

Continual Learning for Visual Search with Backward Consistent Feature Embedding

Timmy S. T. Wan¹

Jun-Cheng Chen²

Tzer-Yi Wu³

Chu-Song Chen^{1,*}

National Taiwan University¹

Academia Sinica²

ucfunnel Co. Ltd.³

{r08944004, chusong}@csie.ntu.edu.tw, pullpull@citi.sinica.edu.tw, kenny.wu@ucfunnel.com

Abstract

In visual search, the gallery set could be incrementally growing and added to the database in practice. However, existing methods rely on the model trained on the entire dataset, ignoring the continual updating of the model. Besides, as the model updates, the new model must re-extract features for the entire gallery set to maintain compatible feature space, imposing a high computational cost for a large gallery set. To address the issues of long-term visual search, we introduce a continual learning (CL) approach that can handle the incrementally growing gallery set with backward embedding consistency. We enforce the losses of inter-session data coherence, neighbor-session model coherence, and intra-session discrimination to conduct a continual learner. In addition to the disjoint setup, our CL solution also tackles the situation of increasingly adding new classes for the blurry boundary without assuming all categories known in the beginning and during model update. To our knowledge, this is the first CL method both tackling the issue of backward-consistent feature embedding and allowing novel classes to occur in the new sessions. Extensive experiments on various benchmarks show the efficacy of our approach under a wide range of setups¹.

1. Introduction

Continual learning (CL) aims to learn new tasks while keeping the functions learned from the old sessions. The technology has been rapidly evolving; nevertheless, the active research area in CL focuses on image classification but ignores the demand for image retrieval (aka visual search). For obtaining powerful feature representations in image retrieval, most works [37, 44, 57, 67] still require a model to be trained on an entire dataset simultaneously instead of in an incremental manner. However, a practical visual search system should be capable of continually learning from new materials while consolidating the old knowledge to cope with the data accumulated with time.

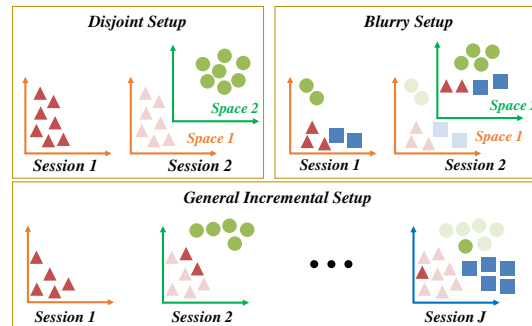


Figure 1. Illustration of the proposed approach with the *general incremental setup*. Our solution allows the new gallery set of seen and unseen classes to be freely and incrementally added to the database with respect to the widely adopted *disjoint* and the recent *blurry setups*. In addition, it also considers backward compatible embedding for a session sequence. This avoids the gallery embeddings of old and new sessions from being separated in incompatible feature spaces. Thus, our approach is more practical for real retrieval applications. Semi-transparent icons represent the data points from the previous sessions collected so far.

As data grows, despite updating the model by simply fine-tuning, many observations [11, 69] reveal that catastrophic forgetting happens. A series of strategies have been developed in CL to address the problem [11, 29, 43]. The methods can make a single deep model capable of updating itself successively while avoiding disappointing overall performance. However, there are still several ongoing issues.

First, many works in CL emphasize the disjoint setup where the data from the old class will not show during training in a new task (or session). The task boundary arising over classes restricts the usage of CL since many retrieval systems need to collect extra data of the seen labels for improving their models in new sessions. Although recent studies (e.g. [2, 4]) allow the class overlapping among the tasks, the blurry setup in these works assumes that all the class labels in the future sessions are pre-given in advance; only the instance ratios in the classes vary with the session (Fig. 1). This kind of data scenario is impractical for most visual-search applications (e.g., in an e-commerce system, the new

* indicates corresponding author.

¹Code: <https://github.com/ivclab/CVS>

product arrives over season). Moving toward a more general setup is desirable to fulfill real-world scenarios.

Second, a model updated from new data will deploy online in retrieval. An essential step is to re-extract the feature embedding from previous gallery images to maintain a consistent feature space on the pairwise distance measurement. For visual search on large-scale data, feature re-extraction is computationally intensive. Thus, an ideal design for continual visual search is that an updated model only extracts features for incoming gallery data while keeping the previously generated feature representations unchanged. However, it leads to the further difficulty of pairwise similarity measurement in uneven feature spaces. Motivated by this, we argue that current CL studies lack consideration for feature compatibility between the ongoing and previous data. Hence, designing a CL algorithm with backward consistent feature embedding is demanded.

To address the above issues, we introduce a novel CL approach, namely, CVS (Continual-learner for Visual Search), for a generally incremental setup of visual search. CVS can learn effective feature representations with backward consistency. For learning new knowledge, our learner obtains discriminating features for the current task. We introduce a cross-task gallery embedding consistency constraint that keeps the currently learned features compatible with the representatives from the obsolete gallery features. For consolidating old knowledge while keeping feature consistency, we develop a metric-based knowledge distillation to draw together the embedding from the different feature spaces. By coordinating the components, CVS can achieve continual visual search with backward-compatible feature embedding effectively. Also, we conduct extensive experiments under various incremental data distributions, especially for the general-incremental setup, to validate the effectiveness of CVS. The main characteristics include:

General Incremental Setup: We introduce a new CL scenario to simulate real-world visual search application systems. It includes the previous disjoint and blurry setups as special cases and tackles the general setting that the classes in a coming session can be either seen or unseen before.

Backward Consistent Feature Space Learning: Our learner can learn discriminative features for unseen classes. It can also maintain the distance metric learned effective for the seen classes in both new and old sessions, where the old-session features can be kept unchanged without the need to be re-extracted every time in a visual search system.

Extensive experiments on multiple datasets under incremental data distributions show that our method achieves state-of-the-art results. Fig. 2 shows our system diagram.

2. Related Work

We briefly review the recent progress of image retrieval in Section 2.1. We then summarize the CL and describe its

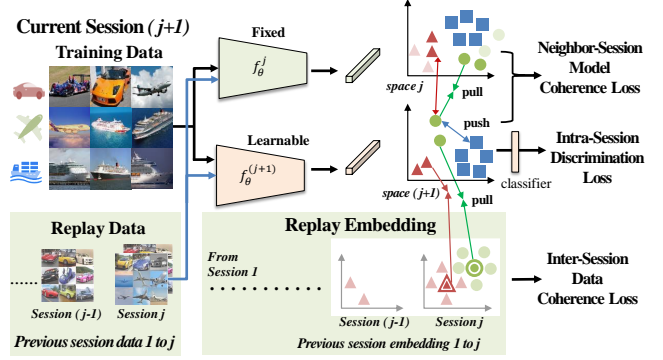


Figure 2. Overview of the proposed approach for CL in General-Incremental setup with backward embedding consistency for a long-term learning. We enforce the three losses for a continual learner: intra-session discrimination loss to learn discriminative representation with mainly the current-session data, neighbor-session model coherence loss to regulate the current model with the previous-session model for backward compatible embedding, and inter-session data coherence of all previous sessions data and replayed embedding of all previous sessions for long-term embedding consistency. Semi-transparent icons represent the data points from the previous sessions collected so far. Note that we omit the replayed data to simplify the illustration.

challenge on similarity-based visual search in Section 2.2.

2.1. Image Retrieval

Image retrieval ranks the gallery images in order given the query image. Previous studies rely on either descriptors by local feature aggregation [23, 39, 49] or low-level visual clues [5, 35, 62] and then perform a nearest neighbor search. The modern way utilizes embeddings from the neural network (*i.e.*, neural mapping) instead because the learnable descriptors show superior accuracy [10, 36, 45, 50, 51, 67] and compact storage [7, 30, 64, 65]. To derive the neural mapping, metric learning optimizes the model in a pairwise manner or pointwise manner. Pairwise methods derive the discriminating space where positive pairs become closer and negative pairs repel each other. A fair benchmark [44] shows the state-of-the-art results of the pairwise method. However, an issue is that the complexity of the sample mining grows in magnitude when the number of samples in a group is beyond two, (*e.g.*, triplet loss [45] and quadruplet [10]). Even though many studies [17, 57–60] delve into the mining issue, sampling the informative pairs is intrinsically hard without clear indicators. On the other hand, pointwise methods [25, 36, 41, 52, 55, 67] treat the optimization as data points against representative samples per class. Instead of meticulous mining, the representative samples could be learnable proxies [25, 36, 41, 52] or randomly sampled points [55, 67], bringing faster convergence for training without pairwise comparisons. For example,

NSoftmax [67] follows the classification training paradigm and applies L2 normalization to the embedding layer output for image retrieval. Such a minor modification yields competitive results compared with pairwise methods in the benchmark [44]. Our method is easily integrated with the strategies mentioned above. For simple and elegant purposes, we optimize with NSoftmax.

Backward Consistency of Feature Embedding: Despite the proliferation of different methods in image retrieval, most work ignores the demand for backward consistency, *i.e.*, making previously frozen features comparable with the newly extracted ones in the gallery set. Feature consistent learning contributes to this goal in two ways. The first line aims to reduce the effort of feature re-extraction. R³AN [9] projects the old features into the new feature space by one-side transform; CMC [56] bridges the multiple feature spaces via a lightweight transformation module. However, they still need re-extraction, which becomes impractical for large galleries and long sessions. The second line maintains the backward consistency without any feature re-extraction. BCT [48] constrains the current feature space by simultaneously enabling gradient flow from both the old and the new classifier. However, it requires the entire previous data seen so far for finetuning and hence is incapable of handling the long session learning scenario. Moreover, the gallery set is assumed fixed in the experimental settings of BCT without addressing the critical issue that the gallery set could be increasingly enlarged in practice. Unlike BCT, our method imposes constraints on both the inter-session feature embeddings and the neighbor-session model; hence it gains a great improvement in backward consistency. Besides, we examine an at most 10-session scenario (in contrast to only 3 sessions at most in BCT’s experiments) for a reality check on the backward feature compatibility.

2.2. Continual Learning

CL aims for a single learner that can sequentially update knowledge without forgetting the previously learned information. The existing works [20, 53] can be categorized into task-incremental and class-incremental CL, where the former assumes the task index presented at the inference time while the later assumes a task-agnostic scenario during inference. Most works assume the categories in the successive tasks (*i.e.*, sessions) to be disjoint to each other.

