



# ED-GCAE: Efficient and Adaptive Disentanglement via Shared Features and Dynamic Noise Injection

Xingshen Zhang<sup>1</sup>, Hong Pan<sup>2</sup>, Bin Chai<sup>2</sup>, Lin Wang<sup>3,1</sup>, Bo Yang<sup>3,1</sup> and Shuangrong Liu<sup>1\*</sup>

<sup>1</sup> Shandong Key Laboratory of Ubiquitous Intelligent Computing, University of Jinan, Jinan 250022, China

<sup>2</sup> Shandong Institute for Product Quality Inspection, Jinan 250022, China

<sup>3</sup> Quan Cheng Laboratory, Jinan 250100, China

liusr.constant@gmail.com

**Abstract.** Autoencoder-based methods have become a dominant framework in disentangled representation learning. However, their reliance on simplistic Gaussian density estimation presents significant limitations to better disentanglement performance. The Gaussian Channel Autoencoder (GCAE) was introduced to address density estimation flexibility yet suffers from high computational costs due to its independent discriminator architecture and sensitivity to noise. To overcome these challenges, we propose ED-GCAE, a novel framework designed to improve the efficiency and dynamic adaptability of GCAE. ED-GCAE incorporates a shared feature extraction backbone into the discriminator architecture, significantly enhancing computational efficiency and training stability. Concurrently, we introduce a dynamic latent-variable-dependent noise injection mechanism to achieve the balance between disentanglement and stability. Experiments demonstrate that ED-GCAE demonstrates superior performance compared to baseline methods, achieving better disentangled representations while exhibiting enhanced training stability and computational efficiency.

**Keywords:** Disentanglement Representation Learning, Representation Learning, Deep Learning.

## 1 Introduction

Disentanglement representation learning [1] has been recognized as a promising paradigm for endowing machines with human-like perception and understanding of the world [2][3]. The objective is to automatically disentangle the data into distinct generative factors within the latent space, which subsequently serve to describe and represent the data's characteristics across different aspects. [1][4][5]. Disentanglement representations have proven beneficial for a mount of downstream tasks, abstract visual reasoning [6][7], image generation and manipulation [8][9], enhanced interpretability [10][11][12], and zero-shot domain adaptation [13][14][15].

To achieve generative factor independence, it is often necessary to ensure the statistical independence of latent variable dimensions. Conventional disentanglement methods, e.g., Variational Autoencoders (VAEs),  $\beta$ -VAE,  $\beta$ -TCVAE, largely operate

\*Corresponding author: Shuangrong Liu

within the information-theoretic framework, imposing constraints on the correlations between latent variable dimensions. Intrinsically, the implementation of correlation constraints within disentanglement learning methods relies fundamentally upon the accurate estimation of the data's probability density. However, these methods typically assume that the posterior and prior distributions of latent variables follow Gaussian distributions. The encoder learns the mean and variance of these Gaussian distributions, and the Kullback-Leibler (KL) divergence loss is employed to constrain the posterior distribution to approximate the prior. However, owing to the inherent diversity of problems, the underlying data distribution does not necessarily conform to a Gaussian assumption. The Gaussian distribution assumption inherent in these methods may lead to posterior approximation errors, consequently limiting model expressiveness and potentially resulting in the loss of crucial data information. This ultimately hinders the achievement of satisfactory disentanglement performance.

To address these challenges, Yeats et al. [17] introduced the Gaussian Channel Autoencoder (GCAE), which adopts a more flexible density estimation approach. GCAE employs multiple discriminators, with each discriminator tasked with estimating the conditional density of one latent dimension given observations of all other latent dimensions. Concurrently, by leveraging the density-ratio trick, GCAE achieves effective density estimation for unknown distributions.

Nevertheless, the GCAE method still has significant limitations. Since each discriminator is only responsible for estimating the conditional density of a single latent variable, the number of discriminators required for training and the data volume needed to train them scale rapidly with the dimensionality of the latent space, leading to prohibitively high training costs. Additionally, like most VAE-based disentanglement models, GCAE also faces instability in training performance, where different initial states lead to significant differences in the model's final disentanglement capability.

Therefore, we propose Efficient and Dynamic GCAE (ED-GCAE) to enhance the GCAE discriminator architecture by incorporating a shared feature extraction backbone mechanism. This integrates the original  $m$  independent discriminators into a multi-head discriminator with  $m$  independent output heads. By sharing the discriminator's backbone, we achieve more efficient and stable density estimation. Concurrently, to more finely balance the quality of disentangled representations and the stability of disentanglement, we design a latent-variable-dependent dynamic noise adjustment mechanism. This mechanism adaptively adjusts the capacity of noise injection based on the characteristics of different latent dimensions, thereby achieving a more nuanced and effective balance between disentanglement capability and stability.

