



Topics in AI (CPSC 532S): Multimodal Learning with Vision, Language and Sound

Lecture 15: Generative Models

Supervised vs. Unsupervised Learning

Supervised Learning

Data: (x, y)

x is data, y is label

Goal: Learn a *function* to map $x \rightarrow y$

Examples: Classification, regression, object detection, semantic segmentation, image captioning, *etc.*



→ Cat

Classification

[This image](#) is [CC0 public domain](#)

Supervised vs. Unsupervised Learning

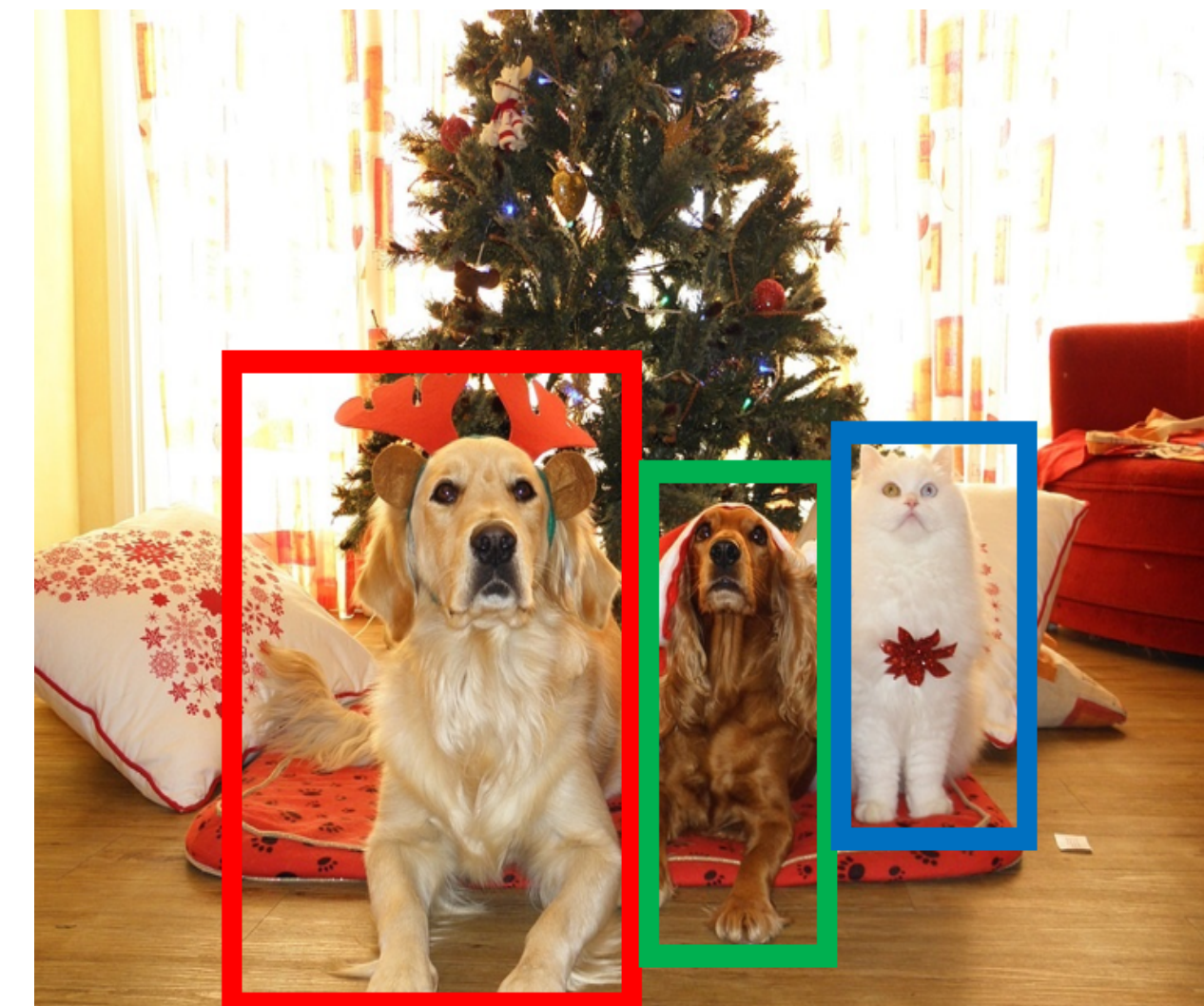
Supervised Learning

Data: (x, y)

x is data, y is label

Goal: Learn a *function* to map $x \rightarrow y$

Examples: Classification, regression, object detection, semantic segmentation, image captioning, *etc.*



DOG, DOG, CAT

Object Detection

[This image is CC0 public domain](#)

Supervised vs. Unsupervised Learning

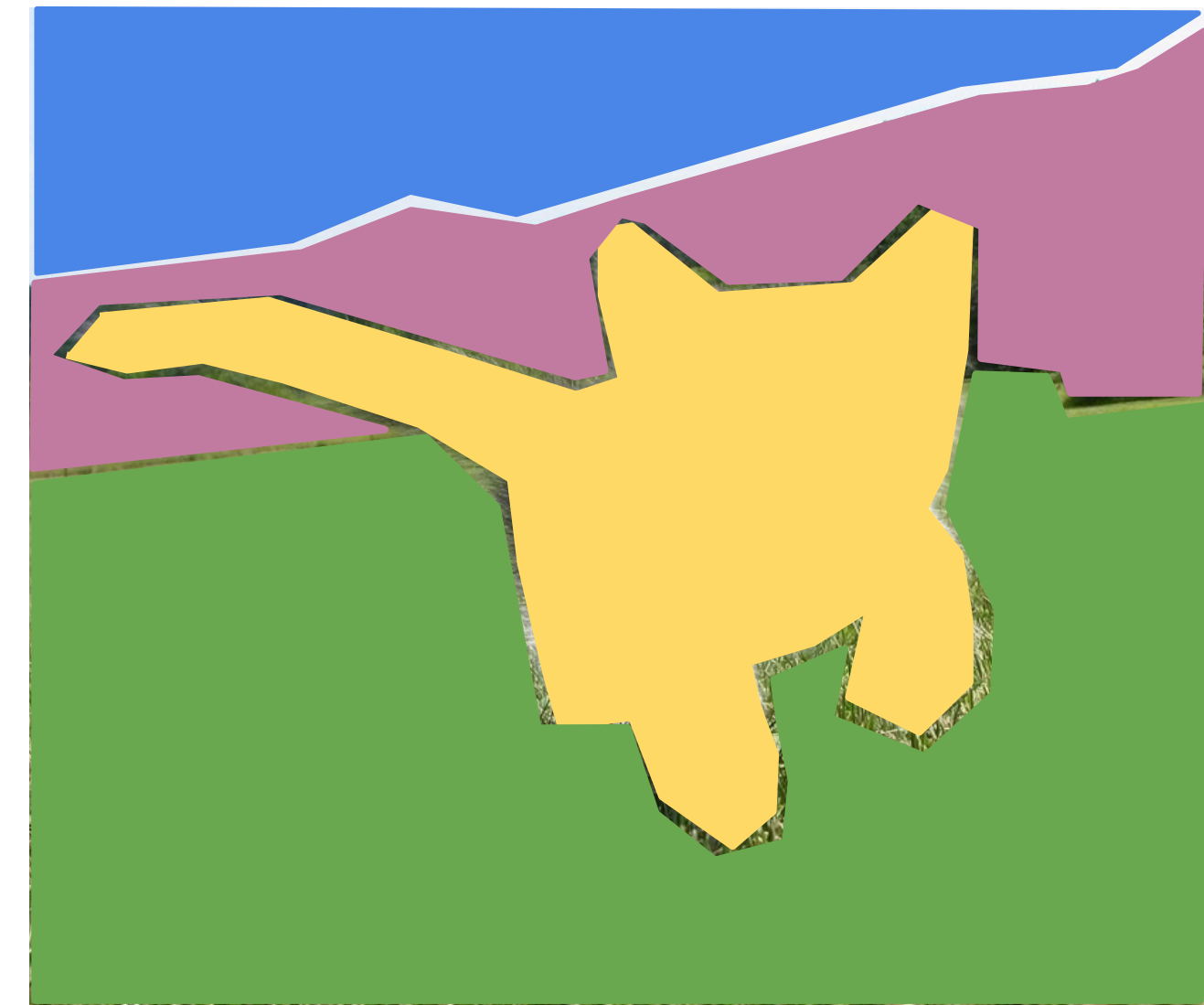
Supervised Learning

Data: (x, y)

x is data, y is label

Goal: Learn a *function* to map $x \rightarrow y$

Examples: Classification, regression, object detection, semantic segmentation, image captioning, *etc.*



GRASS, CAT, TREE, SKY

Semantic Segmentation

[This image is CC0 public domain](#)

Supervised vs. Unsupervised Learning

Supervised Learning

Data: (x, y)

x is data, y is label

Goal: Learn a *function* to map $x \rightarrow y$

Examples: Classification, regression, object detection, semantic segmentation, image captioning, *etc.*



A cat sitting on a suitcase on the floor

Image Captioning

[This image](#) is [CC0 public domain](#)

Supervised vs. Unsupervised Learning

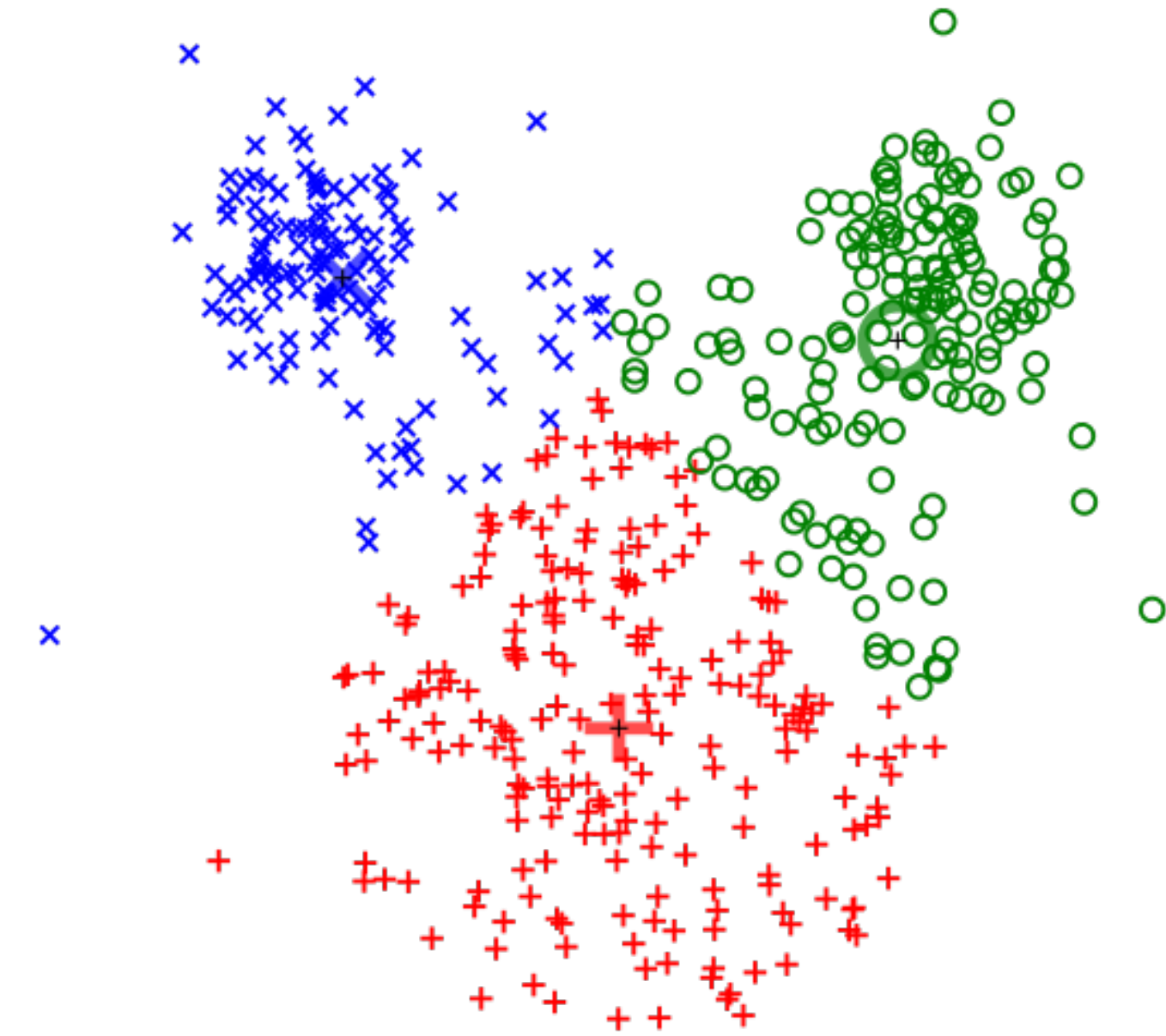
Unsupervised Learning

Data: x

Just data, no labels!

