

# Lifelong Generative Learning via Knowledge Reconstruction

## Abstract

Generative models often incur the catastrophic forgetting problem when they are used to sequentially learning multiple tasks, i.e., lifelong generative learning. Although there are some endeavors to tackle this problem, they suffer from high time-consumptions or error accumulation. In this work, we develop an efficient and effective lifelong generative model based on variational autoencoder (VAE). Unlike the generative adversarial network, VAE enjoys high efficiency in the training process, providing natural benefits with few resources. We deduce a lifelong generative model by expending the intrinsic reconstruction character of VAE to the historical knowledge retention. Further, we devise a feedback strategy about the reconstructed data to alleviate the error accumulation. Experiments on the lifelong generating tasks of MNIST, FashionMNIST, and SVHN verified the efficacy of our approach, where the results were comparable to SOTA.

## 1 Introduction

Lifelong learning is an important yet challenging problem, which is notorious for the catastrophic forgetting phenomenon [Kemker *et al.*, 2018; Parisi *et al.*, 2019]. Such a phenomenon reflects the core challenge that the performance dramatically degrades on the tasks learned before, since the model focuses on learning the current task. As the generative replay methods provide an insightful solution in addressing most types of lifelong discriminative learning problems (i.e., task, domain, class lifelong learning) [Van de Ven and Tolias, 2018; Delange *et al.*, 2021], how to enable the generative models to lifelong learning has drawn much attention [Lesort, 2020; Ramapuram *et al.*, 2020]. This problem is also known as lifelong generative learning.

Some approaches (also known as the pseudo rehearsal) train a new generative model on the mixture data that combines the samples of the current task and the pseudo samples generated from the previous model. These approaches can be fallen into two categories. One is based on the unconditional generative models, while the other is on the conditional generative models. On the one hand, methods based on the

unconditional generative models took the generative adversarial network (GAN) [Shin *et al.*, 2017] or variational autoencoding (VAE) [Van de Ven and Tolias, 2018; van de Ven *et al.*, 2020] to unconditionally accomplish lifelong generative learning. But they are biased towards sampling from the recent task [Seff *et al.*, 2017; Wu *et al.*, 2018]. To relieve the biased-sampling problem, on the other hand, conditional generative models were introduced. In particular, van de Ven *et al.* [2020] modelled the hidden variable in the VAE as a Gaussian mixture distribution, indicating a mixture of the learned tasks. Then, they balanced the learning of historical tasks and the current one to relieve the biased-sampling problem. Ramapuram *et al.* [2020] extended the inference capability of VAE from modeling only one hidden variable into two variables, including a continuous feature variable and a discrete task ID variable. Further, Ye and Bors [2021] added one more discrete label ID variable in each task to achieve a better fine-grained category balance. For the GAN model, Seff *et al.* [2017] only preserved the number of historical labels throughout the lifelong learning process. They creatively embedded the label information of the historical task in the generator model to solve the imbalance problems (i.e., task imbalance and category imbalance). This strategy for GAN is prevalently employed [Zhai *et al.*, 2019; Liu *et al.*, 2020].

Other approaches train the new generative model on the data from the current task only. Referring to [Kirkpatrick *et al.*, 2017], Seff *et al.* [2017] enforced the generator to remember the historical knowledge by constraints to update the parameters that are crucial for historical tasks. Besides, Wu *et al.* [2018] and Liu *et al.* [2020] employed knowledge distillation [Hinton *et al.*, 2015] to transfer learned knowledge from the previous generator to the new generator.

However, all of the above-mentioned pseudo rehearsal methods train the generative model partly on the pseudo data, which own the uncertain quality from the sampling, and thus result in the *error accumulation* problem. In addition, methods trained only on the current data are merely investigated in the GAN framework, and the *time consumption* of generating suitable instances is extensive. Compared with GAN, VAE not only provides a stable training mechanism but keeps satisfactory sample diversity [Ramapuram *et al.*, 2020].

In this paper, motivated by the conditional GAN (CGAN), we first develop a variant of the conditional VAE (CVAE), which enjoys high efficiency in the training process. To en-

sure remembering the previous knowledge, inspired by the intrinsic reconstruction character in the VAE, we particularly propose a knowledge reconstruction loss to guide the decoder training. Further, we devise a feedback strategy about the reconstructed data to encourage CVAE to encode the reconstructed sample consistent with the real one. Finally, we present a lifelong generative learning algorithm via knowledge reconstruction (LGLvKR in short), which is only trained on the current data without error accumulation from the pseudo sampling. Compared with CGAN-based methods, our algorithm LGLvKG obtains comparable results with less computational resources.

## 2 Background and Problem Definition

In this section, we first introduce the concepts used throughout this paper. They include the generalized conditional generative model (GCGM), and a variant of CVAE. We then formalize the definition of our lifelong generative problem.

### 2.1 Genelized Conditional Generative Model

Consider an observed dataset with  $N$  conditional samples  $\mathcal{D} = \{(x_i, y_i) | i = 1, \dots, N\}$  where  $x_i$  and  $y_i$  indicate the image and the label for the  $i^{th}$  sample, respectively. GCGM aims to estimate the distribution of each label in  $y$  based on observations. With estimated distributions, real data with the specified label are well sampled, i.e., conditional sample generation.

