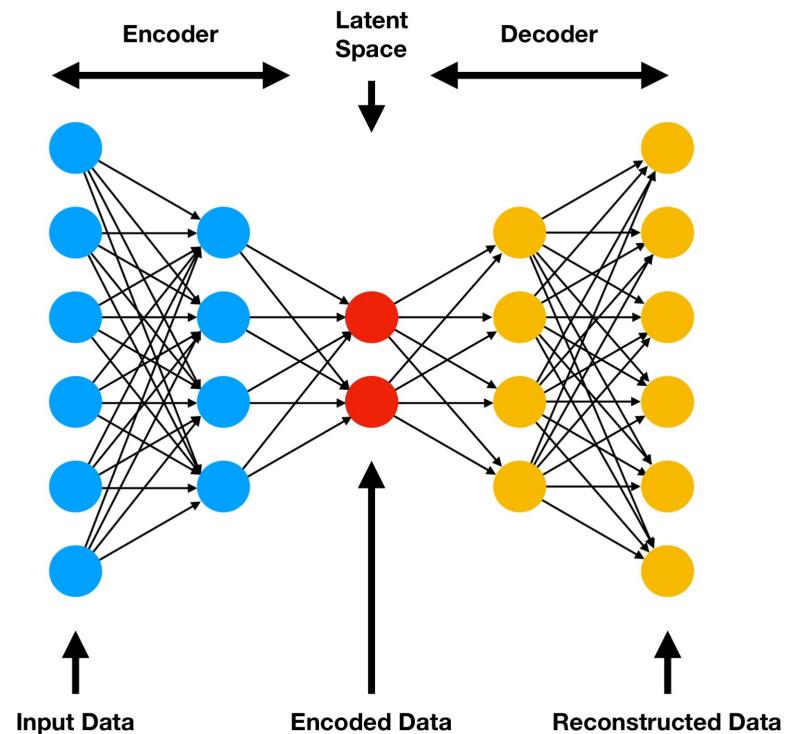
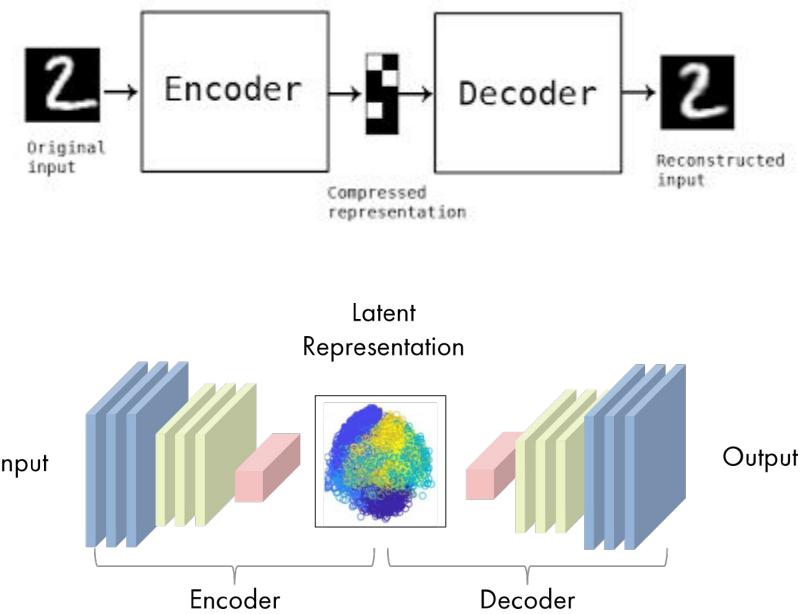
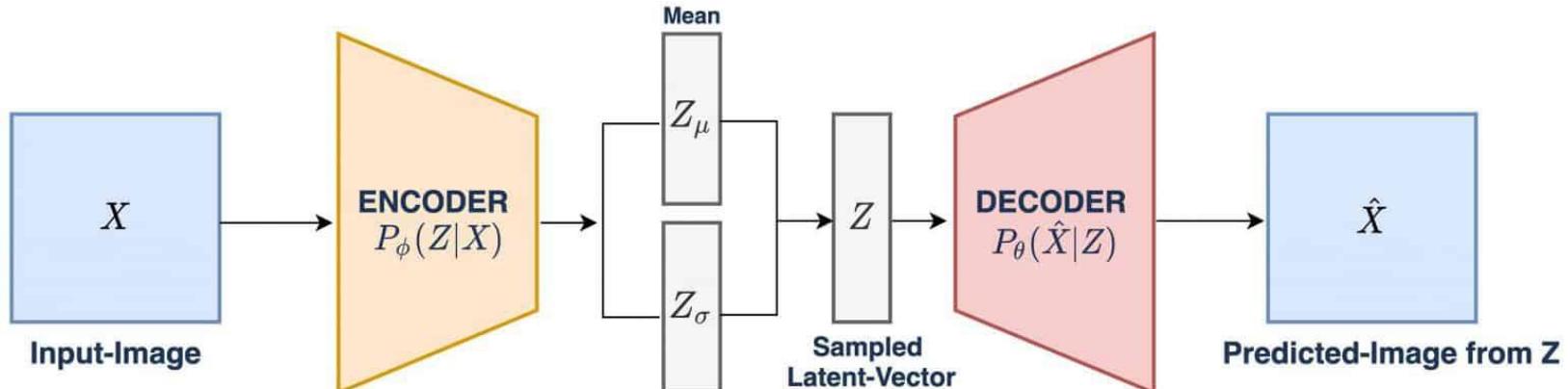


# Deep Learning

Week 6 : Variational AutoEncoders +  
Compression + Perceptual Optimization

# Auto-Encoders Review





Mean  
 $Z_\mu$   
Variance or Standard Deviation  
 $Z_\sigma$

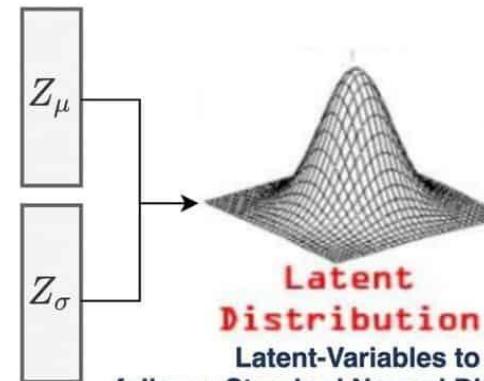
Sampled Latent-Vector  
 $Z$

Predicted-Image from  $Z$   
 $\hat{X}$

Sample a point from  $G(Z_\mu, Z_\sigma)$

$$Z = \mu + \sigma \odot \epsilon$$

$$\epsilon \sim \mathcal{N}(0, 1)$$



Latent Distribution  
Latent-Variables to follow a Standard Normal Distribution

---

# Auto-Encoding Variational Bayes

---

**Diederik P. Kingma**

Machine Learning Group  
Universiteit van Amsterdam  
dpkingma@gmail.com

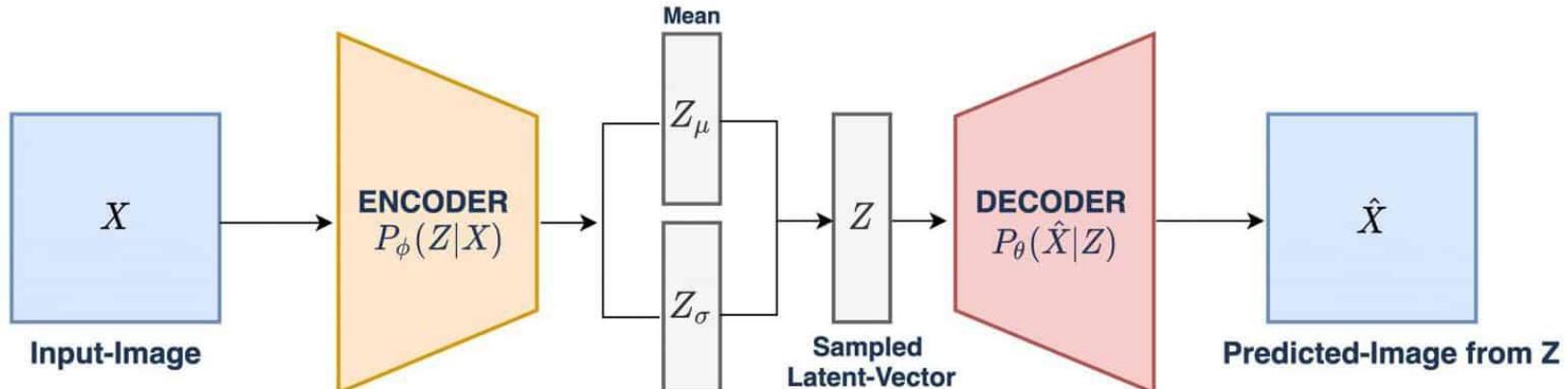
**Max Welling**

Machine Learning Group  
Universiteit van Amsterdam  
welling.max@gmail.com

## Abstract

Can we efficiently learn the parameters of directed probabilistic models, in the presence of continuous latent variables with intractable posterior distributions, and in case of large datasets? We introduce a novel learning and approximate inference method that works efficiently, under some mild conditions, even in the on-line and intractable case. The method involves optimization of a stochastic objective function that can be straightforwardly optimized w.r.t. all parameters, using standard gradient-based optimization methods.

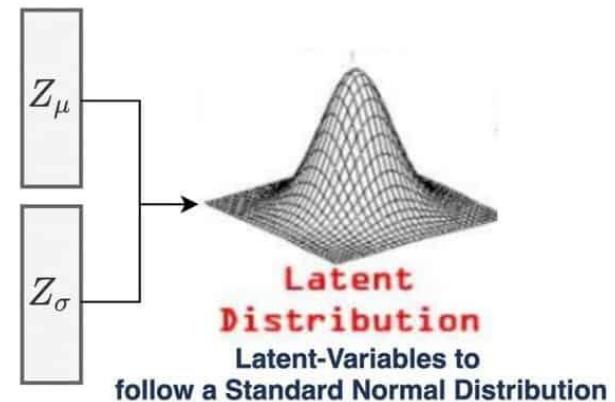
The method does not require the typically expensive sampling loops per data-point required for Monte Carlo EM, and all parameter updates correspond to optimization of the variational lower bound of the marginal likelihood, unlike the wake-sleep algorithm. These theoretical advantages are reflected in experimental results.



Sample a point from  $G(Z_\mu, Z_\sigma)$

$$Z = \mu + \sigma \odot \epsilon$$

$$\epsilon \sim \mathcal{N}(0, 1)$$



# White-Board Derivation of Variational Auto-Encoder

But First : 1) What is KL-Divergence ?

For two discrete pdf's  $p(x)$  and  $q(x)$  sampled over  $x$  we define the Kullback-Leibler Divergence from  $q(x)$  to  $p(x)$  as:

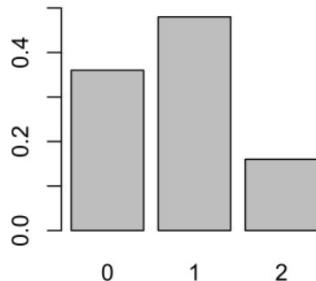
$$D_{\text{KL}}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log \left( \frac{P(x)}{Q(x)} \right),$$

which is equivalent to

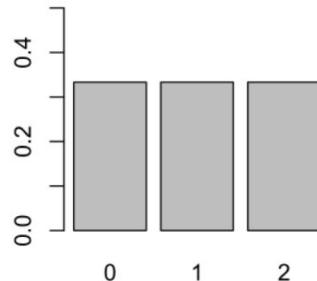
$$D_{\text{KL}}(P \parallel Q) = - \sum_{x \in \mathcal{X}} P(x) \log \left( \frac{Q(x)}{P(x)} \right).$$

# KL - Divergence

**Distribution P**  
Binomial with  $p = 0.4$ ,  $N = 2$



**Distribution Q**  
Uniform with  $p = 1/3$



Two distributions to illustrate relative entropy



Relative entropies  $D_{\text{KL}}(P \parallel Q)$  and  $D_{\text{KL}}(Q \parallel P)$  are calculated as follows. This example uses the natural log with base  $e$ , designated  $\ln$  to get results in nats (see [units of information](#)):

$$\begin{aligned} D_{\text{KL}}(P \parallel Q) &= \sum_{x \in \mathcal{X}} P(x) \ln \left( \frac{P(x)}{Q(x)} \right) \\ &= \frac{9}{25} \ln \left( \frac{9/25}{1/3} \right) + \frac{12}{25} \ln \left( \frac{12/25}{1/3} \right) + \frac{4}{25} \ln \left( \frac{4/25}{1/3} \right) \\ &= \frac{1}{25} (32 \ln(2) + 55 \ln(3) - 50 \ln(5)) \approx 0.0852996, \end{aligned}$$

$$\begin{aligned} D_{\text{KL}}(Q \parallel P) &= \sum_{x \in \mathcal{X}} Q(x) \ln \left( \frac{Q(x)}{P(x)} \right) \\ &= \frac{1}{3} \ln \left( \frac{1/3}{9/25} \right) + \frac{1}{3} \ln \left( \frac{1/3}{12/25} \right) + \frac{1}{3} \ln \left( \frac{1/3}{4/25} \right) \\ &= \frac{1}{3} (-4 \ln(2) - 6 \ln(3) + 6 \ln(5)) \approx 0.097455. \end{aligned}$$

$x$	0	1	2
<b>Distribution <math>P(x)</math></b>	$\frac{9}{25}$	$\frac{12}{25}$	$\frac{4}{25}$
<b>Distribution <math>Q(x)</math></b>	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$

$$\text{JSD}(P \parallel Q) = \frac{1}{2} D(P \parallel M) + \frac{1}{2} D(Q \parallel M),$$

where  $M = \frac{1}{2}(P + Q)$  is a [mixture distribution](#) of  $P$  and  $Q$ .

## Jensen Shannon Divergence

**Note 1:**  $\text{JSD}(P \parallel Q) = \frac{1}{2}D(P \parallel M) + \frac{1}{2}D(Q \parallel M),$

where  $M = \frac{1}{2}(P + Q)$  is a **mixture distribution** of  $P$  and  $Q$ .

