
Audio Barlow Twins: Self-Supervised Audio Representation Learning

Jonah Anton

GLAM, Imperial College London, UK
jonahlouisanton@gmail.com

Harry Coppock

GLAM, Imperial College London, UK
harry.coppock@imperial.ac.uk

Pancham Shukla

Imperial College London, UK
panchamkumar.shukla@imperial.ac.uk

Björn W. Schuller

GLAM, Imperial College London, UK & EIHW, University of Augsburg, Germany
bjoern.schuller@imperial.ac.uk

Abstract

The Barlow Twins self-supervised learning objective requires neither negative samples or asymmetric learning updates, achieving results on a par with the current state-of-the-art within Computer Vision. As such, we present *Audio Barlow Twins*, a novel self-supervised audio representation learning approach, adapting Barlow Twins to the audio domain. We pre-train on the large-scale audio dataset AudioSet, and evaluate the quality of the learnt representations on 18 tasks from the HEAR 2021 Challenge, achieving results which outperform, or otherwise are on a par with, the current state-of-the-art for instance discrimination self-supervised learning approaches to audio representation learning. Code at https://github.com/jonahanton/SSL_audio.

1 Introduction

Inspired by recent successes in Computer Vision (CV) [1, 2, 3] and Natural Language Processing (NLP) [4, 5] in the generation of universal representations¹ through self-supervised learning (SSL) methodologies, much recent interest has been dedicated to using SSL to learn universal representations of audio data [6, 7, 8, 9]. Whilst generative approaches [8, 9] have produced state-of-the-art (SOTA) results for SSL methods in many audio tasks, the current SOTA SSL techniques in CV are dominated by instance discrimination (ID) approaches [10, 11, 12], which build a meaningful representation space through training an encoder network to embed similar instances near one another.

Barlow Twins [13] is one such ID approach, which encourages the empirical cross-correlation matrix between the embeddings of two views of a mini-batch of data samples towards the identity matrix. Through forcing the cross-correlation matrix to the identity, Barlow Twins embeds instances which encode similar semantic content near one another whilst minimising the redundancy between the individual components of the extracted embedding vectors, encouraging the latent representations to be maximally informative. Barlow Twins requires neither negative samples [1] nor asymmetric learning updates [2, 14], instead preventing representational collapse by design. As a result, Barlow Twins i) directly enforces invariances to the applied data augmentations without having to sample

¹A representation is a lower-dimensional and compressed, but highly informative, distillation of an input.

negative pairs, and ii) prevents representational collapse in an intuitive and explainable manner [15], unlike approaches such as BYOL [2] which are theoretically poorly understood (although some attempts have recently been made [16]). Within the audio domain, the sampling of negative pairs is also potentially problematic, since obtaining such a pair from two different audio signals within a mini-batch [6, 17] can lead to low-quality solutions since two signals may share common sounds, such as a chord sequence in music.

It seems reasonable, therefore, that Barlow Twins, when adapted to the audio domain, would produce robust and generalisable audio representations. To this end, we present **Audio Barlow Twins** (ABT), a novel self-supervised audio representation learning method which adapts Barlow Twins [13] to the audio domain. ABT achieves results which outperform, or otherwise are on a par with, the current state-of-the-art for ID self-supervised learning approaches to audio representation learning.

2 Background

Instance Discrimination Instance discrimination (ID) SSL approaches [1, 2, 18] are built on the core idea of similarity: instances which encode similar semantic content should be embedded near one another in representation space. These methods make use of a Siamese network, where each ‘arm’ of the network processes a different view of the data sample. The extracted feature representations of the two views are then pushed together. Solely enforcing representational similarity of positive pairs is vulnerable to mode collapse onto a constant vector for all inputs, a phenomenon known as representational collapse. Contrastive ID approaches, such as SimCLR [1], prevent representational collapse through the use of negative pairs, which are forced apart in representation space. Non-contrastive ID approaches such as BYOL [2] prevent representation collapse, instead, through introduction of asymmetry into the learning framework.

Audio SSL (ID) Many self-supervised learning methods have been proposed to learn generalisable audio representations².Fonseca et al. [17], Saeed et al. [6], Al-Tahan and Mohsenzadeh [20] all adapt SimCLR [1] to the audio domain. Fonseca et al. [17] additionally propose an augmentation which they term *mix-back*, where the incoming spectrogram is mixed with another clip randomly drawn from the training dataset whilst ensuring that the incoming patch remains dominant. Niizumi et al. [7] present BYOL-A, adapting BYOL [2] to the audio domain with minimal modifications from the original learning framework. The key modification they make is their proposed data augmentation module, used to generate the two spectrogram views. BYOL-A also makes use of a lightweight convolutional encoder architecture, based on a network used in the solution of the NTT DCASE2020 Challenge Task 6 (Automated Audio Captioning) [21], which we use in our experiments and term the *AudioNTT* encoder.

3 Method

A schematic depicting ABT’s high-level architecture is detailed in Figure 1.

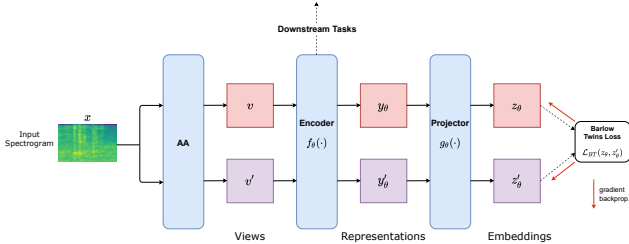


Figure 1: The Audio Barlow Twins learning framework.

Generation of views ABT first produces two views, v, v' of an input spectrogram x by stochastic application of the audio augmentation (AA) module $v, v' \sim \text{AA}(x)$. The audio augmentation module consists of three different augmentation blocks: Mixup, Random Resize Crop (RRC), and Random

²A full and in-depth analysis on the current SOTA audio self-supervised learning methods can be found in the survey produced by [19].

Linear Fader (RLF) [22]. The spectrogram input is first normalised by the dataset mean and standard deviation.

Extraction of embeddings The two views are passed through the encoder to obtain the representations, $y_\theta = f_\theta(v)$, $y'_\theta = f_\theta(v')$. The representations are then passed through the projector network to obtain the embeddings, $z_\theta = g_\theta(y_\theta)$, $z'_\theta = g_\theta(y'_\theta)$.

Barlow Twins objective The Barlow Twins objective, \mathcal{L}_{BT} , is calculated on the embeddings, $\mathcal{L}_{BT}(z_\theta, z'_\theta)$. \mathcal{L}_{BT} , since it uses batch statistics in its calculation of the embeddings’ cross-correlation matrix C , cannot in practice be calculated on an input-by-input basis, but instead must be calculated over a batch of embeddings Z_θ, Z'_θ , with $Z_\theta = [z_\theta^1, \dots, z_\theta^B] \in \mathbb{R}^{B \times d}$, and likewise for Z'_θ . Formally

$$\mathcal{L}_{BT} = \alpha \sum_i (1 - C_{ii})^2 + \lambda \sum_{i \neq j} C_{ij}^2, \quad (1)$$

where the first term enforces representational invariance to the applied audio augmentations, and the second term minimises the redundancy between the individual components of the embedding vectors. The positive constants α and λ control the trade-off between the importance of these two terms, and by default α is set to 1 and λ to 0.005 (as in the original publication [13]). The cross-correlation matrix C is computed between the embeddings within the batch B ,

$$C_{ij} = \sum_{b=1}^B \hat{Z}_{\theta,i}^b \hat{Z}'_{\theta,j}{}^b, \quad (2)$$

where \hat{Z}_θ is the normalised embedding Z_θ along the batch dimension, and $\hat{Z}'_{\theta,i}{}^b$ corresponds to the i^{th} component of the b^{th} batch element of \hat{Z}_θ .

