

Self-supervised Audio Teacher-Student Transformer for Both Clip-level and Frame-level Tasks

Xian Li, Nian Shao, and Xiaofei Li*

Abstract— Self-supervised learning (SSL) has emerged as a popular approach for learning audio representations. One goal of audio self-supervised pre-training is to transfer knowledge to downstream audio tasks, generally including clip-level and frame-level tasks. While frame-level tasks are important for fine-grained acoustic scene/event understanding, prior studies primarily evaluate on clip-level downstream tasks. In order to tackle both clip-level and frame-level tasks, this paper proposes Audio Teacher-Student Transformer (ATST), with a clip-level version (named ATST-Clip) and a frame-level version (named ATST-Frame), responsible for learning clip-level and frame-level representations, respectively. Both methods use a Transformer encoder and a teacher-student training scheme. We have carefully designed the view creation strategy for ATST-Clip and ATST-Frame. Specifically, ATST-Clip uses segment-wise data augmentations, and ATST-Frame integrates frame-wise data augmentations and masking. Experimental results show that our ATST-Frame model obtains state-of-the-art (SOTA) performances on most of the clip-level and frame-level downstream tasks. Especially, it outperforms other models by a large margin on the frame-level sound event detection task. In addition, the performance can be further improved by combining the two models through knowledge distillation. Our code is available online.

Index Terms—Audio self-supervised learning, audio representation learning

I. INTRODUCTION

AUDIO self-supervised learning (SSL), which learns knowledge from a large amount of unlabeled audio data, has emerged as a popular approach for learning audio representations [1]–[10].

The siamese models [3]–[5], [11], [12] maximize the embedding similarity of two augmented views of the same audio clip, having shown a great promise for learning good audio representations. Another promising technical line for audio SSL follows the spirit of BERT (Bidirectional Encoder Representations from Transformers) [13], using Transformer encoder [14] and performing a predictive task for the masked frames [6]–[10].

One goal of audio self-supervised pre-training is to transfer knowledge to downstream audio tasks. Generally speaking, audio tasks are defined within two different ways, i) clip-level tasks are to classify the acoustic scene or event of an entire audio clip, e.g. audio tagging, musical instrument recognition, etc., and ii) frame-level tasks are to detect and recognize event-level timestamps from an audio clip, e.g. sound event detection (SED). Previous studies primarily evaluate their methods on

clip-level audio tasks, leaving the performance on frame-level audio tasks unclear. Clip-level tasks currently account for the majority of the downstream audio tasks. Only a few frame-level tasks have been well-defined in the field, such as speaker diarization and sound event detection. However, the frame-level tasks are more important for fine-grained acoustic scene/event understanding, and they are generally more challenging than clip-level tasks. To handle both clip-level and frame-level tasks, there are several issues to be considered.

In terms of the training criterion, a portion of previous methods focus on learning global representation of an audio clip by using clip-level training criteria [3], [5], [12], while others propose learning local frame-wise or patch-wise representations by using frame-level [6], [7] or patch-level criteria [6], [7], [9], [10], [15], [16]. Most of the clip-level methods allow for extracting frame-wise representations from the intermediate output. However, since the frame-wise representations are not explicitly trained during pre-training, it is questionable whether a model trained by clip-level criterion can perform well on frame-level downstream tasks.

Besides the training criterion, a high temporal resolution for frame-level representations is necessary for frame-level tasks. For Transformer-based methods, the temporal resolution is determined by how the input sequence is organized, either patch-wisely or frame-wisely. SSAST (Self-Supervised Audio Spectrogram Transformer) [6] and MAE-AST (Masked Autoencoding Audio Spectrogram Transformer) [7] have shown that the patch-wise strategy and frame-wise strategy perform differently for different downstream tasks. Other studies [9], [10], [15], [16] only use the patch-wise strategy. Generally, the frame-wise strategy has a better temporal resolution than the patch-wise strategy, and thus may be more suitable for frame-level downstream tasks.

Accounting for learning both clip-level and frame-level audio representations, this paper proposes two models: ATST-Clip and ATST-Frame, where ATST stands for Audio Teacher-Student Transformer. They are developed based on the teacher-student scheme of Bootstrap Your Own Latent (BYOL) [17] and BYOL for Audio (BYOL-A) [5]. They both use Transformer encoder and frame-wise strategy. ATST-Clip and ATST-Frame are responsible for learning global and frame-wise representations by using a clip-level and frame-level training criterion, respectively. This work is a continuation of our previous conference paper [18], in which ATST was first proposed for clip-level representation learning, which is renamed ATST-Clip in this paper to avoid ambiguity. This paper proposes a new ATST-Frame model, and a combination method of the two models based on knowledge distillation.

{lixian, shaonian, lixiaofei}@westlake.edu.cn

¹Westlake Institute for Advanced Study & ²Westlake University, Hangzhou, China

* corresponding author

In addition, the proposed models have been more thoroughly evaluated in this paper.

ATST-Clip draws inspirations from the teacher-student scheme of BYOL [17] and BYOL-A [5], which contains a teacher network and a student network. Given an audio clip, two different views are created through augmentations, e.g. randomly cropping at time dimension. The two views are then separately fed into the teacher network and the student network. Considering the similarity of the two views, the student network weights are updated by maximizing the embedding similarity of the two views. The teacher network weights, on the other hand, are updated by taking exponential moving average (EMA) of the student network weights. ATST-Clip proposes to replace the convolutional neural network (CNN) encoder of BYOL-A with a Transformer encoder, which shows a clear superiority over the CNN encoder, especially for learning the long-term semantic information of speech. More importantly, a new view creation strategy is proposed to fully leverage the capability of Transformer encoder. BYOL-A uses one short segment to create two views. Instead, we propose to use two different long segments to create the two views, which is more fit for Transformer, as the network can learn longer temporal dependencies. The length of segments is carefully studied to control the distinction and overlap of the two segments, which is especially important for rationalizing the difficulty of matching the representations of the two views at latent space.

ATST-Frame extends ATST-Clip to explicitly learn frame-wise representations by maximizing the agreement of student’s frame-level embeddings to the teacher’s frame-level embeddings. Creating proper views for teacher and student branches is the key for achieving a proper difficulty for matching the frame-level embeddings, and thus guiding the model to learn meaningful frame-level representations. Both teacher and student branches process the entire audio clip to maintain the frame-to-frame correspondence between their output sequences. To increase the matching difficulty, data augmentation is applied to one of the teacher and student branches. Moreover, masking is further applied to the student branch to encourage the model to learn semantic relations between frames by accomplishing the prediction of masked frames. Our experiments show that data augmentation and masking are both necessary and are a good combination for frame-level audio pre-training within the teacher-student framework.

Finally, as the training criterion of ATST-Frame and ATST-Clip are totally different, they could learn complementary features. We propose to combine ATST-Frame and ATST-Clip at the fine-tuning stage of downstream tasks, based on cross-model knowledge distillation, which outperforms ATST-Frame or ATST-clip alone.

We use the large-scale AudioSet [19] for pre-training, and evaluate the models with a variety of clip-level downstream tasks and two frame-level downstream tasks. Downstream tasks cover multiple audio domains: environmental sound, speech, and music. Our results show that i) on clip-level tasks, after fine-tuning, the proposed models outperform other state-of-the-art (SOTA) methods for most of the tasks. Especially,

the precision on the AudioSet-2M and AudioSet-20K datasets reach a new SOTA of 49.7% and 40.5% (without model ensembling), respectively. ii) on the frame-level SED task, the proposed ATST-Frame model performs particularly well, outperforming ATST-Clip and other methods by a large margin. We open-source our code online¹ for the research community to replicate and expedite future research.

II. RELATED WORKS

This section introduces related works on audio self-supervised learning.

A. Siamese Models

Siamese models use a two-tower architecture, in which each tower processes a view of the data sample and the embedding similarity of the two views are maximized during training [11], [20]. This idea often confronts the issue of model collapse, e.g. the model can find an easy solution to output a constant value for any inputs. Various training strategies are developed to avoid model collapse. Most of these strategies are originally developed in image SSL pre-training and then are adopted by audio SSL pre-training. One of the strategies is contrastive learning, which introduces negative samples and not only pull the two views close in the latent space but also push them far away from negative samples in the latent space, e.g. SimCLR (Simple Framework for Contrastive Learning of Visual Representations) [21] in image SSL and its audio counterparts COLA (CONtrastive Learning for Audio) [3], [22]. However, the negative samples are possibly similar to positive samples in some scenarios, which will harm the pre-training performance. Due to this reason, some recent works investigated to train siamese models without using negative samples, e.g. BYOL [17] in image SSL pre-training and its audio counterpart BYOL-A [5]. As this work is inspired by the BYOL-style strategy, we will introduce the framework of BYOL-A in Section III-A.

Our ATST-Clip extends BYOL-A to use Transformer encoder, and proposes a new view creation strategy to fit the Transformer encoder. Our ATST-Frame further extends ATST-Clip to explicitly learn frame-wise representations.

B. Masked Audio Modelling

Other methods follow the line of Masked Language Modelling (MLM) [13]. This kind of method has been first applied to speech self-supervised pre-training [23]–[27], and then to audio self-supervised pre-training, e.g. SSAST [28], Conformer-based audio SSL method [8], MAE-AST [7] and Audio-MAE (Audio Masked Autoencoders) [9]. The idea is to mask an arbitrary region of the input, and then perform a prediction task on the masked region. Some of them train the model by reconstructing the masked region [23], [26], while others replace the reconstruction loss with a classification loss. Wav2vec2 [24] and its follower [8], a method based on Conformer (Convolution-augmented Transformer), solve a frame-level contrastive problem by introducing positive and negative

¹<https://github.com/Audio-WestlakeU/audioss1>

frames. SSAST [6] jointly solves a masked reconstruction and a wave2vec-style contrastive problem. HuBERT (Hidden-Unit BERT) [25] creates pseudo classification labels by performing clustering on MFCC (Mel-Frequency Cepstral Coefficient) features or output features of the model trained in the previous iteration. BEATs (Bidirectional Encoder representation from Audio Transformers) [10] proposes an iterative audio pre-training framework, where an acoustic tokenizer and an audio SSL model are iteratively optimized. From the perspective of model architecture, some works [7], [9], [15], [29] follow the asymmetric encoder-decoder structure of masked autoencoders (MAE) [30], in which the encoder encodes the unmasked region, while the decoder processes both the masked and unmasked regions and reconstructs the masked region.