Concretely, as the graphical model shown in Fig.1(a), GCGM expects to estimate a conditional probability distribution with the parameter,  $\theta$ , through the observable variable  $y$  and the unobservable variable  $z$ . In the network architecture, the generative model with the parameter  $\theta$  enables us to generate  $x$  conditioned by  $y$  with a prior  $z$ . To get a proper  $\theta$ , the maximum logarithmic likelihood estimation is typically introduced:

$$\max_{\theta} \log p_{\theta}(x|y) = \max_{\theta} \log \int p_{\theta}(x|y, z) p(z) dz. \quad (1)$$

To make Eq.(1) integrally tractable [Kingma and Welling, 2014], we suggest a novel conditional variational autoencoder shown in Section 2.2.

### 2.2 Conditional Variational Autoencoder

The existing CVAE architecture embeds the conditional information,  $y$ , either in the encoder [Sohn *et al.*, 2015] or in the decoder [Kingma *et al.*, 2014]. Unlike them, we use a variant of CVAE where  $y$  is embedded in both encoder and decoder, parametrized by  $\phi$  and  $\theta$ , respectively, as shown in Fig.1(b). In this way, we utilize the conditional information further. Using the Jensen's inequality [Ramapuram *et al.*, 2020], the evidence lower bound (ELBO) is,

$$\begin{aligned} \log p_{\theta}(x|y) &= \log \int \frac{q_{\phi}(z|y, x)}{q_{\phi}(z|y, x)} p_{\theta}(x|y, z) p(z) dz, \\ &\geq \mathbb{E}_{q_{\phi}(z|y, x)} [\log p_{\theta}(x|y, z)] - \text{KL} [q_{\phi}(z|y, x) || p(z)], \end{aligned} \quad (2)$$

where  $\mathbb{E}_{q_{\phi}(z|y, x)} [\log p_{\theta}(x|y, z)]$  means to take the expectation of  $\log p_{\theta}(x|y, z)$  with respect to the distribution  $q_{\phi}(z|y, x)$ ,

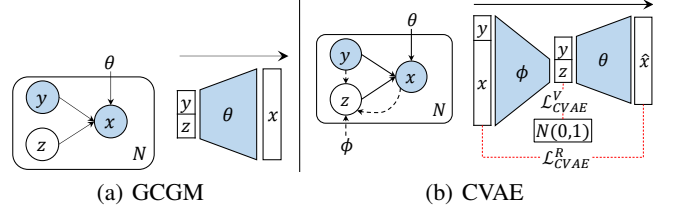


Figure 1: Schematics of the conditional generative networks, i.e., the graphical model (left) and network architecture (right). (a) GCGM shows how the observed variable,  $x$ , is generated by an unobserved latent variable  $z$  and another observed variable  $y$ . GCGM aims to get the parameter  $\theta$  so that one can sample  $x$  with specific characteristics as indicated in  $y$ . (b) In CVAE, first  $x$  and  $y$  are encoded to a latent variable  $z$  which is expected to be standard Gaussian. Next,  $y$  is again used with  $z$  to the decoder to reconstruct a variable  $\hat{x}$ . We expect variables  $x$  and  $\hat{x}$  are as same as possible.

$\text{KL} [q_{\phi}(z|y, x) || p(z)]$  is the Kullback–Leibler divergence between the posteriori distribution  $q_{\phi}(z|y, x)$  and the prior distribution  $p(z)$  [Bishop, 2006]. In this paper, we consider the prior  $p(z)$  as standard Gaussian distribution  $N(0, 1)$ . For more priors, please refer to [Kingma *et al.*, 2014; Kingma and Welling, 2014].

Therefore, a proper conditional data generation model  $p_{\theta}(x|y, z)$  and a conditional encoding model  $q_{\phi}(z|y, x)$  can be obtained by back-propagating the following CVAE loss based on the observed data set  $\mathcal{D}$ ,

$$\begin{aligned} \mathcal{L}_{CVAE}(\phi, \theta) &= \mathcal{L}_{CVAE}^R(\phi, \theta) + \mathcal{L}_{CVAE}^V(\phi), \\ &= -\mathbb{E}_{q_{\phi}(z|y, x)} [\log p_{\theta}(x|y, z)] + \text{KL} [q_{\phi}(z|y, x) || p(z)]. \end{aligned} \quad (3)$$

Here, for convenience, we named  $\mathcal{L}_{CVAE}^R$  and  $\mathcal{L}_{CVAE}^V$  as the reconstruction error and variational error, respectively. By minimizing  $\mathcal{L}_{CVAE}(\phi, \theta)$ , the conditional encoder  $q_{\phi}(z|y, x)$  promotes  $x$  and  $y$  to be the standard Gaussian distributions  $N(0, 1)$  while the conditional decoder  $p_{\theta}(x|y, z)$  reconstructs  $y$  and  $z$  back to the original sample  $x$ .

