

---

# GAN Memory with No Forgetting

---

Yulai Cong\*    Miaoyun Zhao\*    Jianqiao Li    Sijia Wang    Lawrence Carin  
Department of Electrical and Computer Engineering  
Duke University

## Abstract

As a fundamental issue in lifelong learning, catastrophic forgetting is directly caused by inaccessible historical data; accordingly, if the data (information) were memorized perfectly, no forgetting should be expected. Motivated by that, we propose a GAN memory for lifelong learning, which is capable of remembering a stream of datasets via generative processes, with *no* forgetting. Our GAN memory is based on recognizing that one can modulate the “style” of a GAN model to form perceptually-distant targeted generation. Accordingly, we propose to do sequential style modulations atop a well-behaved base GAN model, to form sequential targeted generative models, while simultaneously benefiting from the transferred base knowledge. The GAN memory – that is motivated by lifelong learning – is therefore itself manifested by a form of lifelong learning, via forward transfer and modulation of information from prior tasks. Experiments demonstrate the superiority of our method over existing approaches and its effectiveness in alleviating catastrophic forgetting for lifelong classification problems. Code is available at [https://github.com/MiaoyunZhao/GANmemory\\_LifelongLearning](https://github.com/MiaoyunZhao/GANmemory_LifelongLearning).

## 1 Introduction

Lifelong learning (or continual learning) is a long-standing challenge for machine learning and artificial intelligence systems [76, 28, 73, 11, 14, 60], concerning the ability of a model to continually learn new knowledge without forgetting previously learned experiences. An important issue associated with lifelong learning is the notorious catastrophic forgetting of deep neural networks [48, 36, 87], *i.e.*, training a model with new information severely interferes with previously learned knowledge.

To alleviate catastrophic forgetting, many methods have been proposed, with most focusing on discriminative/classification tasks [36, 65, 95, 55, 94]. Reviewing existing methods, [77] revealed generative replay (or pseudo-rehearsal) [72, 69, 86, 66, 88] is an effective and general strategy for lifelong learning, with this further supported by [40, 78]. That revelation is anticipated, for if the characteristics of previous data are remembered perfectly (*e.g.*, via realistic generative replay), no forgetting should be expected for lifelong learning. Compared with the coreset idea, that saves representative samples of previous data [55, 65, 11], generative replay has advantages in addressing privacy concerns and remembering potentially more complete data information (via the generative process). However, most existing generative replay methods either deliver blurry generated samples [10, 40] or only work well on simple datasets [40, 78, 40] like MNIST; besides, they often don’t scale well to practical situations with high resolution [60] or a long sequence [86], sometimes even with negative backward transfer [96, 82]. Therefore, it’s challenging to continually learn a well-behaved generative replay model [40], even for moderately complex datasets like CIFAR10.

We seek a realistic generative replay framework to alleviate catastrophic forgetting; going further, we consider developing a realistic generative memory with growing (expressive) power, believed to be a fundamental building block toward general lifelong learning systems. We leverage the popular

---

\*Equal Contribution. Correspondence to: Miaoyun Zhao <miaoyun9zhao@gmail.com>.

GAN [25] setup as the key component of that generative memory, which we term GAN memory, because (i) GANs have shown remarkable power in synthesizing realistic high-dimensional samples [9, 51, 33, 34]; (ii) by modeling the generative process of training data, GANs summarize the data statistical information in the model parameters, consequently also protecting privacy (the original data need not be saved); and (iii) a GAN often generates realistic samples not observed in training data, delivering a synthetic data augmentation that potentially benefits better performance of downstream tasks [80, 8, 9, 21, 33, 26, 27]. Distinct from existing methods, our GAN memory leverages transfer learning [6, 16, 46, 92, 58] and (image) style transfer [18, 30, 41]. Its key foundation is a discovery that one can leverage the modified variants of style-transfer techniques [64, 98] to modulate a source generator/discriminator into a powerful generator/discriminator for perceptually-distant target domains (see Section 4.1), with a limited amount of style parameters. Exploiting that discovery, our GAN memory sequentially modulates (and also transfers knowledge from) a well-behaved base/source GAN model to realistically remember a sequence of (target) generative processes with *no* forgetting. Note by “well-behaved” we mean the shape of source kernels is well trained (see Section 4.1 for details); empirically, this requirement can be *readily satisfied* if (i) the source model is pretrained on a (moderately) large dataset (*e.g.*, CelebA [43]; often a dense dataset is preferred [85]) and (ii) it’s sufficiently trained and shows relatively high generation quality. Therefore, many pretrained GANs can be “well-behaved”,<sup>2</sup> showing great flexibility in selecting the base/source model. Our experiments will show that flexibility roughly means source and target data should have the same data *type* (*e.g.*, images).

Our GAN memory serves as a solution to the fundamental memory issue of general lifelong learning, and its construction also leverages a form of lifelong learning. In practice, the GAN memory can be used, for example, as a realistic generative replay to alleviate catastrophic forgetting for challenging downstream tasks with high-dimensional data and a long (and varying) task sequence. Our contributions are as follows.

- Based on FiLM [64] and AdaFM [98], we develop modified variants, termed mFiLM and mAdaFM, to better adapt/transfer the source fully connected (FC) and convolutional (Conv) layers to target domains, respectively. We demonstrate that mFiLM and mAdaFM can be leveraged to modulate the “style” of a source GAN model (including both the generator and discriminator) to form a generative/discriminative model capable of addressing a perceptually-distant target domain.
- Based on the above discovery, we propose our GAN memory, endowed with growing (expressive) generative power, yet with *no* forgetting of existing capabilities, by leveraging a limited amount of task-specific style parameters. We analyze the roles played by those style parameters and reveal their further compressibility.
- We generalize our GAN memory to its conditional variant, followed by empirically verifying its effectiveness in delivering realistic synthesized samples to alleviate catastrophic forgetting for challenging lifelong classification tasks.

## 2 Related work

**Lifelong learning** Existing lifelong learning methods can be roughly grouped into three categories, *i.e.*, regularization-based [36, 69, 95, 55, 68], dynamic-model-based [47, 68, 47], and generative-replay-based methods [72, 42, 86, 66, 88, 78]. Among these methods, generative replay is believed an effective and general strategy for lifelong learning problems [66, 88, 40, 78], as discussed above. However, most existing methods of this type often have blurry/distorted generation (for images) or scalability issues [86, 66, 59, 60, 96, 54, 82]. MeRGAN [86] leverages a copy of the current generator to replay previous data information, showing increasingly blurry historical generations with reinforced generation artifacts [96] as training proceeds. CloGAN [66] uses an auxiliary classifier to filter out a portion of distorted replay but may still suffer from the reinforced forgetting, especially in high-dimensional situations. OCDVAE [54] unifies open-set recognition and VAE-based generative replay, whereas showing blurry generations for high-resolution images. Based on a shared generator, DGMw [59] introduces task-specific binary masks to the generator weights, accordingly suffering from scalability issues when the generator is large and/or the task sequence is long. Lifelong GAN [96] employs cycle consistency (via an auxiliary encoder) and knowledge distillation (via copies of

<sup>2</sup> One can of course expect better performance if a better source model (pretrained on a large-scale dense and diverse dataset) is used.

that encoder and the generator) to remember image-conditioned generation; however, it still shows decreased performance on historical tasks. By comparison, our GAN memory delivers realistic synthesis with *no* forgetting and scales well to high-dimensional situations with a long task-sequence, capable of serving as a realistic generative memory for general lifelong learning systems.

