

Soft-IntroVAE: Analyzing and Improving the Introspective Variational Autoencoder

CVPR 2021 (Oral)

Tal Daniel, Aviv Tamar

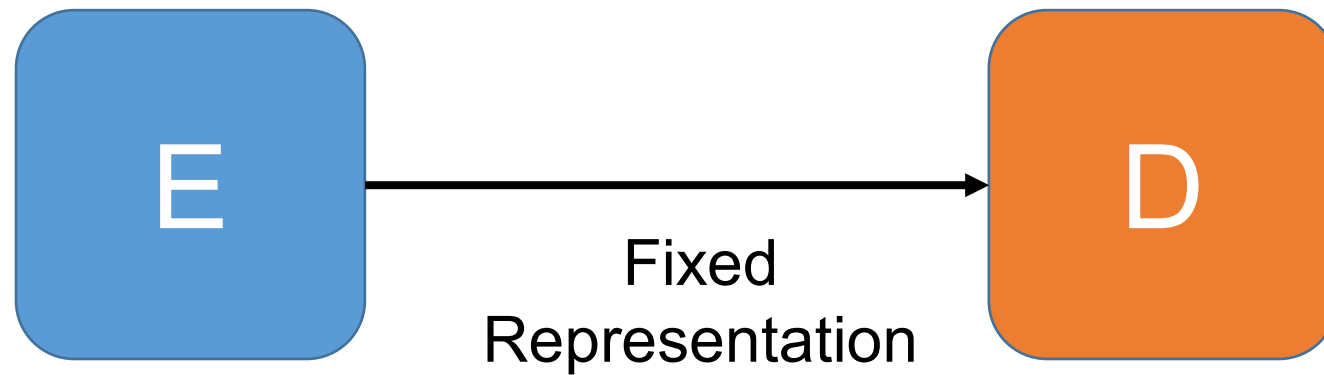
Department of Electrical Engineering Technion, Haifa, Israel

