

# “Creative” Neural Networks

The Hitchhiker’s Guide to ~~Deepfakes~~ *Generative Modelling*

Artem Maevskiy

Research Fellow – Institute for Functional Intelligent  
Materials @ National University of Singapore



MLHEP-2023

Erice, Apr 11 – 18, 2023



Institute for Functional  
Intelligent Materials



National Institute for Nuclear Physics

Practicum



# Outline

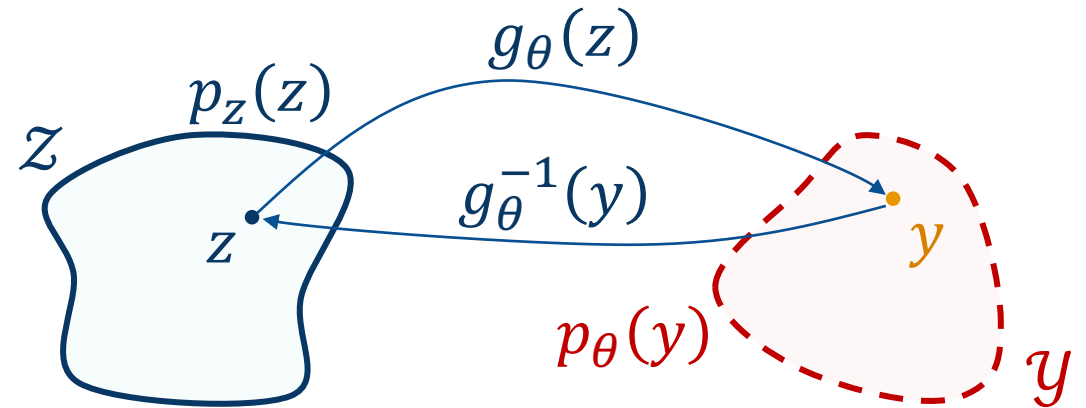
- Intro:
  - Generative modelling: what it is and how it's useful
  - Generative modelling: how it's done
- Generative Adversarial Networks (GAN)
- Variational Autoencoders (VAE)
- Normalizing Flows (NF)
- Evaluating Generative Models

# Normalizing Flows

# Basic idea

- Similarly to GAN and VAE, start with simple distribution for the latent codes  $z \sim p_z$
- Similarly to GAN, map  $z \rightarrow y$  deterministically
- Distinguishing feature: choose mapping  $g_\theta(z)$  to be invertible

$$z = g_\theta^{-1}(y)$$



# Simple max likelihood training

- Define  $f_{\theta}(y) \equiv g_{\theta}^{-1}(y)$

$$\begin{aligned}\log p_{\theta}(y) &= \log \left[ \left| \det \frac{\partial f_{\theta}}{\partial y} \right| \cdot p_z(f_{\theta}(y)) \right] \\ &= \log \left| \det \frac{\partial f_{\theta}}{\partial y} \right| + \log p_z(f_{\theta}(y)) \rightarrow \max_{\theta}\end{aligned}$$

- For a sequence of transformations  $f_{\theta}(y) = f_{\theta}^n \left( \dots f_{\theta}^2 \left( f_{\theta}^1(y) \right) \dots \right)$ :

$$\log p_{\theta}(y) = \sum_i \log \left| \det \frac{\partial f_{\theta}^i}{\partial f_{\theta}^{i-1}} \right| + \log p_z(f_{\theta}(y)) \rightarrow \max_{\theta}$$

# Practical considerations

- Our mappings should:
  - Be invertible
    - This means the dimensionalities of  $z$  and  $y$  are same
  - Be sufficiently expressive
  - Be computationally efficient (both functions and Jacobian determinant)

# Simplest case

- Element-wise bijective function

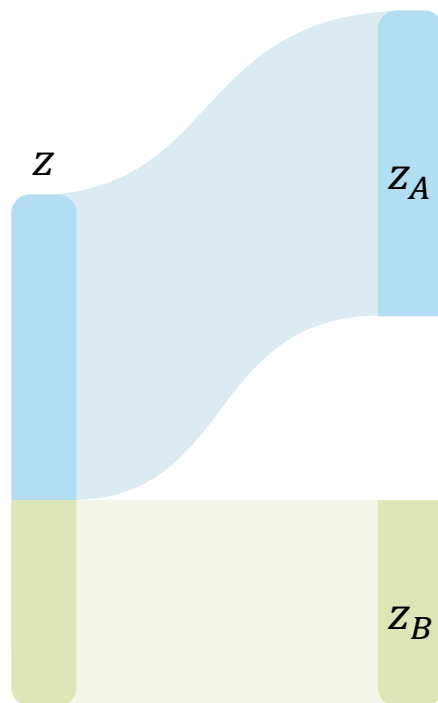
# Linear flows

- For any invertible matrix  $\mathbf{A}$ , function  $f(\mathbf{y}) = \mathbf{A}\mathbf{y} + \mathbf{b}$  is invertible too, e.g.:
  - Triangular
    - Possibly with permutations  $\leq$  these can't be trained though
    - Inversion is  $\sim \mathcal{O}(d^2)$
  - Orthogonal
    - Parameterized as product of reflections  $\mathbf{H} = \mathbb{I} - \frac{2}{\|\mathbf{v}\|^2} \mathbf{v}\mathbf{v}^T$



