

# CSDS 600: Deep Generative Models

## **Autoregressive Models**

Yu Yin ([yu.yin@case.edu](mailto:yu.yin@case.edu))

Case Western Reserve University

# Outline

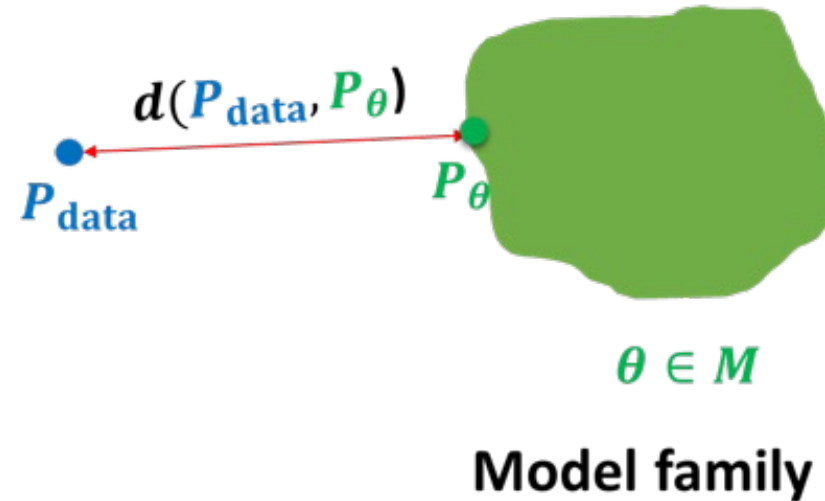
- Basic Representation
- Definition of Autoregressive Models
- Learning and Inference of Autoregressive Models
- Examples
  - FVSBN
  - NADE
  - MADE
  - Pixel RNN
  - Pixel CNN

# Learning a generative model

- Given a training set of examples, e.g., images of dogs

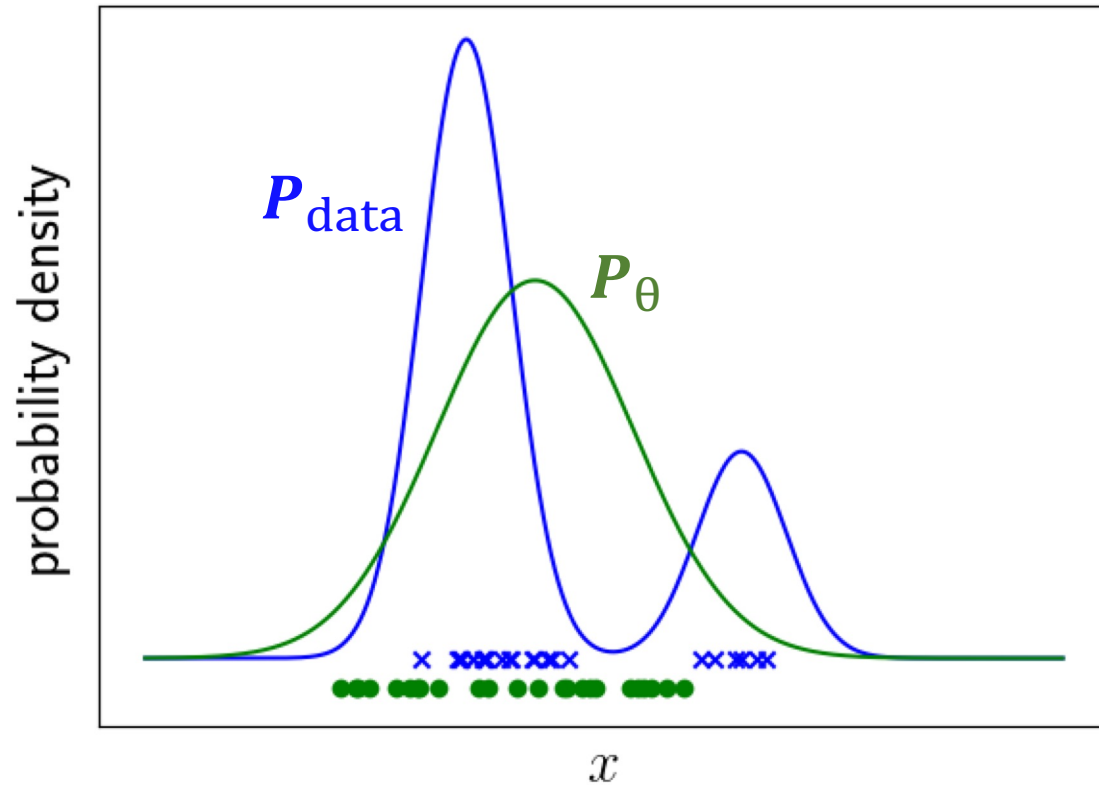


$$\mathbf{x}_i \sim P_{\text{data}} \\ i = 1, 2, \dots, n$$



- Goal: Learn a probability distribution  $p(x)$  over images  $x$

# Example

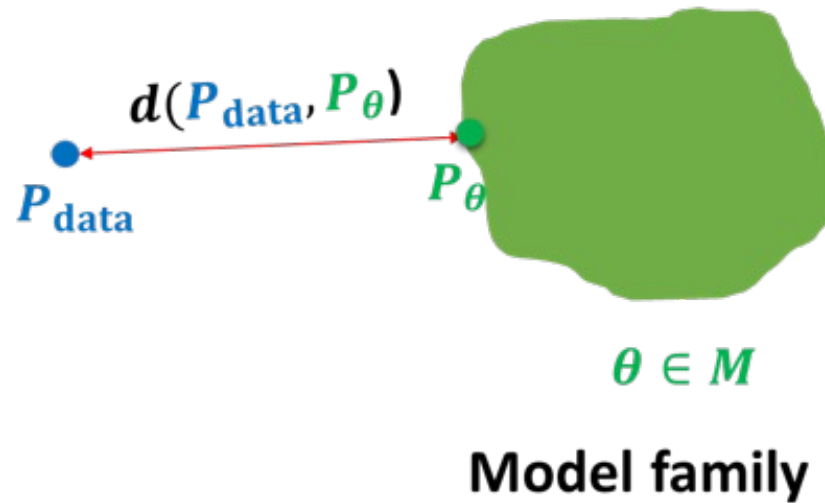


- $P_{\text{data}}$  : data distribution (unknown)
- $\times$  : data samples (known)
- $M$  : mixture of 1d Gaussians
- $P_{\theta}$  : generative model
- $\bullet$  : generated samples

# Challenge of High Dimensionality



$$\mathbf{x}_i \sim P_{\text{data}} \\ i = 1, 2, \dots, n$$



- $P_{\text{data}}$  can contain many modes
- Training samples cover very small region of true support

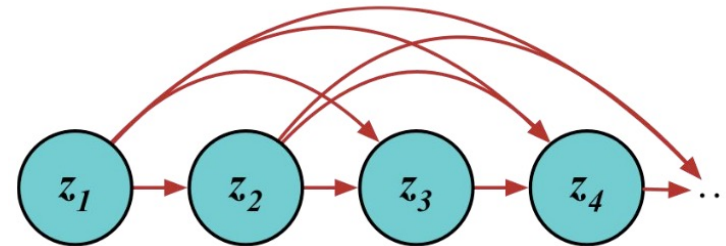
# Learning a generative model

- Three Objectives:
  - **Generation:** If we sample  $x_{new} \sim p(x)$ ,  $x_{new}$  should look like a dog (sampling)
  - **Density estimation:**  $p(x)$  should be high if  $x$  looks like a dog, and low otherwise (anomaly detection)
  - **Unsupervised representation learning:** We should be able to learn what these images have in common, e.g., ears, tail, etc. (features)
- Question: how to **represent** and **learn**  $p(x)$ ?

# Types of Generative Models

- **Fully-observed models**

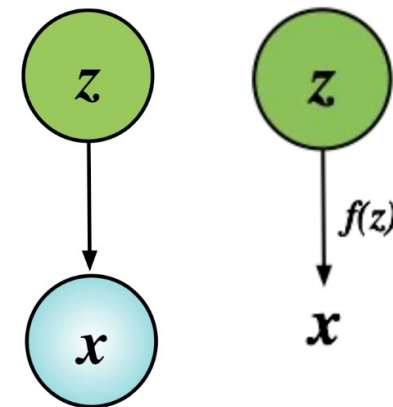
Model observed data directly without introducing any new unobserved local variables.



- **Latent Variable Models**

Introduce an unobserved random variable for every observed data point to explain hidden causes.

- **Prescribed models:** Use observer likelihoods and assume observation noise.
- **Implicit models:** Likelihood-free models (*e.g.*, GANs).



# Motivating Example: MNIST

- **Given:** a dataset  $D$  of handwritten digits (binarized MNIST)



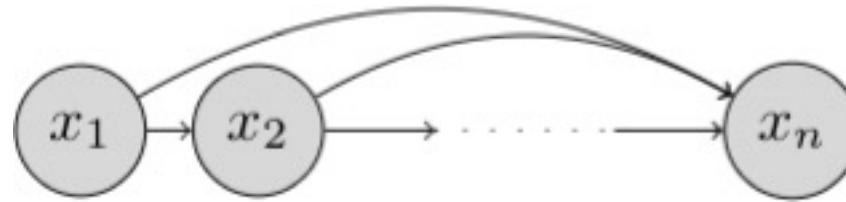
Each image has  $n = 28 \times 28 = 784$  pixels. Each pixel can either be black (0) or white (1).

