

Overcoming Catastrophic Forgetting with Gaussian Mixture Replay

Benedikt Pfülb

Department of Applied Computer Science
 Fulda University of Applied Sciences
 Fulda, Germany
 benedikt.pfuelb@cs.hs-fulda.de

Alexander Gepperth

Department of Applied Computer Science
 Fulda University of Applied Sciences
 Fulda, Germany
 alexander.gepperth@cs.hs-fulda.de

Abstract—We present Gaussian Mixture Replay (GMR), a rehearsal-based approach for continual learning (CL) based on Gaussian Mixture Models (GMM). CL approaches are intended to tackle the problem of catastrophic forgetting (CF), which occurs for Deep Neural Networks (DNNs) when sequentially training them on successive sub-tasks. GMR mitigates CF by generating samples from previous tasks and merging them with current training data. GMMs serve several purposes here: sample generation, density estimation (e.g., for detecting outliers or recognizing task boundaries) and providing a high-level feature representation for classification. GMR has several conceptual advantages over existing replay-based CL approaches. First of all, GMR achieves sample generation, classification and density estimation in a single network structure with strongly reduced memory requirements. Secondly, it can be trained at constant time complexity w.r.t. the number of sub-tasks, making it particularly suitable for life-long learning. Furthermore, GMR minimizes a differentiable loss function and seems to avoid mode collapse. In addition, task boundaries can be detected by applying GMM density estimation. Lastly, GMR does not require access to sub-tasks lying in the future for hyper-parameter tuning, allowing CL under real-world constraints. We evaluate GMR on multiple image datasets, which are divided into class-disjoint sub-tasks.

Index Terms—Gaussian Mixture Models, Continual Learning, Life-long Learning, Pseudo-Rehearsal, Incremental Learning

I. INTRODUCTION

This article is set in the context of continual learning (CL). Continual (or life-long) learning is an important element for successful biological systems. CL denotes the ability to acquire new knowledge without forgetting previously learned knowledge. The CL problem is usually formalized as sequentially training a model (e.g., a DNN) on data that is divided into a sequence of sub-tasks (here denoted as T_1, T_2, \dots). When doing this, a basic problem that occurs with DNNs is catastrophic forgetting (CF) [12].

Catastrophic Forgetting Essentially, the CF effect implies an abrupt and near-complete loss of knowledge about T_1, \dots, T_{t-1} with a few training iterations on T_t (see Fig. 1). In real-world applications, many variations of the CL scenario can be distinguished. This includes, for example, the problem of unavailable sub-task boundaries, which can even be fluid (concept drift/shift). Many approaches (some of which are targeting very specific CL scenarios) for avoiding CF have

been proposed recently, see Sec. I-B. This article proposes a model that is based on *replay*.

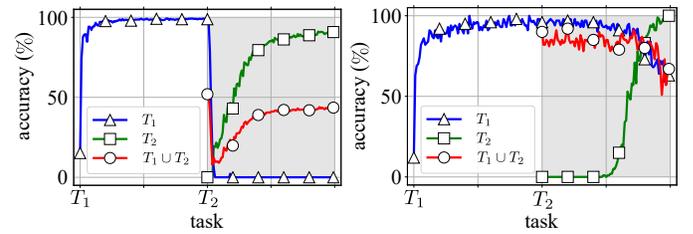


Fig. 1: Visualization of the catastrophic forgetting (CF) effect [30]. Left image shows CF. The blue line drops complete after a few training iterations on a second task. Hence, knowledge of previous task is lost. Right side shows a more controllable linear decline.

Replay An inspiration for avoiding CF is the learning process of students. In school, new topics are continuously taught. The challenge is to retain older topics when studying current ones. In order to solve this problem, students must 1) recognize that they have forgotten something and 2) repeat and revise, if necessary. This real-world association can serve as an inspiration for a replay approach towards CL, see Fig. 2.

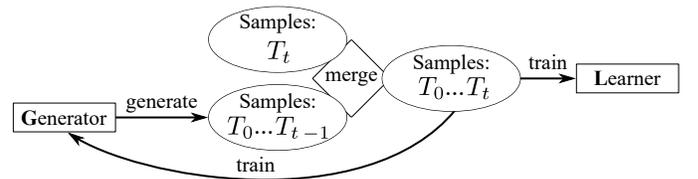


Fig. 2: The replay approach to continual learning: a Learner (L), e.g., a DNN, is trained on several sub-tasks sequentially. To avoid forgetting, a generator (G) is trained to generate samples from past sub-tasks. For training L, G *generates* samples from past sub-tasks, which are merged with current sub-task samples.

A. Gaussian Mixture Replay

GMR is a new approach for lifelong learning based on Gaussian Mixture Models (GMMs). GMMs aim to represent the distribution of training samples, enabling them to assess the probability of unknown samples (density estimation), as well as generate new samples from the learned distribution (sampling). In the context of CL, density estimation is a

useful tool for detecting task boundaries, because samples of a previously unseen nature (e.g., a new visual class) start arriving at task boundaries. Furthermore, GMMs can play the role of a generator in replay-based CL due to their sampling ability see Fig. 2 (e.g., sampling from a multivariate Gaussian distributions). In GMR, the learner and the generator from Fig. 2 are contained in a single structure. The GMM layer provides sampling, density estimation, as well as a feature representation for the linear classifier layer. This layer operates on quantities termed *responsibilities* or posterior probabilities of the GMM layer, see Fig. 3. Responsibilities are normalized and lie in the range $[0 \dots 1]$. Therefore, they are ideal for linear classification. In turn, by approximate inversion of the linear classifier layer, GMM sampling can be informed about the classes of samples that should be generated.