### 2.3 Lifelong Generative Problem Definition

Lifelong generative learning extends the single-distribution estimation task mentioned in Section 2.1 to a sequence of  $T$  tasks<sup>1</sup>. Concretely, each task characterized by a distinct observed dataset  $\mathcal{D}^t = \{(x_i^t, y_i^t) | i = 1, \dots, N_t\}$ ,  $t = 1, \dots, T$ , sampled from the desired distributions,  $p_{\theta_t} = p_{\theta_t}(x^t|y^t)$ . Lifelong generative models aim to estimate the true distribution  $p_{\theta} = \int \prod_{t=1}^T p_{\theta_t}(x^t|y^t, z) p(z) dz$  related the whole tasks learned so far although the tasks occur in an sequential manner.

As shown in Fig.2, during learning the second task with the only accessible dataset  $\mathcal{D}^2$ , the model not only needs to well estimate the distribution of  $\mathcal{D}^2$  but also remember the learned distribution of  $\mathcal{D}^1$ . Following this lifetime accumulated learning, the final estimated distribution  $p_{\theta_T}$  with only accessible dataset  $\mathcal{D}^T$  is required to be the same as the true distribution  $p_{\theta}$ .

<sup>1</sup>For simplicity, we use  $p_{\theta_t}(x^t|y^t)$  to replace  $p_{\theta}(x|y)$  in Eq.(1) if handling the  $t^{th}$  task with the  $t^{th}$  distribution,  $t = 1, \dots, T$ .

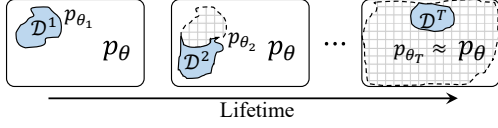


Figure 2: A demo of lifelong generative problems. Here, there is only one accessible dataset,  $\mathcal{D}^i, i = 1, \dots, T$ , indicated by the shaded area for each current task. The grid-filled area represents the distribution that has been learned before. Eventually, the true distribution,  $p_{\theta}$ , is accumulated learned after the lifetime of learning.

### 3 Proposed LGLvKR

In this section, we propose LGLvKR to well solve the lifelong generative problem mentioned in Section 2.3. Concretely, we first extend the intrinsic reconstruction character of VAE to the retention of historical knowledge. Then, the feedback consolidation strategy is conducted on the reconstructed data to ensure the superior performance. Finally, we give the whole training process and its pseudo code about our proposed LGLvKR.

#### 3.1 Knowledge Reconstruction

Recall from Section 2.3 that the key to solving the lifelong generative problem is how to retain the historically learned distributions when estimating the current task. As shown in the left-hand side of Fig.3(a), we claimed that the learned distribution is well retained in the parameter  $\theta_{t-1}$ . Here,  $y^{1:t-1}$  in Fig.3 indicates the labels of tasks learned so far.

Instead of training on the mixture data that are partly from the current task and partly generated from the previous model to retain the historical knowledge [Ramapuram *et al.*, 2020; Ye and Bors, 2021], we extend the intrinsic reconstruction character of VAE to the knowledge reconstruction. Specifically, as shown in Fig.3(b), the historical decoder  $p_{\theta_{t-1}}$  is frozen and saved after training on the dataset  $\mathcal{D}^{t-1}$ , and the knowledge about the learned tasks (i.e., from task 1 to  $t-1$ ) is retained in  $\theta_{t-1}$ . Note that the learned labels are accumulated and also saved. By inputting the current decoder,  $p_{\theta_t}$ , and the historical decoder,  $p_{\theta_{t-1}}$ , with the same latent variable  $z$  and labels  $y^{1:t-1}$ , we restrict their reconstruction output to be as consistent as possible. In this way,  $p_{\theta_t}$  could well reconstruct the historical knowledge retained in the  $p_{\theta_{t-1}}$ . Therefore, we define a knowledge reconstruction loss associated with the lifelong generation as:

$$\begin{aligned} \mathcal{L}_{LG}^R(\theta_t) \\ = -\mathbb{E}_{p_{\theta_{t-1}}(\hat{x}|y,z), y \sim U(1_y, (t-1)_y), z \sim N(0,1)} [\log p_{\theta_t}(\hat{x}|y,z)], \end{aligned} \quad (4)$$

where  $U(a, b)$  is the discrete uniform distribution,  $1_y$  and  $(t-1)_y$  indicate the unique labels contained in  $y^1$  and  $y^{t-1}$ , respectively<sup>2</sup>. It is worth noting that such a knowledge reconstruction loss about each sample has two inputs: a latent  $z$  sampled from the prior distribution, and a label-attribute  $y$  sampled from the distribution  $U(1_y, (t-1)_y)$ . Given  $z$  and

<sup>2</sup>In this paper, we assumed the labels related to the whole  $T$  tasks are well organized and disjointed, which is common in lifelong learning [Belouadah *et al.*, 2021]. Therefore,  $y \sim U(1_y, (t-1)_y)$  means uniformly sampling  $y$  from the labels learned so far.

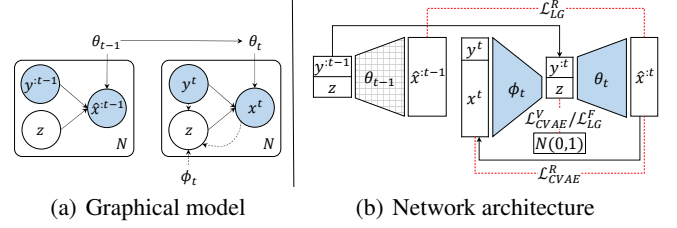


Figure 3: Schematics of the proposed LGLvKR network.

