

# Deep Generative Model

Generative vs Discriminative models, Variational  
Autoencoders, Autoregressive models: RIDE,  
PixelRNN/CNN, PixelCNN++, Conditional Generative  
models

# Success of Discriminative models



→ Cat

Classification



DOG, DOG, CAT

Object Detection



Semantic Segmentation

## ImageNet

- 1 million labeled images

## MS COCO, Pascal VOC

- 0.1 million images with object instances

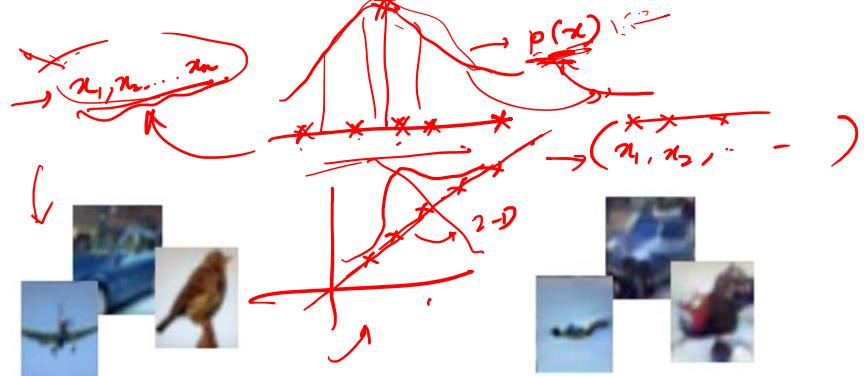
## CityScapes

- 5k HR frames with pixel-wise labels

Need *tons* of data

# Generative models

Goal: To model the data density,  $p(x)$



Training data  $\sim p_{\text{data}}(x)$

Generated samples  $\sim p_{\text{model}}(x)$

So, Why model  $p(x)$ ?

- To improve generalisation
- To create synthetic training data
- Practical tasks like speech synthesis
- To simulate situations
- To understand the data better