Fig. 3: Principal structure of the GMR model, composed of a GMM modeling the distribution of training samples (left), and a linear classifier operating on the posterior probabilities (sometimes termed *responsibilities*) produced by the GMM. The coupled GMM/classifier implements the learner from Fig. 2 (operating from left to right), whereas the GMM implements the generator from Fig. 2. The GMM sampling process is informed by feedback from the classifier (operating from right to left).

B. Related Work

The field of CL is expanding rapidly, see [9], [15], [29], [36] for reviews. Systematic comparisons between different approaches to avoid CF are performed in, e.g., [18], [30]. As discussed in [30], many recently proposed methods demand specific experimental setups, which deviate significantly from application scenarios. For example, some methods require access to samples from *all* sub-tasks for tuning hyper-parameters, whereas others need to store all samples from past tasks. Many proposed methods have a time and/or memory complexity that scales at least linearly with the number of sub-tasks and, thus, may fail if this number is large. Among the proposed remedies to CF, three major directions may be distinguished according to [9]: parameter isolation, regularization and replay.

Parameter Isolation Isolation methods aim at determining (or creating) a group of DNN parameters that are mainly “responsible” for a certain sub-task. CL is then avoided by *protecting* these parameters when training on successive sub-tasks. Representative works are [3], [11], [26], [27], [32], [34].

Regularization Regularization methods mostly propose modifying the DNN loss function, including additional terms that protect knowledge acquired in previous sub-tasks. Current approaches are very diverse: SSL [5] focuses on enhancing sparsity of neural activities, whereas approaches such as LwF [24] rely on knowledge distillation mechanisms. A method that has attracted significant attention is Elastic Weight Consolidation (EWC) [19], which inhibits changes to weights that are important to previous sub-tasks, measuring their

importance based on the Fisher information matrix (FIM). Synaptic intelligence [38] is pursuing a similar goal. An online version of EWC is also published [33]. Incremental Moment Matching (IMM) [21] makes use of the FIM to merge the parameters obtained for different sub-tasks. The Matrix of Squares (MasQ) method [14] is identical to EWC, but relies on the calculus of derivatives to assess the importance of parameters for a given sub-task. It is, therefore, more simple w.r.t its concepts and much more memory-efficient.

Replay Replay-based methods are at the core of this article. They mainly occur in two forms: rehearsal and pseudo-rehearsal. *Rehearsal methods* store a subset of samples from past sub-tasks preventing CF, either by putting constraints on current sub-task training or by adding retained samples to the current sub-task training set. Typical representatives of rehearsal methods are iCaRL [31], (A-)GEM [7], [25], GBSS [4] and TEM [8]. *Pseudo-rehearsal* methods, in contrast, do not store samples but generate them using a dedicated *generator* that is trained along with the learner, see Fig. 2. Typical models used as generators are Generative Adversarial Networks (GANs), Variational Autoencoders (VAEs) and their variants, see [35] and [17]. While being impressive in terms of CL performance, these approaches have problems, as well. First of all, training time scales linearly with the number of sub-tasks. Secondly, even for simple datasets like MNIST, the generators require significant computational resources, in addition to the resource requirements for the learner (CNNs or DNNs). Lastly, GANs can suffer from mode collapse, which is very hard to detect, since GANs do not minimize a differentiable loss function. The GMR model we are proposing belongs to pseudo-rehearsal approaches, as well.

Training and Evaluation Paradigms for CL In the context of continual learning, a variety of trainings and evaluation paradigms are proposed the literature [2], [6], [10], [16], [18], [22]. In addition to a very broad definition of various topics related to continuous learning, [23] addresses the topic of evaluation protocols and metrics. The influence of different difficulties and the correlative order of tasks is addressed. The evaluation of generative models is also highlighted. De Lange et. al. describe the random selection of classes that are divided into tasks of size two [9]. It must be ensured that all tasks are approximately equal in difficulty. Effectiveness is measured by determining the difference between the accuracy and the expected value after each task. Another comprehensive survey is presented by [28]. Here an evaluation strategy that is “more comprehensive” with regard to real-world applications is suggested, among other things. The classical evaluation pipeline is extended by adding queries, e.g., related to the importance of the data. In addition to human activities, statistical methods can be used to influence the training process. Parisi et al. [29] also address the issue of continual learning from a biological perspective. Moreover, the benchmark procedure used is shown. Again, “incremental class learning” is one of the methods used.

C. Contributions

GMR offers several novel contributions to the field of CL:

- Pseudo-rehearsal based CL integrating learner and generator in a single structure
- Easy-to-implement detection of task boundaries that is well-justified by probability theory
- Constant time complexity w.r.t. the number of sub-tasks
- No cross-validation on future sub-tasks for hyper-parameter selection required

To validate our approach, we perform a comparison to the Elastic Weight Consolidation (EWC) model which is viewed as a “standard model” for CL in many recent publications. Furthermore, we provide a public TensorFlow implementation¹.

II. DATASETS

In order to measure the impact of forgetting during continual learning, three public datasets are used (see Tab. I). All datasets consist of gray scale images with dimensions of 784 or 1024. Each of them is divided in the ratio 90% to 10% in training and test data with an almost equal distribution of samples within classes. The used datasets are normalized to the $[0, 1]$ range.

MNIST contains images of handwritten digits (0-9) with a resolution of 28×28 pixels. It is probably the most commonly used benchmark for classification problems. FashionMNIST contains pictures of different types of clothes. This dataset is supposed to be harder to classify compared to MNIST (same resolution) and, thus, leads to lower accuracies. Similar to MNIST, the Devanagari dataset contains written Devanagari letters. It is available in a resolution of 32×32 pixels per image. Since there are more classes included than needed, we randomly select 10 classes.

