

Latent Replay for Real-Time Continual Learning

Lorenzo Pellegrini¹, Gabriele Graffieti¹, Vincenzo Lomonaco¹ and Davide Maltoni¹.

Abstract—Training deep neural networks at the edge on light computational devices, embedded systems and robotic platforms is nowadays very challenging. Continual learning techniques, where complex models are incrementally trained on small batches of new data, can make the learning problem tractable even for CPU-only embedded devices enabling remarkable levels of adaptiveness and autonomy. However, a number of practical problems need to be solved: catastrophic forgetting before anything else. In this paper we introduce an original technique named “Latent Replay” where, instead of storing a portion of past data in the input space, we store activations volumes at some intermediate layer. This can significantly reduce the computation and storage required by native rehearsal. To keep the representation stable and the stored activations valid we propose to slow-down learning at all the layers below the latent replay one, leaving the layers above free to learn at full pace. In our experiments we show that Latent Replay, combined with existing continual learning techniques, achieves state-of-the-art performance on complex video benchmarks such as CORE50 NICv2 (with nearly 400 small and highly non-i.i.d. batches) and OpenLORIS. Finally, we demonstrate the feasibility of nearly real-time continual learning on the edge through the deployment of the proposed technique on a smartphone device.

I. INTRODUCTION

Training on the edge (e.g., on light computing devices such as smartphones, smart cameras, embedded systems and robotic platforms) is highly desirable in several applications where privacy, lack of network connection and fast adaptation are real constraints. While some steps in this direction have been recently moved [1], training on the edge often remains unfeasible. In fact, given the high demand in terms of memory and computation, most machine learning models nowadays are trained on powerful multi-GPUs servers, and only frozen models are deployed to edge devices for inference.

Furthermore, in some applications (e.g., robotic vision, see Fig. 1), training a deep model from scratch as soon as new data becomes available is prohibitive in terms of storage / computation even if performed server side. Continual Learning (CL), that is the ability of continually training existing models using only new data, is gaining a lot of attention and several solutions have been recently proposed to deal with the daunting issue of catastrophic forgetting (i.e., as the model learns new concepts and skills, it tends to forget the old ones) [2]. Recent surveys [3], [4] provide an overview of the CL field. In principle, CL approaches could be exploited not only to control forgetting but also to reduce the training complexity.

¹Department of Computer Science and Engineering (DISI), University of Bologna, Via dell’università, 50, 47521 Cesena FC, Italy. E-mails: {l.pellegrini, gabriele.graffieti, vincenzo.lomonaco, davide.maltoni}@unibo.it

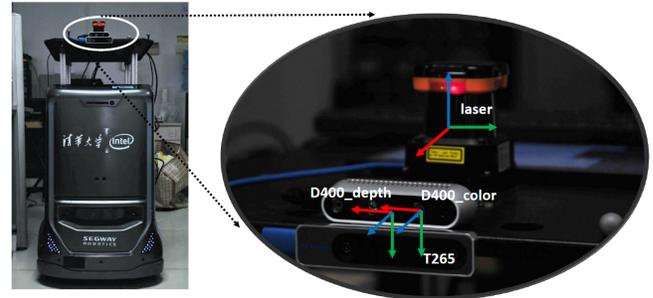


Fig. 1: Wheeled robot used in the *Lifelong Robotic Vision* challenge at IROS 2019 [5] (left) and equipped with multiple sensors including two Real Sense cameras (right).

In this paper we focus on real-time CL and prove that continual training with small batches can be compatible with the limited computing power made available by CPU-only embedded devices and robotic platforms.

In [6] it was shown that some CL approaches can effectively learn to recognize objects (on the CORE50 dataset [7]) even when fed with fine-grained incremental batches. CORE50 NICv2 [6] is a continual learning benchmark where objects from 50 different classes have to be learned incrementally. What makes this benchmark challenging is that classes are discovered a little at a time and the training batches are small and non i.i.d. In particular, in NICv2 - 391 each training batch includes only 300 frames extracted from a short video (15 seconds at 20 fps) of a single object slowly moving in front of the camera: hence, patterns within each batch are highly correlated. Despite these nuisances, in [6] the approaches denoted as CWR* and AR1* proved to be able to learn continually even in absence of replay mechanisms, that is, the periodic refresh of old examples maintained in an external memory. Under this challenging setting both CWR* and AR1* performed significantly better than well known techniques such as LWF [8] and EWC [9]. While these results are encouraging:

- the accuracy gap w.r.t. the cumulative approach (a sort of upper bound obtained by training the model on the entire training set) remains quite relevant (about 20%).
- in NICv2 - 391, the most challenging setup, AR1* was not able to effectively adapt the representation layers during continual learning.

The aim of this work is to *reduce as much as possible the gap w.r.t. the cumulative upper bound* and, at the same time, to *provide an efficient implementation strategy of CL approaches to enable nearly real-time training on the edge*.

