

# Lecture Notes for **Neural Networks and Machine Learning**



Generative Networks  
and  
Auto-Encoding Generators



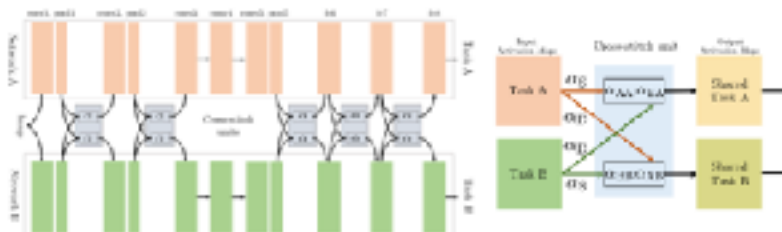
# Logistics and Agenda

- Logistics
  - Lab three due date pushed back (see schedule)
- Agenda
  - A historical perspective of generative Neural Networks
  - Variational Auto-Encoding
  - VAE in Keras Demo (if time)
  - Adversarial Auto-Encoders



# Last Time

## Multi-task: Cross Stitch Networks

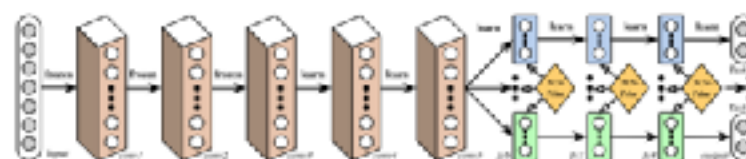


- Only works for simultaneous multi-label problems
  - like semantic segmentation and surface normal segmentation (clustering similarly facing objects)
- Take a learned weighted sum of the activations
- Works a little better than single task, but no worse

<https://arxiv.org/pdf/1804.05501.pdf>

75

## Multi-task: Deep Relationship Networks



- Start training traditionally
- Minimize Kronecker Product between fully connected task specific layers
  - that is, make Covariance between layers close to identity
  - encourages feature maps in each task to be **less correlated** to feature maps of another task

<https://arxiv.org/pdf/1805.05117.pdf>

76



## Multi-Task Learning in Keras with Multi-Label Data

Fashion week, colors and dresses

"finish demo"

Follow Along: <https://www.pyimagesearch.com/2018/06/04/keras-multiple-outputs-and-multiple-losses/>

80



## Multi-Task Learning

School Data, Computer Surveys



Justin Hlop



Luke Wood

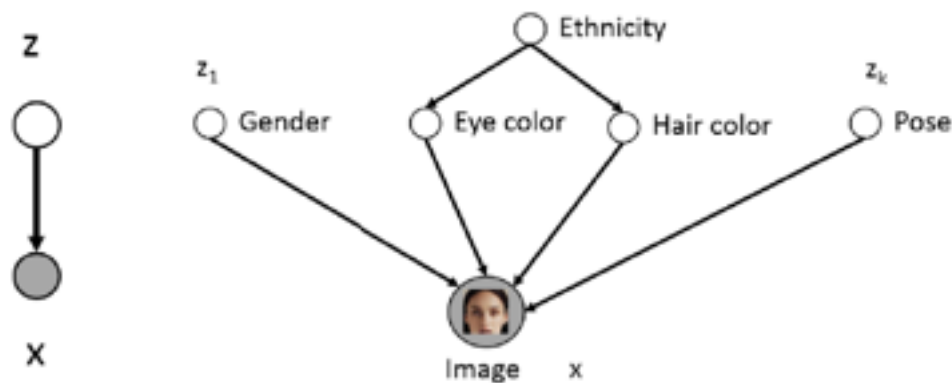
"finish demo"

Follow Along: [LectureNotesMaster/05 LectureMultiTask.ipynb](#)

81

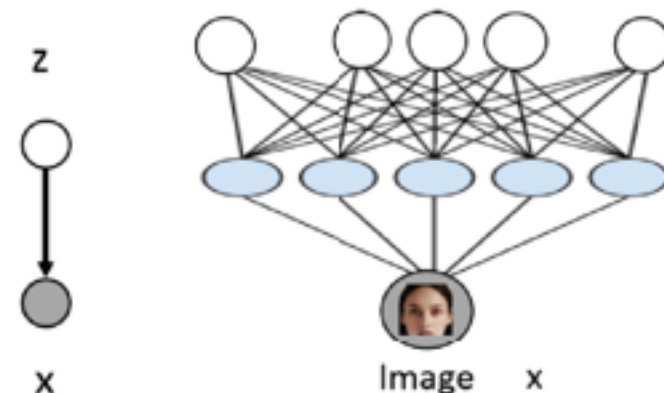


# Motivations: Generative Latent Variables



$$p(\mathbf{x} | \mathbf{z})$$

**Hard:**  $\mathbf{z}$  is expertly chosen



$$p(\mathbf{x} | \mathbf{z})$$

**Not as Hard:**  $\mathbf{z}$  is trained,  
latent variables are uncontrolled

**Want:**

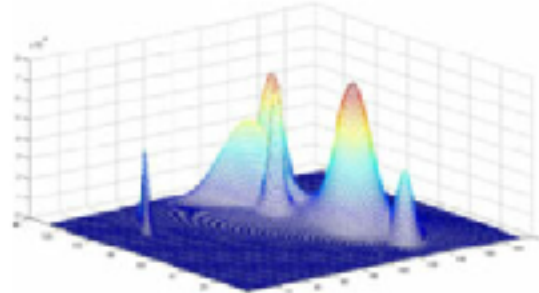
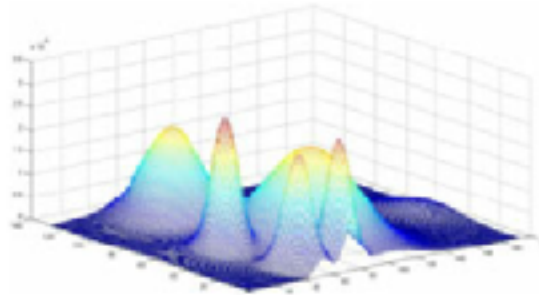
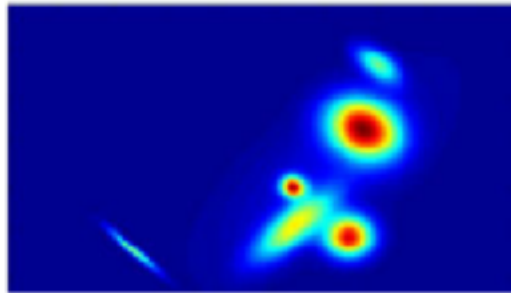
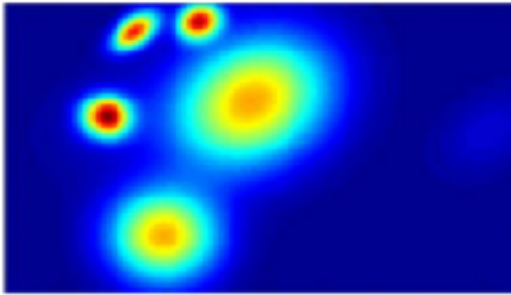
$$p(\mathbf{x}) \approx \sum_{\mathbf{z}} p(\mathbf{z}) p_{\theta}(\mathbf{x} | \mathbf{z})$$



