# Online Continual Learning for Robust Indoor Object Recognition

Umberto Michieli and Mete Ozay

*Abstract*— **Vision systems mounted on home robots need to interact with unseen classes in changing environments. Robots have limited computational resources, labelled data and storage capability. These requirements pose some unique challenges: models should adapt without forgetting past knowledge in a data- and parameter-efficient way. We characterize the problem as few-shot (FS) online continual learning (OCL), where robotic agents learn from a non-repeated stream of few-shot data updating only a few model parameters. Additionally, such models experience variable conditions at test time, where objects may appear in different poses (*e.g.*, horizontal or vertical) and environments (*e.g.*, day or night). To improve robustness of CL agents, we propose RobOCLe, which; 1) constructs an enriched feature space computing high order statistical moments from the embedded features of samples; and 2) computes similarity between high order statistics of the samples on the enriched feature space, and predicts their class labels. We evaluate robustness of CL models to train/test augmentations in various cases. We show that different moments allow RobOCLe to capture different properties of deformations, providing higher robustness with no decrease of inference speed.**

## I. INTRODUCTION

In the last decades, vision models have outperformed human-level accuracy on many recognition benchmarks. Deep learning has produced outstanding results by training all parameters of large models via offline updates over large batches of abundant samples available all at once.

The rise of specialized robots (*e.g.*, indoor domestic and service robots) has driven the recent advancements in automatic scene analyses. However, vision models deployed on robotic agents experience new classes and domains specific to users at test time [1]; and classical models cannot be easily personalized without suffering from catastrophic forgetting. Hence, Continual Learning (CL) paradigm has emerged to address such adaptation-preservation challenge [2], [3].

At the same time, robotic agents are equipped with low computation resources and no storage availability due to hardware and privacy constraints [4], [5], [6]. Therefore, online CL (OCL) methods have gained attention [7] to learn from a non-repeated stream of data and tasks. In particular, OCL has been recently investigated for low-resource embedded devices with no data storage available [4], [8]. However, no extensive study exists at the time of writing. Furthermore, when learning new classes, users will only provide a few labelled samples, as the labeling operation is time-consuming and tedious. To tackle this, few-shot (FS) learning [9] should pair with OCL (FS-OCL) [10].

Finally, after learning new concepts, robotic agents move around the environment (*e.g.*, the user home) and discover the same or other instances of newly learned objects in different conditions. Therefore, vision models should be robust to test-time distortion of objects (*e.g.*, different poses, viewpoints, *etc.*) and/or environment (*e.g.*, different illumination, room, *etc.*) [11], [12]. For example, a user may want to update the robotic agent to recognize, *e.g.*, its own pet. The user will likely provide *a few* images of the pet, *e.g.*, lying on the living room carpet; however, the agent should identify the same pet in other poses (*e.g.*, standing in front of a door) and in other domains (*e.g.*, corridor or kitchen).

In this paper, we introduce the novel scenario of parameter-efficient FS-OCL for low-resource robotic agents for robust test-time performance in variable conditions. We extend previous setups [4] to analyse and improve the robustness of the algorithms to test-time distortions without updating model weights at test time. To our knowledge, this challenging and very practical scenario has not been addressed in the existing literature, and in general, little exploration has been carried out on robustness of CL models to variable test-time conditions [13], [14]. The main details of our scenario are depicted in Fig. 1.

Our contributions are summarized as follows:

- We tackle a new FS-OCL paradigm for low-resource robots where storage is unavailable and computation power is low (*i.e.*, most of model weights cannot be updated).
- We introduce a new parameter-efficient FS-OCL method (RobOCLe) suitable for mobile robotic agents and FS data. RobOCLe extracts high order statistics from embeddings, building more reliable feature representations robust to variable domain and object conditions. RobOCLe shows consistent improvements against 10 OCL baselines on 4 benchmarks and 16 backbones. RobOCLe shows a room-aware relative accuracy gain (RARG, Sec. IV) of average 65.8% on same-domain data and 16.4% on other-domain data compared to the best competing approach, while decreasing inference FPS by less than 0.5%.
- We present a new evaluation paradigm to determine robustness of vision models, and particularly CL models, over a suite of real and synthetic augmentations, resembling practical use cases. RobOCLe shows average RARG of 18.3% on controlled augmentations of other-domain data.

## II. RELATED WORK

**Continual Learning** has many definitions. Classical CL assumes data released at incremental steps and multiple training epochs performed with large batch sizes updating the whole net [3], [15], [16], [17]. Online CL (OCL) assumes, instead, that data comes as a non-repeated stream of small-batch samples [7], [18], [19], [20], [21], [22]. OCL methods

All authors are with Samsung Research UK.
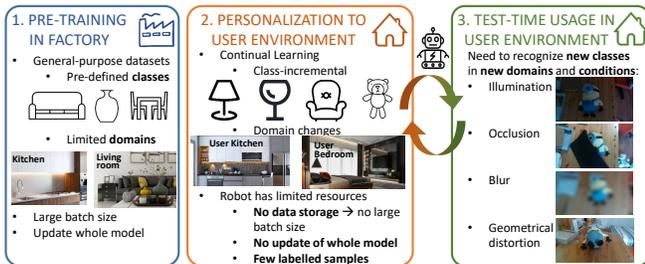`{n.surname}@samsung.com`

Fig. 1. Overview of the considered use-case in 3 stages. 1) A vision recognition model is pre-trained on factory servers on large-scale general-purpose datasets on a pre-defined set of classes and a limited set of acquisition domains. Model pre-training requires large computational power to update the whole model via large batches. 2) The model is embedded into *e.g.*, home robots, and shipped to users. Users customize the recognition model to discover novel classes in their target domain/environment (*e.g.*, user home). Robots do not have large storage due to privacy and storage constraints, and cannot update the whole model due to computational limitations. Large batches are also unavailable, since no storage and low RAM is assumed. Users do not label many samples as it is a tedious operation. 3) The robot navigates the new environment and discovers recently-added classes in unseen domains and conditions, *e.g.*, under variable illumination, occlusion, blurring and geometric distortions. Stages 2) and 3) can alternate at any time.

are mainly based on replay due to the inherent difficulty to prevent forgetting. Recently, OCL for embedded devices [4], [8], [23], [24] further restricts the scope to parameter-efficient updates with unitary batch size due to low-resource devices. In parallel, few-shot CL (FS-CL) [25], [26], [27], [28], [29] and FS-OCL [10] train a model to recognize new classes based on few labelled data only, via fast and efficient model updates with minimal human effort.

Along all these definitions, three main scenarios have been considered [30]: 1) class-incremental (CI) where new classes are introduced to models over time; 2) domain-incremental (DI) where the same problem is learned in different contexts (*e.g.*, with a fixed set of classes); and 3) task-incremental (TI), where new distinct tasks are presented to models.

In our paper, we focus on the most practical and less-explored FS-OCL for low-resource devices and few-shot data in the CI setup, with domain changes happening at both fine-tuning and testing stages. OCL, FSCL and FSOCL methods generally assume replay data [31], [32], [33], [7], violating the important assumption that no storage is guaranteed. Therefore, we rely on the baselines introduced in [4]. Additionally, we evaluate robustness of agents at test time.

**Test-Time Performance** is heavily influenced by the domain gap with respect to the training set [34], [35]. Test-Time Augmentation (TTA) approaches use data augmentation at test time to obtain greater robustness [11], [12], improved accuracy [36], [37], [38], [39], [40], [41], or estimates of uncertainty [42], [43]. TTA pools predictions from several transformed versions of a given test input to obtain a *smoothed* prediction. In particular, [37], [38], [39] average the predictions obtained over augmented views (*e.g.*, cropped, rotated, *etc.*) of input samples. [40] averages embeddings obtained over augmented views of the input samples. [36], [41] propose a learnable aggregation of TTAs. [44], [41] learn a selection policy for TTAs. [45] enhances images to make them recognition-friendly.

Test-Time Training (TTT) is a special case of Unsupervised Domain Adaptation (UDA) where a model trained on a source domain is adapted to an unlabeled target domain without accessing source data [46], [47], [48]. Many approaches design novel loss functions to train the model and reduce domain shift, requiring access to the full target set. Recently, [12] performs TTT on single images coupling TTA with TTT to improve test-time robustness. Continual TTT [13], [14] has been recently considered.

