MgSvF: Multi-Grained Slow vs. Fast Framework for Few-Shot Class-Incremental Learning

Hanbin Zhao, Yongjian Fu, Mintong Kang, Qi Tian, Fei Wu, and Xi Li

Abstract—As a challenging problem, few-shot class-incremental learning (FSCIL) continually learns a sequence of tasks, confronting the dilemma between slow forgetting of old knowledge and fast adaptation to new knowledge. In this paper, we concentrate on this "slow vs. fast" (SvF) dilemma to determine which knowledge components to be updated in a slow fashion or a fast fashion, and thereby balance old-knowledge preservation and new-knowledge adaptation. We propose a multi-grained SvF learning strategy to cope with the SvF dilemma from two different grains: intra-space (within the same feature space) and inter-space (between two different feature spaces). The proposed strategy designs a novel frequency-aware regularization to boost the intra-space SvF capability, and meanwhile develops a new feature space composition operation to enhance the inter-space SvF learning performance. With the multi-grained SvF learning strategy, our method outperforms the state-of-the-art approaches by a large margin.

Index Terms—Few-shot class-incremental learning, multi-grained, class-incremental learning

1 Introduction

R Ecent years have witnessed a great development of class-incremental learning [1]–[8], which aims at enabling a learner to acquire new knowledge from new data while preserving the learned knowledge from previous data. In practice, the new knowledge from new data is often represented in a challenging few-shot learning scenario (i.e., few annotated samples), leading to a problem named fewshot class-incremental learning [9] (FSCIL). FSCIL typically involves the learning stages of the base task (i.e., the first task with large-scale training samples) and the new tasks (with limited samples). In principle, FSCIL is in a dilemma between slow forgetting of old knowledge and fast adaptation to new knowledge. As shown in Figure 2 (a) and (b), slow forgetting typically leads to underfitting on new tasks, while fast adaptation incurs a catastrophic forgetting problem. Hence, a "slow vs. fast" (SvF) learning pipeline is needed to be implemented to determine which knowledge components to be updated in a slow fashion or a fast fashion, keeping a trade-off between slow-forgetting and fastadaptation as shown in Figure 2 (c). In this paper, we focus on investigating the SvF learning performance differences from two different grains: within the same feature space (called intra-space SvF analysis) and between two different feature spaces (called inter-space SvF analysis).

In the literature, a number of approaches only maintain a unified feature space for SvF learning w.r.t. different feature dimensions [9], [11]–[17]. Since the unified feature space has mutually correlated feature dimensions, it is difficult to disentangle the feature for SvF analysis. Besides, the learning directions for old-knowledge preservation and new-knowledge adaptation are usually inconsistent with each

E-mail: tian.qi1@huawei.com.

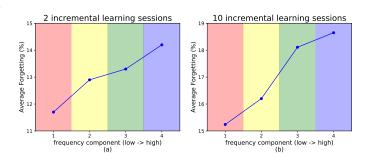


Fig. 1. Analysis of intra-space SvF on CIFAR100. The forgetting of previous tasks is estimated with *average forgetting* [1], [10]. (a) Results with 2 learning sessions. (b) Results with 10 learning sessions. It can be seen that different frequency components appear different characteristics for old-knowledge transfer. In both learning settings, the lower frequency components achieve less forgetting, and the average forgetting increases along with the frequency.

other (even contradictory sometimes). In the case of FSCIL, the unified feature space tends to fit the data of new-task well, but suffers from the degradation of discriminability and generalization ability, as well as catastrophic forgetting.

Motivated by the above observations, we build an intraspace SvF feature disentanglement scheme by Discrete Cosine Transform (DCT), resulting in an orthogonal frequency space with mutually uncorrelated frequency components. Subsequently, we propose to evaluate the SvF performance differences w.r.t. different frequency components. The evaluation results indicate that different frequency components indeed appear to have different characteristics for knowledge transfer. As shown in Figure 1, the low-frequency components contribute more to preserving old knowledge, and the average forgetting increases along with the frequency. Therefore, we turn out to build the new feature space updated by frequency-aware knowledge distillation, which enforces higher weights on the regularization term of approximating the low-frequency components of the old

H. Zhao, Y. Fu, M. Kang, F. Wu, and X. Li are with College of Computer Science and Technology, Zhejiang University, Hangzhou 310027, China. E-mail: {zhaohanbin, yjfu, kangmintong, wufei, xilizju}@zju.edu.cn.

Q. Tian is with Cloud BU, Huawei Technologies, Shenzhen 518129, China.

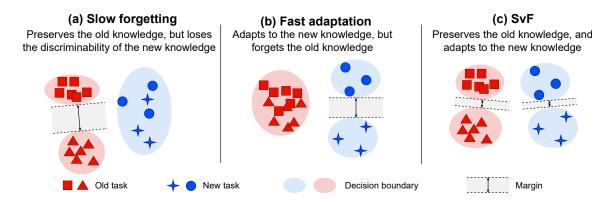


Fig. 2. Illustration of slow vs. fast analysis for few-shot class-incremental learning. (a) mainly pays attention to slow forgetting. Samples of old tasks are separated by a large margin but that of new tasks are mixed up. (b) puts emphasis on fast adaptation. New-task samples are separable while old-task samples are mixed up. (c) keeps a trade-off between slow-forgetting and fast-adaptation and solves all tasks well.

feature space. Experiments show that this simple frequency-aware learning strategy results in an effective performance improvement.

Given the aforementioned frequency space, we also maintain another separate feature space that focuses on preserving the old knowledge to further enhance the discriminability. In this way, we set up an inter-space SvF learning scheme to make the separate feature space updated more slowly than the other. In the inter-space SvF learning scheme, we propose a feature space composition operation to compose the above two spaces. In principle, different composition operations are flexible to be used. From extensive experiments, we find out that even an extremely simple uniform concatenation strategy can result in a dramatic performance improvement.

Overall, the main contributions of this work are summarized in the following three aspects: 1) We propose a multi-grained "slow vs. fast" (SvF) learning framework for FSCIL, which aims to balance old-knowledge preserving and new-knowledge adaptation. To our knowledge, it is the first work to introduce multi-grained SvF into FSCIL. 2) We present two simple yet effective SvF learning strategies for intra-space and inter-space cases, that is, frequency-aware regularization and feature space composition operation. 3) Extensive experiments over all the datasets demonstrate the effectiveness of our approach as our method outperforms state-of-the-art approaches by a large margin.

2 RELATED WORK

2.1 Incremental Learning.

Recently, there has been a large body of research in incremental learning [18]–[41]. These works can be categorized into three major families: 1) architectural strategies, 2) rehearsal strategies, 3) regularization strategies. Architectural strategies [42]–[51] keep the learned knowledge from previous tasks and acquire new knowledge from the current task by manipulating the network architecture, e.g., parameter masking, network pruning. Rehearsal strategies [15]–[17], [52]–[60] replay old tasks information when learning the new task, and the past knowledge is memorized by storing old tasks' exemplars or old tasks data distribution via generative models. Regularization strategies [36], [61]–[69]

alleviate forgetting by regularization loss terms enabling the updated parameters of networks to retain past knowledge. Incremental learning is usually conducted under the task-incremental [70], [71] or the class-incremental learning scenarios [2], [68], [72]–[75]. This paper considers the latter where the task identity is non-available at inference time. Few-shot class-incremental learning [9] is a more practical and challenging problem, where only a few number of samples for new tasks are available. The aforementioned approaches resort to one unified feature space for SvF learning to balance old-task knowledge preserving and new-task knowledge adaptation. In this paper, we investigate the SvF learning performance from intra-space and inter-space grains and propose a multi-grained SvF learning strategy for FSCIL (i.e., frequency-aware regularization and feature space composition operation).

