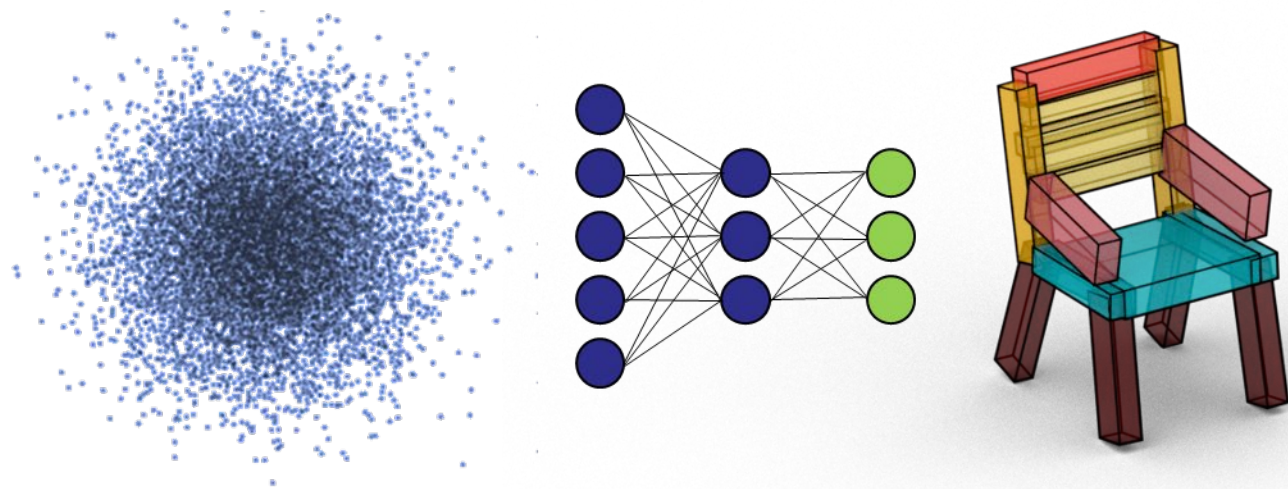# CS348n: Neural Representations and Generative Models for 3D Geometry
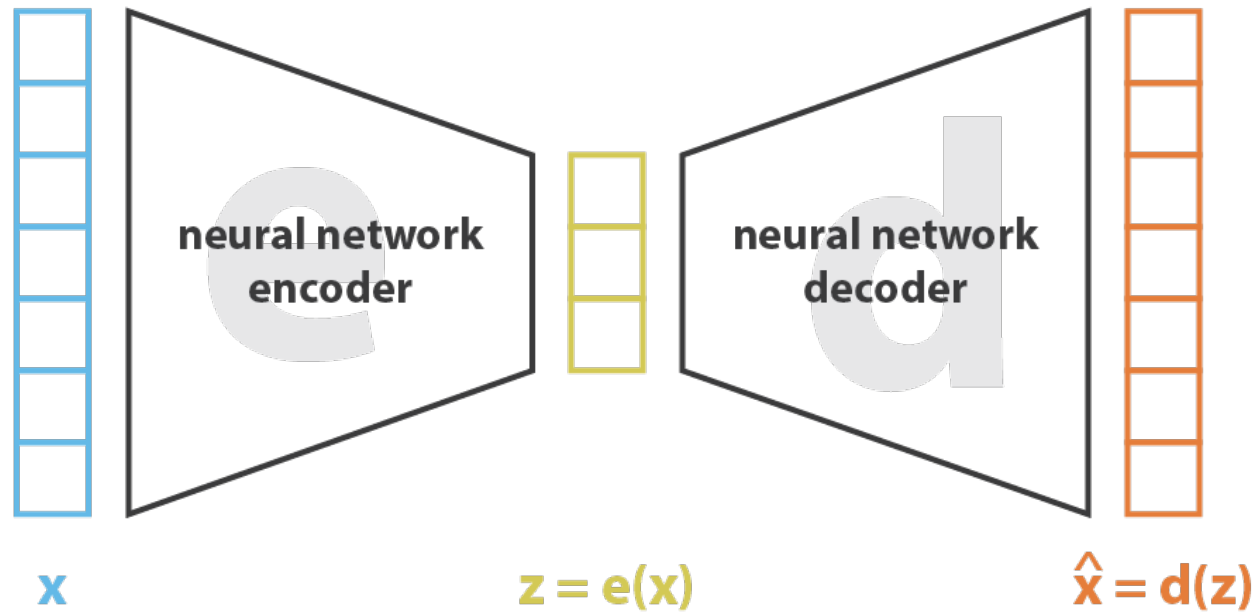
Leonidas Guibas
Computer Science Department
Stanford University

Leonidas Guibas Laboratory

Geometric Computing

# Last Time: Variational AutoEncoders, Neural Implicits

# Autoencoders



$$\text{loss} \; = \; || \, \mathbf{x} - \hat{\mathbf{x}} \, ||^2 \; = \; || \, \mathbf{x} - \mathbf{d}(\mathbf{z}) \, ||^2 \; = \; || \, \mathbf{x} - \mathbf{d}(e(\mathbf{x})) \, ||^2$$

# Autoencoders for Content Generation



**training process**

encoder

**e**

encoded vector
(in latent space)

input

**generation process**

sampler

sampled vector
(from latent space)

decoder

**d**

decoded content

(reconstructed input /
generated content)

Severe overfitting

# Autoencoders for Content Generation



encoder

decoder

point sampled from the one dimensional latent space for new content generation

"training" data for the autoencoder

encoded data can be decoded without loss if the autoencoder has enough degrees of freedom