III. THE GAUSSIAN MIXTURE REPLAY MODEL

In this section, we provide details concerning the proposed GMR model. We put a particular emphasis on extensions/modifications of the GMM component to allow for a pseudo-rehearsal mechanism that performs well in practice. For implementing and training the GMM as core component, we use the approach (and code, in a modified form) presented in [13]. This allows to train GMMs by Stochastic Gradient Descent (SGD) on high-dimensional data. The latter is a hard requirement for the presented work, since all datasets are image-based and contain a significant number of samples. The size of the used datasets makes the Expectation-Maximization (EM) approach to train GMMs unfeasible, which is why

¹<https://gitlab.cs.hs-fulda.de/ML-Projects/gmr>

the efficient SGD training procedure proposed in [13] is advantageous.

GMR is a layered model containing a folding GMM and classification layer, see Fig. 3. This is discussed in the following sections. All layers expect their inputs to be four-dimensional tensors in NHWC format, as it is usual for CNNs.

A. Folding Layer

The folding layer performs a part of what a convolutional layer (in CNNs) would be doing, namely to extract organized slices from its input. Folding layers are parameterized by receptive field/filter sizes f_X, f_Y and strides Δ_X, Δ_Y . For an input tensor of the dimensions $N \times H \times W \times C$, a folding layer produces an output tensor of the dimensions $N \times H' \times W' \times C'$:

$$\begin{aligned} H' &= 1 + \frac{H - f_Y}{\Delta_Y} \\ W' &= 1 + \frac{W - f_X}{\Delta_X} \\ C' &= f_X f_Y C \end{aligned} \quad (1)$$

These quantities result from “dumping” the flattened content of a receptive field into a single channel dimension (a single *column*) of the output. In a CNN, this would be followed by a scalar product with the CNN filters, which is replaced here by the processing steps of the subsequent GMM layer.

B. Gaussian Mixture Model Layer

GMM layers contain a single instance of a Gaussian Mixture Model. GMMs describe data samples \vec{x} as a mixture of K Gaussian densities \mathcal{N} , see Eq. (2).

$$p(\vec{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\vec{x}; \vec{\mu}_k, \Sigma_k) \quad (2)$$

For an input tensor of the dimensions $N \times H \times W \times C$, this GMM instance processes all HW columns of its input tensor independently. Thus, the GMM instance processes input vectors \vec{x} of dimension C . The posterior probability (or responsibility) of the GMM, which we treat as the GMM layer’s output to the classifier layer, is computed as:

$$\gamma_j(\vec{x}) = \frac{\pi_j \mathcal{N}(\vec{x}; \vec{\mu}_j, \Sigma_j)}{\sum_k \pi_k \mathcal{N}(\vec{x}; \vec{\mu}_k, \Sigma_k)}. \quad (3)$$

Thus, responsibilities have the dimensions of $N \times H \times W \times K$. They are normalized and bounded in the interval $[0, 1]$.

To learn the distribution of inputs from data, we use the SGD procedure from [13]. This is, however, a technical detail, and trained GMMs behave like any other GMM trained by

TABLE I: Detailed information to the used datasets (including examples of each class).

dataset	ref.	resolution	number of training samples	number of test samples	random examples (from classes)									
					0	1	2	3	4	5	6	7	8	9
MNIST	[20]	28×28	50 000	10 000										
FashionMNIST	[37]	28×28	60 000	10 000										
Devanagari	[1]	32×32	18 000	2 000										

Expectation-Maximization (EM). In particular, Eq. (2) can be used for detecting outliers whose (log-)probability is low. For sampling, a value is drawn from a component distribution $\mathcal{N}(\vec{x}; \vec{\mu}_k, \Sigma_k)$. k is drawn from a multinomial distribution parameterized by the weights $\vec{\pi}$ of the GMM. In GMR, the downstream linear classifier layer can provide a control signal to the GMM layer \vec{t} for sampling. In this case, the parameter to the multinomial distribution is \vec{t} , and sampling is performed for each column independently, resulting in a sampling result of the same dimension as the input.

C. Linear Classification Layer

Since GMMs are inherently trained in an unsupervised manner, they cannot compute class labels. For the assignment of labels to GMM outputs, we use a linear classifier layer that performs the operation $\vec{\gamma}W + \vec{b}$ on the responsibilities $\vec{\gamma}$. The softmax function, which is usually applied beforehand, is already included in the responsibilities, see Eq. (3). For training, either cross-entropy or mean squared error (MSE) can be used as loss function, in combination with SGD. For generating a control signal to the GMM layer, a one-hot vector \vec{o} representing the desired class is created and transformed as $\vec{t} = \text{norm}_{[0,1]}(\vec{o} - \vec{b})W^\top$. Thus, it is possible to select components from the overlying GMM layer to enable conditional sampling.

D. Classification Process Example

In the following section, the classification process by the GMR is described and also visualized as an example (see Fig. 4). Let us assume that the input tensor has the following format: $N = 1$, $H = 5$, $W = 5$ and $C = 1$. For the sake of simplicity, we omit the dimension N of shape 1. We are left with a single $5 \times 5 \times 1$ gray-scale image as input \vec{x} (see Fig. 4 left). This leads to a given GMR which consists of a folding layer, a GMM layer and a classification layer (central part).

The folding layer unfolds the input signal with a receptive field of the size $f_X = f_Y = 5$. It follows that striding is not used. This shows that the image is transformed line by line into a 1D vector, which is passed to the GMM layer.