2.2 Frequency domain learning.

Recently, a series of research works explore introducing frequency transformation into deep learning to provide auxiliary information [76], [77]. Some of these frequency-aware methods aim to reduce the computing cost with frequency transformation [78], [79], thus improving the network efficiency. Others propose to conduct frequency-aware augmentation [80], [81] to improve the robustness or to solve the problem of domain adaptation. Apart from that, there is also a work that decouples the entangled information and regularizes the model on some specific components [82] (e.g., focus more on details in vehicles recognition task). Our intra-space SvF analysis is more inspired by the work on frequency-based decouple [82], discussing different properties of different frequency components and proposing a frequency-aware regularization method for FSCIL.

3 METHODS

3.1 Few-Shot Class-Incremental Learning

To better understand our representations, we provide detailed explanations of the main notations and symbols used throughout this paper as shown in Table 1.

In a few-shot class-incremental learning setup, a network learns several tasks continually, and each task contains a

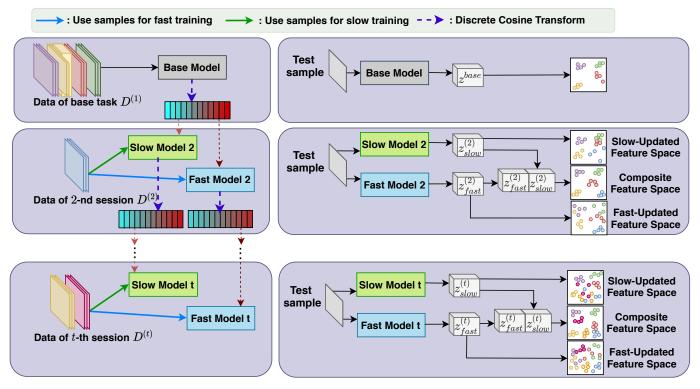


Fig. 3. Illustration of our method. The dash lines show the frequency-aware intra-space regularization, and the solid lines indicate the inter-space composition operation. At the 1-st session, a base embedding model is initially trained on a large-scale training set of base task. At the t-th (t>1) learning session, two embedding models are fast or slowly updated on data of t-th task by intra-space SvF learning, then we composite the slow-updated feature space and the fast-updated feature space, and finally use the composite feature space for classification.

TABLE 1
Main notations and symbols used throughout the paper.

Notation	Definition
$D^{(t)}$	the training set for the <i>t</i> -th task and contains only a few samples
$D^{(1)}$	the first task (termed as a base task) and has a large number of samples
$C^{(t)}$	the set of classes of the <i>t</i> -th task
$\operatorname{dist}(\cdot,\cdot)$	the distance metric (e.g., Euclidean distance)
\mathbf{z}_{j}	the embedding of a given sample x_j in the original entangled feature space
\mathbf{u}_c	the mean of embedding of class c
\hat{y}_{j}	the prediction result of a given sample x_j
$T(\cdot)$	the transformation function
$\overline{\mathbf{z}}_j$	the embedding of a given sample x_j in an orthogonal feature space
$\overline{\mathbf{z}}_{j,q}$	the q -th component of \mathbf{z}_j in the frequency domain
Q	the total number of frequency components
γ_q	the weight on the regularization term of approximating the q -th frequency component of the old embedding space
$\Psi(\cdot,\cdot)$	the composition function (e.g., for a naive implementation, a simple concatenation operation)
$\mathbf{\tilde{z}}_{j}$	the composite feature of sample x_j

batch of new classes [1], [9], [50]. The time interval from the arrival of the current task to that of the next task is considered as a FSCIL session [83]. We suppose the training set for the t-th task is $D^{(t)}$ and each $D^{(t)}$ only contains a few samples, except the first task (termed as a base task) $D^{(1)}$, which has a large number of training samples and classes instead. $C^{(t)}$ is the set of classes of the t-th task, and we consider the generally studied case where there is no overlap between the classes of different tasks: for $i \neq j, C^{(i)} \cap C^{(j)} = \varnothing$. At the t-th session, we only have access to the data $D^{(t)}$, and the model is trained on it.

FSCIL can be formulated in two frameworks [1], [15]. The first framework is composed of a feature extractor and a softmax classifier, and they are trained jointly. The other

one only needs to train an embedding network, mapping samples to a feature space where distance represents the semantic discrepancy between samples [84], and utilizes an nearest class mean (NCM) classifier for classification [1], [15], [85], which is defined as:

$$\hat{y}_j = \operatorname*{argmin}_{c \in \bigcup_i C^{(i)}} \operatorname{dist}(\mathbf{z}_j, \mathbf{u}_c), \tag{1}$$

where $\operatorname{dist}(\cdot, \cdot)$ is the distance metric (e.g., Euclidean distance) and we denote the embedding of a given sample x_j as \mathbf{z}_i and \mathbf{u}_c is the mean of embedding of class c [15].

We follow the latter framework in this paper and thus our goal is to obtain a discriminative feature space \mathbf{z}_j performing well on all the seen tasks, which means balancing

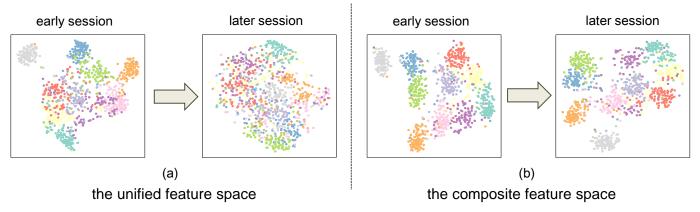


Fig. 4. (Best viewed in color.) Visualization of samples in "the unified feature space" or "the composite feature space" by t-SNE on CIFAR100. Samples of ten classes are from two tasks and each class is represented by one color. (a): Samples in the unified feature space at an early session and a later session; (b): Samples in the composite feature space at an early session and a later session.

old knowledge slow-forgetting and new knowledge fast-adapting well at each session. In FSCIL, It becomes rather more difficult because the feature space tends to overfit the very few number of samples, and suffers from the degradation of discriminability and generalization ability, as well as catastrophic forgetting. Therefore, a method needs to be proposed to disentangle the learned knowledge embedded in \mathbf{z}_j , determining which knowledge components to be updated slowly or fast, and thereby achieving slow old-knowledge forgetting and fast new-knowledge adaptation.

Our multi-grained "slow vs. fast" (SvF) few-shot class-incremental learning pipeline (i.e., how to train our embedding model) is detailed in the rest of this section. An overview of our whole pipeline is shown in Figure 3. The dash lines in Figure 3 illustrate the intra-space grained frequency regularization, which is analyzed in Section 3.2, while the solid lines show inter-space grained space composition operation, elaborated in Section 3.3.