Want to learn  $p_{\text{model}}(x)$  similar to  $p_{\text{data}}(x)$

$y = x + n$

$y = Ax + n$       *Substitution*

$\min_{\hat{x}} \|y - Ax\|^2$

$\max_x p(x|y) = \frac{p(x,y)}{p(y)}$

$\max_x \frac{p(x)p(y|x)}{p(y)}$

$A$   $n \times m$   $n < m$

$x$   $m \times 1$

$y$   $n \times 1$

$A^T$   $m \times n$

$A^T y$

$A^T A$

$A^T A^{-1}$

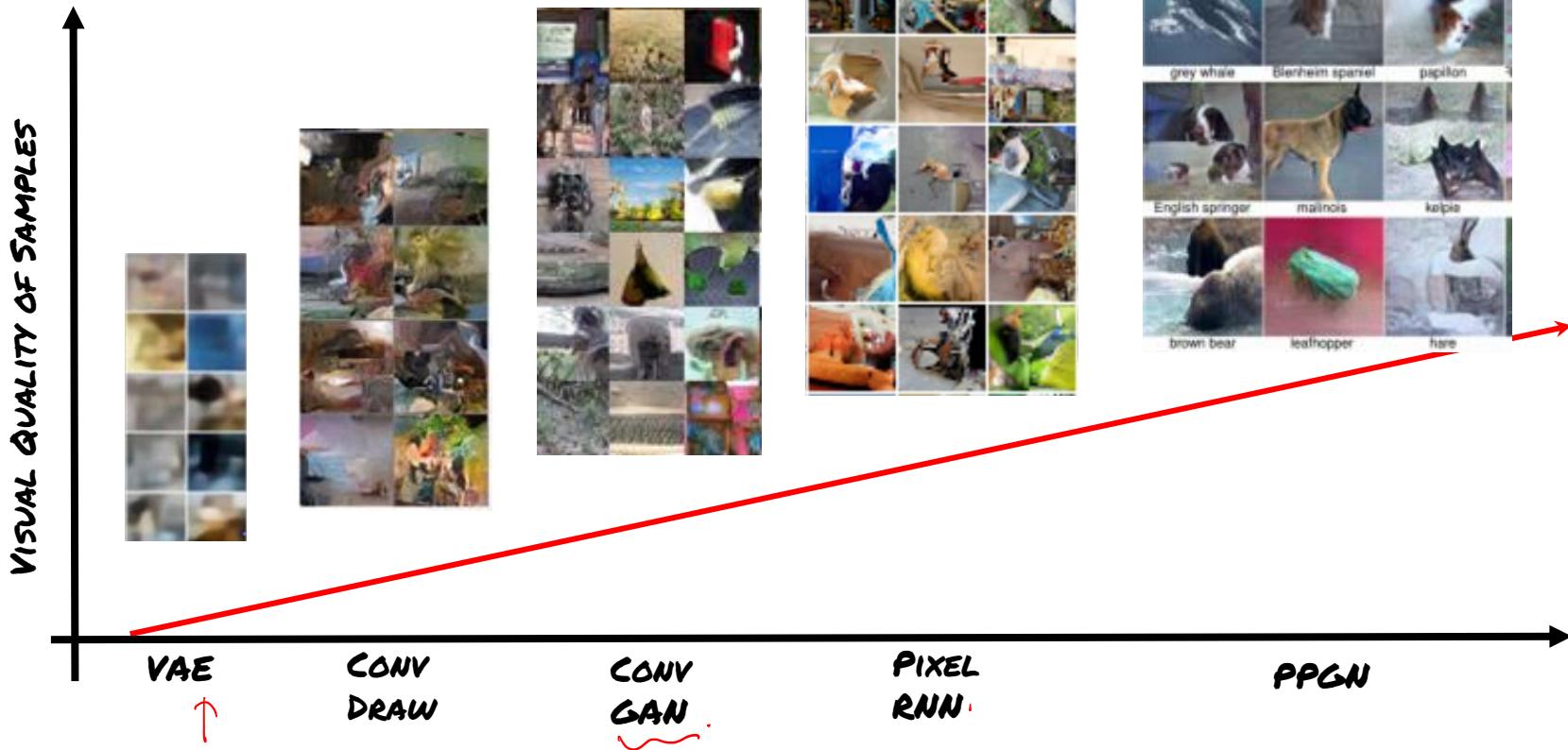
$A^T A$

$A^T A^{-1} A$

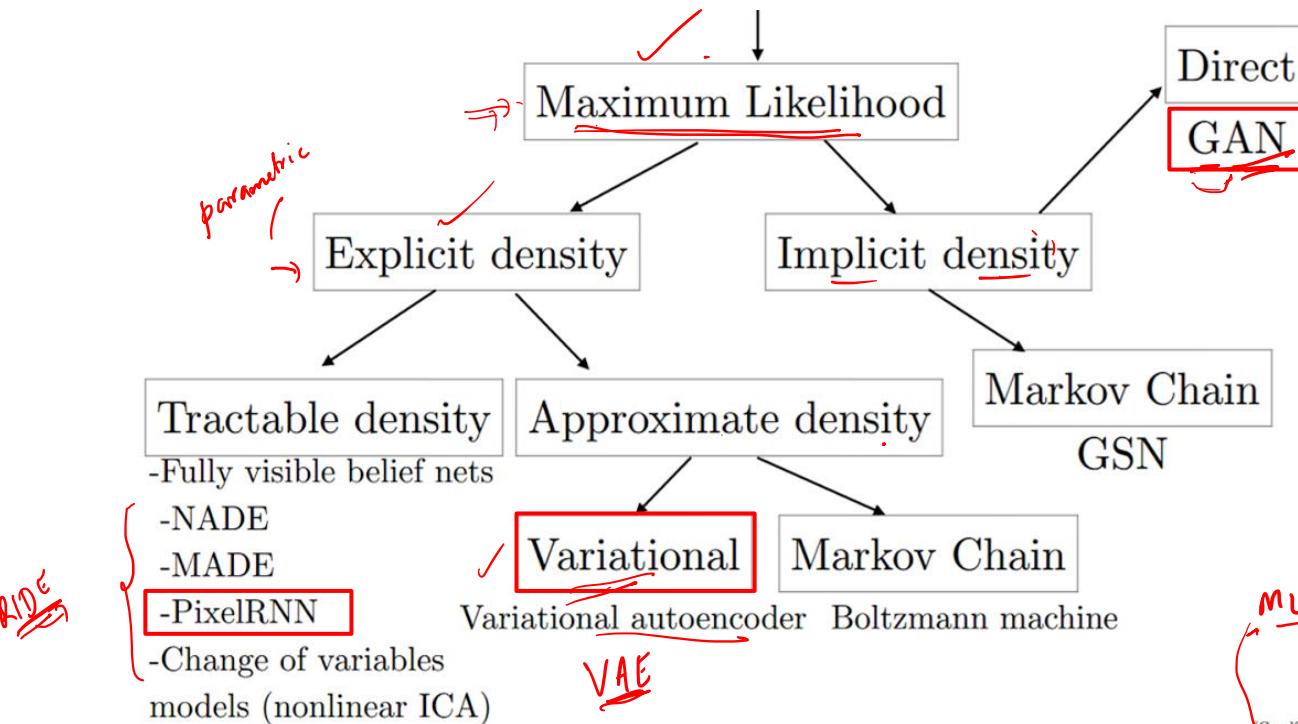
$A^T A^{-1} A^T y$

$x = A^T A^{-1} A^T y$

# Deep Generative models



# What are we going to look at



Handwritten notes on the right side of the slide:

- A bell curve with mean  $\mu$  and variance  $\sigma^2$ . The formula is  $p(x; \mu, \sigma^2)$ .
- A set of data points  $x_1, \dots, x_n$  with mean  $\bar{x}$  and standard deviation  $s$ . The formula is  $\mu = \frac{1}{n} \sum x_i$  and  $\sigma^2 = \frac{1}{n} \sum (x_i - \bar{x})^2$ .
- The formula  $\max_{\mu, \sigma^2} p(x)$ .
- A sequence of data points  $x_1, \dots, x_n$  with a probability distribution  $p(x)$ .
- A bell curve with mean  $\mu$  and variance  $\sigma^2$ . The formula is  $p(x; \mu, \sigma^2)$ .
- The formula  $\max_{\theta} P(x_1, \dots, x_n; \theta)$ .
- The formula  $\max_{\theta} \ln P(\dots; \theta)$ .
- The formula  $NLL \leftarrow -\min_{\theta} \sum_i P(x_i; \theta)$ .

Goodfellow 2016

$$y = Ax + \tilde{w}(x, \epsilon^2) \\ p(y|x) \sim \mathcal{N}(Ax, \sigma^2 I)$$

# Applications: Solving inverse problems

Original Image

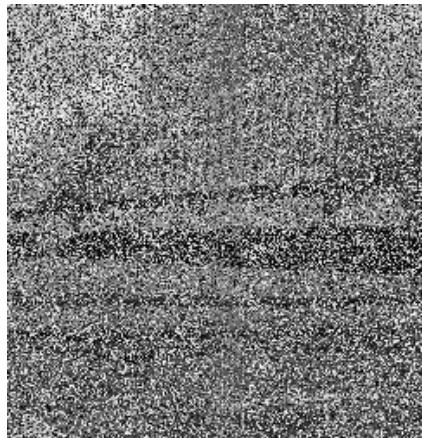


$$\text{MAP} \rightarrow \max_x p(x|y) = \max_x p(x)p(y|x)$$

70% missing pixels



During inference



$$\begin{aligned} & \text{MAP} \rightarrow \max_x p(x|y) = \max_x p(x)p(y|x) \\ & \text{MAP} \rightarrow \min_{\hat{x}} \|y - A\hat{x}\|^2 \\ & \text{rank}(A) \leq \min(m, n) \\ & m < n \quad \text{under-determined} \\ & n > m \quad \text{over-determined} \\ & n \gg m \quad \text{unknown} \\ & \hat{x} = (A^T A)^{-1} A^T y \end{aligned}$$

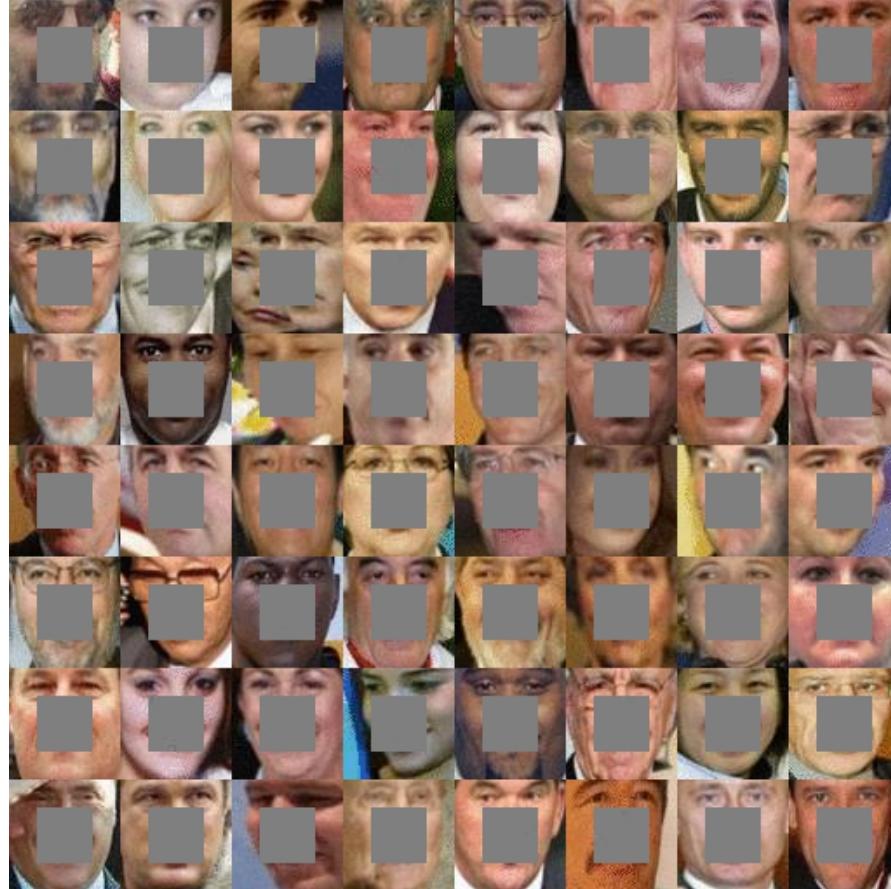
Final inpainted image



$$\begin{aligned} & y = Ax + \tilde{w}(x, \epsilon^2) \\ & p(y|x) \sim \mathcal{N}(Ax, \sigma^2 I) \\ & \text{prior knowledge} \\ & p(x) \sim \mathcal{N}(0, \sigma^2 I) \\ & p(\tilde{w}(x, \epsilon^2)) \sim \mathcal{N}(0, \sigma^2 I) \\ & \text{MAP} \rightarrow \min_{\hat{x}} \|y - A\hat{x}\|^2 + \lambda \|x\|^2 \\ & \hat{x} = (A^T A + \lambda I)^{-1} A^T y \\ & \text{rank}(A) \leq \min(m, n) \\ & m < n \quad \text{under-determined} \\ & n < m \quad \text{over-determined} \\ & y = Ax \Rightarrow \\ & A\hat{x} = y \\ & A\hat{x}_0 = y \\ & \hat{x}_0 = \hat{x} + \lambda x_0 \end{aligned}$$