Our contributions are summarized as follows:

- a) The ED-GCAE is proposed to enhance the disentanglement performance of autoencoder-based approaches, specifically improving training stability and efficiency through the introduction of noise perturbation and information sharing mechanisms.
- b) Dynamic noise injection adjustment mechanism is proposed to effectively balance disentanglement performance and stability.

- c) The proposed method is demonstrated to achieve enhanced disentanglement performance and greater stability by extensive experiments compared to the baseline methods.

## 2 Related Work

To achieve disentanglement in the latent space, current mainstream methodologies predominantly adhere to the information-theoretic framework. These approaches commonly impose various constraints to diminish the statistical dependencies among latent variable dimensions, thereby encouraging the posterior distribution of latent variables to approximate a factorized form. Within these methodologies, the Variational Autoencoder (VAE) [16] framework and its variants have emerged as dominant paradigms, owing to the VAE's provision of an effective generative model learning framework and its capacity to handle intricate posterior distributions through variational inference techniques.

For instance,  $\beta$ -VAE [13] enhances the constraint imposed by the prior distribution on the latent space by introducing a hyperparameter  $\beta$  into the VAE's Evidence Lower Bound (ELBO) loss function. This indirectly promotes the independence of latent variable dimensions [18].  $\beta$ -TCVAE [19], through the decomposition of the ELBO loss function in conjunction with importance weighted sampling techniques, explicitly minimizes the Total Correlation (TC) of latent representations, further bolstering disentanglement efficacy. FactorVAE [20], also committed to reducing the total correlation of latent representations, innovatively employs adversarial training and the density-ratio trick. This approach leverages a discriminator network to estimate total correlation, circumventing the challenges associated with directly computing high-dimensional joint distributions. The DIP-VAE [21] methods, conversely, directly regularizes the covariance matrix of latent representations. By constraining the covariance matrix to approximate a diagonal matrix, DIP-VAE explicitly encourages the independence of latent variable dimensions. DynamicVAE [22] employs an incremental Proportional-Integral (PI) controller alongside moving average techniques to dynamically adjust the regularization coefficient  $\beta$  in  $\beta$ -VAE. This method progressively anneals the  $\beta$  value from an initially larger value ( $\beta > 1$ ) to a smaller value ( $\beta \leq 1$ ), mitigating the issue of degraded reconstruction quality that often accompanies the enhanced disentanglement achieved through setting a larger  $\beta$  value in  $\beta$ -VAE.

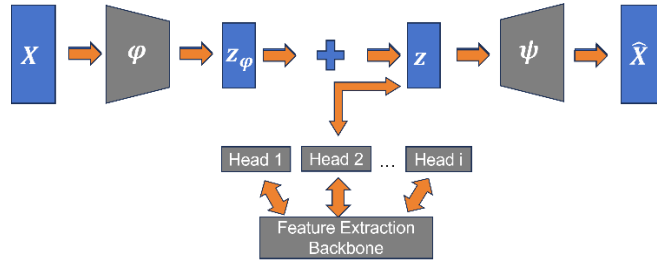
The above VAEs predominantly adopt parametric density estimation strategies, typically assuming that the posterior and prior distributions of latent variables conform to simplistic Gaussian distributions. While Gaussian distributions offer mathematically tractable properties, empirical evidence suggests that data posterior distributions are often significantly more complex than Gaussian approximations. To address the limitations imposed by the Gaussian distribution assumption in traditional VAE methods, Yeats et al. [17] introduced the Gaussian Channel Autoencoder (GCAE). GCAE dispenses with the Gaussian posterior assumption and instead employs a discriminator-based conditional density estimation approach. GCAE leverages multiple discriminators, each dedicated to estimating the conditional density of a single latent dimension

given the observations of all other latent dimensions. Furthermore, by employing the density-ratio trick, GCAE achieves effective density estimation for unknown distributions, thereby affording the model enhanced density estimation flexibility.

### 3 Methodology

In the original GCAE, discriminators were employed to estimate the conditional densities of latent variables, deviating from the prevalent practice in many VAE-based methods of directly using Gaussian distributions as posterior distributions. This departure allowed for more flexible density estimation. However, as this method needs to estimate density for each latent variable, the original GCAE architecture, corresponding to the latent space dimensionality  $m$ , utilized  $m$  mutually independent discriminators, each responsible for estimating the conditional density of a single latent variable. While this design potentially enabled more accurate estimation of conditional densities for different latent variables, it substantially increased the training time and data volume required for discriminators. For instance, according to the original paper's recommended settings, each training step of the entire framework required five separate training steps for the discriminators, resulting in suboptimal overall training efficiency.