- **Goal:** Learn a probability distribution  $p(x) = p(x_1, \dots, x_{784})$  over  $x \in \{0, 1\}^{784}$  such that when  $x \sim p(x)$ ,  $x$  looks like a digit



# Autoregressive Models

- Directed, fully-observed graphical models



**Key idea:** Decompose the joint distribution as a **product of tractable conditionals**

$$p_{\theta}(\mathbf{x}) = \prod_{i=1}^n p_{\theta}(x_i | x_1, x_2, \dots, x_{i-1}) = \prod_{i=1}^n p_{\theta}(x_i | x_{<i})$$

# Autoregressive Models

- We can pick an ordering of all the random variables, i.e., raster scan ordering of pixels from top-left ( $X_1$ ) to bottom-right ( $X_{n=784}$ )
- Without loss of generality, we can use chain rule for factorization

$$p(x_1, \dots, x_{784}) = p(x_1)p(x_2 \mid x_1)p(x_3 \mid x_1, x_2) \cdots p(x_n \mid x_1, \dots, x_{n-1})$$

Too complex!

Chain Rule:  $p(x_1, x_2, x_3, x_4) = p(x_1)p(x_2 \mid x_1)p(x_3 \mid x_1, x_2)p(x_4 \mid x_1, x_2, x_3)$

# Autoregressive Models

- To simplify, we assume

$$p(x_1, \dots, x_{784}) = p_{\text{CPT}}(x_1; \alpha^1) p_{\text{logit}}(x_2 \mid x_1; \alpha^2) p_{\text{logit}}(x_3 \mid x_1, x_2; \alpha^3) \dots p_{\text{logit}}(x_n \mid x_1, \dots, x_{n-1}; \alpha^n)$$

conditional Bernoulli ←

$$p_{\text{CPT}}(X_1 = 1; \alpha^1) = \alpha^1, \quad p(X_1 = 0) = 1 - \alpha^1$$

$$p_{\text{logit}}(X_2 = 1 \mid x_1; \alpha^2) = \sigma(\alpha_0^2 + \alpha_1^2 x_1)$$

$$p_{\text{logit}}(X_3 = 1 \mid x_1, x_2; \alpha^3) = \sigma(\alpha_0^3 + \alpha_1^3 x_1 + \alpha_2^3 x_2)$$

- Parameterized functions (e.g., logistic regression) is used to predict next pixel given all the previous ones.

# Learning and Inference

- **Learning** maximizes the model log-likelihood over the dataset !

$$\max_{\theta} \log p_{\theta}(\mathcal{D}) = \sum_{\mathbf{x} \in \mathcal{D}} \log p_{\theta}(\mathbf{x}) = \sum_{\mathbf{x} \in \mathcal{D}} \sum_{i=1}^n \log p_{\theta}(x_i | x_{<i})$$

- Tractable conditionals allow for **exact likelihood evaluation**
  - Conditional evaluated in parallel during training
- Directed model permits **ancestral sampling**, one variable at a time

$$x_1 \sim p_{\theta}(x_1), x_2 \sim p_{\theta}(x_2 | x_1), \dots, x_n \sim p_{\theta}(x_n | x_{<n})$$

# Learning and inference

- **Inference** samples each variable of one data from estimated conditional distributions step by step, until the whole data is generated.
- Ancestral sampling:
  - A process of producing samples from a probabilistic model.
  - **First** sample variables which has no conditional constraints using their prior distribution.  $x_1 \sim p_\theta(x_1)$
  - **Then** sample child variables using conditional distribution based on their parents and repeat so on.  $x_2 \sim p_\theta(x_2|x_1)$
- The attribute of Autoregressive Models that directly model and output distributions allows for ancestral sampling.

# Learning and inference

Differences between Autoregressive models (AR), VAE and GAN:

- **GAN** model doesn't define any distribution, it adapts discriminator to learn the data distribution implicitly.  $P(X, Z) = P(X|Z)P(Z)$
- **VAE** model believes the data distribution is too complex to model directly, thus it tries to learn the distribution by defining an intermediate distribution and learning the map between the defined simple distribution to the complex data distribution.  $P(X, Z) = P(X|Z)P(Z)$
- **AR** model on the one hand assumes that the data distribution can be learned directly (tractable), then it define its outputs as conditional distributions to solve the generation problem by directly modeling each conditional distribution.

# Learning and inference

## Conclusion:

1. Using complex networks, each step Autoregressive Models output an approximated complex conditional distribution  $\hat{x}_i = p_\theta(x_i|x_1, x_2, \dots, x_{i-1})$
2. Taking in the previous inputs  $x_1, x_2, \dots, x_{i-1}$  and the next input  $x_i$  by sampling previous estimated conditional distribution  $\hat{x}_i$ , Autoregressive Model is able to generate all conditional distribution iteratively

$$x_1 \sim P_\theta(x_1), x_2 \sim P_\theta(x_2|x_1), x_3 \sim P_\theta(x_3|x_1, x_2), \dots x_n \sim P_\theta(x_n|x_1, \dots, x_{n-1})$$

3. Product rule makes sure the generated data that made up of sampled result  $x_i$  from each step follows the data distribution.

$$(x_1, x_2, \dots, x_n) \sim \prod_{i=1}^n p_\theta(x_i|x_{<i})$$

# Parameterizing Autoregressive Models

Linear/multi-layer NN based parameterizations

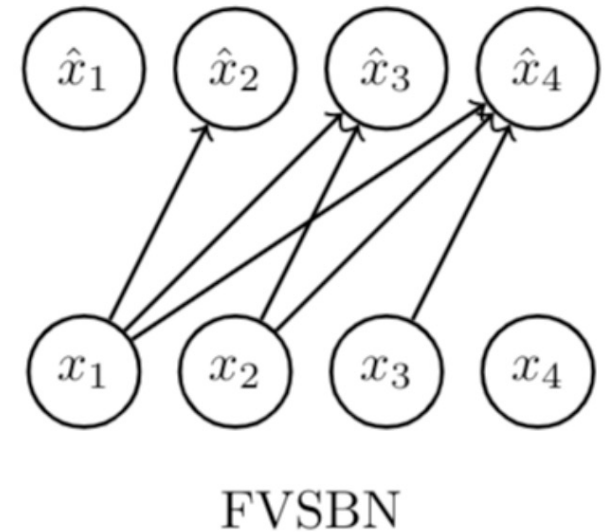
- Fully-visible sigmoid belief networks (FVSBN)
- Neural autoregressive density estimator (NADE): 1 hidden layer NN
- Masked autoencoder distribution estimation (MADE): multilayer NN



# Fully Visible Sigmoid Belief Network (FVSBN)

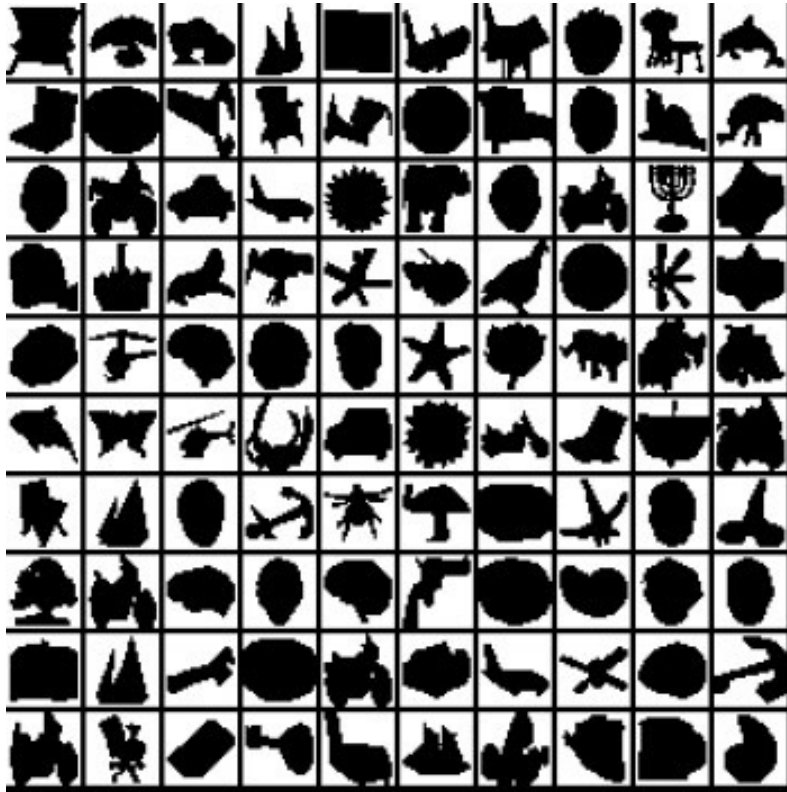
- The fully visible sigmoid belief network **without any hidden** units is denoted FVSBN.
- The conditional variables  $x_i | x_{<i}$  in FVSBN are **Bernoulli** with parameters.
- Some conditionals are too complex. So FVSBN assume:

$$\begin{aligned}\hat{x}_i &= p(x_i = 1 | x_1, x_2, \dots, x_{i-1}) = f_i(x_1, x_2, \dots, x_{i-1}; \alpha^i) \\ &= \sigma(\alpha_0^{(i)} + \alpha_1^{(i)} x_1 + \dots + \alpha_{i-1}^{(i)} x_{i-1})\end{aligned}$$

