

# Deep Learning

## Lecture 9: Generalisation theory

---

Chris G. Willcocks

Durham University



# Lecture overview

## 1 Generalisation theory

---

- universal approximation theorem
- fixed width and arbitrary depth?
- empirical risk minimisation
- noise, training set size, and regularisation
- generalisation of unsupervised models
- a cycle that amplifies biases
- no free lunch theorem and Occam's razor
- bias-variance tradeoff on nasty densities
- increasing capacity through double descent
- neural scaling laws

## 2 Modelling theory

---

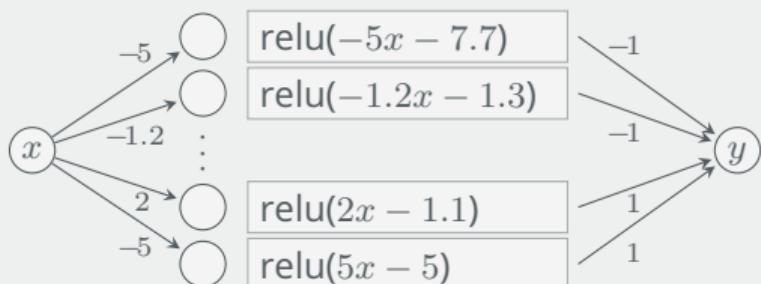
- the trilemma
- the five main modelling equations
- hybrids can satisfy multiple criteria

# Generalisation theory universal approximation theorem

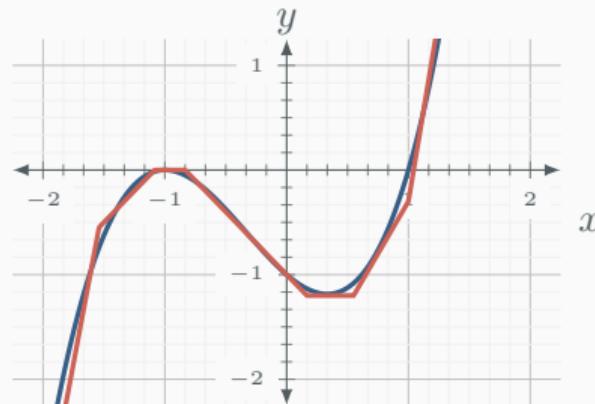
**Theorem:** universal function approximator

## Arbitrary width

A network with a single hidden layer, containing a finite number of neurons, can approximate any continuous function under mild assumptions.



original function  $f(x) = x^3 + x^2 - x - 1$



$$f(x) = -r(-5x - 7.7) - r(-1.2x - 1.3) - r(1.2x + 1) + r(1.2x - 0.2) + r(2x - 1.1) + r(5x - 5)$$

Example ReLU weights by Brendan Fortuner

# Generalisation theory universal approximation theorem

**Theorem:** universal function approximator

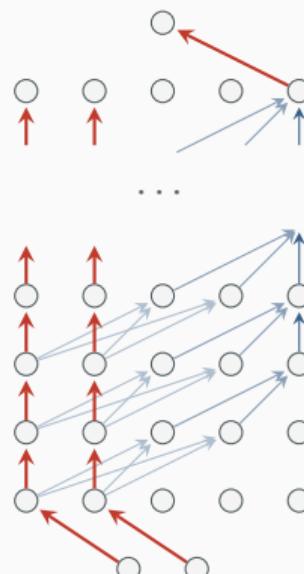
## Arbitrary depth (fixed width)

Does the theorem still hold for fixed width and arbitrary depth? **Yes!**

For a network of  $n$  inputs and  $m$  outputs, [1] show universal approximator holds true for:

- width  $n + m + 2$  for almost any activation function
- width  $n + m + 1$  for most activation functions