4 Experiments

We pre-train on the large-scale audio dataset AudioSet [23] for 100 epochs with a batch size of 128, which corresponds to $\sim 1.3\text{M}$ training iterations. We successfully download 1,629,756 clips (corresponding to $\sim 4,500$ hours of audio) from AudioSet’s unbalanced train subset, which are used for ABT pre-training.

Audio preprocessing All audio samples are converted with a sampling frequency of 16 kHz to (log-scaled) mel-spectrograms using a 64 ms sliding window with a 10 ms step size, extracting $F = 64$ mel frequency bins in the range 60 – 7,800 Hz. By default, during pre-training, we randomly crop $T = 96$ time frames (all clips with shorter duration are padded with zeros), corresponding to 950 ms of audio. This produces a mel-spectrogram of size $F \times T = 64 \times 96$.

Architecture We consider two encoders, the AudioNTT convolutional encoder [21], and the ViT_C encoder [24]. We consider the ViT_C-B(ase)³ model, using a patch size of 16×8 . A learnable [CLS] token is prepended to the sequence of patches, and its output representation, $O_{[\text{CLS}]}$, is taken as representative of the clip as a whole. Fixed sinusoidal positional encodings are added to each patch.

Downstream Tasks We use 18 tasks from the HEAR 2021 Challenge [25] for evaluation. HEAR includes two types of tasks, i) scene-based tasks, corresponding to classification of an entire audio clip, and ii) timestamp-based tasks, corresponding to sound event detection and or transcription over time. For each task, the representations are extracted from the frozen pre-trained model and then evaluated using the hear-eval⁴ toolkit, which trains a shallow Multilayer Perceptron (MLP) classifier on the extracted representations.

5 Results

We compare the performance of the ABT pre-trained models on the 18 HEAR tasks with two baseline models, CREPE [26], wav2vec2.0 [27], and to BYOL-A^{*5} [7]. The results for the scene-based and timestamp-based tasks are detailed in Tables 1,2 and 3.

³The ViT_C-B corresponds to the ViT_C-18GF model proposed in the original publication [24].

⁴<https://github.com/hearbenchmark/hear-eval-kit>

⁵BYOL-A* is a reimplementaion of BYOL-A [7] by Elbanna et al. [28], which we use since Niizumi et al. [7] did not evaluate on the HEAR tasks.

Table 1: Results on HEAR speech and environmental sound scene-based tasks. Top two performing models for each task are shown underlined and **highlighted**. % \uparrow_{RAND} refers to the fractional increase (**not** the absolute increase) from the average score obtained by the random baseline for that model.

| Model | Speech | | | | | | | Environmental Sound | | | |
|---------------------------------|--------------|--------------|--------------|--------------|--------------|--------------|-----------------------------------|---------------------|--------------|--------------|-----------------------------------|
| | CREMA-D | L5C | SPC-5h | SPC-F | Voclm | VoxL | Avg (% \uparrow_{RAND}) | ESC-50 | FSD50K | Gunshot | Avg (% \uparrow_{RAND}) |
| HEAR CREPE | 0.383 | 0.499 | 0.180 | 0.211 | 0.051 | 0.142 | 0.244 | 0.301 | 0.159 | 0.863 | 0.441 |
| HEAR wav2vec2.0 | <u>0.656</u> | 0.692 | 0.838 | 0.879 | 0.080 | <u>0.493</u> | <u>0.606</u> | 0.561 | 0.342 | 0.848 | 0.584 |
| HEAR BYOL-A* | <u>0.623</u> | <u>0.788</u> | <u>0.896</u> | <u>0.924</u> | <u>0.137</u> | <u>0.390</u> | <u>0.626</u> | <u>0.789</u> | <u>0.489</u> | <u>0.875</u> | <u>0.7180</u> |
| ABT AudioNTT | 0.594 | 0.745 | <u>0.882</u> | <u>0.910</u> | <u>0.111</u> | 0.324 | 0.594 (17%) | <u>0.786</u> | <u>0.474</u> | <u>0.905</u> | <u>0.721 (24%)</u> |
| ABT ViT _{C-B} (16 × 8) | 0.581 | <u>0.812</u> | 0.724 | 0.771 | 0.087 | 0.312 | 0.548 (140%) | 0.705 | 0.446 | 0.845 | 0.666 (49%) |

Table 2: Results on HEAR music scene-based tasks.

| Model | Music | | | | | | | |
|---------------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|-----------------------------------|
| | Beijing | GTZAN-Genre | GTZAN-M/S | Mrd-Stroke | Mrd-Tonic | NSynth 5h | NSynth 50h | Avg (% \uparrow_{RAND}) |
| HEAR CREPE | <u>0.928</u> | 0.645 | 0.929 | 0.898 | 0.824 | <u>0.870</u> | <u>0.900</u> | <u>0.856</u> |
| HEAR wav2vec2.0 | 0.907 | 0.780 | 0.946 | 0.943 | 0.828 | 0.402 | 0.653 | 0.780 |
| HEAR BYOL-A* | 0.919 | <u>0.835</u> | <u>0.969</u> | <u>0.970</u> | <u>0.900</u> | 0.290 | 0.642 | 0.789 |
| ABT AudioNTT | <u>0.966</u> | <u>0.818</u> | 0.962 | <u>0.970</u> | <u>0.932</u> | <u>0.476</u> | <u>0.740</u> | <u>0.838 (-1%)</u> |
| ABT ViT _{C-B} (16 × 8) | 0.869 | 0.765 | <u>0.992</u> | 0.952 | 0.897 | 0.280 | 0.632 | 0.769 (15%) |

Table 3: Results on HEAR timestamp-based tasks. Error rate (\downarrow) indicates that a lower error rate is better. Table format adapted from [28].

| Model | DCASE | | MAESTRO | | Avg (% \uparrow_{RAND}) |
|---------------------------------|--------------|-----------------------------|----------------|---------------------|-----------------------------------|
| | Onset FMS | Error rate (\downarrow) | Onset FMS | Onset w/ Offset FMS | Onset FMS |
| HEAR CREPE | 0.552 | 0.420 | <u>0.3910</u> | <u>0.15</u> | <u>0.472</u> |
| HEAR wav2vec2.0 | 0.670 | 0.320 | 0.0328 | 0.009 | 0.351 |
| HEAR BYOL-A* | 0.499 | 0.503 | 0.0028 | 0.00029 | 0.251 |
| ABT AudioNTT | <u>0.761</u> | <u>0.274</u> | <u>0.04801</u> | <u>0.00672</u> | <u>0.405 (27%)</u> |
| ABT ViT _{C-B} (16 × 8) | <u>0.722</u> | <u>0.275</u> | 0.0263 | 0.00429 | 0.374 (-10%) |