Outline

- Background
- Method
- Experiments
- Conclusion

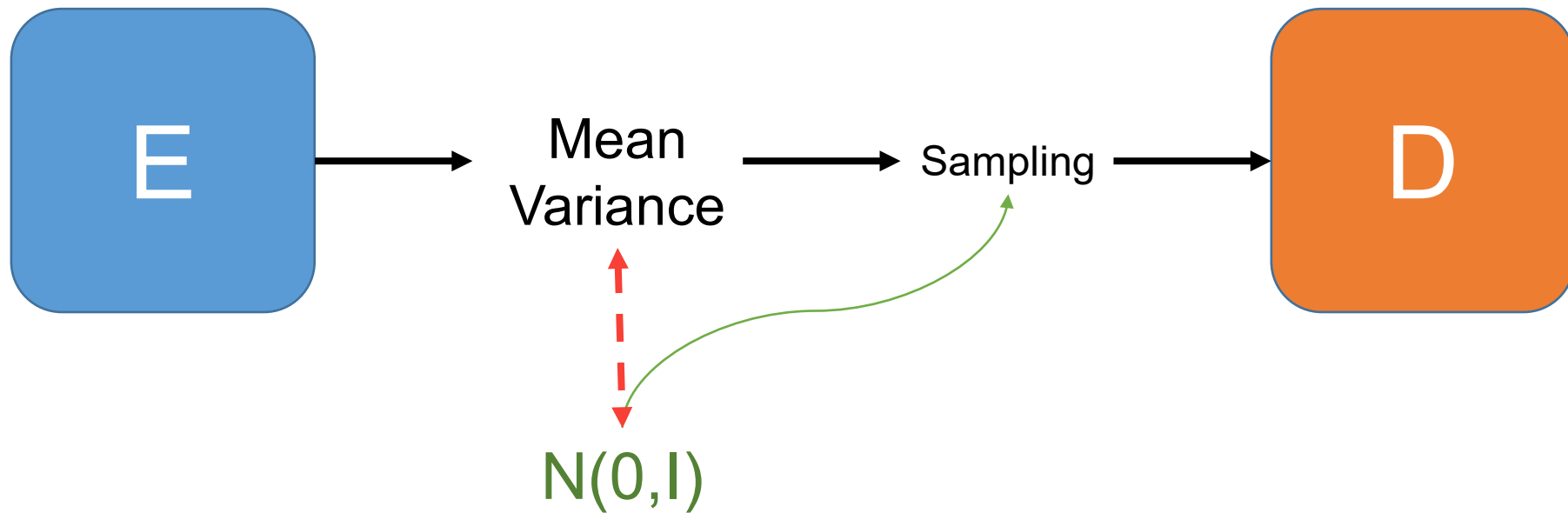
Auto-Encoder

- Compression



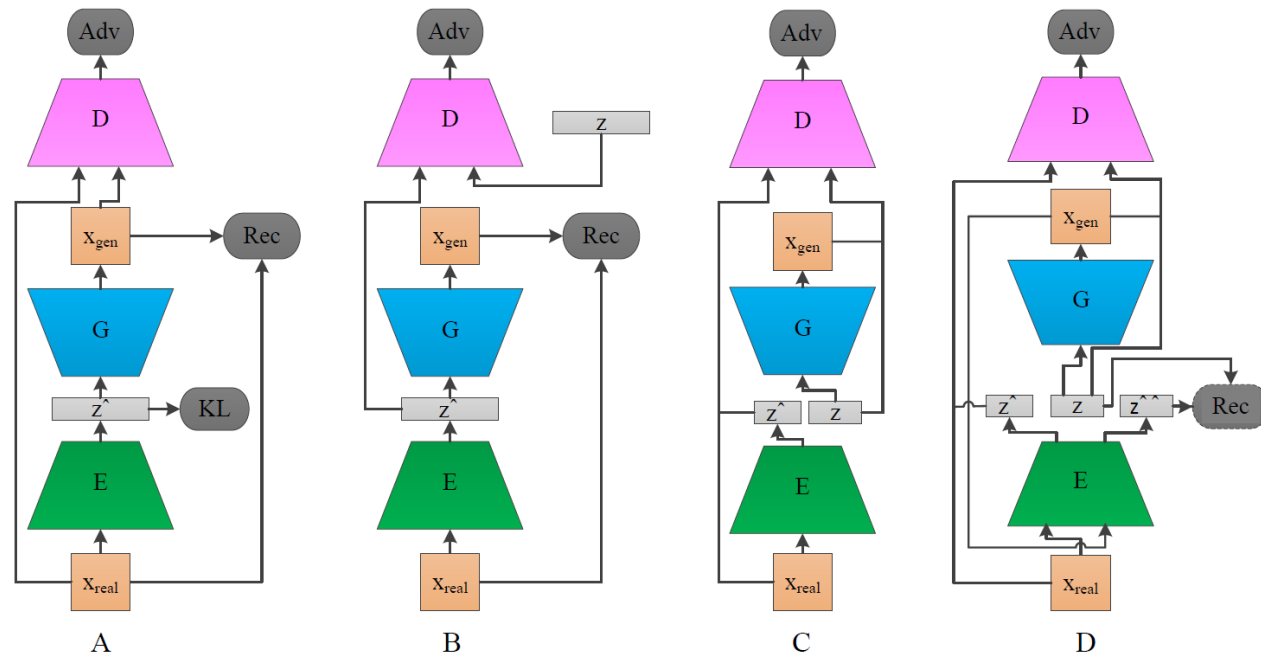
Varitional Auto-Encoder (VAE)

- Generation



Hybrid VAE

- VAE is blurry
 - Learning from GAN



- Too Complex

IntroVAE

- Real \rightarrow minimizing KL ; Generated \rightarrow maximizing KL
 - Encoder: Maximizing KL of generated samples
 - Decoder: Fooling the Encoder \rightarrow Generating samples minimizing KL
 - Nash equilibrium
 - $p_G = p_{data}$

KL between posterior and prior

Encoder $L_E(x, z) = E(x) + [m - E(G(z))]^+ + L_{AE}(x),$

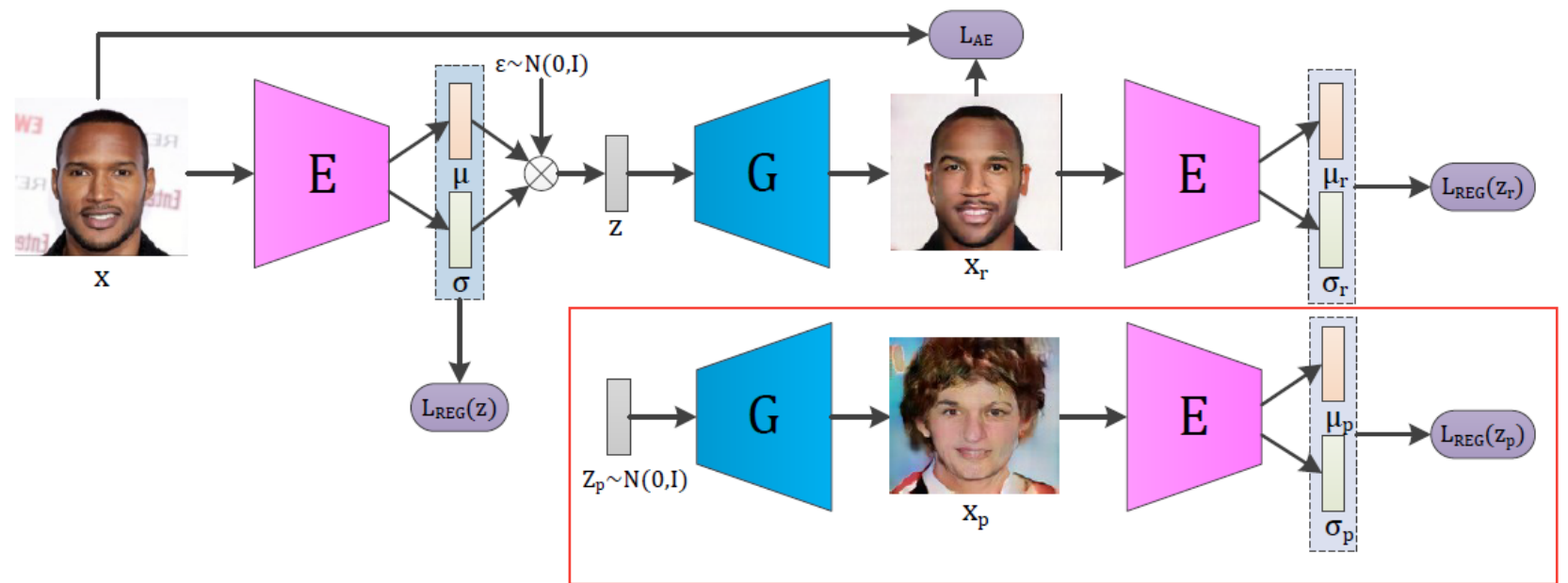
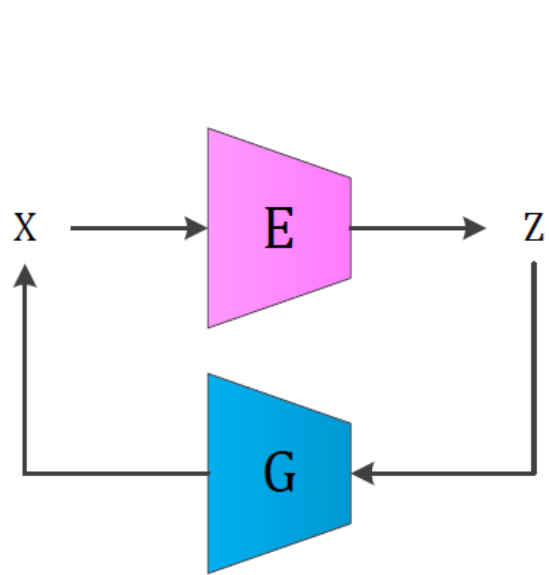
Hard Threshold