Simultaneously, to ensure the continuity and smoothness of the latent space, the original GCAE incorporated Gaussian noise with a fixed parameter sigma for each latent variable. This was intended to prevent the dimensions of the latent variables from collapsing into sharp, peaked distributions, which could impede training. However, this approach presents several limitations. Firstly, the disentanglement performance of this method is sensitive to sigma. As demonstrated in the original experimental results, when sigma is within the range of 0.1-0.3, different sigma values often result in significant variations across various disentanglement metrics, with differences reaching up to approximately fivefold. Secondly, even with a constant sigma value, repeated experiments frequently yielded results with significant gaps in disentanglement effectiveness. The original experiments confirmed this phenomenon across multiple metrics and datasets, revealing considerable variance in disentanglement performance for the same sigma value over several trials, which points to an underlying instability problem with the method.



**Fig. 1.** Framework of ED-GCAE.

To mitigate the computational cost associated with discriminator training, we introduce a shared feature extraction backbone mechanism. This approach decouples the

original discriminators into a backbone component and prediction heads. The density estimation for each latent variable shares the backbone component, while being complemented by  $m$  independent prediction heads, thereby enabling efficient discriminator training. Recognizing the inherent limitations of a fixed sigma, we design a dynamic sigma adjustment strategy. This strategy allows each dimension of the latent variables to dynamically adjust the Gaussian noise magnitude based on its information entropy, thus improving disentanglement performance while maintaining robust reconstruction capabilities.

### 3.1 Shared Feature Extraction Backbone Multi-Head Discriminator

To enhance the efficiency and scalability of discriminator training within the GCAE framework, and to overcome the inherent parameter redundancy and computational bottlenecks of the independent discriminator architecture, we propose a shared feature extraction backbone multi-head discriminator architecture (hereafter referred to as the shared discriminator architecture). In the original GCAE, the conditional density ratio  $D_i(z_i, z_{\setminus i})$  for each latent dimension  $z_i$  was modeled by an independent discriminator  $D_i$ , leading to parameter scaling and computational complexity that increased linearly with the latent space dimensionality  $m$ . To address this, we decouple the discriminator's functionality and introduce a sharing mechanism, Shared Feature Extraction Backbone  $F$ , which serves as a generic feature encoder module for the discriminators. Through a neural network  $F(\cdot)$ , it maps the input latent variable  $z$  to a shared, high-dimensional feature space. The parameters  $\theta_F$  of module  $F$  are shared across all conditional density estimation tasks, realizing feature reuse and knowledge transfer.

For each latent dimension  $z_i$ , we design a lightweight, independent prediction head  $H_i(\cdot)$ , with parameters  $\theta_{H_i}$ . Prediction head  $H_i$  receives the feature representation output by the shared feature extraction backbone  $F$  and further maps it to a scalar output, approximating the conditional density ratio of the  $i$ -th latent dimension. Thus, the function of the  $i$ -th discriminator  $D_i$  can be expressed as:

$$D_i(z) = H_i(F(z); \theta_{H_i}; \theta_F) \quad (1)$$

where  $\theta_F$  are the parameters of the shared feature extraction backbone  $F$ , and  $\theta_{H_i}$  are the parameters of the  $i$ -th prediction head  $H_i$ . Through the parameter-sharing strategy, the parameter count and computational complexity of the discriminator architecture are significantly reduced, enabling efficient and scalable conditional density estimation.

### 3.2 Latent-Variable Dependent Dynamic Sigma Adjustment Strategy

To overcome the hyperparameter sensitivity to the fixed Gaussian noise scale  $\sigma$  in the original GCAE method, and to adaptively balance disentanglement capability and reconstruction quality, we propose a latent-variable dependent dynamic Sigma adjustment strategy. This strategy aims to dynamically adjust the scale  $\sigma_i$  of the injected Gaussian noise based on the information content of each latent dimension  $z_i$ , achieving

fine-grained noise control. Specifically, we integrate simulated annealing and exponential moving average techniques, and introduce an information entropy-guided adaptive adjustment mechanism:

### Information Entropy-Guided Adaptation

The design of the information entropy mapping function  $g(\cdot)$  is important, as it demonstrates how  $\sigma_i$  is adjusted based on  $H(z_i)$ . To implement information entropy-guided adaptive noise scale adjustment, we first perform normalization on the estimated information entropy values of each latent dimension. This ensures that the information entropy values across different dimensions possess comparability and a uniform scale. Specifically, we employ min-max normalization, mapping the original information entropy values to a standardized range of  $[0, 1]$ . For each training batch, we compute the minimum  $H_{min}$  and maximum  $H_{max}$  values of the information entropy estimates  $H(z_i)$  for all latent dimensions. Subsequently, the normalized information entropy  $\hat{H}(z_i)$  for each latent dimension  $z_i$  is calculated according to the following formula:

$$\hat{H}(z_i) = \frac{(H(z_i) - H_{min})}{(H_{max} - H_{min})} \quad (2)$$

where  $H(z_i)$  represents the original information entropy estimate of latent dimension  $z_i$ ,  $H_{min}$  and  $H_{max}$  represent the minimum and maximum values of information entropy estimates for all latent dimensions within the current batch, respectively, and  $\hat{H}(z_i)$  is the normalized information entropy value corresponding to dimension  $z_i$ . This normalization process ensures that the information entropy values are bound within the  $[0, 1]$  range, thereby providing a standardized control signal for subsequent dynamic noise scale adjustment. This normalization strategy eliminates discrepancies in the units and scales of information entropy values across different latent dimensions, enabling us to more effectively leverage information entropy values to guide the adaptive adjustment of noise scale.

### Exponential Moving Average Update

For each latent dimension  $z_i$ , its Gaussian noise scale  $\sigma_i$  is dynamically adjusted in each training iteration  $t$  according to the following EMA update rule:

$$\sigma_i^{(t)} = \alpha_{EMA} \cdot \sigma_i^{(t-1)} + (1 - \alpha_{EMA}) \cdot g\left(H(z_i^{(t-1)})\right) \quad (3)$$

where  $\sigma_i^{(t)}$  denotes the noise scale at iteration  $t$ ,  $\alpha \in [0, 1]$  is the smoothing factor,  $H(z_i^{(t-1)})$  is the estimated information entropy of latent dimension  $z_i$  at iteration  $(t - 1)$ , and  $g(\cdot)$  represents the information entropy mapping function, which maps the information entropy value to an appropriate noise scale adjustment amount.

### Simulated Annealing

Employing a simulated annealing strategy to anneal  $\sigma$ , the overall coefficient  $\gamma$  is initialized to a value approximating 1 and progressively decreases with increasing

training iterations. This annealing strategy imbues  $\sigma_i$  with the capacity for rapid adjustment in the nascent stages of training, helping the identification of optimization directions favorable for disentanglement. As training deepens and  $\gamma$  diminishes, the responsiveness of  $\sigma_i$  to information entropy fluctuations becomes more smoothly, enabling more fine-tuned adaptive adjustments and mitigating abrupt shifts in the optimization trajectory.

The overall formula for dynamically adjusting Sigma is presented below:

$$\sigma_i^{(t)} = \max \left( \sigma_{min}, \left[ \alpha_{EMA} \cdot \sigma_i^{(t-1)} + (1 - \alpha_{EMA}) \cdot \left( \sigma_{max} - (\sigma_{max} - \sigma_{min}) \cdot \frac{H(z_i^{(t)}) - H_{min}^{(t)}}{H_{max}^{(t)} - H_{min}^{(t)}} \right) \right] \cdot \gamma^g \right) \quad (4)$$

## 4 Experiments

In our experimental evaluation, we primarily focus on comparing ED-GCAE against the original GCAE [17],  $\beta$ -VAE [13],  $\beta$ -TCVAE [19], FactorVAE [20], DIP-VAE-II [21] to assess the performance enhancements afforded by our proposed improvements. The disentanglement performance of our proposed method and comparative methods are quantitatively evaluated using four widely recognized supervised metrics: Mutual Information Gap (MIG) [19], Factor Score (FAC) [20], DCI Disentanglement (DCI) [23], and Separated Attribute Predictability (SAP) [24]. These metrics, collectively representing the three major categories of disentanglement assessment delineated by Carbonneau et al. [25], facilitate a rigorous and comprehensive comparative analysis of disentanglement capabilities.

We consider two datasets which cover different data modalities. The Beamsynthesis dataset [26] is a collection of 360 timeseries data from a linear particle accelerator beamforming simulation. The waveforms are 1000 values long and are made of two independent data generating factors: duty cycle (continuous) and frequency (categorical). The dSprites dataset [27] is a collection of 737280 synthetic images of simple white shapes on a black background. Each  $64 \times 64$  pixel image consists of a single shape generated from the following independent factors: shape (categorical), scale (continuous), orientation (continuous), x-position (continuous), and y-position (continuous).