ABT, with the AudioNTT encoder, generally performs on a par with, or outperforms, BYOL-A*, which uses the same AudioNTT encoder architecture, on the scene-based tasks, and consistently outperforms BYOL-A* on the timestamp-based tasks. We see further consistent improvements over the HEAR baseline models (CREPE, wav2vec2.0), except on the type of tasks on which these models have been specialised (music for CREPE, speech for wav2vec2.0). These results demonstrate the robustness of ABT in the generation of general-purpose audio representations. Interestingly, ABT pre-training appears damaging to performance on several of the music tasks, often leading to performance degradation from the random baselines. We find this to be particularly evident for the Mridingam Stroke and Tonic, NSynth (5h and 50h), and MAESTRO tasks. This extends to other ID methods, with BYOL-A*, performing similarly poorly. The aforementioned tasks all require a sound’s pitch to be correctly discerned. However, invariance to pitch perturbations is enforced through RRC, and as such, it is intuitive that a model will consequently struggle to classify pitch. That said, we find through extensive ablation studies, detailed in Appendix C, that RRC does considerably improve the quality of the learnt representations. We therefore observe an issue with transferring ID methods from CV to audio, since such methods rely on applying data augmentations to generate two views, and any given data augmentation may benefit one type of audio task but harm another. This provides support for generative self-supervised methods for learning universal audio representations [8, 9], since they don’t require the use of any data augmentations.

6 Conclusion

In this paper, we presented *Audio Barlow Twins* (ABT), a novel self-supervised audio representation learning method which adapts Barlow Twins [13] to the audio domain. ABT pre-training on AudioSet [23] for 100 epochs with the AudioNTT encoder [21] results in model performance which is on a par with, and in several cases better than, BYOL-A [7]. We found commonly introduced augmentations to be harmful to ABT in certain settings. Future works should consider the effect on different downstream tasks of different augmentations that act directly on raw waveforms, within the ABT learning framework. Applying the augmentations directly on raw waveforms, instead of spectrograms, allows for a) a better control of the strength of these augmentations (as it is possible to *listen* directly to their effect), and b) a greater number of augmentations to be considered (e.g. pitch shift, time masking, time shift, time stretch, fade in/out, compression, etc.). We were unable to apply data augmentations during training directly on raw waveforms in this work due to an I/O bottleneck.

References

- [1] T. Chen, S. Kornblith, M. Norouzi, and G. E. Hinton, “A simple framework for contrastive learning of visual representations,” *CoRR*, vol. abs/2002.05709, 2020. [Online]. Available: <https://arxiv.org/abs/2002.05709> pages 1, 2, 13
- [2] J. Grill, F. Strub, F. Althé, C. Tallec, P. H. Richemond, E. Buchatskaya, C. Doersch, B. Á. 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,” *CoRR*, vol. abs/2006.07733, 2020. [Online]. Available: <https://arxiv.org/abs/2006.07733> pages 1, 2
- [3] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin, “Emerging properties in self-supervised vision transformers,” *CoRR*, vol. abs/2104.14294, 2021. [Online]. Available: <https://arxiv.org/abs/2104.14294> pages 1
- [4] J. Devlin, M. Chang, K. Lee, and K. Toutanova, “BERT: pre-training of deep bidirectional transformers for language understanding,” *CoRR*, vol. abs/1810.04805, 2018. [Online]. Available: <http://arxiv.org/abs/1810.04805> pages 1
- [5] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, “Roberta: A robustly optimized BERT pretraining approach,” *CoRR*, vol. abs/1907.11692, 2019. [Online]. Available: <http://arxiv.org/abs/1907.11692> pages 1
- [6] A. Saeed, D. Grangier, and N. Zeghidour, “Contrastive learning of general-purpose audio representations,” *CoRR*, vol. abs/2010.10915, 2020. [Online]. Available: <https://arxiv.org/abs/2010.10915> pages 1, 2
- [7] D. Niizumi, D. Takeuchi, Y. Ohishi, N. Harada, and K. Kashino, “Byol for audio: Self-supervised learning for general-purpose audio representation,” 2021. [Online]. Available: <https://arxiv.org/abs/2103.06695> pages 1, 2, 3, 4, 12
- [8] Y. Gong, C. J. Lai, Y. Chung, and J. R. Glass, “SSAST: self-supervised audio spectrogram transformer,” *CoRR*, vol. abs/2110.09784, 2021. [Online]. Available: <https://arxiv.org/abs/2110.09784> pages 1, 4
- [9] P.-Y. Huang, H. Xu, J. Li, A. Baevski, M. Auli, W. Galuba, F. Metze, and C. Feichtenhofer, “Masked autoencoders that listen,” 2022. [Online]. Available: <https://arxiv.org/abs/2207.06405> pages 1, 4
- [10] J. Zhou, C. Wei, H. Wang, W. Shen, C. Xie, A. L. Yuille, and T. Kong, “ibot: Image BERT pre-training with online tokenizer,” *CoRR*, vol. abs/2111.07832, 2021. [Online]. Available: <https://arxiv.org/abs/2111.07832> pages 1
- [11] X. Chen, S. Xie, and K. He, “An empirical study of training self-supervised vision transformers,” *CoRR*, vol. abs/2104.02057, 2021. [Online]. Available: <https://arxiv.org/abs/2104.02057> pages 1
- [12] M. Assran, M. Caron, I. Misra, P. Bojanowski, F. Bordes, P. Vincent, A. Joulin, M. Rabbat, and N. Ballas, “Masked siamese networks for label-efficient learning,” 2022. [Online]. Available: <https://arxiv.org/abs/2204.07141> pages 1, 14, 15
- [13] J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny, “Barlow twins: Self-supervised learning via redundancy reduction,” *CoRR*, vol. abs/2103.03230, 2021. [Online]. Available: <https://arxiv.org/abs/2103.03230> pages 1, 2, 3, 4, 8, 9, 11, 13, 14
- [14] X. Chen and K. He, “Exploring simple siamese representation learning,” *CoRR*, vol. abs/2011.10566, 2020. [Online]. Available: <https://arxiv.org/abs/2011.10566> pages 1
- [15] Y.-H. H. Tsai, S. Bai, L.-P. Morency, and R. Salakhutdinov, “A note on connecting barlow twins with negative-sample-free contrastive learning,” 2021. [Online]. Available: <https://arxiv.org/abs/2104.13712> pages 2
- [16] Y. Tian, X. Chen, and S. Ganguli, “Understanding self-supervised learning dynamics without contrastive pairs,” *CoRR*, vol. abs/2102.06810, 2021. [Online]. Available: <https://arxiv.org/abs/2102.06810> pages 2
- [17] E. Fonseca, D. Ortego, K. McGuinness, N. E. O’Connor, and X. Serra, “Unsupervised contrastive learning of sound event representations,” 2020. [Online]. Available: <https://arxiv.org/abs/2011.07616> pages 2

