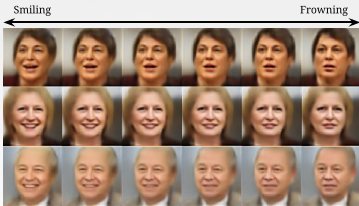


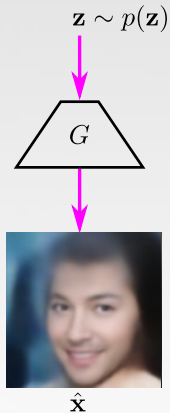
Laplacian Pyramid of Conditional Variational Autoencoders

Garoe Dorta^{1,2} Sara Vicente² Lourdes Agapito³
Neill D.F. Campbell¹ Simon Prince² Ivor Simpson²

¹University of Bath ²Anthropics Technology Ltd. ³University College London



Sampling and editing with generative models [1, 2, 3]



$\hat{\mathbf{x}}$

Sample

\mathbf{z}

Latent vector

$\boldsymbol{\theta}$

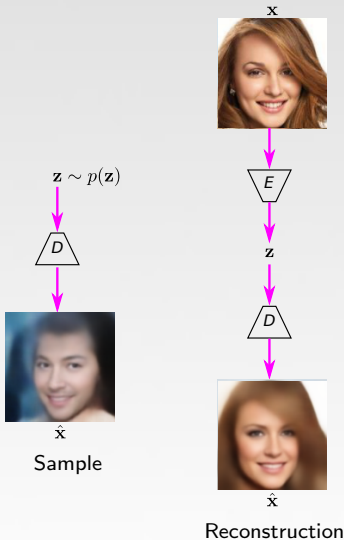
Model parameters

$G(\mathbf{z}; \boldsymbol{\theta})$

Generative function

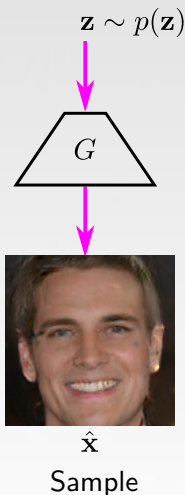
$p(\mathbf{z})$

Simple known distribution



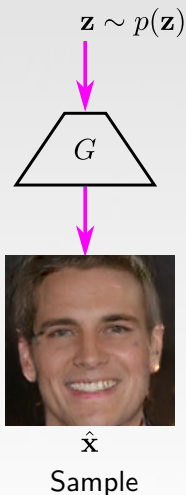
Variational Autoencoders [4, 5]

- Encoder $z \sim E(x)$
- Decoder $\hat{x} \sim D(z)$
- Gaussian likelihood estimation
 - Easy to train
 - Generate blurry images



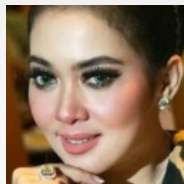
Generative Adversarial Networks [6]

- Generator $\hat{\mathbf{x}} = G(\mathbf{z})$
- Discriminator
- Implicit likelihood estimation
 - Impressive results
 - Unstable training
- Only for sampling



Generative Adversarial Networks [6]

- Generator $\hat{\mathbf{x}} = G(\mathbf{z})$
- Discriminator
- Implicit likelihood estimation
 - Impressive results
 - Unstable training
- Only for sampling



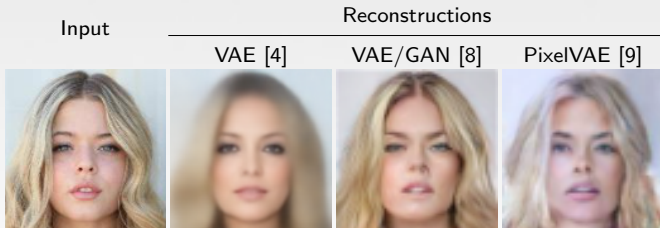
Input



Reconstruction [7]

VAE extensions

- Complex distributions in the latent and output space [8, 9]
- Throw away the simplicity of the Gaussian likelihood estimation.



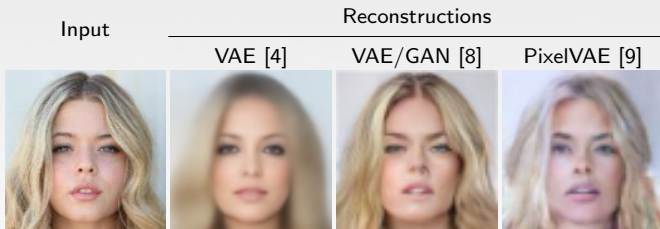
VAE extensions

- Complex distributions in the latent and output space [8, 9]
- Throw away the simplicity of the Gaussian likelihood estimation.



VAE extensions

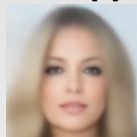
- Complex distributions in the latent and output space [8, 9]
- Throw away the simplicity of the Gaussian likelihood estimation.



