



HAL
open science

Pseudo-Labeling for Class Incremental Learning

Alexis Lechat, Stéphane Herbin, Frédéric Jurie

► **To cite this version:**

Alexis Lechat, Stéphane Herbin, Frédéric Jurie. Pseudo-Labeling for Class Incremental Learning. BMVC 2021: The British Machine Vision Conference, Nov 2021, virtuel, United Kingdom. hal-03464466

HAL Id: hal-03464466

<https://hal.science/hal-03464466>

Submitted on 3 Dec 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Pseudo-Labeling for Class Incremental Learning

Alexis Lechat¹²

alexis.lechat@onera.fr

Stéphane Herbin¹

stephane.herbin@onera.fr

Frédéric Jurie¹³

frederic.jurie@safrangroup.com

¹ DTIS, ONERA

Université Paris-Saclay

FR-91123 Palaiseau, France

² Normandie Univ

UNICAEN, ENSICAEN, CNRS, GREYC

14000 Caen, France

³ Safran Tech

Safran

78114 Magny-les-Hameaux, France

Abstract

Class Incremental Learning (CIL) consists in training a model iteratively with limited amount of data from few classes that will never be seen again, resulting in catastrophic forgetting and lack of diversity. In this paper, we address these phenomena by assuming that, during incremental learning, additional unlabeled data are continually available, and propose a Pseudo-Labeling approach for class incremental learning (PLCiL) that makes use of a new adapted loss. We demonstrate that our method achieves better performance than supervised or other semi-supervised methods on standard class incremental benchmarks (CIFAR-100 and ImageNet-100) even when a self-supervised pre-training step using a large set of data is used as initialization. We also illustrate the advantages of our method in a more complex context with fewer labels. The code is available at <https://github.com/alechat/PLCiL>.

1 Introduction

Natural vision systems learn in a continuous way, benefiting from their constant interaction with the environment. Their skills are dynamically updated and accumulated throughout their life. While artificial models such as Deep Neural Networks (DNNs) have now achieved similar or even better performance in several perceptual tasks, their learning process is fundamentally different and relies mainly on supervised *batch training* which requires a large amount of annotated data.

Incrementally training artificial DNNs from an incoming data stream suffers from *catastrophic forgetting* [1, 2]: previously learned skills tend to be less accurate when new ones are integrated in the system. Continual Learning (CL) explores solutions to alleviate this phenomenon. It can be seen as finding a way to solve a *plasticity-stability dilemma* [3]: the model should be flexible enough to dynamically expand its knowledge (plasticity) while ensuring the integrity of previously accumulated knowledge (stability).

In this work, we focus on the *Class Incremental Learning* (CIL) scenario applied to image recognition: new classes gradually appear from the data stream and the total number of categories is not known beforehand. Many studies emulate a data stream by splitting a classification dataset into disjoint batches of several classes [0, 83]. Using this definition, we can see class incremental as analogous to a succession of small-scale supervised batch sessions, with each one learning a small subset of disjoint classes.

There is a large consensus that learning good visual representations is crucial to competitive classification performance [0]. Optimizing a feature extractor via *representation learning* or pre-training with an auxiliary task becomes mandatory when working with datasets of limited size. In particular, *self-supervised visual representation learning* [0, 8, 88] now achieves performance close to fully supervised methods while requiring less labeled data.

In a CL setting, however, methods struggle to learn good representations: this is especially true at the beginning of the learning process, since only a fraction of the data is available at each session. Training a DNN from scratch in a continuous framework requires the feature extractor to be constantly adjusted, resulting in unstable representations. This explains why some authors have taken the easier option of having a large number of classes at the beginning (e.g. 50 classes out of a total of 100) in order to consolidate the representations even before starting the incremental learning of the remaining classes [10, 18, 30].

In this paper, we aim to learn a model truly from scratch, using a semi-supervised representation learning mechanism, assuming that unlabeled data is available throughout the learning process. Note that the idea of using unlabeled data in this context has already been considered in [25, 45].

The underlying intuition behind this proposition is that, with an ideal representation space, the solution to class incremental reduces to allocating unassigned regions in the representation space to the new classes without needing to strongly modify the previous regions. While the advantages expressed in the representation learning literature directly transfer to CL, we also study how semi-supervision provides an answer to the plasticity-stability dilemma. We propose to exploit a process based on pseudo-labeling as a way to combine self-supervision with CIL.

Our contributions are threefold: i) We introduce a mechanism of *Pseudo-Labeling for Class incremental Learning* (PLCiL) and show its benefit in an original learning scheme adapted to semi-supervised CIL that combines 3 losses: a supervised loss, a self-supervised loss using pseudo-labels and a new distillation loss that ensures prediction consistency during CIL sessions. ii) Using PLCiL, we demonstrate that a self-supervision provides regularization against catastrophic forgetting, reaching state-of-the-art performance on class incremental benchmarks. iii) We propose a new class incremental evaluation protocol with even fewer labeled images. We show that our method can still learn the continual task thanks to the semi-supervision while fully supervised CL approaches can hardly compete on such data-scarce problems.

2 Related Work

Our proposed method borrows ideas from two different areas of the literature: that of continual learning and that of representation learning with self-supervision.

Continual Learning [29] refers to different settings [19, 69] such as CIL, which is the focus of this paper, or Task Incremental (TI) learning. The big challenge of CIL is to propose methods that are as insensitive to catastrophic forgetting during learning. The main

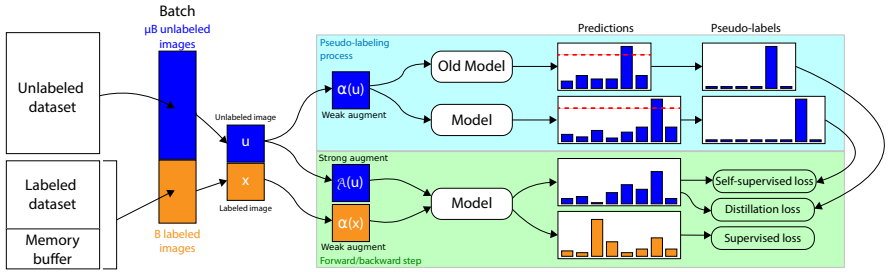


Figure 1: Overview of the proposed incremental training process. The upper part (blue box) represents the pseudo-labeling process dedicated to the automatic generation of a label for an unlabeled u image. The lower part (green box) is a standard supervised forward/backward process, including a 3-term loss function: i) standard supervised cross-entropy, ii) regularization between sessions by distillation, iii) regularization by self-supervision.

approaches are described in the following.

