

Appendix for Lifelong Variational Autoencoder Using Online Adversarial Expansion Strategy

December 1, 2022

Contents

A	Additional information for the proposed OAES	3
B	Theorem 1	4
C	Lemma 2	6
D	Lemma 3	8
E	Theoretical analysis for the expansion threshold	8
F	Theoretical analysis for other VAE models	9
F.1	Importance weighted autoencoders	9
F.2	Hierarchical Variational Inference	10
G	Theoretical analysis for task-known continual learning	11
G.1	Theoretical analysis for the static network architecture	11
G.1.1	Memory-based model	11
G.1.2	GRM-based model	12
G.2	Theoretical analysis for the dynamic expansion model	16
G.3	Theoretical analysis for the existing GRM-based models	18
G.3.1	Lifelong VAEGAN	19
G.3.2	Lifelong infinite mixture model	19
H	Additional information for experiments	20
H.1	Additional information for experiment settings	20
H.2	Experiment setting	21
H.3	The configuration for the classification task.	21

I	Additional information for ablation study	23
I.1	The impact of the threshold β	23
I.2	The impact of the memory buffer size	24
I.3	Knowledge diversity among experts	24
I.4	The impact of batch size change	25
I.5	Theoretical results	26
I.6	Classification on fuzzy task boundaries	26
I.7	Comparison of computational costs	27
I.8	The other performance criterion on the generative modelling	27

A Additional information for the proposed OAES

In this section, we provide the detailed learning procedure of the proposed model in Fig. 1.

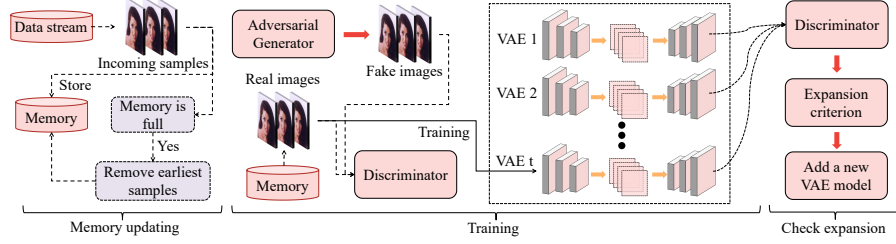


Figure 1: The learning procedure of the proposed Online Adversarial Expansion Strategy (OAES) model. During the memory updating stage, we continually add incoming samples to the memory buffer. When the memory buffer is full, we automatically remove the earliest added samples. In the second stage, if the proposed OAES does not have any VAE components, we build the first VAE component which is then trained until for \mathcal{T}_{100} steps. Then we freeze the first component and treat it as the adversary in the subsequent learning. If we have already learnt t components $\mathbf{V} = \{\mathcal{V}_1^{c_1}, \dots, \mathcal{V}_t^{c_t}\}$ at \mathcal{T}_i , we only train $\mathcal{V}_t^{c_t}$ on \mathcal{M}_i at \mathcal{T}_i by using Eq.(1) of the paper to avoid forgetting previously learnt knowledge. In addition, we train the generator G_{ε^i} and the discriminator D_{ψ^i} on \mathcal{M}_i using adversarial loss (Eq.(14) of the paper). In the third stage, we check the expansion criterion using the evaluation stage of OAES (Eq.(15) of the paper). If Eq.(15) of the paper is satisfied, we add a new component \mathcal{V}_{t+1} into \mathbf{V} and clear up \mathcal{M}_i in order to learn data samples which are statistically non-overlapping in the following training step.

In the following, we provide the pseudocode of the proposed OAES in Algorithm 1, which can be summarized into four stages :

Stage 1 . Memory updating : Let $|\mathcal{M}|^{max}$ be the maximum number of samples in the memory buffer. Since this paper does not focus on the sample selection for the memory buffer, we simply remove the earliest stored samples and add new incoming samples into \mathcal{M}_i at the i -th training step (\mathcal{T}_i), if the maximum memory buffer size $|\mathcal{M}|^{max}$ is reached.

Stage 2 . Training a component : If \mathbf{V} has a single component, then we train \mathcal{V}_1 on \mathcal{M}_i until finishes the training step $\mathcal{T}_i = |\mathcal{M}|^{max}$ in order to preserve the initial information about the data stream. Then we build the second component \mathcal{V}_2 in the subsequent learning. We describe the following training as follows. Let us suppose that we have already learnt t components $\mathbf{V} = \{\mathcal{V}_1^{c_1}, \dots, \mathcal{V}_t^{c_t}\}$ at \mathcal{T}_i , we only train $\mathcal{V}_t^{c_t}$ on \mathcal{M}_i at \mathcal{T}_i by using Eq.(1) of the paper to avoid forgetting previously learnt knowledge. In

addition, we train the generator G_{ε^i} and the discriminator D_{ψ^i} on \mathcal{M}_i using adversarial loss (Eq.(14) of the paper).

Stage 3 . Check the expansion : If $|\mathcal{M}_i|$ reaches $|\mathcal{M}|^{max}$, then we check the expansion criterion using the evaluation stage of OAES (Eq.(15) of the paper). If Eq.(15) of the paper is satisfied, we add a new component \mathcal{V}_{t+1} into \mathbf{V} and clear up the memory \mathcal{M}_i in order to learn statistically non-overlapping samples in the following training step.

Stage 4 . Component selection at the testing phase : For a given sample \mathbf{x} , we choose a component with the highest sample log-likelihood by :

$$s^* = \arg \max_{s=1, \dots, |\mathbf{V}|} \{\mathcal{L}_{ELBO}(\mathbf{x}; \mathcal{V}_s)\}, \quad (1)$$

Algorithm 1 Algorithm for OAES

```

1: (Input: The data stream);
2: for  $i < n$  do
3:   Memory updating
4:    $\mathcal{B}_i \sim W$ 
5:    $\mathcal{M}_i = \mathcal{M}_i \cup \mathcal{B}_i$ 
6:   if  $|\mathcal{M}_i| > |\mathcal{M}|^{max}$  then
7:      $\mathcal{M}_i = \bigcup_{j=10}^{|\mathcal{M}|^{max}+10} \mathcal{M}_i$ 
8:   end if
9:   Training the component
10:  if  $|\mathbf{V}| = 1$  and  $\mathcal{T}_i = |\mathcal{M}|^{max}$  then
11:    Add the second component  $\mathcal{V}_2$ 
12:  end if
13:  Train the current VAE component  $\mathcal{V}_t$  on  $\mathcal{M}_i$  using  $\mathcal{L}_{ELBO}$ 
14:  Train the generator  $G_{\varepsilon^i}$  and the discriminator  $D_{\psi^i}$  on  $\mathcal{M}_i$  using adversarial loss
15:  Check the expansion
16:  if  $|\mathcal{M}_i| > |\mathcal{M}|^{max}$  then
17:    if  $\min \left\{ C_{\psi^i}(\mathbb{P}_{\theta_1^{c_1}}, \mathbb{P}_{\theta_1^i}), \dots, C_{\psi^i}(\mathbb{P}_{\theta_{t-1}^{c_{t-1}}}, \mathbb{P}_{\theta_t^i}) \right\} \geq \beta$  then
18:      Add a new Component  $\mathcal{V}_{t+1}$ 
19:    end if
20:  end if
21: end for
22: for  $i < n'$  do
23:    $\mathbf{x} \sim \mathbf{X}_{test}$ 
24:    $s^* = \arg \max_{s=1, \dots, |\mathbf{V}|} \{\mathcal{L}_{ELBO}(\mathbf{x}; \mathcal{V}_s)\}$ 
25:   Choose  $\mathcal{V}_{s^*}$  for the evaluation.
26: end for

```

B Theorem 1

Theorem 1 Let $p_{\theta_i}(\mathbf{x})$ be a probability density function for a single model \mathcal{V}^i updated at \mathcal{T}_i . Let \mathbb{P}_i^W denote a distribution of all visited data batches $\{\mathcal{B}_1, \dots, \mathcal{B}_i\}$ drawn

from W at \mathcal{T}_i . Let $p_{\mathcal{M}_i}(\mathbf{x})$ and $p_{W^i}(\mathbf{x})$ denote the density function for $\mathbb{P}_{\mathcal{M}_i}$ and \mathbb{P}_i^W , respectively. We then derive an upper bound for a single VAE model trained on \mathcal{M}_i at \mathcal{T}_i as :

$$\begin{aligned}\mathbb{E}_{\mathbb{P}_i^W} [\log p_{\theta^i}(\mathbf{x})] &\geq \mathbb{E}_{\mathbb{P}_{\mathcal{M}_i}} [\log p_{\theta^i}(\mathbf{x})] \\ &\quad - D_{KL}(\mathbb{P}_i^W \parallel \mathbb{P}_{\mathcal{M}_i}) - \mathcal{F}_{\text{DL}}(\mathbb{P}_{\mathcal{M}_i}, \mathbb{P}_i^W, \mathbb{P}_{\theta^i}) \\ &\quad + \mathcal{F}_{\text{dis}}(\mathbb{P}_i^W, \mathbb{P}_{\mathcal{M}_i}),\end{aligned}\tag{2}$$

where $\mathcal{F}_{\text{DL}}(\mathbb{P}_{\mathcal{M}_i}, \mathbb{P}_i^W, \mathbb{P}_{\theta^i})$ is defined as :

$$\begin{aligned}\mathcal{F}_{\text{DL}}(\mathbb{P}_{\mathcal{M}_i}, \mathbb{P}_i^W, \mathbb{P}_{\theta^i}) &\triangleq |D_{KL}(\mathbb{P}_{\mathcal{M}_i} \parallel \mathbb{P}_{\theta^i}) \\ &\quad - D_{KL}(\mathbb{P}_i^W \parallel \mathbb{P}_{\theta^i})|\end{aligned}\tag{3}$$

and $\mathcal{F}_{\text{dis}}(\mathbb{P}_i^W, \mathbb{P}_{\mathcal{M}_i})$ is :

$$\begin{aligned}\mathcal{F}_{\text{dis}}(\mathbb{P}_i^W, \mathbb{P}_{\mathcal{M}_i}) &\triangleq \mathbb{E}_{\mathbb{P}_i^W} [p_{W^i}(\mathbf{x}) \log p_{W^i}(\mathbf{x})] \\ &\quad - \mathbb{E}_{\mathbb{P}_{\mathcal{M}_i}} [p_{\mathcal{M}_i}(\mathbf{x}) \log p_{\mathcal{M}_i}(\mathbf{x})]\end{aligned}\tag{4}$$

We can observe that $\mathcal{F}_{\text{dis}}(\mathbb{P}_i^W, \mathbb{P}_{\mathcal{M}_i})$ is constant if and only if \mathbb{P}_i^W and $\mathbb{P}_{\mathcal{M}_i}$ are fixed. $\mathcal{F}_{\text{dis}}(\mathbb{P}_{W^i}, \mathbb{P}_{\mathcal{M}_i})$ can also be bounded by $|D_{KL}(\mathbb{P}_i^W \parallel \mathbb{P}_{\mathcal{M}_i}) - D_{KL}(\mathbb{P}_{\mathcal{M}_i} \parallel (\mathbb{P}_i^W))|$. Based on Eq. (2), we can estimate the sample log-likelihood of \mathbb{P}_i^W by ELBO :

$$\begin{aligned}\mathbb{E}_{\mathbb{P}_i^W} [\log p_{\theta^i}(\mathbf{x})] &\geq \mathbb{E}_{\mathbb{P}_{\mathcal{M}_i}} [\mathcal{L}_{\text{ELBO}}(\mathbf{x}; \theta^i, \omega^i)] \\ &\quad - D_{KL}(\mathbb{P}_i^W \parallel \mathbb{P}_{\mathcal{M}_i}) - \mathcal{F}_{\text{DL}}(\mathbb{P}_{\mathcal{M}_i}, \mathbb{P}_i^W, \mathbb{P}_{\theta^i}) \\ &\quad + \mathcal{F}_{\text{dis}}(\mathbb{P}_i^W, \mathbb{P}_{\mathcal{M}_i}),\end{aligned}\tag{5}$$

We find that Eq. (5) can be recovered to a standard ELBO (Eq. (1) of the paper) if and only if \mathbb{P}_i^W is equal to $\mathbb{P}_{\mathcal{M}_i}$.

