# Online Cooperative Memorization for Variational Autoencoders Supporting Document

Fei Ye and Adrian G. Bors

## APPENDIX A
### PROOF OF THEOREM 1

When $p_\theta(\mathbf{x} \,|\, \mathbf{z})$ is the Gaussian decoder, the computation of $\log p_\theta(\mathbf{x} \,|\, \mathbf{z})$ involves the noise value $\sigma$ :

$$\log p_\theta\left(\mathbf{x} \,|\, \mathbf{z}\right) = -\frac{1}{2\sigma^2}\|\mathbf{x} - \mu_\theta\left(\mathbf{z}\right)\|^2 - \frac{1}{2}\log 2\pi\sigma^2, \quad (1)$$

where $\mu_\theta(\mathbf{z})$ is the mean of distribution $p_\theta(\mathbf{x} \,|\, \mathbf{z})$. In order to simplify Eq. (1), the noise $\sigma$ is set to $1/\sqrt{2}$, resulting in :

$$\log p_\theta\left(\mathbf{x} \,|\, \mathbf{z}\right) = -\|\mathbf{x} - \mu_\theta\left(\mathbf{z}\right)\|^2 - \frac{1}{2}\log \pi. \quad (2)$$

We substract the KL divergence resulting in :

$$\log p_\theta\left(\mathbf{x} \,|\, \mathbf{z}\right) - D_{KL}(q_\omega(\mathbf{x} \,|\, \mathbf{z}) \,|\, p(\mathbf{z})) =$$
$$- \|\mathbf{x} - \mu_\theta\left(\mathbf{z}\right)\|_2^2 - D_{KL}(q_\omega(\mathbf{x} \,|\, \mathbf{z}) \,|\, p(\mathbf{z})) - \frac{1}{2}\log \pi. \quad (3)$$

Then we consider the expectation in both sides, resulting in :

$$\inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_\mathbf{x}}\mathbb{E}_{q_\omega(\mathbf{z}\,|\,\mathbf{x})}\left[\log p_\theta\left(\mathbf{x} \,|\, \mathbf{z}\right) - D_{KL}(q_\omega(\mathbf{x}\,|\,\mathbf{z})\,|\,p(\mathbf{z}))\right]$$
$$= \inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_\mathbf{x}}\mathbb{E}_{q_\omega(\mathbf{z}\,|\,\mathbf{x})}\left[ - \|\mathbf{x} - \mu_\theta\left(\mathbf{z}\right)\|_2^2\right.$$
$$\left. - D_{KL}(q_\omega(\mathbf{x}\,|\,\mathbf{z})\,|\,p(\mathbf{z})) - \frac{1}{2}\log \pi\right]. \quad (4)$$

where the first term in the right-hand side of Eq. (4) can be rewritten as $\mathcal{L}(\mathbf{x}, \mathrm{G}_i(\mathbf{z}))$, and this relationship becomes :

$$\inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_\mathbf{x}}\mathbb{E}_{q_\omega(\mathbf{z}\,|\,\mathbf{x})}\left[\log p_\theta\left(\mathbf{x}\,|\,\mathbf{z}\right) - D_{KL}(q_\omega(\mathbf{x}\,|\,\mathbf{z})\,|\,p(\mathbf{z}))\right]$$
$$= \inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_\mathbf{x}}\mathbb{E}_{q_\omega(\mathbf{z}\,|\,\mathbf{x})}\left[ - \mathcal{L}(\mathbf{x}, \mathrm{G}_i(\mathbf{z}))\right.$$
$$\left. - D_{KL}(q_\omega(\mathbf{x}\,|\,\mathbf{z})\,|\,p(\mathbf{z})) - \frac{1}{2}\log \pi\right]. \quad (5)$$

where the first term in the left-hand side (LHS) of Eq. (5) is the ELBO, defined in Eq. (1) of the paper. Since the KL divergence $D_{KL}(\cdot)$ is equal or larger than 0, we have the following inequality :