Our work is linked to this, with some key differences:

- We do not assume the test set to be available all at once, as in practical applications where an agent encounters new samples and domains over time.
- Our scope is different. TTA boosts model performance at test time aggregating multiple predictions. TTT updated the model at test time. Our aim is to propose an OCL method with increased robustness over variable test sets.
- We do not compute multiple TTAs of the same input sample, but rather mimic a deployment scenario where robots experience variable domains with no clear delimitation.
- We extract rich information from inputs, to reduce the need for TTA/TTT and maintain same inference time (we use one inference step per sample, opposed to TTA/TTT).

**Pooling** has been one of the key enablers for deep learning models mapping high dimensional features into a summary of active features [49]. Multiple schemes have been proposed to extract the most useful information from the input data. We compare our pooling mechanism with the following schemes: average pooling (AVG) [50]; max pooling (MAX) [51]; a linear combination of AVG and MAX (MIX) [52]; a dropout-inspired probabilistic selection of activations based on a multinomial distribution fit (STOCHASTIC) [53]; a concatenation of top-k% of features to discover information spread over the spatial map (RAP) [54]; a concatenation of AVG and MAX pooling (AVGMAX) [55]; an average of the $p$-th order non-central moment ($L_p$-norm) [56]; and iSQRT-COV, which computes second-order covariance pooling via iterative matrix square root normalization [57]. In our work, we consider high order central moments [58] to discover robust features from limited input data.

## III. METHODS

We investigate OCL from sequential data streams. We evaluate existing approaches and propose a new method (RobOCLe) to overcome some of the challenges that CL methods face when tested in novel conditions.

In our setup, a model learns from streamed examples, seen one at a time with no repetitions, since many applications cannot store samples due to the continuous nature of the data stream and the limited time and storage to process it. Over time, the learner sees new classes but no task identifier is provided, since the agent should infer the class regardless of the specific task information. On resource-constrained devices, CL methods should be both data and computation efficient: they should learn from few-shot labelled data and adhere to strict memory and time constraints. Finally, agents will operate in uncontrolled environments, therefore, CL methods should be robust to (severe) test-time corruptions.

## A. RobOCLe: **Rob**ust **O**nline **C**ontinual **Le**arning

To tackle our scenario, we propose a new **Rob**ust **O**nline **C**ontinual **Le**arning method (RobOCLe), comprising of 3 main blocks: a feature extractor $G(\cdot)$, a pooling scheme $P(\cdot)$ and a classifier $F(\cdot)$ to categorize samples from their extracted features. We incrementally train a neural network $F(P(G(\mathbf{x}_{n,k})))$ via supervised OCL updates over subsequent tasks to generate predictions $\hat{y}_{n,k}$, where $\mathbf{x}_{n,k}$ is the $n$-th input sample, $n \in [N_k]$, of the $k$-th task, $k \in [K]$.

Following traditional transfer learning setups, $G(\cdot)$ is a backbone network (*e.g.*, ResNets [38]) pre-trained on a server on public benchmarks, with large batch sizes over many epochs updating the whole model (point 1 of Fig. 1). When $G(\cdot)$ is embedded in robotic agents, it experiences new domains (*e.g.*, user environments) compared to training ones.

Additionally, users want to personalize models to recognize personal objects. Therefore, a novel classifier $F(\cdot)$ adapts the general features extracted by $G(\cdot)$ to the class set of particular users (point 2 of Fig. 1). Users show few labeled samples of the new classes to recognize, and robots do not have storage or large computation resources. Therefore, $F(\cdot)$ is trained using online updates, seeing a non-repeated stream of few-shot labelled samples using a frozen feature extractor $G(\cdot)$. In the experiments, we consider the most challenging case of a single class being learned per task.

**High order pooling.** In order to apply $F(\cdot)$ over the pre-trained backbone $G(\cdot)$, a pooling scheme is generally inserted to summarize the input information and reduce the feature map size. At test time, the robotic agent should recognize personal classes in a variety of conditions (*e.g.*, illumination, occlusion, blur, pose). Therefore, the overall OCL procedure should be robust to such factors of variations at test time. With this aim in mind, we perform the pooling $P(\cdot)$, employing higher order statistical moments to increase the amount of clues extracted from the input samples. In particular, $P(\cdot)$ computes and concatenates the first $R$ statistical moments from the output of $G(\cdot)$. Such moments characterize the distribution $\mathcal{G}$ of $g \triangleq G(\mathbf{x}_{n,k}) \in \mathbb{R}^{h \times w \times d}$. where $h$ and $w$ are the spatial sizes of features and $d$ is the number of channels, and thus capture rich spatial statistics to improve recognition accuracy. More formally, we employ

$$P(g)=\Big\|\left(\mu, E_{\mathcal{G}}\left[(g-\mu)^2\right]^{\frac{1}{2}}, \Big\|_{r=3}^{R} E_{\mathcal{G}}\left[\frac{g-\mu}{E_{\mathcal{G}}\left[(g-\mu)^2\right]^{\frac{1}{2}}}\right]^r\right),$$

(1)

where $E_{\mathcal{G}}[\cdot]$ is the expectation over $\mathcal{G}$, $\mu$ is the empirical mean, and $\|(\cdot)$ denotes the concatenation operation. That is, $P(g) \in \mathbb{R}^{R \cdot d}$ concatenates the first $R$ moments of $g$. In our case, we feed $R = 3$ moments to the classifier. Fig. 2 shows the distribution of the first statistical moments of features extracted from clean vs. augmented samples. AVG has more variability (higher Wasserstein distance) than higher moments when changing domain. In Sec. V we confirm that high order moments increase robustness of OCL models and their invariance to augmentations on both same-domain and other-domain test data (*i.e.*, when models receive test data in other conditions than at train time).

**Similarity estimation.** The classifier $F(\cdot)$ should be representative of the new classes and extremely lightweight for on-device deployment. We extend the Nearest Classifier Mean (NCM) and the Linear Discriminant Analysis (LDA) to support the streaming setup [20] and the high order pooling. Streaming NCM and LDA are denoted by NCM and SLDA.

RobOCLe$_{\text{NCM}}$ computes a running mean feature vector per class (*i.e.*, the $c$-th class prototype $\mathbf{m}_c \in \mathbb{R}^{R \cdot d}$, $\forall c \in \mathcal{C}$) each with an associated counter denoting the number of samples employed to compute each average value ($\mathbf{t}_c$). Given a new data vector $\mathbf{x}_{n,k}$ with the associated label $y_{n,k}$, we embed it to $\mathbf{z}_{n,k} \triangleq P(G(\mathbf{x}_{n,k})) \in \mathbb{R}^{R \cdot d}$, and update the class mean and associated counter by

$$\mathbf{m}_{y_{n,k}} \leftarrow \frac{\mathbf{t}_{y_{n,k}} \cdot \mathbf{m}_{y_{n,k}} + \mathbf{z}_{n,k}}{\mathbf{t}_{y_{n,k}} + 1}, \quad \mathbf{t}_{y_{n,k}} \leftarrow \mathbf{t}_{y_{n,k}} + 1. \quad (2)$$

For inference on new samples, RobOCLe$_{\text{NCM}}$ assigns the label of the nearest prototype according to its $\ell_2$ distance on the enriched feature space. The original NCM has shown to be a simple yet effective baseline in CL [59], [19], [4], [60].

RobOCLe$_{\text{SLDA}}$ computes one channel-wise covariance matrix of features shared across classes ($\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$) that is updated online via [61]. During inference, SLDA assigns to a new sample the label of the closest Gaussian model in the feature space defined using the running class means and shared $\boldsymbol{\Sigma}$. RobOCLe$_{\text{NCM}}$ is a special case of RobOCLe$_{\text{SLDA}}$ where $\boldsymbol{\Sigma}$ is equal to the identity matrix. We use the implementation from [20] to update the $\boldsymbol{\Sigma}$ and compute predictions. RobOCLe$_{\text{SLDA}}$ runs inference on a new sample[1] $\mathbf{x}_n$ by $F(P(G(\mathbf{x}_n))) = \mathbf{W}\mathbf{z}_n + \mathbf{b}$, where $\mathbf{W} \in \mathbb{R}^{|\mathcal{C}| \times R \cdot d}$ and $\mathbf{b} \in \mathbb{R}^{|\mathcal{C}|}$, where $|\cdot|$ is the set cardinality. Rows of $\mathbf{W}$ are computed by $\mathbf{w}_c = \boldsymbol{\Lambda}\mathbf{m}_c^T$, and elements of $\mathbf{b}$ are computed by $b_c = -0.5(\mathbf{m}_c \cdot \boldsymbol{\Lambda}\mathbf{m}_c) = (\mathbf{m}_c \cdot \mathbf{w}_c)$ with the shrinking approximation $\boldsymbol{\Lambda} = [(1 - \epsilon)\boldsymbol{\Sigma} + \epsilon\mathbf{I}]^{-1}$ with parameter $\epsilon = 10^{-4}$. Running covariance [61] is then computed by

$$\boldsymbol{\Sigma}_{n+1} = \frac{n\boldsymbol{\Sigma}_n + \delta_n}{n + 1}, \quad \delta_n = \frac{n(\mathbf{z}_n - \mathbf{m}_{y_n})(\mathbf{z}_n - \mathbf{m}_{y_n})^T}{n + 1}.$$