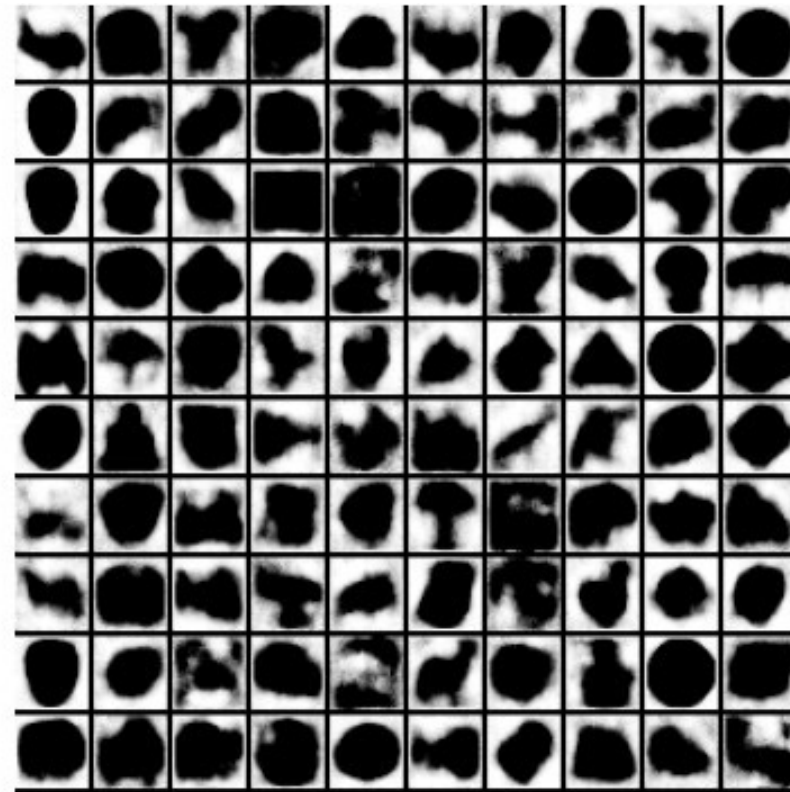


- # of parameters:  $1 + 2 + 3 + \dots + n \approx n^2/2$

# FVSBN Results



Training data from Caltech  
101 Silhouettes



Samples from the model

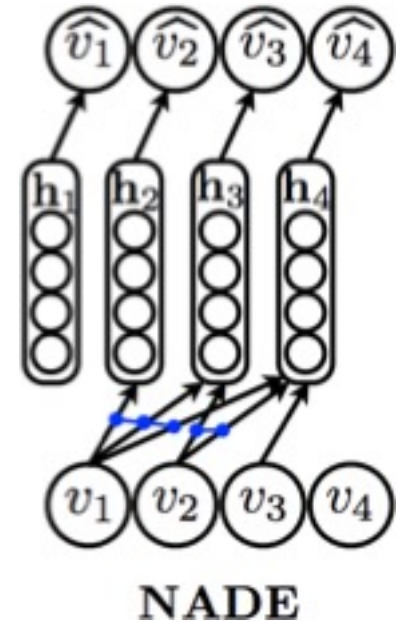
# Neural autoregressive density estimator (NADE)

Improve FVSBN: use one hidden layer neural network instead of logistic regression

$$\mathbf{h}_i = \sigma(A_i \mathbf{x}_{<i} + \mathbf{c}_i)$$

$$\hat{\mathbf{x}}_i = p(x_i | x_1, x_2 \dots x_{i-1}; \underbrace{A_i, \mathbf{c}_i, \boldsymbol{\alpha}_i, b_i}_{\text{parameters}}) = \sigma(\boldsymbol{\alpha}_i \mathbf{h}_i + b_i)$$

$\mathbf{x}_i = \mathbf{v}_i$



$$\mathbf{h}_2 = \sigma \left( \underbrace{\begin{pmatrix} \vdots \\ \mathbf{w}_1 \\ \vdots \end{pmatrix}}_{A_2} x_1 \right) \quad \mathbf{h}_3 = \sigma \left( \underbrace{\begin{pmatrix} \vdots \\ \mathbf{w}_1 & \mathbf{w}_2 \\ \vdots \end{pmatrix}}_{A_3} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \right) \quad \mathbf{h}_4 = \sigma \left( \underbrace{\begin{pmatrix} \vdots \\ \mathbf{w}_1 & \mathbf{w}_2 & \mathbf{w}_3 \\ \vdots \end{pmatrix}}_{A_3} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} \right)$$

Tie weights are shared to *reduce the number of parameters* and *speedup computation* (see **blue** dots in the figure)

# NADE

$$\mathbf{h}_i = \sigma(\mathbf{A}_i \mathbf{x}_{<i} + \mathbf{c}_i)$$

$$\hat{\mathbf{x}}_i = p(x_i | x_1, x_2 \dots x_{i-1}; \mathbf{A}_i, \mathbf{c}_i, \boldsymbol{\alpha}_i, b_i)$$

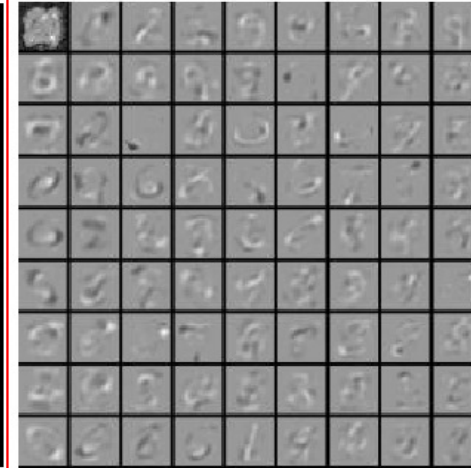
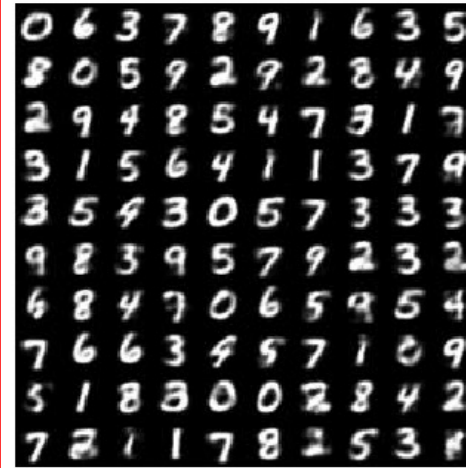
$$= f_i(x_1, x_2, \dots, x_{i-1}) = \sigma(\boldsymbol{\alpha}_i \mathbf{h}_i + b_i)$$

- Parameters:  $\theta_i = \{\mathbf{A}_i \in R^{d \times (i-1)}, \mathbf{c}_i \in R^d, \boldsymbol{\alpha}_i \in R^d, b_i \in R\}$
- The total number of parameters in this model is dominated by the matrices  $\{\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_n\}$  given by  $O(n^2 * d)$ .
- Sharing parameters: Tie weights are shared to reduce the number of parameters and speed up computation.  $\rightarrow O(n * d)$ .