Proof. Firstly, we consider the JS divergence $D_{JS}(\mathbb{P}_i^W \parallel \mathbb{P}_{\mathcal{M}_i})$ and two KL divergences $D_{KL}(\mathbb{P}_{\mathcal{M}_i} \parallel \mathbb{P}_{\theta^i})$ and $D_{KL}(\mathbb{P}_i^W \parallel \mathbb{P}_{\theta^i})$. Then we have :

$$\begin{aligned}D_{KL}(\mathbb{P}_i^W \parallel \mathbb{P}_{\theta^i}) &\leq D_{KL}(\mathbb{P}_{\mathcal{M}_i} \parallel \mathbb{P}_{\theta^i}) + |D_{KL}(\mathbb{P}_{\mathcal{M}_i} \parallel \mathbb{P}_{\theta^i}) - D_{KL}(\mathbb{P}_i^W \parallel \mathbb{P}_{\theta^i})| \\ &\quad + D_{JS}(\mathbb{P}_i^W \parallel \mathbb{P}_{\mathcal{M}_i})\end{aligned}\tag{6}$$

Eq. (6) holds because the sum of last two terms in RHS is larger than LHS while

8 $D_{KL}(\mathbb{P}_i^W \parallel \mathbb{P}_{\mathcal{M}_i}) \geq 0$. In the following, we can rewrite $D_{KL}(\mathbb{P}_i^W \parallel \mathbb{P}_{\theta^i})$ and
 9 $D_{KL}(\mathbb{P}_{\mathcal{M}_i} \parallel \mathbb{P}_{\theta^i})$ as :

$$D_{KL}(\mathbb{P}_i^W \parallel \mathbb{P}_{\theta^i}) = \mathbb{E}_{\mathbb{P}_i^W} [p_{W^i}(\mathbf{x}) \log p_{W^i}(\mathbf{x})] - \mathbb{E}_{\mathbb{P}_i^W} [\log p_{\theta^i}(\mathbf{x})] \quad (7)$$

$$D_{KL}(\mathbb{P}_{\mathcal{M}_i} \parallel \mathbb{P}_{\theta^i}) = \mathbb{E}_{\mathbb{P}_{\mathcal{M}_i}} [p_{\mathcal{M}_i}(\mathbf{x}) \log p_{\mathcal{M}_i}(\mathbf{x})] - \mathbb{E}_{\mathbb{P}_{\mathcal{M}_i}} [\log p_{\theta^i}(\mathbf{x})] \quad (8)$$

10 Then we take Eq. (7) and Eq. (8) into Eq. (6), resulting in :

$$\begin{aligned} \mathbb{E}_{\mathbb{P}_i^W} [p_{W^i}(\mathbf{x}) \log p_{W^i}(\mathbf{x})] - \mathbb{E}_{\mathbb{P}_i^W} [\log p_{\theta^i}(\mathbf{x})] &\leq \mathbb{E}_{\mathbb{P}_{\mathcal{M}_i}} [p_{\mathcal{M}_i}(\mathbf{x}) \log p_{\mathcal{M}_i}(\mathbf{x})] \\ &\quad - \mathbb{E}_{\mathbb{P}_{\mathcal{M}_i}} [\log p_{\theta^i}(\mathbf{x})] \\ &\quad + |D_{KL}(\mathbb{P}_{\mathcal{M}_i} \parallel \mathbb{P}_{\theta^i}) - D_{KL}(\mathbb{P}_i^W \parallel \mathbb{P}_{\theta^i})| \\ &\quad + D_{JS}(\mathbb{P}_i^W \parallel \mathbb{P}_{\mathcal{M}_i}) \end{aligned} \quad (9)$$

11 We rearrange Eq. (9) as :

$$\begin{aligned} -\mathbb{E}_{\mathbb{P}_i^W} [\log p_{\theta^i}(\mathbf{x})] &\leq \mathbb{E}_{\mathbb{P}_{\mathcal{M}_i}} [p_{\mathcal{M}_i}(\mathbf{x}) \log p_{\mathcal{M}_i}(\mathbf{x})] - \mathbb{E}_{\mathbb{P}_i^W} [p_{W^i}(\mathbf{x}) \log p_{W^i}(\mathbf{x})] \\ &\quad - \mathbb{E}_{\mathbb{P}_{\mathcal{M}_i}} [\log p_{\theta^i}(\mathbf{x})] + |D_{KL}(\mathbb{P}_{\mathcal{M}_i} \parallel \mathbb{P}_{\theta^i}) - D_{KL}(\mathbb{P}_i^W \parallel \mathbb{P}_{\theta^i})| \\ &\quad + D_{JS}(\mathbb{P}_i^W \parallel \mathbb{P}_{\mathcal{M}_i}) \end{aligned} \quad (10)$$

12 We then multiply with -1 both LHS and RHS of Eq. (10), resulting in :

$$\begin{aligned} \mathbb{E}_{\mathbb{P}_i^W} [\log p_{\theta^i}(\mathbf{x})] &\geq -\mathbb{E}_{\mathbb{P}_{\mathcal{M}_i}} [p_{\mathcal{M}_i}(\mathbf{x}) \log p_{\mathcal{M}_i}(\mathbf{x})] + \mathbb{E}_{\mathbb{P}_i^W} [p_{W^i}(\mathbf{x}) \log p_{W^i}(\mathbf{x})] \\ &\quad + \mathbb{E}_{\mathbb{P}_{\mathcal{M}_i}} [\log p_{\theta^i}(\mathbf{x})] - |D_{KL}(\mathbb{P}_{\mathcal{M}_i} \parallel \mathbb{P}_{\theta^i}) - D_{KL}(\mathbb{P}_i^W \parallel \mathbb{P}_{\theta^i})| \\ &\quad - D_{JS}(\mathbb{P}_i^W \parallel \mathbb{P}_{\mathcal{M}_i}) \end{aligned} \quad (11)$$

13 □

14 This proves Theorem 1.

15 C Lemma 2

16 **Lemma 2** Let $\{D^T(1, k), \dots, D^T(C(T, k), k)\}$ be several target sets where each tar-
 17 get set $D^T(c, k)$ can be divided into several data batches $\{\mathcal{B}^T(c, 1), \dots, \mathcal{B}^T(c, n(T, c, k))\}$
 18 where $n(T, c, k)$ is the total number of data batches for $D^T(c, k)$. Let $\mathbb{P}_T^{\mathcal{B}}(c, j)$ repre-

sent the probabilistic representation of the data batch $\mathcal{B}^T(c, j)$. We suppose that \mathbf{V} has already learnt t components trained on \mathcal{M}_i at \mathcal{T}_i . The generalization performance on all target sets, achieved by \mathbf{V} at \mathcal{T}_i , is defined as :

$$\sum_{c=1}^{C_k^T} \left\{ \sum_{j=1}^{n_{c,k}^T} \left\{ \mathbb{E}_{\mathbb{P}_{T(c,j)}^{\mathcal{B}}} [\log p_{\Theta^i}(\mathbf{x})] \right\} \right\} \geq \sum_{c=1}^{C_k^T} \left\{ \sum_{j=1}^{n_{c,k}^T} \left\{ \mathcal{F}_s(\mathbb{P}_T^{\mathcal{B}}(c, j), \mathbf{V}) \right\} \right\}, \quad (12)$$

which is Eq. (10) in the paper.

Similar to the conclusion of Theorem 2, increasing the number of components in \mathbf{V} would lead to better generalization performance on all target sets. In practice, we can not use Eq.(10) of the paper for the component selection since it involves several intractable terms. Instead, we perform the component selection by comparing the sample log-likelihood, expressed as :

$$\mathcal{F}_s^r(\mathbb{P}_j^{\mathcal{B}}, \mathbf{V}) \triangleq \arg \max_{i=1, \dots, t} \left\{ \mathbb{E}_{\mathbb{P}_j^{\mathcal{B}}} [\mathcal{L}_{ELBO}(\mathbf{x}; \theta_i^{c_i}, \omega_i^{c_i})] \right\}, \quad (13)$$

We can rewrite Eq. (12) by using the component selection function (Eq. (13)) :

$$\sum_{c=1}^{C_k^T} \left\{ \sum_{j=1}^{n_{c,k}^T} \left\{ \mathbb{E}_{\mathbb{P}_{T(c,j)}^{\mathcal{B}}} [\log p_{\Theta^i}(\mathbf{x})] \right\} \right\} \geq \sum_{c=1}^{C_k^T} \left\{ \sum_{j=1}^{n_{c,k}^T} \left\{ \mathcal{F}_s^r(\mathbb{P}_T^{\mathcal{B}}(c, j), \mathbf{V}) \right\} \right\}, \quad (14)$$

Since Eq. (13) can not guarantee the best performance for Eq. (14) because Eq. (13) does not always select the component that has the highest function value (RHS of Eq. (14)). We can measure this error caused by Eq. (13) as :

$$\sum_{c=1}^{C_k^T} \left\{ \sum_{j=1}^{n_{c,k}^T} \left\{ \mathcal{F}_s(\mathbb{P}_T^{\mathcal{B}}(c, j), \mathbf{V}) - \mathcal{F}_s^r(\mathbb{P}_T^{\mathcal{B}}(c, j), \mathbf{V}) \right\} \right\}, \quad (15)$$

From Eq. (14), we can observe that the generalization performance of \mathbf{V} on all target sets relies on not only the KL divergence term but also the error caused by the component selection. In the following section, we study how component diversity in a mixture system affects performance.

D Lemma 3

In this section, we investigate whether a large number of components in a DEM can always ensure an optimal performance.

Lemma 3 Let $\{D^T(1, k), \dots, D^T(C(T, k), k)\}$ be several target sets and we use $\mathbb{P}_{(j,k)}^T$ to denote the probabilistic representation of $D^T(j, k)$. Learning a large number of components in \mathcal{V} can not always ensure an optimal performance.

Proof We assume that \mathcal{V} has learnt $u > C(T, k)$ components and several components would capture the information corresponding to a unique target data distribution. We assume that a certain target distribution is ignored by all components. According to Theorem 3, Eq. (11) of the paper, we have :

$$\begin{aligned} \sum_{j=1}^{C'} \left\{ \mathbb{E}_{\mathbb{P}_{(j,k)}^T} [\log p_{\Theta^i}(\mathbf{x})] \right\} &\geq \sum_{j=1}^{C'} \left\{ \mathcal{F}_{\text{dis}}(\mathbb{P}_{(j,k)}^T, \mathbb{P}_{\mathcal{M}_{c_1:t}}) \right. \\ &\quad + \mathbb{E}_{\mathbb{P}_{\mathcal{M}_{c_1:t}}} [\mathcal{L}_{ELBO}(\mathbf{x}; \Theta^i, \Omega^i)] - D_{JS}(\mathbb{P}_{(j,k)}^T \parallel \mathbb{P}_{\mathcal{M}_{c_1:t}}) \\ &\quad \left. - \mathcal{F}_{\text{DL}}(\mathbb{P}_{\mathcal{M}_{c_1:t}}, \mathbb{P}_{(j,k)}^T, \mathbb{P}_{\Theta^i}) \right\}, \end{aligned} \quad (16)$$

□

From Eq. (16), we observe that the model \mathcal{V} can not achieve the optimal performance due to the error caused by the JS divergence term.

E Theoretical analysis for the expansion threshold

In this section, we provide the theoretical analysis for the expansion threshold β from Eq. (13) of the paper. From **Theorem 3**, we know that encouraging the knowledge diversity would improve the generalization performance of the DEM with a minimal number of components. A large expansion threshold β in Eq.(15) of the paper can ensure knowledge diversity among components, leading however to losing some information learnt previously. According to **Theorem 3**, we have :

$$\begin{aligned}
\sum_{j=1}^{C'} \left\{ \mathbb{E}_{\mathbb{P}_{(j,k)}^T} [\log p_{\Theta^i}(\mathbf{x})] \right\} &\geq \sum_{j=1}^{C'} \left\{ \mathcal{F}_{\text{dis}}(\mathbb{P}_{(j,k)}^T, \mathbb{P}_{\mathcal{M}_{c_{1:t}}}) \right. \\
&\quad + \mathbb{E}_{\mathbb{P}_{\mathcal{M}_{c_{1:t}}}} [\mathcal{L}_{\text{ELBO}}(\mathbf{x}; \Theta^i, \Omega^i)] - D_{JS}(\mathbb{P}_{(j,k)}^T \parallel \mathbb{P}_{\mathcal{M}_{c_{1:t}}}) \\
&\quad \left. - \mathcal{F}_{\text{DL}}(\mathbb{P}_{\mathcal{M}_{c_{1:t}}}, \mathbb{P}_{(j,k)}^T, \mathbb{P}_{\Theta^i}) \right\}, \tag{17}
\end{aligned}$$

if we choose a very large expansion threshold β , the model \mathbf{V} would use fewer components while some underlying data distributions $\{\mathbb{P}_{j,1}^T, \mathbb{P}_{j,2}^T\}$ would not be captured by the components. Therefore, RHS of Eq. (17) would be increased, resulting in a degenerated performance. In contrast, a minimal expansion threshold β can allow \mathbf{V} to create more components, which would capture all underlying data distributions. However, this would also lead to an increase in the number of parameters. Therefore, a suitable trade-off for the expansion threshold β can ensure good performance with a fair number of components.