### 4.1 Comparison of ED-GCAE with GCAE and Baseline

Table 1 presents a comparative evaluation of ED-GCAE against several VAE baselines and the original GCAE across various disentanglement metrics on dSprites dataset. For GCAE, hyperparameters were set to  $\sigma = 0.2, k = 5$ , while ED-GCAE was configured with  $\sigma_{max} = 0.4, \sigma_{min} = 0.2, \alpha_{EMA} = 0.1, \gamma = 0.999, k = 1$ . Comparison with VAE baselines reveals that ED-GCAE achieves performance superior to or comparable with these baselines across all metrics. Furthermore, in direct comparison to

GCAE, our proposed method demonstrates superior mean performance across nearly all metrics. Notably, the original GCAE exhibits generally higher variance in metric scores, indicating inherent instability. Through the implementation of dynamic noise injection adjustment, our method significantly reduces this variance, thereby achieving enhanced training stability.

Table 2 presents a corresponding comparative analysis on the Beamsynthesis dataset, consistently demonstrating the superior performance of our proposed method, exhibiting enhanced disentanglement capability compared to the original GCAE.

**Table 1.** Disentanglement metric comparison of ED-GCAE with VAE baselines and original GCAE on dSprites dataset.

	MIG	FAC	DCI	SAP
ED-GCAE	<b><math>0.358 \pm 0.011</math></b>	$0.584 \pm 0.007$	$0.220 \pm 0.006$	<b><math>0.586 \pm 0.045</math></b>
GCAE	$0.274 \pm 0.093$	$0.554 \pm 0.073$	$0.173 \pm 0.065$	$0.579 \pm 0.022$
$\beta$ -VAE	$0.352 \pm 0.010$	$0.600 \pm 0.006$	<b><math>0.293 \pm 0.059</math></b>	$0.280 \pm 0.066$
$\beta$ -TCVAE	$0.321 \pm 0.016$	$0.562 \pm 0.060$	$0.260 \pm 0.053$	$0.244 \pm 0.052$
FactorVAE	$0.162 \pm 0.046$	<b><math>0.714 \pm 0.002</math></b>	$0.101 \pm 0.033$	$0.342 \pm 0.026$
DIP-VAE-II	$0.025 \pm 0.001$	$0.424 \pm 0.020$	$0.022 \pm 0.001$	$0.081 \pm 0.046$

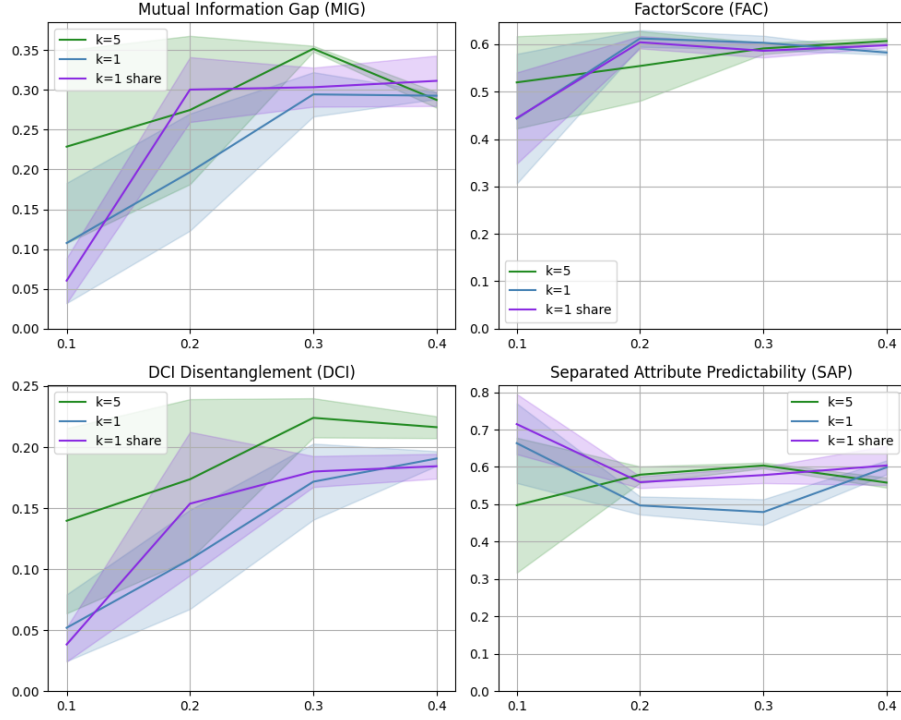
**Table 2.** Disentanglement metric comparison of ED-GCAE with VAE baselines and original GCAE on Beamsynthesis dataset.