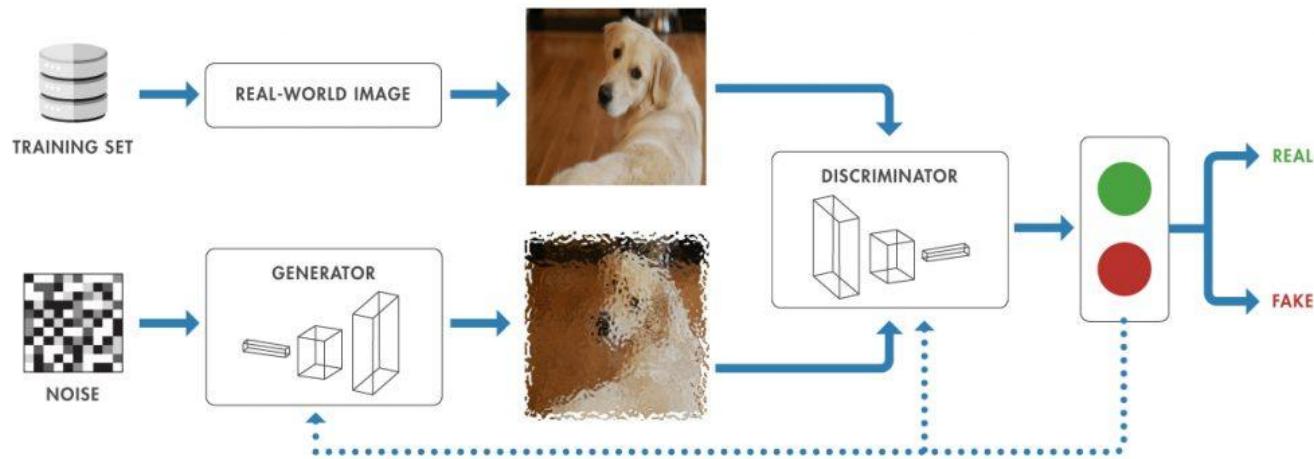
### Note 2:

With base-e logarithm, which is commonly used in statistical thermodynamics, the upper bound is  $\ln(2)$ . In general, the bound in base  $b$  is  $\log_b(2)$ :

$$0 \leq \text{JSD}(P \parallel Q) \leq \log_b(2).$$

**Note 3:  $\text{JSD}(P \parallel Q) = \text{JSD}(Q \parallel P)$**   
**Symmetrical Property!**

# Generative Adversarial Networks

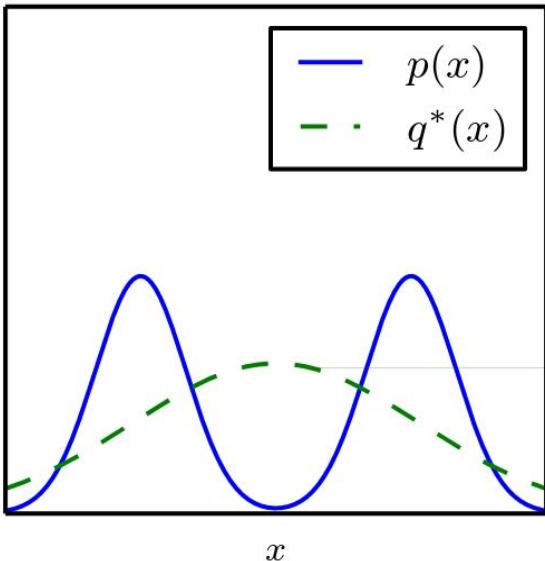


$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))].$$

Consider  $p(x)$  being the real data distribution,  
and  $q^*(x)$  being the generated data distribution

$$q^* = \operatorname{argmin}_q D_{\text{KL}}(p\|q)$$

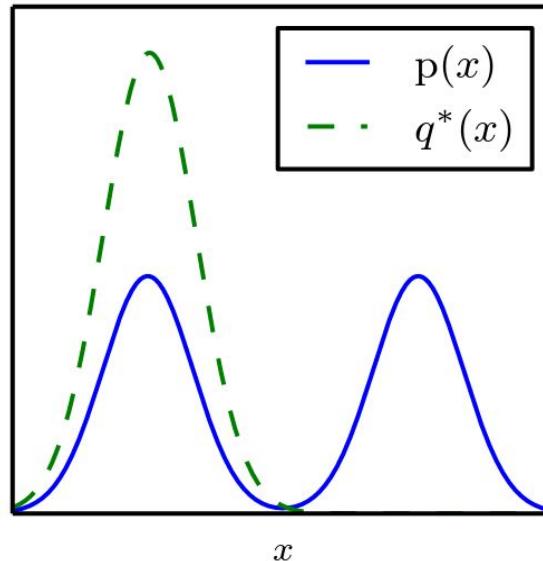
Probability Density



Maximum likelihood

$$q^* = \operatorname{argmin}_q D_{\text{KL}}(q\|p)$$

Probability Density

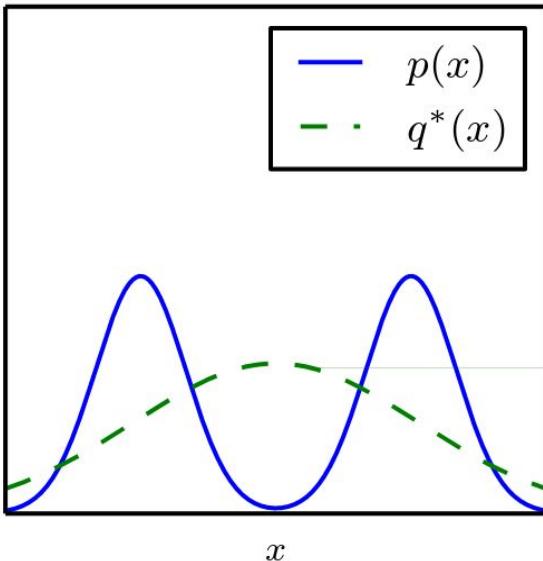


Reverse KL

Consider  $p(x)$  being the real data distribution,  
and  $q^*(x)$  being the generated data distribution

$$q^* = \operatorname{argmin}_q D_{\text{KL}}(p\|q)$$

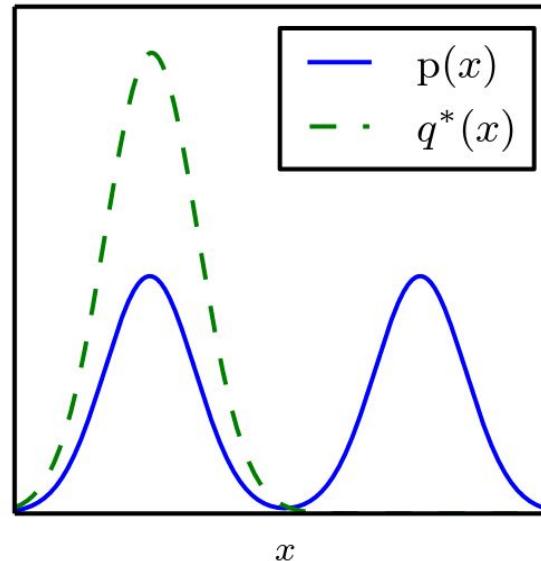
Probability Density



Maximum likelihood

$$q^* = \operatorname{argmin}_q D_{\text{KL}}(q\|p)$$

Probability Density

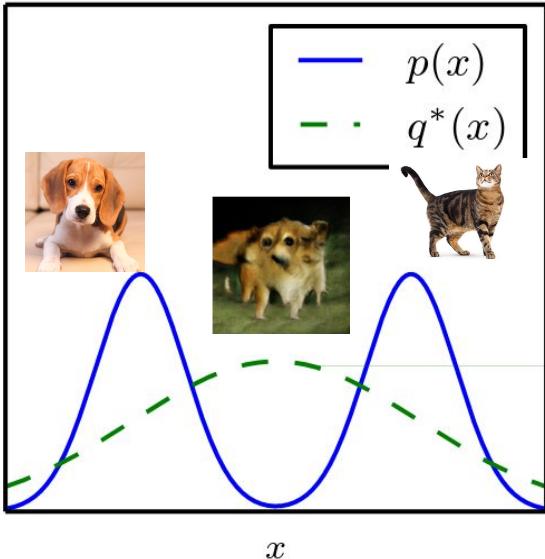


Reverse KL

Consider  $p(x)$  being the real data distribution,  
and  $q^*(x)$  being the generated data distribution

$$q^* = \operatorname{argmin}_q D_{\text{KL}}(p\|q)$$

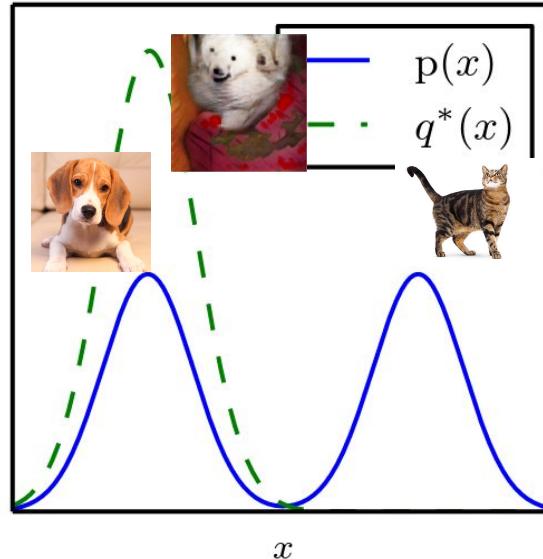
Probability Density



Maximum likelihood

$$q^* = \operatorname{argmin}_q D_{\text{KL}}(q\|p)$$

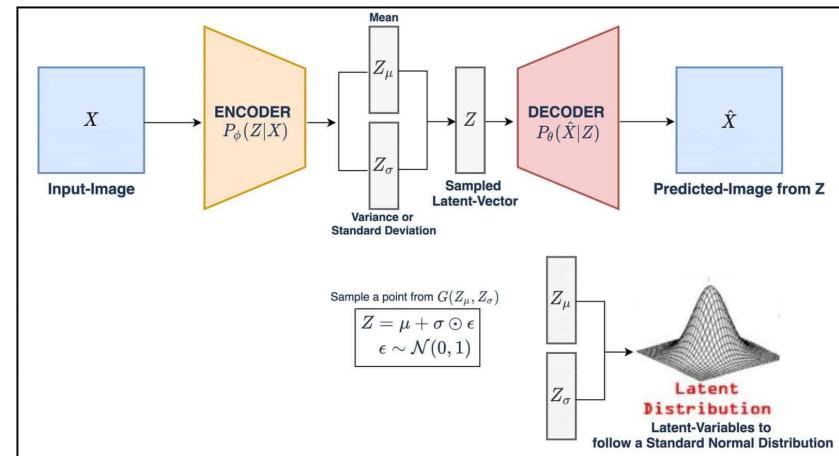
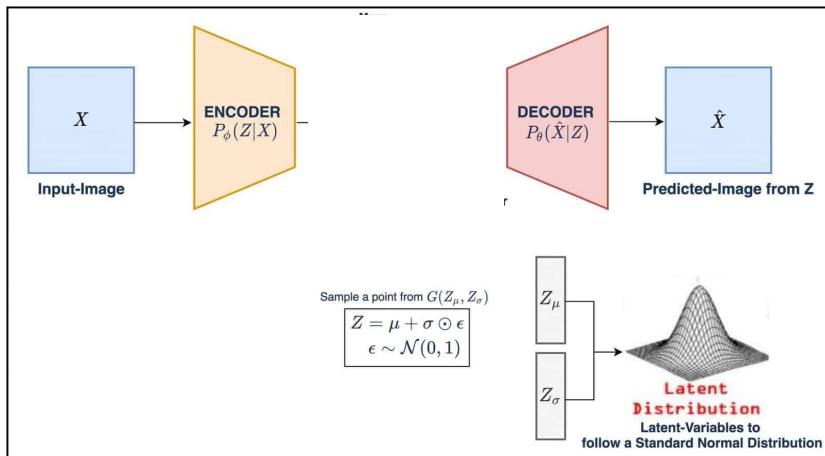
Probability Density