**Transfer learning** Being general and effective, transfer learning has attracted increasing attention recently in various research fields, with a focus on discriminative tasks like classification [6, 70, 90, 93, 16, 24, 44, 46, 74, 92] and those in natural language processing [3, 63, 53, 52, 2]. However for generative tasks, only a few efforts have been made [85, 58, 98]. For example, a GAN pretrained on a large-scale source dataset is used in [85] to initialize the GAN model in a target domain, for efficient training or even better performance; alternatively, [58] freezes the source GAN generator only to modulate its hidden-layer statistics to “add” new generation power with  $L1$ /perceptual loss; observing the general applicability of low-level filters in GAN generator/discriminator, [98] transfers-then-freezes them to facilitate generation with limited data in a target domain. Those methods either only concern synthesis in the target domain (completely forgetting source generation) [85, 98] or deliver blurry target generation [58]. Our GAN memory provides both realistic target generation and no forgetting on source generation.

### 3 Preliminary

We briefly review two building blocks on which our method is constructed: generative adversarial networks (GANs) [25, 33] and style-transfer techniques [64, 98].

**Generative adversarial networks (GANs)** GANs have shown increasing power to synthesize highly realistic observations [32, 45, 51, 9, 33, 34], and have found wide applicability in various fields [39, 1, 19, 81, 83, 84, 12, 89, 37]. A GAN often consists of a generator  $G$  and a discriminator  $D$ , with both trained adversarially with objective

$$\min_G \max_D \mathbb{E}_{\mathbf{x} \sim q_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))] \quad (1)$$

where  $p(\mathbf{z})$  is a simple distribution (*e.g.*, Gaussian) and  $q_{\text{data}}(\mathbf{x})$  is the underlying (unknown) data distribution from which we observe samples.

**Style-transfer techniques** An extensive literature [23, 71, 79, 13, 62, 38] has explored how one can manipulate the style of an image (*e.g.*, the texture [18, 30, 41] or attributes [33, 34]) by modulating the statistics of its latent features. These methods use style-transfer techniques like conditional instance normalization [18] or adaptive instance normalization [30], most of which are related to Feature-wise Linear Modulation (FiLM) [64]. FiLM imposes simple element-wise affine transformations to latent features of a neural network, showing remarkable effectiveness in various domains [30, 17, 33, 58]. Given a  $d$ -dimensional feature  $\mathbf{h} \in \mathbb{R}^d$  from a layer of a neural network,<sup>3</sup> FiLM yields

$$\hat{\mathbf{h}} = \gamma \odot \mathbf{h} + \beta, \quad (2)$$

where  $\hat{\mathbf{h}}$  is forwarded to the next layer,  $\odot$  denotes the Hadamard product, and the scale  $\gamma \in \mathbb{R}^d$  and shift  $\beta \in \mathbb{R}^d$  may be conditioned on other information [18, 64, 17]. Different from FiLM modulating latent features, another technique named adaptive filter modulation (AdaFM) modulates source convolutional (Conv) filters to manipulate its “style” to deliver a boosted transfer performance [98]. Specifically, given a Conv filter  $\mathbf{W} \in \mathbb{R}^{C_{\text{out}} \times C_{\text{in}} \times K_1 \times K_2}$ , where  $C_{\text{in}}/C_{\text{out}}$  denotes the number of input/output channels and  $K_1 \times K_2$  is the kernel size, AdaFM yields

$$\hat{\mathbf{W}} = \mathbf{\Gamma} \odot \mathbf{W} + \mathbf{B}, \quad (3)$$

where the scale matrix  $\mathbf{\Gamma} \in \mathbb{R}^{C_{\text{out}} \times C_{\text{in}}}$ , the shift matrix  $\mathbf{B} \in \mathbb{R}^{C_{\text{out}} \times C_{\text{in}}}$ , and the modulated  $\hat{\mathbf{W}}$  is used to convolve with input feature maps for output ones.

### 4 Proposed method

Targeting the fundamental memory issue of lifelong learning, we propose to exploit popular GANs to design a realistic *generative* memory (named GAN memory) to sequentially remember data-generating processes. Specifically, we consider a lifelong generation problem: the GAN memory

<sup>3</sup>For simplicity, we omit layer-index notation throughout the paper.

sequentially accesses a stream of datasets/tasks  $\{\mathcal{D}_1, \mathcal{D}_2, \dots\}$ <sup>4</sup> (during task  $t$ , only  $\mathcal{D}_t$  is accessible); after task  $t$ , the GAN memory should be able to synthesize realistic samples resembling  $\{\mathcal{D}_1, \dots, \mathcal{D}_t\}$ . Below we first reveal a surprising discovery that lays the foundation of the paper. We then build on top of it our GAN memory followed by a detailed analysis, and finally compression techniques are presented to facilitate our GAN memory for lifelong problems with a long task sequence.

#### 4.1 A surprising discovery

Moving well beyond the style-transfer literature modulating image features to manipulate its style [18, 30, 17], we discover that one can even modulate the “style” of a source generative/discriminative process (e.g., a GAN generator/discriminator trained on a source dataset  $\mathcal{D}_0$ ) to form synthesis power for a perceptually-distant target domain (e.g., a generative/discriminative power on  $\mathcal{D}_1$ ), via manipulating its FC and Conv layers with the style-modulation techniques developed below. Note different from the classical style-transfer literature, the “style” terminology here is associated with the characteristics of a function (e.g., for a GAN generator, its style manifests as the generation content); because of the similarity in mathematics, we reuse that terminology but name our approach as style-modulation techniques.

Before introducing the technical details, we emphasize our basic assumption of well-behaved source FC and Conv parameters; often parameters from a GAN model trained on large-scale datasets satisfy that assumption, as discussed in the Introduction. To highlight our discovery, we choose a moderately sophisticated GP-GAN model [49] trained on the CelebA [43] (containing only faces) as the source,<sup>5</sup> and select perceptually-distant target datasets including Flowers, Cathedrals, Cats, Brain-tumor images, Chest X-rays, and Anime images (see Figure 5 and Section 5.1). With the style-modulation techniques detailed below, we observe realistic generations in all target domains (see Figure 5), even though the generation power is modulated from an entirely different source domain. Alternatively, given a specific target domain, that observation also implies the flexibility in choosing a source model (see also Appendix H), *i.e.*, the source (with well-behaved parameters) should have the same target data type, but it need not be related to the target domain. In the context of image-based data, this implies a certain universal structure to images, that may be captured within a GAN by one (relatively large) image dataset. Via appropriate style modulation of this model, it can be adapted to new and very different image classes, using potentially limited observations from those target domains.

We next present the style-modulation techniques employed here, modified FiLM (mFiLM) and modified AdaFM (mAdaFM), for modulating FC and Conv layers, respectively.