# Coupling flows

$$y = g(z), \quad g:$$
$$z \rightarrow (z_A, z_B)$$

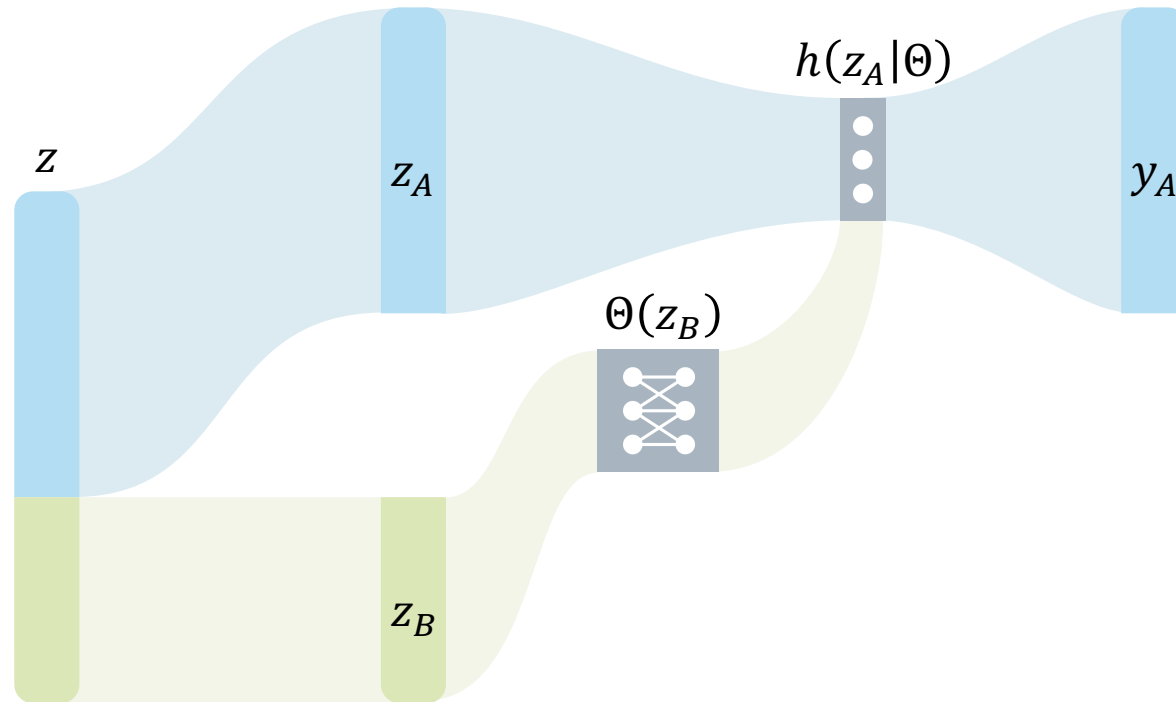


# Coupling flows

$$y = g(z), \quad g:$$

$$z \rightarrow (z_A, z_B)$$

$$y_A = h(z_A | \Theta(z_B))$$



# Coupling flows

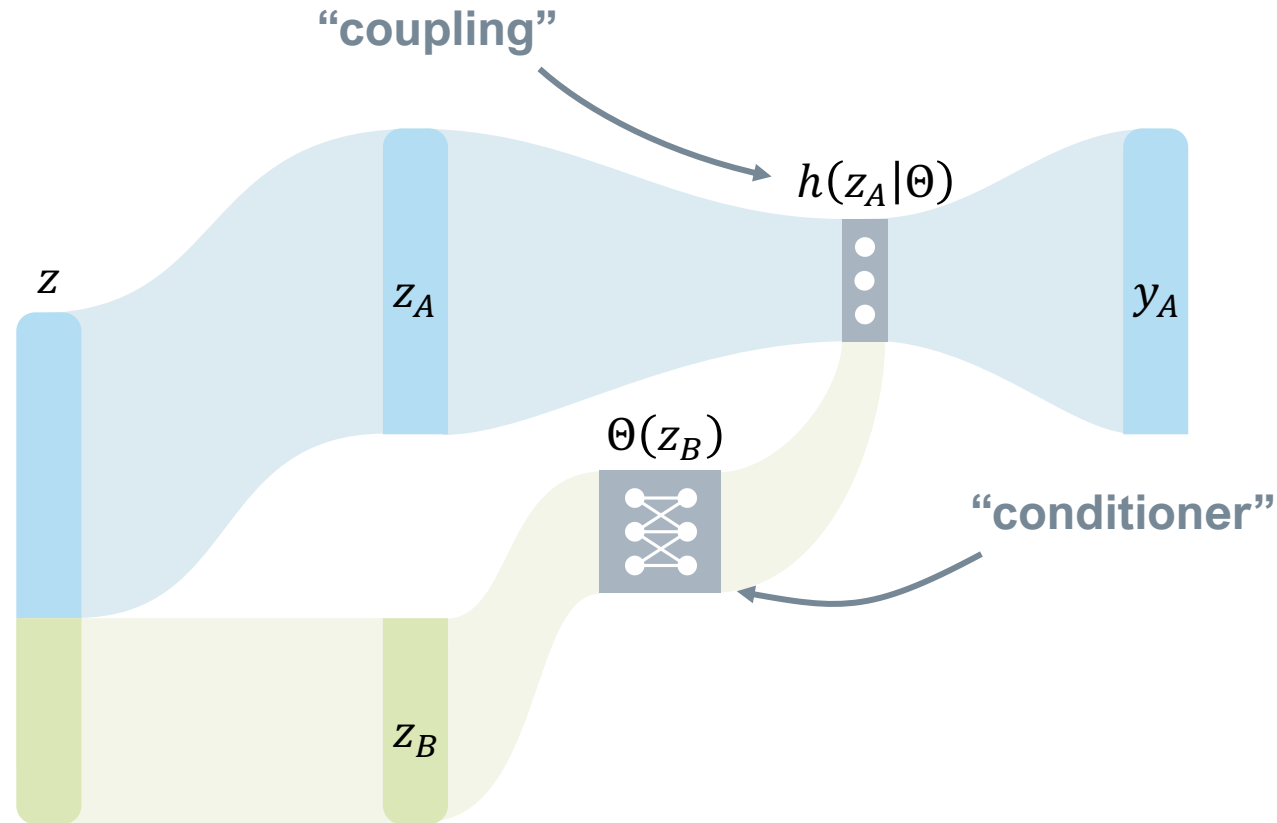
$$y = g(z), \quad g:$$

$$z \rightarrow (z_A, z_B)$$

$$y_A = h(z_A | \Theta(z_B))$$

$h(z_A | \Theta)$  – some  
function, invertible  
wrt  $z_A$

$\Theta(z_B)$  – any function (e.g., neural net)



