

Learning Representations for New Sound Classes With Continual Self-Supervised Learning

Zhepei Wang[#], Cem Subakan^{bb}, Xilin Jiang[†], Junkai Wu[#], Efthymios Tzinis[#], Mirco Ravanelli^{bx^b}, and Paris Smaragdis^{#‡}, *Fellow, IEEE*

Abstract— In this paper, we work on a sound recognition system that continually incorporates new sound classes. Our main goal is to develop a framework where the model can be updated without relying on labeled data. For this purpose, we propose adopting representation learning, where an encoder is trained using unlabeled data. This learning framework enables the study and implementation of a practically relevant use case where only a small amount of the labels is available in a continual learning context. We also make the empirical observation that a similarity-based representation learning method within this framework is robust to forgetting even if no explicit mechanism against forgetting is employed. We show that this approach obtains similar performance compared to several distillation-based continual learning methods when employed on self-supervised representation learning methods.

Index Terms—Continual Learning, Representation Learning, Self-Supervised Learning, Sound Classification.

I. INTRODUCTION

DEEP neural networks for audio classification require large amounts of data to be trained on [1], [2]. However, in many realistic deployments, additional labeled data for new sound classes might present itself after initial training, necessitating a time and resource-consuming retraining process that incorporates both past and new data. This problem could be exacerbated by storage constraints, hardware corruption, or privacy regulations that can limit access to past data.

Continual/lifelong learning [3] has been a field of rising interest to address the aforementioned concerns. The goal is to design learning algorithms that would continuously learn from a sequence of tasks to let the models imitate the learning process of human beings. A major problem that arises when a model is continually trained is called *catastrophic forgetting* [4], which is associated with deteriorating performance on previously learned tasks.

Various solutions have been proposed to combat catastrophic forgetting in supervised learning settings, including storage of previous data [5]–[8], or using generative models [9], regularization of the model parameters between tasks [3], [10]–[12], and the progressive growth of neural networks to combat forgetting [13]–[16]. These continual learning methods are applied in the image domain, however, continual learning has potential applications in speech and audio processing systems as well, including personal voice assistants, and conversational AI agents [17]. There indeed exists a few papers on

supervised continual learning on audio, such as environmental sound classification [18], [19], fake audio detection [20], and audio captioning [21].

In this paper, we propose to adopt continual representation learning (CRL) in the sound event classification domain, and experimentally show its efficacy for class-incremental continual learning. Representation learning decouples the learning into two stages: a) Learning of an encoder with a representation-specific objective (e.g., with a self-supervised learning objective). b) Finetuning of a shallow classifier on top of the encoder for a downstream task. To the best of our knowledge, this is the first time that continual representation learning has been used in the sound event classification domain. We argue that using a representation learning pipeline in continual learning is especially beneficial for practical use-case implications. Namely, i) the majority of the computational burden is passed on to the encoder learning stage, and therefore the shallow output layer can be trained including earlier tasks. ii) Decoupling the output head gives the additional flexibility to add an indefinite number of classes by re-training an ad-hoc output head for each task, which is required in real-life applications. iii) In the case where labeled data is scarce, learning an encoder from unlabeled data is beneficial for generalization performance, as we showcase with our experiments.

As hinted above, we adopt self-supervised learning (SSL) within CRL. SSL is a relatively novel research area that is revolutionizing deep learning due to its impressive ability to learn features that generalize well without labels. In SSL, deep networks are trained with pretext tasks such as clustering [22], mutual information maximization [23], image colorization [24], masked token prediction [25], [26], contrastive learning [27]–[29] and so on. These early works on self-supervised learning focused on images and natural language processing. Only recently has self-supervised learning been extended to audio and speech as well [2], [30]–[34].

In the space of incorporating representation learning within a continual learning setup, recent works include [35]–[38]. Different from this paper, these works are applied to images and involve explicit ways to combat catastrophic forgetting. A continual representation learning method for speech recognition is proposed in [39] by learning more languages over time, but the size of the model grows as more languages are involved, and language information is required during inference. In contrast, we focus on a more challenging setup where the model remains constant in size while progressively learning new sound classes, and no prior information on the input data is needed at test time (e.g., category of the sound).

To summarize, in this paper we propose using the continual

This manuscript is submitted for review on 10/03/2022. This work was supported in part by NIFA award #2020-67021-32799.