F Theoretical analysis for other VAE models

In this section, we extend the proposed theoretical framework to analyze the forgetting behaviour of the existing VAE models.

F.1 Importance weighted autoencoders

The Importance Weighted Autoencoder (IWELBO) [Burda et al. \(2015\)](#) is another VAE model employing a recognition network to generate multiple samples during the optimization to lead to better modelling of the posterior probabilities. The corresponding ELBO for sampling K' samples is defined as :

$$\mathcal{L}_{\text{ELBO}_{K'}}(\mathbf{x}; \theta, \omega) = \mathbb{E}_{\mathbf{z}_1, \dots, \mathbf{z}_{K'} \sim q(\mathbf{z}|\mathbf{x})} \left[\log \frac{1}{K'} \sum_{i=1}^{K'} \frac{p(\mathbf{x}, \mathbf{z}_i)}{q(\mathbf{z}_i|\mathbf{x})} \right] \tag{18}$$

where K' is the number of weighted samples and $K' = 1$ is equivalent to the standard ELBO.

Lemma 4 Let $p_{\theta^i}(\mathbf{x})$ be a probability density function for a single model \mathcal{V}^i updated at \mathcal{T}_i . Let \mathbb{P}_i^W denote a distribution of all visited data batches $\{\mathcal{B}_1, \dots, \mathcal{B}_i\}$ drawn from W at \mathcal{T}_i . Let $p_{\mathcal{M}_i}(\mathbf{x})$ and $p_{W^i}(\mathbf{x})$ denote the density function for $\mathbb{P}_{\mathcal{M}_i}$ and \mathbb{P}_i^W ,

respectively. We then derive an upper bound for a single VAE model trained on \mathcal{M}_i at \mathcal{T}_i as :

$$\begin{aligned} \mathbb{E}_{\mathbb{P}_i^W} [\log p_{\theta^i}(\mathbf{x})] &\geq \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\mathcal{M}_i}} \mathbb{E}_{\mathbf{z}_1, \dots, \mathbf{z}_{K'} \sim q(\mathbf{z}|\mathbf{x})} \left[\log \frac{1}{K'} \sum_{i=1}^{K'} \frac{p(\mathbf{x}, \mathbf{z}_i)}{q(\mathbf{z}_i|\mathbf{x})} \right] \\ &\quad - D_{JS}(\mathbb{P}_i^W \parallel \mathbb{P}_{\mathcal{M}_i}) - \mathcal{F}_{DL}(\mathbb{P}_{\mathcal{M}_i}, \mathbb{P}_i^W, \mathbb{P}_{\theta^i}) \\ &\quad + \mathcal{F}_{dis}(\mathbb{P}_i^W, \mathbb{P}_{\mathcal{M}_i}), \end{aligned} \quad (19)$$

Proof Since we have $\mathcal{L}_{ELBO_{K'}}(\mathbf{x}; \theta^i, \omega^i) < \mathcal{L}_{ELBO_{K'+1}}(\mathbf{x}; \theta^i, \omega^i)$, we can replace $\mathcal{L}_{ELBO}(\mathbf{x}; \theta^i, \omega^i)$ by $\mathcal{L}_{ELBO_{K'+1}}(\mathbf{x}; \theta^i, \omega^i)$ in Eq. (5), resulting in Eq. (19).

From Eq. (19), we can observe that by increasing the number of weighted samples K' can not ensure an optimal performance since RHS of Eq. (19) is also relying on the JS divergence term $D_{JS}(\mathbb{P}_i^W \parallel \mathbb{P}_{\mathcal{M}_i})$.

F.2 Hierarchical Variational Inference

In this section, we extend our theoretical analysis to the Auxiliary Deep Generative Models (ADGM) [Maaløe et al. \(2016\)](#). ADGM is a classical hierarchical latent variable model. ADGM introduces an auxiliary variable \mathbf{a} into the variational distribution $q(\mathbf{a}, \mathbf{z} | \mathbf{x}) = q(\mathbf{z} | \mathbf{a}, \mathbf{x})q(\mathbf{a} | \mathbf{x})$ and its ELBO is expressed as :

$$\begin{aligned} \log p(\mathbf{x}) &= \log \iint p(\mathbf{x}, \mathbf{a}, \mathbf{z}) d\mathbf{a} d\mathbf{z} \geq \mathbb{E}_{q(\mathbf{a}, \mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{a} | \mathbf{z}, \mathbf{x}) p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q(\mathbf{a} | \mathbf{x}) q(\mathbf{z} | \mathbf{a}, \mathbf{x})} \right] \\ &= \mathcal{L}_{ADGM}(\mathbf{x}; \theta, \omega). \end{aligned} \quad (20)$$

Then according to the results from **Lemma 4**, we have :

$$\begin{aligned} \mathbb{E}_{\mathbb{P}_i^W} [\log p_{\theta^i}(\mathbf{x})] &\geq \mathbb{E}_{\mathbb{P}_{\mathcal{M}_i}} \mathcal{L}_{ADGM}(\mathbf{x}; \theta^i, \omega^i) \\ &\quad - D_{KL}(\mathbb{P}_i^W \parallel \mathbb{P}_{\mathcal{M}_i}) - \mathcal{F}_{DL}(\mathbb{P}_{\mathcal{M}_i}, \mathbb{P}_i^W, \mathbb{P}_{\theta^i}) \\ &\quad + \mathcal{F}_{dis}(\mathbb{P}_i^W, \mathbb{P}_{\mathcal{M}_i}), \end{aligned} \quad (21)$$

Similar to **Lemma 4**, maximizing $\mathcal{L}_{ADGM}(\mathbf{x}; \theta^i, \omega^i)$ only can not ensure the optimal performance on the target distribution \mathbb{P}_i^W . These results show that deriving a tight ELBO on the source distribution can not ensure the optimal performance on the target distribution under continual learning.

G Theoretical analysis for task-known continual learning

In this section, we extend our theoretical analysis to a general continual learning where the task boundary is known during the training.

G.1 Theoretical analysis for the static network architecture

In this section, we firstly provide theoretical analysis for the static model and then extend it to the dynamic expansion model in the next section. We provide the necessary notations below.

G.1.1 Memory-based model

Definition 5. Let $\mathcal{C} = \{\mathcal{C}_1, \dots, \mathcal{C}_n\}$ be a sequence of tasks, where each task \mathcal{C}_i is associated with the training dataset D_i^S as well as the testing dataset D_i^T . Let \mathbb{P}_i^T and \mathbb{P}_i^S represent the probabilistic representation of D_i^T and D_i^S , respectively. First, we provide the theoretical analysis for a single VAE model learning a sequence of n tasks.

Definition 6. In the task-known continual learning, let \mathcal{M}_i represent the memory buffer updated at the i -th task learning (does not store the samples from the i -th task).

Lemma 5 Let $\{\mathbb{P}_1^T, \dots, \mathbb{P}_n^T\}$ be the distribution of the testing sets from a sequence of n tasks $\{\mathcal{C}_1, \dots, \mathcal{C}_n\}$. Let $p_{\theta^i}(\mathbf{x})$ be a probability density function for a single model \mathcal{V}^i updated at \mathcal{T}_i . Let $p_{\mathcal{M}_i}(\mathbf{x})$ and $p_{T^i}(\mathbf{x})$ denote the density function for $\mathbb{P}_{\mathcal{M}_i}$ and \mathbb{P}_i^T , respectively. Let $\mathbb{P}_j^S \otimes \mathbb{P}_{\mathcal{M}_j}$ be the combined distribution of D_j^S and \mathcal{M}_j . We then derive a lower bound for a single VAE model trained on \mathcal{M}_j at \mathcal{C}_j as :

$$\begin{aligned} \sum_{i=1}^n \{ \mathbb{E}_{\mathbb{P}_i^T} [\log p_{\theta^i}(\mathbf{x})] \} &\geq \sum_{i=1}^n \{ \mathbb{E}_{\mathbb{P}_{\mathcal{M}_j} \otimes \mathbb{P}_j^S} [\mathcal{L}_{ELBO}(\mathbf{x}; \theta^i, \omega^i)] \\ &\quad - D_{JS}(\mathbb{P}_i^T \parallel \mathbb{P}_{\mathcal{M}_j} \otimes \mathbb{P}_j^S) \\ &\quad - \mathcal{F}_{DL}(\mathbb{P}_{\mathcal{M}_j} \otimes \mathbb{P}_j^S, \mathbb{P}_i^T, \mathbb{P}_{\theta^i}) \\ &\quad + \mathcal{F}_{dis}(\mathbb{P}_i^T, \mathbb{P}_{\mathcal{M}_j} \otimes \mathbb{P}_j^S) \}, \end{aligned} \quad (22)$$

Eq. (22) can explain the forgetting behaviour of the existing memory-based methods [Ye and Bors \(2022a\)](#); [Egorov et al. \(2021\)](#); [Deja et al. \(2021\)](#) when learning a sequence of n tasks. As learning more tasks (j is increased), the model would accumulate more knowledge over time and can thus improve its performance. However, when the mem-

108 ory buffer \mathcal{M}_j has a fixed capacity, the model would suffer from a performance loss on
 109 past tasks since the memory buffer can not preserve all previously learnt knowledge.

110 In the following, we study the forgetting behaviour for the generative replay ap-
 111 proaches that trains a generator to replay past samples. Unlike the memory-based ap-
 112 proach, the generative replay approaches can provide an infinite number of generative
 113 replay samples over time.

114 G.1.2 GRM-based model

115 **Definition 7. (Generative Replay Mechanism (GRM)).** Let $\mathbb{P}_{\hat{\mathbf{x}}^i}$ represent the distribu-
 116 tion of the generative replay samples drawn from a single VAE model that has learnt i
 117 tasks $\{\mathcal{C}_1, \dots, \mathcal{C}_i\}$. Let $F_{\mathcal{C}}: \mathcal{X} \rightarrow \mathcal{C}$ be an optimal task-inference model that always
 118 return the true task label for a given input \mathbf{x} . By using $F_{\mathcal{C}}$, we can form an approximate
 119 distribution $\mathbb{P}_{\hat{\mathbf{x}}^{(j,i-j)}}$ through the sampling process $\mathbf{x} \sim \mathbb{P}_{\hat{\mathbf{x}}^i}$ if $F_{\mathcal{C}}(\mathbf{x}) = j$ where the
 120 superscript $i-j$ denotes that $\mathbb{P}_j^S = \mathbb{P}_{\hat{\mathbf{x}}^{(j,0)}}$ is transformed to $\mathbb{P}_{\hat{\mathbf{x}}^{(j,i-j)}}$, $j < i$ through
 121 $i-j$ GRM processes. Therefore, $\mathbb{P}_{\hat{\mathbf{x}}^{(j,i-j)}}$ is an approximate distribution of the training
 122 set of a certain task (\mathcal{C}_j), generated by a single model that has learnt i tasks.

123 Based on the definition of the GRM process, we derive new theoretical analysis for
 124 the GRM-based methods in the following.