To this purpose we first show that a small amount of pattern replay is sufficient to significantly improve accuracy on NICv2 - 391 (Section III). However, even if in the CORE50 setting the extra memory required by replay is not an issue (we store only 30 patterns for each of the 50 classes), a constant refresh significantly increases the required amount of computation because of the extra forward and backward steps and this makes the resulting training too resource demanding for real-time applications. Therefore, we propose a “*Latent Replay*” approach (Section IV) where old data are injected at some intermediate layer selected according to the desired accuracy-efficiency trade-off.

In Section V we compare our approach with several continual learning algorithms and show its advantages reaching state-of-the-art performances on two different benchmarks: CORE50 and OpenLORIS. Finally, in order to demonstrate the practical applicability of the proposed approach, in Section VI we discuss the implementation of a continual learning application for Android smartphones that, starting from a pre-trained MobileNetV1 model with 10 classes, can incrementally learn (in near real-time) new classes and/or new objects of existing classes.

II. RELATED WORKS

While several works in the continual learning literature focus on Multiple independent Tasks (MT) scenario, in many practical applications such as robotic vision, a Single Incremental Task (SIT) scenario is more appropriate [10]. In particular, a robot should be able to incrementally improve its object recognition capabilities while being exposed to new instances of both known and completely new classes (denoted as NIC setting - New Instances and Classes). CORE50 NICv2 benchmark specifically addresses this problem [6]. Other datasets have been released to study continual learning for robotic vision (e.g., iCub-transformation [11], OpenLORIS [12]) but no NIC benchmarks have been yet defined for them. ImageNet-1K [13] and CIFAR-100 [14] have also been used to evaluate continual learning techniques, but these datasets do not fit well the object recognition task because of the lack of multiple videos of the same objects taken under different poses, lighting and backgrounds.

In [6], two approaches denoted as CWR* and AR1* have been evaluated on CORE50 NICv2: in CWR* the last fully connected layer is implemented as a double memory, and simple initialization and fusions steps are performed before and after each training batch to synchronize the two memories. However, after the first training batch, CWR* freezes all the layers except the last one, thus losing the benefits of a continual adaptation of the underlying representation. AR1* extends CWR* by enabling end-to-end continual training throughout the entire network; to this purpose the Synaptic Intelligence [15] regularization approach (similar to Online-EWC [16]) is adopted to constrain the change of critical weights.

Patterns replay, which is central in the proposed approach, proved to be an effective approach to contrast forgetting in continual learning scenarios [17], [18], [19], [20]. In fact,

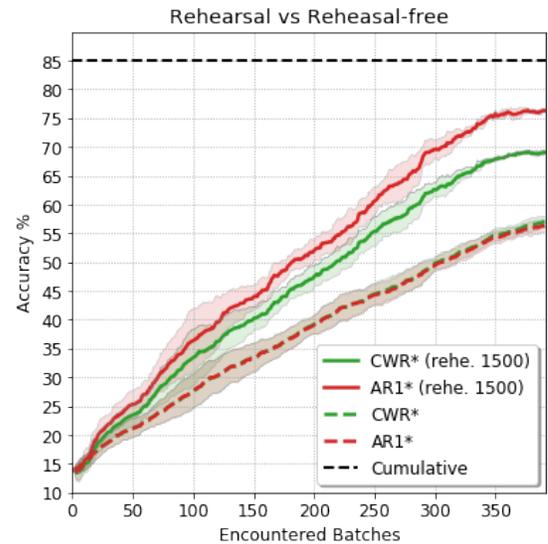


Fig. 2: Comparison of CWR* and AR1* on CORE50 NICv2 - 391 with and without rehearsal ($RM_{size} = 1500$). Each experiment was averaged on 5 runs with different batch ordering: colored areas represent the standard deviation of each curve. The black dashed line denotes the reference accuracy of the cumulative upper bound.

periodically replaying some representative patterns from old data helps the model to retain important information of past tasks / classes while learning new concepts. iCaRL [17] uses well-designed entry / exit criteria (denoted as herding) to maintain a class-balanced set of exemplars that maximize representativeness. A comparison between the proposed technique and iCaRL is reported in Section V-A. Generative Replay (also known as “*Pseudo-rehearsal*” [21]), where surrogates of past data are generated without explicitly storing native patterns, looks very appealing because of the storage saving; however, most of the proposed approaches to date do not allow on-line generation of effective replay patterns.

Another class of relevant techniques for our study are the so called *streaming* continual learning approaches [22], where a model can be incrementally trained with a single pattern at a time. Even if in a robotic vision scenario learning from single frames does not appear necessary (in fact, using short videos of single objects can be more efficient and looks more biologically plausible) efficient streaming learning techniques can be effortlessly applied to NIC setting. Deep Streaming Linear Discriminant Analysis (DSLDA) was recently proposed [23] where an online extension of the LDA classifier works on the top of a fixed deep learning feature extractor. This approach, which achieved state-of-the-art accuracy on (partitioned) ImageNet-1K and CORE50 (10 classes version) was run on NICv2 and compared with other techniques in [6], showing competitive results.