3.2 Intra-Space Level SvF learning

To train our embedding model, we utilize a metric learning loss to ensure that the distance between similar instances is close, and vice versa. The objective function is formulated as:

$$\mathcal{L} = \mathcal{L}_{ML} + \lambda \mathcal{L}_{R},\tag{2}$$

where \mathcal{L}_{ML} is the metric loss term and \mathcal{L}_R is the regularization term for retaining past knowledge [1], [62], [64], [67], λ denotes the trade-off. The triplet loss [86] is often adopted as the metric learning loss \mathcal{L}_{ML} :

$$\mathcal{L}_{ML} = \max(0, d_{+} - d_{-} + r), \tag{3}$$

where d_+ and d_- are the Euclidean distances between the embeddings of the anchor x_a and the positive instance x_p and the negative instance x_n respectively, and r denotes the margin. For regularization term \mathcal{L}_R , previous methods usually retain old knowledge by distilling on a unified feature space (i.e., directly approximating the embedding $\mathbf{z}_j^{(t-1)}$ while updating $\mathbf{z}_j^{(t)}$ at the t-th session):

$$\mathcal{L}_R = \left\| \mathbf{z}_j^{(t)} - \mathbf{z}_j^{(t-1)} \right\|,\tag{4}$$

where $\|\cdot\|$ denotes the Frobenius norm.

To better analyze intra-space SvF learning in an orthogonal space with uncorrelated components, we propose an intra-space SvF feature disentanglement scheme. We utilize $T(\cdot)$ to denote the transformation function, and the embedding in the original entangled feature space \mathbf{z}_j is transformed to the embedding in an orthogonal feature space $\overline{\mathbf{z}}_j$:

$$\overline{\mathbf{z}}_j = T(\mathbf{z}_j),\tag{5}$$

In this paper, we utilize Discrete Cosine Transform (DCT) to conduct the transformation. In this way, the transformed feature $\overline{\mathbf{z}}_i$ is with the same length as \mathbf{z}_i :

$$\overline{\mathbf{z}}_j = [\overline{\mathbf{z}}_{j,1}, \overline{\mathbf{z}}_{j,2}, \cdots, \overline{\mathbf{z}}_{j,Q}], \tag{6}$$

where Q denotes the total number of frequency components (also the length of \mathbf{z}_j), and each component $\overline{\mathbf{z}}_{j,q}$ denotes the q-th component of \mathbf{z}_j in the frequency domain.

As shown in Figure 1, we find that different frequency components indeed appear to have different characteristics for knowledge transfer when separately distilling them in the fashion of Equation 4 to preserve old knowledge. Specifically, we find that the low-frequency components of the feature space often contribute more to preserving old-knowledge. Therefore, we turn out to approximate the old feature space by frequency-aware knowledge distillation, formulated as:

$$\mathcal{L}_{R} = \sum_{q=1}^{Q} \gamma_{q} \left\| \overline{\mathbf{z}}_{j,q}^{(t)} - \overline{\mathbf{z}}_{j,q}^{(t-1)} \right\|, \tag{7}$$

where γ_q denotes the weight on the regularization term of approximating the q-th frequency component of the old embedding space. To obtain a slow-updated space for old-knowledge preservation, we enforce higher γ_q on the low-frequency components in the regularization term for less forgetting and vice versa.

3.3 Inter-Space Level SvF Learning

Apart from the above intra-space SvF, we here introduce our inter-space SvF in detail. Previous works strive to maintain a uniform feature space that balances the old-knowledge slow-forgetting and new-knowledge fast-adapting. However, since the data of old tasks is non-available and the

number of new samples is few, a unified feature space is prone to overfit the new tasks, and the discriminability and generalization ability of the feature space easily degrades. As shown in Figure 4 (a), upon the arrival of new tasks, the samples of different classes which are separated at an early session, turn out to overlap with each other at later sessions. It indicates that the model suffers catastrophic forgetting after several sessions, and the samples which can be discriminated well at an early session are indistinguishable at the later session.

To this end, we propose to maintain another separate feature space and set up an inter-space SvF scheme to update the two feature spaces in different fashions. While the one updated in a fast fashion (e.g., training with a larger learning rate) is prone to new knowledge adaptation, the slowly updated one can better preserve the old knowledge throughout the learning process.

For classification, we propose a feature space composition operation to compose the above two spaces. After feature space composition, we can obtain a discriminative feature space, as shown in Figure 4 (b), where samples are clustered according to their classes and can be well separated. We use $\Psi(\cdot,\cdot)$ to denote the composition function (e.g., for a naive implementation, a simple concatenation operation). The composite feature $\tilde{\mathbf{z}}_j$ for sample x_j is denoted as $\tilde{\mathbf{z}}_j = \Psi(\mathbf{z}_j^{slow}, \mathbf{z}_j^{fast})$, where \mathbf{z}_j^{slow} denotes the embedding in the slow-updated feature space and \mathbf{z}_j^{fast} denotes that in the fast-updated feature space (trained at the current session). We conduct classification in this composite feature space, which is defined as:

$$\hat{y}_j = \underset{c \in \bigcup_i C^{(i)}}{\operatorname{argmin}} (\tilde{\mathbf{z}}_j - \tilde{\mathbf{u}}_c)^\top \mathbf{A} (\tilde{\mathbf{z}}_j - \tilde{\mathbf{u}}_c), \tag{8}$$

$$\tilde{\mathbf{u}}_c = \frac{1}{n_c} \sum_j [y_j = c] \cdot \Psi(\mathbf{u}_j^{slow}, \mathbf{u}_j^{fast}). \tag{9}$$

where **A** is a metric matrix. For a simple formulation, **A** can be a diagonal matrix and thereby indicating the importance of the slow-updated features and fast-updated features, and all features will be concerned equally if **A** is an identity matrix.

4 EXPERIMENTS AND RESULTS

4.1 Datasets

CIFAR100 [87] is a labeled subset of the 80 million tiny image dataset for object recognition. It contains $60000~32\times32$ RGB images in 100 classes, with 500 images per class for training and 100 images per class for testing. CUB200-2011 [88] contains 6000 images for training and 6000 for testing with the resolution of 256×256 , over 200 bird categories. It is originally designed for fine-grained image classification. MiniImageNet is the subset of ImageNet-1k [89] that is utilized by few-shot learning. It contains 100 classes. Each class contains $500~84\times84$ images for training and 100 images for testing.

4.2 Implementation Details

4.2.1 Evaluation Protocol

We conduct experiments under the FSCIL setting. Following [9], we evaluate our method on three datasets (CI-

FAR100, MiniImageNet, and CUB200) with similar evaluation protocols. For each $D^{(t)}, t>1$, if the number of classes $|C^{(t)}|$ is M and the number of training samples per class is K, we denote the setting as M-way K-shot. For CIFAR100 and miniImageNet datasets, we choose 60 and 40 classes for the base task and new tasks, respectively, and adopt the 5-way 5-shot setting, leading to 9 training sessions in total. For CUB200, we adopt the 10-way 5-shot setting, by picking 100 classes into 10 new learning sessions and the base task has 100 classes. For all datasets, we construct the training set of each learning session by randomly selecting 5 training samples per class from the original dataset, and the test set is the same as the original one. After training on a new batch of classes, we evaluate the trained model on test samples of all seen classes.