- [18] M. Caron, I. Misra, J. Mairal, P. Goyal, P. Bojanowski, and A. Joulin, “Unsupervised learning of visual features by contrasting cluster assignments,” *CoRR*, vol. abs/2006.09882, 2020. [Online]. Available: <https://arxiv.org/abs/2006.09882> pages 2
- [19] S. Liu, A. Mallool-Ragolta, E. Parada-Cabeleiro, K. Qian, X. Jing, A. Kathan, B. Hu, and B. W. Schuller, “Audio self-supervised learning: A survey,” 2022. [Online]. Available: <https://arxiv.org/abs/2203.01205> pages 2
- [20] H. Al-Tahan and Y. Mohsenzadeh, “Clar: Contrastive learning of auditory representations,” 2020. [Online]. Available: <https://arxiv.org/abs/2010.09542> pages 2
- [21] Y. Koizumi, D. Takeuchi, Y. Ohishi, N. Harada, and K. Kashino, “The ntt dcase2020 challenge task 6 system: Automated audio captioning with keywords and sentence length estimation,” 2020. [Online]. Available: <https://arxiv.org/abs/2007.00225> pages 2, 3, 4
- [22] D. Niizumi, D. Takeuchi, Y. Ohishi, N. Harada, and K. Kashino, “Byol for audio: Exploring pre-trained general-purpose audio representations,” 2022. [Online]. Available: <https://arxiv.org/abs/2204.07402> pages 3, 13
- [23] J. F. Gemmeke, D. P. W. Ellis, D. Freedman, A. Jansen, W. Lawrence, R. C. Moore, M. Plakal, and M. Ritter, “Audio set: An ontology and human-labeled dataset for audio events,” in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2017, pp. 776–780. pages 3, 4
- [24] T. Xiao, M. Singh, E. Mintun, T. Darrell, P. Dollár, and R. B. Girshick, “Early convolutions help transformers see better,” *CoRR*, vol. abs/2106.14881, 2021. [Online]. Available: <https://arxiv.org/abs/2106.14881> pages 3, 8, 9, 11
- [25] 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,” 2022. [Online]. Available: <https://arxiv.org/abs/2203.03022> pages 3, 8
- [26] J. W. Kim, J. Salamon, P. Li, and J. P. Bello, “Crepe: A convolutional representation for pitch estimation,” 2018. [Online]. Available: <https://arxiv.org/abs/1802.06182> pages 3
- [27] A. Baeovski, H. Zhou, A. Mohamed, and M. Auli, “wav2vec 2.0: A framework for self-supervised learning of speech representations,” *CoRR*, vol. abs/2006.11477, 2020. [Online]. Available: <https://arxiv.org/abs/2006.11477> pages 3
- [28] G. Elbanna, N. Scheidwasser-Clow, M. Kegler, P. Beckmann, K. E. Hajal, and M. Cernak, “Byol-s: Learning self-supervised speech representations by bootstrapping,” 2022. [Online]. Available: <https://arxiv.org/abs/2206.12038> pages 3, 4, 8
- [29] Y. You, I. Gitman, and B. Ginsburg, “Large batch training of convolutional networks,” 2017. [Online]. Available: <https://arxiv.org/abs/1708.03888> pages 8, 11
- [30] I. Loshchilov and F. Hutter, “Fixing weight decay regularization in adam,” *CoRR*, vol. abs/1711.05101, 2017. [Online]. Available: <http://arxiv.org/abs/1711.05101> pages 8
- [31] B. McFee, C. Raffel, D. Liang, D. Ellis, M. Mcvicar, E. Battenberg, and O. Nieto, “librosa: Audio and music signal analysis in python,” 01 2015, pp. 18–24. pages 8
- [32] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” 2014. [Online]. Available: <https://arxiv.org/abs/1412.6980> pages 8
- [33] E. Fonseca, X. Favory, J. Pons, F. Font, and X. Serra, “Fsd50k,” Oct. 2020. [Online]. Available: <https://doi.org/10.5281/zenodo.4060432> pages 9
- [34] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama, “Optuna: A next-generation hyperparameter optimization framework,” *CoRR*, vol. abs/1907.10902, 2019. [Online]. Available: <http://arxiv.org/abs/1907.10902> pages 10
- [35] F. Hutter, H. Hoos, and K. Leyton-Brown, “An efficient approach for assessing hyperparameter importance,” *31st International Conference on Machine Learning, ICML 2014*, vol. 2, pp. 1130–1144, 01 2014. pages 10
- [36] 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. [Online]. Available: <https://arxiv.org/abs/2010.11929> pages 11

- [37] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick, “Masked autoencoders are scalable vision learners,” 2021. [Online]. Available: <https://arxiv.org/abs/2111.06377> pages 14

Appendix A Further Implementations Details

In this appendix we provide further details for Audio Barlow Twins pre-training and HEAR evaluation.

A.1 Pre-training

We pre-train on the unbalanced train subset of AudioSet ($\sim 1.6\text{M}$ audio segments) for 100 epochs with a batch size of 128, which corresponds to $\sim 1.3\text{M}$ training iterations. We use a smaller version of the projector network than that proposed in the original Barlow Twins publication [13], although with the same modular structure. The projector network corresponds to a small MLP with one hidden layer, which has hidden dimension 8192, and an output dimension of 1048. The first layer of the projector is followed by a batch normalisation and Rectified Linear Unit (ReLU) non-linearity. The Barlow Twins loss hyperparameters are set to 1 and 5×10^{-3} , for α and λ respectively. For the AudioNTT encoder, we use the Layer-wise Adaptive Rate Scaling (LARS) optimizer [29], with a learning rate of 0.4 for the weights and 0.0048 for the biases and batch normalisation parameters. We use a weight decay of $1 \cdot 10^{-5}$. Following from [13], LARS adaptation, as well as weight decay, do not apply to the biases and batch normalisation parameters. The choice of the LARS optimizer, over Adam and Stochastic gradient descent (SGD), as well as the optimizer hyperparameters (learning rate weights, learning rate biases, weight decay) are selected after conducting an extensive hyperparameter sweep (Appendix B). We also consider the default values⁶ as used by [13], but find a noticeable degradation in model performance. For the ViT_{C-B} encoder, we use AdamW [30] with the default hyperparameter values as suggested⁷ by [24], using a learning rate of $6.25 \cdot 10^{-5}$ and a weight decay of 0.24. Following [24], weight decay is not applied to the biases and any normalisation parameters. AdamW’s β parameters are set as $\beta_1 = 0.9$ and $\beta_2 = 0.999$. All experiments are run on a single NVIDIA RTX 6000 GPU.

A.2 HEAR evaluation

We use 18 tasks from the HEAR 2021 Challenge [25], derived from 15 datasets, to evaluate the quality of the learned audio representations. HEAR includes two types of tasks, *scene-based*, corresponding to classification (multi-class or multi-label) of an entire audio clip, and *timestamp-based*, corresponding to sound event detection and or transcription over time. Scene-based tasks, following Elbanna et al. [28], can be subdivided into three subcategories: speech, environmental sounds, and music. All datasets are downloaded⁸ at 48 kHz and re-sampled to 16 kHz to align with the sampling frequency of the AudioSet clips used during model pre-training. The re-sampling is implemented using the `librosa` Python library [31].

For each task, the embeddings are first extracted from the frozen pre-trained model. The timestamp-based tasks first require the input audio clips to be divided into fixed-size segments, such that embeddings can be extracted corresponding to specific timestamps. We use a segment size of 950 ms with a hop size of 50 ms for all timestamp-based tasks. The embeddings are then evaluated using the `hear-eval`⁹ toolkit, which trains a shallow MLP classifier on the frozen embeddings (linear evaluation). The MLP is trained for a maximum of 500 epochs with the Adam optimizer [32], implementing early stopping on the validation set, checking every 3 epochs with a patience of 20 (except for with DCASE 2016 Task 2, which is checked every 10 epochs). Model selection is performed over a choice of 8 models each of which uses a different hyperparameter configuration, selecting the optimal model using the validation score. Variations in the number of hidden layers, learning rate, and weight initialization are considered. Full details can be found in the original HEAR publication [25].