The visualized GMM layer consists of $K = 9$ components arranged in a 3×3 grid. Each component is composed of $\vec{\mu}_k$ (the mean) and Σ_k (the variances) with the shape of the input signal (1×25). For simplicity, only the means are displayed, which in turn were reshaped from a 1D vector (same shape an input signal) into a 2D image. We assume that the GMR has already converged and derived the following data distribution from the training data. Note that the prototypes are presented in an exaggerated way. The central component (green highlighted) is the one that has the largest association for the input signal (0 sample), which is why the responsibility (γ_k) for this prototype is the highest (also green). For clarity, all responsibilities ($\vec{\gamma}$) are additionally transformed (3×3) to match the visualization of the GMM components.

The linear layer receives the responsibilities ($\vec{\gamma}$) of the GMM layer. Here, the class label is assigned based on the

learned distribution. The output of this layer, and thus, the GMR, is a class label associated with the input.

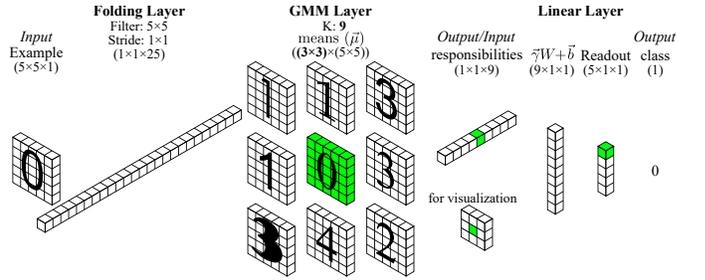


Fig. 4: Visualization of the classification process of a GMR model.

E. Sample Generation Example

The inverse function for classification is the sample generation. As described in Sec. III-D, we assume that the model is already trained. Again, we set $N = 1$ and neglect this dimension for visualization. In Fig. 5 the reverse data flow through the GMR is depicted. To generate samples conditionally, the class from which a sample should be generated is used as input. It is converted as one-hot vector (see Fig. 5 left). The class from which an example is generated is highlighted in green.

The linear layer is used to inverse the mapping of a GMM component to a class by generating a signal which is passed on to the GMM layer above. For this, the inverse transformation is taken as described in Sec. III-C. It is also highlighted in green and arranged for visualization (3×3) as the components of the upper GMM layer.

This time, for each GMM component the means ($\vec{\mu}$) and corresponding variances (Σ) in their true form are given. We again assume that the GMR has already learned the distribution from the data (see Fig. 4 for an impression of the prototypes). The selection strategy also includes the choice of different components that simultaneously represent the same class. This is achieved on the basis of a probability distribution depending on the signal of the classifier layer. In this example, only one sample should be created and only the central component (green highlighted) is responsible, which is why it is very likely to be selected for sampling. A sample x is drawn from the selected GMM component by sampling from a multivariate normal distribution $\mathcal{N}(\mu_k, \Sigma_k)$.

The last step is to restore the original form of the input data (images). The folding layer inverts the flattening process by reshaping x into the origin input shape ($5 \times 5 \times 1$).

IV. ELASTIC WEIGHT CONSOLIDATION

The approach from [19] is a typical regularization-based model for DNNs, see Sec. I-B. EWC stores DNN parameters θ_t^T after training on sub-task T_t . In addition, EWC computes the “importance” of each parameter after training on sub-task T_t by approximating the diagonal of the Fisher Information Matrix² $\vec{F}^T t$. The EWC loss function contains additional

²See [14] for a discussion of this approximation.

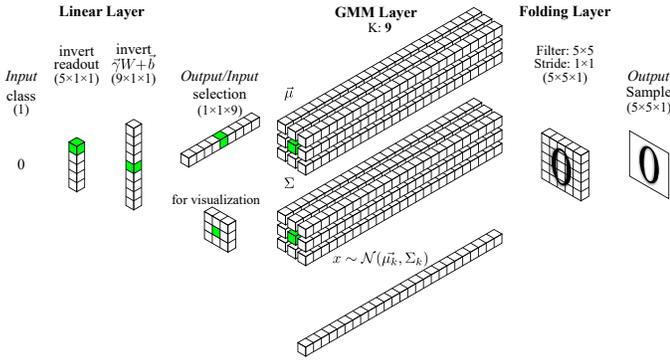


Fig. 5: Visualization of the sampling process of a GMR model.

terms, see Eq. (4) besides the cross-entropy loss computed on the current sub-task T_c . These additional loss terms punish deviations from “important” DNN parameter values obtained after training on past sub-tasks:

$$\mathcal{L}^{EWC} = \mathcal{L}_{T_c}(\theta) + \frac{\lambda}{2} \sum_{t=1}^{c-1} \sum_i F_i^{T_t} (\theta_i - \theta_i^{T_t})^2 \quad (4)$$

EWC is optimized using SGD, usually in conjunction with the Adam optimizer. EWC hyper-parameters comprise the SGD step size ϵ , the regularization constant λ and of course the number and size of layers in the DNN. In [19], it is proposed to set $\lambda = \epsilon^{-1}$, thereby eliminating one hyper-parameter.

V. EXPERIMENTAL SETUP

In this section, the used parameters, tasks and evaluation methods for EWC and GMR are described. Generally, each experiment is repeated 10 times. This ensures that the results are not influenced by the random initialization of the model.