Goal: Learn some underlying hidden *structure* of the data

Examples: Clustering, dimensionality reduction, feature learning, density estimation, *etc.*



k-means clustering

[This image is CC0 public domain](#)

Supervised vs. Unsupervised Learning

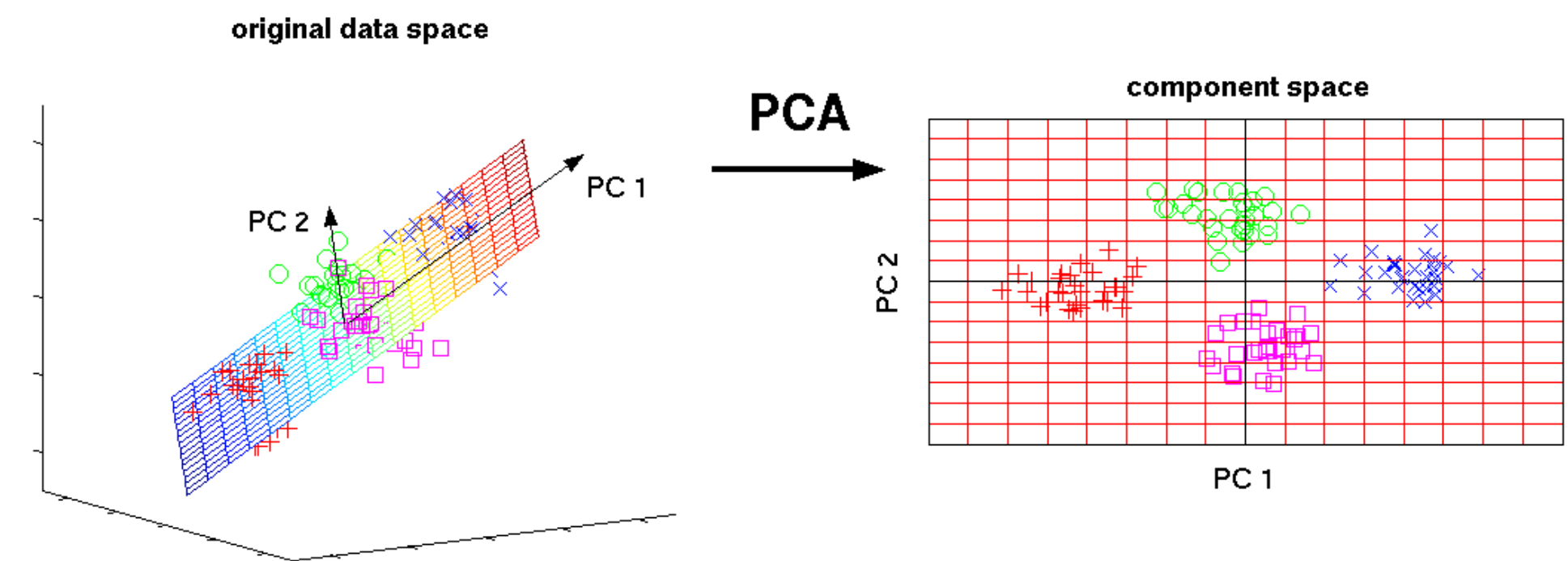
Unsupervised Learning

Data: X

Just data, no labels!

Goal: Learn some underlying hidden *structure* of the data

Examples: Clustering, dimensionality reduction, feature learning, density estimation, *etc.*



dimensionality reduction

[This image is CC0 public domain](#)

Supervised vs. Unsupervised Learning

Unsupervised Learning

Data: X

Just data, no labels!

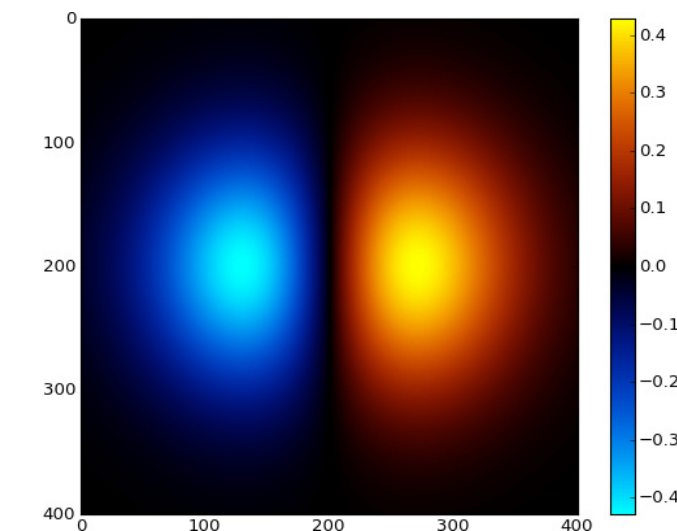
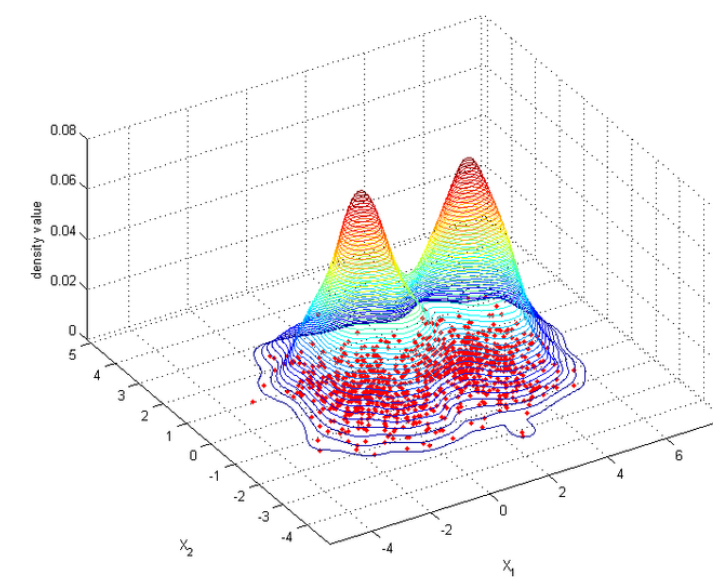
Goal: Learn some underlying hidden *structure* of the data

Examples: Clustering, dimensionality reduction, feature learning, density estimation, *etc.*



Figure copyright Ian Goodfellow, 2016. Reproduced with permission.

1-dim density estimation



2-dim density estimation

2-d density images [left](#) and [right](#) are [CC0 public domain](#)

Supervised vs. Unsupervised Learning

Supervised Learning

Data: (x, y)

x is data, y is label

Goal: Learn a *function* to map $x \rightarrow y$

Examples: Classification, regression, object detection, semantic segmentation, image captioning, *etc.*

Unsupervised Learning

Data: x

Just data, no labels!

Goal: Learn some underlying hidden *structure* of the data

Examples: Clustering, dimensionality reduction, feature learning, density estimation, *etc.*

Generative Models

Given training data, generate new samples from the same distribution



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



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

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

Generative Models

Given training data, generate new samples from the same distribution



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



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

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

Addresses **density estimation**, a core problem in unsupervised learning

- **Explicit** density estimation: explicitly define and solve for $p_{\text{model}}(\mathbf{x})$
- **Implicit** density estimation: learn model that can sample from $p_{\text{model}}(\mathbf{x})$ w/o explicitly defining it

