# Active Class Incremental Learning for Imbalanced Datasets

Eden Belouadah[1,2][0000−0002−3418−1546], Adrian Popescu[1][0000−0002−8099−824X], Umang Aggarwal[1][0000−0002−3982−9284], and Lo Saci[1][0000−0002−6957−7823]

[1] CEA, LIST, F-91191 Gif-sur-Yvette, France
[2] IMT Atlantique, Computer Science Department, F-29238, Brest, France
{eden.belouadah,adrian.popescu,umang.aggarwal,leo.saci}@cea.fr

**Abstract.** Incremental Learning (IL) allows AI systems to adapt to streamed data. Most existing algorithms make two strong hypotheses which reduce the realism of the incremental scenario: (1) new data are assumed to be readily annotated when streamed and (2) tests are run with balanced datasets while most real-life datasets are imbalanced. These hypotheses are discarded and the resulting challenges are tackled with a combination of active and imbalanced learning. We introduce sample acquisition functions which tackle imbalance and are compatible with IL constraints. We also consider IL as an imbalanced learning problem instead of the established usage of knowledge distillation against catastrophic forgetting. Here, imbalance effects are reduced during inference through class prediction scaling. Evaluation is done with four visual datasets and compares existing and proposed sample acquisition functions. Results indicate that the proposed contributions have a positive effect and reduce the gap between active and standard IL performance.

**Keywords:** Incremental Learning, Active Learning, Imbalanced Learning, Computer Vision, Image Classification.

## 1 Introduction

AI systems are often deployed in dynamic settings where data are not all available at once [29]. Examples of applications include: (1) robotics - where the robot evolves in a changing environment and needs to adapt to it, (2) news analysis - where novel entities and events appear at a fast pace and should be processed swiftly, and (3) medical document processing - where parts of the data might not be available due to privacy constraints.

In such cases, incremental learning (IL) algorithms are needed to integrate new data while also preserving the knowledge learned for past data. Following a general trend in machine learning, recent IL algorithms are all built around Deep Neural Networks (DNNs) [2,3,10,17,24,33]. The main challenge faced by such algorithms is catastrophic interference or forgetting [25], a degradation of performance for previously learned information when a model is updated with new data.

IL algorithm design is an open research problem if computational complexity should remain bounded as new data are incorporated and/or if only a limited memory is available to store past data. These two conditions are difficult to satisfy simultaneously and existing approaches address one of them in priority. A first research direction allows for model complexity to grow as new data are added [2,3,24,35,41]. They focus on minimizing the number of parameters added for each incremental update and no memory of past data is allowed. Another research direction assumes that model complexity should be constant across incremental states and implements rehearsal over a bounded memory of past data to mitigate catastrophic forgetting [10,14,20,33,46]. Most existing IL algorithms assume that new data are readily labeled at the start of each incremental step. This assumption is a strong one since data labeling is a time consuming process, even with the availability of crowdsourcing platforms. Two notable exceptions are presented in [3] and [31] where the authors introduce algorithms for self-supervised face recognition. While interesting, these works are applicable only to a specific task and both exploit pretrained models to start the process. Also, a minimal degree of supervision is needed in order to associate a semantic meaning (i.e. person names) to the discovered identities. A second hypothesis made in incremental learning is that datasets are balanced or nearly so. In practice, imbalance occurs in a wide majority of real-life datasets but also in research datasets constructed in controlled conditions. Public datasets such as ImageNet [12], Open Images [21] or VGG-Face2 [9] are all imbalanced. However, most research works related to ImageNet report results with the ILSVRC subset [34] which is nearly balanced.

These two hypotheses limit the practical usability of existing IL algorithms. We replace them by two weaker assumptions to make the incremental learning scenario more realistic. First, full supervision of newly streamed data is replaced by the possibility to annotate only a small subset of these data. Second, no prior assumption is made regarding the balanced distribution of new data in classes. We combine active and imbalanced learning methods to tackle the challenges related to the resulting IL scenario.

The main contribution of this work is to adapt sample acquisition process, which is the core component of active learning (AL) methods, to incremental learning over potentially imbalanced datasets. A two phases procedure is devised to replace the classical acquisition process which uses a single acquisition function. A standard function is first applied to a subset of the active learning budget in order to learn an updated model which includes a suboptimal representation of new data. In the second phase, a balancing-driven acquisition function is used to favor samples which might be associated to minority classes (i.e. those having a low number of associated samples). The data distribution in classes is updated after each sample labeling to keep it up-to-date. Two balancing-driven acquisition functions which exploit the data distribution in the embedding space of the IL model are introduced here. The first consists of a modification of the core-set algorithm [37] to restrain the search for new samples to data points which are closer to minority classes than to majority ones. The second function

prioritizes samples which are close to the poorest minority classes (i.e. those represented by the minimum number of samples) and far from any of the majority classes. The balancing-driven acquisition phase is repeated several times and new samples are successively added to the training set in order to enable an iterative active learning process [38].

A secondary contribution is the introduction of a backbone training procedure which considers incremental learning with memory as an instance of imbalanced learning. The widely used training with knowledge distillation [10,17,20,33,46] is consequently replaced by a simpler procedure which aims to reduce the prediction bias towards majority classes during inference [8]. Following the conclusions of this last work, initial predictions are rectified by using the prior class probabilities from the training set.

Four public datasets designed for different visual tasks are used for evaluation. The proposed balancing-driven sample acquisition process is compared with a standard acquisition process and results indicate that it has a positive effect for imbalanced datasets.

## 2  Related Works

We discuss existing works from incremental, imbalanced and active learning areas and focus on those which are most closely related to our contribution. Class incremental learning witnessed a regain of interest and all recent methods exploit DNNs. One influential class of IL methods build on the adaptation of fine tuning and exploit increasingly sophisticated knowledge distillation techniques to counter catastrophic forgetting [25]. *Learning without Forgetting* ($LwF$) [23] introduced this trend and is an inspiration for a wide majority of further IL works. *Incremental Classifier and Representation Learning* ($iCaRL$) [33] is one such work which uses $LwF$ and also adds a bounded memory of the past to implement replay-based IL efficiently. $iCaRL$ selects past class exemplars using a herding approach. The classification layer of the neural nets is replaced by a nearest-exemplars-mean, which adapts nearest-class-mean [26], to counter class imbalance. *End-to-end incremental learning* [10] uses a distillation component which adheres to the original definition of $LwF$ from [16]. A balanced fine tuning step is added to counter imbalance. As a result, a consequent improvement over $iCaRL$ is reported. *Learning a Unified Classifier Incrementally via Rebalancing* ($LUCIR$) [17] tackles incremental learning problem by combining cosine normalization in the classification layer, a less-forget constraint based on distillation and an inter-class separation to improve comparability between past and new classes. *Class Incremental Learning with Dual Memory* ($IL2M$) [4] and *Bias Correction* ($BiC$) [42] are recent approaches that add an extra layer to the network in order to remove the prediction bias towards new classes which are represented by more images than past classes.