	MIG	FAC	DCI	SAP
ED-GCAE	<b><math>0.369 \pm 0.040</math></b>	<b><math>0.989 \pm 0.008</math></b>	<b><math>0.549 \pm 0.015</math></b>	<b><math>0.395 \pm 0.094</math></b>
GCAE	$0.291 \pm 0.052$	$0.932 \pm 0.062$	$0.312 \pm 0.026$	$0.332 \pm 0.082$
$\beta$ -VAE	$0.142 \pm 0.044$	$0.981 \pm 0.011$	$0.512 \pm 0.057$	$0.152 \pm 0.034$
$\beta$ -TCVAE	$0.238 \pm 0.062$	$0.986 \pm 0.006$	$0.423 \pm 0.087$	$0.225 \pm 0.025$
FactorVAE	$0.153 \pm 0.051$	$0.946 \pm 0.052$	$0.424 \pm 0.041$	$0.162 \pm 0.021$
DIP-VAE-II	$0.082 \pm 0.023$	$0.824 \pm 0.047$	$0.376 \pm 0.035$	$0.163 \pm 0.057$

## 4.2 Effect of Shared Discriminator

Fig 2 presents a comparative analysis of the model incorporating the shared discriminator architecture against the original GCAE. In the original GCAE, we evaluated configurations with  $k=5$  and  $k=1$ , where  $k$  denotes the number of discriminator iterations per training step of the overall framework. The shared discriminator model was configured with  $k=1$ . All models were trained for a total of 20,000 steps. As depicted in Figure 1, when comparing different  $\sigma$  values within the original GCAE framework, a substantial performance gap in disentanglement metrics is evident between the  $k=1$  and  $k=5$  configurations. In contrast, the shared discriminator model with  $k=1$  demonstrates disentanglement performance generally superior to that of the original GCAE with  $k=1$ , and in certain metrics, it approximates or even surpasses the performance of the original

GCAE with  $k=5$ . This observation underscores the capacity of the shared discriminator architecture to maintain robust disentanglement performance while requiring fewer discriminator training iterations per step, thereby significantly enhancing training efficiency.

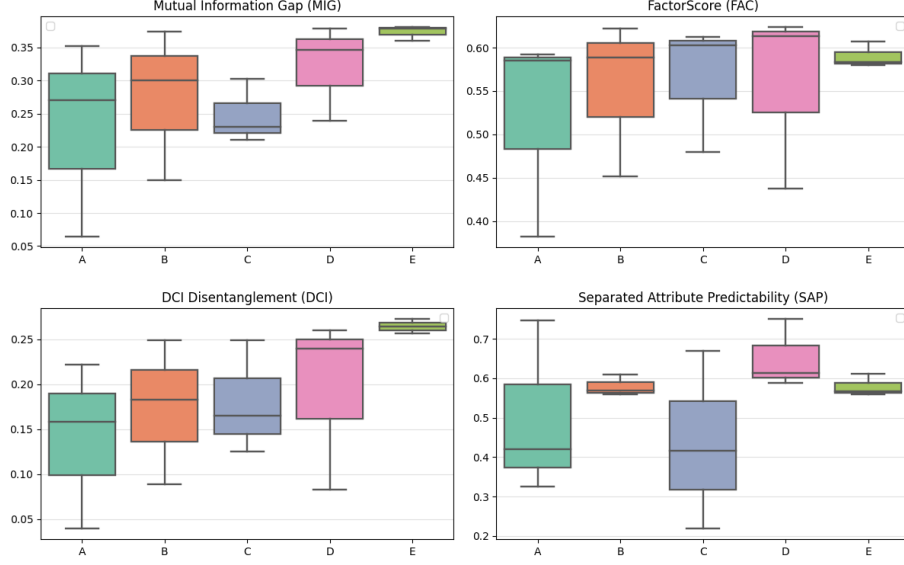


**Fig. 2.** Comparison of original GCAE with varying  $k$  values and whether using a shared discriminator across disentanglement metrics. The x-axis depicts the sigma value employed by the model, while the y-axis depicts the value of the corresponding disentanglement metric.

### 4.3 Effect of Dynamic Sigma Adjustment Strategy

Fig 3 presents a comparative analysis of the dynamic sigma adjustment strategy against the fixed sigma approach. For this comparison, both model configurations were set to  $k=5$ . We observe that, in contrast to  $\sigma=0.1$ , configurations with fixed sigma values within the range of  $[0.1, 0.3]$  exhibit reduced variance and enhanced stability across disentanglement metrics. Furthermore, when comparing  $\sigma=0.2$  with the dynamic sigma configurations ( $\sigma \in [0.2, 0.3]$  and  $\sigma \in [0.2, 0.4]$ ), the models employing dynamic sigma demonstrate a discernible improvement in disentanglement performance, coupled with superior stability. Notably, the configuration with  $\sigma \in [0.2, 0.4]$  particularly

excels, showcasing both exceptional disentanglement efficacy and robustness.