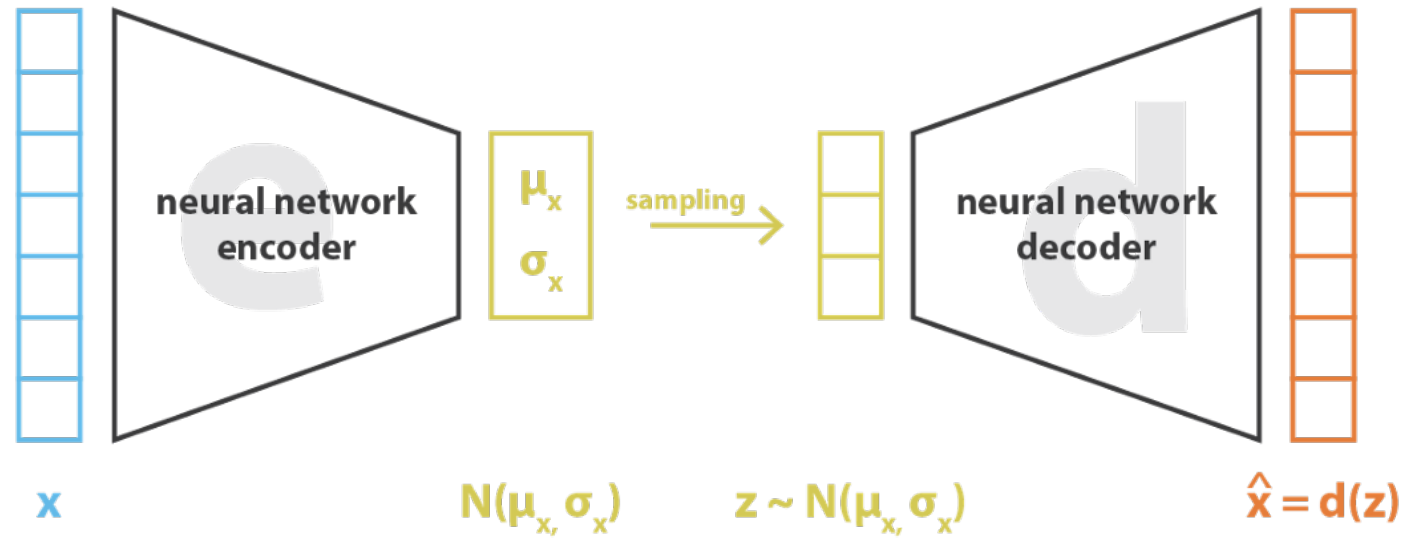without explicit regularisation, some points of the latent space are "meaningless" once decoded

An autoencoder is solely trained to encode and decode with as small loss as possible, no matter how the latent space is organized

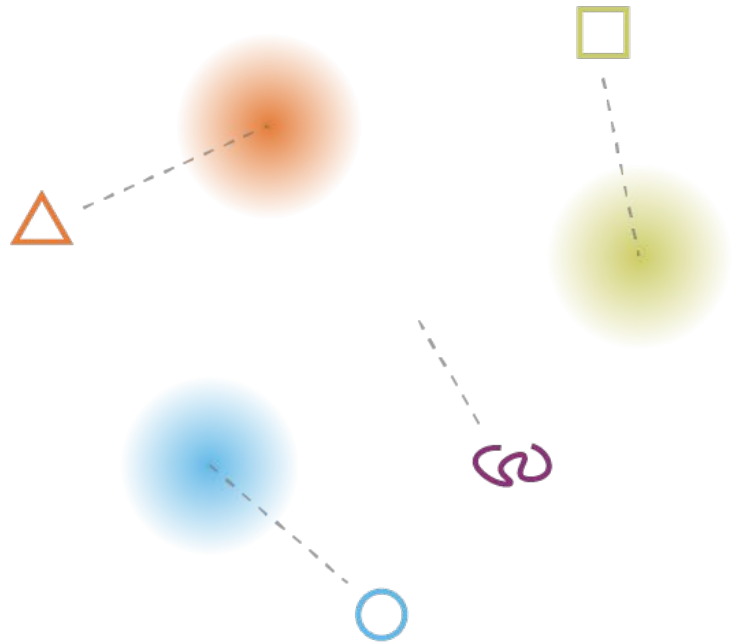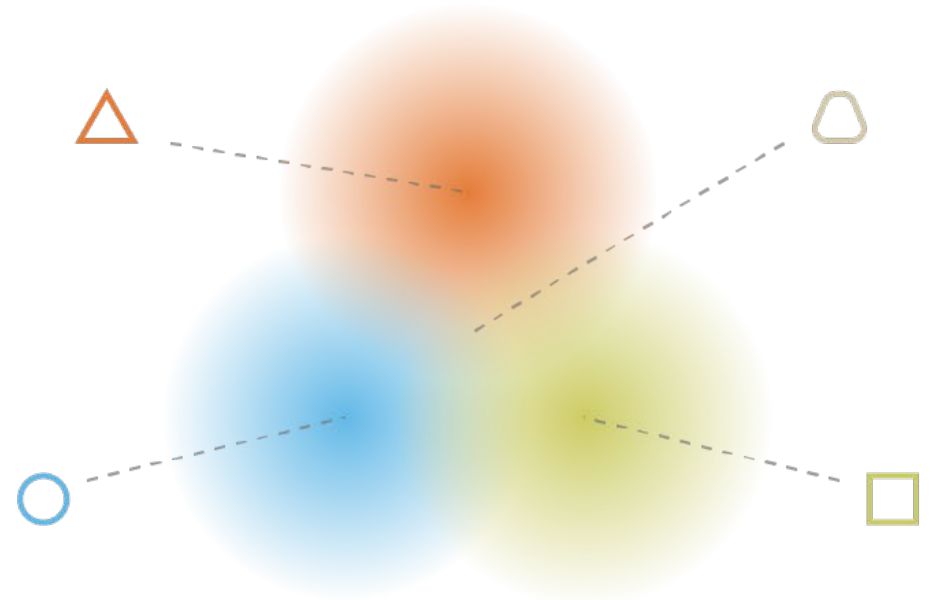# Regularize the Distribution in Latent Space



$$\text{loss} = || x - \hat{x} ||^2 + KL[\, N(\mu_x, \sigma_x), N(0, I)\,] = || x - d(z) ||^2 + KL[\, N(\mu_x, \sigma_x), N(0, I)\,]$$