Classical active learning is thoroughly reviewed in [38]. A first group of approaches exploits informativeness to select items for uncertain regions in the classification space. Uncertainty is often estimated with measures such as en-

tropy [40], least confidence first [11] or min margin among top predictions [36]. Another group of approaches leverages sample representativeness computed in the geometric space defined by a feature extractor. Information density [39] was an early implementation of such an approach. *Core-set*, which rely on the classical *K-centers* algorithm to discover an optimal subset of the unlabeled dataset, was introduced in [37].

Recent active learning works build on the use of deep learning. The labeling effort is examined in [19] to progressively prune labels as labeling advances. An algorithm which learns a loss function specifically for AL was proposed in [44]. While very interesting, such an approach is difficult to exploit in incremental learning since the main challenge here is to counter data imbalance between new and past classes or among new classes. Another line of works proposes to exploit multiple network states to improve the AL process. *Monte Carlo Dropout* [13] uses softmax prediction from a model with random dropout masks. In [6], an ensemble approach which combines multiple snapshots of the same training process is introduced. These methods are not usable in our scenario because they increase the number of parameters due to the use of multiple models. We retain the use of the same deep model through incremental states to provide embeddings and propose a stronger role for them during the sample acquisition process. Recently, [1] proposed a method which focuses on single-stage AL for imbalanced datasets. They exploit a pretrained feature extractor and annotate the unlabeled samples so as to favor minority classes.

Ideally, incremental updates should be done in a fully unsupervised manner [43] in order to remove the need for manual labeling. However, unsupervised algorithms are not mature enough to capture dataset semantics with the same degree of refinement and performance as their supervised or semi-supervised counterparts. Closest to our work are the self-supervision approaches designed for incremental face recognition [3,31]. They are tightly related to unsupervised learning since no manual labeling is needed, except for naming the person. Compared to self-supervision, our approach requires manual labeling for a part of new data and has a higher cost. However, it can be applied to any class IL problem and not only to specific tasks such as face recognition as it is the case for [3,31].

A comprehensive review of imbalanced object-detection problems is provided in [27]. The authors group these problems in a taxonomy depending on their class imbalance, scale imbalance, spatial imbalance or objective imbalance. The study shows the increasing interest of the computer vision community in the imbalanced problems for their usefulness in real life situations.

## 3   Proposed Method

The proposed active learning adaptation to an incremental scenario is motivated by the following observations:

- Existing acquisition functions ($\mathcal{AF}$s) were designed and tested successfully for active learning over balanced datasets. However, a wide majority of real-life datasets are actually imbalanced. Here, no prior assumption is made
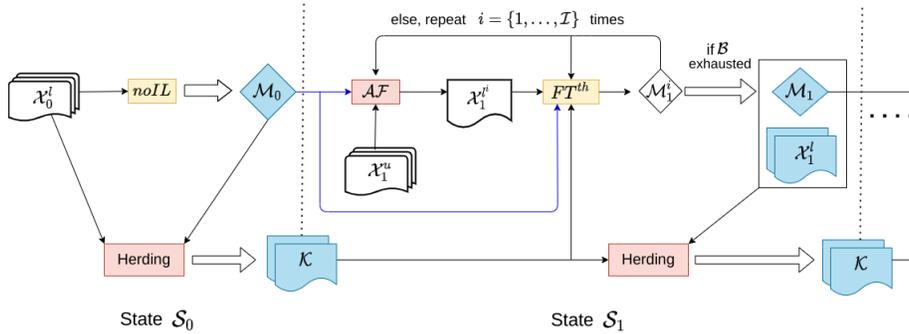
**Fig. 1.** Illustration of the proposed training process with one initial state $\mathcal{S}_0$, and one incremental state $\mathcal{S}_1$. The initial deep model $\mathcal{M}_0$ is trained from scratch on a fully-labeled dataset $\mathcal{X}_0^l$ using $noIL$ (a non-incremental learning). $\mathcal{M}_0$ and $\mathcal{X}_0^l$ are used to prepare the past class memory $\mathcal{K}$ using herding (a mechanism that selects the best representative past class images). State $\mathcal{S}_1$ starts with a sample acquisition function $\mathcal{AF}$ that takes the unlabeled set $\mathcal{X}_1^u$ and the model $\mathcal{M}_0$ as inputs, and provides a part of the budget $\mathcal{B}$ annotated as $\mathcal{X}_1^{l^i}$. The model $\mathcal{M}_0$ is then updated with data from $\mathcal{X}_1^{l^i} \cup \mathcal{K}$ using $FT^{th}$ (a fine tuning followed by a threshold calibration). The updated model $\mathcal{M}_1^i$ is again fed into the acquisition function $\mathcal{AF}$ with the rest of unlabeled examples from $\mathcal{X}_1^u$ to further annotate a part of the budget $\mathcal{B}$ and the model is updated afterwards. This process is repeated $\mathcal{I}$ times until $\mathcal{B}$ is exhausted. The model $\mathcal{M}_1$ is then returned with the annotated dataset $\mathcal{X}_1^l$ and the memory $\mathcal{K}$ is updated by inserting exemplars of new classes from $\mathcal{X}_1^l$ and reducing exemplars of past classes in order to fit its maximum size. Note that the two blue arrows are applicable in the first AL iteration only (when $i = 1$). Best viewed in color.

regarding the imbalanced or balanced character of the unlabeled data which is streamed in IL states. Unlike existing sample acquisition approaches which exploit a single $\mathcal{AF}$, we propose to split the process in two phases. The first phase uses a classical $\mathcal{AF}$ to kick-off the process. The second one implements an $\mathcal{AF}$ which is explicitly designed to target a balanced representation of labeled samples among classes.

– In IL, a single deep model can be stored throughout the process. This makes the application of recent ensemble methods [6] inapplicable. Following the usual AL pipeline, an iterative fine tuning of the model is implemented to incorporate labeled samples from the latest AL iteration.