# Coupling flows

$$y = g(z), \quad g:$$

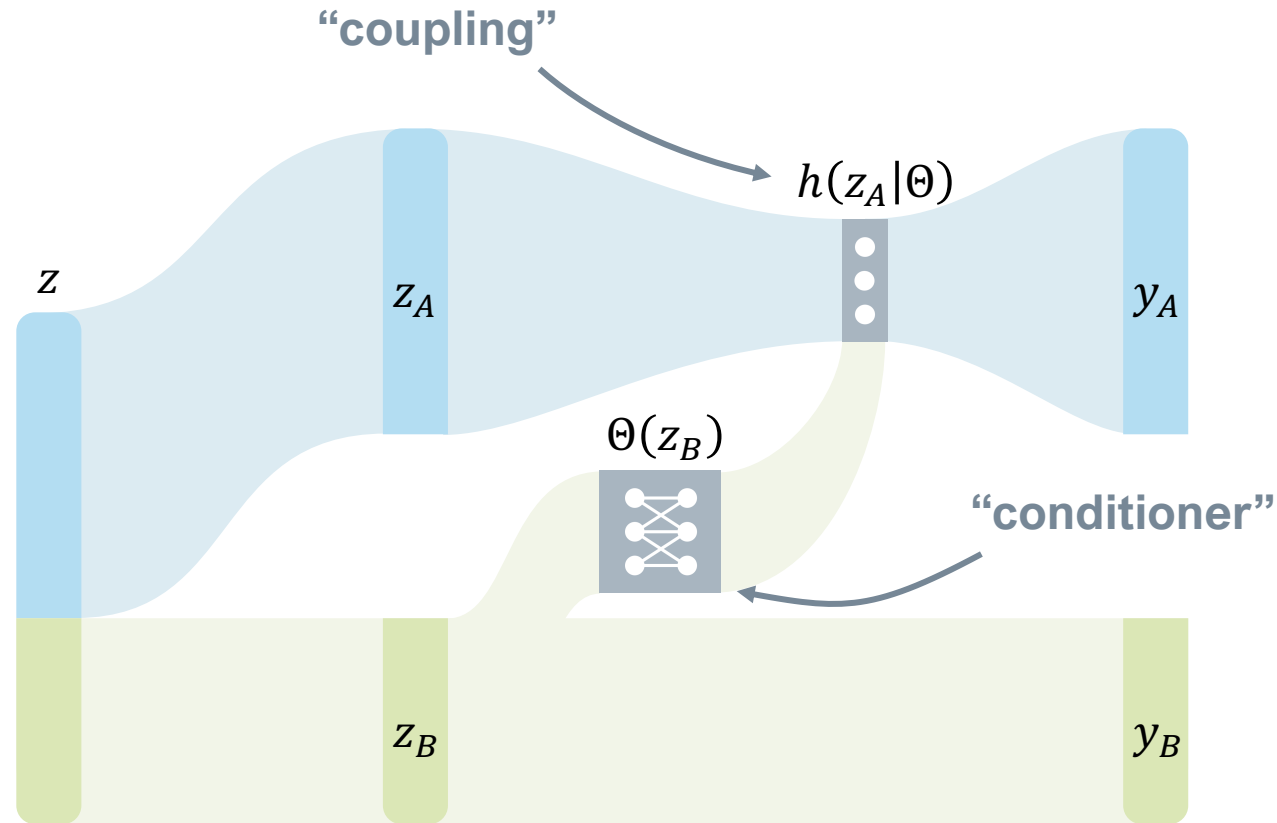
$$z \rightarrow (z_A, z_B)$$

$$y_A = h(z_A | \Theta(z_B))$$

$$y_B = z_B$$

$h(z_A | \Theta)$  – some  
function, invertible  
wrt  $z_A$

$\Theta(z_B)$  – any function (e.g., neural net)



# Coupling flows

$$y = g(z), \quad g:$$

$$z \rightarrow (z_A, z_B)$$

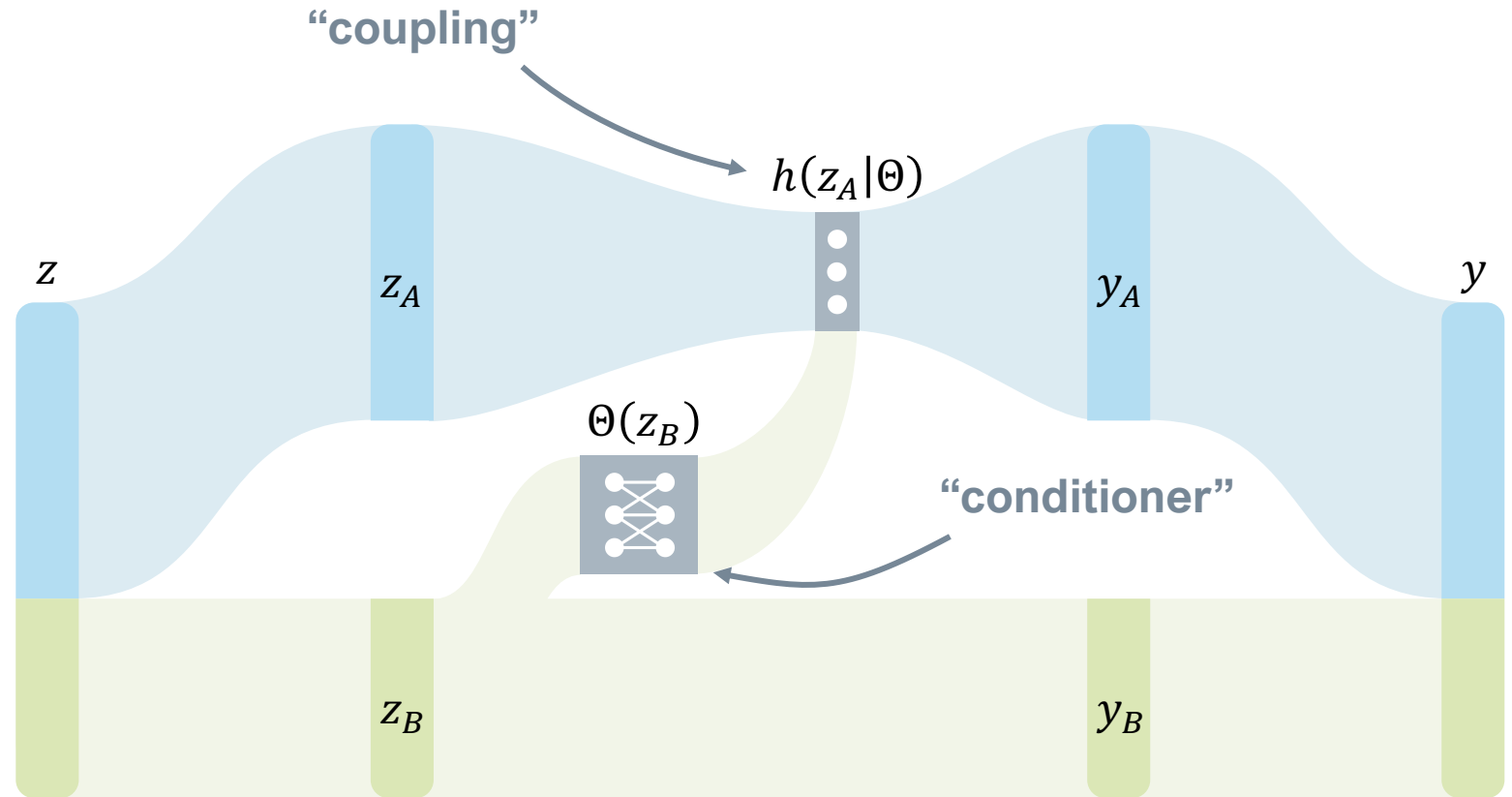
$$y_A = h(z_A | \Theta(z_B))$$

$$y_B = z_B$$

$$(y_A, y_B) \rightarrow y$$

$h(z_A | \Theta)$  – some  
function, invertible  
wrt  $z_A$

$\Theta(z_B)$  – any function (e.g., neural net)