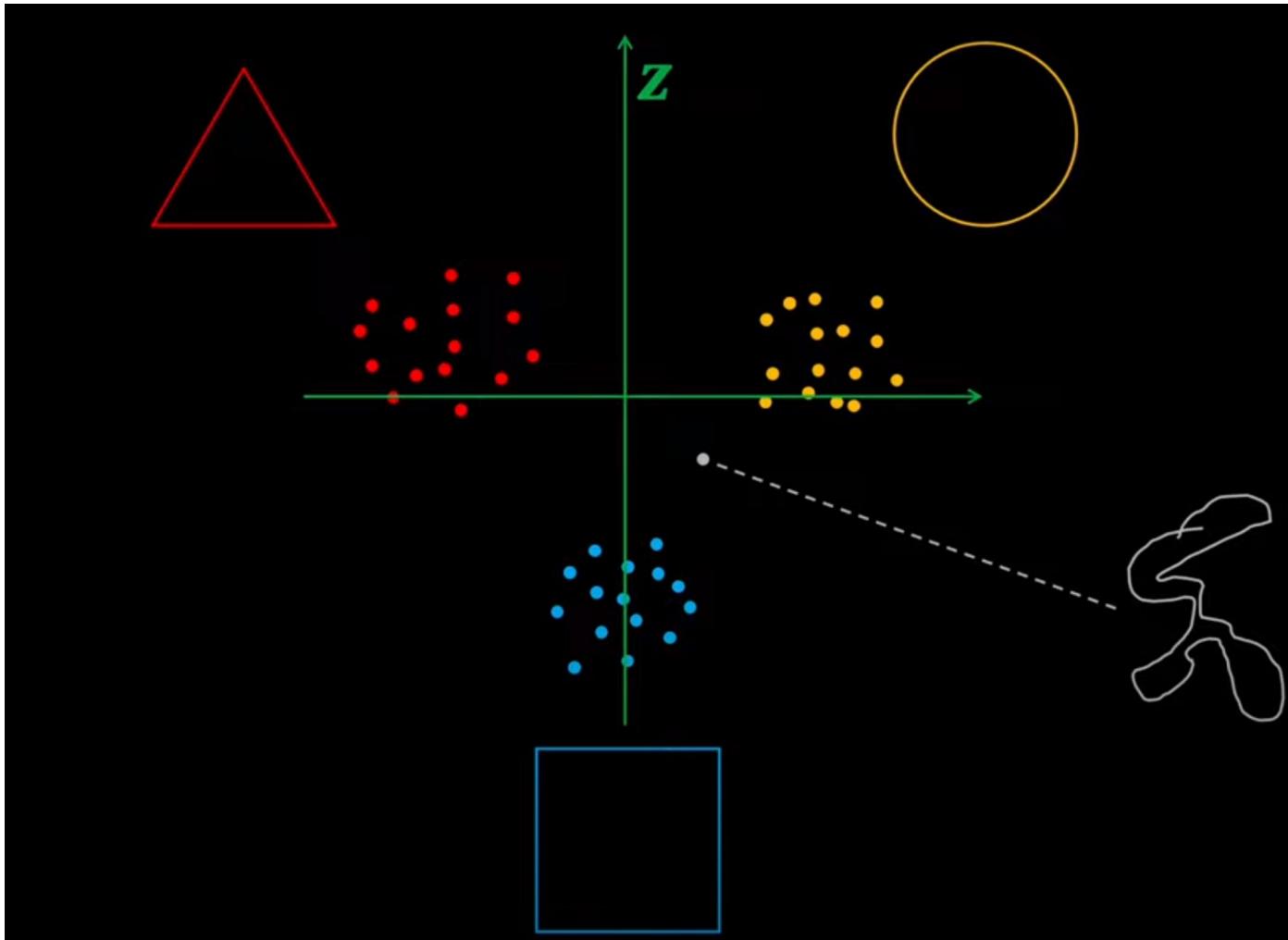


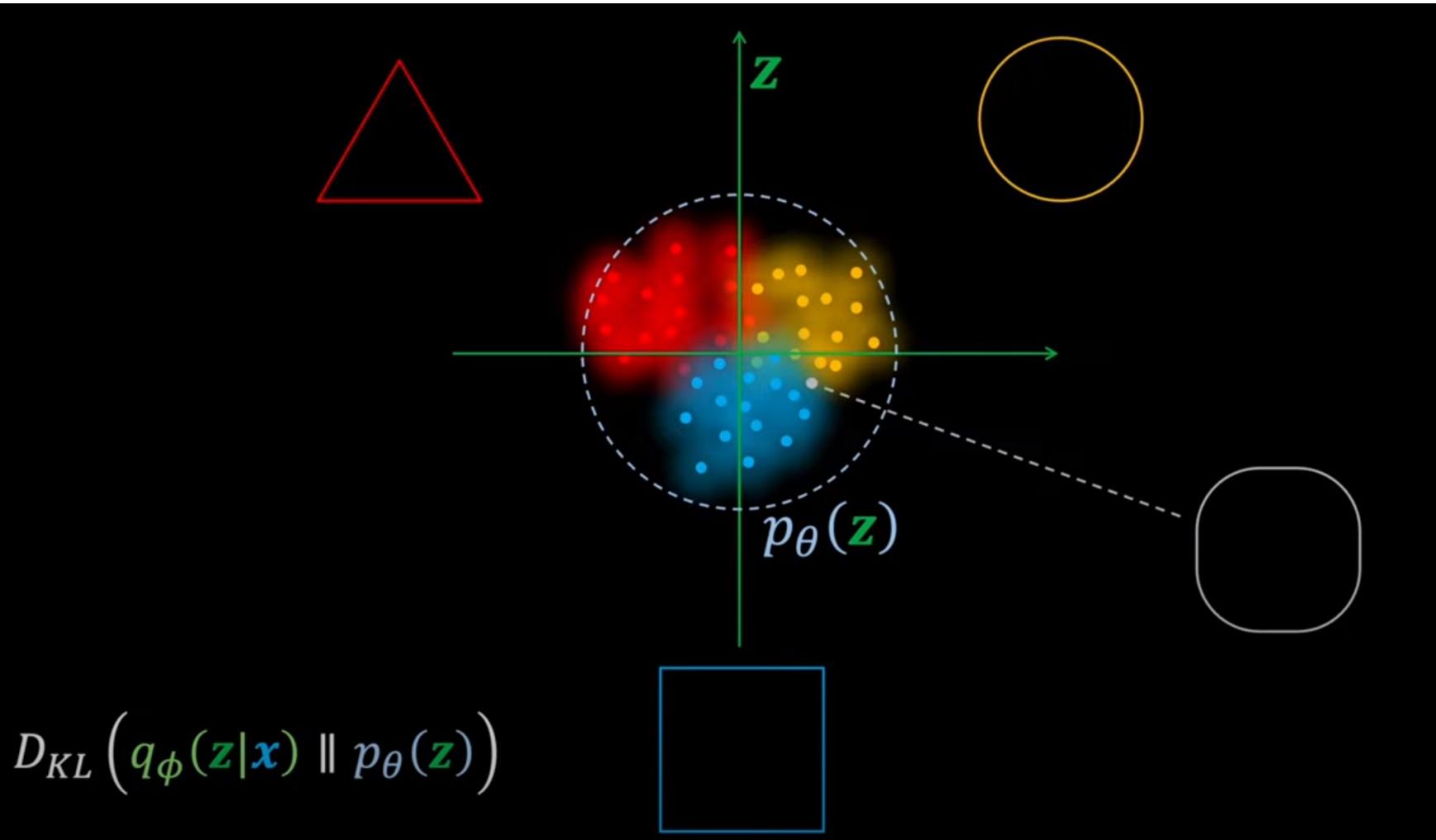
Reverse KL

*“Mode Dropping”*

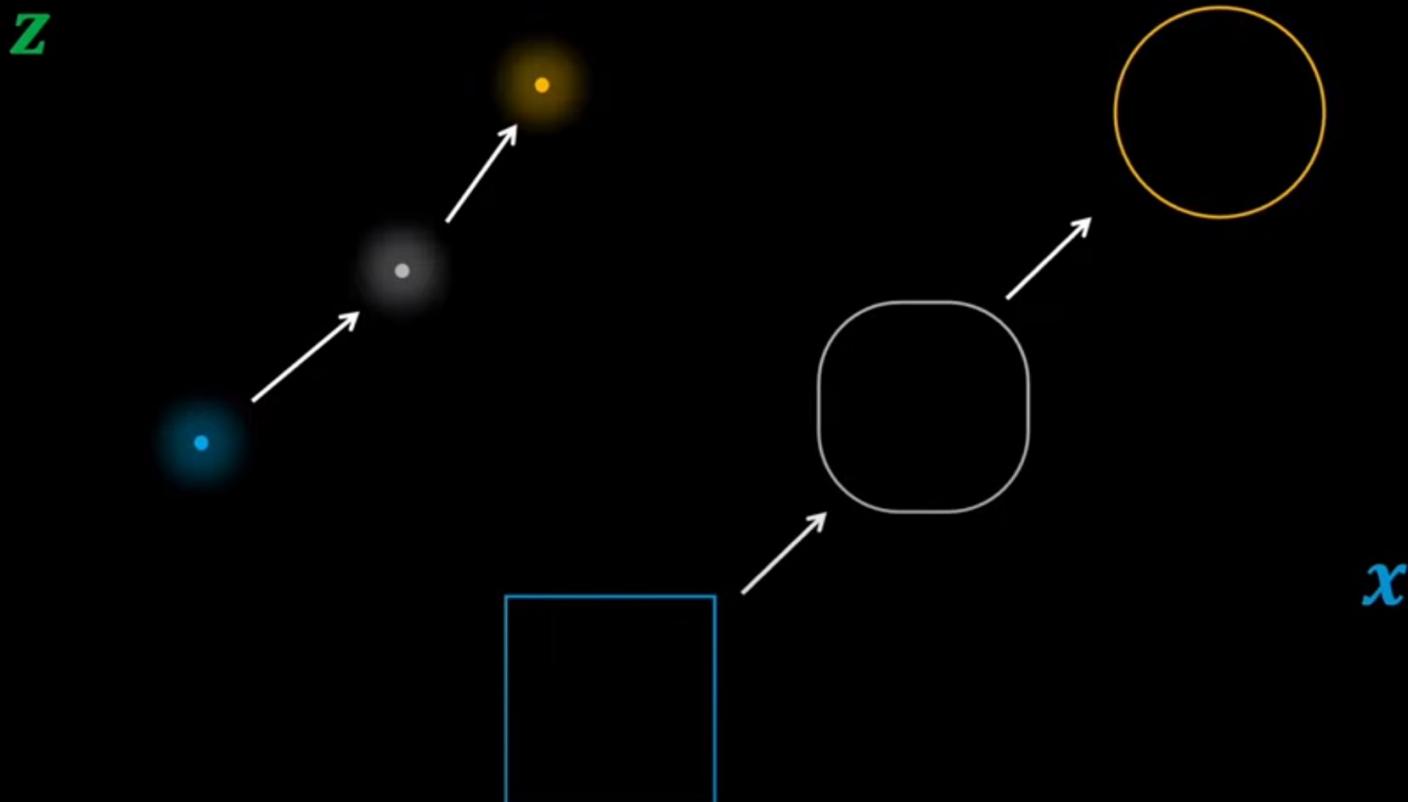
# Differences + Similarities between AE's + VAE's







# Regularized latent space: Continuity



# Loss balance

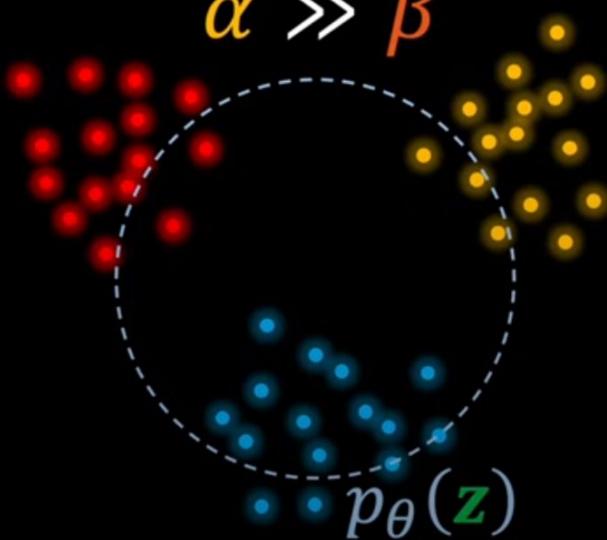
VAE loss

$$\mathcal{L} = \mathcal{L}_{\text{recons}} + \mathcal{L}_{\text{KL}}$$

Weighted loss

$$\mathcal{L} = \alpha \mathcal{L}_{\text{recons}} + \beta \mathcal{L}_{\text{KL}}$$

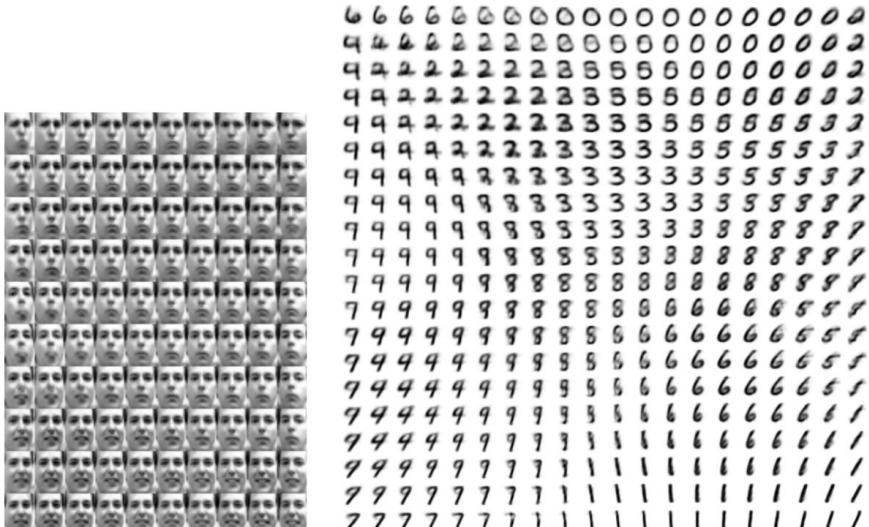
$$\alpha \gg \beta$$



$$\beta \gg \alpha$$



# Variational Auto-Encoders



(a) Learned Frey Face manifold

(b) Learned MNIST manifold

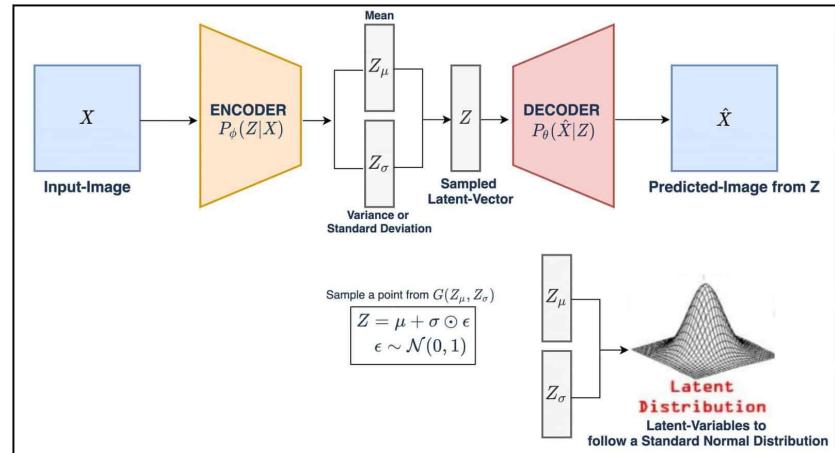
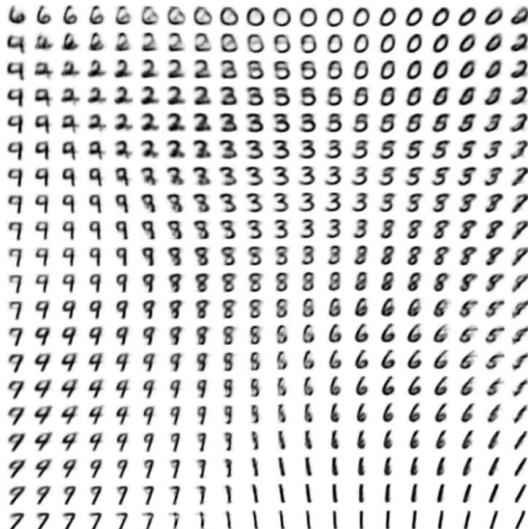


Figure 4: Visualisations of learned data manifold for generative models with two-dimensional latent space, learned with AEVB. Since the prior of the latent space is Gaussian, linearly spaced coordinates on the unit square were transformed through the inverse CDF of the Gaussian to produce values of the latent variables  $\mathbf{z}$ . For each of these values  $\mathbf{z}$ , we plotted the corresponding generative  $p_\theta(\mathbf{x}|\mathbf{z})$  with the learned parameters  $\theta$ .

# Variational Auto-Encoders



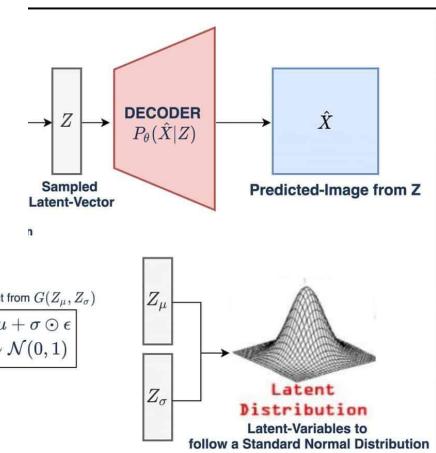
(a) Learned Frey Face manifold



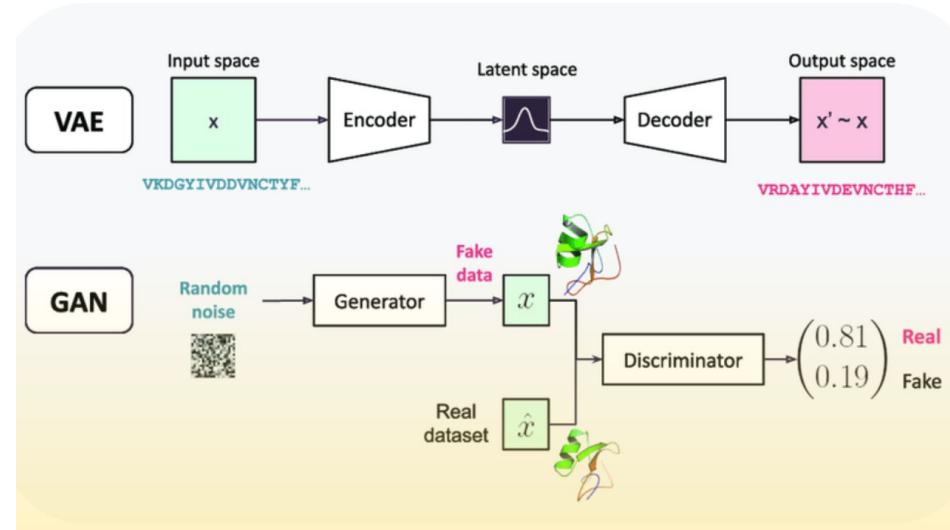
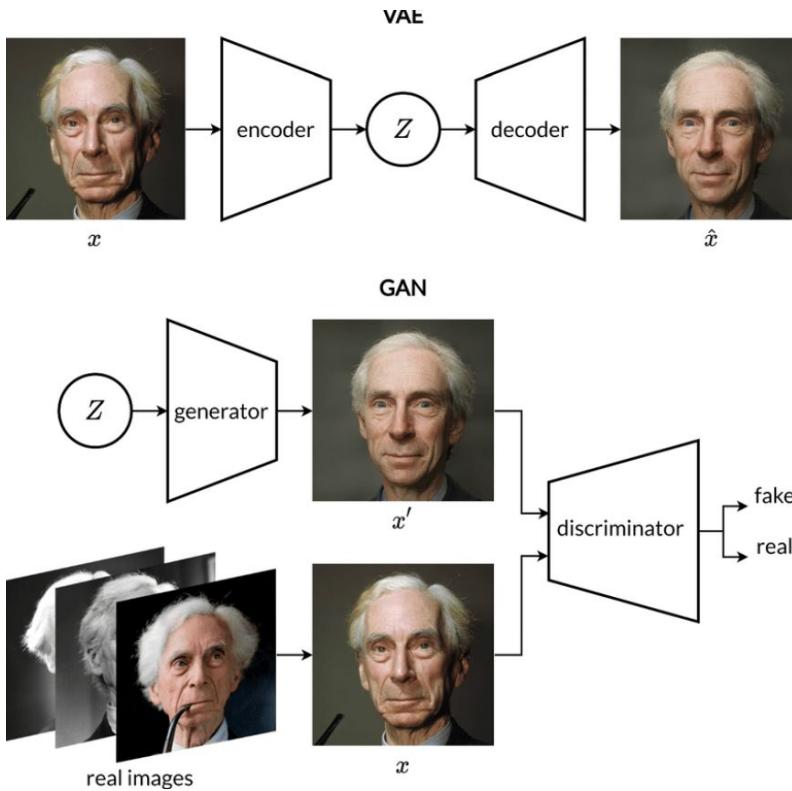
(b) Learned MNIST manifold

Figure 4: Visualisations of learned data manifold for generative models with two-dimensional latent space, learned with AEVB. Since the prior of the latent space is Gaussian, linearly spaced coordinates on the unit square were transformed through the inverse CDF of the Gaussian to produce values of the latent variables  $\mathbf{z}$ . For each of these values  $\mathbf{z}$ , we plotted the corresponding generative  $p_{\theta}(\mathbf{x}|\mathbf{z})$  with the learned parameters  $\theta$ .