To avoid overfitting to the current session, regularization-based methods [8, 26, 46, 66] impose constraints on change of the critical neural weights, distillation-based methods [13, 14, 29, 42, 68] transfer the knowledge from the previous learner to the current one, replay-based methods [4, 6, 15, 31, 43] revisit a small amount of old data to prevent forgetting, and isolation-based methods [21, 33, 34, 47] assign the subnetwork capacity per task. Though an automatic task selector [1] can help extend an

isolation-based method to class-incremental, it still requires the disjoint assumption of the classes among tasks. To break the disjoint-class limitation in CL, recent studies propose the blurry setup (*i.e.*, class-overlapping [2, 4, 40]). However, it requires the presentation of all predefined classes across subsequent sessions, and is thus still impractical for real-world visual search applications.

Besides, most works focus on the classification problem, ignoring the demand for image retrieval. MMD [11] studies retrieval-based CL in a distillation manner. By narrowing down the mean difference between two distributions in reproducing kernel space, the current model distills the knowledge of the old one. However, the method ignores the backward-compatible needs for retrieval; hence the embeddings have to be re-extracted for the gallery data as the task grows in CL. It is thus impractical for a realistic system.

3. Continual Learner for Visual Search (CVS)

Consider a CL retrieval problem with J sessions. In session $j \in \{1 \cdots J\}$, a neural-network model is learned and let f_j be the mapping from the input image to the feature embedding of the model. During CL training, the backbone feature extractor f_{j+1} is initialized from f_j , the neural mapping obtained immediately before the current session. To fulfill the requirement in our General-Incremental setup, the gallery set is incrementally extensible with the newly added embedding. Specifically, let G_j be the images newly added in time session j and \mathbf{g}_j be the feature embedding obtained for G_j (via f_j). The embedding will be added to the gallery set up to the j -th session, $\mathbf{g}_{1:j} = \bigcup_{i=1}^j \mathbf{g}_i$. The feature embedding cumulatively saved (*i.e.*, $\mathbf{g}_{1:j}$) then serve as the gallery set for retrieval for future sessions $l (> j)$. In addition, to learn f_{j+1} capable of handling the new-session data containing possibly both existing and novel classes and also maintaining the performance on all previous-session data, we conduct the following loss terms in our CVS approach, namely, inter-session term of data coherence ($\mathbf{L}_{\{1:j\};j+1}^d$), neighbor-session term for model coherence ($\mathbf{L}_{j;j+1}^m$), and intra-session term for discrimination (\mathbf{L}_{j+1}^c).

3.1. Inter-session data coherence

Each training sample of the previous-session data has been already provided with an embedding by one of the old neural mappings, f_1, \cdots, f_j . To fulfill the condition of not re-extracting the embedding for updating the gallery set every time, they are fixed and kept invariant once having been built upon a session $k \in \{1 \cdots j\}$. That is, after establishing the embedding, they can be used for a series of future tasks in CL. Besides hoping to un-forget the classification power of known labels, we further need to un-forget the capability of retrieval based on the embedding already built.

Replayed Embedding Together with Data: Replayed data are widely adopted to avert forgetting in CL. To resolve stor-

age efficiency and avoid training on all data, a small portion of the data from sessions 1 to j can be stored and replayed for a joint training with the current data in session $j + 1$. In our work, the old data embeddings have been extracted as $\mathbf{g}_i = f_i\{G_i\}_{i=1}^j$. Besides replaying the data, we replay the embedding to facilitate CL, which is more efficient with learning. We call the technique *replayed embedding*, where a small portion of features sampled from the embedding space are replayed for training.

The deep model updated in the current session $j + 1$ should be confined to the stored embedding of seen classes. Therefore, the goal of our learner is to train f_{j+1} to extract features from G_{j+1} ($\mathbf{g}_{j+1} = f_{j+1}\{G_{j+1}\}$) while keeping the existing obsolete features $\mathbf{g}_{1:j}$ unchanged. However, the embeddings are established from different neural mappings varying with sessions. The sessions have separated distributions and could be diverse for a long sequence of sessions. As different feature spaces are not necessarily comparable, to avoid an over-fitting to individual sessions, we aggregate the embedding across the previous sessions by taking the expectations as follows.

$$\mathcal{E}_c = \frac{1}{j} \sum_{i=1}^j \xi(\mathbf{g}_{ic}), c \in \mathcal{C}(i), \quad (1)$$

where $\mathcal{C}(i)$ is the set of class indices appearing in session i , \mathbf{g}_{ic} is the set of embedding extracted from the data of class c in G_i , and $\xi(\cdot)$ denotes the expectation operator. The loss enforcing inter-session data coherence is defined as follows:

$$\mathbf{L}_{\{1:j\};j+1}^{\text{dr}} = \sum_{c \in \Pi_{j+1}} \sum_{x_i \in c} \|f_{j+1}(x_i) - \mathcal{E}_c\|_2^2, \quad (2)$$

where $\Pi_{j+1} = \bigcup_{i=1}^j \mathcal{C}(i) \cap \mathcal{C}(j+1)$ is the indices intersection between the classes in the current and all previous sessions; x_i denotes the data in the current session ($j+1$).

However, the current-session classes may not contain all the previous classes in our General Incremental setup. *E.g.*, for the special case of Disjoint setup, there is no old class sample in the incoming sessions at all. Hence, besides employing the replayed embedding \mathcal{E}_c (in Π_{j+1}), we use the replayed data for the old classes $\Gamma_{j+1} = \bigcup_{i=1}^j \mathcal{C}(i)$ too. Without loss of generality, we use the exemplar mining technique iCaRL [43] to conduct a small portion of data for replay, which searches random neighbors around the mean per class. The replayed embedding and data in Γ_{j+1} are co-used as follows:

$$\mathbf{L}_{\{1:j\};j+1}^{\text{do}} = \sum_{c \in \Gamma_{j+1}} \sum_{\tilde{x}_i \in c} \|f_{j+1}(\tilde{x}_i) - \mathcal{E}_c\|_2^2, \quad (3)$$

where \tilde{x}_i denotes the replayed data. Then, the inter-session data coherence loss is a combination of Eqs. 2 and 3, which helps maintain the backward feature consistency:

$$\mathbf{L}_{\{1:j\};j+1}^{\text{d}} = (\mathbf{L}_{\{1:j\};j+1}^{\text{dr}} + \mathbf{L}_{\{1:j\};j+1}^{\text{do}})/n, \quad (4)$$

where n is the mini-batch size.

The loss combines both the replayed embedding and data to maintain the coherence of the feature spaces of sessions $1, \dots, (j+1)$. Even the aggregated embeddings provide strong constraints, as the freedom of neural networks is often large in the capacity, the learner can still easily adapt the model from f_j to f_{j+1} . In our experience, a nice usage of this precedent constraint can help build an effective CL learner because the fixed embeddings can act as attractors to regularize the deep model training. Our approach only needs $|\Gamma_{j+1}|$ replayed embeddings, and we set a ratio of around 5% samples of the total data in sessions $1 : J$ as a fixed budget shared by all sessions up to now in the replayed data. Since the loss establishes the connection between session $(j+1)$ to all previous ones ($1 : j$), unlike BCT [48] employing only the inter relationship between the current session $(j+1)$ and that immediately before the session (j) , our approach can effectively enhance the backward feature consistency for long-term sessions.

Besides, our replayed-embedding solution coincides with the idea in the cross-batch learning [58, 61], where previously learned embeddings in the past mini-batches guide the training in the current one. The approaches have been demonstrated effective to boost retrieval performance. Our approach is analogously a cross-session learning method. In the experiments, we show the efficacy of our cross-session CL solution and make an in-depth analysis on the results.

3.2. Neighbor-session model coherence

In the above, we utilize the replayed materials of seen classes in all previous sessions and the current session data for training. This section develops the loss using mainly the current-session data while co-using the current model (f_{j+1}) and last-time model (f_j) in training. It enforces the neighbor-session model coherence.

The loss is designed using the distillation principle, where the mapping f_l (learned and fixed) serves as the teacher model to guide the training of the student f_{l+1} . Unlike previous studies of distillation techniques that are mainly for classification [18], we conduct a distillation-based loss to regulate the metric learning of f_{j+1} from the metric space of f_j based on the triplet loss to better fit the nature of CL on retrieval. In a triple of (x_a, x_p, x_n) for an anchor, a positive sample, and a negative sample respectively are selected from mainly the current-session data. The purpose of training is to shrink the distance between x_a and x_p while enlarging that between x_a and x_n .

To distill knowledge from the embedding space already built from the previous session j , we generate the embeddings of positive and negative samples using $f_j(\cdot)$ and obtain $f_j(x_p)$ and $f_j(x_n)$, respectively. Then, we use the teacher-generated embedding to guide the training of the student anchor $f_{j+1}(x_a)$. The current session data are fed

into the current and previous models to constrain the embedding distributions produced by f_{j+1} and f_j , respectively. The loss conducted can thus enforce the model coherence between neighbor sessions.

Besides, we co-use x_a and x_p via setting $x_p = x_a$. Hence, we conduct a 2-sample-3-embedding triplet-loss strategy where the triplet embeddings are $f_{j+1}(x_a)$, $f_j(x_a)$, and $f_j(x_n)$. This is because it can facilitate drawing the embedding distributions of the student and the teacher based on the same x_a for comparison. In the meantime, it also saves the computation of choosing positive samples and enforces efficient sampling. Recent studies have shown that mining easy positive samples (*i.e.*, similar positives) benefits to metric learning [60]. Our approach directly forms the positive sample’s embedding from the anchor, which eliminates the positive mining effort in [60] during pair sampling. For the negative, we follow the hardest negative mining principle [17] to pick the embedding. Denote $d_j^{j+1}(x, y) = \|f_{j+1}(x) - f_j(y)\|_2^2$. The neighbor-session-model-coherence loss is written as:

$$\mathbf{L}_{j;j+1}^m = \frac{1}{n} \sum_{x_a} \left[d_j^{j+1}(x_a, x_a) - d_j^{j+1}(x_a, x_n) + m \right]_+, \quad (5)$$

where m is the margin set as 0.1 by default. Besides, $f_{j+1}(x)$ and $f_j(x)$ are l_2 -normalized to alleviate the issue of forgetting according to previous studies [19, 32]. This loss helps restrict the behavior of the updated model coherent to the previous one.