The most similar works with our ATST-Frame are data2vec [27] and M2D (Masked Modeling Duo) [16]. They both use a teacher-student scheme, in which the student encodes a masked version of the training sample and the teacher provides the training/prediction target for the student. The teacher take as input either the unmasked version of the same training sample (data2vec) or only the masked parts of the student input (M2D). Besides, they both use a frame/patch-level criterion.. However, there exist several major differences: i) data augmentation is applied in our ATST-Frame, but not in data2vec and M2D. Data augmentation is critical for adjusting the prediction difficulty; ii) data2vec constructs the training/prediction target for the student network by taking the average of the last eight Transformer blocks of the teacher encoder, while our ATST-Frame uses the asymmetric structure of the BYOL [17], where an extra predictor network is set for the student branch. iii) M2D organizes spectrograms patch-wisely and uses a MAE structure in the student branch, while ATST-Frame adopts a frame-wise strategy and uses a regular Transformer encoder architecture.

III. THE PROPOSED METHOD

Two models are proposed in this work: ATST-Clip and ATST-Frame. ATST-Clip focuses on learning the global representation of an audio clip, while ATST-Frame focuses on learning frame-wise representations. Both of them use a Transformer encoder [14] to process audio spectrograms. And both of them are trained in a teacher-student scheme [17], in which the teacher model is updated by an exponential moving average (EMA) of the student model, while the student model is updated by maximizing the similarity of its embedding to the embedding of teacher.

We will introduce the baseline teacher-student scheme in Section III-A, the Transformer encoder in Section III-B, and then present ATST-Clip and ATST-Frame in Section III-C and Section III-D, respectively. The combination of ATST-Clip and ATST-Frame is presented in Section III-E.

A. Baseline Teacher-Student Scheme

In this work, we adopt the teacher-student scheme as our baseline framework, which was first proposed by Bootstrap your own latent (BYOL) [17] for image pre-training, and

adopted by BYOL-A [5] for audio pre-training. In BYOL-A, given one augmented view of an audio clip, the student network is trained to predict a data representation being close to the teacher network’s representation on another augmented view of the same audio clip. During training, the teacher network weights are updated by taking the EMA of the student network weights.

Formally, the student network, defined by a set of weights θ , contains an encoder f_θ , a projector g_θ and a predictor q_θ , while the teacher network, defined by a set of weights ϕ , contains only an encoder f_ϕ and a projector g_ϕ . The encoder, a CNN in BYOL and BYOL-A, extracts representations from the augmented views. The projectors and the predictor are multi-layer perceptrons (MLPs) that consist of a linear layer (with output dimension of 4096) followed by batch normalization, rectified linear units (RELU), and a final linear layer (with output dimension of 256). The output representation of encoder is used in downstream tasks. Using projectors in pre-training is shown to improve the representation quality [21]. It has been shown that the additional predictor in the student network (combined with the stop-gradient operation of teacher network) is the key factor for preventing the model from collapsing [20]. During training, ϕ is updated by the EMA of θ as: $\phi \leftarrow m\phi + (1-m)\theta$, where m is a decay rate. θ is updated as follows. Let $(\mathbf{X}, \mathbf{X}')$ be two views created from an audio clip. \mathbf{X} is fed into the teacher network to obtain $\mathbf{h} = f_\phi(\mathbf{X})$ and $\mathbf{z} = g_\phi(\mathbf{h})$. \mathbf{X}' is fed into the student network to obtain $\mathbf{h}' = f_\theta(\mathbf{X}')$, $\mathbf{z}' = g_\theta(\mathbf{h}')$ and $q_\theta(\mathbf{z}')$. \mathbf{z} and $q_\theta(\mathbf{z}')$ are then L2-norm normalized to $\bar{\mathbf{z}}$ and $\bar{q}_\theta(\mathbf{z}')$, and the mean square error (MSE) loss between them is computed:

$$L_\theta = \|\bar{\mathbf{z}} - \bar{q}_\theta(\mathbf{z}')\|_2^2 \quad (1)$$

A symmetric loss L'_θ is also calculated by feeding \mathbf{X} to the student network and \mathbf{X}' to the teacher network. During training, θ is updated by minimizing $L_\theta^{total} = L_\theta + L'_\theta$.

In BYOL-A, the encoder is a CNN. The proposed models will replace the CNN encoder with a Transformer encoder, and use the same projectors and predictors as BYOL-A.

B. Audio Spectrogram Transformer Encoder

Both ATST-Clip and ATST-Frame use the same encoding network architecture. The raw waveform is first converted to log-mel spectrogram $\mathbf{X} \in \mathbb{R}^{L \times C}$, where L and C denote the number of frames and the number of frequency bins, respectively. Since modeling long sequences with Transformer is computationally demanding, four consecutive frames of \mathbf{X} are stacked as one frame to reduce the sequence length. The stacked frames are fed to a linear projection layer (with output dimension of d) to obtain a new embedding sequence $\mathbf{E} \in \mathbb{R}^{\frac{L}{4} \times d}$ as the input sequence of Transformer encoder. The embedding sequence is then added with a trainable absolute lookup table positional embedding $\mathbf{P} \in \mathbb{R}^{(\frac{L}{4}) \times d}$. Eventually, the embedding sequence is processed by a Transformer encoder, obtaining an output embedding sequence of $\mathbf{O} \in \mathbb{R}^{(\frac{L}{4}) \times d}$. Our Transformer encoder architecture is the same as Vision Transformer [31], which is a Pre-LN (Layer Norm) Transformer [32]. To represent the entire clip, ATST-Clip

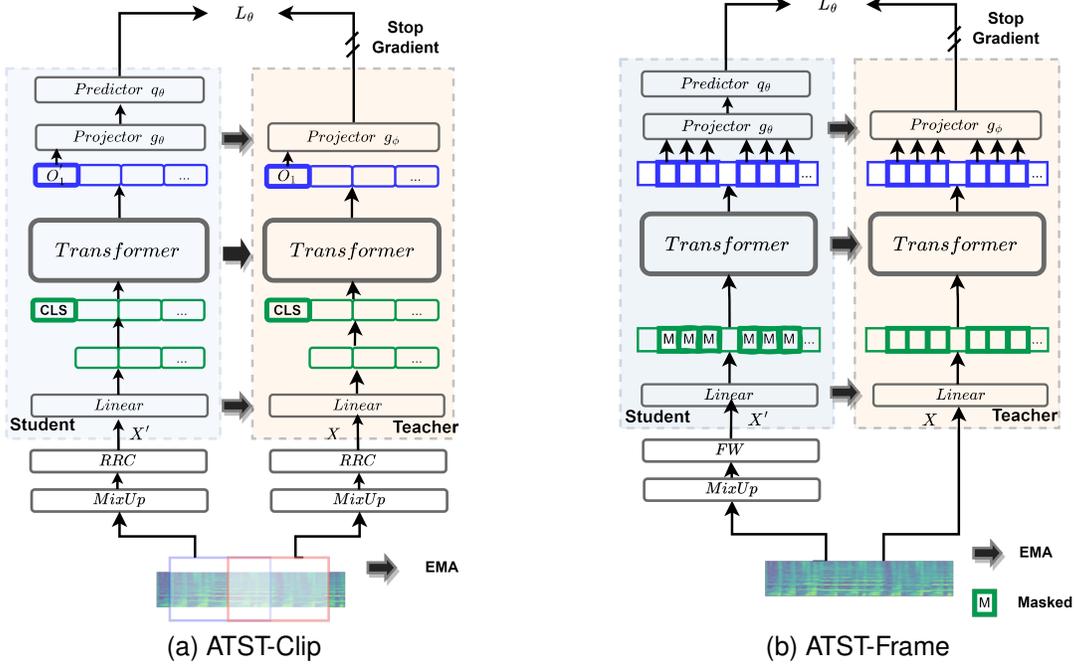


Fig. 1: The proposed methods. (a) ATST-Clip (b) ATST-Frame. The loss L_θ is computed by feeding X to the teacher branch and X' to the student branch. The symmetric loss L'_θ can be computed by swapping X and X' (not shown in the figure).

incorporates an extra trainable class token $[\text{CLS}] \in \mathbb{R}^{1 \times d}$, which will be detailed in III-C.

C. ATST-Clip

The major difference between ATST-Clip and BYOL-A is twofold. ATST-Clip uses a Transformer encoder to leverage its powerful abilities in modeling long-term dependencies and uses a new view creation strategy specifically fit for the Transformer encoder.

First, given an audio clip, it is converted from the waveform domain to the log-mel spectrogram, from which two views are created through a set of augmentations. The two views are then fed into the student and teacher branches respectively, generating a clip-level representation at each branch. In the end, the two clip-level representations are used to calculate the training loss.

1) *Creation of Views*: For siamese model, the two views should be similar with each other to be identified as the same sample yet different enough to increase the difficulty of the identification. BYOL-A [5] randomly crops a single 1-second segment from the input audio and then creates two views by applying different data augmentations to this single segment. It is considered in BYOL-A [5] that different segments may be too different to be identified as the same sample. The work in [4] uses two segments to create two views, however, it uses negative samples to mitigate the problem caused by using two segments.

Our view creation strategy is shown in Fig. 1(a). The time domain input audio clip is first transformed to log-mel spectrogram. We randomly crop two different segments from the log-mel spectrogram. Then, two types of data augmentation

are applied to each of the segments, creating two views of the input audio clip, i.e. (X, X') . The augmentations we employed include Mixup [5] (a modified version of the original Mixup [33], [34]) and Random Resize Cropping (RRC) [5] (adapted from RRC [35] in computer vision to accommodate audio signals).

In order to take full advantage of the Transformer’s ability in modeling long-term dependencies, the proposed method intends to use longer segments, e.g. 6-second segments randomly cropped from 10-second training audio clips in our experiments. The proposed method separately creates two views from two different segments for the purpose of increasing the difficulty of identifying the two views as the same sample, thus leading the model to learn more generalized representations. On the other hand, the two segments cannot be too far away from each other, otherwise, the similarity between them is completely lost. This is guaranteed by properly setting the segment length to make the two segments have a certain portion of overlap. Overall, the proposed strategy does not lose the rationality of identifying two segments as the same sample due to the overlap constraint, and meanwhile increases the task difficulty by using two segments and thus helping to learn more generalized presentation.

2) *Encoding*: The encoding procedure is illustrated in Fig. 1(a). To obtain a representation for the entire clip, an extra class token is used. First, a linear projection layer processes the view, X or X' , obtaining an embedding sequence, at the beginning of which a trainable class token $[\text{CLS}] \in \mathbb{R}^{1 \times d}$ is inserted. The embedding sequence is then added with a trainable absolute lookup table positional embedding sequence, and then fed into the encoder. The class token $[\text{CLS}]$ is widely

used for sentence embedding in neural language processing [13], global image embedding [36], as well as audio segment embedding [28]. It aggregates information from the embedding sequence at every Transformer blocks with the self-attention mechanism. In the output embedding sequence, the class token, denoted as $\mathbf{O}_1 \in \mathbb{R}^{1 \times d}$, is taken as the final clip representation. \mathbf{O}_1 is then processed by the following projector (and predictor).