(3)

## B. OCL Baselines

We describe several OCL methods used to update $F(\cdot)$, mostly focusing on data-free approaches. **Fine-Tuning (FT)** updates a fully-connected output layer using SGD and CE loss with no mechanisms to prevent forgetting [4]. **Online Perceptron (PRCPT)** keeps one weight vector for each class [4], initialized to the first sample of the class. After that, when the model misclassifies a new sample, the class weight vector and the weight vector of the misclassified class with the highest score are updated. During inference, a label is assigned by taking the $\arg\max$ over the dot product between weights and input vector. **Online Centroid-Based Concept Learning (CBCL)** extends NCM with multiple class prototypes [25]. At inference, CBCL searches for the weighted nearest neighbor, where class weights are inversely

---

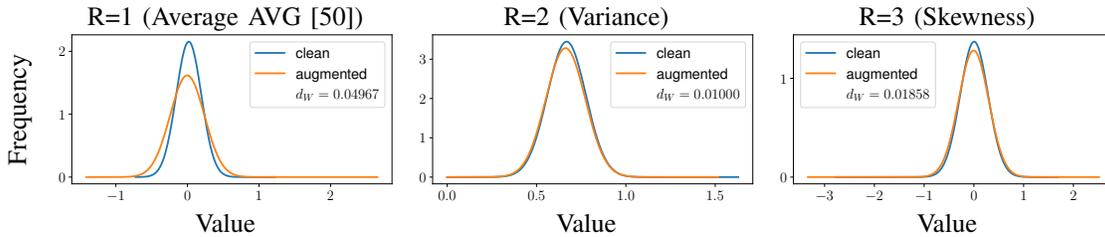[1]To simplify notation, we remove task index $k$.

Fig. 2. Distribution of statistical moments in $\mathbf{z}_n$. $d_W$: Wasserstein distance.



Fig. 3. Datasets employed in our work. Right-hand side table summarizes differences between train and test sets. [Text within squared brackets] specifies how properties change when our augmentations are applied (Sec. IV).



Fig. 4. Synthetic data augmentations employed to analyse different conditions in which objects could be found.

proportional to the number of samples seen so far for that class. We use the default parameters suggested in [26], *i.e.*, a distance threshold of 17, a nearest neighbor value of 1, and a maximum buffer size of 44 prototypes. We noticed no clear difference compared to NCM, as in [4]. **Streaming One-vs-Rest (SOvR)** measures how close a new input is to a class mean vector while also considering its distance to examples from other classes [4], which is reminiscent of SVMs. **SQDA** extends SLDA estimating one covariance matrix for each class (*i.e.*, Quadratic). The drawbacks are increased memory consumption and need for many samples per class to estimate reliable covariance matrices [62]. **Streaming Gaussian Naïve Bayes (SNB)** estimates a running variance vector per class [63] (*i.e.*, diagonal covariance matrices assuming independent features). It requires significantly less memory than SQDA. **Online iCaRL** maintains a class-balanced memory buffer [20], [59] and randomly replaces an example from the most represented class with a new sample when the buffer is full. During training, it randomly draws examples from the buffer, combines them with the new sample and makes a single update. Similar to other methods, we train a linear output layer. While effective, it requires storing sensitive samples on devices. We use a buffer size of 1K samples (iCarl) or of 2 samples per class (iCarl-2pc).

## IV. EXPERIMENTAL SETUP

**Datasets.** We analyse our novel setup on F-SIOL-310 [26] and OpenLORIS(-Objects) [64], which are robot-acquired datasets specifically designed for CL and FS, and their train and test sets have different properties. A visual summary is reported in Fig. 3. We consider a class-IID setup, where all samples are ordered by class and shuffled within each class.

*OpenLORIS* collects 121 instances including 40 categories of daily necessities objects under 20 scenes from variable camera-object distance/angle. The dataset comprises four environmental factors: *illumination* variability during recording, *occlusion* percentage of the objects, object pixel size in each frame (*i.e.*, *scale* of objects), and *clutter* of the scene. Each factor has three levels of intensity, yielding around

550,000 total RGB samples. Following [18], [4], we consider two splits: instance and instance-small.

The instance split (referred to as OpenLORIS) presents videos (*i.e.*, image frames from the same video) of shuffled object instances to learners one at a time. For example, a learner experiences a video of toy#1, then a video of hat#3, then a video of toy#2, and so on. The learner experiences temporally-correlated data with repeating classes. This split uses the original train-val division provided by the authors.

The instance-small split (OpenLORIS-small) evaluates low-shot out-of-domain generalization ability of models. When learning a new class, frames from a single video of an object instance are shown to the learner. When testing, the model is evaluated on all test data containing the classes seen so far. Test data contains same object instance seen under different domains (*i.e.*, factors) as well as other object instances from the same classes seen under different domains. Challenges of this setup are: low-shot labelled learning from correlated data, generalization to different object instances of same class, generalization to unseen domains. This setup is often encountered in practice. For example, given a new pet, users provide a few labels and a classifier should recognize the pet in different domains and conditions.

*F-SIOL-310* is specifically designed for Few-Shot Incremental Object Learning. We employ this dataset to further evaluate the ability of models to generalize from extremely small labelled data. It contains 310 instances of 22 household categories with no object instances overlapping between train and test sets. We consider two splits, respectively with 5 and 10 shots (*i.e.*, 5 or 10 labelled training samples).

**Objects, Conditions & Augmentations.** The considered datasets span a wide range of setups with variable factors. OpenLORIS uses the whole dataset and the same objects are shown in train and test sets. Images are acquired under some different conditions (illumination, occlusion, scale, clutter). However, co-existence of conditions is limited. OpenLORIS-small restricts the set of samples and has only one instance of objects in training. F-SIOL-310 uses 5 or 10 training samples, and has different objects between train and test sets. However, objects are acquired with similar conditions (same background, viewpoint/scale, no occlusion, no clutter).

To cope with the limited variability of existing datasets,

we introduce some augmentations at either train or test splits. A visual summary is shown in Fig. 4: we mimic several conditions as described next.1) Illumination changes via color jittering with parameters chosen uniformly in the following ranges: brightness in $[0.5, 1.5]$, contrast in $[0.5, 1]$, saturation in $[0.5, 1.5]$, and hue in $[-0.1, 0.1]$. 2) Image blur (*e.g.*, due to dust on camera or wrong focus) via Gaussian noise with $11 \times 11$ kernel and standard deviation (std) in $[0.1, 0.5]$. 3) Geometrical changes via a combination of a) affine transformation keeping image center invariant, with degree in $[-30, 30]$, translate shift in $[-0.2, 0.2]$ for both horizontal and vertical axes, scale in $[0.8, 1.2]$, and shear in $[-0.1, 0.1]$); b) random perspective with probability $0.2$ and distortion amount $0.2$, c-d) horizontal and vertical flip with probability $0.5$ and $0.3$, respectively. 4) Last, we also consider all of the above together. To ensure fair comparisons among OCL methods, all samples are augmented with the same parameters in each set of experiments.

All the results are averaged across 5 class orderings, although std is not reported to improve readability. Offline training upper bounds are not considered, because the classifier $F(\cdot)$ is different for each approach.

**Metrics.** We evaluate both plasticity and forgetting in comparison of OCL methods [7]. All metrics are computed on test sets at the end of training phases. We compute accuracy (Acc %, ↑) to assess the final model performance across all classes. Often, we report the room-aware relative gain, RARG, ($\Delta_R$, %) of model$_2$ with accuracy Acc$_2$ compared to a model$_1$ with Acc$_1$ with respect to the available room of improvement, defined by $\Delta_R \triangleq (\text{Acc}_2 - \text{Acc}_1)/(100 - \text{Acc}_1)$. In addition, we also measure: (1) backward transfer (BwT, ↑), which tracks the influence that learning a new task has on the preceding tasks' performance, to measure stability; (2) forgetting (Forg, ↓), which averages difference of class-wise Acc at the last step and the best previous class-wise Acc; (3) plasticity (Pla, ↑) as the average Acc achieved on each task evaluated right after learning that task; (4) training time (TTime [min]); and (5) frames per second (FPS) in inference. Evaluation is carried out using an NVIDIA GeForce RTX 2080Ti GPU supported by an Intel(R) Xeon(R) Gold 5218R CPU@2.10GHz with 80 cores and 252GB RAM. Nonetheless, relative considerations remain the same on hardware mounted on robotic agents. All results are averaged across 5 runs with different seeds: only the mean value is shown since the standard deviation is negligible ($\leq 0.16$ in all cases).

