

## InfoGAN

**InfoGAN** [1] learns disentangled and interpretable representations by maximizing the mutual information between a subset of the latent variables and the GAN generated sample. This is done through the addition of an extra term to the objective function.

### Background: GANs and Mutual Information

**GANs** are trained by a two-player minimax game between Discriminator  $D$  and Generator  $G$  with value function  $V_{GAN}(D, G)$ :

$$\min_G \max_D V_{GAN}(D, G) = \mathbb{E}_{x \sim p_{\text{real}}(x)}[\log D(x)] + \mathbb{E}_{z \sim p_z(z)}[\log(1 - D(G(z)))]$$

The **mutual information** (MI)  $I(X; Y)$  between random variables  $X$  (all images) and  $Y$  (real/fake label) is:

$$I(X; Y) = H(X) - H(X|Y) = D_{KL}(P_{(X,Y)} || P_X \otimes P_Y)$$

As the MI is intractable, we use a **variational lower bound** on  $I(c; G(z, c))$ :

$$I(c; G(z, c)) \geq \mathbb{E}_{x \sim G(z, c)}[\mathbb{E}_{c' \sim P(c|x)}[\log Q(c'|x)]] + H(c) = L_I(G, Q)$$

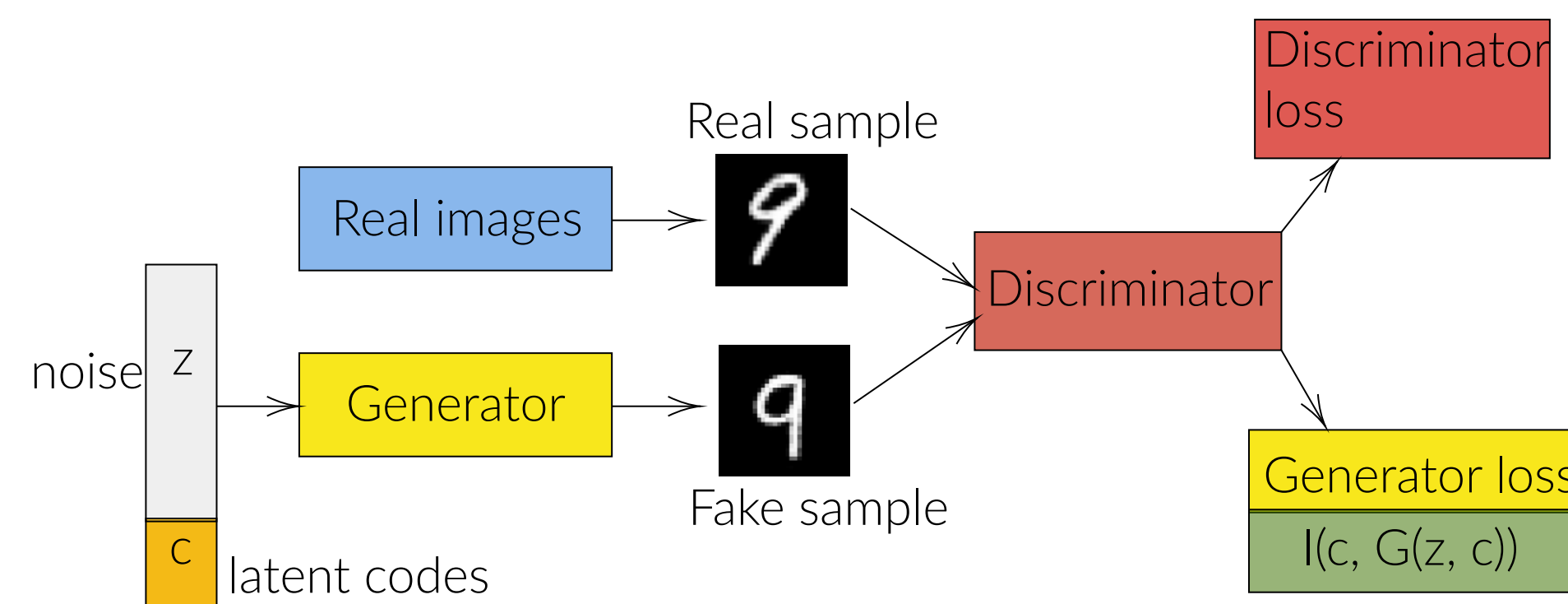
In GANs  $G$  **minimizes the same variational lower bound** on  $I(X; Y)$  [5]:

$$\begin{aligned} I(X; Y) &\geq \mathbb{E}_{x \sim p_{\text{all images}}(x)} \mathbb{E}_{y \sim p_{\text{is real}}(y|x)}[\log q(y|x)] + H(Y) \\ &= \mathbb{E}_{x \sim p_{\text{real}}(x)}[\log q(y=1|x)] + \mathbb{E}_{x \sim p_{\text{fake}}(x, z)}[\log(1 - q(y=1|x))] + H(Y) \end{aligned}$$

