

# Class-Incremental Learning with Cross-Space Clustering and Controlled Transfer

Arjun Ashok<sup>1,2</sup> , K J Joseph<sup>1</sup> , and Vineeth N Balasubramanian<sup>1</sup> 

<sup>1</sup> Indian Institute of Technology Hyderabad

<sup>2</sup> PSG College Of Technology, Coimbatore

arjun.ashok.psg@gmail.com, {cs17m18p100001, vineethnb}@iith.ac.in  
<https://cscct.github.io>

**Abstract.** In class-incremental learning, the model is expected to learn new classes continually while maintaining knowledge on previous classes. The challenge here lies in preserving the model’s ability to effectively represent prior classes in the feature space, while adapting it to represent incoming new classes. We propose two distillation-based objectives for class incremental learning that leverage the structure of the feature space to maintain accuracy on previous classes, as well as enable learning the new classes. In our first objective, termed **cross-space clustering** (CSC), we propose to use the feature space structure of the previous model to characterize directions of optimization that maximally preserve the class: directions that all instances of a specific class should collectively optimize towards, and those directions that they should collectively optimize away from. Apart from minimizing forgetting, such a class-level constraint indirectly encourages the model to reliably cluster all instances of a class in the current feature space, and further gives rise to a sense of “herd-immunity”, allowing all samples of a class to jointly combat the model from forgetting the class. Our second objective termed **controlled transfer** (CT) tackles incremental learning from an important and understudied perspective of inter-class transfer. CT explicitly approximates and conditions the current model on the semantic similarities between incrementally arriving classes and prior classes. This allows the model to learn the incoming classes in such a way that it maximizes positive forward transfer from similar prior classes, thus increasing plasticity, and minimizes negative backward transfer on dissimilar prior classes, whereby strengthening stability. We perform extensive experiments on two benchmark datasets, adding our method (CSCCT) on top of three prominent class-incremental learning methods. We observe consistent performance improvement on a variety of experimental settings.

**Keywords:** Incremental Learning, Continual Learning, Knowledge Distillation, Transfer Learning.

## 1 Introduction

Incremental learning is a paradigm of machine learning where learning objectives are introduced to a model incrementally in the form of *phases* or *tasks*, and the

model must dynamically learn new tasks while maintaining knowledge on previously seen tasks. The differences of this setup from a *static training* scenario is that the model is not allowed to *retrain* from scratch on encountering new tasks, and no task information is available upfront. A fundamental challenge in incremental learning is in the stability-plasticity trade-off [37], where stability relates to maintaining accuracy on the previous tasks, while plasticity relates to learning the current task effectively. In their naive form, deep learning models are too plastic; the model changes significantly during training, and incurs *catastrophic forgetting* [16] of old tasks when exposed to new ones.

Class-incremental learning (CIL) [42, 36] is a specific sub-paradigm of incremental learning where tasks are composed of *new classes* and we seek to learn a *unified* model that can *represent and classify* all classes seen so far equally well. The main challenge in class-incremental learning lies in how knowledge over the long stream of classes can be consolidated at every phase. Regularization-based methods [25, 4, 58, 9] quantify and preserve important parameters corresponding to prior tasks. Another set of approaches [34, 10, 14, 50] focus on modifying the learning algorithm to ensure that gradient updates do not conflict with the previous tasks. In dynamic architecture methods [55, 52, 31, 1, 41], the network architecture is modified by expansion or masking when encountering new tasks during learning. Replay-based methods [42, 19, 13, 51, 6, 7, 33, 8, 30, 2, 47, 46, 54] store a subset of each previous task in a separate memory, and replay the tasks when learning a new one, to directly preserve the knowledge on those tasks. A wide variety of such replay methods have been developed recently, and have attained promising results in the CIL setting. A number of these methods [42, 19, 13, 7, 2, 47, 8, 30] use variants of knowledge distillation [18], where the model and its predictions corresponding to the previous task are utilized to prevent the current task’s model from diverging too much from its previous state.

Our work herein falls under distillation-based methods. Prior work has advocated for utilizing distillation to directly constrain an example’s position or angle using its position in the previous feature space [19], to preserve pooled convolutional outputs of an instance [13], or to maintain the distribution of logits that the model’s classifier outputs on the data [2, 42]. We argue that preserving the features or predictions of a model on independent individual instances are only useful to a certain extent, and do not characterize and preserve properties of a class captured globally by the model as a whole. *Class-level* semantics may be more important to be preserved in the *class-incremental* learning setting, to holistically prevent individual classes from being forgotten. To this end, we develop an objective termed **Cross-Space Clustering** (CSC) that characterizes *entire regions* in the feature space that classes should stay away from, and those that a class should belong to, and distills this information to the model. Our objective indirectly establishes multiple goals at once: (i) it encourages the model to cluster all instances of a given class; (ii) ensures that these clusters are well-separated; and (iii) regularizes to preserve *class cluster positions* as a single entity in the feature space. This provides for a class-consolidated distillation

objective, prodding instances of a given class to “unite” and thus prevent the class from being forgotten.

Next, as part of our second objective, we tackle the class-incremental-learning problem from a different perspective. While all prior distillation objectives seek better ways to preserve properties of the learned representations in the feature space [42, 19, 13, 7, 2, 47, 8, 30], we believe that controlling *inter-class transfer* is also critical for class-incremental learning. This comes from the observation that forgetting often results from *negative backward transfer* from new classes to previous classes, and plasticity is ensured when there is *positive forward transfer* from prior classes to new ones [34]. To this end, we develop an objective called **Controlled Transfer** (CT) that controls and regularizes transfer of features between classes at a fine-grained level. We formulate an objective that approximates the relative similarity between an incoming class and all previous classes, and conditions the current task’s model on these estimated similarities. This encourages *new classes* to be situated optimally in the feature space, ensuring maximal positive transfer and minimal negative transfer.

A unique characteristic of our objectives is their ability to extend and enhance existing distillation-based CIL methodologies, without any change to their methodologies. We verify this by adding our objectives to three prominent and state-of-the-art CIL methods that employ distillation in their formulation: iCARL [42], LUCIR [19] and PODNet [13]. We conduct thorough experimental evaluations on benchmark incremental versions of large-scale datasets like CIFAR-100 and ImageNet subset. We perform a comprehensive evaluation of our method, considering a wide variety of experimental settings. We show that our method consistently improves incremental learning performance across datasets and methods, at no additional cost. We further analyze and present ablation studies on our method, highlighting the contribution of each of our components.

## 2 Related Work

### 2.1 Incremental Learning

In an incremental setting, a model is required to consistently learn new tasks, without compromising performance on old tasks. Incremental learning methodologies can be split into five major categories, each of which we review below.

**Regularization-based** methods focus on quantifying the importance of each parameter in the network, to prevent the network from excessively changing the important parameters pertaining to a task. These methods include EWC [25], SI [58], MAS [4] and RWalk [9]. A recent method ELI [22], introduces an energy based implicit regularizer which helps to alleviate forgetting.

**Algorithm-based** methods seek to avoid forgetting from the perspective of the network training algorithm. They modify gradients such that updates to weights do not deteriorate performance on previously seen tasks. Methods such as GEM [34], A-GEM [10], OGD [14] and NSCL [50] fall under this category. Meta-learning based methods [21, 26] are also useful while learning continually.

**Architecture-based** methods modify the network architecture dynamically to fit more tasks, by expanding the model by adding more weights [55, 52], or masking and allocating subnetworks to specific tasks [45], or by gating the parameters dynamically using a task identifier [1].

