# Self-Supervised Models are Continual Learners

Enrico Fini[*1,2]    Victor G. Turrisi da Costa[*1]    Xavier Alameda-Pineda[2]

Elisa Ricci[1,3]    Karteek Alahari[2]    Julien Mairal[2]

[1] University of Trento    [2] Inria[†]    [3] Fondazione Bruno Kessler

## Abstract

*Self-supervised models have been shown to produce comparable or better visual representations than their supervised counterparts when trained offline on unlabeled data at scale. However, their efficacy is catastrophically reduced in a Continual Learning (CL) scenario where data is presented to the model sequentially. In this paper, we show that self-supervised loss functions can be seamlessly converted into distillation mechanisms for CL by adding a predictor network that maps the current state of the representations to their past state. This enables us to devise a framework for Continual self-supervised visual representation Learning that (i) significantly improves the quality of the learned representations, (ii) is compatible with several state-of-the-art self-supervised objectives, and (iii) needs little to no hyperparameter tuning. We demonstrate the effectiveness of our approach empirically by training six popular self-supervised models in various CL settings. Code: `github.com/DonkeyShot21/cassle`.*

## 1. Introduction

During the last few years, self-supervised learning (SSL) has become the most popular paradigm for unsupervised visual representation learning [3, 7, 8, 13, 14, 26, 28, 58]. Indeed, under certain assumptions (*e.g.*, offline training with large amounts of data and resources), SSL methods are able to extract representations that match the quality of representations obtained with supervised learning, without requiring annotations. However, these assumptions do not always hold in real-world scenarios, *e.g.*, when new unlabeled data are made available progressively over time. In fact, in order to integrate new knowledge into the model, training needs to be repeated on the whole dataset, which is impractical, expensive, and sometimes even impossible when old data is not available. This issue is exacerbated by the fact that SSL models are notoriously computationally expensive to train.
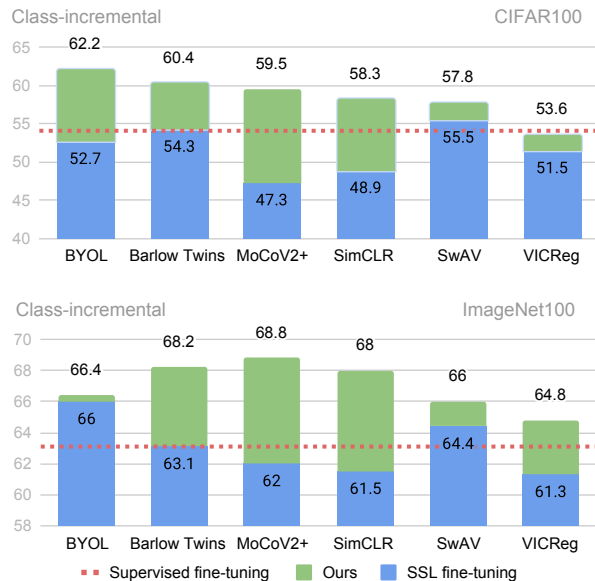


Figure 1. Linear evaluation accuracy of representations learned with different self-supervised methods on class-incremental CIFAR100 and ImageNet100. In blue the accuracy of SSL fine-tuning, in green the improvement brought by CaSSLe. The red dashed line is the accuracy attained by supervised fine-tuning.

Continual learning (CL) studies the ability of neural networks to learn tasks sequentially. Prior art in the field focuses on mitigating catastrophic forgetting [17, 23, 25, 40]. Common benchmarks in the CL literature evaluate the discriminative performance of classifiers learned *with supervision* from non-stationary distributions. In this paper, we tackle the same forgetting phenomenon in the context of SSL. Unsupervised representation learning is indeed appealing for sequential learning since it does not require human annotations, which are particularly hard to obtain when new data is generated on-the-fly. This setup, called Continual Self-Supervised Learning (CSSL), is surprisingly under-investigated in the literature.

In this work, we propose CaSSLe, a simple and effective framework for CSSL of visual representations based on the intuition that SSL models are intrinsically capable of learning continually, and that SSL losses can be seamlessly converted into distillation losses. Our key idea is to train the

---

current model to predict past representations with a prediction head, thus encouraging it to remember past knowledge. CaSSLe has several favourable features: (i) it is compatible with popular state-of-the-art SSL loss functions and architectures, (ii) it is simple to implement, and (iii) it does not require any additional hyperparameter tuning with respect to the original SSL method. Our experiments demonstrate that SSL methods trained continually with CaSSLe significantly outperform all the related methods (CSSL baselines and several methods adapted from supervised CL).

We also perform a comprehensive analysis of the behavior of six popular SSL methods in diverse CL settings (*i.e.*, class, data, and domain incremental). We provide empirical results on small (CIFAR100), medium (ImageNet100), and large (DomainNet) scale datasets. Our study sheds new light on interesting properties of SSL methods that emerge when learning continually. Among other findings, we discover that, in the class-incremental setting, SSL methods typically approach or outperform supervised learning (see Fig.1), while this is not generally true for other settings (data-incremental and domain-incremental) where supervised learning still shows a sizeable advantage.

## 2. Related Work

**Self-Supervised Learning.** Recent SSL approaches have shown performance comparable to their supervised learning equivalents [3,7,8,13,14,26,28,58]. In a nutshell, most of these methods use image augmentation techniques to generate correlated views (positives) from a sample, and then learn a model that is invariant to these augmentations by enforcing the network to output similar representations for the positives. Initially, contrastive learning, based on instance discrimination [56] using noise-contrastive estimation [27,41], was a popular strategy [13,28]. However, this learning paradigm requires large batch sizes or memory banks. A few methods that use a negative-free cosine similarity loss [15,26] have addressed such issues.

Concurrently, clustering-based methods (SwAV [7], DeepCluster v2 [6,7] and DINO [8]) have also been proposed. They do not operate on the features directly, and instead compare positives through a cross-entropy loss using cluster prototypes as a proxy. Redundancy reduction-based methods have also been popular [3,20,58]. Among them, BarlowTwins [58] considers an objective function measuring the cross-correlation matrix between the features, and VicReg [3] uses a mix of variance, invariance and covariance regularizations. Methods such as [19] have explored the use of nearest-neighbour retrieval and divide and conquer [53]. However, none of these works studied the ability of SSL methods to learn continually and adaptively.

**Continual Learning.** A plethora of methods have been developed to counteract catastrophic forgetting [2,4,9–12, 18,22,31,34,36,38,42,44,46–50,55,59]. Following [17],

these works can be organized into three macro-categories: replay-based [4,12,38,42,46,47], regularization-based [2, 9–11,18,22,31,34,36,50,55,59], and parameter isolation methods [48,49]. All these works evaluate the effectiveness of CL methods using a linear classifier learned sequentially over time. However, this evaluation does not reflect an important aspect, *i.e.*, the internal dynamics of the hidden representations. Moreover, most CL methods tend to rely on supervision in order to mitigate catastrophic forgetting. A few of them can be adapted for the unsupervised setting, although their effectiveness is greatly reduced (see discussion in Sec. 5, Sec. 6 and the supplementary material).

Works such as [1,45,51] laid the foundations of unsupervised CL, but their studies are severely limited to digit-like datasets, *e.g.*, MNIST and Omniglot, and the proposed methods are unfit for large-scale scenarios. Recently, [5,24] explored self-supervised pretraining for supervised continual learning with online and few-shot tasks, and [10] presented a supervised contrastive CL approach. Two concurrent works [37,39] have also attempted to address CSSL recently. The former [37] extends [10] to the unsupervised setting, but is specifically designed for contrastive SSL, such as [13,28], and lacks generalizability to other popular SSL paradigms. The latter [39] is also limited as it only shows small-scale experiments in the class-incremental setting and considers just two SSL methods. In contrast, we present a general framework for CSSL with superior performance, conduct large-scale experiments on three challenging settings, thereby presenting a deeper analysis of CSSL.

## 3. Preliminaries

**Self-Supervised Learning.** The training procedure of several state-of-the-art SSL methods [3,7,8,13,19,26,28,58] can be summarized as follows. Given an image $x$ in a batch sampled from a distribution $\mathcal{D}$, two correlated views $x^A$ and $x^B$ are extracted by applying stochastic image augmentations, such as random cropping, color jittering and horizontal flipping. View $x^A$ is fed to an encoder $f_\theta = f_p \circ f_b$, which is parametrized by $\theta$ and has a backbone $f_b$ and a projection head $f_p$, that extracts feature representations $z^A = f_\theta(x^A)$. Similarly, $x^B$ is forwarded into the same networks, or possibly copies thereof, updated with exponential moving average (EMA), to obtain the representation $z^B$. A loss function $\mathcal{L}_{SSL}$ is applied to these representations to learn the parameters $\theta$ as follows:

$$\underset{\theta}{\operatorname{argmin}} \ \mathbb{E}_{x \sim \mathcal{D}} \left[ \mathcal{L}_{SSL} \left( z^A, z^B \right) \right]. \qquad (1)$$

More details on the implementation of $\mathcal{L}_{SSL}$ are provided in Sec. 5.1 and Tab. 1. This procedure turns out to be extremely powerful at extracting visual representations from large unlabeled datasets. The intuition behind the success of these models is that they learn to be invariant to augmentations. Importantly, augmentations are hand-crafted in a way

that the two views $\boldsymbol{x}^A$ and $\boldsymbol{x}^B$ contain roughly the same semantics as $\boldsymbol{x}$, but their overall appearance (geometry, colors, resolution, *etc*.) is different. This forces the model map images with the same semantics to similar regions of the feature space. Interestingly, these augmentations are much stronger, *i.e.*, they distort the image more, than augmentations commonly used to train supervised models.