# Motivation: Mixtures for Simplicity

**Want:** 
$$p(\mathbf{x}) \approx \sum_{\mathbf{z}} p(\mathbf{z}) p_{\theta}(\mathbf{x} | \mathbf{z})$$

- Each latent variable is mostly independent of other latent variables
- The sum of various mixtures can approximate most any distribution
- Good choice for conditional is Normal Distribution
- Can parameterize  $p(\mathbf{x} | \mathbf{z})$  to be a Neural Network



$$p_{\theta}(\mathbf{x} | \mathbf{z} = k) = \mathcal{N}(\mathbf{x}, \mu_k, \Sigma_k)$$

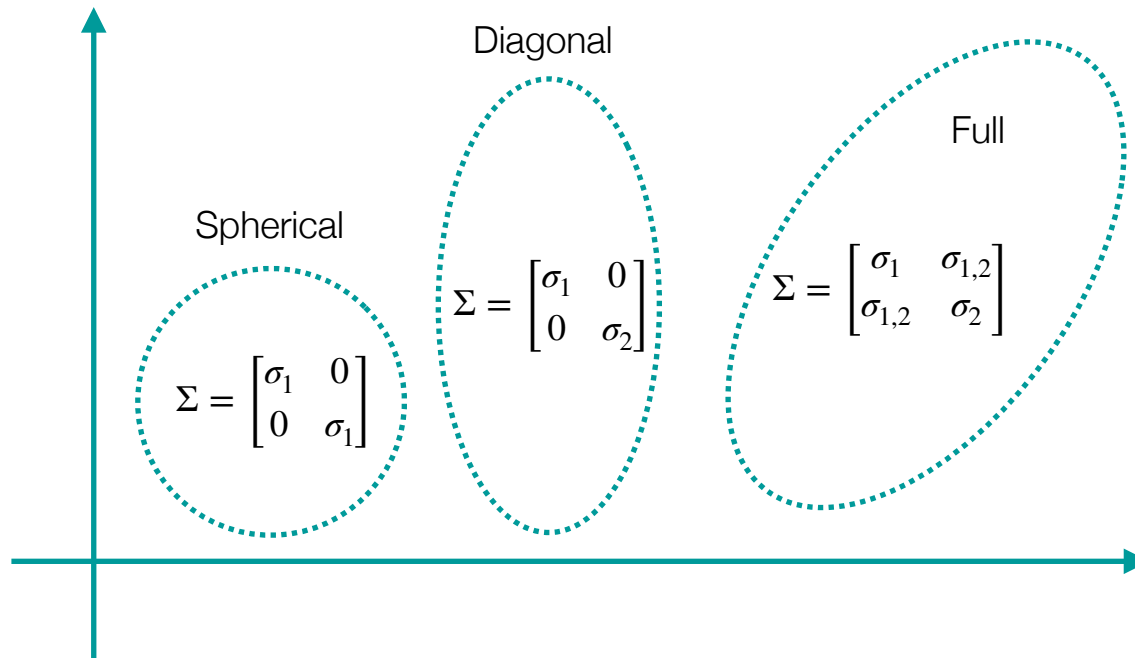
mean and covariance learned



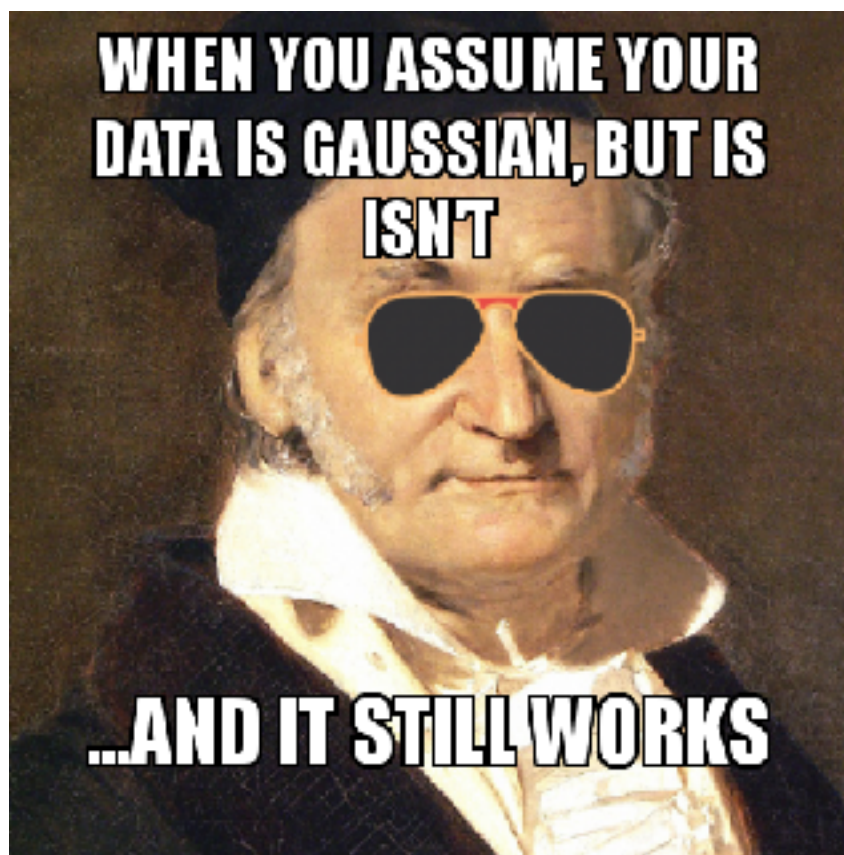
# Motivation: Mixtures for Simplicity

$$= \mathcal{N}(\mathbf{x}, \mu_k, \Sigma_k)$$

mean and covariance learned



# A History of Generative Networks



# Aside: Notation

$$\overset{\text{some function}}{\mathbf{E}_{s \leftarrow q(s|x)}}[f(\cdot)] = \int q(s|x) \cdot f(x) dx \approx \overset{\text{could be neural networks}}{\sum_{\forall i} q(s|x^{(i)}) \cdot f(x^{(i)})}$$