Finally, training on the edge was recently addressed in [1], where an object detection model was incrementally trained based on LWF and pattern replay. While training is not actually real-time (it requires a few minutes on Nvidia Jetson TX2 board) and only few large continual training batches are

presented to the model, the detection problem approached in [1] is more difficult than the classification problem here considered and therefore we cannot make a direct comparison.

III. NATIVE REHEARSAL

In [10] it was shown that a very simple rehearsal implementation (hereafter denoted as *native rehearsal*), where for every training batch a random subset of the batch patterns is added to the external storage to replace a (equally random) subset of the external memory, is not less effective than more sophisticated approaches such as iCaRL. Therefore, in this study, we started by expanding CWR* and AR1* with the trivial rehearsal approach summarized in Algorithm 1. In Figure 2 we compare the learning trend of CWR* and AR1* of a MobileNetV1¹ trained with and without rehearsal on CORe50 NICv2 - 391. We use the same protocol and hyper-parameters introduced in [6] and a rehearsal memory of 1,500 patterns. It is well evident that even a moderate external memory (about 1.27% of the total training set) is very effective to improve the accuracy of both approaches and to reduce the gap with the cumulative upper bound that for this model is $\sim 85\%$.

Algorithm 1 Pseudocode describing how the external memory RM is populated across the training batches. Note that the amount h of patterns to add progressively decreases to maintain a nearly balanced contribution from the different training batches, but no constraints are enforced to achieve a class-balancing.

- 1: $RM = \emptyset$
- 2: $RM_{size} =$ number of patterns to be stored in RM
- 3: **for each** training batch B_i :
- 4: train the model on shuffled $B_i \cup RM$
- 5: $h = \frac{RM_{size}}{i}$
- 6: $R_{add} =$ random sampling h patterns from B_i
- 7: $R_{replace} = \begin{cases} \emptyset & \text{if } i == 1 \\ \text{random sample } h \text{ patterns from } RM & \text{otherwise} \end{cases}$
- 8: $RM = (RM - R_{replace}) \cup R_{add}$

To understand the influence of the external memory size we repeated the experiment with different RM_{size} values: 500, 1,000, 1,500, 3,000². Since rehearsal itself protects the model from forgetting we also run AR1* (where important weights of lower layers are protected from forgetting by using Synaptic Intelligence [15] regularization) without Synaptic Intelligence protection, that is lower layers weights are left totally unconstrained; Hereafter we denote this approach as AR1*free. We found that increasing the rehearsal memory leads to better accuracy for all the algorithms, but the gap between 1500 and 3000 is not large and we found 1,500 is a good trade-off for this dataset. AR1*free works slightly better than AR1* when a sufficient number of rehearsal patterns are provided but, as expected, accuracy is worse with light (i.e. $RM_{size} = 500$) or no rehearsal.

It is worth noting that the best combination (AR1*free with 3000 patterns) is only 5% worse than the cumulative

¹The network was pre-trained on ImageNet-1k.

²More details about this experiment can be found in the extended preprint version of this work [24].

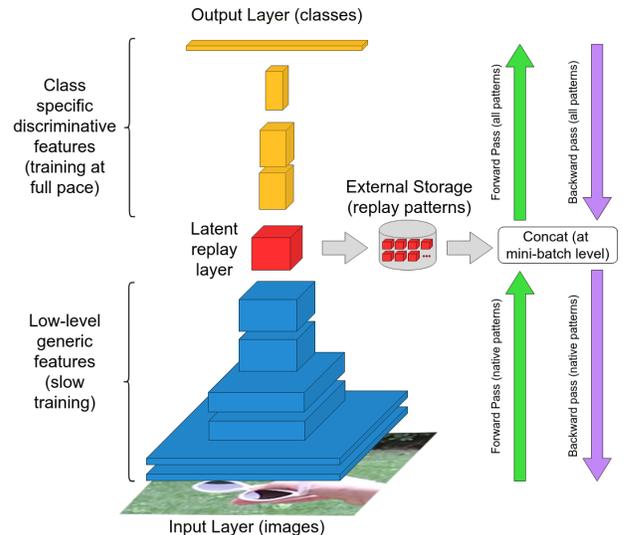


Fig. 3: Architectural diagram of Latent Replay.

upper bound and a better parametrization and exploitation of the rehearsal memory could further reduce this gap.

IV. LATENT REPLAY

In deep neural networks the layers close to the input (often denoted as representation layers) usually perform low-level feature extraction and, after a proper pre-training on a large dataset (e.g., ImageNet), their weights are quite stable and reusable across applications. On the other hand, higher layers tend to extract class-specific discriminant features and their tuning is often important to maximize accuracy.

