

# Yes, we GAN.

VLC =  $\iiint$  Vision  
Learning *aVdLdC*  
Control

## Deep Generative Modelling

Jonathon Hare

Vision, Learning and Control  
University of Southampton

- What is generative modelling and why do we do it?
- Differentiable Generator Networks
- Variational Autoencoders
- Generative Adversarial Networks

## Generative Modelling and Differentiable Generator Networks

## Recap: Generative Models

- Learn models of the data:  $p(\mathbf{x})$
- Learn *conditional* models of the data:  $p(\mathbf{x}|y = y)$
- Some generative models allow the probability distributions to be evaluated explicitly
  - i.e. compute the probability of a piece of data  $\mathbf{x}$ :  $p(\mathbf{x} = \mathbf{x})$
- Some generative models allow the probability distributions to be sampled
  - i.e. draw a sample  $\mathbf{x}$  based on the distribution:  $\mathbf{x} \sim p(\mathbf{x})$
- Some generative models can do both of the above
  - e.g. a Gaussian Mixture Model is an explicit model of the data using  $k$  Gaussians
    - The likelihood of data  $\mathbf{x}$  is the weighted sum of the likelihood from each of the  $k$  Gaussians
    - Sampling can be achieved by sampling the categorical distribution of  $k$  weights followed by sampling a data point from the corresponding Gaussian

## Why do generative modelling?

- Try to understand the processes through which the data was itself generated
  - Probabilistic latent variable models like VAEs or topic models (PLSA, LDA, ...) for text
  - Models that try to disentangle latent factors like  $\beta$ -VAE
- Understand how likely a new or previously unseen piece of data is
  - outlier prediction, anomaly detection, ...
- Make 'new' data
  - Make 'fake' data to use to train large supervised models?
  - 'Imagine' new, but plausible, things?

# Differentiable Generator Networks

- Generative Modelling is not new; we've known how to make arbitrarily complex probabilistic graphical models for many years.
  - ...But difficult to train and scale to real data, relying on MCMC.
- The past few years has seen major progress along four loose strands:
  - Invertible density estimation - A way to specify complex generative models by transforming a simple latent distribution with a series of invertible functions.
  - Autoregressive models - Another way to model  $p(x)$  is to break the model into a series of conditional distributions:  
$$p(x) = p(x_1)p(x_2|x_1)p(x_3|x_2, x_1) \dots$$
  - Variational autoencoders - Latent-variable models that use a neural network to do approximate inference.
  - Generative adversarial networks - A way to train generative models by optimizing them to fool a classifier
- **Common thread in recent advances is that the loss functions are end-to-end differentiable.**

## Differentiable Generator Networks: key idea

- We're interested in models that transform samples of latent variables  $z$  to
  - samples  $x$ , or,
  - distributions over samples  $x$
- The model is a (differentiable) function  $g(z, \theta)$ 
  - typically  $g$  is a neural network.

## Example: drawing samples from $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

- Consider a simple generator network with a single affine layer that maps samples  $\mathcal{N}(\mathbf{0}, \mathbf{I})$  to  $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ :

$$\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \longrightarrow \boxed{g_{\theta}(\mathbf{z})} \longrightarrow \mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

- Note: Exact solution is  $\mathbf{x} = g_{\theta}(\mathbf{z}) = \boldsymbol{\mu} + \mathbf{L}\mathbf{z}$  where  $\mathbf{L}$  is the Cholesky decomposition of  $\boldsymbol{\Sigma}$ :  $\boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}^{\top}$ , lower triangular  $\mathbf{L}$ .

## Generating samples

More generally, we can think of  $g$  as providing a nonlinear change of variables that transforms a distribution over  $\mathbf{z}$  into the desired distribution over  $\mathbf{x}$ :

$$p_{\mathbf{z}}(\mathbf{z}) \longrightarrow \boxed{g(\mathbf{z})} \longrightarrow p_{\mathbf{x}}(\mathbf{x})$$

For any *invertible, differentiable, continuous*  $g$ :

$$p_{\mathbf{z}}(\mathbf{z}) = p_{\mathbf{x}}(g(\mathbf{z})) \left| \det \left( \frac{\partial g}{\partial \mathbf{z}} \right) \right|$$