*Sample any point from the N-Dimensional Gaussian and render the sample.*



# Differences + Similarities between GAN's and VAE's



# Training the VAE

Reverse KL

$$D_{KL}(q_\phi || p_\theta) = \mathbb{E}_{q_\phi} \left[ \log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z}|\mathbf{x})} \right]$$

Minimize the KL Divergence  
[Expected Log likelihood ratio]

# Training the VAE

$$\begin{aligned} D_{KL}(q_\phi || p_\theta) &= \mathbb{E}_{q_\phi} \left[ \log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z}|\mathbf{x})} \right] \\ &= \mathbb{E}_{q_\phi} [\log q_\phi(\mathbf{z}|\mathbf{x})] - \mathbb{E}_{q_\phi} [\log p_\theta(\mathbf{z}|\mathbf{x})] \\ &= \mathbb{E}_{q_\phi} [\log q_\phi(\mathbf{z}|\mathbf{x})] - \mathbb{E}_{q_\phi} \left[ \log \frac{p_\theta(\mathbf{z}, \mathbf{x})}{p_\theta(\mathbf{x})} \right] \\ &= \mathbb{E}_{q_\phi} [\log q_\phi(\mathbf{z}|\mathbf{x})] - \mathbb{E}_{q_\phi} [\log p_\theta(\mathbf{z}, \mathbf{x})] + \mathbb{E}_{q_\phi} [\log p_\theta(\mathbf{x})] \\ &= \mathbb{E}_{q_\phi} [\log q_\phi(\mathbf{z}|\mathbf{x})] - \mathbb{E}_{q_\phi} [\log p_\theta(\mathbf{z}, \mathbf{x})] + \int q_\phi(\mathbf{z}|\mathbf{x}) \log p_\theta(\mathbf{x}) d\mathbf{z} \\ &= \mathbb{E}_{q_\phi} [\log q_\phi(\mathbf{z}|\mathbf{x})] - \mathbb{E}_{q_\phi} [\log p_\theta(\mathbf{z}, \mathbf{x})] + \log p_\theta(\mathbf{x}) \int q_\phi(\mathbf{z}|\mathbf{x}) d\mathbf{z} \\ &= \mathbb{E}_{q_\phi} [\log q_\phi(\mathbf{z}|\mathbf{x})] - \mathbb{E}_{q_\phi} [\log p_\theta(\mathbf{z}, \mathbf{x})] + \log p_\theta(\mathbf{x}) \end{aligned}$$

# Training the VAE

$$\log p_{\theta}(\mathbf{x}) = -\mathbb{E}_{q_{\phi}} [\log q_{\phi}(\mathbf{z}|\mathbf{x})] + \mathbb{E}_{q_{\phi}} [\log p_{\theta}(\mathbf{z}, \mathbf{x})] + D_{KL}(q_{\phi}||p_{\theta})$$

Marginal log Likelihood      Component 1      Component 2  
Log Evidence                    ↓  
intractable

$\geq 0$

# Training the VAE

$$\log p_{\theta}(\mathbf{x}) = -\mathbb{E}_{q_{\phi}} [\log q_{\phi}(\mathbf{z}|\mathbf{x})] + \mathbb{E}_{q_{\phi}} [\log p_{\theta}(\mathbf{z}, \mathbf{x})] + D_{KL}(q_{\phi}||p_{\theta})$$

Marginal log Likelihood  
Log Evidence

Component 1

Component 2

If I maximize this

Indirectly I will minimize this

$\geq 0$



intractable

# Training the VAE

$$\begin{aligned}\text{ELBO} &= -\mathbb{E}_{q_\phi} [\log q_\phi(\mathbf{z}|\mathbf{x})] + \mathbb{E}_{q_\phi} [\log p_{\color{blue}\theta}(\mathbf{z}, \mathbf{x})] \\&= -\mathbb{E}_{q_\phi} [\log q_\phi(\mathbf{z}|\mathbf{x})] + \mathbb{E}_{q_\phi} [\log p_{\color{violet}\theta}(\mathbf{x}|\mathbf{z})] + \mathbb{E}_{q_\phi} [\log p_{\color{blue}\theta}(\mathbf{z})] \\&= \mathbb{E}_{q_\phi} [\log p_{\color{violet}\theta}(\mathbf{x}|\mathbf{z})] - \mathbb{E}_{q_\phi} [\log q_\phi(\mathbf{z}|\mathbf{x})] + \mathbb{E}_{q_\phi} [\log p_{\color{blue}\theta}(\mathbf{z})] \\&= \mathbb{E}_{q_\phi} [\log p_{\color{violet}\theta}(\mathbf{x}|\mathbf{z})] - \mathbb{E}_{q_\phi} \left[ \log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p_{\color{blue}\theta}(\mathbf{z})} \right]\end{aligned}$$

# Training the VAE

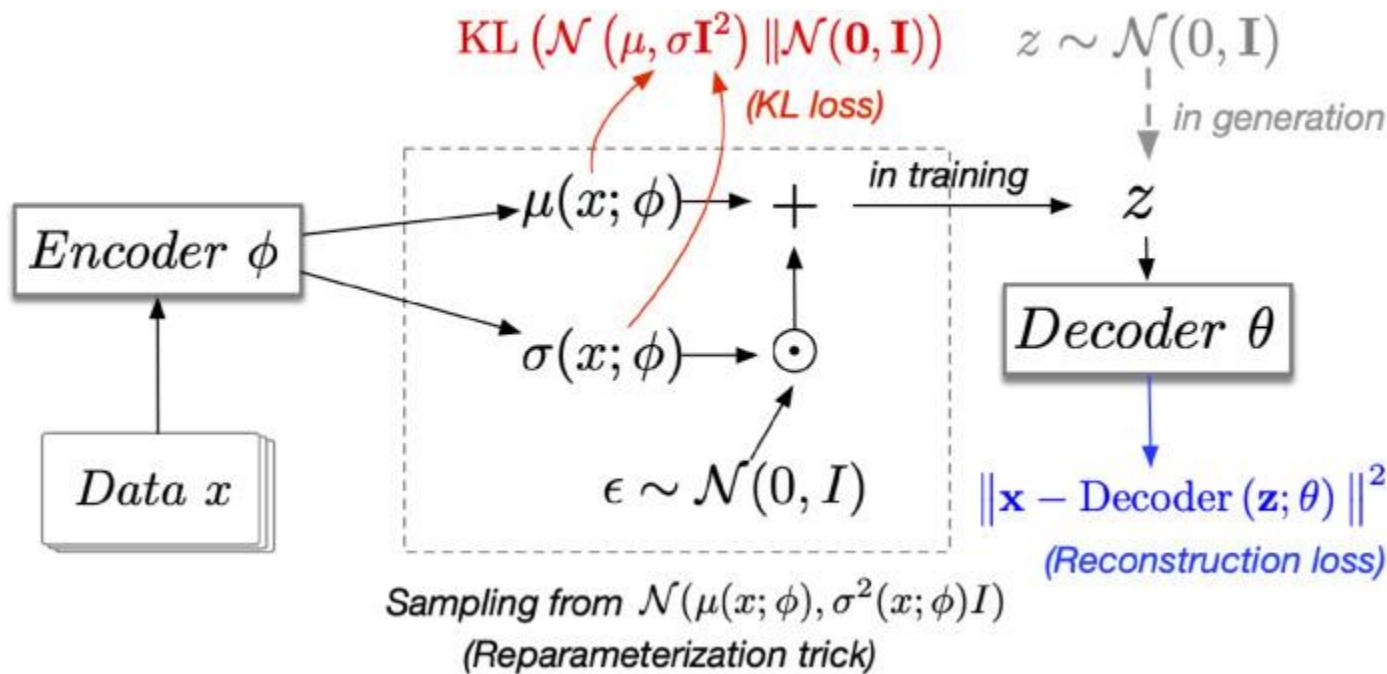
$$\text{ELBO} = \mathbb{E}_{q_\phi} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \mathbb{E}_{q_\phi} \left[ \log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p_{\theta}(\mathbf{z})} \right]$$

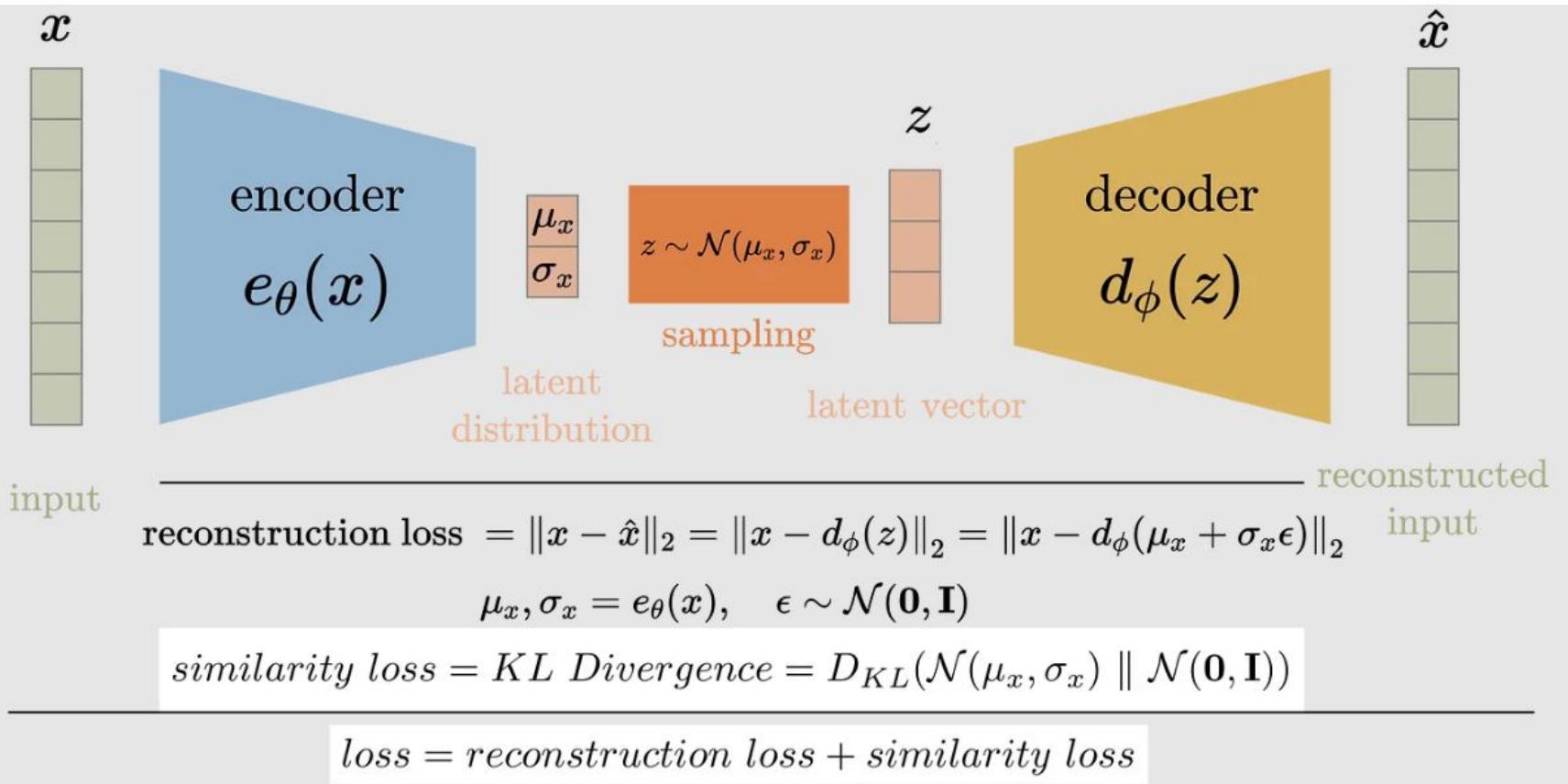


Expected reconstruction error



KL Divergence between  
approx. posterior & the prior





# TOWARDS METAMERISM VIA FOVEATED STYLE TRANSFER

Arturo Deza<sup>1,4</sup>, Aditya Jonnalagadda<sup>3</sup>, Miguel P. Eckstein<sup>1,2,4</sup>

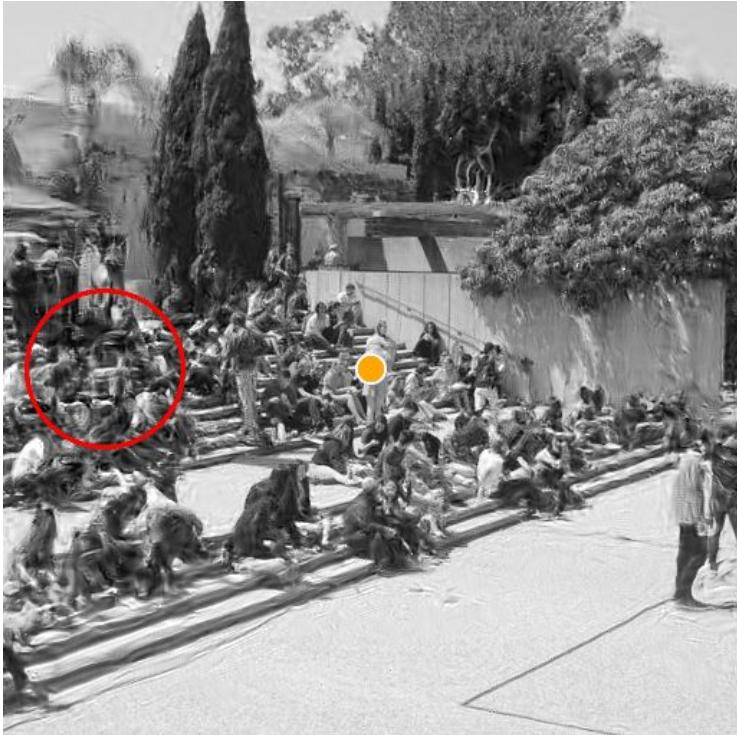
<sup>1</sup> Dynamical Neuroscience, <sup>2</sup>Psychological and Brain Sciences,