**FC layers** Given a source FC layer  $\mathbf{h}^{\text{source}} = \mathbf{W}\mathbf{z} + \mathbf{b}$  with weight  $\mathbf{W} \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$ , bias  $\mathbf{b} \in \mathbb{R}^{d_{\text{out}}}$ , and input  $\mathbf{z} \in \mathbb{R}^{d_{\text{in}}}$ , mFiLM modulates its *parameters* to form a target function  $\mathbf{h}^{\text{target}} = \hat{\mathbf{W}}\mathbf{z} + \hat{\mathbf{b}}$  with

$$\hat{\mathbf{W}} = \gamma \odot \frac{\mathbf{W} - \boldsymbol{\mu}}{\boldsymbol{\sigma}} + \boldsymbol{\beta}, \quad \hat{\mathbf{b}} = \mathbf{b} + \mathbf{b}_{\text{FC}}, \quad (4)$$

where  $\boldsymbol{\mu}, \boldsymbol{\sigma} \in \mathbb{R}^{d_{\text{out}}}$ , with the elements  $\mu_i, \sigma_i$  denoting the mean and standard derivation of the vector  $\mathbf{W}_{i,:}$ , respectively;  $\gamma, \boldsymbol{\beta}, \mathbf{b}_{\text{FC}} \in \mathbb{R}^{d_{\text{out}}}$  are target-specific scale, shift, and bias style parameters trained with target data ( $\mathbf{W}$  and  $\mathbf{b}$  are frozen – from learning on the original source domain – during target training). One may interpret mFiLM as applying FiLM [64] (or batch normalization [31]) to a source FC *weight* to modulate its (row) statistics/style (encoded in  $\boldsymbol{\mu}$  and  $\boldsymbol{\sigma}$ ) to adapt to a target domain.

**Conv layers** Given a source Conv layer  $\mathbf{H}^{\text{source}} = \mathbf{W} * \mathbf{H}' + \mathbf{b}$  with input feature maps  $\mathbf{H}'$ , Conv filters  $\mathbf{W} \in \mathbb{R}^{C_{\text{out}} \times C_{\text{in}} \times K_1 \times K_2}$ , and bias  $\mathbf{b} \in \mathbb{R}^{C_{\text{out}}}$ , we leverage mAdaFM to modulate its parameters to form a target Conv layer as  $\mathbf{H}^{\text{target}} = \hat{\mathbf{W}} * \mathbf{H}' + \hat{\mathbf{b}}$ , where

$$\hat{\mathbf{W}} = \boldsymbol{\Gamma} \odot \frac{\mathbf{W} - \mathbf{M}}{\mathbf{S}} + \mathbf{B}, \quad \hat{\mathbf{b}} = \mathbf{b} + \mathbf{b}_{\text{Conv}}, \quad (5)$$

where  $\mathbf{M}, \mathbf{S} \in \mathbb{R}^{C_{\text{out}} \times C_{\text{in}}}$  with the elements  $M_{i,j}, S_{i,j}$  denoting the mean and standard derivation of the vector  $\text{vec}(\mathbf{W}_{i,j,:,:})$ , respectively. The trainable target-specific style parameters are  $\boldsymbol{\Gamma}, \mathbf{B} \in$

<sup>4</sup>This setup is not limited, as it’s often convenient to use a physical memory buffer to form the dataset stream. Note concerning practical applications, it’s often unnecessary to consider the extreme case where each dataset  $\mathcal{D}_t$  contains only one data sample; accordingly, we assume a moderate number of samples per dataset by default.

<sup>5</sup>See Appendix A for the detailed architectures. Note our method is deemed considerably robust to the (pretrained) source model, as discussed in Appendix H, where additional experiments are conducted based on a different source GAN model pretrained on LSUN Bedrooms [91].

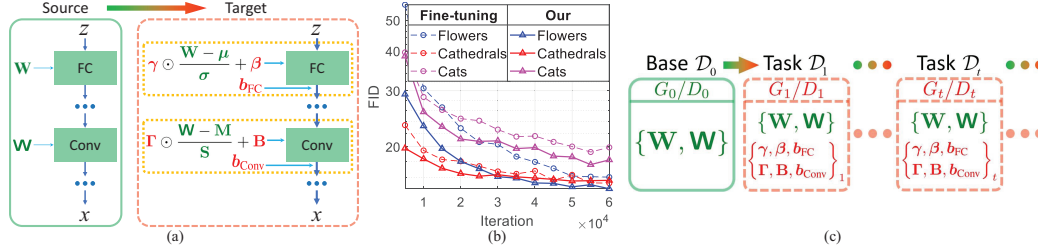


Figure 1: (a) Style modulation of a source GAN model (demonstrated with the generator, but it is also applied to the discriminator). Source parameters (green  $\{\mathbf{W}, \mathbf{W}\}$ ) are frozen, with limited trainable style parameters (*i.e.*, red  $\{\gamma, \beta, b_{FC}, \Gamma, \mathbf{B}, b_{Conv}\}$ ) introduced to form the augmentation to the target domain. (b) Comparing our style modulation to the strong fine-tuning baseline (see Appendix B for details). (c) The architecture of our GAN memory for a stream of target generation tasks.

$\mathbb{R}^{C_{out} \times C_{in}}$  and  $\mathbf{b}_{Conv} \in \mathbb{R}^{C_{out}}$ , with  $\mathbf{W}$  and  $\mathbf{b}$  frozen. Similar to mFiLM, mAdaFM first removes the source style (encoded in  $\mathbf{M}$  and  $\mathbf{S}$ ), followed by leveraging  $\Gamma$  and  $\mathbf{B}$  to learn the target style.

The adopted style modulation process (shown with the generator) is illustrated in Figure 1(a). Given a source GAN model, we transfer and freeze all its parameters to a target domain, followed by using mFiLM/mAdaFM in (4)/(5) to modulate its FC/Conv layers (with style parameters  $\{\gamma, \beta, b_{FC}, \Gamma, \mathbf{B}, b_{Conv}\}$ ) to yield the target generation model. With that style modulation, we transfer the source knowledge (within the frozen parameters) to the (potentially) perceptually-distant target domain to facilitate a realistic generation therein (see Figure 5 for generated samples). For quantitative evaluations, comparisons are made with a strong baseline, *i.e.*, fine-tuning the whole source model (including both generator and discriminator) on target data, which is expected to outperform training from scratch in both efficiency and performance by referring to [85] and the transfer learning literature. The FID scores (lower is better) [29] from both methods are summarized in Figure 1(b), where our method consistently outperforms that strong baseline by a large margin, on both training efficiency and performance,<sup>6</sup> highlighting the valuable knowledge within frozen source parameters (apparently a type of universal information about images) and showing a potentially better way for transfer learning.

Complementing the common knowledge of transfer learning, *i.e.*, low-level filters (those close to observations) are generally applicable, while high-level ones are task-specific [90, 93, 44, 4, 5, 20, 58, 98], the above discovery reveals an orthogonal dimensionality for transfer learning. Specifically, the shape of kernels (*i.e.*, relative relationship among kernel elements  $\{\mathbf{W}_{i,j,1,1}, \mathbf{W}_{i,j,1,2}, \dots\}$ ) may be generally transferable whereas the *statistics* of kernels (the mean  $\mathbf{M}_{i,j}$  or standard derivation  $\mathbf{S}_{i,j}$ ) or among-kernel correlations (*e.g.*, relative relationship among kernel statistics) are task-specific. A similar conjecture on low-level Conv filters was discussed in [98]; we reveal such patterns even hold for the whole GAN model (for both low-level and high-level kernels of the generator/discriminator), which is unanticipated because common experience associates high-level kernels with task-specific information. This insight might reveal a new avenue for transfer learning.