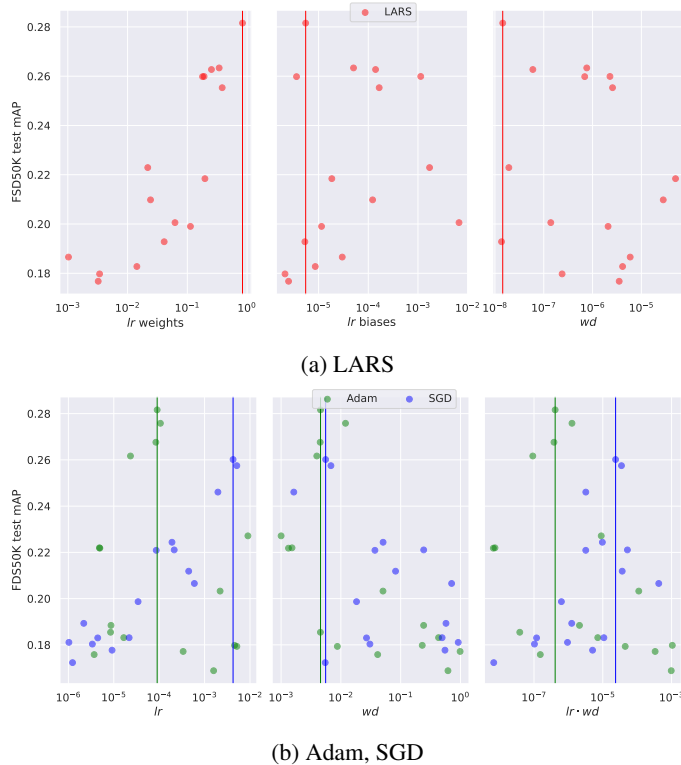


Figure A-1: AudioNTT optimizer hyperparameter sweep. Scatterplots are shown for all three optimizers: a) LARS (red), b) Adam (green), SGD (blue). FSD50K test mAP under linear evaluation is shown on the vertical axis and vertical bars correspond to the optimal values for each optimizer.

Appendix B Hyperparameter Search

Before performing full ABT pre-training on AudioSet, hyperparameter sweeps are conducted over several important variables that we anticipate will most significantly affect optimization. We consider these to be the choice of optimizer, the optimizer learning rate (lr) and weight decay (wd).

For each hyperparameter configuration, ABT pre-training is performed on the FSD50K [33] development subset for a maximum of 20 epochs, corresponding to ~ 6400 training iterations (with a batch size of 128). We measure performance through training a linear classifier¹⁰ (linear evaluation protocol) on the frozen features of the FSD50K train subset, extracted from the pre-trained model, implementing early stopping on the FSD50K validation subset with a patience of 10. The classifier is trained for a maximum of 100 epochs. We report the model performance as the mean Average Precision (mAP) on the FSD50K evaluation subset. We anticipate this metric to be generally indicative of the quality of the learned audio representations for a given set of hyperparameter values. During linear evaluation on FSD50K, all audio clips are first randomly cropped to 96 time frames, the same number as used during ABT pre-training. This is done to reduce the time to extract the embeddings from the pre-trained model. This explains the discrepancy between the FSD50K test scores (mAP)

⁶Appropriately scaled for a batch size of 128 (using linear scaling), [13] use a learning rate of 0.1 for the weights and 0.0024 for the biases, with a weight decay of $1.5 \cdot 10^{-6}$.

⁷Xiao et al. [24] perform an extensive hyperparameter sweep for the ViT_C with a batch size of 2048, using a patch size of 16×16 . They find, with AdamW, a lr of $1 \cdot 10^{-3}$ and a wd of 0.24 to be optimal. We scale this lr linearly by batch size ($0.24 \cdot 128/2048$) to obtain the used value of $6.25 \cdot 10^{-5}$.

⁸<https://zenodo.org/record/5887964>

⁹<https://github.com/hearbenchmark/hear-eval-kit>

¹⁰We train the linear classifier with the following set of hyperparameters: Adam optimizer, batch size = 200, $lr = 1 \cdot 10^{-3}$, $wd = 1 \cdot 10^{-8}$, Adam $\beta_1 = 0.9$, Adam $\beta_2 = 0.999$, Adam $\epsilon = 1 \cdot 10^{-8}$.

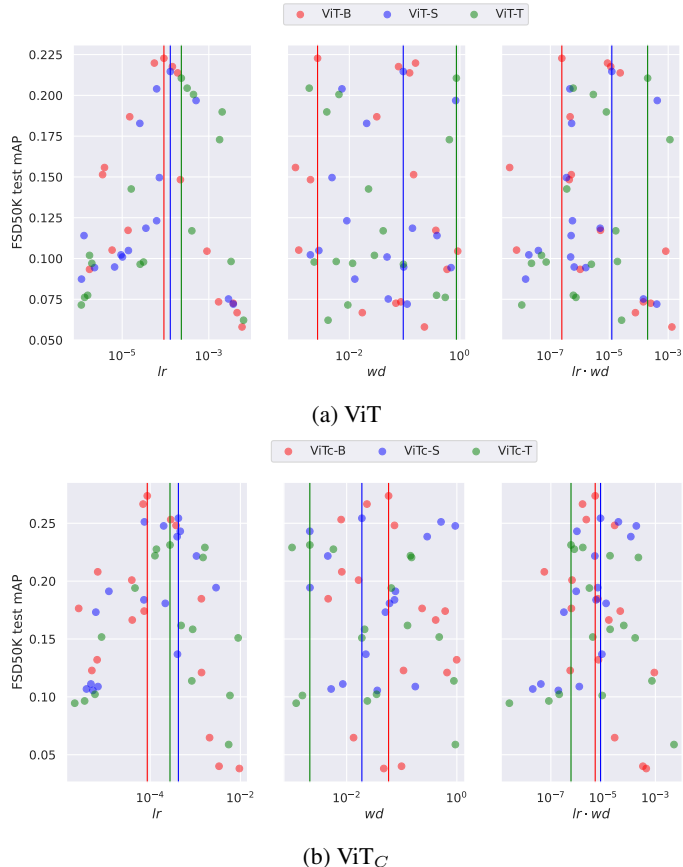


Figure A-2: ViT / ViT_C optimizer hyperparameter sweep. Scatterplots are shown for all three model sizes (-B (red), -S (blue), -T (green)) for both the vanilla ViT (a) and ViT_C (b) models, with FSD50K test mAP under linear evaluation shown on the vertical axis. Vertical bars correspond to the optimal values for each model.

reported here and in Section 5, where linear evaluation is also performed on FSD50K but with no cropping.