Taxonomy of Generative Models

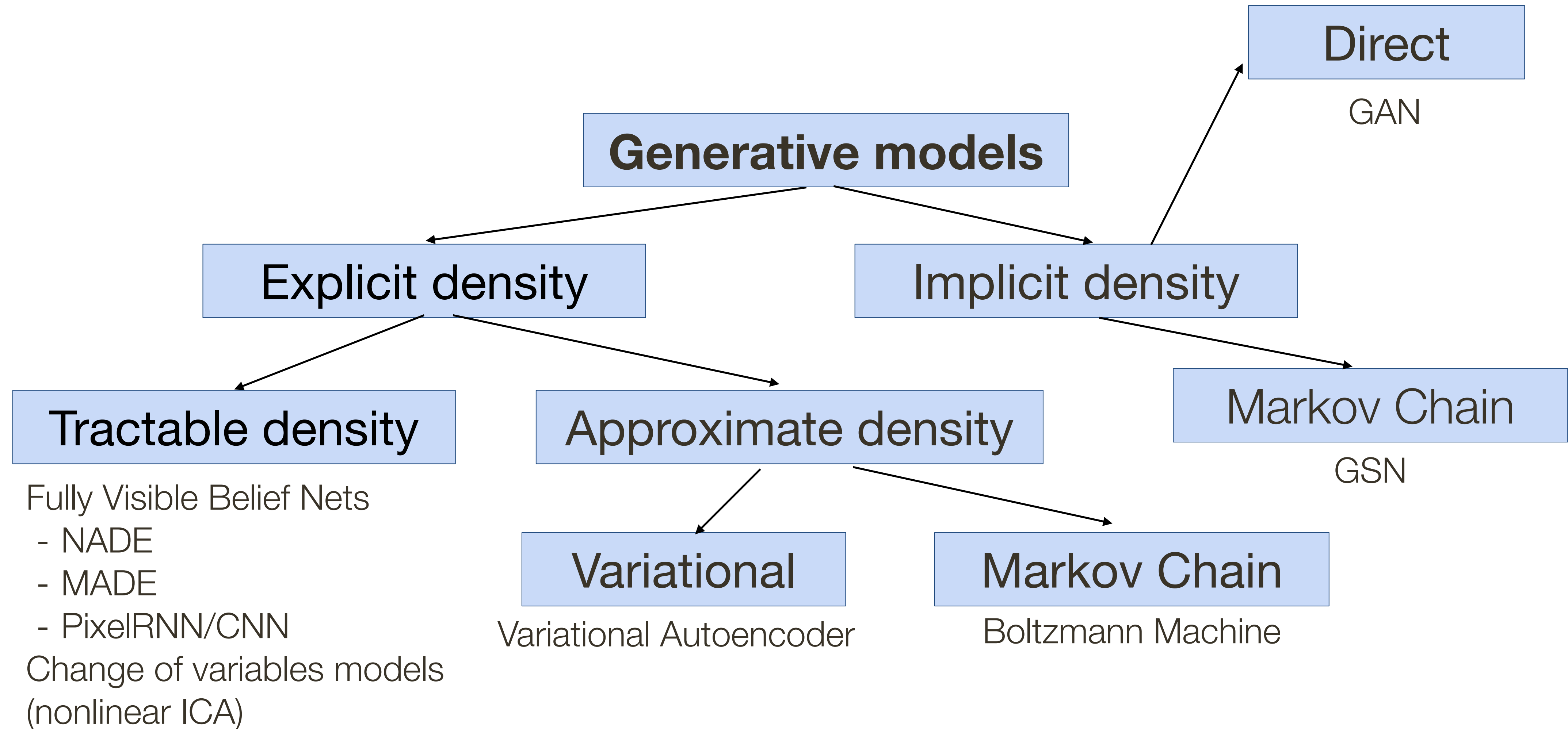


Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

Taxonomy of Generative Models

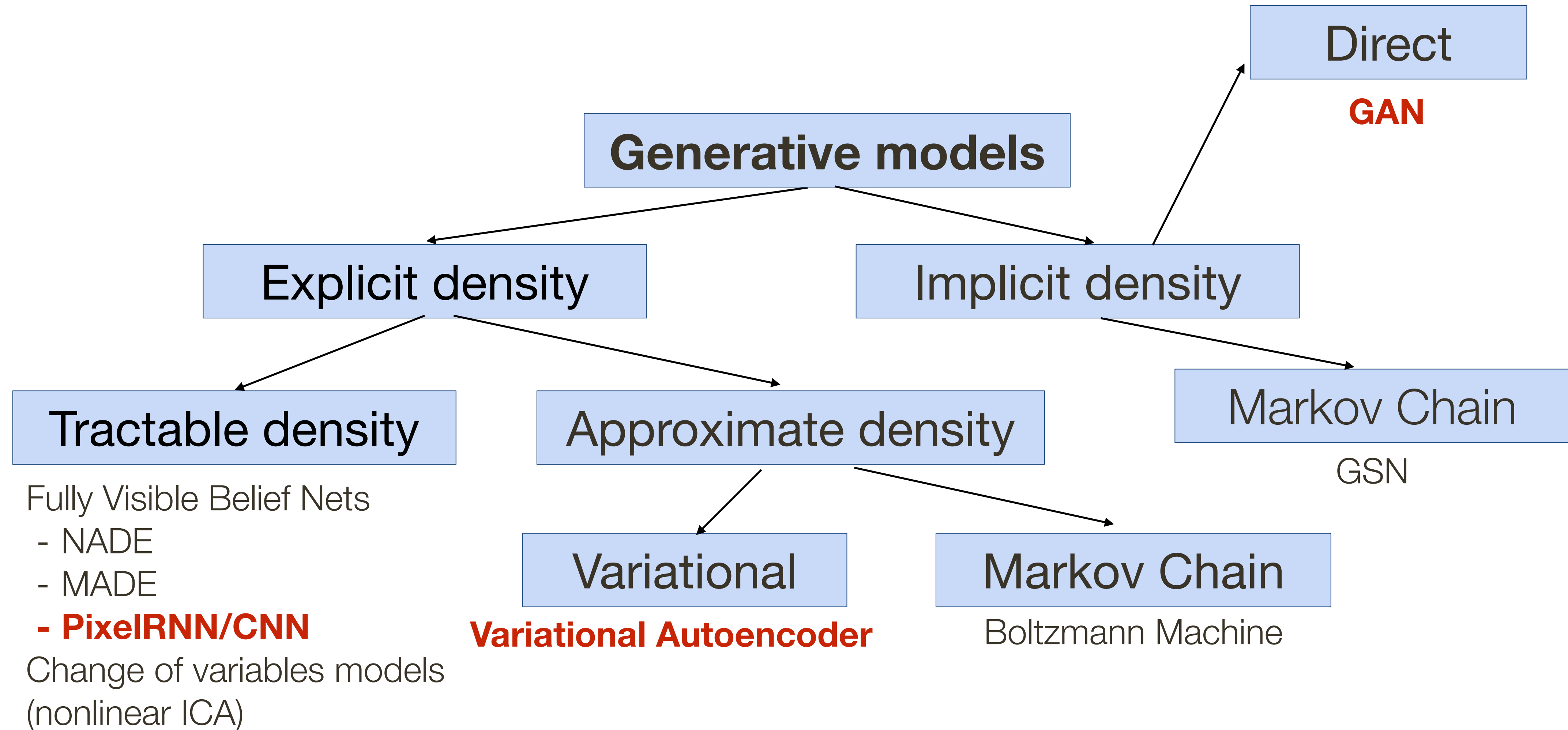
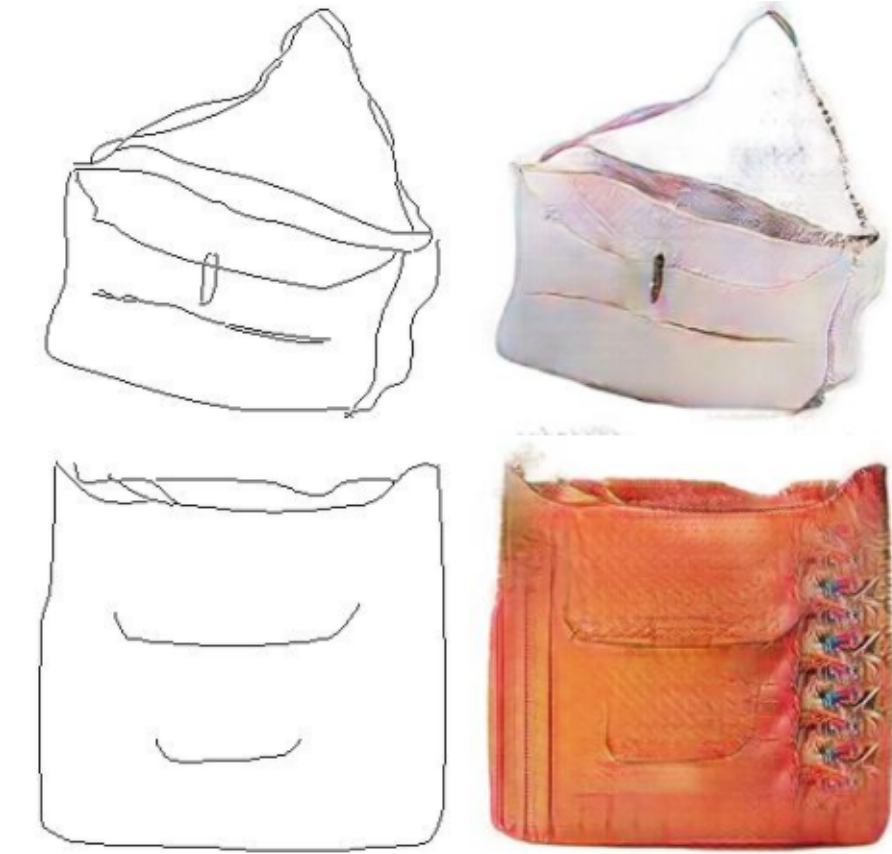
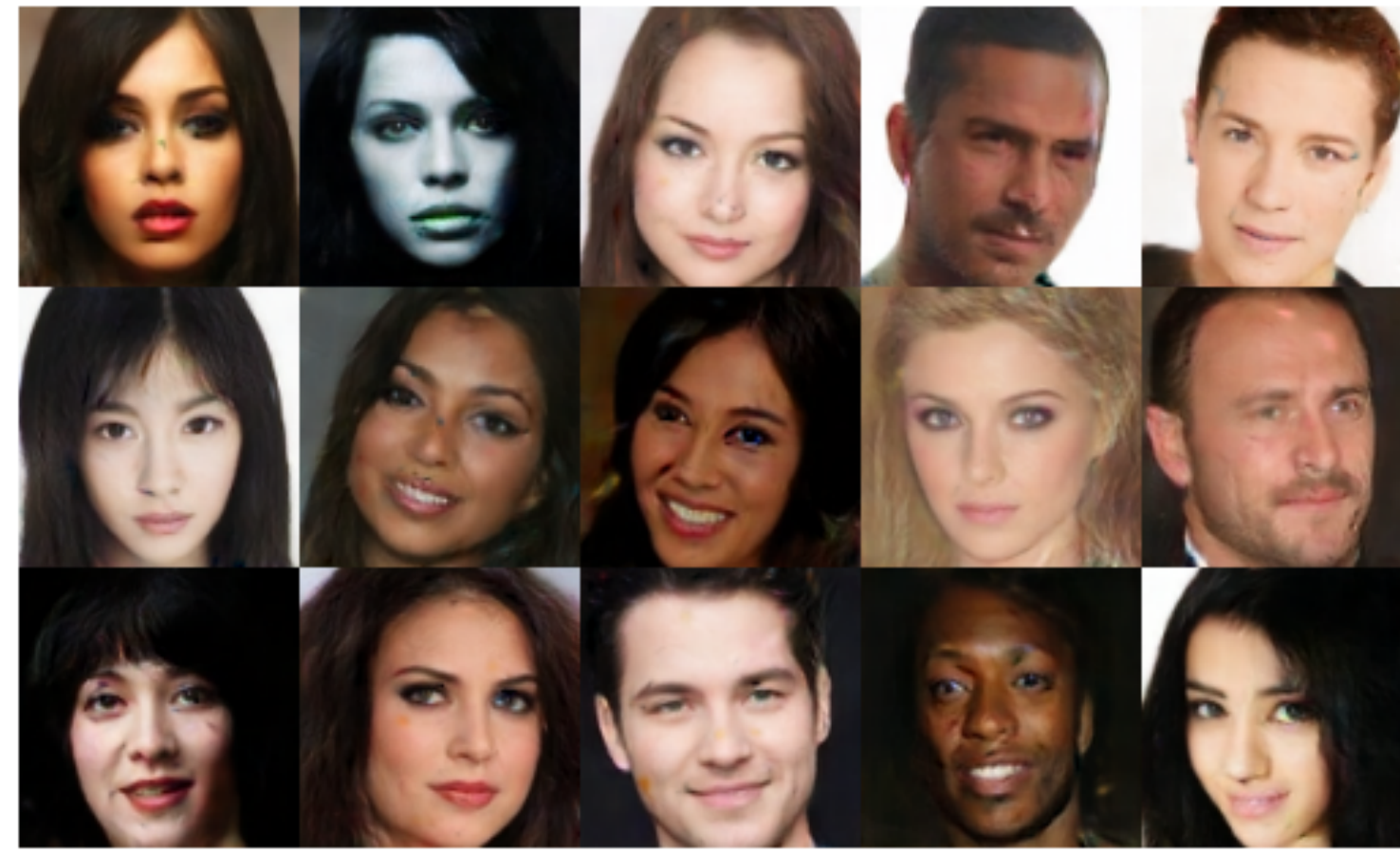


Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

* slide from Fei-Dei Li, Justin Johnson, Serena Yeung, **cs231n Stanford**

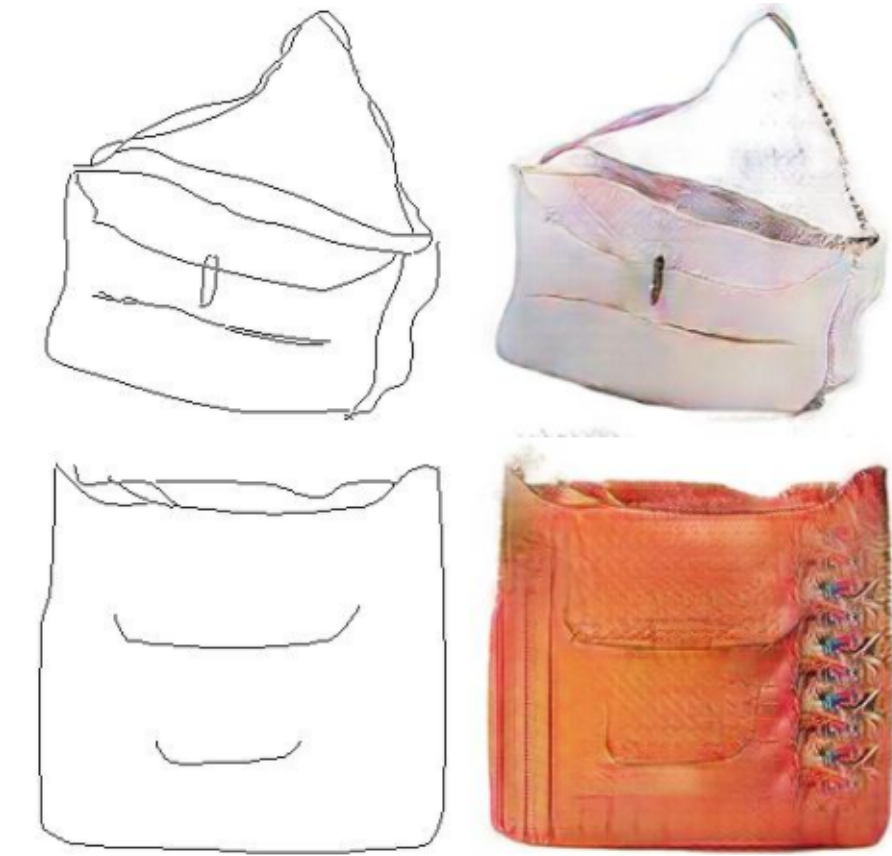
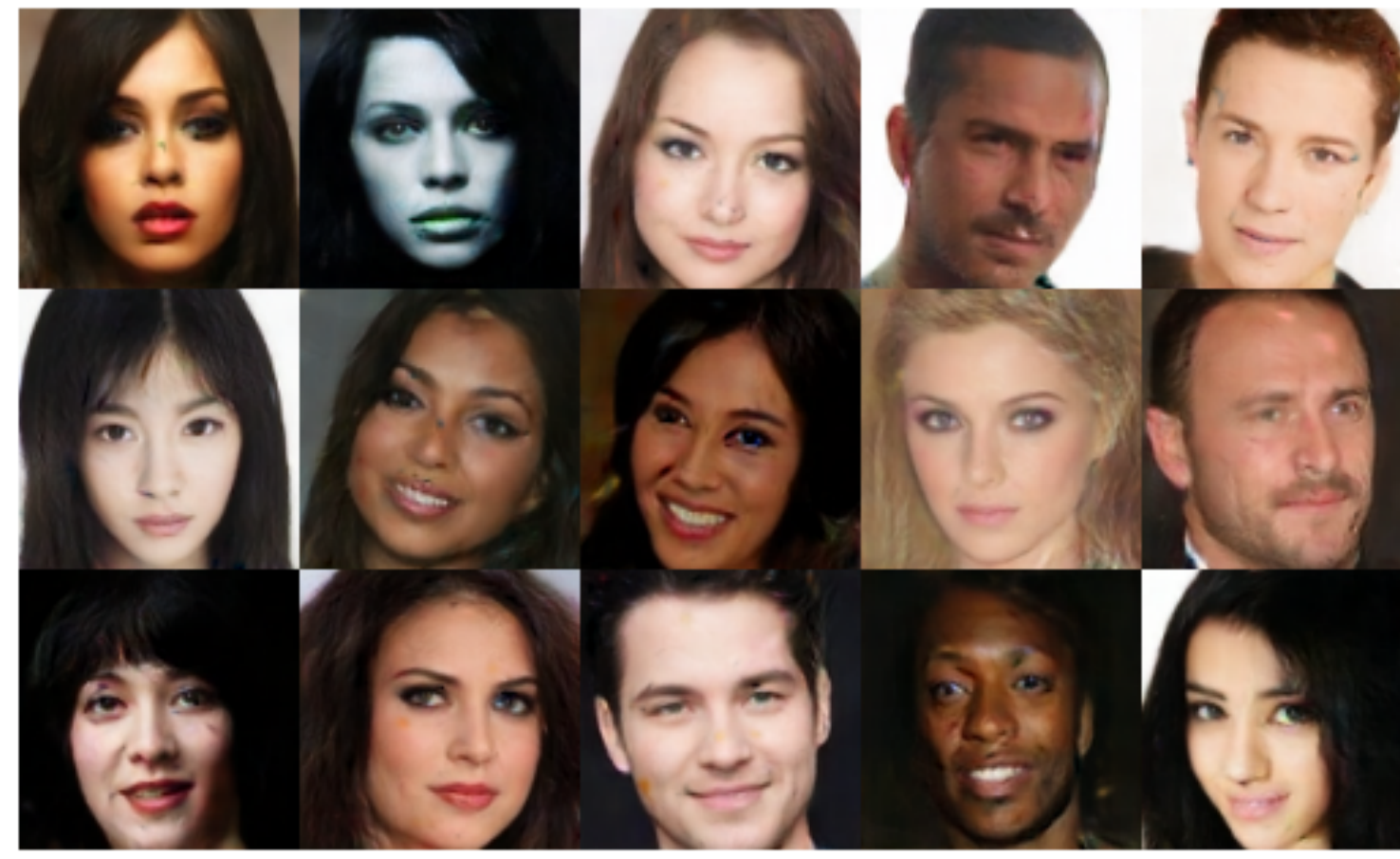
Why **Generative** Models?

- Realistic **samples** for artwork, super-resolution, colorization, *etc.*



Why **Generative** Models?