3) *Loss Function*: The loss function is the same as the one in the baseline scheme described in Section III-A.

D. ATST-Frame

As BYOL-A and ATST-Clip have shown powerful abilities in learning clip-level audio representations with the teacher-student scheme, we further adopt the teacher-student scheme to develop ATST-Frame, which explicitly learns fine-grained frame-wise representation. ATST-Clip creates two different views of the audio clip, and then maximizes the agreement between the clip-level representation of the two views. Instead, ATST-Frame maximizes the agreement between the frame-level representations of two views. The key is to properly design the two views to achieve a good trade off between the difficulty and rationality of the frame-level pretext task.

1) *Creation of Views*: Different from ATST-Clip which randomly crops the audio clip, ATST-Frame processes the entire audio clip. The reasons are i) the frame correspondence of the two views should be preserved for measuring the frame-level agreement; and ii) in order to take full advantage of Transformer in modeling long-term dependencies, the views are set to be as long as possible.

To increase the difference of the two views, and thus increase the task difficulty, data augmentation is first applied to one of the two views. This time, the augmentations should preserve the frame correspondence of the two views. Two augmentations are used: Mixup [5] and Frequency Warping (FW). For computational efficiency, FW is implemented in the spectrogram domain through cropping and then resizing at the frequency axis. Specifically, the input log-mel spectrogram $\mathbf{X} \in \mathbb{R}^{L \times C}$, is first cropped at the frequency axis as $\mathbf{X}_{1:L, 1:a}$, and then resized by bi-cubic interpolation at the frequency axis as $\mathbf{X}^{FW} \in \mathbb{R}^{L \times C}$, where C is the number of the mel frequency bins, and the integer number a is uniformly sampled from the frequency range of $[C * 0.6, C]$. These operations lead to an approximate yet efficient frequency warping.

Due to the constraint of frame-to-frame correspondence for the two views, data augmentation does not bring sufficient task difficulty. Thence, we adopt BERT-like masking [13], which masks/replaces a portion of the frames with a certain trainable mask token and then performs a prediction task to predict the masked frames. The student network takes as input a masked version of the training sample and learns to predict the masked frames, while the prediction target is provided by the teacher network taking as input the unmasked version of the training sample. To prevent the model cheating by simply interpolating, we adopt the group masking strategy [24] that forces N adjacent frames to be masked together. Specifically, we set a probability of 0.65 for masking and

force five adjacent frames to be masked together as a masked block, and the masked blocks are allowed to overlap. With this setting, approximately 50% of the frames are masked. As will be explained later, the pre-training loss will be computed only on the masked frames.

Overall, combining data augmentation and masking is able to create two proper views for the frame-wise pretext task within the teacher-student framework. Although these techniques, i.e. data augmentation, masking and teacher-student scheme, have already been individually (or together with other techniques) used for audio pre-training in the literature, this work carefully integrates them in a different way from other methods, and achieves noticeably better performance. For example, frame-level training of other teacher-student schemes [16], [27], [37] do not use data augmentation. And other masking-based methods [7], [9] reconstruct the masked region.

2) *Encoding*: The encoding procedure is illustrated in Fig. 1(b). The input log-mel spectrogram is first data augmented with Mixup and FW for one view, then processed with a linear projection. After the linear projection, the embedding sequence $\mathbf{E} \in \mathbb{R}^{\frac{L}{4} \times d}$ is randomly masked along the time dimension, only for the student branch. Each masked frame is substituted with a trainable vector $\mathbf{M} \in \mathbb{R}^{1 \times d}$. Subsequently, the embedding sequence is added with a trainable absolute lookup table positional embedding sequence, and then fed into the encoder. After the encoder, the unmasked frames are thrown away, and only the masked frames are further processed by the following projector (and predictor).

3) *Loss Function*: The loss function of ATST-Frame differs from the one of baseline or ATST-Clip in the sense that the loss of ATST-Frame is computed frame-wisely. Feeding the two views, i.e. \mathbf{X} and \mathbf{X}' , to the teacher branch and the student branch, we obtain $\mathbf{z} \in \mathbb{R}^{N_{\text{mask}} \times 256}$ and $q_{\theta}(\mathbf{z}') \in \mathbb{R}^{N_{\text{mask}} \times 256}$, respectively, where N_{mask} denotes the number of masked (stacked-)frames. The MSE loss is calculated on the L2-normalized embeddings $\bar{\mathbf{z}}$ and $\bar{q}_{\theta}(\mathbf{z}')$ as

$$L_{\theta} = \frac{1}{N_{\text{mask}}} \sum_{i=1}^{N_{\text{mask}}} \|\bar{\mathbf{z}}_i - \bar{q}_{\theta}(\mathbf{z}'_i)\|_2^2. \quad (2)$$

A symmetric loss L'_{θ} is also calculated by feeding \mathbf{X} to the student network and \mathbf{X}' to the teacher network. During training, θ is updated by minimizing $L_{\theta}^{\text{total}} = L_{\theta} + L'_{\theta}$.

E. Combine ATST-Clip and ATST-Frame

ATST-Clip and ATST-Frame focus on learning global clip-level representation and local frame-level representations, respectively. Combining them may yield more comprehensive representations. ATST-Clip is trained by a clip-level criterion, but still can extract frame-wise representations with the encoder. However, these representations are only used for storing local information from which the class token [CLS] can aggregate information, but are not optimized specifically to fit frame-level downstream tasks. On the other hand, ATST-Frame can obtain a clip-level representation by applying average pooling to the frame-wise representations. However, this type of average information is not optimized specifically to fit clip-level downstream tasks.

It is possible to jointly train ATST-Clip and ATST-Frame, e.g. add a [CLS] token to ATST-Frame and then perform multi-task learning. Actually, ASiT [37] is trained with a combination of three tasks, including a global task of maximizing agreement, a local task of maximizing agreement, and a reconstruction task. However, our preliminary experiments show that, when creating two views, it is hard to trade off the task difficulty for both ATST-Clip and ATST-Frame. ATST-Clip requires two different randomly cropped segments and the two segments have only a certain portion of overlap, while ATST-Frame asks for a frame-to-frame correspondence between the two views. Besides, ATST-Clip and other clip-level audio siamese methods [5], [12], [22] largely leverage the RRC augmentation [5] to achieve a good performance, but RRC will distort the frame-to-frame correspondence.

Therefore, we leave ATST-Clip and ATST-Frame trained separately to maintain their own advantages and combine them in the evaluation stage. A straightforward way is to ensemble the two models at the inference stage, e.g. to concatenate the outputs of the two models, but this will double the computational cost for inference. Instead, we use knowledge distillation [38] to combine the two models. CMKD [39] has explored cross-model knowledge distillation in the context of supervised audio tagging, e.g. distilling knowledge from EfficientNet-B0 [40] to AST [28]. We follow the principle of CMKD. We first fine-tune ATST-Clip on a downstream task, and then use the fine-tuned ATST-Clip as a teacher to fine-tune ATST-Frame on the same downstream task. Specifically, ATST-Frame is fine-tuned by using two classification losses computed with the ground-truth labels and the ATST-Clip predictions, respectively, and the two losses are weighted with a balance term $\lambda = 0.5$. This strategy is denoted as ATST-C2F. Or we can reverse the fine-tuning order to have the ATST-F2C strategy. These strategies approximately double the fine-tuning time compared with ATST-Frame or ATST-Clip alone, but do not increase the computational cost for inference.

F. Transferring to Downstream Task

For both ATST-Clip and ATST-Frame, after pre-training, we discard the projector in the teacher network, and use the teacher encoder to extract embeddings for downstream tasks.

IV. EXPERIMENTAL SETUP

We conduct extensive experiments using the large-scale AudioSet [19] for pre-training and a variety of downstream tasks for evaluation. The evaluation is performed under the protocol of linear evaluation or fine-tuning. In linear evaluation, the pre-trained encoder is frozen as a feature extractor, on top of which a linear classifier is trained. Whereas in fine-tuning, the pre-trained encoder and linear classifier are fine-tuned together.

A. Pre-training

We use AudioSet [19] for pre-training. The full AudioSet contains 2 million audio clips captured from Youtube videos, with a fixed clip length of 10 seconds. The AudioSet is published with an unbalanced set with 2,042,985 clips and

	#parameters	#blocks	#heads	dimension
ATST-Clip _{small}	22M	12	6	384
ATST-Clip	86M	12	12	768
ATST-Frame _{small}	22M	12	6	384
ATST-Frame	86M	12	12	768

TABLE I: The size of models.

a balanced set with 22,176 clips. We use the unbalanced set of the AudioSet for pre-training. Due to the change of YouTube video availability, the unbalanced set we use contains 1,912,024 clips (AS-1.9M).

For both ATST-Clip and ATST-Frame, a base model is trained using AS-1.9M, which contains 12 Transformer encoder blocks, and 12 heads for each block. The dimension and inner dimension are 768 and 3072, respectively. Besides, we also trained a small model for them for accelerating the development process and conducting ablation studies, using a subset of 200 thousand randomly sampled audio clips (AS-200K). The small model contains 12 Transformer encoder blocks, and 6 heads for each block. The dimension and inner dimension are 384 and 1536, respectively. In the following, we use ATST-Clip and ATST-Frame to represent the base models by default, while ATST-Clip_{small} and ATST-Frame_{small} for the small models.

Audio is re-sampled to 16 kHz. Audio clips are transformed to the log-mel spectrogram domain, with a Hamming window, a window length of 64 ms, a hop size of 10 ms, and 64 mel-frequency bins ranging from 60 Hz to 7800 Hz. The mel-spectrogram feature is min-max normalized, where the minimum and maximum values are calculated globally on the pre-training dataset.

ATST-Clip: We intentionally set the length of two segments (for creating two views) to 6 seconds, which will lead to a segment overlap of at least 2 second, considering that the length of audio clip is 10 seconds. The two randomly sampled segments are augmented by Mixup and RRC with the same configurations used in BYOL-A [5].

ATST-Frame: We use the entire audio clip (10 seconds in AudioSet) for training. We set a probability of 0.65 for masking, and force five adjacent frames to be masked.

Hyper-parameters for pre-training are listed in Table II. We pre-train our models with the AdamW optimizer [41]. The learning rate is warmed up for 10 epochs, and then annealed to 10^{-6} at cosine rate [42]. Similar to DINO [36], the weight decay of Transformer is increased from 0.04 to 0.4 at cosine rate. The EMA decay rate increases from an initial value to 1 at cosine rate.

B. Clip-level Downstream Tasks

1) *Datasets:* Evaluations are carried out on a variety of clip-level downstream tasks, which cover multiple audio domains: environmental sound, speech and music.