3.3. Intra-session discrimination

Recent studies show that classification is a strong baseline for learning effective feature embedding (if an appropriate l_2 -norm normalization layer is added) [67]. Without loss of generality, we employ the method to establish the retrieval capability of f_{j+1} using mainly the current-session data. The loss for intra-session discrimination is

$$\mathbf{L}_{j+1}^c = \frac{1}{n} \sum_{i=1}^n -\log\left(\frac{\exp(w_{y_i}^T f_{j+1}(x_i)/T)}{\sum_k \exp(w_k^T f_{j+1}(x_i)/T)}\right), \quad (6)$$

where $(x_i, y_i)_{i=1}^n$ are the data and labels, w denote the l_2 -normalized weight of the classification layer, $f_{j+1}(x_i)$ is the l_2 -normalized embedding extracted for x_i , and T is the temperature term set as 0.05 by default.

Total Loss. The final learning objective is as follows:

$$\min_{f_{j+1}} \mathbf{L} = \mathbf{L}_{j+1}^c + \alpha \mathbf{L}_{j;j+1}^m + \beta \mathbf{L}_{\{1:j\};j+1}^d. \quad (7)$$

By default, we set α to 10 and β to 1 empirically in our experiments, and cosine distance is used for retrieval in this work. Only the first term \mathbf{L}_1^c is used in session 1. Then, all three terms are used in sessions 2 to J . Fig. 2 gives an illustration of our approach.

In sum, to fulfill the General Incremental setup, the first term \mathbf{L}_{j+1}^c provides a fundamental retrieval ability of the current session ($j + 1$). The second term $\mathbf{L}_{j;j+1}^m$ makes the neural mapping f_{j+1} close to f_j based on both the novel- and seen-label data in the current session. It enforces the prediction to mimic the behavior of the previous feature extractor. The third term $\mathbf{L}_{\{1:j\};j+1}^d$ enforces the backward feature consistency from the replayed embedding and data of all sessions (1 to j) and a joint training with the current-session data based only on the labels seen before (*i.e.*, $\bigcup_{i=1}^j \mathcal{C}(i)$). It helps bias the feature space to align with the one obsolete gallery features lie in. By joining the three losses, our CVS approach can handle the General Incremental setup and its special cases (Blurry, Disjoint) well.

4. Experiments

We conduct extensive experiments across various datasets under different data distributions. Five datasets are used, including two *coarse-grained* datatests, CIFAR100 and Tiny ImageNet, and three *fine-grained* datasets, Stanford Dog, iNaturalist 2017, and Product-10K. We use CIFAR100 for fundamental study and then Tiny Imagenet on a longer sequence of sessions.

The datasets are summarized as follows. **CIFAR100** [27] has 100 categories; each owns 500 32×32 images for training and 100 images for testing. **Tiny ImageNet** [28] has 100,000 training and 10,000 testing images (of size 64×64 sampled from the 200 classes from ImageNet. Due to its compact size and various categories, we conduct a long sequence learning on it. **Stanford Dog** [24] contains 120 dog breed-level categories picked from ImageNet, including 8,580 images for testing and 12,000 images for training. For the training split, each class has 100 images. **iNaturalist 2017** [54] is a large-scale long-tailed image retrieval dataset with 5,089 species-level categories. We sample 527 images per class for the 200 classes; each contains at least 527 samples to avoid very small classes and unbalanced data before partition. We call it **iNat-M** in this work. **Product-10K** [3] is an ultrafine-grained long-tailed dataset covering the top-9691 frequently bought product images from a real e-commerce system; Each class’s images contain diverse appearances of the specific product by collecting offline customer-taken and online in-shop photos. We remove the category with fewer than 20 images in the training set. We then construct the train-test split based on the original training set because the official testing data do no offer labels. Finally, there are 2,743 classes. We call it **Product-M** in this work. Fine-grained images are more challenging as only subtle differences exist.

4.1. Implementation Details

In all of the experiments, the embedding dimension sets as 128 by default. We use ResNet-18 [16] for coarse-

grained datasets, and ResNet-50 [16] for fine-grained datasets. The hyperparameter details are depicted in the supplementary material.

Disjoint setup assumes that previously seen classes are unavailable when learning new session data. As shown in Fig. 3a, we partition the CIFAR100 into five sessions with each containing 20 classes. For Tiny ImageNet, the dataset is divided into ten sessions with 20 classes per session.

Blurry setup [4] provides all categories in the beginning; each session shows a subset of samples from all of them. The data distribution is controlled via a percentage of the data from the major and minor classes per session. In our experiments, CIFAR100 is divided into five sessions; each contains 20 major classes and 80 minor classes, with 90% samples from the major and 10% from the minor, as shown in 3b. For Tiny ImageNet, we simulate a 10-session CL; each has 70% from the major and 30% from the minor.

General-incremental setup: The new-class and novel old-class samples may co-exist. The model is learned from the S classes at initial in a total of L sessions. For the later sessions, we add C classes and assume that $M\%$ of the data are from the old categories and $(1 - M)\%$ are from the new, denoted as (S, C, M, L) . We set $(20, 20, 10, 5)$ for CIFAR100 (Fig. 3c). For the rest, we apply $(20, 20, 30, 10)$ to Tiny ImageNet, $(60, 20, 30, 4)$ to Stanford Dog, $(100, 25, 30, 5)$ to iNat-M, and $(1343, 700, 40, 3)$ to Product-M.

Replayed data are used for all loss terms. Following [4], we set the memory buffer size for replay as 2,000 for CIFAR100. Note that this is a budget-limited buffer (for the replay-based methods) co-used by all sessions. We assume a memory budget of 600 samples for Dog due to the small scale, 4,000 for Tiny ImageNet, 4,000 for iNat-M, and 3,000 for the ultrafine-grained Product-M, respectively. It is pretty challenging for Product-M as it has 2,743 classes; for each session, only 1 \sim 2 images will be stored for replay.

4.2. Evaluation

We follow the train-validation-test protocol. The model achieving the best recall@1 in the validation phase is picked for fair comparisons. Then, we report the final performance by using the original testing set as a query set. For collecting the validation query set for each class in disjoint and general-incremental setups, we sample 5% of current training session data for iNat-M and Tiny ImageNet, and 10% for CIFAR100 and Stanford Dog. We sample 2 images randomly from each class for Product-M due to scarce training data. Note that such a validation query set is accumulated as the session grows in the disjoint and general setup. For the blurry setup, we keep a fixed portion of data from the entire training set as a validation query set (5% for Tiny ImageNet and 10% for CIFAR100) class-wisely since the model can see all predefined categories in the beginning. The embedding extracted for the current-session training data is added

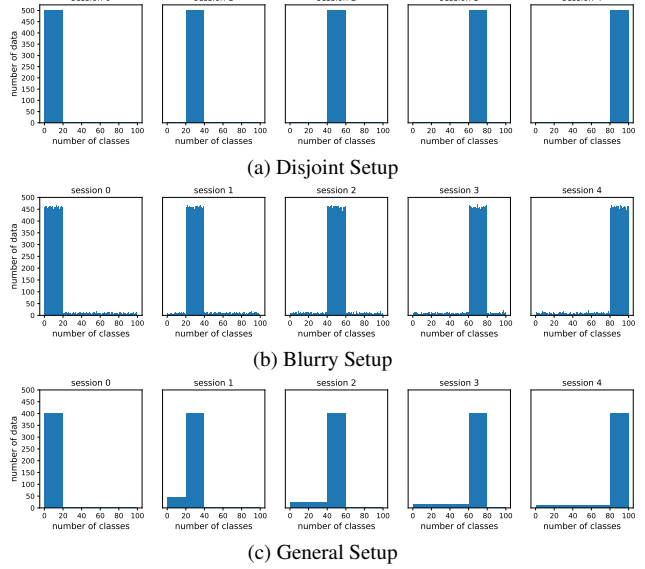


Figure 3. Class distributions of different setups on CIFAR100.

to the gallery set after finishing the session as mentioned before; thus the new-classes images first observed in one session are allowed to be searched in the future sessions. Such an extensible gallery set matches the real-world scenarios.

We report the **recall@k** [22] as it is the most popular metric in fine-grained retrieval [44, 51]. For the testing query, we evaluate all the classes for the blurry setup, and the classes seen so far for both the disjoint and general setup, respectively. Finally, we average the scores over sessions, denoted as **AR@K**. To compare with existing approaches, we re-implement the LWF [29], and MMD [11] and re-run EWC [26], RWalk [8], and Rainbow [4] from the official Github of [4]. Among them, all (except MMD) are for classification since CL is rare to be studied on retrieval yet, and we replace all the plain softmax loss with normalized softmax loss [67] for them for a fair comparison. We also re-implement BCT [48] that is a backward feature compatible approach for retrieval. We apply the same exemplar mining as ours for BCT in disjoint setup as the Naïve BCT is not a CL solution and requires old classes to pass through.

Besides, we provide direct fine-tuning as the **lower bound**; in contrast, we offer joint training of all gallery images seen so far allowing re-extraction as the **upper bound**.

4.3. Results on the Coarse-grained Datasets

Results on CIFAR100: We conduct a 5-session experiment on CIFAR100. The results are shown in Fig. 4 (recall@1) and Table 1(a) (**AR@K**). For all three setups, except for *Joint Train* that is an upper bound of this experiment, our CVS performs the best on both evaluation measures. In the blurry and general-incremental setups, we observe that BCT performs even worse than *Finetune*, a lower bound in this

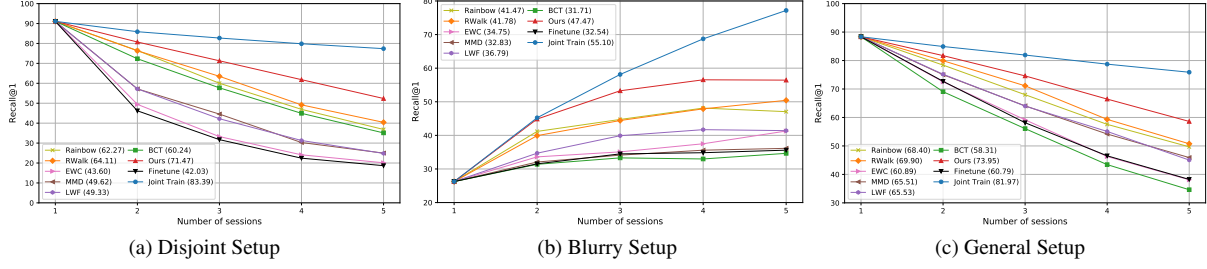


Figure 4. Recall@1 on CIFAR100 for the three setups, where Average Recall@1 across sessions is reported in parentheses.