[#] University of Illinois at Urbana-Champaign, USA; [‡] Concordia University, Canada; ^χ Université de Montréal, Canada; [†] Columbia University; ^b Mila-Quebec AI Institute, Canada; [‡] Amazon Web Services, USA.

representation learning paradigm for class-incremental learning of new audio classes. We list our contributions as follows:

- We show that adopting representation learning for continual learning of new sound classes enables the study of the practically interesting case where only a small amount of labels is present. This is a realistic use case where the system is mainly trained on unlabeled data. We investigate both in-domain and out-of-domain performance. For the evaluation of the partially labeled use case, we propose several evaluation metrics.
- We empirically observe that even if we do not employ an explicit mechanism to combat forgetting, employing similarity-based self-supervised learning within CRL yields better performance than continual supervised representation learning, and comparable performance to distillation-based continual learning methods to combat forgetting, which requires additional storage and computation.

II. PROBLEM STATEMENT

A. Continual, Class-Incremental, Self-Supervised Learning

Continual learning (CL) aims to train a model on a sequence of datasets in a non-stationary training distribution. In this setup, the model is restricted from accessing any past or future data samples. Given an ordered sequence of tasks $\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_T$ with the associated datasets $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_T$, the model is trained sequentially on one task at a time. In particular, we are interested in the class-incremental learning scenario [40] where any two tasks would not share a common class. During inference, the model is required to discriminate samples from all classes it has seen without knowing which task the test data comes from.

In this paper, we employ similarity-based SSL methods [27]–[29], [41], [42], where we maximize the similarity between embeddings from a pair of distorted views of the same input to make representations robust to distortions. We consider SimCLR [28], [41], MoCo [27], and Barlow Twins [42] as candidates for the continual representation learning framework. To create positive pairs for these similarity-based methods, we apply audio data augmentation by first taking fixed-length random segments from each clip, and then applying SpecAugment [43] to the spectrograms.

B. Continual Representation Learning

State-of-the-art representation learning approaches require a large dataset to train high-quality representations. In practical cases where the dataset grows over time, frequently updating the model using all of the data samples is computationally prohibitive. The proposed continual representation learning framework circumvents this bottleneck by training an encoder only for the most recent task, which carries the majority of the computational burden. For instance, suppose we collect new data with new urban sound types (e.g., car horn, gunshot) to train a model that has been previously trained to recognize animal sounds (e.g., dog, rooster). Instead of retraining the network on both datasets, the encoder is updated only on the urban sound classes while preserving the ability to effectively

represent animal classes. Once the encoder is updated with new **unlabeled** data, a shallow output layer is trained from scratch on top of the encoder, with available **labeled** data.

That is, in continual representation learning, we aim to learn low-dimensional embeddings by training an encoder model f_θ on one task at a time. Under the class-incremental setting, in each task, the encoder is trained with data from a set of classes that corresponds to the current task: When entering a new task, f_θ is initialized from the end of the previous task. It is then updated with a representation learning objective with the current dataset $\mathcal{D}_t^{\text{train}}$.

After training, the learned representations are evaluated on a classification task. A shallow classifier h'_ψ is trained on top of the encoder’s output using a small amount of labeled data $\mathcal{D}_t^{\text{eval}}$ with the encoder’s weights fixed. In the case of in-domain evaluation, $\mathcal{D}_t^{\text{eval}}$ contains examples from previous classes (see Figure 1), and in the case of out-of-distribution evaluation, $\mathcal{D}_t^{\text{eval}}$ may come from a different distribution from the dataset used for training the encoder. This continual learning framework is shown in Fig. 1.

We propose using similarity-based self-supervised algorithms for representation learning in this framework, and we call this approach continual self-supervised representation learning (**CSSL**). We compare CSSL with continual supervised representation learning (**CSUP**), which relies on the assumption that annotations are accessible. In CSUP, we train an additional classifier h_ψ jointly with f_θ to propagate label information. This classifier is discarded for downstream tasks or evaluation. For a fair comparison, we apply identical data augmentations and encoder architecture to both approaches.

We hypothesize that similarity-based SSL is beneficial in the context of CRL due to its strong generalization ability. The set of features learned from the current data generalizes well to past data, and consequently, this results in a system that is less prone to forgetting.

TABLE I
OFFLINE CLASSIFICATION ACCURACY

| | CSUP | SimCLR | Barlow Twins | MoCo |
|--------------|-------------|--------|--------------|------|
| UrbanSound8K | 80.9 | 74.3 | 73.3 | 69.3 |
| DCASE TAU19 | 68.2 | 62.5 | 58.8 | 50.3 |

Offline accuracy when representations are trained with the entire dataset.