Expected value of  $f$  under conditional distribution,  $q$   
 $s$  is latent variable,  $x^{(i)}$  is an observation

$$\mathbf{E}_{s \leftarrow q(s|x)}[\log f(\cdot)] = \sum_{\forall i} q(s|x^{(i)}) \cdot \log(f(x^{(i)}))$$

If function is a probability, this is just the negative of cross entropy of distributions:

$$H(q, p) = - \sum_x q(x) \cdot \log(p(x))$$

Recall that KL divergence is a measure of difference in two distribution, and is just:

$$D(p||q) = \sum_x p(x) \cdot \log \left( \frac{p(x)}{q(x)} \right) = \mathbf{E}_p \left[ \log \left( \frac{p(x)}{q(x)} \right) \right]$$





# Taxonomy of Generative Models

## Taxonomy of Generative Models

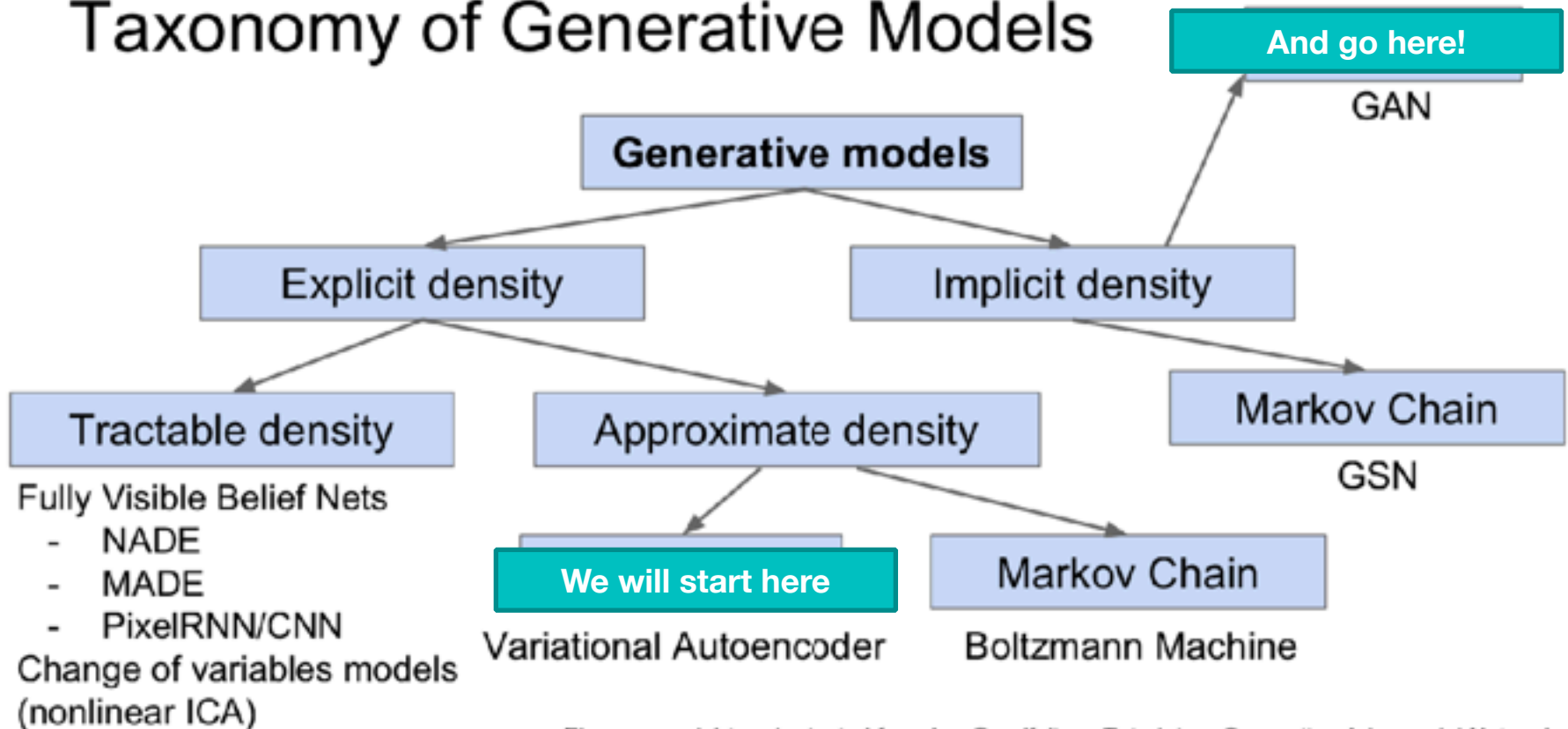


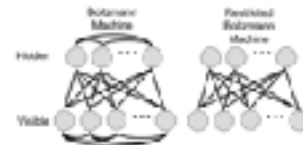
Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.