A. Sequential Learning Tasks

Sequential Learning Tasks (SLTs) simulate a continuous learning scenario by dividing datasets given in Sec. II. The resulting sub-datasets are enumerated and contain only samples of non-overlapping classes. For example, a D_{5-5} task consists of two sub-datasets consisting of 5 classes each. Each sub-task is identified by its order, e.g., T_1, T_2, \dots, T_x . Baseline experiments (D_{10}) contain all available classes to investigate the effect of incremental task-by-task training.

SLTs can be used to perform basic experiments to determine the effect of forgetting under the conditions described above. To measure the impact of the number of classes contained in a task, different combinations and subdivisions are evaluated. Tab. II shows evaluated SLTs and their definition of sub-tasks.

B. Hyper-Parameter

For both investigated models, a fixed batch size of $\mathcal{B} = 100$ is used. Both models are informed about the start and end of each sub-task, so they can react appropriately.

The following hyper-parameters were explored for GMR by performing a grid search. GMR is trained for 50 epochs (\mathcal{E}) for the first task. The number of epochs for retraining is

TABLE II: Definition of Sequential Learning Tasks (SLTs) and the class divisions of their sub-tasks.

SLT	Sub-tasks
D_{10} (baseline)	$T_1(0, 1, 2, 3, 4, 5, 6, 7, 8, 9)$
D_{9-1a}	$T_1(0, 1, 2, 3, 4, 5, 6, 7, 8) \quad T_2(9)$
D_{9-1b}	$T_1(0, 1, 2, 4, 5, 6, 7, 8, 9) \quad T_2(3)$
D_{5-5a}	$T_1(0, 1, 2, 3, 4) \quad T_2(5, 6, 7, 8, 9)$
D_{5-5b}	$T_1(0, 1, 2, 6, 7) \quad T_2(3, 4, 5, 8, 9)$
$D_{3-3-3-1}$	$T_1(0, 1, 2) \quad T_2(3, 4, 5) \quad T_3(6, 7, 8) \quad T_4(9)$
$D_{2-2-2-2-2a}$	$T_1(0, 1) \quad T_2(2, 3) \quad T_3(4, 5) \quad T_4(6, 7) \quad T_5(8, 9)$
$D_{2-2-2-2-2b}$	$T_1(1, 7) \quad T_2(0, 2) \quad T_3(6, 8) \quad T_4(4, 5) \quad T_5(3, 9)$
$D_{1-1-1-1-1-1-1-1-1-1a}$	$T_1(0) \quad T_2(1) \quad T_3(2) \quad T_4(3) \quad T_5(4) \quad T_6(5) \quad T_7(6) \quad T_8(7) \quad T_9(8) \quad T_{10}(9)$
$D_{1-1-1-1-1-1-1-1-1-1b}$	$T_1(7) \quad T_2(1) \quad T_3(2) \quad T_4(0) \quad T_5(6) \quad T_6(8) \quad T_7(4) \quad T_8(5) \quad T_9(9) \quad T_{10}(3)$

adjusted depending on the ratio of the number of classes of the previous and the new task. We do not modify the model itself after each training task. The start of a new sub-task triggers the generation of samples from past sub-tasks, see Fig. 3. The number of GMM components varies: $K \in \{36, 84, 100\}$. As indicated in [13], the principle “the more the better” applies, which makes selection unproblematic. In order to investigate the influence of the loss function of the linear layer, we evaluate cross-entropy and MSE losses. The learning rate of both, the GMM and the linear layer, is fixed after initial tests. This is reasonable due to the GMM learning rate exclusively depending on the data which do not change, and because the inputs to the linear classifier layer are normalized. Thus, the total number of 1800 experiments performed is composed of the three datasets, the conducted SLTs (10), varied parameters (3) and 10 repetitions.

For EWC, we perform a grid search for the parameter ϵ . We vary ϵ as $\epsilon \in \{0.001, 0.0001, 0.00001, 0.000001, 0.0000001\}$ while λ is always set to $\frac{1}{\epsilon}$. We fix the model architecture to a three-layer DNN, each of size 800 (used in most recent studies on EWC). Training epochs \mathcal{E} are empirically set to 10 for each training task. This is done to ensure a fair comparison of GMR and EWC. Due to replay, GMR is insensitive to the duration of training for each sub-task. EWC, however, is very sensible to this parameter as performance tends to deteriorate over (training) time. As longer training times lead to less favorable EWC results, we generally evaluate performance during the training process. Using this evaluation strategy, the high number of training iterations for EWC should not be problematic. This contradicts a real-world scenario evaluation, since the best training duration is explicitly unknown and cannot be measured but it is used here, however, to ensure comparability in the first place. The best hyper-parameters and experiments are selected based on the maximum measured accuracy during training averages over 10 experiment repetitions. From the number of datasets (3), the performed SLTs (10), the different parameters (5) and the repetitions (10) results in 1500 conducted EWC experiments.

C. Task Boundaries Specification

One difficulty of continual learning (CL), which is often not addressed, is the specification of task boundaries. This

TABLE III: Results of the conducted GMR and EWC experiments. The accuracy in % of the baseline experiments on all available classes is shown for each dataset. For each best SLT experiment (defined in Tab. II) the difference to the baseline is given. Therefore, the maximum measured accuracy value is used, averaged over 10 experiment repetitions. To measure the accuracy, the joint test dataset consisting of all tasks (D_{10}) is used.