**Exemplar-based** methods (also called replay-based or rehearsal methods) assume that a small subset of every class can be stored in a memory at every phase. They replay the seen classes later along with the incoming new classes, directly preventing them from being forgotten. One set of works focus on reducing the recency bias due to the new classes being in majority at every phase [51, 6, 19, 7]. Another set of works focus on optimizing which samples to choose as exemplars to better represent the class distributions [33, 5].

**Distillation-based** methods use the model learned until the previous task as a teacher and provide extra supervision to the model learning the current tasks (the student). Since the data of the previous tasks are inaccessible, these methods typically enforce distillation objectives on the current data [32, 12], data from an exemplar memory [19, 13, 42, 7, 2, 27, 23], external data [30] or synthetic data [59]. Since our method falls under this category, we extend our discussion on related methods below.

Early works in this category distill logit scores [32, 42] or attention maps [12] of the previous model. iCARL [42] proposes to enforce distillation on new tasks as well exemplars from old tasks, along with herding selection of exemplars and nearest-mean-of-exemplars (NME) based classification. GD [30] calibrates the confidence of the model’s outputs using external unlabelled data, and propose to distill the calibrated outputs instead. LUCIR [19] introduces a less-forget constraint that encourages the orientation of a sample in the current feature space to be similar to its orientation in the old feature space. Apart from that, LUCIR proposes to use cosine-similarity based classifiers and a margin ranking loss that mines hard negatives from the new classes to better separate the old class, additionally avoiding ambiguities between old and new classes. PODNet [13] preserves an example’s representation throughout the model with a spatial distillation loss. SS-IL [2] show that general KD preserves the recency bias in the distillation, and propose to use task-wise KD. Co2L [8] introduces a contrastive learning based self-supervised distillation loss that preserves the exact feature relations of a sample with its augmentations and other samples from the dataset. GeoDL [47] introduces a term that enhances knowledge distillation by performing KD across low-dimensions path between the subspaces of the two models, considering the gradual shift between models.

The main difference of our cross-space clustering objective from these works is that we do not optimize to preserve the properties of individual examples, and instead preserve the previously learned semantics or properties of each *class* in a *holistic manner*. Our formulation takes into account the global position of a class in the feature space, and optimizes all samples of the class towards the same region, making the model indifferent to instance-level semantics. Further, classes are supervised with specific “negative” regions all over the feature space, also intrinsically giving rise to better separation between class clusters.

Our controlled transfer objective, on the other hand, attempts to regularize transfer between tasks. MER [43], an algorithm-based method is related to our high-level objective. MER works in the online continual learning setup, combining meta-learning [15, 38] with replay. Their method optimizes such that the model receives weight updates are restricted to those directions that agree with the gradients of prior tasks. Our objective proposes to utilize the *structure* of the feature space of the previous model to align the current feature space, in order to maximize transfer. Our novelty here lies in how we explicitly approximate *inter-class semantic similarities* in a *continual task stream*, and utilize them to appropriately position new tasks representations, regularizing transfer.

## 2.2 Knowledge Distillation

Hinton et al. [18] introduced knowledge distillation (KD) in their work as a way to transfer knowledge from an ensemble of teacher networks to a smaller student network. They show that there is dark knowledge in the logits of a trained network that can give more structure about the data, and use them as soft targets to train the student. Since then, a number of other works have explored variants of KD. Attention Transfer [57] focused on the attention maps of the network instead of the logits, while FitNets [44] also deal with intermediate activation maps of a network. Several other papers have enforced criteria based on multi-layer representations of a network [53, 20, 24, 3, 28].

Among these, our controlled transfer objective shares similarities with a few works that propose to exploit the mutual relation between data samples for distillation. Tung and Mori [49] propose a distillation objective that enforces an L2 loss in the student that constraints the similarities between activation maps of two examples to be consistent with those of the teacher. Relational KD [39] additionally propose to preserve the angle formed by the three examples in the feature space by means of a triplet distillation objective. Extending this direction, Correlation Congruence [40] models relations between examples through higher-order kernels, to enforce the same objectives with better relation estimates.

The difference of our controlled transfer objective from these works lies in the high-level objective in the context of the incremental learning setting, as well as the low-level formulation in terms of the loss objective. All the above works propose to use sample relations in the feature space to provide additional supervision to a student model by regularizing the feature relations of the student. The main challenge in incremental learning is how we can reduce the effect that a new class has on the representation space, to minimize forgetting.

Our objective also exploits sample relations in the feature space, however our novelty lies in how we estimate a measure of relative similarity between an *unseen class* and each previously seen class, and utilize them to control where the *new samples* are located in the embedding space, in relation to the old samples. Our specific formulation indirectly promotes forward transfer of features from prior classes similar to the new class, and simultaneously prevents negative backward transfer of features from the new class to dissimilar previous classes .

### 3 Method

We briefly introduce the problem setting in Sec. 3.1. Next, we explain in detail our learning objectives in Sec. 3.2 and Sec. 3.3 respectively, and discuss the overall objective function in Sec. 3.4.

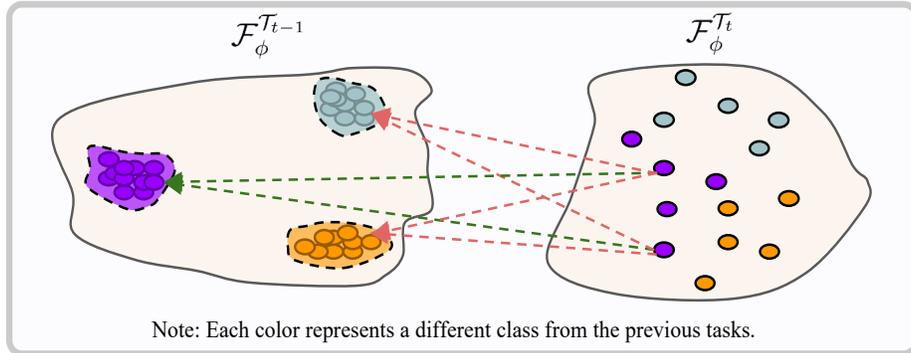
#### 3.1 Problem Setting

In the incremental learning paradigm, a set of tasks  $\mathcal{T}_t = \{\tau_1, \tau_2, \dots, \tau_t\}$  is introduced to the model over time, where  $\mathcal{T}_t$  represents the tasks that the model has seen until time step  $t$ .  $\tau_t$  denotes the task introduced at time step  $t$ , which is composed of images and labels sampled from its corresponding data distribution:  $\tau_t = (\mathbf{x}_i^{\tau_t}, y_i^{\tau_t}) \sim p_{data}^{\tau_t}$ . Each task  $\tau_t$  contains instances from a disjoint set of classes.  $\mathcal{F}^{\mathcal{T}_t}$  denotes the model at time step  $t$ . Without loss of generality,  $\mathcal{F}^{\mathcal{T}_t}$  can be expressed as a composition of two functions:  $\mathcal{F}^{\mathcal{T}_t}(\mathbf{x}) = (\mathcal{F}_\phi^{\mathcal{T}_t} \circ \mathcal{F}_\theta^{\mathcal{T}_t})(\mathbf{x})$ , where  $\mathcal{F}_\phi^{\mathcal{T}_t}$  represents a feature extractor, and  $\mathcal{F}_\theta^{\mathcal{T}_t}$  denotes a classifier. The challenge in incremental learning is to learn a model that can represent and classify all seen classes equally well, at any point in the task stream.