# Coupling flows

$$y = g(z), \quad g:$$

$$z \rightarrow (z_A, z_B)$$

$$y_A = h(z_A | \Theta(z_B))$$

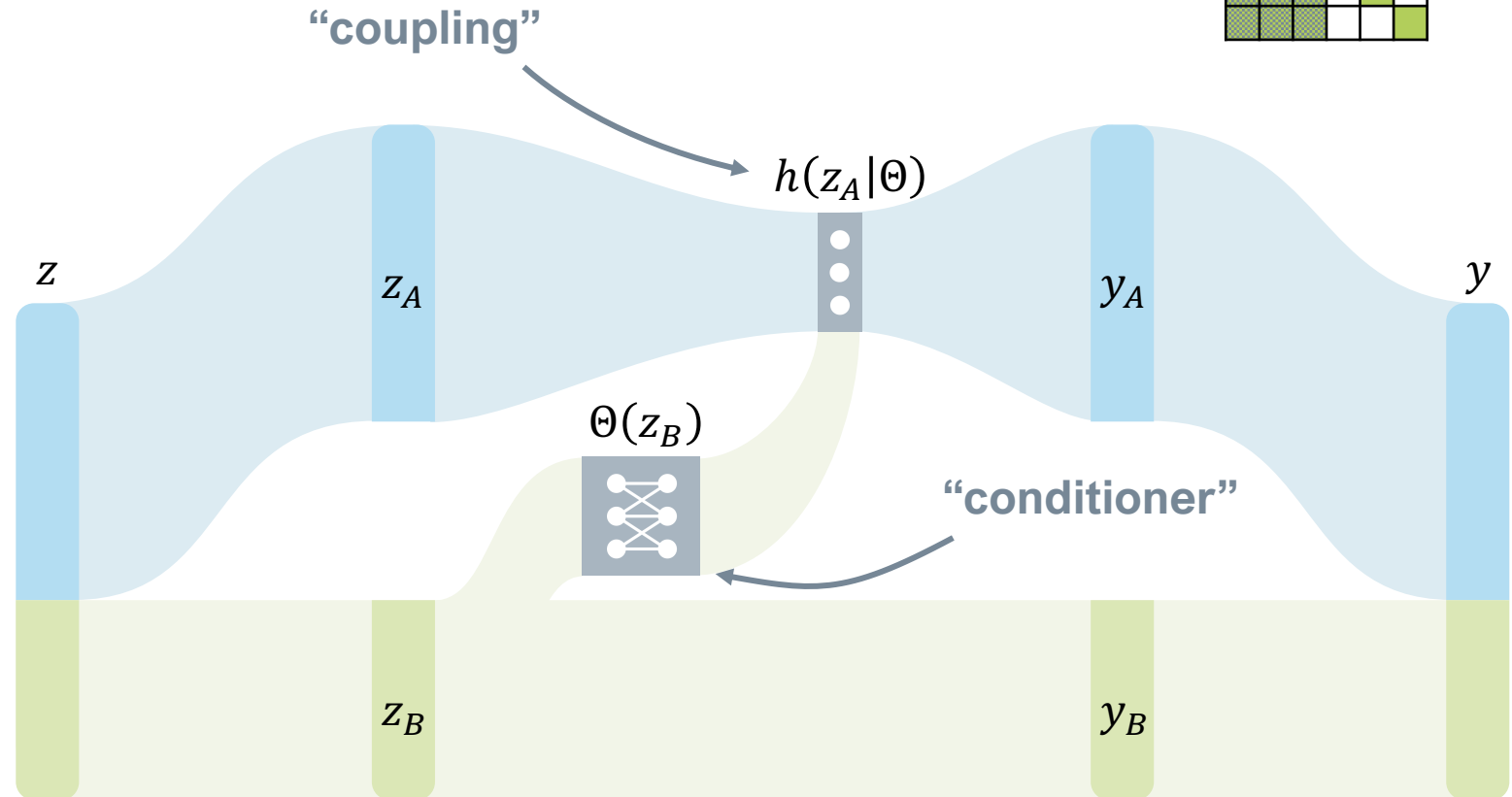
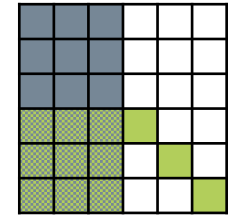
$$y_B = z_B$$

$$(y_A, y_B) \rightarrow y$$

$h(z_A | \Theta)$  – some  
function, invertible  
wrt  $z_A$

$\Theta(z_B)$  – any function (e.g., neural net)

Jacobian structure:



# Particular choices

- NICE (Non-linear Independent Components Estimation):

$$h(z_A|\Theta) = z_A + \Theta$$

- RealNVP (real-valued non-volume preserving transformations):

$$h(z_A|\Theta) = z_A \odot \exp \Theta_s + \Theta_t$$

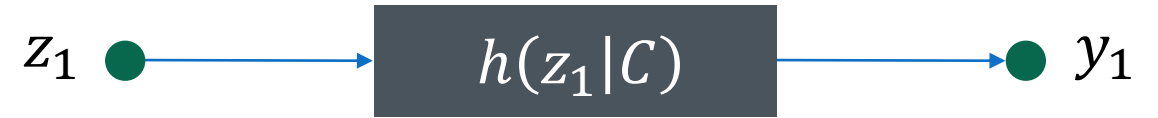
# Autoregressive flows

- $z = (z_1, \dots, z_D)^T \rightarrow y = (y_1, \dots, y_D)^T$
- $h(\cdot | \theta): \mathbb{R} \rightarrow \mathbb{R}$  – invertible function
- $y_i = h(z_i | \theta(z_{1:i-1}))$



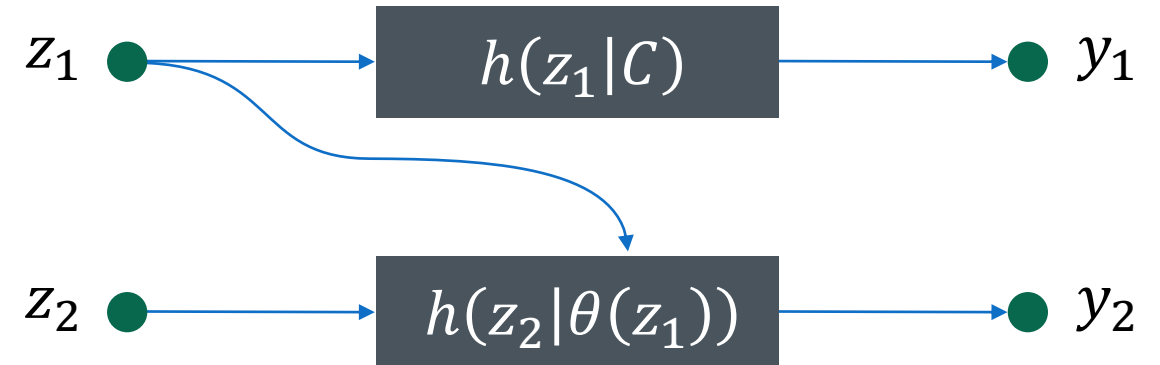
# Autoregressive flows

- $z = (z_1, \dots, z_D)^T \rightarrow y = (y_1, \dots, y_D)^T$
- $h(\cdot | \theta): \mathbb{R} \rightarrow \mathbb{R}$  – invertible function
- $y_i = h(z_i | \theta(z_{1:i-1}))$