**Continual Learning.** The CL problem focuses on training models such as deep neural networks from non-stationary data distributions. More formally, this involves a network $f'_{\theta'} = f'_c \circ f'_b$ with parameters $\theta'$, backbone $f'_b$ and a classifier $f'_c$, that learns from an ordered set of tasks $\{1, \ldots, T\}$, each exhibiting a different data distribution $\mathcal{D}_t$. Usually, an image $\boldsymbol{x}$ sampled i.i.d. from $\mathcal{D}_t$ is processed by $f'$ that predicts a probability distribution $\boldsymbol{p}$ over the set of classes $\mathcal{Y}_t$. The objective is to find parameters $\theta'$ such as:

$$\operatorname*{argmin}_{\theta'} \sum_{t=1}^{T} \mathbb{E}_{(\boldsymbol{x},\boldsymbol{y}) \sim \mathcal{D}_t} \left[ \mathcal{L}_{CL} \left( \boldsymbol{p}, \boldsymbol{y} \right) \right], \qquad (2)$$

where, in most cases, $\mathcal{L}_{CL}$ is the cross-entropy loss. However, during task $t$, the previous data distribution $\mathcal{D}_{t-1}$ is not available and therefore Eq. (2) cannot be minimized directly. Current research focuses on approximating $\theta'$ using indirect approaches. Some of them [18, 36] are based on knowledge distillation [30], *i.e.*, transferring knowledge from one network to another by forcing them to produce the same outputs. We will discuss the applicability of distillation methods in CSSL in Sec. 5.

# 4. Continual Self-Supervised Learning

In this paper, we tackle the problem of Continual Self-Supervised Learning as an extension of both SSL and CL. In practice, a CSSL experiment starts with the first task, where the model is trained as per the specific self-supervised method that it implements, with no difference from offline training. Subsequent tasks are then presented to the model sequentially, and the data from the previous tasks are discarded. No labels are provided during this training phase. For the sake of simplicity and since we are exploring a new, challenging setting, we assume task boundaries to be provided to the model. More formally, the CSSL objective is to learn a strong feature extractor that is invariant to augmentations on all tasks. Following the notation introduced in Sec. 3, we define:

$$\operatorname*{argmin}_{\theta} \sum_{t=1}^{T} \mathbb{E}_{\boldsymbol{x} \sim \mathcal{D}_t} \left[ \mathcal{L}_{SSL} \left( \boldsymbol{z}^A, \boldsymbol{z}^B \right) \right]. \qquad (3)$$

Note the absence of labels $\boldsymbol{y}$ when sampling from $\mathcal{D}_t$, the summation over the set of tasks inherited from Eq. (2) and the SSL loss function in Eq. (1). The expectation is approximated using stochastic gradient descent on minibatches.

**Evaluation.** After each task, it is possible (for evaluation purposes) to train a linear classifier on top of the obtained backbone $f_b$. With this linear classifier we report accuracy on the test set. This protocol is compatible with standard CL metrics, as shown in Sec. 6.1. We explore three CSSL settings in our work.

▶ **Class-incremental**: each task $t$ is represented by a dataset $D_t \sim \mathcal{D}_t$ containing images that belong to a set of classes $\mathcal{Y}_t$ such that $\mathcal{Y}_t \cap \mathcal{Y}_s = \varnothing$ for each other task $s \neq t$. Note that the class labels are only used for splitting the dataset and they are unknown to the model. In practice, the set of classes in the dataset are shuffled and then partitioned into $T$ tasks. Each task contains the same number of classes.

▶ **Data-incremental**: each task $t$ contains a set of images $D_t$ such that $D_t \cap D_s = \varnothing$ for each other task $s \neq t$. No additional constraints are imposed on the classes. In practice, the whole dataset is shuffled and then partitioned into $T$ tasks. Each task can potentially contain all the classes.

▶ **Domain-incremental**: each task $t$ contains a set of images $D_t$ drawn from a different domain. We assume that the set of classes $\mathcal{Y}_t$ in each dataset remains the same for all tasks but the data distribution changes, as if the data were collected from different sources.

# 5. The CaSSLe Framework

We now introduce "CaSSLe", our framework for continual self-supervised learning of visual representations and detail its compatibility with several SSL methods.

**Distillation in CSSL.** From a supervised CL perspective, the concept of invariance is interesting. Here, we would like to learn representations of previously-learned semantic concepts that are invariant to the state of the model's parameters. Indeed, this idea was investigated in prior works [18, 31] that leverage knowledge distillation for CL. However, such approaches are only mildly effective in a CSSL scenario, as we show in Sec. 6. We believe this is due to CSSL being fundamentally different from supervised CL. In CSSL, we aim to extract the best possible representations that can be subsequently reused in a variety of tasks, and maximize the linear separability of features at the end of the CL phase. Thus, the linear classifier does not benefit much from the stability of the representations. Also, forcing the representations not to change may prevent the model from learning new concepts. This is especially critical for SSL methods for two reasons: (i) the performance of the models improve substantially with longer training, implying that the representations continue to get refined, and (ii) they exhibit different losses and feature normalizations that might interfere with distillation and vice-versa (*e.g.*, BarlowTwins uses standardization while [18, 31] use $l2$-normalization). Nonetheless, the features still need to be informative of previous tasks to maximize the separability of the old distribution but the current state might be too different from the pre-
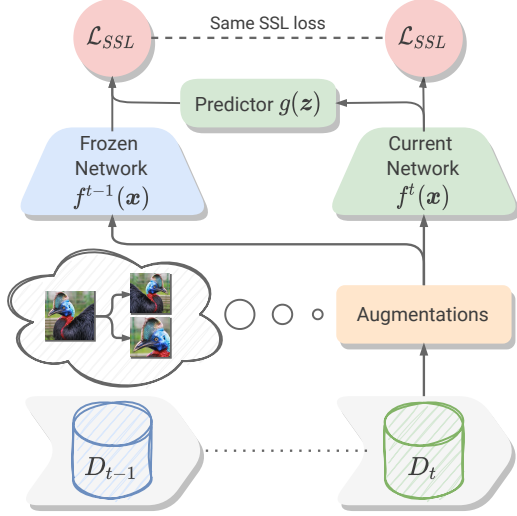
Figure 2. Overview of the CaSSLe framework.

vious one making comparing representations complicated.

**Distillation through SSL losses.** Our framework, shown in Fig. 2, is based on the following ideas: (i) a predictor network that maps the current state of the representations to their past state, by leveraging a distillation through time strategy that satisfies both stability and plasticity principles, and (ii) a family of adaptable distillation losses inherited from the SSL literature that solves the issue of having different objectives interfering with each other.

When a new task is received, we start by making a copy of the current model. This copy does not require gradient computation and will not be updated. We call this the *frozen encoder* $f^{t-1}$. As soon as an image $\boldsymbol{x} \in D_t$ is available we apply our stochastic image augmentations and extract its features $\boldsymbol{z} = f^t(\boldsymbol{x})$. In addition, we also use the frozen encoder to extract another feature vector $\bar{\boldsymbol{z}} = f^{t-1}(\boldsymbol{x})$. Now, our goal is to ensure that $\boldsymbol{z}$ contains at least as much information as (and ideally more than) $\bar{\boldsymbol{z}}$. Instead of enforcing the two feature vectors to be similar, and hence discouraging the new model from learning new concepts, we propose to use a predictor network $g$ to project the representations from the new feature space to the old one. If the predictor is able to perfectly map from one space to the other, then it implies that $\boldsymbol{z}$ is at least as powerful as $\bar{\boldsymbol{z}}$.

We are now ready to perform distillation, but which is the most appropriate distillation loss? Since we want the representations produced by $g$ to be invariant to the state of the model, we propose to use the same SSL loss used to simulate invariance to augmentations. Empirically, we verify that this choice reduces interference and minimizes the need for hyperparameter tuning. We can hence write a generic distillation loss by reusing the definition of $\mathcal{L}_{SSL}$:

$$\mathcal{L}_D(\boldsymbol{z}, \bar{\boldsymbol{z}}) = \mathcal{L}_{SSL}(g(\boldsymbol{z}), \bar{\boldsymbol{z}}). \tag{4}$$

Note that $\bar{\boldsymbol{z}}$ is always detached from the computational graph, such that the frozen encoder does not receive any gradient, and the gradient only flows through the predictor $g$, as prescribed in [15]. On the one hand, if training converges and $\mathcal{L}_D$ is minimized, the features predicted by $g$ will likely be quasi-invariant to the state of the model, which satisfies the stability principle. On the other hand, the current encoder is less bound to its previous state, hence representations $\boldsymbol{z}$ can be more plastic. The loss can be extended to multiple views by applying it to both representations, *i.e.*, $\mathcal{L}_D(\boldsymbol{z}^A, \bar{\boldsymbol{z}}^A) + \mathcal{L}_D(\boldsymbol{z}^B, \bar{\boldsymbol{z}}^B)$, and also swapped distillation, *e.g.*, $\mathcal{L}_D(\boldsymbol{z}^A, \bar{\boldsymbol{z}}^B)$ and vice-versa (see ablation in Tab. 6).

The final loss of an SSL method trained continually with the CaSSLe framework is given by:

$$\begin{aligned}\mathcal{L} &= \mathcal{L}_{SSL}(\boldsymbol{z}^A, \boldsymbol{z}^B) + \mathcal{L}_D(\boldsymbol{z}^A, \bar{\boldsymbol{z}}^A) \\ &= \mathcal{L}_{SSL}(\boldsymbol{z}^A, \boldsymbol{z}^B) + \mathcal{L}_{SSL}(g(\boldsymbol{z}^A), \bar{\boldsymbol{z}}^A).\end{aligned} \tag{5}$$

This loss can be made symmetric by applying it to both the views (swapping $A$ and $B$ in Eq. (5)) and it can also be easily adapted for multi-crop [7]. Note that we do not use any hyperparameter to weight the importance of the distillation loss with respect to the SSL loss.