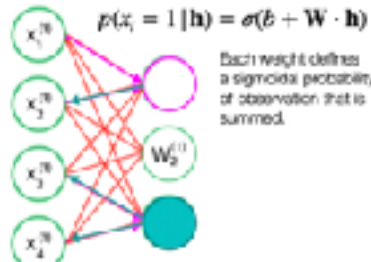
# Abridged History of Generative Networks

## • Restricted Boltzmann Machine

- Forward pass (visible to latent)
- Backward pass (latent to visible)
- Similar to an auto-encoder



AUTOENCODERS



RBM

<https://www.courville.com/blog/restricted-boltzmann-machine-tutorial/>

## 2006 Restricted Boltzmann Machine

## • Deep Boltzmann Machine

$$P(x, h^{(1)}, h^{(2)}) = \frac{1}{Z(\theta)} \exp(-E(x, h^{(1)}, h^{(2)}; \theta)). \quad (26.24)$$

To simplify our presentation, we omit the bias parameters below. The DBM energy function is then defined as follows:

$$E(x, h^{(1)}, h^{(2)}; \theta) = -x^T W^{(1)} h^{(1)} - h^{(1)T} W^{(2)} h^{(2)} - h^{(2)T} W^{(3)} h^{(3)}. \quad (26.25)$$

We now develop the mean field approach for the example with two hidden layers. Let  $Q(h^{(1)}, h^{(2)} | x)$  be the approximation of  $P(h^{(1)}, h^{(2)} | x)$ . The mean field assumption implies that

$$Q(h^{(1)}, h^{(2)} | x) = \prod_i Q(h_i^{(1)} | x) \prod_j Q(h_j^{(2)} | x). \quad (26.26)$$

**Not tractable: Can only optimize the Evidence lower bound, ELBO**

One can consider of many ways of measuring how well  $Q(h | x)$  fits  $P(h | x)$ . The mean field approach is to minimize

$$\text{KL}(Q||P) = \sum_x Q(x) \log \left( \frac{Q(h^{(1)}, h^{(2)} | x)}{P(h^{(1)}, h^{(2)} | x)} \right). \quad (26.27)$$

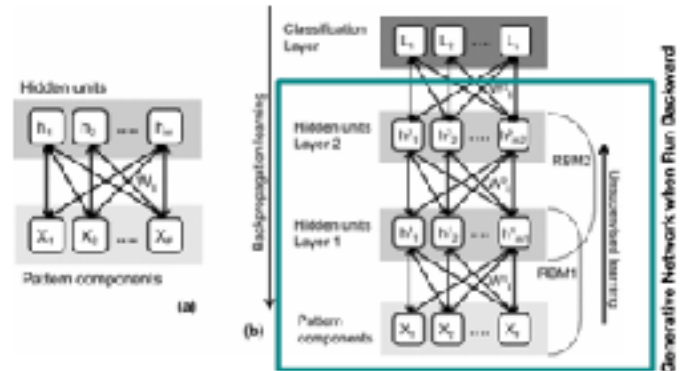
Approximate via MCMC via Gibbs Sampling

Goodfellow, Bengio, Courville, Deep learning, MIT press, 2016

## 2009 Deep Boltzmann Machine Goodfellow, Bengio, Courville

## • Deep Belief Network

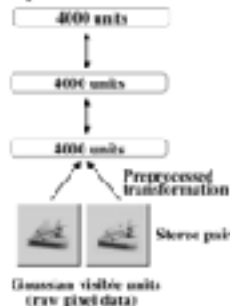
- Many RBM blocks together!



[https://paperswithcode.com/abstract/10/ICCV/978-3-319-63132-7\\_23](https://paperswithcode.com/abstract/10/ICCV/978-3-319-63132-7_23)

## 2007, Deep Belief Networks RBMs with many layers

### Deep Boltzmann Machine



### Training Samples



### Generated Samples



Figure 5: Left: The architecture of deep Boltzmann machine used for NCRB. Right: Random samples from the training set, and samples generated from the deep Boltzmann machines by running the Gibbs sampler for 10,000 steps.

## 2009, Practical Examples Salakhutdinov and Hinton



# Contemporary Modeling

- DBNs and DBMs did not become very popular
  - Mathematics detracts from popular understanding
  - Often methods using sampling are not scalable
  - Cannot directly use Gradients (no Back Prop ) 😓
- Popular method for calculating generative networks with Evidence Lower Bound (ELBO) approximation:
  - Variational Auto Encoding
    - ◆ Guaranteed NOT to find global minimum
    - ◆ But scalable and will converge in finite time



# Variational Auto Encoding

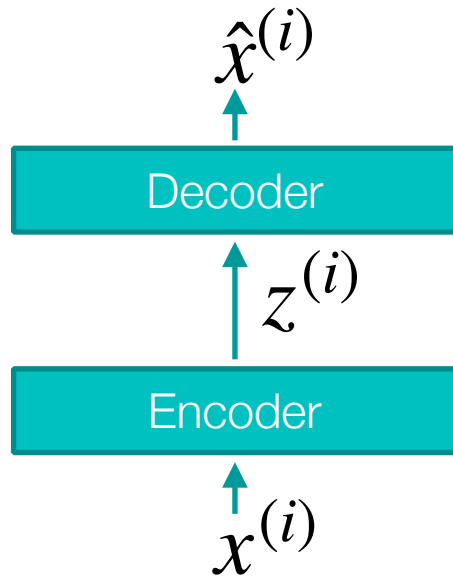
**“Mathematics is the  
Khaleesi of sciences.”**



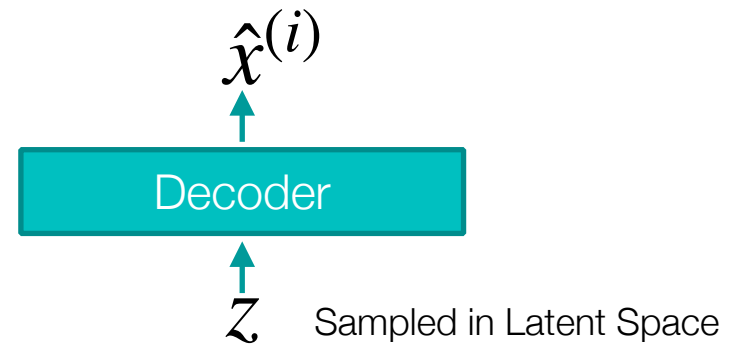
**– Khal Friedrich Gauss**



# Can Auto Encoding Generate Samples?



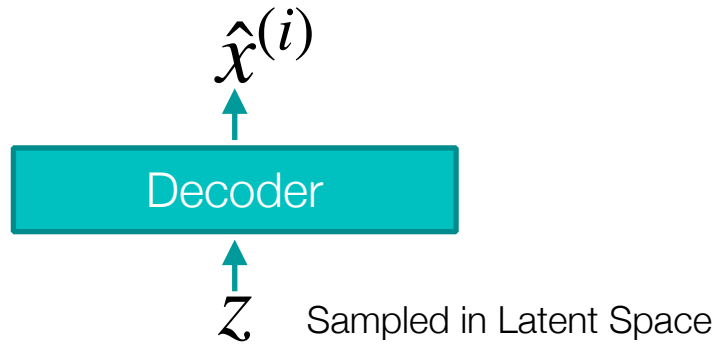
Once trained, is it possible to generate data?