Decoder

$$L_G(z) = E(G(z)) + L_{AE}(x).$$

Reconstruction

IntroVAE



IntroVAE

Algorithm 1 Training IntroVAE model

- 1: $\theta_G, \phi_E \leftarrow$ Initialize network parameters
 - 2: **while** not converged **do**
 - 3: $X \leftarrow$ Random mini-batch from dataset
 - 4: $Z \leftarrow Enc(X)$
 - 5: $Z_p \leftarrow$ Samples from prior $N(0, I)$
 - 6: $X_r \leftarrow Dec(Z), X_p \leftarrow Dec(Z_p)$
 - 7: $L_{AE} \leftarrow L_{AE}(X_r, X)$
 - 8: $Z_r \leftarrow Enc(ng(X_r)), Z_{pp} \leftarrow Enc(ng(X_p))$
 - 9: $L_{adv}^E \leftarrow L_{REG}(Z) + \alpha\{[m - L_{REG}(Z_r)]^+ + [m - L_{REG}(Z_{pp})]^+\}$
 - 10: $\phi_E \leftarrow \phi_E - \eta \nabla_{\phi_E} (L_{adv}^E + \beta L_{AE})$ ▷ Perform Adam updates for ϕ_E
 - 11: $Z_r \leftarrow Enc(X_r), Z_{pp} \leftarrow Enc(X_p)$
 - 12: $L_{adv}^G \leftarrow \alpha\{L_{REG}(Z_r) + L_{REG}(Z_{pp})\}$
 - 13: $\theta_G \leftarrow \theta_G - \eta \nabla_{\theta_G} (L_{adv}^G + \beta L_{AE})$ ▷ Perform Adam updates for θ_G
 - 14: **end while**
-