– A memory $\mathcal{K}$ of past class samples is allowed and, following [4,18], we model IL as an instance of imbalanced learning. The distillation component, which is central to most existing class IL algorithms [10,20,33,42], is removed. Instead, we deal with imbalance by using a simple but efficient post-processing step which modifies class predictions based on their prior probabilities in

the training set. The choice of this method is motivated by its superiority in deep imbalanced learning over a large array of other methods [8].

An illustration of the proposed learning process is provided in Fig. 1. In the next sections, we first formalize the proposed active incremental learning scenario. Then, we introduce the adapted sample acquisition process, with focus on the balancing-driven acquisition functions. Finally, we present the incremental learning backbone which is inspired from imbalanced learning [8].

### 3.1 Problem Formalization

The formalization of the problem is inspired by [10,33] for the incremental learning part and by [6] for the active learning part. We note $\mathcal{T}$ the number of states (including the first non-incremental state), $\mathcal{K}$ - the bounded memory for past classes, $\mathcal{B}$ - the labeling budget available for active learning, $\mathcal{AF}$ - an acquisition function designed to optimize sample selection in active learning, $\mathcal{I}$ the number of iterations done during active learning, $\mathcal{S}_t$ - an incremental state, $N_t$ - the total number of classes recognizable in $\mathcal{S}_t$, $\mathcal{X}_t^u$ - the unlabeled dataset associated to $\mathcal{S}_t$, $\mathcal{X}_t^l$ - a manually labeled subset of $\mathcal{X}_t^u$, $\mathcal{M}_t$ - the deep model learned in $\mathcal{S}_t$. The initial state $\mathcal{S}_0$ includes a dataset $\mathcal{X}_0 = \{X_0^1, X_0^2, ..., X_0^j, ..., X_0^{P_0}\}$ with $N_0 = P_0$ classes. $X_t^j = \{x_1^j, x_2^j, ..., x_{n_j}^j\}$ is the set of $n_j$ training examples for the $j^{th}$ class, $p_t^j$ is its corresponding classification probability in the state $\mathcal{S}_t$.

We assume that all the samples are labeled in the first state. An initial non-incremental model $\mathcal{M}_0 : \mathcal{X}_0 \to \mathcal{C}_0$ is trained to recognize a set $\mathcal{C}_0$ containing $N_0$ classes using all their data from $\mathcal{X}_0$. $P_t$ new classes need to be integrated in each incremental state $\mathcal{S}_t$, with $t > 0$. Each IL step updates the previous model $\mathcal{M}_{t-1}$ into the current model $\mathcal{M}_t$ which recognizes $N_t = P_0 + P_1 + ... + P_t$ classes in the incremental state $\mathcal{S}_t$. Active learning is deployed using $\mathcal{AF}(\mathcal{X}_t^u)$ to obtain $\mathcal{X}_t^l$, a labeled subset from $\mathcal{X}_t^u$.

$\mathcal{X}_t^l$ data of the $P_t$ new classes are available but only a bounded exemplar memory $\mathcal{K}$ for the $N_{t-1}$ past classes is allowed. $\mathcal{M}_t$, the model associated to the state $\mathcal{S}_t$ is trained over the $\mathcal{X}_t^l \cup \mathcal{K}$ training dataset. An iterative AL process is implemented to recognize a set of classes $\mathcal{C}_t = \{c_t^1, c_t^2, ..., c_t^{N_{t-1}}, c_t^{N_{t-1}+1}, ..., c_t^{N_t}\}$

### 3.2 Active Learning in an Incremental Setting

We discuss the two phases of the adapted active learning process below. Classical sampling is followed by a phase which exploits the proposed balancing-driven acquisition functions.

**Classical Sample Acquisition Phase.** At the start of each IL state $\mathcal{S}_t$, an unlabeled dataset $\mathcal{X}_t^u$ is streamed into the system and classical AL acquisition functions are deployed to label $\mathcal{X}_t^l$, a part of $\mathcal{X}_t^u$, for inclusion in the training set. Due to IL constraints, the only model available at the beginning of $\mathcal{S}_t$ is $\mathcal{M}_{t-1}$, learned for past classes in the previous incremental step. It is used to extract

the embeddings needed to implement acquisition functions. A number of acquisition functions were proposed to optimize the active learning process [38], with adaptations for deep learning in [6,13,37,47]. Based on their strong experimental performance [6,36,37,38], four $\mathcal{AF}$s are selected for the initial phase:

- **core-set sampling** [37] (*core* hereafter): whose objective is to extract a representative subset of the unlabeled dataset from the vectorial space defined by the deep embeddings. The method selects samples with:

$$x_{next} = \underset{x_u \in \mathcal{X}_t^u}{\operatorname{argmax}} \{ \min_{1 \le k \le n} \Delta(e(x_u), e(x_k)) \} \tag{1}$$

  where: $x_{next}$ is the next sample to label, $x_u$ is an unlabeled sample left, $x_k$ is one of the $n$ samples which were already labeled, $e()$ is the embedding extracted using $\mathcal{M}_{t-1}$ and $\Delta$ is the Euclidean distance between two embeddings.
- **random sampling** (*rand* hereafter) : a random selection of images for labeling. While basic, random selection remains a competitive baseline in active learning.
- **entropy sampling** [38] (*ent* hereafter): whose objective is to favor most uncertain samples as defined by the set of probabilities given by the model.

$$x_{next} = \underset{x_u \in \mathcal{X}_t^u}{\operatorname{argmax}} \{ -\sum_{j=1}^{J} (p_t^j * \log(p_t^j)) \} \tag{2}$$

  where $p_t^j$ is the prediction score of $x_u$ for the class $j$ and $J$ is the number of detected classes so far by AL.
- **margin sampling** [36] (*marg* hereafter): selects the most uncertain samples based on their top-2 predictions of the model.

$$x_{next} = \underset{x_u \in \mathcal{X}_t^u}{\operatorname{argmax}} \{ max(p_t^1, .., p_t^j, .., p_t^J) - max_2(p_t^1, .., p_t^j, .., p_t^J) \} \tag{3}$$

  where $max(\cdot)$ and $max_2(\cdot)$ provide the top-2 predicted probabilities for the sample $x_u$. This $\mathcal{AF}$ favors samples that maximize the difference between their top two predictions.

This acquisition phase is launched once at the beginning of each incremental state to get an initial labeled subset of the new data. This step is necessary to include the samples for the new classes in the trained model and initiate the iterative AL process.