## 5.1. Compatibility of SSL methods with CaSSLe

The main difference among SSL methods is the loss function that they use. Following the notation defined in Sec. 3, and the loss functions in Tab. 1, we now detail if and how SSL losses can be used in our CaSSLe framework. Full derivation of distillation losses is deferred to the supplementary material.

**InfoNCE-based** methods [13, 28] perform instance discrimination, where positive samples help to build invariance to augmentations. The negatives prevent the model from falling into degenerate solutions. The InfoNCE (*a.k.a.* contrastive) loss can be written as in Eq. (6), where subscript $i$ is the index of a generic sample in the batch, sim is the cosine similarity and $\eta(i)$ is the set of negatives for sample $i$ in the current batch. Distilling knowledge with this loss is equivalent to performing instance discrimination of current task samples but in the feature space learnt in the past. Thus, the predictor $g$ learns to project samples from the present to the past space to maximize the distance with the negative samples, and the similarity with itself in the past.

**MSE-based** approaches [15, 26] enforce consistency among positive samples and ignore the negatives. BYOL [26] uses a momentum encoder and SimSiam [15] performs a stop gradient operation to avoid degenerate solutions. Since the representations are $l2$-normalized, their loss (Eq. 7) can be rewritten as the negative cosine similarity: $-\text{sim}(\boldsymbol{q}^A, \boldsymbol{z}^B) = -\frac{\boldsymbol{q}^A}{||\boldsymbol{q}^A||_2} \cdot \frac{\boldsymbol{z}^B}{||\boldsymbol{z}^B||_2}$, where $\boldsymbol{q}^A = h(\boldsymbol{z}^A)$ and $h$ is a prediction head. The gradient is backpropagated only

Table 1. Overview of state-of-the-art SSL methods and losses. In all tables, highlight colors are coded according to the type of loss.

| Methods | Loss | Equation | |
|---|---|---|---|
| SimCLR [13]<br>MoCo [28]<br>NNCLR [19] | InfoNCE | $-\log \frac{\exp\left(\text{sim}\left(\boldsymbol{z}_i^A, \boldsymbol{z}_i^B\right)/\tau\right)}{\sum_{\boldsymbol{z}_j \in \eta(i)} \exp\left(\text{sim}\left(\boldsymbol{z}_i^A, \boldsymbol{z}_j\right)/\tau\right)}$ | (6) |
| BYOL [26]<br>SimSiam [15]<br>VICReg [3] | MSE | $-\lVert \boldsymbol{q}^A - \boldsymbol{z}^B \rVert_2^2$ | (7) |
| SwAV [7]<br>DCV2 [7]<br>DINO [8] | Cross-entropy | $-\sum_d \boldsymbol{a}_d^B \log \frac{\exp\left(\text{sim}\left(\boldsymbol{z}^A, \boldsymbol{c}_d\right)/\tau\right)}{\sum_k \exp\left(\text{sim}\left(\boldsymbol{z}^A, \boldsymbol{c}_k\right)/\tau\right)}$ | (8) |
| Barlow Twins [58]<br>VICReg [3] | Cross-correlation | $\sum_u \left(1 - \mathcal{C}_{uv}\right)^2 + \lambda \sum_u \sum_{v \neq u} \mathcal{C}_{uv}^2$ | (9) |

through the representations of the first argumentation. A special case of this family of methods is VICReg [3], which uses a combination of multiple losses, where MSE acts as invariance term. Features are not $l2$-normalized in VICReg and its predictor is the identity function. In our framework, this loss encourages the model to predict the past state of representations without additional regularization.

**Cross-entropy-based.** Instead of simply enforcing invariance of the representations to augmentations, cluster prototypes $\mathbf{C} = \{\mathbf{c}_1, \ldots, \mathbf{c}_K\}$ are used as a proxy in these approaches, so that the model learns to predict invariant cluster assignments. Slight variations of this idea result in different methods: SwAV [7], DeepClusterV2 [7] and DINO [8]. Once a probability distribution over the prototypes is predicted, the cross-entropy loss (Eq. 8) is used to compare the two views. Features and cluster prototypes $\boldsymbol{c}$ are $l2$-normalized. The assignments $\boldsymbol{a}^B$ can be calculated in several ways, e.g., k-means in DeepCluster, Sinkhorn-Knopp in SwAV and EMA in DINO. When employed as a distillation loss, cross-entropy encourages $g$ to predict the assignments generated by the frozen encoder with a set of frozen prototypes: $\boldsymbol{a}^B = \frac{\exp\left(\text{sim}\left(\bar{\boldsymbol{z}}^B, \boldsymbol{c}_d^{t-1}\right)/\tau\right)}{\sum_k \exp\left(\text{sim}\left(\bar{\boldsymbol{z}}^B, \boldsymbol{c}_k^{t-1}\right)/\tau\right)}$, where $\mathbf{C}^{t-1} = \{\mathbf{c}_1^{t-1}, \ldots, \mathbf{c}_K^{t-1}\}$.

**Cross-correlation-based**. These methods use a different approach based on decorrelating the components of the feature space, e.g., Barlow Twins [58], VICReg [3] and W-MSE [20]. For our analysis, we will mainly focus on Barlow Twins' implementation of this objective. Extensions to VICReg are left for future work. The cross-correlation based objective function is shown in Eq. 9, where $\lambda$ is an hyperparameter to control the importance of the first and the second terms of the loss, and $\mathcal{C}_{uv} = \frac{\sum_i \boldsymbol{z}_{i,u}^A \boldsymbol{z}_{i,v}^B}{\sqrt{\sum_i \left(\boldsymbol{z}_{i,u}^A\right)^2} \cdot \sqrt{\sum_i \left(\boldsymbol{z}_{i,v}^B\right)^2}}$ is the value of position $(u, v)$ of the cross-correlation matrix computed between the representations of the views along the batch dimension. Note that the representations here are mean centered along the batch dimension, such that each unit has mean output zero over the batch. Performing distillation with this loss has the addi-

tional effect of decorrelating the dimensions of the predicted features $g(\boldsymbol{z}^A)$.

# 6. Experiments

## 6.1. Experimental Protocol

**Evaluation Metrics.** Following previous work [38], we propose the following metrics to evaluate the quality of the representations extracted by our CSSL model:
► **Linear Evaluation Accuracy**: accuracy of a classifier trained on top of the backbone $f_b$ on all tasks (or a subset, e.g., 10% of the data) or a downstream task. For class-incremental and data-incremental, we use the task-agnostic setting, meaning that at evaluation time we do not assume to know the task ID. For the domain-incremental setting, we perform both task-aware and task-agnostic evaluations (the latter is discussed in the supplementary material). To calculate the average accuracy we compute $A = \frac{1}{T} \sum_{i=1}^{T} A_{T,i}$, where $A_{j,k}$ is the linear evaluation accuracy of the model on task $k$ after observing the last sample from task $j$.
► **Forgetting**: a common metric in the CL literature, it quantifies how much information the model has forgotten about previous tasks. It is formally defined as: $F = \frac{1}{T-1} \sum_{i=1}^{T-1} \max_{t \in \{1,\ldots,T\}} \left(A_{t,i} - A_{T,i}\right)$.
► **Forward Transfer**: measures how much the representations that we learned so far are helpful in learning new tasks, namely: $FT = \frac{1}{T-1} \sum_{i=2}^{T} A_{i-1,i} - R_i$ where $R_i$ is the linear evaluation accuracy of a random network on task $i$.

**Datasets.** We perform experiments on 3 datasets: CIFAR100 [35] (class-incremental), a 100-class dataset with 60k 32x32 colour images; ImageNet100 [54] (class- and data-incremental), 100-class subset of the ILSVRC2012 dataset with ≈130k images in high resolution (resized to 224x224); DomainNet [43] (domain-incremental), a 345-class dataset containing roughly 600k high-resolution images (resized to 224x224) divided into 6 domains. We experiment with 5 tasks for the class- and data-incremental settings and with 6 tasks (one for each domain in Domain-Net) in the case of domain-incremental. The supplementary material presents additional results with different number of tasks. For the domain-incremental setting, we order the domains in decreasing number of images.