$y$ , we can generate a sample from the frozen historical decoder  $p_{\theta_{t-1}}(\hat{x}|y, z)$ , and subsequently compare this sample with the one generated by the trainable  $p_{\theta_t}(\hat{x}|y, z)$  with the same inputs. We take consistently the same comparing strategy as the one used in CVAE reconstruction loss.

#### 3.2 Feedback Consolidation

We assume in Section 3.1 that the froze historical decoder  $p_{\theta_{t-1}}$  has optimal parameters. Then, the main reason for  $p_{\theta_t}$  not remembering the considerable knowledge is due to the poor generation by the decoder  $p_{\theta_{t-1}}$ .

To guarantee the generation performance of the decoder, inspired by Verma *et al.* [2018], we introduce a feedback consolidation strategy. Specifically, as shown in Figure 3(b), when training CVAE, the encoder is used twice to ensure that the distribution obtained from the current generated data  $\hat{x}^t$  follows the true sampling distribution. Therefore, we formally design a feedback consolidation loss associated with the lifelong generation as:

$$\mathcal{L}_{LG}^F(\phi_t, \theta_t) = -\mathbb{E}_{p_{\theta_t}(\hat{x}^t|y^t, z)} [\text{KL}[q_{\phi_t}(z|y^t, \hat{x}^t) \| p(z)]]. \quad (5)$$

Note that the feedback consolidation is only devoted to the reconstructed data related to the current task  $t$ .

Hence, accompanying the CVAE loss with the knowledge reconstruction loss as well as the feedback consolidation loss, the overall learning objective is defined as:

$$\min_{\phi_t, \theta_t} \left\{ \mathcal{L}_{CVAE}(\phi_t, \theta_t) + \lambda_t^r \cdot \mathcal{L}_{LG}^R(\theta_t) + \lambda_t^f \cdot \mathcal{L}_{LG}^F(\phi_t, \theta_t) \right\}, \quad (6)$$

where the weights,  $\lambda_t^r > 0$  and  $\lambda_t^f > 0$ , are hyperparameters.

#### 3.3 Training LGLvKR

In this section, we summarize the whole training process of LGLvKR when encountering multiple sequential tasks. Look back to Fig.2, on the one hand, we could only access the observed dataset  $\mathcal{D}^1 = \{(x_i^1, y_i^1) | i = 1, \dots, N_1\}$  when encountering the first task. At this stage, we train the proposed CVAE model with the below objective function augmented only with the feedback consolidation loss,

$$\min_{\phi_1, \theta_1} \left\{ \mathcal{L}_{CVAE}(\phi_1, \theta_1) + \lambda_1^f \cdot \mathcal{L}_{LG}^F(\phi_1, \theta_1) \right\}, \quad (7)$$

where  $\lambda_1^f = 1$ . After the training, the decoder model  $p_{\theta_1}$  is frozen and saved with the unique labels contained in  $y^1$ .

On the other hand, when encountering the next task characterized by the observed dataset  $\mathcal{D}^2 = \{(x_i^2, y_i^2) | i =$

**Algorithm 1** LGL<sub>v</sub>KR

**Input:** A sequence of  $T$  datasets  $\mathcal{D}^t$ ,  $t = 1, \dots, T$ .

**Parameter:**  $\{\lambda_t^f \mid t = 1, \dots, T\}$ .

**Output:** Decoder,  $p_{\theta_T}$ , and the set of unique labels,  $\{y^T\}$ .

- 1: Observe the dataset  $\mathcal{D}_1$ .
- 2:  $\{\phi_1, \theta_1\} \leftarrow$  Update CVAE using Eq.(7) with  $\mathcal{D}_1$ .
- 3: Save the decoder,  $p_{\theta_1}$ , and the labels,  $\{y^1\}$ .
- 4: **for**  $t = 2, \dots, T$  **do**
- 5:    $\{\theta_t\} \leftarrow$  Initialize the current decoder with  $\{\theta_{t-1}\}$ .
- 6:   Observe dataset  $\mathcal{D}_t$ .
- 7:    $\{\phi_t, \theta_t\} \leftarrow$  Update the current CVAE using Eq.(6)  
with  $\mathcal{D}_t$ ,  $p_{\theta_1}$ , and  $\{y^{1:t-1}\}$ .
- 8:   Save the decoder,  $p_{\theta_t}$ , and the labels,  $\{y^t\}$ .
- 9: **end for**
- 10: **return**  $p_{\theta_T}(x^{:T} | y^{:T}, z), \{y^{:T}\}$

$1, \dots, N_2\}$ , a snapshot about  $p_{\theta_1}$  is taken to initialize  $p_{\theta_2}$  before learning. Note that, the historical encoder  $p_{\theta_1}$  is frozen while the current encoder  $p_{\theta_1}$  is trainable throughout the training stage. At this stage, we train the proposed CVAE model using Eq.(6) with  $t = 2$ . Since we have saved the unique labels in  $y^1$ , the knowledge reconstruction loss,  $\mathcal{L}_{LG}^R$ , could be well calculated. After learning this task, again, the current decoder model  $p_{\theta_2}$  is frozen and saved along with the accumulated unique labels contained in  $y^2$  (i.e.,  $y^1 \cup y^2$ ). Along with this iterative way, we train LGLvKR until finishing learning the whole sequential  $T$  tasks. Here, the weights  $\lambda_t^f = 1, t = 1, \dots, T$  while  $\lambda_t^r = t - 1, t = 1, \dots, T$ . We summarize our method LGLvKG in Algorithm 1.