# Applications

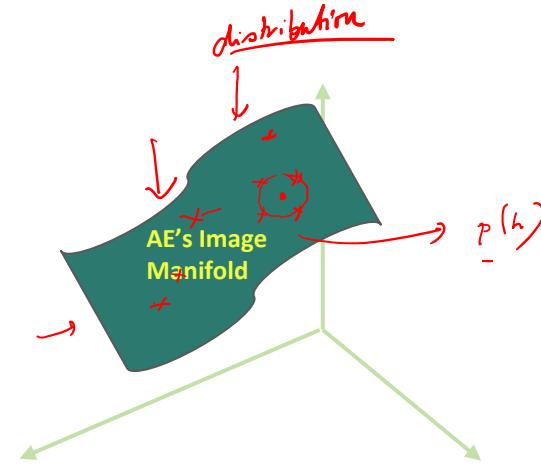
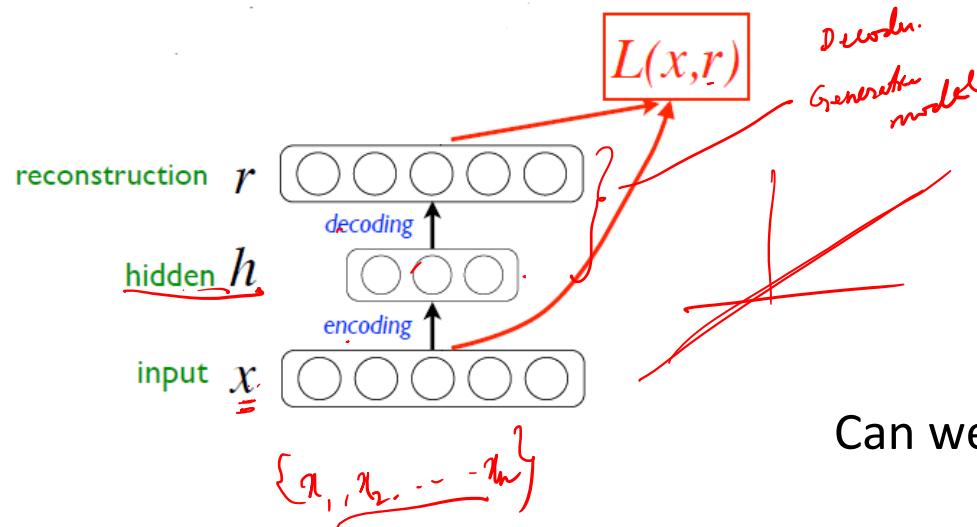
$x^{(n)}$



# Variational Auto Encoder (VAE)

p(?)

# Autoencoders



Can we generate samples from this manifold?

# Training a Decoder to sample

$$\{x_1, x_2, \dots, x_n\}$$

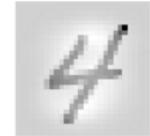
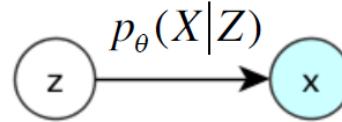
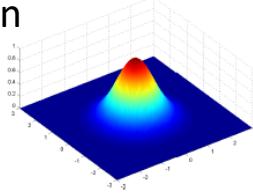
$$\max_{\theta} \prod_i p(x_i; \theta)$$

$$p(x) =$$

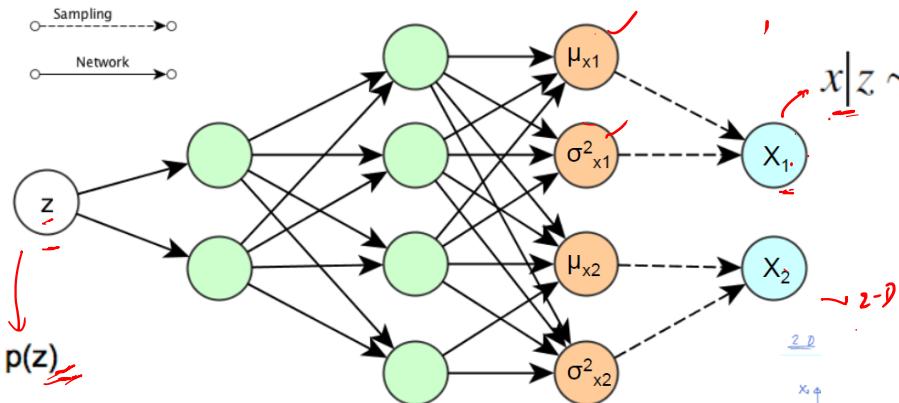
- Specify a distribution on latent variable  $P(z)$

$\Rightarrow P(z)$  could be multivariate Gaussian

$$\begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_n \end{pmatrix} \sim$$



- Learn a decoder distribution  $P(x|z)$



$$p(x)$$

One Example

$$\max_{\theta} \prod_i p(x_i; \theta)$$

$$\max_{\theta} \prod_i p(x_i; \theta) = \prod_i p(x_i; \theta)$$

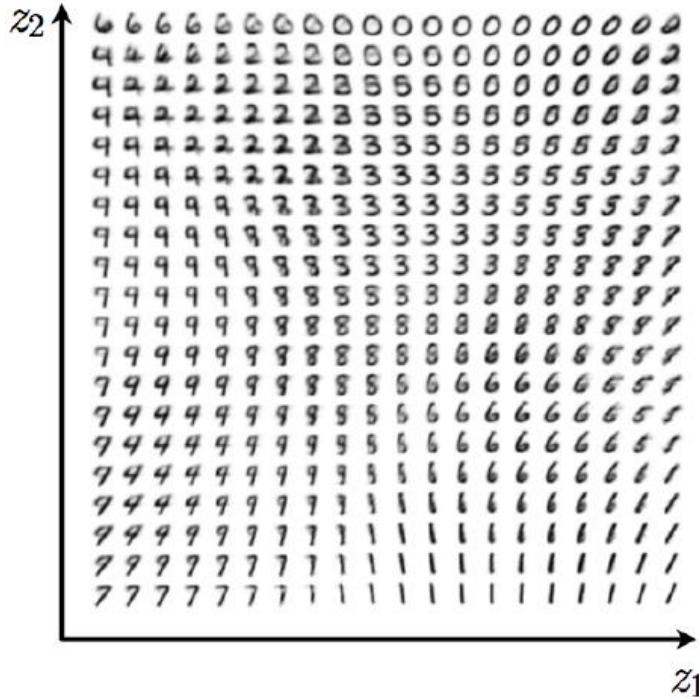
$$\text{Given } p(z), p(x|z)$$

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

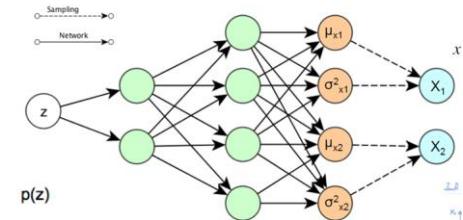
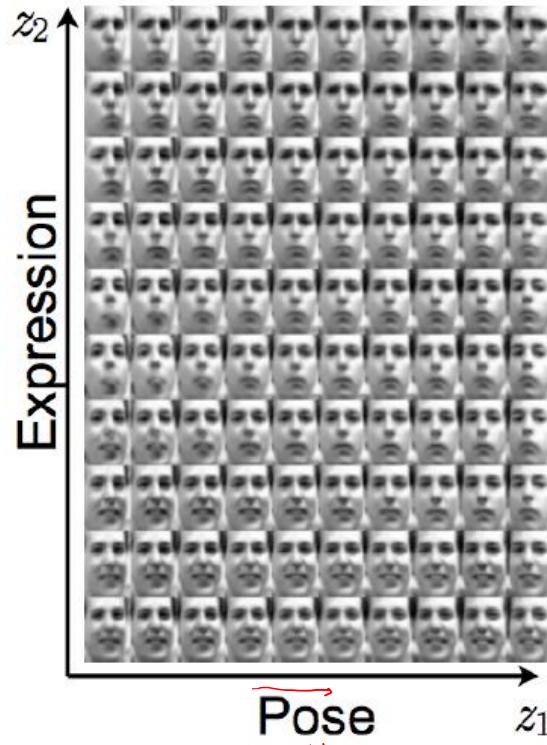
$$p(z) = \int p(z) p(x|z) dz$$

