

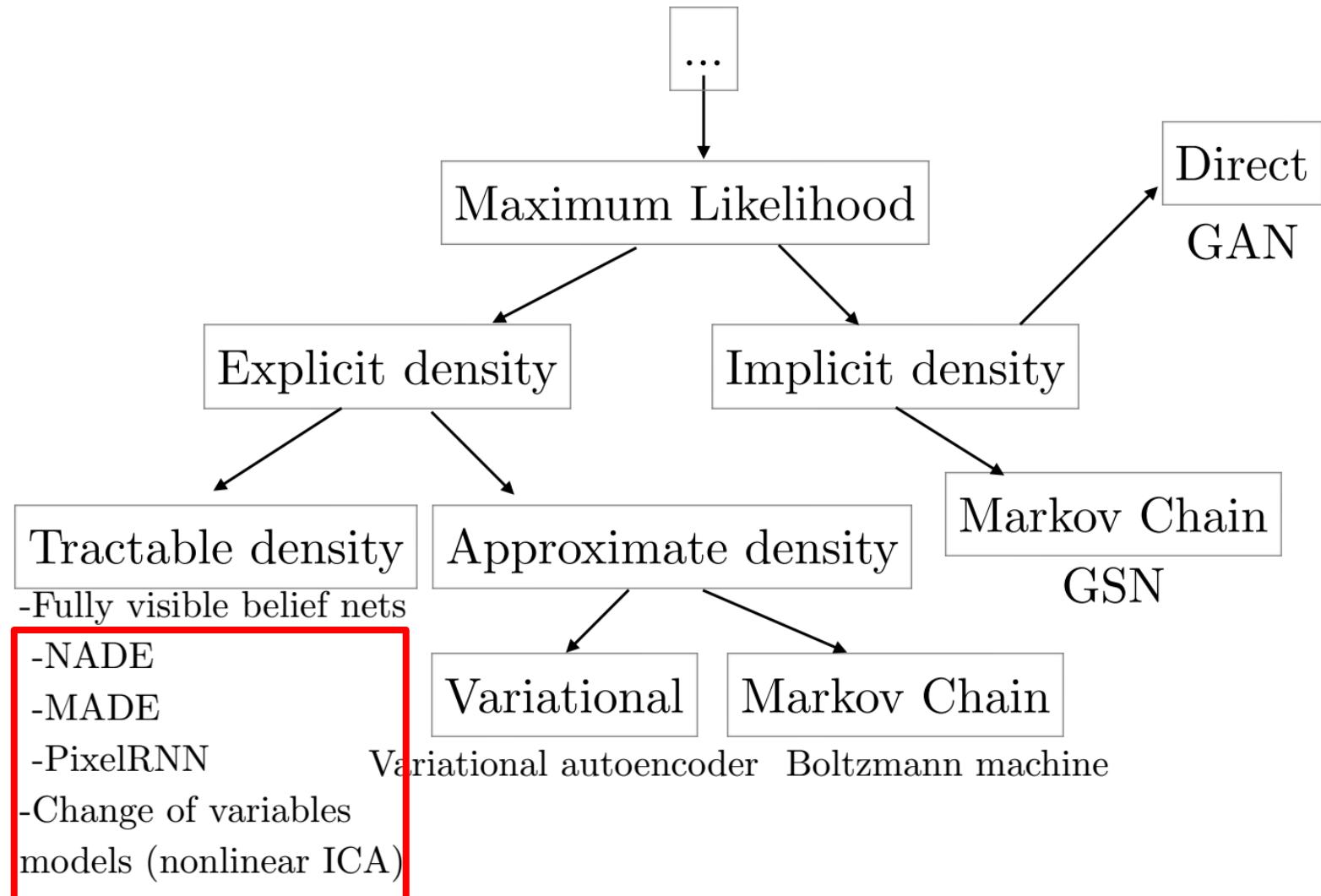
Autoregressive and Flow-based Generative Models

Efstratios Gavves

Lecture overview

- Early autoregressive models
- Modern autoregressive models
- Normalizing flows
- Flow-based models

A map of generative models



Beyond independent dimensions

- Often, in data there is either an order or we can make up an order
 - From a generation point of view, data dimensions depend on each other

Decomposing likelihood of sequential data

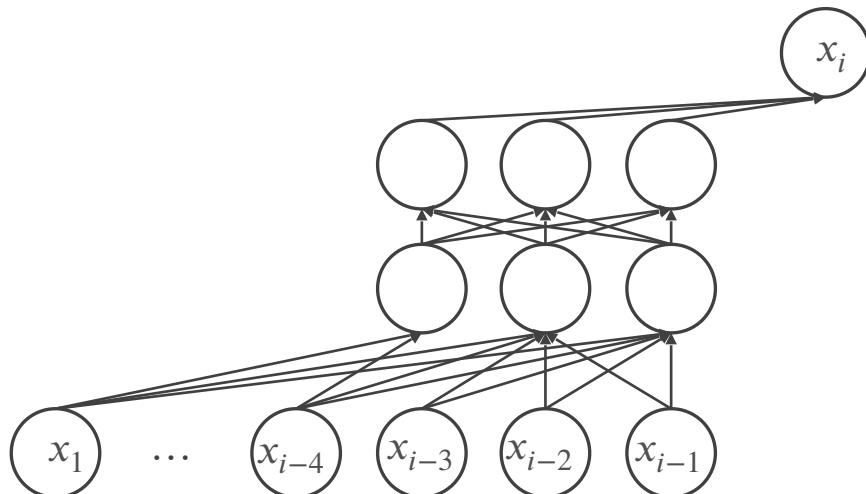
- If $\mathbf{x} = [x_1, x_2, \dots, x_d]$ is sequential, $p(\mathbf{x})$ decomposes with chain rule of probabilities

$$p(\mathbf{x}) = p(x_1) \cdot p(x_2 | x_1) \cdot p(x_3 | x_1, x_2) \cdot \dots \cdot p(x_d | x_1, \dots, x_{d-1}) = \prod_{i=1}^d p(x_i | x_{<i})$$

- If \mathbf{x} is *not* sequential, we can assume an artificial order
 - *e.g.*, the order with which pixels make (generate) an image
 - This can create artificial bias, however

Deep networks to model conditional likelihoods

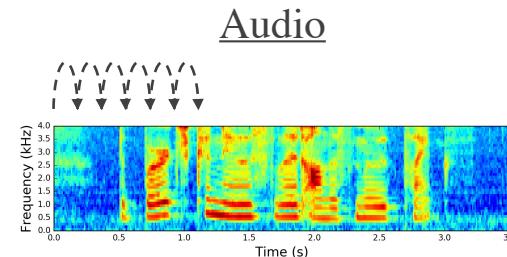
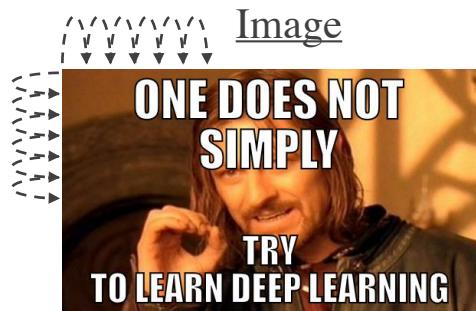
- Model the conditional likelihoods with deep neural networks
 - Logistic regression (Frey et al., 1996), Neural nets (Bengio and Bengio, 2000)
 - E.g., learn a deep net to generate one pixel at a time given past pixels
- The learning objective is to maximize the log-likelihood $\log p(\mathbf{x})$
 - If each conditional is tractable, $\log p(\mathbf{x})$ is tractable
 - Model conditional probabilities directly and with no partition functions Z



Tractability vs Flexibility

$$p(\mathbf{x}) = p(x_1) \cdot p(x_2 | x_1) \cdot p(x_3 | x_1, x_2) \cdot \dots \cdot p(x_d | x_1, \dots, x_{d-1}) = \prod_{i=1}^d p(x_i | x_{<i})$$

- Autoregressivity means there is no integral in our computations per se
- The Bayes rule ensures the final quantity to be a probability for as long as the decompositions are probabilities
- With the caveat of sequential computation, so we literally must compute pixel by pixel
- During the forward pass, the backward pass, inference, sampling



Text

‘One does not simply

Neural Autoregressive Density Estimation

- Inspired by RBMs but with tractable density estimation
 - Each conditional modelled with sigmoidal neural net like in RBMs
- Parameter matrix W maps past inputs $\mathbf{v}_{<i}$ to hidden feature \mathbf{h}_i
- Parameter matrix V generates pixel v_i given the hidden feature \mathbf{h}_i
$$p(v_i | \mathbf{v}_{<i}) = \sigma(b_i + (V^T)_{i,\cdot} \mathbf{h}_i)$$
$$\mathbf{h}_i = \sigma(c + W_{\cdot,<i} \mathbf{v}_{<i})$$
- Teacher forcing
 - During training use ground truth past inputs $\mathbf{v}_{<i}$
 - During testing use predicted past inputs $\hat{\mathbf{v}}_{<i}$

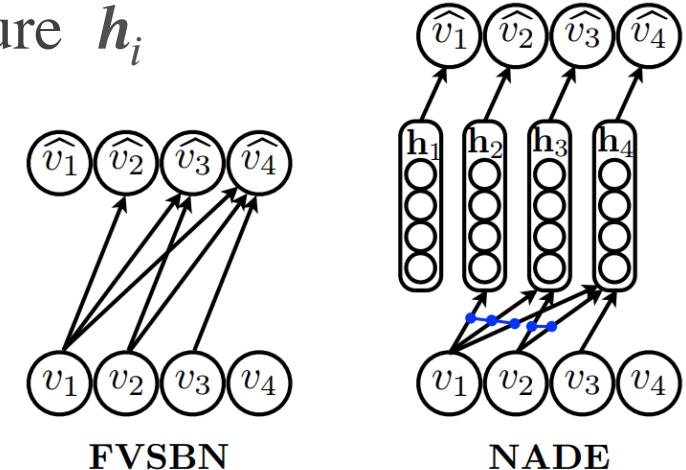


Figure 1: (Left) Illustration of a fully visible sigmoid belief network. (Right) Illustration of a neural autoregressive distribution estimator. \hat{v}_i is used as a shorthand for $p(v_i = 1 | \mathbf{v}_{<i})$. Arrows connected by a blue line correspond to connections with shared or tied parameters.