**Implementation details.** The SSL methods are adapted from `solo-learn` [16], an established SSL library, which is the main code base for all our experiments. The number of epochs per task is as follows: 500 for CIFAR100, 400 for ImageNet100, 200 for DomainNet. The backbone $f_b$ is a ResNet18 [29], with batch size 256. We use LARS [57] for all our experiments. The offline version of each method, that serves as an upper bound, is trained for the same number of epochs as the continual counterpart for a fair comparison. All the results for offline upper bounds are ob-

Table 2. Comparison with state-of-the-art CL methods on CI-FAR100 (5 tasks, class-incremental) using linear evaluation top-1 accuracy, forgetting and forward transfer.

| Strategy | SimCLR | | | Barlow Twins | | | BYOL | | |
|---|---|---|---|---|---|---|---|---|---|
| | A ($\uparrow$) | F ($\downarrow$) | T ($\uparrow$) | A ($\uparrow$) | F ($\downarrow$) | T ($\uparrow$) | A ($\uparrow$) | F ($\downarrow$) | T ($\uparrow$) |
| Fine-tuning | 48.9 | 1.0 | 33.5 | 54.3 | 0.4 | 39.2 | 52.7 | 0.1 | 35.9 |
| EWC [34] | 53.6 | **0.0** | 33.3 | 56.7 | **0.2** | 39.1 | 56.4 | **0.0** | 39.9 |
| ER [47] | 50.3 | 0.1 | 32.7 | 54.6 | 3.0 | 39.4 | 54.7 | 0.4 | 36.3 |
| DER [4] | 50.7 | 0.4 | 33.2 | 55.3 | 2.5 | 39.6 | 54.8 | 1.1 | 36.7 |
| LUMP [39] | 52.3 | 0.3 | 34.5 | 57.8 | 0.3 | 41.0 | 56.4 | 0.2 | 37.9 |
| Less-Forget [31] | 52.5 | 0.2 | 33.8 | 56.4 | **0.2** | 40.1 | 58.6 | 0.2 | 41.1 |
| POD [18] | 51.3 | 0.1 | 33.8 | 55.9 | 0.3 | 40.3 | 57.9 | **0.0** | 41.1 |
| CaSSLe | **58.3** | 0.2 | **36.4** | **60.4** | 0.4 | **42.2** | **62.2** | **0.0** | **43.6** |
| Offline | 65.8 | - | - | 70.9 | - | - | 70.5 | - | - |

Table 3. Comparison with Lin *et al.* [37] on CIFAR100 (2 and 5 tasks, class-incremental setting). MoCoV2+ is an updated version of MoCoV2 that uses a symmetric loss. The difference between the two is ≈1% at convergence [15].

| Strategy | Method | 2 Tasks | 5 Tasks |
|---|---|---|---|
| Lin *et al.* [37] | SimCLR | 55.7 | - |
| | MoCoV2 | 56.1 | 53.8 |
| CaSSLe | SimCLR | **61.8** | **58.3** |
| | MoCoV2+ | **63.3** | **59.5** |

tained using the checkpoints provided in [16]. For some SSL methods, it was necessary to slightly increase the learning rate over the values provided by [16] in order for the methods to fully convergence in the CSSL setting. Although tuning the hyperparameters might be beneficial in some settings, we do **not** perform any hyperparameter tuning for CaSSLe. We also neither change the parameters of the SSL methods, nor use a weight for the distillation loss (as per Eq. (5)).

**Baselines.** Most of the CL methods require labels which makes them unsuitable for CSSL. However, a few works can be adapted for our setting with minimal changes. We choose baselines from three categories [17]: prior-focused regularization (EWC [34]), data-focused regularization (POD [18], Less-Forget [31]), and rehearsal-based replay (ER [47], DER [4]) methods. We also compare with two concurrent works that propose approaches for CSSL (LUMP [39], Lin *et al.* [37]). Finally, we do not consider methods based on VAEs [1,45], since they have been shown to yield poor performance on large scale. Details on how the baselines are selected, implemented and tuned for CSSL can be found in the supplementary material.

## 6.2. Results

**Comparison with the state of the art.** In Tab. 2 we report comparison with CL baselines and fine-tuning in composition with three SSL methods: SimCLR, Barlow Twins and BYOL. We select these three methods for the following reasons: (i) they feature different losses (InfoNCE, Cross-correlation and MSE), (ii) they exhibit different feature normalizations ($l2$, standardization and mean centering), and (iii) they use different techniques to avoid collapse (neg-
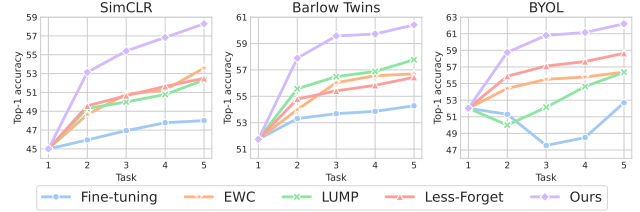


Figure 3. Evolution of top-1 linear evaluation accuracy over tasks on CIFAR100 (5 tasks, class-incremental).

atives, redundancy reduction, momentum encoder). The comparison is performed on class-incremental CIFAR100 with 5 tasks. Offline learning results are reported as upper bound.

First, we notice that CaSSLe produces better representations than all the other strategies, outperforming them by large margins with all SSL methods in terms of top-1 accuracy. Moreover, our framework also shows better forward transfer, meaning that its features are easier to generalize to other tasks (also evident in Tab. 8). CaSSLe appears to reduce catastrophic forgetting with respect to fine-tuning, and is comparable to other methods. In general, SSL methods already have low forgetting with respect to supervised learning on CIFAR100 (see Tab. 4) and therefore there is little margin for improvement. However, on higher resolution images (ImageNet100) CaSSLe actually achieves remarkable results in the mitigation of catastrophic forgetting.

Replay-based methods (ER, DER) clearly do not help against forgetting in CSSL. We found two reasons for this failure. First, in supervised CL, replay-based methods benefit from storing labels, which contain a lot of information about previous tasks and enable the retraining of the linear classifier on old classes. This is not the case in CSSL, where labels are unavailable. Second, SSL models need more training epochs to converge, which means that samples in the buffer are also replayed many more times. This causes severe overfitting on these exemplars, defeating the purpose of the replay buffer. LUMP mitigates this effect by augmenting the buffer using mixup but does not reach too far, surpassing other baselines only with Barlow Twins. EWC holds up surprisingly well, outperforming more recent methods, meaning that the importance of the weights can be calculated accurately with the self-supervised loss. Distillation methods (POD, Less-Forget) show good performance. However, they use $l2$-normalization in their loss, causing loss of information when coupled with Barlow Twins, which decreases accuracy.

Fig. 3 shows the evolution of top-1 linear evaluation accuracy over the whole training trajectory on class-incremental CIFAR100 with 5 tasks. CaSSLe outperforms the other methods, and keeps improving throughout the sequence. We found BYOL to be unstable when simply fine-tuning the model. CaSSLe, EWC and Less-Forget mitigate this instability completely. On the other hand, LUMP first

Table 4. Linear evaluation top-1 accuracy on class-incremental CIFAR100 and ImageNet100 with 5 tasks. CaSSLe is compared to fine-tuning, offline and supervised learning.

| Method | Strategy | CIFAR100 | | | ImageNet100 | | |
|---|---|---|---|---|---|---|---|
| | | A ($\uparrow$) | F ($\downarrow$) | T ($\uparrow$) | A ($\uparrow$) | F ($\downarrow$) | T ($\uparrow$) |
| Barlow Twins | Fine-tuning | 54.3 | **0.4** | 39.2 | 63.1 | 10.7 | 44.4 |
| | CaSSLe | **60.4** | 0.4 | **42.2** | **68.2** | **1.3** | **47.9** |
| | Offline | 70.9 | - | - | 80.4 | - | - |
| SwAV | Fine-tuning | 55.5 | **0.0** | 32.8 | 64.4 | 4.3 | 42.8 |
| | CaSSLe | **57.8** | 0.0 | **34.5** | **66.0** | **0.2** | **43.6** |
| | Offline | 64.9 | - | - | 74.3 | - | - |
| BYOL | Fine-tuning | 52.7 | 0.1 | 35.9 | 66.0 | 2.9 | 43.2 |
| | CaSSLe | **62.2** | **0.0** | **42.2** | **66.4** | **1.1** | **46.6** |
| | Offline | 70.5 | - | - | 80.3 | - | - |
| VICReg | Fine-tuning | 51.5 | 0.9 | 36.4 | 61.3 | 7.9 | 42.0 |
| | CaSSLe | **53.6** | **0.2** | **41.1** | **64.8** | **4.3** | **45.3** |
| | Offline | 68.5 | - | - | 79.4 | - | - |
| MoCoV2+ | Fine-tuning | 47.3 | 0.2 | 33.4 | 62.0 | 8.4 | 41.6 |
| | CaSSLe | **59.5** | **0.0** | **39.6** | **68.8** | **1.5** | **46.8** |
| | Offline | 69.9 | - | - | 79.3 | - | - |
| SimCLR | Fine-tuning | 48.9 | 1.0 | 33.5 | 61.5 | 8.1 | 40.3 |
| | CaSSLe | **58.3** | **0.2** | **36.4** | **68.0** | **2.2** | **45.8** |
| | Offline | 65.8 | - | - | 77.5 | - | - |
| Supervised | Fine-tuning | 54.1 | 6.8 | 36.5 | 63.1 | 5.6 | 42.5 |
| | Offline | 75.6 | - | - | 81.9 | - | - |

Table 5. Training 5 times longer on 1/5 of the data vs. training continually w/ and w/o CaSSLe on ImageNet100 (5 tasks, class- and data-incremental). **Bold** is best, underlined is second best.

| Setting | Method | Fine-tune | Offline 1/5 | CaSSLe |
|---|---|---|---|---|
| Class-inc. | SimCLR | 61.5 | _63.1_ | **68.0** |
| | Barlow Twins | 63.1 | _63.5_ | **68.2** |
| | BYOL | _66.0_ | 60.6 | **66.4** |
| Data-inc. | SimCLR | _68.9_ | 67.2 | **72.1** |
| | Barlow Twins | _71.3_ | 70.2 | **74.9** |
| | BYOL | **74.0** | 66.7 | _73.3_ |

Table 6. Ablation study of design choices in CaSSLe.

| Strategy | Method | Swap | No pred. | Ours |
|---|---|---|---|---|
| CaSSLe | SimCLR | 49.3 | 52.6 | **58.3** |
| | Barlow Twins | 57.4 | 57.3 | **60.4** |
| | BYOL | 52.0 | 58.6 | **62.2** |

drops slightly and then recovers. We believe this is due to some instability introduced by the mixup regularization, to which the model takes time to adapt.