Rehearsal / Replay consists in replaying part of the old data, stored in an episodic memory, and mixes it with new data during learning [5, 12, 21, 27, 33]. The best performing approaches follow this strategy, but their success depends on the number and representativeness of the chosen examples and require a large memory footprint. The limited amount of examples stored induces an imbalance between old classes and newly introduced ones. Models learning with rehearsal are then heavily biased. [15, 18, 41, 46] proposed solutions to compensate for the bias at the classifier level. Rather than simply storing old examples, an alternative is to produce data online by a continuously learned generative model [57, 40].

Another strategy is to make the classifier have outputs close to those of the previous sessions, especially on old data. To this end, Learning without Forgetting [26] has adapted the *Knowledge Distillation (KD)* loss introduced in [10]. The old classifier is the teacher and distills his knowledge to the new classifier, seen as the student. Most of the rehearsal based methods make use of some form of KD [6, 25, 33, 41, 46].

Other interesting strategies are *Parameter Control* [11, 1, 23, 42], *Dynamic Architectures* [51, 56, 42], *Generative Replay* [57, 40] or *Meta-Learning* [21, 52, 52].

Two approaches [25, 45] make the same hypothesis as our method: a large amount of unlabeled data available during the learning sessions. They both implement a self-supervised task: since KD does not require any ground-truth, they leverage unlabeled data by distilling knowledge from two teachers, a model expert on the old classes and a model trained only on the new classes, into a global model performing on all classes. While this kind of self-supervision is an efficient regularization against catastrophic forgetting, these methods do not exploit these additional data to enhance their representations.

Our method belongs to the rehearsal with memory category, but introduces both a new KD scheme and a self-supervised objective that aims to learn better representations throughout the CIL process.

Self-Supervised Representation Learning refers to a particular set of representation learning methods that use *pretext tasks*. Pretext tasks are used to automatically create artificial labels from unlabeled data that can be used to compute an error signal for learning.

A wide range of pretext tasks has been explored in the literature. We refer the reader to [22] for a complete review of the field. We do not detail this literature here because it is less

central to our problem.

Our primary objective is to introduce regularization mechanisms against catastrophic forgetting, and we focus on semi-supervised methods [6, 48] that we believe are good candidates. Among the most powerful recent methods, [58] combines coherence regularization (robustness to transformations) and pseudo-labeling (semi-supervised learning).

Our approach proposes an original way to combine pseudo-labeling with CIL by introducing a specific KD loss ensuring prediction consistency between learning sessions.

3 Pseudo-Labeling for Class incremental learning

The PLCiL method proposed in this paper is based on a standard class incremental paradigm that relies on rehearsal learning with episodic memory [63]. In the classical setting, the available labeled samples during a training session only come from the memory and from the new annotated data. As in [45, 45], we propose to complement this paradigm with the possibility of using unlabeled data, with the motivation that access to unlabeled data at training time is easy (inexpensive) and does not violate the principle of exclusive annotation vocabulary between sessions that typifies class incremental learning.

We present this method in 4 steps: i) its overview (notations, training sessions, prediction model), ii) how the data is organized for each training session, iii) the various loss functions at the heart of our approach, and iv) a discussion and justification of its main components.

3.1 Class Incremental Learning

The class incremental scenario formalizes the incoming labeled data as a stream of subsets $\mathcal{X} = \{X^1, X^2, \dots, X^j, \dots\}$ where each $X^j = \{x_1^j, \dots, x_{n_j}^j\}$ only contains instances of class j . An incremental learning session uses several subsets pooled together and submitted to the network for training. Once the session is done, the pool of data is discarded and the associated classes will not appear again in the data stream. During the i -th session, the model is trained with the pool $T_i = \{X^{(i-1)s+1}, \dots, X^{is}\}$ with s being the incremental step. For practical reasons and without loss of generality, we set s fixed during the whole training process.

In our approach, we also have access to another source of data \mathcal{U} providing *unlabeled* images belonging to the same domain as \mathcal{X} . We make the hypothesis that \mathcal{U} is available to the process with no restriction at any time, although in practice only a limited quantity of data can be exploited during each learning session.

At each session i , the learning process uses a deep neural network with parameters Θ_i capable of predicting the class probability for any element $y \in \mathcal{Y}_i = \{y_1, \dots, y_{i \times s}\}$ and any input sample $x \in \mathbb{R}^n$: $p(y = j|x) = f_j^i(x; \Theta_i)$. The DNN consists in a convolutional part with parameters θ_i , which can be seen as a feature extractor $\phi(x, \theta_i) : \mathbb{R}^n \rightarrow \mathbb{R}^d$, followed by a fully connected layer classifier with parameters $\mathbf{w}_i \in \mathbb{R}^{d \times (i \times s)}$. At each session, the model is initialized with the previous set of parameters $\Theta_{i-1} = (\theta_{i-1}, \mathbf{w}_{i-1})$. s outputs are added to the single-head classifier while the encoder keeps the same parameter structure in θ_i .

3.2 Buffer management and training data

The training data at each session comes from four different sources: the pool of annotated data T_i , the unlabeled data \mathcal{U} , the rehearsal buffer \mathcal{B} and a data augmentation process.

Memory Buffer \mathcal{B} Its role is to store old annotated samples to mimic an episodic memory. It is characterized by a hyperparameter K that defines the number of stored samples. During the learning process, we use random selection to pick the exemplars while ensuring the balance between classes, i.e. after the i -th session, \mathcal{B} contains $\lfloor \frac{K}{|\mathcal{Y}_s} \rfloor$ exemplars per class.

Data Augmentation We will see in the next section that our approach makes use of the generation of pseudo-labels for non-annotated images. This mechanism is based on the idea that images must keep the same pseudo-labels even when a transformation is applied to them.

We define two types of possible transformations: strong and weak, denoted respectively as $\mathcal{A}(\cdot)$ and $\alpha(\cdot)$. In practice, our weak transformations consist of random horizontal flips and translations, as it is practiced in most representation learning methods [63, 44]. Our strong transformations strictly follow the implementation of [68]: they include transformations such as cutout, translation, rotation, color and brightness adjustment, etc. and use the CTAugment sampling strategy described in [9]. The complete list of augmentations is provided in appendix D along with an ablation evaluating the combination of both weak and strong augmentations for consistency regularization.

3.3 Learning process

The parameters Θ_i are learned by stochastic gradient descent (SGD) at each session i . As is standard when applying SGD, each parameter update step makes use of a mini-batch of data which is randomly sampled at each step. In the semi-supervised scheme proposed in our approach, such a mini-batch \mathcal{S} is made of two types of data: a subset \mathcal{S}_l containing B labeled images sampled from $T_i \cup \mathcal{B}$ and another subset \mathcal{S}_u composed of μB images from \mathcal{U} where μ is a scalar hyperparameter.