Larochelle and Murray, Neural Autoregressive Distribution Estimation

Masked Autoencoder for Distribution Estimation

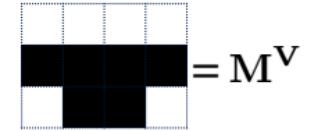
- Make an autoregressive autoencoder by setting each output x_i depend only on previous outputs $\mathbf{x}_{<i}$
 - In autoencoders the output dimensions depend on ‘future’ dimensions also
- Implement this by introducing a masking matrix M to multiply weights

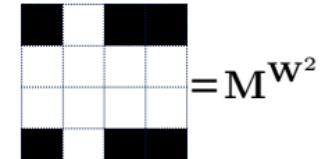
$$h(\mathbf{x}) = g\left(\mathbf{b} + \left(W \odot M^W\right) \cdot \mathbf{x}\right) = \sigma\left(c + \left(V \odot M^V\right) \cdot h(\mathbf{x})\right)$$

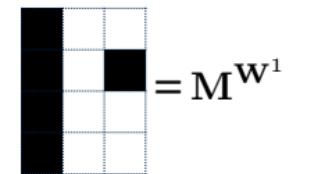
For the k -th neuron the mask column is $M_{k,d} = \begin{cases} 1 & m(k) \geq d \\ 0 & \text{otherwise} \end{cases}$

And $m(k)$ is a integer between 1 and $d - 1$

Masks

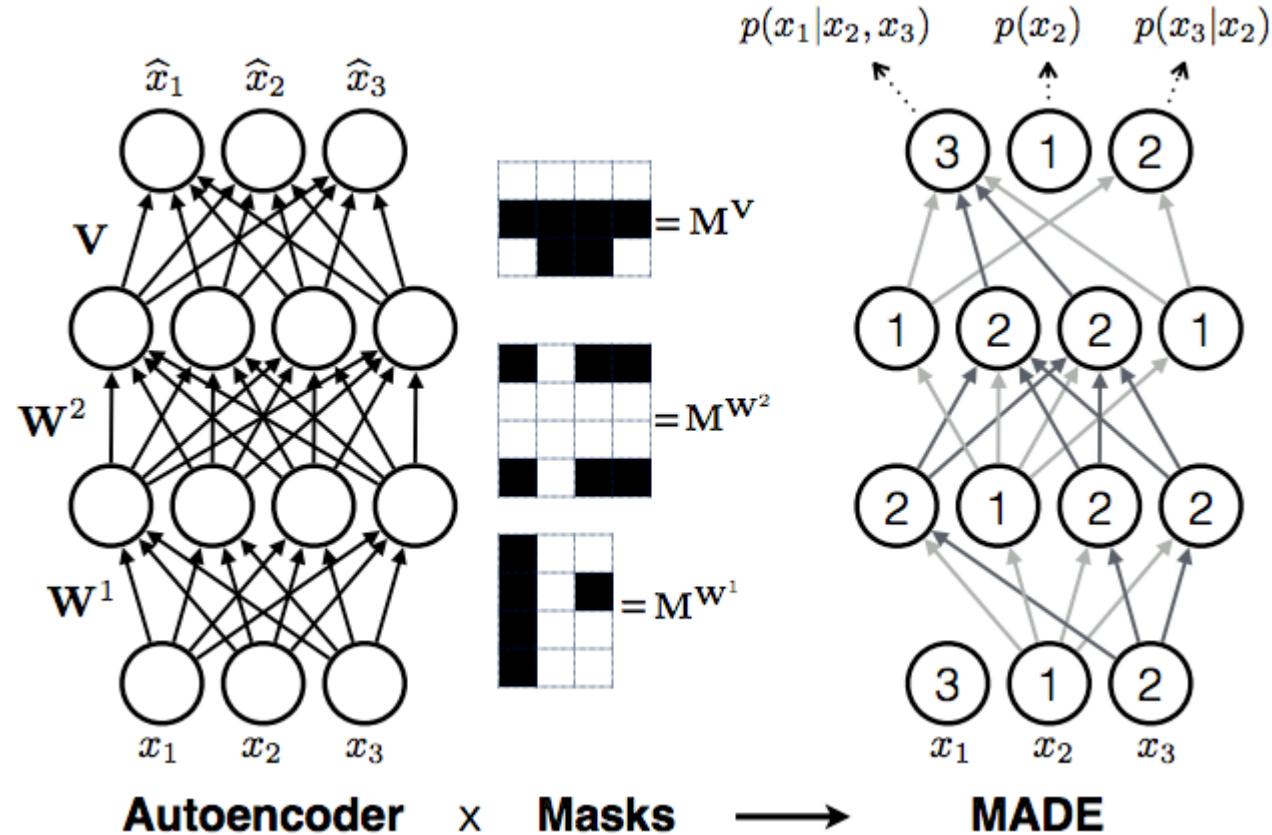

$$= \mathbf{M}^V$$


$$= \mathbf{M}^{W^2}$$


$$= \mathbf{M}^{W^1}$$

Germain, Gregor, Murray, Larochelle, Masked Autoencoder for Distribution Estimation

MADE architecture



Modern autoregressive models



Figure 1. Image completions sampled from a PixelRNN.

PixelRNN

- Decompose the data likelihood of an $n \times n$ image $p(\mathbf{x}) = \prod_{i=1}^{n^2} p(x_i | \mathbf{x}_{<i})$
- Each pixel conditional corresponds to a triplet of colors
→ Further decompose per color (same as above)

$$p(x_i | \mathbf{x}_{<i}) = p(x_{i,R} | \mathbf{x}_{<i}) \cdot p(x_{i,G} | \mathbf{x}_{<i}, x_{i,R}) \cdot p(x_{i,B} | \mathbf{x}_{<i}, x_{i,R}, x_{i,G})$$

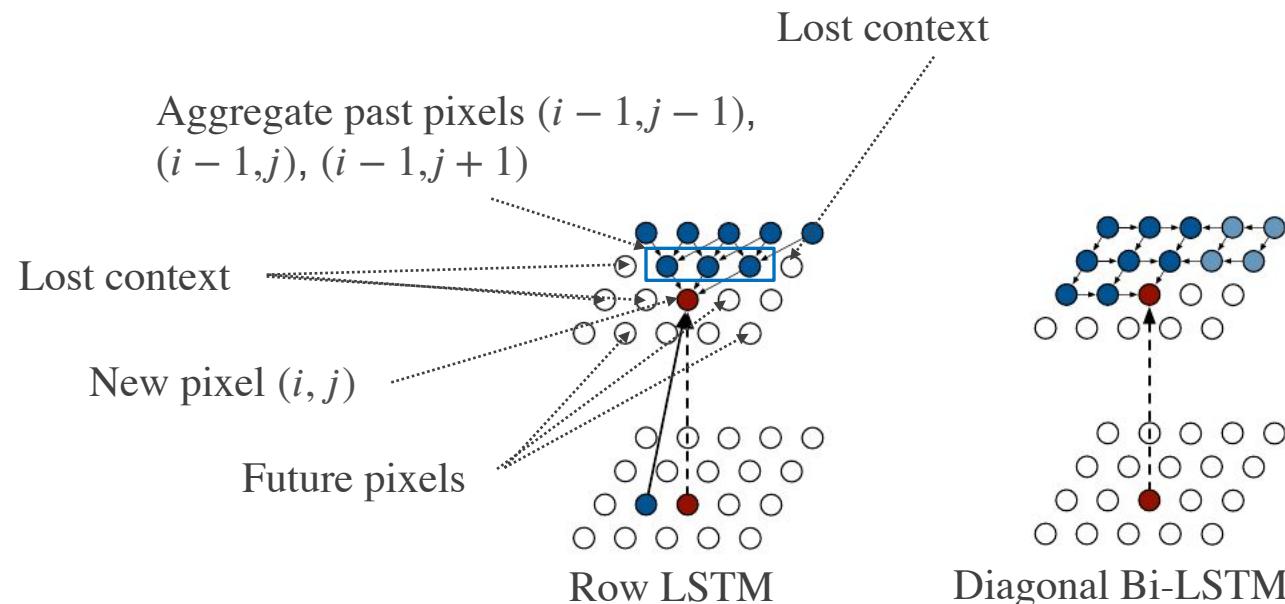


- Model the conditionals $p(x_{i,R} | \mathbf{x}_{<i}), \dots$ with 12-layer convolutional RNN
 - The MLP from NADE cannot easily scale and statistics are not shared
- Model the output as a categorical distribution
 - 256-way softmax