- **AS-20K** for multi-label sound event classification. We use the balanced set of AudioSet, with 527 audio classes. We successfully downloaded 20,886 audio clips for training and 18,886 audio clips for evaluation.

	ATST-Clip _{small}	ATST-Clip	ATST-Frame _{small}	ATST-Frame
Dataset	AS-200K	AS-1.9M	AS-200K	AS-1.9M
Optimizer	AdamW	AdamW	AdamW	AdamW
Batch size	1536	1536	1024	864
Learning rate	5e-4	2e-4	4e-4	8e-5
Warm up (epochs)	10	10	10	10
Epochs	300	200	300	200
Initial EMA Decay Rate	0.99	0.9995	0.997	0.9996
Initial Weight Decay	0.04	0.04	0.04	0.04
Final Weight Decay	0.4	0.4	0.4	0.4
Drop path	0.1	0.1	0.1	0.1
Dropout	0	0	0	0

TABLE II: Hyper-parameters for pre-training.

- **AS-2M** for multi-label sound event classification. We use the unbalanced set and balanced set (1,932,110 clips in total) together for training, and use the 18,886 audio clips for evaluation.
- **US8K** for single-label audio scene classification. We use the Urbansound8k dataset [43] to classify audio clips (less than 4 seconds) into 10 classes. It contains 8,732 audio clips and has ten folds for cross-validation.
- **SPCV2** for spoken command recognition. We use Speech Command V2 [44] to recognize 35 spoken commands for one second of audio. It contains 84,843, 9,981 and 11,005 audio clips for training, validation and evaluation respectively.
- **VOX1** for speaker identification. We use the Voxceleb1 dataset [45], with 1,251 speakers. It contains 13,8361, 6,904 and 8,251 for training, validation and evaluation, respectively.
- **NSYNTH** for musical instrument classification. We use the NSYNTH dataset [46], to recognize 11 musical instrument family classes from 4-second audio clips.
- **FSD50K** for multi-label sound event classification. We use FSD50K dataset [47], which contains 36,796, 4,170 and 10,231 audio clips for training, validation and evaluation, respectively.

2) *Metric*: We take classification accuracy (Acc) as the performance metric for the single-label tasks, including audio scene classification, spoken command recognition, speaker identification and musical instrument classification, and mean average precision (mAP) for the tasks of multi-label sound event classification.

For datasets containing validation set, we use the validation set for hyper-parameters tuning and model selection, and report metric score of the selected model on evaluation set. As for AS-20K and AS-2M, we tune hyper-parameters and report metric score on evaluation set. Note this is a common practice for AudioSet fine-tuning [9], [10], [28], mainly because it is non-trivial to sample a meaningful validation set from AudioSet due to extreme class imbalance and label co-occurrence. For

US8K, we conduct 10-fold cross-validation, and report the average accuracy of the 10 folds.

3) *Clip-level Embedding Extraction*: For clip-level downstream tasks, the pre-trained models need to provide a clip-level representation. ATST-Clip is directly designed for this purpose. Although ATST-Frame is designed for frame-wise learning, the average of frame-level representations could also be a reasonable clip-level representation. Therefore, we evaluate both ATST-Clip and ATST-Frame for the clip-level downstream tasks.

In linear evaluation experiments, we use the output of all 12 encoder blocks to construct the clip-level representation. For ATST-Clip, the embedding for the class token and the average of the rest of embedding sequence are first concatenated for each block. We find that the latter still provides some extra information besides the former. Then, the embeddings are concatenated over blocks. For ATST-Frame, the average of the embedding sequence for all blocks are concatenated. The embedding dimensions for ATST-Clip_{small}, ATST-Clip, ATST-Frame_{small} and ATST-Frame are $384 \times 12 \times 2$, $768 \times 12 \times 2$, 384×12 and 768×12 , respectively.

In fine-tuning experiments, we only use the output of the last block. The embedding dimensions for ATST-Clip_{small}, ATST-Clip, ATST-Frame_{small} and ATST-Frame are 384×2 , 768×2 , 384 and 768 , respectively.

The long audio clips will be split into chunks without overlap, with a chunk length of 6 seconds and 10 seconds for ATST-Clip and ATST-Frame respectively. In pre-training, ATST-Clip (ATST-Frame) processes audio clips with a fixed length of 6 seconds (10 seconds), and accordingly the length of positional embedding sequence is also 6 seconds (10 seconds). To account for the length of positional embedding, ATST-Clip (ATST-Frame) will also process audio clips not longer than 6 seconds (10 seconds) for downstream tasks. Note that, audio clips that are longer than 12 seconds are first centrally cropped with a maximum length of 12 seconds, thus there will be at most two chunks for one clip. The chunks are independently processed by the pre-trained models, and their outputs are averaged to obtain the final clip representation.

4) *Downstream Task Training*: In linear evaluation experiments, we train the linear classifier for 100 epochs with the SGD optimizer. The learning rate is annealed to 10^{-6} at cosine rate during training. The optimal initial learning rate is searched for each task separately. Batch size is set to 1024. Data augmentation is not used.

In fine-tuning experiments, we fine-tune all models with the SGD optimizer. The learning rate is warmed up for 5 epochs, and then annealed to 10^{-6} at cosine rate [42]. The learning rate is also scheduled by a layer-wise learning rate schedule [48], in which the learning rate is multiplied with a scaling factor computed with the scaling function $s(i) = \alpha^{n-i}$, where n and i are the number of layers and the layer index, respectively, and α is usually less than 1 and set to be 0.75 in our experiments. The optimal learning rate is searched for each task separately. Batch size is set to 512 for SPCV2, Vox1, AS-20K, FSD50K, NSYNTH and 1024 for AS-2M. We trained AS-2M for 10 epochs, SPCV2, VOX1 and NSYNTH for 50 epochs, FSD50K for 100 epochs and AS-20K for 200 epochs.

For SPCV2, FSD50K, NSYNTH, AS-20K and AS-2M, we use Mixup [33], [34] and RRC for data augmentation. For VOX1, data augmentation is not applied. For AS-2M, we use balance sampling [40], which is a common strategy to train AudioSet, due to its unbalanced distribution of classes. Note that, the Mixup methods used for pre-training and downstream tasks are different, the former only mixes the audio clips since labels are not involved in training, while the latter mixes both audio clips and labels.

C. Frame-level Downstream Task

1) *Dataset*: The evaluations of frame-level downstream tasks are conducted on the sound event detection (SED) task. SED is a frame-level multi-class classification task, which requires the model to recognize the sound events as well as their corresponding timestamps from the given audio clips. Two datasets are used for evaluation: domestic environment sound event detection (DESED) [49] and strongly-labeled AudioSet [50].

- **DESED** dataset is provided by the Detection and Classification of Acoustic Scenes and Events (DCASE) Challenge 2022, task 4 - Sound Event Detection in Domestic Environment. The DESED dataset provides both labeled and unlabeled recordings for training. Since our goal is to evaluate the SSL methods, we only utilize the labeled audio clips in our experiments for linear evaluation and fine-tuning of the pre-trained models. Specifically, 1,476 clip-level labeled (weakly-labeled) real audio clips and 12,500 frame-level labeled (strongly-labeled) synthetic audio clips are used for training and validation. As for evaluation, the DCASE task 4 development set is used, containing 1,168 frame-level labeled real audio clips. It is worth mentioning that, in DCASE task 4, there are a large amount of unlabelled data used for semi-supervised training, which however are not used in this work.
- **Strongly-labeled AudioSet** is a subset of the AS-2M dataset. It contains 102,561 and 15,958 frame-level labeled real audio clips for training and evaluation, respectively. There are 407 audio classes in total. The audio classes are seriously unbalanced, where the most frequently appeared 10 event classes occupy 50.7% of the total amount of events.

2) *Metric*: The evaluation on the DESED dataset is accomplished by the official metrics of DCASE Challenge 2022, i.e. the polyphonic sound event detection scores (PSDS) [51]. This metric measures the intersection between truth events and detected events. Two different sets of PSDS parameters are used, denoted as $PSDS_1$ and $PSDS_2$, to emphasize the low reaction time (accurate localization of sound event) and the low confusion rate between classes, respectively. For both metrics, the higher the better.

The evaluation methods used in the original work of strongly-labeled AudioSet [50] have a coarse temporal resolution of 960 ms. We think they are not fine-grained enough to represent the detection accuracy, considering that the length of sound events could be as small as tens of milliseconds. In our experiments, the two PSDSs used for the DESED task are used

for the strongly-labeled AudioSet as well. The vanilla PSDS [51] includes an optional penalty term by the performance variance across all classes. Such penalty term evaluates the stability of performance across classes. However, the number of classes of strongly-labeled AudioSet is very high (407) and the classes are heavily unbalanced, such that all the test models in our experiments have a high performance variance across classes. With the penalty term, PSDS scores could be reduced to 0 for many models, which cannot conduct fair comparison. Therefore, we will report the scores without applying the penalty term as well.

3) *Frame-level Embedding Extraction*: For this frame-level downstream task, both ATST-Clip and ATST-Frame only use the frame-level embedding sequence. In both linear evaluation and fine-tuning experiments, the embedding sequence of the last encoder block is used.

The length of all the audio clips in both DCASE and strongly-labeled AudioSet datasets is 10 seconds, which is exactly the same as the length of the positional embedding adopted by ATST-Frame, therefore, ATST-Frame processes these audio clips without splitting. For ATST-Clip, to account for the length of positional embedding, the audio clips are splitted into two chunks and the representations of the two chunks are concatenated along the time dimension.

4) *Downstream Task Training*: On top of the frame-level embedding, a multi-class classifier is added. For the DESED dataset, to account for the weak labels, the frame-level detection results are pooled with a softmax attention linear classifier, following the principle of DCASE challenge baseline method [52]. Note that, our setup is different from the one of the DCASE baseline, as the latter uses some extra networks besides the pre-trained model and uses some unlabelled data. As for the strongly-labeled AudioSet, a simple linear classifier is cascaded behind the frame-level embeddings to generate detection results.

In fine-tuning experiments, considering the limited data size of the DESED dataset, we only unfreeze the last encoder block to avoid over-fitting. The batch size is set to [128, 128] for weakly- and strongly-labeled samples, respectively. For the strongly-labeled AudioSet, we unfreeze the entire model, where the batch size is set to 256. The models are trained with the SGD optimizer for 100 epochs. The learning rate is warmed up for 5 epochs.

V. RESULTS

A. Ablation Study

Ablation experiments are conducted using ATST-Clip_{small} and ATST-Frame_{small} with the linear evaluation protocol, due to their low computational cost.