<sup>3</sup>Electric and Computer Engineering, <sup>4</sup> Institute for Collaborative Biotechnologies  
UC Santa Barbara, CA, USA

deza@dync.ucs.eds, aditya\_jonnalagada@ece.ucs.eds, eckstein@psych.ucs.eds

## ABSTRACT

The problem of *visual metamerism* is defined as finding a family of perceptually indistinguishable, yet physically different images. In this paper, we propose our NeuroFovea metamer model, a foveated generative model that is based on a mixture of peripheral representations and style transfer forward-pass algorithms. Our gradient-descent free model is parametrized by a foveated VGG19 encoder-decoder which allows us to encode images in high dimensional space and interpolate between the content and texture information with adaptive instance normalization anywhere in the visual field. Our contributions include: 1) A framework for computing metamers that resembles a noisy communication system via a foveated feed-forward encoder-decoder network – We observe that metamerism arises as a byproduct of noisy perturbations that partially lie in the perceptual null space; 2) A perceptual optimization scheme as a solution to the hyperparametric nature of our metamer model that requires tuning of the image-texture tradeoff coefficients everywhere in the visual field which are a consequence of internal noise; 3) An ABX psychophysical evaluation of our metamers where we also find that the rate of growth of the receptive fields in our model match V1 for reference metamers and V2 between synthesized samples. Our model also renders metamers at roughly a second, presenting a  $\times 1000$  speed-up compared to the previous work, which allows for tractable data-driven metamer experiments.

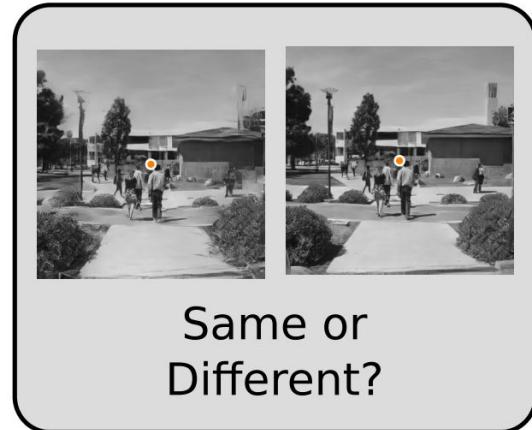




Image



Metamer



Same or  
Different?

Figure 1: Two visual metamers are physically different images that when fixated on the orange dot (center), should remain perceptually indistinguishable to each other for an observer. Colored circles highlight different distortions in the visual field that observers do not perceive in our model.

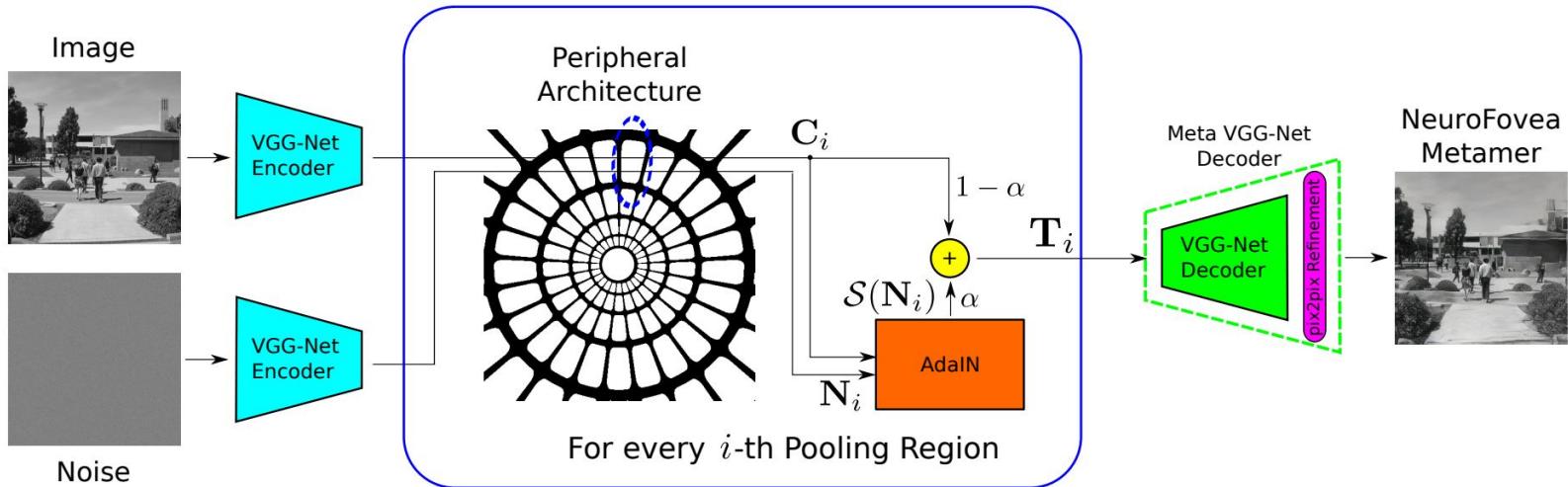
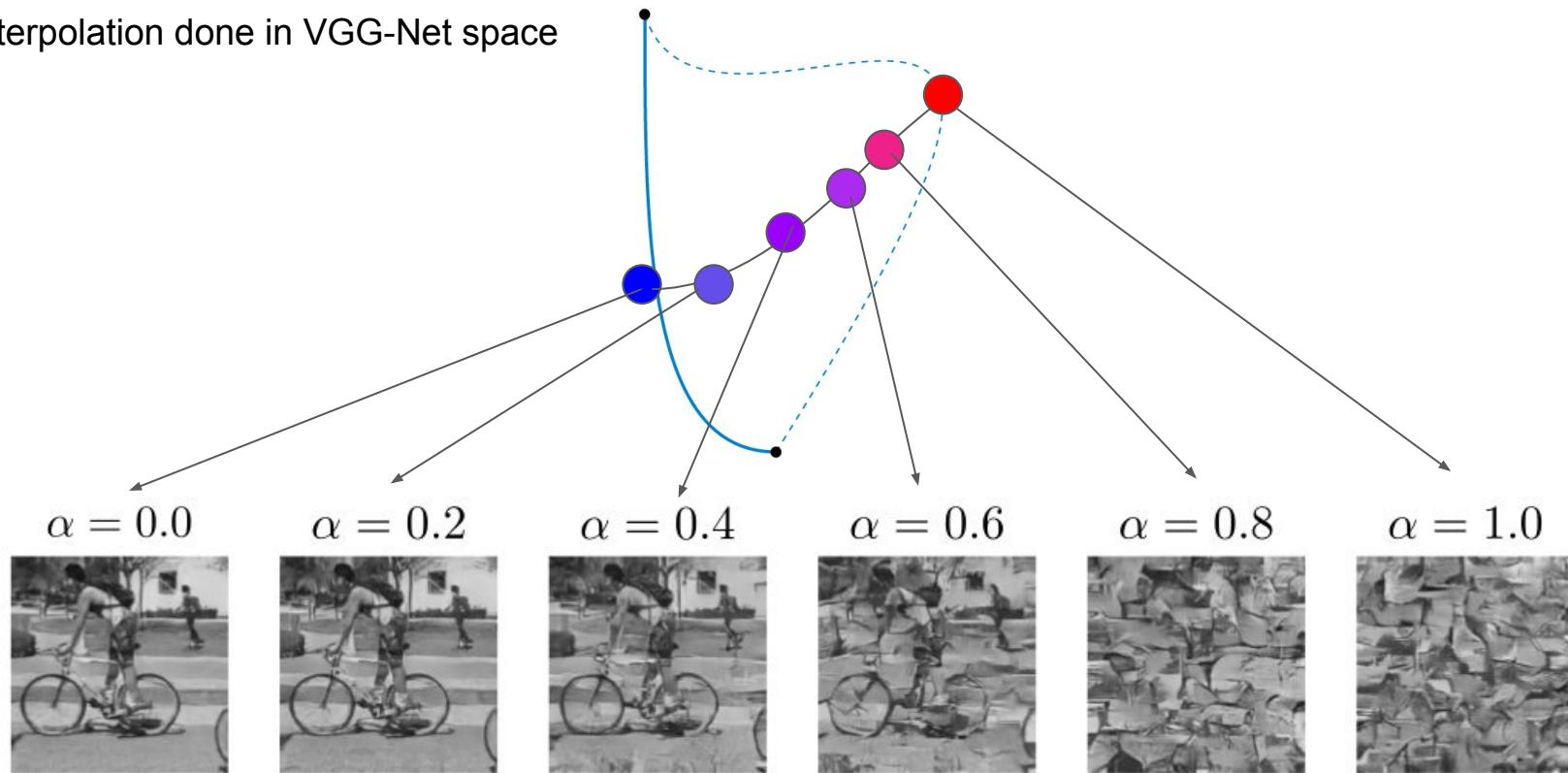


Figure 2: The NeuroFovea metamer generation schematic: An input image and a noise patch are fed through a VGG-Net encoder into a new feature space. Through spatial control we can produce an interpolation for each pooling region in such feature space between the stylized-noise (texture), and the content (the input image). This is how we successfully impose both global image and local texture-like constraints in every pooling region. The metamer is the output of the pooled (and interpolated) feature vector through the Meta VGG-Net Decoder.

$$\mathbf{M}(I|\mathcal{N};\bar{\alpha}) = \mathcal{D}(\mathcal{E}_{\Sigma}(I|\mathcal{N};\bar{\alpha})) = \mathcal{D}\left(\sum_{i=1}^k w_i[(1-\alpha_i)\mathcal{E}_i(I) + \alpha_i\mathcal{S}(\mathcal{E}_i(\mathcal{N}))]\right)$$

Interpolation done in VGG-Net space



## Case Study:

What is the maximum amount of texture-driven distortion we can put in this **receptive field**?

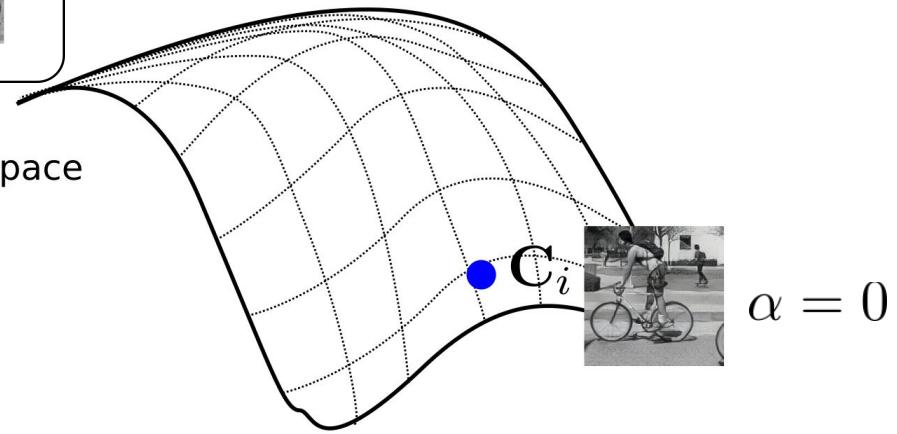


## Case Study:

What is the maximum amount of texture-driven distortion we can put in this **receptive field**?



Encoded Space



## Case Study:

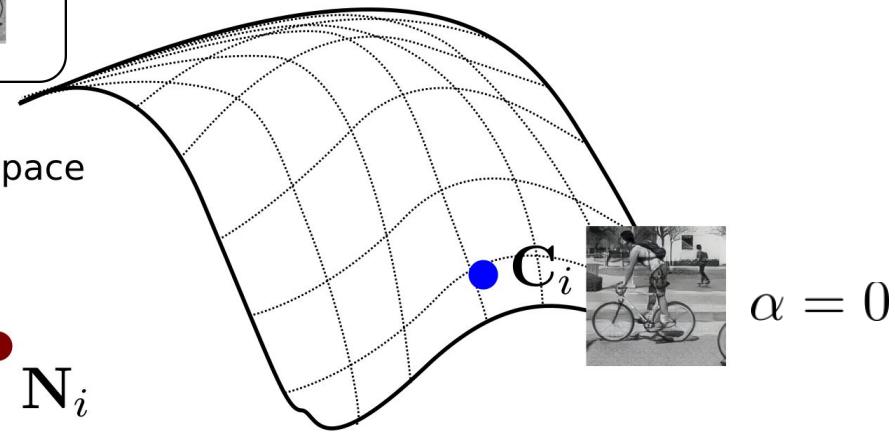
What is the maximum amount of texture-driven distortion we can put in this **receptive field**?



Encoded Space



$\mathbf{N}_i$



## Case Study:

What is the maximum amount of texture-driven distortion we can put in this **receptive field**?