## 4 Experiments

We evaluated the LGLvKR mainly in two aspects: (1) the error accumulation, and (2) the time complexity on different datasets.

**Baseline Methods.** There are three types of training strategies in our baseline, fine tuning, joint training, and the existing lifelong training. Specifically, the fine tuning sequentially trains the model with parameters initialized from the recently fine-tuned model on the previous task. The joint training strategy trains the model on the combined real data from tasks seen so far. As for the lifelong training, we engaged two widely adopted approaches, pseudo rehearsal and the approach trained only on the dataset of the current task. In particular, both of them were built upon the conditional generative model. The pseudo rehearsal approaches are implemented on CGAN and CVAE models [van de Ven *et al.*, 2020; Ramapuram *et al.*, 2020; Ye and Bors, 2021]. For CVAE, only the decoder used the label information. As for the other approaches trained only on the current dataset, we engaged CGAN+EWC [Seff *et al.*, 2017] and CWGAN+RA [Wu *et al.*, 2018]. In addition, we also compared with the methods that only take the feedback consolidation (LGL-noKR) or knowledge reconstruction (LGL-noFC) on our proposed CVAE where both encoder and decoder use the label information.

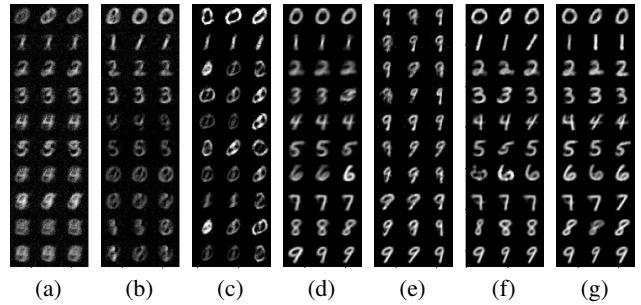


Figure 4: The lifelong generated digits after sequentially learning 10 tasks with methods (a) CGAN, (b) CGAN+EWG, (c) CWGAN+RA, (d) CVAE, (e) LGL-noKR, (f) LGL-noFC, and (g) LGLvKR.

**Quantitative Metrics.** We used the accuracy (ACC) and the Fréchet Inception Distance (FID) [Heusel *et al.*, 2017] as the quantitative metrics. ACC is the accuracy of the classifier trained on the generated dataset and evaluated on real samples (higher ACC indicates better generation results). FID, instead of evaluating the generated samples directly, compares the statistics between the generated dataset and the real one (lower FID indicates higher generation quality).

### 4.1 Forgetting and Error Accumulation Test

**Experimental Setting.** To evaluate the effectiveness of our LGLvKR, we consider the lifelong digits generation problem [Zhai *et al.*, 2019] in the MNIST [LeCun, 1998] dataset. It consists of  $32 \times 32$  pixels images of handwritten digits. Each task contains only one digital category which appears in an ascending order, thus we get 10 separate tasks. Note that each method takes a small amount of training with only 5 epochs on each task. For fairness, we implemented all methods on the network with the same level of parameters (i.e., each model has 0.66 million parameters). Specifically, CGAN and CGAN+EWC use the same architecture [Linder-Norén, 2018] while CWGAN+RA takes the original network [Liu *et al.*, 2020]. LGL-noKR, LGL-noFC, and LGLvKR use the networks as shown in Fig.1(b) while CVAE has one less layer in encoding label than the network used by LGLvKG. We recorded the average ACC and FID for each method after running 10 times with different seeds.

**Analysis.** We analyze the results from three aspects. Firstly, in Tab.1 and Fig.4(a), fine tuning the generative networks (used in CGAN, CWGAN+EWC, CVAE, LGLvKR) will incur the catastrophic forgetting problem since its ACC dramatically degrades from 100% to the level below 50%. Correspondingly, the FID value substantially increases, which means the quality of the generated images decreases.