# What $z$ , latent code, can learn?

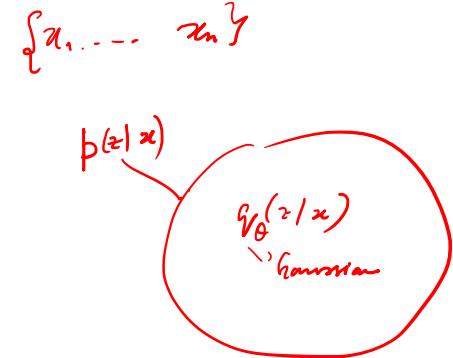
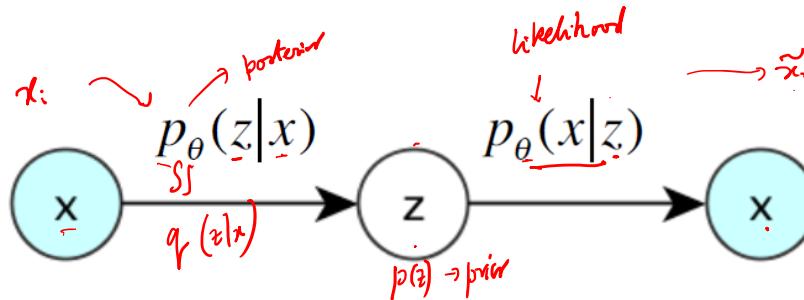
MNIST:



Frey Face dataset:



# Where does $z$ , latent code, come from?



- Marginal,  $p(x)$ , is intractable for complex distributions
- Solution:** approximate  $p(z|x)$  by  $\underline{q(z|x)}$  using variational inference

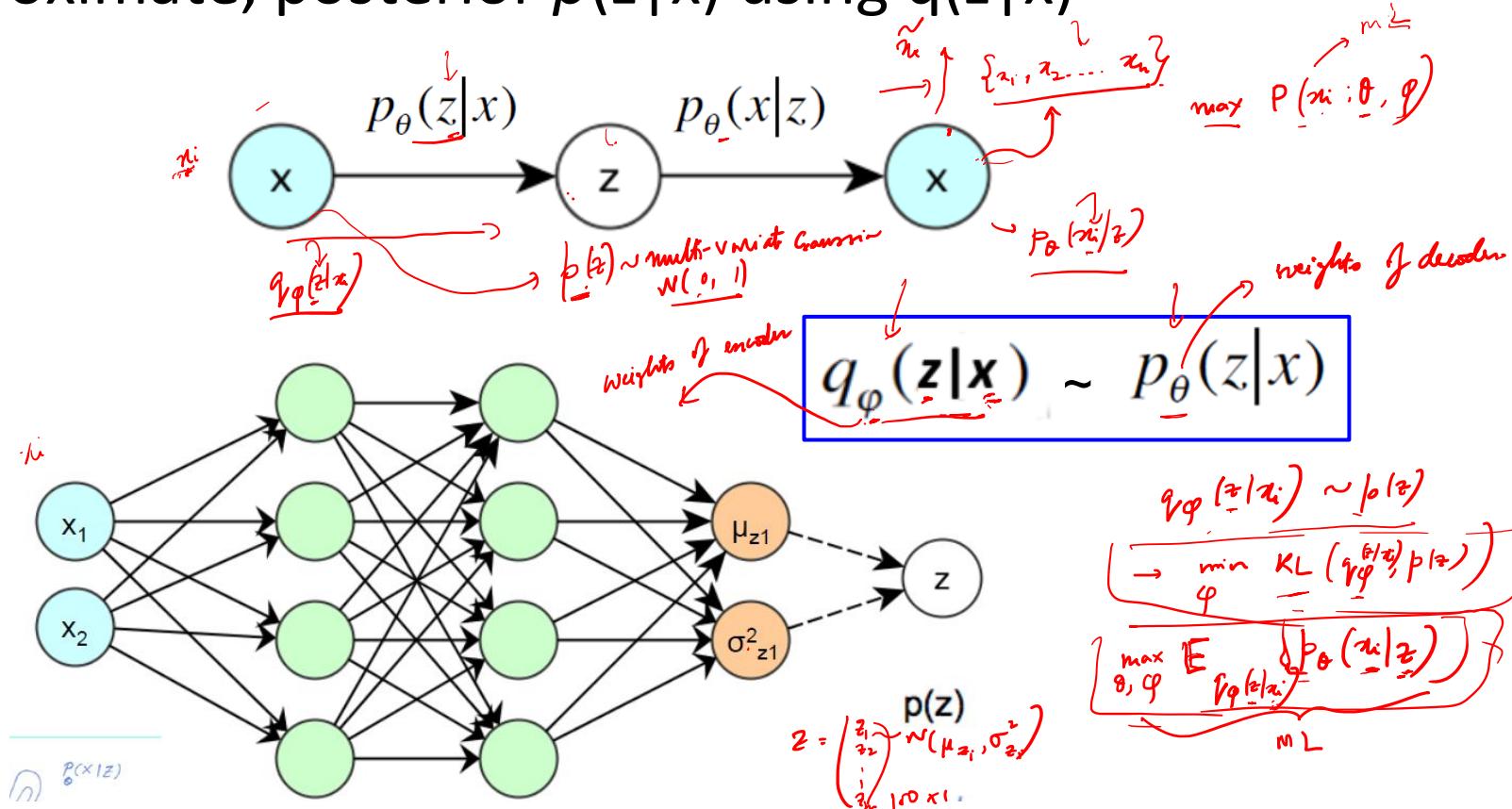
$$p(z|x) = \frac{p(x|z)p(z)}{\int p(x,z)dz}$$

likelihood  
 prior  
 posterior  
 marginal

$\downarrow p(x)$

# Approximate, posterior $p(z|x)$ using $q(\underline{z}|x)$

$$\max_{\theta} p_{\theta}(z_i)$$



# Variational Autoencoder (VAE)

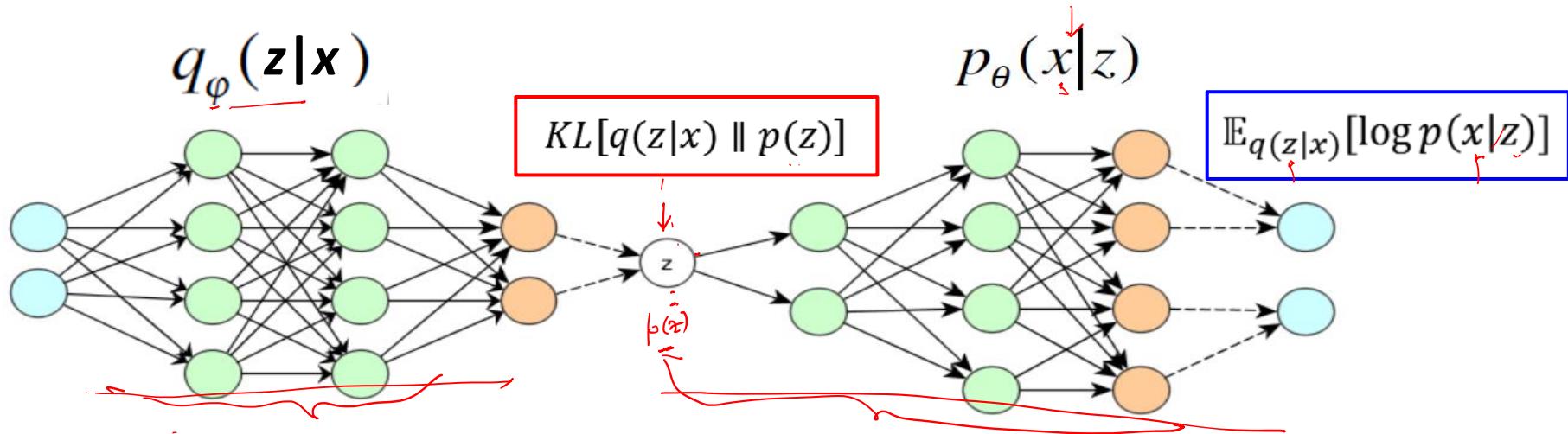
- ❖ Find encoder and decoder parameters such that
  - If we sample from  $\tilde{q}(z|x)$ , then the decoder should reconstruct back  $x$ 
    - Maximize the expected log likelihood value  $\mathbb{E}_{\tilde{q}(z|x)}[\log p(x|z)]$
  - $\tilde{q}(z|x)$  should be close to  $p(z)$ 
    - Minimize KL divergence  $KL[q(z|x) \parallel p(z)]$

- ❖ Thus overall, maximize

$$\mathbb{E}_{\tilde{q}(z|x)}[\log p(x|z)] - KL[q(z|x) \parallel p(z)]$$

# Variational Autoencoder (VAE)

$$\begin{aligned}x &= A(z) \quad w(0, I) \\ \text{Gaussian } E[x^T] &= E[A z^T A^T] \\ w(0, A A^T) &= A A^T\end{aligned}$$



$$\max_{\theta, \phi} \mathbb{E}_{q(z|x)}[\log p(x|z)] - KL[q(z|x) \parallel p(z)]$$

# Autoregressive models

Autoregressive models: RIDE, PixelRNN/CNN,  
PixelCNN++, Conditional Generative models

# Factorization of a distribution

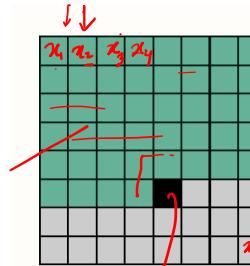
$$p(x_1, x_2, x_3) = p(x_1) p(x_2 | x_1) p(x_3 | x_2, x_1)$$

Is this the only valid factorization?