1) *Ablations on ATST-Clip*: We separately evaluate the effectiveness of the Transformer encoder and the proposed view creation strategy. Table III shows the results. The result of BYOL-A is also given, which uses a CNN encoder and a single 1-second segment. Our models use single or two segments, with a length of 1 second or 6 seconds. For a fair comparison, when the segment length is set to 1 second, we split audio clips into 1-second chunks for downstream tasks.

Method	Segments	length of segment (s)	AS-20K mAP (%)	SPCV2 Acc (%)	VOX1 Acc (%)	NSYNTH Acc (%)	US8K Acc (%)	Average Acc (%)
BYOL-A [5]	single	1	-	92.2	40.1	74.1	79.1	71.4
ATST-Clip _{small}	single	1	21.0	94.3	52.3	73.8	79.3	74.9
	two	1	19.1	91.3	50.0	74.3	76.6	73.1
	single	6	25.7	94.0	57.3	73.8	80.9	76.5
	two	6	27.9	93.6	61.9	75.3	82.0	78.2

TABLE III: Ablation studies on ATST-Clip_{small}. Linear evaluation results are shown. "Average" is taken over the last four tasks.

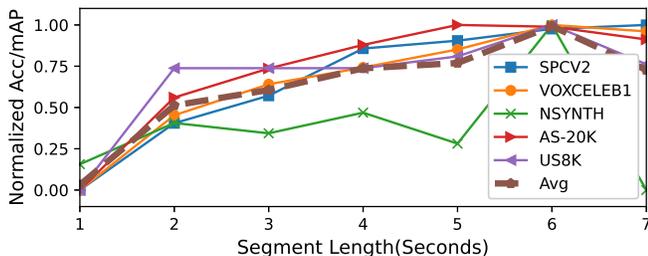


Fig. 2: Acc/mAP of ATST-Clip_{small} as a function of segment length, Acc/mAP of each task is normalized into the range of [0,1]. "Avg" denotes averaging over all tasks.

Transformer Encoder: With the same view creation strategy, i.e. creating two views from a 1-second segment, our model (line 2 in Table III) outperforms BYOL-A, especially for the two speech tasks (SPCV2 and VOX1). Speech involves more long-term semantic information, and Transformer is more suitable than CNN for learning these long-term dependencies.

View Creation Strategy: As shown in Table III, when the segment length is set to 1 second, using one single segment is better than using two segments. This phenomenon is consistent with the claim made in BYOL-A [5] that the two segments may be too different to be identified as the same sample. However, two views created from a single segment may share too much semantic content, thus leading our model to find an easy solution. When the segment length is increased to 6 seconds, the performance measures of AS-20K, VOX1 and US8K are systematically increased, no matter whether using one or two segments. This is partially due to the capability of learning long-term dependencies of the Transformer encoder. In addition, for the 6-second case, using two segments exhibits superior performance over using one segment. The possible reasons are: the two segments can be rationally identified as the same sample as they share a small portion of overlap, and meanwhile they are different enough to increase the task difficulty and thus leads the model to learn a more generalized representation.

Fig. 2 shows the normalized performance of each task as a function of segment length, where two segments are used. We can observe that as the segment length increases, the performance metrics continue to improve until they reach a maximum at 6 seconds. This further verifies our findings: i) when Transformer encoder is used, increasing the segment length helps to learn more information; ii) when two segments

are used, the segment length should be set to make the segments share a proper amount of overlap, and have a proper difficulty for matching them as the same sample.

2) *Ablations on ATST-Frame:* Table IV shows the ablation results on the effectiveness of ATST-Frame components.

Data augmentation: Compared with the no augmentation case (configuration A in Table IV), using data augmentation (configuration B, C, D, E) brings a significant performance improvement on all tasks. This means data augmentation is able to properly increase the task difficulty and to encourage the model to learn more meaningful audio representations. Augmenting two views (for both teacher and student branches, configuration E) leads to a large task difficulty, and achieves better results on three out of five tasks. Only augmenting one view (for either student or teacher branch, configuration D) achieves slightly worse performance than augmenting two views on AS-20K and SPCV2, but much better performance on VOX1, thus has a better average result. This is consistent with our observations in the ablation studies of ATST-Clip that better performance can be achieved with a balanced difficulty of the pre-training task.

Masking: In configuration C, both the student and teacher branches are masked with the same time index, thus the two branches need to predict the masked frames from unmasked frames, and the predictions should be matched. In configuration D, only the student branch is masked, while the teacher branch sees the whole audio clip. We can see that masking both branches performs worse than masking only the student branch. The possible reasons are i) the teacher branch provides more meaningful guidance for the student branch when seeing the frames that are not visible to the student branch; ii) the teacher encoder consistently sees unmasked input for pre-training and downstream tasks. As for the masking strategy, group masking that forces N adjacent frames to be masked together performs better than random masking (Configuration B). This is consistent with the observations in the speech pre-training works [24], [27].

Based on the above analysis, the proposed ATST-Frame is set up with configuration D. Unless noted, the following experiments of ATST-Frame use configuration D by default.

The symmetrical loss: In ATST-Frame, augmentation is applied to one of the two views. As the symmetrical loss is used, both the teacher branch and the student branch see the augmented view during training. We conduct experiments by using only L_θ or L'_θ to evaluate which one of teacher and student branches is more important to see the augmented view. The results are shown in Table V, which shows that it is

Configuration	Augmented views	Mask teacher	Mask student	Mask strategy	Downstream Tasks					Average
					AS-20K mAP (%)	SPCV2 Acc (%)	VOX1 Acc (%)	NSYNTH Acc (%)	US8K Acc (%)	
A	0		✓	Group	8.0	63.3	25.8	56.8	60.1	42.8
B	1		✓	Random	18.0	85.2	43.7	69.9	76.1	58.6
C	1	✓	✓	Group	22.5	88.2	47.0	72.9	72.9	60.7
D	1		✓	Group	28.1	92.3	67.0	72.5	84.0	68.8
E	2		✓	Group	28.5	92.5	59.6	74.7	83.4	67.7

TABLE IV: Ablation studies on ATST-Frame_{small}. Linear evaluation results are shown.

Method	Symmetrical	Loss	Augmented branch	AS-20K mAP (%)	SPCV2 Acc (%)	VOX1 Acc (%)	NSYNTH Acc (%)	US8K Acc (%)	Average
ATST-Frame _{small}	True	$L_\theta + L'_\theta$	Teacher & Student	28.1	92.3	67.0	72.5	84.0	68.8
	False	L'_θ	Teacher	6.4	77.8	18.0	63.0	67.1	46.7
	False	L_θ	Student	22.5	87.2	50.2	68.7	80.1	61.7

TABLE V: Ablation studies on the symmetrical loss of ATST-Frame_{small}. "Augmented branch" denotes the branch taking as input the augmented view. Linear evaluation results are shown.

Method	Strategy	Frequency warping	AS-20K mAP (%)	SPCV2 Acc (%)	VOX1 Acc (%)	NSYNTH Acc (%)	US8K Acc (%)	Average
ATST-Frame _{small}	Frame-wise	True	28.1	92.3	67.0	72.5	84.0	68.8
	Frame-wise	False	23.5	89.4	56.9	69.0	79.7	63.7
	Patch-wise	True	16.3	77.8	35.9	71.6	75.2	55.4
	Patch-wise	False	23.3	80.8	48.6	70.2	79.5	60.5

TABLE VI: Ablation studies on comparison of frame-wise and patch-wise strategy. Linear evaluation results are shown.

Method	Symmetrical	Augmented branch	AS-20K mAP (%)	SPCV2 Acc (%)	VOX1 Acc (%)	NSYNTH Acc (%)	US8K Acc (%)	Average
ATST-Frame _{small}	True	Teacher & Student	28.1	92.3	67.0	72.5	84.0	68.8
ATST-Frame-data2vec* _{small}	False	None	24.1	92.0	58.4	73.0	81.4	65.8
	False	Student	18.6	89.3	43.7	71.5	76.0	59.8
	True	Teacher & Student	21.6	90.7	52.3	70	78.4	62.6

TABLE VII: Ablation studies on comparison with data2vec-style training target. "ATST-Frame-data2vec*_{small}" denotes ATST-Frame_{small} with data2vec-style [27] training target. "Augmented branch" denotes the branch takes as input the augmented view. Linear evaluation results are shown.

more important for the student branch than the teacher branch to see the augmented view, and using the symmetrical loss outperforms the case that only one branch sees the augmented view.

Patch-wise strategy and frequency warping: ATST-Frame uses frame-wise strategy for log-mel spectrogram, while other works [6], [7] have reported that patch-wise strategy exhibits better performance than frame-wise strategy on sound event/scene classification task, as sound events/scenes have complex frequency structure, which can be better captured by the frequency split of patch-wise models [6]. To testify the frame-wise strategy of our ATST-Frame model, we further conduct experiments using patch-wise strategy in the framework of ATST-Frame. Specifically, we organize the log-mel spectrogram into patches in the size of 16 frequency bins \times 16 frames, which leads to the same number of tokens as ATST-Frame for a 64-bin log-mel spectrogram. For patch-wise models, frequency warping conflicts with the principle of the patch-wise loss, as it distorts the patch correspondences. Therefore, we conducted experiments both with or without using frequency warping. Note that Mixup is always used. The

results are shown in VI. Without using frequency warping, the frame-wise model noticeably performs better than the patch-wise model on speech tasks (SPCV2 and VOX1), which is consistent with the observations in other works [6], [7]. However, we do not observe the advantage of patch-wise strategy on sound event/scene classification (AS-20K and US8K), where patch-wise strategy and frame-wise strategy are comparable. The frame-wise model significantly benefits from frequency warping for all tasks whereas the patch-wise model does not. Frequency warping (FW) encourages learning FW-invariant representations, which may help to learn the spectral pattern, even for the complex spectral structure of sound events/scenes. Overall, within the framework of ATST-Frame, the frame-wise strategy is suitable for both speech and sound events/scenes, and frequency warping helps to largely improve the performance.

Using the training target of data2vec: We apply the training target of data2vec [27] to our ATST-Frame model (referred to as ATST-Frame-data2vec). Specifically, the last 8 blocks of the teacher encoder are averaged to form the training target; the projectors are removed; the predictor is replaced

Method	AS-20K mAP (%)	SPCV2 Acc (%)	VOX1 Acc (%)	NSYNTH Acc (%)	US8K Acc (%)	FSD50K mAP (%)
TRILL [53]	-	-	17.9	-	-	-
COLA [3]	-	62.4	29.9	63.4	-	-
BYOL-A [5]	-	92.2	40.1	74.1	79.1	-
BYOL-A-V2 [12]	-	93.1	57.6	73.1	79.7	44.8
SF NFNNet-F0 [54]	-	93.0	64.9	78.2	-	-
M2D [16]	-	95.4	73.1	76.9	87.6	-
ATST-Clip (ours)	33.8	95.1	72.0	76.2	85.8	58.5
ATST-Frame (ours)	33.0	94.9	77.4	75.9	85.8	55.1

TABLE VIII: Linear evaluation results on clip-level downstream tasks. The scores of comparison models are quoted from their papers.

with a linear projection. Although the original data2vec does not use data augmentation and symmetrical loss, we also test how will data augmentation and symmetrical loss perform when used to ATST-Frame-data2vec. The results are reported in Table VII. It can be seen that the data2vec target performs well, but it does not benefit from data augmentation.