TABLE II
IN-DOMAIN EVALUATION FOR CRL

| Method | UrbanSound8K ($T=5$) | | | | DCASE TAU19 ($T=5$) | | | |
|-------------------------------------|------------------------|--------------------|------------------|--------------------|-----------------------|--------------------|------------------|--------------------|
| | LEP | | SLEP | | LEP | | SLEP | |
| | A (\uparrow) | F (\downarrow) | A (\uparrow) | F (\downarrow) | A (\uparrow) | F (\downarrow) | A (\uparrow) | F (\downarrow) |
| No distillation | | | | | | | | |
| CSUP | 65.6 | 19.6 | 48.4 | 33.1 | 48.2 | 27.6 | 32.5 | 38.4 |
| SimCLR | 70.3 | 15.3 | 50.3 | 26.6 | 59.7 | 17.8 | 42.1 | 27.9 |
| Barlow Twins | 68.5 | 14.1 | 49.7 | 20.5 | 55.9 | 19.0 | 41.0 | 23.4 |
| MoCo | 68.4 | 15.6 | 50.3 | 25.8 | 49.5 | 20.2 | 34.8 | 25.8 |
| With distillation | | | | | | | | |
| CSUP + \mathcal{L}_{MSE} | 58.6 | 27.0 | 43.8 | 39.2 | 49.1 | 26.1 | 35.1 | 35.8 |
| CSUP + \mathcal{L}_{sim} | 70.6 | 13.8 | 54.9 | 27.1 | 56.2 | 19.7 | 42.1 | 32.4 |
| CSUP + \mathcal{L}_{KLD} | 69.8 | 15.9 | 55.4 | 27.4 | 55.7 | 19.4 | 42.6 | 30.3 |
| SimCLR + \mathcal{L}_{MSE} | 70.9 | 14.6 | 50.6 | 25.1 | 56.2 | 19.6 | 42.0 | 26.5 |
| SimCLR + \mathcal{L}_{sim} | 70.6 | 14.0 | 51.1 | 25.0 | 60.0 | 17.6 | 42.8 | 25.9 |

Average accuracy (A) and forgetting (F) for CSSL and CSUP methods. Best performances are in bold. Note that rows with gray background correspond to the CSSL approaches where no-distillation against forgetting is used.

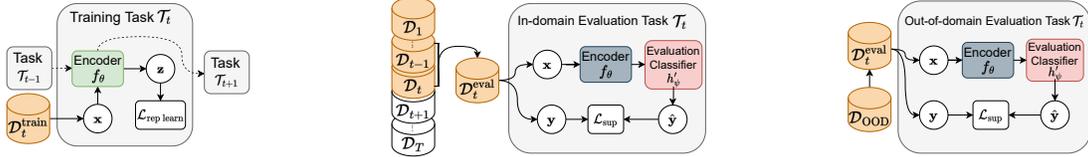


Fig. 1. **(Left) Training pipeline for continual representation learning:** Representations are learned by continually training an encoder f_θ with one task/dataset at a time. **(Middle) In-domain evaluation:** The representations are evaluated by training a classifier h'_ψ using a labeled training set on top of the frozen f_θ . **(Right) Out-of-domain evaluation:** The dataset used for training the output head h'_ψ is different from the dataset for encoder training.

III. EVALUATION PROTOCOL

SSL algorithms are usually trained and evaluated in two steps [27], [28], [42]. Similarly, in the continual learning setting, we first learn the representation with the encoder f_θ only using data from the current task. At the end of each task, we evaluate the representation by training an evaluation classifier, h'_ψ , randomly initialized, on top of the pre-trained frozen encoder (see Fig. 1). Note that h'_ψ is much smaller in size compared to the encoder and is computationally more affordable to train. Also, notice that the classifier h_ψ in CSUP (trained jointly with f_θ) is discarded before evaluation since only the classes for the current task are seen during training while all previously seen classes are used for evaluation. Initializing a new evaluation classifier h'_ψ from scratch ensures a fair comparison between CSSL and CSUP.

Following supervised continual learning metrics for classification, for each task $t \in \{1, 2, \dots, T\}$, we compute the accuracy of the classifier on the test set of task j using the encoder at task t denoted as $A_{t,j}$. We measure the average accuracy at the end of each task: $\bar{A}_t = \frac{1}{t} \sum_i A_{t,i}$ and obtain the average accuracy at the end of the training, $\bar{A} = \bar{A}_T$. Additionally, we compute forgetting, which measures the average decrease in accuracy of each task between its peak and the final performance, defined by $\bar{F} = \frac{1}{T-1} \sum_{j=1}^{T-1} \max_{\tau=1,2,\dots,T} (A_{\tau,j} - A_{T,j})$.

We propose a set of evaluation protocols with different data distributions and types of models for the downstream classification task. When encoder training and evaluation take place on the same dataset (see Fig. 1 Middle), we propose the following in-domain protocols: i) Linear Evaluation Protocol (**LEP**) [37], where a linear classifier h'_ψ is trained with the output of the fixed, pre-trained encoder f_θ using labeled data from past and current task t , $\{\mathcal{D}_\tau\}_{\tau=1}^t$; ii) Subset Linear Evaluation Protocol (**SLEP**), where h'_ψ is trained with a random subset, $\{\mathcal{D}'_\tau\}_{\tau=1}^t$, where $\mathcal{D}'_\tau \subset \mathcal{D}_\tau$. This protocol simulates the scenario with limited labels. We keep a subset of 200 samples per task ($\approx 12\%$ of data). For both LEP and SLEP, classification becomes more challenging over time since the test data involves more classes as more tasks are learned.