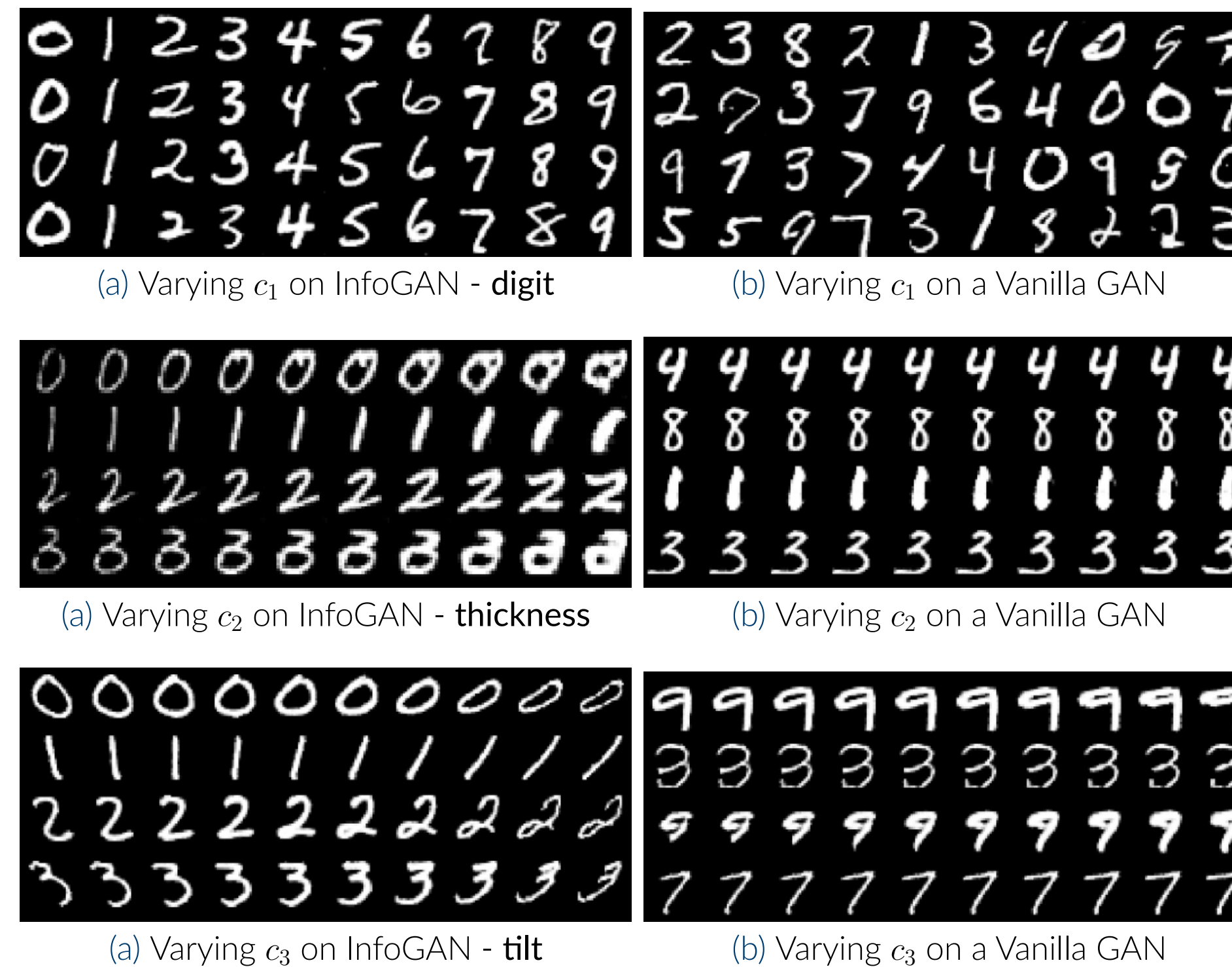
### Method: Information Regularizing GANs