The proposed PLCiL algorithm relies on the optimization of 3 combined losses targeting 3 different objectives: i) supervision coming from the novel labeled data, ii) consistency regularization using pseudo-labels automatically generated on the unlabeled data, iii) self-supervised knowledge distillation adding extra regularization. We describe these 3 losses in the following, removing the session number i in the notations when its reference is not necessary. The overall training process is illustrated in Figure 1.

Supervised loss l_{sup} . Its role is to use labeled data to learn the model. To improve robustness, a data augmentation step is introduced using a weak transformation α . The supervised signal is back propagated to the network using the standard cross-entropy loss between the output of the DNN and the true labels. It can be expressed as:

$$l_{\text{sup}} = \frac{1}{B} \sum_{(x,y) \in \mathcal{S}_l} H(y, f(\alpha(x); \Theta)) \quad (1)$$

where $f(x; \Theta) \in \mathbb{R}^{|\mathcal{Y}|}$ is the predicted class distribution given an input x , $|\mathcal{Y}|$ is the number of classes considered in the current session. H is the cross-entropy defined by:

$$H(y, f(\alpha(x); \Theta)) = - \sum_{j \in \mathcal{Y}} y_j \log(f_j(\alpha(x); \Theta)) \quad (2)$$

where y is a one-hot encoding of the true class label.

Self-Supervised loss l_{self} . Its role is to regularize the image representations by mimicking true annotation using pseudo-labels on unlabeled data u from \mathcal{S}_u and two levels of image transformations: weak and strong. A weakly transformed data $\alpha(u)$ is first fed to the DNN. If the model is confident enough on its output (according to a threshold τ on the scores), this

prediction is used as a pseudo-label for a cross-entropy loss on the strongly augmented image $\mathcal{A}(u)$. Given the prediction on weakly augmented data $q_u = f(\alpha(u), \Theta)$ and the pseudo-label $\hat{q}_u = \operatorname{argmax}(q_u)$, the resulting self-supervised loss is:

$$l_{\text{self}} = \frac{1}{\mu B} \sum_{u \in \mathcal{S}_u} \mathbb{1}_{\max(q_u) > \tau} H(\hat{q}_u, f(\mathcal{A}(u); \Theta_i)) \quad (3)$$

Distillation loss l_{kd} . Distillation is used to ensure prediction consistency between sessions and is thus expected to lower forgetting. We again use pseudo-labeling with confidence thresholding but with the difference that the pseudo-labels are generated using the model from the previous session $f(x; \Theta_{i-1})$. Let $q_{\text{old}} = f(\alpha(u), \Theta_{i-1})$ and $\hat{q}_{\text{old}} = \operatorname{argmax}(q_{\text{old}})$ be respectively the prediction and the associated pseudo-label at the previous session. The knowledge distillation loss is defined as:

$$l_{\text{kd}} = \frac{1}{\mu B} \sum_{u \in \mathcal{S}_u} \mathbb{1}_{\max(q_{\text{old}}) > \tau} H(\hat{q}_{\text{old}}, f(\mathcal{A}(u); \Theta_i)) \quad (4)$$

Note that in our approach, l_{kd} is computed only on unlabeled samples.

Total loss. PLCiL combines the 3 training objectives during the optimization process:

$$loss = l_{\text{sup}} + \lambda(l_{\text{self}} + \eta l_{\text{kd}}) \quad (5)$$

with $\eta = \frac{|\mathcal{Y}_{i-1}|}{|\mathcal{Y}_i|}$ the ratio between the number of classes learned by the old model and the current number of classes. This scalar is used in [26, 41, 46] to balance the distillation loss. Note that with the assumption of s constant, $\eta = \frac{i-1}{i}$. λ is a scalar hyper-parameter balancing supervision and self-supervision.

3.4 Discussion

The design of our algorithm was guided by 3 key objectives: i) the use of additional unlabeled data to improve and make the learned representations more stable, due to the visual diversity they provide; ii) the use of unlabeled data to add self-regularization to the classification head; iii) the use of pseudo-labels, generated by the model of the previous session to distill knowledge between incremental steps, adding additional regularization.

One of our main contributions is therefore knowledge distillation via unlabeled data. This is related to consistency regularization by matching the distribution of outputs of the two models (the old and the new one), while in related work KD is usually based on soft sharpening [26]. Pseudo-labeling KD is very specific to our approach as we use distillation on unlabeled samples: a confidence threshold allows to select the nature of the distilled knowledge, retaining only the relevant examples, while classical distillation [25, 41, 45, 46] blindly transfers a fraction (temperature parameter) of the whole knowledge contained in all examples. This, together with the fact that the classification head grows with each training session, makes our problem very different from that of FixMatch [33].

The argmax-based distillation loss l_{kd} , defined in Eq. (4), takes advantage of both the unlabeled mini-batch \mathcal{S}_u and the old stored model $f(x; \Theta_{i-1})$. In addition to the self-supervision loss (l_{self}) which regularizes the model as such, l_{kd} enforces the consistency between sessions (i.e., between Θ_{i-1} and Θ_i).

During the learning process, the model of the previous session is expert for the classes already seen. During a new session it can produce pseudo-labels for unlabeled images. In

practice, even if these images come from categories never seen, they have links with the categories seen because of the selection process (application of transformation and score thresholding), thus creating bridges between sessions. This is very different from pseudo-labeling for semi-supervised learning, which concerns already known classes.

This distillation mechanism also makes it possible to properly take into account the increase in the number of outputs of the classification head. The sudden increase in the number of logits overwhelms the output distribution, due to the softmax activation. Thus, during the first training epochs of $f(x; \Theta_i)$, almost none of the predictions on unlabeled images will exceed the τ threshold (usually set close to 1), which means that the pseudo-labeling process is reset at each session. The KD loss ensures that consistent pseudo-labeling is maintained over the sessions. The η parameter takes into account the fact that the pseudo-labels can change over time as the number of possible outputs increases. When the teacher (i.e. the previous model) knows only a few classes, many images are likely to get a wrong pseudo-label, so less weight is given to the labels provided by the previous model in favor of the pseudo-labeling done by the current model. However, a ratio η close to 1 means that the two models know about the same amount of classes. We can assume that, in this case, most of the images can be correctly labeled by both and then give the same credit to l_{kd} and l_{self} .

4 Experiments

We experimentally validated our method on the 2 datasets commonly used to evaluate class incremental methods, namely CIFAR-100 [24] and ImageNet-100 [65]. ImageNet-100 is a subset of ImageNet-1000 where only 100 classes are considered (same classes as [6, 63]).