Lemma 6 Let $\{\mathbb{P}_1^T, \dots, \mathbb{P}_n^T\}$ be the distribution of the testing sets from a sequence of
 n tasks $\{\mathcal{C}_1, \dots, \mathcal{C}_n\}$. Let $p_{\theta^i}(\mathbf{x})$ be a probability density function for a single model
 \mathcal{V}^i updated at \mathcal{T}_i . We assume that a VAE model is enabled with the GRM process to
 relieve forgetting. We then derive a lower bound for the j -th task, achieved by a single
 VAE model trained at \mathcal{C}_i as :

$$\begin{aligned} \mathbb{E}_{\mathbb{P}_j^T} [\log p_{\theta^i}(\mathbf{x})] &\geq \mathbb{E}_{\mathbb{P}_{\hat{\mathbf{x}}^{(j,i-j)}} \otimes \mathbb{P}_i^S} [\mathcal{L}_{ELBO}(\mathbf{x}; \theta^i, \omega^i)] \\ &\quad - D_{JS}(\mathbb{P}_j^T \parallel \mathbb{P}_{\hat{\mathbf{x}}^{(j,i-j)}} \otimes \mathbb{P}_i^S) \\ &\quad - \mathcal{F}_{DL}(\mathbb{P}_{\hat{\mathbf{x}}^{(j,i-j)}} \otimes \mathbb{P}_i^S, \mathbb{P}_j^T, \mathbb{P}_{\theta^i}) \\ &\quad + \mathcal{F}_{dis}(\mathbb{P}_j^T, \mathbb{P}_{\hat{\mathbf{x}}^{(j,i-j)}} \otimes \mathbb{P}_i^S), \end{aligned} \quad (23)$$

115 From Eq. (23), it observes that the quality of the approximate distribution $\mathbb{P}_{\hat{\mathbf{x}}^{(j,i-j)}}$
 116 is crucial for the performance on the j -th task, achieved by a single VAE model trained
 117 at \mathcal{C}_i . Specifically, if $\mathbb{P}_{\hat{\mathbf{x}}^{(j,i-j)}}$ can approximate the corresponding data distribution
 118 \mathbb{P}_j^T exactly, RHS of Eq. (23) would be large, resulting in better performance. As the
 119 number of tasks to be increased (i is increased), the model would suffer from more for-

getting because the approximate distribution $\mathbb{P}_{\hat{\mathbf{x}}(j,i-j)}$ would far away from \mathbb{P}_j^T , caused by the frequent GRM process ($i - j$ is large). In the following, we derive a new lower bound that describes how a single VAE model forgets its previously learnt knowledge in each task learning.

Lemma 7 Let $\{\mathbb{P}_1^T, \dots, \mathbb{P}_n^T\}$ be the distribution of the testing sets from a sequence of n tasks $\{\mathcal{C}_1, \dots, \mathcal{C}_n\}$. Let $p_{\theta^i}(\mathbf{x})$ be a probability density function for a single model \mathcal{V}^i updated at \mathcal{T}_i . We assume that a VAE model is enabled with the GRM process to relieve forgetting. We then derive a lower bound for the j -th task, achieved by a single VAE model trained at \mathcal{C}_i as :

$$\begin{aligned} \mathbb{E}_{\mathbb{P}_j^T} [\log p_{\theta^i}(\mathbf{x})] &\geq \mathbb{E}_{\mathbb{P}_{\hat{\mathbf{x}}(j,i-j)} \otimes \mathbb{P}_i^S} [\mathcal{L}_{ELBO}(\mathbf{x}; \theta^i, \omega^i)] - D_{JS}(\mathbb{P}_j^T \parallel \mathbb{P}_{\hat{\mathbf{x}}(j,0)} \otimes \mathbb{P}_i^S) \\ &- \mathcal{F}_{DL}(\mathbb{P}_{\hat{\mathbf{x}}(j,0)} \otimes \mathbb{P}_i^S, \mathbb{P}_j^T, \mathbb{P}_{\theta^i}) + \mathcal{F}_{dis}(\mathbb{P}_j^T, \mathbb{P}_{\hat{\mathbf{x}}(j,0)} \otimes \mathbb{P}_i^S) \\ &- \sum_{c=2}^{i-j} \left\{ D_{JS}(\mathbb{P}_{\hat{\mathbf{x}}(j,c-1)} \parallel \mathbb{P}_{\hat{\mathbf{x}}(j,c)} \otimes \mathbb{P}_i^S) \right. \\ &\left. + \mathcal{F}_{DL}(\mathbb{P}_{\hat{\mathbf{x}}(j,c)} \otimes \mathbb{P}_i^S, \mathbb{P}_{\hat{\mathbf{x}}(j,c-1)}, \mathbb{P}_{\theta^i}) - \mathcal{F}_{dis}(\mathbb{P}_{\hat{\mathbf{x}}(j,c-1)}, \mathbb{P}_{\hat{\mathbf{x}}(j,c)} \otimes \mathbb{P}_i^S) \right\}, \end{aligned} \quad (24)$$

Remark. We have several observations from Lemma 7 :

- As the number of tasks increases (i is increased), the performance on the j -th task, achieved by a single model, would be degenerated caused by the accumulated JS divergence terms (the fifth term in RHS of Eq. (24)).
- Learning early tasks would suffer from more forgetting than the learning of the recent tasks because early tasks cause more accumulated JS divergence terms ($i - j$ is large in RHS of Eq. (24)).

Proof. First, we consider \mathbb{P}_j^T and $\mathbb{P}_{\hat{\mathbf{x}}(j,1)}$ as the target and source distribution, respectively. Based on Lemma 6, we have :

$$\begin{aligned} \mathbb{E}_{\mathbb{P}_j^T} [\log p_{\theta^i}(\mathbf{x})] &\geq \mathbb{E}_{\mathbb{P}_{\hat{\mathbf{x}}(j,0)} \otimes \mathbb{P}_i^S} [\mathcal{L}_{ELBO}(\mathbf{x}; \theta^i, \omega^i)] \\ &- D_{JS}(\mathbb{P}_j^T \parallel \mathbb{P}_{\hat{\mathbf{x}}(j,0)} \otimes \mathbb{P}_i^S) \\ &- \mathcal{F}_{DL}(\mathbb{P}_{\hat{\mathbf{x}}(j,0)} \otimes \mathbb{P}_i^S, \mathbb{P}_j^T, \mathbb{P}_{\theta^i}) \\ &+ \mathcal{F}_{dis}(\mathbb{P}_j^T, \mathbb{P}_{\hat{\mathbf{x}}(j,0)} \otimes \mathbb{P}_i^S), \end{aligned} \quad (25)$$

Then, we consider $\mathbb{P}_{\hat{\mathbf{x}}(j,0)}$ and $\mathbb{P}_{\hat{\mathbf{x}}(j,1)}$ as the target and source data distribution, respectively. This setting is reasonable since $\mathbb{P}_{\hat{\mathbf{x}}(j,0)}$ is more closed to the distribution \mathbb{P}_j^T of the testing set compared with $\mathbb{P}_{\hat{\mathbf{x}}(j,1)}$. We can derive a lower bound between

$\mathbb{P}_{\hat{\mathbf{x}}^{(j,0)}}$ and $\mathbb{P}_{\hat{\mathbf{x}}^{(j,1)}}$.

$$\begin{aligned}
\mathbb{E}_{\mathbb{P}_{\hat{\mathbf{x}}^{(j,0)}}} [\log p_{\theta^i}(\mathbf{x})] &\geq \mathbb{E}_{\mathbb{P}_{\hat{\mathbf{x}}^{(j,1)}} \otimes \mathbb{P}_i^S} [\mathcal{L}_{ELBO}(\mathbf{x}; \theta^i, \omega^i)] \\
&\quad - D_{JS}(\mathbb{P}_{\hat{\mathbf{x}}^{(j,0)}} \parallel \mathbb{P}_{\hat{\mathbf{x}}^{(j,1)}} \otimes \mathbb{P}_i^S) \\
&\quad - \mathcal{F}_{DL}(\mathbb{P}_{\hat{\mathbf{x}}^{(j,1)}} \otimes \mathbb{P}_i^S, \mathbb{P}_{\hat{\mathbf{x}}^{(j,0)}}, \mathbb{P}_{\theta^i}) \\
&\quad + \mathcal{F}_{dis}(\mathbb{P}_{\hat{\mathbf{x}}^{(j,0)}}, \mathbb{P}_{\hat{\mathbf{x}}^{(j,1)}} \otimes \mathbb{P}_i^S), \tag{26}
\end{aligned}$$

Then we treat $\mathbb{P}_{\hat{\mathbf{x}}^{(j,1)}}$ and $\mathbb{P}_{\hat{\mathbf{x}}^{(j,2)}}$ as the target and source distribution, respectively. We have :

$$\begin{aligned}
\mathbb{E}_{\mathbb{P}_{\hat{\mathbf{x}}^{(j,1)}}} [\log p_{\theta^i}(\mathbf{x})] &\geq \mathbb{E}_{\mathbb{P}_{\hat{\mathbf{x}}^{(j,2)}} \otimes \mathbb{P}_i^S} [\mathcal{L}_{ELBO}(\mathbf{x}; \theta^i, \omega^i)] \\
&\quad - D_{JS}(\mathbb{P}_{\hat{\mathbf{x}}^{(j,1)}} \parallel \mathbb{P}_{\hat{\mathbf{x}}^{(j,2)}} \otimes \mathbb{P}_i^S) \\
&\quad - \mathcal{F}_{DL}(\mathbb{P}_{\hat{\mathbf{x}}^{(j,2)}} \otimes \mathbb{P}_i^S, \mathbb{P}_{\hat{\mathbf{x}}^{(j,1)}}, \mathbb{P}_{\theta^i}) \\
&\quad + \mathcal{F}_{dis}(\mathbb{P}_{\hat{\mathbf{x}}^{(j,1)}}, \mathbb{P}_{\hat{\mathbf{x}}^{(j,2)}} \otimes \mathbb{P}_i^S), \tag{27}
\end{aligned}$$

131 According to the summary, we can have the following bounds :

$$\begin{aligned}
&\mathbb{E}_{\mathbb{P}_{\hat{\mathbf{x}}^{(j,2)}}} [\log p_{\theta^i}(\mathbf{x})] \geq \mathbb{E}_{\mathbb{P}_{\hat{\mathbf{x}}^{(j,3)}} \otimes \mathbb{P}_i^S} [\mathcal{L}_{ELBO}(\mathbf{x}; \theta^i, \omega^i)] \\
&\quad - D_{JS}(\mathbb{P}_{\hat{\mathbf{x}}^{(j,2)}} \parallel \mathbb{P}_{\hat{\mathbf{x}}^{(j,3)}} \otimes \mathbb{P}_i^S) \\
&\quad - \mathcal{F}_{DL}(\mathbb{P}_{\hat{\mathbf{x}}^{(j,3)}} \otimes \mathbb{P}_i^S, \mathbb{P}_{\hat{\mathbf{x}}^{(j,2)}}, \mathbb{P}_{\theta^i}) \\
&\quad + \mathcal{F}_{dis}(\mathbb{P}_{\hat{\mathbf{x}}^{(j,2)}}, \mathbb{P}_{\hat{\mathbf{x}}^{(j,3)}} \otimes \mathbb{P}_i^S), \\
&\dots \\
&\mathbb{E}_{\mathbb{P}_{\hat{\mathbf{x}}^{(j,i-j-1)}}} [\log p_{\theta^i}(\mathbf{x})] \geq \mathbb{E}_{\mathbb{P}_{\hat{\mathbf{x}}^{(j,i-j)}} \otimes \mathbb{P}_i^S} [\mathcal{L}_{ELBO}(\mathbf{x}; \theta^i, \omega^i)] \\
&\quad - D_{JS}(\mathbb{P}_{\hat{\mathbf{x}}^{(j,i-j-1)}} \parallel \mathbb{P}_{\hat{\mathbf{x}}^{(j,i-j)}} \otimes \mathbb{P}_i^S) \\
&\quad - \mathcal{F}_{DL}(\mathbb{P}_{\hat{\mathbf{x}}^{(j,i-j)}} \otimes \mathbb{P}_i^S, \mathbb{P}_{\hat{\mathbf{x}}^{(j,i-j-1)}}, \mathbb{P}_{\theta^i}) \\
&\quad + \mathcal{F}_{dis}(\mathbb{P}_{\hat{\mathbf{x}}^{(j,i-j-1)}}, \mathbb{P}_{\hat{\mathbf{x}}^{(j,i-j)}} \otimes \mathbb{P}_i^S),
\end{aligned}$$