$$\inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_\mathbf{x}}[\mathcal{L}_{ELBO}(\mathbf{x}; \theta, \omega)] =$$
$$\inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_\mathbf{x}}\mathbb{E}_{q_\omega(\mathbf{z}\,|\,\mathbf{x})}[-\mathcal{L}(\mathbf{x}, \mathrm{G}_i(\mathbf{z}))$$
$$- D_{KL}(q_\omega(\mathbf{z}\,|\,\mathbf{x})\,||\,p(\mathbf{z})) - \frac{1}{2}\log \pi]$$
$$\leq \inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_\mathbf{x}}\mathbb{E}_{q_\omega(\mathbf{z}\,|\,\mathbf{x})}[-\mathcal{L}(\mathbf{x}, \mathrm{G}_i(\mathbf{z}))] - \frac{1}{2}\log \pi, \quad (6)$$

From the inequality from Eq. (8) from the paper after multiplying with $-1$ :

$$-\mathrm{W}_\mathcal{L}^\star(\mathrm{P}_\mathbf{x}, \mathrm{P}_{\mathbf{G}_i}) \geq \inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_\mathbf{x}}\mathbb{E}_{q_\omega(\mathbf{z}\,|\,\mathbf{x})}[-\mathcal{L}(\mathbf{x}, \mathrm{G}_i(\mathbf{z}))], \quad (7)$$

and then rewrite Eq. (6) by considering Eq. (7), resulting in:

$$\inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_\mathbf{x}}[\mathcal{L}_{ELBO}(\mathbf{x}; \theta, \omega)] \leq -\mathrm{W}_\mathcal{L}^\star(\mathrm{P}_\mathbf{x}, \mathrm{P}_{\mathrm{G}_i}) - \frac{1}{2}\log \pi. \quad (8)$$

Eq. (8) proves Theorem 1 $\qquad\square$

## APPENDIX B
### PROOF OF THEOREM 2

We consider Eq. (8) and add $-\mathrm{W}_\mathcal{L}^\star(\mathrm{P}_{m_i}, \mathrm{P}_{\mathrm{G}_i})$ to both sides of resulting in :

$$\inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_\mathbf{x}}[\mathcal{L}_{ELBO}(\mathbf{x}; \theta, \omega)] - \mathrm{W}_\mathcal{L}^\star(\mathrm{P}_{m_i}, \mathrm{P}_{\mathrm{G}_i}) \leq$$
$$- \mathrm{W}_\mathcal{L}^\star(\mathrm{P}_{m_i}, \mathrm{P}_{\mathrm{G}_i}) - \mathrm{W}_\mathcal{L}^\star(\mathrm{P}_\mathbf{x}, \mathrm{P}_{\mathrm{G}_i}) - \frac{1}{2}\log \pi \quad (9)$$

The first term in the right-hand side (RHS) is bounded, similarly to Eq. (7), but on the memory buffer $m_i$ :

$$\inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_{m_i}}\mathbb{E}_{q_\omega(\mathbf{z}\,|\,\mathbf{x})}[-\mathcal{L}(\mathbf{x}, \mathrm{G}_i(\mathbf{z}))] \leq -\mathrm{W}_\mathcal{L}^\star(\mathrm{P}_{m_i}, \mathrm{P}_{\mathrm{G}_i}), \quad (10)$$

then we have :

$$\inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_{m_i}}\mathbb{E}_{q_\omega(\mathbf{z}\,|\,\mathbf{x})}[-\mathcal{L}(\mathbf{x}, \mathrm{G}_i(\mathbf{z}))]+$$
$$\left| \inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_{m_i}}\mathbb{E}_{q_\omega(\mathbf{z}\,|\,\mathbf{x})}[-\mathcal{L}(\mathbf{x}, \mathrm{G}_i(\mathbf{z}))] - \mathrm{W}_\mathcal{L}^\star(\mathrm{P}_{m_i}, \mathrm{P}_{\mathrm{G}_i})\right|$$
$$\geq -\mathrm{W}_\mathcal{L}^\star(\mathrm{P}_{m_i}, \mathrm{P}_{\mathrm{G}_i}). \quad (11)$$