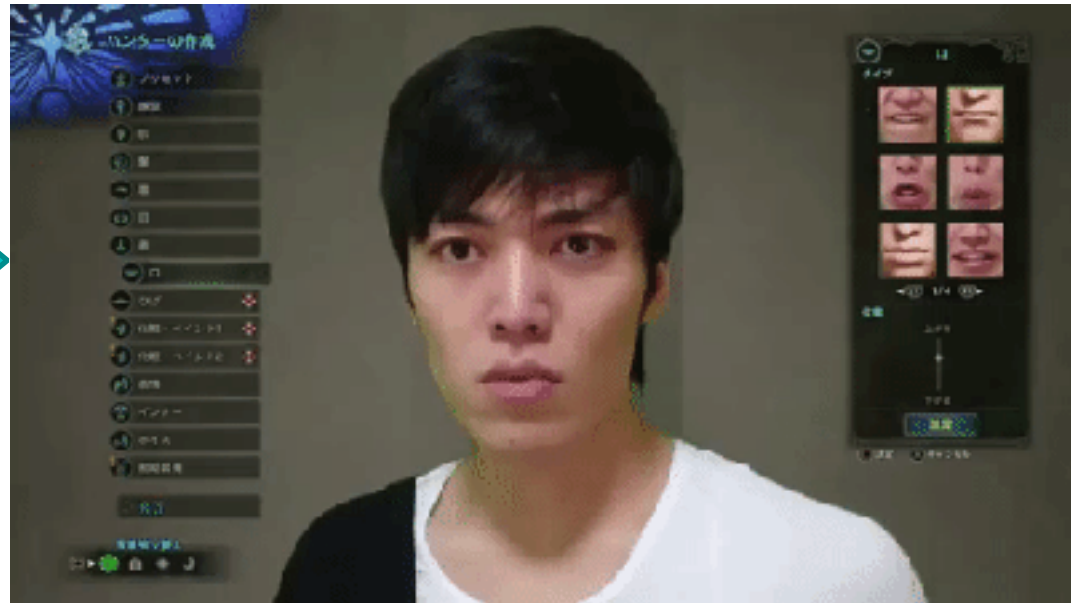
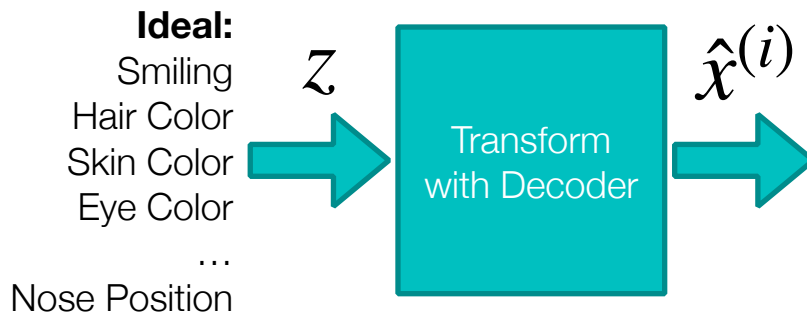
- Does this work for simple auto encoding?
  - Yes, but not satisfactory results
- Learned space is not continuous
- Features could be highly correlated, related in complex ways
  - So, how should we sample from the latent space?
- Need to define some constraints on the latent space...



# Reasonable constraints for $p(z)$ ?



- Should be simple, easy to sample from: **Normal**
- Each component should be i.i.d.:  
**Diag. Covariance**
  - Encourages features that may be semantic, like expert might select