Let's look at one such factorization over an image,  $x$

$$p(x) = \prod_{i=1}^{n^2} p(x_i | x_1, \dots, x_{i-1})$$

$$\begin{aligned} p(x_1, x_2, \dots, x_n) &= p(x_1) p(x_2 | x_1) p(x_3 | x_1, x_2) \\ &= p(x_m) p(x_{m+1} | x_m) \dots p(x_n | x_1, \dots, x_{n-1}) \end{aligned}$$



■  $x_{ij}$   
■  $X_{<ij}$

$$p(x_i | x_{<i})$$

- If you look at the individual factor  $p(x_i/x_{<i})$  - conditional distribution
- Modeling  $p(x_i/x_{<i})$  is equivalent to sequence prediction,
- Current pixel  $x_i$ 's distribution is modeled based on  $x_{<i}$  causal context

# Modeling the image sequence ...

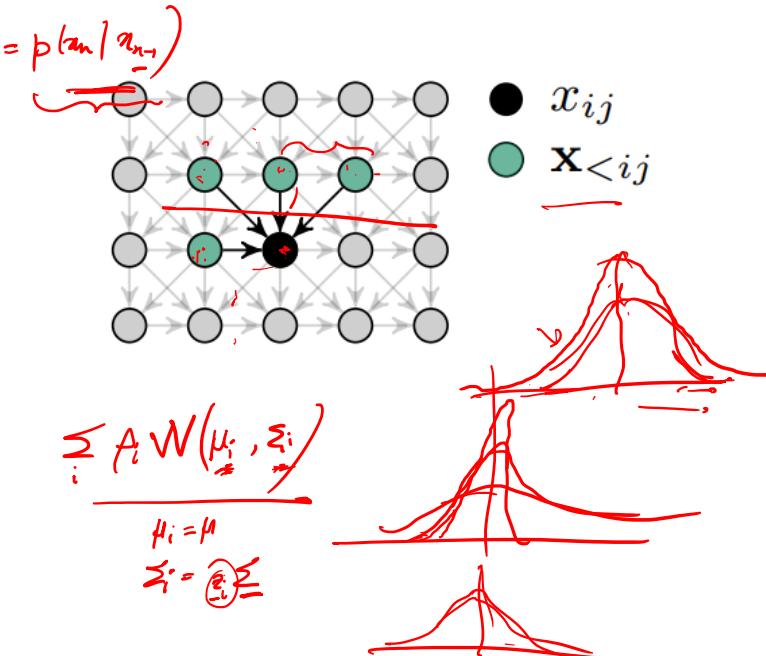
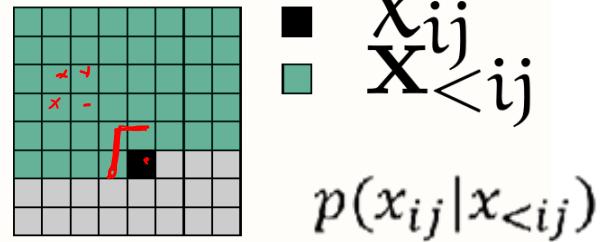
As you can see, depending on the position ' $ij$ ' the length of causal context,  $x_{<ij}$ , changes.

How to handle this?

- One way is to make Markovian assumption
- Model the distribution with Gaussian Scale Mixture (GSM) → mixture of Gaussian

$$p(x) = \int \varphi(z) \mathcal{N}(x; \mu, z\sigma^2) dz,$$

$$= \int \varphi(z) \mathcal{N}(x; \mu, z\sigma^2) dz$$



$$\sum_i A_i \mathcal{N}(\mu_i, \sigma_i^2)$$

$$\mu_i = \mu$$

$$\sigma_i^2 = \sigma^2$$

# Modeling the image sequence ...

- Markovian assumption
- Model the distribution with Mixture of Conditional GSM (MCGSM) by Theis et al.

2012

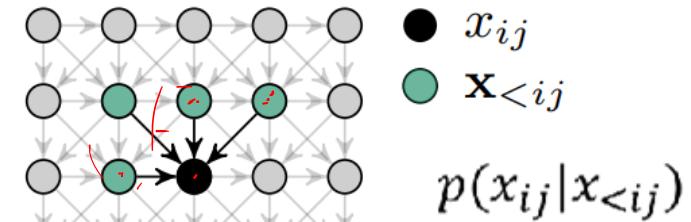
*Gaussian scale mixture*

$$p(x_{ij} | \mathbf{x}_{<ij}, \theta_{ij}) = \sum_{c,s} p(c, s | \mathbf{x}_{<ij}, \theta_{ij}) p(x_{i,j} | \mathbf{x}_{<ij}, c, s, \theta_{ij})$$

↓ gate      ↓ scales      ↓ expert      ↓ scale  
 c              s              c              s  
 covariances      Gaussian      GSM

$\mathcal{N}(x_j | w_k^T x_{<j}, \sigma^2)$

where the sum is over mixture component indices c corresponding to different covariances and scales s corresponding to different variances.



**Note:** length of causal context,  $x_{<ij}$ , is fixed

# How to Train?

$$\max_{\theta} p(\mathbf{x}|\theta) \\ = \min NLL p(\mathbf{x}; \theta)$$

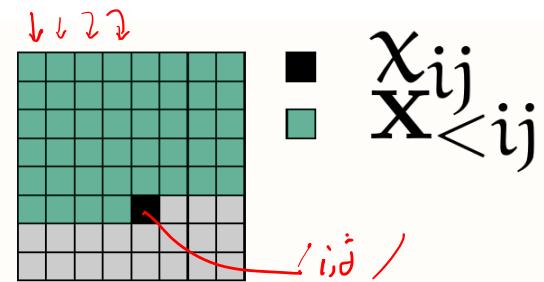
Factorization,  $p(\mathbf{x}) = \prod_{ij} p(x_{ij} | x_{<ij})$

$$NLL: -\log p(\mathbf{x}) = - \sum_{ij} \log p(x_{ij} | x_{<ij}, \theta_{ij})$$

$$= - \sum_{ij} \log p(x_{ij} | x_{<ij}, \theta)$$

$\min_{\theta}$  NLL  
↑  
Training: Maximum likelihood  
parametrized

$$\operatorname{argmin}_{\theta} - \sum_n \log p(\mathbf{x}_n, \theta)$$

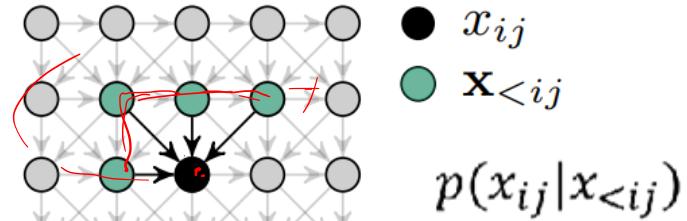


Parametrize,  $p(x_{ij} | x_{<ij})$ ,  
conditional distribution  
with GMMs, GSMs ...