Short YouTube visual proof 



# Generalisation theory empirical risk minimisation

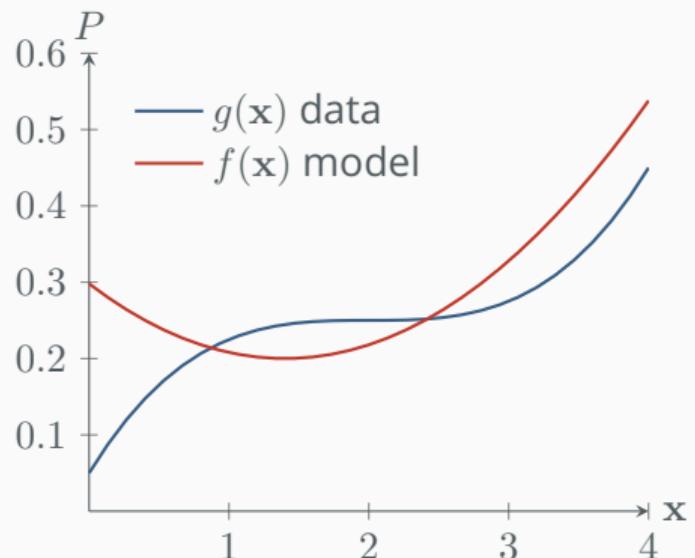
## Finding a model

Lets illustrate this in 1D, but try to imagine it in  $ND$ . Given the target probability distribution of our data

$$g(\mathbf{x}) = P(X, Y)$$

we want to find (design) a model  $f(\mathbf{x}; \theta)$  with parameters  $\theta$  and optimise  $\theta$  such that

$$f(\mathbf{x}; \theta) = g(\mathbf{x})$$



# Generalisation theory empirical risk minimisation

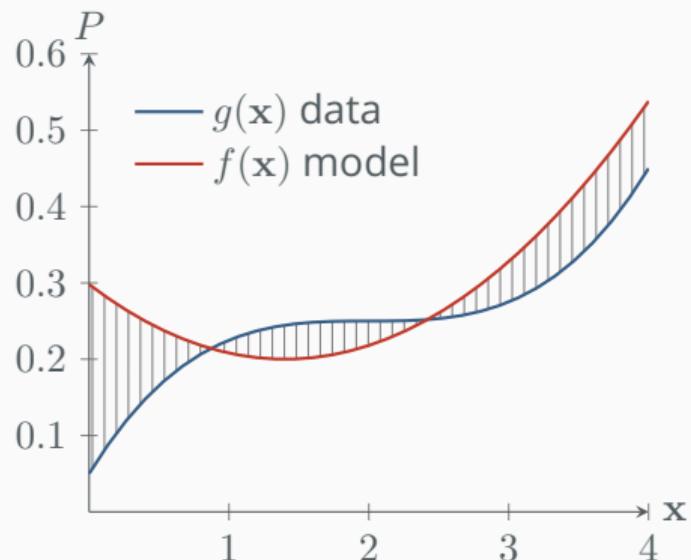
## Optimising the model

The total error is therefore the entire area. Modifying the parameters  $\theta$  will cause the area to change, where we want to find

$$\hat{\theta} = \arg \min_{\theta} \int \mathcal{L}(f(\mathbf{x}; \theta), g(\mathbf{x})) d\mathbf{x}.$$

where  $\mathcal{L}$  is a 'loss function', e.g. a 0-1 loss function  $\mathcal{L}(\hat{x}, x) = \mathbb{I}(\hat{x} \neq x)$  or a mean squared error loss.

The solution is a function  $f$  that has the capacity to exactly represent  $g(\mathbf{x})$