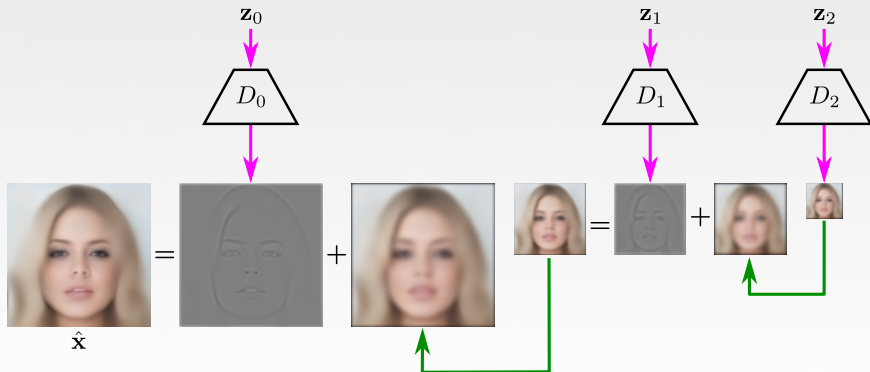
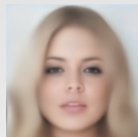
Our method

- Hierarchical approach
- Image generation in tractable steps
- Penalize errors in high-frequency images

VAE [4]

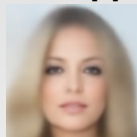


Ours

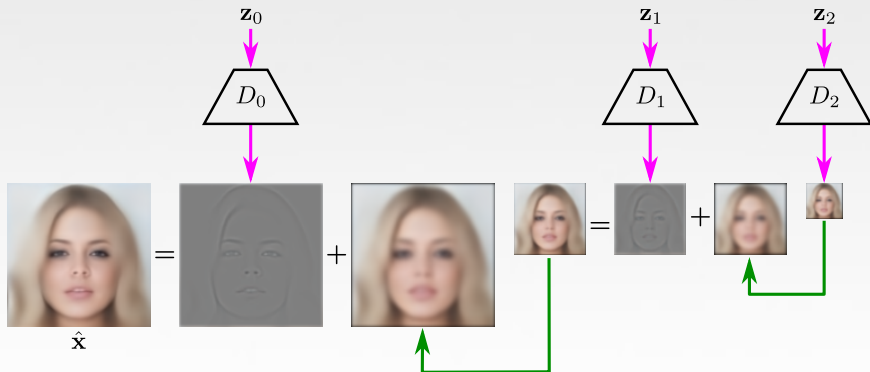


VAE [4]

Ours

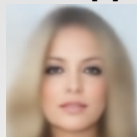


- Hierarchical approach
- Image generation in tractable steps
- Penalize errors in high-frequency images



- Hierarchical approach
- Image generation in tractable steps
- Penalize errors in high-frequency images

VAE [4]



Ours

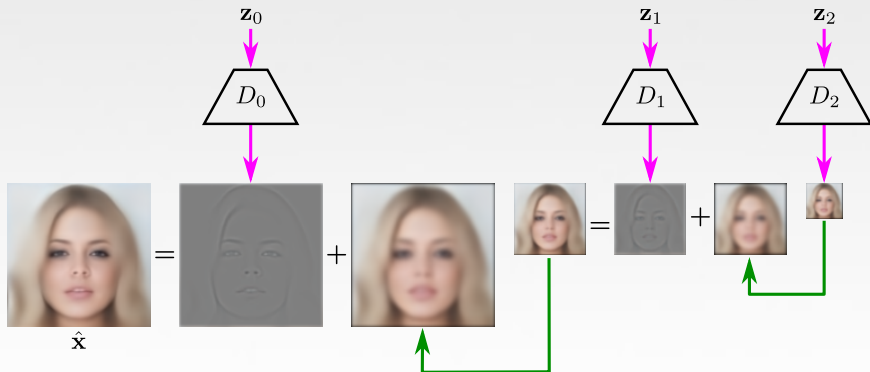
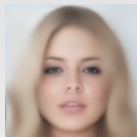