# Optimizing

$$p(z | x) = \frac{p(x | z)p(z)}{p(x)}$$

We need this inference in order to compute latent variable

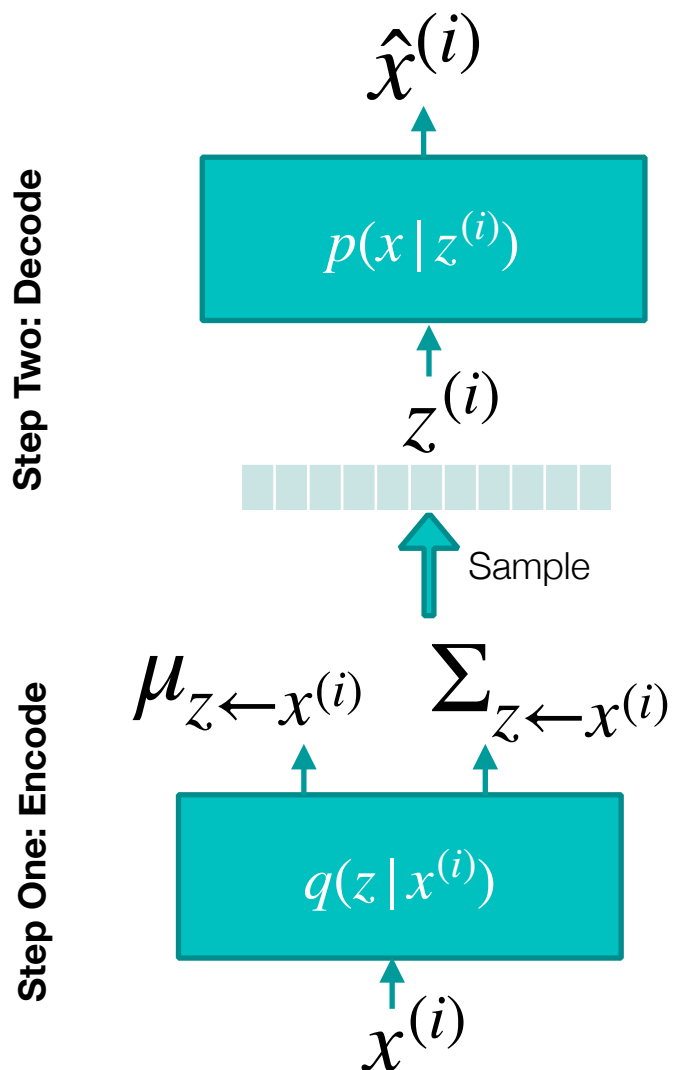
$$p(x) = \int p(x | z)p(z)dz$$

Denominator is of this form

- We can't compute! **Intractable computation** for all “ $z$ ”
- So let's define this with **variational inference**:
  - Only needs to work for  $z$  **with observed**  $x^{(i)}$
  - 1. **Encode** observed  $x^{(i)}$  using network  $q(z | x^{(i)})$ ,
  - 2. Use  $q(z | x^{(i)})$  to sample  $z$  appropriately, then **decode** with another neural network,  $p(x^{(i)} | z^{(i)})$
  - 3. Use bounds to make  $q(z | x^{(i)})$  as close as possible to what we think  $p(z | x)$  might look like



# Need a new formulation



**Step Three: Make conditional  $p$  and  $q$  Similar**

$$D_{KL} [q(z | x^{(i)}) || p(z | x^{(i)})] = \mathbf{E}_{q(z|x)} \left[ \log \left( \frac{q(z | x^{(i)})}{p(z | x^{(i)})} \right) \right]$$

**Step Four: Simplify Optimization**

we can manipulate step three to yield the following approximation:

$$\log p(x^{(i)}) \approx \mathbf{E}_{z \leftarrow q(z|x^{(i)})} [\log p(x^{(i)})]$$

**Intuition MLE:** maximize probability of observed  $x^{(i)}$   
given function  $q$

Output of network,  $q$ , are the mean and covariance for sampling a variable  $z$





# Need a new formulation

$$\log p(x^{(i)}) \approx \mathbf{E}_{z \leftarrow q(z|x)} [\log p(x^{(i)})] \quad \text{Maximize!}$$

$$= \mathbf{E}_q \left[ \log \frac{p(x^{(i)} | z) p(z)}{p(z | x^{(i)})} \frac{q(z | x^{(i)})}{q(z | x^{(i)})} \right] \quad \text{Bayes rule + multiply by one}$$

$$= \mathbf{E}_q [\log p(x^{(i)} | z)] + \mathbf{E}_q \left[ \log \frac{p(z)}{q(z | x^{(i)})} \right] + \mathbf{E}_q \left[ \log \frac{q(z | x^{(i)})}{p(z | x^{(i)})} \right]$$

$$= \mathbf{E}_q [\log p(x^{(i)} | z)] - \mathbf{E}_q \left[ \log \frac{q(z | x^{(i)})}{p(z)} \right] + \mathbf{E}_q \left[ \log \frac{q(z | x^{(i)})}{p(z | x^{(i)})} \right]$$

$$= \mathbf{E}_q [\log p(x^{(i)} | z)] - D_{KL} [q(z | x^{(i)}) || p(z)] + D_{KL} [q(z | x^{(i)}) || p(z | x^{(i)})]$$

always non-negative

$$\log p(x^{(i)}) \geq \mathbf{E}_q [\log p(x^{(i)} | z)] - D_{KL} [q(z | x^{(i)}) || p(z)] \quad \text{Will Maximize Lower Bound}$$

**What have we really done here? Could we have motivated this differently?**



# The Loss Function

Maximize through  
Error of Reconstruction  
Same as minimizing cross entropy

want  $p(z)$  to be  $\mathcal{N}(\mu = 0, \Sigma = I)$   
because it makes nice latent space  
 $q(z|x^{(i)}) \rightarrow (\mu_{z|x}, \Sigma_{z|x}) \quad p(z) \rightarrow \mathcal{N}(0, 1)$

$$\begin{aligned}
 D_{KL}((\mu, \Sigma) \parallel \mathcal{N}(0, 1)) &= \frac{1}{2} \left( \text{tr}(\Sigma) + \mu \cdot \mu^T - \underbrace{k}_{|z|} - \log(\det(\Sigma)) \right) \begin{array}{l} \text{Determinant of diagonal} \\ \text{matrix is simple.} \\ \text{Motivates diagonal} \\ \text{covariance...} \end{array} \\
 \text{Can get this by manipulating} \\
 \text{the KL for normal distribution} \\
 &= \frac{1}{2} \left( \sum_k \Sigma_{k,k} + \sum_k \mu_k^2 - \sum_k 1 - \log \left( \prod_k \Sigma_{k,k} \right) \right) \\
 &\geq \mathbf{E}_{q(z|x^{(i)})} \left[ \log p(x^{(i)} | z) - D_{KL}[q(z|x^{(i)}) \parallel p(z)] \right] \\
 &= \frac{1}{2} \sum_k (\Sigma_{k,k} + \mu_k^2 - 1 - \log \Sigma_{k,k})
 \end{aligned}$$



# The Covariance Output

$$\geq \mathbf{E}_{q(z|x^{(i)})} [\log p(x^{(i)} | z)] - D_{KL} [q(z | x^{(i)}) || p(z)]$$

Maximize through  
Error of Reconstruction  
Same as minimizing cross entropy

want  $p(z)$  to be  $\mathcal{N}(\mu = 0, \Sigma = I)$   
because it makes nice latent space  
 $q(z | x^{(i)}) \rightarrow (\mu_{z|x}, \Sigma_{z|x}) \quad p(z) \rightarrow \mathcal{N}(0, 1)$

$$= \frac{1}{2} \sum_k (\Sigma_{k,k} + \mu_k^2 - 1 - \log \Sigma_{k,k})$$

raw covariance is not numerically stable because of underflow

$$\log \Sigma_{k,k} = \widehat{\Sigma_{k,k}}$$

predicted by  
 $q(z | x^{(i)})$

$$= \frac{1}{2} \sum_k \left( \exp \left( \widehat{\Sigma_{k,k}} \right) + \mu_k^2 - 1 - \widehat{\Sigma_{k,k}} \right)$$