B. Results on Clip-level Downstream Tasks

1) *Linear Evaluation Results:* Table VIII shows the linear evaluation results on six tasks. For a fair comparison, we compare with other methods that also use Audioset for pre-training and have also reported the linear evaluation results in their papers, including TRILL [53], COLA [3], BYOL-A [5], BYOL-A-v2 [12], SF NFNNet-F0 [54] and M2D [16]. The proposed ATST-Clip is developed based on BYOL-A and BYOL-A-V2, using a Transformer encoder and a new view creation strategy. It can be seen that ATST-Clip noticeably outperforms BYOL-A and BYOL-A-V2 on all tasks, which indicates that our modifications are very effective. On average, the proposed models and the recently proposed M2D model perform better than other models. The performance of the proposed models and M2D are comparable, as M2D performs better on SPCV2, NSYNTH and US8K with small advantages, while the proposed ATST-Frame performs better on VOX1. Among the two proposed models, ATST-Clip outperforms ATST-Frame for all the tasks except for VOX1. ATST-Clip is dedicated to learning clip-level representation, its embedding is more representative for the audio clip than the one obtained by averaging the frame-level embeddings of ATST-Frame. However, the drawbacks of ATST-Frame are not significant, which means the average of its frame-level embeddings is still an effective clip-level representation.

2) *Fine-tuning Results:* Linear evaluation cannot fully reflect the capabilities of pre-trained models, as normally the models can be further fine-tuned with the data of downstream tasks. Fine-tuning experiments are conducted on the tasks of multi-label audio event classification (AS-2M, AS-20K and FSD50K), Spoken command recognition (SPCV2), speaker identification (VOX1) and musical instrument classification (NSYNTH). We compare with two groups of prior methods: supervised methods and self-supervised methods. The results are shown in Table IX.

ATST-Frame outperforms ATST-Clip. After fine-tuning, the performance of ATST-Frame is better than ATST-Clip on five out of six tasks. As mentioned above, ATST-Frame

does not explicitly learn clip-level representation during pre-training. However, fine-tuning allows the adjustment of the pre-trained parameters to fit a specific downstream task. ATST-Frame is pre-trained by maximizing the agreement of frame-level embeddings, which is more fine-grained and challenging compared with ATST-Clip. This may help ATST-Frame to learn more sophisticated knowledge and network parameters (such as the self-attention parameters), which happen to be a better initial setting for fine-tuning even on clip-level downstream tasks.

Comparison with supervised methods. The proposed ATST-Frame outperforms the supervised methods on AS-2M, AS-20K and SPCV2. The proposed ATST-C2F model further improves the performance, and outperforms the supervised methods on all tasks. This is encouraging for the field of audio self-supervised learning, as we no longer need to annotate audio data for pre-training when we want to further scale up the dataset. Compared with supervised pre-training, self-supervised pre-training does not suffer from the problem of inaccurate and erroneous labels.

Comparison with other self-supervised methods. Compared with other self-supervised methods, the proposed ATST-Frame achieves comparable or better performance on all tasks. In particular, compared with the recent state-of-the-art self-supervised method BEATs_{iter3+} [10], ATST-Frame achieves the same performance on AS-2M, and better performance on AS-20K. This indicates that ATST-Frame is more effective with less fine-tuning data than BEATs.

Combination through knowledge distillation. ATST-Clip and ATST-Frame learn complementary features in the pre-training stage. Combining ATST-Clip and ATST-Frame through knowledge distillation can further improve the performance, as shown by the results of ATST-C2F in Table IX. Even though ATST-Clip performs worse than ATST-Frame on AS-2M and AS-20K, as a teacher, it still successfully helps to fine-tune ATST-Frame to achieve better performance. Performing knowledge distillation the other way around, i.e. from ATST-Frame to ATST-Clip, does not perform well on most of the tasks. This verifies that ATST-Frame conveys more fine-grained and semantically complicated information than ATST-Clip, and ATST-Frame should be used as the final model. BEATs_{iter3+} [10] also performs knowledge distillation across models at the fine-tuning stage, specifically, it uses the fine-tuned BEATs_{iter2} model as a teacher to fine-tune the final BEATs_{iter3} model. Thence, BEATs_{iter3+} can be regarded as a fair comparison with the proposed ATST-C2F model.

C. Results on Frame-level Downstream Task - DESED

1) *Comparison Methods:* We compare with six SSL pre-trained models: BYOL-A-v2 [12], SSAST [6], MAE-AST [7], Audio-MAE [9], BEATs [10] and M2D [16]. Sound event detection requires to perform frame-level multi-class classification. As mentioned in Section IV-C, the proposed models can be directly used for this task by adding a linear classifier on top of their frame-level representations, with a temporal resolution of 40 ms per frame. The comparison models are pre-trained either frame-wisely or patch-wisely. The frame-wise models, e.g. SSAST and MAE-AST, can also be directly

Method	# Param	Pre-training data	AS-2M mAP (%)	AS-20K mAP (%)	SPCV2 Acc (%)	VOX1 Acc (%)	FSD50K mAP (%)	NSYNTH Acc (%)
Supervised Methods								
PANN [55]	81M		43.9	27.8	-	-	-	-
PSLA [56]	14M		44.4	31.9	-	-	55.4	-
AST [28]	86M		45.9	34.7	98.1	-	-	-
HTS-AT [57]	31M		47.1	-	98.0	-	-	-
PassT [58]	86M		47.1	-	-	-	65.3	-
KD-AST [39]	86M		47.1	-	-	-	62.9	-
Self-supervised Methods								
SSAST-PATCH [6]	89M	AS+LS	-	31.0	98.0	64.2	-	-
SSAST-FRAME [28]	89M	AS+LS	-	29.2	98.1	80.8	-	-
Conformer-Based [8]	88M	67K hours *	41.5	27.6	-	-	-	-
MAE-AST-PATCH [7]	86M	AS+LS	-	30.6	97.9	-	-	-
MAE-AST-FRAME [7]	86M	AS+LS	-	23.0	98.0	63.3	-	-
ASiT [37]	85M	AS	-	35.2	98.8	63.1	-	-
data2vec [27]	94M	AS	-	34.5	-	-	-	-
MaskSpec [29]	86M	AS	47.1	34.7	97.6	-	-	-
MSM-MAE [15] †	86M	AS	-	36.7	98.4	95.3	-	-
Audio-MAE (local) [9]	86M	AS	47.3	37.0	98.3	94.8	-	-
BEAT _{iter3} [10]	90M	AS	48.0	38.3	98.3	-	-	-
BEAT _{iter3+} [10] **	90M	AS	48.6	38.9	98.1	-	-	-
M2D [16]	86M	AS	-	37.4	98.5	94.4	-	-
Ours								
ATST-Clip	86M	AS	45.2	37.9	98.0	95.5	63.4	78.6
ATST-Frame	86M	AS	48.0	39.0	98.1	97.3	61.8	79.2
ATST-C2F **	86M	AS	49.7	40.5	98.4	97.5	65.5	79.2
ATST-F2C **	86M	AS	46.8	39.0	98.1	95.5	64.6	79.8

* Self-hold dataset [8].

† Results are quoted from M2D [16].

** Perform knowledge distillation across two models at the finetuning stage.

TABLE IX: Finetuning results on clip-level downstream tasks. The scores of comparison models are quoted from their papers. AS and LS denote AudioSet and Librispeech [59], respectively.

used for this task, with a temporal resolution of 20 ms per frame. According to the setup of the proposed models, we also evaluate SSAST and MAE-AST with a temporal resolution of 40 ms per frame, by applying average pooling to the frame-level representations. As for the patch-wise models, we average the patch-level representations for each time interval to obtain the interval/frame-level representations, except for M2D, since its authors propose to concatenate instead of average the patch-level representations [16]. Note that, the patch-wise models have a coarser temporal resolution, i.e. 160 ms per frame.

All the pre-trained models are fine-tuned by ourselves using the SED supervised dataset. For linear evaluation, pre-trained models are frozen, and only the two dense layers of the linear classifier are trained. In fine-tuning experiments, for all the Transformer-based models, we unfreeze the last Transformer block. For BYOL-A-v2, we unfreeze the entire model for fair comparison such that the amount of the trainable parameters of each model are similar. The best learning rate for each model has been carefully searched, which is also given in Table X.

2) *Results Analysis*: Table X shows the results. These results are deviated from the results reported in the DCASE challenge, that is because of the different experimental setups as discussed in Sec. IV-C4. The objective of this study is to conduct fair comparison between different SSL models, instead of pursuing the SOTA performances on this dataset. It can be seen that, as expected, relative to linear evaluation, the performance of all models can be improved by fine-tuning with

the SED dataset. As the fine-tuning setup is more practically important than linear evaluation, we mainly analyze the fine-tuning results in the following, and most of the analyses are valid for the linear evaluation results as well.

BYOL-A-V2 does not achieve reasonable performance, possibly due to its limited capacity for sequential processing with a two-layer CNN architecture. For SSAST [6] and MAE-AST [7], increasing the temporal resolution of their frame-wise models from 20 ms to 40 ms largely improves the performance. The possible reasons are that the frame-level representations get more stable when averaging two frames, and meanwhile, the 40 ms temporal resolution is still fine enough for tracking the time variation of sound events. This could also be because the characteristics of one event cannot be well represented without sufficiently long frames. However, the performances of their frame-wise models still largely lag behind their patch-wise models. This is consistent with the observations in [6], [7] that, the patch-wise models are more suitable for sound events, while the frame-wise models are more suitable for speech signals. Sound events have more complex frequency structure, which can be better captured by the frequency split of patch-wise models. Among the comparison models, BEATs performs the best in terms of both PSDS₁ and PSDS₂, and MAE-AST-PATCH achieves close performance with BEATs.