- Realistic **samples** for artwork, super-resolution, colorization, *etc.*



- Generative models of time-series data can be used for **simulation**, **predictions** and planning (reinforcement learning applications)
- Training generative models can also enable inference of latent representation that can be useful as **general features**
- **Dreaming** / hypothesis visualization

PixelRNN and PixelCNN

Explicit Density model

Use chain rule to decompose likelihood of an image \mathbf{x} into product of (many) 1-d distributions

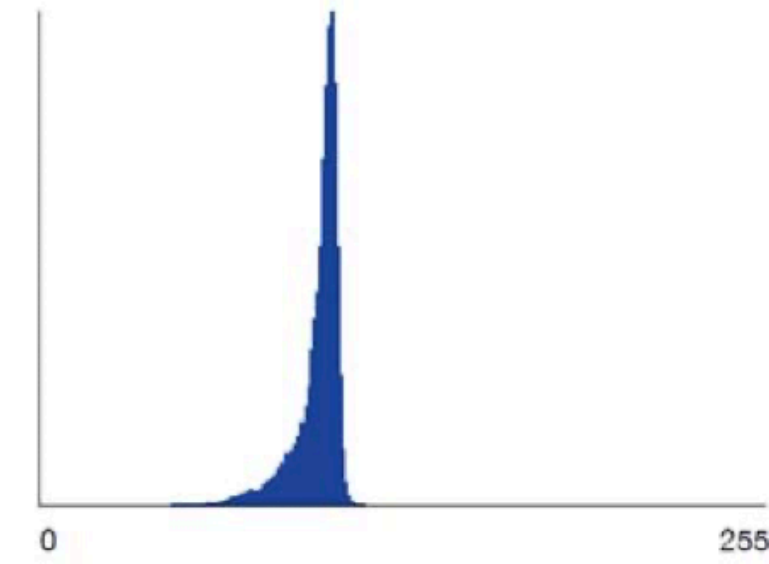
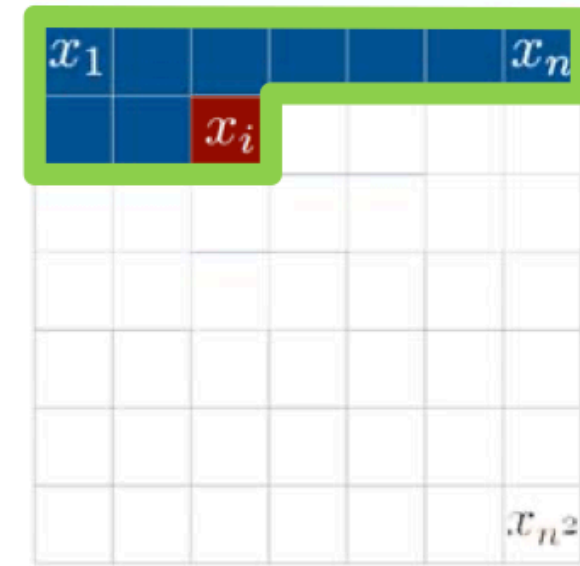
$$p(\mathbf{x}) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1})$$

The diagram illustrates the decomposition of the likelihood of an image \mathbf{x} into a product of 1-d distributions. A green box labeled "Likelihood of image \mathbf{x} " points to $p(\mathbf{x})$. A blue box labeled "Probability of i 'th pixel value given all previous pixels" points to $p(x_i | x_1, \dots, x_{i-1})$.

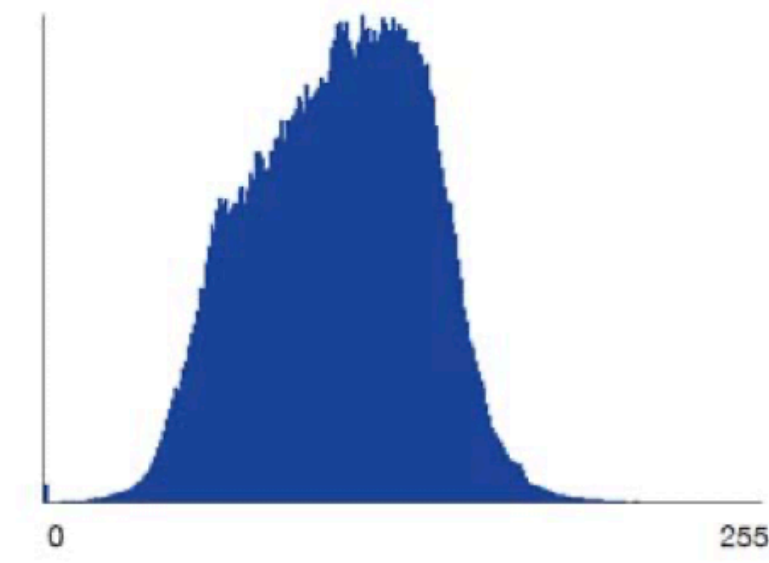
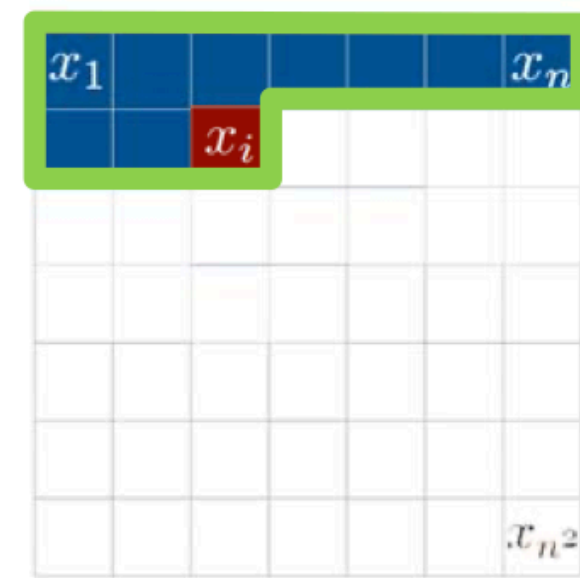
then maximize likelihood of training data

PixelRNN

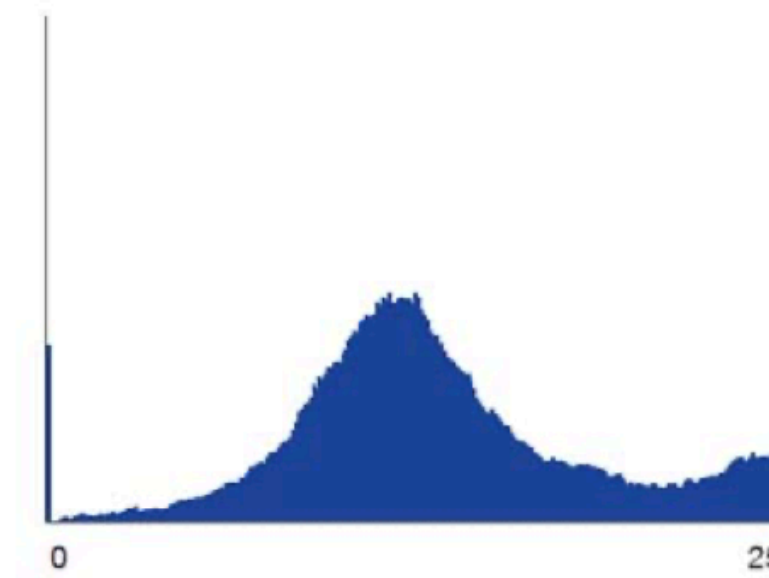
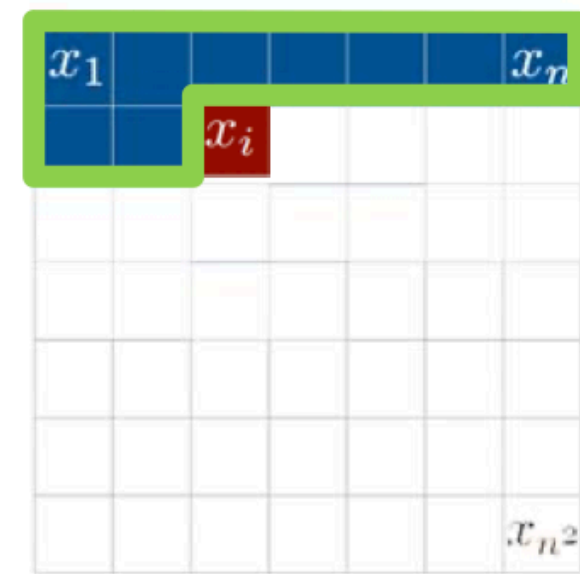
R



G



B



Explicit Density model

Use chain rule to decompose likelihood of an image \mathbf{x} into product of (many) 1-d distributions

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

Likelihood of image \mathbf{x}

Probability of i 'th pixel value given all previous pixels

Complex distribution over pixel values, so lets model using **neural network**

then maximize likelihood of training data

Explicit Density model

Use chain rule to decompose likelihood of an image \mathbf{x} into product of (many) 1-d distributions

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

Likelihood of image \mathbf{x}

Probability of i 'th pixel value given all previous pixels

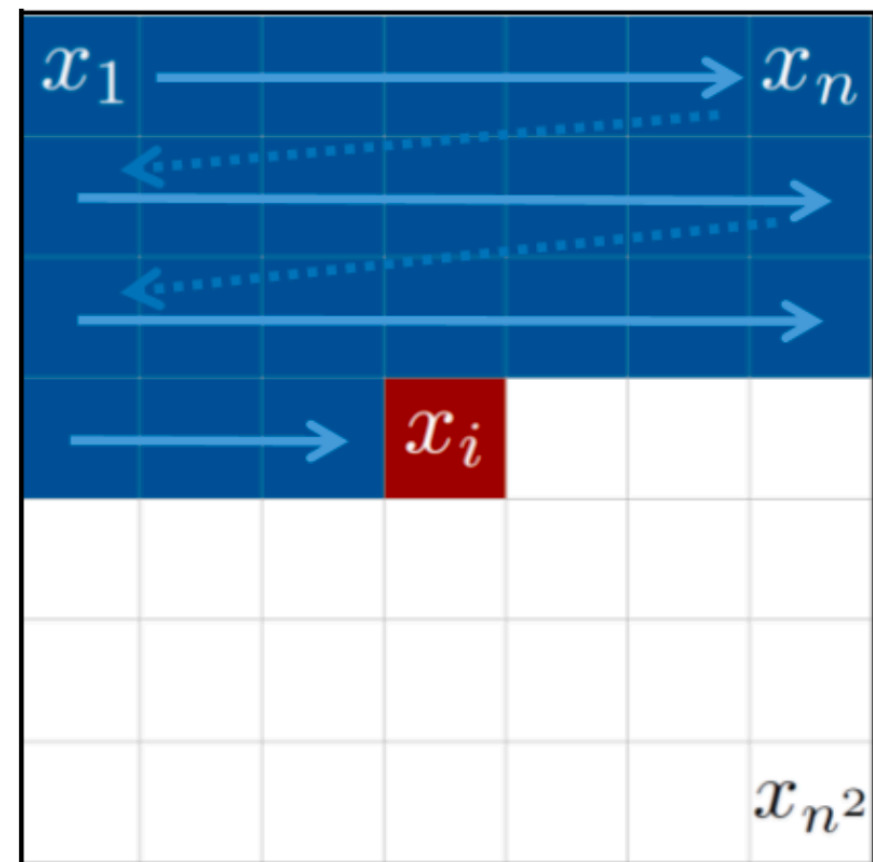
then maximize likelihood of training data

Complex distribution over pixel values, so lets model using **neural network**

Also requires defining **ordering** of “previous pixels”

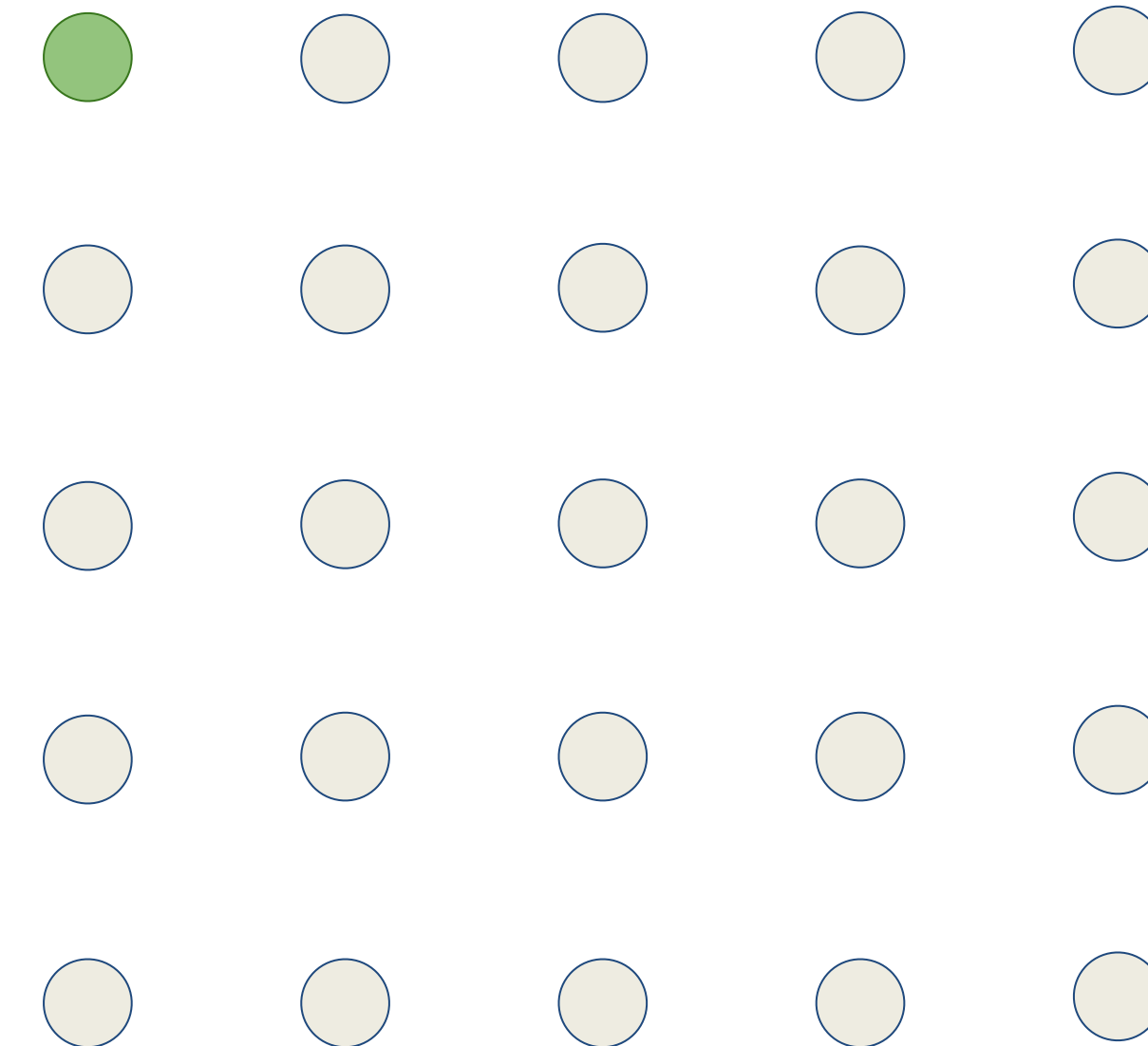
PixelRNN

[van der Oord et al., 2016]