Which implicitly imposes a probability distribution over  $\mathbf{x}$ :

$$p_{\mathbf{x}}(\mathbf{x}) = \frac{p_{\mathbf{z}}(g^{-1}(\mathbf{x}))}{\left| \det \left( \frac{\partial g}{\partial \mathbf{z}} \right) \right|}$$

Note: usually use an indirect means of learning  $g$  rather than minimise  $-\log(p(\mathbf{x}))$  directly

- Rather than use  $g$  to provide a sample of  $\mathbf{x}$  directly, we could instead use  $g$  to define a conditional distribution over  $\mathbf{x}$ ,  $p(\mathbf{x}|\mathbf{z})$ 
  - For example,  $g$  might produce the parameters of a particular distribution - e.g.:
    - means of Bernoulli
    - mean and variance of a Gaussian
- The distribution over  $\mathbf{x}$  is imposed by marginalising  $\mathbf{z}$ :  $p(\mathbf{x}) = \mathbb{E}_{\mathbf{z}} p(\mathbf{x}|\mathbf{z})$

## Distributions vs Samples

- In both cases ( $g$  generates samples and  $g$  generates distributions) we can use the reparameterisation tricks we saw last lecture to train models.
- Generating distributions:
  - + works for both continuous and discrete data
  - - need to specify the form of the output distribution
- Generating samples:
  - + works for continuous data
    - + discrete data is recently possible - we need the STargmax
  - + don't need to specify the distribution in explicit form

- In classification both input and output are given
  - Optimisation only needs to learn the mapping
- Generative modelling is more complex than classification because
  - learning requires optimizing intractable criteria
  - data does not specify both input  $\mathbf{z}$  and output  $\mathbf{x}$  of the generator network
  - learning procedure needs to determine how to arrange  $\mathbf{z}$  space in a useful way and how to map  $\mathbf{z}$  to  $\mathbf{x}$

## Variational Autoencoders

# Variational Autoencoder (VAE)

- VAEs architecturally similar to autoencoders (AEs).
- VAEs (vs AEs) significantly different in their goal and mathematical formulation.
- AEs map the input into a fixed vector.
- However, VAEs map the input into a distribution.
- VAEs are a combination of neural networks (AEs) and **graphical models**.

## Graphical Models (Background)

- A graphical model is a probabilistic model for which a graph expresses the conditional dependence structure between random variables.
- Graphical models are commonly used in probability theory, statistics—particularly Bayesian statistics—and machine learning.<sup>1</sup>

---

<sup>1</sup>Definition taken from Wikipedia



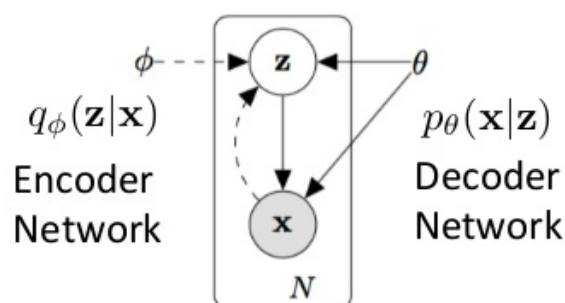
# KL Divergence (Background)

- Kullback–Leibler divergence,  $D_{KL}(P \parallel Q)$ : a measure of how one probability distribution  $Q$  is different from a second, reference probability distribution  $P$ .<sup>2</sup>
- A simple interpretation of the divergence of  $P$  from  $Q$  is the expected excess surprise from using  $Q$  as a model when the actual distribution is  $P$ .
- While it is a distance, it is not a metric, the most familiar type of distance: it is asymmetric in the two distributions.