# Generalisation theory empirical risk minimisation

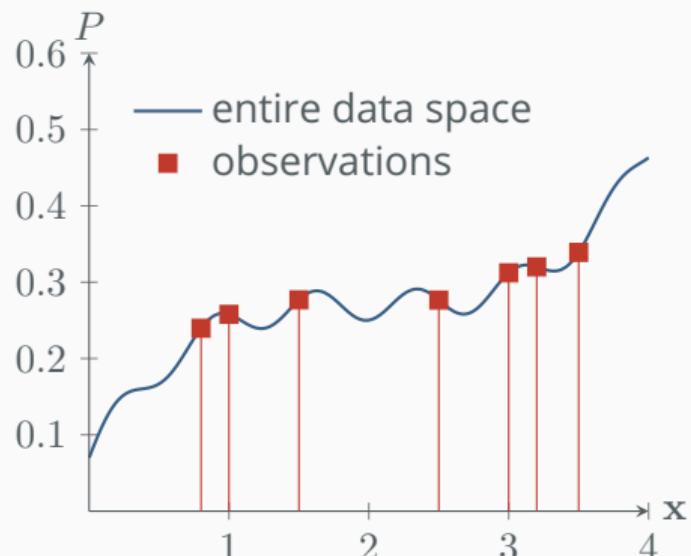
## The generalisation problem

However there's a big problem:

**In practice, we can't observe all of  $g(x)$**

This means:

1. We don't know how smooth the function is between the observations
2. Noise can be difficult to interpret
3. Optimisation is highly sensitive to the sampling process



# Generalisation theory empirical risk minimisation

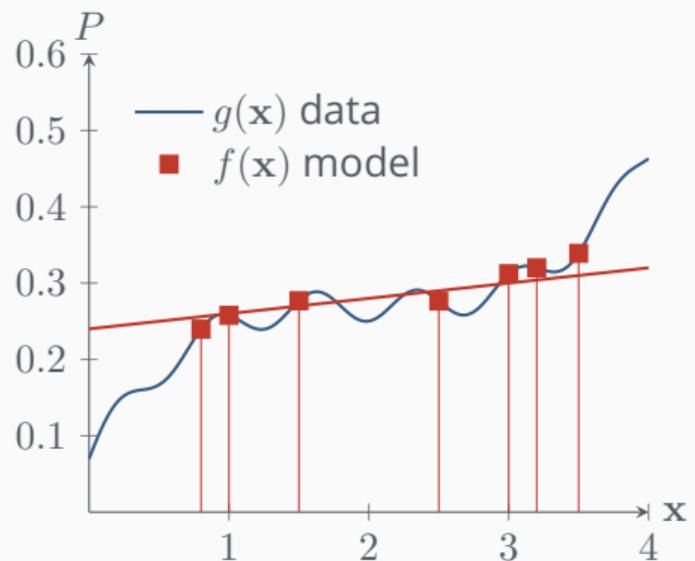
## Definition: empirical risk

The expected error (or risk) is the average error over the entire space, which we can't compute:

$$\mathbb{E}[\mathcal{L}(f(\mathbf{x}; \theta), g(\mathbf{x}))] = \int \mathcal{L}(f(\mathbf{x}; \theta), g(\mathbf{x})) d\mathbf{x}.$$

Therefore we minimise the **empirical estimate** of the risk as an average over the samples:

$$\mathbb{E}[\mathcal{L}(f(\mathbf{x}; \theta), g(\mathbf{x}))] \approx \frac{1}{N} \sum_{i=1}^N \mathcal{L}(f(\mathbf{x}_i; \theta), g(\mathbf{x}_i)).$$

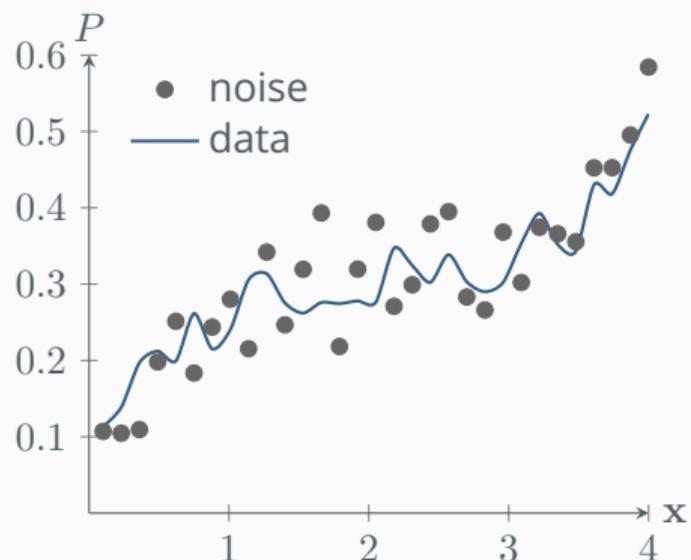


# Generalisation theory empirical risk minimisation

## Noise and regularisation

Given that the shape of the distribution outside of the observations is unknown, it is easy to overfit to noise, especially when you have limited data.