van den Oord, Kalchbrenner and Kavukcuoglu, Pixel Recurrent Neural Networks

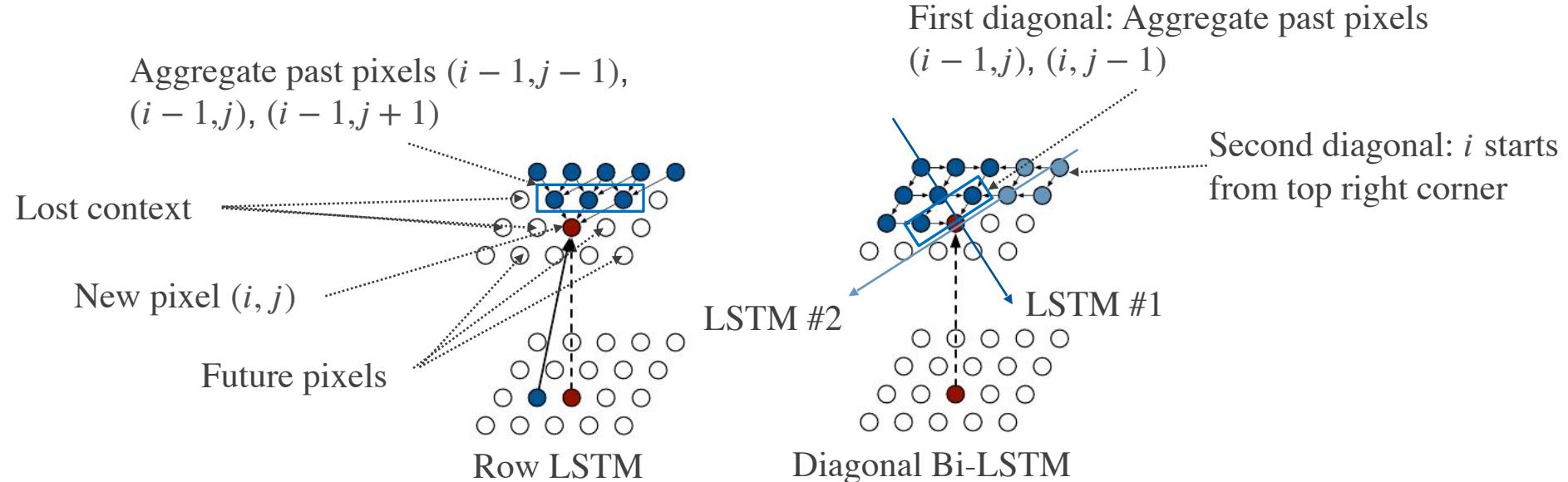
Row LSTM

- Row LSTM with ‘causal’ triangular receptive field
 - Per new pixel (row i) use 1-d conv (size 3) to aggregate pixels above $(i - 1)$
 - The effective receptive field spans a triangle
 - Convolution only on ‘past’ pixels $(i - 1)$, not ‘future pixels’ \rightarrow causal
 - Loses some context



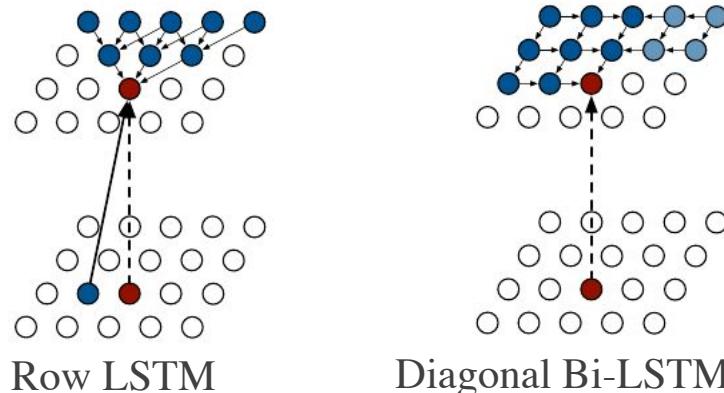
Diagonal BiLSTM

- Have two LSTMs moving on opposite diagonals
 - First diagonal: the convolution past is $(i - 1, j), (i, j - 1)$
- Combine the two LSTMs
 - recursively the entirety of past context is captured



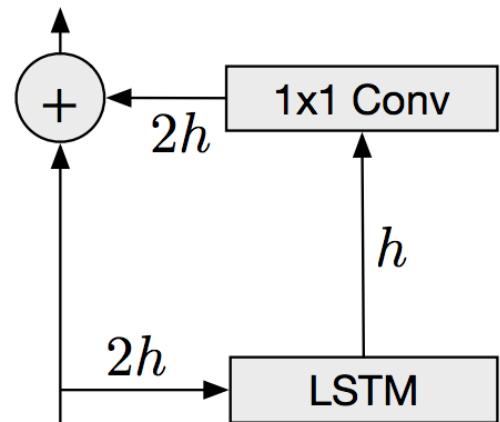
Why not a regular LSTM?

- It would require sequential, pixel-wise computations
 - Less parallelization
 - Slower training
- With Row LSTM and Diagonal BiLSTM we process one row at a time
 - Parallelization possible



Deep LSTMs with Residual connections

- Use 12 layers of LSTMs
- Add residual connections to speed up learning
- Although good modelling of $p(x) \rightarrow$ nice image generation
- Slow training because of LSTM, slow generation



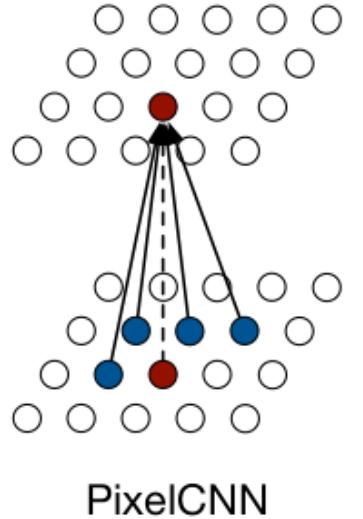
PixelRNN - Generations



Figure 1. Image completions sampled from a PixelRNN.

PixelCNN

- Replace LSTMs with fully convolutional networks
 - 15 layers
 - No pooling layers to preserve spatial resolution
- Use masks to mask out future pixels in convolutions
 - Otherwise ‘access to future’ → no ‘autoregressiveness’
- Faster training as no recurrent steps required
 - Better parallelization
 - Pixel generation still sequential and thus slow



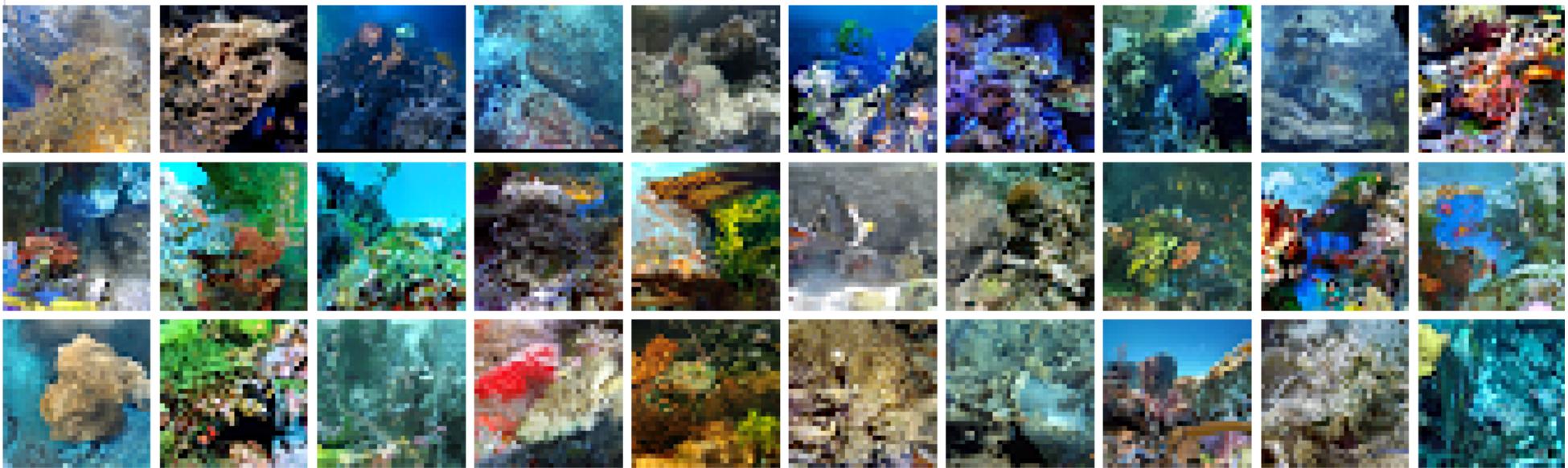
1	1	1	1	1
1	1	1	1	1
1	1	0	0	0
0	0	0	0	0
0	0	0	0	0