<sup>2</sup>Definition taken from Wikipedia

## Variational Autoencoders (VAEs)

### Variational Autoencoder



Minimize:  $D_{KL}[q_\phi(z|x) \parallel p_\theta(z|x)]$

Intractable:  $p_\theta(z|x) = \frac{p_\theta(x|z)p_\theta(z)}{p_\theta(x)}$

<sup>3</sup>Auto-Encoding Variational Bayes <https://arxiv.org/abs/1312.6114>

# Variational Autoencoders (VAEs)

The distance loss just defined is expanded as

$$\begin{aligned} D_{KL}(q_{\Phi}(\mathbf{z} | \mathbf{x}) || p_{\theta}(\mathbf{z} | \mathbf{x})) &= \int q_{\Phi}(\mathbf{z} | \mathbf{x}) \log \frac{q_{\Phi}(\mathbf{z} | \mathbf{x})}{p_{\theta}(\mathbf{z} | \mathbf{x})} d\mathbf{z} \\ &= \int q_{\Phi}(\mathbf{z} | \mathbf{x}) \log \frac{q_{\Phi}(\mathbf{z} | \mathbf{x}) p_{\theta}(\mathbf{x})}{p_{\theta}(\mathbf{z}, \mathbf{x})} d\mathbf{z} \\ &= \int q_{\Phi}(\mathbf{z} | \mathbf{x}) \left( \log(p_{\theta}(\mathbf{x})) + \log \frac{q_{\Phi}(\mathbf{z} | \mathbf{x})}{p_{\theta}(\mathbf{z}, \mathbf{x})} \right) d\mathbf{z} \\ &= \log(p_{\theta}(\mathbf{x})) + \int q_{\Phi}(\mathbf{z} | \mathbf{x}) \log \frac{q_{\Phi}(\mathbf{z} | \mathbf{x})}{p_{\theta}(\mathbf{z}, \mathbf{x})} d\mathbf{z} \\ &= \log(p_{\theta}(\mathbf{x})) + \int q_{\Phi}(\mathbf{z} | \mathbf{x}) \log \frac{q_{\Phi}(\mathbf{z} | \mathbf{x})}{p_{\theta}(\mathbf{x} | \mathbf{z}) p_{\theta}(\mathbf{z})} d\mathbf{z} \\ &= \log(p_{\theta}(\mathbf{x})) + E_{\mathbf{z} \sim q_{\Phi}(\mathbf{z} | \mathbf{x})} \left( \log \frac{q_{\Phi}(\mathbf{z} | \mathbf{x})}{p_{\theta}(\mathbf{z})} - \log(p_{\theta}(\mathbf{x} | \mathbf{z})) \right) \\ &= \log(p_{\theta}(\mathbf{x})) + D_{KL}(q_{\Phi}(\mathbf{z} | \mathbf{x}) || p_{\theta}(\mathbf{z})) - E_{\mathbf{z} \sim q_{\Phi}(\mathbf{z} | \mathbf{x})} (\log(p_{\theta}(\mathbf{x} | \mathbf{z}))) \end{aligned}$$

At this point, it is possible to rewrite the equation as

$$\log(p_{\theta}(\mathbf{x})) - D_{KL}(q_{\Phi}(\mathbf{z} | \mathbf{x}) || p_{\theta}(\mathbf{z} | \mathbf{x})) = E_{\mathbf{z} \sim q_{\Phi}(\mathbf{z} | \mathbf{x})} (\log(p_{\theta}(\mathbf{x} | \mathbf{z}))) - D_{KL}(q_{\Phi}(\mathbf{z} | \mathbf{x}) || p_{\theta}(\mathbf{z}))$$

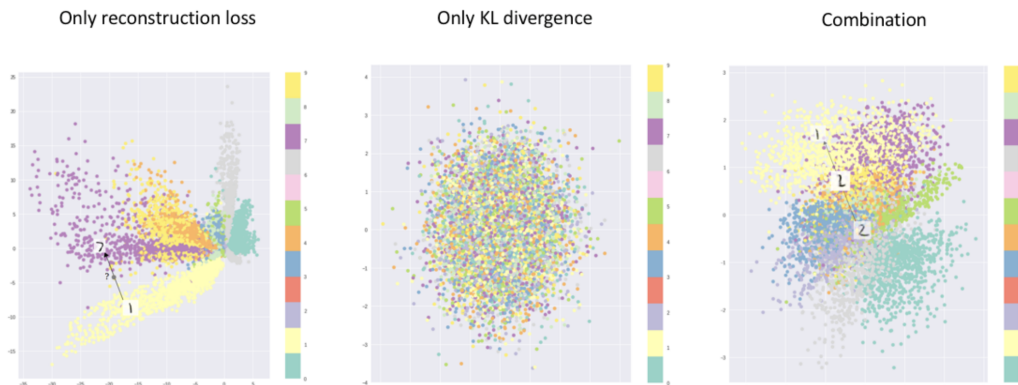
## Evidence Lower Bound (ELBO) Loss

$$L_{VAE}(\theta, \phi) = -\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z} | \mathbf{x})} \log(p_{\theta}(\mathbf{x} | \mathbf{z})) + D_{KL}(q_{\phi}(\mathbf{z} | \mathbf{x}) || p_{\theta}(\mathbf{z}))$$

- We are trying to minimize the ELBO loss with respect to the model parameters.