With latent replay (see Figure 3) we denote an approach where, instead of maintaining copies of input patterns in the external memory in the form of raw data, we store the activations volumes at a given layer (denoted as *Latent Replay layer*). To keep the representation stable and the stored activations valid we propose to slow-down the learning at all the layers below the latent replay one and to leave the layers above free to learn at full pace. In the limit case where lower layers are completely frozen (i.e., slow-down to 0) latent replay is functionally equivalent to rehearsal from the input, but achieves a computational and storage saving thanks to the smaller fraction of patterns that need to flow forward and backward across the entire network and the typical information compression that networks perform at higher layers.

In the general case where the representation layers are not completely frozen, the activations stored in the external memory suffer from an aging effect (i.e., as the time passes they tend to increasingly deviate from the activations that the same pattern would produce if feed-forwarded from the input layer). However, if the training of these layers is sufficiently slow, the aging effect is not disruptive since the external memory has enough time to be rejuvenated with fresh patterns. When latent replay is implemented with mini-batch SGD training: (i) in the forward step, a concatenation is performed at the replay layer (on the mini-batch dimension) to join

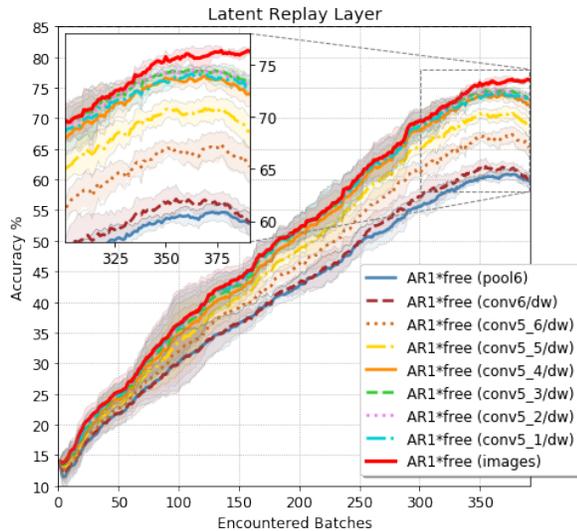


Fig. 4: AR1*free with latent replay ($RM_{size} = 1500$) for different choices of the latent replay layer. Setting the replay layer at the `pool6` layer makes AR1*free equivalent to CWR*. Setting the replay layer at the “images” layer corresponds to native rehearsal. The saturation effect which characterizes the last training batches is due to the data distribution in NICv2 - 391 (see [6]): in particular, the lack of new instances for some classes (that already introduced all their data) slows-down the accuracy trend and intensifies the effect of activations aging.

patterns coming from the input layer with activations coming from the external storage; (ii) the backward step is stopped just before the replay layer for the replay patterns.

V. EXPERIMENTS AND RESULTS

Hereafter, while the proposed latent replay approach is architecture agnostic, we discuss its specific design with several continual learning algorithms CWR*, AR1*, AR1*free and LWF over a MobileNet [25] pre-trained on ImageNet-1K. We compare our latent replay approach with other state-of-the-art techniques on CORE50 and OpenLORIS.

A. Experiments on CORE50 NICv2 - 391

For the CORE50 experiments:

- we use a MobileNetV1 and focus on CWR*, AR1* which have been already proved to be competitive on this benchmark.
- for all the methods the output layer (`fc7`) must be implemented as a double memory with proper (pre)initialization and (post)fusion for each training batch (for details see CWR* pseudocode in Algorithm 2 of [6]);
- for CWR* the latent replay layer is the second-last layer (i.e., `pool6`);
- for AR1* and AR1*free the latent replay layer can be pushed down and selected according to the accuracy-efficiency trade-off discussed below;
- for AR1*free the Synaptic Intelligence regularization is switched off.

To simplify the network design and training we keep the proportion of original and replay patterns fixed: for example, if the training batches contain 300 patterns and the external memory 1500 patterns, in a mini-batch of size 128 we concatenate 21 ($128 \times 300/1800$) original patterns (of the current batch) with 107 ($128 \times 1500/1800$) replay patterns. In this case only 21 patterns (over 128) need to travel across the blue layers in Figure 3.

Concerning the learning slow-down in the representation layers we found that an effective (and efficient) strategy is blocking the weight changes after the first batch (i.e., learning rate set to 0), while leaving the batch normalization moments free to adapt to the statistics of the input patterns across all the batches. Batch Normalization (BN) [26] is widely used in modern deep neural networks (including MobileNets) to control internal covariate shift thus making learning faster and more robust. Replacing BN with Batch Renormalization (BRN) [27] was proved to be a very important step for effective continual learning with fine-grained non-i.d.d. batches [6], so in the MobileNetV1 here adopted BN layers have been replaced with BRN layers. In the context of latent replay, if we leave the BRN moments free to adapt, the activations stored in the external memory suffer the aging effect described in Section IV. However, we experimentally verified that, upon proper setting of the global moment mobile windows (more details are provided in the additional material), the accuracy drop due to the aging effect is quite limited and in any case the final accuracy is higher w.r.t. the case where BRN moments in the representation layers are frozen. On the computational side, blocking the weight changes in the representation layers allows to skip the backward pass in the lower part of the network also for native patterns, since updating the BRN moments only relies on the forward pass.