4.2.2 Training Details

We implement our models with Pytorch and use Adam [90] for optimization. Following [9], we use the ResNet18 [91] as the backbone network. The model is trained with the same strategy in [9] for a fair comparison. We use an embedding network as the feature extractor and the dimensions of the feature extractor are 512. The base model is trained with the same strategy in [9]. For inter-space SvF learning, we train our slow-updated models with learning rate 1e-6 for 50 epochs and train our fast-updated models with learning rate 1e-5 for 50 epochs. We follow the setting discussed in the part "Inter-space SvF analysis" of Section 4.4 and choose the scalar a = 0.5 as default for the feature space composition operation. For intra-space SvF learning, we obtain 512 frequency components from the original feature by DCT, and the overall spectrum is divided into 8 groups. The DCT operation is implemented with torch-dct project on github. The time for training a model on MiniImageNet only increases by around 2% (memory 4%), which are almost negligible. We conduct frequency-aware regularization by adjusting the weights of these 8 frequency groups (i.e. $\{\gamma_q\}_{q=1}^8$). For the slow-updated models, when updating on each new task, we set the weight $\gamma_q = 1$ if q = 1 otherwise $\gamma_q = 0$. For the fast-updated models, we set the weight $\gamma_q = 0$ if q = 1 otherwise $\gamma_q = 1$. We compute the centers of old classes after training one task and fix them in following tasks, which is standard as SDC [1]. For data augmentation, we use standard random cropping and flipping as in [9].

4.3 Comparison to State-of-the-Art Methods

In this section, we compare our proposed method in the few-shot class-incremental learning scenario, with existing state-of-the-art methods. They include iCaRL [15], EEIL [16], LUCIR [17], TOPIC [9], SDC [1], POD [4]. We report the results of iCaRL and LUCIR with both a softmax classifier and an NCM classifier, which are denoted as iCaRL-CNN and iCaRL-NCM (as well as LUCIR-CNN and LUCIR-NCM) respectively. Using different classifiers achieves different accuracy at the first task (the base task). Table 2 reports the test accuracy of different methods on CUB-200 dataset. Compared with the others, our method shows clear superiority (more than 10%). As shown in Figure 5 (a) and Figure 5 (b), we can observe that our method outperforms all other methods at each encountered learning session on CIFAR100 and MiniImageNet datasets.

TABLE 2
Comparison results on CUB200 with ResNet18 using the 10-way 5-shot FSCIL setting. Our method outperforms others at all learning sessions.

Method	learning sessions								our relative			
Method	1	2	3	4	5	6	7	8	9	10	11	improvements
iCaRL-CNN [15]	68.68	52.65	48.61	44.16	36.62	29.52	27.83	26.26	24.01	23.89	21.16	+33.17
EEIL [16]	68.68	53.63	47.91	44.20	36.30	27.46	25.93	24.70	23.95	24.13	22.11	+32.22
LUCIR-CNN [17]	68.68	57.12	44.21	28.78	26.71	25.66	24.62	21.52	20.12	20.06	19.87	+34.46
TOPIC [9]	68.68	62.49	54.81	49.99	45.25	41.40	38.35	35.36	32.22	28.31	26.28	+28.05
iCaRL-NCM [15]	72.29	56.67	50.76	43.29	40.99	34.07	30.01	28.83	26.56	23.76	23.32	+31.01
LUCIR-NCM [17]	72.29	58.70	50.68	49.82	45.59	43.10	34.77	31.35	28.53	25.73	22.91	+31.42
SDC [1]	72.29	68.22	61.94	61.32	59.83	57.30	55.48	54.20	49.99	48.85	42.58	+11.75
POD [4]	72.29	59.77	51.23	48.78	47.83	44.22	39.76	37.79	35.23	31.92	31.27	+23.06
Ours	72.29	70.53	67.00	64.92	62.67	61.89	59.63	59.15	57.73	55.92	54.33	

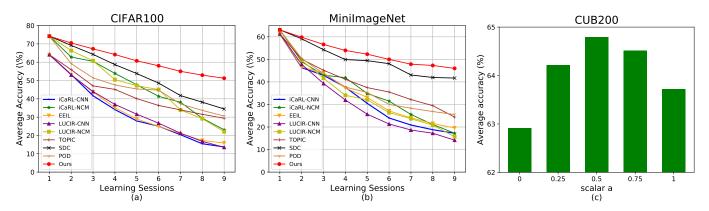


Fig. 5. (a): Comparison results on CIFAR100 with ResNet18 using the 5-way 5-shot FSCIL setting. (b): Comparison results on MiniImageNet with ResNet18 using the 5-way 5-shot FSCIL setting. Our method shows clear superiority and outperforms all other methods at each encountered learning session. (c): Change of average accuracy when varying a on CUB200. The performance peaks on an intermediate value, which indicates the importance and complementarity of both the slow-updated space and the fast-updated space.

TABLE 3

Validation of our intra-space frequency-aware knowledge distillation and inter-space composition operation on CUB200, CIFAR100 and MiniImageNet. The performance at the last learning session (i.e., Last) and the average results over all the learning sessions (i.e., Average) are reported here.

Method	Cī	JB200	CIF	AR100	MiniImageNet		
Metriod	Last	Average	Last	Average	Last	Average	
Baseline	42.58	57.45	34.46	53.68	41.71	49.96	
Baseline+intra-space Baseline+inter-space	49.96 53.16	59.70 60.68	50.60 49.11	61.09 60.41	44.46 44.72	51.93 51.90	

4.4 Ablation Study

In this section, we carry out ablation experiments to validate our inter-space grain and intra-space grain SvF strategies. Extensive experiments are conducted to show the effect of training samples' number, feature space composition method, etc. We also explore the properties of fast-updated space, slow-updated space, as well as different frequency components when transferring old knowledge and further elaborated on the inter-space and intra-space SvF analysis, respectively.

4.4.1 Baseline

As described in Section 3.1, we follow SDC [1] and implement our method with an embedding network and an NCM classifier. We adopt the triplet loss [86] as the metric loss, and conduct knowledge distillation in the embedding space [62] to retain old knowledge.

4.4.2 Effect of Intra-space and Inter-space SvF Strategies

We first conduct ablation experiments to show the effectiveness of our inter-space and intra-space SvF strategies, respectively. As shown in Table 3, both the inter-space grain and intra-space grain SvF strategies improve the performance of our baseline. Inter-space SvF strategy achieves better performance than the intra-space SvF strategy on CUB200 (around 1%) which only contains the samples of bird categories. On the datasets which contain more diverse classes (e.g., CIFAR100 and MiniImageNet), the intra-space SvF strategy achieves comparable performance to the interspace SvF strategy. For inter-space SvF, we also evaluate the performance of only using the slow-updated feature space alone for classification. As shown in Figure 6, a large performance gap on the data of the current task exists between the slow-updated feature space and the fast-updated one, which shows the importance of the fast-updated space to fit new knowledge. It indicates the reasonability of inter-space SvF. In this way, it is difficult to keep a good balance between

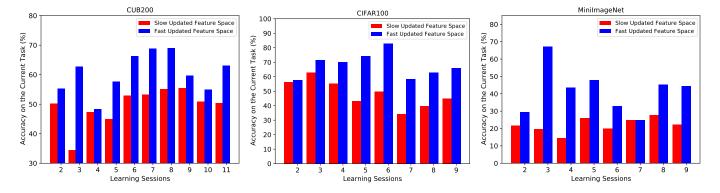


Fig. 6. Accuracy for the current task while only utilizing the slow-updated feature space or the fast-updated feature space on CUB200, CIFAR100, and MinilmageNet. A large performance gap on the data of the current task exists between the slow-updated space and fast-updated space.