Evaluation is done using the standard incremental accuracy metric: the model is tested at the end of each session on all the classes seen so far. The Last Accuracy is measured on all classes once the class incremental process is completed, and the Average Accuracy is the mean of all incremental accuracies, excluding the first session which cannot be considered as incremental. For ImageNet experiments, we reported the top-5 accuracy. We averaged 3 runs for CIFAR-100 and 1 for ImageNet, as it is commonly done in the literature [41, 46].

Note that exploiting additional unlabeled data – which is the main objective of the proposed approach – is not directly feasible with existing methods since they are designed to work under full supervision. Only DMC+ [45] is designed to work under the same scenario. Thus, the results, more than a direct raw comparison between the methods, should be seen as a showcase of the advantages of using semi-supervision when training a CIL model.

We nevertheless compare our method to the following (fully supervised) rehearsal-based solutions which are known to work in large scale CIL scenarios: GDumb [60], iCaRL [63], BiC [40], Weight Aligning (WA) [46] and, as mentioned above, with the semi-supervised DMC+ method [45]. To give fully supervised methods a fairer chance, we propose some experiments in which they are pre-trained with the same unlabeled data. We have re-implemented all the competing methods so as to have exactly the same backbone network and the same DL framework for all the methods being compared. Due to encountered difficulties in reproducing Global Distillation (GD) [25] and its costly requirements for unlabeled data sampling, we directly report results from the original paper when possible.

All methods are tested using the same DNN : a Wide-ResNet-28-8 [43] (WRN28-8) for CIFAR and a ResNet-18 [16] for ImageNet. For reference, WRN28-8 achieves 82.8% on CIFAR-100 and ResNet-18 top-5 accuracy on ImageNet-100 is 94.4%. All hyperparameters used are detailed in appendix A.1 with an in depth study of the sensitivity of μ , τ and λ

Method	Last (%)	Avg (%)
GDumb [60]	27.8	42.0
iCaRL [63]	53.9	63.9
BiC [41]	55.9	67.1
WA [46]	50.8	64.4
DMC+ [15]	50.4	62.8
GD [29]	54.9*	68.1*
Ours	61.5	74.0

* Results from [29] with WRN16-2 backbone and Tiny Images (80M images) as unlabeled data.

Table 1: CIL comparison on CIFAR-100-full.

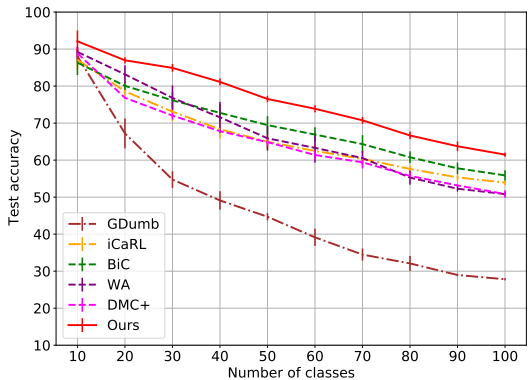


Figure 2: Plot of incremental accuracy over the successive sessions on CIFAR-100-full. Best viewed in PDF.

Method	CIFAR-100-20%				ImageNet-100-10%			
	Last (%)		Avg (%)		Last (%)		Avg (%)	
	Random Init	RotNet Init	Random Init	RotNet Init	Random Init	RotNet Init	Random Init	RotNet Init
GDumb [60]	28.2	42.2	25.3	40.4	40.6	59.6	43.7	62.0
iCaRL [63]	42.7	48.9	43.9	51.1	45.4	57.8	51.9	59.5
BiC [41]	43.3	49.8	43.5	57.3	50.7	62.4	52.2	68.8
WA [46]	40.5	49.7	45.5	55.4	30.2	54.7	40.9	64.5
DMC+ [15]	36.4	42.8	39.3	49.8	56.2	68.1	57.5	69.6
Ours	59.8	66.5	59.5	67.6	61.3	73.8	61.2	75.0

Table 2: Experiments with fewer labels available: CIFAR-100-20% and ImageNet-100-10%.

in appendix C. For CIFAR-100, WRN28-8 contains many more parameters than ResNet-32 commonly used by CIL methods. This choice of backbone is discussed in appendix A.2.

4.1 Class Incremental Results

We conducted experiments following the standard class incremental protocol [6, 63] on CIFAR-100. We set $s = 10$, i.e. 10 sessions of 10 classes. The memory size for rehearsal K is set to 2000. Since we use the whole labeled dataset here, we refer to this experimental setting as CIFAR-100-full. For DMC+ and PLCiL, we emulate an unlabeled data stream by randomly sampling 100K unlabeled images from ImageNet-1000 (subsampled to 32×32).

The results obtained are available in Table 1. They show the value of exploiting unlabeled data in a joint CIL framework. Indeed, our PLCiL method consistently outperforms other state-of-the-art methods on CIFAR-100-full, gaining +5.6% on final accuracy. In the following section, we demonstrate that the semi-supervised nature of our method is able to address even more challenging scenarios where less supervision is available.

Scenario	Last (%)	Avg (%)
a. ImageNet-900	61.3	73.8
b. ImageNet-100	76.9	83.3
c. ImageNet-1000	65.1	76.2
d. Places365	59.6	72.7

Table 3: Class incremental performance on ImageNet-100-10% with 4 different unlabeled datasets.

Loss	Last (%)	Avg (%)
l_{sup}	44.1	59.8
$l_{sup} + \lambda \eta l_{kd}$	61.9	63.9
$l_{sup} + \lambda l_{self}$	50.3	65.15
$l_{sup} + \lambda (l_{self} + \eta l_{kd})$	61.5	74.0
$l_{sup} + \lambda (l_{self} + \eta l_{standkd})$	52.0	65.6
$l_{sup} + \lambda (l_{soft} + \eta l_{standkd})$	50.9	63.4

Table 4: CIL on CIFAR-100-full with only specific components of the loss enabled.

4.2 Semi-Supervised Class Incremental Results

The following experiments study the behavior of our method in scenarios even closer to the non semi-supervised scenarios: only a very limited set of labeled data is available, while a large amount of cheap unlabeled data is accessible.

Following this principle, we keep the previous settings with $s = 10$ but reduce the size of the labeled dataset. For CIFAR-100-20%, we randomly pick 100 samples per class out of the 500, while for ImageNet-100-10%, 130 labeled samples per class are retained. This is in line with the amount of labels commonly used in semi-supervised works [8]. The memory budget remains unchanged with $K = 2000$.

The process of collecting the unlabeled data is the same as before: we sample 100,000 data from ImageNet-1000 at the beginning of each session. To avoid any data leakage, images belonging to classes learned incrementally are removed from the unlabeled data.

