
Memory Replay GANs: learning to generate images from new categories without forgetting

Chenshen Wu, Luis Herranz, Xialei Liu, Yaxing Wang,
Joost van de Weijer, Bogdan Raducanu

Computer Vision Center

Universitat Autònoma de Barcelona, Spain

{chenshen, lherranz, xialei, yaxing, joost, bogdan}@cvc.uab.es

Abstract

Previous works on sequential learning address the problem of forgetting in discriminative models. In this paper we consider the case of generative models. In particular, we investigate generative adversarial networks (GANs) in the task of learning new categories in a sequential fashion. We first show that sequential fine tuning renders the network unable to properly generate images from previous categories (i.e. forgetting). Addressing this problem, we propose *Memory Replay GANs* (MeRGANs), a conditional GAN framework that integrates a memory replay generator. We study two methods to prevent forgetting by leveraging these replays, namely *joint training with replay* and *replay alignment*. Qualitative and quantitative experimental results in MNIST, SVHN and LSUN datasets show that our memory replay approach can generate competitive images while significantly mitigating the forgetting of previous categories.¹

1 Introduction

Generative adversarial networks (GANs) [6] are a popular framework for image generation due to their capability to learn a mapping between a low-dimensional latent space and a complex distribution of interest, such as natural images. The approach is based on an adversarial game between a generator that tries to generate good images and a discriminator that tries to discriminate between real training samples and generated. The original framework has been improved with new architectures [23, 10] and more robust losses [2, 7, 17].

GANs can be used to sample images by mapping a randomly sampled latent vector. While providing diversity, there is little control over the semantic properties of what is being generated. Conditional GANs [19] enable the use of semantic conditions as inputs, so the semantic properties and the inherent diversity can be decoupled. The simplest condition is just the category label, allowing to control the category of the generated image [22].

As most machine learning problems, image generation models have been studied in the conventional setting that assumes all training data is available at training time. This assumption can be unrealistic in practice, and modern neural networks face scenarios where tasks and data are not known in advance, requiring to continuously update their models upon the arrival of new data or new tasks. Unfortunately, neural networks suffer from severe degradation when they are updated in a sequential manner without revisiting data from previous tasks (known as *catastrophic forgetting* [18]). Most strategies to prevent forgetting in neural networks rely on regularizing weights [4, 15] or activations [14], keeping a small set of exemplars from previous categories [24, 16], or memory replay mechanisms [25, 27, 11].

¹The code is available at <https://github.com/WuChenshen/MeRGAN>

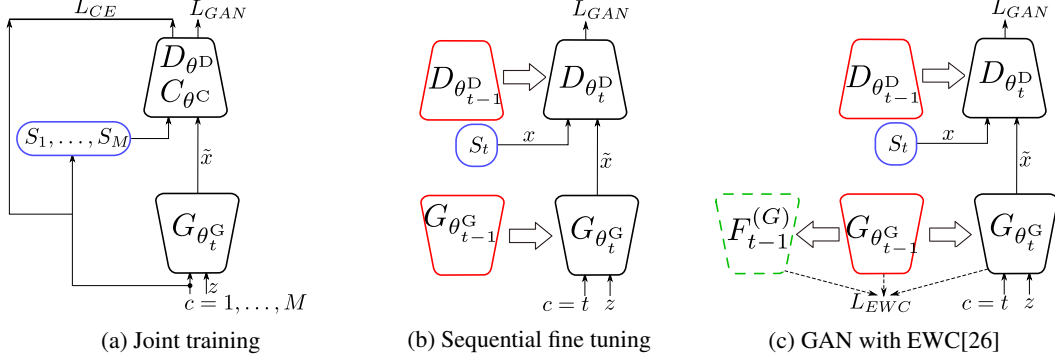


Figure 1: Baseline architectures.

While previous works study forgetting in discriminative tasks, in this paper we focus on forgetting in generative models (GANs in particular) through the problem of generating images when categories are presented sequentially as disjoint tasks. The closest related work is [26], that adapts elastic weight consolidation (EWC) [4] to GANs. In contrast, our method relies on memory replay and we describe two approaches to prevent forgetting by joint retraining and by aligning replays. The former includes replayed samples in the training process, while the latter forces to synchronize the replays of the current generator with those generated by an auxiliary generator (a snapshot taken before starting to learn the new task). An advantage of studying forgetting in image generation is that the dynamics of forgetting and consolidation can be observed visually through the generated images themselves.

2 Sequential learning in GANs

2.1 Joint learning

We first introduce our conditional GAN framework in the non-sequential setting, where all categories are learned jointly. In particular, this first baseline is based on the AC-GAN framework [22] combined with the WGAN-GP loss for robust training [7]. Using category labels as conditions, the task is to learn from a training set $S = \{S_1, \dots, S_M\}$ to generate images given an image category c . Each set S_c represents the training images for a particular category.

The framework consists of three components: generator, discriminator and classifier. The discriminator and classifier share all layers but the last ones (task-specific layers). The conditional generator is parametrized by θ^G and generates an image $\tilde{x} = G_{\theta^G}(z, c)$ given a latent vector z and a category c . In our case the conditioning is implemented via conditional batch normalization [3], that dynamically switches between sets of batch normalization parameters depending on c . Note that, in contrast to unconditional GANs, the latent vector is completely agnostic to the category, and the same latent vector can be used to generate images of different categories just by using a different c .