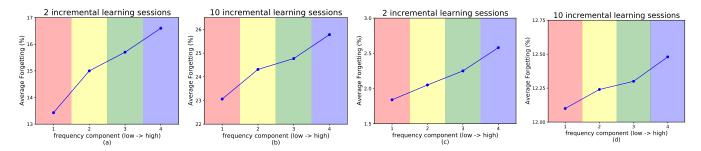


Fig. 7. Analysis of intra-space SvF. (a) Results with 2 learning sessions on CUB200. (b) Results with 10 learning sessions on CUB200. (c) Results with 2 learning sessions on MinilmageNet. (d) Results with 10 learning sessions on MinilmageNet.

TABLE 4
Results with different inter-space composition operations on CUB200, CIFAR100, and MiniImageNet. The sophisticated composition operation implemented with PCA outperforms the simple version by around 1%.

Method	Cī	JB200	CIF	AR100	MiniImageNet		
Wettod	Last	Average	Last	Average	Last	Average	
Baseline	42.58	57.45	34.46	53.68	41.71	49.96	
Baseline+inter-space-simple Baseline+inter-space-pca	53.16 53.76	60.68 61.15	49.11 50.51	60.41 60.78	44.72 45.38	51.90 53.24	

slow-forgetting for old knowledge and fast-adapting for new knowledge by adjusting the updating fashion of a unified feature space.

4.4.3 Intra-space SvF Analysis

We here introduce our intra-space SvF analysis, examining the effect of different frequency components for old-knowledge preservation. We divide the overall spectrum into four groups and regularize the model with one of them separately. It can be observed in Figure 1 that different groups appear different characteristics for old-knowledge transferring, and a pretty clear trend can be found that the forgetting rate of old tasks increases along with the frequency. It indicates that the regularization of low-frequency components is more conducive to preserving old knowledge than that on high-frequency components. Therefore, we enforce higher weights on those low-frequency components in the regularization term to make these knowledge components of the feature space updating more slowly. The results on CUB200 and MiniImageNet are shown in Figure 7.

4.4.4 Inter-Space SvF analysis

We here briefly show that the fast-updated space and the slow-updated space are complementary, even with the simplest implementation strategy for space composition, where

 $a\mathbf{I}$ A is constructed as A =with a scalar a $(1-a){\bf I}$ (I is an identity matrix with dimension half of A's). a = 0means that only using the slow-updated feature space and a = 1 for only using the fast-updated feature space at the current session. The change of accuracy, with respect to a, is shown in Figure 5(c). Using a slow-updated feature space achieves the lowest accuracy since it contains limited newtask knowledge. The performance of using a fast-updated feature space independently is also lower than that of the composite space, because of forgetting old knowledge. The performance peaks on an intermediate value, which indicates the complementarity. More sophisticated forms of the metric matrix A can also be constructed (e.g., in data-driven learning), and the discussion and analysis of another space composition strategy are detailed in the next paragraph.

4.4.5 Different Inter-space Composition Operations

We evaluate another more sophisticated inter-space composition method which first reduces the dimensions of the features from the fast-updated space and slow-updated space by principal component analysis (PCA), and then concatenates them. The results are shown in Table 4, where the simplest composition strategy is denoted by "inter-space-simple", and the sophisticated one is denoted by "inter-space-pca". Specifically, we compute the PCA transformation matrix \mathbf{P}_1 and \mathbf{P}_2 with embeddings of samples obtained by the slowly updated model and the fast updated

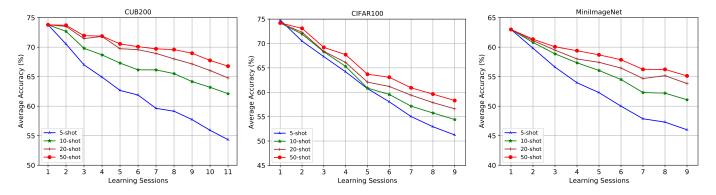


Fig. 8. Results on CUB200, CIFAR100, and MiniImageNet under different FSCIL settings (5-shot, 10-shot, 20-shot, and 50-shot). The performance of our method increases as the number of training samples increases.

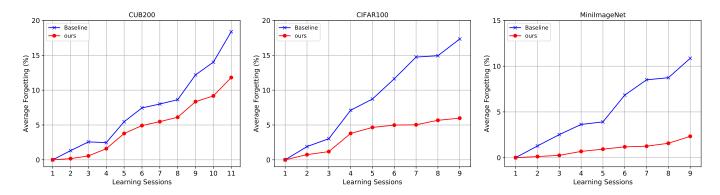


Fig. 9. Average forgetting of our method and the baseline with a unified feature space on CUB200, CIFAR100, and MiniImageNet. Our method outperforms the baseline by a large margin at the last learning session.

model respectively. Another choice is learning \mathbf{P}_1 and \mathbf{P}_2 incrementally [92], [93] with samples from the subsequent tasks. Then we reduce the dimension to, e.g., half of the original, and therefore the size of \mathbf{P}_1 and \mathbf{P}_2 are both 256×512 . For concatenating features in the composite space, the transformation matrix can be constructed as $\mathbf{Q} = \begin{bmatrix} \mathbf{P}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{P}_2 \end{bmatrix}$. Given a feature $\tilde{\mathbf{z}}_j$, and a center $\tilde{\mathbf{u}}_c$ for class c, the dimension-reduced data can be $\mathbf{Q}\tilde{\mathbf{z}}_j$ and $\mathbf{Q}\tilde{\mathbf{u}}_c$, respectively. And the classification can be formulated as:

$$\hat{y}_{j} = \underset{c \in \bigcup_{i} C^{(i)}}{\operatorname{argmin}} (\mathbf{Q}\tilde{\mathbf{z}}_{j} - \mathbf{Q}\tilde{\mathbf{u}}_{c})^{\top} (\mathbf{Q}\tilde{\mathbf{z}}_{j} - \mathbf{Q}\tilde{\mathbf{u}}_{c})$$

$$= \underset{c \in [], C^{(i)}}{\operatorname{argmin}} (\tilde{\mathbf{z}}_{j} - \tilde{\mathbf{u}}_{c})^{\top} \mathbf{Q}^{\top} \mathbf{Q} (\tilde{\mathbf{z}}_{j} - \tilde{\mathbf{u}}_{c}).$$
(10)

Compared with Equation (8), matrix $\mathbf{Q}^{\top}\mathbf{Q}$ can be considered as a specific data-driven metric matrix \mathbf{A} . Also, the transformation matrix \mathbf{Q} can be viewed as a part of composition function $\Psi(\cdot,\cdot)$. We set $\mathbf{P}_1 = \mathbf{P}_2$ in the fewshot scenario because such few new-task samples are not able to estimate a reasonable \mathbf{P}_2 . As shown in Table 4, "Baseline+inter-space-pca" achieves 1% higher accuracy.