## V. EXPERIMENTAL VALIDATION

We evaluate our RobOCLe in multiple scenarios.
**Same-domain data (OpenLORIS).** In the first analyses, we consider 16 backbones, including both CNN-based (2x MobileNets-V3 [65], 2x EfficientNets [66], and 5x ResNets [38]) and transformer-based (3x Swins [67] and 4x ViTs [68]) of variable sizes. These are commonly used baselines and allow us to examine variability of results depending on model size and learned visual features. All architectures have been pre-trained on the ImageNet [37] dataset.



Fig. 5. Per-step Acc of RN152 on OpenLORIS-small. Average gains of RobOCLe over NCM and SLDA are 5.81% and 1.20%, respectively.

Accuracy is reported in Tab. I. Most naïve OCL methods (*i.e.*, PRCPT, SNB, SOvR) are unable to outperform FT, since, in this case, data is sufficient and correlated, with no severe change of domains/objects at test time. Replay data (*i.e.*, iCaRL) shows benefits, however, contradicts the important OCL assumption on data availability: robots should not store data due to privacy and storage constraints.

SLDA outperforms all other competitors thanks to the Gaussian modeling of features, while its extension (SQDA) can achieve competitive performance only on certain backbones: we argue that this is because SQDA cannot create reliable class-wise correlation matrices, especially for larger networks. RobOCLe brings significant improvements with respect to the best approach (SLDA), reaching an outstanding average accuracy of $99.21\%$ with a room-aware relative gain ($\Delta_R$) of $65.8\%$ over it. This first analysis justifies the robustness of RobOCLe to multiple backbones on same-domain test data when sufficient training data is available. This experimental case is the main benchmark employed in the existing literature [7]; however, we argue that it may often not be the case in practice, as we discuss next.

**Real other-domain few-shot data (OpenLORIS-small, F-SIOL-310).** When deploying robots to users and updating their AI models to recognize new objects, we typically encounter two simultaneous problems: 1) users do not label a large amount of data as it is a tedious time-consuming operation, and 2) training and testing sets have different conditions in terms of, *e.g.*, object pose and background environment. For example, users show some pictures of their pet lying on the living room carpet during daytime, while at a later stage, the dog's pose and the environment will change.

In Tab. II, we analyse this challenging scenario of FS-OCL with a different domain (other-domain) at test time. We consider three benchmarks based on OpenLORIS-small and F-SIOL-310 (with 5 and 10 shots), while we restrict the evaluation to ResNets and ViTs, being the most popular CNNs and Transformers at the time of writing (also, we confirmed their effectiveness in Tab. I).

Also in this case, FT and naïve OCL methods show low accuracy. SQDA achieves competitive performance only in certain cases. iCaRL improves significantly classical FT, however, it shows limited results. We observe that when data is scarce, NCM reduces the gap with respect to SLDA. In some cases (*e.g.*, ResNets on OpenLORIS-small) NCM outperforms SLDA. We argue that this reflects two properties: 1) SLDA optimizes a larger number of parameters (prototypes and covariance matrix) compared to NCM (prototypes matrix only); and 2) feature dimensionality of ResNets (*e.g.* 2048 for RN152) is higher than that of ViTs (*e.g.* 1024 for ViT-

TABLE I

ACC ON SAME-DOMAIN OPENLORIS DATA ON 16 BACKBONES AND 10 OCL BASELINES. MN: MOBILENET, EN: EFFICIENTNET, RN: RESNET.

| | MN-S | MN-L | EN-B0 | EN-B1 | RN18 | RN34 | RN50 | RN101 | RN152 | Swin-T | Swin-S | Swin-B | ViT-B16 | ViT-B32 | ViT-L16 | ViT-L32 | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| FT | 84.22 | 93.38 | 96.14 | 97.89 | 83.35 | 83.84 | 94.19 | 95.43 | 94.33 | 95.68 | 96.07 | 96.04 | 59.07 | 67.27 | 92.47 | 62.85 | 87.01 |
| PRCPT [4] | 74.49 | 89.75 | 94.02 | 95.94 | 80.47 | 78.58 | 92.18 | 91.24 | 92.35 | 93.40 | 94.21 | 94.46 | 48.70 | 58.51 | 85.62 | 55.19 | 82.44 |
| SNB [63] | 31.12 | 37.84 | 77.96 | 83.91 | 1.51 | 0.84 | 0.00 | 0.00 | 0.00 | 87.30 | 86.94 | 86.14 | 3.51 | 4.60 | 42.90 | 5.76 | 34.40 |
| SOvR [4] | 37.42 | 60.17 | 73.92 | 80.65 | 34.65 | 31.77 | 71.54 | 64.57 | 67.68 | 80.03 | 79.01 | 77.92 | 4.91 | 18.83 | 62.61 | 22.02 | 54.23 |
| NCM/CBCL [60], [25] | 72.89 | 81.83 | 85.94 | 88.34 | 79.69 | 80.47 | 84.62 | 83.98 | 84.12 | 87.77 | 87.54 | 87.94 | 64.99 | 63.38 | 79.47 | 68.09 | 80.07 |
| SQDA [62] | 77.71 | 55.84 | 2.45 | 6.16 | 81.59 | 81.36 | 5.66 | 24.24 | 1.64 | 83.84 | 62.53 | 63.27 | 8.91 | 15.74 | 3.53 | 7.50 | 36.37 |
| iCaRL [59] | 91.76 | 95.54 | 97.60 | 98.06 | 92.76 | 93.21 | 97.18 | 97.47 | 97.63 | 97.65 | 97.57 | 97.98 | 86.43 | 89.61 | 95.52 | 89.61 | 94.72 |
| iCaRL (2pc) [59] | 89.29 | 95.09 | 97.07 | 96.43 | 91.21 | 92.34 | 96.50 | 96.99 | 96.79 | 97.13 | 97.13 | 97.68 | 79.41 | 82.49 | 93.34 | 81.78 | 92.54 |
| SLDA [20] | 95.57 | 97.93 | 98.83 | 98.98 | 95.01 | 95.47 | 99.00 | 99.10 | 99.13 | 98.25 | 98.18 | 98.85 | 96.74 | 96.20 | 98.69 | 97.04 | 97.69 |
| RobOCLe$_{SLDA}$ (ours) | **98.20** | **99.37** | **99.70** | **99.72** | **97.65** | **97.96** | **99.69** | **99.78** | **99.78** | **99.28** | **99.31** | **99.65** | **99.16** | **99.02** | **99.73** | **99.33** | **99.21** |
| ($\Delta_R$ [%]) | (+59.5) | (+69.4) | (+74.8) | (+72.3) | (+52.8) | (+54.9) | (+69.1) | (+75.6) | (+75.2) | (+58.6) | (+62.0) | (+69.9) | (+74.2) | (+74.3) | (+79.6) | (+77.4) | (+65.8) |

TABLE II

ACCURACY ON REAL OTHER-DOMAIN FEW-SHOT DATASETS ON RESNETS AND VITS BACKBONES.