Similarly, the discriminator (parametrized by θ^D) tries to discern whether an input image x is real (i.e. from the training set) or generated, while the generator tries to fool it by generating more realistic images. In addition, AC-GAN uses an auxiliary classifier C with parameters θ^C to predict the label $\tilde{c} = C_{\theta^C}(x)$, and thus forcing the generator to generate images that can be classified in the same way as real images. This additional task improves the performance in the original task [22]. For convenience we represent all the parameters in the conditional GAN as $\theta = (\theta^G, \theta^D, \theta^C)$.

During training, the network is trained to solve both the adversarial game (using the WGAN with gradient penalty loss [7]) and the classification task by alternating the optimization of the generator, and the discriminator and classifier. The generator optimizes the following problem:

$$\min_{\theta^G} (L_{\text{GAN}}^G(\theta, S) + L_{\text{CLS}}^G(\theta, S)) \quad (1)$$

$$L_{\text{GAN}}^G(\theta, S) = -\mathbb{E}_{z \sim p_z, c \sim p_c} [D_{\theta^D}(G_{\theta^G}(z, c))] \quad (2)$$

$$L_{\text{CLS}}^G(\theta, S) = -\mathbb{E}_{z \sim p_z, c \sim p_c} [y_c \log C_{\theta^C}(G_{\theta^G}(z, c))] \quad (3)$$

where $L_{\text{GAN}}^{\text{G}}(\theta, S)$ and $\lambda_{\text{CLS}}L_{\text{CLS}}^{\text{G}}(\theta, S)$ are the corresponding GAN and cross-entropy loss for classification, respectively, S is the training set, $p_c = \mathcal{U}\{1, M\}$, $p_z = \mathcal{N}(0, 1)$ are the sampling distributions (uniform and Gaussian, respectively), and y_c is the one-hot encoding of c for computing the cross-entropy. The GAN loss uses the WGAN formulation with gradient penalty. Similarly, the optimization problem in the discriminator and classifier is

$$\min_{\theta^{\text{D}}, \theta^{\text{C}}} (L_{\text{GAN}}^{\text{D}}(\theta, S) + \lambda_{\text{CLS}}L_{\text{CLS}}^{\text{D}}(\theta, S)) \quad (4)$$

$$L_{\text{GAN}}^{\text{D}}(\theta, S) = -\mathbb{E}_{(x,c) \sim S} [D_{\theta^{\text{D}}}(x)] + \mathbb{E}_{z \sim p_z, c \sim p_c} [D_{\theta^{\text{D}}}(G_{\theta^{\text{G}}}(z, c))] \quad (5)$$

$$+ \lambda_{\text{GP}} \mathbb{E}_{x \sim S, z \sim p_z, c \sim p_c, \epsilon \sim p_{\epsilon}} \left[(\|\nabla D_{\theta^{\text{D}}}(\epsilon x + (1 - \epsilon)G_{\theta^{\text{G}}}(z, c))\|_2 - 1)^2 \right]$$

$$L_{\text{CLS}}^{\text{D}}(\theta, S) = -\mathbb{E}_{(x,c) \sim S} [C_{\theta^{\text{C}}}(G_{\theta^{\text{G}}}(z, c))] \quad (6)$$

where ϵ are parameters of the gradient penalty term, sampled as $p_{\epsilon} = \mathcal{U}(0, 1)$. The last term of $L_{\text{GAN}}^{\text{D}}$ is the gradient penalty.

2.2 Sequential fine tuning

Now we modify the previous framework to address the sequential learning scenario. We define a sequence of tasks $\mathbf{T} = (1, \dots, M)$, each of them corresponding to learning to generate images from a new training set S_t . For simplicity, we restrict each S_t to contain only images from a particular category c , i.e. $t = c$.

The joint training problem can be adapted easily to the sequential learning scenario as

$$\min_{\theta_t^{\text{G}}} L_{\text{GAN}}^{\text{G}}(\theta_t, S_t) \quad (7)$$

$$\min_{\theta_t^{\text{D}}} L_{\text{GAN}}^{\text{D}}(\theta_t, S_t) \quad (8)$$

where $\theta_t = (\theta_t^{\text{G}}, \theta_t^{\text{D}})$ are the parameters during task t , which are initialized as $\theta_t = \theta_{t-1}$, i.e. the current task t is learned immediately after finishing the previous task $t - 1$. Note that there is no classifier in this case since there is only data of the current category.

Unfortunately, when the network learns to adjust its parameters to generate images of the new domain via gradient descent, that very drifting away from the original solution for the previous task will cause catastrophic forgetting [18]. This has also been observed in GANs [26, 29] (shown later in Figures 3, 5 and 7 in the experiments section).