Image editing using the coarse-to-fine model structure

Reconstruction



Image editing using the coarse-to-fine model structure

Reconstruction



Coarse reconstruction

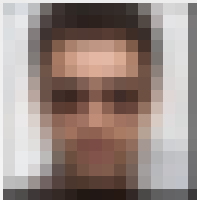


Image editing using the coarse-to-fine model structure

Reconstruction



Coarse painted



Image editing using the coarse-to-fine model structure

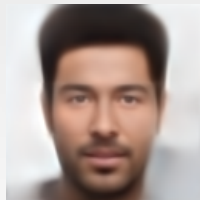
Reconstruction

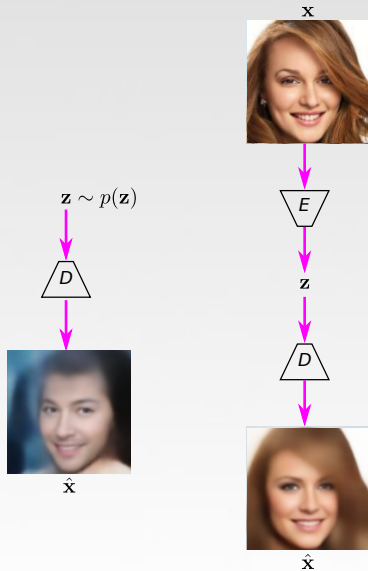


Coarse painted



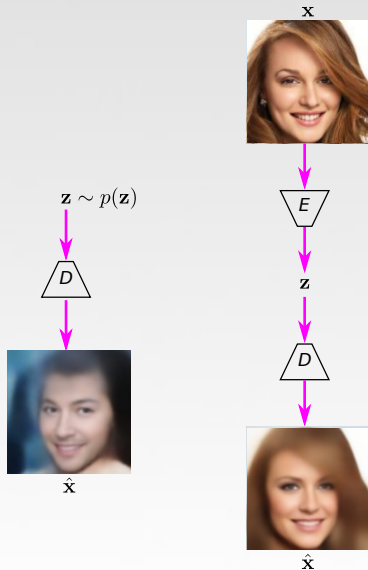
Edited





$$p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma} \mathbf{I})$$

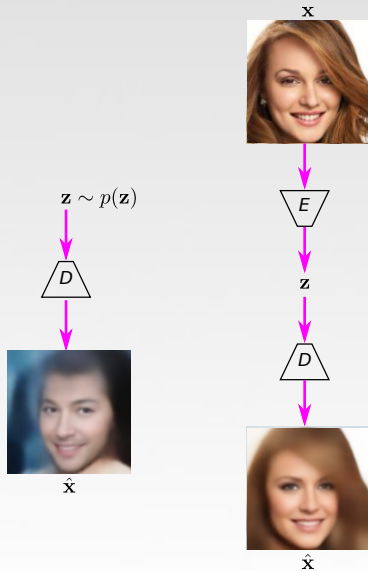
$$D(\mathbf{z}) = \{\boldsymbol{\mu}, \boldsymbol{\sigma}\}$$



$$p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma} \mathbf{I})$$

$$q(\mathbf{z} | \mathbf{x}; \boldsymbol{\phi}) = \mathcal{N}(\boldsymbol{\rho}, \boldsymbol{\omega} \mathbf{I})$$

$$D(\mathbf{z}) = \{\boldsymbol{\mu}, \boldsymbol{\sigma}\}, E(\mathbf{x}) = \{\boldsymbol{\rho}, \boldsymbol{\omega}\}$$

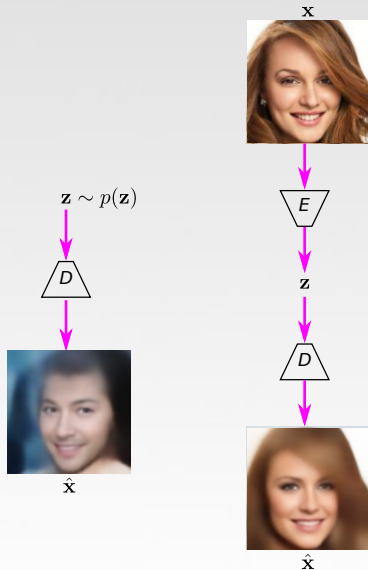


$$p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma} \mathbf{I})$$

$$q(\mathbf{z} | \mathbf{x}; \boldsymbol{\phi}) = \mathcal{N}(\boldsymbol{\rho}, \boldsymbol{\omega} \mathbf{I})$$

$$D(\mathbf{z}) = \{\boldsymbol{\mu}, \boldsymbol{\sigma}\}, E(\mathbf{x}) = \{\boldsymbol{\rho}, \boldsymbol{\omega}\}$$

$$L = \underbrace{-\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z} | \mathbf{x}; \boldsymbol{\phi})} [\log p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta})]}_{\text{Reconstruction loss}} + \underbrace{D_{KL} [q(\mathbf{z} | \mathbf{x}; \boldsymbol{\phi}) || p(\mathbf{z})]}_{\text{Latent space loss}}$$



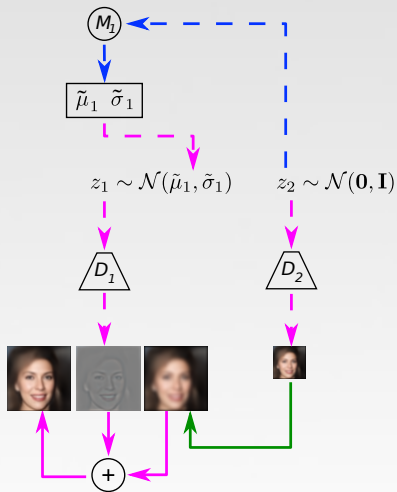
$$p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma} \mathbf{I})$$

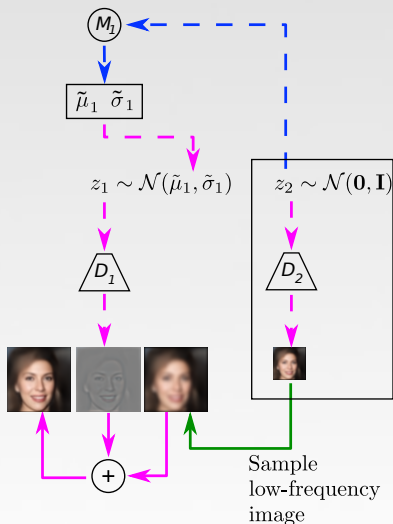
$$q(\mathbf{z} | \mathbf{x}; \boldsymbol{\phi}) = \mathcal{N}(\boldsymbol{\rho}, \boldsymbol{\omega} \mathbf{I})$$