**Balancing-driven Sample Acquisition Phase.** The second acquisition phase tries to label samples so as to tend toward a balanced distribution among new classes. The distribution of the number of samples per class is computed after each sample labeling to be kept up-to-date. The average number of samples per class is used to divide classes into minority and majority ones. These two sets of classes are noted $\mathcal{C}_t^{mnr}$ and $\mathcal{C}_t^{maj}$ for incremental state $\mathcal{S}_t$. Two functions are proposed to implement the balancing-driven acquisition:

- **balanced core-set sampling** ($b - core$ hereafter) is a modified version of
  $core$ presented in Equation 1. $b - core$ acts as a filter which keeps candidate
  samples for labeling only if they are closer to a minority class than to any
  majority class. We write the relative distance of an unlabeled image w.r.t.
  its closest minority and majority classes as:

$$\Delta_{\frac{mnr}{maj}}(x_u) = \min_{c_t^{mnr} \in \mathcal{C}_t^{mnr}} \Delta(e(x_u), \mu(c_t^{mnr})) - \min_{c_t^{maj} \in \mathcal{C}_t^{maj}} \Delta(e(x_u), \mu(c_t^{maj})) \quad (4)$$

where: $x_u$ is an unlabeled sample, $c_t^{mnr}$ and $c_t^{maj}$ are classes from the minor-
ity and majority sets $\mathcal{C}_t^{mnr}$ and $\mathcal{C}_t^{maj}$ respectively, $e(x_u)$ is the embedding of
$x_u$ extracted from the latest deep model available, $\mu(c_t^{mnr})$ and $\mu(c_t^{maj})$ are
the centroids of minority and majority classes $c_t^{mnr}$ and $c_t^{maj}$ computed over
the embeddings of their labeled samples.
The next sample to label is chosen by using the core-set definition from
Equation 1 but after filtering remaining unlabeled samples with Equation 4:

$$x_{next} = \underset{x_u \in \mathcal{X}_t^u \ and \ \Delta_{\frac{mnr}{maj}}(x_u) < 0}{\operatorname{argmax}} \{ \min_{1 \leq k \leq n} \Delta(e(x_u), e(x_k)) \} \quad (5)$$

- **poorest class first sampling** ($poor$) is an acquisition function which gives
  priority to the class represented by the minimum number of labeled samples
  associated to it at a given moment during active learning. If there are sev-
  eral such classes, one of them is selected randomly. The method translates
  the hypothesis that samples which are close to a poor class and far from
  any majority class should be favored in order to achieve a more balanced
  distribution. The next candidate for labeling is selected with:

$$x_{next} = \underset{x_u \in \mathcal{X}_t^u}{\operatorname{argmin}} \{ \Delta(e(x_u), \mu(c_t^{poor})) - \min_{\forall c_t^{maj} \in \mathcal{C}_t^{maj}} \Delta(e(x_u), \mu(c_t^{maj})) \} \quad (6)$$

where $c_t^{poor}$ is a minority class from $\mathcal{C}_t^{mnr}$ which has the lowest number of
samples in the current labeled subset.
$poor$ is similar in spirit to $b - core$ but has a stronger drive towards balancing
because an individual class with poorest representation is targeted instead
of samples which are close to any minority class.

In an iterative active learning scenario, the balancing-driven acquisition can
be repeated several time until the AL budget $\mathcal{B}$ is exhausted.

### 3.3   Imbalance-driven Incremental Learning

The model update within each incremental state is inspired by a usual iterative
AL approach [38] which includes a classical acquisition phase at the beginning
and several balancing-driven iterations. For a total of $\mathcal{I}$ active learning iterations
in each state $\mathcal{S}_t$, intermediate models $\mathcal{M}_t^1$, ..., $\mathcal{M}_t^i$, ..., $\mathcal{M}_t^{\mathcal{I}-1}$ are created while
annotating $\mathcal{X}_t^{l^1}$, ..., $\mathcal{X}_t^{l^i}$, ..., $\mathcal{X}_t^{l^{\mathcal{I}-1}}$ during the first $\mathcal{I}-1$ iterations before obtaining
the final $\mathcal{M}_t$. The number of iterations $\mathcal{I}$ and the size of each iteration can take

different values. The choice of a particular setting is done empirically so as to: (1) have enough new samples in the initial iteration in order for the new classes to be trainable in $\mathcal{M}_t^1$, i.e. the model $\mathcal{M}_t$ in the first iteration, (2) have enough candidates left for the balancing-driven iterations and (3) do not repeat the fine tuning process too many times to keep the incremental update timely. $\mathcal{M}_{t-1}$ is used to extract embeddings if *core* is used in the initial AL iteration. Note that while iterative training increases the level of forgetting in IL [4,5], it is needed in AL to update model representation while annotating the images [38].

As we mentioned, we depart from the usual modeling of the IL problem [10,18,33,42] which exploits knowledge distillation to counter catastrophic forgetting. Following the recent observation that a simpler fine tuning based approach gives interesting results [4], we use an IL backbone inspired from imbalance learning results presented in [8]. This backbone is called fine tuning with thresholding ($FT^{th}$ below), also known as threshold moving or post scaling [8]. Thresholding adjusts the decision threshold of the model. It consists in the addition of a calibration layer at the end of the model during inference to compensate the prediction bias in favor of majority classes. This layer rectifies the class prediction $p_t^j$ of the $j^{th}$ class in the state $\mathcal{S}_t$ as follows:

$$p_t^{j'} = p_t^j \times \frac{|\mathcal{X}_t^l \cup \mathcal{K}|}{|X_t^j|} \tag{7}$$

where $|X_t^j|$ is the number of training examples for the $j^{th}$ class in the state $\mathcal{S}_t$ and $|\mathcal{X}_t^l \cup \mathcal{K}|$ is the total number of training examples in state $\mathcal{S}_t$.

$FT^{th}$ boosts the scores of classes with a lower number of associated samples. The method has the interesting property of dealing with imbalance in IL in a uniform manner. It does not matter whether imbalance comes from the distribution of newly streamed data or from the fact that only a bounded memory of past classes is available. This stands in contrast with knowledge distillation which handles imbalance for past classes but not among new ones. $FT^{th}$ is competitive against state-of-the-art algorithms. In a classical (i.e. fully supervised) IL setting, it has 59.59 top-1 accuracy for Cifar100, compared to $iCaRL$ [33] (57.35) and $LUCIR$ [18] (55.36). More results are provided in the next section.

## 4   Experiments

### 4.1   Datasets

Experiments are run with four public datasets, out of which three are imbalanced and one is balanced. We provide a brief description of the datasets below:

- **ImageNet100** - dataset for fine grained object recognition consisting of a subset of 100 randomly selected leaf classes from ImageNet [12] which have at least 50 training images and are not present in the ILSVRC subset [34].
- **Faces100** - face recognition dataset consisting of a subset of randomly selected 100 identities from VGG-Face2 [9] with at least 30 training images.