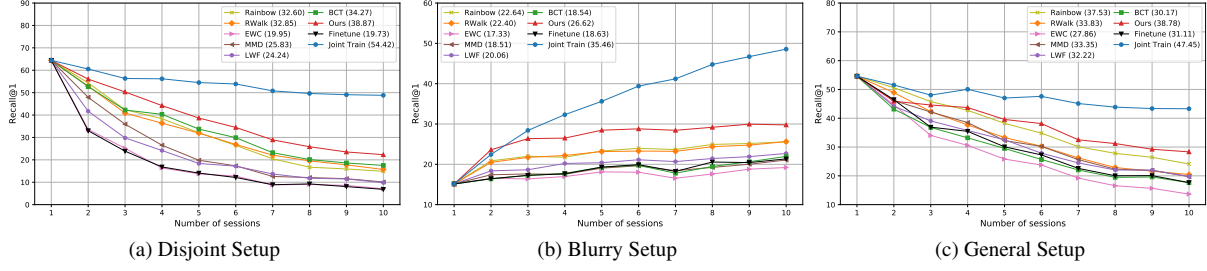


Figure 5. Recall@1 on Tiny ImageNet for the three setups, where Average Recall@1 across sessions is reported in parentheses.

	Disjoint			Blurry			General		
	AR@1	AR@2	AR@4	AR@1	AR@2	AR@4	AR@1	AR@2	AR@4
Joint Train	83.39	85.69	87.64	55.1	57.74	60.33	81.97	84.6	86.82
Finetune	42.03	43.23	44.74	32.54	36.8	41.11	60.79	64.73	68.14
BCT	60.24	64.38	68.37	31.71	35.69	39.61	58.31	62.2	65.66
LWF	49.33	53.51	58.01	36.79	41.63	46.52	65.53	70.91	75.31
MMD	49.62	53.58	57.87	32.83	36.73	40.64	65.51	70.22	74.33
EWC	43.6	45.65	48.12	34.75	40.85	46.71	60.89	64.86	68.2
RWalk	64.11	67.33	70.56	41.78	45.67	49.4	69.9	73.37	76.39
Rainbow	62.27	65.09	67.43	41.47	44.99	48.27	68.4	71.56	74.2
Ours (CVS)	71.47	74.8	77.51	47.47	49.86	52.17	73.95	76.73	78.84

(a)

	Disjoint			Blurry			General		
	AR@1	AR@2	AR@4	AR@1	AR@2	AR@4	AR@1	AR@2	AR@4
Joint Train	54.42	58.62	62.31	35.46	39.29	43.19	47.45	51.94	55.99
Finetune	19.73	20.92	21.81	18.63	23.17	28.11	31.11	35.87	40.8
BCT	34.27	37.92	41.22	18.54	23	27.66	30.17	34.64	39.29
LWF	24.24	27.5	30.67	20.06	25.11	30.52	32.22	38.23	44.56
MMD	25.83	29.22	32.75	18.51	22.65	27.23	33.35	38.21	43.06
EWC	19.95	21.89	23.88	17.33	21.42	26.07	27.86	32.67	37.78
RWalk	32.85	37.09	41.09	22.4	26.11	29.94	33.83	38.43	42.97
Rainbow	32.6	35.51	38.35	22.64	26.47	30.4	37.53	41.64	45.44
Ours (CVS)	38.87	42.03	45.08	26.62	29.19	31.86	38.78	42.38	45.89

(b)

	Dog			iNat-M			Product-M		
	AR@1	AR@2	AR@4	AR@1	AR@2	AR@4	AR@1	AR@2	AR@4
Joint Train	86.98	91.08	93.99	75.85	79.97	83.65	79.36	83.83	87.73
Finetune	82.7	88.32	92.23	67.75	72.68	77	70.99	76.26	81.16
BCT	81.73	87.49	91.55	67.34	72.07	75.99	70.48	75.67	80.47
LWF	83.25	89.29	93.04	68.51	73.71	78.12	72.95	78.38	83.1
MMD	83.2	88.98	92.71	68.58	73.48	77.79	72.89	78.21	82.96
EWC	81.64	88.7	92.82	66.3	71.4	75.82	66.01	72.53	78.24
RWalk	82.41	88.31	92.43	68.77	73.82	78.16	68.71	74.64	80.13
Rainbow	82.78	89.04	93.23	68.68	73.22	77.21	69.39	75.19	80.34
Ours (CVS)	84.71	89.4	92.61	72.57	76.39	79.87	75.47	80.36	84.68

(c)

Table 1. Results on (a) CIFAR100 and (b) Tiny ImageNet, and (c) Fine-grained datasets, where *Joint Train* and *Finetune* specify the theoretically upper and lower bounds, respectively.

work. Even though both old and new classes are shown for each incremental session, BCT would still require all samples seen previously for training to obtain satisfied results. On the other hand, we find that EWC attains almost the same performance as the lower bound. We attribute the bad

performance to the uncertainty of importance weight estimation. The runner-up is RWalk, which improves EWC by imposing constraints on the parameter space while avoiding forgetting via a replay-based mechanism. Such a mixed strategy works for classification but is not sufficiently well for backward-compatible retrieval. Instead, our CVS utilizes the inter- and neighbor-session information in extra, showing the efficacy under all setups.

Results on Tiny ImageNet: We use Tiny ImageNet for simulating the 10-session CL scenario in this experiment. To our knowledge, we’re the first to perform such a long session sequence for backward-compatible retrieval in CL. According to Fig. 5 and Table 1(b), our method consistently outperforms the existing competitors on the three setups for this long sequence setting on retrieval. The overall runner-up in this dataset becomes Rainbow, a data-replay method for classification. Our CVS employs the replayed embedding to summarize a class across sessions for retrieval in addition to the data, and performs more favorably for all cases on recall@1 and $AR@K$ particularly when k is small. For a large $k = 4$ in the general setup, the performance is tie. We attribute the gap is reduced since only one of the four retrieved labels has to be correct.

Comparisons to Other Embedding Distillations: We conduct a comparison of our CVS to the techniques of metric learning from teacher [12, 38, 63] and knowledge distillation [18] by replacing $L^m_{j:j+1}$ in Eq. 7 with different losses for neighbor-session model coherence based on CIFAR-100. For a fair comparison, we carefully tune hyperparameter α at $\{1, 10\}$ for the better results at $AR@1$. Our 2-sample-3-embedding solution demonstrates the best

position among competitors for all setups. Compared to the runner-up results, we obtain CVS (71.47) v.s. Anglewise-RKD [38] (70.45) in disjoint setup, CVS(47.47) v.s. Absolute MLKD [63] (46.9) in blurry setup, and CVS(73.95) v.s. Dark Knowledge [18](73.18) in general incremental setup. We detail the results in the supplementary material.

Comparisons on Classification: As our CVS can produce the classification results via the NSoftmax layer too, we further compare the results with the CL classifier Rainbow [4] based on CIFAR-100 after finishing all sessions. We have already followed the setting of [4] on both the class distribution ratios and replayed buffer size in the blurry setup. The classification accuracy presented in [4] is 41.35% with an online learning protocol (*i.e.*, only a single epoch is allowed in learning). We rerun the learner to converge and get the accuracy of 50.2%. Surprisingly, CVS can achieve the 54.04% classification accuracy that is even higher. We owe the promising classification results to the replayed embedding that can serve as useful exemplars (like the principle of cross-batch learning [58, 61]) to further guide the training in our CVS. We conduct extra experiments on the other two setups (disjoint, general-incremental) and one additional CL classifier RWalk and obtain the results: CVS (50.62) v.s. Rainbow (46.69) v.s. RWalk (46.85) in disjoint setup, CVS (54.04) v.s. Rainbow (50.2) v.s. RWalk(50.89) in blurry setup, and CVS (55.49) v.s. Rainbow (52.26) v.s. RWalk(51.92) in general incremental setup. The results demonstrate the efficacy of CVS for CL.

4.4. Results on the Fine-grained Datasets

As for the fine-grained benchmarks, we consider the general-incremental setup only due to its practical usefulness. We initialize the model with ImageNet pretrained weights as it is a common practice in the fine-grained retrieval benchmark. We assume the first session presents samples from half of the classes from the dataset. It is a practical setting for a modern visual search system because a robust service should be well-trained to a certain extent before going online. Due to the fine-grained limits on the data amount, we partition the Stanford Dog dataset into four, iNat-M into five, and Product-M into three sessions, respectively. The results are shown in Table 1(c). Our CVS consistently outperforms the other approaches for all the fine-grained benchmarks and measures, except for AR@ k for $k = 4$ on Stanford Dog Dataset. In sum, the evaluation results with other state-of-the-art approaches demonstrate that our approach is effective for both coarse- and fine-grained datasets on the general-incremental setup because we not only consider backward compatible embedding of the neighboring sessions but also long-term consistency with all the data of past sessions via replay embedding and data. In addition, this backward compatible feature also saves the CL from expensive computational costs of gallery

	CIFAR100			Dog			Product-M		
	AR@1	AR@2	AR@4	AR@1	AR@2	AR@4	AR@1	AR@2	AR@4
L^c	60.79	64.73	68.14	82.7	88.32	92.23	70.99	76.26	81.16
$L^c + L^m$	63.79	68.13	72.02	82.65	88.43	92.29	71.77	76.92	81.78
$L^c + L^d$ w/o replay	64.5	68.12	71.14	83.29	88.77	92.38	73.96	78.69	83.02
$L^c + L^d$	72.16	74.67	76.7	84.37	89.07	92.46	75.6	80.27	84.33
$L^c + L^m + L^d$	73.95	76.73	78.74	84.71	89.4	92.61	75.47	80.36	84.68

Table 2. Ablation study on each component, where L^c , L^m , and L^d are the losses of intra-session discrimination, neighbor-session model coherence, and inter-session data coherence, respectively.

feature re-extraction as model updates during each session.

Discussion: A *limitation* we find is that almost all methods (including ours) stagnate as the session expands, especially in the blurry setup. The margin between ours and upper bound enlarges gradually even though our method surpasses the remaining approaches. Therefore, there is still room for improvement in the long-term challenge.