**Fig. 3.** Comparison between fixed  $\sigma$  and dynamic  $\sigma$ . A on x-axis means fixed  $\sigma = 0.1$ , B means  $\sigma = 0.2$ , C means  $\sigma$  changes in range of  $[0.1, 0.3]$ , D means  $\sigma$  changes in  $[0.2, 0.3]$ , E means  $\sigma$  changes in  $[0.2, 0.4]$

Conversely, the fixed sigma approach generally displays higher variance across most scenarios, highlighting an inherent instability in disentanglement performance. We posit that a plausible explanation for this instability lies in the fact that, during later stages of training when the information entropy of latent dimensions is diminished, the fixed sigma approach continues to inject a substantial level of Gaussian noise. This persistent noise injection may inadvertently introduce excessive perturbation to the increasingly stable training dynamics, ultimately contributing to the observed disentanglement instability. In contrast, the dynamic sigma adjustment strategy, through the integration of simulated annealing and Exponential Moving Average (EMA), ensures a more gradual and tempered variation in  $\sigma$ , while concurrently facilitating a progressive reduction in  $\sigma$  during later training phases. This dynamic adaptation contributes to the enhanced stability of the model.

#### 4.4 Ablation Experiments

To further indicate the individual contributions of the simulated annealing strategy and Exponential Moving Average (EMA) to the overall training efficacy, we conducted a comparative analysis. Focusing on the previously identified efficacious sigma range of  $[0.2, 0.4]$  and maintaining  $k=5$ , we performed two distinct experimental groups: one implementing the GCAE architecture augmented with EMA (hereafter referred to as GCAE+EMA), and the other employing GCAE integrated with simulated annealing

(GCAE+Annealing). The baseline for this comparison was the original GCAE configuration with a fixed  $\sigma=0.2$ . Table 3 reveals that, in comparison to the original GCAE baseline, GCAE+EMA demonstrates robust disentanglement capability, maintaining a high level of disentanglement performance while concurrently exhibiting minimal variance, indicative of enhanced stability. GCAE+Annealing, in contrast, showcases even more pronounced disentanglement efficacy, albeit with a slightly diminished stability compared to GCAE+EMA. Notably, ED-GCAE, which integrates both EMA and simulated annealing, achieves a compelling synthesis, exhibiting both superior disentanglement capability and sustained stability.

**Table 3.** Ablation study on Exponential Moving Average (EMA) and Simulated Annealing

	MIG	FAC	DCI	SAP
GCAE	$0.274 \pm 0.092$	$0.554 \pm 0.073$	$0.173 \pm 0.066$	$0.579 \pm 0.022$
GCAE+EMA	$0.333 \pm 0.002$	<b><math>0.604 \pm 0.015</math></b>	$0.242 \pm 0.004$	$0.604 \pm 0.009$
GCAE+Annealing	$0.363 \pm 0.022$	$0.598 \pm 0.020$	$0.253 \pm 0.012$	$0.605 \pm 0.027$
ED-GCAE	<b><math>0.373 \pm 0.009</math></b>	$0.590 \pm 0.012$	<b><math>0.264 \pm 0.006</math></b>	<b><math>0.613 \pm 0.024</math></b>

## 5 Conclusion

In conclusion, we have presented ED-GCAE, an innovative approach that leverages a shared backbone discriminator to alleviate the data and training steps requirements inherent in density estimation architectures. Furthermore, through dynamic noise injection refinement, ED-GCAE achieves superior disentanglement capability and enhanced stability compared to GCAE. Concurrently, ED-GCAE consistently outperforms existing baseline models across a comprehensive suite of disentanglement metrics.

**Acknowledgement.** This study was supported by National Natural Science Foundation of China under Grant No. 61872419, No. 62072213. Shandong Provincial Natural Science Foundation No. ZR2022JQ30, No. ZR2022ZD01, No. ZR2023LZH015. Taishan Scholars Program of Shandong Province, China, under Grant No. tsqn201812077. The "New 20 Rules for University" Program of Jinan City under Grant No. 2021GXRC077. Key Research Project of Quancheng Laboratory, China under Grant No. QCLZD202303. Research Project of Provincial Laboratory of Shandong, China under Grant No. SYS202201. University of Jinan Disciplinary Cross-Convergence Construction Project 2023 No. XKJC-202303, University of Jinan Disciplinary Cross-Convergence Construction Project 2024 No. XKJC-202402. National "Challenge-Based" Major Science and Technology Project in the Building Materials Industry Project No. 2023JBGS11-03. Key Research and Development Program of Shandong Province Grant No. 2024CXPT084. High-performance Computing Platform at University of Jinan.