Secondly, we see that the joint training gets good results, which are the upper bounds of lifelong learning. It confirms that the decline of fine tuning results is not limited by the network capacity but by the catastrophic forgetting problem. Besides, by comparing the results of VAE-based models (CVAE, LGLvKR) and GAN-based models (CGAN, CWGAN+RA), we found that the methods based on VAE relatively outperform methods based on GAN. That is mainly because the

Strategy	Task	CGAN		CGAN+EWC		CWGAN+RA		CVAE		LGL-noKR		LGL-noFC		LGLvKR	
		ACC	FID	ACC	FID	ACC	FID	ACC	FID	ACC	FID	ACC	FID	ACC	FID
Fine Tuning	1	100	0.62	-	-	100	0.52	100	0.34	-	-	-	-	100	0.29
	2	91.02	0.67	-	-	91.16	0.58	93.99	0.53	-	-	-	-	88.53	0.55
	5	53.09	0.84	-	-	71.58	0.44	62.42	0.39	-	-	-	-	64.67	0.38
	8	19.00	0.71	-	-	54.96	0.50	47.68	0.44	-	-	-	-	35.80	0.41
	10	12.90	0.74	-	-	48.45	0.58	46.99	0.43	-	-	-	-	31.99	0.41
Joint Training	1	100	0.69	-	-	100	0.53	100	0.39	-	-	-	-	100	0.37
	2	98.68	0.79	-	-	99.80	0.40	99.80	0.41	-	-	-	-	99.85	0.41
	5	80.90	0.59	-	-	94.45	0.35	95.39	0.34	-	-	-	-	95.29	0.34
	8	58.51	0.60	-	-	88.97	0.35	90.34	0.37	-	-	-	-	90.89	0.37
	10	51.66	0.68	-	-	81.98	0.32	85.85	0.33	-	-	-	-	86.01	0.32
Lifelong Training	1	100	0.74	100	0.73	100	0.52	100	0.38	100	0.32	100	0.32	100	0.30
	2	99.71	0.80	99.32	0.73	99.76	0.53	<b>99.90</b>	0.40	99.51	0.51	99.76	0.37	99.85	0.39
	3	94.73	0.85	94.76	0.66	63.09	0.61	96.42	0.38	84.90	0.42	96.88	0.36	<b>97.36</b>	0.36
	4	91.77	0.86	88.09	0.73	67.77	0.63	93.87	0.44	69.70	0.55	94.53	0.40	<b>95.48</b>	0.38
	5	82.19	0.84	79.59	0.55	67.68	0.64	93.24	0.40	68.75	0.42	94.75	0.35	<b>94.96</b>	0.35
	6	79.14	0.89	72.61	0.71	62.25	0.66	90.21	0.39	41.69	0.35	91.42	0.35	<b>91.44</b>	0.34
	7	68.29	0.93	65.05	0.67	63.08	0.67	89.61	0.41	52.04	0.42	90.12	0.35	<b>90.54</b>	0.32
	8	55.24	0.95	60.36	0.70	49.16	0.68	88.90	0.46	50.71	0.48	89.82	0.36	<b>89.84</b>	0.35
	9	53.06	0.93	34.99	0.77	46.82	0.68	85.49	0.49	37.83	0.57	86.18	0.36	<b>86.58</b>	0.33
	10	46.44	0.90	32.65	0.71	44.20	0.69	82.67	0.51	43.06	0.47	82.27	0.33	<b>83.94</b>	0.31

Table 1: The ACC (%) and FID values of different methods with three training strategies.

GAN-based models are not well trained with 5 epochs while it is enough for the VAE-based model, indicating the high efficiency in the training process of VAE. On the other hand, by comparing CVAE and LGLvKR, a relative improvement could be obtained by embedding  $y$  in both encoder and decoder. It is prominent when there are more categories, such as 86.01 for LGLvKR versus 85.85 for CVAE after learning 10 tasks.

Finally, for the lifelong training, LGLvKR achieved the best ACC, almost reaching its upper bound. In contrast, the CVAE method using pseudo rehearsal also obtained considerable ACC, but its FID results are not good, which was also verified by the generated images. As shown in Fig.4(d) and Fig.4(g), the generated digits (from 1 to 5) of CVAE are more blur compared with those of LGLvKR. The GAN-based methods performed worse. First, they could not be well trained in 5 epochs. Besides, error accumulation heavily influences their results. Especially for the CWGAN+RA, its ACC dramatically degraded to 44.20 while its model reached 81.98 with joint training. In Tab.1, Fig.4(a), and Fig.4(b), CGAN and CGAN+EWC obtained appropriate performance against a few task sequences but failed when the task length exceeds 5.

We also conducted ablation studies. As shown in Tab.1 and Fig.4(e), LGL-noKR showed the similar catastrophic forgetting problem since it trained with CVAE loss of Eq.(3) and the feedback consolidation loss of Eq.(5). By expending the intrinsic reconstruction character of VAE to the knowledge reconstruction, LGL-FC well retained the historical tasks seen so far. And it was further improved with the help of feedback consolidation, which came to the LGLvKR.

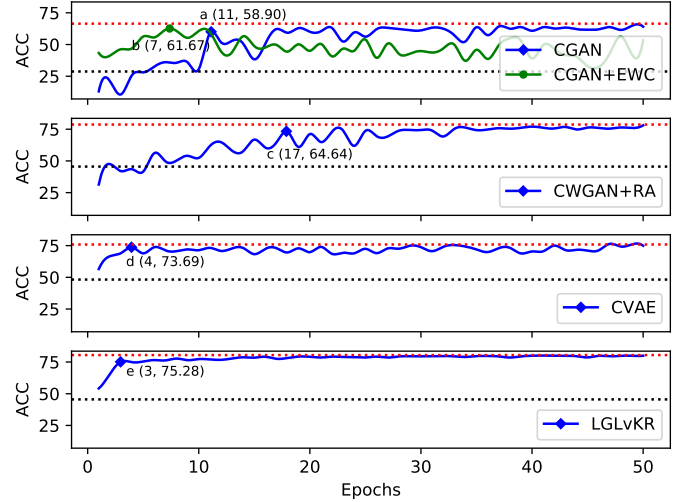


Figure 5: ACC of various lifelong generative methods after sequentially training 10 separate fashion-MNIST tasks. The black and red dash lines indicate, respectively, their corresponding lower and upper bounds.