Make the encoder probabilistic  with a latent space distribution like a simple Gaussian
Add a second loss measuring distribution distance (via the Kulback-Leibler divergence)
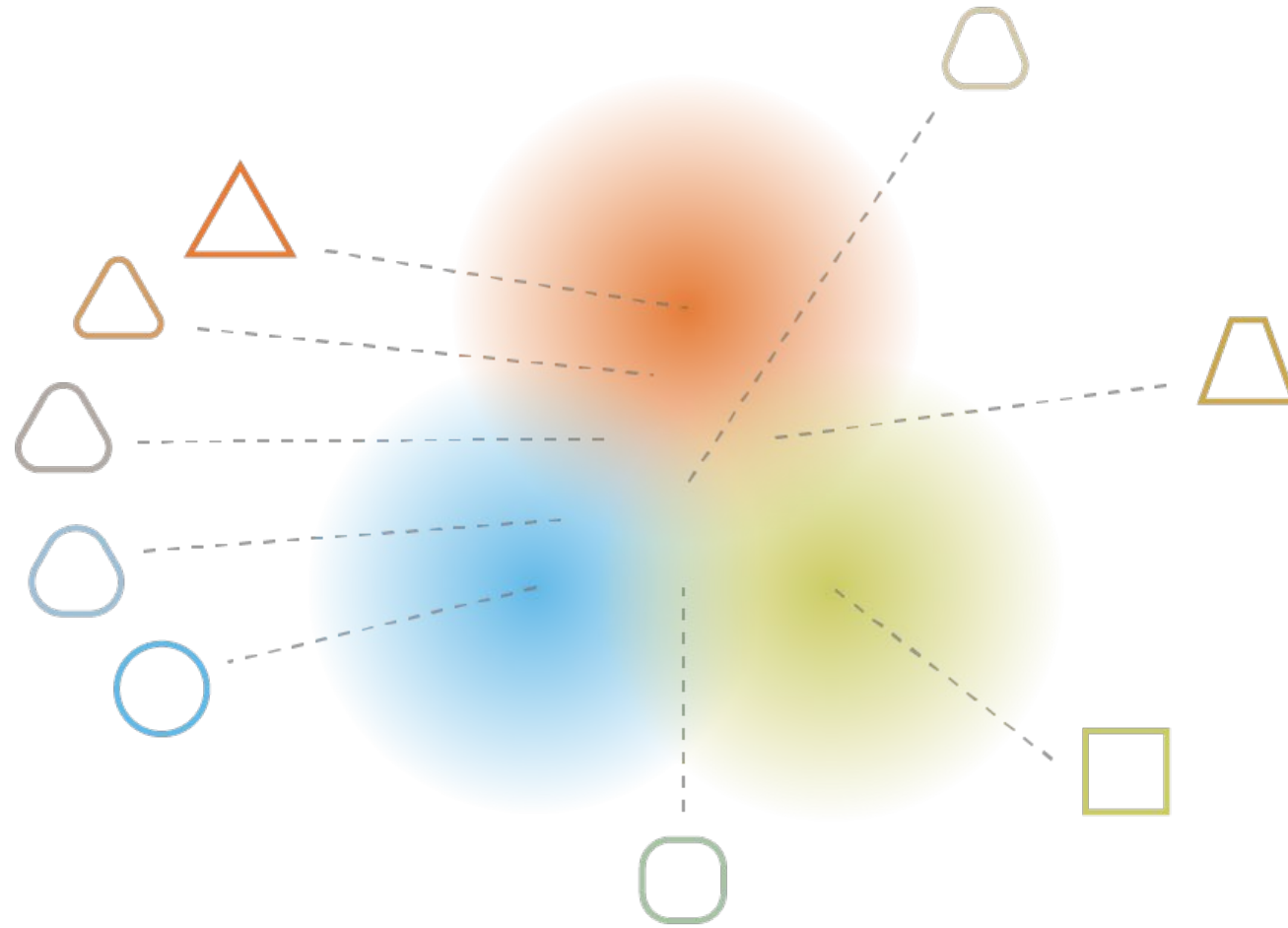
# The Effect of Regularization



what can happen without regularisation ❌

✅ what we want to obtain with regularisation
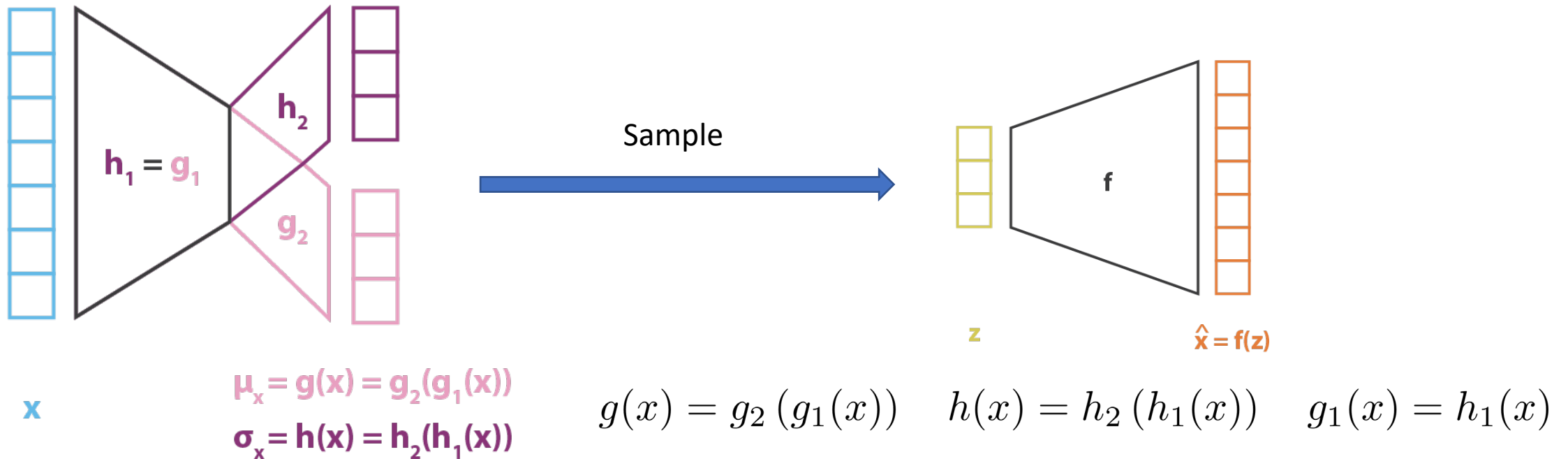
Overfitting with "punctual" distributions