Since this setting is very difficult for fully supervised methods due to the scarcity of data, we also compare with a self-supervised initialization using the same amount of unlabeled data (1M images). WRN28-8 and ResNet-18 were trained using RotNet [17] on ImageNet-1000 (with classes to be learned incrementally excluded). Self-supervised pre-training is the most immediate way to allow fully supervised methods to access as much unlabeled data as PLCiL or DMC+.

The results are presented in Table 2. The lack of data is clearly perceptible for supervised approaches and only GDumb is stable due to the fact that its performance only depends on the buffer size. Meanwhile, PLCiL maintains the level of performance obtained during the first protocol, proving that in scenarios with very few labels, our method can efficiently exploit unlabeled data. Although the representations obtained by self-supervision are better, the exploitation of unlabeled data throughout is more efficient with PLCiL, which is outperforming all pretrained fully supervised methods.

PLCiL is consistently better than DMC+, which however exploits the same quantity of data as PLCiL. We believe this is because DMC+ only exploits unlabeled data to distill knowledge from a learned model with full supervision. From the plasticity-stability dilemma, such approach focuses on leveraging unlabeled data to improve stability. The performance on new classes is dependent on the fully supervised training phase of the teacher model. While PLCiL has a similar behavior with l_{kd} , it also aims to improve plasticity by also implementing self-supervised representation learning with l_{self} and its data augmentation strategy. The ablation study gives more details on the contribution of each loss.

4.3 Ablation study

Impact of the unlabeled data source. In the experiments reported in tables 1 and 2, we use an unlabeled dataset which is semantically similar to the target while excluding all samples belonging to the learned classes (as it is common in self-supervised literature). This emulates an application where both the unlabeled data \mathcal{U} and labeled data \mathcal{X} are collected from the same environment.

In order to evaluate the influence of the unlabeled source, we repeat the same experiment on ImageNet-100-10% with random initialization using 4 different unlabeled datasets for \mathcal{U} : (a) samples from the remaining 900 classes as in section 4.2. (b) samples from the remaining images of ImageNet-100 (which is feasible since only 10% of the labeled images are used in \mathcal{X}). (c) samples from all unused images from ImageNet-1000 so that \mathcal{U} contains about one tenth of images belonging to classes in \mathcal{X} using uniform sampling. (d) samples from a semantically unrelated dataset (Places-365 [47]).

Results are presented in Table 3. (c) shows that a reasonable class leakage between \mathcal{U} and \mathcal{X} is efficiently leveraged by our model, improving the performance from our baseline (a). We also tried the idealistic scenario of (b) with the exact same classes in both dataset. This provides an upper bound of our method given a fully related unlabeled dataset. At last, using a semantically unrelated dataset in (d) slightly lowers the performance compared to (a) but still outperforms the other approaches shown in table 2.

Contribution of each loss component We ran several variants of PLCiL on CIFAR-100-full to evaluate the contribution of each loss term (see Table 4). We noticed two complementary behaviors when l_{kd} and l_{self} are used separately. The version with pseudo-labeling KD focuses on the stability of the model, keeping the most consistent accuracy from session to session and achieves the best final accuracy despite a lower average. This is due to the fact that it struggles to learn new classes, especially during the early stages where the proportion of old classes is low, making KD less relevant. The version with only l_{self} enhances the plasticity of the model, allowing to learn the new classes with better accuracy as it is shown in the first sessions. However, this variant still lacks regularization to alleviate the catastrophic forgetting during the later stages. By combining both, PLCiL optimizes the *plasticity-stability* trade-off and gives satisfactory results during all sessions. The full version of Table 4 with accuracy at all sessions is provided in Appendix B.1.

Efficiency of Pseudo-Labeling for Knowledge Distillation. We replaced our l_{kd} by a standard KD, as in [41, 46], but still applied on unlabeled data only. The results are given in the penultimate line of Table 4. Standard KD ($l_{standkd}$) has little or no effect and gives results similar to BiC and WA. This comparison highlights the efficiency of our custom distillation with only unlabeled data and pseudo-labeling.

We also experimented with the removing of the Pseudo-Labeling step. As for KD, the *soft* alternative consists in using directly the output distribution as target. Using a similar loss to the standard distillation l_{soft} , we apply consistency regularization between weakly and strongly augmented images. In the last line of Table 4 with both l_{soft} and $l_{standkd}$, i.e. with the thresholding removed in all losses, the performance further drop. This result validates the usefulness of a hard thresholding for leveraging unlabeled data.

5 Conclusion

We introduce a simple class incremental method combining learning by rehearsal and self-supervised learning that takes advantage of complementary unlabeled data during the learning process. With this simple approach, we demonstrate that semi-supervision is a valuable aid to the problems of catastrophic forgetting and scarcity of data at training time. Experimental validation on both CIFAR-100 and ImageNet-100 shows that PLCiL has better learning capacities than existing methods, thanks to finer and more stable representations.

Acknowledgments. Frederic Jurie is partly supported by ANR grant #ANR-19-CHIA-0017.

References

- [1] Rahaf Aljundi, Francesca Babiloni, Mohamed Elhoseiny, Marcus Rohrbach, and Tinne Tuytelaars. Memory Aware Synapses: Learning what (not) to forget. In *European Conference on Computer Vision, ECCV*, 2018.
- [2] Yoshua Bengio, Aaron C. Courville, and Pascal Vincent. Representation Learning: A Review and New Perspectives. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(8):1798–1828, 2013.
- [3] David Berthelot, Nicholas Carlini, Ekin D. Cubuk, Alex Kurakin, Kihyuk Sohn, Han Zhang, and Colin Raffel. ReMixMatch: Semi-Supervised Learning with Distribution Matching and Augmentation Anchoring. In *8th International Conference on Learning Representations, ICLR*, 2020.
- [4] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised Learning of Visual Features by Contrasting Cluster Assignments. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems, NeurIPS*, 2020.
- [5] Francisco M. Castro, Manuel J. Marín-Jiménez, Nicolás Guil, Cordelia Schmid, and Karteek Alahari. End-to-End Incremental Learning. In *European Conference on Computer Vision, ECCV*, 2018.
- [6] Olivier Chapelle, Bernhard Schölkopf, and Alexander Zien. *Semi-Supervised Learning*. The MIT Press, 2006. ISBN 9780262033589.
- [7] Arslan Chaudhry, Puneet Kumar Dokania, Thalaiyasingam Ajanthan, and Philip H. S. Torr. Riemannian Walk for Incremental Learning: Understanding Forgetting and Interscience. In *European Conference on Computer Vision, ECCV*, 2018.
- [8] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey E. Hinton. A Simple Framework for Contrastive Learning of Visual Representations. In *37th International Conference on Machine Learning, ICML*, 2020.
- [9] Ekin D. Cubuk, Barret Zoph, Dandelion Mané, Vijay Vasudevan, and Quoc V. Le. AutoAugment: Learning Augmentation Strategies From Data. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, pages 113–123, 2019.

