

---

# Generative Adversarial Networks and Continual Learning\*

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1        There is a strong emphasis in the continual learning literature on sequential classifi-  
2        cation experiments, where each task bears little semblance to previous ones. While  
3        certainly a form of continual learning, such tasks do not accurately represent many  
4        continual learning problems of the real-world, where the data distribution often  
5        evolves slowly over time. We propose using Generative Adversarial Networks  
6        (GANs) as a potential source for generating potentially unlimited datasets of this  
7        nature. We also identify that the dynamics of GAN training naturally constitute a  
8        continual learning problem, and show that leveraging continual learning methods  
9        can improve performance. As such, we show that techniques from both continual  
10       learning and GAN, typically studied separately, can be used to each other's benefit.

## 11    1 Introduction

12    The ability to learn new things continually while retaining previously acquired knowledge is a  
13    desirable attribute of an intelligent system. Humans and other forms of life do this well, but  
14    neural networks are known to exhibit a phenomenon known as *catastrophic forgetting* [12, 19]: the  
15    gradients that adapt a neural network's parameters to perform a new task tend to also clobber the  
16    model's ability to perform old ones. Because of its broad importance to the general field of machine  
17    learning, recent years have seen increased interest in approaches that enable *continual learning* (e.g.  
18    [9, 27, 11, 16, 20, 25]). These methods focus on improving the model architecture, objective, or  
19    training procedure to preserve knowledge of prior tasks while still enabling learning of new ones.

20    However, many of these works tend to conduct experiments that focus on learning a sequence of  
21    disparate tasks, which while certainly a continual learning task, does not capture the dynamics of  
22    a setting in which the data slowly evolves over time, as opposed to making abrupt discontinuous  
23    jumps. Such situations are common in many real-world applications, as deployed systems must  
24    maintain performance in an ever-evolving environment. It is therefore desirable for experiments in  
25    the literature to reflect this setting, but datasets that evolve over time are not readily available, which  
26    makes applying continual learning methods to such circumstances difficult.

27    On the other hand, recent years have seen an enormous amount of progress made in generative  
28    models, specifically with the advent of Generative Adversarial Networks (GANs) [3]. GANs have  
29    demonstrated the ability to learn impressively complex distributions [8, 1] from data samples alone.  
30    Interestingly, since GANs are capable of learning conditional distributions [14], and because the  
31    distribution of the generator's outputs smoothly evolves as training progresses, GANs represent an  
32    opportunity for producing a labeled dataset that varies through time.

33    Importantly though, the implications of the generator's distribution varying through time go beyond  
34    the potential for new sequential task benchmarks for continual learning. GANs are known to be  
35    somewhat challenging to train, with mode collapse a common problem. Inspection of a collapsed

---

\*Part of submission to the International Conference on Learning Representations (ICLR) 2019

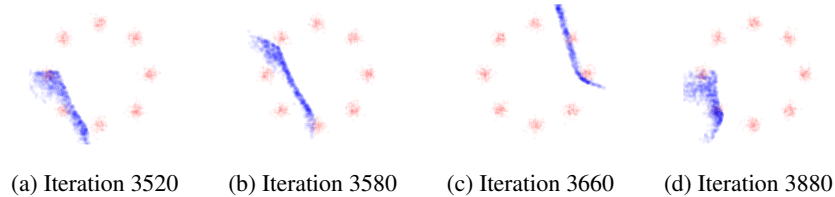


Figure 1: Real samples from a mixture of eight Gaussians in red; generated samples in blue. (a) The generator is currently mode collapsed to the bottom left. (b) The discriminator learns to recognize the generator oversampling this region and pushes the generator away. (c) The generator quickly settles at some new mode. (d) Upon being unseated from this new mode by the discriminator, the generator returns to the previous space. Such oscillations are common while training a vanilla GAN.

36 generator over subsequent training iterations reveal that rather than converging to a stationary distri-  
 37 bution, mode-collapsed generators tend to oscillate wildly, oftentimes revisiting previous locations  
 38 of the data space—modes that the discriminator presumably had previously learned to recognize  
 39 as fake (see Figure 1). We conjecture this phenomenon is at least in part enabled by catastrophic  
 40 forgetting in the discriminator: during training, synthesized fakes are presented to the discriminator  
 41 in a sequential manner reminiscent of the way tasks are learned in continual learning literature. Since  
 42 the discriminator is typically not refreshed with earlier synthesized samples, it loses its ability to  
 43 recognize them, allowing the generator to oscillate back to previous locations.

44 With these perspectives in mind, we make the following observations and contributions:

- 45 • Experiments in continual learning focus on sequences of disjoint tasks and do not cover the more  
 46 realistic scenario where a model encounters an evolving data distribution. GANs represent an  
 opportunity to fill this gap by synthesizing datasets that have the requisite time component.
- The training of a GAN discriminator is a continual learning problem. We show that augmenting  
 GAN models with continual learning methods improves performance on benchmark datasets.