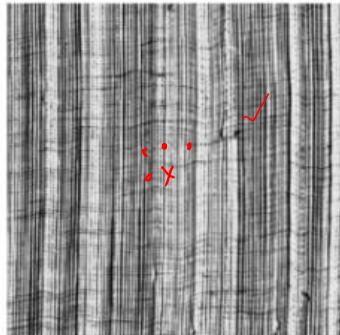
# Limitation

- MCGSM by Theis et al. 2012

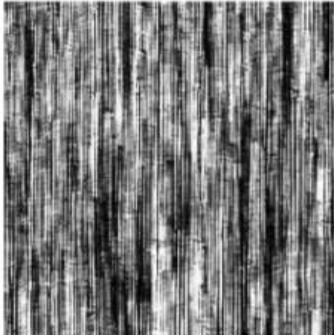
$$p(x_{ij} | \mathbf{x}_{<ij}, \theta_{ij}) = \sum_{c,s} \underbrace{p(c, s | \mathbf{x}_{<ij}, \theta_{ij})}_{\text{gate}} \underbrace{p(x_{i,j} | \mathbf{x}_{<ij}, c, s, \theta_{ij})}_{\text{expert}}$$



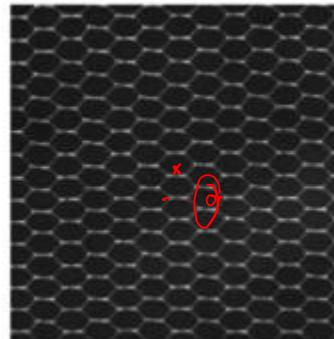
Fails to capture the long range dependencies in an image



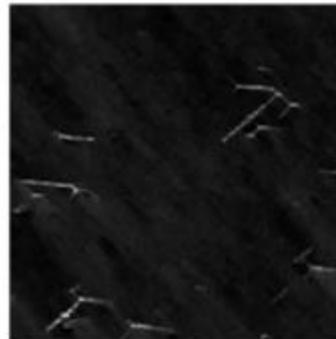
Actual image



Sample



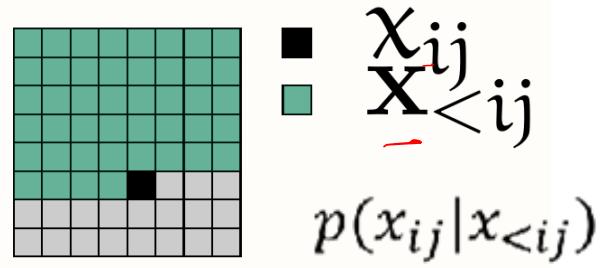
Actual image



Sample

# How to overcome this limitation?

As you can see, depending on the position ' $ij$ ' the length of causal context,  $x_{<ij}$ , changes.



How to handle this?

- Markovian assumption - fails to capture long range dependencies
- Can we use RNNs to model them via  $\underline{x_{<ij}}$  ?

How do we use RNNs to model image sequence?

# RIDE

- Recurrent Image Density Estimator by Theis et al. 2015
- SLSTM: a 2D Spatial LSTM summarizes the causal context,  $x_{<ij}$  as  $h_{ij}$

→ 1-D LSTM

$$g_t = \tanh(W_h h_{t-1} + U_x x_t)$$

$$c_t = g_t \odot i_t + c_{t-1} \odot f_t$$

$$h_t = \tanh(O_t \odot c_t)$$

$h_{ij} = f(x_{<ij}, h_{i-1,j}, h_{i,j-1})$

$\begin{pmatrix} g_{ij} \\ o_{ij} \\ i_{ij} \\ f_{ij}^r \\ f_{ij}^c \end{pmatrix} = \begin{pmatrix} \tanh \\ \sigma \\ \sigma \\ \sigma \\ \sigma \end{pmatrix} T_{A,b} \begin{pmatrix} x_{<ij} \\ h_{i-1,j-1} \\ h_{i-1,j} \end{pmatrix}$

$c_{ij} = g_{ij} \odot i_{ij} + c_{i,j-1} \odot f_{ij}^c + c_{i-1,j} \odot f_{ij}^r$ ,       $h_{ij} = \tanh(c_{ij} \odot o_{ij})$ ,

$p(x_{ij} | h_{<ij}, \theta) = \sum_{c,s} p(c, s | h_{<ij}, \theta) p(x_{ij} | c, s, h_{<ij}, \theta)$

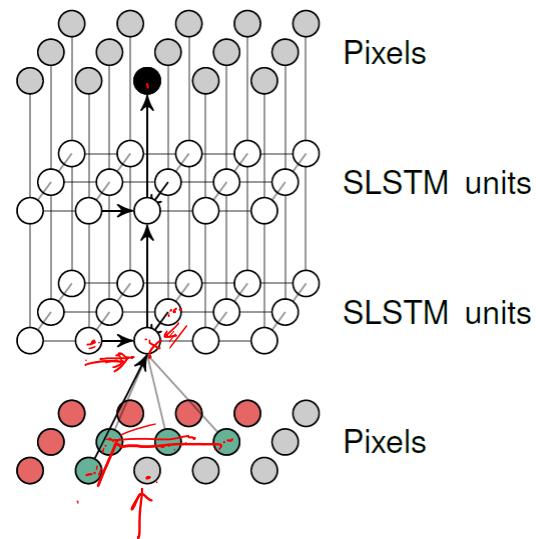
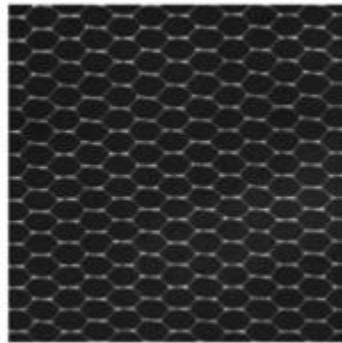
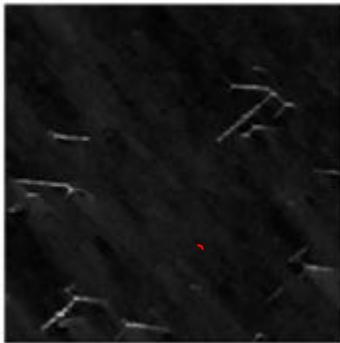


Image courtesy: Theis et al., 2015

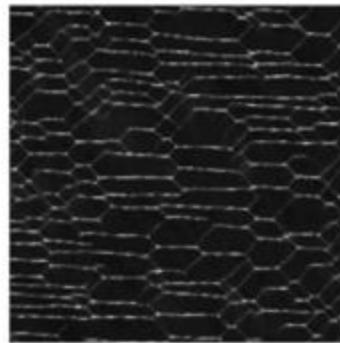
# RIDE



Actual image

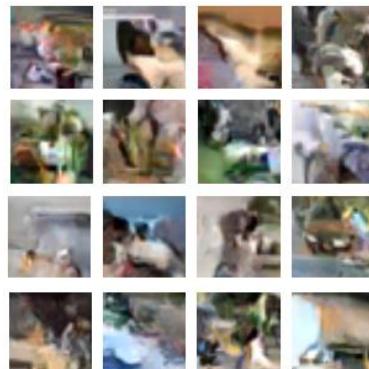


MCGSM sample



RIDE sample

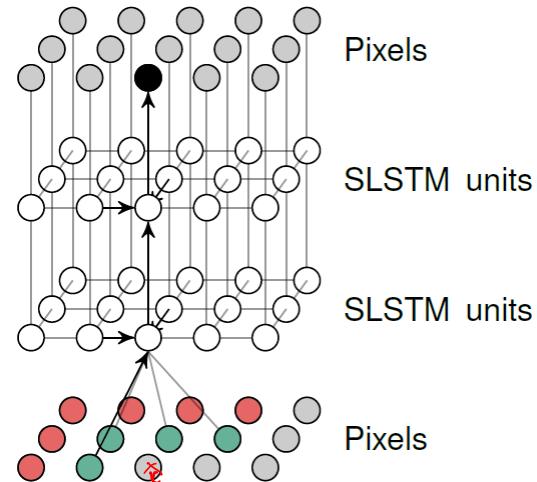
SLSTM improves in  
modeling *long-term*  
dependencies



RIDE samples trained on  
CIFAR dataset

# RIDE summary

- SLSTM captures long-range dependencies in images well
- Sampling and training is **time consuming** because of sequential nature