Additionally, we propose an out-of-domain (OOD) protocol (Fig. 1 Right) where the representations are evaluated on a different dataset from which they are learned. We introduce the Full Linear Evaluation Protocol (**FLEP**), where at each task the linear evaluation classifier is trained on the full downstream dataset. As opposed to the in-domain protocols, here the test set does not grow with more tasks.

IV. EXPERIMENTAL SETUP

A. Datasets

We use UrbanSound8K [44], the TAU Urban Acoustic Scenes 2019 (DCASE TAU19) [45], and VGGSound [46] to continually train the encoder. We partition each dataset into class-disjoint subsets. Following [40], we set the number of tasks to $T = 5$ or $T = 10$, and the classes are evenly split between tasks. Namely, each task corresponds to a set of classes. For example, in a $T = 5$ case for a 10-class dataset, we have two distinct classes in each task. The ordering of classes is randomly shuffled before the split. The UrbanSound8K and the TAU Urban Acoustic Scenes 2019 (DCASE TAU19) are used for in-domain evaluation, and VGGSound [46] is used to learn representations for OOD experiments. We refer the readers to the appendix for more details on the datasets.

B. Model Training

We use the CNN14 architecture as the backbone encoder model with 78M trainable parameters [1]. The output representation has 2048 dimensions. We train the encoder for 50 epochs per task on UrbanSound8K and VGGSound, and 100 epochs per task on DCASE TAU19. The linear classifiers for evaluation are trained for 30 epochs per task on UrbanSound8K and 100 epochs on DCASE TAU19.

Results are averaged across 10 folds on UrbanSound8K and across three runs with different seeds on DCASE TAU19 and VGGSound. Additional details can be found in our SpeechBrain [47] implementation¹.

V. RESULTS AND DISCUSSIONS

A. Offline Evaluation

We first train and evaluate the representations with supervised and SSL algorithms under the ideal offline scenario in which the entire training set is available for training (i.e., there is only $T = 1$ task). We consider three SSL algorithms, namely, SimCLR, MoCo, and Barlow Twins. For evaluation, we train a new classifier h'_ψ on the output of the fixed, pre-trained encoder f_θ following Section III.

As we show in Table I on both datasets, CSUP outperforms the CSSL approaches. This indicates that the supervised classification loss provides a strong baseline for representation learning when trained with the entire corpus at once without continual learning constraints.

B. Comparison between CSSL and CSUP

We compare representations learned from CSSL and CSUP in terms of their classification performance and resilience to

¹https://github.com/zhepei/w/cssl_sound

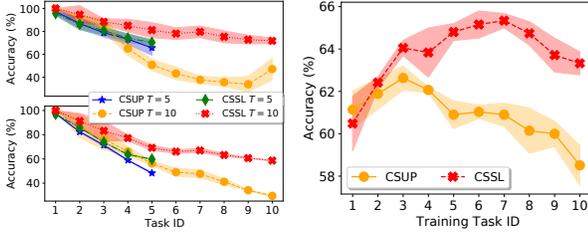


Fig. 2. **(Left) The In-Domain experiment:** The classification performance obtained with CSSL and CSUP on UrbanSound8K (top) and DCASE TAU19 (bottom). **(Right) The OOD experiment:** Performance of CSSL and CSUP trained on VGGSound and evaluated with full-set linear evaluation protocol (FLEP) on DCASE TAU19.

forgetting. Similar to the offline experiments, we consider SimCLR, MoCo, and Barlow Twins for CSSL. In this experiment, we perform an in-domain evaluation with $T = 5$ tasks, and we measure both values for average accuracy and forgetting at the completion of all tasks.

Table II (No distillation) summarizes the average accuracy and forgetting for the in-domain protocols. First, the accuracy values for all learning methods for $T = 5$ are lower than the corresponding offline accuracy in Section V-A; this indicates the existence of catastrophic forgetting for continual representation learning regardless of the approaches. While supervised learning beats SSL in offline evaluation, the trend is reversed when evaluating in a continual learning setup. All three CSSL methods (highlighted in gray) outperform CSUP on both datasets across all evaluation protocols and metrics. The consistent advantage of SSL in both accuracy and forgetting indicates that **the similarity-based objectives learn features that generalize better and are less affected by forgetting than CSUP.**

In Fig. 2 (Left), we also show the trajectories of accuracies obtained with CSSL and CSUP for five tasks ($T = 5$) and ten tasks ($T = 10$) using SimCLR as the CSSL method. CSSL has higher average accuracy, and generally, the gap widens for both $T = 5$ and $T = 10$ as training progresses.