Masking convolutions

van den Oord, Kalchbrenner and Kavukcuoglu, Pixel Recurrent Neural Networks

PixelCNN - Generations

Coral reef



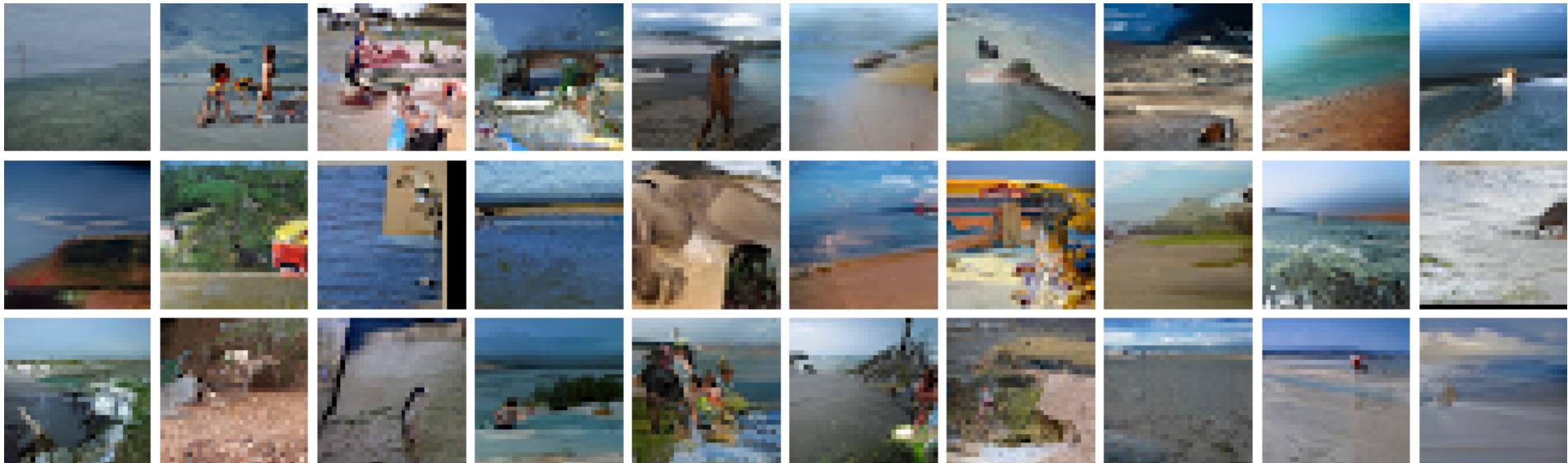
PixelCNN - Generation

Sorrel horse



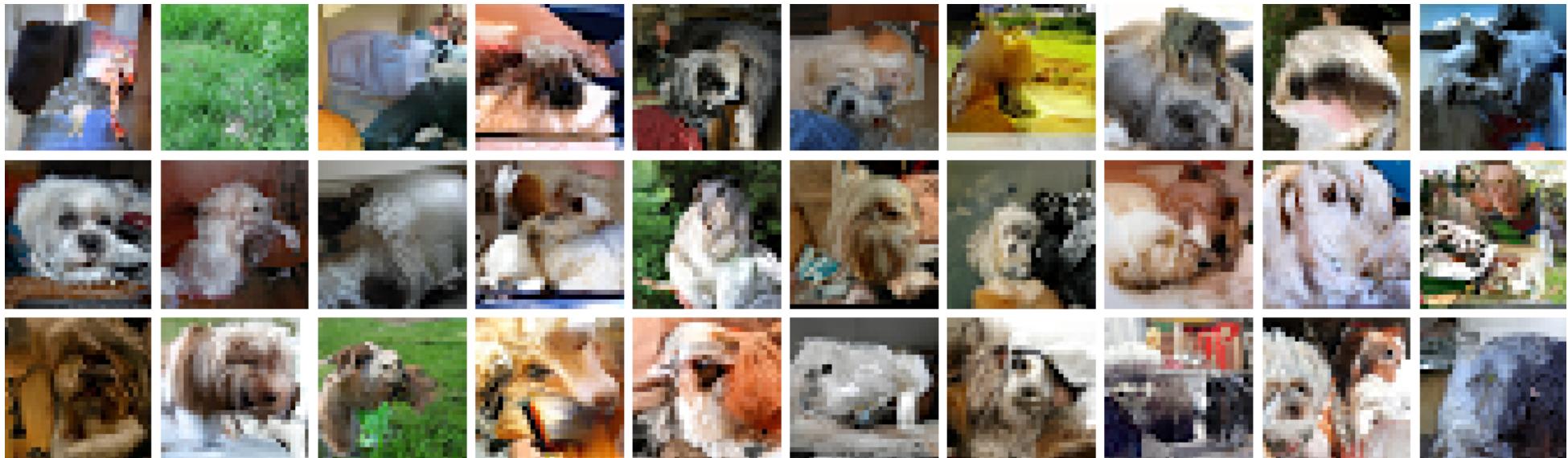
PixelCNN - Generation

Sandbar



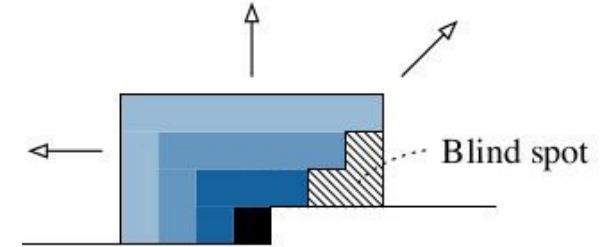
PixelCNN - Generation

Lhasa Apso



PixelCNN: Pros and Cons

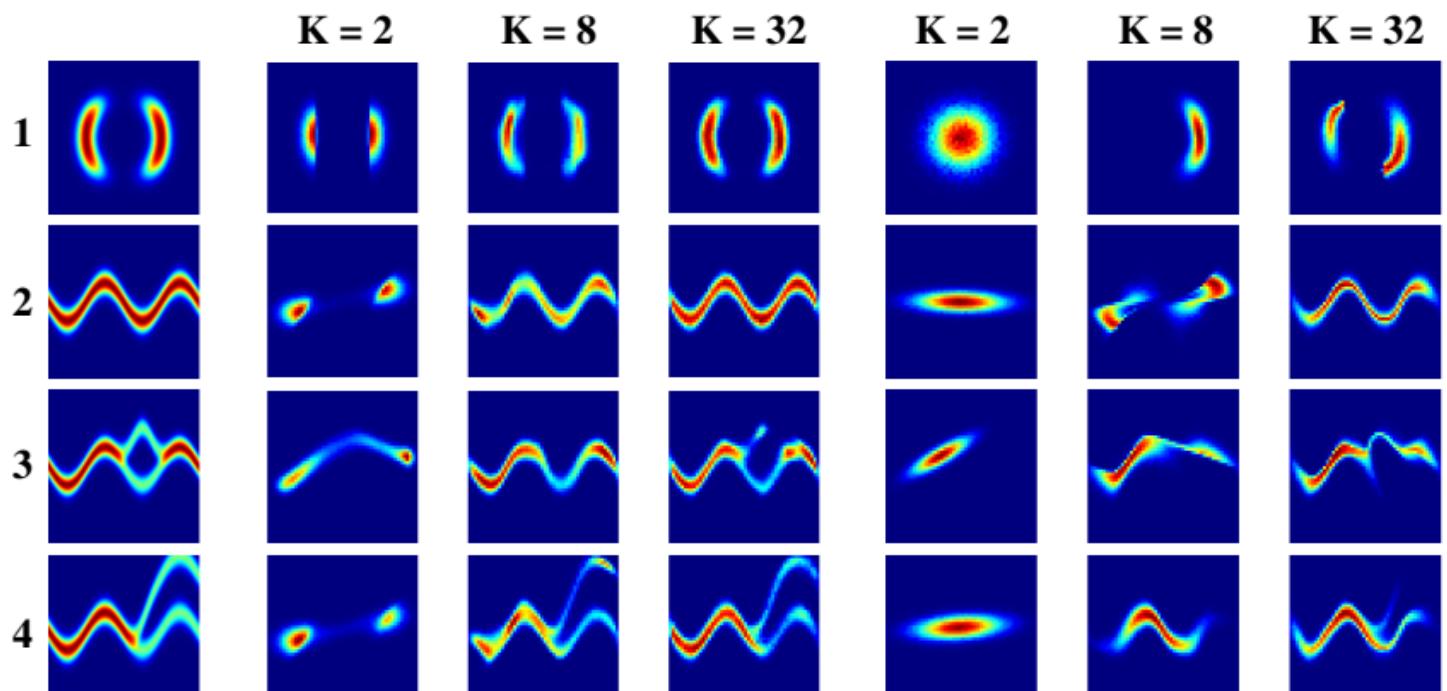
- Faster training
- Performance is worse than PixelRNN as context is discarded
- The cascaded convolutions create a ‘blind spot’
 - Use Gated PixelCNN to fix
- No latent space
- PixelCNN++ improves PixelCNN by (Salimans et al.)
 - Model output by discretized logistic mixture likelihood \leftarrow Softmax requires many parameters and yields very sparse gradients
 - Condition on whole pixels, not colors
 - Architectural innovations



Autoregressive models: pros and cons

- Top density estimation
- They take into account complex co-dependencies
 - Potentially, better generations and more accurate likelihoods
- Autoregressive models are not necessarily latent variable models
 - They neither have necessarily an encoder nor learn representations
- Slow in learning, inference and generation
 - Computations are sequential (one at a time) → limited parallelism
 - E.g., to generate the next word we must generate past words first
- They may introduce artificial bias when assumed order is imposed

Normalizing flows



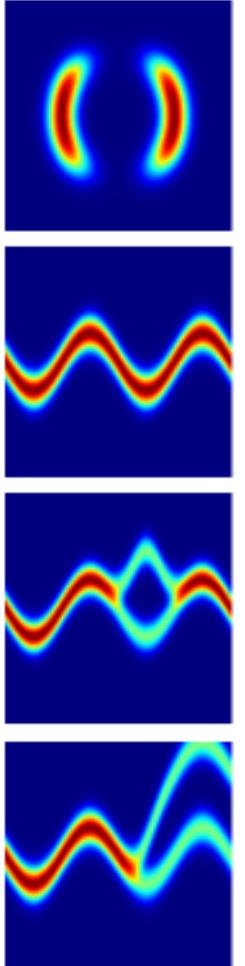
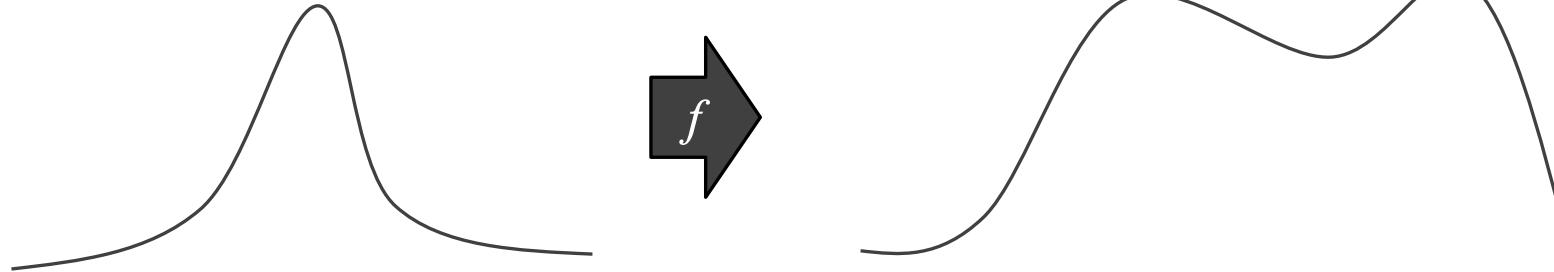
(a)