Then, by using Eq. (9), we derive :

$$\inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_\mathbf{x}}[\mathcal{L}_{ELBO}(\mathbf{x}; \theta, \omega)] - \mathrm{W}_\mathcal{L}^\star(\mathrm{P}_{m_i}, \mathrm{P}_{\mathrm{G}_i})$$
$$\leq \inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_{m_i}}\mathbb{E}_{q_\omega(\mathbf{z}\,|\,\mathbf{x})}[-\mathcal{L}(\mathbf{x}, \mathrm{G}_i(\mathbf{z}))] - \mathrm{W}_\mathcal{L}^\star(\mathrm{P}_\mathbf{x}, \mathrm{P}_{\mathbf{G}_i})$$
$$+ \left| \inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_{m_i}}\mathbb{E}_{q_\omega(\mathbf{z}\,|\,\mathbf{x})}[-\mathcal{L}(\mathbf{x}, \mathrm{G}_i(\mathbf{z}))] - \mathrm{W}_\mathcal{L}^\star(\mathrm{P}_{m_i}, \mathrm{P}_{\mathbf{G}_i})\right|$$
$$- \frac{1}{2}\log \pi. \quad (12)$$

We then add the negative KL divergence term in both sides of Eq. (12) :

$$\inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_\mathbf{x}}[\mathcal{L}_{ELBO}(\mathbf{x};\theta,\omega)] - W_\mathcal{L}^\star(\mathrm{P}_{m_i},\mathrm{P}_{\mathrm{G}_i})$$
$$- \inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_{m_i}}[D_{KL}(q_\omega(\mathbf{z}\,|\,\mathbf{x})\,||\,p(\mathbf{z}))] \le$$
$$\inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_{m_i}}\mathbb{E}_{q_\omega(\mathbf{z}\,|\,\mathbf{x})}[-\mathcal{L}(\mathbf{x},\mathrm{G}_i(\mathbf{z}))]$$
$$- \inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_{m_i}}[D_{KL}(q_\omega(\mathbf{z}\,|\,\mathbf{x})\,||\,p(\mathbf{z}))] - \frac{1}{2}\log\pi \quad (13)$$
$$- W_\mathcal{L}^\star(\mathrm{P}_\mathbf{x},\mathrm{P}_{\mathbf{G}_i}) + \Big| \inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_{m_i}}\mathbb{E}_{q_\omega(\mathbf{z}\,|\,\mathbf{x})}[$$
$$- \mathcal{L}(\mathbf{x},\mathrm{G}_i(\mathbf{z}))] - W_\mathcal{L}^\star(\mathrm{P}_{m_i},\mathrm{P}_{\mathbf{G}_i})\Big|,$$

According to the definition of ELBO, this can be rewritten as :

$$\inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_\mathbf{x}}[\mathcal{L}_{ELBO}(\mathbf{x};\theta,\omega)] - W_\mathcal{L}^\star(\mathrm{P}_{m_i},\mathrm{P}_{\mathrm{G}_i})$$
$$- \inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_{m_i}}[D_{KL}(q_\omega(\mathbf{z}\,|\,\mathbf{x})\,||\,p(\mathbf{z}))] \le$$
$$\inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_{m_i}}[\mathcal{L}_{ELBO}(\mathbf{x};\theta,\omega)] - W_\mathcal{L}^\star(\mathrm{P}_\mathbf{x},\mathrm{P}_{\mathbf{G}_i})$$
$$+ \Big| \inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_{m_i}}\mathbb{E}_{q_\omega(\mathbf{z}\,|\,\mathbf{x})}[-\mathcal{L}(\mathbf{x},\mathrm{G}_i(\mathbf{z}))] - W_\mathcal{L}^\star(\mathrm{P}_{m_i},\mathrm{P}_{\mathrm{G}_i})\Big|,$$
$$(14)$$