2.3 Preventing forgetting with Elastic Weight Consolidation

Catastrophic forgetting can be alleviated using samples from previous tasks [24, 16] or different types of regularization that result in penalizing large changes in parameters or activations [4, 14]. In particular, the elastic weight consolidation (EWC) regularization [4] has been adapted to prevent forgetting in GANs [26] and included as an augmented objective when training the generator as

$$\min_{\theta_t^{\text{G}}} L_{\text{GAN}}^{\text{G}}(\theta_t, S_t) + \sum_i \frac{\lambda_{\text{EWC}}}{2} F_{t-1,i} (\theta_{t,i}^{\text{G}} - \theta_{t-1,i}^{\text{G}})^2 \quad (9)$$

where $F_{t-1,i}$ is the Fisher information matrix that somewhat indicates how sensitive the parameter $\theta_{t,i}^{\text{G}}$ is to forgetting, and λ_{EWC} is a hyperparameter. We will use this approach as a baseline.

3 Memory replay generative adversarial networks

Rather than regularizing the parameters to prevent forgetting, we propose that the generator has an active role by replaying memories of previous tasks (via generative sampling), and using them during the training of current task to prevent forgetting. Our framework is extended with a replay generator, and we describe two different methods to leverage memory replays.

This replay mechanism (also known as pseudorehearsal [25]) resembles the role of the hippocampus in replaying memories during memory consolidation [5], and has been used to prevent forgetting in classifiers [11, 27], but to our knowledge has not been used to prevent forgetting in image generation. Note also that image generation is a generative task and typically more complex than classification.

3.1 Joint retraining with replayed samples

Our first method to leverage memory replays creates an extended dataset $S'_t = S_c \cup_{c \in \{1, \dots, t-1\}} \tilde{S}_c$ that contains both real training data for the current tasks and memory replays from previous tasks. The replay set \tilde{S}_c for a given category c typically samples a fixed number for replays $\hat{x} = G_{\theta_{t-1}^G}(z, c)$.

Once the extended dataset is created, the network is trained using joint training (see Fig. 2a) as

$$\min_{\theta_t^G} (L_{\text{GAN}}^G(\theta_t, S'_t) + \lambda_{\text{CLS}} L_{\text{CLS}}^G(\theta_t, S'_t)) \quad (10)$$

$$\min_{\theta_t^D} (L_{\text{GAN}}^D(\theta_t, S'_t) + \lambda_{\text{CLS}} L_{\text{CLS}}^D(\theta_t, S'_t)) \quad (11)$$

This method could be related to the deep generative replay in [27], where the authors use an unconditional GAN and the category is predicted with a classifier. In contrast, we use a conditional GAN where the category is an input, allowing us finer control of the replay process, with more reliable sampling of (x, c) pairs since we avoid potential classification errors and biased sampling towards recent categories.

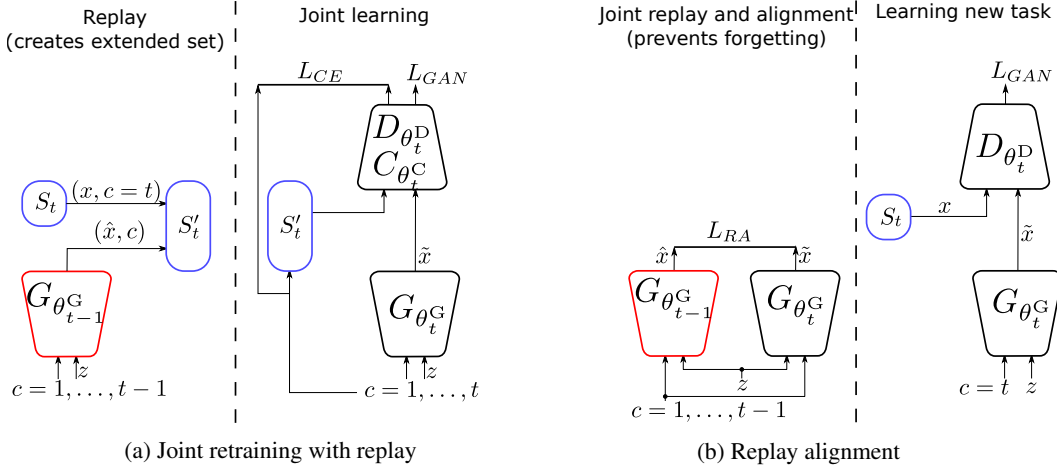


Figure 2: Memory Replay GANs and mechanisms to prevent forgetting (for a given current task t).

3.2 Replay alignment

We can also take advantage of the fact that the current generator and the replay generator share the same architecture, inputs and outputs. Their condition spaces (i.e. categories), and, critically, their latent spaces (i.e. latent vector z) and parameter spaces are also initially aligned, since the current generator is initialized with the same parameters of the replay generator. Therefore, we can synchronize both the replay generator and current one to generate the same image by the same category c and latent vector z as inputs (see Fig. 2b). In these conditions, the generated images \hat{x} and x should also be aligned pixelwise, so we can include a suitable pixelwise loss to prevent forgetting (we use L_2 loss).

In contrast to the previous method, in this case the discriminator is only trained with images of the current task, and there is no classification task. The problem optimized by the generator includes a replay alignment loss

$$\min_{\theta_t^G} L_{\text{GAN}}^G(\theta_t, S_t) + \lambda_{\text{RA}} L_{\text{RA}}(\theta_t, S_t) \quad (12)$$

$$L_{\text{RA}}(\theta_t, S_t) = \mathbb{E}_{x \sim S, z \sim p_z, c \sim \mathcal{U}\{1, t-1\}} \left[\left\| G_{\theta_t^G}(z, c) - G_{\theta_{t-1}^G}(z, c) \right\|^2 \right] \quad (13)$$

Note that in this case both generators engage in memory replay for all previous tasks. The corresponding problem in the discriminator is simply $\min_{\theta_t^D} L_{\text{GAN}}^D(\theta_t, S_t)$.