| | OpenLORIS-small | | | | | | | | F-SIOL-310 (5-shots) | | | | | | | | F-SIOL-310 (10-shots) | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | RN50 | RN101 | RN152 | ViT-B16 | ViT-B32 | ViT-L16 | ViT-L32 | Avg | RN50 | RN101 | RN152 | ViT-B16 | ViT-B32 | ViT-L16 | ViT-L32 | Avg | RN50 | RN101 | RN152 | ViT-B16 | ViT-B32 | ViT-L16 | ViT-L32 | Avg |
| FT | 13.64 | 14.81 | 16.05 | 2.08 | 2.08 | 2.51 | 2.08 | 7.61 | 40.72 | 27.65 | 31.96 | 8.56 | 9.87 | 19.67 | 8.63 | 21.01 | 30.25 | 26.75 | 32.83 | 5.00 | 7.33 | 14.42 | 7.25 | 17.69 |
| PRCPT [4] | 27.45 | 21.68 | 21.36 | 2.08 | 2.18 | 6.64 | 2.08 | 11.92 | 32.16 | 24.12 | 34.31 | 4.84 | 7.52 | 7.39 | 4.84 | 16.15 | 38.25 | 24.42 | 32.75 | 5.17 | 5.17 | 9.75 | 4.92 | 17.20 |
| SNB [63] | 0.57 | 0.58 | 0.25 | 3.10 | 4.19 | 23.52 | 5.16 | 5.34 | 4.90 | 2.94 | 6.47 | 11.31 | 9.28 | 24.90 | 7.84 | 9.66 | 6.00 | 0.83 | 1.92 | 6.92 | 6.42 | 20.58 | 15.42 | 8.30 |
| SOvR [4] | 46.72 | 45.96 | 43.76 | 9.33 | 10.31 | 29.21 | 17.24 | 28.93 | 80.39 | 63.01 | 64.97 | 18.69 | 8.17 | 54.58 | 14.51 | 43.47 | 76.33 | 63.67 | 67.50 | 17.75 | 7.00 | 50.08 | 10.00 | 41.76 |
| SQDA [62] | 37.65 | 47.08 | 46.10 | 37.34 | 36.20 | 41.41 | 36.41 | 40.31 | 93.73 | 94.18 | 91.90 | 96.27 | 94.44 | 95.03 | 93.40 | 94.11 | 94.67 | 96.00 | 95.50 | 96.00 | 96.92 | 96.42 | 95.99 | 95.99 |
| iCaRL [59] | 51.29 | 51.04 | 50.33 | 33.60 | 33.57 | 42.76 | 32.58 | 42.17 | 58.95 | 66.21 | 63.53 | 34.90 | 40.39 | 38.50 | 32.16 | 47.81 | 68.83 | 77.75 | 76.50 | 47.08 | 48.75 | 51.42 | 44.08 | 59.20 |
| iCaRL (2pc) [59] | 49.08 | 47.28 | 47.93 | 31.83 | 28.92 | 40.89 | 31.73 | 39.67 | 58.04 | 67.12 | 64.12 | 37.25 | 37.65 | 38.69 | 32.88 | 47.96 | 69.33 | 78.33 | 76.33 | 46.17 | 49.83 | 50.50 | 45.33 | 59.40 |
| NCM/CBCL [60], [25] | 50.59 | 48.31 | 47.68 | 34.58 | 32.20 | 40.57 | 34.28 | 41.17 | 94.90 | 93.53 | 93.66 | 91.96 | 92.22 | 94.84 | 91.70 | 93.26 | 93.83 | 94.42 | 94.33 | 93.58 | 96.17 | 96.67 | 94.17 | 94.74 |
| RobOCLe$_{NCM}$ (ours) | **55.67** | **54.84** | **54.89** | **36.39** | **32.74** | **41.23** | **35.12** | **44.41** | **95.90** | **95.37** | **94.85** | **94.33** | **93.91** | **95.70** | **92.99** | **94.72** | **95.86** | **96.81** | **96.51** | **95.58** | **97.00** | **97.69** | **95.71** | **96.45** |
| ($\Delta_R$ [%]) | (+10.3) | (+12.6) | (+13.8) | (+2.8) | (+0.8) | (+1.1) | (+1.3) | (+5.5) | (+19.7) | (+28.5) | (+18.8) | (+29.5) | (+21.6) | (+16.) | (+15.6) | (+21.7) | (+32.9) | (+42.8) | (+38.5) | (+31.2) | (+21.7) | (+30.6) | (+26.4) | (+32.6) |
| SLDA [20] | 50.26 | 49.91 | 50.41 | 43.96 | 41.50 | 45.51 | 42.93 | 46.35 | 94.38 | 94.77 | 93.99 | 96.34 | 95.62 | 94.64 | 95.23 | 95.00 | 97.00 | 97.92 | 97.97 | 96.25 | 99.00 | 98.17 | 98.33 | 97.90 |
| RobOCLe$_{SLDA}$ (ours) | **51.33** | **51.44** | **52.42** | **44.73** | **42.86** | 45.22 | **43.07** | **47.29** | **96.12** | **96.98** | **95.84** | **96.96** | **95.80** | **95.38** | 94.73 | **95.97** | **97.42** | **98.33** | 97.47 | **99.18** | 98.31 | **98.40** | **98.78** | **98.27** |
| ($\Delta_R$ [%]) | (+2.1) | (+3.1) | (+4.1) | (+1.4) | (+2.3) | (-0.5) | (+0.2) | (+1.8) | (+31.0) | (+42.2) | (+30.8) | (+17.0) | (+4.0) | (+13.7) | (-10.5) | (+19.5) | (+13.9) | (+20.0) | (+32.6) | (+17.9) | (+8.0) | (+4.3) | (+8.4) | (+17.5) |

TABLE III

ACCURACY ON OTHER-DOMAIN DATA GENERATED VIA CONTROLLED AUGMENTATIONS ON THE OPENLORIS WITH RESNET152 AND VIT-L16.

RESULTS HIGHLIGHTED IN YELLOW: SAME AUGMENTATIONS BETWEEN TRAIN AND TEST. TR: TRAIN, TE:TEST.

| | | NCM | | | | | | RobOCLe$_{NCM}$ (ours) | | | | | | | | | NCM | | | | | | | RobOCLe$_{NCM}$ (ours) | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | tr↓ te→ | clean | illum | geom | noise | all | Avg OD | clean | illum | geom | noise | all | Avg OD | $\Delta_R^{OD}$ | | tr↓ te→ | clean | illum | geom | noise | all | Avg OD | clean | illum | geom | noise | all | Avg OD | $\Delta_R^{OD}$ |
| ResNet152 | clean | 84.12 | 34.75 | 77.37 | 60.82 | 13.40 | 46.59 | **88.15** | 41.26 | 82.10 | 66.16 | 18.16 | 51.92 | (+10.0) | ViT-L16 | clean | 79.47 | 35.11 | 73.43 | 66.57 | 19.49 | 48.65 | **80.31** | 33.13 | 73.69 | 68.38 | 19.74 | 48.73 | (+0.2) |
| | illum | 56.91 | 53.91 | 54.77 | 45.75 | 30.29 | 46.93 | 62.65 | **56.83** | 59.92 | 51.26 | 31.91 | 51.44 | (+8.5) | | illum | 30.78 | 40.58 | 33.00 | 29.12 | 30.00 | 30.73 | 30.82 | **35.00** | 35.17 | 30.65 | 31.22 | 31.96 | (+1.8) |
| | geom | 82.01 | 40.49 | 81.07 | 64.05 | 19.67 | 51.55 | 86.21 | 45.11 | **84.39** | 67.92 | 23.90 | 55.79 | (+8.7) | | geom | 78.70 | 37.05 | 79.19 | 68.73 | 23.31 | 51.67 | 78.00 | 37.05 | **78.64** | 68.87 | 24.32 | 52.24 | (+1.2) |
| | noise | 77.27 | 35.60 | 71.99 | 74.31 | 20.08 | 51.23 | 82.78 | 43.99 | 78.67 | **78.82** | 24.84 | 57.57 | (+13.0) | | noise | 73.67 | 36.37 | 66.27 | 75.64 | 25.41 | 50.43 | 73.84 | 36.58 | 68.55 | **76.51** | 25.56 | 51.13 | (+1.4) |
| | all | 40.45 | 43.49 | 38.22 | 32.35 | 35.02 | 38.63 | 43.07 | 44.33 | 41.64 | 35.95 | 37.33 | 41.25 | (+4.3) | | all | 15.32 | 32.27 | 21.58 | 15.25 | 31.79 | 21.11 | 17.05 | 32.26 | 20.98 | 18.21 | 24.77 | 22.13 | (+1.3) |