(b) Norm. Flow

(c) NICE

Intuition

- Often our posterior approximation is not enough
 - Imagine data generated by two modes
→ approximating with a standard Gaussian would be problematic
- If we applying a transformation to a simple density input
 - *e.g.*, a Gaussian
 - We can morph it into a more complicated density function
- Doing it many times → model any complex density



Change of variables

- For 1-d variables we know that

$$\int f(g(x))g'(x)dx = \int f(u)du, \text{ where } u = g(x)$$

- This is called change of variables (or [integration by substitution](#))
- The density is $du = g'(x)dx$
- For multivariate cases

$$\int f(g(\mathbf{x})) \left| \det \frac{dg}{d\mathbf{x}} \right| d\mathbf{x} = \int f(\mathbf{u})d\mathbf{u}, \text{ where } \mathbf{u} = g(\mathbf{x})$$

And \mathbf{u} and \mathbf{x} have the same dimensionality

Normalizing flows

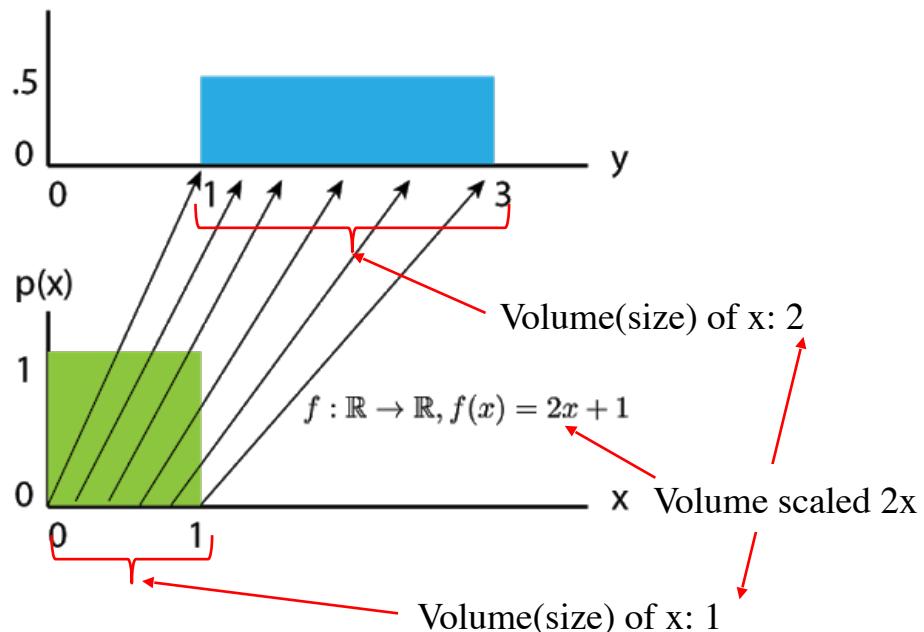
- Our model is an encoder $f: \mathbf{x} \rightarrow \mathbf{z}$
 - It maps the input \mathbf{x} with density $p(\mathbf{x})$ to the latent \mathbf{z} with density $p(\mathbf{z})$
- The inverse model is the decoder $f^{-1}: \mathbf{z} \rightarrow \mathbf{x}$
 - It maps the latent \mathbf{z} back to the input \mathbf{x}
- For our forward model we have

$$p(\mathbf{z}) \left| \det \frac{df}{d\mathbf{x}} \right| = p(\mathbf{x}) \Leftrightarrow \log p(\mathbf{z}) + \log \left| \det \frac{df}{d\mathbf{x}} \right| = \log p(\mathbf{x})$$

Rezende and Mohamed, Variational Inference with Normalizing Flows

Normalizing flows geometrically

- The determinant shows how the volume of the input and output probability spaces changes
 - The volumes (sizes) must change so that $\int p(z) = 1$ and $\int p(x) = 1$
 - Normalizing flows expand or contract the density



Stacking normalizing flows

- The change of variables can be applied recursively

$$\mathbf{x} = \mathbf{z}_K = f_K \circ \dots \circ f_2 \circ f_1(\mathbf{z}_0)$$

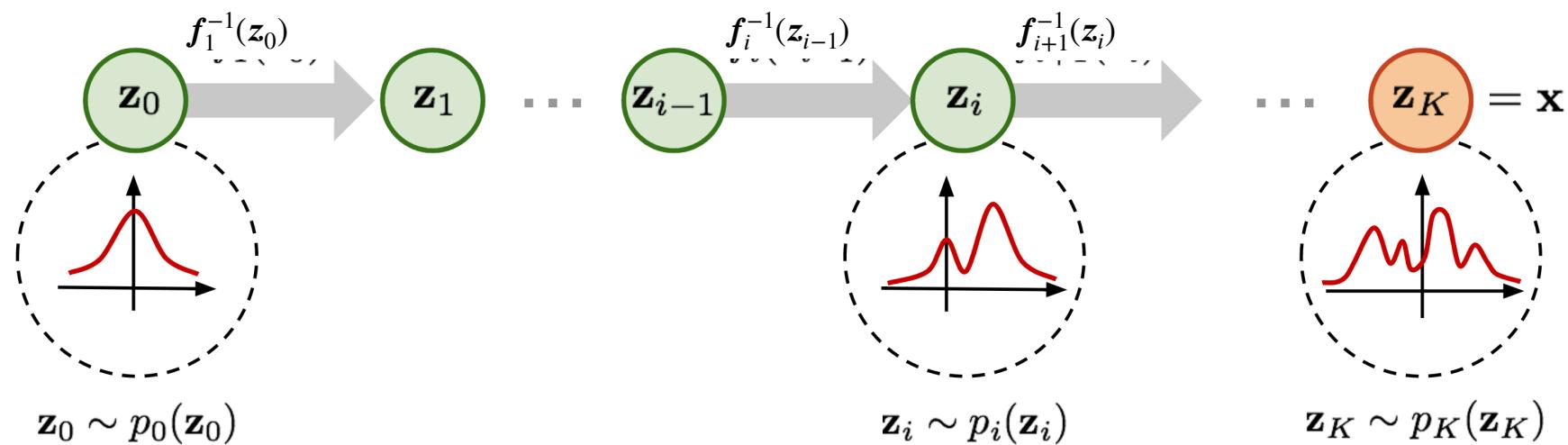
(e.g., $\mathbf{z}_0 = f_1(\mathbf{z}_1) \Leftrightarrow \mathbf{z}_1 = f_1^{-1}(\mathbf{z}_0)$)

- The log density of our data is

$$\log p(\mathbf{x}) = \log p_K(\mathbf{z}_K) = \log p_0(\mathbf{z}_0) - \sum_{k=1}^K \log \det \left| \frac{\partial f_k}{\partial \mathbf{z}_k} \right|$$

- Optimize with maximum likelihood

Stacking normalizing flows



What transformations?

$$\log p_K(\mathbf{z}_K) = \log p_0(\mathbf{z}_0) - \sum_{k=1}^K \text{logdet} \left| \frac{\partial f_k}{\partial \mathbf{z}_k} \right|$$

- We want smooth, differentiable transformations f_k
 - For which it is easy to compute inverse f_k^{-1}
 - and determinant of the Jacobian $\det \frac{df_k}{d\mathbf{z}_k}$
- Example transformations
 - Planar flows
 - Radial flows
 - Coupling layers

Planar flow

- The transformation is

$$f(\mathbf{z}) = \mathbf{z} + \mathbf{u} \mathbf{h}(\mathbf{w}^T \mathbf{z} + \mathbf{b})$$

- \mathbf{u} , \mathbf{w} , \mathbf{b} are free parameters
- \mathbf{h} is an element-wise non-linearity (element-wise so that it is easy to invert)
- The log-determinant of the Jacobian is

$$\begin{aligned}\psi(\mathbf{z}) &= h'(\mathbf{w}^T \mathbf{z} + \mathbf{b}) \mathbf{w} \\ \det \left| \frac{\partial f}{\partial \mathbf{z}} \right| &= |1 + \mathbf{u}^T \psi(\mathbf{z})|\end{aligned}$$

Radial flow

- The transformation is

$$f(\mathbf{z}) = \mathbf{z} + \beta h(\alpha, r)(\mathbf{z} - \mathbf{z}_0)$$

Where $h(\alpha, r) = 1/(\alpha + r)$

- The log-determinant of the Jacobian is

$$\det \left| \frac{\partial f}{\partial \mathbf{z}} \right| = [1 + \beta h(\alpha, r)]^{d-1} [1 + \beta h(\alpha, r) + h'(\alpha, r)r]$$