# Why Autoencoder?

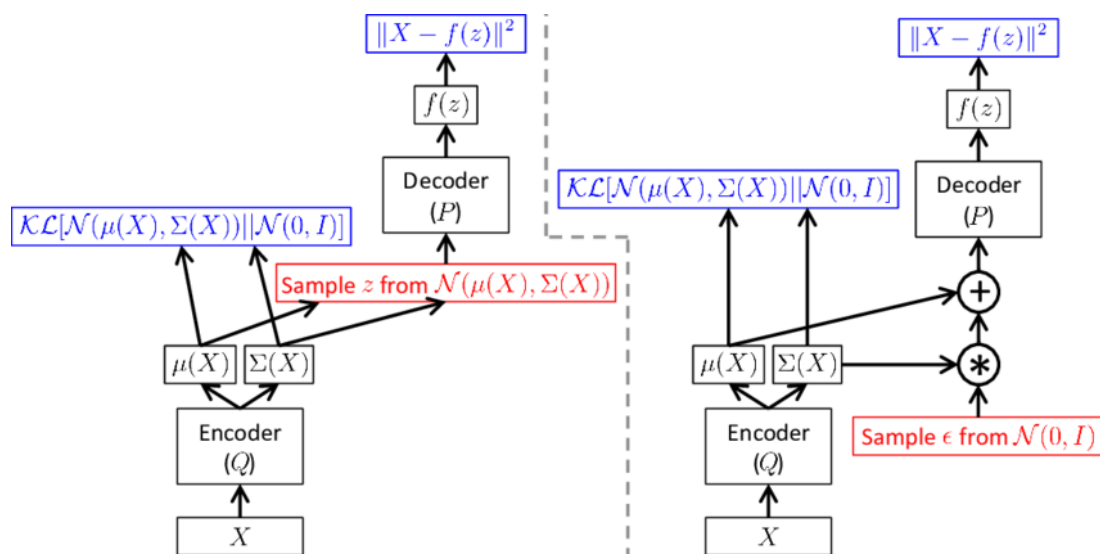
- The reconstruction term, forces each image to be unique and spread out.
- The KL term will push all the images towards the same prior.



4

<sup>4</sup>Figure taken from <https://towardsdatascience.com/intuitively-understanding-variational-autoencoders-1bfe67eb5daf>