In the final, we sum up the above bounds, resulting in :

$$\begin{aligned}
\mathbb{E}_{\mathbb{P}_j^T} [\log p_{\theta^i}(\mathbf{x})] &\geq \mathbb{E}_{\mathbb{P}_{\hat{\mathbf{x}}^{(j,i-j)}} \otimes \mathbb{P}_i^S} [\mathcal{L}_{ELBO}(\mathbf{x}; \theta^i, \omega^i)] \\
&\quad - D_{JS}(\mathbb{P}_j^T \parallel \mathbb{P}_{\hat{\mathbf{x}}^{(j,0)}} \otimes \mathbb{P}_i^S) \\
&\quad - \mathcal{F}_{DL}(\mathbb{P}_{\hat{\mathbf{x}}^{(j,0)}} \otimes \mathbb{P}_i^S, \mathbb{P}_j^T, \mathbb{P}_{\theta^i}) + \mathcal{F}_{dis}(\mathbb{P}_j^T, \mathbb{P}_{\hat{\mathbf{x}}^{(j,0)}} \otimes \mathbb{P}_i^S) \\
&\quad - \sum_{c=2}^{i-j} \{D_{JS}(\mathbb{P}_{\hat{\mathbf{x}}^{(j,c-1)}} \parallel \mathbb{P}_{\hat{\mathbf{x}}^{(j,c)}} \otimes \mathbb{P}_i^S) \\
&\quad - \mathcal{F}_{DL}(\mathbb{P}_{\hat{\mathbf{x}}^{(j,c)}} \otimes \mathbb{P}_i^S, \mathbb{P}_{\hat{\mathbf{x}}^{(j,c-1)}}, \mathbb{P}_{\theta^i}) \\
&\quad + \mathcal{F}_{dis}(\mathbb{P}_{\hat{\mathbf{x}}^{(j,c-1)}}, \mathbb{P}_{\hat{\mathbf{x}}^{(j,c)}} \otimes \mathbb{P}_i^S)\}, \tag{28}
\end{aligned}$$

□

Based on Lemma 7, we can derive a lower bound for all tasks in the following.

Lemma 8 Let $\{\mathbb{P}_1^T, \dots, \mathbb{P}_n^T\}$ be the distribution of the testing sets from a sequence of n tasks $\{\mathcal{C}_1, \dots, \mathcal{C}_n\}$. Let $p_{\theta^i}(\mathbf{x})$ be a probability density function for a single model \mathcal{V}^i updated at \mathcal{T}_i . We assume that a VAE model is enabled with the GRM process to relieve forgetting. We then derive a lower bound for all tasks, achieved by a single VAE model trained at \mathcal{C}_i as :

$$\begin{aligned}
\sum_{j=1}^i \left\{ \mathbb{E}_{\mathbb{P}_j^T} [\log p_{\theta_j}(\mathbf{x})] \right\} &\geq \sum_{j=1}^i \left\{ \mathbb{E}_{\mathbb{P}_{\hat{\mathbf{x}}^{(j,i-j)}} \otimes \mathbb{P}_i^S} [\mathcal{L}_{ELBO}(\mathbf{x}; \theta^i, \omega^i)] \right. \\
&\quad - D_{JS}(\mathbb{P}_j^T \parallel \mathbb{P}_{\hat{\mathbf{x}}^{(j,0)}} \otimes \mathbb{P}_i^S) - \mathcal{F}_{DL}(\mathbb{P}_{\hat{\mathbf{x}}^{(j,0)}} \otimes \mathbb{P}_i^S, \mathbb{P}_j^T, \mathbb{P}_{\theta^i}) \\
&\quad + \mathcal{F}_{dis}(\mathbb{P}_j^T, \mathbb{P}_{\hat{\mathbf{x}}^{(j,0)}} \otimes \mathbb{P}_i^S) - \sum_{c=2}^{i-j} \left\{ D_{JS}(\mathbb{P}_{\hat{\mathbf{x}}^{(j,c-1)}} \parallel \mathbb{P}_{\hat{\mathbf{x}}^{(j,c)}} \otimes \mathbb{P}_i^S) \right. \\
&\quad \left. \left. + \mathcal{F}_{DL}(\mathbb{P}_{\hat{\mathbf{x}}^{(j,c)}} \otimes \mathbb{P}_i^S, \mathbb{P}_{\hat{\mathbf{x}}^{(j,c-1)}}, \mathbb{P}_{\theta^i}) - \mathcal{F}_{dis}(\mathbb{P}_{\hat{\mathbf{x}}^{(j,c-1)}}, \mathbb{P}_{\hat{\mathbf{x}}^{(j,c)}} \otimes \mathbb{P}_i^S) \right\} \right\}, \tag{29}
\end{aligned}$$

From Eq. (29), a robust generative replay network plays a vital role in the GRM-based model's performance. If the generative replay network can approximate each target data distribution precisely, the RHS of Eq. (29) would be increased, and thus, the model can achieve optimal performance on all target distributions. However, In practice, the generator such as VAE or GAN, can not produce reasonable generative replay samples when learning a sequence of different data domains. In the following section, we extend the proposed theoretical analysis to the dynamic expansion model that can overcome the limitations of the static models.

142 G.2 Theoretical analysis for the dynamic expansion model

143 In this section, we study the forgetting behaviour of the dynamic expansion model
 144 in a general continual learning setting where the task boundary is known during the
 145 training.

Lemma 9 Let $\mathbf{V} = \{\mathcal{V}_1, \dots, \mathcal{V}_t\}$ be a mixture model with t components. We consider
 an idea solution that the number of components is equal to the number of tasks ($t = b$).
 We then derive a lower bound for \mathbf{V} as :

$$\begin{aligned} \sum_{j=1}^n \left\{ \mathbb{E}_{\mathbb{P}_j^T} [\log p_{\theta_j}(\mathbf{x})] \right\} &\geq \sum_{j=1}^n \left\{ \mathbb{E}_{\mathbb{P}_j^S} [\mathcal{L}_{ELBO}(\mathbf{x}; \theta^i, \omega^i)] - D_{JS}(\mathbb{P}_j^T || \mathbb{P}_j^S) \right. \\ &\quad \left. - \mathcal{F}_{DL}(\mathbb{P}_j^S, \mathbb{P}_j^T, \mathbb{P}_{\theta_j}) + \mathcal{F}_{dis}(\mathbb{P}_j^T, \mathbb{P}_j^S) \right\}, \end{aligned} \quad (30)$$

146 From Eq. (30), it finds that the mixture model \mathbf{V} can achieve optimal performance
 147 on all tasks since it does not cause forgetting during the training. In practice, the
 148 mixture model \mathbf{V} would dynamically change its network architecture according to the
 149 complexity of the tasks. In the following, we derive new lower bound for \mathbf{V} that can
 150 dynamically change the number of components during the training.

151 **Lemma 10** Let $B = \{b_1, \dots, b_j\}$ denote a set of labels where each one b_i represents the
 152 task that is only trained once. The corresponding distribution for b_i is defined as $\mathbb{P}_{\mathbf{x}^{(i,0)}}$.
 153 Let $C = \{c_1, \dots, c_j\}$ be a set where each c_i denotes the index of the component that
 154 learns the b_i -th task. Let $B' = \{b'_1, \dots, b'_n\}$ represent a set of labels where each one b'_j
 155 indicates the b'_j -th task was used for re-training more than once. Let $C' = \{c'_1, \dots, c'_n\}$
 156 represent a set where each c'_i denotes that the index of the component that learns the
 157 b'_i -th task. We also define a set $\hat{B} = \{\hat{b}_1, \dots, \hat{b}_n\}$ where each one $\hat{b}_i > 1$ denotes
 158 that the b'_i -th task has been learnt for \hat{b}_i times $\mathbb{P}_{\mathbf{x}^{(b'_i,0)}} \rightarrow \mathbb{P}_{\mathbf{x}^{(b'_i, \hat{b}_i-1)}}$ where $\mathbb{P}_{\mathbf{x}^{(b'_i, \hat{b}_i-1)}}$
 159 represents the corresponding probabilistic representation. We derive a lower bound for
 160 \mathbf{V} at the t -th task learning :

$$\begin{aligned}
& \sum_{j=1}^{|B|} \left\{ \mathbb{E}_{\mathbb{P}_{b_j}^T} [\log p_{\theta_{c_j}}(\mathbf{x})] \right\} + \sum_{j=1}^{|B'|} \left\{ \mathbb{E}_{\mathbb{P}_{b'_j}^T} [\log p_{\theta_{c'_j}}(\mathbf{x})] \right\} \geq \\
& \sum_{j=1}^{|B|} \left\{ \mathbb{E}_{\mathbb{P}_{b_j}^S} [\mathcal{L}_{ELBO}(\mathbf{x}; \theta_{c_j}, \omega_{c_j})] - D_{JS}(\mathbb{P}_{b_j}^T \parallel \mathbb{P}_{c_j}^S) \right. \\
& \quad \left. - \mathcal{F}_{DL}(\mathbb{P}_{b_j}^S, \mathbb{P}_{b_j}^T, \mathbb{P}_{\theta_{c_j}}) + \mathcal{F}_{dis}(\mathbb{P}_{b_j}^T, \mathbb{P}_{b_j}^S) \right\} \\
& + \sum_{j=1}^{|B'|} \left\{ \mathbb{E}_{\mathbb{P}_{\hat{\mathbf{x}}^{(j, \hat{b}_j)}} \otimes \mathbb{P}_{c'_j}^S} [\mathcal{L}_{ELBO}(\mathbf{x}; \theta_{c'_j}, \omega_{c'_j})] - D_{JS}(\mathbb{P}_{b'_j}^T \parallel \mathbb{P}_{\hat{\mathbf{x}}^{(j, 0)}} \otimes \mathbb{P}_{c'_j}^S) \right. \\
& \quad \left. - \mathcal{F}_{DL}(\mathbb{P}_{\hat{\mathbf{x}}^{(j, 0)}} \otimes \mathbb{P}_{c'_j}^S, \mathbb{P}_{b'_j}^T, \mathbb{P}_{\theta_{c'_j}}) + \mathcal{F}_{dis}(\mathbb{P}_{b'_j}^T, \mathbb{P}_{\hat{\mathbf{x}}^{(j, 0)}} \otimes \mathbb{P}_{c'_j}^S) \right. \\
& \quad \left. - \sum_{c=2}^{\hat{b}_j} \left\{ D_{JS}(\mathbb{P}_{\hat{\mathbf{x}}^{(j, c-1)}} \parallel \mathbb{P}_{\hat{\mathbf{x}}^{(j, c)}} \otimes \mathbb{P}_{b'_j}^S) \right. \right. \\
& \quad \left. \left. + \mathcal{F}_{DL}(\mathbb{P}_{\hat{\mathbf{x}}^{(j, c)}} \otimes \mathbb{P}_{b'_j}^S, \mathbb{P}_{\hat{\mathbf{x}}^{(j, c-1)}}, \mathbb{P}_{\theta_{c'_j}}) - \mathcal{F}_{dis}(\mathbb{P}_{\hat{\mathbf{x}}^{(j, c-1)}}, \mathbb{P}_{\hat{\mathbf{x}}^{(j, c)}} \otimes \mathbb{P}_{b'_j}^S) \right\} \right\}, \quad (31)
\end{aligned}$$

Remark From Eq. (31), we have several observations :

- Increasing the number of components in \mathbf{V} can improve the performance. Specifically, when the number of components matches the number of tasks, Eq. (31) is transformed to Eq. (30) which does not have forgetting.
- Reducing the number of components would lead to degenerated performance because the accumulated errors is increased in RHS of Eq. (31). $|B'| = 1$ indicates that \mathbf{V} only learns a single component and would suffer from more forgetting due to the more accumulated errors.