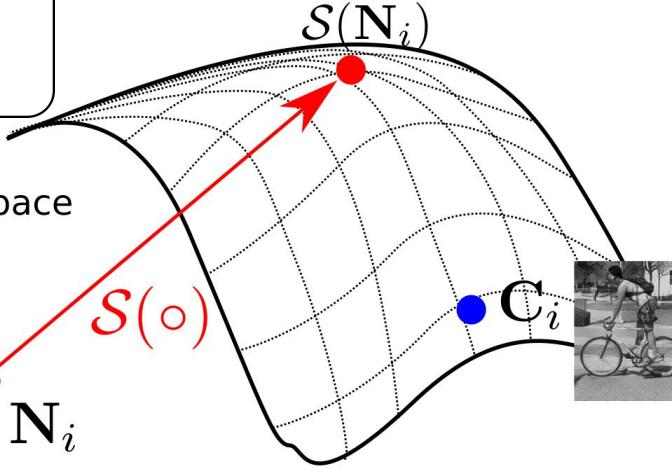
$$\alpha = 1$$

Encoded Space



$N_i$

$S(\circ)$



$$\alpha = 0$$

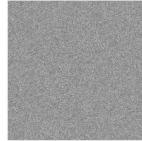
## Case Study:

What is the maximum amount of texture-driven distortion we can put in this **receptive field**?



$$\alpha = 1$$

Encoded Space



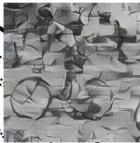
$\mathbf{N}_i$

$\mathcal{S}(\circ)$

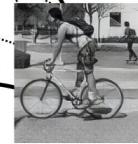
$\mathcal{S}(\mathbf{N}_i)$

$\mathbf{T}_i$

$\mathbf{C}_i$



$$\alpha = 0.4$$



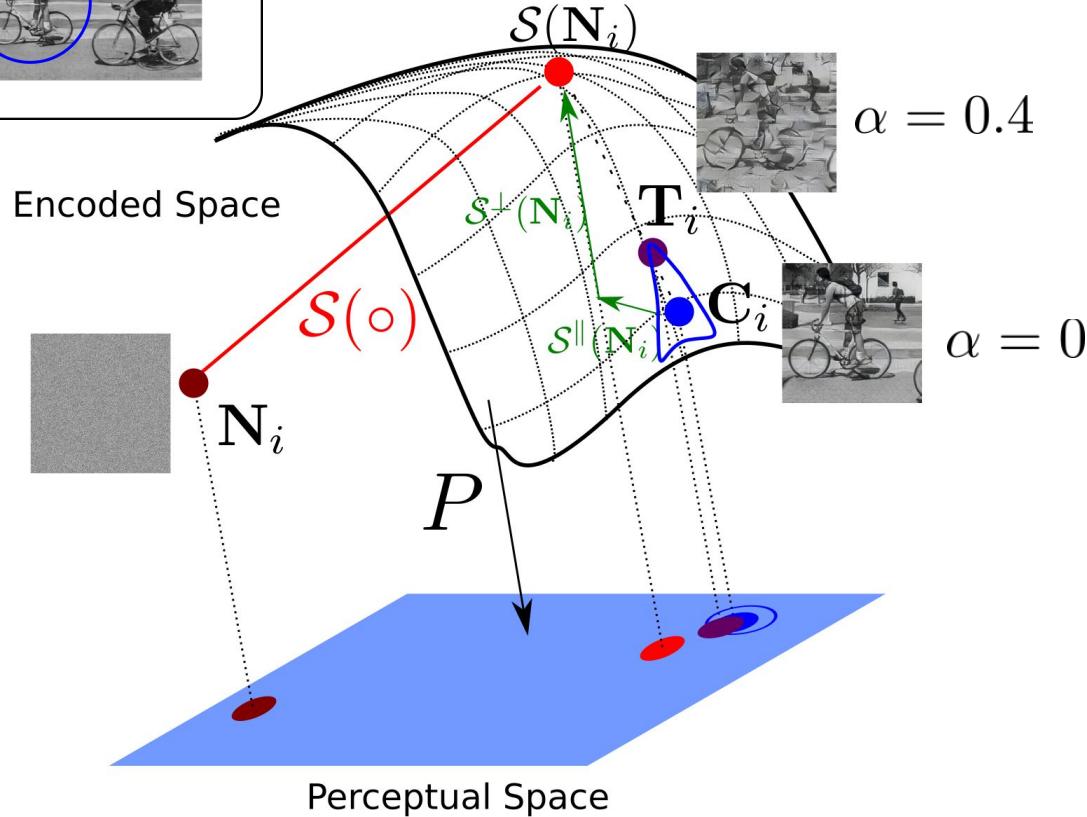
$$\alpha = 0$$

## Case Study:

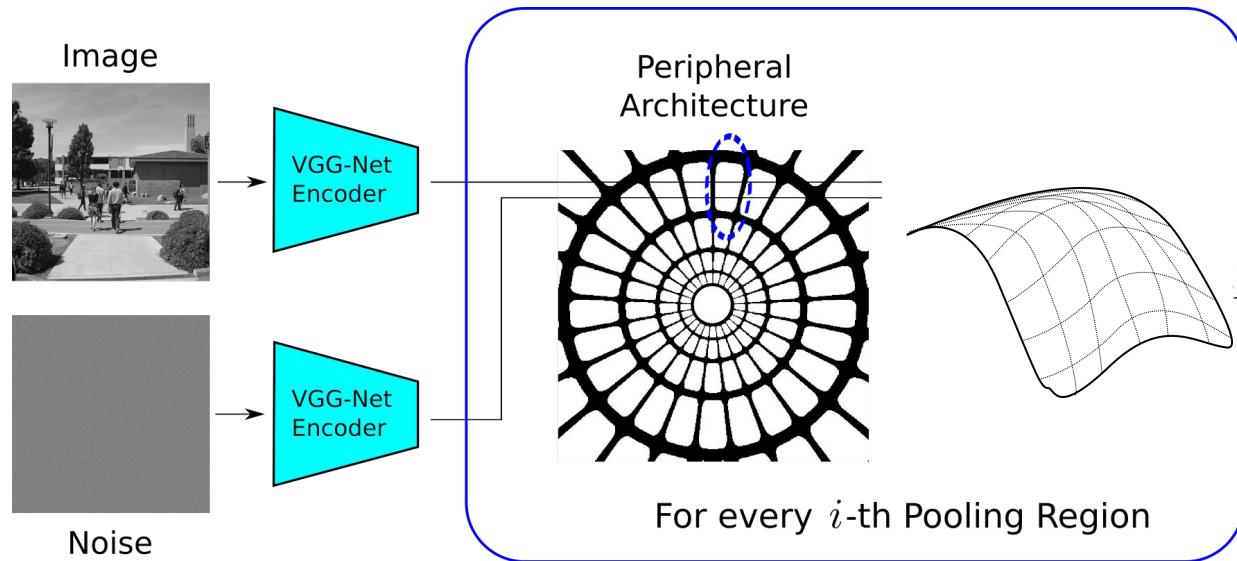
What is the maximum amount of texture-driven distortion we can put in this **receptive field**?



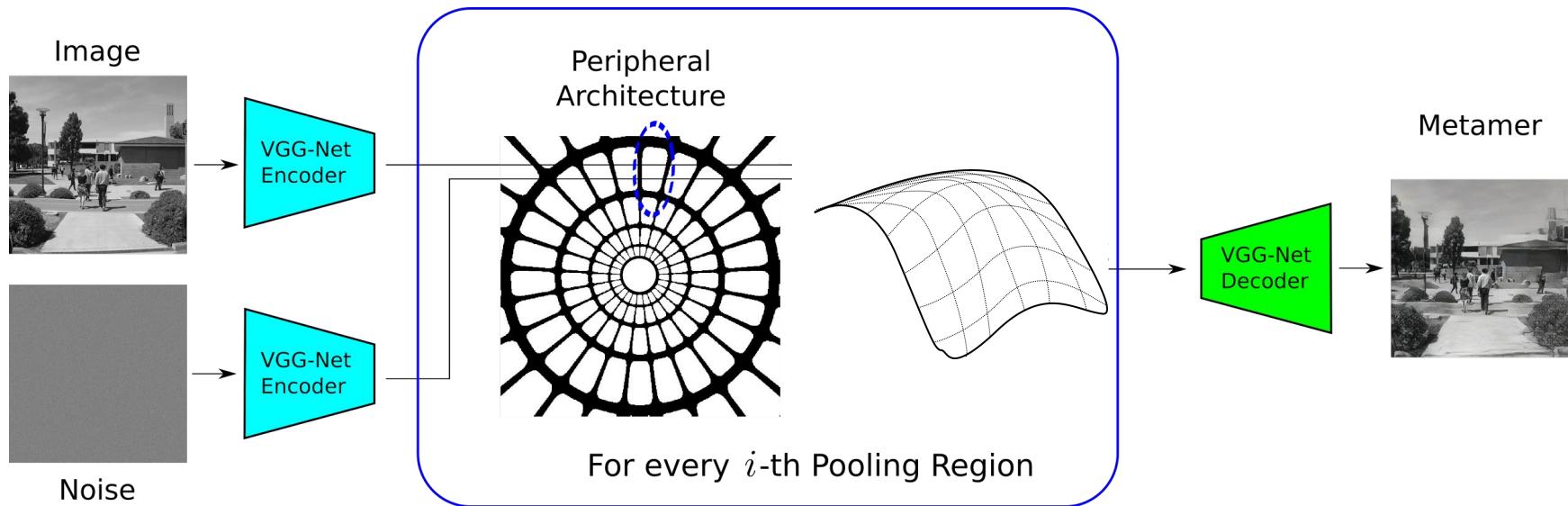
$$\alpha = 1$$



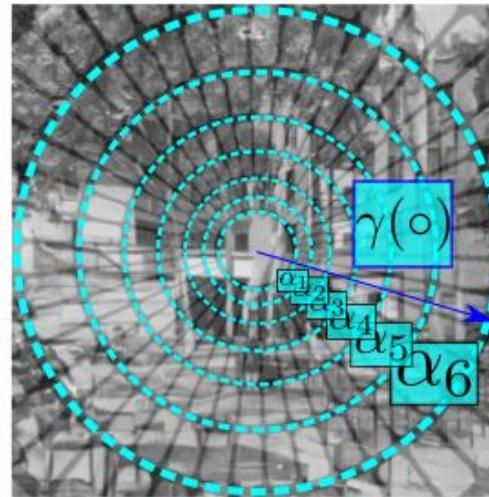
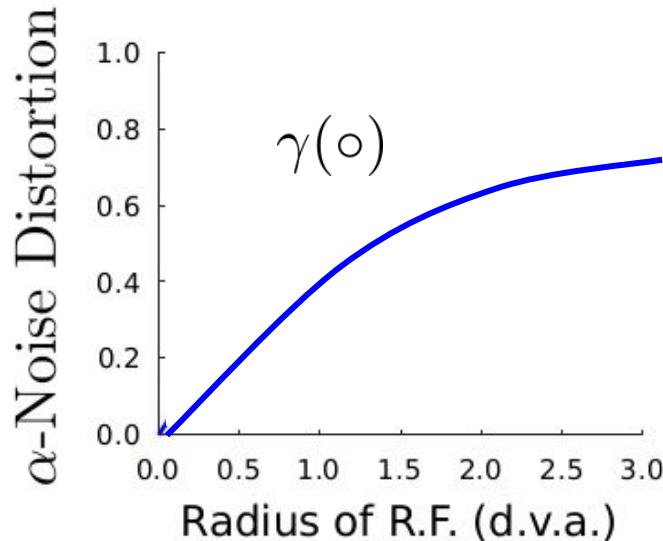
# A new metamer rendering pipeline