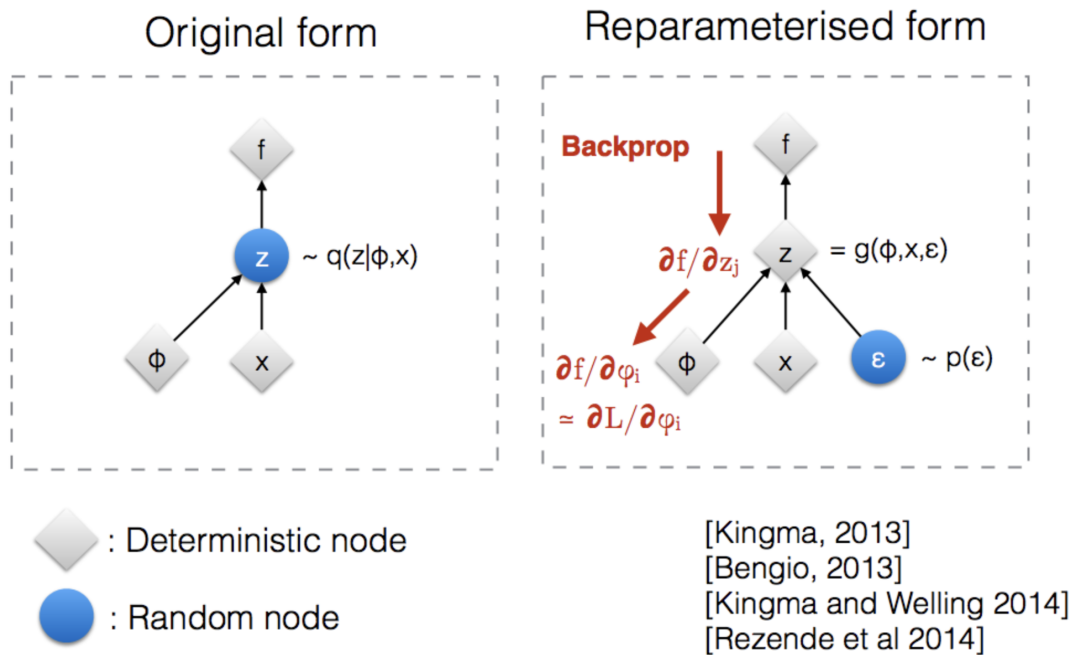
## Training Procedure



5

<sup>5</sup>Figure taken from Carl Doersch tutorial

# Reparametrization Trick Visualisation



## VAE Models and Performance

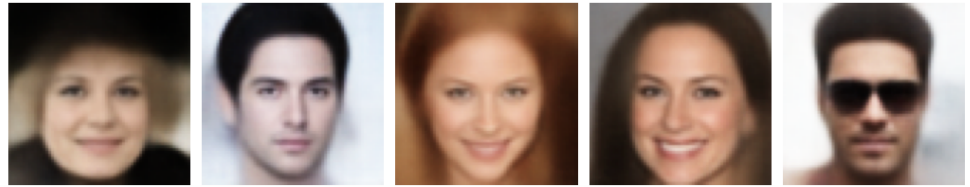
- VAEs can be used with any kind of data
  - the distributions and network architecture just needs to be set accordingly
  - e.g. it's common to use convolutions in the encoder and transpose convolutions in (Gaussian) decoder for image data
- VAEs have nice learning dynamics; they tend to be easy to optimise with stable convergence
- VAEs have a reputation for producing blurry reconstructions of images
  - Not fully understood why, but most likely related to a side effect of maximum-likelihood training
- VAEs tend to only utilise a small subset of the dimensions of  $z$

# Reconstructions Example

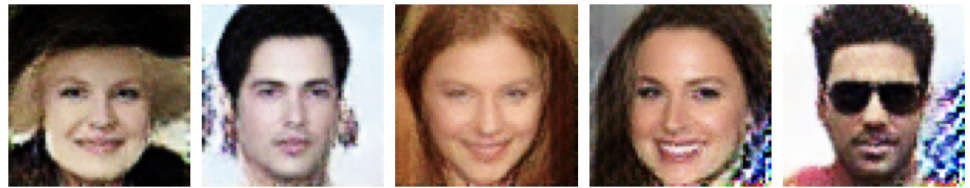
**Input**



**VAE**



**VAE<sub>Dis<sub>l</sub></sub>**



**VAE/GAN**



## Generative Adversarial Networks

# Generative Adversarial Networks (GANs)

- New (old?!<sup>6</sup>) method of training deep generative models
- Idea: pitch a generator and a discriminator against each other
  - Generator tries to draw samples from  $p(x)$
  - Discriminator tries to tell if sample came from the generator (fake) or the real world
- Both discriminator and generator are deep networks (differentiable functions)
- LeCun quote 'GANs, the most interesting idea in the last ten years in machine learning'

---

<sup>6</sup>c.f. Schmidhuber

## Aside: Adversarial Learning vs. Adversarial Examples

The approach of GANs is called adversarial since the two networks have *antagonistic* objectives.