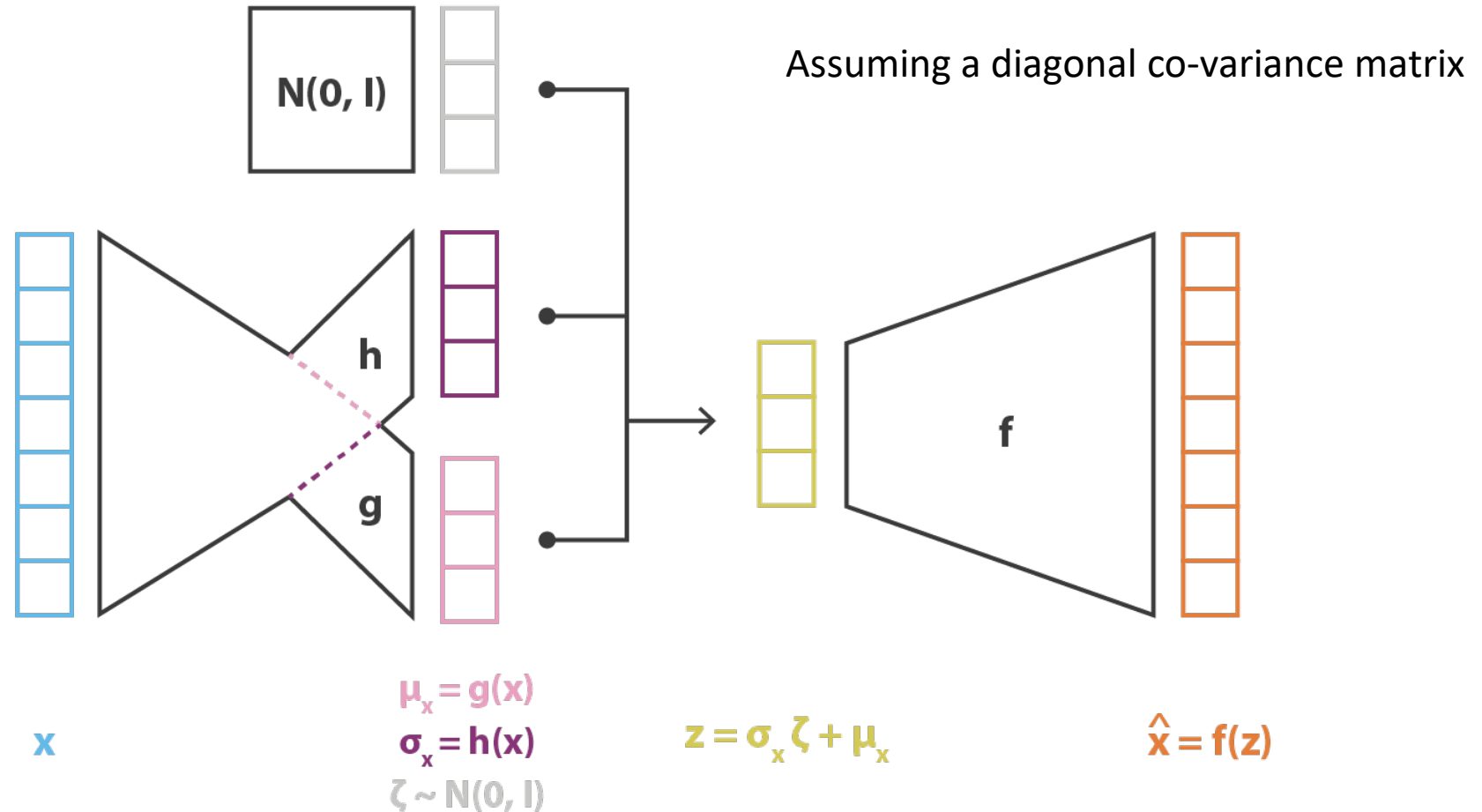Create smooth gradients over the information encoded in the latent space

# Variational AutoEncoder (VAE)



Sample

$\mu_x = g(x) = g_2(g_1(x))$

$\sigma_x = h(x) = h_2(h_1(x))$

$$g(x) = g_2\left(g_1(x)\right) \qquad h(x) = h_2\left(h_1(x)\right) \qquad g_1(x) = h_1(x)$$

How to train:
Sampling is a problem w. back propagation

Assuming a diagonal co-variance matrix

$N(0, I)$

$h$

$g$

$f$

$x$

$\mu_x = g(x)$
$\sigma_x = h(x)$
$\zeta \sim N(0, I)$

$z = \sigma_x \zeta + \mu_x$

$\hat{x} = f(z)$

$$loss = C \, || \, x - \hat{x} \, ||^2 + KL[\, N(\mu_x, \sigma_x), N(0, I)\,] = C \, || \, x - f(z) \, ||^2 + KL[\, N(g(x), h(x)), N(0, I)\,]$$

# Signed Distance Function

NN

$(x,y,z)$

SDF

(b)

# Can then Reconstruct via Marching Cubes



Lorensen et al., 1987

Code

(x,y,z)

NN

SDF

# Auto-Decoder



Auto-Encoder

Auto-Decoder

**Backpropagate**

$z_i$

NN

Code

$(x,y,z)$

$f_\theta$

$\|$ SDF ▬ GT $\|$

$$\arg\min_{\theta,\{\boldsymbol{z}_i\}_{i=1}^N} \sum_{i=1}^N \left( \sum_{j=1}^K \mathcal{L}(f_\theta(\boldsymbol{z}_i, \boldsymbol{x}_j), s_j) \right)$$

20

Our
Reconstruction

Octree Based

Ground Truth          Our Reconstruction          Atlasnet (25 Patches)          Atlasnet (1 Patch)

5D Input
Position + Direction

Output
Color + Density

$(x, y, z, \theta, \phi) \rightarrow F_\Theta \rightarrow (RGB\sigma)$

Ray 1

Ray 2

(a)

(b)