| | | SLDA | | | | | | RobOCLe$_{SLDA}$ (ours) | | | | | | | | | SLDA | | | | | | | RobOCLe$_{SLDA}$ (ours) | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | clean | 99.13 | 47.19 | 92.55 | 77.09 | 19.74 | 59.14 | **99.78** | 48.69 | 93.31 | 80.17 | 20.32 | 60.63 | (+3.6) | | clean | 98.69 | 38.30 | 92.94 | 90.09 | 21.49 | 60.71 | **99.73** | 40.69 | 95.04 | 93.03 | 22.53 | 62.82 | (+5.4) |
| | illum | 86.68 | 79.65 | 77.76 | 61.88 | 39.41 | 66.43 | 89.02 | **82.21** | 79.29 | 64.58 | 40.89 | 68.45 | (+6.0) | | illum | 83.68 | 77.13 | 72.86 | 72.39 | 50.58 | 69.88 | 86.13 | **80.65** | 77.00 | 76.19 | 52.56 | 72.97 | (+10.3) |
| | geom | 98.05 | 49.50 | 97.53 | 78.08 | 22.59 | 62.06 | 99.81 | 51.32 | **98.40** | 78.05 | 23.12 | 62.83 | (+2.0) | | geom | 97.72 | 46.72 | 96.86 | 89.44 | 27.98 | 65.46 | 98.74 | 45.87 | **98.10** | 92.25 | 28.08 | 66.23 | (+2.2) |
| | noise | 97.81 | 50.14 | 90.12 | 96.83 | 28.12 | 66.55 | 98.72 | 51.82 | 92.85 | **98.12** | 28.57 | 67.99 | (+4.3) | | noise | 97.42 | 40.08 | 89.87 | 97.19 | 26.21 | 63.40 | 98.82 | 41.90 | 93.41 | **98.68** | 26.69 | 65.21 | (+4.9) |
| | all | 70.16 | 68.79 | 70.85 | 66.35 | 62.60 | 69.04 | 74.78 | 71.70 | 74.88 | 68.73 | **65.31** | 72.52 | (+11.3) | | all | 71.27 | 68.15 | 66.03 | 68.67 | 65.09 | 68.53 | 72.81 | 71.36 | 71.58 | 71.59 | **68.43** | 71.84 | (+10.5) |

TABLE IV

ACCURACY ON OTHER-DOMAIN DATA WITH CONTROLLED AUGMENTATIONS ON RN152 ON THE FEW-SHOT BENCHMARKS.

| | | SLDA | | | | | | RobOCLe$_{SLDA}$ (ours) | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | tr↓ te→ | clean | illum | geom | noise | all | Avg$_{OD}$ | clean | illum | geom | noise | all | Avg$_{OD}$ | $\Delta_R^{OD}$ |
| OLORIS-small | clean | 50.41 | 20.06 | 44.77 | 32.80 | 10.42 | 27.01 | **52.42** | 20.23 | 44.27 | 33.56 | 10.29 | 27.09 | (+0.1) |
| | illum | 30.24 | 27.71 | 20.88 | 19.28 | 12.99 | 20.85 | 32.36 | **28.75** | 24.58 | 17.26 | 12.56 | 21.69 | (+1.8) |
| | geom | 51.01 | 20.14 | 49.61 | 33.76 | 10.11 | 28.76 | 52.51 | 20.71 | **49.75** | 35.82 | 11.13 | 30.04 | (+1.8) |
| | noise | 46.29 | 18.83 | 39.03 | 42.31 | 11.35 | 28.87 | 46.69 | 19.63 | 40.16 | **43.31** | 11.04 | 29.38 | (+0.7) |
| | all | 23.63 | 26.08 | 22.27 | 21.40 | 24.49 | 23.34 | 23.85 | 28.75 | 22.88 | 23.73 | **25.80** | 24.81 | (+1.9) |
| FSIOL-310-5s | clean | 93.99 | 36.86 | 78.89 | 58.50 | 26.08 | 50.08 | **95.82** | 41.80 | 80.20 | 64.28 | 27.97 | 53.46 | (+6.8) |
| | illum | 41.50 | 60.33 | 39.48 | 33.40 | 21.11 | 33.87 | 41.70 | **62.48** | 40.82 | 32.52 | 22.88 | 34.48 | (+0.9) |
| | geom | 79.74 | 34.71 | 83.59 | 40.98 | 26.54 | 45.49 | 79.93 | 43.53 | **84.38** | 43.53 | 30.07 | 49.26 | (+6.9) |
| | noise | 72.55 | 34.44 | 52.68 | 83.66 | 24.12 | 45.95 | 75.69 | 38.10 | 58.56 | **85.36** | 25.69 | 49.51 | (+6.6) |
| | all | 24.25 | 25.49 | 21.31 | 17.39 | 19.61 | 22.11 | 26.27 | 26.86 | 19.67 | 20.59 | **22.03** | 23.35 | (+1.6) |
| FSIOL-310-10s | clean | 96.25 | 40.17 | 85.50 | 62.75 | 29.33 | 54.44 | **97.58** | 49.08 | 87.25 | 66.75 | 32.92 | 59.00 | (+10.0) |
| | illum | 50.33 | 68.83 | 39.48 | 36.08 | 28.08 | 38.40 | 53.50 | **69.42** | 39.67 | 39.17 | 29.50 | 40.46 | (+3.3) |
| | geom | 88.33 | 43.25 | 90.92 | 42.25 | 32.58 | 51.60 | 91.08 | 49.92 | **92.17** | 52.17 | 33.50 | 56.67 | (+10.5) |
| | noise | 81.92 | 36.58 | 61.00 | 92.08 | 24.08 | 50.90 | 86.67 | 44.17 | 65.83 | **94.92** | 27.50 | 56.04 | (+10.5) |
| | all | 27.25 | 24.42 | 27.67 | 22.42 | 43.17 | 25.44 | 28.58 | 25.17 | 29.83 | 23.00 | **43.83** | 26.65 | (+1.6) |

in the F-SIOL-310 dataset using lower labelled samples. Remarkably, the gains on the other-domain datasets reported in this paragraph are even larger than the gains reported on same-domain data on OpenLORIS in Tab. I.

Overall, these analyses certify the robustness of our method to both low-shot labelled training samples (outperforming competitors with as little as 5-shots) and datasets (improving on both the OpenLORIS-small and F-SIOL-310 with different objects and/or conditions seen at test time).

In Fig. 5, we also show robustness to per-step accuracy. RobOCLe consistently outperforms the baselines at every step, making it particularly suitable for real-world applications where the number of classes is not defined *a priori*.

**Controlled augmentations on other-domain data (OpenLORIS, OpenLORIS-small, F-SIOL-310).** We anticipated in Sec. IV and Fig. 3 that the existing datasets are limited in having either partial (OpenLORIS) or no (F-SIOL-310) different conditions at test time. To cope with this limitation, we apply further augmentations at train and test time.

First, we perform experiments on OpenLORIS using the

L32), thus increasing the size of the covariance matrix. Therefore, on few-shot data, both NCM and SLDA can achieve competitive results. Our RobOCLe improves both NCM and SLDA in almost every backbone and dataset, with just 2 exceptions out of 42 scenarios. The average gain $\Delta_R$ ranges from 1.8% to 32.6%, with the highest gains shown

TABLE V

MULTIPLE POOLING MECHANISMS ON OPENLORIS-SMALL AND SLDA.
$\Delta_f$: THE AVERAGE FEATURE SIZE MULTIPLIER COMPARED TO AVG.

| | $\Delta_f$ | Accuracy | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | RN50 | RN101 | RN152 | ViT-B16 | ViT-B32 | ViT-L16 | ViT-L32 | Avg |
| AVG [50] | 1.0 | 50.26 | 49.91 | 50.41 | 43.96 | 41.50 | 45.51 | 42.93 | 46.35 |
| [CLS] [69] | 1.0 | - | - | - | 43.10 | 40.78 | 45.57 | 41.83 | - |
| MAX [51] | 1.0 | 50.87 | 50.01 | 50.75 | 41.27 | 38.67 | 43.23 | 39.79 | 44.94 |
| AVGMAX [55] | 2.0 | 50.74 | 50.16 | 50.51 | 42.21 | 39.49 | 44.42 | 40.19 | 45.39 |
| MIX (50%) [52] | 1.0 | 50.04 | 49.48 | 50.14 | 43.07 | 40.86 | 45.15 | 40.84 | 45.65 |
| STOCHASTIC [53] | 1.0 | 42.15 | 40.37 | 38.48 | 32.78 | 29.45 | 31.97 | 30.68 | 35.13 |
| L2 [56] | 1.0 | 44.12 | 43.71 | 44.09 | 35.14 | 30.88 | 36.19 | 32.27 | 38.06 |
| L3 [56] | 1.0 | 45.12 | 43.15 | 43.80 | 36.44 | 31.11 | 37.86 | 32.57 | 38.58 |
| RAP (1%) [54] | 2.8 | 50.39 | 50.14 | 50.01 | 43.11 | 41.98 | 44.48 | 42.61 | 46.10 |
| RAP (10%) [54] | 139.0 | 39.48 | 32.14 | 31.15 | 27.48 | 24.41 | 26.10 | 24.13 | 29.27 |
| iSQRT-COV [57] | 1.0 | 50.34 | 48.50 | 48.75 | 42.98 | 41.66 | **45.75** | 42.12 | 45.73 |
| RobOCLe$_{\text{SLDA}}$ (ours) | 3.0 | **51.33** | **51.44** | **52.42** | **44.73** | **42.86** | 45.22 | **43.07** | **47.29** |