- [10] Arthur Douillard, Matthieu Cord, Charles Ollion, Thomas Robert, and Eduardo Valle. PODNet: Pooled Outputs Distillation for Small-Tasks Incremental Learning. In *European Conference on Computer Vision, ECCV*, 2020.
- [11] Robert M. French. Catastrophic forgetting in connectionist networks. *Trends in Cognitive Sciences*, 3(4):128 – 135, 1999. ISSN 1364-6613.
- [12] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised Representation Learning by Predicting Image Rotations. In *6th International Conference on Learning Representations, ICLR*, 2018.
- [13] Ian J. Goodfellow, Mehdi Mirza, Da Xiao, Aaron Courville, and Yoshua Bengio. An Empirical Investigation of Catastrophic Forgetting in Gradient-Based Neural Networks. *arXiv:1312.6211*, 2013.
- [14] Tyler L. Hayes, Kushal Kafle, Robik Shrestha, Manoj Acharya, and Christopher Kanan. REMIND Your Neural Network to Prevent Catastrophic Forgetting. In *European Conference on Computer Vision, ECCV*, 2020.
- [15] Jiangpeng He, Runyu Mao, Zeman Shao, and Fengqing Zhu. Incremental Learning in Online Scenario. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, 2020.
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 770–778, 2016.
- [17] Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the Knowledge in a Neural Network. *arXiv:1503.02531*, 2015.
- [18] Saihui Hou, Xinyu Pan, Chen Change Loy, Zilei Wang, and Dahua Lin. Learning a Unified Classifier Incrementally via Rebalancing. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2019.
- [19] Yen-Chang Hsu, Yen-Cheng Liu, and Zsolt Kira. Re-evaluating Continual Learning Scenarios: A Categorization and Case for Strong Baselines. *arXiv:1810.12488*, 2018.
- [20] Ahmet Iscen, Jeffrey Zhang, Svetlana Lazebnik, and Cordelia Schmid. Memory-Efficient Incremental Learning Through Feature Adaptation. In *European Conference on Computer Vision, ECCV*, 2020.
- [21] Khurram Javed and Martha White. Meta-Learning Representations for Continual Learning. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems, NeurIPS*, 2019.
- [22] Longlong Jing and Yingli Tian. Self-supervised Visual Feature Learning with Deep Neural Networks: A Survey. *arXiv:1902.06162*, 2019.
- [23] 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. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.

- [24] Alex Krizhevsky. Learning Multiple Layers of Features from Tiny Images. Technical report, 2009.
- [25] Kibok Lee, Kimin Lee, Jinwoo Shin, and Honglak Lee. Overcoming Catastrophic Forgetting With Unlabeled Data in the Wild. In *International Conference on Computer Vision, ICCV*, 2019.
- [26] Zhizhong Li and Derek Hoiem. Learning without Forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947, 2017.
- [27] Yaoyao Liu, Yuting Su, An-An Liu, Bernt Schiele, and Qianru Sun. Mnemonics Training: Multi-Class Incremental Learning Without Forgetting. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, 2020.
- [28] Ilya Loshchilov and Frank Hutter. SGDR: stochastic gradient descent with warm restarts. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, 2017.
- [29] German I Parisi, Ronald Kemker, Jose L Part, Christopher Kanan, and Stefan Wermter. Continual Lifelong Learning with Neural Networks: A Review. *Neural Networks*, 113: 54–71, 2019.
- [30] Ameya Prabhu, Philip Torr, and Puneet Dokania. GDumb: A Simple Approach that Questions Our Progress in Continual Learning. In *The European Conference on Computer Vision (ECCV)*, 2020.
- [31] Jathushan Rajasegaran, Munawar Hayat, Salman H. Khan, Fahad Shahbaz Khan, and Ling Shao. Random Path Selection for Continual Learning. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems, NeurIPS*, 2019.
- [32] Jathushan Rajasegaran, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, and Mubarak Shah. iTAML: An Incremental Task-Agnostic Meta-learning Approach. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, 2020.
- [33] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. iCaRL: Incremental Classifier and Representation Learning. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [34] Matthew Riemer, Ignacio Cases, Robert Ajemian, Miao Liu, Irina Rish, Yuhai Tu, and Gerald Tesauero. Learning to Learn without Forgetting by Maximizing Transfer and Minimizing Interference. In *7th International Conference on Learning Representations, ICLR*, 2019.
- [35] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015.
- [36] 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 *35th International Conference on Machine Learning, ICML*, 2018.

- [37] Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual Learning with Deep Generative Replay. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems, NeurIPS*, 2017.
- [38] Kihyuk Sohn, David Berthelot, Chun-Liang Li, Zizhao Zhang, Nicholas Carlini, Ekin D. Cubuk, Alex Kurakin, Han Zhang, and Colin Raffel. FixMatch: Simplifying Semi-Supervised Learning with Consistency and Confidence. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS*, 2020.
- [39] Guido M. van de Ven and Andreas S. Tolias. Three scenarios for continual learning. *arXiv:1904.07734*, 2019.
- [40] Yue Wu, Yinpeng Chen, Lijuan Wang, Yuancheng Ye, Zicheng Liu, Yandong Guo, Zhengyou Zhang, and Yun Fu. Incremental Classifier Learning with Generative Adversarial Networks. *arXiv:1802.00853*, 2018.
- [41] Yue Wu, Yinpeng Chen, Lijuan Wang, Yuancheng Ye, Zicheng Liu, Yandong Guo, and Yun Fu. Large Scale Incremental Learning. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2019.
- [42] Shipeng Yan, Jiangwei Xie, and Xuming He. DER: Dynamically Expandable Representation for Class Incremental Learning. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2021.
- [43] Sergey Zagoruyko and Nikos Komodakis. Wide Residual Networks. In *Proceedings of the British Machine Vision Conference 2016, BMVC 2016, York, UK, September 19-22, 2016*, 2016.
- [44] Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual Learning Through Synaptic Intelligence. In *34th International Conference on Machine Learning, ICML*, 2017.
- [45] Junting Zhang, Jie Zhang, Shalini Ghosh, Dawei Li, Serafettin Tasci, Larry P. Heck, Heming Zhang, and C.-C. Jay Kuo. Class-incremental Learning via Deep Model Consolidation. In *IEEE Winter Conference on Applications of Computer Vision, WACV*, 2020.
- [46] Bowen Zhao, Xi Xiao, Guojun Gan, Bin Zhang, and Shu-Tao Xia. Maintaining Discrimination and Fairness in Class Incremental Learning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, 2020.
- [47] Bolei Zhou, Àgata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 Million Image Database for Scene Recognition. *IEEE Trans. Pattern Anal. Mach. Intell.*, 40(6):1452–1464, 2018.
- [48] Xiaojin Zhu and Andrew B. Goldberg. *Introduction to Semi-Supervised Learning*. Synthesis Lectures on Artificial Intelligence and Machine Learning. Morgan & Claypool Publishers, 2009.