While training the model  $\mathcal{F}^{\mathcal{T}_t}$  on the task  $\tau_t$ , the model does not have access to all the data from previous tasks. Exemplar-based methods [42, 5, 51, 33, 6] sample a very small subset of each task data  $e_t \in \tau_t$  at the end of every task  $\tau_t$  and store it in a memory buffer  $\mathcal{M}_t = \{e_1, e_2, \dots, e_t\}$ , which contains a subset of data from all tasks seen until time  $t$ . When learning a new task at time step  $t$ , the task’s data  $\tau_t$  is combined with samples from the memory containing exemplars of each previous task  $\mathcal{M}_{t-1}$ . Therefore, the dataset that the model has at time step  $t$  is  $\mathcal{D}_t = \tau_t \cup \mathcal{M}_{t-1}$ . In distillation-based methods, we assume access to the previous model  $\mathcal{F}^{\mathcal{T}_{t-1}}$  which has learned the stream of tasks  $\mathcal{T}_{t-1}$ . The model  $\mathcal{F}^{\mathcal{T}_{t-1}}$  is frozen and not updated, and is instead used to guide the learning of the current model  $\mathcal{F}^{\mathcal{T}_t}$ . Effectively utilising the previous model is key to balancing stability and plasticity. Excess constraints tied to the model can prevent the current task from being learned, and poor constraints can lead to easy forgetting of previous tasks.

#### 3.2 Cross-Space Clustering

Our *cross-space clustering (CSC)* objective alleviates forgetting by distilling class-level semantics and inducing tight clusters in the feature space. CSC leverages points across the entire feature space of the previous model  $\mathcal{F}^{\mathcal{T}_{t-1}}$ , to identify regions that a class is optimized to stay within, and other harmful regions that it is prevented from drifting towards. We illustrate our cross-space clustering objective in Fig. 1.



**Fig. 1.** We illustrate our cross-space clustering (CSC) objective. We show instances from 3 classes from the stream  $\mathcal{T}_{t-1}$ , and their positions in  $\mathcal{F}_\phi^{\mathcal{T}_{t-1}}$  and  $\mathcal{F}_\phi^{\mathcal{T}_t}$  respectively. Classes are well represented in  $\mathcal{F}_\phi^{\mathcal{T}_{t-1}}$ , however their representations are dispersed in the  $\mathcal{F}_\phi^{\mathcal{T}_t}$ . Here we illustrate the constraint imposed on an instance of the **violet** class, based on the cluster position of its own class (indicated by the **green** arrows) and the positions of every other class (indicated by the **red** arrows). Note how the exact same constraint is applied on all instances of a class (illustrated here with 2 instances of the **violet** class). Best viewed in color.

Consider that the model  $\mathcal{F}^{\mathcal{T}_t}$  is trained on mini-batches  $\mathcal{B} = \{x_i, y_i\}_{i=1}^k$  sampled from  $D_t$ . Our cross-space clustering objective enforces the following loss on the model:

$$L_{Cross-Cluster} = \frac{1}{k^2} \sum_{i=1}^k \sum_{j=1}^k \left( 1 - \cos(\mathcal{F}_\phi^{\mathcal{T}_t}(x_i), \mathcal{F}_\phi^{\mathcal{T}_{t-1}}(x_j)) \right) * \text{ind}(y_i == y_j) \quad (1)$$

where  $\text{ind}$  is an indicator function that returns **1** when its inputs are equal and **-1** otherwise<sup>3</sup>, and  $\cos(a, b)$  denotes the cosine similarity between two vectors  $a$  and  $b$ .

**Physical Interpretation:** For pairs of samples  $x_i$  and  $x_j$ ,  $\mathcal{F}^{\mathcal{T}_t}(x_i)$  is enforced to minimize cosine distance with  $\mathcal{F}^{\mathcal{T}_{t-1}}(x_j)$  when they are of the same class ( $y_i == y_j$ ), and maximize cosine distance with  $\mathcal{F}^{\mathcal{T}_{t-1}}(x_j)$  when they are of different classes ( $y_i \neq y_j$ ). We expand upon the objective and its implications separately below.

**Explanation:** Consider that there are  $l$  examples of class  $n$  in the considered batch of size  $k$ , and hence  $k - l$  samples belonging to classes other than  $n$ .

Sample  $x_i$  from class  $n$  is allowed to see the previous feature positions of all of the  $l$  samples of the same class in  $\mathcal{B}$ , and is regularized to be *equally close* to all these positions. Since multiple positions are used in the previous feature space and equal constraints are applied, points only see an approximation of its

<sup>3</sup> Note how this is different from a typical indicator function that returns 0 when the inputs are not equal.

class (cluster) in the previous feature space, and do not see individual feature positions. This inherently removes the dependency of the distillation on the specific position of the sample within its class, and instead optimizes the sample to move towards a point that can preserve the class as a whole.

Next, *every sample*  $x_i$  belonging to a class  $n$  in the batch is given the *exact same constraints* with no difference. In our case, all samples belonging to a class are optimized towards the mean of the class’s embeddings in the previous space. This leads to all of them being optimized *jointly* to a single stark region belonging to their class. Repeating this process for several iterations implicitly leads to model to implicitly *cluster* all samples of a class in the *current* feature space  $\mathcal{F}^{\mathcal{T}_t}$  in the specific characterized regions. With respect to clustering, an important point is that our loss is *cross-space* in the sense, it does not encourage clustering of features of a class using features from the same model [8], that would neither exploits prior knowledge about the inter-class distances, nor imposes any constraints on the location of the classes. Our formulation instead encourages a model to keep all these clusters at specific points provided by the previous feature space, thereby directly distilling and preserving the *cluster positions* in the feature space as well. Hence, our objective uses approximate cluster positions from  $\mathcal{F}^{\mathcal{T}_{t-1}}$  to in-turn cluster samples at specific positions in  $\mathcal{F}^{\mathcal{T}_t}$ . Since all samples are optimized towards the same region, all points of the class are optimized to *unite* and jointly protect and preserve the class. Such a formulation gives rise to a sense of *herd-immunity* of classes from forgetting, which better preserves the classes as the model drifts.

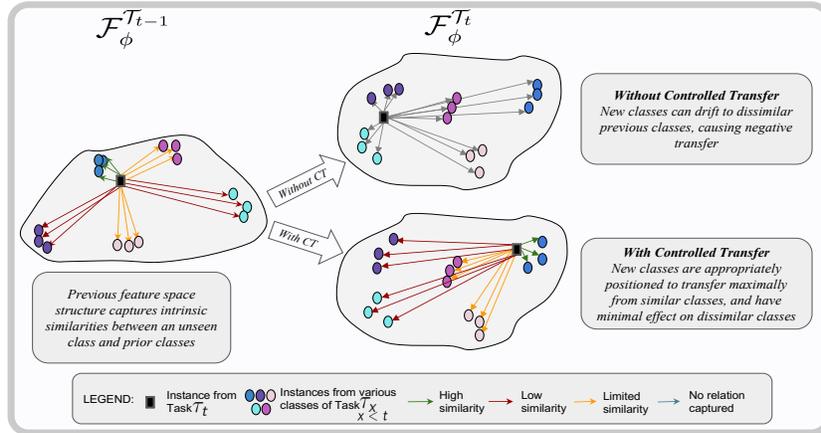
Finally, with very few exemplars stored per-class in the memory, our objective proposes to maximally utilize the memory<sup>4</sup> as well as the current task, leveraging them to identify *negative regions* that an instance is maintained to lie away from. Through multiple iterations of optimization,  $x_i$  belonging to class  $n$  is enforced to stay equally away from the positions of all other  $k - l$  examples from the entire previous space. This indirectly tightens the cluster of class  $n$  in  $\mathcal{F}^{\mathcal{T}_t}$  along multiple directions in the feature space.