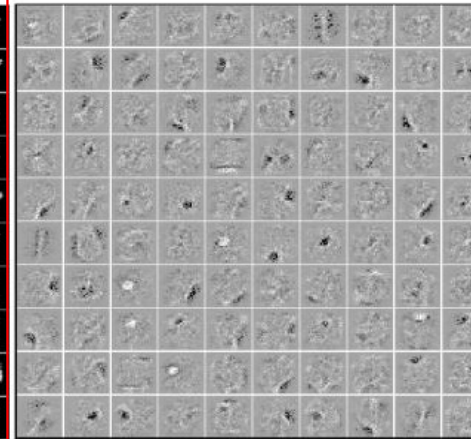


# Results on MNIST

FVSBN



NADE



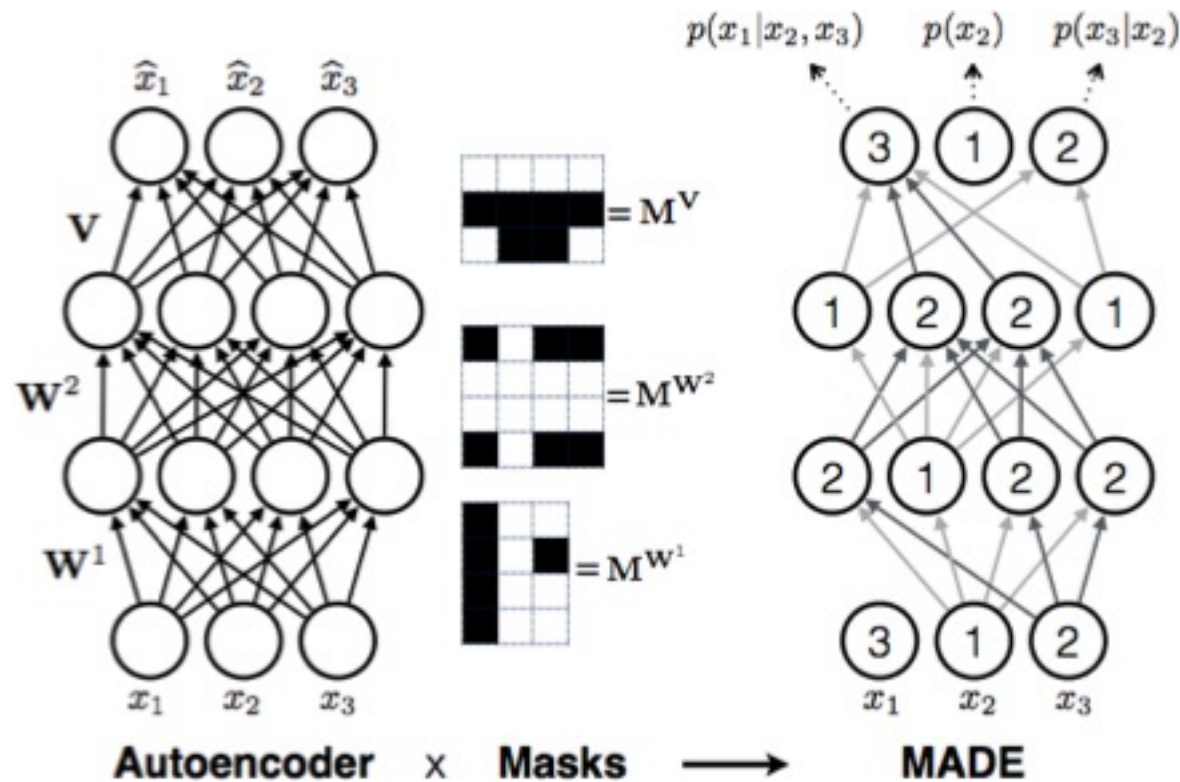
Training data

Generated  
results

Learned Features of the  
bottom layer

# Masked autoencoder distribution estimation (MADE)

- Extension to multilayer NN



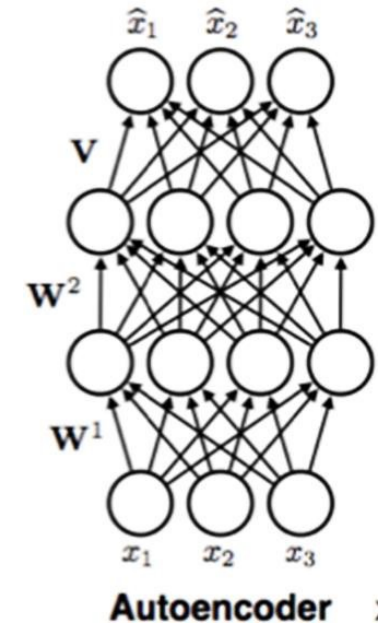
Germain et al., 2015

# Autoregressive model vs. Autoencoders

- FVSBN and NADE look similar to an **autoencoder**.
- An **encoder**  $e(\cdot)$ . E.g.,  $e(x) = \sigma(W^2(W^1x + b^1) + b^2)$
- A **decoder** such that  $d(e(x)) \approx x$ . E.g.,  $d(h) = \sigma(Vh + c)$ .
- Loss function: (if continuous)

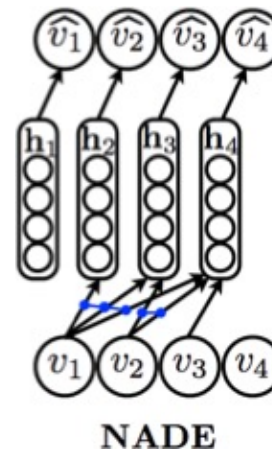
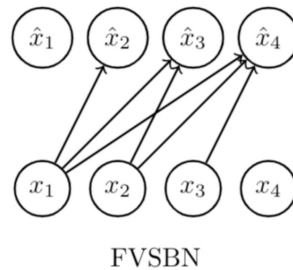
$$\min_{W^1, W^2, b^1, b^2, V, c} \sum_{x \in D} \sum_i (x_i - \hat{x}_i)^2$$

- Encoder: feature learning
- A vanilla autoencoder is not a generative model: it does not define a distribution over  $x$  we can sample from to generate new data points.

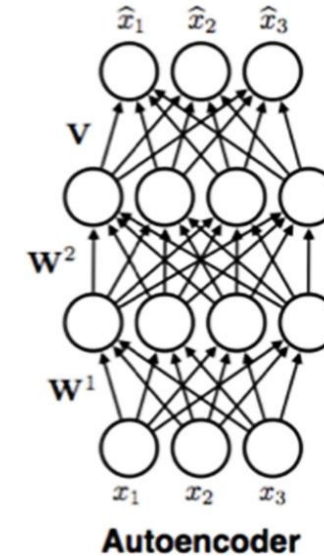


# Autoregressive model vs. Autoencoders

- FVSBN and NADE look similar to an autoencoder.
- Can we get a generative model from an Autoencoder?



VS



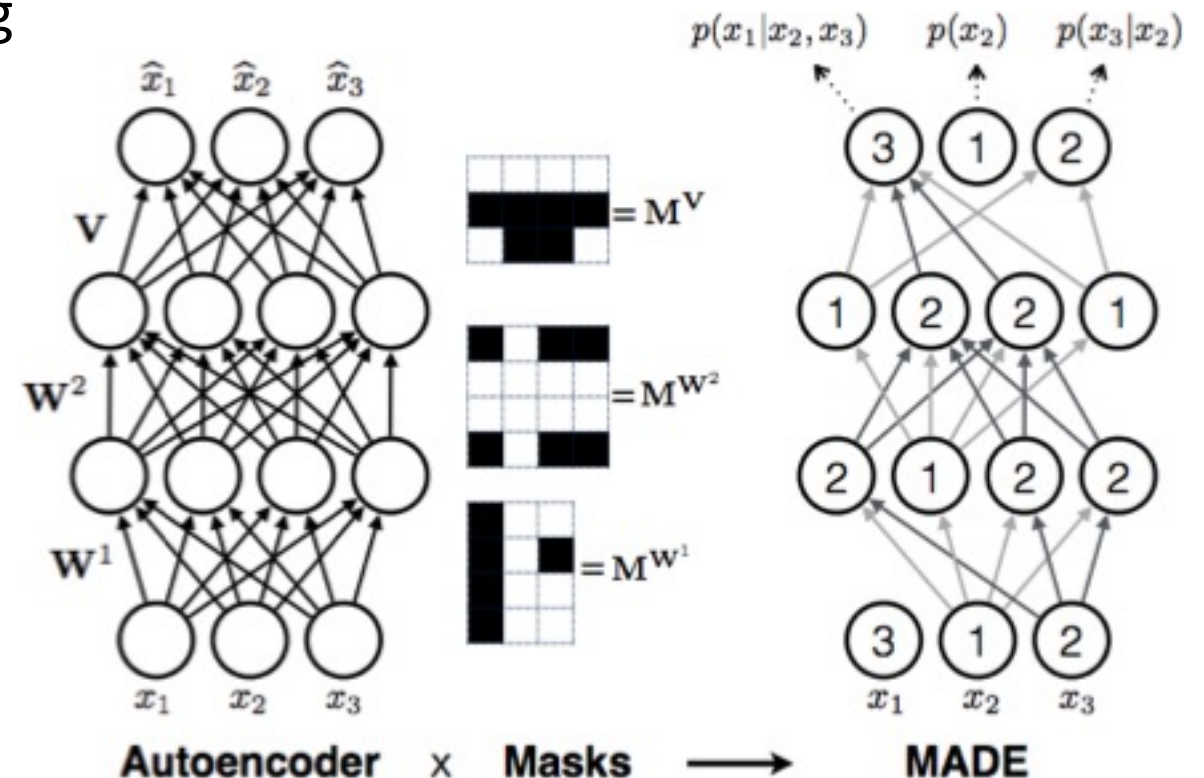
- A dependency order constraint is required for Autoencoder to make it a Bayesian Network.



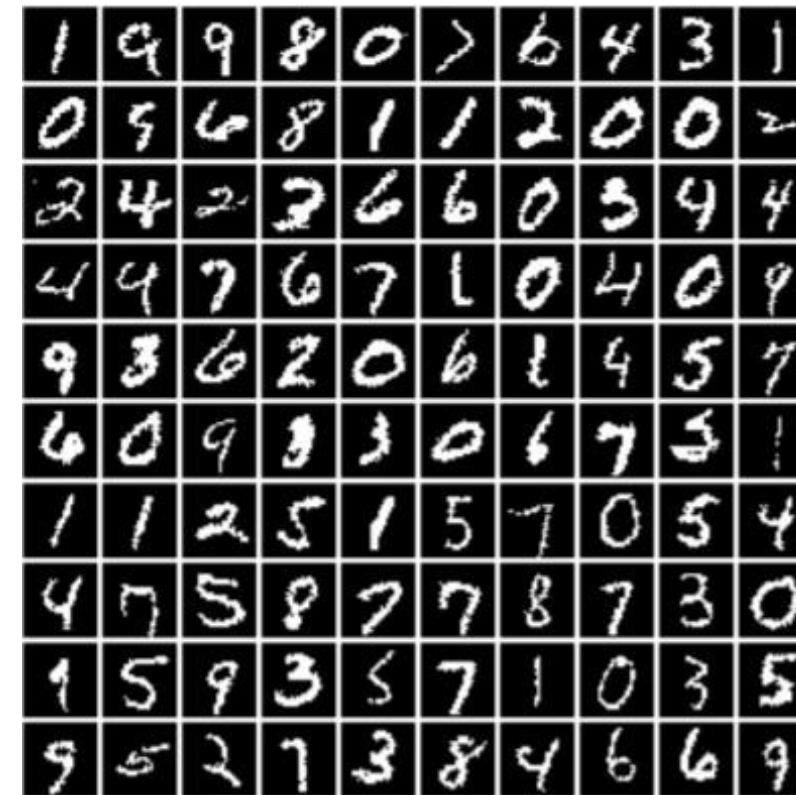
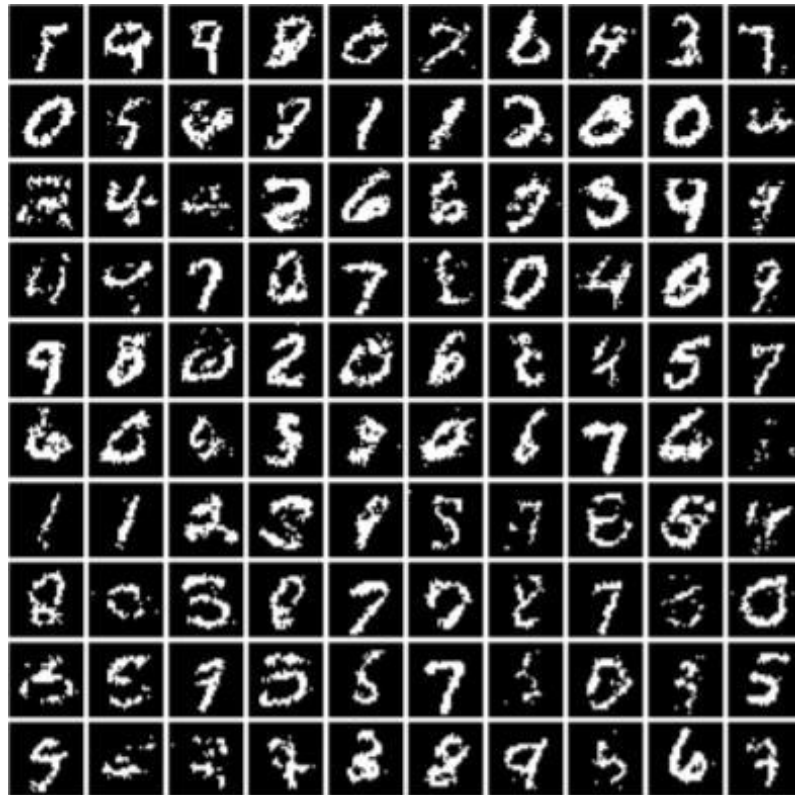
# Masked autoencoder distribution estimation (MADE)