Maximize the MI between Generator output  $G(z, c)$  and latent codes  $c$  as a regularizer:

$$\min_G \max_D V_{I\text{-GAN}}(D, G) = V_{GAN}(D, G) - \lambda L_I(G, Q)$$



## InfoGAN vs Vanilla GAN



## Stability Analysis and Mutual Information

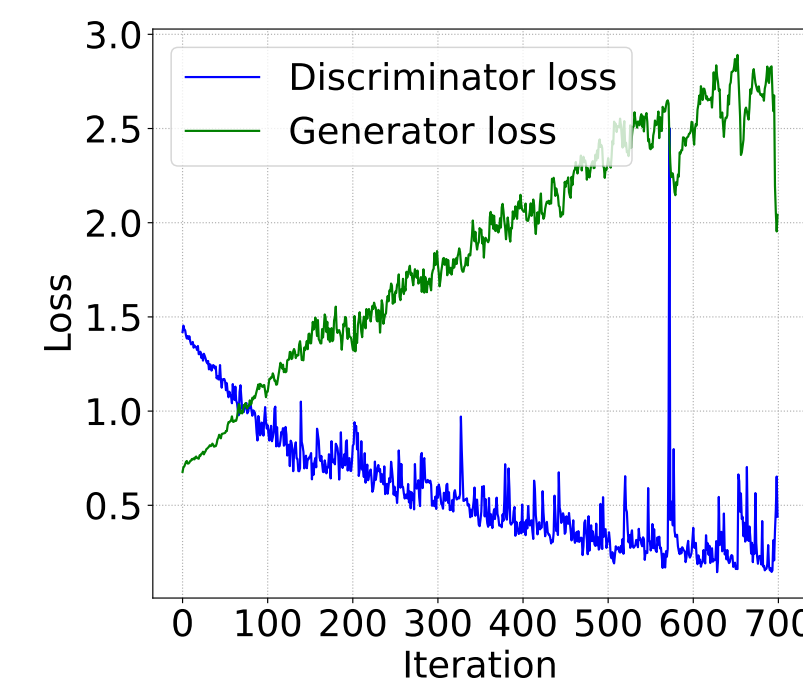


Figure 1. Training loss curves for both networks in our InfoGAN for MNIST.

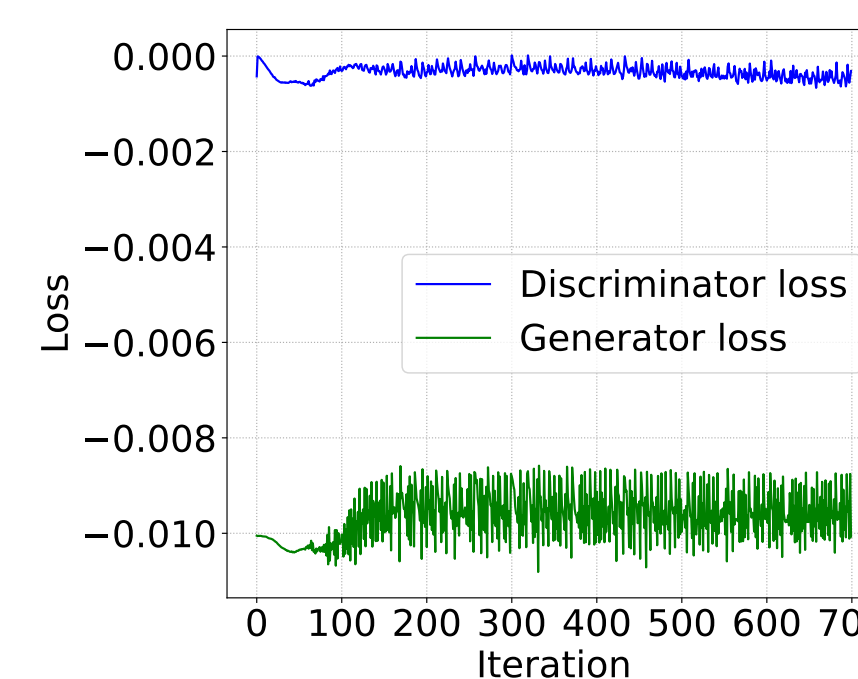


Figure 2. Training loss curves for both networks in our Info-WGAN for MNIST.

