10 709.
- [5] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, "Learning transferable visual models from natural language supervision," in *Proceedings of the 38th International Conference on Machine Learning*, 2021, pp. 8748–8763.
- [6] L. Yuan, D. Chen, Y.-L. Chen, N. C. F. Codella, X. Dai, J. Gao, H. Hu, X. Huang, B. Li, C. Li, C. Liu, M. Liu, Z. Liu, Y. Lu, Y. Shi, L. Wang, J. Wang, B. Xiao, Z. Xiao, J. Yang, M. Zeng, L. Zhou, and P. Zhang, "Florence: A new foundation model for computer vision," *ArXiv*, vol. abs/2111.11432, 2021.
- [7] C. Jia, Y. Yang, Y. Xia, Y.-T. Chen, Z. Parekh, H. Pham, Q. V. Le, Y. Sung, Z. Li, and T. Duerig, "Scaling up visual and vision-language representation learning with noisy text supervision," 2021.
- [8] A. Guzhov, F. Raue, J. Hees, and A. Dengel, "Audioclip: Extending clip to image, text and audio," 2021.
- [9] B. Elizalde, S. Deshmukh, M. A. Ismail, and H. Wang, "Clap: Learning audio concepts from natural language supervision," *arXiv preprint arXiv:2206.04769*, 2022.
- [10] H.-H. Wu, P. Seetharaman, K. Kumar, and J. P. Bello, "Wav2clip: Learning robust audio representations from clip," 2022.
- [11] Y. Zhao, J. Hessel, Y. Yu, X. Lu, R. Zellers, and Y. Choi, "Connecting the dots between audio and text without parallel data through visual knowledge transfer," in *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Seattle, United States: Association for Computational Linguistics, July 2022, pp. 4492–4507.
- [12] Y. Wu\*, K. Chen\*, T. Zhang\*, Y. Hui\*, T. Berg-Kirkpatrick, and S. Dubnov, "Large-scale contrastive language-audio pre-training with feature fusion and keyword-to-caption augmentation," in *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP*, 2023.
- [13] 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 *Proc. IEEE ICASSP 2017*, New Orleans, LA, 2017.
- [14] A. Andonian, S. Chen, and R. Hamid, "Robust cross-modal representation learning with progressive self-distillation," in *CVPR 2022*, 2022.
- [15] Y. Gao, J. Liu, Z.-H. Xu, T. Wu, W. Liu, J. jin Yang, K. Li, and X. Sun, "Softclip: Softer cross-modal alignment makes clip stronger," *ArXiv*, vol. abs/2303.17561, 2023.
- [16] E. Fonseca, X. Favory, J. Pons, F. Font, and X. Serra, "FSD50K: an open dataset of human-labeled sound events," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 30, pp. 829–852, 2022.
- [17] K. Drossos, S. Lipping, and T. Virtanen, "Clotho: an audio captioning dataset," *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 736–740, 2019.
- [18] C. D. Kim, B. Kim, H. Lee, and G. Kim, "AudioCaps: Generating captions for audios in the wild," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Minneapolis, Minnesota: Association for Computational Linguistics, June 2019, pp. 119–132. [Online]. Available: <https://aclanthology.org/N19-1011>
- [19] I. Martín-Morató and A. Mesaros, "What is the ground truth? reliability of multi-annotator data for audio tagging," *2021 29th European Signal Processing Conference (EUSIPCO)*, pp. 76–80, 2021.
- [20] 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*, 2020.
- [21] Z. Wang, C. Subakan, X. Jiang, J. Wu, E. Tzinis, M. Ravanelli, and P. Smaragdis, "Learning representations for new sound classes with continual self-supervised learning," *IEEE Signal Processing Letters*, vol. 29, pp. 2607–2611, 2022.
- [22] K. J. Piczak, "ESC: Dataset for Environmental Sound Classification," in *Proceedings of the 23rd Annual ACM Conference on Multimedia*. ACM Press, 2015, pp. 1015–1018. [Online]. Available: <http://dl.acm.org/citation.cfm?doid=2733373.2806390>
- [23] J. Salamon, C. Jacoby, and J. P. Bello, "A dataset and taxonomy for urban sound research," in *ACM International Conference on Multimedia*, 2014.
- [24] 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 *Proceedings of the Detection and Classification of Acoustic Scenes and Events 2017 Workshop (DCASE2017)*, November 2017, pp. 85–92.