4.4.6 The Effect of The Number of Training Samples

To examine the effect of the number of training samples, we evaluate our method with different shots of training samples, which are 5-shot, 10-shot, 20-shot, and 50-shot settings. As shown in Figure 8, we can see that the performance of our method increases as the number of training

samples increases. It can also be noticed that the number of samples matters more on latter training sessions, and the performance gap grows more rapidly when the number of samples gets larger.

4.4.7 Average Forgetting of Previous Tasks

The forgetting of previous tasks is estimated with *average* forgetting [1], [10]. We illustrate the forgetting curves of our method across 11 learning sessions on CUB200 (9 learning sessions on CIFAR100 and MiniImageNet), shown in Figure 9. On these three datasets, we can observe that our method outperforms "Baseline" by a large margin at the last learning session (more than 5%). These results indicate the stability of our composite feature space against the continuous arrival of new tasks.

4.4.8 Analysis of How to Choose Hyperparameters

The hyperparameters (i.e. gammas and learning rates) are chosen according to the following rules: 1) for gammas (e.g. $\{\gamma_q\}_{q=1}^8$ correspond to 8 frequency groups from low-frequency to high-frequency), we set higher weights $\gamma_i=1$ on low-frequency groups (the first group) and lower weights $\gamma_j=0$ on high-frequency groups (the other groups) when training slow-updated model, and vice versa for fast-updated one. Considering the first group as the low-frequency group usually leads to best performance. 2) for learning rates, the learning rate of fast-updated model is usually 10 times than that of slow-updated one. In our

experiments, we choose 1e-5 and 1e-6 for fast-updated and slow-updated model respectively.

We have also evaluated our method with different numbers of frequency groups (i.e. $N_Q=2,4,8,16$) on MiniImageNet. The average results over all learning sessions are 52.76%, 53.03%, 53.24%, 52.78% respectively. Our method is stable as the number of frequency groups changes and achieves the best performance when $N_Q=8$. In our experiments, we choose $N_Q=8$ as the number of frequency groups.

5 CONCLUSION

In this paper, we propose a novel few-shot class-incremental learning scheme based on a "slow vs. fast" (SvF) learning pipeline. In the pipeline, we design an intra-space frequency-aware regularization to enforce a SvF constraint on different frequency components, and present an interspace composition operation to well balance the slow forgetting of old knowledge and fast adaptation to new knowledge. Comprehensive experimental results demonstrate that the proposed approach significantly outperforms other state-of-the-art approaches by a large margin.

ACKNOWLEDGMENT

The authors would like to thank Xuewei Li, Songyuan Li and Hui Wang for their valuable comments and suggestions.

REFERENCES

- [1] L. Yu, B. Twardowski, X. Liu, L. Herranz, K. Wang, Y. Cheng, S. Jui, and J. van de Weijer, "Semantic drift compensation for class-incremental learning," in *Proceedings of the IEEE conference* on computer vision and pattern recognition (CVPR), 2020.
- [2] X. Tao, X. Chang, X. Hong, X. Wei, and Y. Gong, "Topology-preserving class-incremental learning." ECCV, 2020.
 [3] A. Iscen, J. Zhang, S. Lazebnik, and C. Schmid, "Memory-efficient
- [3] A. Iscen, J. Zhang, S. Lazebnik, and C. Schmid, "Memory-efficient incremental learning through feature adaptation." ECCV, 2020.
- [4] A. Douillard, M. Cord, C. Ollion, T. Robert, and E. Valle, "Small-task incremental learning." ECCV, 2020.
- [5] B. Zhao, X. Xiao, G. Gan, B. Zhang, and S.-T. Xia, "Maintaining discrimination and fairness in class incremental learning," in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020, pp. 13 208–13 217.
- [6] J. Rajasegaran, S. Khan, M. Hayat, F. S. Khan, and M. Shah, "itaml: An incremental task-agnostic meta-learning approach," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 13588–13597.
- [7] T. L. Hayes, K. Kafle, R. Shrestha, M. Acharya, and C. Kanan, "Remind your neural network to prevent catastrophic forgetting." ECCV, 2020.
- [8] M. Masana, X. Liu, B. Twardowski, M. Menta, A. D. Bagdanov, and J. van de Weijer, "Class-incremental learning: survey and performance evaluation," arXiv preprint arXiv:2010.15277, 2020.
- [9] X. Tao, X. Hong, X. Chang, S. Dong, X. Wei, and Y. Gong, "Few-shot class-incremental learning," in *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, 2020.
- [10] A. Chaudhry, P. K. Dokania, T. Ajanthan, and P. H. Torr, "Riemannian walk for incremental learning: Understanding forgetting and intransigence," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 532–547.
- [11] R. M. French, "Catastrophic forgetting in connectionist networks," Trends in cognitive sciences, vol. 3, no. 4, pp. 128–135, 1999.
- [12] I. J. Goodfellow, M. Mirza, D. Xiao, A. Courville, and Y. Bengio, "An empirical investigation of catastrophic forgetting in gradient-based neural networks," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2014.