Soft-IntroVAE

- IntroVAE is **hard to train**

- **Can't reproduce**

- **Soft-IntroVAE**

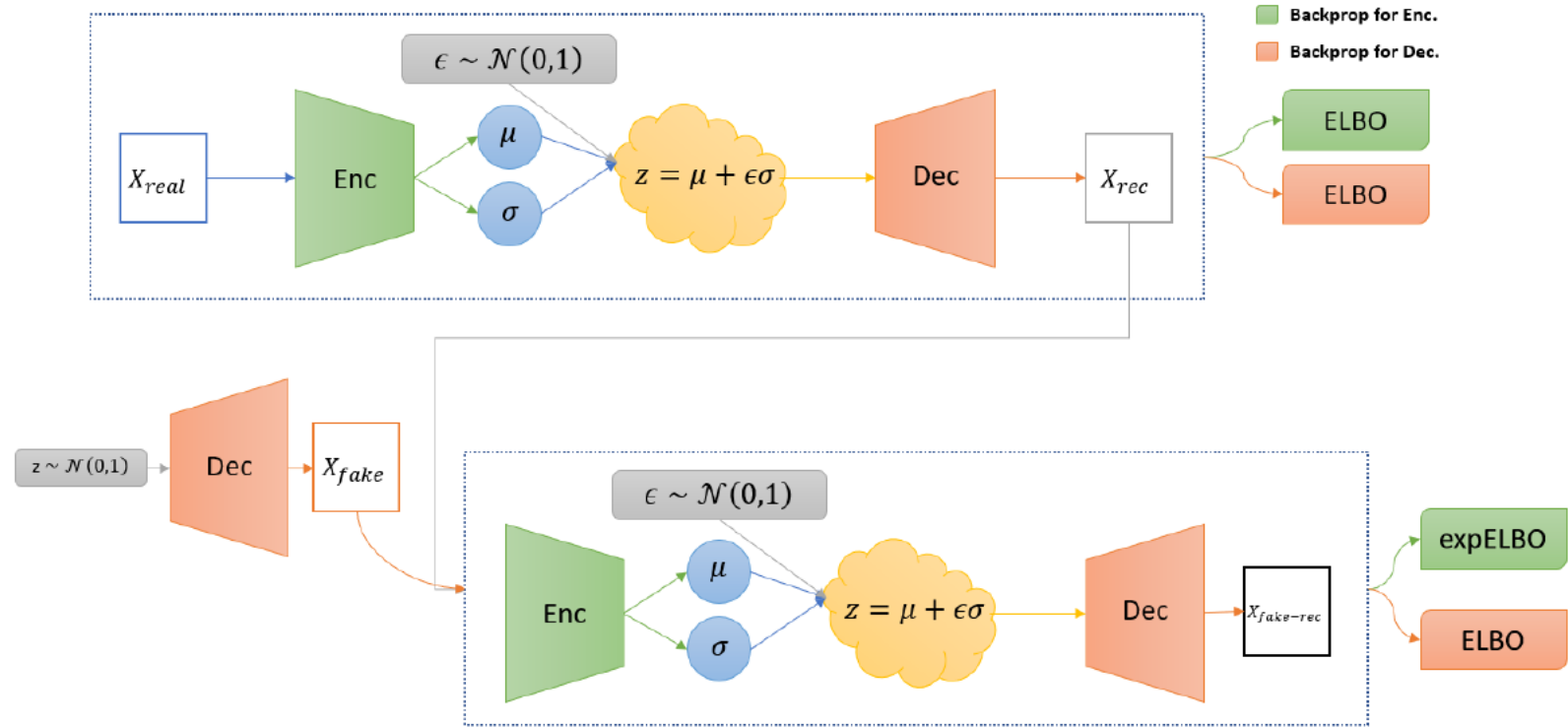
- Utilizing the **complete ELBO** term instead of just the KL

- Replacing the hard threshold with a **soft exponential function** over the ELBO

$$\mathcal{L}_{E_\phi}(x, z) = ELBO(x) - \frac{1}{\alpha} \exp(\alpha ELBO(D_\theta(z))),$$

$$\mathcal{L}_{D_\theta}(x, z) = ELBO(x) + \gamma ELBO(D_\theta(z)),$$

Soft-IntroVAE



$$\mathcal{L}_{E_\phi}(x, z) = ELBO(x) - \frac{1}{\alpha} \exp(\alpha ELBO(D_\theta(z))),$$

$$\mathcal{L}_{D_\theta}(x, z) = ELBO(x) + \gamma ELBO(D_\theta(z)),$$

Soft-IntroVAE

Algorithm 1 Training Soft-IntroVAE (pseudo-code)

Require: $\beta_{rec}, \beta_{kl}, \beta_{neg}, \gamma_r$

```
1:  $\phi_E, \theta_D \leftarrow$  Initialize network parameters
2:  $s \leftarrow 1/\text{input dim}$   $\triangleright$  Scaling constant
3: while not converged do
4:    $X \leftarrow$  Random mini-batch from dataset
5:    $Z \leftarrow E(X)$   $\triangleright$  Encode
6:    $Z_f \leftarrow$  Samples from prior  $N(0, I)$ 
7:   procedure UPDATEENCODER( $\phi_E$ )
8:      $X_r \leftarrow D(Z), X_f \leftarrow D(Z_f)$   $\triangleright$  Decode
9:      $Z_{ff} \leftarrow E(X_f)$ 
10:     $X_{ff} \leftarrow D(Z_{ff})$ 
11:     $\text{ELBO} \leftarrow s \cdot \text{ELBO}(\beta_{rec}, \beta_{kl}, X, X_r, Z)$ 
12:     $\text{ELBO}_f \leftarrow \text{ELBO}(\beta_{rec}, \beta_{neg}, X_f, X_{ff}, Z_{ff})$ 
13:     $\text{expELBO}_f \leftarrow 0.5 \exp(2s \cdot \text{ELBO}_f)$ 
14:     $L_E \leftarrow \text{ELBO} - \text{expELBO}_f$   $\triangleright$  Eq. 4
15:     $\phi_E \leftarrow \phi_E + \eta \nabla_{\phi_E} (L_E)$   $\triangleright$  Adam update
16:  end procedure
```

```
17:   procedure UPDATEDECODER( $\theta_D$ )
18:      $X_r \leftarrow D(Z), X_f \leftarrow D(Z_f)$   $\triangleright$  Decode
19:      $Z_{ff} \leftarrow E(X_f)$ 
20:      $X_{ff} \leftarrow \text{sg}(D(Z_{ff}))$   $\triangleright$  sg: stop-gradient
21:      $\text{ELBO} \leftarrow \beta_{rec} L_{rec}(X, X_r)$ 
22:      $\text{ELBO}_f \leftarrow \text{ELBO}(\gamma_r \cdot \beta_{rec}, \beta_{kl}, X_f, X_{ff}, Z_{ff})$ 
23:      $L_D \leftarrow s \cdot (\text{ELBO} + \text{ELBO}_f)$   $\triangleright$  Eq. 4
24:      $\theta_D \leftarrow \theta_D + \eta \nabla_{\theta_D} (L_D)$   $\triangleright$  Adam update
25:   end procedure
26: end while
```