Ablation Study: Our CVS consists of three loss terms, L^c , L^m , and L^d . To verify their effectiveness, we conduct the ablation study as shown in Table 2. First, we find that $L^c + L^d$ has the greatest impact on the overall performance. It brings the gain in the range of 1.67% to 11.37% in AR@1 compared with L^c alone. Second, the effect of $L^c + L^m$ is weaker despite a slight increase on CIFAR100 and Product-M compared to the one of L^c . Our finding suggests that adopting classification loss with distillation loss alone is insufficient. Therefore, seeking unification of additional information like L^d is practical in our setting. To examine the performance gain from the consistency loss in more depth, we detach the exemplar replay technique from our method. We denote $L^c + L^d$ w/o replay as the aforementioned case. The replay-based trick improves the plain version of the consistency loss on all datasets by a margin, especially for CIFAR100. Without reviewing exemplar data, the performance drop ranges between 1.08% and 7.66% in AR@1. Hence, integrating this design is essential for overall performance.

5. Conclusion

In this work, we present a novel general incremental setup which allows the new gallery set of both seen and unseen classes incrementally added to the database and is closer to the real-world retrieval setup than the widely adopted disjoint and the recent blurry setups. Besides, we also propose a CL method for long-term visual search with backward consistent feature embedding. Our method acts as an extension of cross-batch memory to the cross-session memory for feature embedding learning in CL. We introduce a 2-sample-3-embedding strategy in a triplet for distillation learning across neighbor sessions to enforce the model coherence. The extensive experiments show that our approach achieves the state of the art performance in two coarse-grained classification and three fine-grained datasets under different incremental data distributions.

References

- [1] Davide Abati, Jakub Tomczak, Tijmen Blankevoort, Simone Calderara, Rita Cucchiara, and Babak Ehteshami Bejnordi. Conditional channel gated networks for task-aware continual learning. In *CVPR*, pages 3930–3939, 2020.
- [2] Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient based sample selection for online continual learning. In *NeurIPS*, pages 11816–11825, 2019.
- [3] Yalong Bai, Yuxiang Chen, Wei Yu, Linfang Wang, and Wei Zhang. Products-10k: A large-scale product recognition dataset. *arXiv preprint arXiv:2008.10545*, 2020.
- [4] Jihwan Bang, Heesu Kim, YoungJoon Yoo, Jung-Woo Ha, and Jonghyun Choi. Rainbow memory: Continual learning with a memory of diverse samples. In *CVPR*, pages 8218–8227, 2021.
- [5] S. Belongie, J. Malik, and J. Puzicha. Shape matching and object recognition using shape contexts. *PAMI*, 24(4):509–522, 2002.
- [6] P. P. Brahma and Adrienne Othon. Subset replay based continual learning for scalable improvement of autonomous systems. In *CVPRW*, pages 1179–11798, 2018.
- [7] Zhangjie Cao, Mingsheng Long, Jianmin Wang, and Philip S. Yu. Hashnet: Deep learning to hash by continuation. In *ICCV*, pages 5609–5618, 2017.
- [8] Arslan Chaudhry, Puneet K. Dokania, Thalaiyasingam Ajanthan, and Philip H. S. Torr. Riemannian walk for incremental learning: Understanding forgetting and intransigence. In *ECCV*, pages 556–572, 2018.
- [9] Ken Chen, Yichao Wu, Haoyu Qin, Ding Liang, Xuebo Liu, and Junjie Yan. R³ adversarial network for cross model face recognition. In *CVPR*, pages 9860–9868, 2019.
- [10] Weihua Chen, Xiaotang Chen, Jianguo Zhang, and Kaiqi Huang. Beyond triplet loss: A deep quadruplet network for person re-identification. In *CVPR*, pages 1320–1329, 2017.
- [11] Wei Chen, Yu Liu, Weiping Wang, Tinne Tuytelaars, Erwin M Bakker, and Michael Lew. On the exploration of incremental learning for fine-grained image retrieval. In *BMVC*, 2020.
- [12] Yuntao Chen, Naiyan Wang, and Zhaoxiang Zhang. Dark-rank: Accelerating deep metric learning via cross sample similarities transfer. In *AAAI*, pages 2852–2859, 2018.
- [13] Prithviraj Dhar, Rajat Vikram Singh, Kuan-Chuan Peng, Ziyang Wu, and Rama Chellappa. Learning without memorizing. In *CVPR*, pages 5133–5141, 2019.
- [14] Arthur Douillard, Matthieu Cord, Charles Ollion, Thomas Robert, and Eduardo Valle. Podnet: Pooled outputs distillation for small-tasks incremental learning. In *ECCV*, pages 86–102, 2020.
- [15] Tyler L Hayes, Nathan D Cahill, and Christopher Kanan. Memory efficient experience replay for streaming learning. In *ICRA*, pages 9769–9776, 2019.
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016.
- [17] Alexander Hermans, L. Beyer, and B. Leibe. In defense of the triplet loss for person re-identification. *arXiv preprint arXiv:1703.07737*, 2017.
- [18] Geoffrey Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural network. In *NeurIPS*, 2015.
- [19] Saihui Hou, Xinyu Pan, Chen Change Loy, Zilei Wang, and Dahua Lin. Learning a unified classifier incrementally via rebalancing. In *CVPR*, pages 831–839, 2019.
- [20] Yen-Chang Hsu, Yen-Cheng Liu, Anita Ramasamy, and Zsolt Kira. Re-evaluating continual learning scenarios: A categorization and case for strong baselines. In *NeurIPS*, 2018.
- [21] Ching-Yi Hung, Cheng-Hao Tu, Cheng-En Wu, Chien-Hung Chen, Yi-Ming Chan, and Chu-Song Chen. Compacting, picking and growing for unforgetting continual learning. In *NeurIPS*, pages 13647–13657, 2019.
- [22] H. Jégou, M. Douze, and C. Schmid. Product quantization for nearest neighbor search. *PAMI*, 33:117–128, 2011.
- [23] Hervé Jégou, Matthijs Douze, Cordelia Schmid, and Patrick Pérez. Aggregating local descriptors into a compact image representation. In *CVPR*, pages 3304–3311, 2010.
- [24] Aditya Khosla, Nityananda Jayadevaprakash, Bangpeng Yao, and Li Fei-Fei. Novel dataset for fine-grained image categorization. In *CVPRW*, 2011.
- [25] Sungyeon Kim, Dongwon Kim, Minsu Cho, and Suha Kwak. Proxy anchor loss for deep metric learning. In *CVPR*, pages 3235–3244, 2020.
- [26] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *PNAS*, 114(13):3521–3526, 2017.
- [27] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical report, 2009.
- [28] Ya Le and X. Yang. Tiny imagenet visual recognition challenge. 2015.
- [29] Zhizhong Li and Derek Hoiem. Learning without forgetting. *PAMI*, 40:2935–2947, 2018.
- [30] Kevin Lin, Huei-Fang Yang, Jen-Hao Hsiao, and Chu-Song Chen. Deep learning of binary hash codes for fast image retrieval. In *CVPRW*, pages 27–35, 2015.
- [31] Yaoyao Liu, Yuting Su, An-An Liu, Bernt Schiele, and Qianru Sun. Mnemonics training: Multi-class incremental learning without forgetting. In *CVPR*, pages 12242–12251, 2020.
- [32] Chunjie Luo, Jianfeng Zhan, Xiaohe Xue, Lei Wang, Rui Ren, and Qiang Yang. Cosine normalization: Using cosine similarity instead of dot product in neural networks. In *ICANN*, pages 382–391, 2018.
- [33] Arun Mallya, Dillon Davis, and Svetlana Lazebnik. Piggyback: Adapting a single network to multiple tasks by learning to mask weights. In *ECCV*, pages 72–88, 2018.
- [34] Arun Mallya and Svetlana Lazebnik. Packnet: Adding multiple tasks to a single network by iterative pruning. In *CVPR*, pages 7765–7773, 2018.
- [35] B.S. Manjunath and W.Y. Ma. Texture features for browsing and retrieval of image data. *PAMI*, 18(8):837–842, 1996.
- [36] Yair Movshovitz-Attias, A. Toshev, Thomas Leung, S. Ioffe, and Saurabh Singh. No fuss distance metric learning using proxies. In *ICCV*, pages 360–368, 2017.