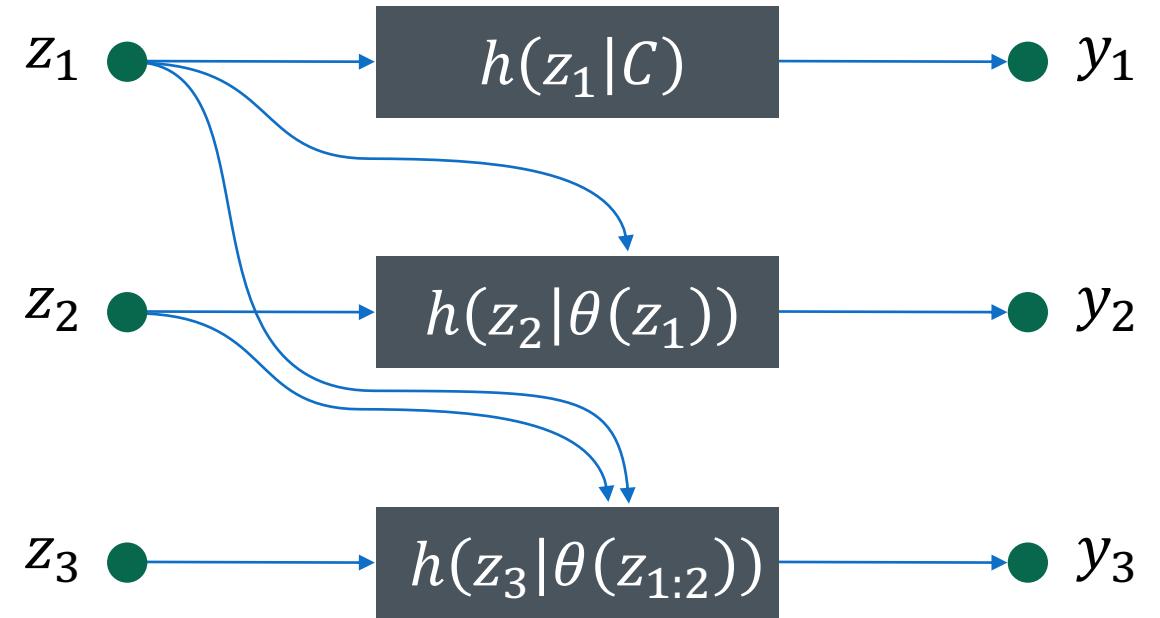
# Autoregressive flows

- $z = (z_1, \dots, z_D)^T \rightarrow y = (y_1, \dots, y_D)^T$
- $h(\cdot | \theta): \mathbb{R} \rightarrow \mathbb{R}$  – invertible function
- $y_i = h(z_i | \theta(z_{1:i-1}))$



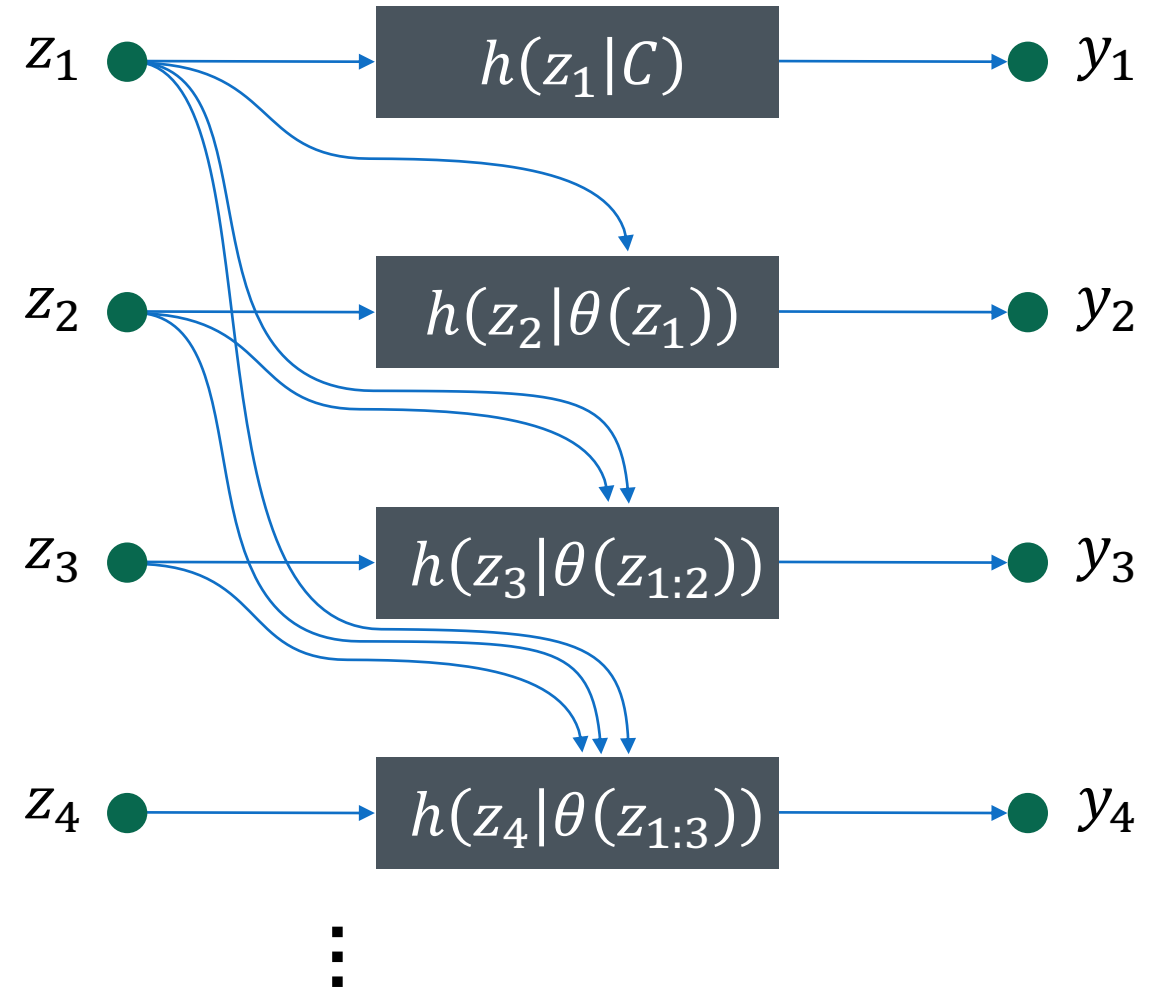
# Autoregressive flows

- $z = (z_1, \dots, z_D)^T \rightarrow y = (y_1, \dots, y_D)^T$
- $h(\cdot | \theta): \mathbb{R} \rightarrow \mathbb{R}$  – invertible function
- $y_i = h(z_i | \theta(z_{1:i-1}))$