  – **Food101** - dataset for fine-grained food recognition [7]. Since the initial
    dataset is perfectly balanced, an imbalance induction procedure was applied
    by removing a variable number of training samples keeping at least 25 images
    per class.
  – **Cifar100** - dataset for object recognition used in its original version [22]
    which is perfectly balanced.

The main statistics of the experimental datasets are provided in Table 1. We
provide the coefficient of variation $cv = \frac{\sigma}{\mu}$, with $\sigma$ the standard deviation and $\mu$
the mean of the distribution of samples per class. $cv$ provides information about
the degree of imbalance associated to each dataset. The larger this value is, the
more imbalanced the dataset will be.

| Dataset | Classes | Train | Test | Mean train ($\mu$) | Std train ($\sigma$) | $cv$ |
|---|---|---|---|---|---|---|
| ImageNet100 | 100 | 50000 | 5K | 500.0 | 376.17 | 0.7523 |
| Faces100 | 100 | 23237 | 5K | 232.37 | 167.68 | 0.7216 |
| Food101 | 101 | 22374 | 10K | 223.74 | 177.66 | 0.7940 |
| Cifar100 | 100 | 50000 | 10K | 500.0 | 0.0 | 0.0 |

**Table 1.** Dataset statistics. $cv$ is the coefficient of variation defined as $cv = \frac{\sigma}{\mu}$.

### 4.2   Methodology

**Incremental Learning Setting.** We run the experiments with $\mathcal{T} = 10$ states
for each dataset[3]. This setting is classically used in class incremental learn-
ing [10,33]. A total of $\mathcal{K}$ images of past classes are kept at any time during
incremental learning. $\mathcal{K}$ approximates 2% of the full training sets. Memory sizes
are thus $\mathcal{K} = 1000$ for ImageNet100 and Cifar100, $\mathcal{K} = 465$ for Faces100 and
$\mathcal{K} = 450$ for Food101. At the end of each incremental state, memory is updated
by inserting exemplars of new classes and reducing exemplars of past classes in
order to fit its maximum size. Note that since $\mathcal{K}$ is constant and the number
of past classes grows, the imbalance in favor of new classes grows for later in-
cremental states and the problem becomes more challenging. The exemplars are
chosen using the herding mechanism introduced in [33]. The herding procedure
consists in choosing the set of images that approximates the best the real mean
of the class.

**Active Learning Process.** Three active learning budgets are tested covering
$\mathcal{B} = \{20\%, 10\%, 5\%\}$ of the unlabeled dataset $\mathcal{X}_t^u$ streamed in state $\mathcal{S}_t$. These
different values are used to get a comprehensive view of each configuration's
behavior. Active learning is implemented with a usual iterative approach [37,38]

---

[3] The initial non-incremental state of Food101 includes 11 classes while the initial
states for the other datasets include 10 classes each.

including $\mathcal{I} = 4$ iterations, 40% of $\mathcal{B}$ are used for classical acquisition and three times 20% of $\mathcal{B}$ for balancing-driven acquisition (values were experimentally chosen). Classical and balancing-driven acquisition phases are independent of one another and we test all their combinations. For completeness, we include results with a baseline in which both phases are implemented with random sampling. Note that the proposed acquisition functions are non-deterministic and experiments are run five times for each configuration in order to have a robust estimation of its performance. To improve comparability of configurations which use the same initial $\mathcal{AF}$, the same initial models are used for all subsequent balancing-driven $\mathcal{AF}$s.

**Training Details.** The experimental setup is inspired by the one proposed in $iCaRL$ [33]. $FT^{th}$ is implemented in Pytorch [30] using a ResNet-18 architecture [15] and an SGD optimizer. The first non-incremental state is run for 100 epochs with $batch\ size = 128$, $lr = 0.1$, $momentum = 0.9$, $weight\ decay = 0.0005$. The learning rate is divided by 10 when the error plateaus for 15 consecutive epochs. Fine tuning is run for 80 epochs, 20 epochs for each active learning iteration with $batch\ size = 32$, $lr = 0.1$, $momentum = 0.9$, $weight\ decay = 0.0005$. The learning rate is initialized at the beginning of the AL process and then divided by 10 when the error plateaus for 10 consecutive epochs.

Training images are preprocessed using randomly resized $224 \times 224$ crops and horizontal flipping and are normalized afterwards. While more advanced data augmentation is known to slightly improve performance [10], we did not apply other image transformations. For Faces100, face cropping is done with MTCNN [45] before further processing.

**Upper Bound Methods.** In addition to the active learning configurations, we present results with:

- $sIL$ - usual supervised incremental learning in which all samples are labeled (equivalent to $\mathcal{B} = 100\%$).
- $noIL$ - classical non-incremental learning in which all samples are provided at once.

For comparability, $sIL$ and $noIL$ are both trained using threshold calibration. $sIL$ is an incremental upper bound for active learning configurations since it is fully supervised. $noIL$ is an upper bound for $sIL$ since all the data are labeled and available at once. These upper bounds are useful insofar they provide information about the performance gap due to a partial labeling of streamed data.

### 4.3   Results and Discussion

**$FT^{th}$ in supervised mode** - Instead of handling catastrophic forgetting [25] as previous works did [10,17,32,42], we address IL with bounded past memory as an

imbalanced learning problem. We use threshold calibration [8] to rectify scores in order to give more chances to minority classes to be selected during inference. The comparison to recent IL methods in supervised mode from Table 2 indicates that $FT^{th}$ is competitive. It clearly outperforms $iCaRL$ [33] and $IL2M$ [4] and is better than LUCIR [17] for three datasets out of four. We also provide the results of vanilla fine tuning before threshold calibration to underline the usefulness of thresholding. It has a positive effect for all four datasets, a finding which validates its usefulness in our scenario.

| Dataset | $FT$ | $FT^{th}$ [8] | $LUCIR$ [18] | $iCaRL$ [33] | $IL2M$ [4] |
|---|---|---|---|---|---|
| Imagenet100 | 54.80 | **61.42** | 60.77 | 52.40 | 57.68 |
| Faces100 | 69.11 | 73.26 | **78.44** | 60.48 | 70.33 |
| Food101 | 30.21 | **34.79** | 25.70 | 21.99 | 32.20 |
| Cifar100 | 50.98 | **59.59** | 55.36 | 57.35 | 54.24 |

**Table 2.** Top-1 average supervised IL accuracy (%). Best results are in bold.

**Active Learning** - The experimental results obtained with $FT^{th}$ for the proposed active incremental learning scenario are presented in Table 3. The comparison of classical $\mathcal{AF}$s ($rand$ -$rand$ and $core - core$ in Table 3) indicates that random sampling outperforms the $core - set$ sampling in a majority of cases. This result is at odds with the one reported in [37] but is in line with the findings of [6,13] that random sampling in AL is a strong baseline and is actually better than the recent core-set method from [35]. The authors of this last paper also report that random sampling has better performance for lower active learning budgets which are studied here. Consequently, improving over random sampling for imbalanced datasets is an interesting result.