The proposed ATST-Clip does not work well, as it is trained for learning global representation, which does not automatically lead to good frame-level representations. By

Method	τ (ms)	learning rate	PSDS ₁	PSDS ₂
Linear Evaluation				
BYOL-A-V2-40ms [12]	40	0.01	0.024	0.181
SSAST-FRAME-20ms [6]	20	0.1	0.028	0.166
SSAST-FRAME-40ms [6]	40	0.1	0.096	0.266
SSAST-PATCH [6]	160	0.1	0.179	0.315
MAE-AST-FRAME-20ms [7]	20	0.1	0.031	0.234
MAE-AST-FRAME-40ms [7]	40	0.1	0.081	0.293
MAE-AST-PATCH [7]	160	0.1	0.225	0.442
Audio-MAE (local) [9]	160	0.1	0.218	0.401
BEAT _{iter3} [10]	160	0.1	0.177	0.358
M2D [16]	160	0.1	0.234	0.438
Ours				
ATST-Clip	40	0.1	0.115	0.293
ATST-Frame	40	0.1	0.304	0.507
Finetuning				
BYOL-A-V2-40ms [12]	40	0.01	0.030	0.219
SSAST-FRAME-20ms [6]	20	0.05	0.046	0.235
SSAST-FRAME-40ms [6]	40	0.1	0.132	0.325
SSAST-PATCH [6]	160	0.1	0.236	0.459
MAE-AST-FRAME-20ms [7]	20	0.05	0.123	0.346
MAE-AST-FRAME-40ms [7]	40	0.1	0.235	0.418
MAE-AST-PATCH [7]	160	0.05	0.281	0.573
Audio-MAE (local) [9]	160	0.1	0.254	0.509
BEAT _{iter3} [10]	160	0.1	0.282	0.584
M2D [16]	160	0.1	0.267	0.500
Ours				
ATST-Clip	40	0.05	0.223	0.422
ATST-Frame	40	0.01	0.361	0.581
ATST-C2F	40	0.1	0.357	0.607
ATST-F2C	40	0.1	0.259	0.445

TABLE X: Results on the frame-level downstream task, DESED. τ stands for temporal resolution.

leveraging the proposed frame-level training criterion and thus learning better frame-level representations, ATST-Frame largely improves the performance over ATST-Clip. Compared with the best comparison model, i.e. BEATs, ATST-Frame achieves much better PSDS₁, and similar PSDS₂. PSDS₁ emphasizes the time localization accuracy of sound events, thence the better PSDS₁ of ATST-Frame means a better temporal detection performance, which is possibly due to the finer temporal resolution of ATST-Frame compared with BEATs, i.e. 40 ms versus 160 ms. PSDS₂ emphasizes the recognition accuracy of sound events. The similar PSDS₂ of ATST-Frame and BEATs reflect the similar representation quality of them. This is consistent with the results on the clip-level AS-2M task, ATST-Frame also performs similarly with BEATs as shown in Table VIII. It is important to note that, the good performance of ATST-Frame conflicts with the observations in [6], [7] that patch-wise models are more suitable for sound events than frame-wise models. As discussed in the ablation study, the success of ATST-Frame is possibly attributed to the frequency warping operation, which helps to capture the complex frequency structure of sound events.

Knowledge distillation is also applied to combine ATST-Clip and ATST-Frame. The results of ATST-C2F show that, taking ATST-Clip as a teacher for fine-tuning ATST-Frame, PSDS₂ can be further improved, while PSDS₁ is slightly decreased. This means the knowledge learned by ATST-Clip is still complementary for improving the accuracy of frame-level representations, but will slightly blur the time localization.

Method	Learning rate	PSDS ₁		PSDS ₂	
		w/o var-pen	with var-pen	w/o var-pen	with var-pen
Linear Evaluation					
BYOL-A-V2-40ms [12]	0.5	0.087	0.0	0.083	0.0
SSAST-PATCH [6]	0.5	0.048	0.0	0.067	0.0
MAE-AST-PATCH [7]	0.5	0.116	0.0	0.185	0.0
Audio-MAE (local) [9]	0.5	0.073	0.0	0.107	0.0
BEAT _{iter3} [10]	0.5	0.034	0.0	0.062	0.0
M2D [16]	0.5	0.182	0.0	0.301	0.039
Ours					
ATST-Clip	0.5	0.120	0.0	0.201	0.001
ATST-Frame	0.5	0.207	0.008	0.304	0.048
Finetuning					
BYOL-A-V2-40ms [12]	0.5	0.110	0.0	0.243	0.027
SSAST-PATCH [6]	0.5	0.243	0.017	0.411	0.122
MAE-AST-PATCH [7]	0.5	0.274	0.039	0.481	0.187
Audio-MAE (local) [9]	0.5	0.276	0.038	0.476	0.182
BEAT _{iter3} [10]	0.5	0.290	0.045	0.491	0.186
M2D [16]	0.1	0.292	0.042	0.509	0.199
Ours					
ATST-Clip	0.5	0.328	0.083	0.478	0.178
ATST-Frame	0.5	0.347	0.069	0.538	0.152
ATST-C2F	0.5	0.374	0.125	0.572	0.266
ATST-F2C	0.5	0.323	0.075	0.470	0.163

TABLE XI: Results on the frame-level downstream task, SED of the strongly-labeled AudioSet. ‘var-pen’ stands for the performance variance penalty term.

D. Results on Frame-level Downstream Task - Strongly-labeled AudioSet

Table XI shows the results on the strongly-labeled AudioSet. According to the DESED performances, only the best-performing model for each comparison method is evaluated. Considering the large data size of strongly-labeled AudioSet, we only search over 3 different learning rates for each model, i.e. 0.05, 0.1 and 0.5. As mentioned in Sec. IV-C2, we evaluate the models by the PSDSs with or without applying the performance variance penalty term.

For all the models, the scores with variance penalty are much lower than the ones without variance penalty, which means the performance variance across classes for all models are very large. The scores with variance penalty for linear evaluation could be reduced to 0 for most of the models. This reflects the data imbalance and task difficulty of the strongly-labeled AudioSet.

After finetuning, The BYOL-A-v2 model has a large performance gap comparing with other Transformer-based models. With better learned frame-level representations, the proposed ATST-Frame model has an obvious advantage over the comparison models and ATST-Clip, when variance penalty is not applied. However, ATST-Clip has a better stability of performance across classes, and thus outperforms ATST-Frame when applying variance penalty. When combining ATST-Frame and ATST-Clip, the performance measures are largely improved by ATST-C2F, and the model is improved in terms of both classification accuracy and performance stability.

	ATST-Clip	ATST-Frame	Best
Beehive	58.3	64.6	87.8
Beijing Opera	95.3	95.8	97.5
CREMA-D	<u>76.0</u>	<u>76.7</u>	75.2
DCASE 2016 *	<u>93.7</u>	<u>95.7</u>	92.5
ESC-50	91.2	89.0	96.1
FSD50K	59.5	55.7	64.1
Gunshot	<u>98.8</u>	94.3	96.7
GTZAN Genre	87.7	88.3	90.8
GTZAN Music/Speech	99.2	<u>100.0</u>	99.2
Libricount1	78.2	78.1	78.5
Maestro 5h *	18.9	24.4	46.9
Mridangam Stroke	<u>97.7</u>	97.5	97.5
Mridangam Tonic	<u>96.7</u>	<u>96.9</u>	94.1
NSynth Pitch 5h	67.8	68.6	87.8
Speech command 5h	93.1	92.6	97.6
Speech command full	95.5	95.1	97.8
Vocal Imitation	18.5	<u>22.3</u>	21.5
VoxLingua107 top 10	53.9	66.9	72.2

* frame-level task

TABLE XII: Results on the HEAR benchmark. ‘Best’ denotes the best result in the HEAR leaderboard. Underlined scores denote better performance than ‘Best’.

E. Results on HEAR benchmark

We also evaluate the proposed models on the HEAR benchmark [60], which includes 17 clip-level and 2 frame-level tasks. We successfully downloaded 18 tasks. The HEAR benchmark trains a shallow MLP classifier on top of frozen embeddings. We use the official hear-eval-kit¹, and our embeddings are extracted in the same way as we did in our linear evaluation experiments except that for frame-level tasks, we concatenate outputs of all the blocks. Table XII shows the results. As a baseline, we quote the best result for each task from the HEAR leaderboard², denoted as ‘Best’ in the table. It is worth noting that there are two frame-level tasks, i.e. DCASE 2016 and Maestro 5h. On DCASE 2016, both ATST-Clip and ATST-Frame perform better than the best baseline. However, the best baseline performs much better than the proposed models for Maestro 5h. Maestro 5h is a piano music transcription task, aiming to extract pitch and onset from raw audio. The data augmentation of RRC and frequency warping in ATST encourage the model to learn frequency-changing-invariant representations, which may be not suitable for pitch learning, as pitch is sensitive to frequency changing. This phenomenon is also observed on the clip-level pitch estimation task, i.e. NSynth Pitch 5h. Overall, both the proposed ATST-Clip and ATST-Frame achieve better performance than the best baseline on five tasks. This is remarkable considering the fact that the best baseline results quoted here for different tasks are achieved by 14 different submissions. Moreover, some of the best baseline results are obtained by the model especially trained for the specific tasks, as HEAR benchmark does not limit the pre-training methods (supervised or unsupervised) and pre-training datasets.

¹<https://github.com/hearbenchmark/hear-eval-kit>

²<https://hearbenchmark.com/hear-leaderboard.html>

VI. CONCLUSION

In this paper, based on the teacher-student scheme of BYOL, we have proposed two effective self-supervised audio pre-training methods, ATST-Clip and ATST-Frame, specifically crafted to learn clip-level and frame-level representations, respectively, enabling effective audio understanding.

The proposed methods have been extensively evaluated on a variety of downstream tasks, including seven clip-level tasks and two frame-level tasks, covering multiple audio domains: environmental sound, speech and music. Both ATST-Clip and ATST-Frame demonstrated their outstanding capabilities of learning audio representations compared with previous state-of-the-art methods. Furthermore, ATST-Clip offers complementary knowledge to ATST-Frame, and these knowledge can be effectively distilled to ATST-Frame at the fine-tuning stage. Especially, the proposed methods achieve new SOTA scores on the AudioSet-2M and AudioSet-20K datasets, with the precision of 49.7% and 40.5% (without model ensembling), respectively. Furthermore, this work also provides a new benchmark for applying pre-trained models to frame-level downstream tasks, on two sound event detection datasets. The frame-level downstream tasks have rarely been studied in the field, and hopefully this work would fill this gap.

We have open-sourced our code online for the research community to replicate and expedite future research. As the scope of this study is limited to audio classification tasks, future work may extend our models to audio generation tasks.