## 47 2 Methods

### 48 2.1 GAN-generated datasets for continual learning

49 Consider distribution  $p_{\text{real}}(\mathbf{x})$ , from which we have data samples  $\mathcal{D}^{\text{real}}$ . We seek to learn a mapping  
 50 from an easy-to-sample distribution  $p(\mathbf{z})$  (e.g. standard normal) to a data distribution  $p_{\text{gen}}(\mathbf{x})$ ,  
 51 which we want to match  $p_{\text{real}}(\mathbf{x})$ . This mapping is parameterized as a neural network  $G_{\phi}(\mathbf{z})$  with  
 52 parameters  $\phi$ , termed the *generator*. The synthesized data are drawn  $\mathbf{x} = G_{\phi}(\mathbf{z})$ , with  $\mathbf{z} \sim p(\mathbf{z})$ . In  
 53 the GAN [3] set-up, we simultaneously learn another neural network  $D_{\theta}(\mathbf{x}) \in [0, 1]$  with parameters  
 54  $\theta$ , termed the *discriminator*, which provides feed-back to  $G_{\phi}(\mathbf{z})$ . Trained by a min-max objective in  
 55 conjunction with the generator, the generator gradually evolves: initial generations resemble random  
 56 noise, but eventually grow to resemble  $\mathcal{D}^{\text{real}}$ . At any point during training, an unlimited number of  
 57 samples can be drawn from  $G_{\phi}(\mathbf{z})$ . Therefore, at any training iteration  $t$ , we can generate a dataset  
 58  $\mathcal{D}_t^{\text{gen}}$ , and because  $p_{\text{gen}}(\mathbf{x})$  smoothly evolves with  $t$ , so does the sequence of datasets  $\mathcal{D}_1^{\text{gen}}, \dots, \mathcal{D}_T^{\text{gen}}$ .

59 As an example, we can train a DCGAN [18] on MNIST and generate an entire “fake” dataset of  
 60 70K samples every 50 training iterations of the DCGAN generator. We propose performing learning  
 61 on each of these generated datasets as individual tasks for continual learning. Selected samples are  
 62 shown in Figure 3 of Appendix A from the datasets  $\mathcal{D}_t^{\text{gen}}$  for  $t \in \{5, 10, 15, 20\}$ , each generated  
 63 from the same 100 samples of  $\mathbf{z}$  for all  $t$ . By conditioning the GAN [14] on randomly generated  
 64 labels, we have a mechanism for generating *labeled* datasets. With the success of large-scale GANs  
 65 [1], a similar method can be used to generate time-varying ImageNet datasets.

### 66 2.2 Continual learning for GAN discriminators

67 The traditional continual learning methods like Elastic Weight Consolidation (EWC) [9] or Intelligent  
 68 Synapses (IS) [27]<sup>1</sup> are designed for certain canonical benchmarks, commonly consisting of a  
 69 small number of clearly defined tasks (e.g., classification datasets in sequence). In GANs, the  
 70 discriminator is trained on dataset  $\mathcal{D}_t = \{\mathcal{D}^{\text{real}}, \mathcal{D}_t^{\text{gen}}\}$  at each iteration  $t$ . However, because of  
 71 the evolution of the generator, the distribution  $p_{\text{gen}}(\mathbf{x})$  from which  $\mathcal{D}_t^{\text{gen}}$  comes changes over time.  
 72 As such, we argue that different instances in time of the generator should be viewed as separate

<sup>1</sup>Summary of both of these methods can be found in Appendix B

73 tasks. Specifically, in the parlance of continual learning, the training data are to be regarded as  
 74  $\mathcal{D} = \{(\mathcal{D}^{\text{real}}, \mathcal{D}_1^{\text{gen}}), (\mathcal{D}^{\text{real}}, \mathcal{D}_2^{\text{gen}}), \dots\}$ . Thus motivated, we would like to apply continual learning  
 75 methods to the discriminator, but doing so is not straightforward for the following reasons:

- 76 • **Definition of a task:** EWC and IS were originally proposed for discrete, well-defined tasks. For  
 GAN, there is no such precise definition as to what a “task” is, and as discriminators are not  
 typically trained to convergence at every iteration, it is also unclear how long a task should be.
- 77 • **Computational memory:** While Equations 3 and 5 are for two tasks, they can be extended to  
 $K$  tasks by adding an additional loss term for each of the  $K - 1$  prior tasks. As each loss term  
 requires saving both a historical reference term  $\theta_k^*$  and either a diagonal Fisher Information  
 matrix  $F_k$  or importance weights  $\omega_k$  (all of which are the same size as the model parameters  
 $\theta$ ) for each task  $k$ , employing these techniques naively quickly becomes impractical for bigger  
 models when  $K$  gets large, especially if  $K$  is set to the number of training iterations  $T$ .
- 78 • **Continual *not* learning:** Early iterations of the discriminator are likely to be non-optimal, and  
 without a forgetting mechanism, EWC and IS may forever lock the discriminator to a poor  
 initialization. Additionally, the unconstrained addition of a large number of loss terms will cause  
 the continual learning regularization term to grow unbounded, which can disincentivize any  
 further changes in  $\theta$ .

79 To address these issues, we build upon EWC and IS by proposing several changes:

80 **Number of tasks as a rate:** We choose the total number of tasks  $K$  as a function of a constant rate  
 81  $\alpha$ , which denotes the number of iterations before the conclusion of a task, as opposed to arbitrarily  
 82 dividing the GAN training iterations into some set number of segments. Given  $T$  training iterations,  
 83 this means a rate  $\alpha$  yields  $K = \frac{T}{\alpha}$  tasks.