Then we rewrite Eq. (14), resulting in :

$$\inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_\mathbf{x}}[\mathcal{L}_{ELBO}(\mathbf{x};\theta,\omega)] \le$$
$$\inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_{m_i}}[\mathcal{L}_{ELBO}(\mathbf{x};\theta,\omega)] + W_\mathcal{L}^\star(\mathrm{P}_{m_i},\mathrm{P}_{\mathrm{G}_i})$$
$$- W_\mathcal{L}^\star(\mathrm{P}_\mathbf{x},\mathrm{P}_{\mathrm{G_i}}) + \inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_{m_i}}[D_{KL}(q_\omega(\mathbf{z}\,|\,\mathbf{x})\,||\,p(\mathbf{z}))]$$
$$+ \Big| \inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_{m_i}}\mathbb{E}_{q_\omega(\mathbf{z}\,|\,\mathbf{x})}[-\mathcal{L}(\mathbf{x},\mathrm{G}_i(\mathbf{z}))] - W_\mathcal{L}^\star(\mathrm{P}_{m_i},\mathrm{P}_{\mathrm{G}_i})\Big|.$$
$$(15)$$

We consider that $\mathcal{L}(\cdot)$ satisfies the triangle inequality :

$$W_\mathcal{L}^\star(\mathrm{P}_{m_i},\mathrm{P}_{\mathrm{G}_i}) + W_\mathcal{L}^\star(\mathrm{P}_\mathbf{x},\mathrm{P}_{\mathrm{G}_i}) \ge W_\mathcal{L}^\star(\mathrm{P}_\mathbf{x},\mathrm{P}_{m_i}) \quad (16)$$

We move the second term from the LHS of Eq. (16) in the RHS :

$$W_\mathcal{L}^\star(\mathrm{P}_\mathbf{x},\mathrm{P}_{\mathrm{G}_i}) \ge W_\mathcal{L}^\star(\mathrm{P}_\mathbf{x},\mathrm{P}_{m_i}) - W_\mathcal{L}^\star(\mathrm{P}_{m_i},\mathrm{P}_{\mathrm{G}_i}) \quad (17)$$

Then we replace $W_\mathcal{L}^\star(\mathrm{P}_\mathbf{x},\mathrm{P}_{\mathrm{G}_i})$ from Eq. (15) by the expression of Eq. (17), resulting in :

$$\inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_\mathbf{x}}[\mathcal{L}_{ELBO}(\mathbf{x};\theta,\omega)] \le$$
$$\inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_{m_i}}[\mathcal{L}_{ELBO}(\mathbf{x};\theta,\omega)] + 2W_\mathcal{L}^\star(\mathrm{P}_{m_i},\mathrm{P}_{\mathrm{G}_i}) \quad (18)$$
$$- W_\mathcal{L}^\star(\mathrm{P}_\mathbf{x},\mathrm{P}_{m_i}) + \tilde{\mathrm{F}}(\mathrm{P}_{\mathrm{G}_i},\mathrm{P}_{m_i}),$$

where $\tilde{\mathrm{F}}(\mathrm{P}_{\mathrm{G}_i},\mathrm{P}_{m_i})$ is expressed as :

$$\tilde{\mathrm{F}}(\mathrm{P}_{\mathrm{G}_i},\mathrm{P}_{m_i}) = \inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_{m_i}}[D_{KL}(q_\omega(\mathbf{z}\,|\,\mathbf{x})\,||\,p(\mathbf{z}))]$$
$$+ \Big| \inf_{q_\omega(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathrm{P}_{m_i}}\mathbb{E}_{q_\omega(\mathbf{z}\,|\,\mathbf{x})}[-\mathcal{L}(\mathbf{x},\mathrm{G}_i(\mathbf{z}))]$$
$$- W_\mathcal{L}^\star(\mathrm{P}_{m_i},\mathrm{P}_{\mathbf{G}_i})\Big|$$
$$(19)$$