TABLE VI

ADDITIONAL METRICS ON THE OPENLORIS-SMALL AND RESNET152.
TTIME: TRAINING TIME [MIN]. FPS: FRAMES PER SEC AT TEST TIME.

| | Acc ↑ | BwT ↑ | Forg ↓ | Pla ↑ | TTime ↓ | FPS ↑ |
|---|---|---|---|---|---|---|
| FT | 16.05 | 23.37 | 78.77 | 92.82 | 5.57 | 165.8 |
| NCM | 47.68 | 57.34 | 46.54 | 61.10 | 3.93 | 164.8 |
| RobOCLe$_{\text{NCM}}$ (ours) | 54.89 (+13.8%) | 63.85 | 39.95 | 69.18 | 4.24 | 164.2 (-0.36%) |
| SLDA | 50.41 | 59.14 | 45.69 | 66.19 | 4.71 | 165.0 |
| RobOCLe$_{\text{SLDA}}$ (ours) | 52.42 (+4.1%) | 60.75 | 43.54 | 70.38 | 7.14 | 164.2 (-0.45%) |

best CNN (RN152) and the best transformer (ViT-L16) on the two best OCL methods (NCM and SLDA) with variable train and test augmentations. Results are reported in Tab. III. For each sub-table, we report train augmentations on the rows and test augmentation on the columns, and we compute the average other-domain accuracy (Avg OD). For each group of results (RobOCLe vs. baseline), we also report the RARG of RobOCLe in terms of Avg OD. On each scenario, RobOCLe shows superior results on severe and controlled synthetic augmentations when applied in training and/or in testing stages. Employing either color, geometric or noise augmentation generally improves Avg OD. Their combination improves Avg OD when using SLDA only, as it has more capacity than NCM to capture and disentangle input-level domain variations. When non-clean same augmentation is applied to both train and test sets, instead, we obtain lower accuracy compared to the original accuracy obtained using clean train/test sets. This is due to the tasks becoming much harder as, in the OCL setup, the agent experiences data only once. Therefore, convergence of models to learn robust representations is hindered. Also in this case, we confirm how RobOCLe is more effective than the baseline counterpart, especially on other-domain setups.

Next, we select RN152, and analyse the effect of controlled augmentations on the few-shot datasets using SLDA in Tab. IV. For each dataset and for each training augmentation, our RobOCLe robustly outperforms the baseline by a large margin on both same-domain and other-domain data.

**Ablation studies.** So far, we presented evaluation of results on the basis of 1) mostly replay-free OCL methods, and 2) accuracy. Here, we complement the analyses. In Tab. V, we report accuracy of other pooling methods described in Sec. III. We observe that [CLS] pooling performs slightly worse than AVG, due to the limited input data. Concatenating the top-k% features (RAP) increases the feature size but

shows lower accuracy compared to the proposed method: using larger feature size is not enough to improve accuracy and robustness. All other approaches do not achieve competitive results. Our method slightly increases the feature size, however, it brings robust recognition results.

As we introduced in Sec. IV, multiple metrics have been defined to characterize performance of CL agents. Tab. VI summarizes the main results obtained on the OpenLORIS-small. RobOCLe improves performance with respect to baseline competitors with higher Acc (overall final performance), BwT (knowledge transfer to past classes), and lower Forg (average forgetting of all classes), by finding a better trade-off between stability and Pla (which is lower than, *e.g.*, FT). For completeness, FwT is always 0, since incremental steps contain disjoint classes. FT implements the classifiers as a fully-connected layer and shows longer TTime than NCM and SLDA, which only requires computation of distance between embedded features and class prototypes. RobOCLe adds some computational overhead which is mostly affecting TTime, while inference FPS remains practically unchanged with a small decrease of less than 0.5% in all cases. On the other hand, the slight increase of TTime is not a relevant concern, since training is done on few-shot data only.

## VI. CONCLUSION

In this paper, we tackled the practical task of FS-OCL targeting robust test-time object recognition for low-resource robots with limited labelled data and computational/storage capability. We introduced RobOCLe that promotes invariance to augmentation via high order statistical moments of the embedded features of input samples. We proved that RobOCLe achieves robust recognition in a variety of scenarios, using several backbones, low-shot setups, per-step accuracy, and controlled train/test augmentation on both same-domain and other-domain data. Overall, FS-OCL task in low-resource devices is far from being solved. We hope that our problem formulation, approach and extensive comparison with previous methods will encourage future works on this direction.

## REFERENCES

[1] N. Churamani, S. Kalkan, and H. Gunes, "Continual learning for affective robotics: Why, what and how?" in *RO-MAN*, 2020.

[2] T. Lesort, V. Lomonaco, A. Stoian, D. Maltoni, D. Filliat, and N. Díaz-Rodríguez, "Continual learning for robotics: Definition, framework, learning strategies, opportunities and challenges," *Inf. fus.*, 2020.

[3] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska, D. Hassabis, C. Clopath, D. Kumaran, and R. Hadsell, "Overcoming catastrophic forgetting in neural networks," *PNAS*, vol. 114, 2017.

[4] T. L. Hayes and C. Kanan, "Online Continual Learning for Embedded Devices," *CoLLAs-2022*, 2022.

[5] L. Pellegrini, V. Lomonaco, G. Graffieti, and D. Maltoni, "Continual learning at the edge: Real-time training on smartphone devices," *arXiv:2105.13127*, 2021.

[6] M. Knauer, M. Denninger, and R. Triebel, "Recall: Rehearsal-free continual learning for object classification," in *IROS*, 2022.

[7] Z. Mai, R. Li, J. Jeong, D. Quispe, H. Kim, and S. Sanner, "On-line continual learning in image classification: An empirical survey," *Neurocomputing*, vol. 469, pp. 28–51, 2022.

[8] G. Borghi, G. Graffieti, and D. Maltoni, "On the challenges to learn from natural data streams," *arXiv:2301.03495*, 2023.

[9] J. Snell, K. Swersky, and R. Zemel, "Prototypical networks for few-shot learning," *NeurIPS*, vol. 30, 2017.

[10] M. Lunayach, J. Smith, and Z. Kira, "Lifelong wandering: A realistic few-shot online continual learning setting," *CVPRW*, 2022.

[11] J. Cohen, E. Rosenfeld, and Z. Kolter, "Certified adversarial robustness via randomized smoothing," in *ICML*. PMLR, 2019, pp. 1310–1320.

[12] M. M. Zhang, S. Levine, and C. Finn, "MEMO: Test Time Robustness via Adaptation and Augmentation," in *NeurIPS*, 2022.

[13] Q. Wang, O. Fink, L. Van Gool, and D. Dai, "Continual test-time domain adaptation," in *CVPR*, 2022, pp. 7201–7211.

[14] M. Toldo, U. Michieli, and P. Zanuttigh, "Learning with style: Continual semantic segmentation across tasks and domains," *arXiv:2210.07016*, 2022.

[15] F. Zenke, B. Poole, and S. Ganguli, "Continual learning through synaptic intelligence," in *ICML*. PMLR, 2017.

[16] M. De Lange, R. Aljundi, M. Masana, S. Parisot, X. Jia, A. Leonardis, G. Slabaugh, and T. Tuytelaars, "A continual learning survey: Defying forgetting in classification tasks," *TPAMI*, 2021.

[17] D. Lopez-Paz and M. Ranzato, "Gradient episodic memory for continual learning," *NeurIPS*, vol. 30, 2017.

[18] T. L. Hayes, N. D. Cahill, and C. Kanan, "Memory efficient experience replay for streaming learning," in *ICRA*. IEEE, 2019, pp. 9769–9776.

[19] Z. Mai, R. Li, H. Kim, and S. Sanner, "Supervised contrastive replay: Revisiting the nearest class mean classifier in online class-incremental continual learning," in *CVPRW*, 2021.

[20] T. L. Hayes and C. Kanan, "Lifelong machine learning with deep streaming linear discriminant analysis," in *CVPRW*, 2020.

[21] C. R. Wolfe and A. Kyrillidis, "Cold start streaming learning for deep networks," *arXiv:2211.04624*, 2022.