84 **Online Memory:** Seeking a way to avoid storing extra  $\theta_k^*$ ,  $F_k$ , or  $\omega_k$ , we observe that the sum of two  
 85 or more quadratic forms is another quadratic, which gives the classifier loss with continual learning  
 86 the following form for the  $(k + 1)^{\text{th}}$  task:

$$\mathcal{L}(\theta) = \mathcal{L}_{k+1}(\theta) + \mathcal{L}^{\text{CL}}(\theta), \quad \text{with} \quad \mathcal{L}^{\text{CL}}(\theta) \triangleq \frac{\lambda}{2} \sum_i S_{k,i} (\theta_i - \bar{\theta}_{k,i}^*)^2, \quad (1)$$

87 where  $\bar{\theta}_{k,i}^* = \frac{P_{k,i}}{S_{k,i}}$ ,  $S_{k,i} = \sum_{\kappa=1}^k Q_{\kappa,i}$ ,  $P_{k,i} = \sum_{\kappa=1}^k Q_{\kappa,i} \theta_{\kappa,i}^*$ , and  $Q_{\kappa,i}$  is either  $F_{\kappa,i}$  or  $\omega_{\kappa,i}$ ,  
 88 depending on the method. We name models with EWC and IS augmentations EWC-GAN and  
 89 IS-GAN, respectively.

90 **Controlled forgetting:** To provide a mechanism for forgetting earlier non-optimal versions of the  
 91 discriminator and to keep  $\mathcal{L}^{\text{CL}}$  bounded, we add a discount factor  $\gamma$ :  $S_{k,i} = \sum_{\kappa=1}^k \gamma^{k-\kappa} Q_{\kappa,i}$  and  
 92  $P_{k,i} = \sum_{\kappa=1}^k \gamma^{k-\kappa} Q_{\kappa,i} \theta_{\kappa,i}^*$ . Together,  $\alpha$  and  $\gamma$  determine how far into the past the discriminator  
 93 remembers previous generator distributions, and  $\lambda$  controls how important memory is relative to the  
 94 discriminator loss. Note, the terms  $S_k$  and  $P_k$  can be updated every  $\alpha$  steps in an online fashion:

$$S_{k,i} = \gamma S_{k-1,i} + Q_{k,i}, \quad P_{k,i} = \gamma P_{k-1,i} + Q_{k,i} \theta_{k,i}^* \quad (2)$$

95 This allows the EWC or IS loss to be applied without necessitating storing either  $Q_k$  or  $\theta_k^*$  for every  
 96 task  $k$ , which would quickly become too costly to be practical. Only a single variable to store a  
 97 running average is required for each of  $S_k$  and  $P_k$ , making this method space efficient.

98 Note that the training of the generator remains the same. Here we have shown two methods to  
 99 mitigate catastrophic forgetting for the original GAN; however, the proposed framework is applicable  
 100 to almost all of the wide range of GAN setups. Similarly, while we focus on EWC and IS here, any  
 101 continual learning method can be applied in a similar way.

### 102 3 Related work

103 There has been previous work investigating continual learning within the context of GANs. Improved  
 104 GAN [21] introduced historical averaging, which regularizes the model with a running average of  
 105 parameters of the most recent iterations. Simulated+Unsupervised training [23] proposed replacing  
 106 half of each minibatch with previous generator samples during training of the discriminator, as  
 107 previous generations should always be considered fake. However, this necessitates a historical buffer  
 108 of samples and halves the number of current samples that can be considered. Continual Learning  
 109 GAN [22] applies EWC to GAN, as we have, but uses it in the context of the class-conditioned

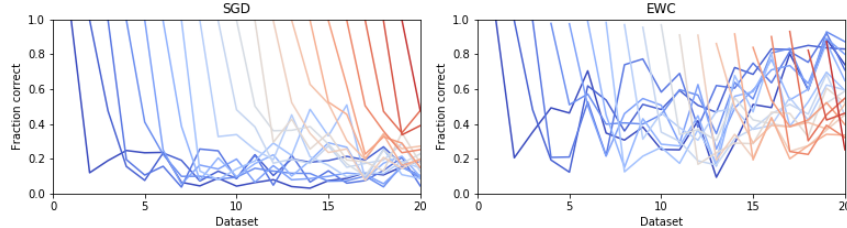


Figure 2: Each line represents the discriminator’s test accuracy on the fake GAN datasets. Note the sharp decrease in the discriminator’s ability to recognize previous fake samples upon fine-tuning on the next dataset using SGD (left). Forgetting still occurs with EWC (right), but is less severe.

Table 1: Image generation quality on CelebA and CIFAR-10

Method	CelebA	CIFAR-10	
	FID ↓	FID ↓	ICP ↑
DCGAN	12.52	41.44	6.97 ± 0.05
DCGAN + EWC	10.92	34.84	7.10 ± 0.05
WGAN-GP	-	30.23	7.09 ± 0.06
WGAN-GP + EWC	-	29.67	7.44 ± 0.08
SN-DCGAN	-	27.21	7.43 ± 0.10
SN-DCGAN + EWC	-	25.51	7.58 ± 0.07

110 generator that learns classes sequentially, as opposed to all at once, as we propose. [24] independently  
 111 makes a similar observation on the continual learning nature of GAN training, but propose momentum  
 112 and gradient penalty solutions instead and restrict themselves to experiments on toy examples.