Proof. Firstly, let we consider the tasks that are only trained once. We have :

$$\begin{aligned}
& \sum_{j=1}^{|B|} \left\{ \mathbb{E}_{\mathbb{P}_{b_j}^T} [\log p_{\theta_{c_j}}(\mathbf{x})] \right\} \geq \sum_{j=1}^{|B|} \left\{ \mathbb{E}_{\mathbb{P}_{b_j}^S} [\mathcal{L}_{ELBO}(\mathbf{x}; \theta_{c_j}, \omega_{c_j})] - D_{JS}(\mathbb{P}_{b_j}^T \parallel \mathbb{P}_{c_j}^S) \right. \\
& \quad \left. - \mathcal{F}_{DL}(\mathbb{P}_{b_j}^S, \mathbb{P}_{b_j}^T, \mathbb{P}_{\theta_{c_j}}) + \mathcal{F}_{dis}(\mathbb{P}_{b_j}^T, \mathbb{P}_{b_j}^S) \right\}, \quad (32)
\end{aligned}$$

Secondly, we consider the tasks that are trained more than once. We have :

$$\begin{aligned}
& \sum_{j=1}^{|B'|} \left\{ \mathbb{E}_{\mathbb{P}_{b'_j}^T} [\log p_{\theta_{c'_j}}(\mathbf{x})] \right\} \geq \sum_{j=1}^{|B'|} \left\{ \mathbb{E}_{\mathbb{P}_{\hat{\mathbf{x}}^{(j, \hat{b}_j)}}} \otimes \mathbb{P}_{c'_j}^S [\mathcal{L}_{ELBO}(\mathbf{x}; \theta_{c'_j}, \omega_{c'_j})] \right. \\
& - D_{JS}(\mathbb{P}_{b'_j}^T, \mathbb{P}_{\hat{\mathbf{x}}^{(j, 0)}} \otimes \mathbb{P}_{c'_j}^S) - \mathcal{F}_{DL}(\mathbb{P}_{\hat{\mathbf{x}}^{(j, 0)}} \otimes \mathbb{P}_{c'_j}^S, \mathbb{P}_{b'_j}^T, \mathbb{P}_{\theta_{c'_j}}) \\
& + \mathcal{F}_{dis}(\mathbb{P}_{b'_j}^T, \mathbb{P}_{\hat{\mathbf{x}}^{(j, 0)}} \otimes \mathbb{P}_{b'_j}^S) - \sum_{c=2}^{\hat{b}_j} \left\{ D_{JS}(\mathbb{P}_{\hat{\mathbf{x}}^{(j, c-1)}} \parallel \mathbb{P}_{\hat{\mathbf{x}}^{(j, c)}} \otimes \mathbb{P}_{b'_j}^S) \right. \\
& \left. + \mathcal{F}_{DL}(\mathbb{P}_{\hat{\mathbf{x}}^{(j, c)}} \otimes \mathbb{P}_{b'_j}^S, \mathbb{P}_{\hat{\mathbf{x}}^{(j, c-1)}}, \mathbb{P}_{\theta_{c'_j}}) - \mathcal{F}_{dis}(\mathbb{P}_{\hat{\mathbf{x}}^{(j, c-1)}}, \mathbb{P}_{\hat{\mathbf{x}}^{(j, c)}} \otimes \mathbb{P}_{b'_j}^S) \right\} \Big\}, \quad (33)
\end{aligned}$$

Then we sum up Eq. (32) and Eq. (33), resulting in :

$$\begin{aligned}
& \sum_{j=1}^{|B|} \left\{ \mathbb{E}_{\mathbb{P}_{b_j}^T} [\log p_{\theta_{c_j}}(\mathbf{x})] \right\} + \sum_{j=1}^{|B'|} \left\{ \mathbb{E}_{\mathbb{P}_{b'_j}^T} [\log p_{\theta_{c'_j}}(\mathbf{x})] \right\} \geq \\
& \sum_{j=1}^{|B|} \left\{ \mathbb{E}_{\mathbb{P}_{b_j}^S} [\mathcal{L}_{ELBO}(\mathbf{x}; \theta_{c_j}, \omega_{c_j})] - D_{JS}(\mathbb{P}_{b_j}^T \parallel \mathbb{P}_{c_j}^S) \right. \\
& - \mathcal{F}_{DL}(\mathbb{P}_{b_j}^S, \mathbb{P}_{b_j}^T, \mathbb{P}_{\theta_{c_j}}) + \mathcal{F}_{dis}(\mathbb{P}_{b_j}^T, \mathbb{P}_{b_j}^S) \Big\} \\
& + \sum_{j=1}^{|B'|} \left\{ \mathbb{E}_{\mathbb{P}_{\hat{\mathbf{x}}^{(j, \hat{b}_j)}}} \otimes \mathbb{P}_{c'_j}^S [\mathcal{L}_{ELBO}(\mathbf{x}; \theta_{c'_j}, \omega_{c'_j})] - D_{JS}(\mathbb{P}_{b'_j}^T \parallel \mathbb{P}_{\hat{\mathbf{x}}^{(j, 0)}} \otimes \mathbb{P}_{c'_j}^S) \right. \\
& - \mathcal{F}_{DL}(\mathbb{P}_{\hat{\mathbf{x}}^{(j, 0)}} \otimes \mathbb{P}_{c'_j}^S, \mathbb{P}_{b'_j}^T, \mathbb{P}_{\theta_{c'_j}}) + \mathcal{F}_{dis}(\mathbb{P}_{b'_j}^T, \mathbb{P}_{\hat{\mathbf{x}}^{(j, 0)}} \otimes \mathbb{P}_{b'_j}^S) \\
& - \sum_{c=2}^{\hat{b}_j} \left\{ D_{JS}(\mathbb{P}_{\hat{\mathbf{x}}^{(j, c-1)}} \parallel \mathbb{P}_{\hat{\mathbf{x}}^{(j, c)}} \otimes \mathbb{P}_{b'_j}^S) \right. \\
& \left. + \mathcal{F}_{DL}(\mathbb{P}_{\hat{\mathbf{x}}^{(j, c)}} \otimes \mathbb{P}_{b'_j}^S, \mathbb{P}_{\hat{\mathbf{x}}^{(j, c-1)}}, \mathbb{P}_{\theta_{c'_j}}) - \mathcal{F}_{dis}(\mathbb{P}_{\hat{\mathbf{x}}^{(j, c-1)}}, \mathbb{P}_{\hat{\mathbf{x}}^{(j, c)}} \otimes \mathbb{P}_{b'_j}^S) \right\} \Big\}, \quad (34)
\end{aligned}$$

169 \square

170 **G.3 Theoretical analysis for the existing GRM-based models**

171 In this section, we employ the proposed theoretical framework to analyze the forgetting
172 behaviour of the existing GRM-based models in a general continual learning setting
173 where task label is given during the training.

G.3.1 Lifelong VAEGAN

174

Lifelong VAEGAN [Ye and Bors \(2020a\)](#) is a popular GRM-based model. Unlike other GRM-based approaches that lack an inference mechanism, Lifelong VAEGAN introduces to combine the powerful inference mechanism of VAE and the robust generation capacity of GAN into a unified framework. Lifelong VAEGAN produces not only high-quality generative replay samples but also captures meaningful latent representations across domains over time. In this section, we extend our theoretical analysis to Lifelong VAEGAN. According to **Lemma 7**, we can derive a lower bound for Lifelong VAEGAN as :

$$\begin{aligned} \mathbb{E}_{\mathbb{P}_j^T} [\log p_{\theta^i}(\mathbf{x})] &\geq \mathbb{E}_{\mathbb{P}_{\mathbf{x}(j,i-j)} \otimes \mathbb{P}_i^S} [\mathcal{L}_{ELBO}(\mathbf{x}; \theta^i, \omega^i)] - D_{JS}(\mathbb{P}_j^T \parallel \mathbb{P}_{\mathbf{x}(j,0)} \otimes \mathbb{P}_i^S) \\ &\quad - \mathcal{F}_{DL}(\mathbb{P}_{\mathbf{x}(j,0)} \otimes \mathbb{P}_i^S, \mathbb{P}_j^T, \mathbb{P}_{\theta^i}) + \mathcal{F}_{dis}(\mathbb{P}_j^T, \mathbb{P}_{\mathbf{x}(j,0)} \otimes \mathbb{P}_i^S) \\ &\quad - \sum_{c=2}^{i-j} \left\{ D_{JS}(\mathbb{P}_{\mathbf{x}(j,c-1)} \parallel \mathbb{P}_{\mathbf{x}(j,c)} \otimes \mathbb{P}_i^S) \right. \\ &\quad \left. + \mathcal{F}_{DL}(\mathbb{P}_{\mathbf{x}(j,c)} \otimes \mathbb{P}_i^S, \mathbb{P}_{\mathbf{x}(j,c-1)}, \mathbb{P}_{\theta^i}) - \mathcal{F}_{dis}(\mathbb{P}_{\mathbf{x}(j,c-1)}, \mathbb{P}_{\mathbf{x}(j,c)} \otimes \mathbb{P}_i^S) \right\}, \end{aligned} \quad (35)$$

Eq. (35) describes the forgetting behaviour of Lifelong VAEGAN on the j -th task when it is trained on the i -th task where $i > j$. Since Lifelong VAEGAN employs GAN as the generative replay network, it can provide high-quality generative replay samples that approximate the target distribution exactly. Furthermore, compared with the VAE-based models, Lifelong VAEGAN can have the small JS divergence terms in RHS of Eq. (35) and thus can achieve better performance.

G.3.2 Lifelong infinite mixture model

181

Lifelong Infinite Mixture (LIMix) [Ye and Bors \(2021a\)](#) is a popular dynamic expansion model which can dynamically change its network architectures to deal with new tasks. LIMix also introduces to evaluate the knowledge similarity between each trained component and a new task, which aims to reuse an appropriate component to learn several similar tasks. According to **Lemma 10**, the forgetting behaviour of LIMix can be

described by :

$$\begin{aligned}
& \sum_{j=1}^{|B|} \left\{ \mathbb{E}_{\mathbb{P}_{b_j}^T} [\log p_{\theta_{c_j}}(\mathbf{x})] \right\} + \sum_{j=1}^{|B'|} \left\{ \mathbb{E}_{\mathbb{P}_{b'_j}^T} [\log p_{\theta_{c'_j}}(\mathbf{x})] \right\} \geq \\
& \sum_{j=1}^{|B|} \left\{ \mathbb{E}_{\mathbb{P}_{b_j}^S} [\mathcal{L}_{ELBO}(\mathbf{x}; \theta_{c_j}, \omega_{c_j})] - D_{JS}(\mathbb{P}_{b_j}^T || \mathbb{P}_{c_j}^S) \right. \\
& \quad \left. - \mathcal{F}_{DL}(\mathbb{P}_{b_j}^S, \mathbb{P}_{b_j}^T, \mathbb{P}_{\theta_{c_j}}) + \mathcal{F}_{dis}(\mathbb{P}_{b_j}^T, \mathbb{P}_{b_j}^S) \right\} \\
& + \sum_{j=1}^{|B'|} \left\{ \mathbb{E}_{\mathbb{P}_{\hat{\mathbf{x}}^{(j, \hat{b}_j)}} \otimes \mathbb{P}_{c'_j}^S} [\mathcal{L}_{ELBO}(\mathbf{x}; \theta_{c'_j}, \omega_{c'_j})] - D_{JS}(\mathbb{P}_{b'_j}^T || \mathbb{P}_{\hat{\mathbf{x}}^{(j, 0)}} \otimes \mathbb{P}_{c'_j}^S) \right. \\
& \quad \left. - \mathcal{F}_{DL}(\mathbb{P}_{\hat{\mathbf{x}}^{(j, 0)}} \otimes \mathbb{P}_{c'_j}^S, \mathbb{P}_{b'_j}^T, \mathbb{P}_{\theta_{c'_j}}) + \mathcal{F}_{dis}(\mathbb{P}_{b'_j}^T, \mathbb{P}_{\hat{\mathbf{x}}^{(j, 0)}} \otimes \mathbb{P}_{b'_j}^S) \right. \\
& \quad \left. - \sum_{c=2}^{\hat{b}_j} \left\{ D_{JS}(\mathbb{P}_{\hat{\mathbf{x}}^{(j, c-1)}} || \mathbb{P}_{\hat{\mathbf{x}}^{(j, c)}} \otimes \mathbb{P}_{b'_j}^S) \right. \right. \\
& \quad \left. \left. + \mathcal{F}_{DL}(\mathbb{P}_{\hat{\mathbf{x}}^{(j, c)}} \otimes \mathbb{P}_{b'_j}^S, \mathbb{P}_{\hat{\mathbf{x}}^{(j, c-1)}}, \mathbb{P}_{\theta_{c'_j}}) - \mathcal{F}_{dis}(\mathbb{P}_{\hat{\mathbf{x}}^{(j, c-1)}}, \mathbb{P}_{\hat{\mathbf{x}}^{(j, c)}} \otimes \mathbb{P}_{b'_j}^S) \right\} \right\}, \quad (36)
\end{aligned}$$

From Eq. (36), it observes that increasing the number of components in LIMix can significantly improve the generalization performance since the JS divergence terms in the RHS of Eq. (36) are reduced. In addition, if a particular component is reused to learn several similar tasks, the JS divergence terms caused by the GRM process would be small, resulting in better relieving forgetting and reducing the total number of components. However, LIMix can not be applied in TFCL since it requires access to the task labels in order to perform the component selection and expansion strategy.