C. Knowledge Distillation

We additionally study the effects of regularization-based CL algorithms by adapting the Learning Without Forgetting (LwF) [3] framework to the continual representation learning setting. At each task $t > 1$, the representation learning objective is optimized jointly with a knowledge distillation loss $\mathcal{L}_{\text{dist}}(f_{\theta}(\mathbf{x}), f_{\theta(t-1)}(\mathbf{x}))$, where $f_{\theta(t-1)}$ is the snapshot of the encoder at the completion of task $t - 1$ whose weights are frozen at the current task t . We consider the following candidates for the distillation loss $\mathcal{L}_{\text{dist}}$: i) mean squared error, \mathcal{L}_{MSE} [48], ii) similarity-based objective that is used in SimCLR [49], \mathcal{L}_{sim} , and iii) KL-Divergence loss, \mathcal{L}_{KLD} (used for CSUP only) applied to the logits obtained after the output head [3]. Table II (With distillation) shows the in-domain evaluation results of CSSL with SimCLR and CSUP with these loss options. To verify the consistency of these results for different SSL algorithms, we also tried Barlow Twins and MoCo and observed similar performance. We omit these results due to space constraints.

We find that distillation with \mathcal{L}_{MSE} does not improve the results for CSUP on UrbanSound8K nor SimCLR on DCASE

TAU19. \mathcal{L}_{KLD} outperforms \mathcal{L}_{MSE} , but it is still beaten by the plain SimCLR except for the SLEP evaluation protocol, even though the plain SimCLR does not require storing any models or labels. Both \mathcal{L}_{MSE} and \mathcal{L}_{KLD} explicitly restrict the output to mitigate forgetting, but they reduce the plasticity of learning knowledge from new tasks and hence do not improve the overall performance. With the similarity-based loss \mathcal{L}_{sim} , the accuracy of CSUP increases significantly on both datasets. The performance gain of CSUP once again highlights **the generalization ability of the similarity-based self-supervised objective and its resilience against forgetting in continual representation learning.** The marginal improvement with \mathcal{L}_{sim} on SimCLR shows that the similarity-based framework alone learns features that generalize across tasks as well as distillation-based methods. For CSSL, distillation does not provide more advantages of generalization despite the cost of additional computation during training and storage for model saving. **Given the computational benefits, we hence conclude that CSSL without explicit methods for combating forgetting is preferable over the alternatives that use distillation.**

D. Out-of-domain (OOD) Evaluation

In addition to the in-domain evaluation, we compare CSSL and CSUP when the representation learning and downstream evaluation are performed on different datasets. Fig. 2 (Right) shows the trajectories of the average accuracy using the FLEP protocol on DCASE TAU19 when the encoder is trained on VGGSound. We consider SimCLR for CSSL using $T = 10$ tasks. When evaluating the representation using a linear classifier with the FLEP protocol, we see that CSSL outperforms CSUP after the second task. The performance gap widens as the encoder learns more tasks. We observe a decreasing trend in the accuracy curve of CSUP, indicating that the transfer of knowledge is overwhelmed by catastrophic forgetting. For CSSL, the curve keeps increasing for the first seven tasks. The results from the OOD evaluation demonstrate the effectiveness of CSSL to learn continually without labels and its ability to generalize to a different downstream task.

VI. CONCLUSION

In this work, we propose a continual representation learning framework for sound classes. In this framework, the encoder training does not rely on labels and therefore is suitable for the practically relevant case where only a subset of labels is used to finetune an output classifier. Additionally, the framework is flexible enough to continually incorporate novel classes without the apriori knowledge of the total number of classes. With the continually pre-trained representations, the computational burden is significantly alleviated by simply finetuning a shallow classifier for downstream tasks. We showed, for the first time, that continual self-supervised learning gets competitive performance even if we do not use any mechanism against forgetting, which helps to reduce computational complexity. In future work, we plan to integrate continual self-supervised learning on more challenging audio tasks such as multi-label classification.