- Use Masks to constraint dependency paths!
- Each output unit is an estimated distribution, it only depends on the inputs with orderings that before its chosen ordering

- With the order  $x_2, x_3, x_1$ :
  1.  $p(x_2)$  doesn't depends on any input
  2.  $p(x_3|x_2)$  depends on input  $x_2$
  3.  $p(x_1|x_2, x_3)$  depends on input  $x_2, x_3$



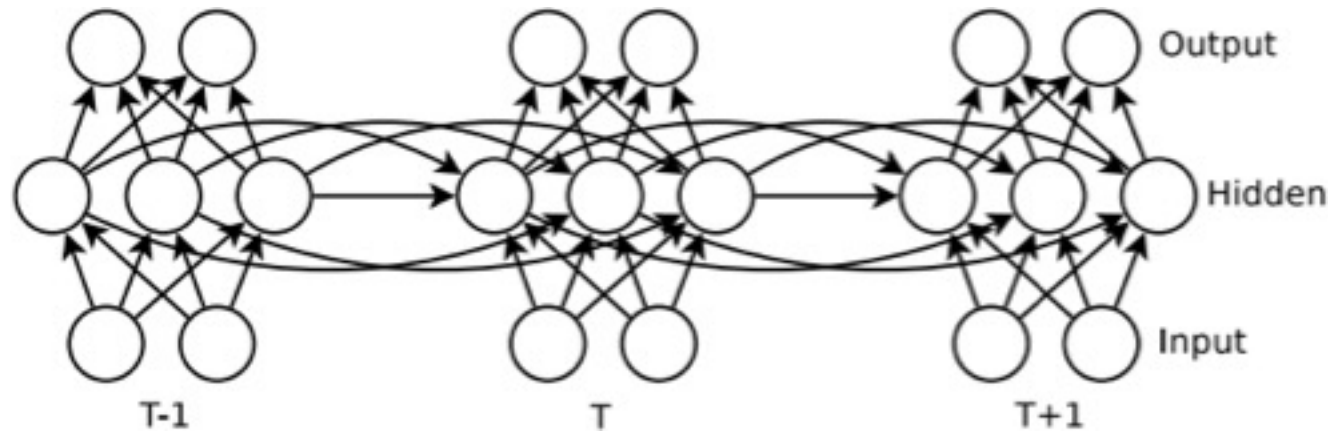
# MADE Results on MNIST



Samples from a 2 hidden layer MADE      Nearest neighbor in dataset

# RNN-based parameterizations

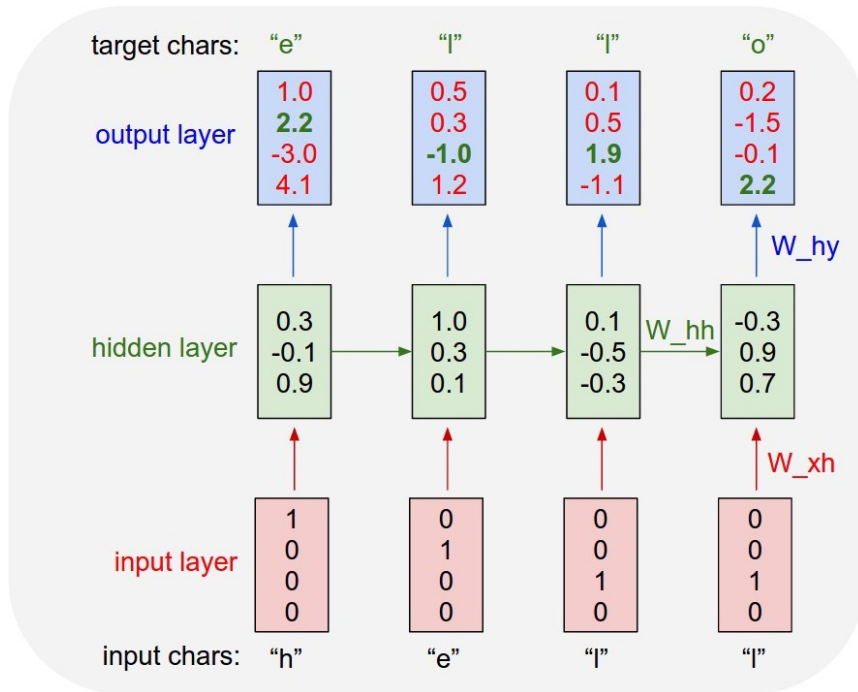
Effective architecture for learning long-range sequential dependencies



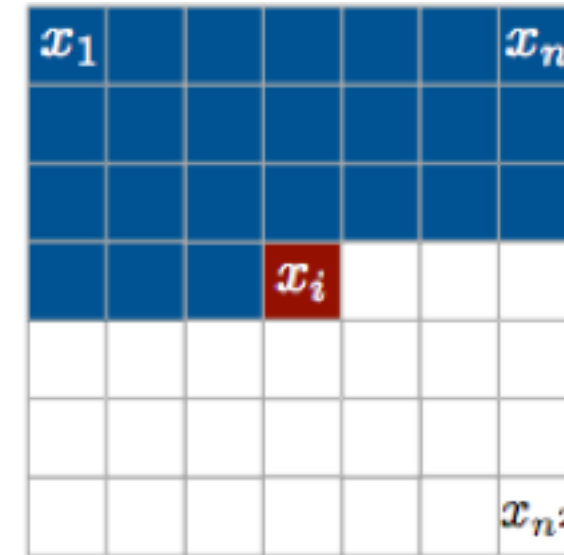
Hidden layer  $h_{\theta}(t)$  is a summary of the inputs seen till  $t - 1$

Output layer  $o_{\theta}(t)$  specifies parameters for conditionals  $p_{\theta}(x_t | x_{<t})$

# RNN-based parameterizations



Char-RNN



Pixel-RNN

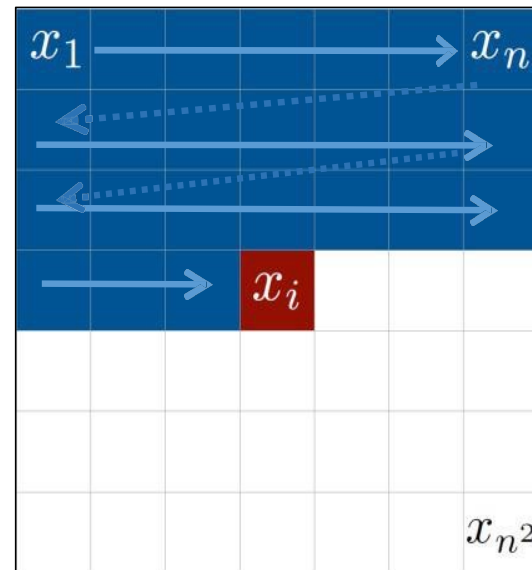
Sutskevar et al., 2011, Karpathy, 2015, Theis & Bethge, 2015, van den Oord, 2016a

# Pixel RNN (Oord et al., 2016)

- Model images pixel by pixel using raster scan order

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

- Sequentially predict pixels rather than predicting the whole image at once (like as GAN, VAE)

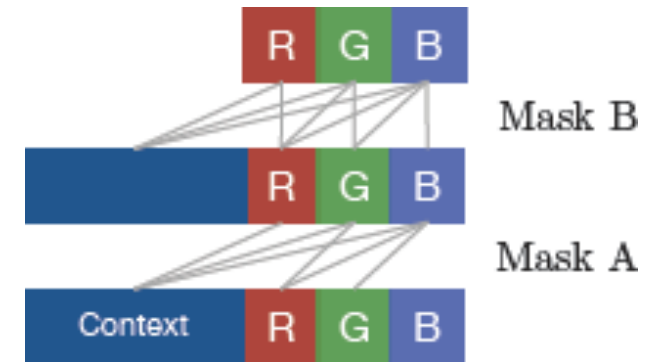
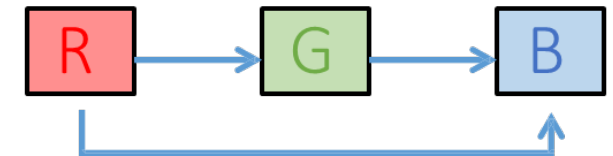