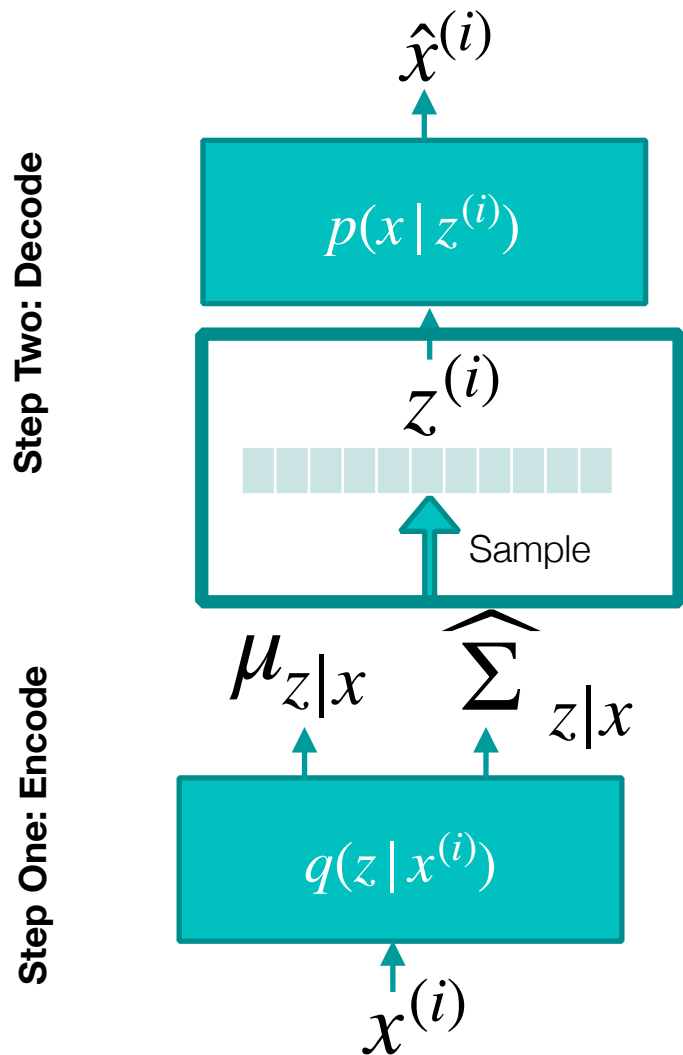
so we will have the neural network output log variance

Also, remember we assume **diagonal covariance**, so  $z$ 's are not correlated

This means covariance is only a vector of variances (the diagonal of  $\Sigma$ )



# Back Propagating



$$\geq \mathbf{E}_{q(z|x^{(i)})} [\log p(x^{(i)} | z)] - D_{KL} [q(z | x^{(i)}) \| p(z)]$$

This is partially differentiable by chain rule...

$$\begin{aligned} \mathcal{N}(\mu_{z|x}, \exp(\widehat{\Sigma_{z|x}})) &= z \\ &= \mu(x^{(i)}) + \exp(\widehat{\Sigma(x^{(i)})}) \cdot \mathcal{N}(0,1) \end{aligned}$$

**To update  $q$ ,  
we need to back propagate  
through sampling layer. How?**



# The Loss Function Implementation

```
# Encode the input into a mean and variance parameter
z_mean, z_log_variance = encoder(input_img)
# Draw a latent point using a small random epsilon
z = z_mean + exp(z_log_variance) * epsilon

# Then decode z back to an image
reconstructed_img = decoder(z)

# Instantiate a model
model = Model(input_img, reconstructed_img)
```

$$z = \mu(x^{(i)}) + \exp(\widehat{\Sigma(x^{(i)})}) \cdot \mathcal{N}(0,1)$$

```
def vae_loss(self, x, z_decoded):
    x = K.flatten(x)
    z_decoded = K.flatten(z_decoded)
    xent_loss = keras.metrics.binary_crossentropy(x, z_decoded)
    kl_loss = -5e-4 * K.mean(
        1 + z_log_var - K.square(z_mean) - K.exp(z_log_var), axis=-1)
    return K.mean(xent_loss + kl_loss)
```

$$-\mathbf{E}_{q(z|x^{(i)})} [\log p(x^{(i)} | z)]$$

$$-\lambda \sum_k 1 + \widehat{\Sigma(x^{(i)})} - \mu(x^{(i)})^2 - \exp(\widehat{\Sigma(x^{(i)})})$$

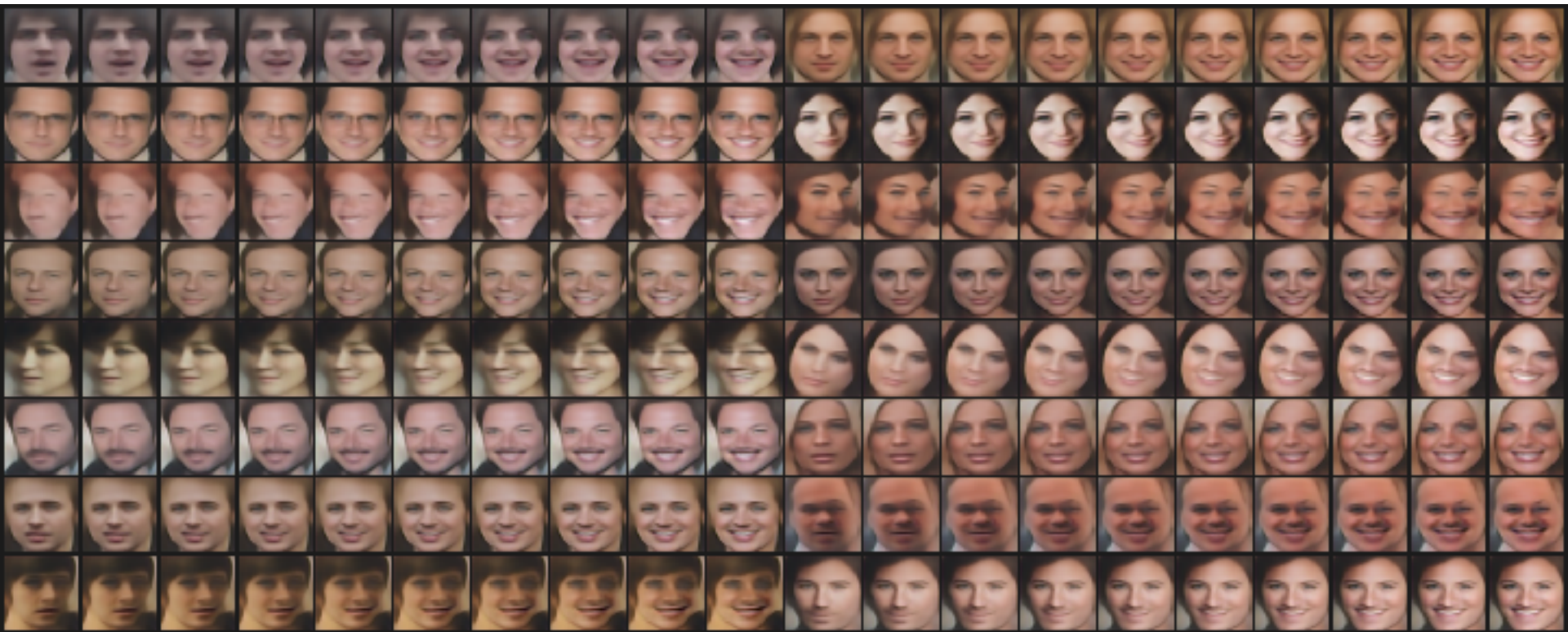
## Note:

Flipped from maximization to minimization  
and added lambda for tradeoff in reconstruction, normal latent space

$$= -\mathbf{E}_{q(z|x^{(i)})} [\log p(x^{(i)} | z)] - \lambda \sum_k 1 + \widehat{\Sigma(x^{(i)})} - \mu(x^{(i)})^2 - \exp(\widehat{\Sigma(x^{(i)})})$$



# VAE Examples



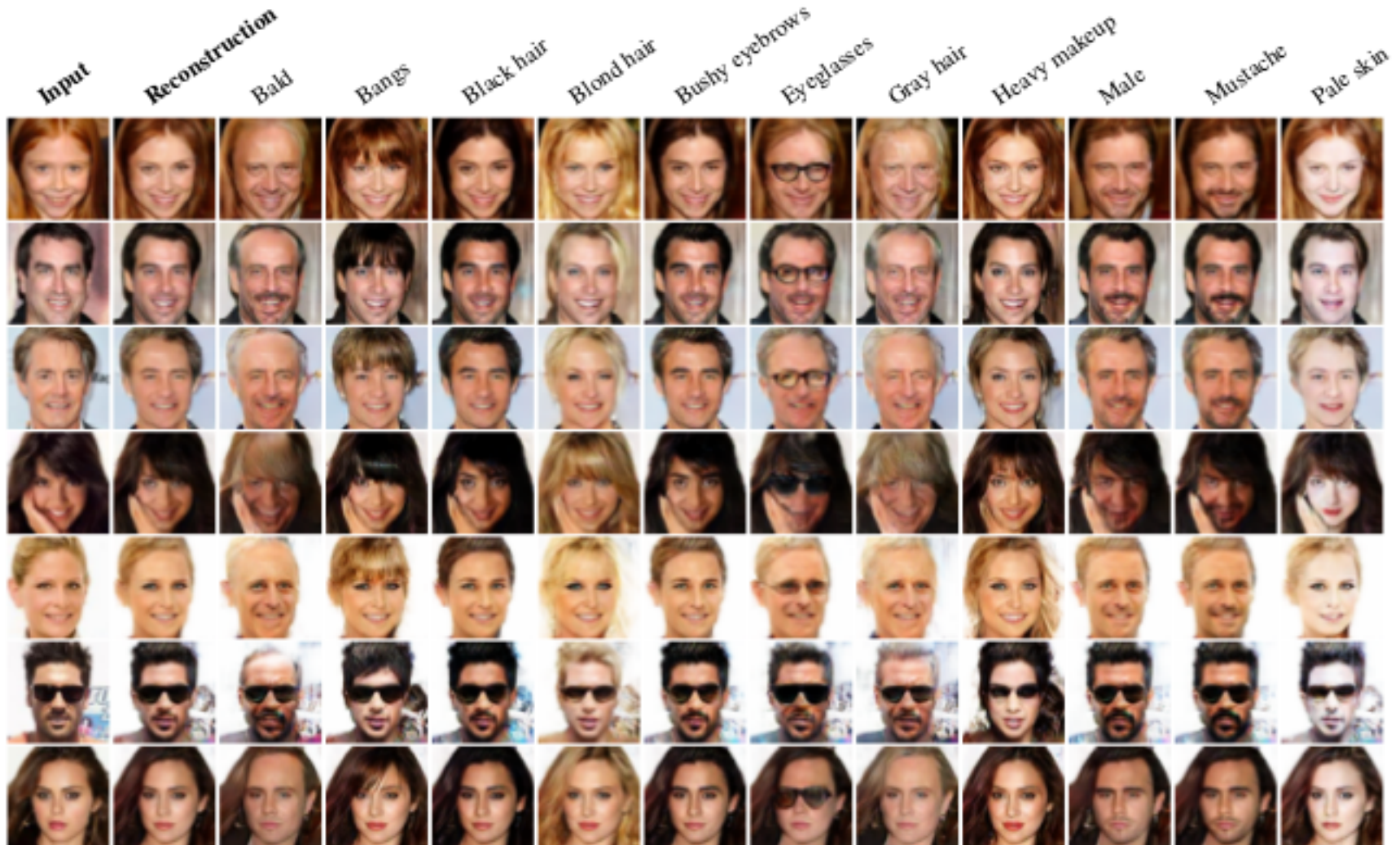
Encoding faces, then finding the “z” that relates to smiling.





# VAE Examples

Different, automatically found  $z$ , latent variables





# VAEs in Keras

Sampling from variational auto encoder  
using MNIST



Demo by Francois Chollet

In Master Repo: [07a VAEs in Keras.ipynb](#)

Follow Along: <https://github.com/fchollet/deep-learning-with-python-notebooks/blob/master/8.4-generating-images-with-vaes.ipynb>