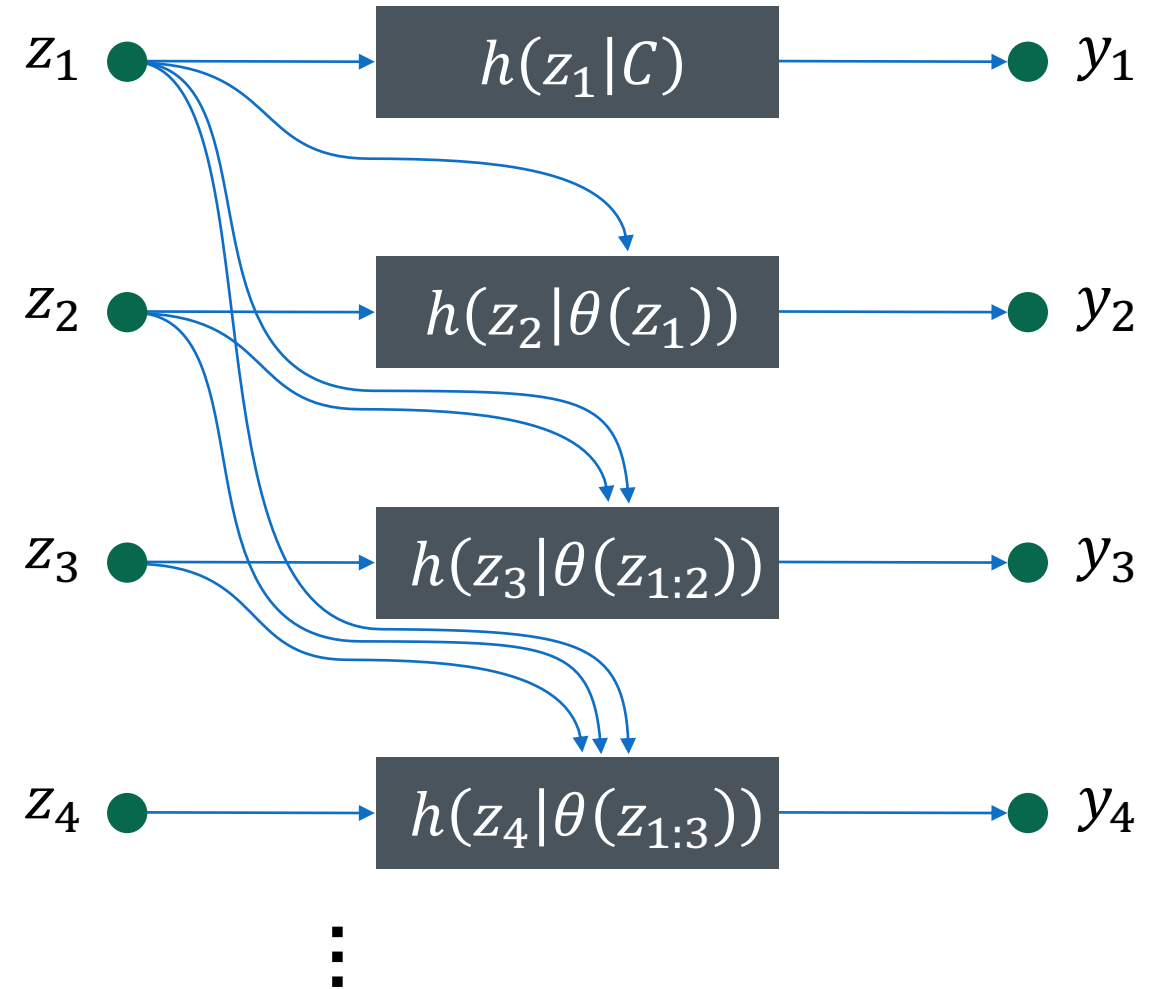
# Autoregressive flows

- $z = (z_1, \dots, z_D)^T \rightarrow y = (y_1, \dots, y_D)^T$
- $h(\cdot | \theta): \mathbb{R} \rightarrow \mathbb{R}$  – invertible function
- $y_i = h(z_i | \theta(z_{1:i-1}))$



# Autoregressive flows

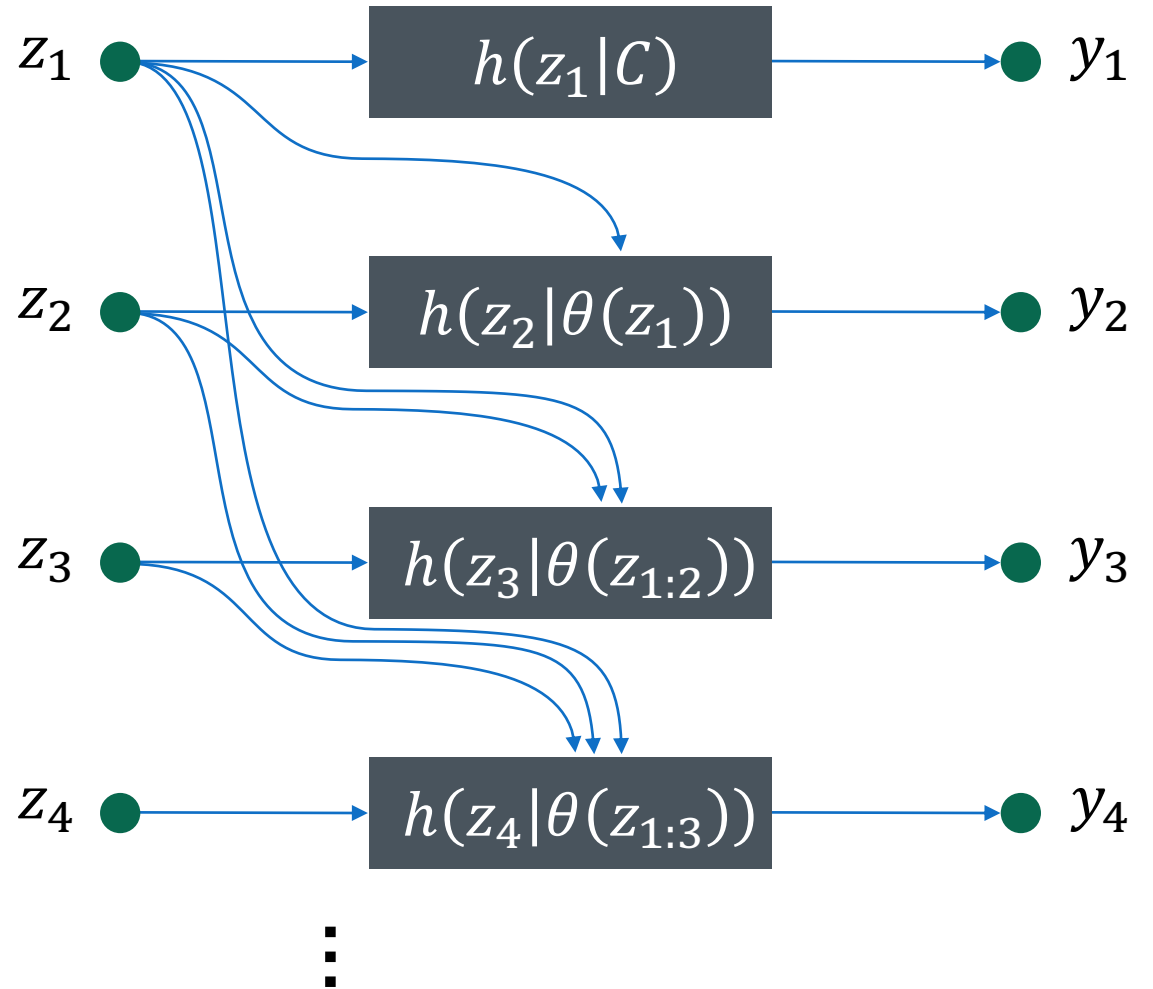
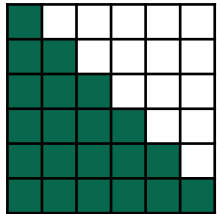
- $z = (z_1, \dots, z_D)^T \rightarrow y = (y_1, \dots, y_D)^T$
- $h(\cdot | \theta): \mathbb{R} \rightarrow \mathbb{R}$  – invertible function
- $y_i = h(z_i | \theta(z_{1:i-1}))$
- Inverse involves recursion