# Pixel RNN (Oord et al., 2016)

- For color image, 3 channels are generated successive conditioning, **blue** given **red** and **green**, **green** given **red**, and **red** given only the pixels above and to the left of all channels

$$p(x_t | x_{1:t-1}) = p(x_t^{red} | x_{1:t-1})p(x_t^{green} | x_{1:t-1}, x_t^{red})p(x_t^{blue} | x_{1:t-1}, x_t^{red}, x_t^{green})$$

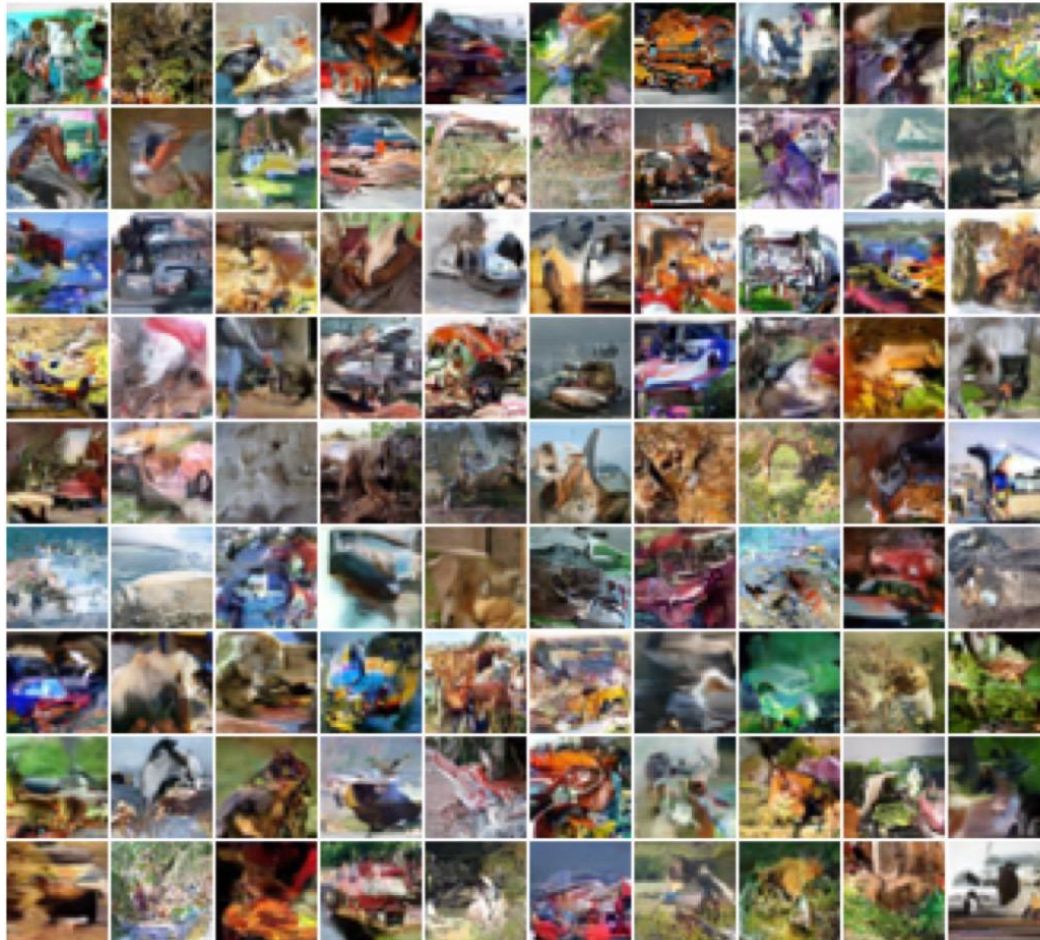
- Conditionals modeled using RNN variants. LSTMs + masking (like MADE)