## 4.2 Time Consumption Test

**Experimental Setting.** We evaluate the efficiency of LGLvKR on the fashion-MNIST dataset [Xiao *et al.*, 2017]. Specifically, we tested the ACC to the number of epochs (from 1 to 50) trained on each task and the time required to achieve comparable results. Similar to Section 4.1, we implemented each method on the network with the same level



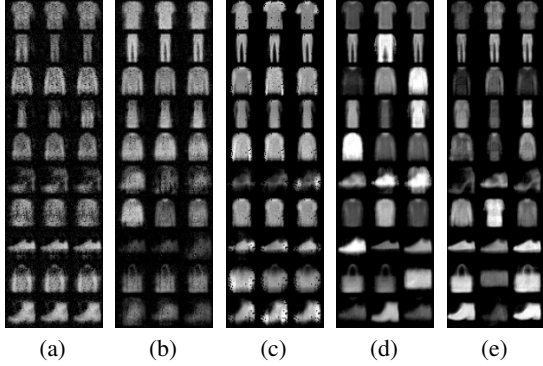


Figure 6: The generated images from methods (a) CGAN (700s), (b) CGAN+EWC (388s), (c) CWGAN+RA (1230s), (d) CVAE (217s), and (e) LGLvKR (157s), with their average training time (seconds). We took the models to generate images when the methods first reached the best ACC, while the best ACC was stressed in Fig.5.

of parameters. And the fashion-MNIST dataset was divided into 10 separate tasks. We took the best results of fine tuning and joint training over 50 epochs as the lower and upper bounds of ACC, respectively. With the fixed epoch  $i$ , each model trains  $i$  epochs on each task after sequentially training 10-tasks. We saved the trained generative model. Then, a classifier was trained on the samples generated from the saved model until convergence. We finally took the classifier to calculate the ACC on the real test dataset. Each experiment runs 10 times with different seeds.

**Analysis.** Since CGAN and CGAN+EWC engaged the same network, they shared the same upper and lower bounds. As shown in Fig.5, CGAN and CGAN+EWC first reached their optimal ACC, with 11 and 7 epochs training, respectively. Due to the nature drawback of unstable training [Arjovsky *et al.*, 2017], CGAN was unstable as the number of epochs grows. In particular, this unstableness pronounced for CGAN+EWC owing to its poor performance on the long length of sequent tasks. By introducing the Wasserstein metric and emphasizing the importance of label in the discriminator [Wu *et al.*, 2018; Liu *et al.*, 2020], CWGAN+RA obtained a higher upper bound of ACC (66.45% for CGAN/CGAN+EWC versus 78.66% for CWGAN+RA) and more stable results on lifelong generating learning. However, CWGAN+RA conducted 17 epochs to reach the best ACC for the first time. Compared with CWGAN+RA, although CVAE got an inferior upper bound in joint training (75.88% for CVAE versus 78.66% for CWGAN+RA), it obtained more stable results benefiting from its stable training mechanism. CVAE took 4 epochs to reach the optimal ACC for the first time. By embedding the label information in both the encoder and decoder, LGLvKR achieved a higher ACC upper bound (80.47%) compared with CVAE (75.88%) with the joint training strategy. On the other hand, instead of pseudo rehearsals, LGLvKR took knowledge reconstruction to retain the historical knowledge and was enhanced by a feedback consolidation. It got both the optimal ACC with only 3 epochs and a considerable improvement in stability.

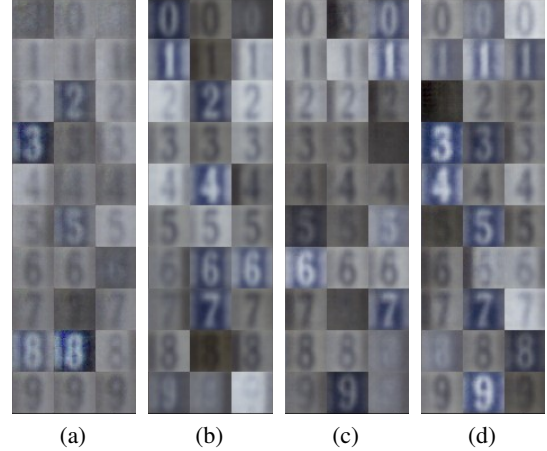


Figure 7: The generated colored images from CVAE with lifelong training (a) and joint training strategy (b), along with LGLvKR with lifelong training (c) and joint training (d).

As demonstrated in Fig.6, overall, LGLvKR achieved lifelong learning with the shortest time while getting a good picture of diversity. CVAE in Fig.6(e) whereas confused Sandal (the sixth class) and Ankle boot (the last class).