**Differences from prior work:** Prior distillation-based methods [19, 13, 47] only apply *sample-to-sample* cross-space constraints, to preserve the representational properties of the previous space. The core difference of our method from all others lies in how it applies *class-to-region* constraints. Here, *class* denotes how all samples of a class are jointly optimized with the same constraints, and *region* denotes how the samples are optimized towards and away from specific *regions* instead of towards individual points.

### 3.3 Controlled Transfer

Catastrophic forgetting has been previously characterized to arise due to a variety of reasons - the inability to access enough data of previous tasks [51, 19, 6],

<sup>4</sup> A batch of sufficient size typically contains at least one sample from each previous class, serving as a rough approximation of the memory



**Fig. 2.** We illustrate our **controlled transfer** objective. We show the positions of instances from five random classes taken from previous tasks  $\tau_x; x < t$ , and one unseen incoming class from the current task  $\tau_t$ , in  $\mathcal{F}_\phi^{\mathcal{T}_{t-1}}$  and  $\mathcal{F}_\phi^{\mathcal{T}_t}$  respectively. With our objective, the new task instances are regularized to position themselves appropriately, to prevent negative transfer to dissimilar classes, and to encourage positive transfer from similar classes (best viewed in color)

change in important parameters [25, 58, 4], representation drift [32, 13, 8], conflicting gradients [34, 10, 14, 50] and insufficient capacity of models [55, 52, 1]. However, all these works ignore the semantic similarities between tasks and their relation to forgetting. We argue that knowing the degree of semantic similarity between two classes can, in fact, be very useful in incremental learning: When a previous class is *dissimilar* to the class currently being learned, the model must learn to treat the previous class distinctively and minimally impact it, so that the semantic specialities of that class are not erased. Conversely, when there is a previous class which is *similar* to the class currently being learned, the model must maximally transfer features from that class, to learn the current class in the best possible way. With these goals, we propose an incremental learning objective that *explicitly quantifies* inter-class similarities, and *controls* transfer between classes in every phase. Fig. 2 illustrates our controlled transfer objective.

**Notation:** We first describe the general notation that we use to denote the similarity between samples in a space. Consider two samples  $x_i$  and  $x_j$  from a dataset  $D_k$ , and a model  $\mathcal{F}^{\mathcal{T}_k}$ . We denote the similarity between  $x_i$  and  $x_j$  computed on the feature space of  $\mathcal{F}^{\mathcal{T}_k}$  as  $z_{x_i, x_j}^{\mathcal{T}_k} = \cos(\mathcal{F}_\phi^{\mathcal{T}_k}(x_i), \mathcal{F}_\phi^{\mathcal{T}_k}(x_j))$  where  $\cos(a, b)$  denotes cosine similarity between two vectors  $a$  and  $b$ . We denote the the normalized distribution of similarities that an individual sample  $x_i$  has with *every sample* in  $D_k$ , in the feature space of  $\mathcal{F}^{\mathcal{T}_k}$  as

$$H_{x_i, D_k, T}^{\mathcal{T}_k} = \left\{ \frac{(z_{x_i, x_j}^{\mathcal{T}_k} / T)}{\sum_{g=1}^{|D_k|} (z_{x_i, x_g}^{\mathcal{T}_k} / T)} \right\}_{j=1}^{|D_k|} \quad (2)$$

where  $T$  is the temperature that is used to control the entropy of the distribution.  $H_{x_i, D_k, T}^{\mathcal{T}_k}$  is a row matrix, where the value in each column  $j$  of the matrix denotes the normalized similarity between  $x_i$  and  $x_j$ , relative to every sample in the dataset  $D_k$ .

**Formulation:** We first aim to estimate the similarities between a *new class*  $\mathcal{C}_{new} \in \tau_t$  and every previously seen class  $\mathcal{C}_{old} \in \tau_k \in \mathcal{T}_{t-1}$ .  $\mathcal{C}_{old}$  is well represented the model  $\mathcal{F}^{\mathcal{T}_{t-1}}$ ; the new class  $\mathcal{C}_{new}$  has not yet been learned by any model. It is not possible to use the drifting feature space of  $\mathcal{F}^{\mathcal{T}_t}$  to represent  $\mathcal{C}_{new}$ ; even representing  $\mathcal{C}_{new}$  once it has been learned by  $\mathcal{F}^{\mathcal{T}_t}$  would heavily bias the representations towards  $\mathcal{C}_{new}$  due to the well-known recency bias [51, 2]. Our formulation instead proposes to utilize the *dark knowledge* that the *previous model* possesses about an *unseen class*: if the representations of an unseen class  $\mathcal{C}_{new}$  lie relatively close to the class representations of a previous class in  $\mathcal{F}^{\mathcal{T}_{t-1}}$ , it indicates that the two classes share semantic similarities. On the other hand, if the representations of an unseen class  $\mathcal{C}_{new}$  lie relatively far from a previous class in  $\mathcal{F}^{\mathcal{T}_{t-1}}$ , it indicates that the two classes do not have any semantic features in common. Note how the similarities for  $\mathcal{C}_{new}$  given by  $\mathcal{F}^{\mathcal{T}_{t-1}}$  are *unbiased* as the model has never seen  $\mathcal{C}_{new}$ , and hence can be used as approximations to the semantic similarities. We propose to use these approximate similarities captured by  $\mathcal{F}^{\mathcal{T}_{t-1}}$  in our objective explained below.

Consider a mini-batch of  $B_n^{\mathcal{T}_t}$  of size  $s$  that contains samples  $\{(x_i^{\mathcal{T}_t}, y_i^{\mathcal{T}_t})\}$  randomly sampled from  $D_t$ . This mini-batch  $B_n^{\mathcal{T}_t}$  is composed of  $p$  samples from the current task denoted by  $P = (x_i^{\mathcal{T}_t}, y_i^{\mathcal{T}_t})_{i=1}^p$ , and  $q$  samples taken from the memory, denoted by  $Q = (x_i^{\mathcal{T}_k}, y_i^{\mathcal{T}_k})_{i=1}^q$ , where  $k < t$ . In an effort to control the transfer between a new and an old sample, our objective regularizes the normalized similarity (closeness) that a sample from the current task  $(x_i^{\mathcal{T}_t}, y_i^{\mathcal{T}_t}) \in \tau_t$  has with every sample from any previous class  $(x_i^{\mathcal{T}_k}, y_i^{\mathcal{T}_k})$ , where  $k < t$ . This is enforced by minimizing the KL Divergence of the similarity distribution of  $x_i^{\mathcal{T}_t} \in P$  over  $Q$ , in the *current space*  $\mathcal{F}_\phi^{\mathcal{T}_t}$ , with the similarity distribution computed in the *previous space*  $\mathcal{F}_\phi^{\mathcal{T}_{t-1}}$ , as follows

$$L_{Transfer} = \frac{1}{p} \sum_{i=1}^p KL(H_{x_i, Q, T}^{\mathcal{T}_t} || H_{x_i, Q, T}^{\mathcal{T}_{t-1}}) \quad (3)$$

This loss modifies the position of the current classes in the current feature space  $\mathcal{F}_\phi^{\mathcal{T}_t}$  such that they have *high similarity* with (lie close to) prior classes that are *very similar*. This encourages *positive forward transfer* of features to the current classes from selected previous classes that are similar, as both their embeddings are optimized to have high similarity in the current space. This helps the model learn the current task better by leveraging transferred features, and lessens the impact that the new task has on the representation space. Conversely, as the embeddings are optimized to have *low similarity* with (lie far from) those previous classes that are *dissimilar* to it, it discourages (negative) backward transfer from the current classes to those dissimilar classes. Consequently, the features of these specific classes are further shielded from being erased, leading to

**Table 1.** The table shows results on **CIFAR100** when our method is added to three top-performing approaches [42, 19, 13]. The red subscript highlights the relative improvement.  $\mathcal{B}$  denotes the number of classes in the first task.  $\mathcal{C}$  denotes the number of classes in every subsequent task.