Mildenhall et al. 2020

# Neural Parametrics: AtlasNet

# Implicit Curves and Surfaces via Functions

- Kernel of a scalar function $f : \mathbb{R}^m \to \mathbb{R}$

  - Curve in 2D: $S = \{x \in \mathbb{R}^2 | f(x) = 0\}$
  - Surface in 3D: $S = \{x \in \mathbb{R}^3 | f(x) = 0\}$

$$x := (a, b)$$
$$a^2 + b^2 - 1$$

- Space partitioning

  $\{x \in \mathbb{R}^m | f(x) > 0\}$ Outside
  $\{x \in \mathbb{R}^m | f(x) = 0\}$ Curve/Surface
  $\{x \in \mathbb{R}^m | f(x) < 0\}$ Inside

$f(x) > 0$

$f(x) < 0$

$v$

$u$

$F$

$R^3$

Embedding

$F$

$z$

$x$

$Y$

$v$

$u$

$F(u,v) \equiv (X(u,v), Y(u,v), Z(u,v))$

$F : R^2 \longrightarrow R^3 \longleftarrow$ ambient space

$\uparrow$ parameter space

These parametric mappings can be explicit functions, but can also be neural networks

$$Y = x^2$$

$$Y - x^2 = 0$$

implicit

$$F(t) = (t, t^2)$$

Parametric

$$X(t) \quad Y(t)$$

$t$

$\mathbb{R}$

# Training setup for 3D reconstruction

# Generating points

Test Shape

AtlasNet: A Papier-Mâché Approach to Learning 3D Surface Generation
Thibault Groueix, Matthew Fisher, Vladimir G. Kim, Bryan C. Russell, Mathieu Aubry
https://arxiv.org/abs/1802.05384

Thibault Groueix, Pierre-Alain Langlois, 2019

# Generating points

Latent shape representation

Generated 3D points

MLP

Test Shape

# Generating points

Latent shape representation

Generated 3D points

MLP

Issue: no idea of surfaces

Test Shape

32

# Limitation of PointSetGen [Fan2017]

➔ **Generate a fixed number of points**
➔ Points connectivity is missing
➔ Generated points are not correlated enough to belong to an implicit surface



Latent shape representation

MLP

Generated 3D points

3

n

# Limitation of PointSetGen [Fan2017]

➔ Generate a fixed number of points
➔ **Points connectivity is missing**
➔ Generated points are not correlated enough to belong to an implicit surface

# Limitation of PointSetGen [Fan2017]

➔ Generate a fixed number of points
➔ **Points connectivity is missing**
➔ Generated points are not correlated enough to belong to an implicit surface

# Limitation of PointSetGen [Fan2017]

➔ Generate a fixed number of points
➔ **Points connectivity is missing**
➔ Generated points are not correlated enough to belong to an implicit surface

# Limitation of PointSetGen [Fan2017]

➜ Generate a fixed number of points

➜ **Points connectivity is missing**

➜ Generated points are not correlated enough to belong to an implicit surface

# Limitation of PointSetGen [Fan2017]

➔ Generate a fixed number of points
➔ Points connectivity is missing
➔ **Generated points are not correlated enough to belong to an implicit surface**



Reconstructing the mesh from a pointcloud :
Poisson Surface Reconstruction [**Kazhdan2013**]

# Training setup for 3D reconstruction



Partial Data **X**

**RGB Image(s)**
**RGBD Image(s)**
**PointCloud**
**PointCloud (Voxelized)**

**...**

**Encoder E**

Feature vector

**Decoder D**

Choice of representation

Volumetric (OctNet)
PointClouds (PointSetGen)
**Surfaces (AtlasNet)**
Signed Distance Function
Geometric Primitives

**Loss**

$$L(D(E(X)),Y)$$

3D Object **Y**

# Deform a surface [Groueix2018]

Latent shape
representation

**x**
**y**

MLP

Sampled
2D point

Generated
3D point

# Deform a surface : space mapping trick [Groueix2018]

Latent shape
representation

x
y

Sampled
2D point

MLP

Generated
3D point

$v$

$F$

$u$

# Deform a surface [Groueix2018]

Latent shape
representation

x
y

Sampled
2D point

MLP

Generated
3D point

$v$

$F$

$u$

# Deform a surface [Groueix2018]

Latent shape representation

Sampled 2D point

**x**
**y**

MLP

Generated 3D point

$F$

$v$

$u$

43

# Deform a surface [Groueix2018]

Latent shape
representation

x
y

Sampled
2D point

MLP

Generated
3D point

# Deform a surface [Groueix2018]

Latent shape
representation

x
y

Sampled
2D point

MLP

Generated
3D point

# Deform a surface [Groueix2018]

Latent shape
representation

x
y

Sampled
2D point

MLP

Generated
3D point

# Deform a surface [Groueix2018]

Test Shape

Latent shape representation

x
y

Sampled 2D point

MLP 2

Generated 3D point

# Deform a surface [Groueix2018]

Latent shape representation

x
y

MLP 2

Sampled 2D point

Generated 3D point

Test Shape

48

# Deform a surface [Groueix2018]

Generated 3D point

MLP 1

Latent shape representation

**x**
**y**

Sampled 2D point

Test Shape

# Deform a surface [Groueix2018]



Encoder E

Decoder D

Latent shape representation

Generated 3D point

MLP 1

MLP 2

MLP 3

Sampled 2D point

x
y

Test Shape

Piecewise Parametric Surfaces

50

Thibault Groueix, Pierre-Alain Langlois, 2019

# Deform a surface [Groueix2018]

Latent shape representation

Sampled 2D point

x
y

MLP 1

MLP 2

MLP 3

Generated 3D point

Test Shape

51

# Results : Single View Reconstruction



(a) Input    (b) 3D-R2N2    (c) HSP    (d) PSG    (e) Ours

52

# Direct application : mesh parametrization

(a) Input Shape    (b) Template    (c) Deformed template

3D-CODED : 3D Correspondences by Deep Deformation

Thibault Groueix, Matthew Fisher, Vladimir G. Kim, Bryan C. Russell, Mathieu Aubry https://arxiv.org/abs/1806.05228

(a) **Network training**

(b) **Local optimization of feature x**

(c) **Correspondences**

# State-of-the-art correspondences of FAUST [Groueix2018b]

# Generative Adversarial Networks

$$CDF_U(u) = \mathbb{P}(U \leq u) = u \quad \forall u \in [0, 1]$$

**Cumulative Distribution Function (CDF)**

$$CDF_X(x) = \mathbb{P}(X \leq x) \quad \in [0, 1]$$



$U$

$CDF_X^{-1}$

$Y = CDF_X^{-1}(U)$

$X$

$$CDF_Y(y) = \mathbb{P}(Y \leq y) = \mathbb{P}\left(CDF_X^{-1}(U) \leq y\right) = \mathbb{P}\left(U \leq CDF_X(y)\right) = CDF_X(y)$$

Conceptually, the purpose of the "transform function" is to deform/reshape the initial probability distribution: the transform function takes from where the initial distribution is too high compared to the targeted distribution and puts it where it is too low.