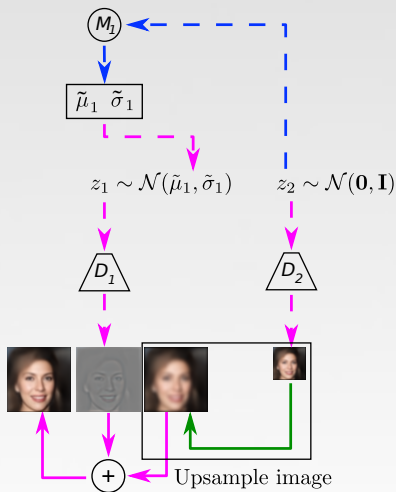
$$D(\mathbf{z}) = \{\boldsymbol{\mu}, \boldsymbol{\sigma}\}, E(\mathbf{x}) = \{\boldsymbol{\rho}, \boldsymbol{\omega}\}$$

$$L = \underbrace{-\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z} | \mathbf{x}; \boldsymbol{\phi})} [\log p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta})]}_{\text{Reconstruction loss}} + \underbrace{D_{KL} [q(\mathbf{z} | \mathbf{x}; \boldsymbol{\phi}) || p(\mathbf{z})]}_{\text{Latent space loss}}$$

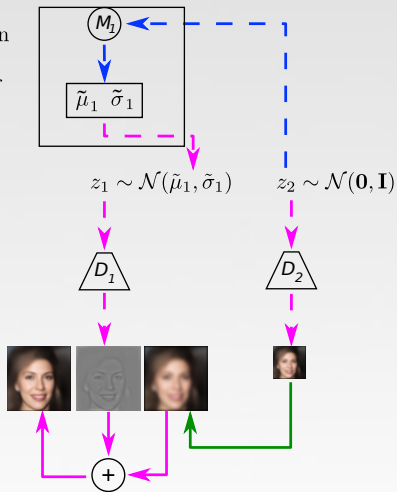
$$p(\mathbf{z}) = \underbrace{\mathcal{N}(\mathbf{0}, \mathbf{I})}_{\text{Prior}}$$