- [13] M. McCloskey and N. J. Cohen, "Catastrophic interference in connectionist networks: The sequential learning problem," in *Psychology of learning and motivation*. Elsevier, 1989, vol. 24, pp. 109– 165.
- [14] B. Pfülb and A. Gepperth, "A comprehensive, application-oriented study of catastrophic forgetting in dnns," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.
- [15] S.-A. Rebuffi, A. Kolesnikov, G. Sperl, and C. H. Lampert, "icarl: Incremental classifier and representation learning," in *Proceedings* of the IEEE conference on computer vision and pattern recognition (CVPR), 2017.
- [16] F. M. Castro, M. J. Marín-Jiménez, N. Guil, C. Schmid, and K. Alahari, "End-to-end incremental learning," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 233–248.
- [17] S. Hou, X. Pan, C. C. Loy, Z. Wang, and D. Lin, "Learning a unified classifier incrementally via rebalancing," in *Proceedings of the IEEE* conference on computer vision and pattern recognition (CVPR), 2019.
- [18] M. De Lange, R. Aljundi, M. Masana, S. Parisot, X. Jia, A. Leonardis, G. Slabaugh, and T. Tuytelaars, "A continual learning survey: Defying forgetting in classification tasks," arXiv preprint arXiv:1909.08383, 2019.
- [19] G. I. Parisi, R. Kemker, J. L. Part, C. Kanan, and S. Wermter, "Continual lifelong learning with neural networks: A review," Neural Networks, 2019.
- [20] Y. Li, L. Zhao, K. Church, and M. Elhoseiny, "Compositional continual language learning," in *Proceedings of the International* Conference on Learning Representations (ICLR), 2020.
- [21] S. Lee, J. Ha, D. Zhang, and G. Kim, "A neural dirichlet process mixture model for task-free continual learning," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020.
- [22] T. Adel, H. Zhao, and R. E. Turner, "Continual learning with adaptive weights (claw)," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020.
- [23] J. von Oswald, C. Henning, J. Sacramento, and B. F. Grewe, "Continual learning with hypernetworks," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020.
- [24] R. Kurle, B. Cseke, A. Klushyn, P. van der Smagt, and S. Günnemann, "Continual learning with bayesian neural networks for non-stationary data."
- [25] M. K. Titsias, J. Schwarz, A. G. d. G. Matthews, R. Pascanu, and Y. W. Teh, "Functional regularisation for continual learning with gaussian processes," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.
- [26] S. Ebrahimi, M. Elhoseiny, T. Darrell, and M. Rohrbach, "Uncertainty-guided continual learning with bayesian neural networks," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020.
- [27] G. Zeng, Y. Chen, B. Cui, and S. Yu, "Continual learning of context-dependent processing in neural networks," *Nature Machine Intelligence*, vol. 1, no. 8, pp. 364–372, 2019.
- [28] J. N. Kundu, R. M. Venkatesh, N. Venkat, A. Revanur, and R. V. Babu, "Class-incremental domain adaptation," 2020.
- [29] O. Ostapenko, M. Puscas, T. Klein, P. Jahnichen, and M. Nabi, "Learning to remember: A synaptic plasticity driven framework for continual learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 11321–11329.
- [30] R. Aljundi, K. Kelchtermans, and T. Tuytelaars, "Task-free continual learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 11254–11263.
- [31] J. Lee, H. G. Hong, D. Joo, and J. Kim, "Continual learning with extended kronecker-factored approximate curvature," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9001–9010.
- [32] J. He, R. Mao, Z. Shao, and F. Zhu, "Incremental learning in online scenario," in *Proceedings of the IEEE/CVF Conference on Computer* Vision and Pattern Recognition, 2020, pp. 13926–13935.
- [33] S. Ebrahimi, F. Meier, R. Calandra, T. Darrell, and M. Rohrbach, "Adversarial continual learning." ECCV, 2020.
- [34] Y. Liu, S. Parisot, G. Slabaugh, X. Jia, A. Leonardis, and T. Tuytelaars, "More classifiers, less forgetting: A generic multi-classifier paradigm for incremental learning." ECCV, 2020.
- [35] M. Caccia, P. Rodriguez, O. Ostapenko, F. Normandin, M. Lin, L. Page-Caccia, I. H. Laradji, I. Rish, A. Lacoste, D. Vázquez et al., "Online fast adaptation and knowledge accumulation (osaka): a new approach to continual learning," Advances in Neural Information Processing Systems, vol. 33, 2020.

- [36] C. V. Nguyen, Y. Li, T. D. Bui, and R. E. Turner, "Variational continual learning," arXiv preprint arXiv:1710.10628, 2017.
- [37] J. Chen, S. Wang, L. Chen, H. Cai, and Y. Qian, "Incremental detection of remote sensing objects with feature pyramid and knowledge distillation," *IEEE Transactions on Geoscience and Remote Sensing*, 2020.
- [38] S. Dang, Z. Cao, Z. Cui, Y. Pi, and N. Liu, "Class boundary exemplar selection based incremental learning for automatic target recognition," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 58, no. 8, pp. 5782–5792, 2020.
- [39] D. Shim, Z. Mai, J. Jeong, S. Sanner, H. Kim, and J. Jang, "Online class-incremental continual learning with adversarial shapley value," arXiv e-prints, pp. arXiv–2009, 2020.
- [40] L. Liu, Z. Kuang, Y. Chen, J.-H. Xue, W. Yang, and W. Zhang, "Incdet: in defense of elastic weight consolidation for incremental object detection," *IEEE transactions on neural networks and learning* systems, 2020.
- [41] H. Li, P. Barnaghi, S. Enshaeifar, and F. Ganz, "Continual learning using bayesian neural networks," *IEEE Transactions on Neural Networks and Learning Systems*, 2020.
- [42] J. Yoon, E. Yang, J. Lee, and S. J. Hwang, "Lifelong learning with dynamically expandable networks," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018.
- [43] X. Li, Y. Zhou, T. Wu, R. Socher, and C. Xiong, "Learn to grow: A continual structure learning framework for overcoming catastrophic forgetting," *International Conference on Machine Learning* (ICML), 2019.
- [44] C.-Y. Hung, C.-H. Tu, C.-E. Wu, C.-H. Chen, Y.-M. Chan, and C.-S. Chen, "Compacting, picking and growing for unforgetting continual learning," in *Advances in Neural Information Processing* Systems, 2019, pp. 13 647–13 657.
- [45] A. Mallya and S. Lazebnik, "Packnet: Adding multiple tasks to a single network by iterative pruning," in *Proceedings of the IEEE* conference on computer vision and pattern recognition (CVPR), 2018.
- [46] A. Mallya, D. Davis, and S. Lazebnik, "Piggyback: Adapting a single network to multiple tasks by learning to mask weights," in Proceedings of the European Conference on Computer Vision (ECCV), 2018, pp. 67–82.
- [47] J. Serrà, D. Surís, M. Miron, and A. Karatzoglou, "Overcoming catastrophic forgetting with hard attention to the task," arXiv preprint arXiv:1801.01423, 2018.
- [48] A. A. Rusu, N. C. Rabinowitz, G. Desjardins, H. Soyer, J. Kirk-patrick, K. Kavukcuoglu, R. Pascanu, and R. Hadsell, "Progressive neural networks," arXiv preprint arXiv:1606.04671, 2016.
- [49] R. Aljundi, P. Chakravarty, and T. Tuytelaars, "Expert gate: Lifelong learning with a network of experts," in *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, 2017.
- [50] J. Rajasegaran, M. Hayat, S. H. Khan, F. S. Khan, and L. Shao, "Random path selection for continual learning," in *Advances in Neural Information Processing Systems*, 2019, pp. 12 648–12 658.
- [51] D. Abati, J. Tomczak, T. Blankevoort, S. Calderara, R. Cucchiara, and B. E. Bejnordi, "Conditional channel gated networks for taskaware continual learning," in *Proceedings of the IEEE/CVF Confer*ence on Computer Vision and Pattern Recognition, 2020, pp. 3931– 3940.
- [52] A. Chaudhry, M. Ranzato, M. Rohrbach, and M. Elhoseiny, "Efficient lifelong learning with a-gem," in Proceedings of the International Conference on Learning Representations (ICLR), 2019.
- [53] D. Lopez-Paz and M. Ranzato, "Gradient episodic memory for continual learning," in Advances in Neural Information Processing Systems, 2017, pp. 6467–6476.
- [54] H. Shin, J. K. Lee, J. Kim, and J. Kim, "Continual learning with deep generative replay," in *Advances in Neural Information Process*ing Systems, 2017, pp. 2990–2999.
- [55] R. Aljundi, E. Belilovsky, T. Tuytelaars, L. Charlin, M. Caccia, M. Lin, and L. Page-Caccia, "Online continual learning with maximal interfered retrieval," in *Advances in Neural Information Processing Systems*, 2019, pp. 11849–11860.
- [56] M. Zhai, L. Chen, F. Tung, J. He, M. Nawhal, and G. Mori, "Life-long gan: Continual learning for conditional image generation," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 2759–2768.
- [57] C. Wu, L. Herranz, X. Liu, J. van de Weijer, B. Raducanu et al., "Memory replay gans: Learning to generate new categories without forgetting," in Advances In Neural Information Processing Systems, 2018, pp. 5962–5972.