# Complex Output Distributions



Use a neural network as the transform function

The problem of generating a new image of dog is equivalent to the problem of generating
a new vector following the "dog probability distribution" over the N dimensional vector space.
So we are, in fact, facing a problem of generating a random variable with respect to a specific probability distribution.

**GENERATIVE NETWORK**

Input random variable (drawn from a simple distribution, for example uniform).

The generative network transforms the simple random variable into a more complex one.

Output random variable (should follow the targeted distribution, after training the generative network).

The output of the generative network once reshaped.

# Loss: Comparing Distributions Based on Samples

- generate some uniform inputs
- make these inputs go through the network and collect the generated outputs
- compare the true "dog probability distribution" and the generated one based on the available samples
- use backpropagation to make one step of gradient descent to lower the distance between true and generated distributions

- This is a very hard problem
  - Maximum Mean Discrepancy (MMD)
  - compute the MMD distance between the sample of true dog images and the sample of generated ones

Forward transform of the initial random variables to generate data

**GENERATIVE NETWORK**

Backpropagation of the matching error to train the network

Input random variables (drawn from a uniform).

Generative network to be trained.

The generated distribution is compared to the true distribution and the "matching error" is backpropagated to train the network.

Extremely expensive!

# An Alternative: Compare on a Downstream Task

- An indirect loss

- **Generative Adversarial Networks (GANs)**: compare distributions through a downstream task

- Use the loss of that task to improve the generator

- Make that task itself be a trainable neural network

## Baseline: the direct approach

The distribution in blue is the true one while the generated distribution is depicted in orange. Iteration by iteration, we compare the two distributions and adjust the networks weights through gradient descent steps. Here the comparison is done over the mean and the variance (similar to a truncated moments matching method). Notice that (obviously) this example is so simple that it doesn't require an iterative approach: the purpose is only to illustrate the intuition given above.

Indirect: use a discriminator

differentiate samples between
the two distributions

The blue distribution is the true one, the orange is the generated one. In grey, with corresponding y-axis on the right, we displayed the probability to be true for the discriminator if it chooses the class with the higher density in each point (assuming "true" and "generated" data are in equal proportions). The closer the two distributions are, the more often the discriminator is wrong. When training, the goal is to "move the green area" (generated distribution is too high) towards the red area (generated distribution is too low).

# How to Get a Discriminator?

- Learn a discriminator through another neural network

- the goal of the <span style="color:red">generator</span> is to fool the discriminator, so the generative neural network is trained to maximize the final classification error (between true and generated data)

- the goal of the <span style="color:red">discriminator</span> is to detect fake generated data, so the discriminative neural network is trained to minimize the final classification error

At each iteration of the training process, the weights of the generative network are updated in order to increase the classification error (error gradient ascent over the generator's parameters) whereas the weights of the discriminative network are updated so that to decrease this error (error gradient descent over the discriminator's parameters).

# An Adversarial Discriminator



Forward propagation (generation and classification)

Backward propagation (adversarial training)

GENERATIVE NETWORK

DISCRIMINATIVE NETWORK

Input random variables.

The generative network is trained to **maximise** the final classification error.

The generated distribution and the true distribution are not compared directly.

The discriminative network is trained to **minimise** the final classification error.

The classification error is the basis metric for the training of both networks.

- a generative network $G(.)$ that takes a random input $z$ with density $p_z$ (the "noise vector") and returns an output $x_g = G(z)$ that should follow (after training) the targeted probability distribution

- a discriminative network $D(.)$ that takes an input $x$ that can be a "true" one ($x_t$, whose density is denoted $p_t$) or a "generated" one ($x_g$, whose density $p_g$ is the density induced by the density $p_z$ going through $G$) and that returns the probability $D(x)$ of to be a "true" data

Error
$$E(G, D) = \frac{1}{2}\mathbb{E}_{x \sim p_t}[1 - D(x)] + \frac{1}{2}\mathbb{E}_{z \sim p_z}[D(G(z))]$$
$$= \frac{1}{2}\left(\mathbb{E}_{x \sim p_t}[1 - D(x)] + \mathbb{E}_{x \sim p_g}[D(x)]\right)$$

$$\max_G \left(\min_D E(G, D)\right)$$

A minimax Nash equilibrium

# Direct vs. Indirect Losses

In an idealized setting of unlimited capacity generator and discriminator and smoothness of the underlying distributions:

it can be shown that the learned generator produces the same density as the true density and the learned discriminator can't do better than being true in one case out of two.

## Deep Convolutional GANs (DCGANs)



Generator Architecture

**Key ideas**:

- Replace FC hidden layers with Convolutions
  - **Generator:** Fractional-Strided convolutions

- Use Batch Normalization after each layer

- **Inside Generator**
  - Use ReLU for hidden layers
  - Use Tanh for the output layer

Radford, Alec, Luke Metz, and Soumith Chintala. "Unsupervised representation learning with deep convolutional generative adversarial networks." arXiv:1511.06434 (2015).

DCGAN: Bedroom images

Radford, Alec, Luke Metz, and Soumith Chintala. "Unsupervised representation learning with deep convolutional generative adversarial networks." arXiv:1511.06434 (2015).