# Autoregressive flows

- $z = (z_1, \dots, z_D)^T \rightarrow y = (y_1, \dots, y_D)^T$
- $h(\cdot | \theta): \mathbb{R} \rightarrow \mathbb{R}$  – invertible function
- $y_i = h(z_i | \theta(z_{1:i-1}))$
- Inverse involves recursion

Jacobian structure:



# Evaluating Generative Models

# Evaluating generative models

- No single guide to follow
- Approaches are very problem-specific
  - E.g. perceived visual quality of generated images vs. quality of a generated dental crown model (<https://arxiv.org/abs/1804.00064>)
  - Most solutions are adapted to or invented for a given particular task
- We'll mention some approaches



# Evaluating generative models

(most obvious thing to do)

- By-eye comparison (if your data allows that)
  - Compare individual objects or whole distributions (e.g. in projections)
  - There might be no need to do any complicated evaluation if the model results simply look bad

# Evaluating generative models

(simple things to do)

- Compare meaningful physical characteristics (if applicable)
  - Means, medians, standard deviations, etc.
  - Correlations
- Statistical tests ( $\chi^2$ , Kolmogorov-Smirnov, etc.)
  - between individual dimensions or projections

# Additional classifier

- Train an independent model (e.g. xgboost) to distinguish real and fake samples
- Evaluate your GAN by checking the classifier's score (e.g. ROC AUC)
- Pros:
  - An objective quality measure
- Cons:
  - Resource consuming
  - Requires hyper-parameter tuning
  - May get picky to things that are not important

# Inception score

- Introduced in <https://arxiv.org/abs/1606.03498>
- Apply the Inception model (pre-trained image classifier) to obtain the conditional label distribution  $p(y|x)$  for each image  $x$ 
  - this should be low-entropy (the classifier should be certain)
- Calculate marginal  $p(y) = \int p(y|x = G(z))p(z)dz$ 
  - this should be high-entropy (diversity of samples)
- Combining these two requirements:

$$\text{IS} = \exp \left[ \mathbb{E}_x \left[ \text{KL}(p(y|x) || p(y)) \right] \right]$$

# Fréchet inception distance (FID score)

- Introduced in <https://arxiv.org/abs/1706.08500>
- One of the drawbacks of IS is that it doesn't care about the true distribution
- Instead one can compare distributions of activations at some Inception layer (originally – last pooling layer)
- The authors proposed calculating the Fréchet (aka Wasserstein-2) distance
- Distance between multivariate Gaussian approximations:

$$\text{FID} = \|\mu_r - \mu_g\|^2 + \text{Tr} \left[ \Sigma_r + \Sigma_g - 2(\Sigma_r \Sigma_g)^{1/2} \right]$$

# Precision and Recall distance (PRD)

**Definition 1.** For  $\alpha, \beta \in (0, 1]$ , the probability distribution  $Q$  has precision  $\alpha$  at recall  $\beta$  w.r.t.  $P$  if there exist distributions  $\mu, \nu_P$  and  $\nu_Q$  such that

Decomposition with  
a common part

$$P = \beta\mu + (1 - \beta)\nu_P \quad \text{and} \quad Q = \alpha\mu + (1 - \alpha)\nu_Q. \quad (3)$$

**Definition 2.** The set of attainable pairs of precision and recall of a distribution  $Q$  w.r.t. a distribution  $P$  is denoted by  $\text{PRD}(Q, P)$  and it consists of all  $(\alpha, \beta)$  satisfying Definition 1 and the pair  $(0, 0)$ .

- The authors provide an algorithm to calculate it for discrete distributions
- They convert Inception activations to discrete distribution using  $k$ -means clustering

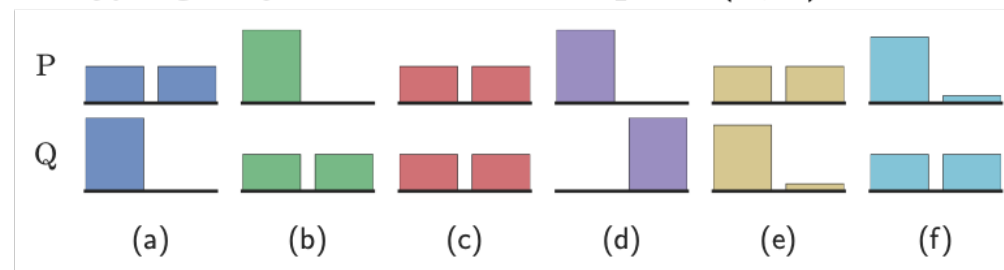


Figure 2: Intuitive examples of  $P$  and  $Q$ .

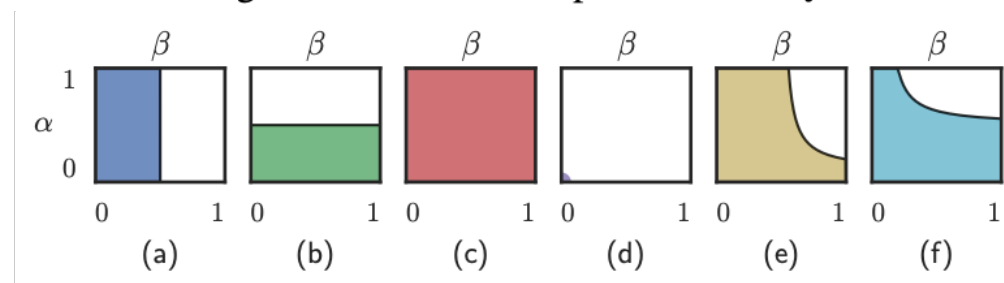


Figure 3:  $\text{PRD}(Q, P)$  for the examples above.