$\square$

# APPENDIX C
# PROOF OF THEOREM 3

Let us firstly consider a certain component ($a_i$-th component) that has been trained only once. From Theorem 2 we derive the bound as follows :

$$\mathbb{E}_{\mathrm{P}_{\hat{\mathbf{x}}^{\tilde{a}_i}}}[\mathcal{L}_{ELBO}(\mathbf{x};\theta,\omega)] \le \mathbb{E}_{\mathrm{P}_{\mathbf{x}^{\tilde{a}_i}}}[\mathcal{L}_{ELBO}(\mathbf{x};\theta,\omega)]$$
$$+ 2W_\mathcal{L}^\star(\mathrm{P}_{\mathbf{x}^{\tilde{a}_i}},\mathrm{P}_{\mathrm{G}^{a_i}}) - W_\mathcal{L}^\star(\mathrm{P}_{\hat{\mathbf{x}}^{\tilde{a}_i}},\mathrm{P}_{\mathbf{x}^{\tilde{a}_i}}) \quad (20)$$
$$+ \tilde{\mathrm{F}}(\mathrm{P}_{\mathrm{G}^{a_i}},\mathrm{P}_{\mathbf{x}^{\tilde{a}_i}}),$$

Eq. (20) holds because we treat $\mathrm{P}_{\hat{\mathbf{x}}^{\tilde{a}_i}}$ and $\mathrm{P}_{\mathbf{x}^{\tilde{a}_i}}$ as the target and source domain respectively. In the following, we consider a component ($b_i$-th component) that has been trained more than once. Since the $b_i$-th component would learn more than one task, we particularly focus on a certain task ($\tilde{b}_i^q$-th task). We firstly consider to treat $\mathrm{P}_{\hat{\mathbf{x}}_i^{\tilde{b}q}}$ and $\mathrm{P}_{\mathbf{x}_i^{\tilde{b}q}}$ as the target and source domain respectively. Then we derive the bound as :

$$\mathbb{E}_{\mathrm{P}_{\hat{\mathbf{x}}_i^{\tilde{b}q}}}[\mathcal{L}_{ELBO}(\mathbf{x};\theta,\omega)] \le \mathbb{E}_{\mathrm{P}_{\mathbf{x}_i^{\tilde{b}q}}}[\mathcal{L}_{ELBO}(\mathbf{x};\theta,\omega)]$$
$$+ 2W_\mathcal{L}^\star(\mathrm{P}_{\mathbf{x}_i^{\tilde{b}q}},\mathrm{P}_{\mathrm{G}^{b_i}}) - W_\mathcal{L}^\star(\mathrm{P}_{\hat{\mathbf{x}}_i^{\tilde{b}q}},\mathrm{P}_{\mathbf{x}_i^{\tilde{b}q}}) \quad (21)$$
$$+ \tilde{\mathrm{F}}(\mathrm{P}_{\mathrm{G}^{b_i}},\mathrm{P}_{\mathbf{x}_i^{\tilde{b}q}}),$$

We do not specify the state (the number of retraining processes) of each generator distribution $\mathrm{P}_{\mathrm{G}i}$ in order to simplify the notation. We have the empirical distribution $\mathrm{P}_{\hat{\mathbf{x}}^{(\tilde{b}_i^q,1)}}$ for one time of the generative replay processes (see Definition 6). We treat $\mathrm{P}_{\hat{\mathbf{x}}^{(\tilde{b}_i^q,0)}} = \mathrm{P}_{\hat{\mathbf{x}}_i^{\tilde{b}q}}$ and $\mathrm{P}_{\hat{\mathbf{x}}^{(\tilde{b}_i^q,1)}}$ as the target and source domain, respectively. We then derive the bound between $\mathrm{P}_{\hat{\mathbf{x}}^{(\tilde{b}_{(i,q)},0)}}$ and $\mathrm{P}_{\hat{\mathbf{x}}^{(\tilde{b}_{(i,q)},1)}}$ as follows :