REFERENCES

- [1] Q. Kong, Y. Cao, T. Iqbal, Y. Wang, W. Wang, and M. D. Plumbley, "Panns: Large-scale pretrained audio neural networks for audio pattern recognition," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2020.
- [2] Y. Gong, C.-I. J. Lai, Y.-A. Chung, and J. Glass, "Ssast: Self-supervised audio spectrogram transformer," *arXiv preprint arXiv:2110.09784*, 2021.
- [3] Z. Li and D. Hoiem, "Learning without forgetting," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018.
- [4] R. French, "Catastrophic forgetting in connectionist networks," *Trends in Cognitive Sciences*, no. 4, 1999.
- [5] S.-A. Rebuffi, A. Kolesnikov, G. Sperl, and C. H. Lampert, "icarl: Incremental classifier and representation learning," *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [6] E. Belouadah and A. Popescu, "Il2m: Class incremental learning with dual memory," in *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019.
- [7] F. M. Castro, M. J. Marín-Jiménez, N. Guil, C. Schmid, and K. Alahari, "End-to-End Incremental Learning," in *European Conference on Computer Vision (ECCV)*, 2018.
- [8] A. Douillard, M. Cord, C. Ollion, T. Robert, and E. Valle, "Small-task incremental learning," *ArXiv*, 2020.
- [9] H. Shin, J. K. Lee, J. Kim, and J. Kim, "Continual learning with deep generative replay," in *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [10] J. Kirkpatrick, R. Pascanu, N. C. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska, D. Hassabis, C. Clopath, D. Kumaran, and R. Hadsell, "Overcoming catastrophic forgetting in neural networks," *Proceedings of the National Academy of Sciences (NAS)*, 2017.
- [11] D. Lopez-Paz and M. Ranzato, "Gradient episodic memory for continual learning," in *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [12] H. Ritter, A. Botev, and D. Barber, "Online structured laplace approximations for overcoming catastrophic forgetting," in *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- [13] A. Mallya and S. Lazebnik, "Packnet: Adding multiple tasks to a single network by iterative pruning," *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [14] X. Tao, X. Hong, X. Chang, S. Dong, X. Wei, and Y. Gong, "Few-shot class-incremental learning," *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [15] A. A. Rusu, N. C. Rabinowitz, G. Desjardins, H. Soyer, J. Kirkpatrick, K. Kavukcuoglu, R. Pascanu, and R. Hadsell, "Progressive neural networks," 2016.
- [16] R. Aljundi, P. Chakravarty, and T. Tuytelaars, "Expert gate: Lifelong learning with a network of experts," *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [17] M. F. McTear, "Conversational ai: Dialogue systems, conversational agents, and chatbots," *Synthesis Lectures on Human Language Technologies*, vol. 13, pp. 1–251, 2020.
- [18] Z. Wang, C. Subakan, E. Tzinis, P. Smaragdis, and L. Charlin, "Continual learning of new sound classes using generative replay," in *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, 2019.
- [19] Y. Wang, N. J. Bryan, M. Cartwright, J. Pablo Bello, and J. Salamon, "Few-shot continual learning for audio classification," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021.
- [20] H. Ma, J. Yi, J. Tao, Y. Bai, Z. Tian, and C. Wang, "Continual learning for fake audio detection," in *Interspeech*, 2021.
- [21] J. van den Berg and K. Drossos, "Continual learning for automated audio captioning using the learning without forgetting approach," in *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*, 2021.
- [22] M. Caron, P. Bojanowski, A. Joulin, and M. Douze, "Deep clustering for unsupervised learning of visual features," in *European Conference on Computer Vision (ECCV)*, 2018.
- [23] R. D. Hjelm, A. Fedorov, S. Lavoie-Marchildon, K. Grewal, P. Bachman, A. Trischler, and Y. Bengio, "Learning deep representations by mutual information estimation and maximization," in *International Conference on Learning Representations (ICLR)*, 2019.
- [24] R. Zhang, P. Isola, and A. A. Efros, "Colorful image colorization," in *European Conference on Computer Vision (ECCV)*, 2016.
- [25] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. B. Girshick, "Masked autoencoders are scalable vision learners," *ArXiv*, 2021.
- [26] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics*, 2019.
- [27] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [28] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *Proceedings of the International Conference on Machine Learning (ICML)*, 2020.
- [29] M. Caron, I. Misra, J. Mairal, P. Goyal, P. Bojanowski, and A. Joulin, "Unsupervised learning of visual features by contrasting cluster assignments," in *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- [30] A. van den Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," *ArXiv*, vol. abs/1807.03748, 2018.
- [31] M. Ravanelli, J. Zhong, S. Pascual, P. Swietojanski, J. Monteiro, J. Trmal, and Y. Bengio, "Multi-task self-supervised learning for robust speech recognition," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020.
- [32] A. Baevski, H. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," in *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- [33] D. Niizumi, D. Takeuchi, Y. Ohishi, N. Harada, and K. Kashino, "Byol for audio: Self-supervised learning for general-purpose audio representation," in *International Joint Conference on Neural Networks (IJCNN)*, 2021.
- [34] P.-H. Chi, P.-H. Chung, T.-H. Wu, C.-C. Hsieh, S.-W. Li, and H. yi Lee, "Audio albert: A lite bert for self-supervised learning of audio representation," *IEEE Spoken Language Technology Workshop (SLT)*, 2021.
- [35] E. Fini, V. G. T. da Costa, X. Alameda-Pineda, E. Ricci, K. Alahari, and J. Mairal, "Self-supervised models are continual learners," *arXiv preprint arXiv:2112.04215*, 2021.
- [36] D. Madaan, J. Yoon, Y. Li, Y. Liu, and S. J. Hwang, "Representational continuity for unsupervised continual learning," in *International Conference on Learning Representations*, 2022.
- [37] Z. Ni, S. Tang, and Y. Zhuang, "Self-supervised class incremental learning," *arXiv preprint arXiv:2111.11208*, 2021.
- [38] S. Zhang, G. Shen, and Z. Deng, "Self-supervised learning aided class-incremental lifelong learning," *ArXiv*, 2020.
- [39] S. Kessler, B. Thomas, and S. Karout, "An adapter based pre-training for efficient and scalable self-supervised speech representation learning," in *International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2022.
- [40] G. M. van de Ven and A. S. Tolias, "Three scenarios for continual learning," *arXiv preprint arXiv:1904.07734*, 2019.
- [41] T. Chen, S. Kornblith, K. Swersky, M. Norouzi, and G. E. Hinton, "Big self-supervised models are strong semi-supervised learners," in *Advances in Neural Information Processing Systems*, 2020.
- [42] J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny, "Barlow twins: Self-supervised learning via redundancy reduction," in *Proceedings of the International Conference on Machine Learning*, M. Meila and T. Zhang, Eds. PMLR, 2021.
- [43] D. S. Park, W. Chan, Y. Zhang, C.-C. Chiu, B. Zoph, E. D. Cubuk, and Q. V. Le, "SpecAugment: A simple augmentation method for automatic speech recognition," in *INTERSPEECH*, 2019.
- [44] J. Salamon, C. Jacoby, and J. P. Bello, "A dataset and taxonomy for urban sound research," in *ACM International Conference on Multimedia*, 2014.
- [45] A. Mesaros, T. Heittola, and T. Virtanen, "A multi-device dataset for urban acoustic scene classification," in *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*, 2018.
- [46] H. Chen, W. Xie, A. Vedaldi, and A. Zisserman, "Vggsound: A large-scale audio-visual dataset," in *International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2020.
- [47] M. Ravanelli, T. Parcollet, P. Plantinga, A. Rouhe, S. Cornell, L. Lugosch, C. Subakan, N. Dawlatabad, A. Heba, J. Zhong, J.-C. Chou, S.-L. Yeh, S.-W. Fu, C.-F. Liao, E. Rastorgueva, F. Grondin, W. Aris, H. Na, Y. Gao, R. D. Mori, and Y. Bengio, "Speechbrain: A general-purpose speech toolkit," *ArXiv*, 2021.
- [48] J. Smith, J. Tian, Y.-C. Hsu, and Z. Kira, "A closer look at rehearsal-free continual learning," *ArXiv*, 2022.
- [49] K. Sohn, "Improved deep metric learning with multi-class n-pair loss objective," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2016.