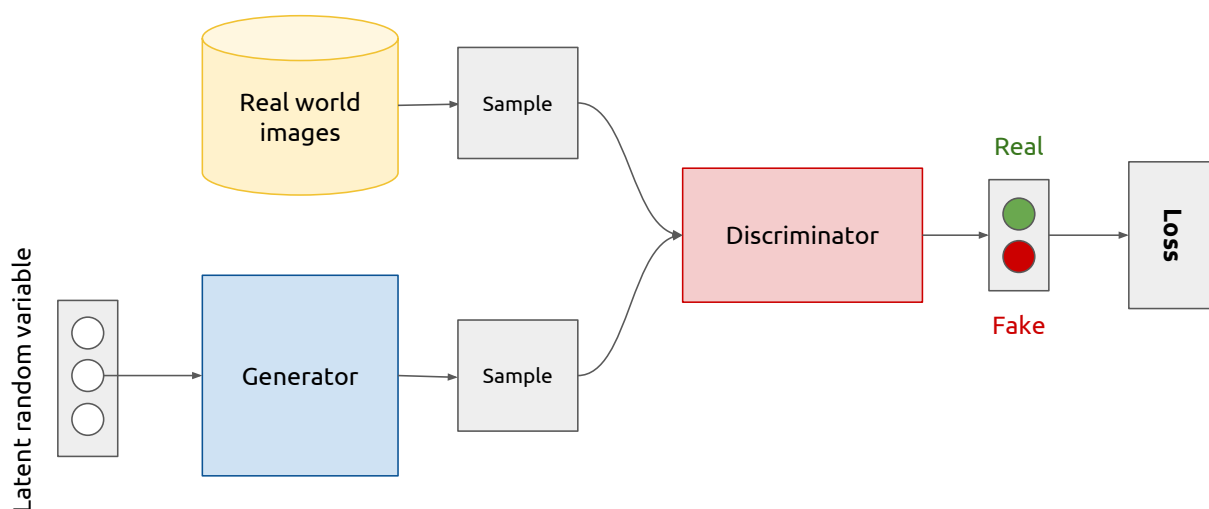
This is not to be confused with *adversarial examples* in machine learning.

See these two papers for more details:

<https://arxiv.org/pdf/1412.6572.pdf>

<https://arxiv.org/pdf/1312.6199.pdf>

# Generative adversarial networks (conceptual)



Picture Credit: Xavier Giro-i-Nieto

## More Formally

- The **generator**

$$\mathbf{x} = g(\mathbf{z})$$

is trained so that it gets a random input  $\mathbf{z} \in \mathbb{R}^n$  from a distribution (typically  $\mathcal{N}(\mathbf{0}, \mathbf{I})$  or  $\mathcal{U}(\mathbf{0}, \mathbf{I})$ ) and produces a sample  $\mathbf{x} \in \mathbb{R}^d$  following the data distribution as output (ideally). Usually  $n \ll d$ .

- The **discriminator**

$$y = d(\mathbf{x})$$

gets a sample  $\mathbf{x}$  as input and predicts a probability  $y \in [0, 1]$  (or real-valued logit of a Bernoulli distribution) determining if it is real or fake.

- Training a standard GAN is difficult and often results in two undesirable behaviours
  - Oscillations without convergence. No guarantee that the loss will actually decrease...
    - It has been shown that a GAN has saddle-point solution, rather than a local minima.
  - The **mode collapse** problem, when the generator models very well a small sub-population, concentrating on a few modes.
- Additionally, performance is hard to assess and often boils down to heuristic observations.

## Deep Convolutional Generative Adversarial Networks (DCGANs)

- Motivates the use of GANS to learn reusable feature representations from large unlabelled datasets.
- GANs known to be unstable to train, often resulting in generators that produce “nonsensical outputs”.
- Model exploration to identify architectures that result in **stable** training across datasets with higher resolution and deeper models.





- Replace pooling layers with strided convolutions in the discriminator and fractional-strided (transpose) convolutions in the generator.
  - This will allow the network to learn its own spatial downsampling.
- Use batchnorm in both the generator and the discriminator.
  - This helps deal with training problems due to poor initialisation and helps the gradient flow.
- Eliminate fully connected hidden layers for deeper architectures.
- Use ReLU activation in the generator for all layers except for the output, which uses tanh.
- Use LeakyReLU activation in the discriminator for all layers.

## Summary

- Generative modelling is a massive field with a long history
- Differentiable generators have had a profound impact in making models that work with real data at scale
- VAEs and GANs are currently the most popular approaches to training generators for spatial data
- We've only scratched the surface of generative modelling
  - Auto-regressive approaches are popular for sequences (e.g. language modelling).
    - But also for images (e.g. PixelRNN, PixelCNN)
  - typically RNN-based
  - but not necessarily - e.g. WaveNet is a convolutional auto-regressive generative model