# Also for 3D

## 3D Chairs



(a) Rotation       (b) Width

Chen, X., Duan, Y., Houthooft, R., Schulman, J., Sutskever, I., & Abbeel, P. InfoGAN: Interpretable Representation Learning by Information Maximization Generative Adversarial Nets, NIPS (2016).

man with glasses − man without glasses + woman without glasses = woman with glasses

Radford, Alec, Luke Metz, and Soumith Chintala. "Unsupervised representation learning with deep convolutional generative adversarial networks." arXiv:1511.06434 (2015).

# GAN Advantages

- Generation is straightforward

- Mode detail is captured

- Training does not require MLE estimation

- Robust to overfitting (generator never sees the training data)

- Impressive empirical results

# GAN Issues

- Learned probability distribution is implicit
  - Vanilla GANS only good for sampling/generation

- Training is difficult and often unstable
  - Non-convergence
  - Vanishing gradients
  - Mode collapse

- Non-convergence
  - Stochastic gradient descent was not designed to find the Nash equilibria in multi-player games
  - Competition between generator and discriminator can cause instabilities
  - A black art, addressed through adding noise to discriminator inputs, toying with learning rates, and various regularizations

- Vanishing gradients
  - Generator training can fail if the discriminator is too good -- an optimal discriminator doesn't provide enough information for the generator to make progress
  - As gradients flow backwards, they can become so small that the early generator layers stop changing
  - A key contribution to address this was the Wasserstein loss, using transportation metrics

$$\min_x \max_y \ xy$$

- $\frac{\partial}{\partial x} = -y \qquad ... \quad \frac{\partial}{\partial y} = x$

- $\frac{\partial^2}{\partial y^2} = \frac{\partial}{\partial x} = -y$

- diff eq has sinusoidal terms
- will not converge, even with small leaning rate

Wasserstein GAN, Martin Arjovsky, Soumith Chintala, Léon Bottou, 2017 https://arxiv.org/abs/1701.07875

- Mode collapse
  - Generator fails to produce diverse-enough samples



Metz, Luke, et al. **"Unrolled Generative Adversarial Networks."** arXiv preprint arXiv:1611.02163 (2016).

- Remedy: let the discriminator looks at the entire batch, not just a single sample – mark as "fake" if there is lack of diversity

# Disentanglement: InfoGAN

- Disentanglement means that individual latent dimensions capture independent key attribute of the output

- How to achieve disentanglement without explicit supervision?

- InfoGAN approach:
  - partition the noise vector into 2 parts
    - a $z$ vector that will capture slight/local variations in the output
    - a small $c$ vector will capture the main attributes of the output
  - maximize <span style="color:red">mutual information</span> between $c$ and the generated data

$$I(X;Y) = \sum_{x,y} p(x,y) \log \left( \frac{p(x,y)}{p(x)p(y)} \right)$$

$$I(X;Y) = H(X) - H(X \mid Y) = H(Y) - \mathrm{H}(Y \mid X)$$

$$\min_G \max_D E_I(D,G) = E(D,G) - \lambda I(c;G(z,c))$$

  - For MNIST, $c$ dimensions correlate with digit class, thickness, slope, etc.

Chen, X., Duan, Y., Houthooft, R., Schulman, J., Sutskever, I., & Abbeel, P. InfoGAN: Interpretable Representation Learning by Information Maximization Generative Adversarial Nets https://arxiv.org/abs/1606.03657

# Disentanglement: StyleGAN

- An alternative architecture based on the style transfer literature

- Allows unsupervised separation of high-level attributes (e.g., pose vs identity for human faces) as well as stochastic variation (freckles, facial hair)

- Allows control of the synthesis

(a) Traditional

(b) Style-based generator

*A* = learned affine transform

*B* = per-channel scaling factors

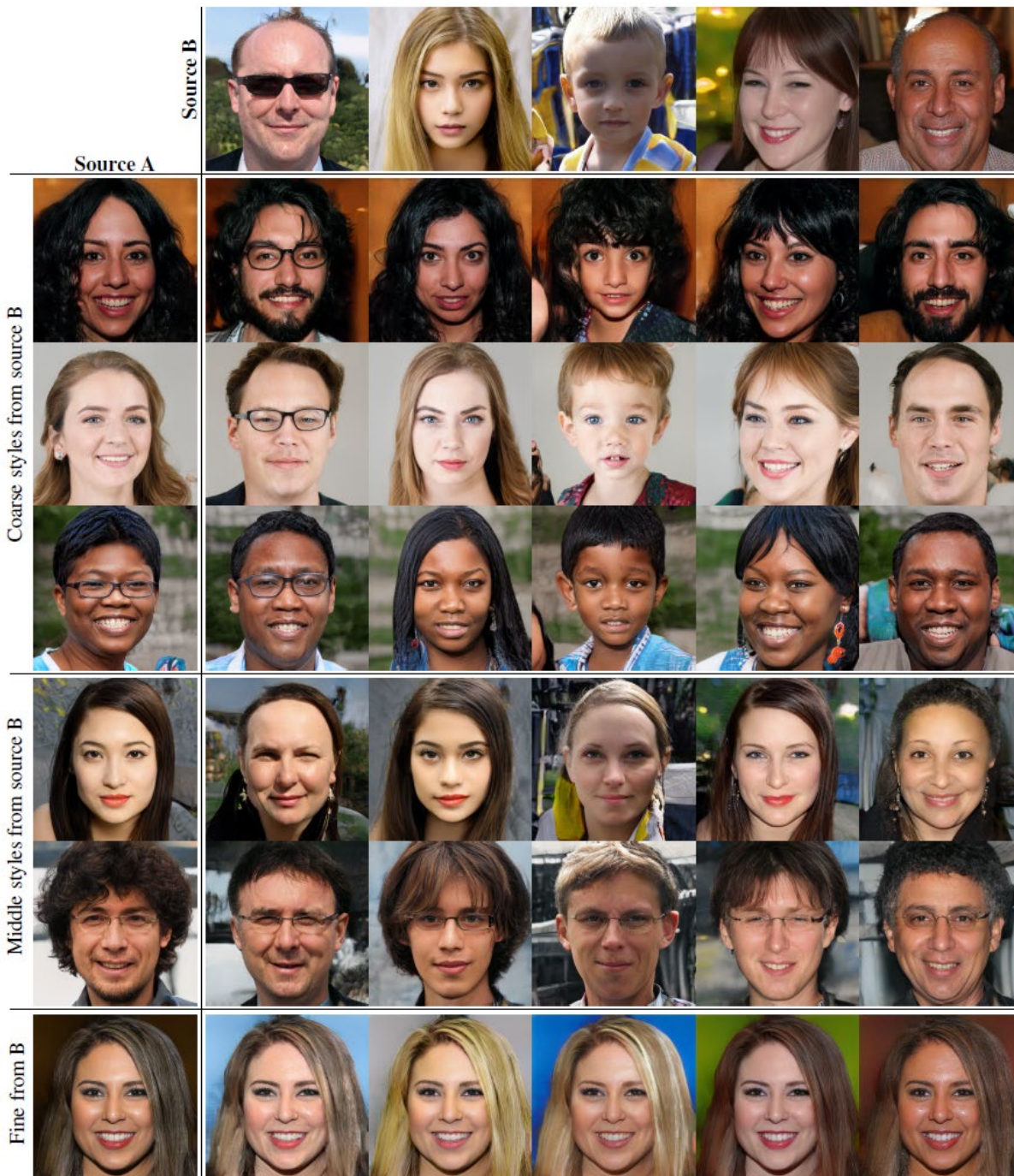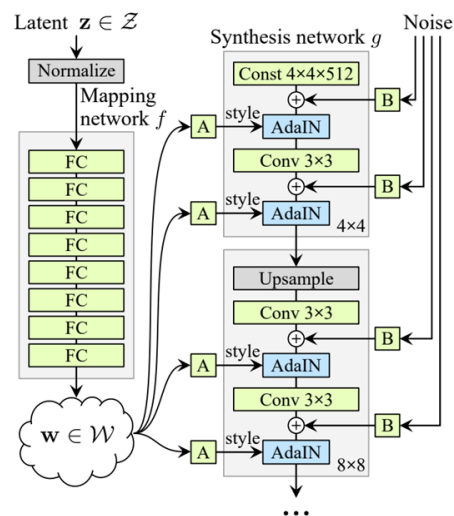AdaIN = adaptive instance normalization