APPENDIX

A. Dataset Details

1) *UrbanSound8K*: The UrbanSound8K dataset [44] is a dataset for sound event recognition that contains 8,732 recordings with a total duration of 8.75 hours, where the length of each clip is limited to 4 seconds. Each recording is labeled with a single event class from the 10 possible classes, where each class has no more than 1,000 clips. All files are pre-sorted into 10 folds for cross-validation.

2) *TAU Urban Acoustic Scenes 2019*: The TAU Urban Acoustic Scenes 2019 dataset is used for the DCASE 2019 Task-1(A) challenge [45] for acoustic scene classification. The development set consists of 40 hours of audio collected from ten cities, with 9,185 recordings for training and 4,185 recordings for evaluation. Each clip is 10 seconds long and is labeled with one of 10 possible classes.

3) *VGGSound*: The VGGSound dataset [46] is an audio-visual dataset with 200,000 clips with a total duration of 560 hours from more than 300 sound classes such as instruments, horns, and city sounds. The recordings are scraped from YouTube and are labeled by pre-trained image and sound classifiers. However, these data samples may not be accurately labeled due to i) the classifiers are pre-trained on a different dataset, and ii) the recordings may contain interfering sound events while each clip is assigned a single label. We consider learning the representation using VGGSound while evaluating with a different downstream dataset with an OOD setup.