# A new metamer rendering pipeline



# How much can we distort in each pooling region?



# Perceptual Optimization Procedure

FS Model

Reference Images



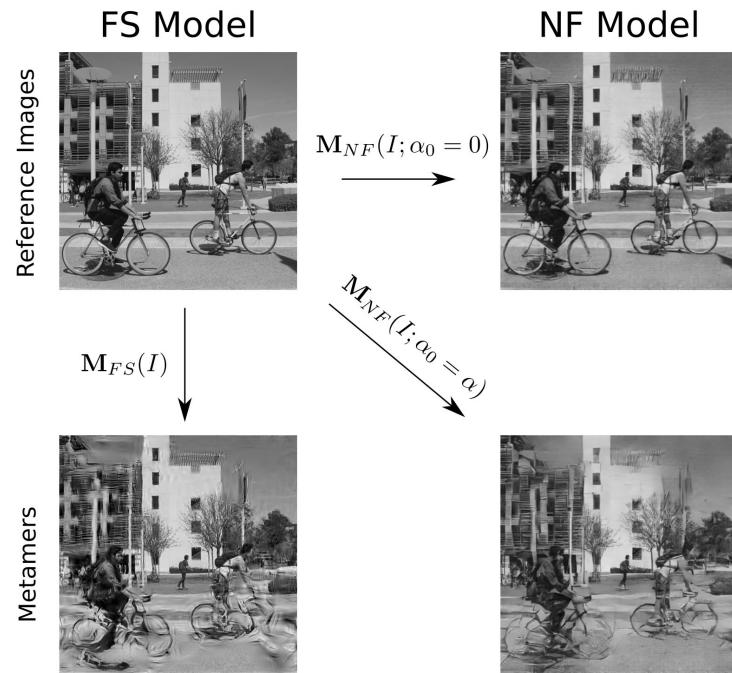
$\mathbf{M}_{FS}(I)$



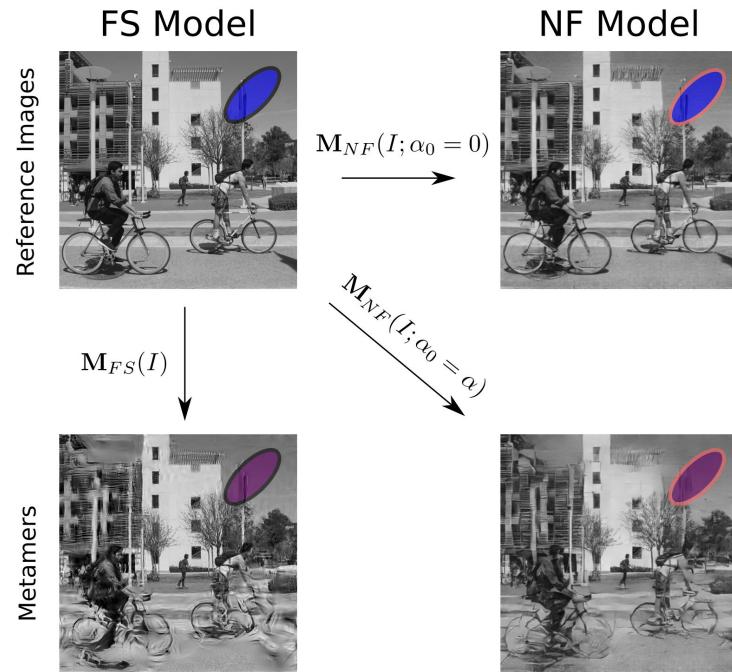
Metamers



# Perceptual Optimization Procedure



# Perceptual Optimization Procedure



# Perceptual Optimization Procedure

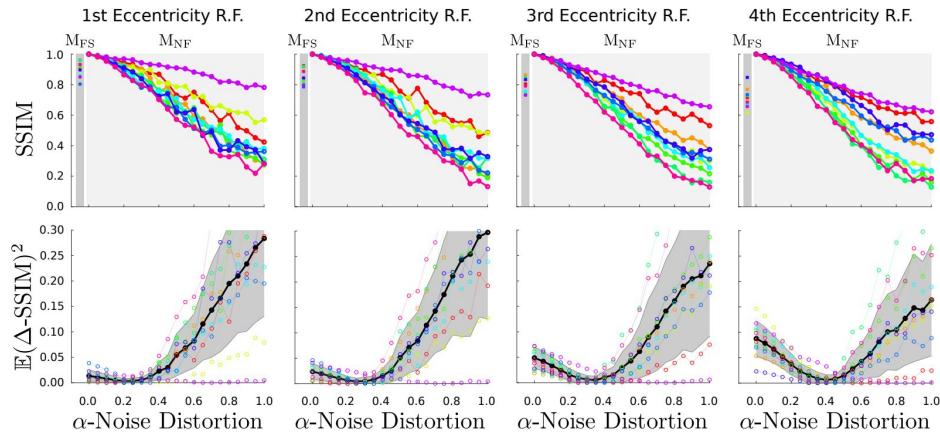
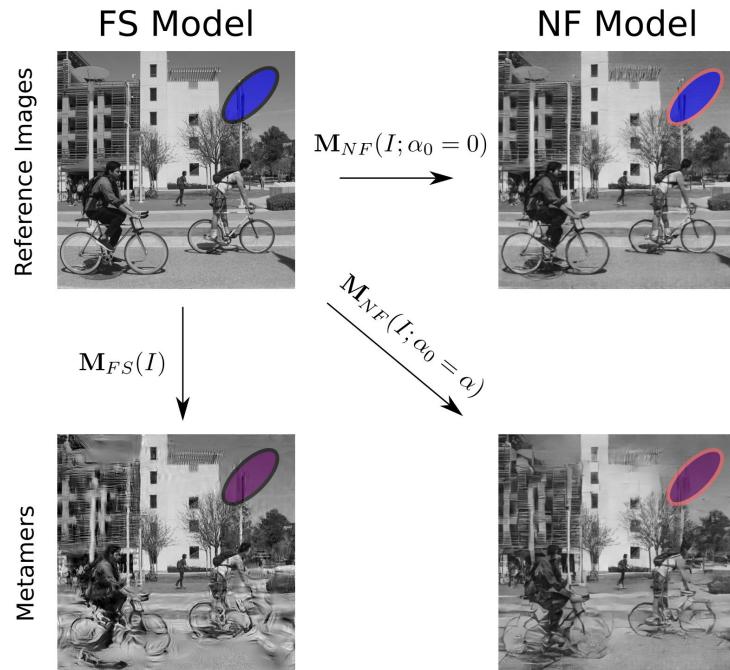
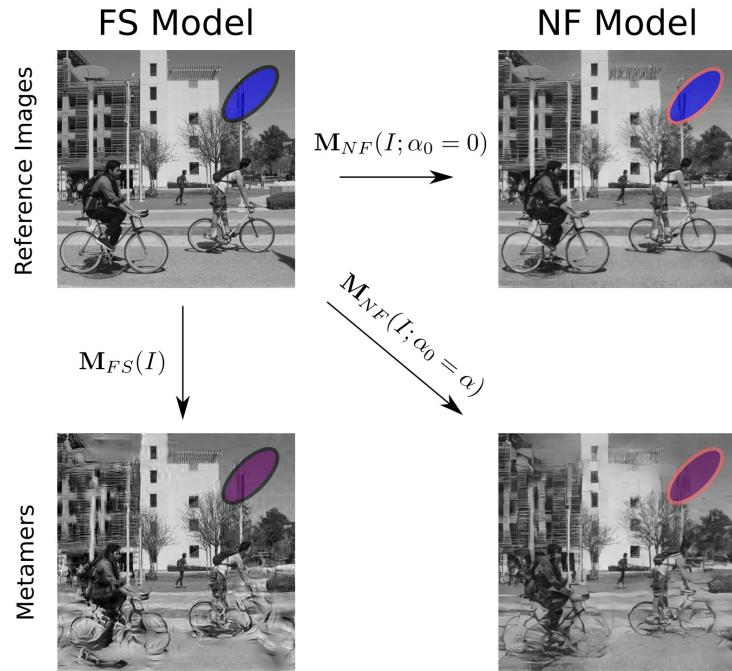


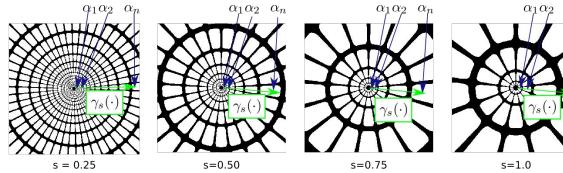
Figure 8: The result of each SSIM (top) for Experiment 1 for a scale of  $s = 0.3$  where we find the critical  $\alpha$  for each receptive field ring as we minimize  $\mathbb{E}(\Delta\text{-SSIM})^2$  (bottom).  $\mathbb{E}(\Delta\text{-SSIM})^2$  is minimized by matching the perceptual distortion of the Freeman & Simoncelli (2011) ( $M_{FS}$ ) and NeuroFovea ( $M_{NF}$ ) metamers in Eq. 9. Each color represents a different  $512 \times 512$  image trajectory, the black line (bottom) shows the average. Only the first 4 eccentricity dependent receptive fields are shown.

# Perceptual Optimization Procedure



Minimize :  $(\text{SSIM}(\text{blue blob}, \text{purple blob}) - \text{SSIM}(\text{blue blob}, \text{pink blob}))^2$   
 for every pooling region in the visual field

Estimation of maximum distortion in the visual field



$$s_0, \bar{\alpha}_0 = \underset{s, \bar{\alpha}}{\operatorname{argmax}} \mathbb{E}[d'(s, \bar{\alpha} | \theta_{obs})]$$

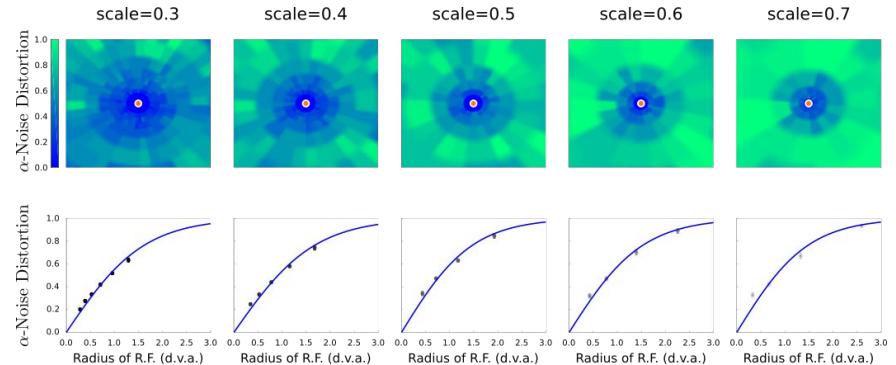
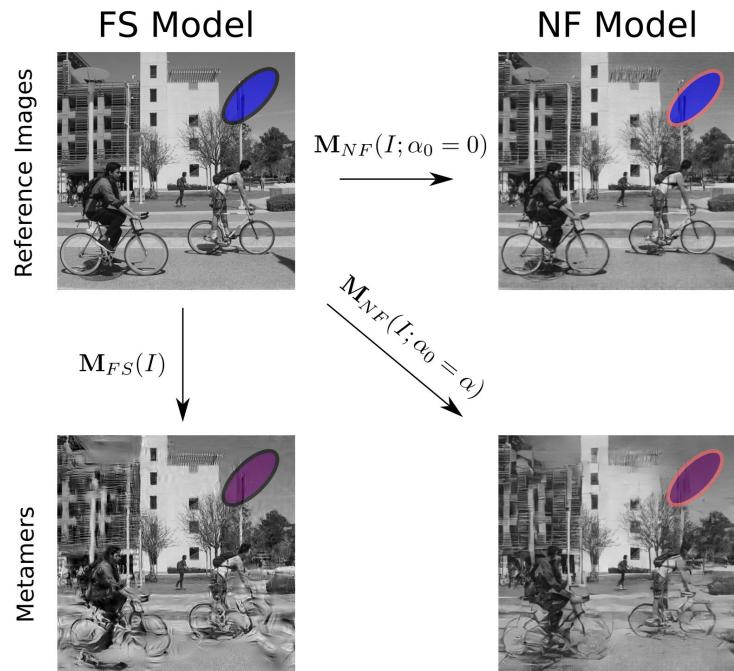
s.t.  $0 < d'(s, \bar{\alpha} | \theta_{obs}) < \epsilon$

$\downarrow \alpha = \gamma(\circ; s)$

$$s_0 = \underset{s}{\operatorname{argmax}} \mathbb{E}[d'(s, \gamma(\circ; s) | \theta_{obs})]$$

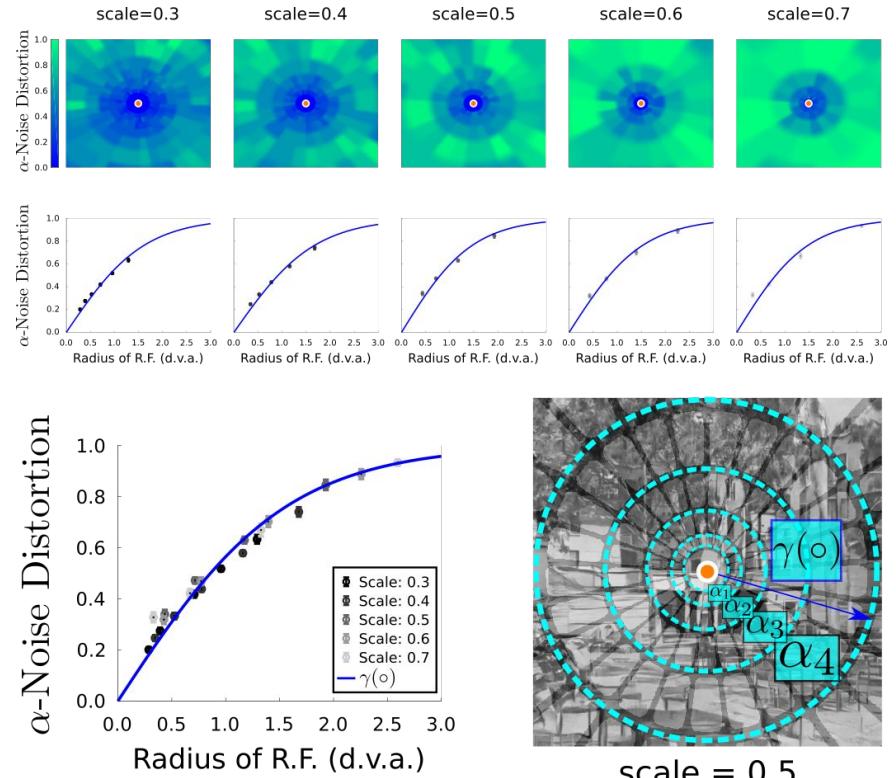
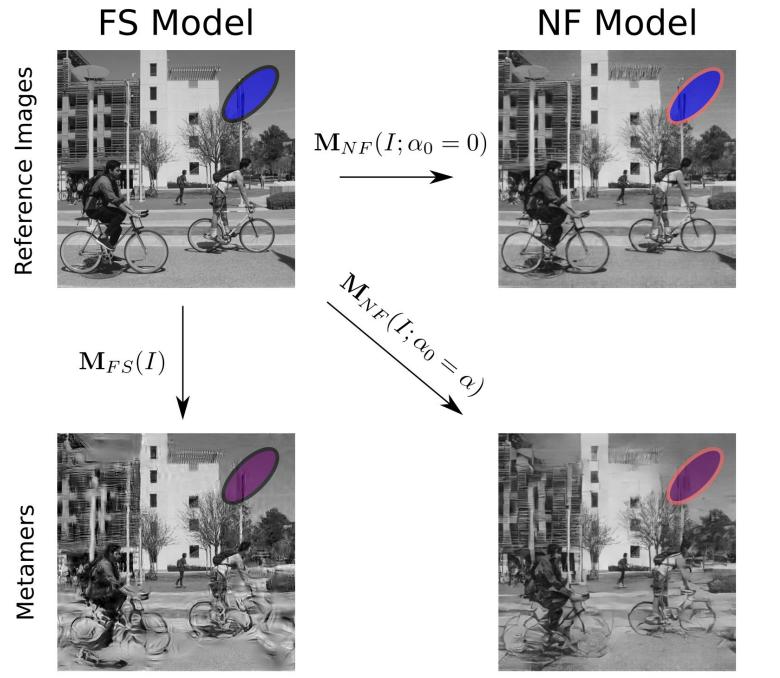
s.t.  $0 < d'(s, \gamma(\circ; s) | \theta_{obs}) < \epsilon$

# Perceptual Optimization Procedure



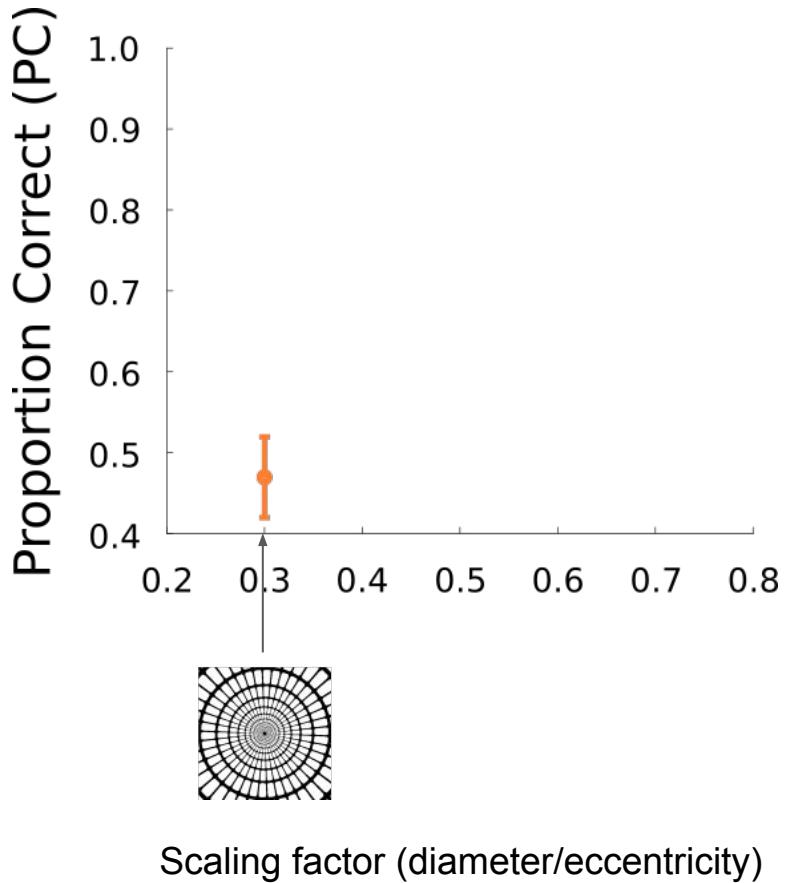
Minimize :  $(\text{SSIM}(\text{blue blob}, \text{purple blob}) - \text{SSIM}(\text{blue blob}, \text{pink blob}))^2$   
for every pooling region in the visual field

# Perceptual Optimization Procedure

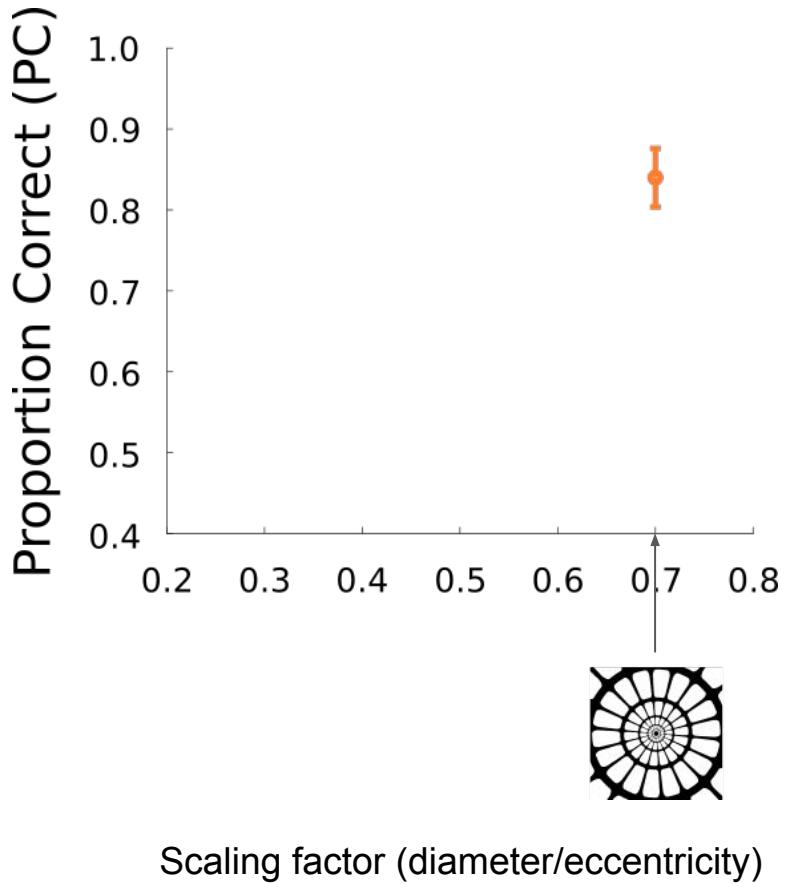


Deza, Jonnalagadda & Eckstein. ICLR 2019.