In *Instance Normalization*, mean and variance are calculated for each individual channel for each individual sample across both spatial dimensions

$$\mathrm{AdaIN}\left(\mathbf{x}_i, \mathbf{y}\right) = \mathbf{y}_{s,i} \frac{\mathbf{x}_i - \mu\left(\mathbf{x}_i\right)}{\sigma\left(\mathbf{x}_i\right)} + \mathbf{y}_{b,i}$$
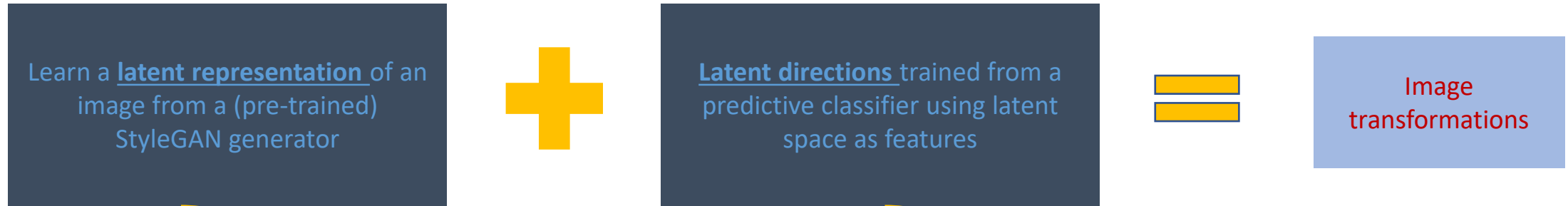
80

# Example Generations

Two sets of images were generated from their respective latent codes (sources A and B); the rest of the images were generated by copying a specified subset of styles from source B and taking the rest from source A. Copying the styles corresponding to coarse spatial resolutions ($4^2 - 8^2$) brings high-level aspects such as pose, general hair style, face shape, and eyeglasses from source B, while all colors (eyes, hair, lighting) and finer facial features resemble A. If we instead copy the styles of middle resolutions ($16^2 - 322$) from B, we inherit smaller scale facial features, hair style, eyes open/closed from B, while the pose, general face shape, and eyeglasses from A are preserved. Finally, copying the fine styles ($64^2 - 1024^2$) from B brings mainly the color scheme and microstructure.
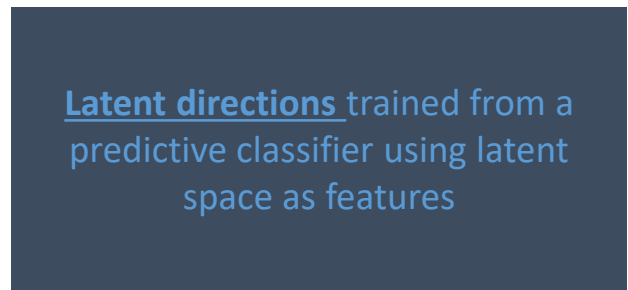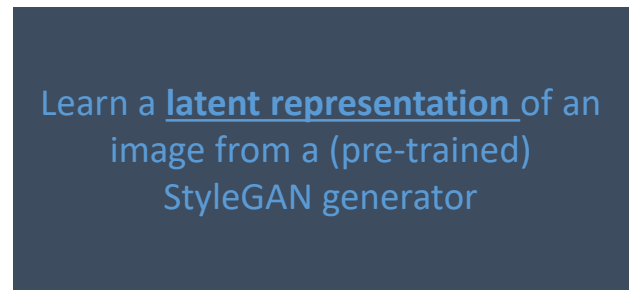
# StyleGAN Encoder

Learn a **latent representation** of an image from a (pre-trained) StyleGAN generator

**+**

**Latent directions** trained from a predictive classifier using latent space as features

**=**

Image transformations

How does the image look like in latent space?

To manipulate the images to e.g. smile

Learn a **latent representation** of an image from a (pre-trained) StyleGAN generator

**Latent directions** trained from a predictive classifier using latent space as features

Image

**+**

Pre-trained StyleGAN generator

Image

**+**

A handful of labelled attributes e.g. is the image smiling? Gender? Age?