## 113 4 Experiments

### 114 4.1 Sequential discrimination

115 While Figure 1 implies catastrophic forgetting in a GAN discriminator, we can show this concretely.  
 116 Using the DCGAN-generated MNIST datasets  $\mathcal{D}_1^{gen}, \dots, \mathcal{D}_T^{gen}$  described in Section 2.1, we now  
 117 train a discriminator to convergence on each  $\mathcal{D}_t^{gen}$  in sequence. Importantly, we do *not* include  
 118 samples from  $\mathcal{D}_{<t}^{gen}$  while fine-tuning on  $\mathcal{D}_t^{gen}$ . After fine-tuning on the train split of dataset  $\mathcal{D}_t^{gen}$ ,  
 119 the percentage of generated examples correctly identified as fake by the discriminator is evaluated  
 120 on the test splits of  $\mathcal{D}_{\leq t}^{gen}$ , with and without EWC (Figure 2). The catastrophic forgetting effect of  
 121 the discriminator trained with SGD is clear, with a steep drop-off in discriminating ability on  $\mathcal{D}_{t-1}^{gen}$   
 122 after fine-tuning on  $\mathcal{D}_t^{gen}$ ; this is unsurprising, as  $p_{gen}(x)$  has evolved specifically to deteriorate  
 123 discriminator performance. While there is still a dropoff with EWC, forgetting is less severe. On the  
 124 other hand, there is certainly room to improve, demonstrating the value of considering such kinds of  
 125 datasets for continual learning methods.

### 126 4.2 Augmenting GAN with continual learning

127 We augment the discriminators of various popular GANs implementations with EWC to preserve  
 128 recognition of previously seen generations, testing on two image datasets, CelebA and CIFAR-10.  
 129 Comparisons are made with the TTUR [6] variants of DCGAN [18] and WGAN-GP [4], as well as  
 130 an implementation of a spectral normalized [15] DCGAN (SN-DCGAN). Without modifying the  
 131 learning rate or model architecture, we show results with and without the EWC loss term added  
 132 to the discriminator for each. Performance is quantified with the Fréchet Inception Distance (FID)  
 133 [6] for both datasets. Since labels are available for CIFAR-10, we also report ICP for that dataset.  
 134 Best values are reported in Table 1. In each model, we see improvement in both FID and ICP from  
 135 the addition of EWC to the discriminator. Additional experiments improving GAN with continual  
 136 learning can be found in Appendix C.

## 137 5 Conclusion

138 We have identified the connections between GANs and continual learning: the training dynamics  
 139 of GAN naturally form a continual learning problem. This perspective allows us to show that (1)  
 140 GAN-generated datasets provide opportunities to form more realistic continual learning benchmarks;  
 141 (2) Existing continual learning methods can be adjusted to improve GAN training. Extensive  
 142 experimental results have demonstrated the proposed observations and solutions.

143 **References**

- 144 [1] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large Scale GAN Training for High Fidelity Natural  
145 Image Synthesis. *arXiv preprint*, 2018.
- 146 [2] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dolí, and C Lawrence  
147 Zitnick. Microsoft COCO Captions: Data Collection and Evaluation Server. *arXiv preprint*, 2015.
- 148 [3] Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron  
149 Courville, and Yoshua Bengio. Generative Adversarial Networks. *Advances In Neural Information  
150 Processing Systems*, 2014.
- 151 [4] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron Courville. Improved  
152 Training of Wasserstein GANs. *Advances In Neural Information Processing Systems*, 2017.
- 153 [5] Jiaxian Guo, Sidi Lu, Han Cai, Weinan Zhang, Yong Yu, and Jun Wang. Long Text Generation via  
154 Adversarial Training with Leaked Information. *AAAI Conference on Artificial Intelligence*, 2018.
- 155 [6] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs  
156 Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium. *Advances In Neural  
157 Information Processing Systems*, 2017.
- 158 [7] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *arXiv  
159 preprint arXiv:1611.01144*, 2016.
- 160 [8] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive Growing of GANs for Improved  
161 Quality, Stability, and Variation. *International Conference on Learning Representations*, 2018.
- 162 [9] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu,  
163 Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis, Claudia  
164 Clopath, Dharshan Kumaran, and Raia Hadsell. Overcoming Catastrophic Forgetting in Neural Networks.  
165 *Proceedings of the National Academy of Sciences*, 2017.
- 166 [10] Kevin Lin, Dianqi Li, Xiaodong He, Zhengyou Zhang, and Ming-Ting Sun. Adversarial ranking for  
167 language generation. *Advances in Neural Information Processing Systems*, 2017.
- 168 [11] David Lopez-Paz and Marc ' Aurelio Ranzato. Gradient Episodic Memory for Continual Learning.  
169 *Advances In Neural Information Processing Systems*, 2017.
- 170 [12] Michael McCloskey and Neal J Cohen. Catastrophic Interference in Connectionist Networks: The  
171 Sequential Learning Problem. *The Psychology of Learning and Motivation*, 1989.
- 172 [13] Luke Metz, Ben Poole, David Pfau, and Jascha Sohl-Dickstein. Unrolled Generative Adversarial Networks.  
173 *International Conference on Learning Representations*, 2017.
- 174 [14] Mehdi Mirza and Simon Osindero. Conditional Generative Adversarial Nets. *arXiv preprint*, 2014.
- 175 [15] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral Normalization for  
176 Generative Adversarial Networks. *International Conference on Learning Representations*, 2018.
- 177 [16] Cuong V Nguyen, Yingzhen Li, Thang D Bui, and Richard E Turner. Variational Continual Learning.  
178 *International Conference on Learning Representations*, 2017.
- 179 [17] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. BLEU: a Method for Automatic  
180 Evaluation of Machine Translation. *Annual Meeting of the Association for Computational Linguistics*,  
181 2002.
- 182 [18] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised Representation Learning with Deep  
183 Convolutional Generative Adversarial Networks. *International Conference on Learning Representations*,  
184 2016.
- 185 [19] Roger Ratcliff. Connectionist Models of Recognition Memory: Constraints Imposed by Learning and  
186 Forgetting Functions. *Psychology Review*, 1990.
- 187 [20] Hippolyt Ritter, Aleksandar Botev, and David Barber. Online Structured Laplace Approximations For  
188 Overcoming Catastrophic Forgetting. *Advances In Neural Information Processing Systems*, 2018.
- 189 [21] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved  
190 Techniques for Training GANs. *Advances In Neural Information Processing Systems*, 2016.