Dataset	CIFAR100					
Settings	$\mathcal{B} = 50$			$\mathcal{B} = \mathcal{C}$		
Methods	$\mathcal{C} = 1$	$\mathcal{C} = 2$	$\mathcal{C} = 5$	$\mathcal{C} = 1$	$\mathcal{C} = 2$	$\mathcal{C} = 5$
iCaRL [42]	43.39	48.31	54.42	30.92	36.80	44.19
iCaRL + CSCCT	46.15 <sub>+2.76</sub>	51.62 <sub>+3.31</sub>	56.75 <sub>+2.33</sub>	34.02 <sub>+3.1</sub>	<b>39.60</b> <sub>+2.8</sub>	46.45 <sub>+2.26</sub>
LUCIR [19]	50.26	55.38	59.40	25.40	31.93	42.28
LUCIR + CSCCT	52.95 <sub>+2.69</sub>	56.49 <sub>+1.13</sub>	62.01 <sub>+2.61</sub>	28.12 <sub>+2.72</sub>	34.96 <sub>+3.03</sub>	44.03 <sub>+1.55</sub>
PODNet [13]	56.88	59.98	62.66	33.58	36.68	45.27
PODNet + CSCCT	<b>58.80</b> <sub>+1.92</sub>	<b>61.10</b> <sub>+1.12</sub>	<b>63.72</b> <sub>+1.06</sub>	<b>36.23</b> <sub>+2.65</sub>	39.3 <sub>+2.62</sub>	<b>47.8</b> <sub>+2.53</sub>

the semantics of those classes being preserved more in the current space, which directly results in lesser forgetting of those classes.

### 3.4 Final Objective Function

The independent nature of our objectives make them suitable to be applied on top of any existing method to improve its performance. Our final objective combines  $L_{Cross-Cluster}$  (1) and  $L_{Transfer}$  (3) with appropriate coefficients:

$$L_{CSCCT} = L_{method} + \alpha * L_{Cross-Cluster} + \beta * L_{Transfer} \quad (4)$$

where  $L_{method}$  denotes the objective function of the base method used, and  $\alpha$  and  $\beta$  are loss coefficients for each of our objectives respectively. We term our method **CSCCT**, indicating **Cross-Space Clustering and Controlled Transfer**.

## 4 Experiments and Results

We conduct extensive experiments adding our method to three prominent methods in class-incremental learning [42, 19, 13].

**Protocols:** Prior work has experimented with two CIL protocols: **a)** training with half the total number of classes in the first task, and equal number of classes in each subsequent task [19, 13, 51], and **b)** training with the same number of classes in each task, including the first [42, 2, 7]. The first setting has the advantage of gaining access to strong features in the first task, while the second tests an extreme continual learning setting; both these are plausible in a real-world incremental classification scenario. We experiment with both these protocols to demonstrate the applicability of our method. On CIFAR100, classes are grouped into 1, 2 and 5 classes per task. On ImageNet-Subset, the classes are split into 2, 5 and 10 classes per task. Hence, we experiment on both *long streams of small tasks*, as well as *short streams of large tasks*.

**Datasets and Evaluation Metric:** Following prior works [19, 13, 42, 7], we test on the incremental versions of CIFAR-100 [29] and ImageNet-Subset [42]. CIFAR100 contains 100 classes, with 500 images per class, and each of

**Table 2.** The table shows results on **ImageNet-Subset** when our method is added to three top-performing approaches [42, 19, 13]. The red subscript highlights the relative improvement.  $\mathcal{B}$  denotes the number of classes in the first task.  $\mathcal{C}$  denotes the number of classes in every subsequent task.

Dataset	ImageNet-Subset					
Settings	$\mathcal{B} = 50$			$\mathcal{B} = \mathcal{C}$		
Methods	$\mathcal{C} = 2$	$\mathcal{C} = 5$	$\mathcal{C} = 10$	$\mathcal{C} = 2$	$\mathcal{C} = 5$	$\mathcal{C} = 10$
iCaRL [42]	55.81	57.34	65.97	40.75	55.92	60.93
iCaRL + CSCCT	57.01 <sub>+1.2</sub>	58.37 <sub>+1.03</sub>	66.82 <sub>+0.8</sub>	42.46 <sub>+1.71</sub>	57.45 <sub>+1.53</sub>	62.60 <sub>+1.67</sub>
LUCIR [19]	60.44	66.55	70.18	36.84	46.40	56.78
LUCIR + CSCCT	61.52 <sub>+1.08</sub>	67.91 <sub>+1.36</sub>	71.33 <sub>+1.15</sub>	37.86 <sub>+1.02</sub>	47.55 <sub>+1.15</sub>	58.07 <sub>+1.29</sub>
PODNet [13]	67.27	73.01	75.32	44.94	58.23	66.24
PODNet + CSCCT	68.91 <sub>+1.64</sub>	74.35 <sub>+1.34</sub>	76.41 <sub>+1.09</sub>	46.06 <sub>+1.12</sub>	59.43 <sub>+1.2</sub>	67.49 <sub>+1.25</sub>

dimension  $32 \times 32$ . ImageNet-Subset is a subset of the ImageNet-1k dataset [11], and contains 100 classes, with over 1300 images per class. Each image is of size  $224 \times 224$ . All our results denote average incremental accuracy. We follow the original papers in their inference methodology: On LUCIR [19] and PODNet [13], classification is performed as usual using the trained classifier, while on iCaRL [42], classification is based on nearest-mean-of-exemplars.

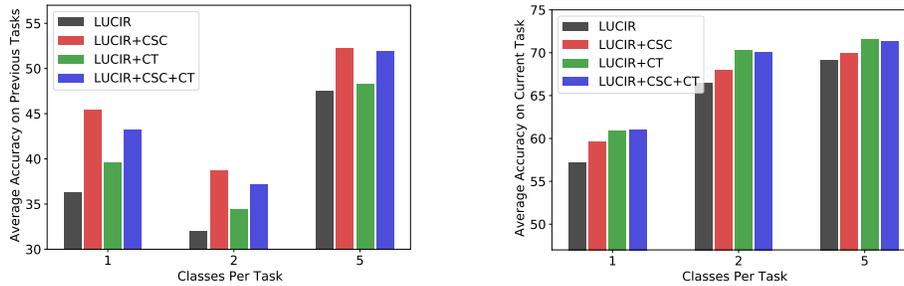
## 5 Implementation Details

Following prior work [19, 13], we use a ResNet-32 and ResNet-18 on CIFAR100 and ImageNet-Subset respectively. On CIFAR100, we use a batch size of 128 and train for 160 epochs, with an initial learning rate of  $4e^{-1}$  that is decayed by 0.1 at the 80<sup>th</sup> and 120<sup>th</sup> epochs respectively. On ImageNet-Subset, we use a batch size of 64 and train for 90 epochs, with an initial learning rate of  $2e^{-2}$  that is decayed by 0.1 at the 30<sup>th</sup> and 60<sup>th</sup> epochs respectively. We use *herding selection* [42] for exemplar sampling, and an exemplar memory size of 20.