A Implementation details

A.1 Training details

In this section, we provide all details needed for training the PLCiL. All parameters were selected through cross-validation.

On CIFAR-100, each session has 150 epochs. The labeled mini-batch size B is 32 with $\mu = 7$. The confidence threshold τ is set to 0.8 and λ to 1. On ImageNet-100, training takes 70 epochs per session and uses the following set of parameters: $\{B = 32, \mu = 7, \tau = 0.7, \lambda = 3\}$.

All experiments are trained with SGD. The learning rate is initialized to 0.03 with Nesterov momentum set to 0.9. We use a weight decay of 10^{-4} , a learning rate decay with cosine annealing and warm restart [28] with an initial period $T_0 = 10$ multiplied by a factor $T_{mult} = 2$ after each restart.

The sensitivity of the main hyperparameters μ , τ and λ are detailed in the following Appendix C. In particular, μ controls the ratio unlabeled/labeled data of each batch. Our experiments have shown that a larger μ is generally better at the expense of increased computational costs. We kept μ quite low so our PLCiL could run on limited hardware with computation time similar to the other methods. τ and λ control the pseudo-labeling process and were tuned for each dataset.

A.2 Continual Training of large DNNs

Our experimentation on CIFAR-100 uses a WRN-28-8 backbone architecture (23M trainable parameters), as [33]. Incidentally, self-supervision – or semi-supervision – is used to train large DNNs when the application is data-scarce. Those models, and even larger ones with hundreds of millions of weights, are now prevailing on classification benchmarks.

Despite that, most CL methods are only evaluated on smaller DNNs: *e.g.* ResNet-32 for CIFAR with only 460K parameters, far from state-of-the-art accuracy in batch training. This questions the quality of the representation that such small models can learn and therefore, their plasticity potential. For reference, the batch accuracy is 80.6% for WRN28-8 and 72.3% for ResNet-32.

We have conducted more experiments on CIFAR-100-full, similar to those of Section 4.1 but with a smaller ResNet-32, to observe the influence of the size of the network. We present performance comparison between these two backbones in Table 5.

With the ResNet-32, our PLCiL falls behind the compared state-of-the-art methods that were designed with this particular backbone. However, our purpose is to see how the methods scale to larger models. Apart from WA, all the methods become more accurate with larger models. The gap stands out with our method with an overall +12.4% obtained when using a larger architecture. Using WRN28-8 instead of ResNet-32 seems less profitable to other methods, and even damages the accuracy of WA. This can be justified by the limited amount of data they can use at each session, making it harder to learn a large number of weights. Our PLCiL addresses this issue and makes larger architectures trainable thanks to the extensive visual diversity submitted to the model through the self-supervised signal.

Method	ResNet-32		WRN28-8	
	Last (%)	Avg (%)	Last (%)	Avg (%)
GDumb [60]	20.7	35.0	27.8	42.0
iCaRL [63]	47.1	57.9	53.9	63.9
BiC [41]	50.8	62.7	55.9	67.1
WA [46]	52.1	65.6	50.8	64.4
DMC+ [45]	43.9	58.4	50.4	62.8
Ours	46.9	62.1	61.5	74.0

Table 5: Comparison between all tested approaches on CIFAR-100-full with two backbones: ResNet-32 (460K parameters) and WRN28-8 (24M parameters).

loss	10	20	30	40	50	60	70	80	90	100	Avg acc
l_{sup}	91.7	82.2	74.2	66.0	62.9	58.3	53.4	50.9	46.4	44.1	59.8
$l_{sup} + \lambda \eta l_{kd}$	91.7	47.6	59.1	68.9	71.5	68.7	67.6	66.1	63.5	61.9	63.9
$l_{sup} + \lambda l_{self}$	91.7	85.6	78.7	72.1	67.7	63.6	59.5	56.2	52.5	50.3	65.15
$l_{sup} + \lambda (l_{self} + \eta l_{kd})$	91.7	86.9	84.9	81.1	76.5	73.9	70.8	66.7	63.8	61.5	74.0
$l_{sup} + \lambda (l_{self} + \eta l_{standardkd})$	91.7	84.2	77.8	72.2	66.6	64.0	62.8	57.6	52.9	52.0	65.6

Table 6: CIL on CIFAR-100-full with only specific components of the loss enabled. Accuracy (%) are computed at the end of each session on all the classes learned so far. Average accuracy does not take into account the first session.

B Complete ablation results

B.1 Contribution of each loss component

In Table 6, we present the full version of Table 2 from our ablation study (Sec. 4.3). With the accuracy at each session, we can clearly distinguish the difference between l_{kd} and l_{self} . The former focuses on stability with a consistent accuracy across all sessions while the later enhance the learning of new classes during the early steps but is still very prone to forgetting.

These results corroborate the importance of the weight η in the complete loss as discussed in section 3.4. When the proportion of new classes is still high compared to the number of classes already learned (η low), the performance is more dependent on the ability to learn new things, KD should then have a low impact on the training. However, during the later sessions, when the number of classes to retain is very high compared to the novelty (η close to 1), KD becomes crucial against the catastrophic forgetting.

C Hyperparameters sensitivity

C.1 Ratio labeled-unlabeled data: μ

The hyperparameter μ defines the amount of unlabeled data sampled in each mini-batch, i.e. for each labeled mini-batch size B , our algorithm considers μB unlabeled data. Thus, increasing μ directly increase the visual diversity seen by the model in a self-supervised fashion. However, this also means larger mini-batches for training which can be timely

μ	0	1	2	3	7	15	31
Avg acc	59.8	69.2	71.2	72.1	74.0	74.4	75.0
Last acc	44.1	55.5	58.9	59.5	61.5	62.2	63.6

Table 7: Comparison of Last Accuracy and Average Accuracy for different values of μ on CIFAR-100-full.

and computationally expensive. In table 7, with B set to 32, we see a consistent increase in performance with larger unlabeled mini-batches. We chose to keep $\mu = 7$ for all our experimentation since it gives satisfactory results and keep the training time and hardware requirement comparable to others CI methods (e.g. $B + \mu B$ gives a total mini-batch size of 256 with $\mu = 7$ and $B = 32$). Higher values of μ only yield minor improvements despite being way more costly.

C.2 Selectivity of the threshold τ and weight of the pseudo-labeling λ

In table 8, we display the results for several combination of the hyperparameters τ and λ . The trends indicate that λ should be kept close to 1, meaning that giving too much weight to the self-supervised part of the loss has a negative impact on the learning of classes.