We can often regularise our function to be invariant to this.

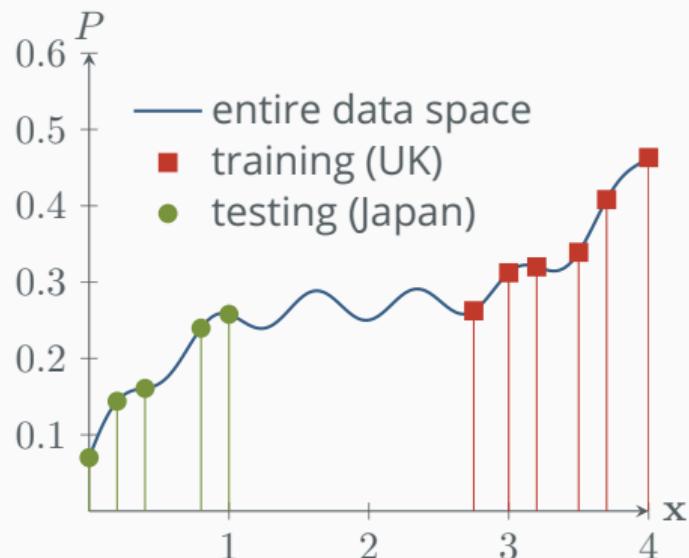


# Generalisation theory empirical risk minimisation

## Out-of-distribution data

Usually the training dataset collection process draws samples from the data space in a way that is not **independent and identically distributed** (abbreviated i.i.d.) to the expected testing (operational) conditions of the model.

Sampling data in a way that is representative of the task/testing/operational distribution is extraordinarily difficult to do properly. It is often a worthwhile investment though!



# Generalisation theory generalisation of unsupervised models

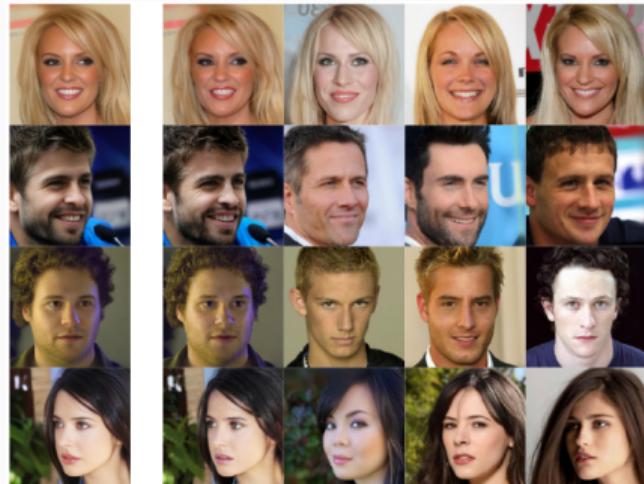
## Similarity measures

How can we evaluate the generalisation of unsupervised generative models?

Deep pre-trained network activations work surprisingly well (e.g. FID, LPIPS [2] for perceptual similarity between images):

```
import lpips  
distance_fn = lpips.LPIPS(net="alex")
```

## Example code usage ↗



Model samples (leftmost column) and nearest neighbours in training set (LPIPS distance [2]) increasing to right.