In Tab. 3 we also compare with Lin *et al.* [37] on class-incremental CIFAR100. Although our method is not specifically designed for contrastive learning, it substantially outperforms Lin *et al.* with 2 and 5 tasks. It is worth nothing that MoCoV2+ is slightly better than MoCoV2 ($\approx$1% difference), whereas our gains are much larger ($\approx$7%).

**Ablation study.** We ablate the most critical design choices we adopt in CaSSLe: (i) distillation without swapped views, and (ii) the presence of a prediction head $g$. These results are reported in Tab. 6. Our full framework clearly outperforms its variants with swapped views and without predictor. This validates our hypothesis that a predictor to map new features to the old feature space is crucial. The result that swapping views does not help is likely due to the frozen encoder not being invariant to the current task.

**Class-incremental.** In Tab. 4 we report a study of CSSL with 6 SSL methods in composition with the CaSSLe framework on class-incremental CIFAR100 and ImageNet100. Fine-tuning and Offline SSL results are reported as lower and upper bounds. The accuracy of supervised learning is also reported. CaSSLe always improves with respect to fine-tuning. In particular, our framework produces higher forward transfer and lower forgetting, especially on ImageNet100, where methods tend to forget more. Notably, CaSSLe outperforms supervised fine-tuning, ex-

cept when coupled with VICReg on CIFAR100. On average, SSL methods trained continually with CaSSLe improve by 6.8% on CIFAR100 and 4% on ImageNet100.

**Data-incremental.** Tab. 7 presents results for linear evaluation top-1 accuracy on ImageNet100 with 5 tasks in a data-incremental scenario. While no SSL method is better than supervised fine-tuning, Barlow Twins coupled with CaSSLe is competitive. CaSSLe improves performance in all cases by 2% on average, except for BYOL. This is likely due to the fact that in the data-incremental scenario remembering past knowledge is less important than in other scenarios, and BYOL already has a momentum encoder that provides some information about the past. This hypothesis is validated by the fact that MoCoV2+ (that uses a momentum encoder) improves less than SimCLR when coupled with CaSSLe. We believe that, by tuning the EMA schedule, improvement could also be achieved for BYOL. In addition, BYOL already shows impressive performance with fine-tuning, outperforming all the other methods by more than 2%. Interestingly, SwAV comes closest to its offline upper bound, with only a 3% decrease in performance when coupled with CaSSLe.

**Domain-incremental.** We also examine the capability of CaSSLe to learn continually when the domain from which the data is drawn changes. Tab. 7 shows the average top-1 accuracy of a linear classifier trained on top of the frozen feature extractor on all domains separately (domain-aware). Domain-agnostic evaluation and results for each domain are presented in the supplementary material. Again, CaSSLe improves every method by 4.4% on average, showing that our distillation strategy is robust to domain shift, and although the data distribution is really different, information transfer is still performed. Interestingly, most of the methods, when trained with CaSSLe get very close to their offline accuracy.

**Long training vs continual training.** We also analyze the following question: is it worth training continually or is it better to train for longer on a small dataset? This

Table 7. Linear evaluation accuracy on ImageNet100 (5 tasks, data-incremental) and DomainNet (6 tasks, domain-incremental).

| Method | Strategy | ImageNet100 (Data-inc.) | DomainNet (Domain-inc.) |
|---|---|---|---|
| Barlow Twins | Fine-tuning | 71.3 | 50.3 |
| | CaSSLe | **74.9** | **55.5** |
| | Offline | 80.4 | 57.2 |
| SwAV | Fine-tuning | 70.8 | 49.6 |
| | Knowledge | **71.3** | **54.3** |
| | Offline | 74.3 | 54.6 |
| BYOL | Fine-tuning | **74.0** | 50.6 |
| | CaSSLe | 73.3 | **55.1** |
| | Offline | 80.3 | 56.6 |
| VICReg | Fine-tuning | 70.2 | 49.3 |
| | CaSSLe | **72.3** | **52.9** |
| | Offline | 79.4 | 56.7 |
| MoCoV2+ | Fine-tuning | 69.5 | 43.2 |
| | CaSSLe | **71.9** | **46.7** |
| | Offline | 78.2 | 53.7 |
| SimCLR | Fine-tuning | 68.9 | 45.1 |
| | CaSSLe | **72.1** | **50.0** |
| | Offline | 77.5 | 52.6 |
| Supervised | Fine-tuning | 75.9 | 55.9 |
| | Offline | 81.9 | 66.4 |

Table 8. Downstream performance with different SSL methods trained on Imagenet-100 and evaluated on DomainNet (Real).

| Strategy | Barlow Twins | SwAV | BYOL | VICReg | MoCoV2+ | SimCLR | Supervised |
|---|---|---|---|---|---|---|---|
| Fine-tune | 56.2 | 55.9 | 55.0 | 54.0 | 52.4 | 51.6 | 54.3 |
| CaSSLe | **60.3** | **56.9** | **56.9** | **56.3** | **58.7** | **56.5** | |

Table 9. Top-1 linear accuracy on Imagenet-100 with different SSL methods, semi-supervised setting with 10% and 1% of labels.

| Percentage | Strategy | Barlow Twins | SwAV | BYOL | VICReg | MoCoV2+ | SimCLR | Supervised |
|---|---|---|---|---|---|---|---|---|
| 10% | Fine-tune | 56.6 | 57.6 | 55.7 | 53.6 | 54.9 | 52.5 | 60.8 |
| | CaSSLe | **60.3** | **58.2** | **56.5** | **56.5** | **61.7** | **58.9** | |
| 1% | Fine-tune | 42.6 | 42.5 | 42.3 | 40.4 | 40.9 | 39.7 | 48.1 |
| | CaSSLe | **47.0** | **43.1** | **43.4** | **43.2** | **47.8** | **46.8** | |

improves all SSL methods when compared with fine-tuning.

depends on two factors: (i) the SSL method, and (ii) the CSSL setting. For SimCLR and Barlow Twins in the class-incremental setting it seems to be better to train offline on 1/5 of the classes instead of training continually with 5 tasks. In this setting, offline BYOL seems to suffer from instability, ending up lower than fine-tuning. On the other hand, on the data-incremental setting, fine-tuning outperforms longer training, especially for BYOL, which also outperforms CaSSLe (as explained previously). Apart from this exception, CaSSLe always produces better representations than other strategies, making it the go-to option.

**Downstream and semi-supervised.** In Tab. 8, we present the downstream performance of CaSSLe compared with fine-tuning when trained on ImageNet100 and evaluated on DomainNet (Real). Barlow Twins, SwAV and BYOL show higher performance than the supervised model, even when considering a fine-tuning strategy. This is probably due to the fact that SSL methods tend to learn more general features than their supervised counterparts. CaSSLe improves performance on all the SSL methods, making them surpass the supervised baseline. Lastly, when compared with fine-tuning, CaSSLe improves the performance of SSL methods by 3.4% on average. Tab. 9 contains the top-1 accuracy on ImageNet100 when training a linear classifier on a frozen backbone with limited amount of labels (10% and 1%). First, we can observe that no SSL method with fine-tune surpasses the performance of supervised learning. When using CaSSLe, MoCoV2+ outperforms supervised with 10% labels and, in general, Barlow Twins and MoCoV2+ work best in both semi-supervised settings. CaSSLe

## 7. Conclusion

In this work, we study Continual Self-Supervised Learning (CSSL), the problem of learning a set of tasks without labels continually. We make two important contributions for the SSL and CL communities: (i) we present CaSSLe, a simple and effective framework for CSSL that shows how SSL methods and losses can be seamlessly reused to learn continually, and (ii) we perform a comprehensive analysis of CSSL, leading to the emergence of interesting properties of SSL methods.

**Limitations.** Although CaSSLe shows exciting performance, it has some limitations. First, it is applicable in settings where task boundaries are provided. Second, our framework increases the amount of computational resources needed for training by roughly 30%, both in terms of memory and time. Finally, CaSSLe does not perform clustering, meaning that it is unable to directly learn a mapping from data to latent classes, and thus needs either a linear classifier trained with supervision, or some clustering algorithm.

**Broader impact.** The capabilities of supervised CL agents are bounded by the need for human-produced annotations. CSSL models can potentially improve without the need for human supervision. This facilitates the creation of powerful AIs that may be used for malicious purposes such as discrimination and surveillance. Also, since in CSSL the data is supposed to come from a non-curated stream, the model may be affected by biases in the data. This is problematic because biases are then be transferred to downstream tasks.

# Appendix

## A. PyTorch-like pseudo-code

We provide a PyTorch-like pseudo-code of our method. As you can see, CaSSLe is simple to implement and does not add much complexity to the base SSL method. In this snippet, the losses are made symmetric by summing the two contributions. In some cases, the two losses are averaged instead. In CaSSLe, we symmetrize in the same way as the base SSL method we are considering.

---

**Algorithm 1** PyTorch-like pseudo-code for CaSSLe.

```
# aug: stochastic image augmentation
# f: backbone and projector
# frozen_f: frozen backbone and projector
# g: CaSSLe's predictor
# loss_fn: any SSL loss in Tab. 1 (main paper)

# PyTorchLightning handles loading and optimization
def training_step(x):

    # correlated views
    x1, x2 = aug(x), aug(x)

    # forward backbone and projector
    z1, z2 = f(x1), f(x2)

    # optionally forward predictor...

    # compute SSL loss (symmetric)
    ssl_loss = loss_fn(z1, z2) \\
             + loss_fn(z2, z1)

    # forward frozen backbone and projector
    z1_bar, z2_bar = frozen_f(x1), frozen_f(x2)

    # compute distillation loss (symmetric)
    distill_loss = loss_fn(g(z1), z1_bar) \\
                 + loss_fn(g(z2), z2_bar)

    # no hyperparameter for loss weighting
    return ssl_loss + distill_loss
```

---

## B. Derivation of distillation losses

In this section, we derive distillation losses from the SSL losses in Tab. 1 of the main paper, starting from the definition of our distillation loss:

$$\mathcal{L}_D(\boldsymbol{z}, \bar{\boldsymbol{z}}) = \mathcal{L}_{SSL}(g(\boldsymbol{z}), \bar{\boldsymbol{z}}), \quad (10)$$

where $z$ and $\bar{z}$ are the representations of the current and frozen encoder, and $g$ is CaSSLe's predictor network implemented as a two layer MLP with 2048 hidden neurons and ReLU activation.