VAE with Normalizing Flows

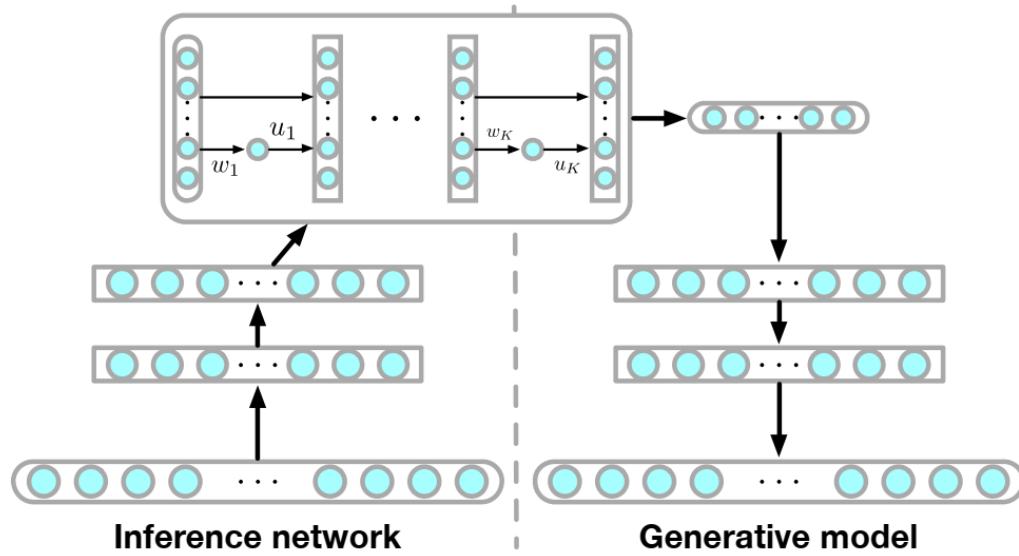


Figure 2. Inference and generative models. Left: Inference network maps the observations to the parameters of the flow; Right: generative model which receives the posterior samples from the inference network during training time. Round containers represent layers of stochastic variables whereas square containers represent deterministic layers.

Learning better posteriors with variational inference

- Again: the evidence lower bound is

$$\text{ELBO}_{\theta,\varphi}(\mathbf{x}) = \log p(\mathbf{x}) - \text{KL}(q_{\varphi}(\mathbf{z} | \mathbf{x}) || p(\mathbf{z} | \mathbf{x}))$$

- Replace the simple approximate posterior by normalizing flows

$$\mathbb{E}_{q_0(\mathbf{z}_0 | \mathbf{x})} \left[\log p_{\theta}(\mathbf{x} | \mathbf{z}_K) \right] - \text{KL}(q_0(\mathbf{z}_0 | \mathbf{x}) || p(\mathbf{z})) + \mathbb{E}_{q_0(\mathbf{z}_0 | \mathbf{x})} \left[\sum_{k=1}^K \log \left| \det \frac{df_k}{d\mathbf{z}_k} \right|^{-1} \right]$$

Algorithm 1 Variational Inf. with Normalizing Flows

Parameters: ϕ variational, θ generative

while not converged **do**

$\mathbf{x} \leftarrow \{\text{Get mini-batch}\}$

$\mathbf{z}_0 \sim q_0(\bullet | \mathbf{x})$

$\mathbf{z}_K \leftarrow f_K \circ f_{K-1} \circ \dots \circ f_1(\mathbf{z}_0)$

$\mathcal{F}(\mathbf{x}) \approx \mathcal{F}(\mathbf{x}, \mathbf{z}_K)$

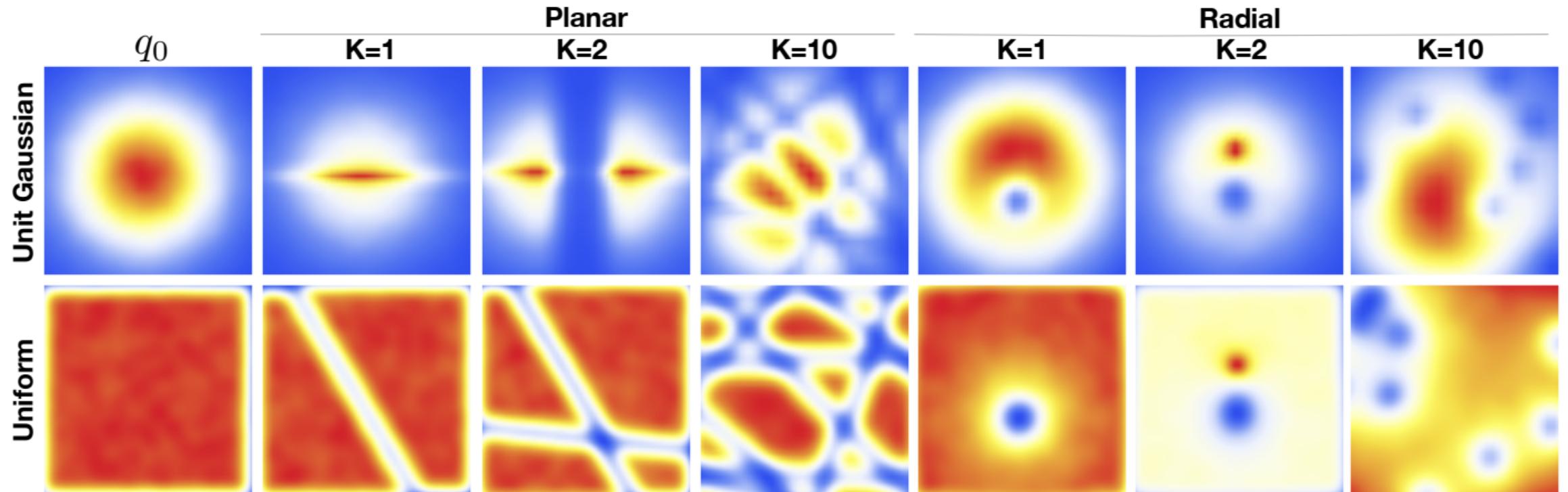
$\Delta\theta \propto -\nabla_{\theta} \mathcal{F}(\mathbf{x})$

$\Delta\phi \propto -\nabla_{\phi} \mathcal{F}(\mathbf{x})$

end while

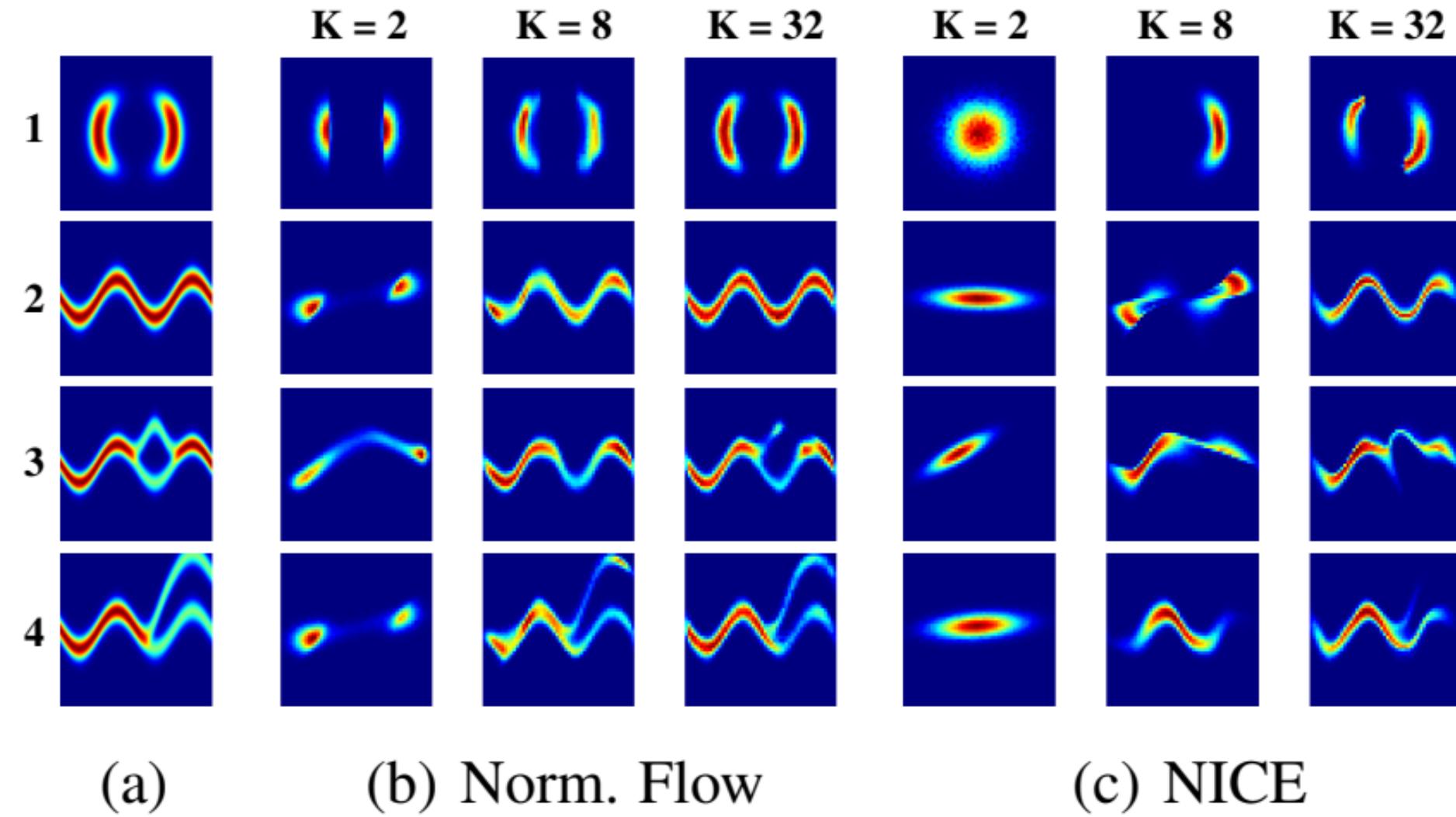
Rezende and Mohamed, Variational Inference with Normalizing Flows