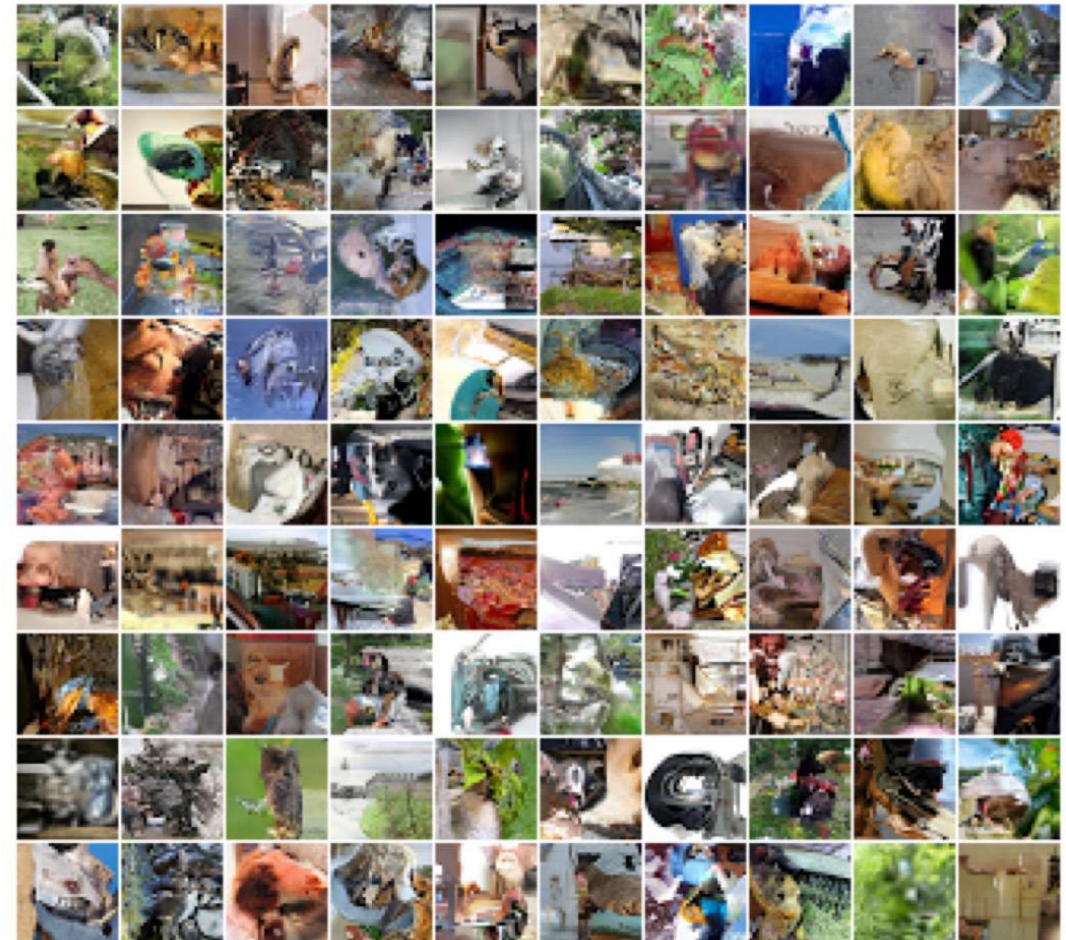


# Results of Pixel RNN

CIFAR-10



ImageNet 32x32





# Results of Pixel RNN



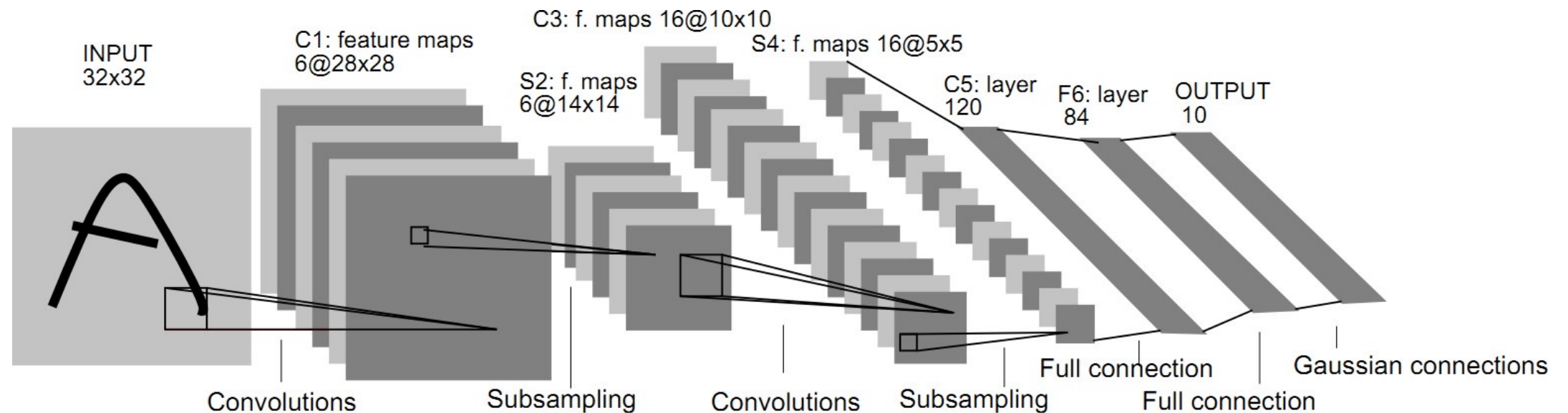
Results on downsampled ImageNet

**Very slow!** (sequential likelihood evaluation)



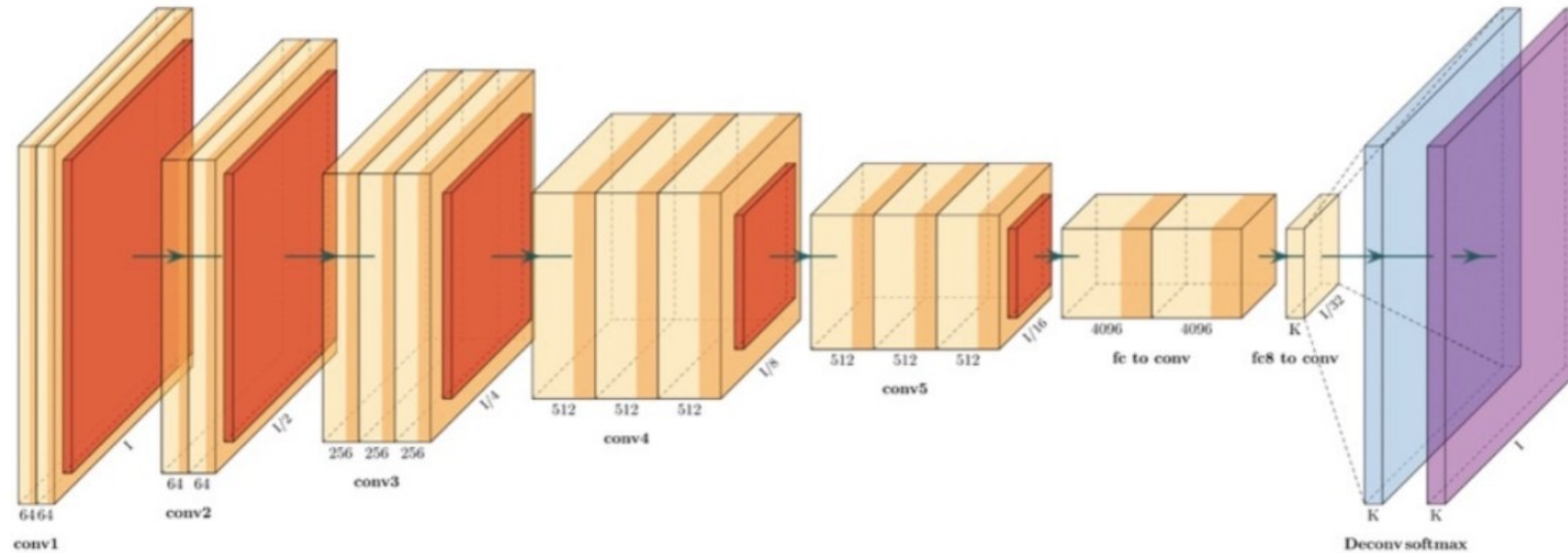
# CNN-based parameterizations

Convolutions are amenable for parallelization on modern hardware



Lecunn et al., 1998

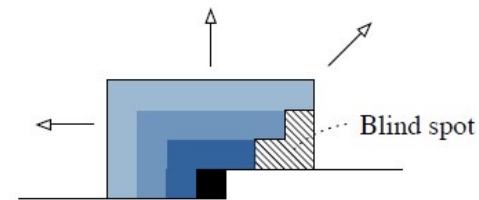
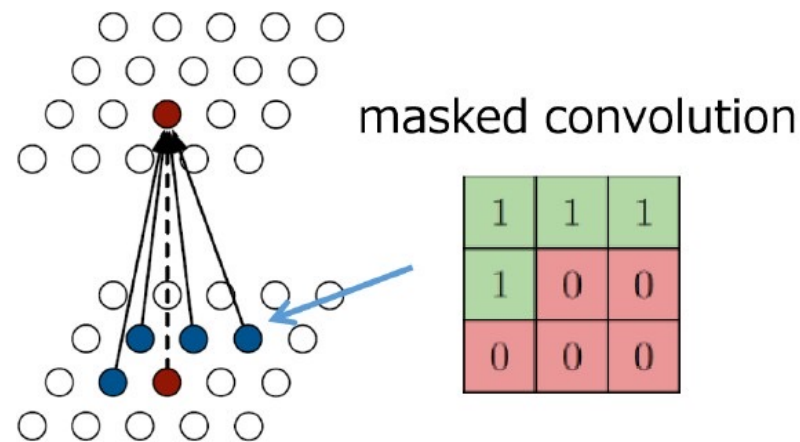
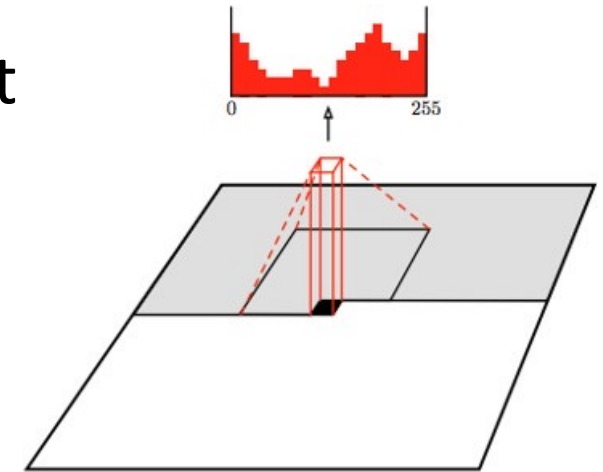
# Convolutional Architectures



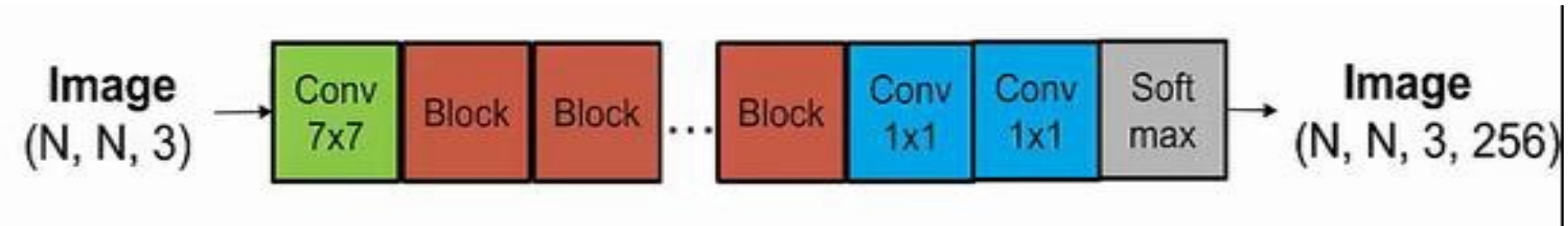
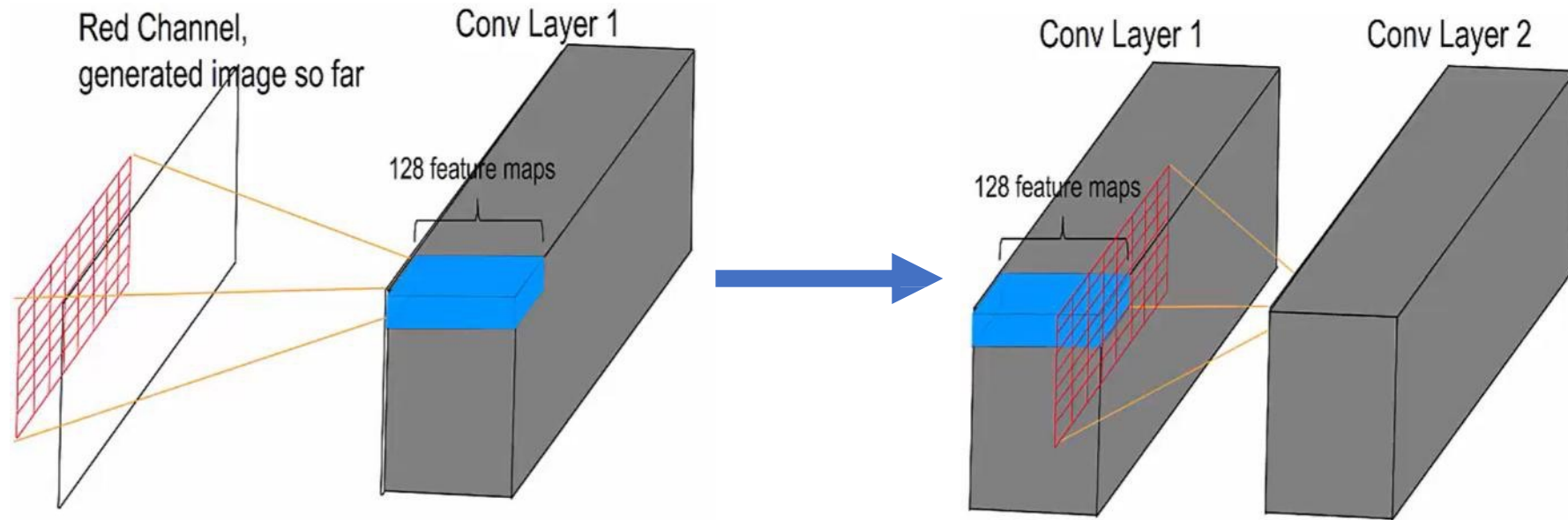
Convolutions are natural for image data and easy to parallelize on modern hardware.