In Figure 4 we report the accuracy of AR1*free with latent replay ($RM_{size} = 1500$) for different choices of the rehearsal layer (reported between parenthesis). As expected, when the replay layer is pushed down the corresponding accuracy increases, proving that a continual tuning of the representation layers is important. However, after `conv5_4/dw` there is a sort of saturation and the model accuracy is no longer improving. The residual gap ($\sim 4\%$) with respect to native rehearsal is not due to the weights freezing of the lower part of the network but to the aging effect introduced above. This can be simply proved by implementing an “intermediate” approach that always feeds the replay pattern from the input and stops the backward at `conv5_4`: such intermediate approach achieved an accuracy at the end of the training very close to the native rehearsal. We believe that the accuracy drop due to the aging effect can be further reduced with better tuning of BNR hyper-parameters and/or with the introduction of a scheduling policy making the global moment mobile windows wider as the continual learning progresses (i.e., more plasticity in the early stages and more stability later); however, such fine optimization is application specific and beyond the scope of this study.

TABLE I: Computation, storage, and accuracy trade-off with Latent Replay at different layers of a MobileNetV1 ConvNet trained continually on NICv2 - 391 with $RM_{size} = 1500$. Computation and pattern size can be easily extrapolated from Table 1 in the supplementary materials where the network architecture is exploded by reporting neurons, connections and weights at each layer.

Layer	Computation % vs Native Rehearsal	Pattern Size	Final Accuracy %	Δ Accuracy % vs Native Rehearsal
Images	100.00%	49152	77.30%	0.00%
conv5_1/dw	59.261%	32768	72.82%	-4.49%
conv5_2/dw	50.101%	32768	73.21%	-4.10%
conv5_3/dw	40.941%	32768	73.22%	-4.09%
conv5_4/dw	31.781%	32768	72.24%	-5.07%
conv5_5/dw	22.621%	32768	68.59%	-8.71%
conv5_6/dw	13.592%	8192	65.24%	-12.06%
conv6/dw	9.012%	16384	59.89%	-17.42%
pool6	0.027%	1024	59.76%	-17.55%

A.1 On the Computation, Storage and Accuracy Trade-off

To better evaluate the latent replay w.r.t. native rehearsal we report in Table I the relevant dimensions: (i) computation refers to the percentage cost in terms of ops of a partial forward (from the latent replay layer on) relative to a full forward step from the input layer; (ii) pattern size is the dimensionality of the pattern to be stored in the external memory (considering that we are using a MobileNetV1 with $128 \times 128 \times 3$ inputs to match CORE50 image size); (iii) accuracy and Δ accuracy quantify the absolute accuracy at the end of the training and the gap with respect to a native rehearsal, respectively. For example, conv5_4/dw exhibits an interesting trade-off because the computation is about 32% of the native rehearsal one, the storage is reduced to 66% and the accuracy drop is mild (5,07%). CWR* (i.e. AR1* with latent replay layer \equiv pool6) has a really negligible computational cost (0.027%) with respect to native rehearsal and still provides an accuracy improvement of $\sim 4\%$ w.r.t. the non-rehearsal case ($\sim 60\%$ vs $\sim 56\%$). It is worth noting that the memory overhead of latent activations volumes can be further reduced with sparsification and compression techniques (see subsection 5.2 of [24] where we report additional experiments on specific topic).

A.2 Comparison with Other Approaches

While the accuracy improvement of the proposed approach w.r.t. state-of-the-art rehearsal-free techniques have been already discussed in the previous sections, a further comparison with other state-of-the-art continual learning techniques may be beneficial for better appreciating its practical impact and advantages. In particular, while AR1* and CWR* have been already proved to be substantially better than LWF and EWC on the NICv2 - 391 benchmark, a comparison with iCaRL, one of the best know rehearsal-based technique, is worth to be considered.

Unfortunately, iCaRL was conceived for incremental class learning scenario and its porting to NIC (whose batches also include patterns of know classes) is not trivial. To avoid subjective modifications, we started from the code shared by the authors and emulated a NIC setting by: (i) always creating new virtual classes from patterns in the coming batches; (ii) fusing virtual classes together when evaluating accuracies. For

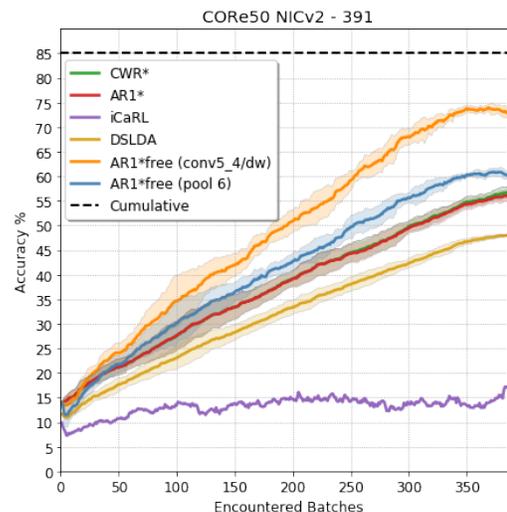


Fig. 5: Accuracy results on the CORE50 NICv2 - 391 benchmark of CWR*, AR1*, DSLDA, iCaRL, AR1*free (conv5_4), AR1*free (pool6). Results are averaged across 10 runs in which the batches order is randomly shuffled. Colored areas indicate the standard deviation of each curve. As an exception, iCaRL was trained only on a single run given its extensive run time (~ 14 days).