# Using Generative Models for Image Processing

Image inpainting, Compressive imaging, LiSens and  
FlatCam

$$\max_x p(x|y) = \max_x \frac{p(x)p(y|x)}{p(y)} = \max_x \frac{p(x)p(y|x)}{\int p(y)dx} = \max_x \frac{p(x)p(y|x)}{\int e^{-\frac{\|x-y\|^2}{2\sigma^2}}dx}$$

$$= \max_x e^{-\frac{\|x-y\|^2}{2\sigma^2}}$$

$$= \min_x \|x-y\|^2$$

$$\min_x \|x-y\|^2 + \lambda \|x\|^2$$

$$\text{L}_2\text{-regularization}$$

# Discriminative models - learns a mapping

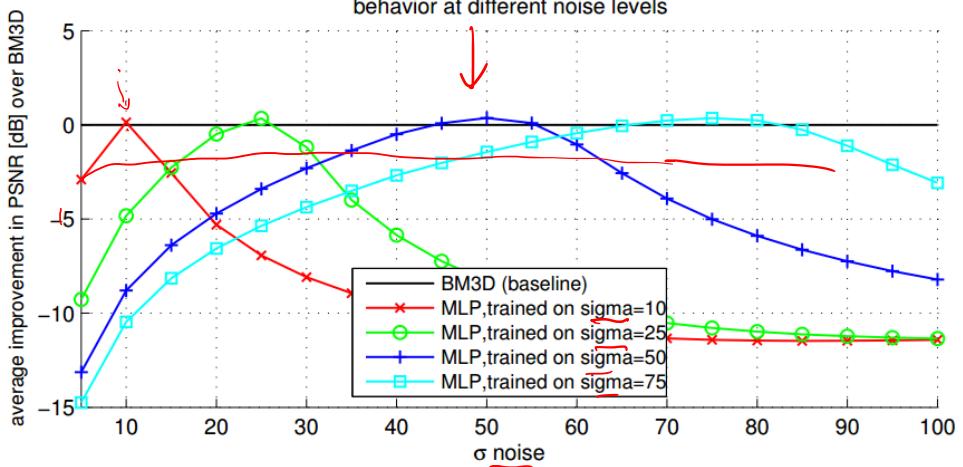


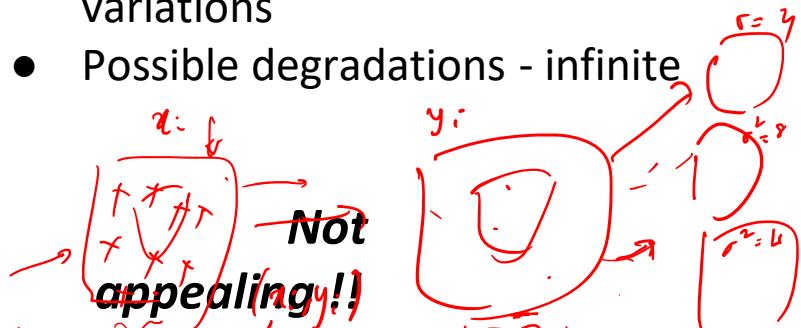
Image denoising using MLPs by HC Burger et al., 2012

Handwritten notes explaining the posterior distribution and MAP estimation:

Given  $p(y|x)$  and  $p(x)$ , the posterior is  $p(x|y) = \frac{p(x)p(y|x)}{p(y)}$ . The MAP estimate is  $\hat{x} = \arg \max_x p(x|y)$ .

Assume  $y = x + n \sim N(\mu_n, \sigma_n^2)$ . Then  $p(y|x) = N(x; \mu_n, \sigma_n^2)$ .

MAP estimate:  $\hat{x} = \mathbb{E}[x|y]$  (MMSE)



$$\hat{x} = \mathbb{E}[x|y]$$

# Generative modeling - *versatility*

$\sim \mathcal{C} + \mathcal{N}$

$$Y = AX + n \quad w(0, r^{\mathbb{I}})$$

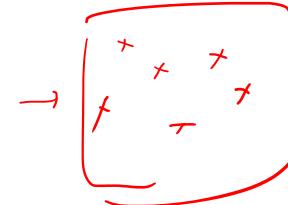
Corrupted Image      Clean Image

$$p(Y|X) = \mathcal{N}(AX, r^{\mathbb{I}})$$

MAP inference

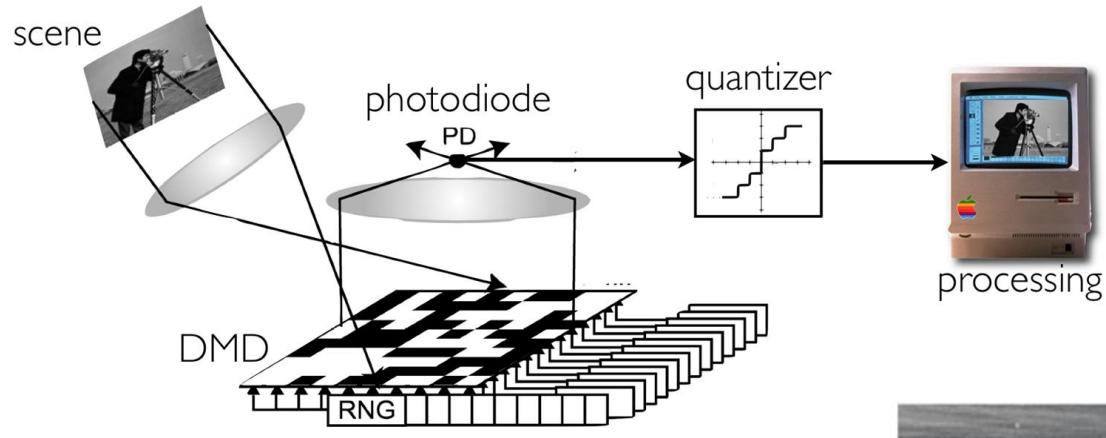
$$\arg \max_X p(Y|X)p(X)$$

Problem Specific      Learned by generative model



One **deep prior** for all problems

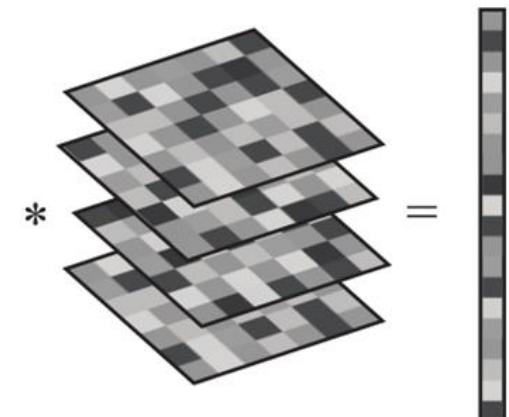
# Single Pixel Camera



SPC by DSP lab, Rice university



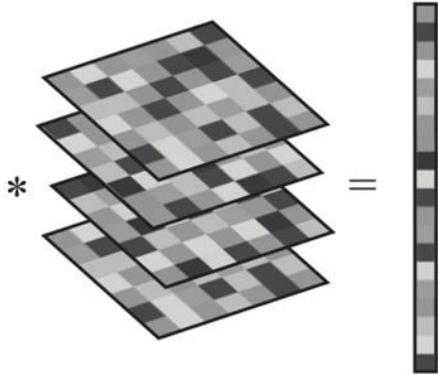
Measurement acquisition simulation



# RIDE for SPC reconstruction



Measurement acquisition



## Inference

While  $k < \text{max\_iter}$ :