The results from Table 3 indicate that the balancing-driven acquisition phase is useful for all three imbalanced datasets and active learning budgets tested. The gains for ImageNet100 and Faces100 are usually between 1 and 2 points compared to the classical acquisition processes implemented here ($rand$ - $rand$ or $core$ - $core$). The gains are low for Food101, the third imbalanced dataset tested. This is probably due to the fact that Food101 is a more difficult task, as shown by $sIL$. More labeled samples per class would probably be needed for an efficient training.

$poor$ strategy is better than $b - core$ for ImageNet100 while more mixed results are obtained for Faces100 and Food101 datasets. Interestingly, the best results are always obtained on top of a $rand$ initial sampling, even when $core$-$core$ baseline is better than a $rand$-$rand$ one, as it is the case for Faces100 with $\mathcal{B} = 20\%$ and $\mathcal{B} = 5\%$.

When applied without balancing, $ent$ and $marg$ have poorer performance compared to that of $rand$ and $core$. Balancing significantly improves results

| Dataset | $\mathcal{B}$ | rand | | | core | | | ent | | | marg | | | sIL | noIL |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | rand | poor | b−core | core | poor | b−core | ent | poor | b−core | marg | poor | b−core | | |
| ImageNet100 | 20% | 57.48 ±0.50 | **58.65** ±0.23 | 58.08 ±0.40 | 56.85 ±0.23 | 57.25 ±0.84 | 57.46 ±0.52 | 45.07 ±0.58 | 56.53 ±0.27 | 56.23 ±0.22 | 54.13 ±0.66 | 56.39 ±0.53 | 56.26 ±0.45 | 61.42 | 72.48 |
| | 10% | 52.61 ±0.45 | **54.89** ±0.53 | 53.40 ±0.26 | 52.09 ±0.41 | 53.55 ±1.21 | 52.22 ±1.13 | 42.15 ±0.43 | 51.91 ±0.51 | 51.81 ±0.76 | 46.26 ±1.24 | 51.40 ±0.89 | 50.61 ±0.36 | | |
| | 5% | 47.72 ±0.69 | **48.71** ±0.97 | 48.18 ±0.56 | 46.01 ±0.51 | 47.39 ±0.85 | 46.45 ±0.30 | 37.95 ±0.61 | 45.10 ±1.46 | 44.70 ±1.02 | 36.74 ±1.05 | 44.18 ±1.09 | 43.93 ±0.93 | | |
| Faces100 | 20% | 65.91 ±0.94 | 66.41 ±0.10 | **67.24** ±0.36 | 66.41 ±0.66 | 66.46 ±0.46 | 66.94 ±1.37 | 48.62 ±0.95 | 63.30 ±0.33 | 64.99 ±0.80 | 59.51 ±1.17 | 62.85 ±1.75 | 65.27 ±0.53 | 73.26 | 93.62 |
| | 10% | 58.40 ±0.71 | **59.13** ±1.66 | 58.92 ±1.05 | 55.82 ±4.70 | 58.76 ±2.93 | 57.26 ±2.90 | 42.12 ±1.38 | 54.93 ±1.07 | 55.45 ±1.47 | 49.32 ±1.40 | 54.82 ±2.19 | 58.69 ±1.52 | | |
| | 5% | 48.38 ±1.27 | 50.09 ±2.30 | **50.12** ±0.96 | 48.74 ±1.21 | 47.71 ±1.28 | 50.04 ±2.04 | 35.61 ±0.51 | 45.79 ±0.83 | 45.39 ±1.13 | 38.37 ±1.02 | 45.90 ±2.31 | 45.64 ±1.56 | | |
| Food101 | 20% | 28.67 ±0.42 | **28.89** ±0.43 | 28.60 ±0.52 | 28.24 ±0.42 | 27.88 ±0.34 | 27.98 ±0.46 | 23.72 ±1.00 | 27.99 ±0.41 | 27.51 ±0.54 | 28.13 ±0.59 | 28.18 ±0.35 | 27.56 ±0.24 | 34.79 | 62.53 |
| | 10% | 24.12 ±0.47 | **24.17** ±0.56 | 24.07 ±0.68 | 22.91 ±0.63 | 23.46 ±0.12 | 23.07 ±0.31 | 19.41 ±0.96 | 22.32 ±0.73 | 22.25 ±0.64 | 23.35 ±0.64 | 22.68 ±0.81 | 22.84 ±0.52 | | |
| | 5% | 20.51 ±0.61 | 19.10 ±0.68 | **20.63** ±0.46 | 19.22 ±0.36 | 19.17 ±0.58 | 18.79 ±0.64 | 16.80 ±0.75 | 18.66 ±0.31 | 18.57 ±0.47 | 18.62 ±0.48 | 18.79 ±0.75 | 18.41 ±0.82 | | |
| Cifar100 | 20% | **49.47** ±0.16 | 49.36 ±0.33 | 48.46 ±0.75 | 46.75 ±0.40 | 46.87 ±0.19 | 46.87 ±0.36 | 39.76 ±1.30 | 46.66 ±0.29 | 47.69 ±0.23 | 46.07 ±0.31 | 45.37 ±0.39 | 46.68 ±0.44 | 59.59 | 76.98 |
| | 10% | **45.49** ±0.61 | 45.23 ±1.17 | 44.83 ±0.29 | 41.76 ±0.54 | 42.60 ±0.77 | 42.04 ±0.77 | 34.87 ±0.66 | 42.64 ±0.50 | 43.76 ±0.55 | 39.92 ±0.43 | 40.94 ±0.31 | 41.82 ±0.35 | | |
| | 5% | **41.58** ±0.29 | 40.69 ±0.23 | 39.69 ±0.46 | 35.23 ±0.64 | 37.70 ±0.46 | 35.72 ±0.67 | 31.74 ±0.74 | 37.68 ±0.34 | 38.02 ±0.66 | 31.88 ±0.58 | 35.96 ±0.42 | 35.69 ±0.64 | | |

**Table 3.** Top-1 average accuracy (%). Following [10], accuracy is averaged only for incremental states (i.e. excluding the initial, non-incremental state). Results are averaged over 5 runs for all AL configurations. $sIL$ is the result obtained in a fully supervised IL scenario. $noIL$ is the non-incremental upper-bound performance obtained with all data available. *Best results for each active learning configuration (row) are in bold.*

for both of uncertainty-based methods, but their overall performance still lags behind that of random followed by balancing. This reinforces the finding that a random selection is a competitive acquisition function in our active incremental learning over imbalanced datasets scenario.