example, let us suppose to encounter 300 patterns of class 5 in batch 2 and other 300 patterns of the same class in batch 7; while two virtual classes are created by iCaRL during training, when evaluating accuracy both classes point to the real class 5. The hereby modified iCaRL implementation, with an external memory of 8000 patterns (much more than the 1500 used by the proposed latent replay, but in line with the settings proposed in the original paper [17]), was run on NICv2 - 391, but we were not able to obtain satisfactory results. In Figure 5 we report the iCaRL accuracy over time and compare it with AR1*free (conv5_4/dw), AR1* (pool6) as well as the top three performing rehearsal-free strategies introduced before: CWR*, AR1* and DSLDA. While iCaRL exhibits better performance than LWF and EWC (as reported in [6]), it is far from DSLDA, CWR* and AR1*.

Furthermore, when the algorithm has to deal with a so large number of classes (including virtual ones) and training batches its efficiency becomes very low (as also reported in [10]). In Table II we also report the total run time (training

TABLE II: Summary of the computation, memory, and accuracy trade-off for each strategy. Memory overhead include both the data used for replay purposes as well as additional trainable parameters needed for continual learning (Data + Params, in Mega Bytes). Run time (in Minutes) includes both training and testing. Each metric is averaged across 10 runs.

C0Re50 NICv2 - 391				
Strategy	Run Time	Mem. Overhead	Final Acc. %	Δ Acc. % vs Cumul.
CWR*	21.4	0 + 0.2	56.99%	-28.27%
AR1*free (pool6)	23.7	5.8 + 12.4	59.75%	-25.51%
AR1*	39.9	0 + 12.4	56.32%	-28.94%
AR1*free (conv5_4/dw)	41.2	48 + 0	72.23%	-13.03%
DSLDA	79.1	0 + 0.2	48.02%	-37.24%
iCaRL	20185.0	375 + 0	15.65%	-69.61%

and testing), memory overhead and accuracy difference with respect to the cumulative upper bound. We believe AR1*free (conv5_4/dw) represents a good trade-off in terms of efficiency-efficacy with a limited computational-memory overhead and only at $\sim 13\%$ distance from the cumulative upper bound. For iCaRL the total training time was 14 days compared to a training time of less than 1 hour for the other techniques.

B. Experiments on OpenLORIS

In order to show the general applicability of latent replay in different continual learning settings, we report and compare its performance also on the OpenLORIS dataset, which has been used as the main benchmark in the recent IROS 2019 “Lifelong Robotic Vision” competition [5].

OpenLORIS is particularly interesting for continual learning in the context of robotic vision since its video sessions have been recorded on a real wheeled robot exploring its environment (see Fig. 1). In this case, however, the scenario is quite different from C0Re50 NICv2 - 391: it is based on a sequence of 12 relatively large batches ($\sim 14,000$ samples each) containing only new examples of the same 69 classes made available in the first batch (also known as the “New Instances” scenario, NI). For this scenario:

- We use a MobileNetV2 [28] pre-trained on ImageNet-1k.
- We apply our *latent replay* approach to LWF [8], whose distillation steps proved to be effective to continually learn over a moderate number or large batches.

Seven finalists passed the first competition stage, and submitted their solutions to the organizers who finally produced the scoreboard reported in Table III. The accuracy of the proposed approach is just slightly lower than the top 1, but its inference time, replay memory and model size are significantly better. Since the challenge evaluation criteria did not include specific metrics on training efficiency, unfortunately from this experiment we cannot appreciate the training efficiency of our solution.

VI. REAL-WORLD DEPLOYMENT ON EMBEDDED DEVICES

The feasibility of continual learning at the edge on embedded devices is demonstrated through the development

TABLE III: Accuracy results at the end of the training and other metrics used for the OpenLORIS challenge benchmark.

OpenLORIS Challenge Results (7 finalists)				
Strategy	Final Acc. %	Inference Time (s)	Replay Size (Sample)	Model Size (MB)
SDU_BFA_PKU	99.56%	2,444.01	28,500	171.40
UniBo-Team (ours)	97.68%	22.41	1,500	5.90
HIK_ILG	96.86%	25.42	0	16.30
Vidit98	96.16%	112.2	13,000	9.40
NTU_LL	93.56%	4,213.76	0	467.10
Neverforget	92.93%	89.15	0	342.90
Guinness	72.9%	346.02	0	9.40

of an Android app (called C0Re) for Android smartphone (see Figure 6). While the app will be open-sourced and uploaded in the Google Play store upon publication of this manuscript, a video showing its functions is already available at <http://bit.ly/core50-app>.