The hyperparameter sweeps are implemented using the Optuna framework [34], evaluating 16 hyperparameter configurations (*trials*) per sweep. Each set of hyperparameter values are sampled each trial using Optuna’s Tree-structured Parzen Estimator (TPE) sampler, and we evaluate model performance at the end of each of the 20 pre-training epochs to allow pruning¹¹ of unpromising trials at intermediate stages (before the maximum of 20 epochs). All hyperparameter sweeps use the same audio preprocessing and ABT architectural defaults as in full AudioSet pre-training. That is, except for the projector output dimension, which has a slightly smaller value of 256¹². We measure the importance of individual hyperparameters (i.e. optimization sensitivity) for a given sweep using the fANOVA [35] hyperparameter importance evaluation algorithm (available within the Optuna framework). fANOVA fits a random forest regression model to the scores of the completed (unpruned) trials for the trial hyperparameter configurations. The variance of the fitted model is then decomposed into additive components, each of which are associated with a specific hyperparameter. The fractional variance associated with a hyperparameter is taken as its importance.

AudioNTT For the AudioNTT encoder, we consider pre-training with Adam, SGD and LARS. For Adam and SGD, each trial a (lr , wd) pair is sampled (on a log scale), considering $lr = \{10^{-6}, 10^{-2}\}$

¹¹We use Optuna’s Hyperband pruner to prune unpromising trials.

¹²We use a projector output dimension of 256 both in the hyperparameter sweeps and all ablation studies as initial experimentation suggests that this value is optimal. Extensive ablation studies, however, reveal that a larger value of 1048 is preferable, and as such this value is used in full AudioSet pre-training.

and $wd = \{10^{-3}, 10^0\}$, with all other optimizer hyperparameters set to their PyTorch defaults¹³. Bias and batch normalisation parameters are excluded from weight decay. For LARS¹⁴, we consider separate lrs for the weights and for the biases, considering $lr\ weights = \{10^{-3}, 10^0\}$, $lr\ biases = \{10^{-6}, 10^{-2}\}$, and $wd = \{10^{-8}, 10^{-4}\}$. A $(lr\ weights, lr\ biases, wd)$ triplet is sampled (on a log scale) each trial. For all three optimizers a total of 16 trials are evaluated. Figures A-1a (LARS) and A-1b (Adam, SGD) show scatterplots for the FSD50K test scores (mAP) for the models trained with the three optimizers. We observe that LARS and Adam both attain higher optimal performance (~ 0.28) than SGD (~ 0.26). LARS is also less sensitive to the wd value than Adam and SGD, and as such LARS is used as the default optimizer for pre-training with the AudioNTT encoder¹⁵. From Fig.A-1a, we see that LARS tends to prefer a larger value for $lr\ weights \sim O(10^0)$, showing less sensitivity to $lr\ biases$ and wd , with importance values of 0.81, 0.10, and 0.09, respectively. The optimal trial, which achieves a score of 0.282, uses $(lr\ weights, lr\ biases, wd)$ values of $(0.84, 5.5 \cdot 10^{-6}, 1 \cdot 10^{-8})$. We find in general that $(lr\ weights, lr\ biases, wd)$ values of $(0.4, 4.8 \cdot 10^{-3}, 1 \cdot 10^{-6})$ are effective in general, and as such these values are used with the LARS optimizer by default.

ViT_C We perform hyperparameter sweeps for both the ViT_C [24] and vanilla ViT [36] encoders, considering only the AdamW¹⁶ optimizer. We use a patch size of 16×16 . Each trial, a (lr, wd) pairs is sampled (on a log scale), considering $lr = \{10^{-6}, 10^{-2}\}$ and $wd = \{10^{-3}, 10^0\}$, with all other hyperparameters set to their PyTorch defaults¹⁷. Bias and batch normalisation parameters are excluded from weight decay. Hyperparameter sweeps are conducted for all ViT and ViT_C model sizes ($-B, -S, -T$), evaluating a total of 16 trials for each model. Figures A-2a (ViT) and A-2b (ViT_C) show scatterplots for the FSD50K test scores (mAP) for the ViT and ViT_C models. We observe that for all ViT and ViT_C models, a lr of $\sim 1 \cdot 10^{-4}$ is optimal, although the smaller models ($-S, -T$) tend to prefer a slightly larger value. The models show less sensitivity to the wd value. In general, a lr/wd of $1 \cdot 10^{-4}/0.06$ is effective for all ViT and ViT_C models, and as such these values are used by default during the ablation studies (Appendix C). However, we find that the tuned values used by Xiao et al. [24] for the ViT_C models, a lr/wd of $6.25 \cdot 10^{-5}/0.24$, lead to better model performance¹⁸ when pre-training on the full AudioSet unbalanced train subset. We anticipate that this is because a slightly lower lr is preferred for a significantly longer training schedule¹⁹. We also note three further salient points: 1) The ViT_C models reach a higher optimal performance than their corresponding vanilla ViT models for all model sizes ($-B$: 0.27 vs 0.22, $-S$: 0.25 vs 0.21, $-T$: 0.23 vs 0.21), 2) The ViT_C models show less sensitivity to the exact lr value used than the ViT models, with wider peaks for all three model sizes, supporting the stability claims made by [24], 3) Both ViT and ViT_C optimal performance scales with model size. Point 1) motivates only using the ViT_C, and not vanilla ViT, encoders for full AudioSet ABT pre-training, since the ViT_C encoders significantly outperform them and full AudioSet pre-training with a Transformer encoder takes several days to complete²⁰.

Appendix C Ablation Studies

We perform extensive ablation studies to investigate the contributions of each of the different components of the ABT learning framework. For all ablation studies (except for those in Appendix C.3, which use a ViT_C- B encoder), we perform ABT pre-training for 100 epochs with the AudioNTT encoder on the FSD50K development subset, which corresponds to $\sim 32k$ training iterations (with a batch size of 128). Besides from the training duration and training dataset, all ablation studies use the same experimental set-up as used for full AudioSet pre-training, except for the projector output dimension, which is set by default to 256. For all ablations considered, we evaluate model

¹³For SGD: Nesterov momentum = False. For Adam: $\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 1 \cdot 10^{-8}$.

¹⁴The other LARS hyperparameters (see [29]) are set as: momentum $m = 0.9$, LARS coefficient $\eta = 1 \cdot 10^{-3}$

¹⁵The choice of LARS as the default optimizer (with a convolutional encoder) is further motivated by that LARS is used in the original Barlow Twins publication [13]

¹⁶[24] find that the ViT_C models are also stable when trained with SGD. However, initial experimentation reveals that training ViT_C encoders with SGD (with ABT pre-training) is unstable, with the loss frequently going to NaN.

¹⁷For AdamW: $\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 1 \cdot 10^{-8}$

¹⁸Model performance as measured by tracking FSD50K linear evaluation score every 5 epochs during training.

¹⁹100 epochs on AudioSet unbalanced train segments corresponds to $\sim 250 \times$ training iterations as with 20 epochs on the FSD50K development subset.

²⁰ABT pre-training on the unbalanced train subset of AudioSet with a ViT_C- B on a single NVIDIA RTX 6000 GPU, using mixed precision, takes approximately 120 hours.