The PLCiL is less sensitive to the threshold value. For $\lambda \leq 2$, our model reach at least 72.6% average accuracy for all τ tested here. This is due to the fact that our model answer confidently for the majority of the unlabeled data seen, outputting high values that goes beyond most threshold values. This is probably due to the curated nature of our unlabeled data pool (ImageNet) which contains visual information easily transferable to CIFAR (close domains). We did not happen to experiment on it, but we believe that τ could be crucial to filter noisy information when dealing with non-curated unlabeled data or if the domains between the labeled and the unlabeled data were further apart.

D Data Augmentation

Our approach makes use of data augmentation strategies to leverage the information provided by the unlabeled data. In this section, we study the impact of the two types of augmentation used: weak and strong, and demonstrate that their combination is useful to improve the performance, especially for large scale datasets and scarce annotations.

D.1 List of Augmentations

In this paper, we used the combination of two sets of augmentation: weak \mathcal{A} and strong \mathcal{A} .

Weak augmentations consist in random vertical and horizontal translations followed by an horizontal flip occurring with a probability of 0.5. A resizing is applied when needed in order to fit the input requirement of the model.

Strong augmentations are applied according to the CTAugment [9] algorithm. It samples two transformations from the following list of 18: autocontrast, brightness adjustment, color adjustment, contrast adjustment, cutout, histogram equalization, pixel inversion, identity, posterizing, rescaling, rotation, sharpness adjustment, horizontal shear, vertical shear,

τ	λ	Avg acc (%)	Last acc (%)
0.5	1	73.6	61.9
0.5	2	73.5	61.0
0.5	5	71.0	60.9
0.5	7	69.4	58.5
0.5	10	62.2	56.9
0.7	1	73.6	61.7
0.7	2	73.6	61.5
0.7	5	69.6	62.3
0.7	7	69.1	60.8
0.7	10	67.3	57.5
0.8	1	74.7	62.3
0.8	2	74.4	62.4
0.8	5	73.4	62.3
0.8	7	72.7	59.3
0.8	10	67.3	59.3
0.9	1	72.6	59.3
0.9	2	73.6	61.6
0.9	5	73.4	58.5
0.9	7	72.3	59.4
0.9	10	68.5	50.9

Table 8: Evaluation of the PLCiL for different combinations of τ and λ . For this experiment, results are reported on only one permutation of classes instead of the usual 3 runs.

Augmentation	Last (%)	Avg (%)
weak + CTAugment	61.5	74.0
weak + weak	2.8	8.4
CTAugment + CTAugment	45.8	62.3
RandAugment + CTAugment	53.8	68.2

Table 9: Class incremental performance for several data argumentation strategies on incremental CIFAR-100-full.

smoothing, solarizing, horizontal translation, vertical translation. Each selected transformation is applied with a magnitude sampled from a learned range.

D.2 Augmentation Strategies

Pseudo-Labeling allows to regularize the output consistency of the model when using weak and strong augmentation of the same image. This choice of using both weak and strong augmentation is similar to [38] and motivated by the ablation study shown in [9, 33].

In this study, we try different combinations of augmentation on the incremental CIFAR-100 benchmark. We compare the following settings: the standard weak + strong CTAugment, only weak augmentation, only strong augmentation using CTAugment and strong augmentation using RandAugment for Pseudo-Labeling and CTAugment for prediction. Results are presented in Table 9.

When using only weak augmentation, at each session, the model quickly reaches an accuracy close to 100% on the training data but shows very bad results on the validation set. This behavior suggests an overfitting situation. Strong augmentation with either CTAugment or RandAugment also yields lower performance compared to the original setting. In [38], the authors mentioned the fact that their model trained with only strong augmentations did not converge. Our small ablation study corroborates the findings of [9, 33]: mixing two families of augmentation is essential for consistency regularization and generating labels from weakly augmented data is more consistent for predicting pseudo-labels.

D.3 Impact of Strong Augmentation on Concurrent Methods

In the main experiments, we evaluate each method given the optimization and parameters provided in the original papers or code repository. Each of these methods preprocess their data using random translation and random horizontal flip which is the same as the weak augment used in our approach.

Since strong augmentation is a core component of the PLCiL, we also evaluate here if other methods can benefit from a wider variety of image transformations. In table 10, we report the performance of each method on the 3 benchmarks of our study using strong augmentation.

The results show that strong augmentation has mitigated results on the other methods. GDumb and WA are significantly improved on all 3 datasets. BiC and iCaRL on the other hand have some inconsistencies which can lead to a drop in accuracy on some benchmarks. For instance, iCaRL reaches top performance on CIFAR-100-full with strong augmentation, even outclassing our PLCiL in final accuracy, but has the opposite behavior on ImageNet-100-10% where the performance is strongly lowered. DMC+ is the most negatively impacted

Method	Weak Augment		Strong Augment	
	Last (%)	Avg (%)	Last (%)	Avg (%)
CIFAR-100-full				
GDumb [50]	27.8	42.0	33.9	48.4
iCaRL [53]	53.9	63.9	62.3	69.7
BiC [41]	55.9	67.1	46.9	68.9
WA [46]	50.8	64.4	51.3	65.8
DMC+ [45]	50.4	62.8	No Convergence	
Ours			61.5	74.0
CIFAR-100-20%				
GDumb [50]	28.2	42.2	34.4	50.4
iCaRL [53]	42.7	48.9	47.2	53.9
BiC [41]	43.3	49.8	46.9	60.8
WA [46]	40.5	49.7	52.6	62.1
DMC+ [45]	36.4	42.8	No Convergence	
Ours			59.8	66.5
ImageNet-100-10%				
GDumb [50]	40.6	59.6	49.1	67.0
iCaRL [53]	45.4	57.8	40.1	52.2
BiC [41]	50.7	62.4	38.6	44.9
WA [46]	30.2	54.7	40.0	50.9
DMC+ [45]	56.2	68.1	44.8	63.0
Ours			61.3	73.8

Table 10: Effect of adding strong augmentation in concurrent baselines. The *Weak Augment* column reports the performance from the original implementations and are the results presented in the section 4.1 and 4.2 of our paper. *Strong Augment* lists all new results obtained when applying strong augmentations. We also report the performance of our PLCiL, which uses both, as a reference.

method with an important degradation of its performance. The model did not converge for both CIFAR-100 experiments.

These experiments suggest that more data augmentation is a simple solution to increase the incremental performance but does not fit some CI methods which have been mostly designed with only weak augmentation. The nature of the dataset is also an important factor. Our PLCiL on the other hand seems to behave consistently with strongly augmented data across the 3 benchmarks evaluated here.