H Additional information for experiments

H.1 Additional information for experiment settings

The release of the code. We have provided the detailed implementation of the proposed Online Adversarial Expansion Strategy (OAES) model. We also provide the source code in the supplemental material. In addition, We will organize the source code of the OAES model for the sake of easy understanding and for facilitating the re-implementation and we will release it publicly on <https://github.com/> if the paper is accepted.

H.2 Experiment setting

The hyperparameter configuration and GPU hardware. To perform the density estimation task, we use Adam [Kingma and Ba \(2015\)](#) with a learning rate of 0.0001 and its default hyperparameters. We set the batch size and the number of epochs for each training step as 64 and 1, respectively. Following from [Ye and Bors \(2022a,b,c, 2021b,c\)](#); [Aljundi et al. \(2019a\)](#); [Ye and Bors \(2022d, 2021d, 2020a,b, 2022e, 2021e,a, 2020c\)](#), the experiments are constructed in the computer with the operating system (Ubuntu 18.04.5). We also use the GPU (GeForce GTX 1080) for all our experiments.

The configuration of the network architecture for log-likelihood estimation task. We adapt the network architecture from [Burda et al. \(2015\)](#) where two fully connected layers implement the inference and generator models. Each layer has 200 hidden units. The shared modules use the expansion mechanism as a single fully-connected neural network with a layer (200 hidden units). A single layer also implements each individual component with 200 hidden units for both the generator and inference models.

Hyperparameter setting. The batch size is of 64 images, and we consider 1 epochs for each training stage. The maximum memory size for Split MNIST, Split Fashion, Split MNIST-Fashion, Cross-domain is 1.5K, 1.5K, 1.9K and 2.0K, respectively.

Additional information for the evaluation. All results reported in the paper are evaluated on the testing datasets after task-free continual learning.

H.3 The configuration for the classification task.

First, we introduce the details about the datasets used in our classification task as follows. The threshold β for Split MNIST, Split CIFAR10, Split CIFAR100 and Split MiniImageNe is 4.2, 4 and 4.5 and 4.8 respectively. The number of components of OAES for Split MNIST, Split CIFAR10, Split CIFAR100 and Split MiniImageNe is 6, 6, 7 and 6, respectively. In the following, we describe the detailed information for each dataset that is used in our classification tasks.

Split MNIST. We divide MNIST which contains 60k training samples into five tasks, each consisting of images from two classes, in consecutive order of their displayed digits, while increasing the numbers represented in the images [De Lange and Tuytelaars \(2021\)](#).

Split CIFAR10. We split CIFAR10 into five tasks where each task consists of samples from two different classes [De Lange and Tuytelaars \(2021\)](#).

Split CIFAR100. We split CIFAR100 into 20 tasks where each task has 2500 examples

from five different classes [Lopez-Paz and Ranzato \(2017\)](#).

Split MiniImageNet. We divide the MiniImageNet into 20 tasks [Vinyals et al. \(2016\)](#), where each task collects the images of five classes [Aljundi et al. \(2019b\)](#).

In the following, we describe the detailed information of the network architecture used in our classification task.

We adapt ResNet 18 [He et al. \(2016\)](#) for Split CIFAR10 and Split CIFAR100. We use an MLP network with 2 hidden layers of 400 units each [De Lange and Tuytelaars \(2021\)](#) for Split MNIST. The maximum memory size for Split MNIST, Split CIFAR10, Split CIFAR100 are 2000, 1000 and 5000, respectively. To enable the proposed OAES for the classification, we train an individual classifier for each component, similar to [Ye and Bors \(2022a\)](#). At the testing phase, we make the component selection by comparing the sample log-likelihood and the classifier of the selected component is used for prediction.

We introduce the baselines used for the classification task but which are not mentioned in the paper.

Finetune trains a single model directly on a new batch of images during the online continual learning.

Gradient Episodic Memory (GEM) [Lopez-Paz and Ranzato \(2017\)](#) is a memory-based approach that would use the memory to store past samples. GEM is also required to access both the task label and class label during the training.

Dynamic-OCM [Ye and Bors \(2022a\)](#) is a dynamic expansion model which proposes an online cooperative memorization (OCM) approach. OCM manages two memory buffers, aiming to store short- and long-term knowledge during training. In addition, Dynamic-OCM detects the change of the loss value as expansion signals, which does not have theoretical guarantees.

Incremental Classifier and Representation Learning (iCARL) [Rebuffi et al. \(2017\)](#) is a standard memory-based method used in a class incremental setup.

reservoir* [Vitter \(1985\)](#) is a memory-based approach that stores the observed sample into a memory buffer \mathcal{M} with probability $|\mathcal{M}|/n$ where n is the number of stored samples, and $|\cdot|$ represents the cardinality of a set.

MIR [Aljundi et al. \(2019b\)](#) introduces a retrieval strategy for the sample selection in the memory during the Online Continual Learning (OCL). However, the retrieval strategy in MIR requires evaluating the loss in each training session. This means that MIR requires modifying the retrieval strategy for different tasks such as classification or

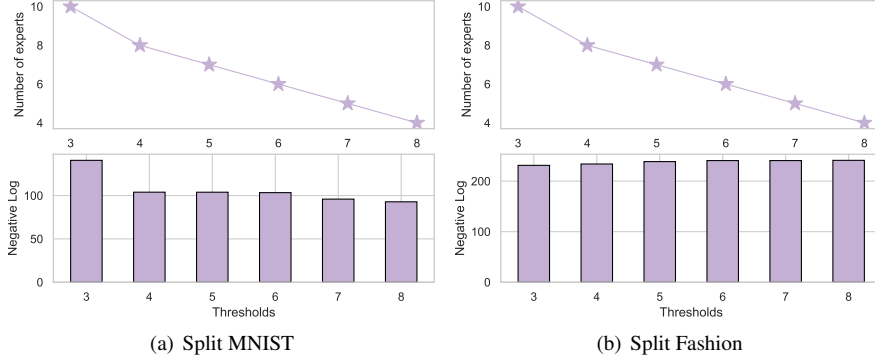


Figure 2: The performance and the number of components of OAES when changing the threshold β .

generation tasks. The proposed OCM does not change the sample selection strategy for different tasks since we evaluate the sample similarity in the given feature space using the kernel function from Eq. (16) from the paper.

GSS [Aljundi et al. \(2019c\)](#) formulates the sample selection process as a constraint reduction problem. GSS stores samples in a buffer based on the gradient information which requires to access the class labels and can not be applied in the unsupervised learning setting.

I Additional information for ablation study

In this section, we provide more ablation studies in order to investigate the effectiveness of each modules in the proposed OAES.

I.1 The impact of the threshold β

We investigate the effect of β by training OAES on Split MNIST with different threshold values and the results are reported in Fig. 2a. It observes that a small β leads to learning more components while improving the performance. In contrast, a large β allows OAES to employ fewer components. We train OAEs on Split Fashion with different threshold values and the results are reported in Fig. 2a. We observe that increasing the number of components can improve the performance of Split Fashion.

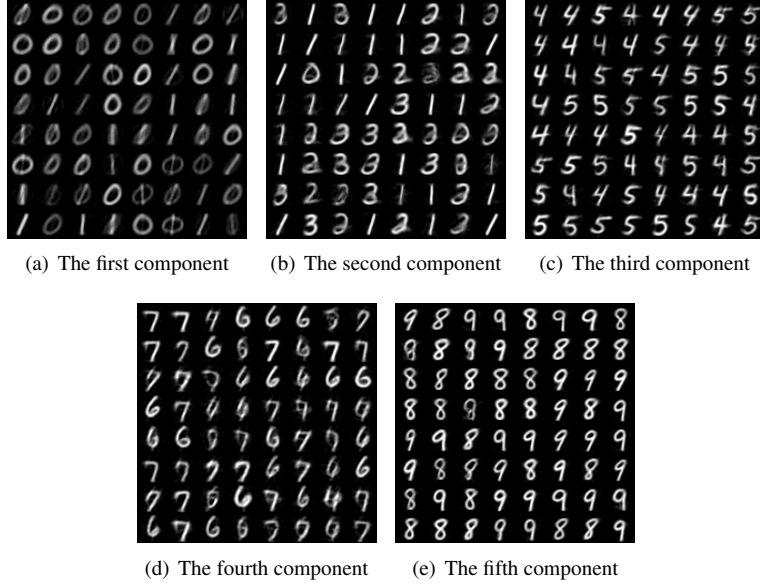


Figure 3: The generation results from each learned VAE component under Split MNIST.

I.2 The impact of the memory buffer size

We train various models under Split MNIST by using different memory buffer sizes and the results are reported in Fig. 4a. These results show that a large-scale memory buffer can improve the performance for all DEM models. We also train various models on Split Fashion with different memory configurations and the results are reported in Fig. 4b. The proposed OAES outperforms other baselines, especially when the memory buffer size is very small (500).

I.3 Knowledge diversity among experts

We train OAES with $\beta = 1.7$ on Split MNIST for the classification task where we evaluate the discrepancy score (the left-hand-side (LHS) of Eq.(15) of the paper) at each training step. We plot the results in Fig. 5, which shows that there are five peak-to-valley discrepancy scores, corresponding to the five different underlying distributions (tasks). We also show the generation results from each learned component of the proposed model in Fig. 3. It observes that the first component learns the same underlying data distribution while other components modelling a unique underlying data distribution.

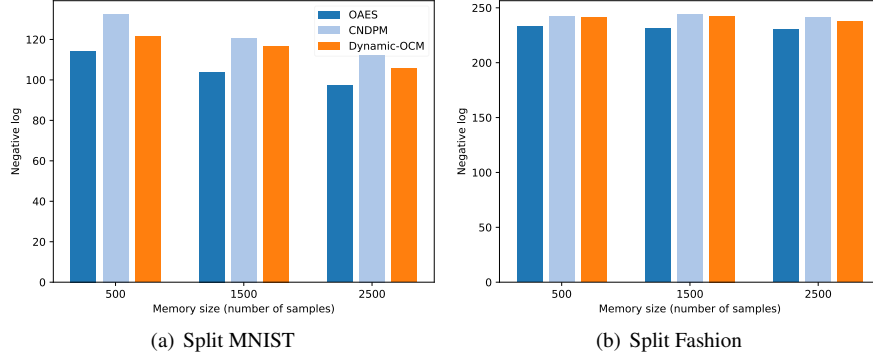


Figure 4: The performance of various models when changing the memory buffer size.

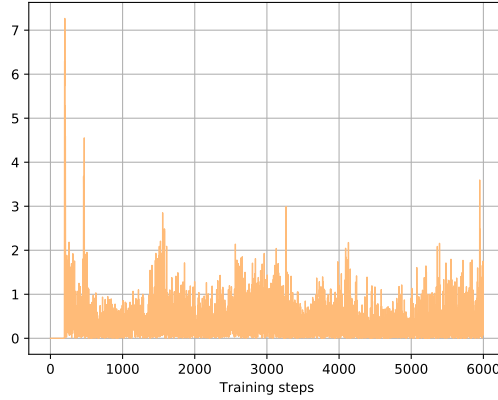


Figure 5: The expansion evaluation (the left-hand-side (LHS) of Eq.(15) of the paper) evaluated by OAES at each training step.

bution. These results show that the proposed OAES can accurately detect the data distribution shift and provide appropriate signals for the model expansion.

I.4 The impact of batch size change

In this section, we investigate the impact of the batch size change of the proposed OAES. We train the OAES under Split MNIST using the different batch sizes, and the empirical results are reported in Fig. 6. It observes that changing the batch size has little effect on the performance and the model size. This shows that the proposed OAES is robust to the change in the batch size.