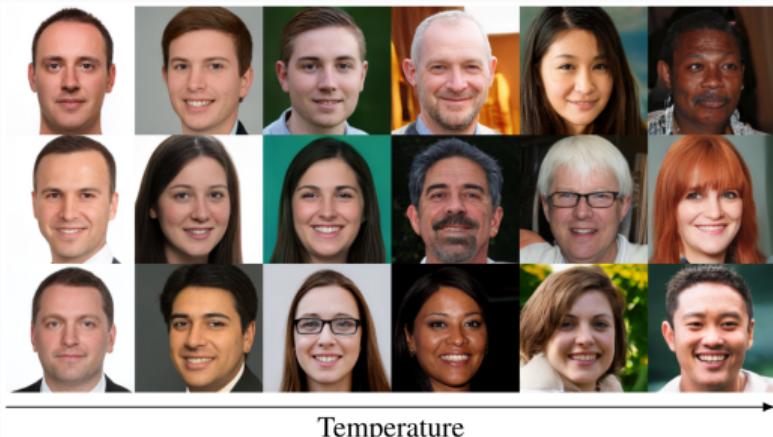


# Generalisation theory a cycle that amplifies biases

## Bias amplification

Judging models solely based on pre-trained measures such as FID can amplify the inherited biases.

FID (currently how the majority of image-based generative model papers are ranked) can favour less diversity (with focus shifted towards the distribution mode).



Sampling away from the mode of the data distribution

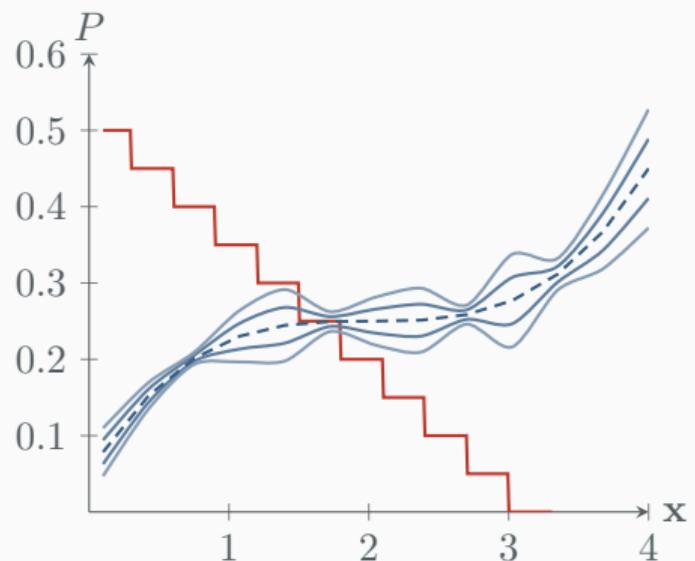
# Generalisation theory no free lunch theorem and Occam's razor

## Definition: no free lunch theorem

We can use methods such as cross validation to empirically choose the best method for our particular problem. However, there is no universally best model — this is sometimes called the no free lunch theorem [3].

## Definition: Occam's razor

'Prefer the simplest hypothesis  $\mathcal{H}$  that fits the data.' In the case of deep learning, this implies the smoothest function that fits the data.



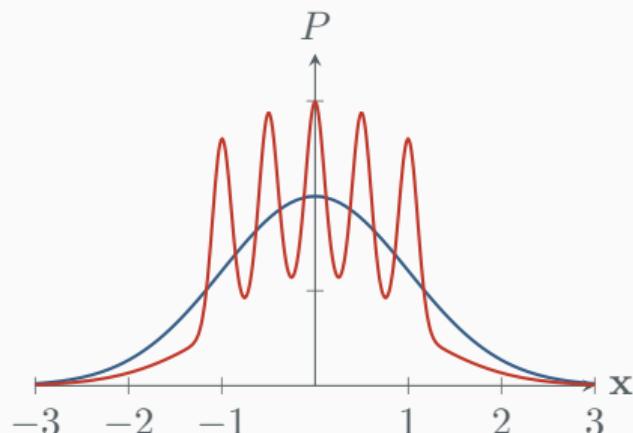
# Generalisation theory bias-variance tradeoff and nasty densities

## Example: problematic densities

In a more complex setting, we may have a density such as the 'Bart Simpson' density, as Nando de Freitas likes to call it

$$p(x) = \frac{1}{2}\phi(x; 0, 1) + \frac{1}{10} \sum_{j=0}^4 \phi(x; (j/2) - 1, 1/10)$$

where  $\phi$  is the normal density with mean  $\mu$  and standard deviation  $\sigma$ . This density cannot be sufficiently estimated with a normal distribution, as the result is over-smoothed (blue).