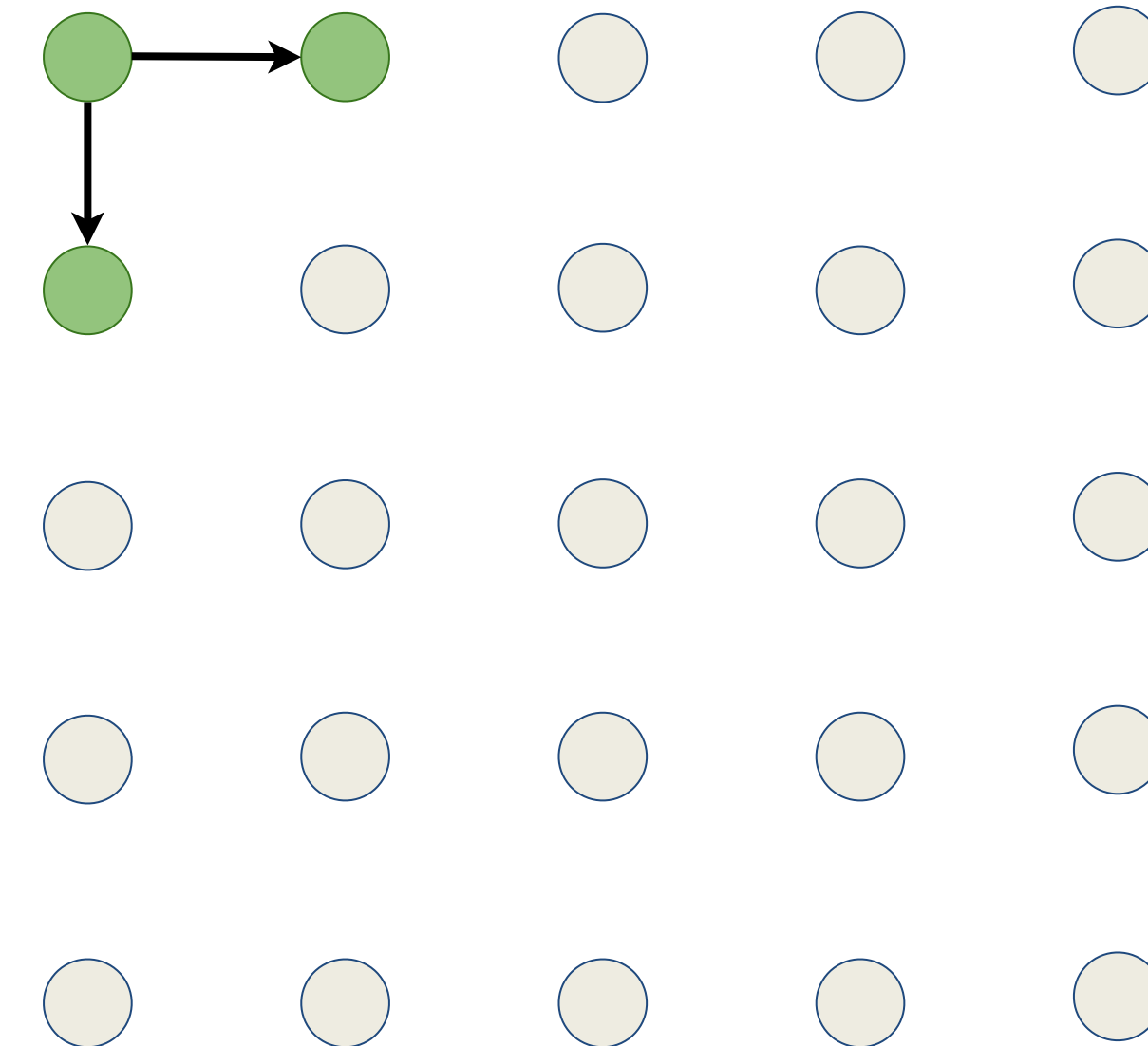
Generate image pixels starting from the corner

Dependency on previous pixels model using an RNN (LSTM)



Generate image pixels starting from the corner

Dependency on previous pixels model using an RNN (LSTM)

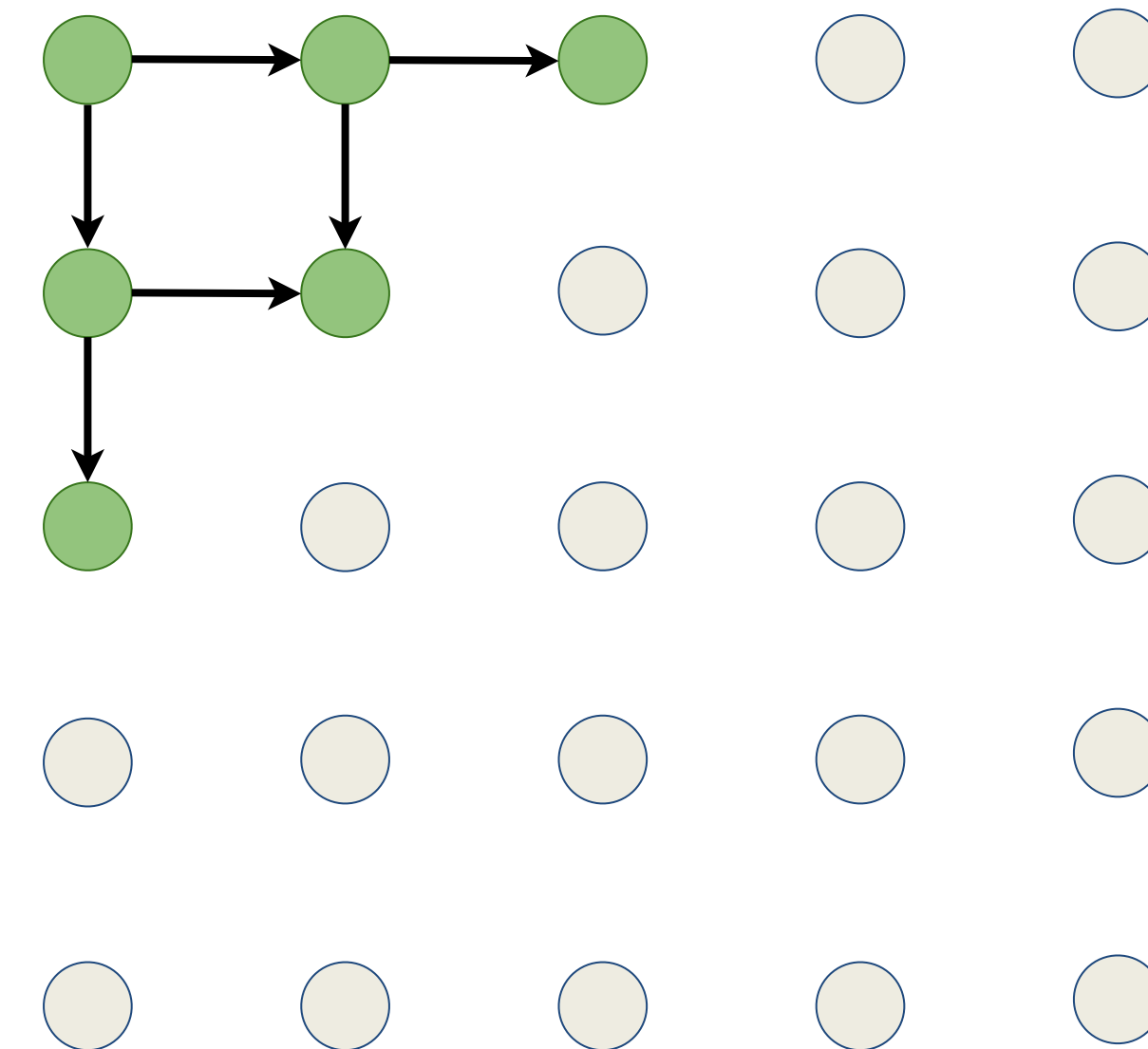


PixelRNN

[van der Oord et al., 2016]

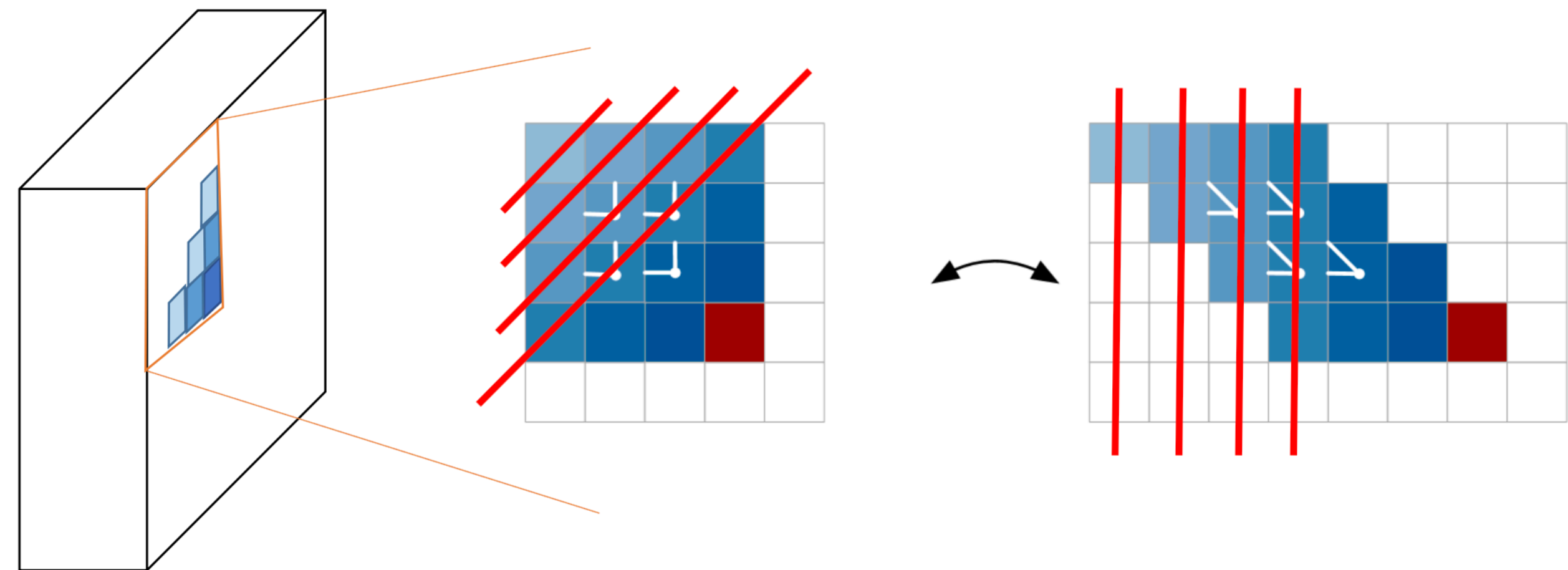
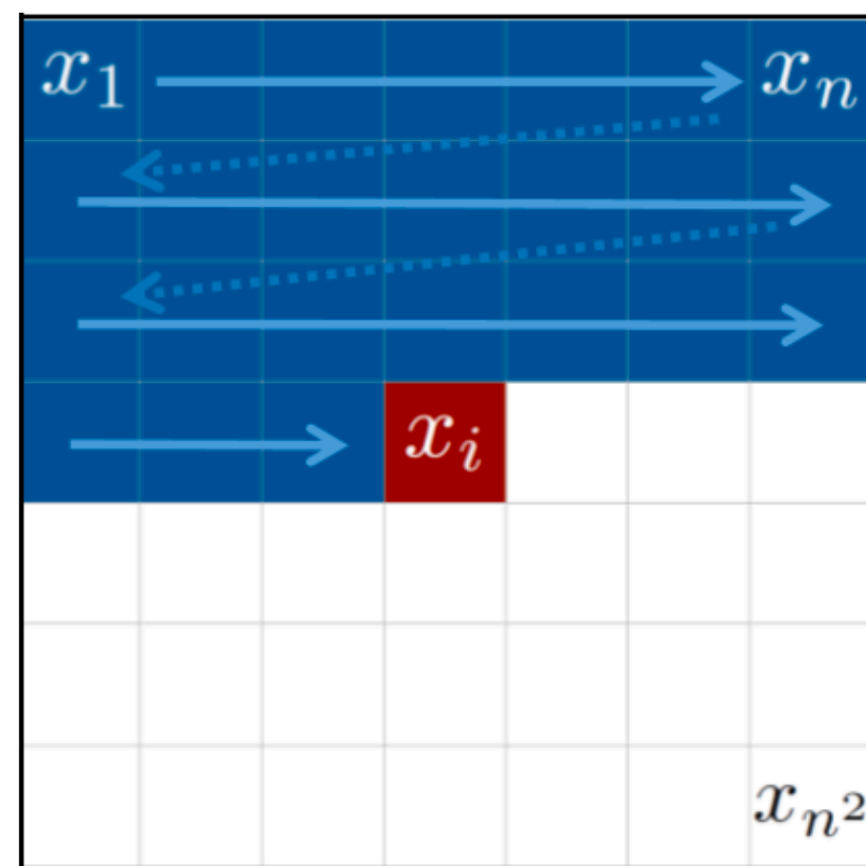
Generate image pixels starting from the corner

Dependency on previous pixels model using an RNN (LSTM)



PixelRNN

[van der Oord et al., 2016]