Analysis

- Nash equilibrium

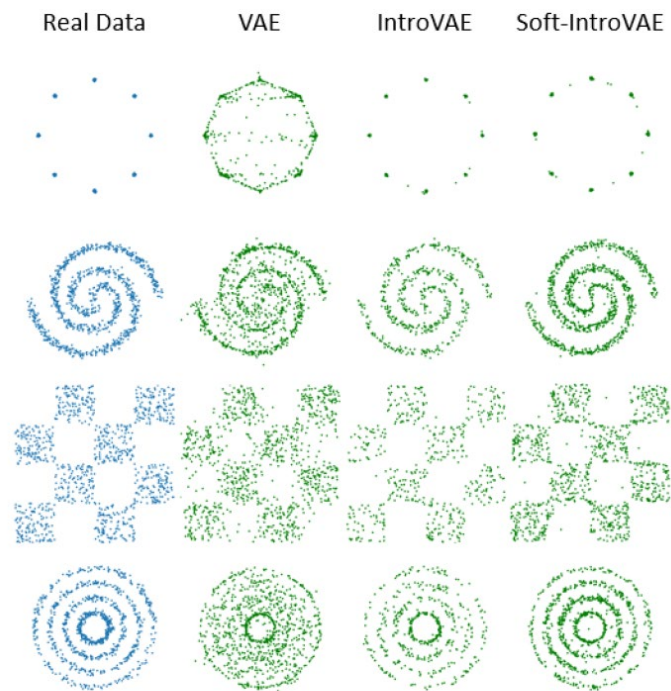
- $d^* \in \arg \min_d \{KL(p_{data} || p_d) + \gamma H(p_d(x))\}.$

- No longer converging to `p_data`, but regularized by entropy

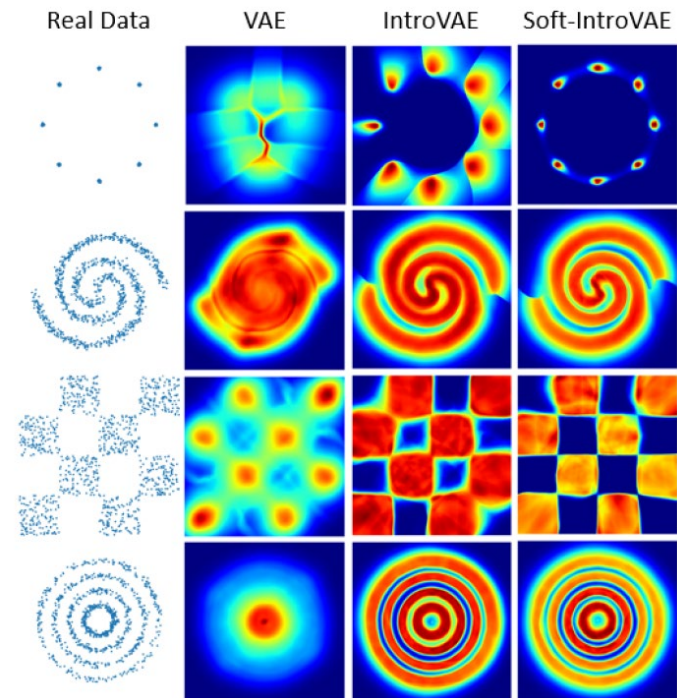
- Math is hard...

Experiments

■ 2D Toy Dataset



(a) Samples from the trained models.