## 4.2 GAN memory to sequentially remember a stream of generative processes

Based on the above observations, we propose a GAN memory that has the power to realistically remember a stream of generative processes with *no* forgetting. The key observation here is that when modulating a source GAN model for a target domain, the source model is frozen (thus no forgetting of the source) with a limited set of target-specific style parameters (*i.e.*, no influence among tasks) introduced to form a realistic target generation. Accordingly, we can use a well-behaved source GAN model as a base, followed by sequentially modulating its “style” to deliver realistic generation power on a stream of target datasets,<sup>7</sup> as illustrated in Figure 1(c). We use the same settings as in Section

<sup>6</sup> The newly-introduced style parameters  $\{\gamma, \beta, b_{FC}, \Gamma, \mathbf{B}, b_{Conv}\}$  of our method are only about 20% of those of the fine-tuning baseline, yet they deliver better efficiency and performance. This is likely because the target data are not sufficient enough to train well-behaved parameters like the source ones, when performing fine-tuning alone. Note that with the techniques from Section 4.3, one can use much less style parameters (*e.g.*, 7.3% of those of the fine-tuning baseline) to yield a comparable performance.

<sup>7</sup> Our GAN memory is amenable to streaming training, parallel training, and even their flexible combinations, thanks to the frozen base model and task-specific style parameters.

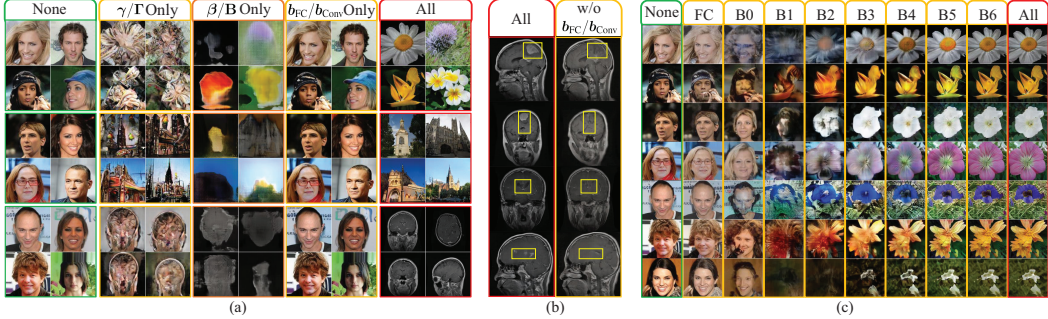


Figure 2: (a) Style parameters modulate different generation perspectives. (b) The biases model sparse objects. (c) Modulations in different blocks have different strength/focus over the generation.  $B_m$  is the  $m$ th residual block. B0/B6 is closest to the noise/observation. See Appendix C for details.

4.1. As style parameters are often limited (and can be further compressed as in Section 4.3), one could expect from our GAN memory a substantial compression of a stream of datasets, while not forgetting realistic generative replay (see Figure 5). To help better understand how/why our GAN memory works, we next reveal five of its properties.

*Each group of style parameters modulates a different generation perspective.* Style parameters consist of three groups, *i.e.*, scales  $\{\gamma, \Gamma\}$ , shifts  $\{\beta, B\}$ , and biases  $\{b_{FC}, b_{Conv}\}$ . Taking as examples the style parameters trained on the Flowers, Cathedrals, and Brain-tumor images, Figure 2(a) demonstrates the generation perspective modulated by each group: (i) when none/all groups are applied, GAN memory generates realistic source/target images (see the first/last column, respectively); (ii) when only modulating via the scales (denoted as  $\gamma/\Gamma$  Only, the second column), the generated samples show textural/structural information from target domains (like the textures on petals or the contours of buildings/skulls); (iii) as shown in the third column, the shifts if solely applied principally control the low-frequency color information from target domains; (iv) finally the biases (see the forth column) control the illumination and localized objects (not obvious here). To clearly reveal the role played by the biases, we keep both scales and shifts fixed, and compare the generated samples with/without biases on the Brain-tumor dataset; Figure 2(b) shows the biases are important in modeling localized objects like the tumor or tissue details, which may be valuable in pathological analysis. Also the following Figure 3(a) shows that the biases help with a better training efficiency.

*Style parameters within different blocks show different strength/focus over the generation.* Figure 2(c) shows the generated samples on the Flowers dataset when gradually and accumulatively adding modulations to each block (from FC to B6). To begin with, the FC modulation changes the overall contrast and illumination of the generation; then the style parameters in B0-B3 make the most effort to modulate the face manifold into a flower, followed by modulations in B4-B6 refining the generation details. Such patterns are somewhat consistent with existing practice, in the sense that high-level/low-level *kernel statistics* are more task-specific/generally-applicable. It’s therefore interesting to consider combining the two orthogonal dimensions, *i.e.*, the existing low-layer/high-layer split and the revealed kernel-shape/kernel-statistics split, for potentially better transfer learning.

*Normalization contributes significantly to better training efficiency and performance.* To investigate how the weight normalization and the biases in (4)/(5) contribute, ablation studies are conducted, with the results shown in Figure 3(a). With the weight normalization to remove the source “style,” our method shows both an improved training efficiency and a boosted performance; the biases contribute to a better efficiency but with minimal influence on the performance.

*GAN memory enables smooth interpolations among generative processes,* by dynamically combining two (or more) sets of style parameters. Figure 3(b) demonstrates smooth interpolations between flower and cat generative processes. Such a property can be used to deliver versatile data augmentation among different domains, which may, for example, benefit a downstream robust classification [61, 7].

*GAN memory readily generalizes to label-conditioned generation problems,* where each task dataset  $\mathcal{D}_t$  contains both observations  $\{x\}$  and the paired labels  $\{y\}$  with  $y \in \{1, \dots, C_t\}$ . Mimicking the conditional GAN [50], we specify one FC bias per class to model the label information; accordingly, the style parameters for the  $t$ th task are  $\{\gamma, \beta, \{b_{FC}\}_{i=1}^{C_t}, \Gamma, B, b_{Conv}\}$ . Figure 3(c) shows the realistic



Figure 3: (a) Ablation study on Flowers. (b) Smooth interpolations between flower and cat generative processes. (c) Realistic replay from our conditional-GAN memory. See Appendix D for details and more demonstrations.

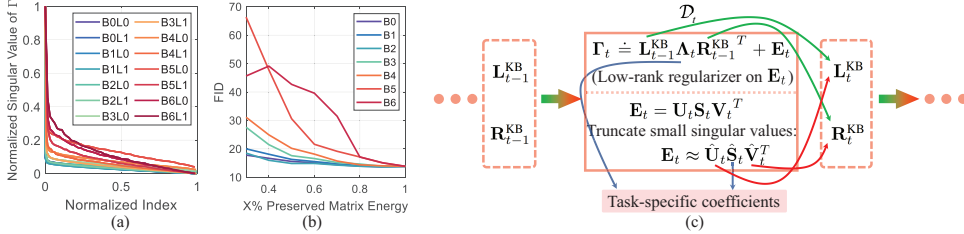


Figure 4: (a) Normalized singular values of  $\Gamma$  at all blocks/layers on Flowers (see Appendix Figure 15 for  $\mathbf{B}$ ). Maximum normalization is applied to both axes.  $BmLn$  stands for the  $n$ th Conv layer of the  $m$ th block. (b) Influence of truncation (via preserving  $X\%$  matrix energy [56]) on generation. (c) GAN memory with further compression and knowledge sharing. See Appendix G for details.

replay from our conditional-GAN memory after sequential training on five tasks (exemplated with bird, dog, and butterfly; see Section 5.2 for details).