The performance drop between active learning configurations and fully supervised IL naturally grows as $\mathcal{B}$ is reduced. The drop between $sIL$ and the best AL configuration is of 3, 6 and 5 points for $\mathcal{B} = 20\%$ for ImageNet100, Faces100, and Food101 respectively. When the AL budget is reduced to only 5% of new data, the corresponding performance losses go to 12.5, 23 and 14 points. Even when as little as 5% of the new data are annotated, suboptimal models are trainable and usable if the IL system needs to be operational quickly.

While the focus is on imbalanced datasets, we also report results with Cifar100, a perfectly balanced dataset for completeness. In this case, the balancing-driven sampling has a slightly negative effect when applied over $rand$ and a slightly positive effect over $core$. It is however notable that $core$ lags consistently behind $rand$ for Cifar100. The best strategy for all $\mathcal{B}$ sizes is $rand$-$rand$, with $rand$ - $poor$ being a close second best configuration.

The gap between active IL and supervised IL is still notable, especially for smaller AL budgets. In practice, active IL is useful when the system needs to be operational quickly after new data are streamed but at the expense of sub-optimal performance. If a longer delay is permitted, it is naturally preferable to annotate all new data before updating the incremental model. The gap is even higher between incremental and classical learning, even though $FT^{th}$ has competitive performance compared to existing IL algorithms. Globally, our results provide further confirmation that the use of incremental learning vs. classical learning should be weighted depending on the time, memory and/or computation constraints associated to an AI system's operation.

## 5   Conclusion

We proposed a more realistic incremental learning scenario which does not assume that streamed data are readily annotated and that they are evenly distributed among classes. An adaptation of the active learning sampling process is proposed in order to obtain a more balanced labeled subset. This adaptation has a positive effect for imbalanced datasets and a slightly negative effect for the balanced dataset evaluated here. Both proposed acquisition functions improve results compared to a classical acquisition process. Also interesting, experiments show that the random baseline outperforms the $core - set$ function. The strong performance of random sampling indicates that this method should be consistently used as a baseline for future works in active incremental learning. As a secondary contribution, we introduce $FT^{th}$, a IL backbone which provides competitive results when compared to state-of-the-art methods. The code is publicly available to facilitate reproducibility[4].

The proposed approach brings the IL scenario closer to practical needs. It reduces the time needed for an IL system to become operational upon receiving new data. The obtained results are encouraging but further investigation is needed to reduce the gap between active and supervised IL. Future work aims to: (1) run experiments with semi-supervised learning methods to automatically expand the labeled dataset and improve overall performance. While appealing, not all semi-supervised methods prove efficient in practice [28] and their usefulness for imbalanced datasets needs to be studied. (2) complement the proposed balancing-driven acquisition functions with a component which pushes the sampling process towards a better coverage of the manifold of each modeled class. This could be done, for instance, by taking inspiration from the herding mechanism [33] already used to select past exemplars. (3) render the IL scenario even more realistic by testing incremental steps of variable size to account for the fact that data might arrive at variable pace and considering that newly streamed data might belong both to unseen and past classes.

---

[4] `https://github.com/EdenBelouadah/class-incremental-learning/`

# References

1. Aggarwal, U., Popescu, A., Hudelot, C.: Active learning for imbalanced datasets. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV) (March 2020)
2. Aljundi, R., Chakravarty, P., Tuytelaars, T.: Expert gate: Lifelong learning with a network of experts. In: Conference on Computer Vision and Pattern Recognition. CVPR (2017)
3. Aljundi, R., Kelchtermans, K., Tuytelaars, T.: Task-free continual learning. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019)
4. Belouadah, E., Popescu, A.: Il2m: Class incremental learning with dual memory. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 583–592 (2019)
5. Belouadah, E., Popescu, A.: Scail: Classifier weights scaling for class incremental learning. In: IEEE Winter Conference on Applications of Computer Vision, WACV 2020, Snowmass Village, CO, USA, March 2-5, 2020 (2020)
6. Beluch, W.H., Genewein, T., Nürnberger, A., Köhler, J.M.: The power of ensembles for active learning in image classification. In: 2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018. pp. 9368–9377 (2018)
7. Bossard, L., Guillaumin, M., Van Gool, L.: Food-101 – mining discriminative components with random forests. In: European Conference on Computer Vision (2014)
8. Buda, M., Maki, A., Mazurowski, M.A.: A systematic study of the class imbalance problem in convolutional neural networks. Neural Networks **106**, 249–259 (2018)
9. Cao, Q., Shen, L., Xie, W., Parkhi, O.M., Zisserman, A.: Vggface2: A dataset for recognising faces across pose and age. In: 13th IEEE International Conference on Automatic Face & Gesture Recognition, FG 2018, Xi'an, China, May 15-19, 2018. pp. 67–74 (2018)
10. Castro, F.M., Marín-Jiménez, M.J., Guil, N., Schmid, C., Alahari, K.: End-to-end incremental learning. In: Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part XII. pp. 241–257 (2018)
11. Culotta, A., McCallum, A.: Reducing labeling effort for structured prediction tasks. In: Proceedings, The Twentieth National Conference on Artificial Intelligence and the Seventeenth Innovative Applications of Artificial Intelligence Conference, July 9-13, 2005, Pittsburgh, Pennsylvania, USA. pp. 746–751 (2005)
12. Deng, J., Dong, W., Socher, R., Li, L., Li, K., Li, F.: Imagenet: A large-scale hierarchical image database. In: 2009 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2009), 20-25 June 2009, Miami, Florida, USA. pp. 248–255 (2009)
13. Gal, Y., Islam, R., Ghahramani, Z.: Deep bayesian active learning with image data. In: Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017. pp. 1183–1192 (2017)