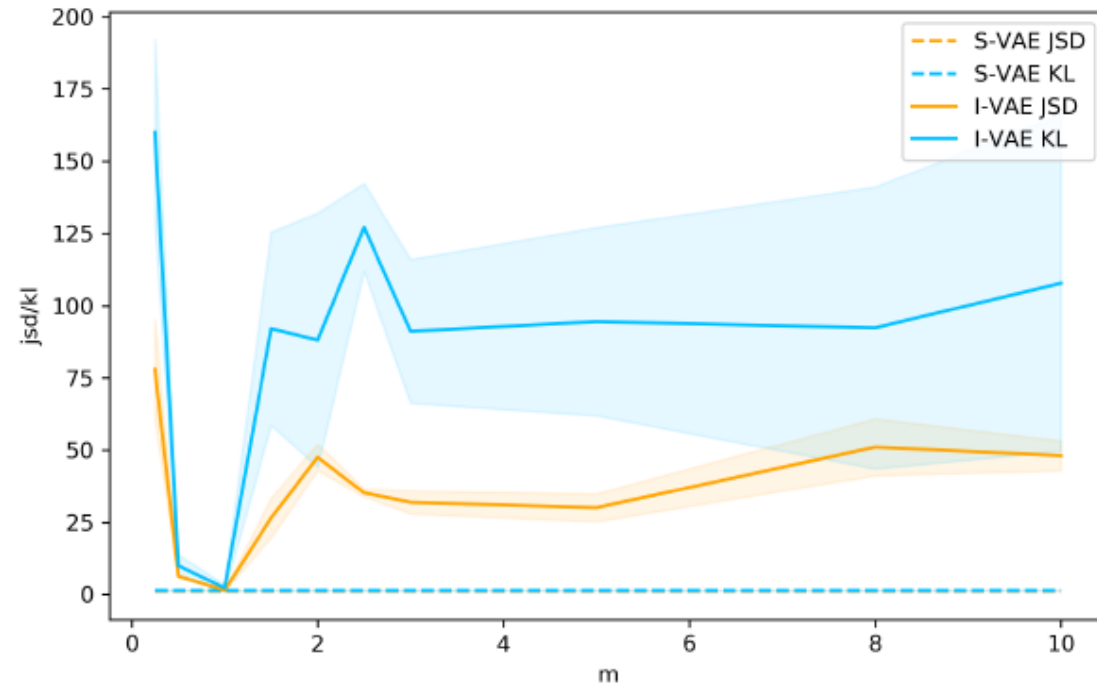


(b) Density estimation with the trained models.

Experiments

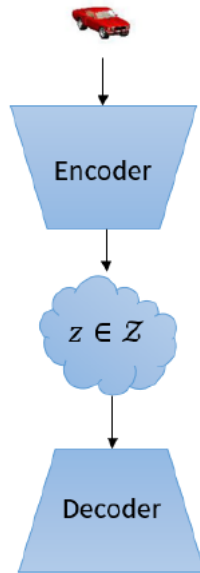
■ Training Stability

- Probably due to the choice of m is sensitive

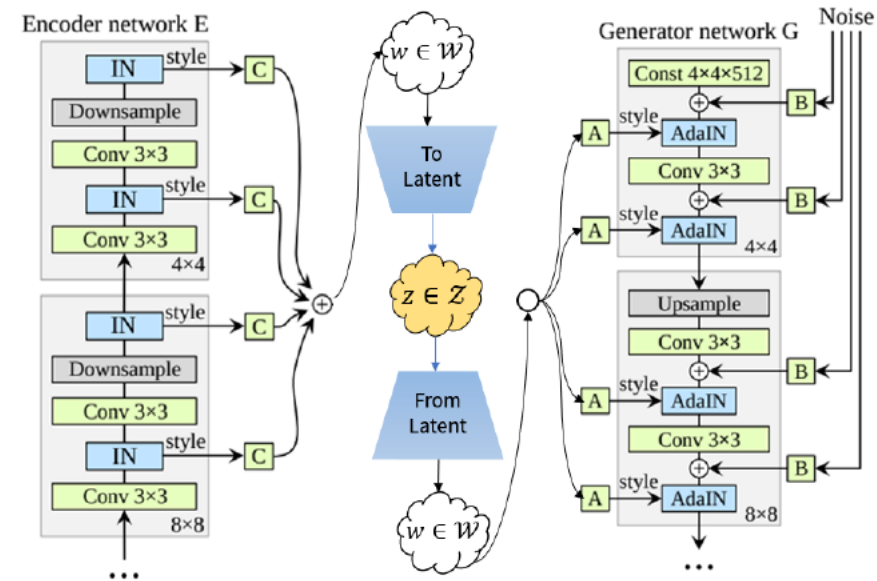


Experiments

- Image Generation
 - Architectures



Standard Architecture

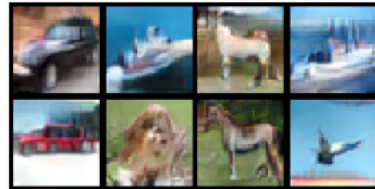


(c) Style-based Architecture inspired by [46]

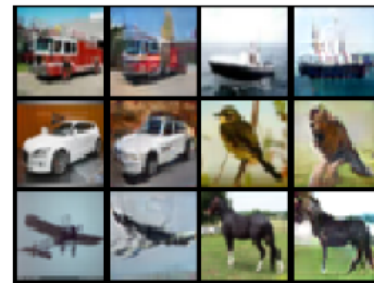
Experiments

- Image Generation

- Cifar-10



(a) Generated samples (FID: 4.6).



(b) Reconstructions on test data: Left: real, right: reconstruction.

Experiments

- Image Generation
 - CelebA-HQ and FFHQ datasets



(a) FFHQ dataset – samples from S-IntroVAE (FID: 17.55).



(b) FFHQ – reconstructions.