**Contrastive based.** Our distillation loss based on contrastive learning is implemented as follows:

$$\mathcal{L}(z_i, \bar{z}_i) = -\log \frac{\exp\left(\text{sim}\left(\boldsymbol{z}_i, \bar{\boldsymbol{z}}_i\right)/\tau\right)}{\sum_{\boldsymbol{z}_j \in \bar{\eta}(i)} \exp\left(\text{sim}\left(\boldsymbol{z}_i, \boldsymbol{z}_j\right)/\tau\right)}, \quad (11)$$

where $\bar{\eta}(i)$ is the set of negatives for the sample with index $i$ in the batch. Note that the negatives are drawn both from the predicted and frozen features.

**MSE based.** This distillation loss is simply the MSE between the predicted features and the frozen features:

$$\mathcal{L}(z, \bar{z}) = -||g(\boldsymbol{z}) - \bar{\boldsymbol{z}}||_2^2. \quad (12)$$

It can be implemented with the cosine similarity as stated in the main manuscript.

**Cross-entropy based.** The cross-entropy loss, when used for distillation in an unsupervised setting, makes sure that the current encoder is able to assign samples to the frozen centroids (or prototypes) consistently with the frozen encoder:

$$\mathcal{L}(z, \bar{z}) = -\sum_d \bar{\boldsymbol{a}}_d \log \frac{\exp\left(\text{sim}\left(g(\boldsymbol{z}), \boldsymbol{c}_d^{t-1}\right)/\tau\right)}{\sum_k \exp\left(\text{sim}\left(g(\boldsymbol{z}), \boldsymbol{c}_k^{t-1}\right)/\tau\right)} \quad (13)$$

where:

$$\bar{\boldsymbol{a}} = \frac{\exp\left(\text{sim}\left(\bar{\boldsymbol{z}}, \boldsymbol{c}_d^{t-1}\right)/\tau\right)}{\sum_k \exp\left(\text{sim}\left(\bar{\boldsymbol{z}}, \boldsymbol{c}_k^{t-1}\right)/\tau\right)}, \quad (14)$$

and the set of frozen prototypes is denoted as follows: $\mathbf{C}^{t-1} = \left\{\mathbf{c}_1^{t-1}, \ldots, \mathbf{c}_K^{t-1}\right\}$.

**Cross-correlation based.** We consider Barlow Twins' [58] implementation of this objective. For VICReg [3] we only consider the invariance term. As a distillation loss, the cross-correlation matrix is computed with the predicted and frozen features:

$$\mathcal{L}(z, \bar{z}) = \sum_u \left(1 - \bar{\mathcal{C}}_{uv}\right)^2 + \lambda \sum_u \sum_{v \neq u} \bar{\mathcal{C}}_{uv}^2, \quad (15)$$

where:

$$\bar{\mathcal{C}}_{uv} = \frac{\sum_i g(\boldsymbol{z}_{i,u}) \bar{\boldsymbol{z}}_{i,v}}{\sqrt{\sum_i g\left(\boldsymbol{z}_{i,u}\right)^2} \cdot \sqrt{\sum_i \left(\bar{\boldsymbol{z}}_{i,v}\right)^2}}. \quad (16)$$

## C. Further discussion and implementation details of the baselines

**Selection.** When evaluating our framework, we try to compare with as many existing related methods as possible. However, given that SSL models are computationally intensive, it was not possible to run all baselines and methods in all the CL settings we considered. As mentioned in the main manuscript, we choose eight baselines (seven related methods + fine-tuning) belonging to three CL macro-categories, and test them on CIFAR100 (class-incremental) in combination with three SSL methods. The selection was based on the ease of adaptation to CSSL and the similarity to our framework.

The most similar to CaSSLe are data-focused regularization methods. Among them, a large majority leverage knowledge distillation using the outputs of a classifier learned with supervision *e.g.* [9, 22, 36], while a few works employ feature distillation [18, 31] which is viable even without supervision. [32] is also related to CaSSLe, but it focuses on memory efficiency which is less interesting in our setting. Also, [32] explicitly uses the classifier after feature adaptation, hence it is unclear how to adapt it

Table A. Linear evaluation top-1 accuracy on DomainNet (6 tasks, domain-incremental setting) w/ and w/o CaSSLe. The sequence of tasks is Real→Quickdraw→Painting→Sketch→Infograph→Clipart. "Aw." stands for task-aware, "Ag," for task-agnostic.

| Method | Strategy | Real | | Quickdraw | | Painting | | Sketch | | Infograph | | Clipart | | Avg. | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Aw. | Ag. | Aw. | Ag. | Aw. | Ag. | Aw. | Ag. | Aw. | Ag. | Aw. | Ag. | Aw. | Ag. |
| Barlow Twins | Finetuning | 56.3 | 50.9 | 54.1 | 45.8 | 42.7 | 35.9 | 49.0 | 41.9 | 22.0 | 17.4 | 59.0 | 52.5 | 50.3 | 43.7 |
| | CaSSLe | **62.7** | **57.1** | **59.1** | **50.6** | **49.2** | **42.1** | **53.8** | **47.7** | **25.5** | **20.6** | **61.9** | **55.6** | **55.5** | **48.9** |
| | Offline | 67.1 | 63.0 | 60.3 | 53.9 | 52.4 | 46.3 | 51.9 | 46.9 | 25.9 | 21.0 | 58.8 | 52.6 | 57.2 | 51.8 |
| SwAV | Finetuning | 57.7 | 52.3 | 53.2 | 43.5 | 43.0 | 35.9 | 46.1 | 39.0 | 21.6 | 16.5 | 53.4 | 46.6 | 49.6 | 42.5 |
| | CaSSLe | **62.8** | **57.8** | **59.5** | **50.2** | **47.5** | **41.2** | **49.5** | **42.5** | **22.5** | **17.9** | **56.5** | **49.6** | **54.3** | **47.5** |
| | Offline | 64.1 | 59.5 | 60.6 | 53.6 | 47.6 | 42.9 | 47.7 | 42.1 | 23.3 | 18.9 | 53.6 | 47.3 | 54.6 | 49.1 |
| BYOL | Finetuning | 58.7 | 53.2 | 51.7 | 41.6 | 44.0 | 37.4 | 49.6 | 43.9 | 23.5 | 19.0 | 58.6 | 53.5 | 50.6 | 43.8 |
| | CaSSLe | **63.7** | **60.5** | **59.3** | **50.9** | **48.6** | **44.1** | **50.4** | **45.2** | **24.1** | **19.4** | **59.0** | **54.4** | **55.1** | **49.7** |
| | Offline | 67.2 | 64.0 | 60.2 | 53.3 | 51.5 | 47.3 | 50.4 | 46.2 | 24.5 | 20.8 | 57.0 | 51.5 | 56.6 | 51.9 |
| VICReg | Finetuning | 54.7 | 49.6 | 53.0 | 44.9 | 42.1 | 34.7 | 49.0 | 41.9 | 21.1 | 16.4 | 58.5 | 52.6 | 49.3 | 42.8 |
| | CaSSLe | **59.0** | **53.2** | **56.4** | **47.8** | **46.0** | **38.9** | **52.3** | **45.6** | **23.9** | **18.5** | **60.9** | **55.3** | **52.9** | **46.1** |
| | Offline | 66.4 | 62.7 | 59.2 | 53.5 | 52.4 | 47.2 | 53.2 | 48.1 | 25.3 | 20.7 | 58.3 | 53.2 | 56.7 | 51.9 |
| SimCLR | Finetuning | 52.5 | 47.6 | 48.2 | 38.1 | 37.5 | 31.7 | 42.8 | 35.7 | 18.8 | 14.4 | 50.9 | 46.8 | 45.1 | 38.4 |
| | CaSSLe | **58.4** | **43.4** | **54.2** | **44.7** | **43.9** | **37.7** | **47.6** | **41.9** | **22.0** | **17.8** | **54.9** | **50.5** | **50.0** | **44.2** |
| | Offline | 62.1 | 59.5 | 58.3 | 52.9 | 46.1 | 42.5 | 45.6 | 41.3 | 22.1 | 18.8 | 51.0 | 45.9 | 52.6 | 48.6 |
| MoCoV2+ | Finetuning | 50.9 | 45.5 | 45.8 | 37.5 | 36.0 | 29.3 | 39.5 | 32.1 | 17.9 | 13.5 | 50.3 | **44.5** | 43.2 | 36.7 |
| | CaSSLe | **56.0** | **50.3** | **48.7** | **40.0** | **40.4** | **33.6** | **42.0** | **35.0** | **19.9** | **15.2** | **51.7** | 44.5 | **46.7** | **38.8** |
| | Offline | 65.2 | 61.3 | 57.9 | 51.3 | 48.7 | 43.1 | 44.7 | 39.1 | 23.4 | 19.0 | 51.3 | 44.8 | 53.7 | 48.4 |
| Supervised Contrastive | Finetuning | 57.7 | 52.6 | 55.3 | 45.5 | 44.9 | 38.0 | 51.7 | 45.0 | 22.6 | 18.3 | 64.0 | 60.0 | 52.1 | 45.4 |
| | CaSSLe | **63.4** | **58.8** | **59.7** | **51.3** | **50.1** | **44.7** | **55.9** | **50.3** | **26.9** | **22.4** | **65.0** | **61.3** | **56.7** | **50.9** |
| | Offline | 67.4 | 65.3 | 65.8 | 63.0 | 53.6 | 50.9 | 56.0 | 53.1 | 28.0 | 25.7 | 62.8 | 59.6 | 60.0 | 57.4 |
| Supervised | Finetuning | 63.0 | 58.2 | 56.9 | 47.6 | 49.1 | 44.0 | 55.7 | 50.3 | 27.7 | 23.3 | 68.6 | 63.5 | 55.9 | 49.8 |
| | Offline | 74.7 | 73.2 | 68.5 | 67.8 | 62.0 | 59.3 | 65.7 | 63.7 | 33.7 | 34.5 | 72.3 | 69.3 | 66.4 | 65.0 |