- [37] Kevin Musgrave, Serge J. Belongie, and Ser-Nam Lim. A metric learning reality check. In *ECCV*, pages 681–699, 2020.
- [38] Wonpyo Park, Dongju Kim, Yan Lu, and Minsu Cho. Relational knowledge distillation. In *CVPR*, pages 3962–3971, 2019.
- [39] Florent Perronnin, Yan Liu, Jorge Sánchez, and Hervé Poirier. Large-scale image retrieval with compressed fisher vectors. In *CVPR*, pages 3384–3391, 2010.
- [40] Ameya Prabhu, Philip Torr, and Puneet Dokania. Gdumb: A simple approach that questions our progress in continual learning. In *ECCV*, pages 524–540, 2020.
- [41] Qi Qian, Lei Shang, Baigui Sun, Juhua Hu, H. Li, and Rong Jin. Softtriple loss: Deep metric learning without triplet sampling. In *ICCV*, pages 6449–6457, 2019.
- [42] Amal Rannen, Rahaf Aljundi, Matthew B. Blaschko, and Tinne Tuytelaars. Encoder based lifelong learning. In *ICCV*, pages 1329–1337, 2017.
- [43] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H. Lampert. icarl: Incremental classifier and representation learning. In *CVPR*, pages 5533–5542, 2017.
- [44] Karsten Roth, Timo Milbich, Samarth Sinha, Prateek Gupta, Bjorn Ommer, and Joseph Paul Cohen. Revisiting training strategies and generalization performance in deep metric learning. In *ICML*, pages 8242–8252, 2020.
- [45] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *CVPR*, pages 815–823, 2015.
- [46] Jonathan Schwarz, Wojciech Czarnecki, Jelena Luketina, Agnieszka Grabska-Barwinska, Yee Whye Teh, Razvan Pascanu, and Raia Hadsell. Progress & compress: A scalable framework for continual learning. In *ICML*, pages 4528–4537, 2018.
- [47] Joan Serra, Didac Suris, Marius Miron, and Alexandros Karatzoglou. Overcoming catastrophic forgetting with hard attention to the task. In *ICML*, pages 4548–4557, 2018.
- [48] Yantao Shen, Yuanjun Xiong, Wei Xia, and Stefano Soatto. Towards backward-compatible representation learning. In *CVPR*, pages 6367–6376, 2020.
- [49] Josef Sivic and Andrew Zisserman. Video google: a text retrieval approach to object matching in videos. In *ICCV*, pages 1470–1477, 2003.
- [50] Kihyuk Sohn. Improved deep metric learning with multi-class n-pair loss objective. In *NeurIPS*, pages 1857–1865, 2016.
- [51] Hyun Oh Song, Yu Xiang, S. Jegelka, and S. Savarese. Deep metric learning via lifted structured feature embedding. In *CVPR*, pages 4004–4012, 2016.
- [52] Eu Wern Teh, Terrance Devries, and Graham W. Taylor. Proxynca++: Revisiting and revitalizing proxy neighborhood component analysis. In *ECCV*, pages 448–464, 2020.
- [53] Gido M Van de Ven and Andreas S Tolias. Three scenarios for continual learning. In *NeurIPS*, 2018.
- [54] Grant Van Horn, Oisín Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro Perona, and Serge Belongie. The inaturalist species classification and detection dataset. In *CVPR*, pages 8769–8778, 2018.
- [55] Ji Wan, Dayong Wang, Steven Chu Hong Hoi, Pengcheng Wu, Jianke Zhu, Yongdong Zhang, and Jintao Li. Deep learning for content-based image retrieval: A comprehensive study. In *ACM MM*, pages 157–166, 2014.
- [56] Chi Wang, Ya-Liang Chang, Shang-Ta Yang, Dong Chen, and Shang-Hong Lai. Unified representation learning for cross model compatibility. In *BMVC*, 2020.
- [57] Xun Wang, Xintong Han, Weilin Huang, Dengke Dong, and Matthew R. Scott. Multi-similarity loss with general pair weighting for deep metric learning. In *CVPR*, pages 5017–5025, 2019.
- [58] Xun Wang, Haozhi Zhang, Weilin Huang, and Matthew R. Scott. Cross-batch memory for embedding learning. In *CVPR*, pages 6387–6396, 2020.
- [59] Chao-Yuan Wu, R. Manmatha, Alexander J. Smola, and Philipp Krahenbuhl. Sampling matters in deep embedding learning. In *ICCV*, pages 2859–2867, 2017.
- [60] Hong Xuan, Abby Stylianou, and Robert Pless. Improved embeddings with easy positive triplet mining. In *WACV*, pages 2463–2471, 2020.
- [61] Huei-Fang Yang, Kevin Lin, and Chu-Song Chen. Cross-batch reference learning for deep classification and retrieval. In *ACM MM*, pages 1237–1246, 2016.
- [62] Hui Yu, Mingjing Li, Hong-Jiang Zhang, and Jufu Feng. Color texture moments for content-based image retrieval. In *ICIP*, pages 929–932, 2002.
- [63] Lu Yu, Vacit Oguz Yazici, Xialei Liu, Joost van de Weijer, Yongmei Cheng, and Arnau Ramisa. Learning metrics from teachers: Compact networks for image embedding. In *CVPR*, pages 2902–2911, 2019.
- [64] Li Yuan, Tao Wang, Xiaopeng Zhang, Francis EH Tay, Zequn Jie, Wei Liu, and Jiashi Feng. Central similarity quantization for efficient image and video retrieval. In *CVPR*, pages 3080–3089, 2020.
- [65] Xin Yuan, Liangliang Ren, Jiwen Lu, and Jie Zhou. Relaxation-free deep hashing via policy gradient. In *ECCV*, pages 141–157, 2018.
- [66] Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. In *ICML*, pages 3987–3995, 2017.
- [67] Andrew Zhai and Hao-Yu Wu. Classification is a strong baseline for deep metric learning. In *BMVC*, 2019.
- [68] Junting Zhang, J. Zhang, Shalini Ghosh, Dawei Li, S. Tasci, Larry Heck, H. Zhang, and C.-C. Jay Kuo. Class-incremental learning via deep model consolidation. In *WACV*, pages 1120–1129, 2020.
- [69] Bo Zhao, S. Tang, Dapeng Chen, Hakan Bilen, and Rui Zhao. Continual representation learning for biometric identification. In *WACV*, pages 1197–1207, 2021.

Supplementary Material

Continual Learning for Visual Search with Backward Consistent Feature Embedding

Timmy S. T. Wan¹ Jun-Cheng Chen² Tzer-Yi Wu³ Chu-Song Chen^{1,*}
 National Taiwan University¹ Academia Sinica² ucfunnel Co. Ltd.³

{r08944004, chusong}@csie.ntu.edu.tw, pullpull@citi.sinica.edu.tw, kenny.wu@ucfunnel.com

A. Contribution Review

General Incremental Setup: Unlike previous works, we investigate a general case for CL (Fig. A.1). Our setup considers the incremental classes of the disjoint setup and also covers overlapped classes of the blurry setup. In a broad sense, it offers a more common situation for the CL research community on both classification and retrieval.

Backward Consistent Feature Space Learning: We propose a novel continual learner for visual search allowing acquiring knowledge for unseen classes and making both the previous and the current feature space comparable without backfilling (*i.e.*, re-extraction) of the previously processed gallery images. We bridge this gap in three loss terms:

- An inter-session data coherence loss learns from the history of all sessions by taking the extensible replayed embedding as a free supervision signal for guidance.
- A neighbor-session model coherence loss preserves the distance metric for the seen classes in both new and old sessions; it leverages a revised triplet loss with a new sampling strategy for distillation.
- An intra-session discrimination loss grasps knowledge from the novel categories using pointwise metric learning without loss of flexibility.

An enlarged Fig. 2 of the main paper is provided for a more precise illustration, as shown in Fig. A.2. In the following, we complement more implementation details in Section B and ablation studies in Section C.

B. Hyperparameter Details

As presented in Section 4.1 of the main paper, we show the hyperparameter details as follows. We implement all models with Pytorch [13] using NVIDIA V100 GPUs. To control the embedding size, we insert the fully connected layer with dimension 128 before the final softmax layer given the network architecture. The full experiments on the general incremental setup are shown in Tables C.1 and C.2.

* indicates corresponding author.

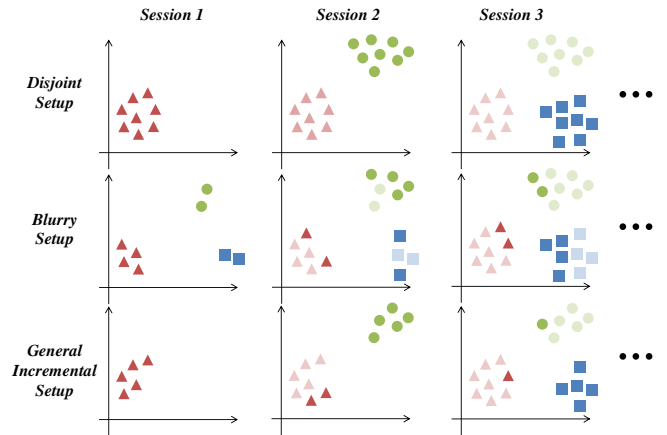


Figure A.1. The widely adopted Disjoint setup (upper row) assumes the image categories mutually disjoint among sessions. The middle row shows the recent Blurry setup, where different sessions allow overlapping classes but all the classes are given initially; every session has a specific data distribution over the known classes. The bottom row shows our General Incremental setup, where the classes in a new session can be either old or novel.

B.1. Details on the Coarse-grained Datasets

The hyperparameters almost follow those in Rainbow [1] but with different batch sizes for Tiny ImageNet. We train ResNet-18 [6] over 256 epochs with the batch size of 16 and 64 for CIFAR100 and Tiny ImageNet, respectively. The networks are optimized using SGD with an initial learning rate of 0.03, a momentum of 0.9, and a weight decay of 0.0001. We adjust the learning rate in the range between 0.03 and 0.0003 by the cosine learning rate scheduler [10]. About data augmentation, training images from CIFAR100 are padded by 4 pixels on all borders and then preprocessed through randomly cropping at 32×32 , randomly horizontal flipping followed by AutoAug [4]. For Tiny ImageNet, we follow the similar augmentation process but use randomly cropping at 64×64 and RandAug [5] instead.

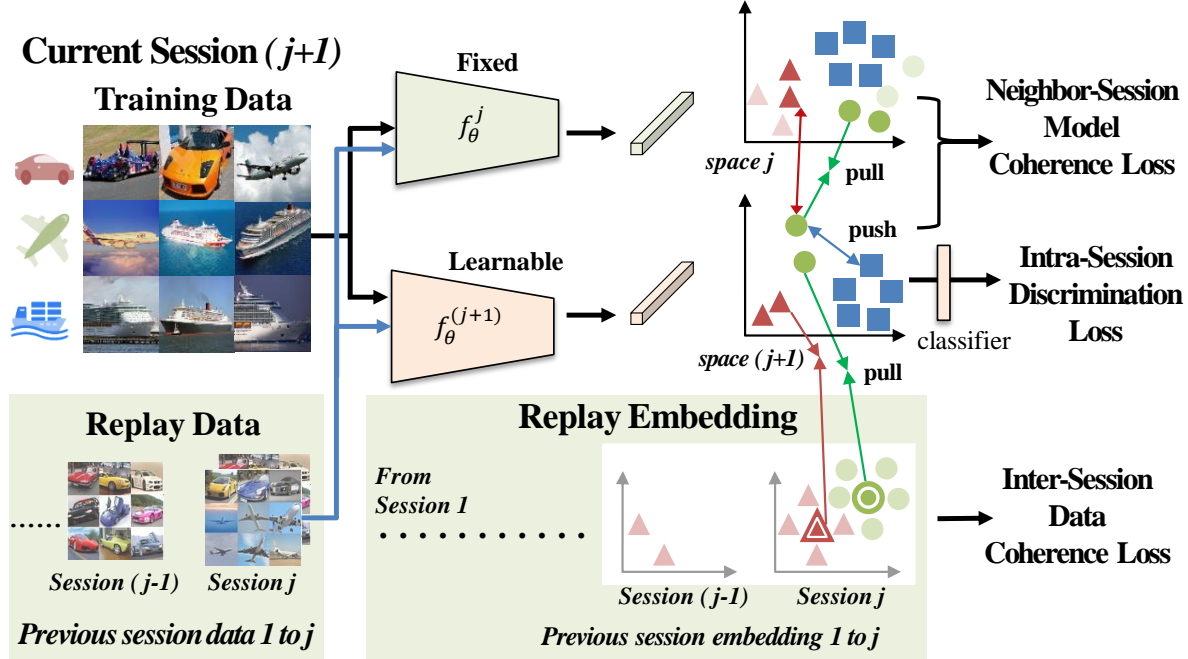


Figure A.2. Overview of our CVS method for CL in General-Incremental setup with long-term backward embedding consistency. In the current session $j + 1$, the training data of the session (together with the replayed data under a budget control) are used in the three loss terms. In addition, the replayed embeddings summarized from previous sessions $1 : j$ serve as historically concentrated attractors to guide the inter-session data-coherence training; it acts as an extension of cross-batch memory [17, 18] to cross-session memory in CL. We introduce a 2-sample-3-embedding strategy in a triplet for distillation learning across neighbor sessions to enforce the model coherence. Note that we omit the replayed data to simplify the illustration. We use a L2-normalized embedding in classification [20] to provide the intra-session discriminating capability, and normalized embedding is adopted in all three loss terms. Our approach is simple but effective in all three CL setups, and we provide the first study on general-incremental setup in CL.