### 5.1 Quantitative Results

We add our method to three state-of-the-art class-incremental learning methodologies: iCaRL [42], LUCIR [19] and PODNet [13]. Table 1 showcases results on CIFAR100, and Table 2 showcases results on ImageNet-Subset. We see a consistent improvement across all settings and methods when CSCCT is added to them. Specifically, on CIFAR100, adding CSCCT to iCaRL [42], LUCIR [19] and PODNet [13] provides strong relative improvement of 2.76%, 2.28% and 1.99% respectively averaged across all settings, while on the much more high-dimensional ImageNet-Subset, adding our method to the respective baselines provides consistent relative improvements of 1.32%, 1.17% and 1.35%.

Evaluating iCaRL [42], LUCIR [19] and PODNet [13] on the equal class protocol show that LUCIR [19] suffers from a severe performance degradation due to its inherent reliance on a large initial task, while iCaRL [42] and PODNet [13] do not. On CIFAR100, simply adding our method to iCaRL [42] gives it



**Fig. 3.** Average accuracy on previous tasks (APT) and average accuracy on the current task (ACT), plotted across various settings on the CIFAR-100 dataset

strong boosts of 2.2% – 3.1% in this setting, bringing it much closer to the state-of-the-art PODNet [13]. Overall, our method improves performance consistently across both settings, showing that our formulation *does not rely on a large initial task to learn strong representations*.

## 6 Ablation Study and Analysis

### 6.1 Effect of Each Component on Average Incremental Accuracy

In Table 3, we ablate each component of our objective. Each of our objectives can improve accuracy independently. In particular, CSC is more effective when the number of classes per task is extremely low, while CT stands out in the improvement it offers when there are more classes per task. Overall, combining our objectives achieves the best performance across all settings.

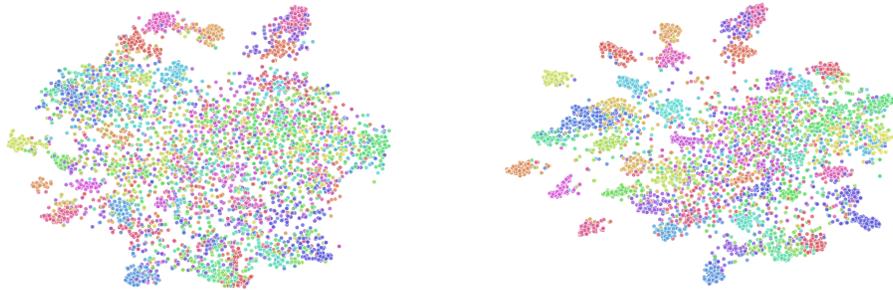
**Table 3.** Ablating each objective on CIFAR100. Maroon denotes  $2^{nd}$  best result.

Settings	$\mathcal{B} = 50$			$\mathcal{B} = \mathcal{C}$		
	Methods	$\mathcal{C} = 1$	$\mathcal{C} = 2$	$\mathcal{C} = 5$	$\mathcal{C} = 1$	$\mathcal{C} = 2$
LUCIR [19]	50.26	55.38	59.4	25.4	31.93	42.28
LUCIR + CSC	<b>52.04</b>	<b>55.95</b>	60.45	<b>27.16</b>	32.89	42.98
LUCIR + CT	51.5	55.87	<b>61.97</b>	26.53	<b>33.98</b>	<b>43.69</b>
LUCIR + CSCCT	<b>52.95</b>	<b>56.49</b>	<b>62.01</b>	<b>28.12</b>	<b>34.96</b>	<b>44.03</b>

### 6.2 Effect of Each Component on Stability/Plasticity

To further investigate how each component is useful specifically in the incremental learning setup, we look into how each component improves the stability and plasticity of the model under various settings. The left plot of Fig. 3 shows the average accuracy on previous tasks (denotes as APT). This serves as an indicator of the **stability** of the model. Mathematically, APT can be expressed as

$$APT = \frac{\sum_{t=2}^T \left( \frac{\sum_{k=1}^{t-1} Acc(\mathcal{F}^T, \tau_k)}{t-1} \right)}{T-1} \quad (5)$$



**Fig. 4.** T-SNE [35] visualizations of the base 50 classes of CIFAR100 in the embedding space, after all 100 classes have been learned (**Left:** LUCIR [19], **Right:** LUCIR+CSC)

where  $Acc(\mathcal{F}^{\mathcal{T}_t}, \tau_k)$  denotes accuracy of model  $\mathcal{F}^{\mathcal{T}_t}$  on the test set of task  $k$ .

The right plot of Fig. 3 shows the average accuracy on the current task (denoted as ACT). This serves as an indicator of the **plasticity** of the model. ACT is expressed as

$$ACT = \frac{\sum_{t=1}^T Acc(\mathcal{F}^{\mathcal{T}_t}, \tau_t)}{T} \quad (6)$$

Across all considered settings, *both of our objectives* increase *stability* as well as *plasticity* of the base model. However, one can see that the effect of the **CSC objective** is much more pronounced on the **stability** of the model. This aligns with intuition that the CSC helps in preserving previous classes better in the representation space. At the same time, the **CT objective** impacts the **plasticity** consistently more than the CSC objective, as it mainly aims at appropriately positioning the current task samples to maximize transfer.

### 6.3 Embedding Space Visualization

In Fig. 4, we present T-SNE [35] visualizations of the embedding space, without and with our Cross-Space Clustering (CSC) objective (1). The 50 classes learned in the initial task are plotted in the embedding spaces of both models, once all the 100 classes have been learned. It is seen that the CSC objective results in better clusters of prior classes in the feature space, compared to the baseline. The number of overlapping classes are reduced to a significant extent, as our objective ensures that the clusters are well-separated.

## 7 Conclusion

We introduced two complementary distillation-based objectives for class incremental learning. Our first objective called *cross-space clustering* positions classes appropriately in the embedding space and enables classes to counteract forgetting jointly. Our second objective called *controlled transfer* controls the positive and negative transfer between classes by estimating and utilizing inter-class relationships. We perform extensive experiments across a wide range of experimental settings to showcase the effectiveness of our objectives.

**Acknowledgements:** We are grateful to the Department of Science and Technology, India, as well as Intel India for the financial support of this project through the IMPRINT program (IMP/2019/000250) as well as the DST ICPS Data Science Cluster program. KJJ thanks TCS for their PhD Fellowship. We also thank the anonymous reviewers and Area Chairs for their valuable feedback in improving the presentation of this paper.