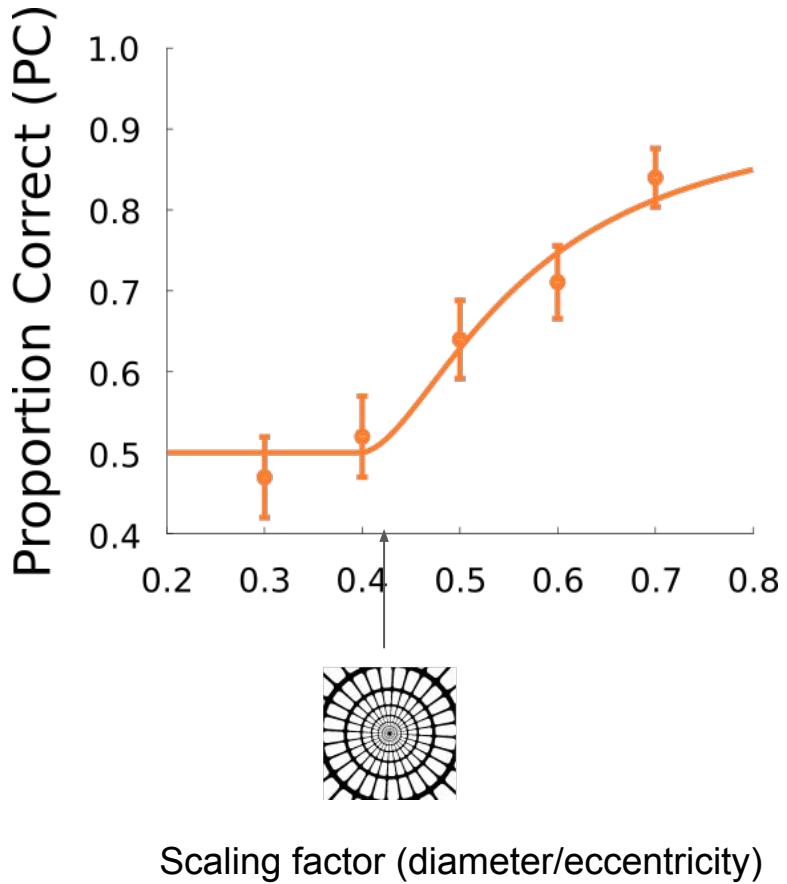
# How do we test this?



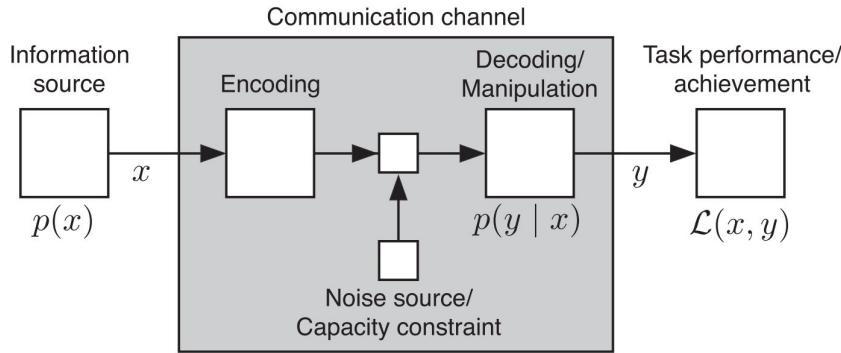
# How do we test this?



# How do we test this?



# Fundamentals of Rate-Distortion Theory



$p(x)$ : Statistics of information source

$p(y|x)$ : Capacity-limited information channel

$\mathcal{L}(x,y)$ : Cost of communication error

**Fig. 1.** The core constructs of rate-distortion theory. An information source is described by a probability distribution over its alphabet. Samples from this source are communicated over a noisy or capacity-limited channel, resulting in a conditional probability distribution over the channel output. A cost function defines the consequences of error in communication.

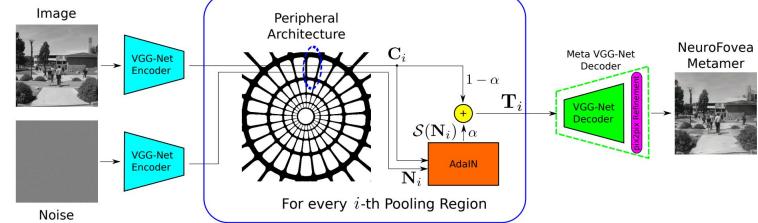
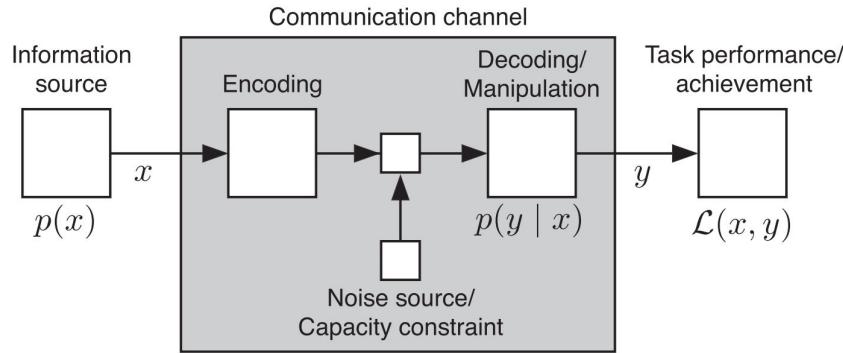


Figure 2: The NeuroFovea metamer generation schematic: An input image and a noise patch are fed through a VGG-Net encoder into a new feature space. Through spatial control we can produce an interpolation for each pooling region in such feature space between the stylized-noise (texture), and the content (the input image). This is how we successfully impose both global image and local texture-like constraints in every pooling region. The metamer is the output of the pooled (and interpolated) feature vector through the Meta VGG-Net Decoder.

# Fundamentals of Rate-Distortion Theory



$p(x)$ : Statistics of information source

$p(y | x)$ : Capacity-limited information channel

$\mathcal{L}(x, y)$ : Cost of communication error

**Fig. 1.** The core constructs of rate-distortion theory. An information source is described by a probability distribution over its alphabet. Samples from this source are communicated over a noisy or capacity-limited channel, resulting in a conditional probability distribution over the channel output. A cost function defines the consequences of error in communication.

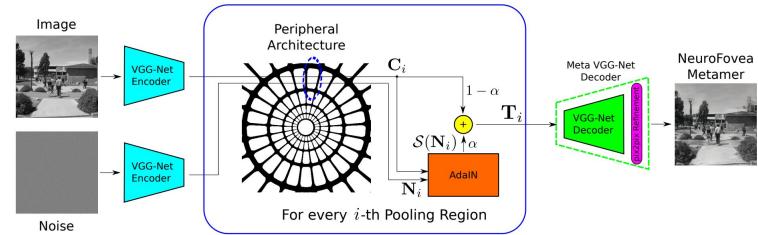
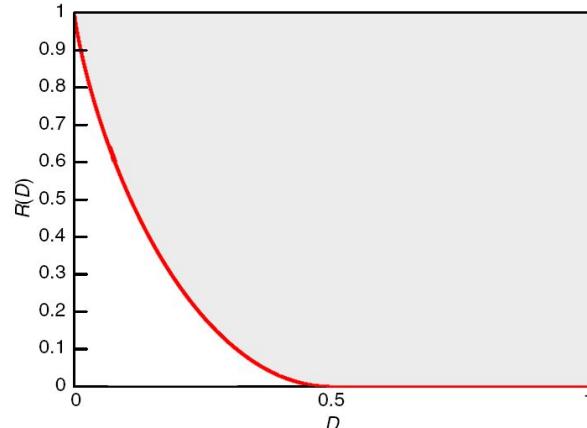


Figure 2: The NeuroFovea metamer generation schematic: An input image and a noise patch are fed through a VGG-Net encoder into a new feature space. Through spatial control we can produce an interpolation for each pooling region in such feature space between the stylized-noise (texture), and the content (the input image). This is how we successfully impose both global image and local texture-like constraints in every pooling region. The metamer is the output of the pooled (and interpolated) feature vector through the Meta VGG-Net Decoder.

Example of a RD curve:



# $\beta$ -VAE: LEARNING BASIC VISUAL CONCEPTS WITH A CONSTRAINED VARIATIONAL FRAMEWORK

Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot,  
Matthew Botvinick, Shakir Mohamed, Alexander Lerchner

Google DeepMind

{irinah, lmatthey, arkap, cpburgess, glorotx,  
botvinick, shakir, lerchner}@google.com

## ABSTRACT

Learning an interpretable factorised representation of the independent data generating factors of the world without supervision is an important precursor for the development of artificial intelligence that is able to learn and reason in the same way that humans do. We introduce  $\beta$ -VAE, a new state-of-the-art framework for automated discovery of interpretable factorised latent representations from raw image data in a completely unsupervised manner. Our approach is a modification of the variational autoencoder (VAE) framework. We introduce an adjustable hyperparameter  $\beta$  that balances latent channel capacity and independence constraints with reconstruction accuracy. We demonstrate that  $\beta$ -VAE with appropriately tuned  $\beta > 1$  qualitatively outperforms VAE ( $\beta = 1$ ), as well as state of the art unsupervised (InfoGAN) and semi-supervised (DC-IGN) approaches to disentangled factor learning on a variety of datasets (*celebA*, *faces* and *chairs*). Furthermore, we devise a protocol to quantitatively compare the degree of disentanglement learnt by different models, and show that our approach also significantly outperforms all baselines quantitatively. Unlike InfoGAN,  $\beta$ -VAE is stable to train, makes few assumptions about the data and relies on tuning a single hyperparameter  $\beta$ , which can be directly optimised through a hyperparameter search using weakly labelled data or through heuristic visual inspection for purely unsupervised data.

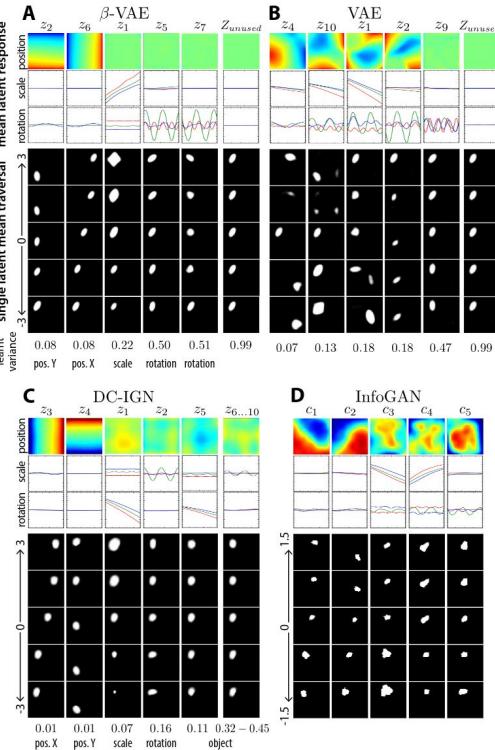
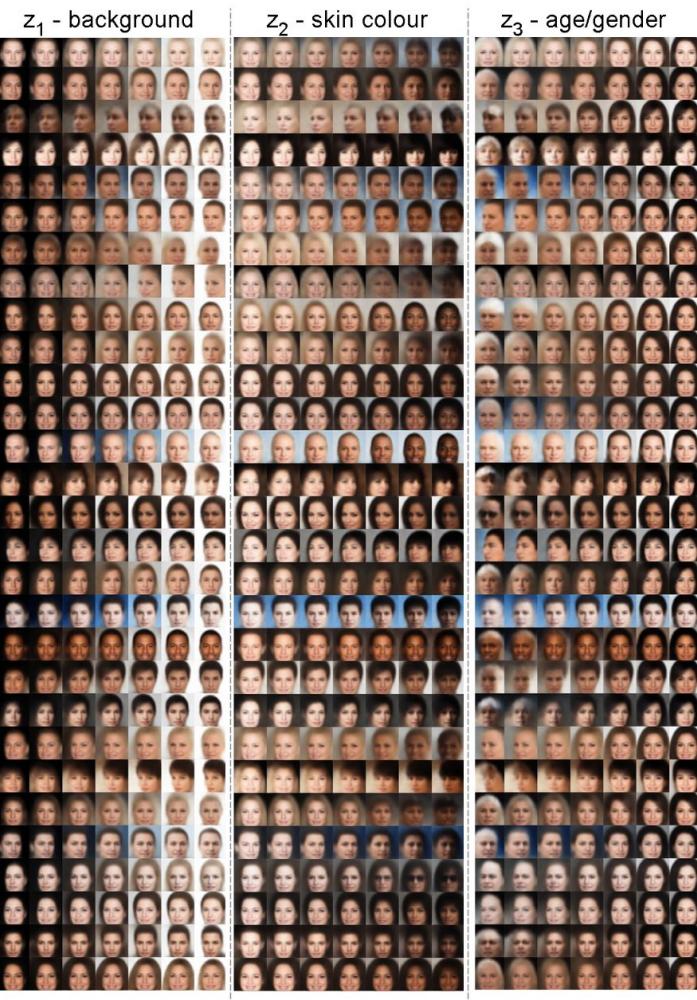


Figure 7: **A:** Representations learnt by a  $\beta$ -VAE ( $\beta = 4$ ). Each column represents a latent  $z_i$ , ordered according to the learnt Gaussian variance (last row). Row 1 (position) shows the mean activation (red represents high values) of each latent  $z_i$  as a function of all 32x32 locations averaged across objects, rotations and scales. Row 2 and 3 show the mean activation of each unit  $z_i$  as a function of scale (respectively rotation), averaged across rotations and positions (respectively scales and positions). *Square* is red, *oval* is green and *heart* is blue. Rows 4-8 (second group) show reconstructions resulting from the traversal of each latent  $z_i$  over three standard deviations around the unit Gaussian prior mean while keeping the remaining 9/10 latent units fixed to the values obtained by running inference on an image from the dataset. **B:** Similar analysis for VAE ( $\beta = 1$ ). **C:** Similar analysis for DC-IGN, clamping a single latent each for scale, positions, orientation and 5 for shape. **D:** Similar analysis for InfoGAN, using 5 continuous latents regularized using the mutual information cost, and 5 additional unconstrained noise latents (not shown).



$$\mathcal{F}(\theta, \phi, \beta; \mathbf{x}, \mathbf{z}) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x} \mid \mathbf{z})] - \beta[D_{KL}(\log q_\theta(\mathbf{z} \mid \mathbf{x}) \parallel p(\mathbf{z})) - \epsilon]$$

Figure 12: Latent traversal plots from  $\beta$ -VAE that learnt disentangled representations on the CelebA dataset.