### 4.3 GAN memory with further compression

Delivering realistic sequentially learned generations with no forgetting, the GAN memory presented above is expected to be sufficient for many practical applications with a moderate number of tasks. However, for challenging situations with many tasks, to save a set of style parameters for each task might become prohibitive.<sup>8</sup> For that problem, we reveal below (i) style parameters (*i.e.*, the expensive matrices  $\Gamma$  and  $\mathbf{B}$ ) can be compressed to lower the memory cost for each task; (ii) one can exploit sharing parameters among tasks to further enhance memory savings.<sup>9</sup>

We first investigate the singular values of  $\Gamma$  and  $\mathbf{B}$  learned at different blocks/layers, and summarize them in Figure 4(a). It's clear the  $\Gamma$  and  $\mathbf{B}$  parameters are in general low-rank; moreover, the closer a  $\Gamma$  or  $\mathbf{B}$  is to the noise (often with a larger matrix size and thus more expensive), the stronger its low-rank property (yielding better compressibility). Accordingly, we truncate/zero-out small singular values at each block/layer to test the corresponding performance decline. Figure 4(b) summarizes the results, where keeping 80% matrix energy [56] ( $\approx 35\%$  top singular values) of  $\Gamma$  and  $\mathbf{B}$  in B0-B4 almost has the same performance, verifying the compressibility of  $\Gamma$  and  $\mathbf{B}$ .

Based on the compressibility of  $\Gamma$  and  $\mathbf{B}$ , we next reveal parameter sharing among tasks can be exploited for further memory saving. Specifically, we propose to leverage matrix factorization and low-rank regularization to form a lifelong knowledge base mimicking [68]. Taking the  $t$ th task as an example, instead of optimizing over a task-specific  $\Gamma_t$  (similarly for  $\mathbf{B}_t$ ), we alternatively optimize over its parameterization  $\Gamma_t \doteq \mathbf{L}_{t-1}^{KB} \Lambda_t (\mathbf{R}_{t-1}^{KB})^T + \mathbf{E}_t$ , where  $\mathbf{L}_{t-1}^{KB}$  and  $\mathbf{R}_{t-1}^{KB}$  are respectively the existing left/right knowledge base,  $\Lambda_t = \text{Diag}(\lambda_t)$ , and  $\lambda_t$  and  $\mathbf{E}_t$  are task-specific trainable parameters. The nuclear norm  $\|\mathbf{E}_t\|_*$  is added to the loss to encourage a low-rank property. After training on the  $t$ th task, we apply singular value decomposition to  $\mathbf{E}_t$ , keep the top singular values to preserved  $X\%$  matrix energy, and use the corresponding left/right singular vectors to update the left/right knowledge base to  $\mathbf{L}_t^{KB}$  and  $\mathbf{R}_t^{KB}$ . The overall procedure is demonstrated in Figure 4(c).

<sup>8</sup> A compromise may save limited task-specific style parameters to hard disks and only load them when used.

<sup>9</sup> The cheap vector parameters  $\{\gamma, \beta, \mathbf{b}_{FC}, \mathbf{b}_{Conv}\}$  can be similarly processed in a dictionary learning manner. As they are often quite inexpensive to retain, we consider them being task-specific for simplicity.

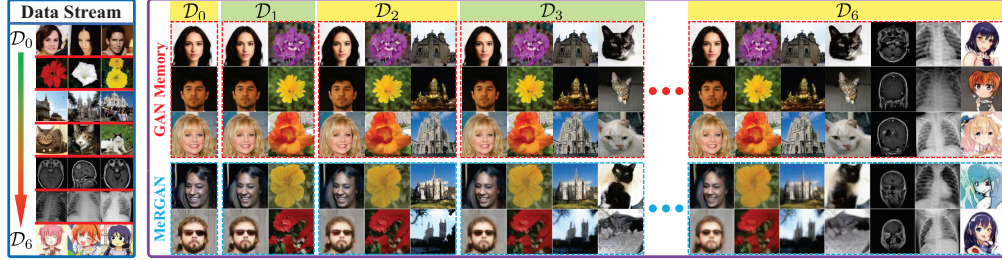


Figure 5: The task/dataset stream (left) and generated samples after training on each task/dataset (right).

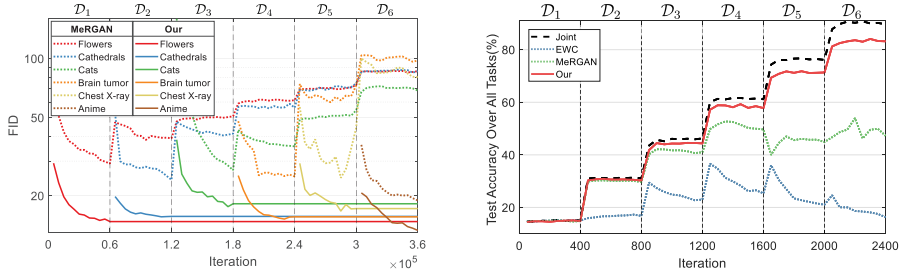


Figure 6: (Left) FID curves on the lifelong generation problem of Section 5.1. (Right) Classification accuracy on the lifelong classification problem of Section 5.2. The curve labeled “Joint” denotes the upper-bound, where the classifier is trained jointly on all the data from the current and historical tasks.

## 5 Experiments

Experiments on high-dimensional image datasets from diverse fields are conducted to demonstrate the effectiveness of the proposed techniques. Specifically, we first test our GAN memory to realistically remember a stream of generative processes; we then show that our (conditional) GAN memory can be used to form realistic pseudo rehearsal (synthesis of data from prior tasks) to alleviate catastrophic forgetting for challenging lifelong classification tasks; and finally for long task sequences, we reveal the techniques from Section 4.3 enable significant memory savings but with comparable performance. Detailed experimental settings are given in Appendix A.

### 5.1 GAN memory on a stream of generation tasks

To demonstrate the superiority of our GAN memory over existing replay-based methods, we design a challenging lifelong generation problem consisting of 6 perceptually-distant tasks/datasets (see Figure 5): Flowers [57], Cathedrals [99], Cats [97], Brain-tumor images [15], Chest X-rays [35], and Anime faces.<sup>10</sup> The GP-GAN [49] trained on the CelebA [43] ( $D_0$ ) is selected as the base; other well-behaved GAN models may readily be considered. We compare our GAN memory with the memory replay GAN (MeRGAN) [86], which keeps another copy of the generator in memory to replay historical generations, to mitigate catastrophic forgetting. Qualitative comparisons between both methods along the sequential training are shown in Figure 5. It’s clear our GAN memory delivers realistic generations with no forgetting on historical tasks, whereas MeRGAN shows increasingly blurry historical generations with reinforced generation artifacts [85] as training proceeds. For quantitative comparisons, the FID scores [29] along the training are summarized in Figure 6 (left), highlighting the advantages of our GAN memory, *i.e.*, realistic generations with no forgetting. Also revealed is that our method even performs better with a better efficiency for the current task, likely thanks to the transferred common knowledge within frozen base parameters.

### 5.2 Conditional-GAN memory as pseudo rehearsal for lifelong classifications

Witnessing the success of our (conditional) GAN memory in continually remembering realistic generations, we next utilize it as a pseudo rehearsal to assist downstream lifelong classifications. Specifically, we design 6 streaming tasks by selecting fish, bird, snake, dog, butterfly, and insect

<sup>10</sup><https://github.com/jayleicn/animeGAN>



images from the ImageNet [67]; each task is then formalized as a 6-classification problem (*e.g.*, for task bird, the goal is to classify 6 categories of birds). We employ the challenging class-incremental learning setup [78], *i.e.*, the classifier (after task  $t$ ) is expected to accurately classify all observed (first  $6t$ ) categories. For comparisons, we employ the regularization-based EWC [36] and the generative-replay-based MeRGAN [86]. For replay-based methods (*i.e.*, the MeRGAN and our conditional-GAN memory), at task  $t$ , we train the classifier with a combined dataset that contains both the current observed data (from task  $t$ ) and the generated/replayed historical samples (mimicking the data from task  $1 \sim t - 1$ ); after that, the MeRGAN/conditional-GAN-memory is updated to remember the current data generative process [72, 86]. See Appendix F for more details.