Our approach can be seen as *aligned distillation*, where distillation requires spatially aligned data. Note that in that way it could be related to the *learning without forgetting* approach [14] to prevent forgetting. However, we want to emphasize several subtle yet important differences:

Different tasks and data Our task is image generation where outputs have a spatial structure (i.e. images), while in [14] the task is classification and the output is a vector of category probabilities.

Spatial alignment Image generation is a one-to-many task with many latent factors of variability (e.g. pose, location, color) that can result in completely different images yet sharing the same input category. The latent vector z somewhat captures those factors and allows a unique solution for a given (z, c) . However, pixelwise comparison of the generated images requires that not only the input but also the output representations are aligned, which is ensured in our case since at the beginning of training both have the same parameters. Therefore we can use a pixelwise loss.

Task-agnostic inputs and seen categories In [14], images of the current classification task are used as inputs to extract output features for distillation. Note that this implicitly involves a domain shift, since a particular input image is always linked to an unseen category (by definition, in the sequential learning problem the network cannot be presented with images of previous tasks), and therefore the outputs for the old task suffer from domain shift. In contrast, our approach does not suffer from that problem since the inputs are not real data, but a category-agnostic latent vector z and a category label c . In addition, we only replay seen categories for both generators, i.e. 1 to $t - 1$.

4 Experimental results

We evaluated the proposed approaches in different datasets with different level of complexity. The architecture and settings are set accordingly. We use the Tensorflow [1] framework with Adam optimizer [12], learning rate $1e-4$, batch size 64 and fixed parameters for all experiments: $\lambda_{EWC} = 1e9$, $\lambda_{RA} = 1e-3$ and $\lambda_{CLS} = 1$ except for $\lambda_{RA} = 1e-2$ on SVHN dataset.

4.1 Digit generation

We first consider the digit generation problem in two standard digit datasets. Learning to generate a digit category is considered as a separate task. MNIST [13] consists of images of handwritten digits which are resized 32×32 pixels in our experiment. SVHN [21] contains cropped digits of house numbers from real-world street images. The generation task is more challenging since SVHN contains much more variability than MNIST, with diverse backgrounds, variable illumination, font types, rotated digits, etc.

The architecture used in the experiments is based on the combination of AC-GAN and Wasserstein loss described in Section 2.1. We evaluated the two variants of the proposed memory replay GANs: joint training with replay (MeRGAN-JTR) and replay alignment (MeRGAN-RA). As upper and lower bounds we also evaluated joint training (JT) with all data (i.e. non-sequential) and sequential fine tuning (SFT). We also implemented two additional methods based on related works: the adaptation of EWC to conditional GANs proposed by [26], and the deep generative replay (DGR) module of [27], implemented as an unconditional GAN followed by a classifier to predict the label. For experiments with memory replay we generate one batch of replayed samples (including all tasks) for every batch of real data. We use a three layer DCGAN [23] for both datasets. In order to compare the methods in a more challenging setting, we keep the capacity of the network relatively limited for SVHN.

Figure 3 compares the images generated by the different methods after sequentially training the ten tasks. Since DGR is unconditional, the category for visualization is the one predicted by its classifier. We observe that SFT completely forgets previous tasks in both datasets, while the other methods show different degrees of forgetting. The four methods are able to generate MNIST digits properly, although both MeRGANs show sharper ones. In the more challenging setting of SVHN (note that the JT baseline also struggles to generate realistic images), the digits generated by EWC are hardly recognizable, while DGR is more unpredictable, sometimes generates good images but often generating images with ambiguous digits. Those generated by MeRGANs are in general clear and more recognizable, but still showing some degradation due to the limited capacity of the network.

We also trained a classifier with real data, using classification accuracy as a proxy to evaluate forgetting. The rationale behind is that in general bad quality images will confuse the classifier and