## Reference

1. Bengio Y.: Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence* 35(8), 1798-1828 (2013)
2. Bengio Y.: Scaling learning algorithms towards AI. *Large-scale kernel machines* 34(5), 1-41 (2007)
3. Schmidhuber J.: Learning factorial codes by predictability minimization. *Neural Computation*, 4(6), 863-879 (1992)
4. Kulkarni, T. D., Whitney, W. F.: Deep convolutional inverse graphics network. *Advances in neural information processing systems*, pp. 2539-2547, (2015)
5. Chen X, Duan Y, Houthoofd R, et al.: Infogan: Interpretable representation learning by information maximizing generative adversarial nets. *Advances in neural information processing systems*, 29, (2016)
6. Van Steenkiste S, Locatello F, Schmidhuber J, et al.: Are disentangled representations helpful for abstract visual reasoning?. *Advances in neural information processing systems*, 32, (2019)
7. Suter R, Miladinovic D, Schölkopf B, et al.: Robustly disentangled causal mechanisms: Validating deep representations for interventional robustness. *International Conference on Machine Learning*, PMLR, pp. 6056-6065. (2019)
8. Lample G, Zeghidour N, Usunier N, et al.: Fader networks: Manipulating images by sliding attributes. *Advances in neural information processing systems*, 30, (2017)
9. Kulkarni T D, Whitney W F, Kohli P, et al.: Deep convolutional inverse graphics network. *Advances in neural information processing systems*, 28 (2015)
10. Tran L, Yin X, Liu X.: Disentangled representation learning gan for pose-invariant face recognition. *Proceedings of the IEEE conference on computer vision and pattern recognition*, 1415-1424 (2017)
11. Lee H Y, Tseng H Y, Huang J B, et al.: Diverse image-to-image translation via disentangled representations. *Proceedings of the European conference on computer vision (ECCV)*, 35-51, (2018)
12. Xing X, Han T, Gao R, et al.: Unsupervised disentangling of appearance and geometry by deformable generator network. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 10354-10363, (2019)
13. Higgins I, Matthey L, Pal A, et al.: beta-vae: Learning basic visual concepts with a constrained variational framework. *International conference on learning representations* (2017)
14. Cao J, Katzir O, Jiang P, et al.: DiDA: Iterative boosting of disentangled synthesis and domain adaptation. *2021 11th International Conference on Information Technology in Medicine and Education (ITME)*, IEEE, 201-208, (2021)
15. Peng X, Huang Z, Sun X, et al.: Domain agnostic learning with disentangled representations. *International conference on machine learning*. PMLR, 5102-5112, (2019)
16. Kingma D P, Welling M.: Auto-encoding variational bayes, (2013-12-20)
17. Yeats E, Liu F Y, Li H.: Disentangling Learning Representations with Density Estimation, *The Eleventh International Conference on Learning Representations*, (2023)
18. Burgess C P, Higgins I, Pal A, et al.: Understanding disentangling in beta-VAE. *arXiv preprint arXiv:1804.03599*, (2018)
19. Chen R T Q, Li X, Grosse R B, et al.: Isolating sources of disentanglement in variational autoencoders. *Advances in neural information processing systems*, 31, (2018)
20. Kim H, Mnih A.: Disentangling by factorizing. *International conference on machine learning*, PMLR, 2649-2658, (2018)



21. Kumar A, Sattigeri P, Balakrishnan A.: Variational Inference of Disentangled Latent Concepts from Unlabeled Observations. International Conference on Learning Representations, (2018)
22. Shao H, Lin H, Yang Q, et al.: Dynamicvae: Decoupling reconstruction error and disentangled representation learning. arXiv preprint arXiv:2009.06795, (2020)
23. Eastwood C, Williams C K I.: A framework for the quantitative evaluation of disentangled representations. 6th International Conference on Learning Representations, (2018)
24. Kumar A, Sattigeri P, Balakrishnan A.: Variational Inference of Disentangled Latent Concepts from Unlabeled Observations. International Conference on Learning Representations, (2018)
25. Carbonneau M A, Zaidi J, Boilard J, et al.: Measuring disentanglement: A review of metrics. IEEE transactions on neural networks and learning systems, 35(7), 8747-8761, (2022)
26. Yeats E, Liu F, Womble D, et al.: Nashae: Disentangling representations through adversarial covariance minimization. European Conference on Computer Vision, Cham: Springer Nature Switzerland, 36-51, (2022)
27. Loic Matthey, Irina Higgins, Demis Hassabis, and Alexander Lerchner: dsprites: Disentanglement testing sprites dataset, <https://github.com/deepmind/dsprites-dataset/>, 2017.