# More metrics...

- An extensive comparison of a large variety of measures:

- <https://arxiv.org/abs/1802.03446>

	Measure	Description
Quantitative	1. Average Log-likelihood [18, 22]	• Log likelihood of explaining realworld held out/test data using a density estimated from the generated data ( <i>e.g.</i> using KDE or Parzen window estimation). $L = \frac{1}{N} \sum_i \log P_{model}(\mathbf{x}_i)$
	2. Coverage Metric [33]	• The probability mass of the true data “covered” by the model distribution $C := P_{data}(dP_{model} > t)$ with $t$ such that $P_{model}(dP_{model} > t) = 0.95$
	3. Inception Score (IS) [3]	• KLD between conditional and marginal label distributions over generated data. $\exp(\mathbb{E}_{\mathbf{x}}[\mathbb{KL}(p(y \mathbf{x}) \  p(y))])$
	4. Modified Inception Score (m-IS) [34]	• Encourages diversity within images sampled from a particular category. $\exp(\mathbb{E}_{\mathbf{x}_i}[\mathbb{E}_{\mathbf{x}_j}[(\mathbb{KL}(P(y \mathbf{x}_i) \  P(y \mathbf{x}_j)))]])$
	5. Mode Score (MS) [35]	• Similar to IS but also takes into account the prior distribution of the labels over real data. $\exp(\mathbb{E}_{\mathbf{x}}[\mathbb{KL}(p(y \mathbf{x}) \  p(y^{train}))]) - \mathbb{KL}(p(y) \  p(y^{train}))$
	6. AM Score [36]	• Takes into account the KLD between distributions of training labels vs. predicted labels, as well as the entropy of predictions. $\mathbb{KL}(p(y^{train}) \  p(y)) + \mathbb{E}_{\mathbf{x}}[H(y \mathbf{x})]$
	7. Fréchet Inception Distance (FID) [37]	• Wasserstein-2 distance between multi-variate Gaussians fitted to data embedded into a feature space $FID(r, g) = \ \mu_r - \mu_g\ _2^2 + Tr(\Sigma_r + \Sigma_g - 2(\Sigma_r \Sigma_g)^{\frac{1}{2}})$
	8. Maximum Mean Discrepancy (MMD) [38]	• Measures the dissimilarity between two probability distributions $P_r$ and $P_g$ using samples drawn independently from each distribution. $M_k(P_r, P_g) = \mathbb{E}_{\mathbf{x}, \mathbf{x}' \sim P_r}[k(\mathbf{x}, \mathbf{x}')] - 2\mathbb{E}_{\mathbf{x} \sim P_r, \mathbf{y} \sim P_g}[k(\mathbf{x}, \mathbf{y})] + \mathbb{E}_{\mathbf{y}, \mathbf{y}' \sim P_g}[k(\mathbf{y}, \mathbf{y}')]$
	9. The Wasserstein Critic [39]	• The critic ( <i>e.g.</i> an NN) is trained to produce high values at real samples and low values at generated samples $\hat{W}(\mathbf{x}_{test}, \mathbf{x}_g) = \frac{1}{N} \sum_{i=1}^N \hat{f}(\mathbf{x}_{test}[i]) - \frac{1}{N} \sum_{i=1}^N \hat{f}(\mathbf{x}_g[i])$
	10. Birthday Paradox Test [27]	• Measures the support size of a discrete (continuous) distribution by counting the duplicates (near duplicates)
	11. Classifier Two Sample Test (C2ST) [40]	• Answers whether two samples are drawn from the same distribution ( <i>e.g.</i> by training a binary classifier)
	12. Classification Performance [1, 15]	• An indirect technique for evaluating the quality of unsupervised representations ( <i>e.g.</i> feature extraction; FCN score). See also the GAN Quality Index (GQI) [41].
	13. Boundary Distortion [42]	• Measures diversity of generated samples and covariate shift using classification methods.
	14. Number of Statistically-Different Bins (NDB) [43]	• Given two sets of samples from the same distribution, the number of samples that fall into a given bin should be the same up to sampling noise
	15. Image Retrieval Performance [44]	• Measures the distributions of distances to the nearest neighbors of some query images ( <i>i.e.</i> diversity)
	16. Generative Adversarial Metric (GAM) [31]	• Compares two GANs by having them engaged in a battle against each other by swapping discriminators or generators. $p(\mathbf{x} y=1; M_1)/p(\mathbf{x} y=1; M_2) = (p(y=1 \mathbf{x}; D_1)p(\mathbf{x}; G_2))/(p(y=1 \mathbf{x}; D_2)p(\mathbf{x}; G_1))$
	17. Tournament Win Rate and Skill Rating [45]	• Implements a tournament in which a player is either a discriminator that attempts to distinguish between real and fake data or a generator that attempts to fool the discriminators into accepting fake data as real.
	18. Normalized Relative Discriminative Score (NRDS) [32]	• Compares $n$ GANs based on the idea that if the generated samples are closer to real ones, more epochs would be needed to distinguish them from real samples.
	19. Adversarial Accuracy and Divergence [46]	• Adversarial Accuracy. Computes the classification accuracies achieved by the two classifiers, one trained on real data and another on generated data, on a labeled validation set to approximate $P_g(y \mathbf{x})$ and $P_r(y \mathbf{x})$ . Adversarial Divergence: Computes $\mathbb{KL}(P_g(y \mathbf{x}), P_r(y \mathbf{x}))$
	20. Geometry Score [47]	• Compares geometrical properties of the underlying data manifold between real and generated data.
	21. Reconstruction Error [48]	• Measures the reconstruction error ( <i>e.g.</i> $L_2$ norm) between a test image and its closest generated image by optimizing for $z$ ( <i>i.e.</i> $\min_z \ G(\mathbf{z}) - \mathbf{x}^{(test)}\ ^2$ )
	22. Image Quality Measures [49, 50, 51]	• Evaluates the quality of generated images using measures such as SSIM, PSNR, and sharpness difference
	23. Low-level Image Statistics [52, 53]	• Evaluates how similar low-level statistics of generated images are to those of natural scenes in terms of mean power spectrum, distribution of random filter responses, contrast distribution, etc.
	24. Precision, Recall and $F_1$ score [23]	• These measures are used to quantify the degree of overfitting in GANs, often over toy datasets.
Qualitative	1. Nearest Neighbors	• To detect overfitting, generated samples are shown next to their nearest neighbors in the training set
	2. Rapid Scene Categorization [18]	• In these experiments, participants are asked to distinguish generated samples from real images in a short presentation time ( <i>e.g.</i> 100 ms); <i>i.e.</i> real v.s fake
	3. Preference Judgment [54, 55, 56, 57]	• Participants are asked to rank models in terms of the fidelity of their generated images ( <i>e.g.</i> pairs, triples)
	4. Mode Drop and Collapse [58, 59]	• Over datasets with known modes ( <i>e.g.</i> a GMM or a labeled dataset), modes are computed as by measuring the distances of generated data to mode centers
	5. Network Internals [1, 60, 61, 62, 63, 64]	• Regards exploring and illustrating the internal representation and dynamics of models ( <i>e.g.</i> space continuity) as well as visualizing learned features