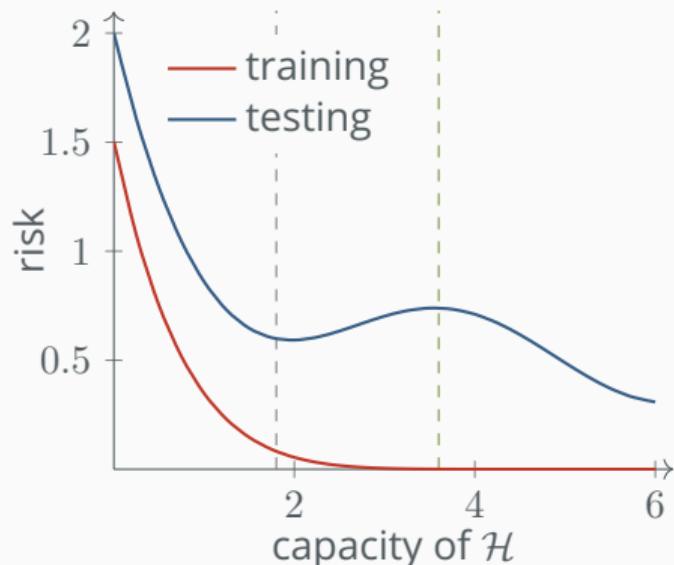


# Generalisation theory increasing capacity to double descent

## Definition: double descent

Traditionally, we know that increasing the parameters lowers the bias (fitting), but the variance (test risk) will eventually reach a 'sweet spot' (first dashed line) and start to increase again.

The full story has a double descent curve [4], as higher capacity functions past the interpolation threshold (second dashed line) lead again to smoother fitting (Occam's razor).





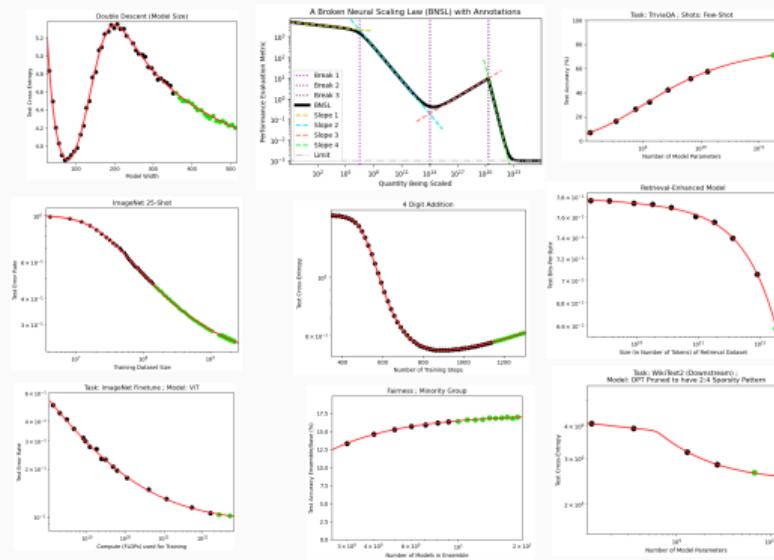
# Generalisation theory neural scaling laws

## Neural scaling laws

Can we predict all this?

- Yes, yes and yes! Extremely well!  
(Images, LLMs, RL, Bias, ...)

Read David Krueger's "Broken neural scaling laws" [5], ICLR 2023.  
Extrapolates many different criteria when considering this behaviour. See also his unifying grokking and double descent paper [6].