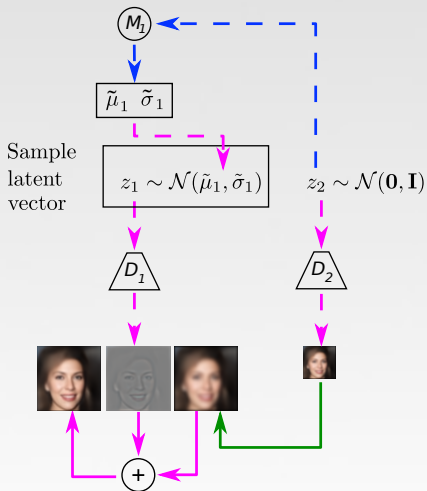


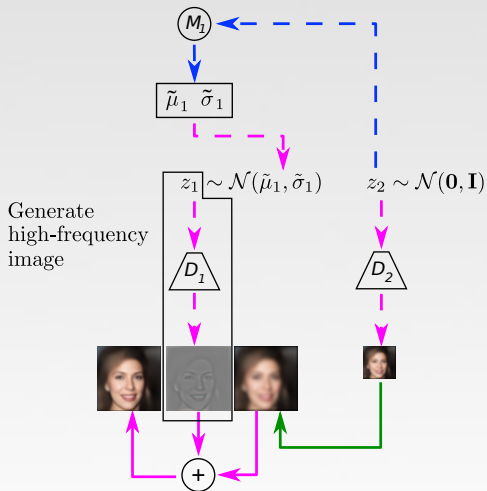


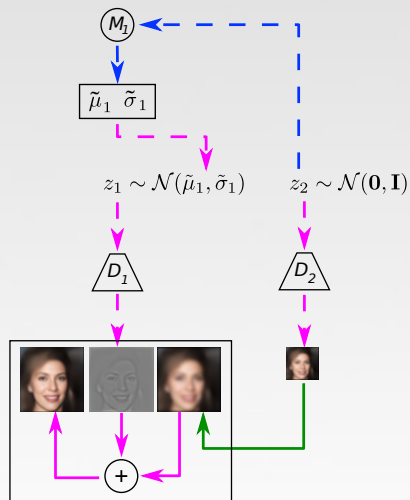


Predict the mean and variance of the latent vector

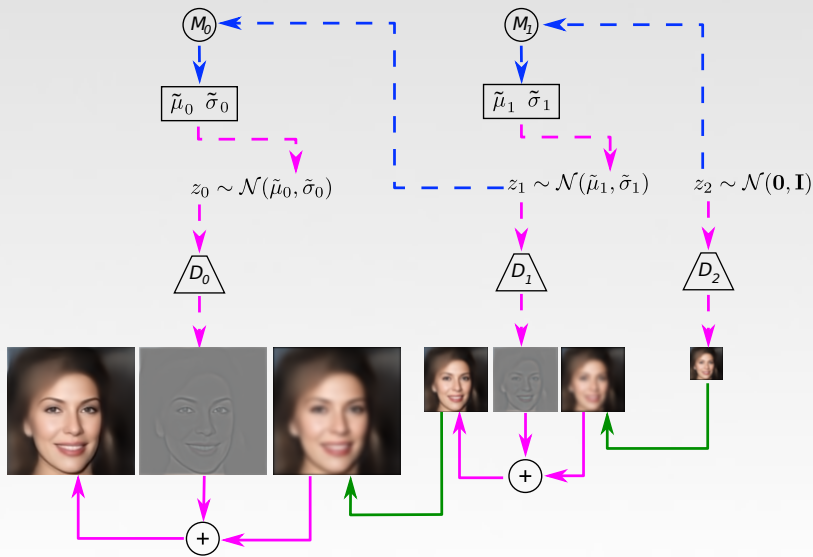


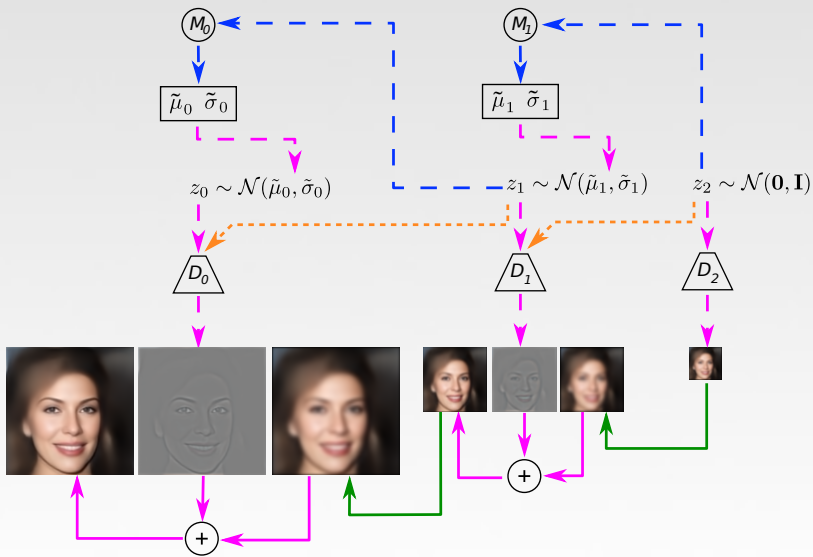






Add high-frequency image
to low-frequency image

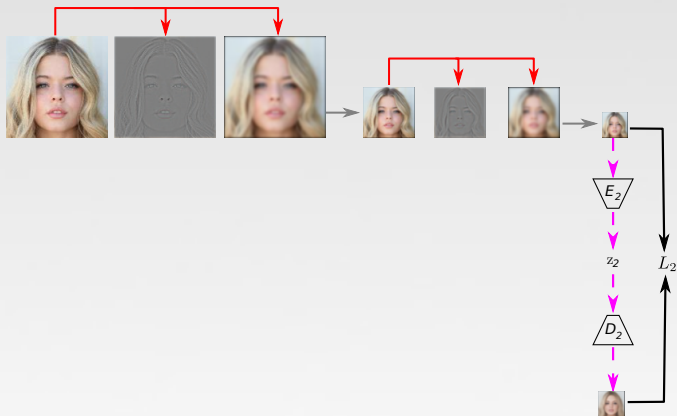


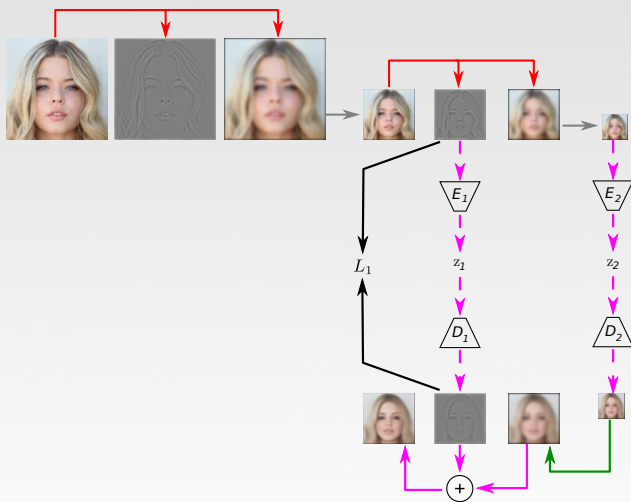


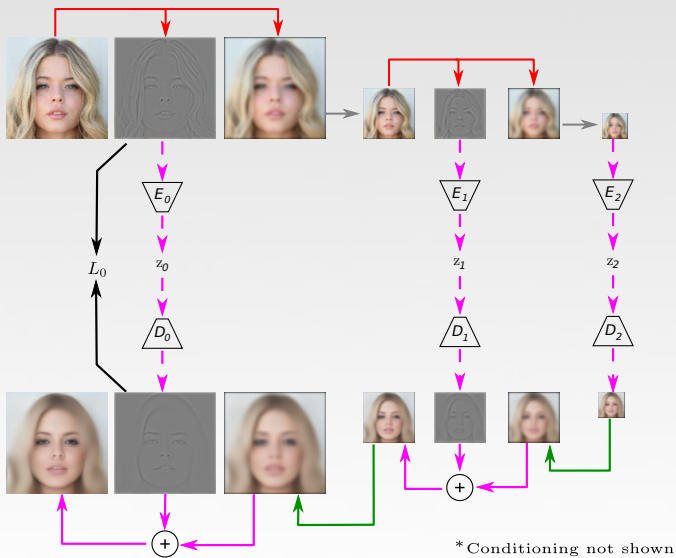












Input image

 \mathbf{x}_k

High frequency

 \mathbf{h}_k

=

+

Low frequency

 $u(\mathbf{x}_{k+1})$

Reconstruction loss

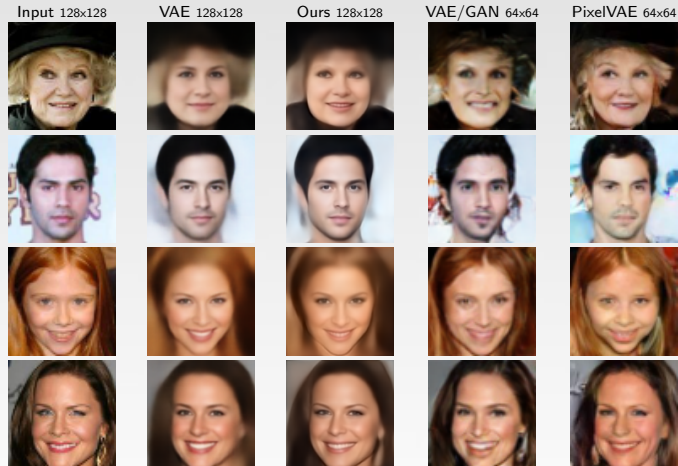
$$L_k = \underbrace{-\mathbb{E}_{\mathbf{z}_k \sim q_k(\mathbf{z}_k | \mathbf{h}_k, u(\mathbf{x}_k); \phi_k)} [\log p_k(\mathbf{h}_k | \mathbf{z}_k, \mathbf{z}_{k+1}, \dots, \mathbf{z}_K; \theta_k)]}_{\text{Reconstruction loss}} +$$

Latent space loss

$$\underbrace{\lambda_k D_{KL} [q_k(\mathbf{z}_k | \mathbf{h}_k, u(\mathbf{x}_{k+1}); \phi_k) || p(\mathbf{z}_k)]}_{\text{Latent space loss}}$$

Prior

$$p(\mathbf{z}_k) = \mathcal{N} \left(R_k(\boldsymbol{\mu}_{k+1}; \boldsymbol{\xi}_k), S_k(\boldsymbol{\sigma}_{k+1}; \boldsymbol{\omega}_k) \right), \quad M_k = \underbrace{\{R_k, S_k\}}_{\text{Prior network}}$$