$$\mathbb{E}_{\mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,0)}}}[\mathcal{L}_{ELBO}(\mathbf{x};\theta,\omega)] \le \mathbb{E}_{\mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,1)}}}[\mathcal{L}_{ELBO}(\mathbf{x};\theta,\omega)]$$
$$+ 2W_\mathcal{L}^\star(\mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_{(i,q)},1)}},\mathrm{P}_{\mathrm{G}^{b_i}}) - W_\mathcal{L}^\star(\mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,1)}},\mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,0)}})$$
$$+ \tilde{\mathrm{F}}(\mathrm{P}_{\mathrm{G}^{b_i}},\mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,1)}}),$$
$$(22)$$

Through mathematical induction, we have the bounds :

$$\mathbb{E}_{\mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,1)}}}[\mathcal{L}_{ELBO}(\mathbf{x};\theta,\omega)] \le \mathbb{E}_{\mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,2)}}}[\mathcal{L}_{ELBO}(\mathbf{x};\theta,\omega)]$$
$$+ 2W_\mathcal{L}^\star(\mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,2)}},\mathrm{P}_{\mathrm{G}^{b_i}})$$
$$- W_\mathcal{L}^\star(\mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,2)}},\mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,1)}})$$
$$+ \tilde{\mathrm{F}}(\mathrm{P}_{\mathrm{G}^{b_i}},\mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,2)}})$$
$$\cdots$$
$$\mathbb{E}_{\mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,c_i^q-1)}}}[\mathcal{L}_{ELBO}(\mathbf{x};\theta,\omega)] \le \mathbb{E}_{\mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,c_i^q)}}}[\mathcal{L}_{ELBO}(\mathbf{x};\theta,\omega)]$$
$$+ 2W_\mathcal{L}^\star(\mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,c_i^q)}},\mathrm{P}_{\mathrm{G}^{b_i}})$$
$$- W_\mathcal{L}^\star(\mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,c_i^q)}},\mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,c_i^q-1)}})$$
$$+ \tilde{\mathrm{F}}(\mathrm{P}_{\mathrm{G}^{b_i}},\mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,c_i^q)}},)$$
$$(23)$$

where $c_i^q$ denotes the number of generative replay processes for the $\tilde{b}_i^q$-th task, achieved by the $b_i$-th component. We then sum up all above inequalities, resulting in :

$$\mathbb{E}_{\mathrm{P}_{\hat{\mathbf{x}}_i^{\tilde{b}q}}}[\mathcal{L}_{ELBO}(\mathbf{x};\theta,\omega)] \le \mathbb{E}_{\mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,c_i^q)}}}[\mathcal{L}_{ELBO}(\mathbf{x};\theta,\omega)]$$
$$+ \sum_{s=0}^{\tilde{c}_{(i,q)}} \Big\{ 2W_\mathcal{L}^\star(\mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,s)}},\mathrm{P}_{\mathrm{G}^{b_i}}) - W_\mathcal{L}^\star(\mathrm{P}_{\hat{\mathbf{x}}^{(\tilde{b}_iq,s-1)}},\mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,s)}})$$
$$+ \tilde{\mathrm{F}}(\mathrm{P}_{\mathrm{G}^{b_i}},\mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,s)}}) \Big\}.$$
$$(24)$$

Fig. 1. The performance of the model when changing $\alpha$ in Eq. (20) of the paper.

| Methods | Log | Memory | N |
|---|---|---|---|
| VAE-ELBO-OCM-COS | -137.92 | 1.6K | 1 |
| VAE-ELBO-OCM | -132.07 | 1.6K | 1 |
| VAE-IWVAE50-OCM | -127.11 | 1.6K | 1 |
| Dynamic-ELBO-OCM | **-115.89** | 1.1K | 5 |