- [58] C. He, R. Wang, S. Shan, and X. Chen, "Exemplar-supported generative reproduction for class incremental learning." in BMVC, 2018, p. 98.
- [59] Y. Wu, Y. Chen, L. Wang, Y. Ye, Z. Liu, Y. Guo, and Y. Fu, "Large scale incremental learning," in *Proceedings of the IEEE conference on* computer vision and pattern recognition (CVPR), 2019.
- [60] Y. Liu, A.-A. Liu, Y. Su, B. Schiele, and Q. Sun, "Mnemonics training: Multi-class incremental learning without forgetting," arXiv preprint arXiv:2002.10211, 2020.
- [61] S.-W. Lee, J.-H. Kim, J. Jun, J.-W. Ha, and B.-T. Zhang, "Overcoming catastrophic forgetting by incremental moment matching," in *Advances in Neural Information Processing Systems*, 2017, pp. 4652–4662.
- [62] Z. Li and D. Hoiem, "Learning without forgetting," IEEE transactions on pattern analysis and machine intelligence, vol. 40, no. 12, pp. 2935–2947, 2017.
- [63] H. Ritter, A. Botev, and D. Barber, "Online structured laplace approximations for overcoming catastrophic forgetting," in Advances in Neural Information Processing Systems, 2018, pp. 3738–3748.
- [64] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska et al., "Overcoming catastrophic forgetting in neural networks," Proceedings of the national academy of sciences, vol. 114, no. 13, pp. 3521–3526, 2017.
- [65] F. Zenke, B. Poole, and S. Ganguli, "Continual learning through synaptic intelligence," *International Conference on Machine Learning* (ICML), 2017.
- [66] X. Liu, M. Masana, L. Herranz, J. Van de Weijer, A. M. Lopez, and A. D. Bagdanov, "Rotate your networks: Better weight consolidation and less catastrophic forgetting," in 2018 24th International Conference on Pattern Recognition (ICPR). IEEE, 2018, pp. 2262– 2268.
- [67] R. Aljundi, F. Babiloni, M. Elhoseiny, M. Rohrbach, and T. Tuyte-laars, "Memory aware synapses: Learning what (not) to forget," in Proceedings of the European Conference on Computer Vision (ECCV), 2018, pp. 139–154.
- [68] P. Dhar, R. V. Singh, K.-C. Peng, Z. Wu, and R. Chellappa, "Learning without memorizing," in *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, 2019.
- [69] S. I. Mirzadeh, M. Farajtabar, R. Pascanu, and H. Ghasemzadeh, "Understanding the role of training regimes in continual learning," Advances in Neural Information Processing Systems, vol. 33, 2020.
- [70] M. Kanakis, D. Bruggemann, S. Saha, S. Georgoulis, A. Obukhov, and L. Van Gool, "Reparameterizing convolutions for incremental multi-task learning without task interference," arXiv preprint arXiv:2007.12540, 2020.
- [71] A. Rosenfeld and J. K. Tsotsos, "Incremental learning through deep adaptation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 42, no. 3, pp. 651–663, 2018.
- [72] J.-M. Perez-Rua, X. Zhu, T. M. Hospedales, and T. Xiang, "Incremental few-shot object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 13.846–13.855
- [73] F. Cermelli, M. Mancini, S. R. Bulo, E. Ricci, and B. Caputo, "Modeling the background for incremental learning in semantic segmentation," in *Proceedings of the IEEE/CVF Conference on Com*puter Vision and Pattern Recognition, 2020, pp. 9233–9242.
- [74] E. Belouadah and A. Popescu, "Deesil: Deep-shallow incremental learning." in Proceedings of the European Conference on Computer Vision (ECCV), 2018, pp. 0–0.
- [75] Y. Xiang, Y. Fu, P. Ji, and H. Huang, "Incremental learning using conditional adversarial networks," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 6619–6628.
- [76] H. Su, J. Su, D. Wang, W. Gan, W. Wu, M. Wang, J. Yan, and Y. Qiao, "Collaborative distillation in the parameter and spectrum domains for video action recognition," arXiv preprint arXiv:2009.06902, 2020.
- [77] Z. Wang, Y. Yang, A. Shrivastava, V. Rawal, and Z. Ding, "To-wards frequency-based explanation for robust cnn," arXiv preprint arXiv:2005.03141, 2020.
- [78] L. Gueguen, A. Sergeev, B. Kadlec, R. Liu, and J. Yosinski, "Faster neural networks straight from jpeg," in *Advances in Neural Infor*mation Processing Systems, 2018, pp. 3933–3944.
- [79] K. Xu, M. Qin, F. Sun, Y. Wang, Y.-K. Chen, and F. Ren, "Learning in the frequency domain," in *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition, 2020, pp. 1740–1749.

- [80] J. Kim, S. Cha, D. Wee, S. Bae, and J. Kim, "Regularization on spatio-temporally smoothed feature for action recognition," in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020, pp. 12103–12112.
- [81] Y. Yang and S. Soatto, "Fda: Fourier domain adaptation for semantic segmentation," in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020, pp. 4085–4095.
- [82] P. Khorramshahi, N. Peri, J.-c. Chen, and R. Chellappa, "The devil is in the details: Self-supervised attention for vehicle reidentification," arXiv preprint arXiv:2004.06271, 2020.
- [83] R. Kemker and C. Kanan, "Fearnet: Brain-inspired model for incremental learning," in *Proceedings of the International Conference* on Learning Representations (ICLR), 2018.
- [84] J. Bromley, I. Guyon, Y. LeCun, E. Säckinger, and R. Shah, "Signature verification using a" siamese" time delay neural network," in Advances in neural information processing systems, 1994, pp. 737–744.
- [85] T. Mensink, J. Verbeek, F. Perronnin, and G. Csurka, "Distance-based image classification: Generalizing to new classes at near-zero cost," *IEEE transactions on pattern analysis and machine intelligence*, vol. 35, no. 11, pp. 2624–2637, 2013.
- [86] J. Wang, Y. Song, T. Leung, C. Rosenberg, J. Wang, J. Philbin, B. Chen, and Y. Wu, "Learning fine-grained image similarity with deep ranking," in *Proceedings of the IEEE conference on computer* vision and pattern recognition (CVPR), 2014.
- [87] A. Krizhevsky and G. Hinton, "Learning multiple layers of features from tiny images," Citeseer, Tech. Rep., 2009.
- [88] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie, "The caltech-ucsd birds-200-2011 dataset. technical report cns-tr-2011-001," California Institute of Technology, 2011.
- [89] O. Vinyals, C. Blundell, T. Lillicrap, D. Wierstra et al., "Matching networks for one shot learning," in Advances in Neural Information Processing Systems, 2016, pp. 3630–3638.
- Processing Systems, 2016, pp. 3630–3638.

 [90] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in Proceedings of the International Conference on Learning Representations (ICLR), 2014.
- [91] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, 2016.
- [92] P. M. Hall, A. D. Marshall, and R. R. Martin, "Incremental eigenanalysis for classification." in *BMVC*, vol. 98. Citeseer, 1998, pp. 286–295.
- [93] J. Weng, Y. Zhang, and W.-S. Hwang, "Candid covariance-free incremental principal component analysis," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 25, no. 8, pp. 1034– 1040, 2003.