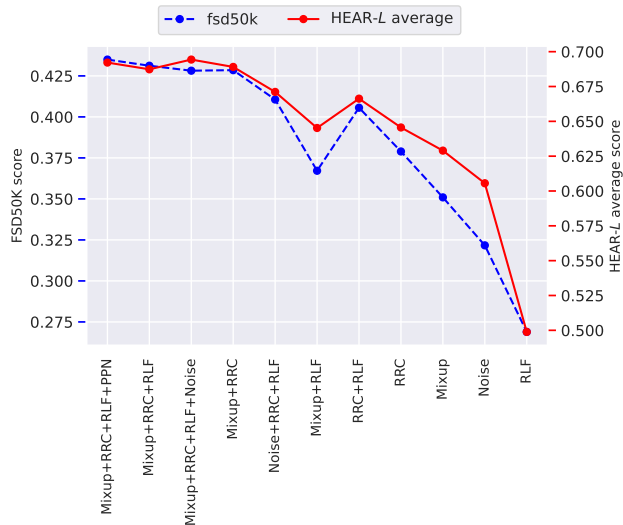


Figure A-3: We compare the effect of pre-training with different combinations of the components of the Audio Barlow Twins audio augmentation (AA) module. Results are shown both evaluated on FSD50K (blue) and the average score on the 5 HEAR-*L* tasks (red).

performance through linear evaluation, using the `hear-eval` toolkit, on a lightweight version of the HEAR Challenge, which we term HEAR-*L*. HEAR-*L* consists of five HEAR tasks covering all three of the scene-based task subcategories: CREMA-D (speech), LibriCount (speech), FSD50K (environmental sound), ESC-50 (environmental sound), and GTZAN Genre (music).

C.1 Audio Augmentations

We consider using different combinations of the components of the audio augmentation (AA) module, which by default consists of Mixup, Random Resize Crop (RRC), and Random Linear Fader (RLF). We further consider two different variations, namely Pre-Post-Norm (*PPN*) and *Noise*. PPN refers to removal of the normalisation block, which standardises input spectrograms by the dataset mean and standard deviation, and replacing it with the Pre- and Post-Normalisation blocks proposed by Niizumi et al. [7] in BYOL-A. Specifically, the pre-normalisation block normalises input spectrograms by batch (and not dataset) statistics, and the post-normalisation block does the same, but after the application of the audio augmentations (Mixup, RRC, RLF). Niizumi et al. [7] argue the post-normalisation corrects the statistical drifts caused by the applied augmentations. Noise refers to the addition of random noise²¹ to an incoming spectrogram. We implement this for direct comparison with Mixup, which interpolates the incoming spectrogram with a natural background signal randomly sampled from the training dataset.

From Figure A-3 we note four salient points:

1. Strong audio augmentations are essential to learn high-quality representations

When all the augmentations are removed from the baseline except RLF (remove RRC and Mixup), Audio Barlow Twins performance drops significantly, by 19 points from 69% to 50% average on the HEAR-*L* tasks.

2. Mixing with natural background signals is more effective than with noise

Mixup improves 13 points from the RLF average on HEAR-*L* to 63%, whereas addition of Gaussian noise results in a smaller improvement of only 10%. Further, using Noise as well as Mixup+RRC+RLF leads to no significant additional performance improvements.

3. RRC is the most effective audio augmentation

RRC, which approximates pitch shift and time stretch, attains the highest performance when any of the audio augmentations are applied alone, achieving a HEAR-*L* average score of 65% (compared with 63% for Mixup, 60% for Noise, and

²¹We sample the noise from a Gaussian distribution $\mathcal{N}(0, \lambda)$, where $\lambda \sim U(0, \alpha)$, $\alpha = 0.2$.

50% for RLF). RLF is by far the least effective augmentation. These findings are consistent with previous results found by [22].

4. PPN shows minimal improvement over dataset normalisation Mixup+RRC+RLF+PPN results in almost identical model performance as Mixup+RRC+RLF (with normalisation by dataset statistics), both having a HEAR-*L* average score of $\sim 69\%$.

C.2 Learning Framework

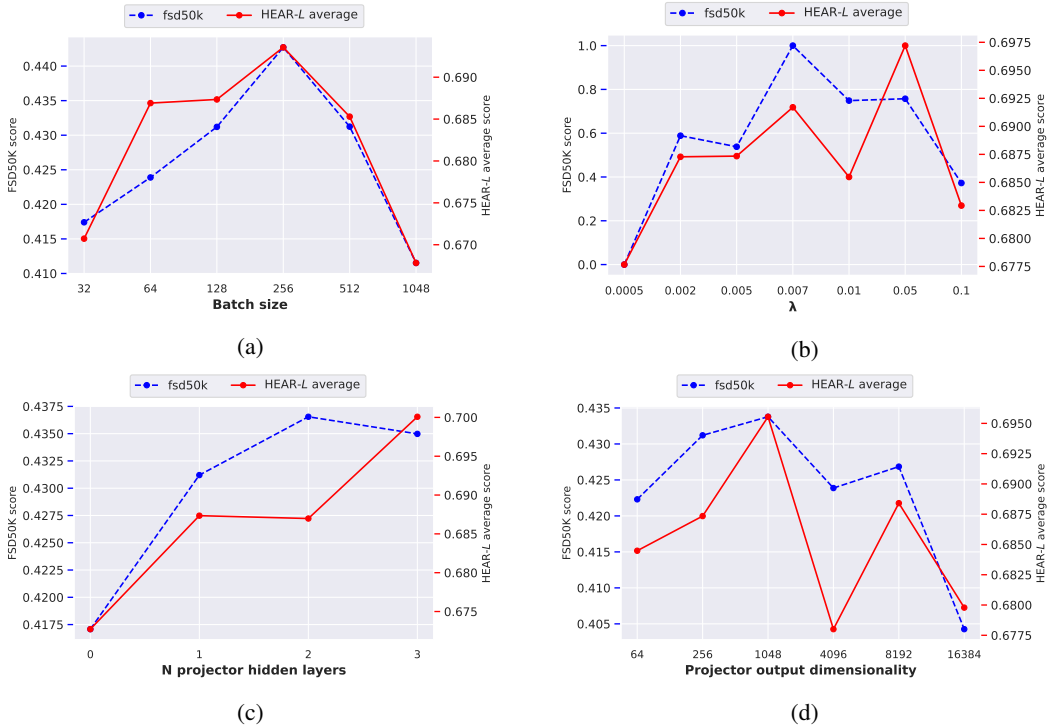


Figure A-4: Ablation studies for (a) batch size, (b) Barlow Twins objective hyperparameter λ , (c) projector depth, and (d) projector output dimensionality. Results are shown both evaluated on FSD50K (red) and the average score on the 5 HEAR-*L* tasks (blue).

Batch Size Figure A-4a shows the sensitivity of Audio Barlow Twins to batch size. The size of the batch is expected to influence training, and therefore downstream performance, since batch dynamics contribute significantly to the Barlow Twins objective (Eqn.1) through the empirical cross-correlation matrix, which is computed across the batch embeddings (Eqn.2: $C_{ij} = \sum_{b=1}^B \hat{Z}_{\theta,i}^b \hat{Z}'_{\theta,j}{}^b$). We observe that ABT exhibits reasonable sensitivity to batch size, with strongest performance with a medium-sized batch containing 64 – 512 samples. This is contrary to methods in CV such as SimCLR [1], which prefer much larger batch sizes (SimCLR requires a batch size of at least 1048 for strong performance). However, we note that all models are trained with the learning rates (*lr weights* and *lr biases*) tuned with a batch size of 128, applying linear scaling: $lr = lr_{128} \times \text{BatchSize}/128$. Re-tuning LARS learning rates for each batch size is beyond the scope of this project. We are unable to consider batch sizes above 1048 due to GPU memory restrictions.