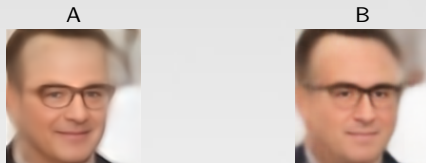
Comparison of image reconstructions

Model	Error ($\sqrt{\text{MSE}}$)
VAE [4] 64×64	22.78 ± 4.64
VAE/GAN [8] 64×64	30.49 ± 7.32
Ours 64×64	20.60 ± 4.81
VAE [4] 128×128	20.75 ± 4.40
Ours 128×128	20.61 ± 5.15

Quantitative model comparison of image reconstructions

Model	Error ($\sqrt{\text{MSE}}$)
VAE [4] 64×64	22.78 ± 4.64
VAE/GAN [8] 64×64	30.49 ± 7.32
Ours 64×64	20.60 ± 4.81
VAE [4] 128×128	20.75 ± 4.40
Ours 128×128	20.61 ± 5.15

Quantitative model comparison of image reconstructions

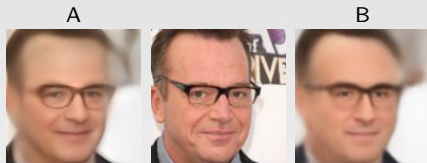


Model	Preference %
	Without original
VAE [4]	15.61 ± 8.14
Ours	84.39 ± 8.14

User study: evaluation of pairs of reconstructions

Model	Error ($\sqrt{\text{MSE}}$)
VAE [4] 64×64	22.78 ± 4.64
VAE/GAN [8] 64×64	30.49 ± 7.32
Ours 64×64	20.60 ± 4.81
VAE [4] 128×128	20.75 ± 4.40
Ours 128×128	20.61 ± 5.15

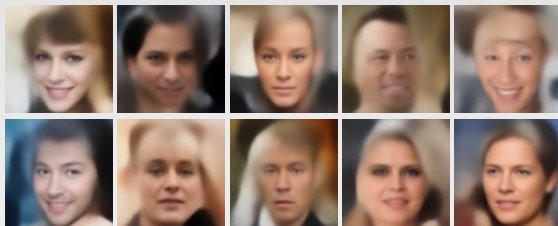
Quantitative model comparison of image reconstructions



Model	Preference %	
	Without original	With original
VAE [4]	15.61 ± 8.14	26.30 ± 7.35
Ours	84.39 ± 8.14	73.70 ± 7.35

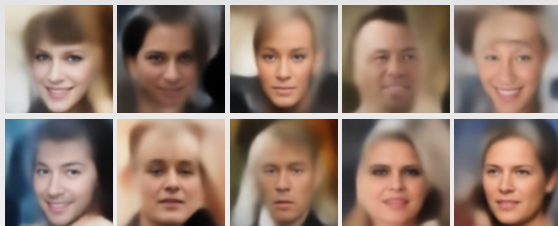
User study: evaluation of pairs of reconstructions