# PixelCNN (Oord et al., 2016)

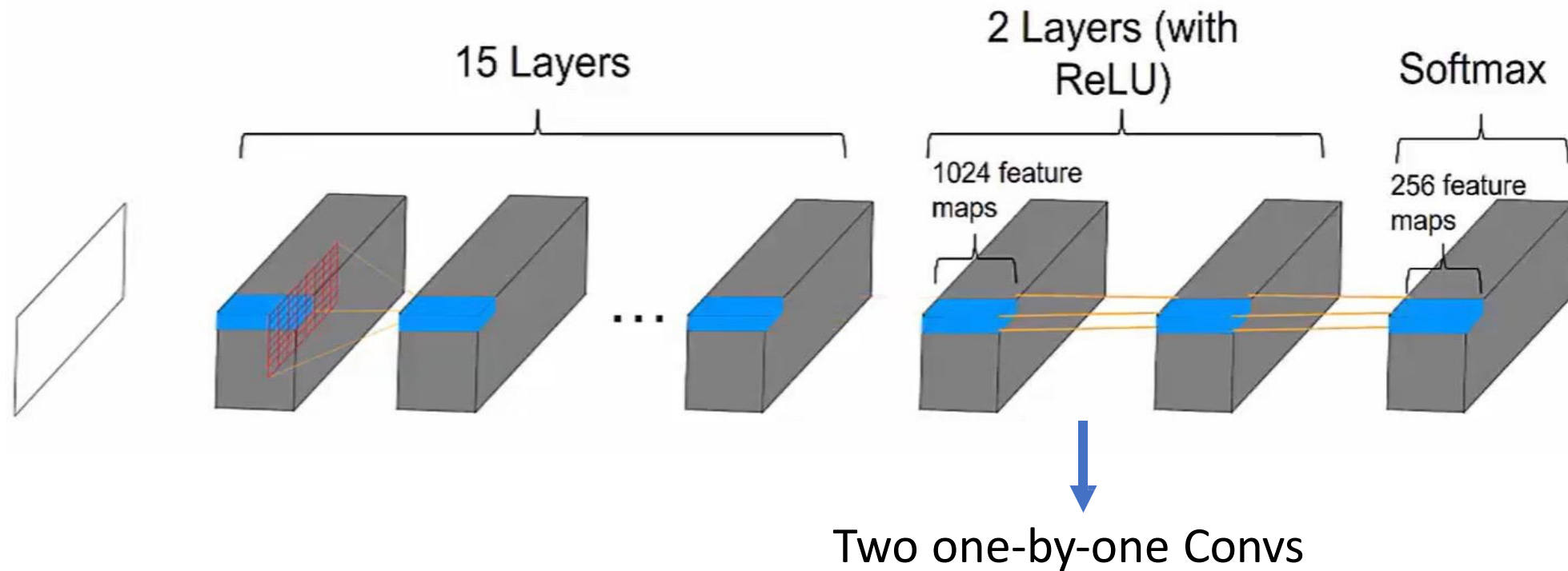
- **Idea:** Use convolutional architecture to predict next pixel given context (a neighborhood of pixels).
- **Challenge:** Has to be autoregressive. Masked convolutions preserve raster scan order. Additional masking for colors order.



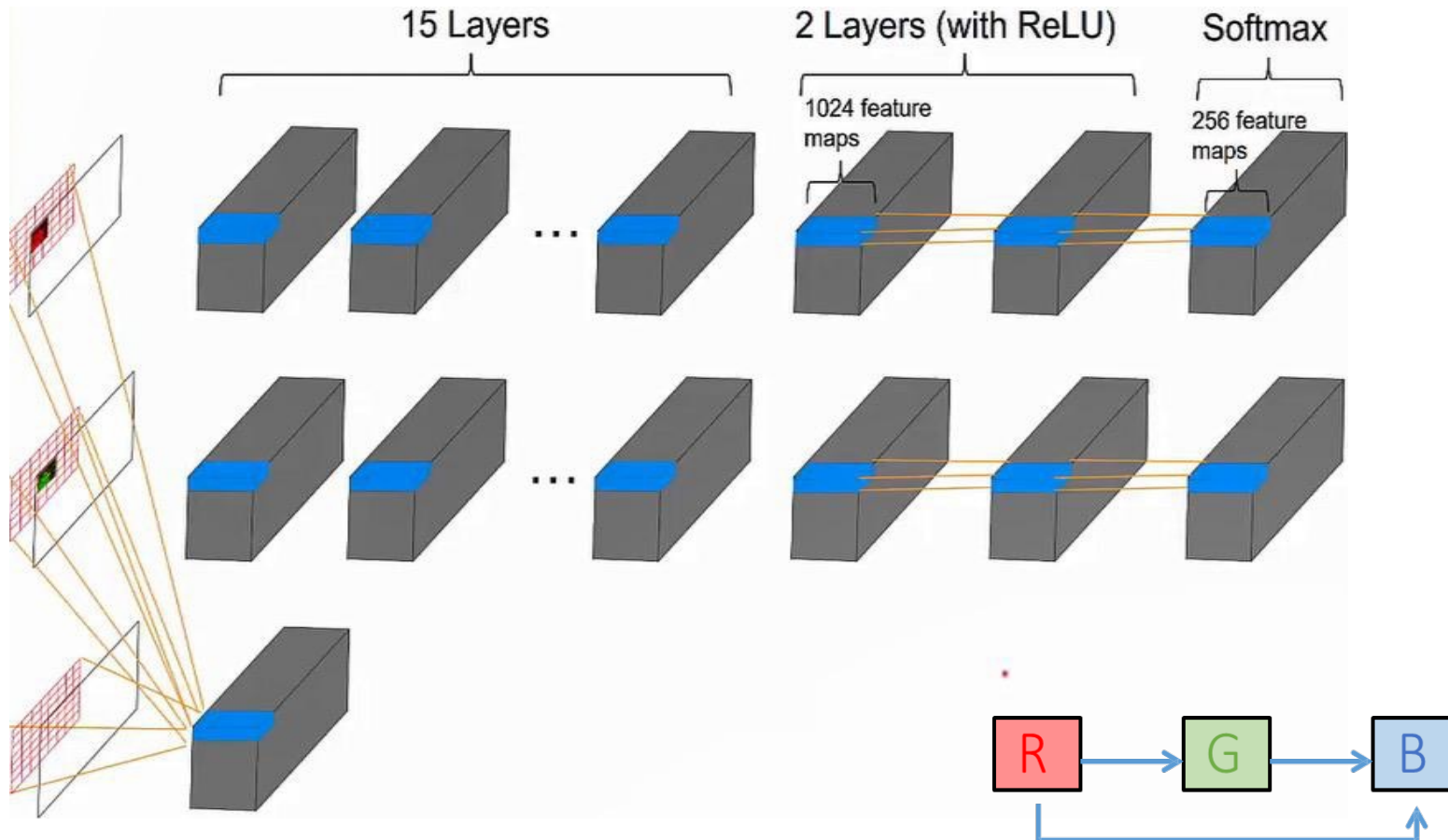
# PixelCNN Network Structure



# PixelCNN Network Structure



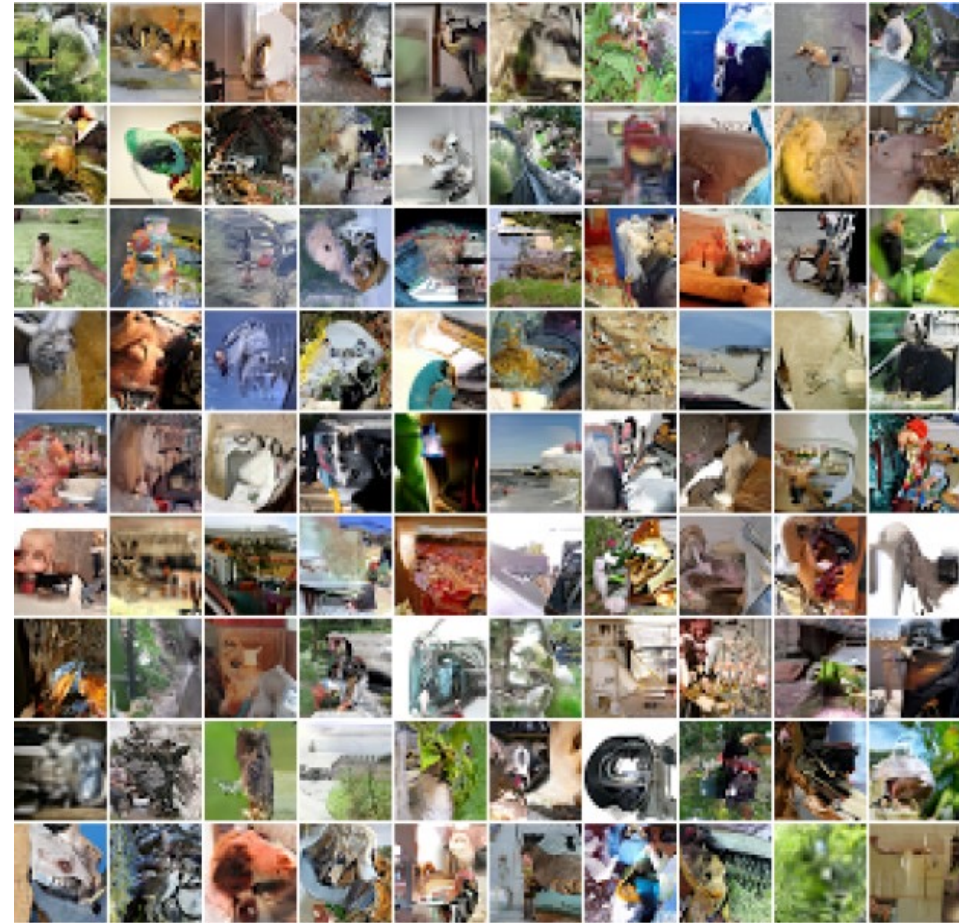
# PixelCNN Network Structure





# PixelCNN

Samples from the model  
trained on Imagenet ( $32 \times 32$   
pixels). Similar performance  
to PixelRNN, but much faster.



# Summary

- Autoregressive models:
  - Chain rule based factorization is fully general
  - Compact representation via conditional independence and/or neural parameterizations
- Autoregressive models Pros:
  - Easy to evaluate likelihoods
  - Easy to train
- Autoregressive models Cons:
  - Requires an ordering
  - Generation is sequential
  - Cannot learn features in an unsupervised way



# Thank You

- Questions?
- Email: [yu.yin@case.edu](mailto:yu.yin@case.edu)