Experiments

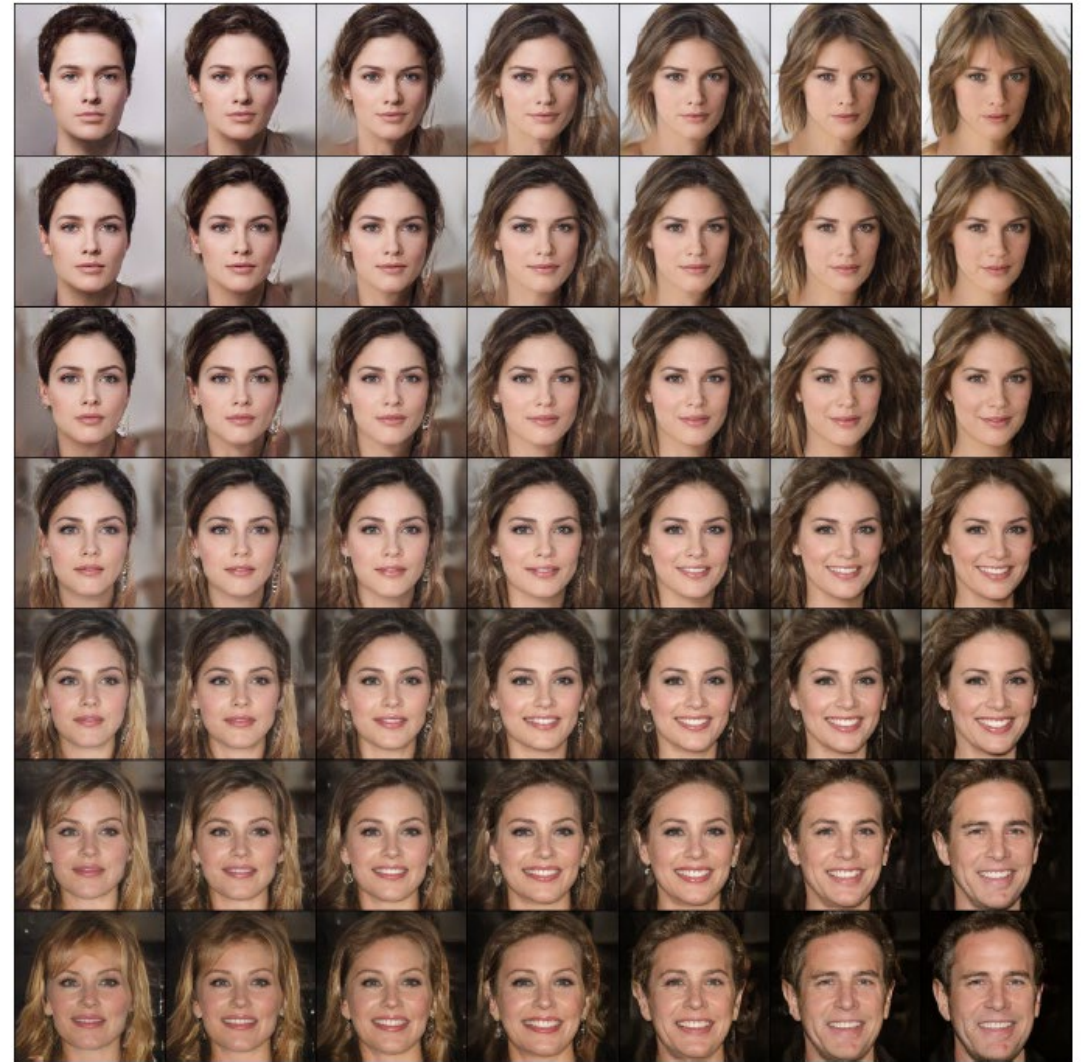
- Image Generation
 - Interpolation in the latent space



Figure 5: Interpolation in the latent space between two samples from a model trained on CelebA-HQ.