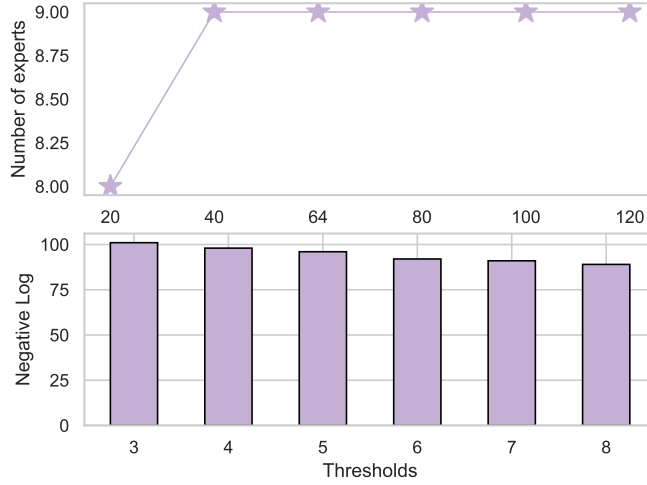


Figure 6: The performance and the number of experts of the proposed OAES under Split MNIST when changing the batch size.

I.5 Theoretical results

We train the proposed OAES and the static model (ER) under Split MNIST in which we evaluate the source risk (the LHS of Eq. (2)) and the target risk (the first term in RHS of Eq. (2)). We plot the results evaluated by each training step in Fig. 7, where 'OAES-source' and 'OAES-target' denote the source and target risk evaluated by the OAES. It observes that ER firstly gets good performance on the initial training phase and gradually losses performance as the number of training steps increases. This is because the memory-based approaches can not store all previously learnt information during the training. In contrast, the proposed OAES gradually improves the performance of the target samples when performing more training steps. These results demonstrate that the proposed OAES can accumulate knowledge without forgetting it and can achieve better generalization performance than the static model.

I.6 Classification on fuzzy task boundaries

In this section, we apply our model in a real-world and more challenging TFCL setting, called fuzzy task boundaries [Lee et al. \(2020\)](#). In this setting, we randomly swap samples between two tasks for each data stream. The results for the fuzzy task boundaries are reported in Tab. 1. These results show that the proposed OAES still outperforms other baselines in this challenging setting.

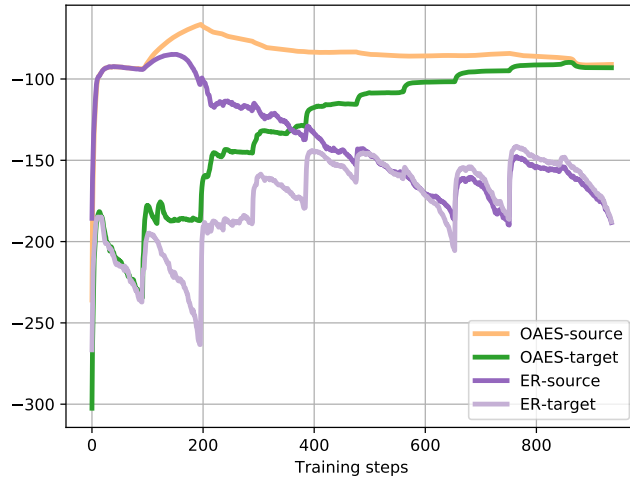


Figure 7: The source and target risks (sample log-likelihood) of the OAES and the static model under Split MNIST.

Methods	Split MNIST	Split CIFAR10	Split MImageNet
Vanilla	21.53 ± 0.1	20.69 ± 2.4	3.05 ± 0.6
ER	79.74 ± 4.0	37.15 ± 1.6	26.47 ± 2.3
MIR	84.80 ± 1.9	38.70 ± 1.7	25.83 ± 1.5
ER + GMED	82.73 ± 2.6	40.57 ± 1.7	28.20 ± 0.6
MIR+GMED	86.17 ± 1.7	41.22 ± 1.1	26.86 ± 0.7
OAES	90.23 ± 1.2	44.26 ± 1.1	29.63 ± 0.8

Table 1: The classification accuracy of five independent runs for various models over data streams with fuzzy task boundaries.

I.7 Comparison of computational costs

In this section, we compare our model (OAES) with Dynamic-OCM in terms of computational costs. The training times (minutes) of various models under the density estimation task are reported in Tab. 2. It observes that the proposed OAES requires less computational costs than Dynamic-OCM while achieving better performance than Dynamic-OCM.

I.8 The other performance criterion on the generative modelling

In this section, we employ the other performance criteria including the Inception Score (IS) and Fréchet Inception Distance (FID) for the generative modelling task. We train our model (OAES) under Split CIFAR10 and report the results in Tab. 3 where the

Methods	Split MNIST	Split Fashion	Split MNIST-Fashion	Cross domain
OAES	30.2	35.6	71.2	105.4
Dynamic-OCM	46.6	52.6	82.7	120.6

Table 2: The training times (minutes) for various models under the continual generative modelling task.

Methods	IS	FID	Memory	N
VAE-ELBO-Random	3.84	116.26	1.0K	1
CNDPM Lee et al. (2020)	4.12	95.23	1.0K	30
LIMix Ye and Bors (2021a)	3.02	156.46	1.0K	30
VAE-ELBO-OCM	4.13	98.76	1.0K	1
Dynamic-ELBO-OCM	4.16	92.99	0.9K	3
OAES	4.25	88.62	0.9K	3

Table 3: IS and FID scores under Split CIFAR10.

332 results of other baselines are taken from [Ye and Bors \(2022a\)](#). These results show that
333 the proposed OAES outperforms other baselines in terms of IS and FID.

334 References

- 335 Yuri Burda, Roger Grosse, and Ruslan Salakhutdinov. Importance weighted autoen-
336 coders. In *Proc. Int. Cont. of Learning Representations (ICLR)*, *arXiv preprint*
337 *arXiv:1509.00519*, 2015. 9, 21
- 338 Lars Maaløe, Casper Kaae Sønderby, Søren Kaae Sønderby, and Ole Winther. Auxil-
339 iary deep generative models. In *Proc. Int. Conf. on Machine Learning (ICML)*, vol.
340 *PMLR 48*, pages 1445–1453, 2016. 10
- 341 Fei Ye and Adrian G. Bors. Continual variational autoencoder learning via online
342 cooperative memorization, 2022a. 11, 21, 22, 28
- 343 Evgenii Egorov, Anna Kuzina, and Evgeny Burnaev. Boovae: Boosting approach for
344 continual learning of vae. *Advances in Neural Information Processing Systems*, 34:
345 17889–17901, 2021. 11
- 346 Kamil Deja, Paweł Wawrzyński, Daniel Marczak, Wojciech Masarczyk, and Tomasz

Trzciński. Multiband vae: Latent space partitioning for knowledge consolidation in	347
continual learning. <i>arXiv preprint arXiv:2106.12196</i> , 2021. 11	348
Fei Ye and Adrian G.Bors Bors. Learning latent representations across multiple data	349
domains using lifelong vaegan. In <i>Proc. of European Conference on Computer Vi-</i>	350
<i>sion (ECCV)</i> , vol. LNCS 12365, pages 777–795, 2020a. 19, 21	351
Fei Ye and Adrian G. Bors. Lifelong infinite mixture model based on knowledge-	352
driven dirichlet process. In <i>Proceedings of the IEEE/CVF International Conference</i>	353
<i>on Computer Vision (ICCV)</i> , pages 10695–10704, October 2021a. 19, 21, 28	354
D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In <i>Proc. Int.</i>	355
<i>Conf. on Learning Representations (ICLR)</i> , <i>arXiv preprint arXiv:1412.6980</i> , 2015.	356
21	357
Fei Ye and Adrian G Bors. Task-free continual learning via online discrepancy distance	358
learning. <i>arXiv preprint arXiv:2210.06579</i> , 2022b. 21	359
Fei Ye and Adrian G Bors. Dynamic self-supervised teacher-student network learning.	360
<i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 2022c. 21	361
Fei Ye and Adrian G. Bors. Lifelong twin generative adversarial networks. In <i>Proc.</i>	362
<i>IEEE Int. Conf. on Image Processing (ICIP)</i> , pages 1289–1293, 2021b. doi: 10.	363
1109/ICIP42928.2021.9506116. 21	364
Fei Ye and Adrian G. Bors. Lifelong mixture of variational autoencoders. <i>IEEE</i>	365
<i>Transactions on Neural Networks and Learning Systems</i> , pages 1–14, 2021c. doi:	366
10.1109/TNNLS.2021.3096457. 21	367
Rahaf Aljundi, Klaas Kelchtermans, and Tinne Tuytelaars. Task-free continual learn-	368
ing. In <i>Proc. of IEEE/CVF Conf. on Computer Vision and Pattern Recognition</i> , pages	369
11254–11263, 2019a. 21	370
Fei Ye and Adrian G Bors. Learning an evolved mixture model for task-free continual	371
learning. In <i>2022 IEEE International Conference on Image Processing (ICIP)</i> , pages	372
1936–1940. IEEE, 2022d. 21	373
Fei Ye and Adrian Bors. Lifelong teacher-student network learning. <i>IEEE Transactions</i>	374
<i>on Pattern Analysis and Machine Intelligence</i> , 2021d. 21	375

- 376 Fei Ye and Adrian G. Bors. Lifelong learning of interpretable image representations. In
 377 *Proc. Int. Conf. on Image Processing Theory, Tools and Applications (IPTA)*, pages
 378 1–6, 2020b. doi: 10.1109/IPTA50016.2020.9286663. 21
- 379 Fei Ye and Adrian G. Bors. Lifelong generative modelling using dynamic expan-
 380 sion graph model. In *Thirty-Sixth AAAI Conference on Artificial Intelligence*,
 381 AAAI, pages 8857–8865. AAAI Press, 2022e. URL [https://ojs.aaai.org/
 382 index.php/AAAI/article/view/20867](https://ojs.aaai.org/index.php/AAAI/article/view/20867). 21
- 383 Fei Ye and Adrian G. Bors. Deep mixture generative autoencoders. *IEEE Trans-*
 384 *actions on Neural Networks and Learning Systems*, pages 1–15, 2021e. doi:
 385 10.1109/TNNLS.2021.3071401. 21
- 386 Fei Ye and Adrian G Bors. Mixtures of variational autoencoders. In *Proc. Int. Conf.*
 387 *on Image Processing Theory, Tools and Applications (IPTA)*, pages 1–6, 2020c. 21
- 388 Matthias De Lange and Tinne Tuytelaars. Continual prototype evolution: Learning
 389 online from non-stationary data streams. In *Proc. of the IEEE/CVF International*
 390 *Conference on Computer Vision*, pages 8250–8259, 2021. 21, 22
- 391 David Lopez-Paz and Marc’Aurelio Ranzato. Gradient episodic memory for continual
 392 learning. In *Advances in Neural Information Processing Systems*, pages 6467–6476,
 393 2017. 22
- 394 Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Koray Kavukcuoglu, and Daan
 395 Wierstra. Matching networks for one shot learning. *Advances in neural information*
 396 *processing systems (NIPS)*, 29:3637–3645, 2016. 22
- 397 Rahaf Aljundi, Eugene Belilovsky, Tinne Tuytelaars, Laurent Charlin, Massimo Caccia,
 398 Min Lin, and Lucas Page-Caccia. Online continual learning with maximal interfered
 399 retrieval. In *Advances in Neural Information Processing Systems (NeurIPS)*,
 400 pages 11872–11883, 2019b. 22
- 401 K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition.
 402 In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages
 403 770–778, 2016. 22
- 404 Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert.
 405 icarl: Incremental classifier and representation learning. In *Proceedings of the*

IEEE conference on Computer Vision and Pattern Recognition, pages 2001–2010, 2017. 22

Jeffrey S Vitter. Random sampling with a reservoir. *ACM Transactions on Mathematical Software (TOMS)*, 11(1):37–57, 1985. 22

Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient based sample selection for online continual learning. In *Proc. Neural Inf. Proc. Systems (NeurIPS)*, pages 11817–11826, 2019c. 23

Soochan Lee, Junsoo Ha, Dongsu Zhang, and Gunhee Kim. A neural Dirichlet process mixture model for task-free continual learning. In *Int. Conf. on Learning Representations (ICLR)*, *arXiv preprint arXiv:2001.00689*, 2020. 26, 28