Testing classification accuracy on all 36 categories from the compared methods along training are summarized in Figure 6 (right). It’s clear that EWC barely works in the class-incremental learning scenario [77, 40, 78]. MeRGAN doesn’t work well when the task sequence is long, because of its increasingly blurry rehearsal as shown in Figure 5. By comparison, our conditional-GAN memory with no forgetting succeeds in stably maintaining an increasing accuracy as new tasks come and shows performance close to the joint-training upper-bound, highlighting its practical value in alleviating catastrophic forgetting for general lifelong learning problems. See Appendix F for evolution of the performance on each task along the training process.

### 5.3 GAN memory with parameter compression and sharing

Table 1: Comparisons of GAN memory with (Compr) or without (Naive) compression techniques. #Params denotes the number of newly-introduced style parameters for each task. The number of the frozen source parameters is 52.2M.

Task	$\mathcal{D}_1$	$\mathcal{D}_2$	$\mathcal{D}_3$	$\mathcal{D}_4$	$\mathcal{D}_5$	$\mathcal{D}_6$
#ParamsNaive	10.6M	10.6M	10.6M	10.6M	10.6M	10.6M
#ParamsCompr	3.8M	1.7M	0.9M	0.8M	0.3M	0.3M
#ParamsCompr/#ParamsNaive	36.4%	16.0%	9.0%	7.6%	2.7%	2.9%
FID (Compr)	<b>27.67</b>	23.49	<b>28.90</b>	<b>31.07</b>	32.19	49.28
FID (Naive)	28.89	<b>22.80</b>	34.36	35.72	<b>29.50</b>	<b>40.03</b>

To verify the effectiveness of the compression techniques presented in Section 4.3, which are believed valuable for lifelong learning with many tasks, we design another lifelong generation problem based on the ImageNet for better demonstration.<sup>11</sup> Specifically, we select 6 categories of butterfly images to form a stream of 6 generation tasks/datasets (one category per task), among which similarity/knowledge-sharability is expected. The procedure shown in Figure 4(c) is employed for our method. See Appendix G for details. We compare our GAN memory with compression techniques (denoted as Compr) to its naive implementation with task-specific style parameters (Naive), with the results summarized in Table 1. It’s clear that (*i*) even for task  $\mathcal{D}_1$  (with an empty knowledge base), the low-rank property of  $\Gamma/\mathbf{B}$  enables a significant parameter compression; (*ii*) based on the existing knowledge base, much less new parameters are necessary to form a new generation model (*e.g.*, for task  $\mathcal{D}_2$  or  $\mathcal{D}_3$ ), confirming the reusability/sharability of existing knowledge; and (*iii*) though it has significant parameter compression, Compr delivers comparable performance to Naive, confirming the effectiveness of the presented compression techniques.

## 6 Conclusions

We reveal that one can modulate the “style” of a GAN model to accurately synthesize the statistics of data from perceptually-distant targets. Based on that recognition, we propose our GAN memory with growing generation power, but with no forgetting. We then analyze our GAN memory in detail, reveal for it new compression techniques, and empirically verify its advantages over existing methods. Concerning a better base model, one may leverage the generative replay ability of the GAN memory to form a long-period update, mimicking human behavior during rapid-eye-movement sleep [75, 22].

<sup>11</sup> The datasets from Section 5.1 are too perceptually-distant to illustrate parameter sharability among tasks.

## Broader impact

Capable of remembering a stream of data generative processes with no forgetting, our GAN memory has the following potential positive impact in the society: (i) it may serve as a powerful generative replay for challenging lifelong applications such as self-driving; (ii) as no original data are saved, the concerns on data privacy may be well addressed; (iii) GAN memory enables flexible control over the replayed contents, which is of great value to practical applications, where training data are unbalanced, or where one needs to flexibly select which model capability to maintain/forget during training; (iv) the counter-intuitive discovery that lays the foundation of our GAN memory may disclose another dimension for transfer learning, *i.e.*, the kernel shape is generally applicable while the corresponding kernel statistics/style is task-specific; similar patterns may also apply to classifiers. Since our GAN memory is built on top of GANs, it may inherit their ethical and societal impact. Despite being versatility, GANs may be improperly used to synthesize fake images/news/videos, resulting in negative consequences. Furthermore, we should be cautious of the failure of adversarial training due to mode collapse, which may compromise the generative capability on the current task. Note that training failure, if it happens, will not hurt the performance on other tasks, showing certain robustness of our GAN memory.

## Acknowledgements

We thank the anonymous reviewers for their constructive comments. The research was supported in part by DARPA, DOE, NIH, NSF, ONR and SOC R&D lab of Samsung Semiconductor Inc. The Titan Xp GPU used was donated by the NVIDIA Corporation.

## References

- [1] K. Ak, J. Lim, J. Tham, and A. Kassim. Attribute manipulation generative adversarial networks for fashion images. In *ICCV*, pages 10541–10550, 2019.
- [2] Y. Arase and J. Tsujii. Transfer fine-tuning: A BERT case study. *arXiv preprint arXiv:1909.00931*, 2019.
- [3] X. Bao and Q. Qiao. Transfer learning from pre-trained bert for pronoun resolution. In *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*, pages 82–88, 2019.
- [4] D. Bau, B. Zhou, A. Khosla, A. Oliva, and A. Torralba. Network dissection: Quantifying interpretability of deep visual representations. In *CVPR*, pages 6541–6549, 2017.
- [5] D. Bau, J. Zhu, H. Strobelt, B. Zhou, J. Tenenbaum, W. Freeman, and A. Torralba. GAN dissection: Visualizing and understanding generative adversarial networks. *arXiv preprint arXiv:1811.10597*, 2018.
- [6] Y. Bengio. Deep learning of representations for unsupervised and transfer learning. In *Proceedings of ICML workshop on unsupervised and transfer learning*, pages 17–36, 2012.
- [7] D. Bertsimas, J. Dunn, C. Pawlowski, and Y. Zhuo. Robust classification. *INFORMS Journal on Optimization*, 1(1):2–34, 2019.
- [8] C. Bowles, L. Chen, R. Guerrero, P. Bentley, R. Gunn, A. Hammers, D. Dickie, M. Hernández, J. Wardlaw, and D. Rueckert. GAN augmentation: augmenting training data using generative adversarial networks. *arXiv preprint arXiv:1810.10863*, 2018.
- [9] A. Brock, J. Donahue, and K. Simonyan. Large scale GAN training for high fidelity natural image synthesis. In *ICLR*, 2019.
- [10] H. Caselles-Dupré, M. Garcia-Ortiz, and D. Filliat. Continual state representation learning for reinforcement learning using generative replay. *arXiv preprint arXiv:1810.03880*, 2018.
- [11] F. Castro, M. Marín-Jiménez, N. Guil, C. Schmid, and K. Alahari. End-to-end incremental learning. In *ECCV*, pages 233–248, 2018.
- [12] C. Chan, S. Ginosar, T. Zhou, and A. Efros. Everybody dance now. In *CVPR*, pages 5933–5942, 2019.
- [13] D. Chen, L. Yuan, J. Liao, N. Yu, and G. Hua. Stylebank: An explicit representation for neural image style transfer. In *CVPR*, pages 1897–1906, 2017.