Figure A.3. Sample images of fine-grained datasets. The first pair is from Stanford Dog, the latter pair is from iNaturalist 2017, and the final four images are from Product-10K.

B.2. Details on the Fine-grained Datasets

Some fine-grained samples (Stanford Dog, iNat-M, and Product-M datasets) are shown in Fig. A.3. They have only subtle changes between classes and are more demanding for retrieval. Following similar experimental settings mentioned in [11, 15], we finetune the ImageNet-pretrained ResNet-50 for 100 epochs using SGD with a small fixed learning rate of 0.0001. The batch size is 64 by default but 32 for Product-M. We follow [15] to preprocess the training images by randomly resizing and cropping them to 224×224 with random horizontal flipping. At testing time, we emphasize the object by central cropping of 224×224 from the 256×256 resized image for feature extraction.

B.3. Reimplementation details

We re-implement the LWF [9], MMD [2], and BCT [16] in our experiments. All loss terms are equal weighting to meet the balance between previously learned information and new knowledge. For MMD, the maximum mean discrepancy loss is solely used for blurry setup, and an additional knowledge distillation loss is applied to the novel class data according to the original definition under the disjoint and general-incremental setups. For BCT, we make some modifications to suit for different setups as it is not a CL solution and requires all old class samples collected so far. We use all the old samples from the seen classes for the blurry and general-incremental setups, and employ the replayed data mined by iCaRL [14] for the disjoint setup.

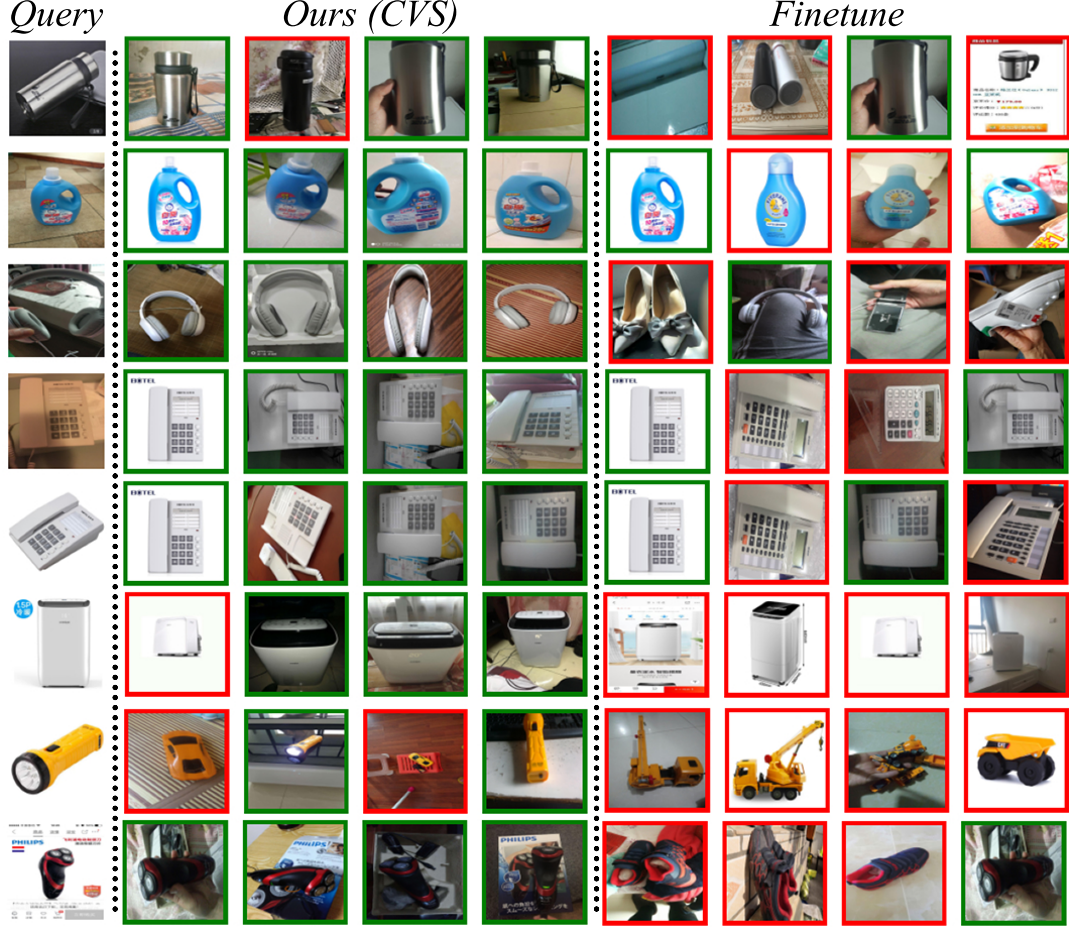


Figure A.4. **Qualitative comparison** of the top-4 results using our method (CVS) and Finetune on the Product-M dataset. The correct and incorrect matches are highlighted in green and red, respectively.

B.4. Qualitative Study

We demonstrate some qualitative results in Fig. A.4. Our method maintains the backward consistency and captures fine-grained characteristics of the particular object. *E.g.*, the first row shows that our method can retrieve visually similar bottles, but Fine-tune yields the results with perturbation.

C. Ablation study on different losses for neighbor-session model coherence

As referenced in Section 4.3 of the main paper, we provide a complete experiment result as follows.

Comparisons to Other Embedding Distillations: We perform an in-depth analysis of different metric losses for the loss term $\mathbf{L}_{j:j+1}^m$, as shown in Table 3. Dark Knowledge [8] minimizes the KL divergence on the classifier side. Absolute MLKD [19] performs distillation at the penultimate layer output. By estimating relational structural information given a mini-batch, RKD [12] reduces the Huber loss using

pairwise euclidean distance difference between two models (i.e., Distancewise RKD) and angle from three points (i.e., Anglewise RKD). Following the same spirit, Relative MLKD [19] uses the difference of Frobenius norm instead, and DarkRank [3] re-estimates the similarity ranks using listwise relationships. Except for Dark Knowledge, we use their official Github implementation for fair comparisons. We tune hyperparameter α at $\{1, 10\}$ for the best AR@1 value at the validation phase for importance weighting. As can be seen in Table C.3, our method provides the most favorable results at AR@ k when k is small (1 or 2) on all setups, and only slightly inferior to Relative MLKD and DarkRank at AR@4 on the blurry and general-incremental setups, respectively.

Comparisons to Other Sample Mining Strategy: We examine different mining techniques because $L_{j;j+1}^m$ is computed based on the sampled triplets. Instead of our easiest positive mining, we use the hardest positive mining (*i.e.*, **BatchHard**, a hardest-positive-hardest-negative online mining strategy mentioned in [7]) for fair comparisons (Table. C.4). We implement BatchHard with a balanced batch sampler to enforce each batch containing at least two samples per class for forming sufficient valid triplets. Unfortunately, BatchHard obtains the worst result or even lower than the one with $L_{j;j+1}^m$ disabled. We attribute this result to the misleading guidance due to a large variation between feature spaces; thus, mining triplets according to the cross-session distance is unreliable. On the other hand, our strategy gains the best result by forming the positive samples using the outputs from two models without explicit mining. Therefore, our triplet mining design is simple but effective, and easy to implement.

D. More Technical Analysis

Results with different memory budgets: The experiment about the influence of the buffer size is provided in Table. C.5. We observe that CVS consistently beats other replay-based methods on AR@1 using half the budget under the general-incremental setup.

Multiple run results: We present the results of averaging 3 runs on all selected datasets. With Table. C.6, most methods show no significant deviation ($< 1\%$) except for RWalk in Tiny ImageNet. But, this doesn't change any conclusions in our main paper.

Acknowledgement

This work was supported in part by the MOST (Ministry of Science and Technology) under the grant MOST 110-2634-F-002-050, MOST 110-2634-F-006-022, and MOST 110-2221-E-002-185 -MY2.

References

- [1] Jihwan Bang, Heesu Kim, YoungJoon Yoo, Jung-Woo Ha, and Jonghyun Choi. Rainbow memory: Continual learning with a memory of diverse samples. In *CVPR*, pages 8218–8227, 2021.
- [2] Wei Chen, Yu Liu, Weiping Wang, Tinne Tuytelaars, Erwin M Bakker, and Michael Lew. On the exploration of incremental learning for fine-grained image retrieval. In *BMVC*, 2020.
- [3] Yuntao Chen, Naiyan Wang, and Zhaoxiang Zhang. Dark-rank: Accelerating deep metric learning via cross sample similarities transfer. In *AAAI*, pages 2852–2859, 2018.
- [4] Ekin D. Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V. Le. Autoaugment: Learning augmentation strategies from data. In *CVPR*, pages 113–123, 2019.
- [5] Ekin D. Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V. Le. Randaugment: Practical automated data augmentation with a reduced search space. In *CVPRW*, pages 3008–3017, 2020.
- [6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016.
- [7] Alexander Hermans, L. Beyer, and B. Leibe. In defense of the triplet loss for person re-identification. *arXiv preprint arXiv:1703.07737*, 2017.
- [8] Geoffrey Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural network. In *NeurIPS*, 2015.
- [9] Zhizhong Li and Derek Hoiem. Learning without forgetting. *PAMI*, 40:2935–2947, 2018.
- [10] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. In *ICLR*, 2017.
- [11] Kevin Musgrave, Serge J. Belongie, and Ser-Nam Lim. A metric learning reality check. In *ECCV*, pages 681–699, 2020.
- [12] Wonpyo Park, Dongju Kim, Yan Lu, and Minsu Cho. Relational knowledge distillation. In *CVPR*, pages 3962–3971, 2019.
- [13] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. In *NeurIPS*, 2017.
- [14] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H. Lampert. icarl: Incremental classifier and representation learning. In *CVPR*, pages 5533–5542, 2017.
- [15] Karsten Roth, Timo Milbich, Samarth Sinha, Prateek Gupta, Bjorn Ommer, and Joseph Paul Cohen. Revisiting training strategies and generalization performance in deep metric learning. In *ICML*, pages 8242–8252, 2020.
- [16] Yantao Shen, Yuanjun Xiong, Wei Xia, and Stefano Soatto. Towards backward-compatible representation learning. In *CVPR*, pages 6367–6376, 2020.
- [17] Xun Wang, Haozhi Zhang, Weilin Huang, and Matthew R. Scott. Cross-batch memory for embedding learning. In *CVPR*, pages 6387–6396, 2020.
- [18] Huei-Fang Yang, Kevin Lin, and Chu-Song Chen. Cross-batch reference learning for deep classification and retrieval. In *ACM MM*, pages 1237–1246, 2016.
- [19] Lu Yu, Vacit Oguz Yazici, Xialei Liu, Joost van de Weijer, Yongmei Cheng, and Arnau Ramisa. Learning metrics from teachers: Compact networks for image embedding. In *CVPR*, pages 2902–2911, 2019.
- [20] Andrew Zhai and Hao-Yu Wu. Classification is a strong baseline for deep metric learning. In *BMVC*, 2019.