model dataset	GMR						EWC					
	MNIST		FashionMNIST		Devanagari		MNIST		FashionMNIST		Devanagari	
SLT	acc. %	std	acc. %	std	acc. %	std	acc. %	std	acc. %	std	acc. %	std
D_{10} baseline	87.4	0.59	73.9	0.26	74.1	0.73	97.57	0.26	87.55	0.38	95.58	0.56
	diff.	std	diff.	std	diff.	std	diff.	std	diff.	std	diff.	std
D_{9-1a}	-1.3	0.59	-2.7	0.26	-3.2	0.73	-41.8	0.26	-9.6	0.38	-56.6	0.56
D_{9-1b}	-3.5	2.19	-1.5	0.87	-1.4	0.88	-50.7	7.77	-20.1	2.52	-29.7	13.34
D_{5-5a}	-0.6	1.53	-1.2	1.53	-6.8	1.38	-35.3	6.65	-32.7	4.22	-46.0	15.38
D_{5-5b}	-1.3	1.92	-1.9	0.49	-4.7	1.59	-35.0	1.83	-36.0	2.72	-47.1	0.11
$D_{3-3-3-1}$	-8.3	2.68	-8.3	1.46	-15.5	1.83	-59.4	2.99	-52.7	2.57	-63.9	0.11
$D_{2-2-2-2-2a}$	-9.5	3.83	-8.5	0.91	-22.5	2.71	-72.2	7.43	-55.6	4.05	-72.1	2.75
$D_{2-2-2-2-2b}$	-10.4	5.28	-5.7	2.37	-14.7	2.94	-72.6	3.22	-57.3	4.99	-73.2	2.31
$D_{1-1-1-1-1-1-1-1-1a}$	-23.3	4.10	-19.2	1.83	-27.5	7.52	-74.5	3.79	-55.2	5.64	-71.1	2.26
$D_{1-1-1-1-1-1-1-1-1b}$	-16.3	2.06	-19.6	2.06	-35.6	2.69	-76.6	3.40	-56.0	4.17	-77.2	2.74

This effect of forgetting can be recognized in all EWC experiments (see also Figs. 9, 11 and 12). For the best D_{5-5a} experiment (see Fig. 9) we want to give more insights into the (re-)training process, especially the GMM layer. After

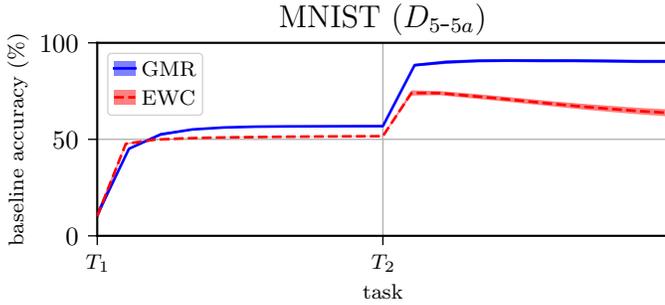
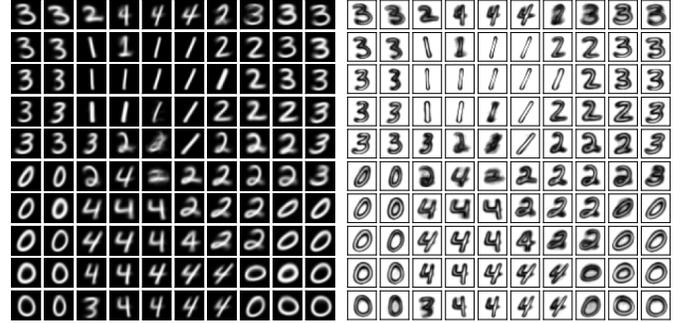


Fig. 9: Visualization of accuracy trend for SLT D_{5-5a} .

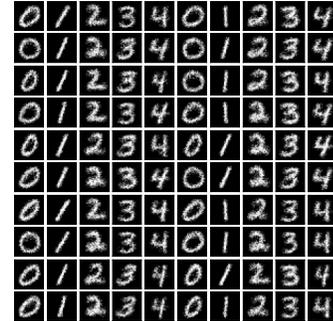
training T_1 for 50 epochs with the MNIST dataset, the 100 components K respectively the means/prototypes ($\bar{\mu}$) are visualized in Fig. 10a. The corresponding variances are displayed in Fig. 10b. Note that the topological order of the prototypes caused by the initial annealing process is no longer preserved after T_2 . After finishing the training on T_1 a dataset G_1 is generated which is used for replay (see Fig. 10c). The examples shown here are noisy, which is due to the high variances present in the original data. After training with $G_1 \otimes T_2$ the prototypes and variances can be observed in Fig. 10d and Fig. 10e. This is the only experiment for which the cross-entropy loss for the linear layer give the best results. It can be seen that the use of MSE objective function leads to a higher variance within the other GMR experiment groups.

VII. DISCUSSION

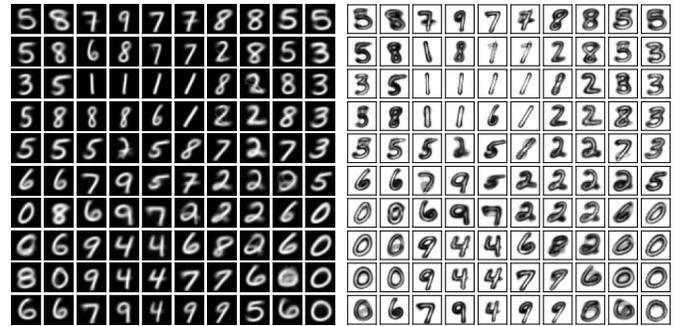
The experimental results underline that GMR shows consistently better continual learning performance on all tested datasets, for all sequential learning tasks. This effect is remarkable since GMR is significantly less complex in terms of free model parameters and memory requirements. In addition, GMR has been shown to be capable of detecting task



(a) Prototypes (μ) after training T_1 . (b) Variances (Σ) after training T_1 .



(c) Batch of generated samples (G_1).