14. He, C., Wang, R., Shan, S., Chen, X.: Exemplar-supported generative reproduction for class incremental learning. In: British Machine Vision Conference 2018, BMVC 2018, Northumbria University, Newcastle, UK, September 3-6, 2018. p. 98 (2018)
15. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Conference on Computer Vision and Pattern Recognition. CVPR (2016)
16. Hinton, G.E., Vinyals, O., Dean, J.: Distilling the knowledge in a neural network. CoRR **abs/1503.02531** (2015)
17. Hou, S., Pan, X., Loy, C.C., Wang, Z., Lin, D.: Learning a unified classifier incrementally via rebalancing. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019)
18. Hou, S., Pan, X., Loy, C.C., Wang, Z., Lin, D.: Learning a unified classifier incrementally via rebalancing. In: IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019. pp. 831–839 (2019)
19. Hu, P., Lipton, Z.C., Anandkumar, A., Ramanan, D.: Active learning with partial feedback. In: 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019 (2019)
20. Javed, K., Shafait, F.: Revisiting distillation and incremental classifier learning. CoRR **abs/1807.02802** (2018)
21. Krasin, I., Duerig, T., Alldrin, N., Ferrari, V., Abu-El-Haija, S., Kuznetsova, A., Rom, H., Uijlings, J., Popov, S., Kamali, S., Malloci, M., Pont-Tuset, J., Veit, A., Belongie, S., Gomes, V., Gupta, A., Sun, C., Chechik, G., Cai, D., Feng, Z., Narayanan, D., Murphy, K.: Openimages: A public dataset for large-scale multi-label and multi-class image classification. Dataset available from https://storage.googleapis.com/openimages/web/index.html (2017)
22. Krizhevsky, A.: Learning multiple layers of features from tiny images. Tech. rep., University of Toronto (2009)
23. Li, Z., Hoiem, D.: Learning without forgetting. IEEE Trans. Pattern Anal. Mach. Intell. **40**(12), 2935–2947 (2018)
24. Mallya, A., Lazebnik, S.: Packnet: Adding multiple tasks to a single network by iterative pruning. In: 2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018. pp. 7765–7773 (2018)
25. Mccloskey, M., Cohen, N.J.: Catastrophic interference in connectionist networks: The sequential learning problem. The Psychology of Learning and Motivation **24**, 104–169 (1989)
26. Mensink, T., Verbeek, J.J., Perronnin, F., Csurka, G.: Distance-based image classification: Generalizing to new classes at near-zero cost. IEEE Trans. Pattern Anal. Mach. Intell. **35**(11), 2624–2637 (2013)
27. Oksuz, K., Cam, B.C., Kalkan, S., Akbas, E.: Imbalance problems in object detection: A review. CoRR **abs/1909.00169** (2019)
28. Oliver, A., Odena, A., Raffel, C.A., Cubuk, E.D., Goodfellow, I.: Realistic evaluation of deep semi-supervised learning algorithms. In: Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., Garnett, R. (eds.) Advances in Neural Information Processing Systems 31, pp. 3235–3246. Curran Associates, Inc. (2018)
29. Parisi, G.I., Kemker, R., Part, J.L., Kanan, C., Wermter, S.: Continual lifelong learning with neural networks: A review. Neural Networks **113**, 54–71 (2019)
30. Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., Lerer, A.: Automatic differentiation in pytorch. In: Advances in Neural Information Processing Systems Workshops. NIPS-W (2017)

31. Pernici, F., Bartoli, F., Bruni, M., Del Bimbo, A.: Memory based online learning of deep representations from video streams. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018)
32. Rebuffi, S., Bilen, H., Vedaldi, A.: Learning multiple visual domains with residual adapters. In: Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA. pp. 506–516 (2017)
33. Rebuffi, S., Kolesnikov, A., Sperl, G., Lampert, C.H.: icarl: Incremental classifier and representation learning. In: Conference on Computer Vision and Pattern Recognition. CVPR (2017)
34. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M.S., Berg, A.C., Li, F.: Imagenet large scale visual recognition challenge. International Journal of Computer Vision **115**(3), 211–252 (2015)
35. Rusu, A.A., Rabinowitz, N.C., Desjardins, G., Soyer, H., Kirkpatrick, J., Kavukcuoglu, K., Pascanu, R., Hadsell, R.: Progressive neural networks. CoRR **abs/1606.04671** (2016)
36. Scheffer, T., Decomain, C., Wrobel, S.: Mining the web with active hidden markov models. In: Proceedings of the 2001 IEEE International Conference on Data Mining, 29 November - 2 December 2001, San Jose, California, USA. pp. 645–646 (2001)
37. Sener, O., Savarese, S.: Active learning for convolutional neural networks: A coreset approach. In: 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings (2018)
38. Settles, B.: Active learning literature survey. Tech. rep., University of Winsconsin (2010)
39. Settles, B., Craven, M.: An analysis of active learning strategies for sequence labeling tasks. In: 2008 Conference on Empirical Methods in Natural Language Processing, EMNLP 2008, Proceedings of the Conference, 25-27 October 2008, Honolulu, Hawaii, USA, A meeting of SIGDAT, a Special Interest Group of the ACL (2008)
40. Shannon, C.E.: A mathematical theory of communication **27**, 379–423 (1948)
41. Wang, Y., Ramanan, D., Hebert, M.: Growing a brain: Fine-tuning by increasing model capacity. In: Conference on Computer Vision and Pattern Recognition. CVPR (2017)
42. Wu, Y., Chen, Y., Wang, L., Ye, Y., Liu, Z., Guo, Y., Fu, Y.: Large scale incremental learning. In: IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019. pp. 374–382 (2019)
43. Yang, J., Parikh, D., Batra, D.: Joint unsupervised learning of deep representations and image clusters. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2016)
44. Yoo, D., Kweon, I.S.: Learning loss for active learning. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019)
45. Zhang, K., Zhang, Z., Li, Z., Qiao, Y.: Joint face detection and alignment using multitask cascaded convolutional networks. IEEE Signal Process. Lett. **23**(10), 1499–1503 (2016)
46. Zhou, P., Mai, L., Zhang, J., Xu, N., Wu, Z., Davis, L.S.: M2KD: multi-model and multi-level knowledge distillation for incremental learning. CoRR **abs/1904.01769** (2019)

47. Zhou, Z., Shin, J.Y., Zhang, L., Gurudu, S.R., Gotway, M.B., Liang, J.: Fine-tuning convolutional neural networks for biomedical image analysis: Actively and incrementally. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017. pp. 4761–4772 (2017)