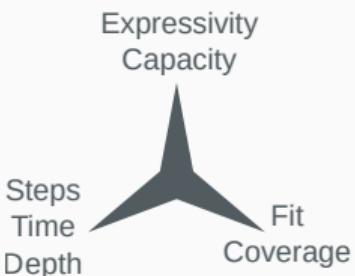


BNSL predictions (extrapolations) [5]

# Modelling theory the trilemma

## The trilemma

- The five main modelling approaches  $p(\mathbf{x}) \approx ?$
- Each vary in
  - Mode coverage
  - Expressivity
  - Time or depth



## Review

| Method                              | Train Speed | Sample Speed | Num. Params. | Resolution Scaling | Free-form Jacobian | Exact Density | FID    | NLL (in BPD)   |
|-------------------------------------|-------------|--------------|--------------|--------------------|--------------------|---------------|--------|----------------|
| Generative Adversarial Networks     |             |              |              |                    |                    |               |        |                |
| DCGAN [182]                         | *****       | *****        | *****        | ****               | ✓                  | ✗             | 37.11  | -              |
| ProGAN [114]                        | *****       | *****        | *****        | ****               | ✓                  | ✗             | 15.52  | -              |
| BigGAN [19]                         | *****       | *****        | *****        | ****               | ✓                  | ✗             | 14.73  | -              |
| StyleGAN2 + ADA [115]               | ****        | *****        | *****        | *****              | ✓                  | ✗             | 2.42   | -              |
| Energy Based Models                 |             |              |              |                    |                    |               |        |                |
| IGEBM [46]                          | *****       | *****        | *****        | ****               | ✓                  | ✗             | 37.9   | -              |
| Denoising Diffusion [87]            | *****       | *****        | *****        | ****               | ✓                  | (✓)           | 3.17   | $\leq 3.75$    |
| DDPM++ Continuous [206]             | *****       | *****        | *****        | ****               | ✓                  | (✓)           | 2.20   | -              |
| Flow Contrastive (EBM) [55]         | *****       | *****        | *****        | ****               | ✓                  | ✗             | 37.30  | $\approx 3.27$ |
| VAEBM [247]                         | *****       | *****        | *****        | ****               | ✓                  | ✗             | 12.19  | -              |
| Variational Autoencoders            |             |              |              |                    |                    |               |        |                |
| Convolutional VAE [123]             | *****       | *****        | *****        | ****               | ✓                  | (✓)           | 106.37 | $\leq 4.54$    |
| Variational Lossy AE [29]           | *****       | *****        | *****        | ****               | ✗                  | (✓)           | -      | $\leq 2.95$    |
| VQ-VAE [184], [235]                 | ***         | *****        | *****        | *****              | ✗                  | (✓)           | -      | $\leq 4.67$    |
| VD-VAE [31]                         | *****       | *****        | *****        | ****               | ✓                  | (✓)           | -      | $\leq 2.87$    |
| Autoregressive Models               |             |              |              |                    |                    |               |        |                |
| PixelRNN [234]                      | ****        | ****         | ****         | ****               | ✗                  | ✓             | -      | 3.00           |
| Gated PixelCNN [233]                | ****        | ****         | ****         | ****               | ✗                  | ✓             | 65.93  | 3.03           |
| PixelIQN [173]                      | ****        | ****         | ****         | ****               | ✗                  | ✓             | 49.46  | -              |
| Sparse Trans. + DistAug [32], [110] | *****       | ****         | ****         | ****               | ✗                  | ✓             | 14.74  | 2.66           |
| Normalizing Flows                   |             |              |              |                    |                    |               |        |                |
| RealNVP [43]                        | *****       | *****        | *****        | ****               | ✗                  | ✓             | -      | 3.49           |
| GLOW [124]                          | *****       | *****        | *****        | ****               | ✗                  | ✓             | 45.99  | 3.35           |
| FFJORD [62]                         | ****        | ****         | ****         | ****               | ✓                  | (✓)           | -      | 3.40           |
| Residual Flow [26]                  | *****       | *****        | *****        | ****               | ✓                  | (✓)           | 46.37  | 3.28           |