The effect of number of transformations/flows



Rezende and Mohamed, Variational Inference with Normalizing Flows

Some results



Rezende and Mohamed, Variational Inference with Normalizing Flows

Flow-based models



Figure 4: Random samples from the model, with temperature 0.7

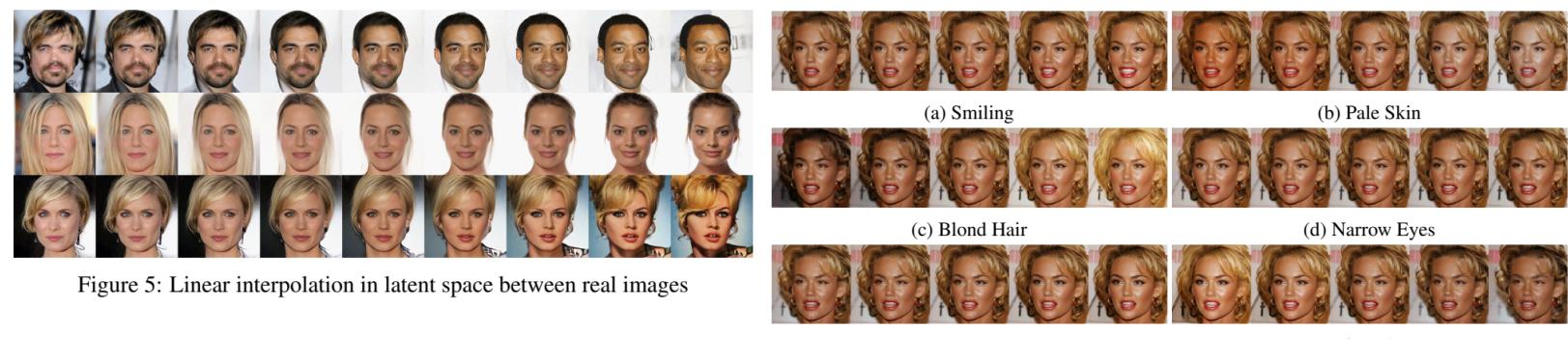
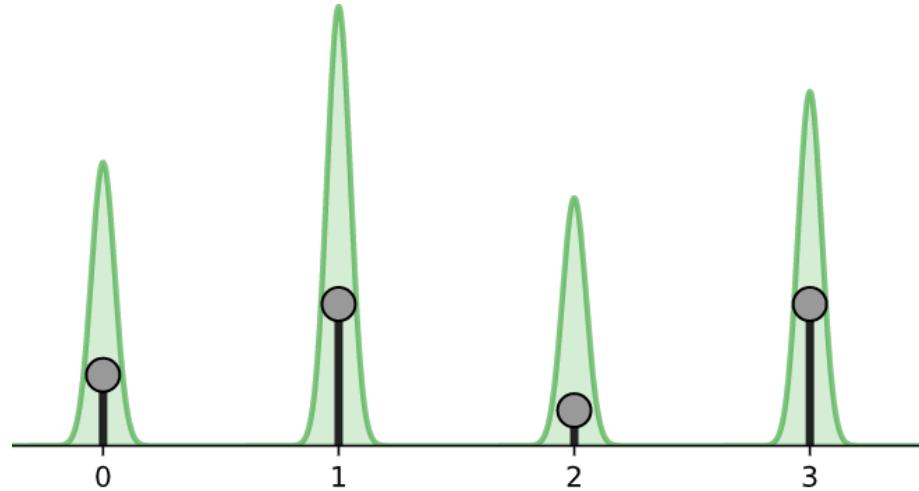


Figure 5: Linear interpolation in latent space between real images

Normalizing flows on images

- Normalizing flows are continuous transformations
- Images contain discrete values
 - The model will assign δ -peak probabilities on integer (pixel) values only
 - These probabilities will be nonsensical, there is no smoothness



[UVADLC tutorial](#)

(Variational) dequantization

- Add (continuous) noise $u \sim q(u | x)$ to input variables $v = x + u$
- The data log-likelihood then is

$$\log p(x) = \log \int p(x + u) du = \log \mathbb{E}_{u \sim q(u | x)} \left[\frac{p(x + u)}{q(u | x)} \right] \geq \mathbb{E}_{u \sim q(u | x)} \log \left[\frac{p(x + u)}{q(u | x)} \right]$$

- If $q(u | x)$ is the uniform distribution the standard dequantization
 - Probability between two consecutive values is fixed
→ resemble boxy boundaries between values
- Better learn $q(u | x)$ in a variational manner
→ Variational dequantization

Coupling layers

- Given input \mathbf{z} the output of the transformation is

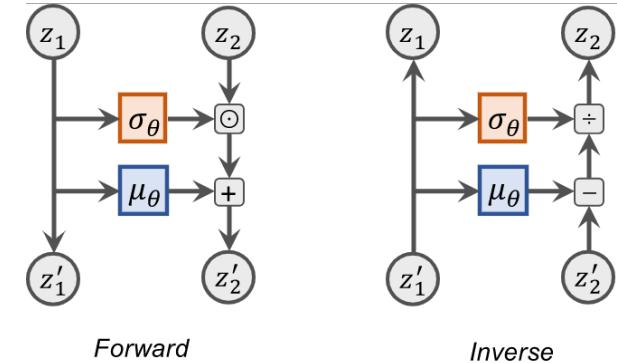
$$\mathbf{z}' = \begin{bmatrix} \mathbf{z}'_{1:j} \\ \mathbf{z}'_{j+1:d} \end{bmatrix} = \begin{bmatrix} \mathbf{z}_{1:j} \\ \mu_\theta(\mathbf{z}_{1:j}) + \sigma_\theta(\mathbf{z}_{1:j}) \odot \mathbf{z}_{j+1:d} \end{bmatrix}$$

- $\mu_\theta, \sigma_\theta$ are neural networks with shared parameters

$$\text{Easy inverse: } \mathbf{z} = \begin{bmatrix} \mathbf{z}_{1:j} \\ (\mathbf{z}'_{j+1:d} - \mu_\theta(\mathbf{z}_{1:j})) \\ \hline \sigma_\theta(\mathbf{z}_{1:j}) \end{bmatrix}$$

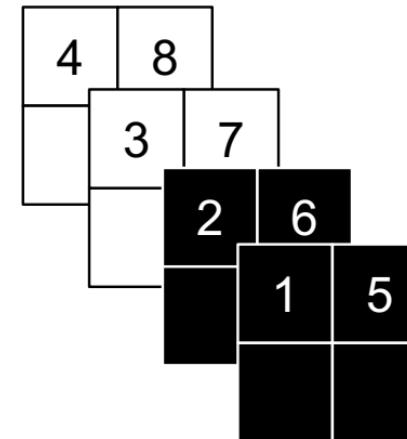
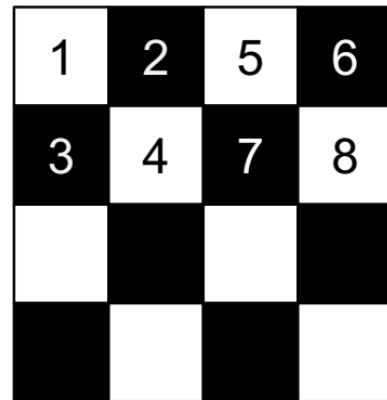
- Easy triangular Jacobian $\frac{\partial \mathbf{z}'}{\partial \mathbf{z}} = \begin{bmatrix} \mathbb{I}_d & 0 \\ \frac{\partial \mathbf{z}'_{j+1:d}}{\partial \mathbf{z}_{1:j}} & \text{diag}\left(\sigma_\theta(\mathbf{z}_{1:j})\right) \end{bmatrix}$

- The log determinant is $\sum_j \log \sigma_\theta(\mathbf{z}_j)$



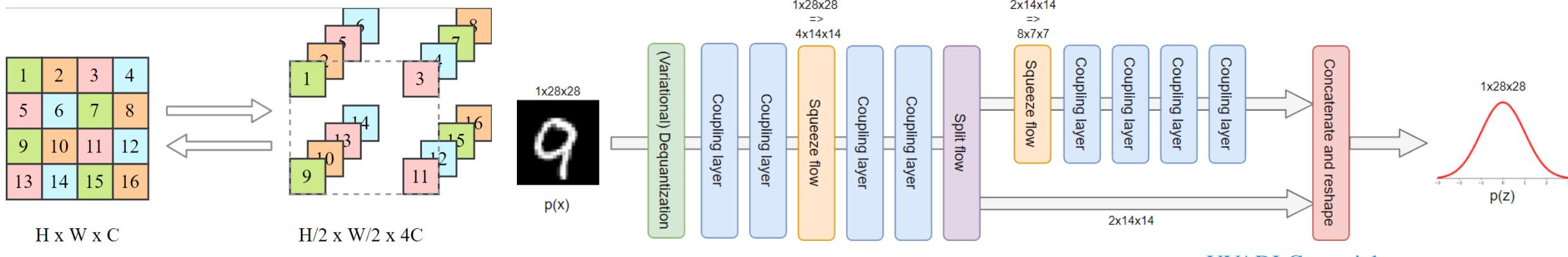
Splitting dimensions in images

- Use masking
 - Checkers pattern
 - Splitting across channels
- Alternate dimensions between consecutive layers
→ not always the same $1:d$ dimensions remain untouched