- [14] Z. Chen and B. Liu. Lifelong machine learning. *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 12(3):1–207, 2018.
- [15] J. Cheng, W. Yang, M. Huang, W. Huang, J. Jiang, Y. Zhou, R. Yang, J. Zhao, Y. Feng, and Q. Feng. Retrieval of brain tumors by adaptive spatial pooling and fisher vector representation. *PloS one*, 11(6), 2016.
- [16] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell. Decaf: A deep convolutional activation feature for generic visual recognition. In *ICML*, pages 647–655, 2014.
- [17] V. Dumoulin, E. Perez, N. Schucher, F. Strub, H. Vries, A. Courville, and Y. Bengio. Feature-wise transformations. *Distill*, 3(7):e11, 2018.
- [18] V. Dumoulin, J. Shlens, and M. Kudlur. A learned representation for artistic style. *arXiv preprint arXiv:1610.07629*, 2016.
- [19] W. Fedus, I. Goodfellow, and A. Dai. MaskGAN: better text generation via filling in the  $\_$ . *arXiv preprint arXiv:1801.07736*, 2018.
- [20] Y. Frégier and J. Gouray. Mind2mind: transfer learning for GANs. *arXiv preprint arXiv:1906.11613*, 2019.
- [21] M. Frid-Adar, I. Diamant, E. Klang, M. Amitai, J. Goldberger, and H. Greenspan. GAN-based synthetic medical image augmentation for increased CNN performance in liver lesion classification. *Neurocomputing*, 321:321–331, 2018.
- [22] S. Gais, G. Albouy, M. Boly, T. Dang-Vu, A. Darsaud, M. Desseilles, G. Rauchs, M. Schabus, V. Sterpenich, G. Vandewalle, et al. Sleep transforms the cerebral trace of declarative memories. *PNAS*, 104(47):18778–18783, 2007.
- [23] G. Ghiasi, H. Lee, M. Kudlur, V. Dumoulin, and J. Shlens. Exploring the structure of a real-time, arbitrary neural artistic stylization network. *arXiv preprint arXiv:1705.06830*, 2017.
- [24] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR*, pages 580–587, 2014.
- [25] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *NIPS*, pages 2672–2680, 2014.
- [26] C. Han, K. Murao, T. Noguchi, Y. Kawata, F. Uchiyama, L. Rundo, H. Nakayama, and S. Satoh. Learning more with less: conditional PGGAN-based data augmentation for brain metastases detection using highly-rough annotation on MR images. *arXiv preprint arXiv:1902.09856*, 2019.
- [27] C. Han, L. Rundo, R. Araki, Y. Furukawa, G. Mauri, H. Nakayama, and H. Hayashi. Infinite brain MR images: PGGAN-based data augmentation for tumor detection. In *Neural Approaches to Dynamics of Signal Exchanges*, pages 291–303. Springer, 2020.
- [28] D. Hassabis, D. Kumaran, C. Summerfield, and M. Botvinick. Neuroscience-inspired artificial intelligence. *Neuron*, 95(2):245–258, 2017.
- [29] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter. GANs trained by a two time-scale update rule converge to a local Nash equilibrium. In *NIPS*, pages 6626–6637, 2017.
- [30] X. Huang and S. Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*, pages 1501–1510, 2017.
- [31] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- [32] T. Karras, T. Aila, S. Laine, and J. Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017.
- [33] T. Karras, S. Laine, and T. Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, June 2019.
- [34] T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen, and T. Aila. Analyzing and improving the image quality of stylegan. *arXiv preprint arXiv:1912.04958*, 2019.

- [35] D. S Kermany, M. Goldbaum, W. Cai, C. CS Valentim, H. Liang, S. L Baxter, A. McKeown, G. Yang, X. Wu, F. Yan, et al. Identifying medical diagnoses and treatable diseases by image-based deep learning. *Cell*, 172(5):1122–1131, 2018.
- [36] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *PNAS*, page 201611835, 2017.
- [37] K. Kumar, R. Kumar, T. de Boissiere, L. Gestin, W. Teoh, J. Sotelo, A. de Brebisson, Y. Bengio, and A. Courville. MelGAN: Generative adversarial networks for conditional waveform synthesis. *arXiv preprint arXiv:1910.06711*, 2019.
- [38] L. Kurzman, D. Vazquez, and I. Laradji. Class-based styling: Real-time localized style transfer with semantic segmentation. In *ICCV*, pages 0–0, 2019.
- [39] C. Ledig, L. Theis, F. Huszár, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. In *CVPR*, pages 4681–4690, 2017.
- [40] T. Lesort, H. Caselles-Dupré, M. Garcia-Ortiz, A. Stoian, and D. Filliat. Generative models from the perspective of continual learning. In *IJCNN*, pages 1–8. IEEE, 2019.
- [41] Y. Li, C. Fang, J. Yang, Z. Wang, X. Lu, and M. Yang. Universal style transfer via feature transforms. In *NIPS*, pages 386–396, 2017.
- [42] Z. Li and D. Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947, 2017.
- [43] Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. In *ICCV*, pages 3730–3738, 2015.
- [44] M. Long, Y. Cao, J. Wang, and M. Jordan. Learning transferable features with deep adaptation networks. *arXiv preprint arXiv:1502.02791*, 2015.
- [45] M. Lucic, K. Kurach, M. Michalski, S. Gelly, and O. Bousquet. Are GANs created equal? a large-scale study. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *NeurIPS*, pages 700–709. Curran Associates, Inc., 2018.
- [46] Z. Luo, Y. Zou, J. Hoffman, and L. Fei-Fei. Label efficient learning of transferable representations across domains and tasks. In *NIPS*, pages 165–177, 2017.
- [47] N. Masse, G. Grant, and D. Freedman. Alleviating catastrophic forgetting using context-dependent gating and synaptic stabilization. *Proceedings of the National Academy of Sciences*, 115(44):E10467–E10475, 2018.
- [48] M. McCloskey and N. Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pages 109–165. Elsevier, 1989.
- [49] L. Mescheder, A. Geiger, and S. Nowozin. Which training methods for GANs do actually converge? In *ICML*, pages 3478–3487, 2018.
- [50] M. Mirza and S. Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.
- [51] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida. Spectral normalization for generative adversarial networks. *arXiv preprint arXiv:1802.05957*, 2018.
- [52] M. Moradshahi, H. Palangi, M. Lam, P. Smolensky, and J. Gao. HUBERT untangles BERT to improve transfer across NLP tasks. *arXiv preprint arXiv:1910.12647*, 2019.
- [53] M. Mozafari, R. Farahbakhsh, and N. Crespi. A BERT-based transfer learning approach for hate speech detection in online social media. *arXiv preprint arXiv:1910.12574*, 2019.
- [54] M. Mundt, S. Majumder, I. Pliushch, and V. Ramesh. Unified probabilistic deep continual learning through generative replay and open set recognition. *arXiv preprint arXiv:1905.12019*, 2019.
- [55] C. Nguyen, Y. Li, T. Bui, and R. Turner. Variational continual learning. *arXiv preprint arXiv:1710.10628*, 2017.
- [56] V. Nikiforov. The energy of graphs and matrices. *JMAA*, 326(2):1472–1475, 2007.