Eq. (24) describes the bound for a single task. We then extend this bound to the components learning more than one task:

$$
\sum_{i=1}^{|\mathcal{B}|} \left\{ \sum_{q=1}^{|\tilde{b}_i|} \left\{ \mathbb{E}_{\mathrm{P}_{\hat{\mathbf{x}}^{\tilde{b}_i^q}}} [\mathcal{L}_{ELBO}(\mathbf{x};\theta,\omega)] \right\} \right\} \leq
$$
$$
\sum_{i=1}^{|\mathcal{B}|} \left\{ \sum_{q=1}^{|\tilde{b}_i|} \left\{ \mathbb{E}_{\mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,c_i^q)}}} [\mathcal{L}_{ELBO}(\mathbf{x};\theta,\omega)] \right. \right.
$$
$$
+ \sum_{s=0}^{c_i^q} \left\{ 2\mathrm{W}_{\mathcal{L}}^{\star}(\mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,s)}}, \mathrm{P}_{\mathrm{G}^{b_i}})
$$
$$
\left. \left. - \mathrm{W}_{\mathcal{L}}^{\star}(\mathrm{P}_{\hat{\mathbf{x}}^{(\tilde{b}_i^q,s-1)}}, \mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,s)}}) + \tilde{\mathrm{F}}(\mathrm{P}_{\mathrm{G}^{b_i}}, \mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,s)}}) \right\} \right\} \right\}, \tag{25}
$$

We also extend the bound from Eq. (20) to components that would only learn one task each :

$$
\sum_{i=1}^{|\mathcal{A}|} \left\{ \mathbb{E}_{\mathrm{P}_{\hat{\mathbf{x}}^{\tilde{a}_i}}} [\mathcal{L}_{ELBO}(\mathbf{x};\theta,\omega)] \right\} \leq
$$
$$
\sum_{i=1}^{|\mathcal{A}|} \left\{ \mathbb{E}_{\mathrm{P}_{\mathbf{x}^{\tilde{a}_i}}} [\mathcal{L}_{ELBO}(\mathbf{x};\theta,\omega)] + 2\mathrm{W}_{\mathcal{L}}^{\star}(\mathrm{P}_{\mathbf{x}^{\tilde{a}_i}}, \mathrm{P}_{\mathrm{G}^{a_i}}) \right.
$$
$$
\left. - \mathrm{W}_{\mathcal{L}}^{\star}(\mathrm{P}_{\hat{\mathbf{x}}^{\tilde{a}_i}}, \mathrm{P}_{\mathbf{x}^{\tilde{a}_i}}) + \tilde{\mathrm{F}}(\mathrm{P}_{\mathrm{G}^{a_i}}, \mathrm{P}_{\mathbf{x}^{\tilde{a}_i}}) \right\}, \tag{26}
$$

Eventually, the bound for all components is defined by con-sidering both Eq. (25) and (26), resulting in :