Multi-scale architecture

- Invertibility → number of dimensions before and after f is the same
 - High computational complexity for large images
- Apply new transformations to half the input only
 - For the other half use the prior (previous) transformations
- Use squeeze to turn spatial to channel dimensions
 - And split for halving the input



GLOW, FLOW, FLOW++

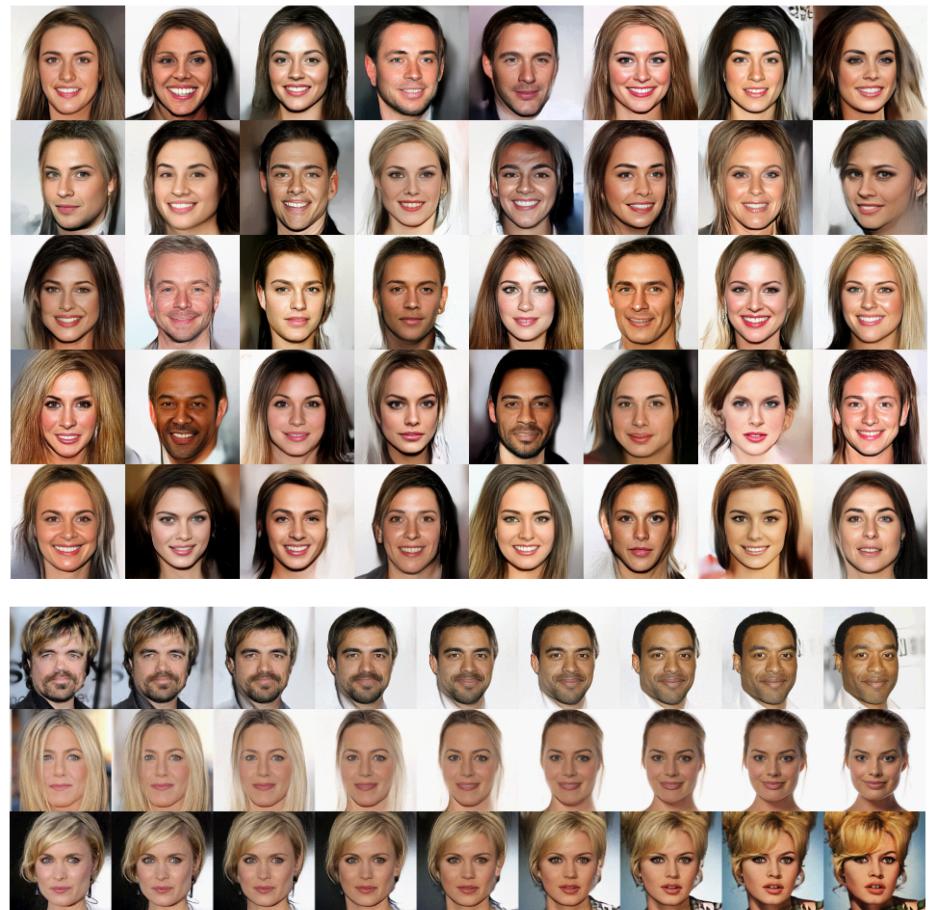


Figure 5: Linear interpolation in latent space between real images

Kingma, Dhariwal, Glow: Generative Flow with Invertible 1x1 Convolutions

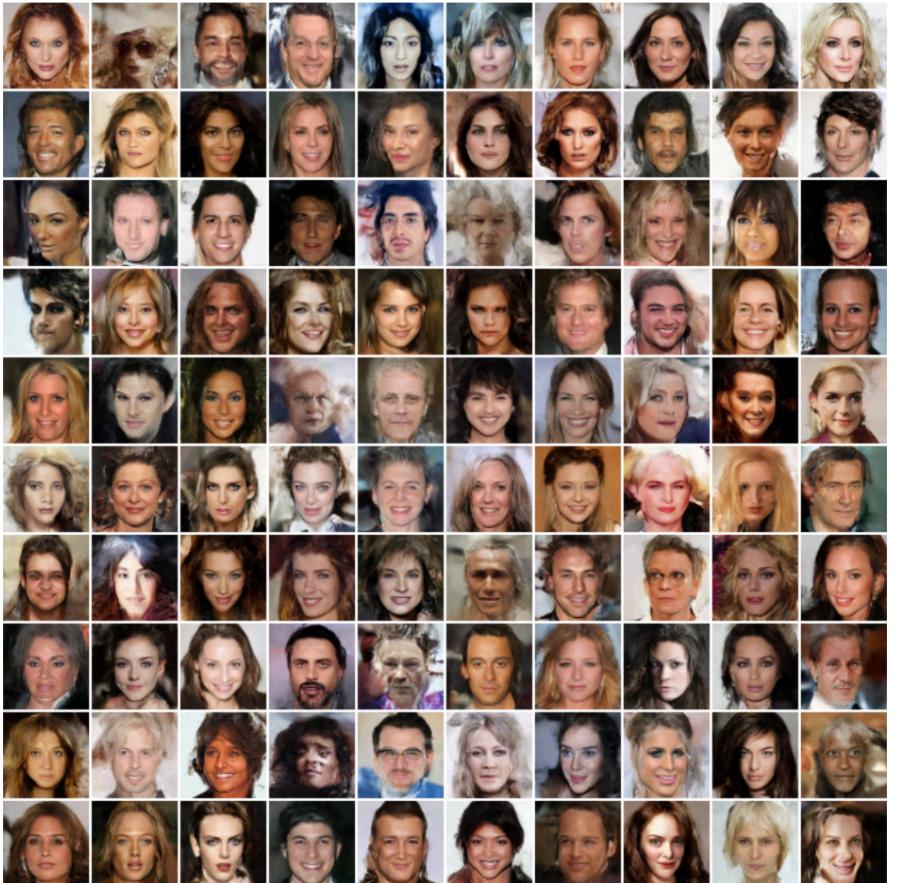


Figure 4. Samples from Flow++ trained on 5-bit 64x64 CelebA, without low-temperature sampling.

Kingma, Dhariwal, Flow++: Improving Flow-Based Generative Models with Variational Dequantization and Architecture Design

Categorical normalizing flows

- Normalizing flows with variational inference to learn representations of categorical data on continuous space
 - Learnable, smooth, support for higher dimensions
- Learning must ensure no loss of information
 - the volumes that represent categorical data must not-overlap
 - Otherwise, to which category does the representation correspond to?

$$p(\mathbf{x}) \geq \mathbb{E}_{\mathbf{z} \sim q(\cdot | \mathbf{x})} \left[\frac{\prod_i p(x_i | z_i)}{q(\mathbf{z} | \mathbf{x})} p(\mathbf{z}) \right]$$

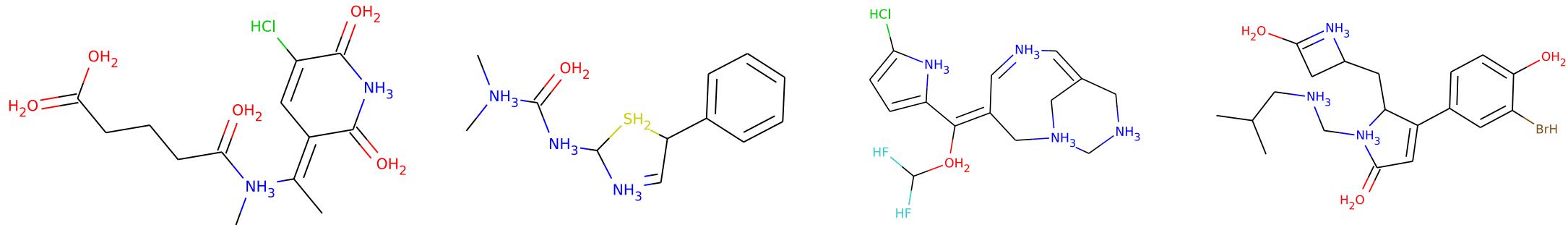
- Factorized posterior $\prod_i p(x_i | z_i)$ to encourage learning non-overlapping z_i

Lippe and Gavves, Categorical Normalizing Flows via Continuous Transformations, in submission to ICLR 2021

Graph generation with categorical normalizing flows

Results on the Zinc250k dataset (224k examples)

Method	Validity	Uniqueness	Novelty	Reconstruction	Parallel	General
JT-VAE	100%	100%	100%	71%	✗	✗
GraphAF	68%	99.10%	100%	100%	✗	✓
R-VAE	34.9%	100%	—	54.7%	✓	✓
GraphNVP	42.60%	94.80%	100%	100%	✓	✓
GraphCNF	83.41% (± 2.88)	99.99% (± 0.01)	100% (± 0.00)	100% (± 0.00)	✓	✓
+ Sub-graphs	96.35% (± 2.21)	99.98% (± 0.01)	99.98% (± 0.02)	100% (± 0.00)	✓	✓



Normalizing flows: pros and cons

- Starting from a simple density like a unit Gaussian we can obtain any complex density that match our data without even knowing its analytic form
- Tractable density estimation
- Efficient parallel sampling and learning
- Often very many transformations required → Very large networks needed
- Constrained to invertible transformations with tractable determinant
- Tied encoder and decoder weights
- Transformations cannot easily introduce bottlenecks

[UVADLC tutorial](#)

A summary of properties

	Training	Likelihood	Sampling	Compression
Autoregressive models (e.g., PixelCNN)	Stable	Yes	Slow	No
Flow-based models (e.g., RealNVP)	Stable	Yes	Fast/Slow	No
Implicit models (e.g., GANs)	Unstable	No	Fast	No
Prescribed models (e.g., VAEs)	Stable	Approximate	Fast	Yes

17

J. Tomczak's lecture from April, 2019

Summary

- Early autoregressive models
- Modern autoregressive models
- Normalizing flows
- Flow-based models