- 191 [22] Ari Seff, Alex Beatson, Daniel Suo, and Han Liu. Continual Learning in Generative Adversarial Nets.  
192 *arXiv preprint*, 2018.
- 193 [23] Ashish Shrivastava, Tomas Pfister, Oncel Tuzel, Josh Susskind, Wenda Wang, and Russ Webb. Learning  
194 from Simulated and Unsupervised Images through Adversarial Training. *Conference on Computer Vision  
195 and Pattern Recognition*, 2017.
- 196 [24] Hoang Thanh-Tung, Truyen Tran, and Svetha Venkatesh. On catastrophic forgetting and mode collapse in  
197 Generative Adversarial Networks. *arXiv preprint*, 2018.
- 198 [25] Ju Xu and Zhanxing Zhu. Reinforced Continual Learning. *Advances In Neural Information Processing  
199 Systems*, 2018.
- 200 [26] Lantao Yu, Weinan Zhang, Jun Wang, and Yong Yu. SeqGAN: Sequence Generative Adversarial Nets with  
201 Policy Gradient. *AAAI Conference on Artificial Intelligence*, 2017.
- 202 [27] Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual Learning Through Synaptic Intelligence.  
203 *International Conference on Machine Learning*, 2017.
- 204 [28] Yizhe Zhang, Zhe Gan, Kai Fan, Zhi Chen, Ricardo Henao, Dinghan Shen, and Lawrence Carin. Adversarial  
205 feature matching for text generation. *ICML*, 2017.