(d) Prototypes after training $G_1 \otimes T_1$. (e) Variances after training $G_1 \otimes T_1$.

Fig. 10: Visualization of GMM prototype, GMM variances and generated samples for SLT D_{5-5a} : $T_1(0,1,2,3,4)$ and $T_2(5,6,7,8,9)$.

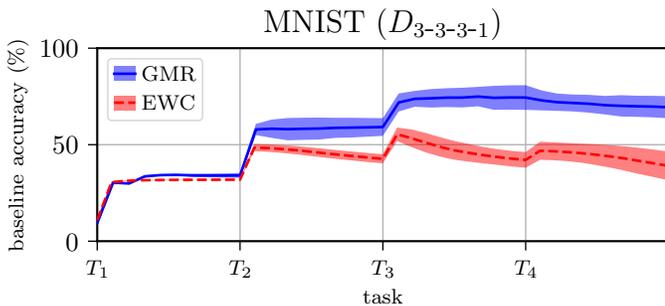


Fig. 11: Visualization of accuracy trend for SLT $D_{3-3-3-1}$.

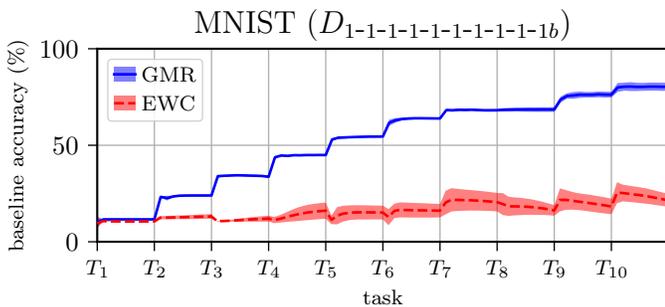


Fig. 12: Visualization of accuracy trend for SLT $D_{1-1-1-1-1-1-1-1-1-1b}$.

boundaries. While GMR and EWC are informed about sub-task boundaries in the current training setup, this is not really a requirement for GMR. Moreover, EWC has no possibility to determine this, as DNNs are known to be incapable of out-of-sample recognition or outlier detection. A similar argument distinguishes GMR from generative replay approaches as discussed in [35]. Furthermore, we find that GMR can be re-trained with a constant number of current and generated samples and, thus, has constant time complexity w.r.t. the number of previous sub-tasks. Lastly, GMR is a compact model in comparison to other replay approaches such as [35], because it integrates learner and generator in a single structure. **Baseline performance** is the first issue that needs to be discussed. Due to GMR’s simple structure, baseline (i.e., non-continual) performance on MNIST is only 87.4%, as compared to EWC (97.5%). First of all, it is possible to overcome this limitation by including a *deep* GMM into GMR, which boosts baseline performance to about 97% for a two-layer convolutional GMM (this will be discussed in a separate article). More importantly, this article is about continual learning, and about continual learning performance. The results of Tab. III plainly show that GMR is superior to EWC by a large margin, even though baseline performance is rather weak.

Suitability for applications is another important point (see [30]) in which GMR surpasses EWC. While the only real free parameter is the number of GMM components K , which can be set according to a *the-more-the-better* strategy without reference to data. On the contrary, the EWC parameter λ has a strong impact on performance and must be carefully adapted to the problem by cross-validation. This requires using data from all sub-tasks, which contradicts the principle

of continual learning. This includes that the sub-tasks are processed sequentially.

VIII. FUTURE WORK

Future work will address the following issues:

Weak baseline performance will be resolved using deep convolutional GMMs, consisting of multiple GMM and folding layers. Since this improves the quality of sampling, a further increase in CL performance can be expected.

More refined replay strategies might include strategies in which the model itself decides which and how many data samples to generate. A guiding principle could be to generate samples whose internal representation would otherwise be overwritten by current sub-task data. In comparison, this task is easier for GMMs than for DNNs.

Systematic comparison to other CL models should be a logical next step once a sufficiently high baseline performance can be reported. In particular, a comparison to generative replay models like [35] could be interesting.

REFERENCES

- [1] S. Acharya, A. K. Pant, and P. K. Gyawali. Deep learning based large scale handwritten Devanagari character recognition. In *2015 9th International Conference on Software, Knowledge, Information Management and Applications (SKIMA)*, pages 1–6, 2015.
- [2] Tameem Adel, Cuong V. Nguyen, Richard E. Turner, Zoubin Ghahramani, and Adrian Weller. Interpretable Continual Learning. 2019.
- [3] R. Aljundi, P. Chakravarty, and T. Tuytelaars. Expert gate: Lifelong learning with a network of experts. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7120–7129, July 2017.
- [4] Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient based sample selection for online continual learning. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 11816–11825, 2019.
- [5] Rahaf Aljundi, Marcus Rohrbach, and Tinne Tuytelaars. Selfless Sequential Learning. In *International Conference on Learning Representations*, 2019.
- [6] Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark experience for general continual learning: a strong, simple baseline. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- [7] Arslan Chaudhry, Marc’Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient Lifelong Learning with A-GEM. In *International Conference on Learning Representations*, 2019.
- [8] Arslan Chaudhry, Marcus Rohrbach, Mohamed Elhoseiny, Thalaiyasingam Ajanthan, Puneet K. Dokania, Philip H. S. Torr, and Marc’Aurelio Ranzato. On tiny episodic memories in continual learning.
- [9] M. Delange, R. Aljundi, M. Masana, S. Parisot, X. Jia, A. Leonardis, G. Slabaugh, and T. Tuytelaars. A continual learning survey: Defying forgetting in classification tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2021.
- [10] Sebastian Farquhar and Yarin Gal. Towards robust evaluations of continual learning. *arXiv*, 2018.
- [11] Chrisantha Fernando, Dylan Banarse, Charles Blundell, Yori Zwols, David Ha, Andrei A. Rusu, Alexander Pritzel, and Daan Wierstra. PathNet: Evolution Channels Gradient Descent in Super Neural Networks. 2017.
- [12] Robert M. French. Catastrophic forgetting in connectionist networks. *Trends in Cognitive Sciences*, 3(4):128–135, 1999.
- [13] Alexander Gepperth and Benedikt Pfülb. Gradient-Based Training of Gaussian Mixture Models in High-Dimensional Spaces. *arXiv*, 2019.
- [14] Alexander Gepperth and Florian Wiech. Simplified computation and interpretation of fisher matrices in incremental learning with deep neural networks. In *Artificial Neural Networks and Machine Learning – ICANN 2019: Deep Learning*, pages 481–494, Cham, 2019. Springer International Publishing.