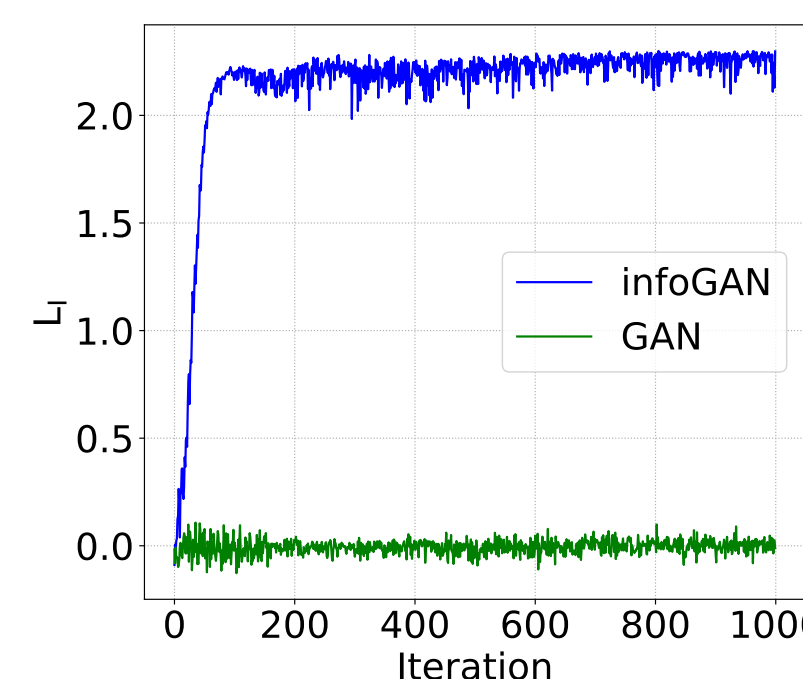


Figure 3.  $L_I$  for discrete code  $c_1$ .

## Info-WGANs

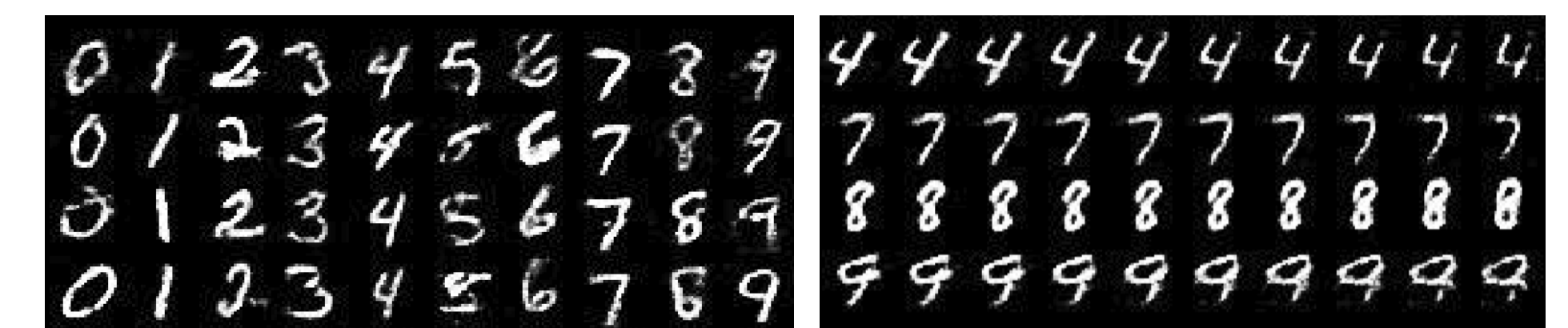
**Wasserstein GANs** (WGANs) [2] optimise Wasserstein distance. Benefits:

1. This metric is less prone to model collapse and vanishing gradients.
2. Introduction of weight clipping to enforce Lipschitz constraint.

$$\min_G \max_D V_{WGAN}(D, G) = \mathbb{E}_{x \sim p_{\text{real}}}[D(x)] - \mathbb{E}_{z \sim p(z)}[D(G(z))]$$

↓

$$\min_G \max_D V_{I\text{-WGAN}}(D, G) = V_{WGAN}(D, G) - \lambda L_I(G, Q)$$



(a) Varying  $c_1$  on Info-WGAN - digit (b) Varying  $c_{\text{cont}}$  on Info-WGAN - ?

**Findings:** Info-WGAN performs satisfactorily with discrete latent codes although it finds difficulties interpreting the continuous ones.

## MINE + GANs

**MINE** (Mutual Information Neural Estimator) [3] is a lower bound on the MI, obtained from the Donsker-Varadhan representation of the KL divergence by restricting function  $T$  to be parametrized by a neural net.

$$D_{KL}(P_{(X,Y)} || P_X \otimes P_Y) \geq \sup_{\theta \in \Theta} \mathbb{E}_{P_{X,Y}}[T_\theta] - \log \mathbb{E}_{P_X \otimes P_Y}[e^{T_\theta}]$$



(a) Varying  $c_2$  on MineGAN - thickness (b) Varying  $c_3$  on MineGAN - tilt

**Findings:** MineGAN is only able to learn the continuous codes and is harder to train: tricks were needed to stabilize the MI.

## Future work: Disentanglement of VAEs

**Variational Autoencoders (VAEs)** [4]:

- Have a more continuous and smooth latent space.
- Provide a more structured and interpretable latent space.
- Can perform interpolation in the latent space.

$$L_{VAE+c} = \mathbb{E}_{c \sim Q(c|x), z \sim Q(z|x)} \left[ \ln \frac{P(z)P(c)P(x|z, c)}{Q(c|x)Q(z|x)} \right] \rightarrow L_{I\text{-VAE}} = L_{VAE} + \lambda L_I$$

- Can we predict the appropriate number of latent codes in an unsupervised manner for different datasets?