REFERENCES

- [1] A. v. d. Oord, Y. Li, and O. Vinyals, “Representation Learning with Contrastive Predictive Coding,” *arXiv:1807.03748 [cs, stat]*, 2019.
- [2] 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.
- [3] A. Saeed, D. Grangier, and N. Zeghidour, “Contrastive Learning of General-Purpose Audio Representations,” *arXiv:2010.10915 [cs, eess]*, 2020.
- [4] E. Fonseca, D. Ortego, K. McGuinness, N. E. O’Connor, and X. Serra, “Unsupervised Contrastive Learning of Sound Event Representations,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 371–375.
- [5] D. Niizumi, D. Takeuchi, Y. Ohishi, N. Harada, and K. Kashino, “BYOL for Audio: Self-Supervised Learning for General-Purpose Audio Representation,” *arXiv:2103.06695 [cs, eess]*, 2021.
- [6] Y. Gong, C.-I. J. Lai, Y.-A. Chung, and J. Glass, “SSAST: Self-Supervised Audio Spectrogram Transformer,” *arXiv:2110.09784 [cs, eess]*, 2022.
- [7] A. Baade, P. Peng, and D. Harwath, “MAE-AST: Masked Autoencoding Audio Spectrogram Transformer,” Mar. 2022.
- [8] S. Srivastava, Y. Wang, A. Tjandra, A. Kumar, C. Liu, K. Singh, and Y. Saraf, “Conformer-Based Self-Supervised Learning for Non-Speech Audio Tasks,” *arXiv:2110.07313 [cs, eess]*, 2022.
- [9] P.-Y. Huang, H. Xu, J. Li, A. Baevski, M. Auli, W. Galuba, F. Metze, and C. Feichtenhofer, “Masked Autoencoders that Listen,” Jan. 2023.
- [10] S. Chen, Y. Wu, C. Wang, S. Liu, D. Tompkins, Z. Chen, and F. Wei, “BEATs: Audio Pre-Training with Acoustic Tokenizers,” Dec. 2022.
- [11] S. Liu, A. Mallol-Ragolta, E. Parada-Cabeleiro, K. Qian, X. Jing, A. Kathan, B. Hu, and B. W. Schuller, “Audio Self-supervised Learning: A Survey,” Mar. 2022, arXiv:2203.01205 [cs, eess].
- [12] D. Niizumi, D. Takeuchi, Y. Ohishi, N. Harada, and K. Kashino, “BYOL for Audio: Exploring Pre-Trained General-Purpose Audio Representations,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 31, pp. 137–151, 2023.
- [13] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding,” *arXiv:1810.04805 [cs]*, 2019.

- [14] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention Is All You Need," *arXiv:1706.03762 [cs]*, 2017.
- [15] D. Niizumi, D. Takeuchi, Y. Ohishi, N. Harada, and K. Kashino, "Masked Spectrogram Modeling using Masked Autoencoders for Learning General-purpose Audio Representation."
- [16] —, "Masked modeling duo: Learning representations by encouraging both networks to model the input," in *ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2023, pp. 1–5.
- [17] J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. H. Richemond, E. Buchatskaya, C. Doersch, B. A. Pires, Z. D. Guo, M. G. Azar, B. Piot, K. Kavukcuoglu, R. Munos, and M. Valko, "Bootstrap your own latent: A new approach to self-supervised Learning," *arXiv:2006.07733 [cs, stat]*, 2020.
- [18] X. LI and X. Li, "ATST: Audio Representation Learning with Teacher-Student Transformer," in *Proc. Interspeech 2022*, 2022, pp. 4172–4176.
- [19] 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 *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2017, pp. 776–780.
- [20] X. Chen and K. He, "Exploring Simple Siamese Representation Learning," *arXiv:2011.10566 [cs]*, 2020.
- [21] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A Simple Framework for Contrastive Learning of Visual Representations," *arXiv:2002.05709 [cs, stat]*, 2020.
- [22] E. Fonseca, D. Ortego, K. McGuinness, N. E. O'Connor, and X. Serra, "Unsupervised contrastive learning of sound event representations," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 371–375.
- [23] A. T. Liu, S.-w. Yang, P.-H. Chi, P.-c. Hsu, and H.-y. Lee, "Mockingjay: Unsupervised Speech Representation Learning with Deep Bidirectional Transformer Encoders," *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6419–6423, 2020.
- [24] A. Baevski, H. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A Framework for Self-Supervised Learning of Speech Representations," *arXiv:2006.11477 [cs, eess]*, 2020.
- [25] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhota, R. Salakhutdinov, and A. Mohamed, "HuBERT: Self-Supervised Speech Representation Learning by Masked Prediction of Hidden Units," *arXiv:2106.07447 [cs, eess]*, 2021.
- [26] A. T. Liu, S.-W. Li, and H.-y. Lee, "Tera: Self-Supervised Learning of Transformer Encoder Representation for Speech," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 2351–2366, 2021.
- [27] A. Baevski, W.-N. Hsu, Q. Xu, A. Babu, J. Gu, and M. Auli, "data2vec: A General Framework for Self-supervised Learning in Speech, Vision and Language," Oct. 2022.
- [28] Y. Gong, Y.-A. Chung, and J. Glass, "AST: Audio Spectrogram Transformer," *arXiv:2104.01778 [cs]*, 2021.
- [29] D. Chong, H. Wang, P. Zhou, and Q. Zeng, "Masked spectrogram prediction for self-supervised audio pre-training," in *ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2023, pp. 1–5.
- [30] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick, "Masked autoencoders are scalable vision learners," *arXiv preprint arXiv:2111.06377*, 2021.
- [31] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," *CoRR*, vol. abs/2010.11929, 2020.
- [32] Q. Wang, B. Li, T. Xiao, J. Zhu, C. Li, D. F. Wong, and L. S. Chao, "Learning deep transformer models for machine translation," *CoRR*, vol. abs/1906.01787, 2019.
- [33] Y. Tokozume, Y. Ushiku, and T. Harada, "Learning from between-class examples for deep sound recognition," *CoRR*, vol. abs/1711.10282, 2017.
- [34] H. Zhang, M. Cissé, Y. N. Dauphin, and D. Lopez-Paz, "mixup: Beyond empirical risk minimization," *CoRR*, vol. abs/1710.09412, 2017. [Online]. Available: <http://arxiv.org/abs/1710.09412>
- [35] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1–9.
- [36] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin, "Emerging Properties in Self-Supervised Vision Transformers," in *arXiv:2104.14294 [cs]*, 2021.
- [37] S. Aïto, M. Awais, W. Wang, M. D. Plumbley, and J. Kittler, "ASiT: Audio Spectrogram vSion Transformer for General Audio Representation," Nov. 2022.
- [38] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," 2015.
- [39] Y. Gong, S. Khurana, A. Rouditchenko, and J. Glass, "CMKD: CNN/Transformer-Based Cross-Model Knowledge Distillation for Audio Classification," Mar. 2022.
- [40] Y. Gong, Y.-A. Chung, and J. Glass, "PSLA: Improving Audio Tagging with Pretraining, Sampling, Labeling, and Aggregation," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 3292–3306, 2021.
- [41] I. Loshchilov and F. Hutter, "Decoupled weight decay regularization," *arXiv:1711.05101 [cs]*, 2017.
- [42] —, "Sgdr - stochastic gradient descent with warm restarts," *arXiv:1608.03983*, 2016.
- [43] J. Salamon, C. Jacoby, and J. P. Bello, "A dataset and taxonomy for urban sound research," in *Proceedings of the 22nd ACM international conference on Multimedia*, 2014, pp. 1041–1044.
- [44] P. Warden, "Speech commands: A dataset for limited-vocabulary speech recognition," *arXiv preprint arXiv:1804.03209*, 2018.
- [45] A. Nagrani, J. S. Chung, and A. Zisserman, "Voxceleb: a large-scale speaker identification dataset," *arXiv preprint arXiv:1706.08612*, 2017.
- [46] 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.
- [47] E. Fonseca, X. Favory, J. Pons, F. Font, and X. Serra, "FSD50K: An Open Dataset of Human-Labeled Sound Events," Apr. 2022.
- [48] H. Bao, L. Dong, S. Piao, and F. Wei, "BEiT: BERT Pre-Training of Image Transformers," Sep. 2022, *arXiv:2106.08254 [cs]*.
- [49] N. Turpault, R. Serizel, A. Shah, and J. Salamon, "Sound event detection in domestic environments with weakly labeled data and soundscape synthesis," in *Acoustic Scenes and Events 2019 Workshop (DCASE2019)*, 2019, p. 253.
- [50] S. Hershey, D. P. W. Ellis, E. Fonseca, A. Jansen, C. Liu, R. Channing Moore, and M. Plakal, "The benefit of temporally-strong labels in audio event classification," in *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 366–370.
- [51] Ç. Bilen, G. Ferroni, F. Tuveri, J. Azcarreta, and S. Krstulović, "A framework for the robust evaluation of sound event detection," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 61–65.
- [52] L. JiaKai, "Mean teacher convolution system for dcase 2018 task 4," *DCASE2018 Challenge*, Tech. Rep., June 2018.
- [53] J. Shor, A. Jansen, R. Maor, O. Lang, O. Tuval, F. d. C. Quitry, M. Tagliasacchi, I. Shavitt, D. Emanuel, and Y. Haviv, "Towards learning a universal non-semantic representation of speech," *arXiv preprint arXiv:2002.12764*, 2020.
- [54] L. Wang, P. Luc, Y. Wu, A. Recasens, L. Smaira, A. Brock, A. Jaegle, J.-B. Alayrac, S. Dieleman, J. Carreira, and A. v. d. Oord, "Towards Learning Universal Audio Representations," Jun. 2022.
- [55] 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.
- [56] S. Pascual, M. Ravanelli, J. Serrà, A. Bonafonte, and Y. Bengio, "Learning Problem-agnostic Speech Representations from Multiple Self-supervised Tasks," *arXiv:1904.03416 [cs, eess, stat]*, 2019.
- [57] K. Chen, X. Du, B. Zhu, Z. Ma, T. Berg-Kirkpatrick, and S. Dubnov, "HTS-AT: A Hierarchical Token-Semantic Audio Transformer for Sound Classification and Detection," in *ICASSP 2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. Singapore, Singapore: IEEE, May 2022, pp. 646–650.
- [58] K. Koutini, J. Schlüter, H. Eghbal-zadeh, and G. Widmer, "Efficient Training of Audio Transformers with Patchout," in *Interspeech 2022*. ISCA, Sep. 2022, pp. 2753–2757.
- [59] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: an asr corpus based on public domain audio books," in *IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2015, pp. 5206–5210.

- [60] J. Turian, J. Shier, H. R. Khan, B. Raj, B. W. Schuller, C. J. Steinmetz, C. Malloy, G. Tzanetakis, G. Velarde, K. McNally, M. Henry, N. Pinto, C. Noufi, C. Clough, D. Herremans, E. Fonseca, J. Engel, J. Salamon, P. Esling, P. Manocha, S. Watanabe, Z. Jin, and Y. Bisk, "HEAR: Holistic Evaluation of Audio Representations," May 2022.