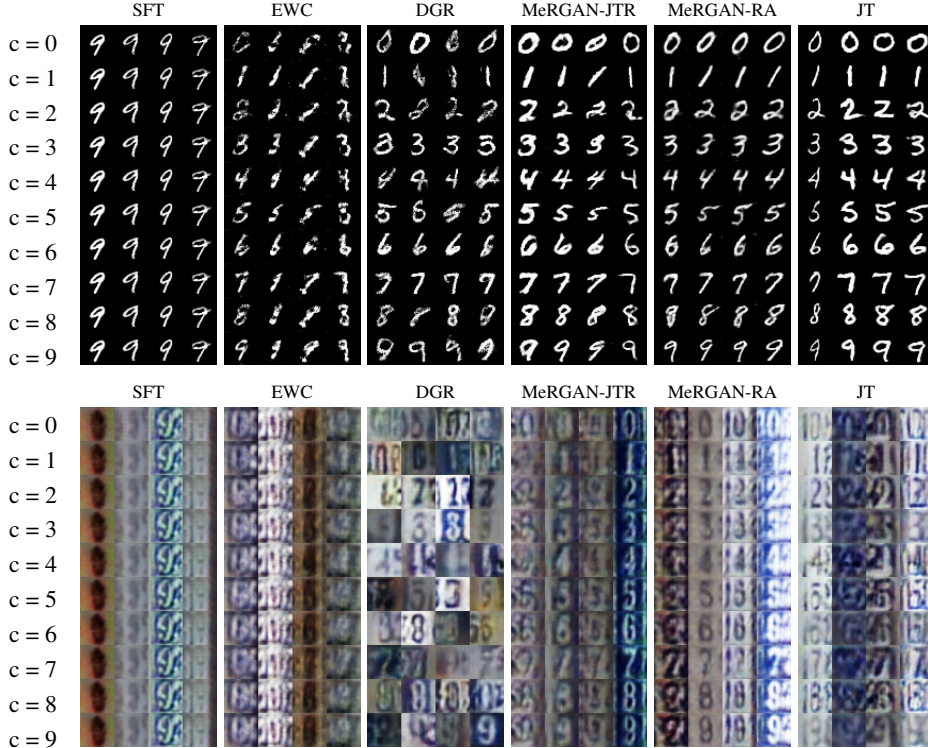


Figure 3: Images generated for MNIST and SVHN after learning the ten tasks. Rows are different conditions (i.e. categories), and columns are different latent vectors.

Table 1: Average classification accuracy (%) in digit generation (ten sequential tasks).

	5 tasks (0-4)										10 tasks (0-9)			
	Baselines		Others		MeRGAN		Baselines		Others		MeRGAN			
	JT	SFT	EWC[26]	DGR[27]	JTR	RA	JT	SFT	EWC[26]	DGR[27]	JTR	RA		
MNIST	97.66	19.87	70.62	90.39	97.93	98.19	96.92	10.06	77.03	85.40	97.00	97.01		
SVHN	85.30	19.35	39.84	61.29	80.90	76.05	84.82	10.10	33.02	47.28	66.50	66.78		

result in lower classification rates. Table 1² shows the classification accuracy after the first five tasks (digits 0 to 4) and after the ten tasks. SFT forgets previous tasks so the accuracy is very low. As expected, EWC performs worse than DGR since it does not leverage replays, however it significantly mitigates the phenomenon of catastrophic forgetting by increasing the accuracy from 19.87 to 70.62 on MNIST, and from 19.35 to 39.84 on SVHN compared to SFT in the case of 5 tasks. The same conclusion can be drawn in the case of 10 tasks. By using the memory replay mechanism, MeRGANs obtain significant improvement compared to the baselines and the others related methods. Especially, our approach performs about 8% better on MNIST and about 21% better on SVHN compared to the strong baseline DGR in the case of 5 tasks. Note that our approach achieves about 12% gain in the case of 10 tasks, which shows that our approach is much more stable with increasing number of tasks. In the more challenging SVHN dataset, all methods decrease in terms of accuracy, however MeRGAN are able to mitigate forgetting and obtain comparable results to JT.

Another interesting way to compare the different methods is through t-SNE visualizations. We use a classifier trained with real digits to extract embeddings of the methods to compare. Fig. 4a shows real 0s from MNIST and generated 0s from the different methods after training 10 tasks (i.e. the first task, and therefore the most difficult to remember). In contrast to SFT and EWC, the distributions of 0s generated by MeRGANs greatly overlap with the distribution of real 0s (in red) and no isolated clusters of real samples are observed, which suggests that MeRGANs prevent forgetting better while

²The reverse classification accuracy can be found in Appendix A.

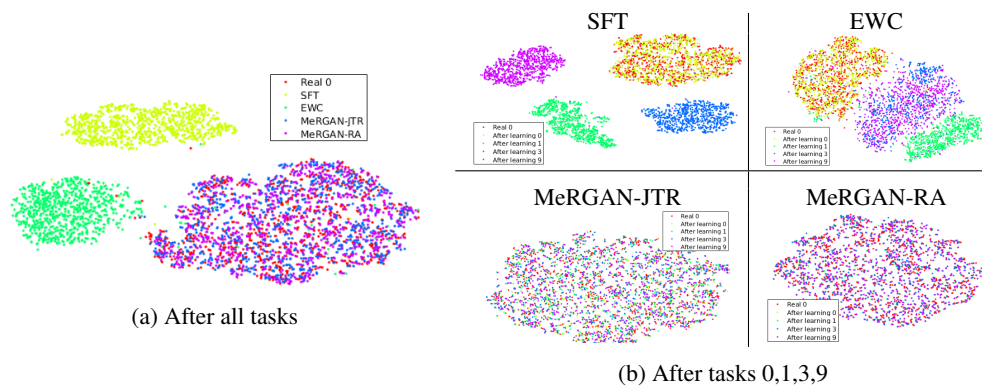


Figure 4: t-SNE visualization of generated 0s. Real 0s correspond to red dots. Please view in electronic format with zooming.