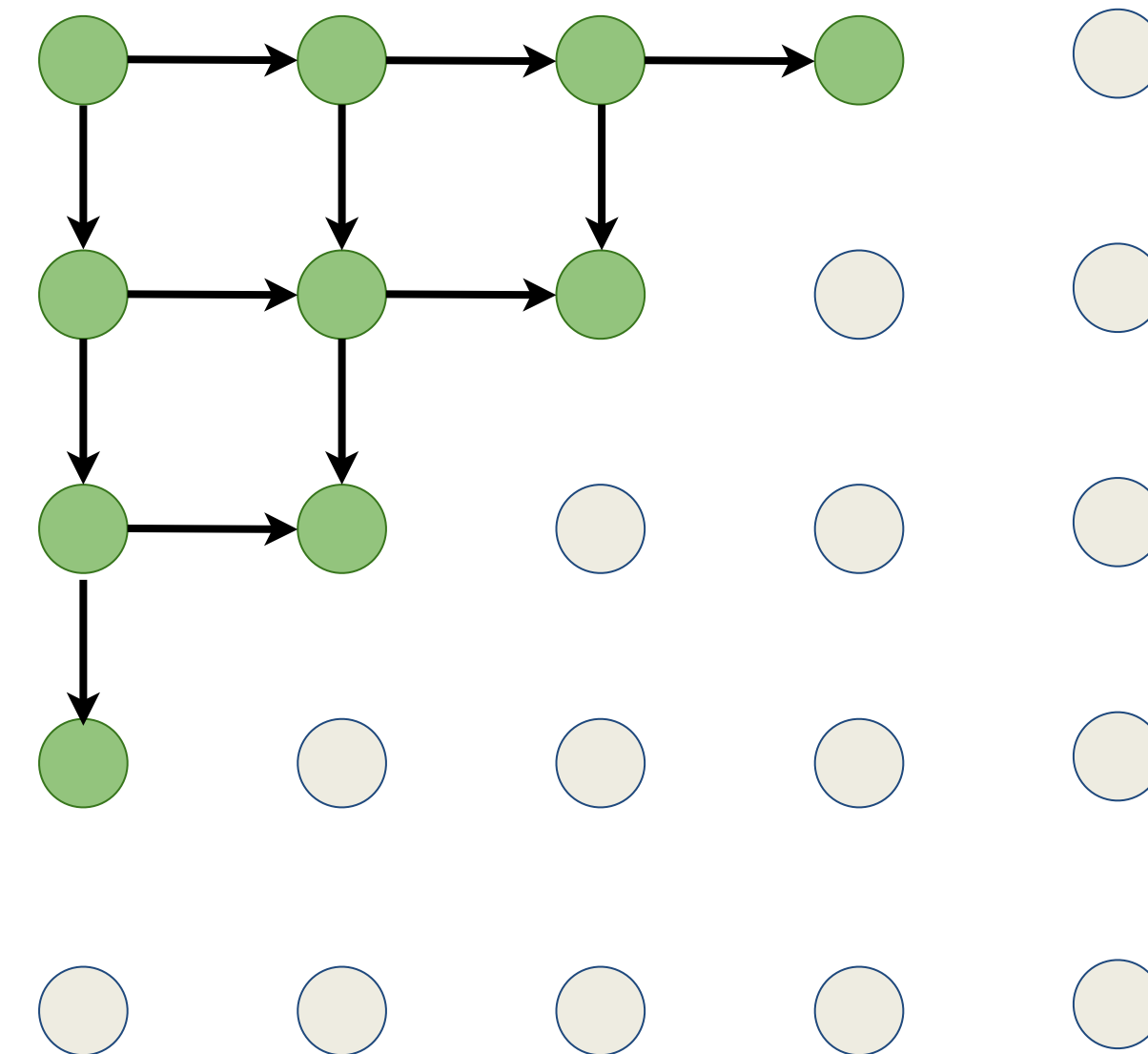


PixelRNN

[van der Oord et al., 2016]

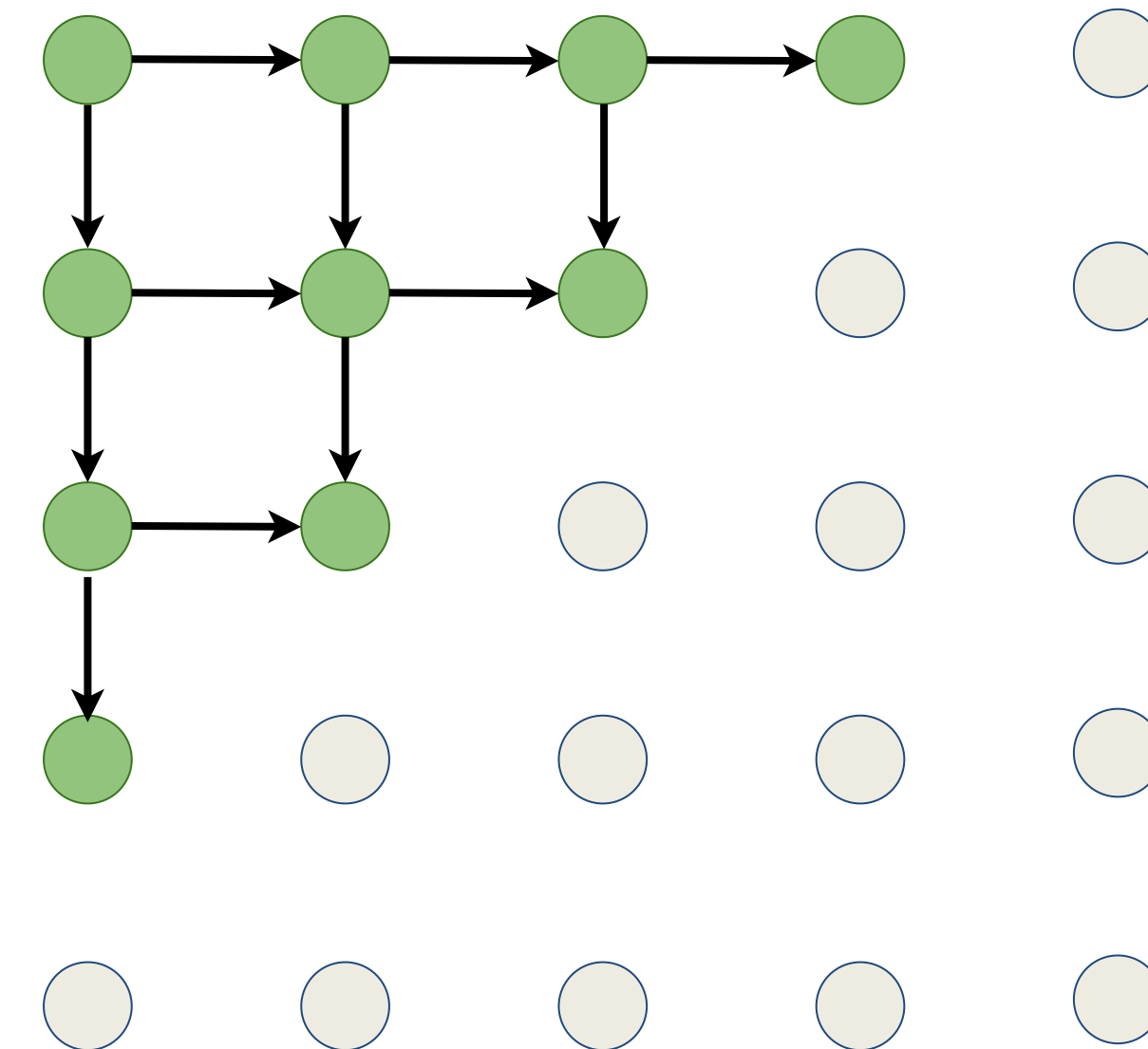
Generate image pixels starting from the corner

Dependency on previous pixels
model using an RNN (LSTM)



Generate image pixels starting from the corner

Dependency on previous pixels
model using an RNN (LSTM)



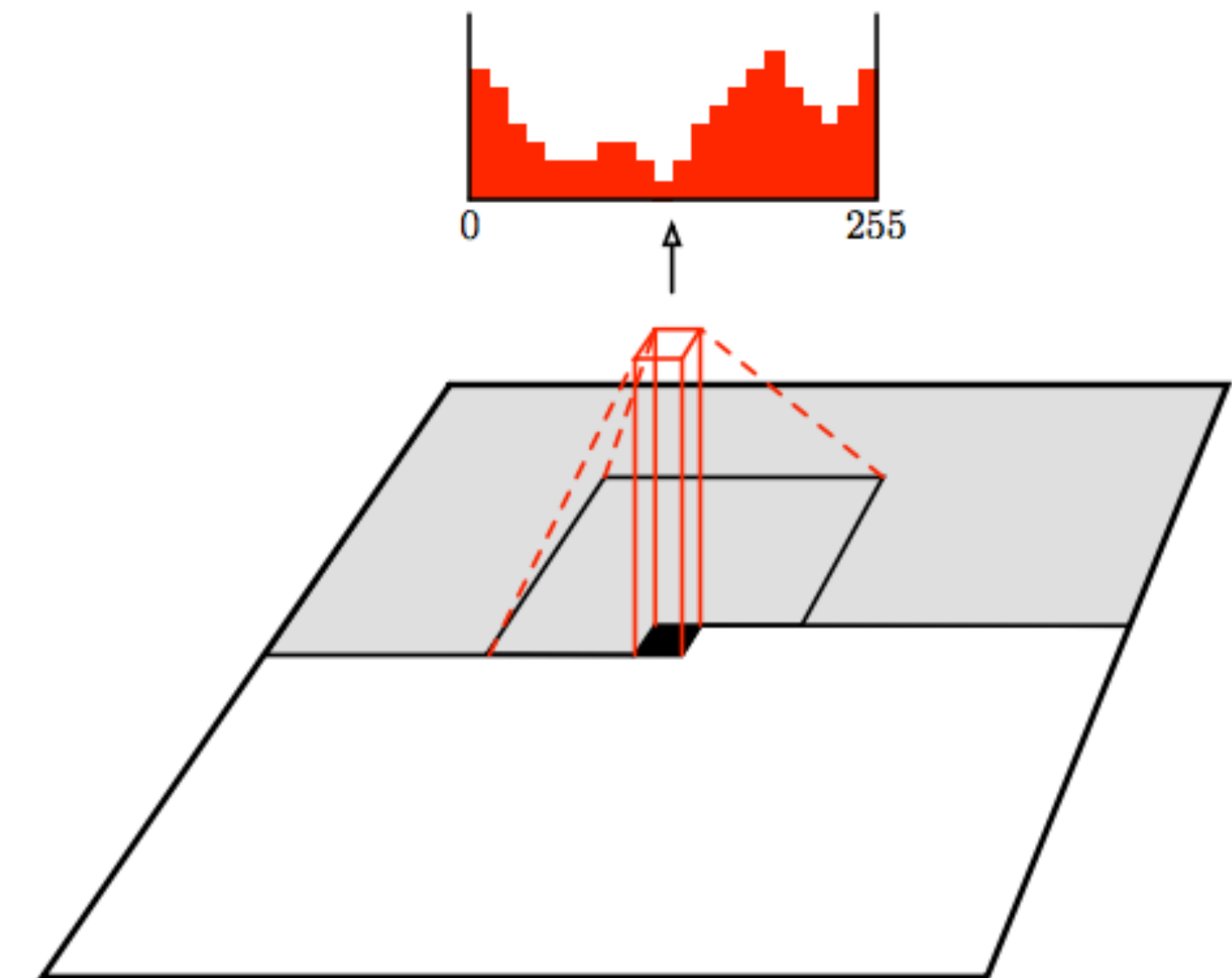
Problem: sequential generation is slow

PixelCNN

[van der Oord et al., 2016]

Still generate image pixels
starting from the corner

Dependency on previous pixels
now modeled using a CNN over
context region



PixelCNN

[van der Oord et al., 2016]

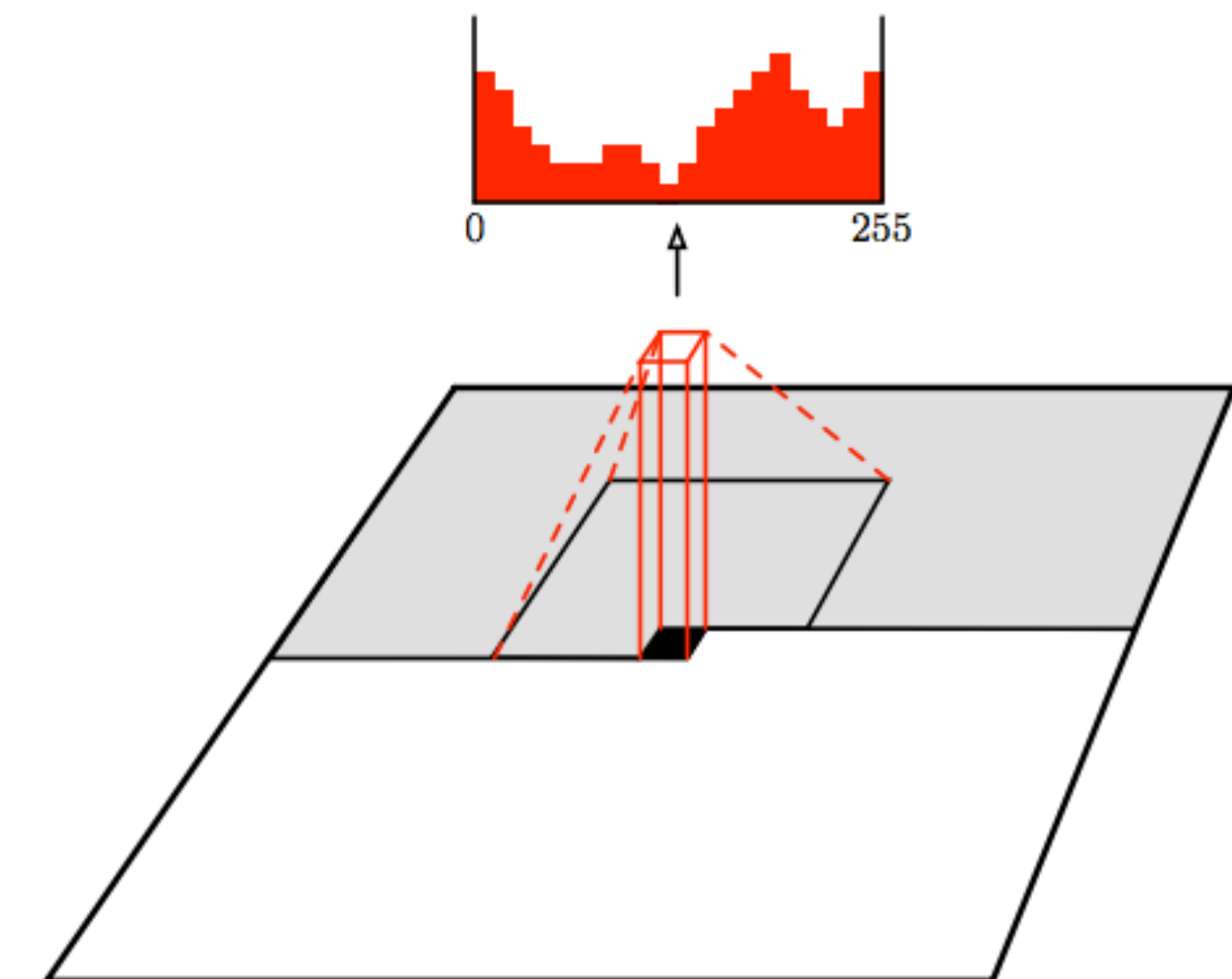
Still generate image pixels starting from the corner