[22] L. Pellegrini, G. Graffieti, V. Lomonaco, and D. Maltoni, "Latent replay for real-time continual learning," in *IROS*. IEEE, 2020.

[23] L. Ravaglia, M. Rusci, D. Nadalini, A. Capotondi, F. Conti, and L. Benini, "A tinyml platform for on-device continual learning with quantized latent replays," *JESTCS*, vol. 11, no. 4, pp. 789–802, 2021.

[24] G. Demosthenous and V. Vassiliades, "Continual learning on the edge with tensorflow lite," *arXiv:2105.01946*, 2021.

[25] A. Ayub and A. R. Wagner, "Cognitively-inspired model for incremental learning using a few examples," in *CVPRW*, 2020, pp. 222–223.

[26] ——, "F-SIOL-310: A Robotic Dataset and Benchmark for Few-Shot Incremental Object Learning," in *ICRA*. IEEE, 2021.

[27] ——, "Tell me what this is: Few-shot incremental object learning by a robot," in *IROS*. IEEE, 2020, pp. 8344–8350.

[28] P. Mazumder, P. Singh, and P. Rai, "Few-shot lifelong learning," in *AAAI*, vol. 35, no. 3, 2021, pp. 2337–2345.

[29] J. Von Oswald, D. Zhao, S. Kobayashi, S. Schug, M. Caccia, N. Zucchet, and J. Sacramento, "Learning where to learn: Gradient sparsity in meta and continual learning," *NeurIPS*, vol. 34, 2021.

[30] G. van de Ven, T. Tuytelaars, and A. Tolias, "Three types of incremental learning," in *NMI*, 2022.

[31] A. Chaudhry, P. K. Dokania, T. Ajanthan, and P. H. Torr, "Riemannian walk for incremental learning: Understanding forgetting and intransigence," in *ECCV*, 2018, pp. 532–547.

[32] D. Shim, Z. Mai, J. Jeong, S. Sanner, H. Kim, and J. Jang, "Online class-incremental continual learning with adversarial shapley value," in *AAAI*, vol. 35, no. 11, 2021, pp. 9630–9638.

[33] A. Chaudhry, M. Ranzato, M. Rohrbach, and M. Elhoseiny, "Efficient lifelong learning with a-gem," *ICLR*, 2019.

[34] M. Wang and W. Deng, "Deep visual domain adaptation: A survey," *Neurocomputing*, vol. 312, pp. 135–153, 2018.

[35] M. Toldo, A. Maracani, U. Michieli, and P. Zanuttigh, "Unsupervised domain adaptation in semantic segmentation: a review," *Technologies*, vol. 8, no. 2, p. 35, 2020.

[36] D. Shanmugam, D. Blalock, G. Balakrishnan, and J. Guttag, "Better aggregation in test-time augmentation," in *ICCV*, 2021.

[37] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," *Communications of the ACM*, vol. 60, no. 6, pp. 84–90, 2017.

[38] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016, pp. 770–778.

[39] H. Lu, D. Shanmugam, H. Suresh, and J. Guttag, "Improved text classification via test-time augmentation," *arXiv:2206.13607*, 2022.

[40] A. Ashukha, A. Atanov, and D. Vetrov, "Mean embeddings with test-time data augmentation for ensembling of representations," *arXiv:2106.08038*, 2021.

[41] S. Chun, J. Y. Lee, and J. Kim, "Cyclic test time augmentation with entropy weight method," in *UAI*. PMLR, 2022, pp. 433–442.

[42] G. Wang, W. Li, M. Aertsen, J. Deprest, S. Ourselin, and T. Vercauteren, "Aleatoric uncertainty estimation with test-time augmentation for medical image segmentation with convolutional neural networks," *Neurocomputing*, vol. 338, pp. 34–45, 2019.

[43] P. Conde, C. Premebida, and P. Coimbra, "Adaptive-tta: accuracy-consistent weighted test time augmentation method for the uncertainty calibration of deep learning classifiers," in *BMVC*, 2022.

[44] I. Kim, Y. Kim, and S. Kim, "Learning loss for test-time augmentation," *NeurIPS*, vol. 33, pp. 4163–4174, 2020.

[45] S. Enomoto, M. R. Busto, and T. Eda, "Augnet: Dynamic test-time augmentation via differentiable functions," *arXiv:2212.04681*, 2022.

[46] S. Goyal, M. Sun, A. Raghunathan, and Z. Kolter, "Test-time adaptation via conjugate pseudo-labels," *arXiv:2207.09640*, 2022.

[47] M. J. Mirza, P. J. Soneira, W. Lin, M. Kozinski, H. Possegger, and H. Bischof, "Actmad: Activation matching to align distributions for test-time-training," *arXiv:2211.12870*, 2022.

[48] Y. Su, X. Xu, and K. Jia, "Revisiting realistic test-time training: Sequential inference and adaptation by anchored clustering," in *NeurIPS*, 2022.

[49] C. Oriet and K. Hozempa, "Incidental statistical summary representation over time," *Journal of Vision*, 2016.

[50] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *IEEE*, 1998.

[51] M. Riesenhuber and T. Poggio, "Hierarchical models of object recognition in cortex," *Nature NN*, vol. 2, no. 11, pp. 1019–1025, 1999.

[52] Q. Zhou, Z. Qu, and C. Cao, "Mixed pooling and richer attention feature fusion for crack detection," *PRL*, vol. 145, pp. 96–102, 2021.

[53] M. D. Zeiler and R. Fergus, "Stochastic pooling for regularization of deep convolutional neural networks," *ICLR*, 2013.

[54] S. Bera and V. K. Shrivastava, "Effect of pooling strategy on convolutional neural network for classification of hyperspectral remote sensing images," *IET IP*, vol. 14, no. 3, pp. 480–486, 2020.

[55] J. Monteiro, M. J. Alam, and T. Falk, "On the performance of time-pooling strategies for end-to-end spoken language identification," in *LREC*, 2020, pp. 3566–3572.

[56] J. Feng, B. Ni, Q. Tian, and S. Yan, "Geometric $l_p$-norm feature pooling for image classification," in *CVPR*, 2011.

[57] P. Li, J. Xie, Q. Wang, and Z. Gao, "Towards faster training of global covariance pooling networks by iterative matrix square root normalization," in *CVPR*, 2018, pp. 947–955.

[58] U. Michieli, P. P. Parada, and M. Ozay, "Online continual learning in keyword spotting for low-resource devices via pooling high-order temporal statistics," *INTERSPEECH*, 2023.

[59] S.-A. Rebuffi, A. Kolesnikov, G. Sperl, and C. H. Lampert, "iCaRL: Incremental classifier and representation learning," in *CVPR*, 2017.

[60] T. Mensink, J. Verbeek, F. Perronnin, and G. Csurka, "Distance-based image classification: Generalizing to new classes at near-zero cost," *TPAMI*, vol. 35, no. 11, pp. 2624–2637, 2013.

[61] S. Dasgupta and D. Hsu, "On-line estimation with the multivariate Gaussian distribution," in *ICCLT*. Springer, 2007, pp. 278–292.

[62] C. Anagnostopoulos, D. K. Tasoulis, N. M. Adams, N. G. Pavlidis, and D. J. Hand, "Online linear and quadratic discriminant analysis with adaptive forgetting for streaming classification," *SADM*, 2012.

[63] B. Welford, "Note on a method for calculating corrected sums of squares and products," *Technometrics*, vol. 4, no. 3, pp. 419–420, 1962.

[64] Q. She, F. Feng, X. Hao, Q. Yang, C. Lan, V. Lomonaco, X. Shi, Z. Wang, Y. Guo, Y. Zhang, F. Qiao, and R. H. M. Chan, "OpenLORIS-Object: A robotic vision dataset and benchmark for lifelong deep learning," in *ICRA*. IEEE, 2020, pp. 4767–4773.

[65] A. Howard, M. Sandler, G. Chu, L.-C. Chen, B. Chen, M. Tan, W. Wang, Y. Zhu, R. Pang, V. Vasudevan, Q. V. Le, and H. Adam, "Searching for MobileNetv3," in *ICCV*, 2019, pp. 1314–1324.

[66] M. Tan and Q. Le, "EfficientNet: Rethinking model scaling for convolutional neural networks," in *ICML*, 2019, pp. 6105–6114.

[67] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," in *ICCV*, 2021, pp. 10012–10022.

[68] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," in *ICLR*, 2021.

[69] J. D. M.-W. C. Kenton and L. K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *NAACL-HLT*, 2019, pp. 4171–4186.