for CSSL, especially since in SSL positives are generated using image augmentations, which are not applicable to a memory bank of features. On the contrary, augmentations can be used in replay methods, among which we select the most common (ER [47]) and one of the most recent (DER [4]). Regarding prior-focused regularization methods, we choose EWC [34] over others (SI [59], MAS [2], *etc*.) as it is considered the most influential and it works best with task boundaries. We also consider two CSSL baselines: LUMP [39] and Lin *et al*. [37]. Finally, we do not consider methods based on VAEs [1, 45], since they have been shown to yield poor performance in the large and medium scale. For instance, as found by [21], a VAE trained offline on CIFAR10 reaches an accuracy of 57.2%, which is lower than any method (except VICReg) trained continually on CIFAR100 with CaSSLe.

**Implementation.** For EWC, we use the SSL loss instead of the supervised loss to estimate importance weights. For POD and Less-Forget, we only re-implement the feature distillation without considering the parts of their methods that explicitly use the classifier. For DER, we replace the logits of the classifier with the projected features in the buffer. We re-implement all these baselines by adapting them from the official implementation (POD), or from the Mammoth framework provided with [4] (DER, ER, EWC), or from the paper (Less-Forget). We also compare with two concurrent works that propose approaches for CSSL (LUMP [39], Lin *et al*. [37]). LUMP uses k-NN evaluation, therefore we adapt the code provided by the authors to run in our code base. For Lin *et*

*al*., we compare directly with their published results, since they use the same evaluation protocol. We perform hyperparameter tuning for all baselines, searching over 5 values for the distillation loss weights of POD and Less-Forget, 3 values for the weight of the regularization in EWC and 3 replay batch sizes for replay methods. The size of the replay buffer is 500 samples for all replay based methods.

## D. Additional results

**Continual supervised contrastive with CaSSLe.** After the popularization of contrastive learning [13,28] for unsupervised learning of representations, [33] proposed a supervised version of the contrastive loss. Here, we show that CaSSLe is easily extendable to support supervised contrastive learning. The implementation is basically the same as for our vanilla contrastive-based distillation loss. In Tab. B, we show the improvement that CaSSLe brings with respect to fine-tuning, which is sizeable in the class-incremental setting. We also report the same comparison on DomainNet in Tab. A, showing interesting results in both task-aware and task-incremental evaluation.

**Task-agnostic evaluation and domain-wise accuracy on DomainNet.** In the main manuscript, we showed that CaSSLe significantly improved performance in the domain-incremental setting using task-aware evaluation. Here, "task-aware" refers to

Table B. Linear evaluation top-1 accuracy on ImageNet100 (5 tasks, class- and data-incremental).

| Method | Strategy | ImageNet100 | |
|---|---|---|---|
| | | Class-inc. | Data-inc. |
| Supervised Contrastive | Fine-tuning | 61.6 | 74.3 |
| | CaSSLe | **69.6** | **76.9** |

Table C. k-NN evaluation on ImageNet100 (5 tasks, class-incremental) performed on backbone and projected features.

| Method | Strategy | k-NN accuracy ($\uparrow$) | |
|---|---|---|---|
| | | Backbone ($f_b$) | Projector ($f_p$) |
| Barlow Twins | Fine-tuning | 59.1 | 34.4 |
| | CaSSLe | **63.4** | **53.2** |
| SwAV | Fine-tuning | **60.0** | 53.9 |
| | CaSSLe | 59.7 | **61.3** |
| BYOL | Fine-tuning | 57.1 | 33.0 |
| | CaSSLe | **61.2** | **60.8** |
| VICReg | Fine-tuning | 56.7 | 35.3 |
| | CaSSLe | **59.5** | **43.4** |
| MoCoV2+ | Fine-tuning | 54.5 | 39.0 |
| | CaSSLe | **61.5** | **53.1** |
| SimCLR | Fine-tuning | 54.8 | 40.1 |
| | CaSSLe | **61.7** | **53.2** |

the fact that linear evaluation is performed on each domain separately, *i.e.* a different linear classifier is learned for each domain. However, it might also be interesting to check the performance of the model when the domain is unknown at test time. For this reason, we report the performance of our model when evaluated in a task-agnostic fashion. In addition, we also show the accuracy on each task (*i.e.* domain). All this information is presented in Tab. A. CaSSLe **always** outperforms fine-tuning with both evaluation protocols. The accuracy of CaSSLe on "Clipart" is also higher than offline. This is probably due to a combination of factors: (i) Clipart is the last task, therefore it probably benefits in forward transfer and (ii) a similar effect to the one found in [53], where dividing data in subgroups tends to enable the learning of better representations. Also, we notice that task-agnostic accuracy is lower than the task-aware counterpart. This is expected and means that the class conditional distributions are not perfectly aligned in different domains. As in the main paper, the colors are related to the type of SSL loss.

**Additional results with k-NN evaluation.** For completeness, in this supplementary material, we also show that CaSSLe yields superior performance when evaluated with a k-NN classifier instead of linear evaluation. We use weighted k-NN with l2-normalization (cosine similarity) and temperature scaling as in [8]. Since since k-NN is much faster than linear evaluation we could also assess the quality of the projected representations, instead of just using the backbone. The results can be inspected in Tab. C. Three interesting phenomena arise: (i) CaSSLe always improves with respect to fine-tuning, (ii) the features of the backbone $f_b$ are usually better than the features of the projector $f_p$ and (iii) CaSSLe causes information retention in the projector, which significantly increases the performance of the projected features. An excep-

Table D. Linear evaluation top-1 accuracy on CIFAR100 (10 tasks, class-incremental).

| Method | Strategy | A ($\uparrow$) |
|---|---|---|
| SimCLR | Fine-tuning | 39.3 |
| | CaSSLe | **52.7** |
| Barlow Twins | Fine-tuning | 49.9 |
| | CaSSLe | **53.7** |

Table E. Linear evaluation top-1 accuracy on ImageNet100 (5 tasks, class- and data-incremental) with ResNet50 [29].

| Method | Strategy | A ($\uparrow$) | |
|---|---|---|---|
| | | Class-inc. | Data-inc. |
| SimCLR | Fine-tuning | 70.7 | 75.6 |
| | CaSSLe | **74.0** | **77.2** |
| Barlow Twins | Fine-tuning | 71.2 | 75.8 |
| | CaSSLe | **74.8** | **78.1** |

tion is represented by SwAV [7], that seems to behave differently to other methods. First, the accuracy of the projected features in SwAV is much higher than other methods. This might be due to the fact that it uses prototypes, which bring the representations 1 layer away from the loss, making them less specialized in the SSL task. Second, it seems that CaSSLe only improves the projected features when coupled with SwAV. However, this is probably an artifact of the evaluation procedure, as the l2-normalization probably causes loss of information. Indeed, although the overall performance is lower, SwAV + CaSSLe outperforms SwAV + fine-tuning (58.7% vs 56.9%) if the euclidean distance is used in place of the cosine similarity for the backbone features. We leave a deeper investigation of this phenomenon for future work.

**Different number of tasks.** The analysis of CSSL settings that we show in the main manuscript is limited to the 5 task scenario. However, it is interesting to run the same benchmarks with a longer task sequence. Nonetheless, one should also remember that SSL methods are data hungry, hence the less data is available per task, the higher the instability of the SSL models. In Tab. D, we present additional results with 10 tasks on CIFAR100 (class-incremental). Barlow Twins seems to hold up surprisingly well, finishing up at roughly 50% accuracy, while SimCLR suffers in the low data regime. Nonetheless, CaSSLe outperforms fine-tuning with Barlow Twins, and to a very large extent with SimCLR.

**Deeper architectures.** The experiments we propose in the main manuscript feature a ResNet18 network. This is a common choice in CL. However, in SSL, it is more common to use ResNet50. For this reason, in Tab. E we show that the same behavior observed with smaller networks is also obtained with deeper architectures. More specifically, CaSSLe outperforms fine-tuning in both class- and data-incremental settings by large margins.

**The role of the predictor.** In the main manuscript, we provided an intuitive explanation of the role of the predictor network that maps the current feature space to the frozen feature space.