λ , Projector Depth, Projector Output Dimensionality Figures A-4b, A-4c, and A-4d show the variation of model performance with the Barlow Twins objective hyperparameter λ , the number of hidden dimensions of the projector network, and the dimensionality of the embeddings (over which the Barlow Twins objective is calculated). We observe that ABT shows minimal sensitivity to the exact value of λ , as found by Zbontar et al. [13] in the original Barlow Twins publication, although a value in the approximate range $0.002 < \lambda < 0.05$ is preferred, allowing for both the invariance and

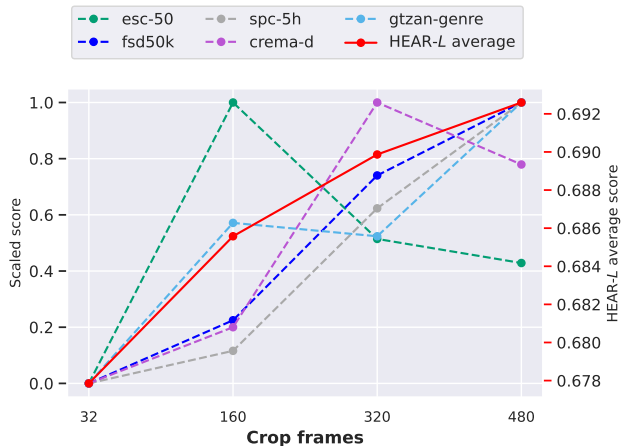


Figure A-5: We compare the effect of pre-training with a different length of the input audio used during pre-training, considering cropping the input spectrograms to 32, 160, 320, and 480 frames. Results are shown evaluated on the individual HEAR-L tasks (ESC-50 (green), Speech Commands 5h (grey), GTZAN Genre (light blue), FSD50K (dark blue), CREMA-D (purple), as well as the HEAR-L average (red). Since the scores on the individual HEAR-L tasks have different scales, we show the MinMax scaled scores, where the best performing input length for each task is set to 1 and the worst performing to 0.

redundancy reduction terms of the Barlow Twins objective (Eqn.1) to contribute. We further observe that a deeper projector is preferred, although performance does not significantly rise above a depth of 2 (1 hidden layer). Contrary to the observations of Zbontar et al. [13], we don’t find that model performance continues to improve as projector output dimensionality grows, with saturation at an output size of 1048, and considerable performance degradation observed with a dimensionality of 16,384.

Input Audio Duration We consider variations in the length of the input audio used during ABT pre-training. Specifically, we consider cropping the input spectrograms to 32, 160, 320, and 480 frames, which correspond to ~ 320 ms, 1.6s, 3.2s, and 4.8s of audio, respectively. We choose not to show the results with the default 96 crop frames as all hyperparameters have been tuned using this value, and as a result it is not considered to be a fair comparison. As shown in Figure A-5, almost all HEAR-L tasks benefit from a longer training window, with the exception of ESC-50. For speech tasks, such as Speech Commands or CREMA-D, this is intuitive, since an element of speech may last several seconds, and as such using only a short segment of under one second in pre-training can result in sounds that don’t retain the original semantic content of the clip (e.g. a word may be cropped to only a syllable or single character). However, ESC-50 clearly seems to benefit from a shorter input duration, showing optimal performance with 160 input frames. The environmental sound dataset ESC-50 contains many sounds categories which consist of short, sharp noises, such as the categories mouse click and door knock, and as such using a short segment during pre-training may better align with the actual sound duration of this dataset. The optimal input duration during pre-training therefore appears to be dependent on the downstream dataset being evaluated on, although there exists a general trend that longer clips are beneficial.

C.3 View Masking

In their recent work *Masked Siamese Networks* (MSN), Assran et al. [12] randomly drop a subset of the patches of one of the two image views before being processed by a Siamese Network architecture using a ViT encoder, matching in feature space the masked view with the unmasked view and thereby performing “implicit denoising” [12] at the representation level. MSN achieves SOTA results whilst simultaneously reducing computational and memory requirements, since the masked patches can be dropped before input into the ViT encoder. It is therefore of great interest to see whether adapting this approach to the audio domain can be beneficial in the pursuit of universal audio representations. Similarly to MSN, we consider adding the step of randomly masking patches from one of the two spectrogram views before input into the ViT_C encoder to ABT’s augmentation module. We implement random patch masking using the algorithm proposed by He et al. [37], where the list of extracted

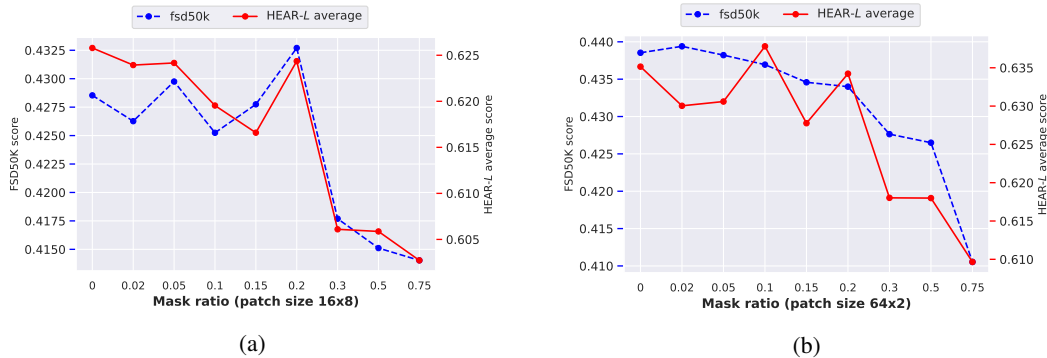


Figure A-6: We consider the effect of partial view masking of one of the two spectrogram views, inspired by recent work by [12]. We pre-train with different masking ratios, considering masking with both (a) 16×8 patches and (b) 64×2 patches. Results are shown both evaluated on FSD50K (blue) and the average score on the 5 HEAR-L tasks (red).

patches is randomly shuffled, and the last M patches from the list are removed, where $M = rN$ (rounded to the nearest integer), with r being the masking ratio and N the initial number of patches. We consider partial view masking with both a patch size of 16×8 and 64×2 , which both correspond to a total of $N = 48$ patches (with 64×96 spectrogram inputs).

Disappointingly, as shown in Figures A-6a and A-6b, partial view masking seems to harm the quality of the learned audio representations, with a clear trend that over a minimum threshold for r ($r \sim 0.2$, corresponding to $M = 10$ masked patches), model performance is significantly reduced. Below this threshold we generally see slight degradation in model performance, although minimal variation (expected as only very few patches have been masked). We anticipate that masking a large number of spectrogram patches may fundamentally change the semantic content of the audio clip, such that matching the representations of the masked and unmasked views encourages the model to embed together audio samples in representation space which no longer share the same semantic content, thereby damaging the quality of the learned representations. This is different to in CV, where strong masking doesn't visually appear to change the overall semantic content contained within an image (e.g. a heavily masked picture of a dog is still recognisable as a dog).

We additionally consider whether, instead of using a fixed masking ratio, slowly increasing the masking ratio during pre-training leads to improved representation quality. Starting the masking ratio at 0, we increase it to a value β at epoch 100 following a sinusoidal schedule with a warm up period of 10 epochs. However, initial experimentation with $\beta = 0.3$ suggests that this also results in a degradation of model performance, although extensive analysis has not been performed.