Table 2: FID and average classification accuracy (%) on LSUN after the 4th task

	SFT	EWC	DGR	MeRGAN-JTR	MeRGAN-RA
Acc.(%)	15.02	14.28	15.40	79.19	81.03
Rev acc.(%)	28.0	63.35	26.17	70.00	83.62
FID	110.12	178.05	93.70	49.69	37.73

keeping diversity (at least in the t-SNE visualizations). Fig. 4b shows the t-SNE visualizations of real and 0s generated after learning 0,1,3 and 9, with similar conclusions.

4.2 Scene generation

We also evaluated MeRGANs in a more challenging domain and on higher resolution images (64×64 pixels) using four scene categories of the LSUN dataset [30]. The experiment consists of a sequence of tasks, each one involving learning the generative distribution of a new category. The sequence of categories is *bedroom*, *kitchen*, *church (outdoors)* and *tower*, in this order. This sequence allows us to have two indoor and outdoor categories, and transitions between relatively similar categories (*bedroom* to *kitchen* and *church* to *tower*) and also a transition between very different categories (*kitchen* to *church*). Each category is represented by a set of 100000 training images, and the network is trained during 20000 iterations for every task. The architectures are based on [7] with 18-layer ResNet[8] generator and discriminator, and for every batch of training data for the new category we generate a batch of replayed images per category.

Figure 5 shows examples of generated images. Each block column corresponds to a different method, and inside, each row shows images generated for a particular condition (i.e. category) and each column corresponds to images generated after learning a particular task, using the same latent vector. Note that we excluded DGR since the generation is not conditioned on the category. We can observe that SFT completely forgets the previous task, and essentially ignores the category condition. EWC generates images that have characteristics of both new and previous tasks (e.g. bluish outdoor colors, indoor shapes), being unable to neither successfully learn new tasks nor remember previous ones. In contrast both variants of MeRGAN are able to generate competitive images of new categories while still remembering to generate images of previous categories.

In addition to classification accuracy (using a VGG [28] trained over the ten categories in LSUN), for this dataset we add two additional measurements. The first one is reverse accuracy measured by a classifier trained with generated data and evaluated with real data. The second one is the Frechet Inception Distance (FID), which is widely used to evaluate the images generated by GANs. Note that FID is sensitive to both quality and diversity[9]. Table 2 shows these metrics after the four tasks are learned. MeRGANs perform better in this more complex and challenging setting, where EWC and DGR are severely degraded.

Figure 6 shows the evolution of these metrics during the whole training process, including transitions to new tasks (the curves have been smoothed for easier visualization). We can observe not only that

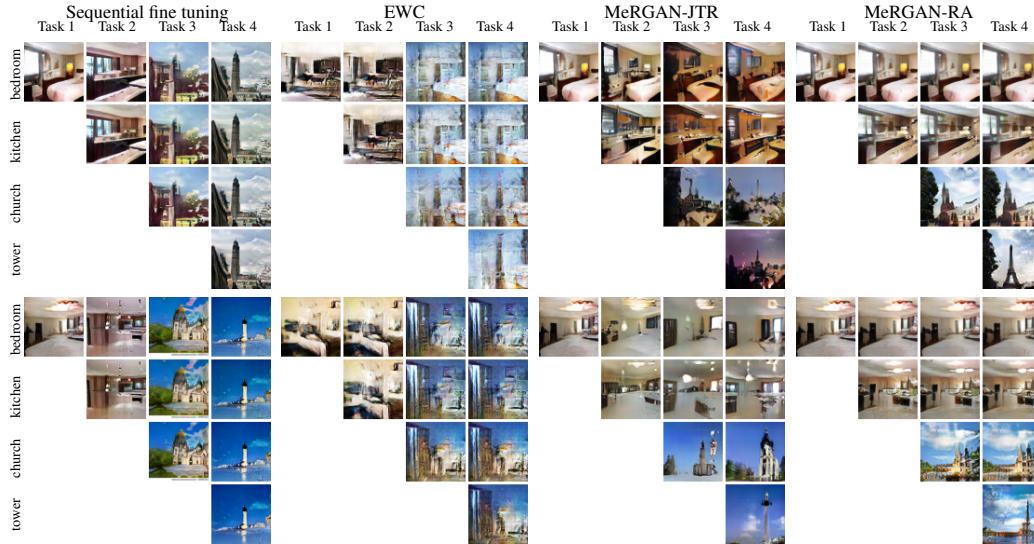


Figure 5: Images generated after sequentially learning each task (column within each block) for different methods (block column), two different latent vectors z (block row) and different conditions c (row within each block). The network learned after the first task is the same in all methods. Note that fine tuning forgets previous tasks completely, while the proposed methods still remember them.