*Bond-Taylor, Willcocks et al., "Deep Generative Modelling: A Comparative Review of VAEs, GANs, Normalizing Flows, Energy-Based and Autoregressive Models" in IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI) [7]*

# Modelling theory the five main modelling equations

## Summary

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

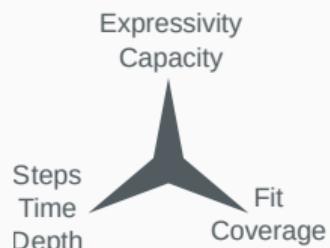
$$p(\mathbf{x}) \approx \frac{e^{-E(\mathbf{x})}}{\int_{\tilde{\mathbf{x}} \in \mathcal{X}} e^{-E(\tilde{\mathbf{x}})}} \text{ e.g. } \nabla_{\mathbf{x}} \log p(\mathbf{x})$$

$$\log p(\mathbf{x}) \geq \mathcal{L}_{\text{recon}}^{\text{pixel}} - D_{\text{KL}}[q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})]$$

$$\log p(\mathbf{x}) \neq \log D(\mathbf{x}) \quad (\text{in GAN})$$

$$p(\mathbf{x}) = p_Z(f^{-1}(\mathbf{x})) \left| \det \left( \frac{\partial f^{-1}(\mathbf{x})}{\partial \mathbf{x}} \right) \right|$$

See our hybrids “Unleashing transformers” [8] ↗ (ECCV22) & “Megapixel image generation” [9].





# Hybrids [9] 2 seconds generation, 2 days training, single GTX 1080Ti





# Take Away Points

## Summary

In summary:

- deep learning isn't magic
- there's lots of formal results [10] (also see new book "understanding deep learning")
- understand the various trade-offs (bias-variance, linear transformers, modelling trilemma)
- categorise the contributions into the three spaces (data, function, modelling) at the end of lecture 1
- ...predictions can be made well by the *very few* researchers who make predictions...



# References I

- [1] Patrick Kidger and Terry Lyons. "Universal approximation with deep narrow networks". In: Conference on Learning Theory. 2020, pp. 2306–2327.
- [2] Richard Zhang et al. "The unreasonable effectiveness of deep features as a perceptual metric". In: Proceedings of the IEEE conference on computer vision and pattern recognition. 2018, pp. 586–595.
- [3] Kevin P Murphy. Machine learning: a probabilistic perspective. MIT press, 2012.
- [4] Mikhail Belkin et al. "Reconciling modern machine-learning practice and the classical bias–variance trade-off". In: Proceedings of the National Academy of Sciences 116.32 (2019), pp. 15849–15854.
- [5] Ethan Caballero et al. "Broken Neural Scaling Laws". In: The Eleventh International Conference on Learning Representations. 2023.
- [6] Xander Davies, Lauro Langosco, and David Krueger. "Unifying Grokking and Double Descent". In: arXiv preprint arXiv:2303.06173 (2023).



## References II

- [7] Sam Bond-Taylor et al. "Deep Generative Modelling: A Comparative Review of VAEs, GANs, Normalizing Flows, Energy-Based and Autoregressive Models". In: IEEE Transactions on Pattern Analysis and Machine Intelligence (2021), pp. 1–1. DOI: 10.1109/TPAMI.2021.3116668.
- [8] Sam Bond-Taylor et al. "Unleashing Transformers: Parallel Token Prediction with Discrete Absorbing Diffusion for Fast High-Resolution Image Generation from Vector-Quantized Codes". In: European Conference on Computer Vision (ECCV) (2022). DOI: 10.1007/978-3-031-20050-2\_11.
- [9] Alex F McKinney and Chris G Willcocks. "Megapixel Image Generation with Step-Unrolled Denoising Autoencoders". In: arXiv preprint arXiv:2206.12351 (2022).
- [10] Ovidiu Calin. Deep learning architectures: a mathematical approach. Springer, 2020.