Finally, we showed the generated images by CVAE and LGLvKR on the colored Street View House Number [Netzer *et al.*, 2011]. In particular, we trained each model with lifelong training and joint training strategies. We referred [Subramanian, 2021] to implement the models. Other experimental settings are similar to those in fashion-MNIST while we fixed the training epochs with 10. As shown in Fig.7, the FID values of all the images were 0.40 (Fig.7(a)), 0.36 (Figs.7(b)), 0.35 (Fig.7(c)), and 0.32 (Fig.7(d)), respectively. The training time of the corresponding methods are, respectively, 1185, 4870, 1059, and 4885 seconds. Overall, LGLvGR achieved considerable results with the shortest training time.

## 5 Conclusion

We studied the lifelong generative learning problem based on a variant of the CVAE. By expending the intrinsic reconstruction character of VAE to reconstruct the historical knowledge learned before, LGLvKR could well handle the catastrophic forgetting problem. In addition, to alleviate the error accumulation, we further developed a feedback strategy for LGLvKR. Experiments on MNIST, FashionMnist, and SVHN verified that the proposed LGLvKR addressed the lifelong generating problem effectively and efficiently.

This work mainly focuses on verifying that the intrinsic reconstruction character of VAE is good enough to retain the historical knowledge and we extend it to lifelong generating learning. However, the lifelong generating output of a more complex image is still an open challenge [Lesort, 2020]. Thus, in the future, it is desirable to introduce a more generalized architecture of the network in the LGLvKR framework. Also, one alternative direction is to generate the high-level features from the complex images in lifelong learning.

## References

- Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International conference on machine learning*, pages 214–223. PMLR, 2017.
- Eden Belouadah, Adrian Popescu, and Ioannis Kanellos. A comprehensive study of class incremental learning algorithms for visual tasks. *Neural Networks*, 135:38–54, 2021.
- Christopher M. Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006.
- Matthias Delange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Ales Leonardis, Greg Slabaugh, and Tinne Tuytelaars. A continual learning survey: Defying forgetting in classification tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- Ronald Kemker, Marc McClure, Angelina Abitino, Tyler Hayes, and Christopher Kanan. Measuring catastrophic forgetting in neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *Proceedings of the 2nd International Conference on Learning Representations*, 2014.
- Diederik P Kingma, Shakir Mohamed, Danilo Jimenez Rezende, and Max Welling. Semi-supervised learning with deep generative models. In *Advances in neural information processing systems*, pages 3581–3589, 2014.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- Yann LeCun. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- Timothée Lesort. Continual learning: Tackling catastrophic forgetting in deep neural networks with replay processes. *arXiv preprint arXiv:2007.00487*, 2020.
- Erik Linder-Norén. Pytorch-gan. <https://github.com/eriklindernoren/PyTorch-GAN/>, 2018. Accessed: 2021-7-29.
- Xialei Liu, Chenshen Wu, Mikel Menta, Luis Herranz, Bogdan Raducanu, Andrew D Bagdanov, Shangling Jui, and Joost van de Weijer. Generative feature replay for class-incremental learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 226–227, 2020.
- Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. *Workshop on Deep Learning and Unsupervised Feature Learning, Neural Information Processing System*, 2011.
- German I Parisi, Ronald Kemker, Jose L Part, Christopher Kanan, and Stefan Wermter. Continual lifelong learning with neural networks: A review. *Neural Networks*, 113:54–71, 2019.
- Jason Ramapuram, Magda Gregorova, and Alexandros Kalousis. Lifelong generative modeling. *Neurocomputing*, 404:381–400, 2020.
- Ari Seff, Alex Beatson, Daniel Suo, and Han Liu. Continual learning in generative adversarial nets. *arXiv preprint arXiv:1705.08395*, 2017.
- Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual learning with deep generative replay. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 2994–3003, 2017.
- Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning structured output representation using deep conditional generative models. *Advances in neural information processing systems*, 28:3483–3491, 2015.
- Anand K Subramanian. Pytorch-vae. <https://github.com/AntixK/PyTorch-VAE>, 2021. Accessed: 2021-03-01.
- Gido M Van de Ven and Andreas S Tolias. Three continual learning scenarios. *Continual Learning Workshop, Neural Information Processing Systems*, 2018.
- Gido M van de Ven, Hava T Siegelmann, and Andreas S Tolias. Brain-inspired replay for continual learning with artificial neural networks. *Nature communications*, 11(1):1–14, 2020.
- Vinay Kumar Verma, Gundeep Arora, Ashish Mishra, and Piyush Rai. Generalized zero-shot learning via synthesized examples. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4281–4289, 2018.
- Chenshen Wu, Luis Herranz, Xialei Liu, Yaxing Wang, Joost Van de Weijer, and Bogdan Raducanu. Memory replay gans: learning to generate images from new categories without forgetting. *arXiv preprint arXiv:1809.02058*, 2018.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- Fei Ye and Adrian Bors. Lifelong teacher-student network learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- Mengyao Zhai, Lei Chen, Frederick Tung, Jiawei He, Megha Nawhal, and Greg Mori. Lifelong gan: Continual learning for conditional image generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2759–2768, 2019.