Table F. Combinations of SSL methods and distillation losses on CIFAR100 (class-incremental, 2 tasks).

| Distillation Loss | SimCLR | Barlow Twins | BYOL |
|---|---|---|---|
| InfoNCE | **61.8** | 64.5 | 64.8 |
| Cross-correlation | 60.1 | **67.2** | 65.8 |
| MSE | 61.3 | 64.6 | **66.7** |

This intuition is corroborated by extensive experimentation and ablation studies. However, one more thing that is worth mentioning is that the success of the predictor network might also be related to the findings in SimSiam [15], BYOL [26] and Direct-Pred [52]. Moreover, we perform additional ablations on the design of CaSSLe's predictor for SimCLR on CIFAR100 (5 tasks): adding BatchNorm after the hidden layer does not make any difference in terms of performance, and removing the non-linearity only causes a 0.3% drop in accuracy.

**Combinations of SSL methods and distillation losses.** For computational reasons, it was not feasible to perform experiments combing all SSL methods with all possible distillation losses. However, in Tab. F we provide a subset of the possible combinations to validate our strategy that uses the same SSL loss for distillation.

# References

[1] Alessandro Achille, Tom Eccles, Loic Matthey, Christopher P Burgess, Nick Watters, Alexander Lerchner, and Irina Higgins. Life-long disentangled representation learning with cross-domain latent homologies. *NeurIPS*, 2018. 2, 6, 10

[2] Rahaf Aljundi, Francesca Babiloni, Mohamed Elhoseiny, Marcus Rohrbach, and Tinne Tuytelaars. Memory aware synapses: Learning what (not) to forget. In *ECCV*, 2018. 2, 10

[3] Adrien Bardes, Jean Ponce, and Yann LeCun. Vicreg: Variance-invariance-covariance regularization for self-supervised learning. *arXiv preprint arXiv:2105.04906*, 2021. 1, 2, 5, 9

[4] Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark experience for general continual learning: a strong, simple baseline. In *NeurIPS*, 2020. 2, 6, 10

[5] Lucas Caccia and Joelle Pineau. Special: Self-supervised pretraining for continual learning. *arXiv preprint arXiv:2106.09065*, 2021. 2

[6] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for unsupervised learning of visual features. In *ECCV*, pages 132–149, 2018. 2

[7] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *NeurIPS*, 2020. 1, 2, 4, 5, 11

[8] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. *ICCV*, 2021. 1, 2, 5, 11

[9] Francisco M Castro, Manuel J Marin-Jimenez, Nicolas Guil, Cordelia Schmid, and Karteek Alahari. End-to-end incremental learning. In *ECCV*, 2018. 2, 9

[10] Hyuntak Cha, Jaeho Lee, and Jinwoo Shin. Co2l: Contrastive continual learning. In *CVPR*, pages 9516–9525, 2021. 2

[11] Arslan Chaudhry, Puneet K Dokania, Thalaiyasingam Ajanthan, and Philip HS Torr. Riemannian walk for incremental learning: Understanding forgetting and intransigence. In *ECCV*, 2018. 2

[12] Arslan Chaudhry, Marc'Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient lifelong learning with a-gem. In *ICLR*, 2019. 2

[13] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, 2020. 1, 2, 4, 5, 10

[14] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020. 1, 2

[15] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *ICCV*, 2021. 2, 4, 5, 6, 12

[16] Victor Guilherme Turrisi da Costa, Enrico Fini, Moin Nabi, Nicu Sebe, and Elisa Ricci. solo-learn: A library of self-supervised methods for visual representation learning. *Journal of Machine Learning Research*, 23(56):1–6, 2022. 5, 6

[17] Matthias De Lange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Ales Leonardis, Gregory Slabaugh, and Tinne Tuytelaars. Continual learning: A comparative study on how to defy forgetting in classification tasks. *IEEE TPAMI*, 2(6), 2019. 1, 2, 6

[18] Arthur Douillard, Matthieu Cord, Charles Ollion, Thomas Robert, and Eduardo Valle. Podnet: Pooled outputs distillation for small-tasks incremental learning. In *ECCV*, 2020. 2, 3, 6, 9

[19] Debidatta Dwibedi, Yusuf Aytar, Jonathan Tompson, Pierre Sermanet, and Andrew Zisserman. With a little help from my friends: Nearest-neighbor contrastive learning of visual representations. *ICCV*, 2021. 2, 5

[20] Aleksandr Ermolov, Aliaksandr Siarohin, Enver Sangineto, and Nicu Sebe. Whitening for self-supervised representation learning. In *ICML*, 2021. 2, 5

[21] William Falcon, Ananya Harsh Jha, Teddy Koker, and Kyunghyun Cho. Aavae: Augmentation-augmented variational autoencoders. *arXiv preprint arXiv:2107.12329*, 2021. 10

[22] Enrico Fini, Stèphane Lathuilière, Enver Sangineto, Moin Nabi, and Elisa Ricci. Online continual learning under extreme memory constraints. In *ECCV*, 2020. 2, 9

[23] Robert M. French. Catastrophic forgetting in connectionist networks. *Trends in Cognitive Sciences*, 3(4):128–135, 1999. 1

[24] Jhair Gallardo, Tyler L Hayes, and Christopher Kanan. Self-supervised training enhances online continual learning. *arXiv preprint arXiv:2103.14010*, 2021. 2

[25] I. J. Goodfellow, M. Mirza, D. Xiao, A. Courville, and Y. Bengio. An empirical investigation of catastrophic forgetting in gradient-based neural networks. *arXiv preprint arXiv:1312.6211*, 2013. 1

[26] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent: A new approach to self-supervised learning. *NeurIPS*, 2020. 1, 2, 4, 5, 12

[27] Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *AISTATS*, 2010. 2

[28] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, 2020. 1, 2, 4, 5, 10

[29] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 5, 11

[30] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. 3

[31] Saihui Hou, Xinyu Pan, Chen Change Loy, Zilei Wang, and Dahua Lin. Learning a unified classifier incrementally via rebalancing. In *CVPR*, pages 831–839, 2019. 2, 3, 6, 9

[32] Ahmet Iscen, Jeffrey Zhang, Svetlana Lazebnik, and Cordelia Schmid. Memory-efficient incremental learning through feature adaptation. In *European Conference on Computer Vision*, pages 699–715. Springer, 2020. 9

[33] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *NeurIPS*, 2020. 10

[34] 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. *Proc. of the national academy of sciences*, 2017. 2, 6, 10

[35] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, Univ. Toronto, 2009. 5

[36] Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE TPAMI*, 2017. 2, 3, 9

[37] Zhiwei Lin, Yongtao Wang, and Hongxiang Lin. Continual contrastive self-supervised learning for image classification. *arXiv preprint arXiv:2107.01776*, 2021. 2, 6, 7, 10

[38] David Lopez-Paz and Marc-Aurelio Ranzato. Gradient episodic memory for continual learning. In *NeurIPS*, 2017. 2, 5

[39] Divyam Madaan, Jaehong Yoon, Yuanchun Li, Yunxin Liu, and Sung Ju Hwang. Rethinking the representational continuity: Towards unsupervised continual learning. *arXiv preprint arXiv:2110.06976*, 2021. 2, 6, 10

[40] Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pages 109–165. Elsevier, 1989. 1

[41] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 2

[42] Oleksiy Ostapenko, Mihai Puscas, Tassilo Klein, Patrick Jahnichen, and Moin Nabi. Learning to remember: A synaptic plasticity driven framework for continual learning. In *CVPR*, 2019. 2

[43] Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching for multi-source domain adaptation. In *ICCV*, 2019. 5

[44] Ameya Prabhu, Philip HS Torr, and Puneet K Dokania. Gdumb: A simple approach that questions our progress in continual learning. In *ECCV*, pages 524–540. Springer, 2020. 2

[45] Dushyant Rao, Francesco Visin, Andrei A Rusu, Yee Whye Teh, Razvan Pascanu, and Raia Hadsell. Continual unsupervised representation learning. *NeurIPS*, 2019. 2, 6, 10

[46] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental classifier and representation learning. In *CVPR*, 2017. 2

[47] Anthony Robins. Catastrophic forgetting, rehearsal and pseudorehearsal. *Connection Science*, 1995. 2, 6, 10

[48] A. A. Rusu, N. C. Rabinowitz, G. Desjardins, H. Soyer, J. Kirkpatrick, K. Kavukcuoglu, R. Pascanu, and R. Hadsell. Progressive neural networks. *arXiv 1606.04671*, 2016. 2

[49] Joan Serra, Didac Suris, Marius Miron, and Alexandros Karatzoglou. Overcoming catastrophic forgetting with hard attention to the task. *ICML*, 2018. 2

[50] Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual learning with deep generative replay. In *NeurIPS*, 2017. 2

[51] James Smith, Cameron Taylor, Seth Baer, and Constantine Dovrolis. Unsupervised progressive learning and the stam architecture. *arXiv preprint arXiv:1904.02021*, 2019. 2

[52] Yuandong Tian, Xinlei Chen, and Surya Ganguli. Understanding self-supervised learning dynamics without contrastive pairs. In *International Conference on Machine Learning*, pages 10268–10278. PMLR, 2021. 12

[53] Yonglong Tian, Olivier J Henaff, and Aaron van den Oord. Divide and contrast: Self-supervised learning from uncurated data. *arXiv preprint arXiv:2105.08054*, 2021. 2, 11

[54] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. In *ECCV*, pages 776–794. Springer, 2020. 5

[55] Yue Wu, Yinpeng Chen, Lijuan Wang, Yuancheng Ye, Zicheng Liu, Yandong Guo, and Yun Fu. Large scale incremental learning. In *CVPR*, 2019. 2

[56] Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *CVPR*, 2018. 2

[57] Yang You, Igor Gitman, and Boris Ginsburg. Large batch training of convolutional networks. *arXiv preprint arXiv:1708.03888*, 2017. 5

[58] Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised learning via redundancy reduction. *ICML*, 2021. 1, 2, 5, 9

[59] Friedeman Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. In *ICML*, 2017. 2, 10