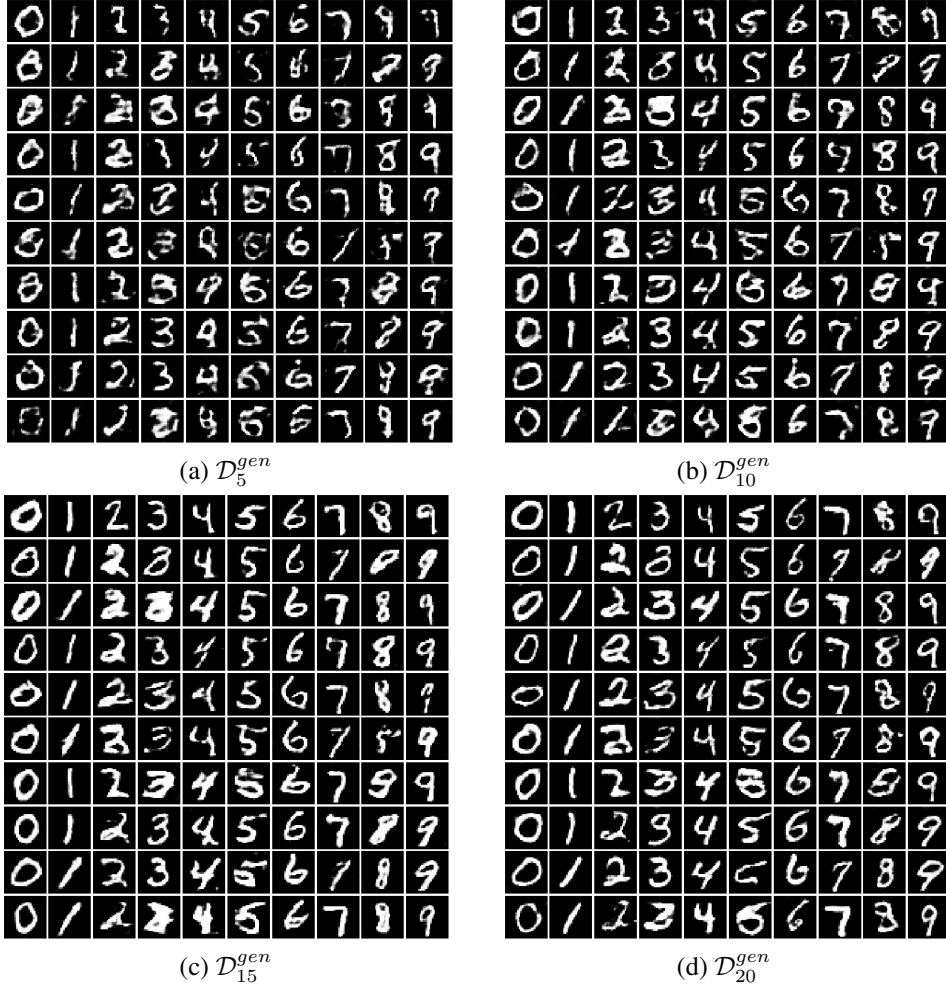


Figure 3: Image samples from generated “fake MNIST” datasets

## 207 B Continual learning methods summary

## 208 B.1 Elastic weight consolidation (EWC)

209 To derive the EWC loss, [9] frames training a model as finding the most probable values of the parameters  $\theta$   
 210 given the data  $\mathcal{D}$ . For two tasks, the data are assumed partitioned into independent sets according to the task,  
 211 and the posterior for Task 1 is approximated as a Gaussian with mean centered on the optimal parameters for  
 212 Task 1  $\theta_1^*$  and diagonal precision given by the diagonal of the Fisher information matrix  $F_1$  at  $\theta_1^*$ . This gives the  
 213 EWC loss the following form:

$$\mathcal{L}(\theta) = \mathcal{L}_2(\theta) + \mathcal{L}^{EWC}(\theta), \quad \text{with} \quad \mathcal{L}^{EWC}(\theta) \triangleq \frac{\lambda}{2} \sum_i F_{1,i} (\theta_i - \theta_{1,i}^*)^2, \quad (3)$$

214 where  $\mathcal{L}_2(\theta) = \log p(\mathcal{D}_2|\theta)$  is the loss for Task 2 individually,  $\lambda$  is a hyperparameter representing the  
 215 importance of Task 1 relative to Task 2,  $F_{1,i} = \left(\frac{\partial \mathcal{L}_1(\theta)}{\partial \theta_i}\right)_{\theta=\theta_1^*}^2$ ,  $i$  is the parameter index, and  $\mathcal{L}(\theta)$  is the  
 216 new loss to optimize while learning Task 2. Intuitively, the EWC loss prevents the model from straying too  
 217 far away from the parameters important for Task 1 while leaving less crucial parameters free to model Task 2.  
 218 Subsequent tasks result in additional  $\mathcal{L}^{EWC}(\theta)$  terms added to the loss for each previous task. By protecting  
 219 the parameters deemed important for prior tasks, EWC as a regularization term allows a single neural network  
 220 (assuming sufficient parameters and capacity) to learn new tasks in a sequential fashion, without forgetting how  
 221 to perform previous tasks.

222 **B.2 Intelligent synapses (IS)**

223 While EWC makes a point estimate of how essential each parameter is at the conclusion of a task, IS [27]  
 224 protects the parameters according to their importance along the task’s entire training trajectory. Termed synapses,  
 225 each parameter  $\theta_i$  of the neural network is awarded an *importance measure*  $\omega_{1,i}$  based on how much it reduced  
 226 the loss while learning Task 1. Given a loss gradient  $\mathbf{g}(t) = \nabla_{\theta} \mathcal{L}(\theta)|_{\theta=\theta_t}$  at time  $t$ , the total change in loss  
 227 during the training of Task 1 then is the sum of differential changes in loss over the training trajectory. With the  
 228 assumption that parameters  $\theta$  are independent, we have:

$$\int_{t^0}^{t^1} \mathbf{g}(t) d\theta = \int_{t^0}^{t^1} \mathbf{g}(t) \theta' dt = \sum_i \int_{t^0}^{t^1} g_i(t) \theta'_i dt \triangleq - \sum_i \omega_{1,i}, \quad (4)$$

229 where  $\theta' = \frac{d\theta}{dt}$  and  $(t^0, t^1)$  are the start and finish of Task 1, respectively. Note the added negative sign, as  
 230 importance is associated with parameters that decrease the loss.

231 The importance measure  $\omega_{1,i}$  can now be used to introduce a regularization term that protects parameters  
 232 important for Task 1 from large parameter updates, just as the Fisher information matrix diagonal terms  $F_{1,i}$   
 233 were used in EWC. This results in an IS loss very reminiscent in form:

$$\mathcal{L}(\theta) = \mathcal{L}_2(\theta) + \mathcal{L}^{\text{IS}}(\theta), \quad \text{with } \mathcal{L}^{\text{IS}}(\theta) \triangleq \frac{\lambda}{2} \sum_i \omega_{1,i} (\theta_i - \theta_{1,i}^*)^2. \quad (5)$$

234 Note that [27] instead consider  $\Omega_{1,i} = \frac{\omega_{1,i}}{(\Delta_{1,i})^2 + \xi}$ , where  $\Delta_{1,i} = \theta_{1,i} - \theta_{0,i}$  and  $\xi$  is a small number for  
 235 numerical stability. We however found that the inclusion of  $(\Delta_{1,i})^2$  can lead to the loss exploding and then  
 236 collapsing as the number of tasks increases and so omit it. We also change the hyperparameter  $c$  into  $\frac{\lambda}{2}$ .