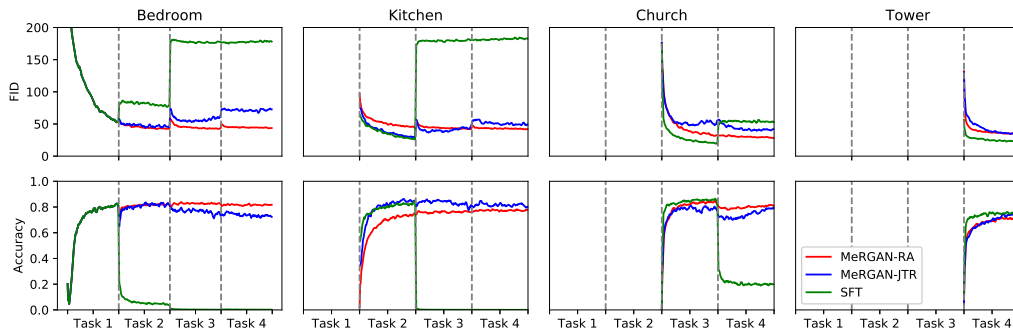


Figure 6: Evolution of FID and classification accuracy (%). Best viewed in color.

sequential fine tuning forgets the task completely, but also that it happens early during the first few iterations. This also allows the network to exploit its full capacity to focus on the new task and learn it quickly. MeRGANs experience forgetting during the initial iterations of a new task but then tend to recover during the training process. In this experiment MeRGAN-RA seems to be more stable and slightly more effective than MeRGAN-JTR.

Figure 6 provides useful insight about the dynamics of learning and forgetting in sequential learning. The evolution of generated images also provides complementary insight, as in the *bedroom* images shown in Figure 7, where we pay special attention to the first iterations. The transition between task 2 to 3 (i.e. *kitchen* to *church*) is particularly revealing, since this new task requires the network to learn to generate many completely new visual patterns found in outdoor scenes. The most clear example is the need to develop filters that can generate the blue sky regions, that are not found in the previous indoor categories seen during task 1 and 2. Since the network is not equipped with knowledge to generate the blue sky, the new task has to reuse and adapt previous one, interfering with previous tasks and causing forgetting. This interference can be observed clearly in the first iterations of task 3 where the walls of bedroom (and kitchen) images turn blue (also related with the peak in forgetting observed at the same iterations in Figure 6). MeRGANs provide mechanisms that penalize forgetting, forcing the network to develop separate filters for the different patterns (e.g. separated filters for wall and sky). MeRGAN-JTR seems to effectively decouple both patterns, since we do not observe the same "blue walls" interference during task 4. Interestingly, the same interference seems to be

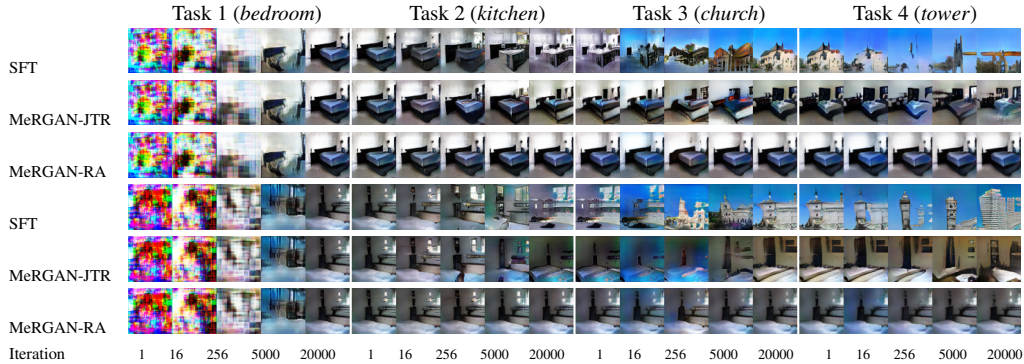


Figure 7: Evolution of the generated images (category *bedroom* and two different values of z) during the sequential learning process (rows). Sequential fine tuning forgets the previous task after just a few iterations (iterations within each task are sampled in a logarithmic fashion). Note that fine tuning forgets previous tasks completely, while the MeRGANs still remember them.

milder in MeRGAN-RA, but recurrent, since it also appears again during task 4. Nevertheless, the interference is still temporary and disappears after a few iterations more.

Another interesting observation from Figures 5 and 7 is that MeRGAN-RA remembers the *same* bedroom (e.g. same point of view, colors, objects), which is related to the replay alignment mechanism that enforces remembering the instance. On the other hand, MeRGAN-JTR remembers bedrooms *in general* as the generated image still resembles a bedroom but not exactly the same one as in previous steps. This can be explained by the fact that the classifier and the joint training mechanism enforce the not-forgetting constraint at the category level.

Conclusions

We have studied the problem of sequential learning in the context of image generation with GANs, where the main challenge is to effectively address catastrophic forgetting. MeRGANs incorporate memory replay as the main mechanism to prevent forgetting, which is then enforced through either joint training or replay alignment. Our results show their effectiveness in retaining the ability to generate competitive images of previous tasks even after learning several new ones. In addition to the application in pure image generation, we believe MeRGANs and generative models robust to forgetting in general, could have important application in many other tasks. We also showed that image generation provides an interesting way to visualize the interference between tasks and potential forgetting by directly observing generated images.

Acknowledgements

C. Wu, X. Liu, and Y. Wang, acknowledge the Chinese Scholarship Council (CSC) grant No.201709110103, No.201506290018 and No.201507040048. Luis Herranz acknowledges the European Union research and innovation program under the Marie Skłodowska-Curie grant agreement No. 6655919. This work was supported by TIN2016-79717-R, and the CHISTERA project M2CR (PCIN-2015-251) of the Spanish Ministry, the ACCIO agency and CERCA Programme / Generalitat de Catalunya, and the EU Project CybSpeed MSCA-RISE-2017-777720. We also acknowledge the generous GPU support from NVIDIA.