## References

1. Abati, D., Tomczak, J.M., Blankevoort, T., Calderara, S., Cucchiara, R., Bejnordi, B.E.: Conditional channel gated networks for task-aware continual learning. *CVPR (2020)* [2](#), [4](#), [9](#)
2. Ahn, H., Kwak, J., Lim, S.F., Bang, H., Kim, H., Moon, T.: Ss-il: Separated softmax for incremental learning. *ICCV (2021)* [2](#), [3](#), [4](#), [10](#), [11](#)
3. Ahn, S., Hu, S.X., Damianou, A.C., Lawrence, N.D., Dai, Z.: Variational information distillation for knowledge transfer. *CVPR (2019)* [5](#)
4. Aljundi, R., Babiloni, F., Elhoseiny, M., Rohrbach, M., Tuytelaars, T.: Memory aware synapses: Learning what (not) to forget. In: *ECCV (2018)* [2](#), [3](#), [9](#)
5. Aljundi, R., Lin, M., Goujaud, B., Bengio, Y.: Gradient based sample selection for online continual learning. In: *NeurIPS (2019)* [4](#), [6](#)
6. Belouadah, E., Popescu, A.D.: Il2m: Class incremental learning with dual memory. *ICCV (2019)* [2](#), [4](#), [6](#), [8](#)
7. Castro, F.M., Marín-Jiménez, M.J., Mata, N.G., Schmid, C., Karteek, A.: End-to-end incremental learning. *ECCV (2018)* [2](#), [3](#), [4](#), [11](#)
8. Cha, H., Lee, J., Shin, J.: Co2l: Contrastive continual learning. *ICCV (2021)* [2](#), [3](#), [4](#), [8](#), [9](#)
9. Chaudhry, A., Dokania, P.K., Ajanthan, T., Torr, P.H.: Riemannian walk for incremental learning: Understanding forgetting and intransigence. In: *ECCV (2018)* [2](#), [3](#)
10. Chaudhry, A., Ranzato, M., Rohrbach, M., Elhoseiny, M.: Efficient lifelong learning with a-GEM. In: *ICLR (2019)* [2](#), [3](#), [9](#)
11. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: *CVPR (2009)* [12](#)
12. Dhar, P., Singh, R.V., Peng, K.C., Wu, Z., Chellappa, R.: Learning without memorizing. *CVPR (2019)* [4](#)
13. Douillard, A., Cord, M., Ollion, C., Robert, T., Valle, E.: Podnet: Pooled outputs distillation for small-tasks incremental learning. In: *ECCV (2020)* [2](#), [3](#), [4](#), [8](#), [9](#), [11](#), [12](#), [13](#), [18](#)
14. Farajtabar, M., Azizan, N., Mott, A., Li, A.: Orthogonal gradient descent for continual learning. In: *AISTATS (2020)* [2](#), [3](#), [9](#)
15. Finn, C., Abbeel, P., Levine, S.: Model-agnostic meta-learning for fast adaptation of deep networks. In: *ICML (2017)* [5](#)
16. French, R.M.: Catastrophic forgetting in connectionist networks. *Trends in Cognitive Sciences (1999)* [2](#)
17. Goyal, P., Dollár, P., Girshick, R., Noordhuis, P., Wesolowski, L., Kyrola, A., Tulloch, A., Jia, Y., He, K.: Accurate, large minibatch sgd: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677 (2017)* [18](#)
18. Hinton, G.E., Vinyals, O., Dean, J.: Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531 (2015)* [2](#), [5](#)

19. Hou, S., Pan, X., Loy, C.C., Wang, Z., Lin, D.: Learning a unified classifier incrementally via rebalancing. *CVPR (2019)* [2](#), [3](#), [4](#), [8](#), [11](#), [12](#), [13](#), [14](#), [18](#), [19](#)
20. Huang, Z., Wang, N.: Like what you like: Knowledge distill via neuron selectivity transfer. arXiv preprint arXiv:abs/1707.01219 (2017) [5](#)
21. Javed, K., White, M.: Meta-learning representations for continual learning. *NeurIPS (2019)* [3](#)
22. Joseph, K., Khan, S., Khan, F.S., Anwer, R.M., Balasubramanian, V.N.: Energy-based latent aligner for incremental learning. In: *CVPR (2022)* [3](#)
23. Joseph, K., Khan, S., Khan, F.S., Balasubramanian, V.N.: Towards open world object detection. In: *CVPR (2021)* [4](#)
24. Kim, J., Park, S., Kwak, N.: Paraphrasing complex network: Network compression via factor transfer. *NeurIPS (2018)* [5](#)
25. Kirkpatrick, J., Pascanu, R., Rabinowitz, N., Veness, J., Desjardins, G., Rusu, A.A., Milan, K., Quan, J., Ramalho, T., Grabska-Barwinska, A., Hassabis, D., Clopath, C., Kumaran, D., Hadsell, R.: Overcoming catastrophic forgetting in neural networks. *PNAS (2017)* [2](#), [3](#), [9](#)
26. KJ, J., N Balasubramanian, V.: Meta-consolidation for continual learning. *NeurIPS (2020)* [3](#)
27. KJ, J., Rajasegaran, J., Khan, S., Khan, F.S., Balasubramanian, V.N.: Incremental object detection via meta-learning. *IEEE TPAMI (2021)* [4](#)
28. Koratana, A., Kang, D., Bailis, P., Zaharia, M.A.: Lit: Learned intermediate representation training for model compression. In: *ICML (2019)* [5](#)
29. Krizhevsky, A., Hinton, G.: Learning multiple layers of features from tiny images. In: *Citeseer (2009)* [11](#)
30. Lee, K., Lee, K., Shin, J., Lee, H.: Overcoming catastrophic forgetting with unlabeled data in the wild. arXiv preprint arXiv:1903.12648 (2019) [2](#), [3](#), [4](#)
31. Li, X., Zhou, Y., Wu, T., Socher, R., Xiong, C.: Learn to grow: A continual structure learning framework for overcoming catastrophic forgetting. In: *ICML (2019)* [2](#)
32. Li, Z., Hoiem, D.: Learning without forgetting. *IEEE TPAMI (2018)* [4](#), [9](#)
33. Liu, Y., Liu, A., Su, Y., Schiele, B., Sun, Q.: Mnemonics training: Multi-class incremental learning without forgetting. *CVPR (2020)* [2](#), [4](#), [6](#)
34. Lopez-Paz, D., Ranzato, M.: Gradient episodic memory for continual learning. *NeurIPS (2017)* [2](#), [3](#), [9](#)
35. van der Maaten, L., Hinton, G.: Visualizing data using t-sne. *JMLR (2008)* [14](#)
36. Masana, M., Liu, X., Twardowski, B., Menta, M., Bagdanov, A.D., van de Weijer, J.: Class-incremental learning: survey and performance evaluation on image classification. arXiv preprint arXiv:2010.15277 (2021) [2](#)
37. Mermillod, M., Bugaiska, A., Bonin, P.: The stability-plasticity dilemma: investigating the continuum from catastrophic forgetting to age-limited learning effects. *Frontiers in Psychology (2013)* [2](#)
38. Nichol, A., Achiam, J., Schulman, J.: On first-order meta-learning algorithms. arXiv preprint arXiv:1803.02999 (2018) [5](#)
39. Park, W., Kim, D., Lu, Y., Cho, M.: Relational knowledge distillation. *CVPR (2019)* [5](#)
40. Peng, B., Jin, X., Liu, J., Zhou, S., Wu, Y., Liu, Y., Li, D., Zhang, Z.: Correlation congruence for knowledge distillation. *ICCV (2019)* [5](#)
41. Rajasegaran, J., Hayat, M., Khan, S.H., Khan, F.S., Shao, L.: Random path selection for continual learning. In: *NeurIPS (2019)* [2](#)
42. Rebuffi, S.A., Kolesnikov, A., Sperl, G., Lampert, C.H.: icarl: Incremental classifier and representation learning. *CVPR (2017)* [2](#), [3](#), [4](#), [6](#), [11](#), [12](#)