VAE

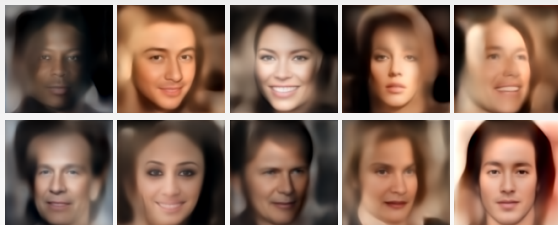


Samples from VAE and our model

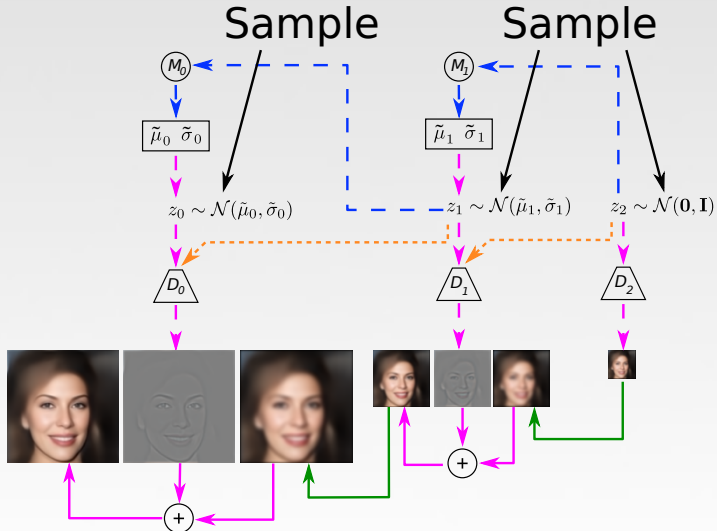
VAE



Ours

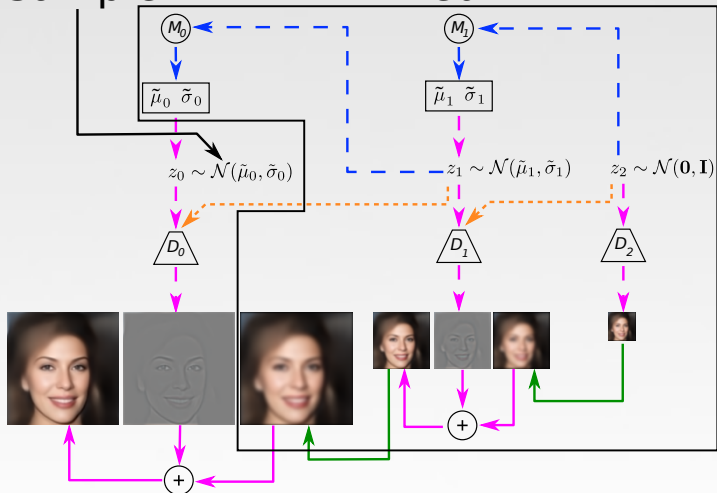


Samples from VAE and our model



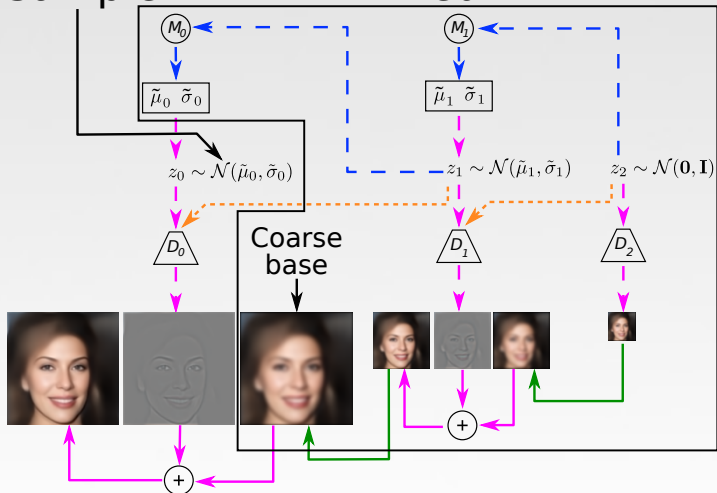
Sample

Fixed

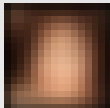


Sample

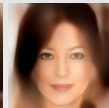
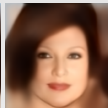
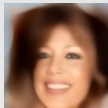
Fixed



Coarse base

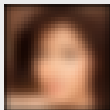
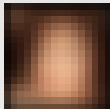


Samples

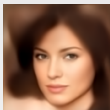
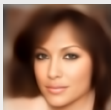
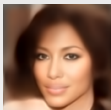
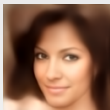
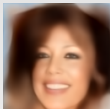


Sampling with $\mathbf{z}_{k \dots K}$ fixed at different levels of the pyramid

Coarse base



Samples



Sampling with $\mathbf{z}_{k \dots K}$ fixed at different levels of the pyramid

Coarse base

Samples



Sampling with $\mathbf{z}_{k \dots K}$ fixed at different levels of the pyramid

Input

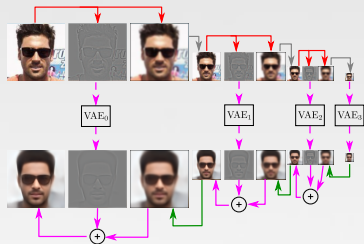


Laplacian pyramid of input

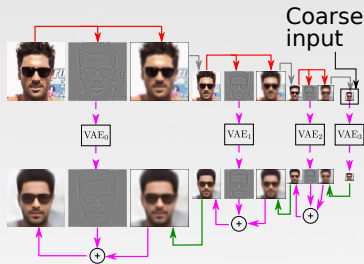
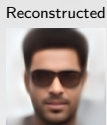
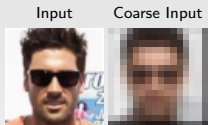
Input