References

- [1] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, et al. Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *arXiv preprint arXiv:1603.04467*, 2016.
- [2] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein gan. *arXiv preprint arXiv:1701.07875*, 2017.
- [3] Vincent Dumoulin, Jonathon Shlens, and Manjunath Kudlur. A learned representation for artistic style. In *ICLR*, 2017.
- [4] James Kirkpatrick et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences of the United States of America*, 14(13):3521–3526, 2017.
- [5] Steffen Gais et al. Sleep transforms the cerebral trace of declarative memories. *Proceedings of the National Academy of Sciences*, 104(47):18778–18783, 2007.
- [6] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NIPS*, 2014.
- [7] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of wasserstein gans. In *NIPS*, 2017.
- [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [9] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, Günter Klambauer, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a nash equilibrium. In *NIPS*, 2017.
- [10] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. In *ICLR*, 2018.
- [11] Ronald Kemker and Christopher Kanan. Fearnnet: Brain-inspired model for incremental learning. In *ICLR*, 2018.
- [12] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- [13] Yann LeCun. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- [14] Zhizhong Li and Derek Hoiem. Learning without forgetting. In *ECCV*, 2016.
- [15] Xialei Liu, Marc Masana, Luis Herranz, Joost Van de Weijer, Antonio M Lopez, and Andrew D Bagdanov. Rotate your networks: Better weight consolidation and less catastrophic forgetting. In *ICPR*, 2018.
- [16] David Lopez-Paz et al. Gradient episodic memory for continual learning. In *NIPS*, 2017.
- [17] Xudong Mao, Qing Li, Haoran Xie, Raymond YK Lau, Zhen Wang, and Stephen Paul Smolley. Least squares generative adversarial networks. In *ICCV*, 2017.
- [18] Michael McCloskey and Neal J. Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. *The psychology of learning and motivation*, 24:109–165, 1989.
- [19] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.
- [20] Takeru Miyato and Masanori Koyama. Cgans with projection discriminator. *arXiv preprint arXiv:1802.05637v1*, 2018.
- [21] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y. Ng. Reading digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning*, 2011.
- [22] Augustus Odena, Christopher Olah, and Jonathon Shlens. Conditional image synthesis with auxiliary classifier gans. In *ICML*, 2017.
- [23] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. In *ICLR*, 2016.
- [24] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Cristoph H. Lampert. icarl: Incremental classifier and representation learning. In *CVPR*, 2017.
- [25] Anthony Robins. Catastrophic forgetting, rehearsal and pseudorehearsal. *Connection Science*, 7(2):123–146, 1995.
- [26] Ari Seff, Alex Beatson, Daniel Suo, and Han Liu. Continual learning in generative adversarial nets. *arXiv preprint arXiv:1705.08395v1*, 2017.
- [27] Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual learning with deep generative replay. In *NIPS*, 2017.

- [28] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.
- [29] Yaxing Wang, Chenshen Wu, Luis Herranz, Joost van de Weijer, Abel Gonzalez-Garcia, and Bogdan Raducanu. Transferring GANs: generating images from limited data. In *ECCV*, 2018.
- [30] Fisher Yu, Yinda Zhang, Shuran Song, Ari Seff, and Jianxiong Xiao. LSUN: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365*, 2015.

Table 3: Direct and reverse classification accuracy (%) in digit generation (ten sequential tasks).

	5 tasks (0-4)				10 tasks (0-9)			
	MeRGAN-JTR		MeRGAN-RA		MeRGAN-JTR		MeRGAN-RA	
	direct	reverse	direct	reverse	direct	reverse	direct	reverse
MNIST	97.93	99.15	98.19	98.55	97.00	96.83	97.01	93.86
SVHN	80.90	40.71	76.05	74.83	66.50	20.07	66.78	49.43
SVHN(ResNet-18)	82.30	81.12	81.74	82.30	73.29	51.00	78.56	70.98

A Reverse Classification Accuracy

In Table 1, we showed the result of direct classification accuracy in MNIST and SVHN dataset³. The accuracies are measured by training a classifier with the real data and then testing it with the generated image. In addition to that, we also evaluated the reverse accuracy by training a classifier with the data generated by the MeRGAN and testing it on the real dataset. The result are shown in table 3.

On the MNIST dataset, both method, MeRGAN-JTR and MeRGAN-RA, produce very high accuracy in 5 tasks setting (digits 0 to 4), 0.992 and 0.985 respectively, and keep relatively high accuracy in 10 tasks setting, 0.968 and 0.939 respectively. However, on the SVHN dataset there is a huge drop compared with the direct classification accuracy. In the 10 tasks setting, the reverse classification accuracy of JTR is only 0.201.

We found that this drop in performance can be reduced by improving the architecture of the GAN. Therefore, we also trained a GAN with ResNet-18 network[8]. This network was further improved by replacing the AC-GAN with cGAN with projection discriminator [20] and using one-hot conditioning [19]. The results of JTR and RA are 51.0 and 71.0 respectively. It shows that result of the reverse classification accuracy can be improved by using a better GAN. It also shows that in a more complex case, the RA works better than JTR.

³This Appendix has been added in 9-2019.