- Gradient ascent

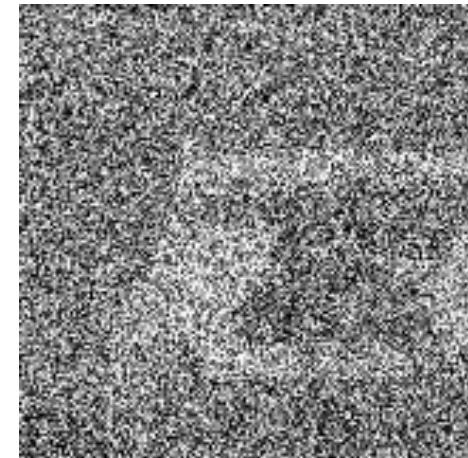
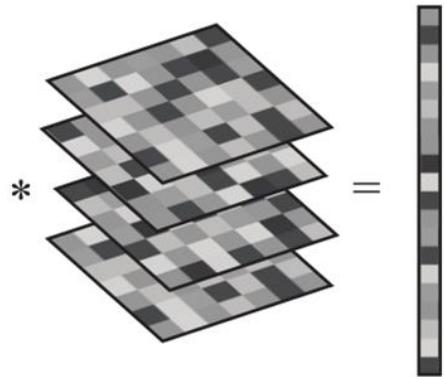
$$\hat{\mathbf{x}}_k = \mathbf{x}_{k-1} + \eta \nabla_{\mathbf{x}_{k-1}} p(\mathbf{x})$$

- Project the solution

$$\mathbf{x}_k = \hat{\mathbf{x}}_k - \Phi^T (\Phi \Phi^T)^{-1} (\Phi \hat{\mathbf{x}}_k - \mathbf{y})$$

$$\hat{\mathbf{x}} = \arg \max_x p(\mathbf{x}) \text{ s.t. } \mathbf{y} = \Phi \mathbf{x}$$

# RIDE for SPC reconstruction

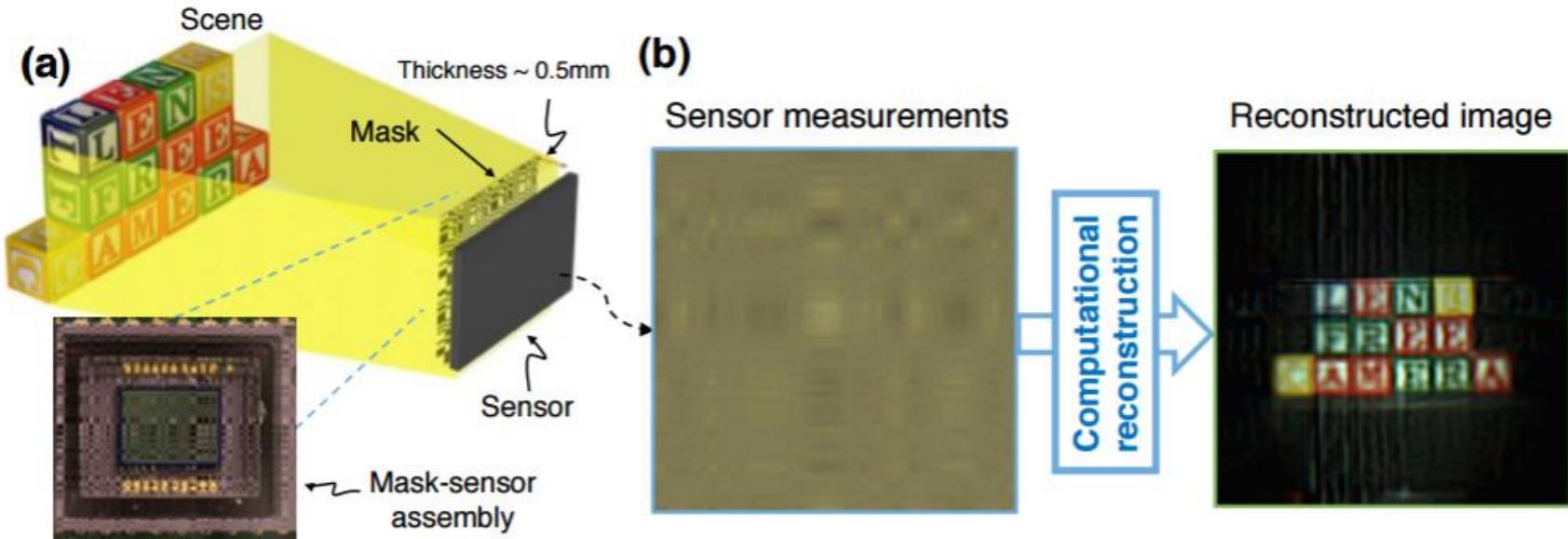
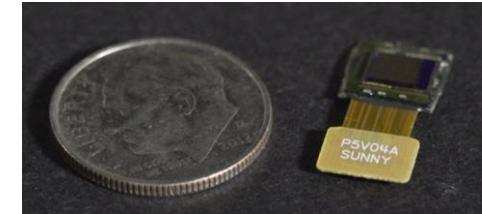


Reconstruction at 30% measurements  
as the iterations progress

$$\hat{\mathbf{x}} = \arg \max_{\mathbf{x}} p(\mathbf{x}) \text{ s.t. } \mathbf{y} = \Phi \mathbf{x}$$

# FlatCam - lens less imaging framework

Thin, flexible and cheap



# FlatCam: results



Original



BM3D

11.82 dB,  
0.199



Ours

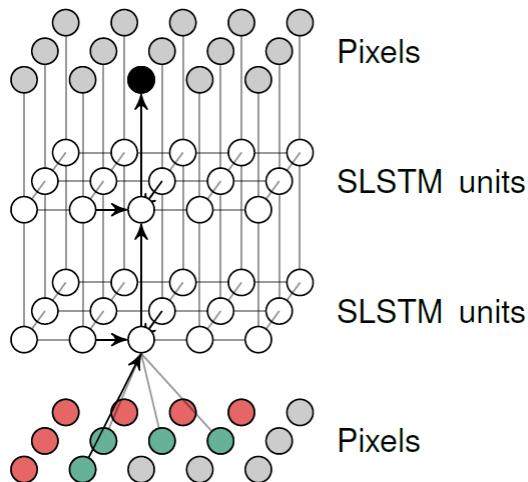
**28.61 dB,**  
**0.325**



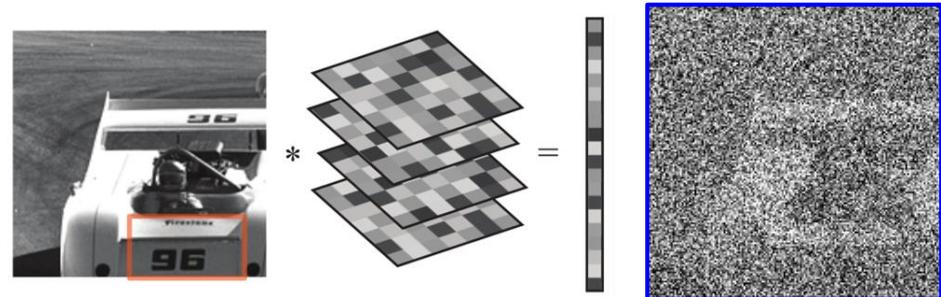
TVAL

11.79 dB,  
0.226

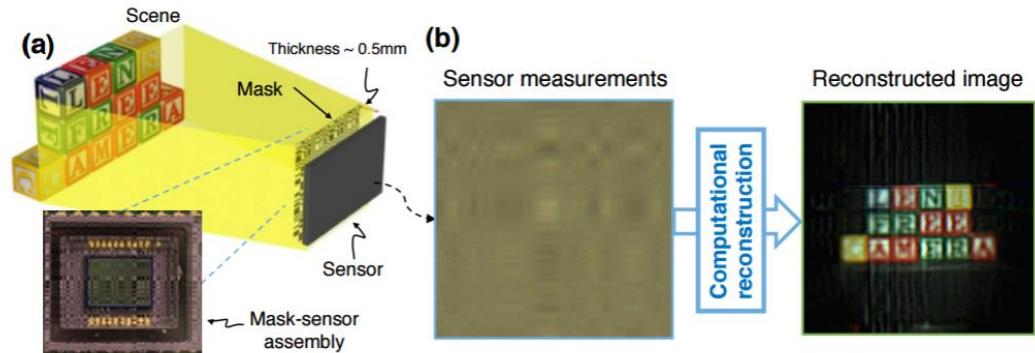
# Summary



Autoregressive models: RIDE,  
PixelCNN/RNN, PixelCNN++,  
Conditional PixelCNN



Single Pixel Camera - reconstructions



FlatCam - reconstructions