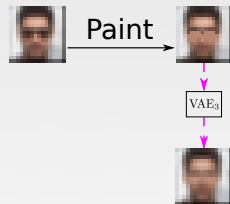
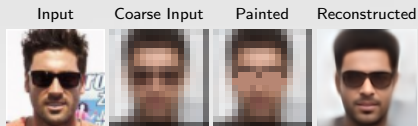
Reconstructed



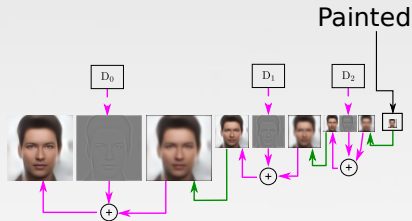
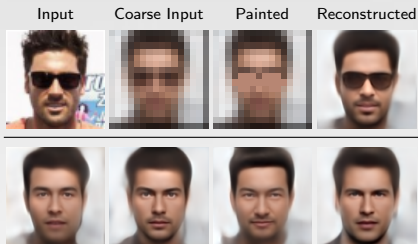
Reconstruct input



Select coarse level



Paint coarse level



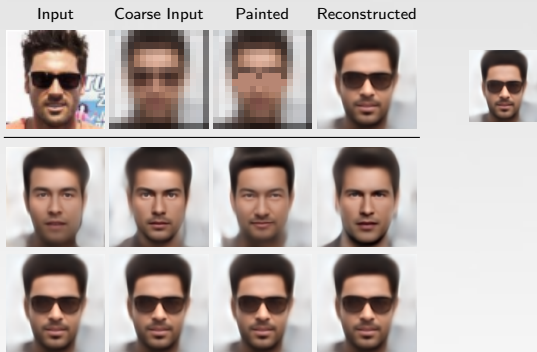
Sample from painted coarse image

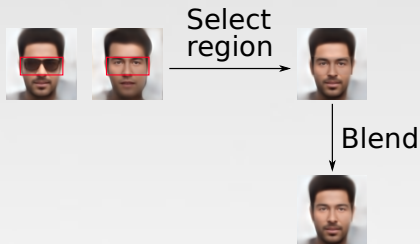
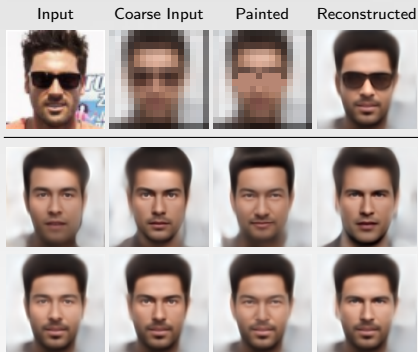
Editing: removing glasses



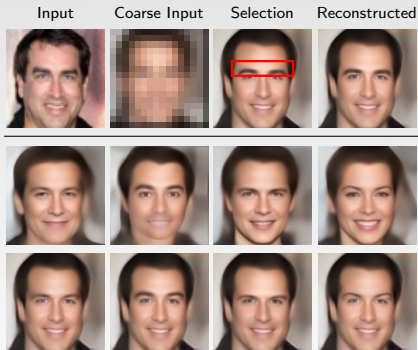
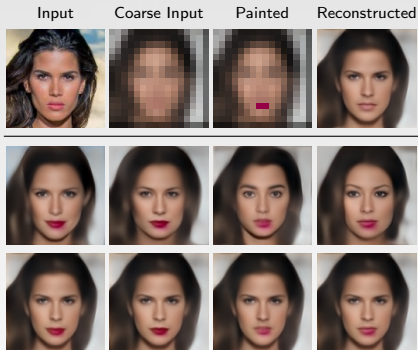
UNIVERSITY OF
BATH

UCL





Blend samples



Add lipstick and adjust eyebrows

Conclusions

- Presented a conditional multi-scale extension of VAE
- Reconstructions and samples are sharper than VAE
- Model allows partial sampling

Limitations and extensions

- Greedy learning
 - End-to-end training strategies
- Gaussian likelihood
 - Complex distributions: perceptual loss or PixelCNN layers



Thank you

Questions?



- [1] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017.
- [2] Xinchun Yan, Jimei Yang, Kihyuk Sohn, and Honglak Lee. Attribute2image: Conditional image generation from visual attributes. *Proceedings of European Conference on Computer Vision (ECCV)*, 2016.
- [3] Jun-Yan Zhu, Philipp Krähenbühl, Eli Shechtman, and Alexei A. Efros. Generative visual manipulation on the natural image manifold. In *Proceedings of European Conference on Computer Vision (ECCV)*, 2016.
- [4] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *International Conference on Learning Representations*, 2014.
- [5] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. *ICML*, 2014.
- [6] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, pages 2672–2680, 2014.
- [7] David Berthelot, Tom Schumm, and Luke Metz. Began: Boundary equilibrium generative adversarial networks. *CoRR*, 2017.
- [8] Anders Boesen Lindbo Larsen, Søren Kaae Sønderby, and Ole Winther. Autoencoding beyond pixels using a learned similarity metric. In *Proceedings of The 33rd International Conference on Machine Learning*, volume 48, pages 1558–1566. JMLR, 2016.
- [9] Ishaan Gulrajani, Kundan Kumar, Faruk Ahmed, Adrien Ali Taiga, Francesco Visin, David Vazquez, and Aaron Courville. PixelVAE: A Latent Variable Model for Natural Images. In *International Conference on Learning Representations (ICLR)*, 2017.