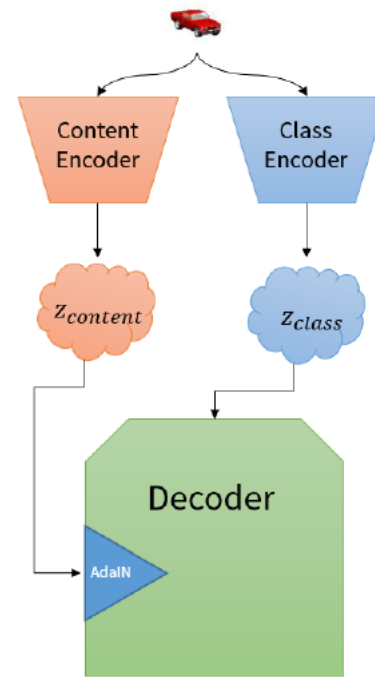
Experiments

- Image Generation
 - Interpolation in the latent space



Experiments

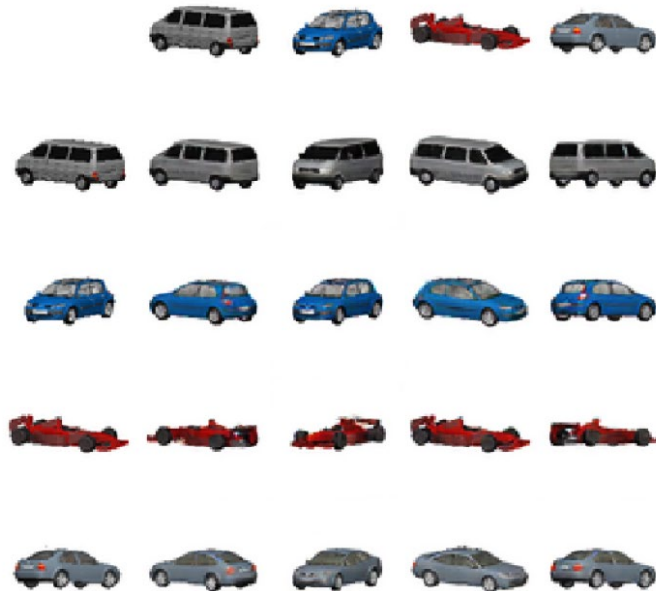
- Image Translation
 - Architectures



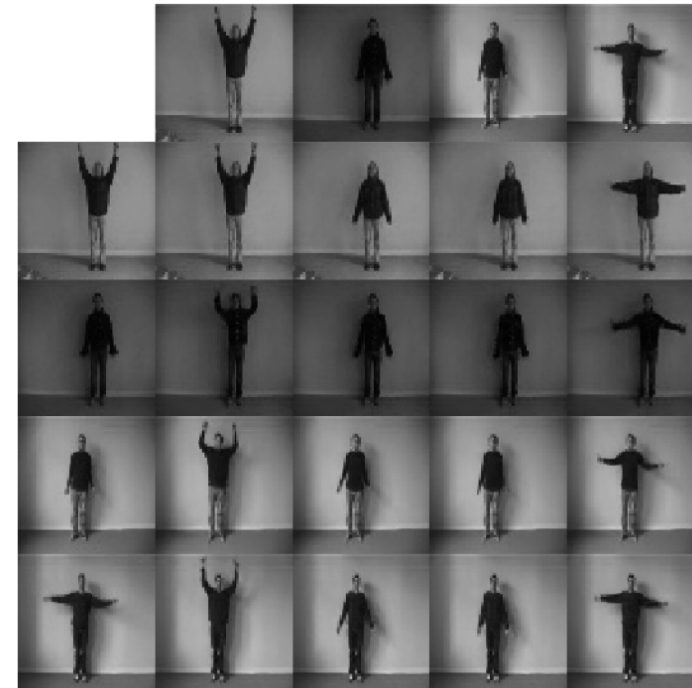
Disentanglement Architecture

Experiments

- Image Translation
 - Cars3D and KTH



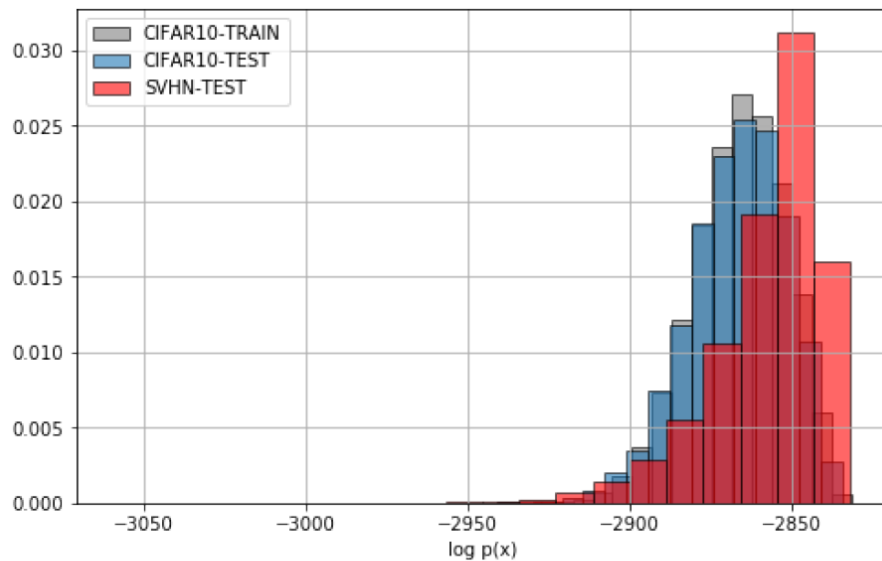
Cars3D



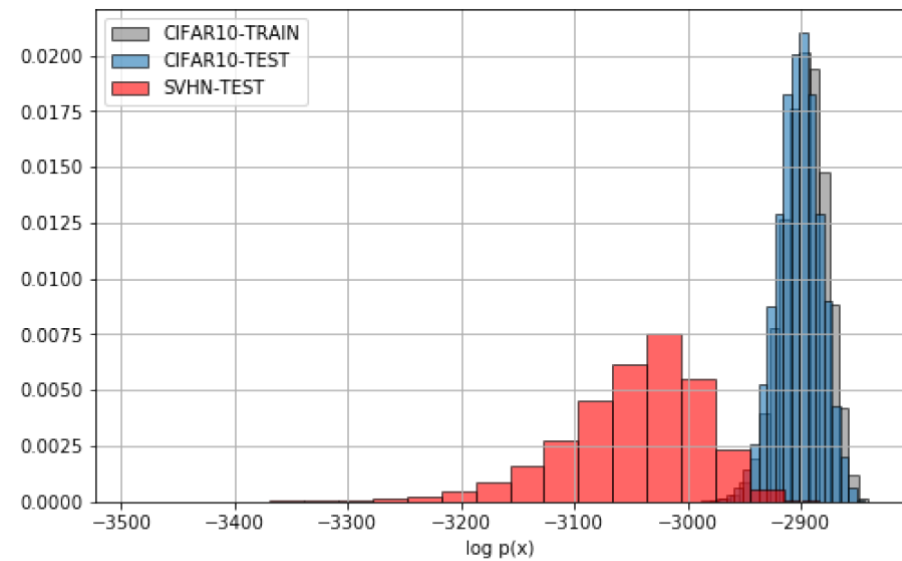
KTH

Experiments

- Out-of-Distribution (OOD) Detection
 - Cifar-10 & SVHN



(a) VAE



(b) Soft-IntroVAE

Conclusion

- Improve the training of IntroVAE
- A deeper theoretical understanding of IntroVAE

Thanks