	CIFAR100			Tiny ImageNet			Dog			iNat-M			Product-M		
	AR@1	AR@2	AR@4	AR@1	AR@2	AR@4	AR@1	AR@2	AR@4	AR@1	AR@2	AR@4	AR@1	AR@2	AR@4
Joint Train	81.97	84.6	86.82	47.45	51.94	55.99	86.98	91.08	93.99	75.85	79.97	83.65	79.36	83.83	87.73
Finetune	60.79	64.73	68.14	31.11	35.87	40.8	82.7	88.32	92.23	67.75	72.68	77	70.99	76.26	81.16
BCT	58.31	62.2	65.66	30.17	34.64	39.29	81.73	87.49	91.55	67.34	72.07	75.99	70.48	75.67	80.47
LWF	65.53	70.91	75.31	32.22	38.23	44.56	83.25	<u>89.29</u>	<u>93.04</u>	68.51	73.71	78.12	72.95	78.38	83.1
MMD	65.51	70.22	74.33	33.35	38.21	43.06	83.2	88.98	92.71	68.58	73.48	77.79	72.89	78.21	82.96
EWC	60.89	64.86	68.2	27.86	32.67	37.78	81.64	88.7	92.82	66.3	71.4	75.82	66.01	72.53	78.24
RWalk	<u>69.9</u>	<u>73.37</u>	<u>76.39</u>	33.83	38.43	42.97	82.41	88.31	92.43	68.77	73.82	78.16	68.71	74.64	80.13
Rainbow	68.4	71.56	74.2	<u>37.53</u>	<u>41.64</u>	<u>45.44</u>	82.78	89.04	93.23	68.68	73.22	77.21	69.39	75.19	80.34
Ours (CVS)	73.95	76.73	78.84	38.78	42.38	45.89	84.71	89.4	92.61	72.57	76.39	79.87	75.47	80.36	84.68
CVS w/o replay	67.61	71.3	74.39	36.39	40.32	44.2	<u>83.51</u>	88.96	92.45	<u>71.31</u>	<u>75.25</u>	<u>78.56</u>	<u>74.19</u>	<u>78.95</u>	<u>83.27</u>

Table C.1. Results on our general incremental setup. The champion is highlighted in bold and the runner-up is underlined in red.

	CIFAR100			Tiny ImageNet			Dog			iNat-M			Product-M		
	AR@1	AR@2	AR@4	AR@1	AR@2	AR@4	AR@1	AR@2	AR@4	AR@1	AR@2	AR@4	AR@1	AR@2	AR@4
CVS: $L^c + L^m + L^d$	73.95	76.73	78.74	38.78	42.38	45.89	84.71	89.4	92.61	<u>72.57</u>	76.39	79.87	<u>75.47</u>	80.36	84.68
CVS w/o replay	67.61	71.3	74.39	36.39	40.32	<u>44.2</u>	83.51	88.96	92.45	71.31	75.25	78.56	74.19	78.95	83.27
$L^c + L^d$	<u>72.16</u>	<u>74.67</u>	<u>76.7</u>	<u>38.4</u>	<u>41.22</u>	43.87	<u>84.37</u>	<u>89.07</u>	<u>92.46</u>	72.61	<u>76.27</u>	<u>79.45</u>	75.6	<u>80.27</u>	<u>84.33</u>
$L^c + L^d$ w/o replay	64.5	68.12	71.14	32.82	36.2	39.55	83.29	88.77	92.38	70.94	74.62	77.86	73.96	78.69	83.02
$L^c + L^m$	63.79	68.13	72.02	32.1	37.15	42.78	82.65	88.43	92.29	67.79	72.64	76.97	71.77	76.92	81.78
L^c	60.79	64.73	68.14	31.11	35.87	40.8	82.7	88.32	92.23	67.75	72.68	77	70.99	76.26	81.16

Table C.2. Ablation results on our general incremental setup. The champion is highlighted in bold and the runner-up is underlined in red.

	Disjoint			Blurry			General		
	AR@1	AR@2	AR@4	AR@1	AR@2	AR@4	AR@1	AR@2	AR@4
None	69.85	72.67	74.58	44.33	46.14	47.72	72.16	74.67	76.7
Dark Knowledge	69.41	73.09	76.05	45.89	47.76	49.5	73.18	76.34	78.68
Absolute MLKD	66.04	70.36	74.27	46.9	48.89	50.77	72.43	75.72	78.47
Relative MLKD	66.64	71.2	75.35	44.1	48.25	52.33	70.91	75.19	78.92
Anglewise RKD	70.45	72.65	74.67	45.07	46.99	48.67	72.65	74.67	76.33
Distancewise RKD	70.1	72.29	74.11	45.09	47.01	48.81	72.83	75.17	77.14
Hard DarkRank	66.45	71.03	75.11	44.8	48.2	51.35	72.33	76.15	79.56
Ours (CVS)	71.47	74.8	77.51	47.47	49.86	52.17	73.95	76.73	78.84

Table C.3. Replace $L^m_{j:j+1}$ with different metric distillation losses on CIFAR100. *None* disables $L^m_{j:j+1}$ for a simple baseline.

	Disjoint			Blurry			General		
	AR@1	AR@2	AR@4	AR@1	AR@2	AR@4	AR@1	AR@2	AR@4
None	69.85	72.67	74.58	44.33	46.14	47.72	72.16	74.67	76.7
BatchHard	60.83	63.91	66.67	39.24	42.59	45.7	68.1	71.02	73.35
Ours (CVS)	71.47	74.8	77.51	47.47	49.86	52.17	73.95	76.73	78.84

Table C.4. Compute $L^m_{j:j+1}$ with different mining strategies on CIFAR100.

	CIFAR100			Tiny ImageNet		
	AR@1	AR@2	AR@4	AR@1	AR@2	AR@4
RWalk (0.5× budget)	66.9	70.32	73.59	31.42	36.39	41.41
Rainbow (0.5× budget)	65.21	69.33	72.85	34.78	39.8	44.83
CVS (0.5× budget)	71.99	74.94	77.25	37.72	41.58	45.59
RWalk (1× budget)	69.9	73.37	76.39	33.83	38.43	42.97
Rainbow (1× budget)	68.4	71.56	74.2	37.53	41.64	45.44
CVS (1× budget)	73.95	76.73	78.84	38.78	42.38	45.89

Table C.5. Results with different memory budgets on our general incremental setup. 1× budget represents the continual learner with a memory buffer size of 2000 samples in CIFAR100 and 4000 samples in TinyImageNet. 0.5× indicates using half the budget.

	CIFAR100			Tiny ImageNet			Dog			iNat-M			Product-M		
	AR@1	AR@2	AR@4	AR@1	AR@2	AR@4	AR@1	AR@2	AR@4	AR@1	AR@2	AR@4	AR@1	AR@2	AR@4
BCT	58.01±0.26	62.13±0.12	65.84±0.17	30.11±0.43	34.89±0.36	39.72±0.51	81.77±0.07	87.33±0.23	91.34±0.24	67.17±0.31	71.88±0.22	75.95±0.16	70.55±0.1	75.83±0.17	80.67±0.4
LWF	65.39±0.24	70.6±0.36	75.06±0.29	32.17±0.59	38.3±0.41	44.81±0.26	83.45±0.33	89.19±0.21	93.01±0.15	68.34±0.3	73.71±0.04	78.2±0.07	73.11±0.14	78.42±0.04	83.01±0.19
MMD	65.15±0.32	69.98±0.22	74.18±0.21	33.3±0.3	38.47±0.27	43.57±0.46	83.53±0.32	89.01±0.15	92.65±0.06	68.63±0.31	73.6±0.26	77.85±0.06	72.81±0.27	78.16±0.17	82.84±0.23
EWC	60.9±0.14	65.1±0.55	68.76±0.92	28.01±0.88	33.21±0.66	38.55±0.66	82.17±0.77	89.1±0.68	93.21±0.54	66.18±0.1	71.53±0.12	76.02±0.17	66.33±0.28	72.53±0.04	78.04±0.17
RWalk	69.69±0.27	73.16±0.19	76.23±0.14	34.05±1.17	38.77±1.07	43.48±0.98	82.53±0.27	88.67±0.33	92.81±0.34	68.85±0.17	73.9±0.11	78.15±0.05	69.05±0.31	74.45±0.18	79.36±0.67
Rainbow	68.18±0.2	71.32±0.27	73.96±0.23	37.47±0.51	41.86±0.38	45.89±0.44	83.06±0.29	89.1±0.1	93.15±0.09	68.85±0.17	73.43±0.33	77.41±0.36	69.66±0.25	74.91±0.3	79.76±0.56
Ours (CVS)	73.81±0.15	76.48±0.26	78.58±0.27	38.56±0.33	42.03±0.07	46.01±0.23	84.64±0.07	89.5±0.09	92.85±0.21	72.7±0.13	76.47±0.13	79.74±0.23	75.77±0.35	80.49±0.21	84.7±0.19

Table C.6. Average results of 3 runs on our general incremental setup.