237 **C Additional Experiments**

238 **C.1 Mixture of eight Gaussians**

239 We show results on a toy dataset consisting of a mixture of eight Gaussians, as in the example in Figure 1.  
 240 Following the setup of [13], the real data are evenly distributed among eight 2-dimensional Gaussian distributions  
 241 arranged in a circle of radius 2, each with covariance  $0.02I$  (see Figure 4). We evaluate our model with Inception  
 242 Score (ICP) [21], which gives a rough measure of diversity and quality of samples; higher scores imply better  
 243 performance, with the true data resulting in a score of around 7.870. For this simple dataset, since we know the  
 244 true data distribution, we also calculate the symmetric Kullback–Leibler divergence (Sym-KL); lower scores  
 245 mean the generated samples are closer to the true data. We show computation time, measured in numbers of  
 246 training iterations per second (Iter/s), averaged over the full training of a model on a single Nvidia Titan X  
 247 (Pascal) GPU. Each model was run 10 times, with the mean and standard deviation of each performance metric  
 248 at the end of 25K iterations reported in Table 2.

249 The performance of EWC-GAN and IS-GAN were evaluated for a number of hyperparameter settings. We  
 250 compare our results against a vanilla GAN [3], as well as a state-of-the-art GAN with spectral normalization

Table 2: Iterations per second, inception score, and symmetric KL divergence comparison on a mixture of eight Gaussians.

Model						
Method	$\alpha$	$\lambda$	$\gamma$	Iter/s $\uparrow$	ICP $\uparrow$	Sym-KL $\downarrow$
GAN	-	-	-	$87.59 \pm 1.45$	$2.835 \pm 2.325$	$19.55 \pm 3.07$
GAN + $\ell_2$ weight	1	0.01	0		$5.968 \pm 1.673$	$15.19 \pm 2.67$
GAN + historical avg.	1	0.01	0.995		$7.305 \pm 0.158$	$13.32 \pm 0.88$
GAN + SN	-	-	-	$49.70 \pm 0.13$	$6.762 \pm 2.024$	$13.37 \pm 3.86$
GAN + IS	1000	100	0.8	$42.26 \pm 0.35$	$7.039 \pm 0.294$	$15.10 \pm 1.51$
GAN + IS	100	10	0.98	$42.29 \pm 0.10$	$7.500 \pm 0.147$	$11.85 \pm 0.92$
GAN + IS	10	100	0.99	$41.07 \pm 0.07$	$7.583 \pm 0.242$	$11.88 \pm 0.84$
GAN + SN + IS	10	100	0.99	$25.69 \pm 0.09$	$7.699 \pm 0.048$	$11.10 \pm 1.18$
GAN + EWC	1000	100	0.8	$82.78 \pm 1.55$	$7.480 \pm 0.209$	$13.00 \pm 1.55$
GAN + EWC	100	10	0.98	$80.63 \pm 0.39$	$7.488 \pm 0.222$	$12.16 \pm 1.64$
GAN + EWC	10	10	0.99	$73.86 \pm 0.16$	$7.670 \pm 0.112$	$11.90 \pm 0.76$
GAN + SN + EWC	10	10	0.99	$44.68 \pm 0.11$	$7.708 \pm 0.057$	$11.48 \pm 1.12$



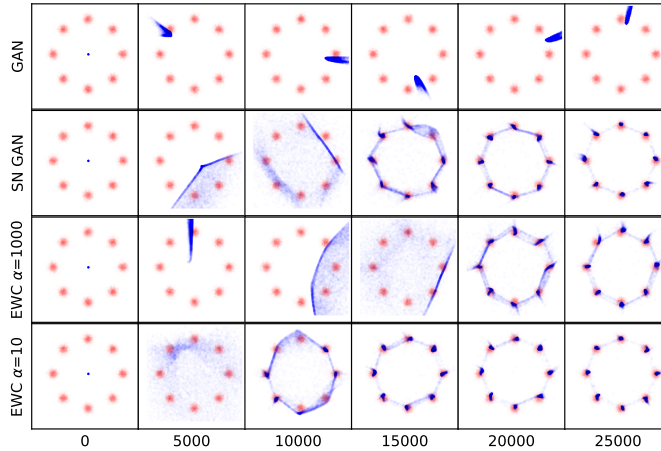


Figure 4: Each row shows the evolution of generator samples at 5000 training step intervals for GAN, SN-GAN, and EWC-GAN for two  $\alpha$  values. The proposed EWC-GAN models have hyperparameters matching the corresponding  $\alpha$  in Table 2. Each frame shows 10000 samples drawn from the true eight Gaussians mixture (red) and 10000 generator samples (blue).

251 (SN) [15] applied to the discriminator. As spectral normalization augments the discriminator loss in a way  
 252 different from continual learning, we can combine the two methods; this variant is also shown.