B. Continual Representation Learning with Full Replay Buffer

We also consider continually training the encoder f_θ and the output head h_ψ with a full replay (FR) buffer by storing data samples from the current and all previous tasks with $T = 5$ tasks. We provide the final average accuracy using the in-domain linear evaluation protocol for UrbanSound8K and DCASE TAU19 in Table III for CSUP and CSSL using SimCLR, denoted by CSUP-FR and SimCLR-FR, respectively. We also include the corresponding offline systems (denoted by CSUP-O, SimCLR-O) from Table I and the continually trained systems without using replay buffer (trained only on the current task) - (denoted by CSUP-NR, SimCLR-NR) taken from Table II. Both FR systems achieve the best performance compared to the offline and NR systems by beating the offline systems by a small margin (possibly due to a curriculum effect) and significantly exceeding the accuracy of the NR systems. When the encoder observes data from all tasks in both offline and full replay scenarios, the supervised algorithm slightly outperforms SimCLR. However, this performance gain comes at the extra computational cost of storing and revisiting past data samples during the encoder training. Also notice that the gap between CSUP-FR and CSUP-NR is significantly larger than the one between SimCLR-FR and SimCLR-NR, and this further indicates that representations learned with supervised objectives are more prone to performance degradation when access to previous data is restricted.

C. Assessing the Impact of Self-Supervised Objectives

To analyze the influence of the SSL objective on continual representation learning, we first perform an ablation study

TABLE III
IN-DOMAIN EVALUATION FOR DIFFERENT ENCODER TRAINING DATASET

| | CSUP-O | SimCLR-O | CSUP-NR | SimCLR-NR | CSUP-FR | SimCLR-FR |
|-------|--------|----------|---------|-----------|-------------|-----------|
| US8K | 80.9 | 74.3 | 65.6 | 70.3 | 82.3 | 77.2 |
| DCASE | 68.2 | 62.5 | 48.2 | 59.7 | 69.1 | 67.4 |

Final average accuracy with the linear evaluation protocol when representations are trained with the full dataset at once (CSUP-O, SimCLR-O), continually trained without using data buffer (CSUP-NR, SimCLR-NR), and using full data buffer (CSUP-FR, SimCLR-FR) for $T = 5$ tasks.

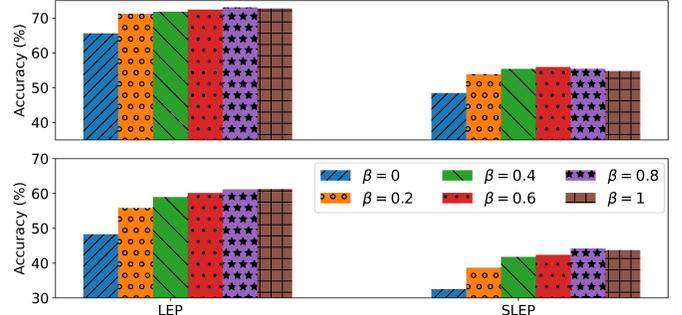


Fig. 3. Average accuracy on the linear combination of CSSL and CSUP with the supervised weight $\alpha = 1$ fixed on in-domain linear evaluation (LEP) and subset evaluation (SLEP) protocols on UrbanSound8K (upper) and DCASE TAU19 (lower).

by controlling the amount of self-supervision in the learning framework. For this experiment, we assume the existence of annotations and train the encoder jointly with both classification loss and the similarity-based SSL loss [49], given by $\mathcal{L}_{\text{joint}} = \alpha \mathcal{L}_{\text{sup}} + \beta \mathcal{L}_{\text{ssl}}$, where $\alpha, \beta \geq 0$ control the weight of each objective. We fix the weight of the supervised loss as $\alpha = 1$ and increase the weight of the SSL objective β from 0 to 1. Notice that $\alpha = 1, \beta = 0$ is equivalent to CSUP. We compute the final average accuracy at the end of $T = 5$ tasks using the in-domain protocols on two different datasets, as shown in Figure 3. In general, higher weights on SSL leads to better performance. In particular, there is a significant performance gain between $\beta = 0$ and $\beta = 0.2$ across both protocols. These observations imply the benefits of incorporating similarity-based SSL objectives in continual representation learning even if label information is available.

D. Impact from the Number of Tasks

We are interested in whether the relative performance between CSSL and CSUP is consistent with a varying total number of tasks. As mentioned in Section V-B, Fig. 2 (Right) displays the trajectories of the average accuracy at the end of each task for CSSL using SimCLR and CSUP with $T = 5, 10$ total tasks. CSSL has a higher average accuracy at the completion for both $T = 5, 10$ on both datasets than CSUP. The performance gap is significantly larger in $T = 10$ than in $T = 5$. Furthermore, in $T = 10$, this gap widens as training progresses. Since both UrbanSound8K and DCASE TAU19 contains 10 classes, $T = 10$ creates an extreme setting where the encoder learns only one class at a time. If training without replay data using cross-entropy, the supervised framework may optimize by biasing the weight towards the current class without learning useful representations. In contrast, without label information, the similarity-based CSSL framework is more robust to the shift in data distribution across tasks.