43. Riemer, M., Cases, I., Ajemian, R., Liu, M., Rish, I., Tu, Y., , Tesauro, G.: Learning to learn without forgetting by maximizing transfer and minimizing interference. In: ICLR (2019) 5
44. Romero, A., Ballas, N., Kahou, S.E., Chassang, A., Gatta, C., Bengio, Y.: Fitnets: Hints for thin deep nets. ICLR (2015) 5
45. Serrà, J., Surís, D., Miron, M., Karatzoglou, A.: Overcoming catastrophic forgetting with hard attention to the task. ICML (2018) 4
46. Shin, H., Lee, J.K., Kim, J., Kim, J.: Continual learning with deep generative replay. In: NeurIPS (2017) 2
47. Simon, C., Koniusz, P., Harandi, M.: On learning the geodesic path for incremental learning. CVPR (2021) 2, 3, 4, 8
48. Smith, S.L., Kindermans, P.J., Le, Q.V.: Don't decay the learning rate, increase the batch size. In: ICLR (2018) 18
49. Tung, F., Mori, G.: Similarity-preserving knowledge distillation. ICCV (2019) 5
50. Wang, S., Li, X., Sun, J., Xu, Z.: Training networks in null space of feature covariance for continual learning. In: CVPR (2021) 2, 3, 9
51. Wu, Y., Chen, Y., Wang, L., Ye, Y., Liu, Z., Guo, Y., Fu, Y.R.: Large scale incremental learning. CVPR (2019) 2, 4, 6, 8, 10, 11
52. Yan, S., Xie, J., He, X.: Der: Dynamically expandable representation for class incremental learning. CVPR (2021) 2, 4, 9
53. Yim, J., Joo, D., Bae, J.H., Kim, J.: A gift from knowledge distillation: Fast optimization, network minimization and transfer learning. CVPR (2017) 5
54. Yin, H., Molchanov, P., Li, Z., Álvarez, J.M., Mallya, A., Hoiem, D., Jha, N.K., Kautz, J.: Dreaming to distill: Data-free knowledge transfer via deepinversion. CVPR (2020) 2
55. Yoon, J., Yang, E., Lee, J., Hwang, S.J.: Lifelong learning with dynamically expandable networks. In: ICLR (2018) 2, 4, 9
56. You, Y., Li, J., Reddi, S., Hseu, J., Kumar, S., Bhojanapalli, S., Song, X., Demmel, J., Keutzer, K., Hsieh, C.J.: Large batch optimization for deep learning: Training bert in 76 minutes. arXiv preprint arXiv:1904.00962 (2019) 18
57. Zagoruyko, S., Komodakis, N.: Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. ICLR (2017) 5
58. Zenke, F., Poole, B., Ganguli, S.: Continual learning through synaptic intelligence. In: ICML (2017) 2, 3, 9
59. Zhai, M., Chen, L., Tung, F., He, J., Nawhal, M., Mori, G.: Lifelong gan: Continual learning for conditional image generation. ICCV (2019) 4

# Class-Incremental Learning with Cross-Space Clustering and Controlled Transfer: Supplementary Material

## 8 Results on Different Batch Sizes

Our method leverages examples from the entire batch in order to facilitate training. All our results on CIFAR100 in the original paper are reported using a batch size of 128, following prior work [19, 13]. Here, we analyze the efficacy of our method on different batch sizes used during training. Prior work has observed that increasing the batch size tends to slightly reduce the test accuracy, and have proposed various methods to reduce the drop in accuracy [48, 17, 56]. In our experiments, we adopt the linear scaling rule [17], which scales the learning rate in proportion to the batch size.

Table S1 showcases the results on two different experimental settings, across various batch sizes. It can be seen that increasing the batch size increases the performance of our method. This is because larger batch sizes allow using more samples from the memory as well as the current task, providing a broader view of the feature space, which directly benefits our objectives.

**Table S1.** Ablation studies on the Batch Size

Settings		$\mathcal{B} = 50$	$\mathcal{B} = \mathcal{C}$
Methods	Batch Size	$\mathcal{C} = 5$	
LUCIR [19]	128	59.4	42.28
LUCIR + CSCCT		60.01 <b>+2.61</b>	44.03 <b>+1.55</b>
LUCIR [19]	256	57.06	39.58
LUCIR + CSCCT		59.71 <b>+2.91</b>	41.48 <b>+1.90</b>
LUCIR [19]	512	56.19	38.16
LUCIR + CSCCT		59.02 <b>+2.83</b>	40.2 <b>+2.04</b>
LUCIR [19]	1024	54.83	37.97
LUCIR + CSCCT		57.8 <b>+2.97</b>	40.27 <b>+2.3</b>

## 9 Results on Different Exemplar Memory Sizes

The memory size specifies the *exemplars-per-class* that the model can store at the end of each phase. Here, we report results varying the exemplar memory size.

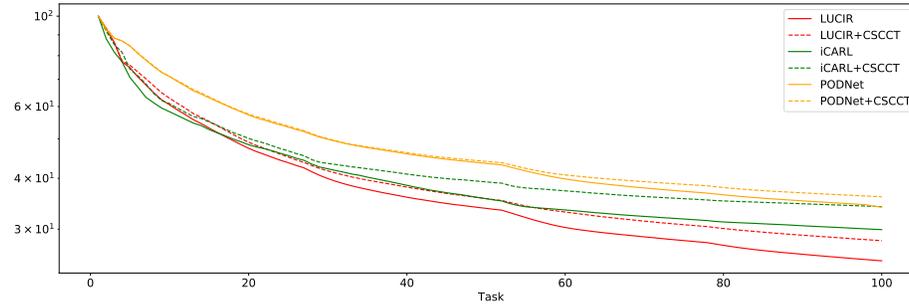
Table S2 showcases the results. Enforcing stricter memory constraints causes a performance drop in LUCIR [19], however, our method still provides strong relative improvements across settings. As the memory size increases, our method offers greater relative improvements.

**Table S2.** Ablation studies on the Memory Size

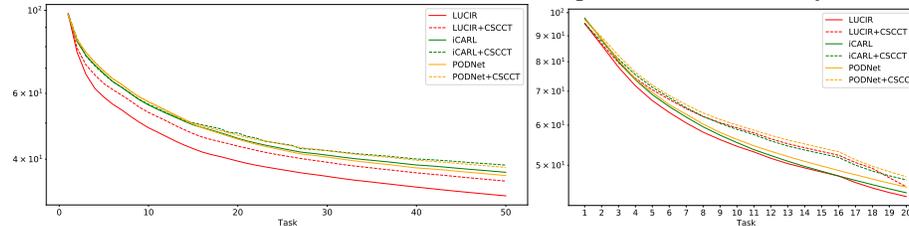
Settings		$B = 50$	$B = C$
Methods	Memory Size	$C = 5$	
LUCIR [19]	10	55.83	39.45
LUCIR + CSCCT	10	57.54 +1.71	40.72 +1.14
LUCIR [19]	20	59.4	42.28
LUCIR + CSCCT	20	60.01 +2.61	44.03 +1.55
LUCIR [19]	30	62.52	46.35
LUCIR + CSCCT	30	65.05 +2.53	48.34 +1.99
LUCIR [19]	40	63.60	50.16
LUCIR + CSCCT	40	66.49 +2.89	52.28 +2.12

## 10 Phase-Wise Plots

Figures A1 and A2 showcase phase-wise plots on three class-incremental learning settings on CIFAR100, on top of multiple baseline methods.



**Fig. A1.** Phase-wise average incremental accuracies on CIFAR100, on the 100 Task Setting with 1 Class per Task. The y-axis is set to log scale for visual clarity.



**Fig. A2.** Phase-wise average incremental accuracies on CIFAR100, on the 50 Task Setting with 2 classes per task (Left) and the 20 Task Setting with 5 classes per task (Right). The y-axis is set to log scale for visual clarity.