253 Note that a discounted version of discriminator historical averaging [21] can be recovered from the EWC and IS  
 254 losses if the task rate  $\alpha = 1$  and  $Q_{k,i} = 1$  for all  $i$  and  $k$ , a poor approximation to both the Fisher information  
 255 matrix diagonal and importance measure. If we also set the historical reference term  $\hat{\theta}_k^*$  and the discount factor  $\gamma$   
 256 to zero, then the EWC and IS losses become  $\ell_2$  weight regularization. These two special cases are also included  
 257 for comparison.

258 We observe that augmenting GAN models with EWC and IS consistently results in generators that better  
 259 match the true distribution, both qualitatively and quantitatively, for a wide range of hyperparameter settings.  
 260 EWC-GAN and IS-GAN result in a better ICP and FID than  $\ell_2$  weight regularization and discounted historical  
 261 averaging, showing the value of prioritizing protecting important parameters, rather than all parameters equally.  
 262 EWC-GAN and IS-GAN also outperform a state-of-the-art method in SN-GAN. In terms of training time,  
 263 updating the EWC loss requires forward propagating a new minibatch through the discriminator and updating  
 264  $S$  and  $P$ , but even if this is done at every step ( $\alpha = 1$ ), the resulting algorithm is only slightly slower than  
 265 SN-GAN. Moreover, doing so is unnecessary, as higher values of  $\alpha$  also provide strong performance for a much  
 266 smaller time penalty. Combining EWC with SN-GAN leads to even better results, showing that the two methods  
 267 can complement each other. IS-GAN can also be successfully combined with SN-GAN, but it is slower than  
 268 EWC-GAN as it requires tracking the trajectory of parameters at each step. Sample generation evolution over  
 269 time is shown in Figure 4.

## 270 C.2 Text generation of COCO Captions

271 We also consider the text generation on the MS COCO Captions dataset [2], with the pre-processing in [5].  
 272 Quality of generated sentences is evaluated by BLEU score [17]. Since BLEU- $b$  measures the overlap of  $b$   
 273 consecutive words between the generated sentences and ground-truth references, higher BLEU scores indicate  
 274 better fluency. Self BLEU uses the generated sentences themselves as references; lower values indicate higher  
 275 diversity.

276 We apply EWC and IS to textGAN [28], a recently proposed model for text generation in which the discriminator  
 277 uses feature matching to stabilize training. This model’s results (labeled “EWC” and “IS”) are compared to a  
 278 Maximum Likelihood Estimation (MLE) baseline, as well as several state-of-the-art methods: SeqGAN [26],  
 279 RankGAN [10], GSGAN [7] and LeakGAN [5]. Our variants of textGAN outperforms the vanilla textGAN for  
 280 all BLEU scores (see Table 3), indicating the effectiveness of addressing the forgetting issue for GAN training in  
 281 text generation. EWC/IS + textGAN also demonstrate a significant improvement compared with other methods,  
 282 especially on BLEU-2 and 3. Though our variants lag slightly behind LeakGAN on BLEU-4 and 5, their self  
 283 BLEU scores (Table 4) indicate it generates more diverse sentences. Sample sentence generations can be found  
 284 in Table 5.

Table 3: Test BLEU  $\uparrow$  results on MS COCO

Method	MLE	SeqGAN	RankGAN	GSGAN	LeakGAN	textGAN	EWC	IS
BLEU-2	0.820	0.820	0.852	0.810	0.922	0.926	0.934	0.933
BLEU-3	0.607	0.604	0.637	0.566	0.797	0.781	0.802	0.791
BLEU-4	0.389	0.361	0.389	0.335	0.602	0.567	0.594	0.578
BLEU-5	0.248	0.211	0.248	0.197	0.416	0.379	0.400	0.388

Table 4: Self BLEU  $\downarrow$  results on MS COCO

Method	MLE	SeqGAN	RankGAN	GSGAN	LeakGAN	textGAN	EWC	IS
BLEU-2	0.754	0.807	0.822	0.785	0.912	0.843	0.854	0.853
BLEU-3	0.511	0.577	0.592	0.522	0.825	0.631	0.671	0.655
BLEU-4	0.232	0.278	0.288	0.230	0.689	0.317	0.388	0.364

Table 5: Sample sentence generations from EWC + textGAN

---

a couple of people are standing by some zebras in the background  
 the view of some benches near a gas station  
 a brown motorcycle standing next to a red fence  
 a bath room with a broken tank on the floor  
 red passenger train parked under a bridge near a river  
 some snow on the beach that is surrounded by a truck  
 a cake that has been perform in the background for takeoff  
 a view of a city street surrounded by trees  
 two giraffes walking around a field during the day  
 crowd of people lined up on motorcycles  
 two yellow sheep with a baby dog in front of other sheep  
 an intersection sits in front of a crowd of people  
 a red double decker bus driving down the street corner  
 an automobile driver stands in the middle of a snowy park  
 five people at a kitchen setting with a woman  
 there are some planes at the takeoff station  
 a passenger airplane flying in the sky over a cloudy sky  
 three aircraft loaded into an airport with a stop light  
 there is an animal walking in the water  
 an older boy with wine glasses in an office  
 two old jets are in the middle of london  
 three motorcycles parked in the shade of a crowd  
 group of yellow school buses parked on an intersection  
 a person laying on a sidewalk next to a sidewalk talking on a cell phone  
 a chef is preparing food with a sink and stainless steel appliances

---