- [57] M. Nilsback and A. Zisserman. Automated flower classification over a large number of classes. In *ICCV, Graphics & Image Processing*, pages 722–729. IEEE, 2008.
- [58] A. Noguchi and T. Harada. Image generation from small datasets via batch statistics adaptation. In *ICCV*, pages 2750–2758, 2019.
- [59] O. Ostapenko, M. Puszcz, T. Klein, P. Jahnichen, and M. Nabi. Learning to remember: A synaptic plasticity driven framework for continual learning. In *CVPR*, pages 11321–11329, 2019.
- [60] G. Parisi, R. Kemker, J. Part, C. Kanan, and S. Wermter. Continual lifelong learning with neural networks: A review. *Neural Networks*, 2019.
- [61] N. Park, M. Mohammadi, K. Gorde, S. Jajodia, H. Park, and Y. Kim. Data synthesis based on generative adversarial networks. *Proceedings of the VLDB Endowment*, 11(10):1071–1083, 2018.
- [62] T. Park, M. Liu, T. Wang, and J. Zhu. Semantic image synthesis with spatially-adaptive normalization. In *CVPR*, pages 2337–2346, 2019.
- [63] Y. Peng, S. Yan, and Z. Lu. Transfer learning in biomedical natural language processing: An evaluation of BERT and ELMo on ten benchmarking datasets. *arXiv preprint arXiv:1906.05474*, 2019.
- [64] E. Perez, F. Strub, H. De Vries, V. Dumoulin, and A. Courville. FiLM: Visual reasoning with a general conditioning layer. In *AAAI*, 2018.
- [65] S. Rebuffi, A. Kolesnikov, G. Sperl, and C. Lampert. iCaRL: Incremental classifier and representation learning. In *CVPR*, pages 2001–2010, 2017.
- [66] A. Rios and L. Itti. Closed-loop memory GAN for continual learning. *arXiv preprint arXiv:1811.01146*, 2018.
- [67] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *IJCV*, 115(3):211–252, 2015.
- [68] J. Schwarz, J. Luketina, W. Czarnecki, A. Grabska-Barwinska, Y. Teh, R. Pascanu, and R. Hadsell. Progress & compress: A scalable framework for continual learning. *arXiv preprint arXiv:1805.06370*, 2018.
- [69] A. Seff, A. Beatson, D. Suo, and H. Liu. Continual learning in generative adversarial nets. *NeurIPS*, 2017.
- [70] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. LeCun. Overfeat: Integrated recognition, localization and detection using convolutional networks. *arXiv preprint arXiv:1312.6229*, 2013.
- [71] F. Shen, S. Yan, and G. Zeng. Neural style transfer via meta networks. In *CVPR*, pages 8061–8069, 2018.
- [72] H. Shin, J. Lee, J. Kim, and J. Kim. Continual learning with deep generative replay. In *NIPS*, pages 2990–2999, 2017.
- [73] K. Shmelkov, C. Schmid, and K. Alahari. Incremental learning of object detectors without catastrophic forgetting. In *ICCV*, pages 3400–3409, 2017.
- [74] Q. Sun, B. Schiele, and M. Fritz. A domain based approach to social relation recognition. In *CVPR*, pages 3481–3490, 2017.
- [75] P. Taupin and F. Gage. Adult neurogenesis and neural stem cells of the central nervous system in mammals. *Journal of neuroscience research*, 69(6):745–749, 2002.
- [76] S. Thrun and T. Mitchell. Lifelong robot learning. *Robotics and autonomous systems*, 15(1-2):25–46, 1995.
- [77] G. van de Ven and A. Tolias. Generative replay with feedback connections as a general strategy for continual learning. *arXiv preprint arXiv:1809.10635*, 2018.
- [78] G. van de Ven and A. Tolias. Three scenarios for continual learning. *arXiv preprint arXiv:1904.07734*, 2019.
- [79] H. Wang, X. Liang, H. Zhang, D. Yeung, and E. Xing. Zm-net: Real-time zero-shot image manipulation network. *arXiv preprint arXiv:1703.07255*, 2017.
- [80] J. Wang and L. Perez. The effectiveness of data augmentation in image classification using deep learning. *Convolutional Neural Networks Vis. Recognit*, 2017.

- [81] K. Wang and X. Wan. SentiGAN: Generating sentimental texts via mixture adversarial networks. In *IJCAI*, pages 4446–4452, 2018.
- [82] R. Wang, D. Tang, N. Duan, Z. Wei, X. Huang, C. Cao, D. Jiang, and M. Zhou. K-adapter: Infusing knowledge into pre-trained models with adapters. *arXiv preprint arXiv:2002.01808*, 2020.
- [83] R. Wang, D. Zhou, and Y. He. Open event extraction from online text using a generative adversarial network. *arXiv preprint arXiv:1908.09246*, 2019.
- [84] T. Wang, M. Liu, J. Zhu, G. Liu, A. Tao, J. Kautz, and B. Catanzaro. Video-to-video synthesis. *arXiv preprint arXiv:1808.06601*, 2018.
- [85] Y. Wang, C. Wu, L. Herranz, J. van de Weijer, A. Gonzalez-Garcia, and B. Raducanu. Transferring GANs: generating images from limited data. In *ECCV*, pages 218–234, 2018.
- [86] C. Wu, L. Herranz, X. Liu, J. van de Weijer, B. Raducanu, et al. Memory replay GANs: Learning to generate new categories without forgetting. In *NeurIPS*, pages 5962–5972, 2018.
- [87] Y. Wu, Y. Chen, L. Wang, Y. Ye, Z. Liu, Y. Guo, and Y. Fu. Large scale incremental learning. In *CVPR*, pages 374–382, 2019.
- [88] Y. Xiang, Y. Fu, P. Ji, and H. Huang. Incremental learning using conditional adversarial networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 6619–6628, 2019.
- [89] R. Yamamoto, E. Song, and J. Kim. Probability density distillation with generative adversarial networks for high-quality parallel waveform generation. *arXiv preprint arXiv:1904.04472*, 2019.
- [90] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson. How transferable are features in deep neural networks? In *NIPS*, pages 3320–3328, 2014.
- [91] F. Yu, Y. Zhang, S. Song, A. Seff, and J. Xiao. LSUN: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365*, 2015.
- [92] A. Zamir, A. Sax, W. Shen, L. Guibas, J. Malik, and S. Savarese. Taskonomy: Disentangling task transfer learning. In *CVPR*, pages 3712–3722, 2018.
- [93] M. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *ECCV*, pages 818–833. Springer, 2014.
- [94] G. Zeng, Y. Chen, B. Cui, and S. Yu. Continual learning of context-dependent processing in neural networks. *Nature Machine Intelligence*, 1(8):364–372, 2019.
- [95] F. Zenke, B. Poole, and S. Ganguli. Continual learning through synaptic intelligence. In *ICML*, pages 3987–3995. JMLR. org, 2017.
- [96] M. Zhai, L. Chen, F. Tung, J. He, M. Nawhal, and G. Mori. Lifelong GAN: Continual learning for conditional image generation. In *ICCV*, pages 2759–2768, 2019.
- [97] W. Zhang, J. Sun, and X. Tang. Cat head detection-how to effectively exploit shape and texture features. In *ECCV*, pages 802–816. Springer, 2008.
- [98] M. Zhao, Y. Cong, and L. Carin. On leveraging pretrained GANs for generation with limited data. In *ICML*, 2020.
- [99] B. Zhou, A. Lapedriza, J. Xiao, A. Torralba, and A. Oliva. Learning deep features for scene recognition using places database. In *NIPS*, pages 487–495, 2014.