- [15] Tyler L. Hayes, Ronald Kemker, Nathan D. Cahill, and Christopher Kanan. New metrics and experimental paradigms for continual learning. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 2112–2123, jun 2018.
- [16] Joseph K J and Vineeth N Balasubramanian. Meta-Consolidation for Continual Learning. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 14374–14386. Curran Associates, Inc., 2020.
- [17] Nitin Kamra, Umang Gupta, and Yan Liu. Deep generative dual memory network for continual learning. *CoRR*, abs/1710.10368, 2017.
- [18] Ronald Kemker, Marc McClure, Angelina Abitino, Tyler Hayes, and Christopher Kanan. Measuring catastrophic forgetting in neural networks.
- [19] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis, Claudia Clopath, Dharshan Kumaran, and Raia Hadsell. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, 114(13):3521–3526, mar 2017.
- [20] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, Nov 1998.
- [21] Sang-Woo Lee, Jin-Hwa Kim, Jaehyun Jun, Jung-Woo Ha, and Byoung-Tak Zhang. Overcoming Catastrophic Forgetting by Incremental Moment Matching. In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17*, page 4655–4665, Red Hook, NY, USA, 2017. Curran Associates Inc.
- [22] Timothée Lesort, Hugo Caselles-Dupré, Michael Garcia-Ortiz, Andrei Stoian, and David Filliat. Generative Models from the perspective of Continual Learning. *2019 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8, 2019.
- [23] Timothée Lesort, Vincenzo Lomonaco, Andrei Stoian, Davide Maltoni, David Filliat, and Natalia Díaz-Rodríguez. Continual learning for robotics: Definition, framework, learning strategies, opportunities and challenges. *Information Fusion*, 58:52–68, 2020.
- [24] Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40:2935–2947, 2018.
- [25] David Lopez-Paz and Marc’Aurelio Ranzato. Gradient episodic memory for continual learning. In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17*, page 6470–6479, Red Hook, NY, USA, 2017. Curran Associates Inc.
- [26] Arun Mallya, Dillon Davis, and Svetlana Lazebnik. Piggyback: Adapting a Single Network to Multiple Tasks by Learning to Mask Weights. In *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part IV*, volume 11208 of *Lecture Notes in Computer Science*, pages 72–88. Springer, 2018.
- [27] Arun Mallya and Svetlana Lazebnik. Packnet: Adding multiple tasks to a single network by iterative pruning. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7765–7773. IEEE, June 2018.
- [28] Martin Mundt, Yong Won Hong, Iuliia Pliushch, and Visvanathan Ramesh. A Wholistic View of Continual Learning with Deep Neural Networks: Forgotten Lessons and the Bridge to Active and Open World Learning. *CoRR*, abs/2009.01797:1–32, 2020.
- [29] German I. Parisi, Ronald Kemker, Jose L. Part, Christopher Kanan, and Stefan Wermter. Continual lifelong learning with neural networks: A review. *Neural Networks*, 113:54–71, 2019.
- [30] Benedikt Pfülb and Alexander Gepperth. A comprehensive, application-oriented study of catastrophic forgetting in DNNs. In *International Conference on Learning Representations*, 2019.
- [31] S. Rebuffi, A. Kolesnikov, G. Sperl, and C. H. Lampert. iCaRL: Incremental Classifier and Representation Learning. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5533–5542, July 2017.
- [32] Andrei A. Rusu, Neil C. Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. *CoRR*, abs/1606.04671, 2016.
- [33] Jonathan Schwarz, Jelena Luketina, Wojciech M. Czarnecki, Agnieszka Grabska-Barwinska, Yee Whye Teh, Razvan Pascanu, and Raia Hadsell. Progress & compress: A scalable framework for continual learning. *35th International Conference on Machine Learning, ICML 2018*, 10:7199–7208, 2018.
- [34] Joan Serra, Didac Suris, Marius Miron, and Alexandros Karatzoglou. Overcoming Catastrophic Forgetting with Hard Attention to the Task. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 4548–4557. PMLR, 10–15 Jul 2018.
- [35] Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual learning with deep generative replay. *Advances in Neural Information Processing Systems*, 2017-December(Nips):2991–3000, 2017.
- [36] Andrea Soltoggio, Kenneth O. Stanley, and Sebastian Risi. Born to Learn: the Inspiration, Progress, and Future of Evolved Plastic Artificial Neural Networks. *Neural Networks*, 108:48–67, 2018.
- [37] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms. pages 1–6, 2017.
- [38] Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual Learning Through Synaptic Intelligence. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 3987–3995. PMLR, 06–11 Aug 2017.