Dependency on previous pixels now modeled using a CNN over context region

Training: maximize likelihood of training images

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

Softmax loss at each pixel



PixelCNN

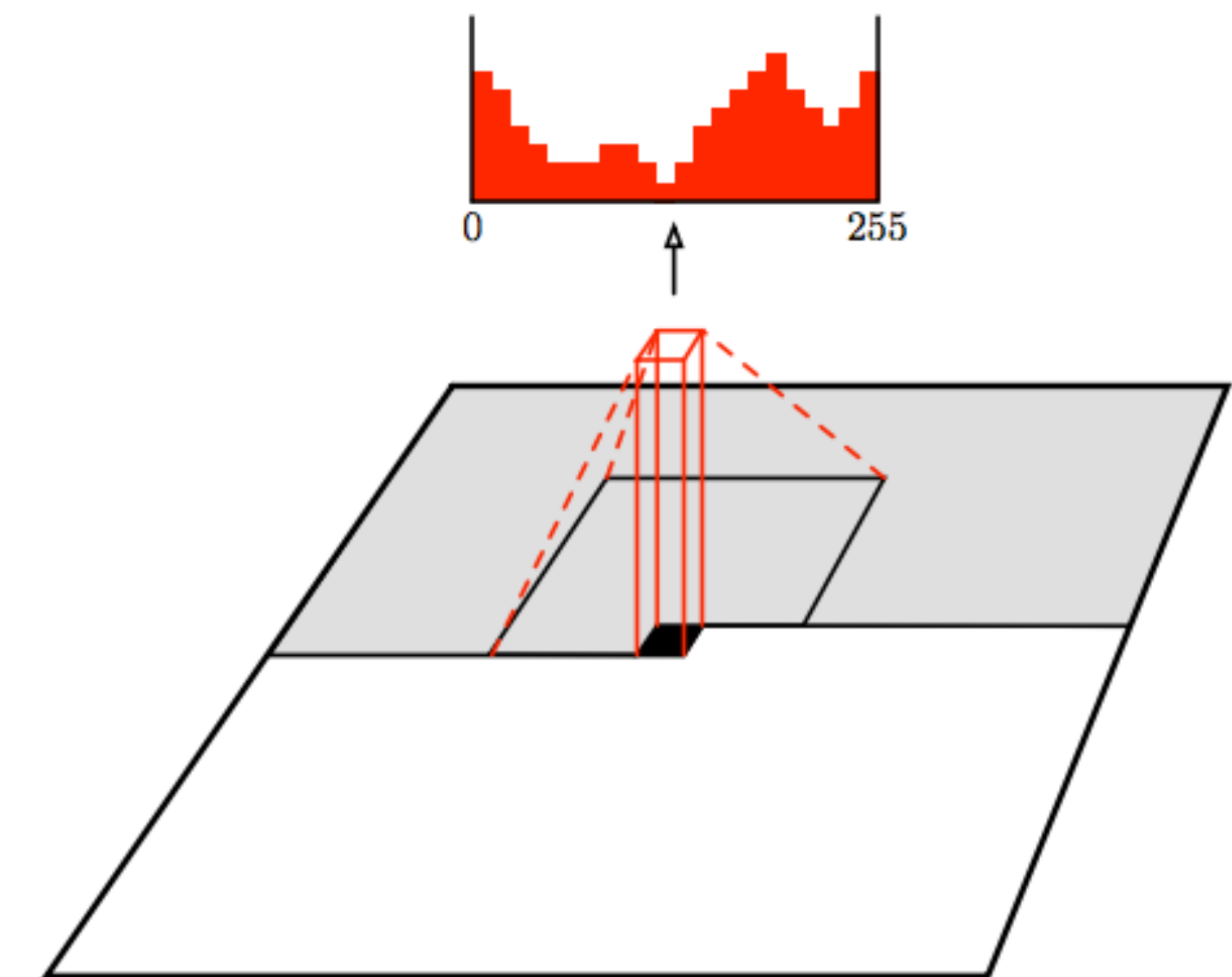
[van der Oord et al., 2016]

Still generate image pixels starting from the corner

Dependency on previous pixels now modeled using a CNN over context region

Training: maximize likelihood of training images

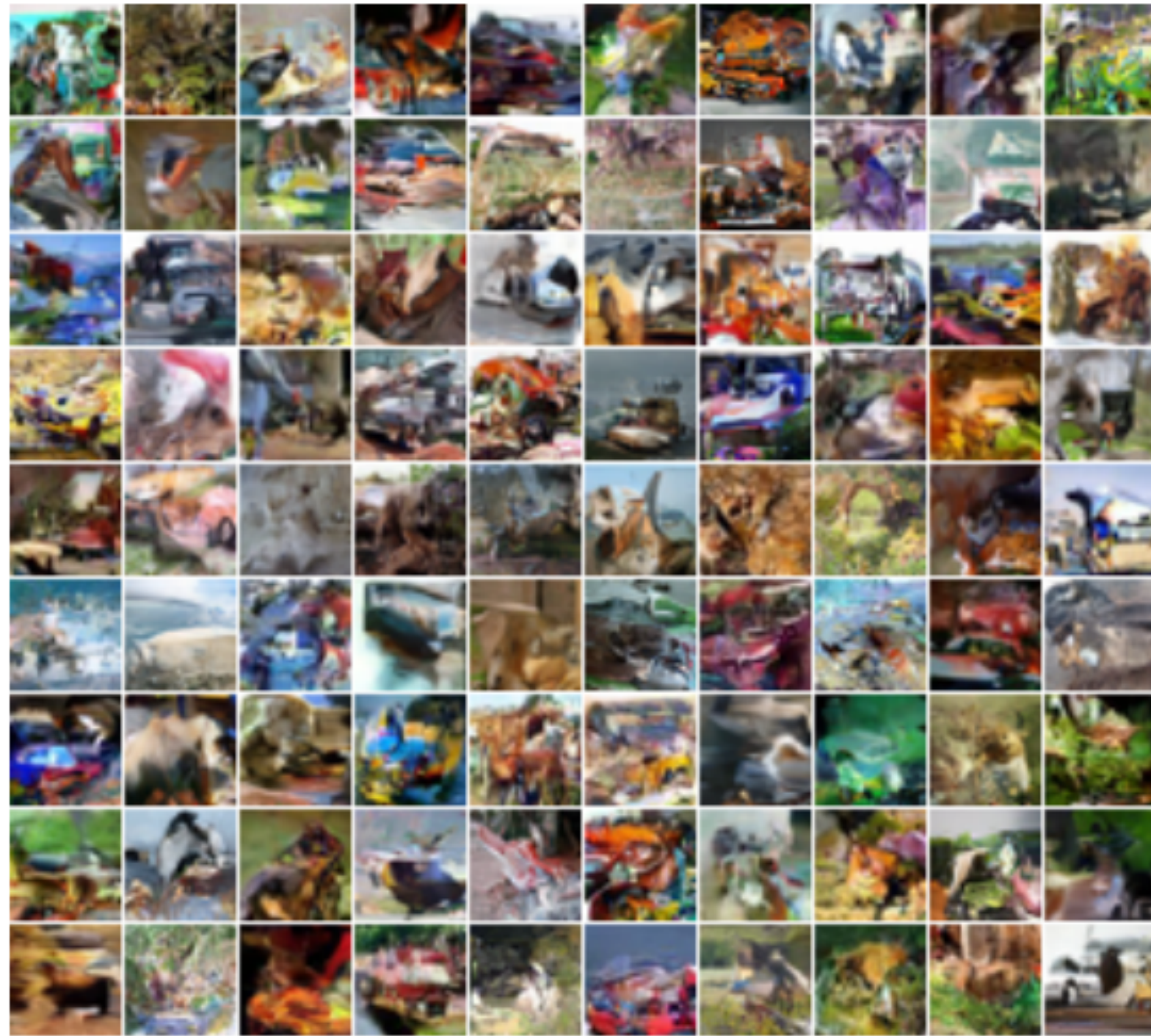
$$p(\mathbf{x}) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1})$$



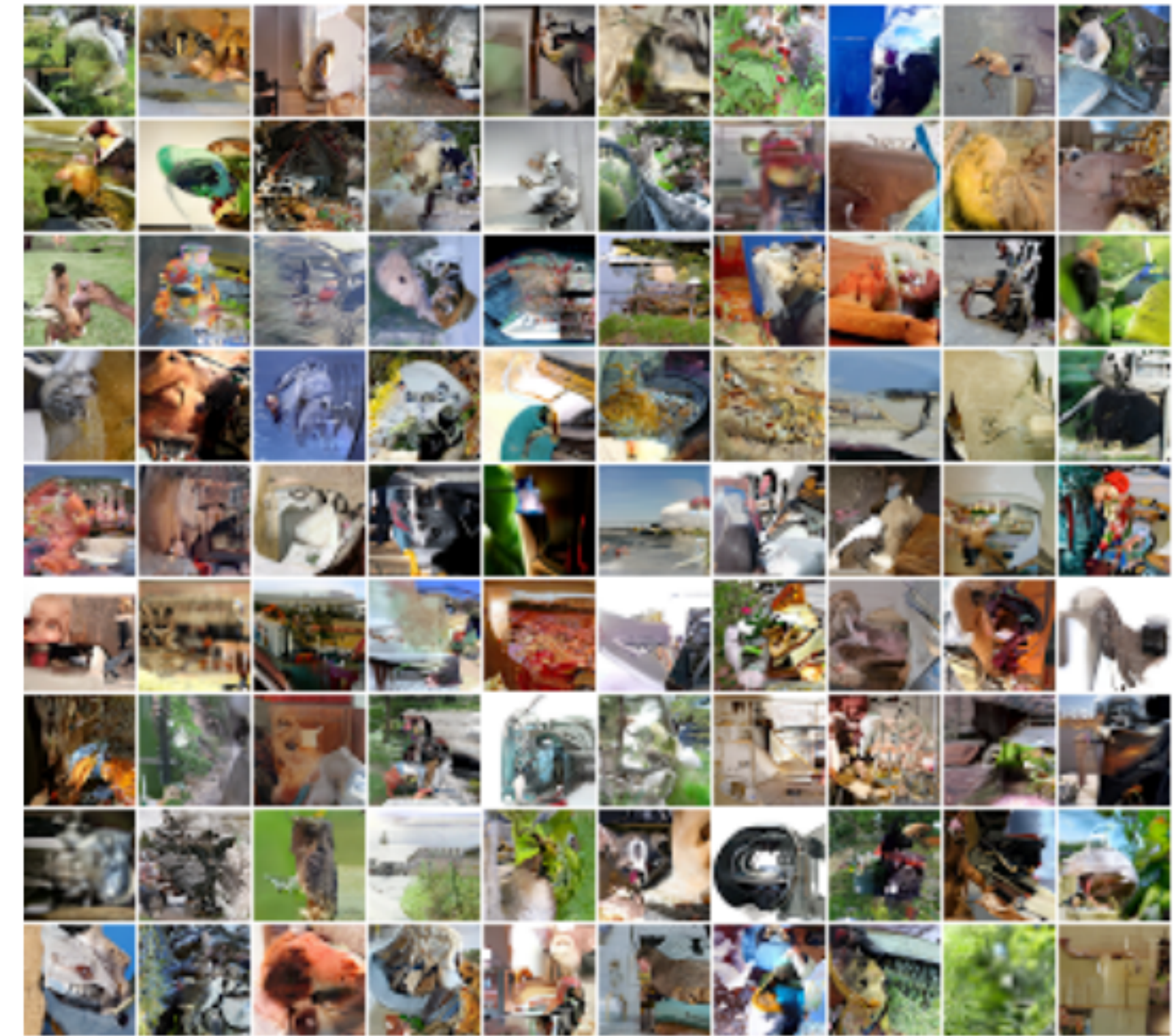
Generation is still slow (sequential),
but learning is faster

Generated Samples

[van der Oord et al., 2016]



32x32 **CIFAR-10**



32x32 **ImageNet**

PixelRNN and PixelCNN

Pros:

- Can explicitly compute likelihood $p(x)$
- Explicit likelihood of training data gives good evaluation metric
- Good samples

Con:

- Sequential generation => slow

Improving PixelCNN performance

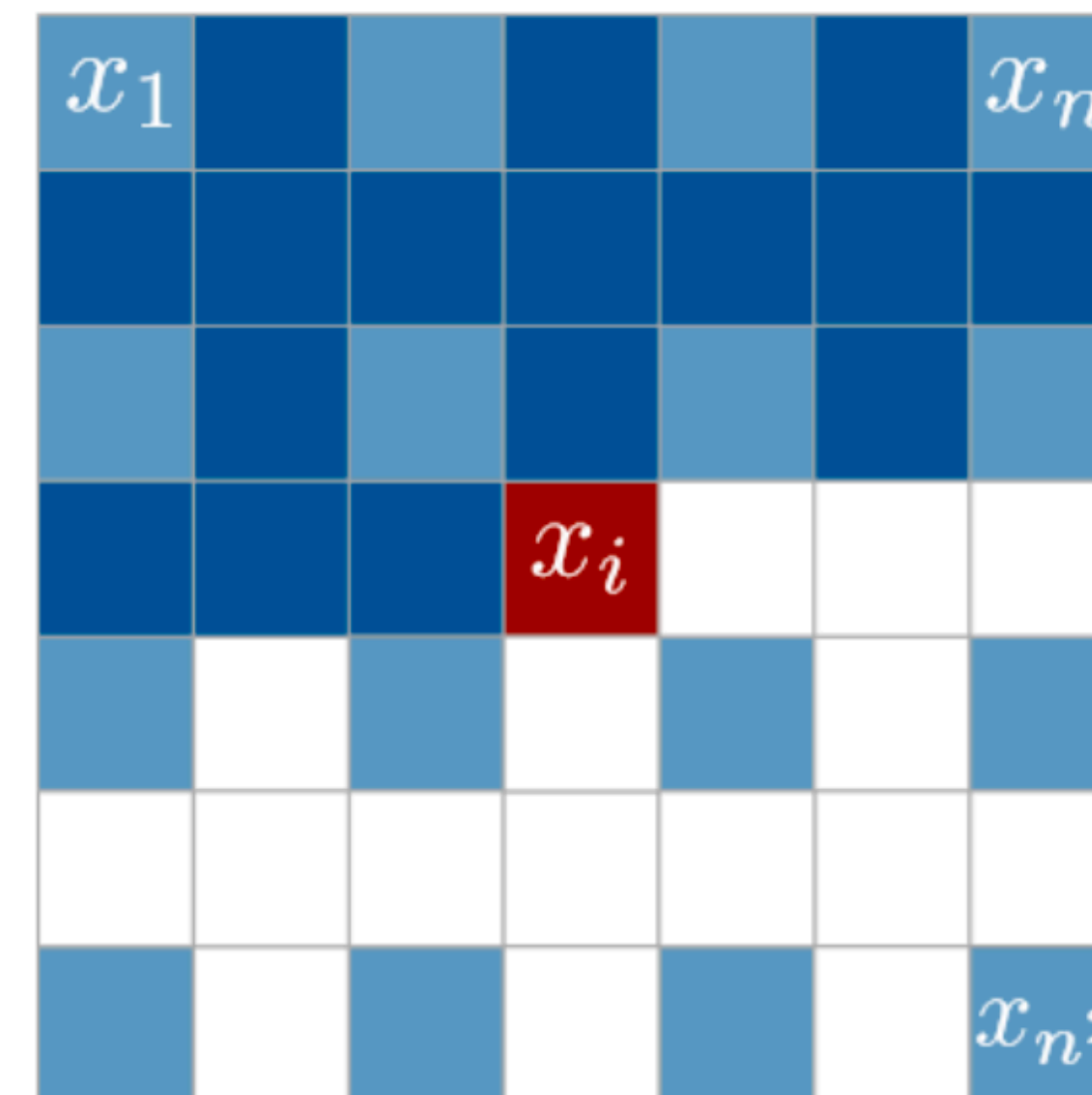
- Gated convolutional layers
- Short-cut connections
- Discretized logistic loss
- Multi-scale
- Training tricks
- Etc...

Multi-scale PixelRNN

[van der Oord et al., 2016]

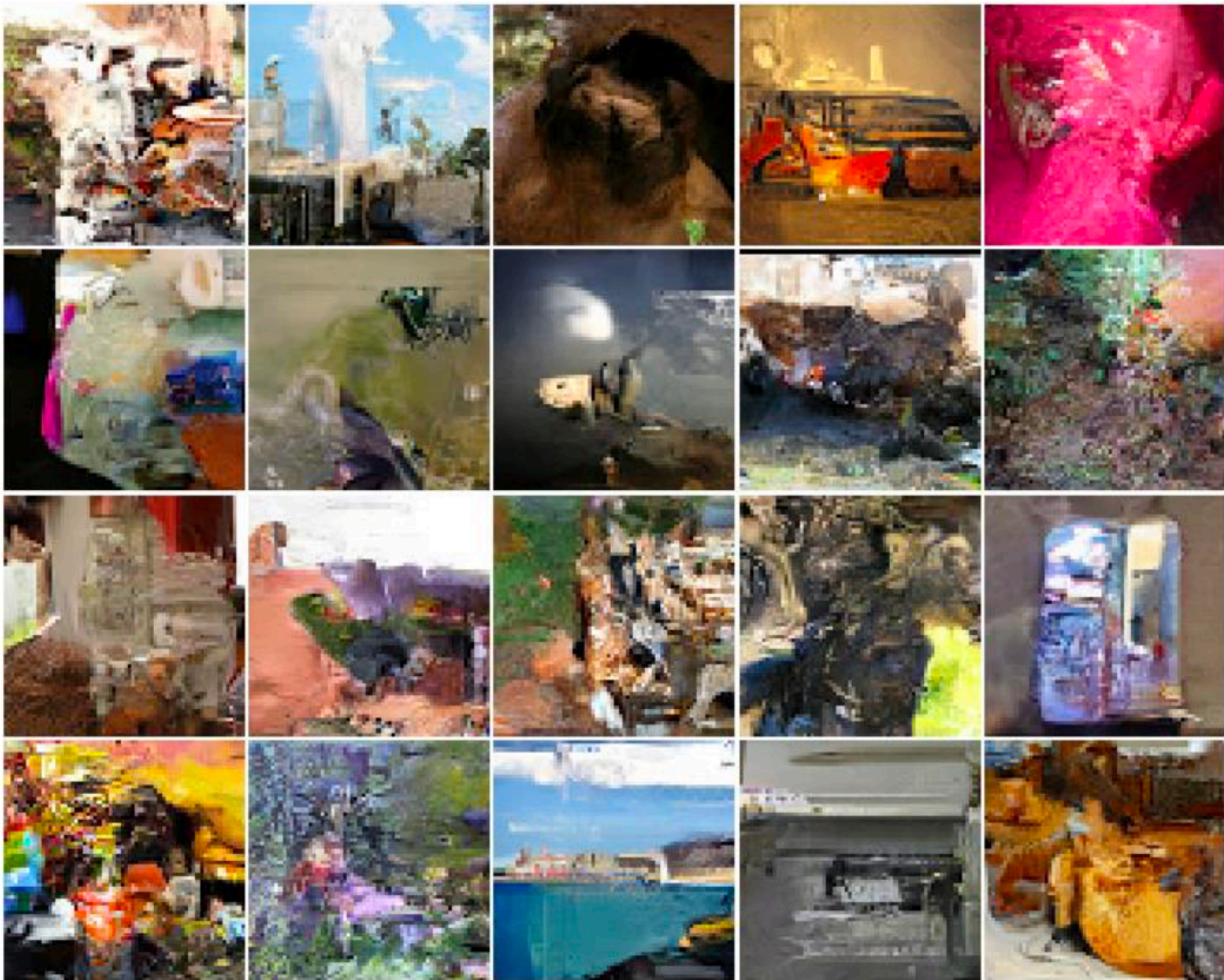
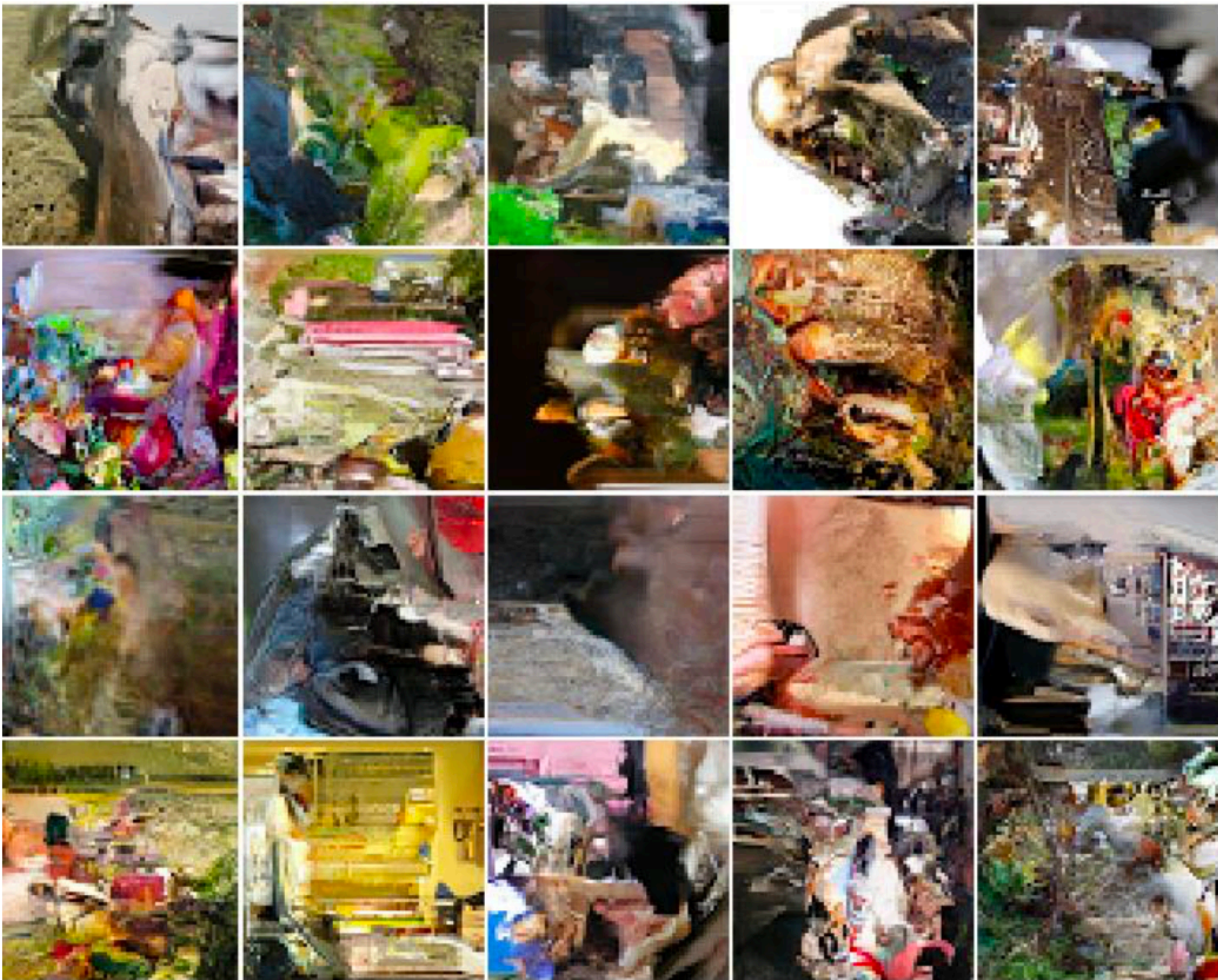
Take sub-sampled pixels as additional input pixels

Can capture better global information (more visually coherent)



Multi-scale PixelRNN

[van der Oord et al., 2016]



* slide from Hsiao-Ching Chang, Ameya Patil, Anand Bhattad

Conditional Image Generation

[van der Oord et al., 2016]

Similar to PixelRNN/CNN but conditioned on a high-level image description vector \mathbf{h}

$$p(\mathbf{x}) = p(x_1, x_2, \dots, x_{n^2})$$



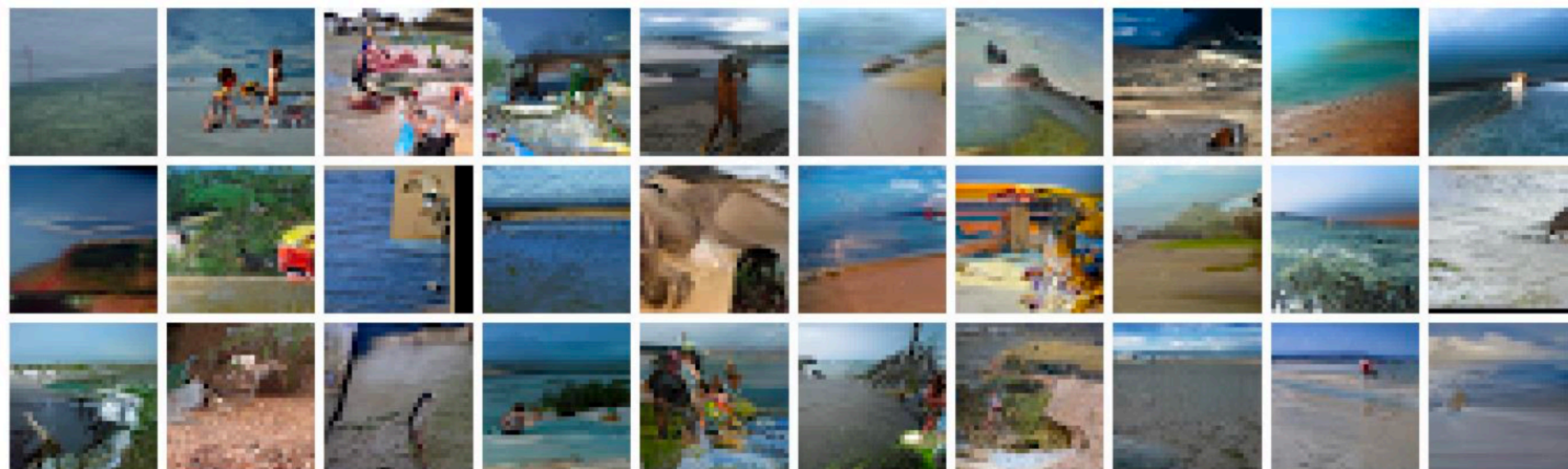
$$p(\mathbf{x}|\mathbf{h}) = p(x_1, x_2, \dots, x_{n^2}|\mathbf{h})$$

Conditional Image Generation

[van der Oord et al., 2016]



African elephant



Sandbar

* slide from Hsiao-Ching Chang, Ameya Patil, Anand Bhattad