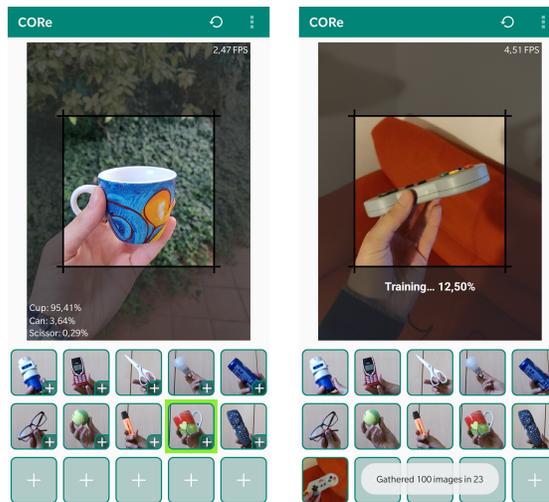


Fig. 6: The user interface of C0Re app. The camera field of view is partially grayed to highlight the central area where the image is cropped, resized to 128×128 and passed to the CNN. The top three categories are returned for each image (closed set classification) and a green frame is placed around the icon of the most likely class. A training session is triggered by tapping the icon of one of the 10 existing classes or one of the (initially) five empty classes.

The app comes pre-trained with 10 classes (corresponding to the 10 C0Re50 categories) and allows to: (i) continually train existing classes (by learning new object/poses) and (ii) to learn up to 5 brand new classes. As the app is started, it switches to inference mode and classifies the framed objects with an inference efficiency of about 5 fps (CPU-only with no hardware acceleration). When learning is triggered, a short video of 20 seconds (at 5 fps) is acquired and the resulting 100 frames are used for continual learning that completes in less than 1 seconds after the end of the acquisition.

Behind the scenes the app is running a customized Caffe version cross-compiled for Android and using the same MobileNetV1 architecture introduced in Section V, here initialized to work with 15 classes. Latent replay in this

case is implemented at the `pool6` layer³ with an external memory of 500 patterns. Low level code is written in C++ and the app interface in Java. A training session consists of 8 epochs, 5 iterations per epoch with a mini-batch size of 120 patterns: each mini-batch includes 20 original frames (from the current batch) and 100 replay patterns. In order to speed up training during the video acquisition, a second thread immediately moves available frames forward in the CNN and caches activations at latent replay layer so that, when the acquisition is concluded, we can directly train the class specific discriminative layers. Further details and precise timing of different phases are provided in the extended preprint version of this work [24].

VII. CONCLUSIONS

In this paper we showed that latent replay is an efficient technique to continually learn new classes and new instances of known classes even from small and non i.i.d. batches. State-of-the-art CL approaches such as AR1*, extended with latent replay, are able to learn efficiently and, at the same time, the achieved accuracy is not far from the cumulative upper bound (about 5% in some cases). The computation-storage-accuracy trade-off can be defined according to both the target application and the available resources so that even edge devices with no GPUs can learn continually from short videos, as we proved through the development of an Android application. In the future we intend to investigate: (i) the design of more sophisticated pattern replacing strategies for the external memory to contrast the aging effect; (ii) replacing the external memory with a generative model trained in the loop and capable of providing pseudo activations on demand.

REFERENCES

- [1] D. Li, S. Tasci, S. Ghosh, J. Zhu, J. Zhang, and L. Heck, "RILOD: Near Real-Time Incremental Learning for Object Detection at the Edge," in *Proceedings of the 4th ACM/IEEE Symposium on Edge Computing (SEC)*. New York, New York, USA: ACM Press, 2019, pp. 113–126.
- [2] M. McCloskey and N. J. Cohen, "Catastrophic Interference in Connectionist Networks: The Sequential Learning Problem," *Psychology of Learning and Motivation - Advances in Research and Theory*, vol. 24, no. C, pp. 109–165, 1989.
- [3] G. I. Parisi, R. Kemker, J. L. Part, C. Kanan, and S. Wermter, "Continual lifelong learning with neural networks: A review," *Neural Networks*, vol. 113, pp. 54–71, may 2019.
- [4] 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," *Information Fusion*, vol. 58, pp. 52–68, jun 2020.
- [5] H. Bae, E. Brophy, R. H. Chan, B. Chen, F. Feng, G. Graffieti, V. Goel, X. Hao, H. Han, S. Kanagarajah, S. Kumar, S.-K. Lam, T. L. Lam, C. Lan, Q. Liu, V. Lomonaco, L. Ma, D. Maltoni, G. I. Parisi, L. Pellegrini, D. Piyasena, S. Pu, Q. She, D. Sheet, S. Song, Y. Son, Z. Wang, T. E. Ward, J. Wu, M. Wu, D. Xie, Y. Xu, L. Yang, Q. Yang, Q. Zhong, and L. Zhou, "IROS 2019 Lifelong Robotic Vision: Object Recognition Challenge [Competitions]," *IEEE Robotics & Automation Magazine*, vol. 27, no. 2, pp. 11–16, jun 2020.
- [6] V. Lomonaco, D. Maltoni, and L. Pellegrini, "Rehearsal-Free Continual Learning over Small Non-I.I.D. Batches," *Continual Learning Workshop at CVPR2020, 2019*, pp. 1–10.
- [7] V. Lomonaco and D. Maltoni, "CORE50: a New Dataset and Benchmark for Continuous Object Recognition," in *Proceedings of the 1st Annual Conference on Robot Learning (CoRL)*, vol. 78, 2017, pp. 17–26.