$$
\sum_{i=1}^{|\mathcal{A}|} \left\{ \mathbb{E}_{\mathrm{P}_{\hat{\mathbf{x}}^{\tilde{a}_i}}} [\mathcal{L}_{ELBO}(\mathbf{x};\theta,\omega)] \right\}
$$
$$
+ \sum_{i=1}^{|\mathcal{B}|} \left\{ \sum_{q=1}^{|\tilde{b}_i|} \left\{ \mathbb{E}_{\mathrm{P}_{\hat{\mathbf{x}}^{\tilde{b}_i^q}}} [\mathcal{L}_{ELBO}(\mathbf{x};\theta,\omega)] \right\} \right\} \leq
$$
$$
\sum_{i=1}^{|\mathcal{A}|} \left\{ \mathbb{E}_{\mathrm{P}_{\mathbf{x}^{\tilde{a}_i}}} [\mathcal{L}_{ELBO}(\mathbf{x};\theta,\omega)] + 2\mathrm{W}_{\mathcal{L}}^{\star}(\mathrm{P}_{\mathbf{x}^{\tilde{a}_i}}, \mathrm{P}_{\mathrm{G}^{a_i}}) \right.
$$
$$
\left. - \mathrm{W}_{\mathcal{L}}^{\star}(\mathrm{P}_{\hat{\mathbf{x}}^{\tilde{a}_i}}, \mathrm{P}_{\mathbf{x}^{\tilde{a}_i}}) + \tilde{\mathrm{F}}(\mathrm{P}_{\mathrm{G}^{a_i}}, \mathrm{P}_{\mathbf{x}^{\tilde{a}_i}}) \right\} \tag{27}
$$
$$
+ \sum_{i=1}^{|\mathcal{B}|} \left\{ \sum_{q=1}^{|\tilde{b}_i|} \left\{ \mathbb{E}_{\mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,c_i^q)}}} [\mathcal{L}_{ELBO}(\mathbf{x};\theta,\omega)] \right. \right.
$$
$$
+ \sum_{s=0}^{c_i^q} \left\{ 2\mathrm{W}_{\mathcal{L}}^{\star}(\mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,s)}}, \mathrm{P}_{\mathrm{G}^{b_i}}) \right.
$$
$$
\left. \left. \left. - \mathrm{W}_{\mathcal{L}}^{\star}(\mathrm{P}_{\hat{\mathbf{x}}^{(\tilde{b}_i^q,s-1)}}, \mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,s)}}) + \tilde{\mathrm{F}}(\mathrm{P}_{\mathrm{G}^{b_i}}, \mathrm{P}_{\tilde{\mathbf{x}}^{(\tilde{b}_i^q,s)}}) \right\} \right\} \right\}
$$
$$
\square
$$

## APPENDIX D
### ADDITIONAL RESULTS ON ABLATION STUDY

**RBF kernel scale.** We investigate the performance of the pro-posed OCM framework when changing the hyperparameters of the RBF kernel in Eq. (20) from the paper. We vary the RBF scale $\alpha = \{5, 10, 20, 30, 50, 70, 100\}$ for the lifelong training a single VAE model trained with OCM under Split MNIST. The results presented in Fig. 1 indicate that OCM with $\alpha = 10$ achieves the best results.

**Using the cosine distance for sample selection.** We consider the cosine distance for evaluating the similarity in the proposed sample selection approach for LTM, instead of the graph based distance from Eq. (22) from the paper, defined as :

$$
\mathrm{R}^{C}(\mathbf{x}_{i,j}^e, \mathbf{x}_{i,u}^l) := \frac{\mathbf{z}_{i,j}^e \cdot \mathbf{z}_{i,u}^l}{\|\mathbf{z}_{i,j}^e\| \|\mathbf{z}_{i,u}^l\|}
$$
$$
= \frac{\sum_{i=1}^{d_z} \mathbf{z}_{i,j}^e(i) \mathbf{z}_{i,u}^l(i)}{\sqrt{\sum_{i=1}^{d_z} (\mathbf{z}_{i,j}^e(i))^2} \sqrt{\sum_{i=1}^{d_z} (\mathbf{z}_{i,u}^l(i))^2}}, \tag{28}
$$

where the evaluation of similarity is based on the latent features $\mathbf{z}_{i,u}^l$ and $\mathbf{z}_{i,j}^e$, corresponding to the data $\mathbf{x}_{i,j}^l$, $\mathbf{x}_{i,u}^e$, from LTM and STM, respectively.

We use "VAE-ELBO-OCM-COS" to represent a single VAE model trained with OCM, where the cosine distance is used as the criterion for the sample selection. Since a small measure in Eq. (26) means that $\mathbf{x}_{i,j}^e$ is far away from $\mathbf{x}_{i,u}^l$, we replace Eq. (23) by considering :

$$\mathrm{R}^C(\mathbf{x}_{i,j}^e, \mathbf{x}_{i,u}^l) < \lambda \Rightarrow \mathcal{M}_i^l = \mathcal{M}_i^l \cup \mathbf{x}_{i,j}^e, \qquad (29)$$

where we set $\lambda = 0$. The results of various models trained under Split MNIST are provided in Table I, showing that the proposed kernel from Eq. (22) for sample selection outperforms the cosine distance.

### REFERENCES