- [8] Z. Li and D. Hoiem, "Learning without forgetting," in *14th European Conference on Computer Vision (ECCV)*, vol. 9908 LNCS, Amsterdam, Netherlands, 2016, pp. 614–629.
- [9] 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," in *Proceedings of the National Academy of Sciences*, vol. 114, 2017, pp. 3521–3526.
- [10] D. Maltoni and V. Lomonaco, "Continuous learning in single-incremental-task scenarios," *Neural Networks*, vol. 116, pp. 56–73, aug 2019.
- [11] G. Pasquale, C. Ciliberto, F. Odone, L. Rosasco, and L. Natale, "Are we done with object recognition? The iCub robot's perspective," *Robotics and Autonomous Systems*, vol. 112, pp. 260–281, 2019.
- [12] 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 *2020 International Conference on Robotics and Automation (ICRA)*, 2020, pp. 4767–4773.
- [13] L.-J. Li, K. Li, F. F. Li, J. Deng, W. Dong, R. Socher, and L. Fei-Fei, "ImageNet: a Large-Scale Hierarchical Image Database Shrimp Project View project hybrid intrusion detection systems View project ImageNet: A Large-Scale Hierarchical Image Database," *2009 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 248–255, 2009.
- [14] A. Krizhevsky, "Learning Multiple Layers of Features from Tiny Images," Ph.D. dissertation, 2009.
- [15] F. Zenke, B. Poole, and S. Ganguli, "Continual Learning Through Synaptic Intelligence," in *Proceedings of the 34th International Conference on Machine Learning (ICML)*, vol. 70, Sydney, Australia, 2017, pp. 3987–3995.
- [16] J. Schwarz, J. Luketina, W. M. Czarnecki, A. Grabska-Barwinska, Y. W. Teh, R. Pascanu, and R. Hadsell, "Progress & Compress: A scalable framework for continual learning," in *Proceedings of the 35th International Conference on Machine Learning (ICML)*, 2018, pp. 4528–4537.
- [17] S. Rebuffi, A. Kolesnikov, G. Sperl, and C. H. Lampert, "iCaRL: Incremental Classifier and Representation Learning," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Honolulu, Hawaii, 2017.
- [18] D. Lopez-paz and M. Ranzato, "Gradient Episodic Memory for Continuum Learning," in *Advances in Neural Information Processing Systems (NIPS 2017)*, 2017.
- [19] D. Rolnick, A. Ahuja, J. Schwarz, T. P. Lillicrap, and G. Wayne, "Experience replay for continual learning," *arXiv preprint arXiv:1811.11682v1*, pp. 1–16, 2018.
- [20] J. Pomponi, S. Scardapane, V. Lomonaco, and A. Uncini, "Efficient continual learning in neural networks with embedding regularization," *Neurocomputing*, vol. 397, pp. 139 – 148, 2020.
- [21] A. Robins, "Catastrophic Forgetting, Rehearsal and Pseudorehearsal," *Connection Science*, vol. 7, no. 2, pp. 123–146, 1995.
- [22] T. L. Hayes, N. D. Cahill, and C. Kanan, "Memory Efficient Experience Replay for Streaming Learning," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, may 2019, pp. 9769–9776.
- [23] T. L. Hayes and C. Kanan, "Lifelong Machine Learning with Deep Streaming Linear Discriminant Analysis," pp. 1–15, 2019.
- [24] L. Pellegrini, G. Graffieti, V. Lomonaco, and D. Maltoni, "Latent Replay for Real-Time Continual Learning," *arXiv preprint arXiv:1912.01100*, 2020.
- [25] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications," *arXiv preprint arXiv:1704.04861*, pp. 1–12, 2017.
- [26] S. Ioffe and C. Szegedy, "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift," *International Conference on Machine Learning (ICML)*, vol. 10, no. 6, pp. 448–456, aug 2015.
- [27] S. Ioffe, "Batch Renormalization: Towards Reducing Minibatch Dependence in Batch-Normalized Models," in *Advances in Neural Information Processing Systems (NIPS)*, 2017, pp. 1945–1953.
- [28] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L. C. Chen, "MobileNetV2: Inverted Residuals and Linear Bottlenecks," *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pp. 4510–4520, 2018.

³Placing latent replay layer at `pool6` corresponds to extending CWR* with latent replay, and leads to maximum efficiency on a CPU-only device.