

Lecture 17:

Generative AI Part 2

GANs & Flow Matching

Administrative

- A3 and midterm grades are out. Midterm solution is out on website. Regrade requests will be accepted between May 27 and June 1
- A5 will be out today/tomorrow and due June 5th
- Final report due June 8th

Taxonomy of Generative Models

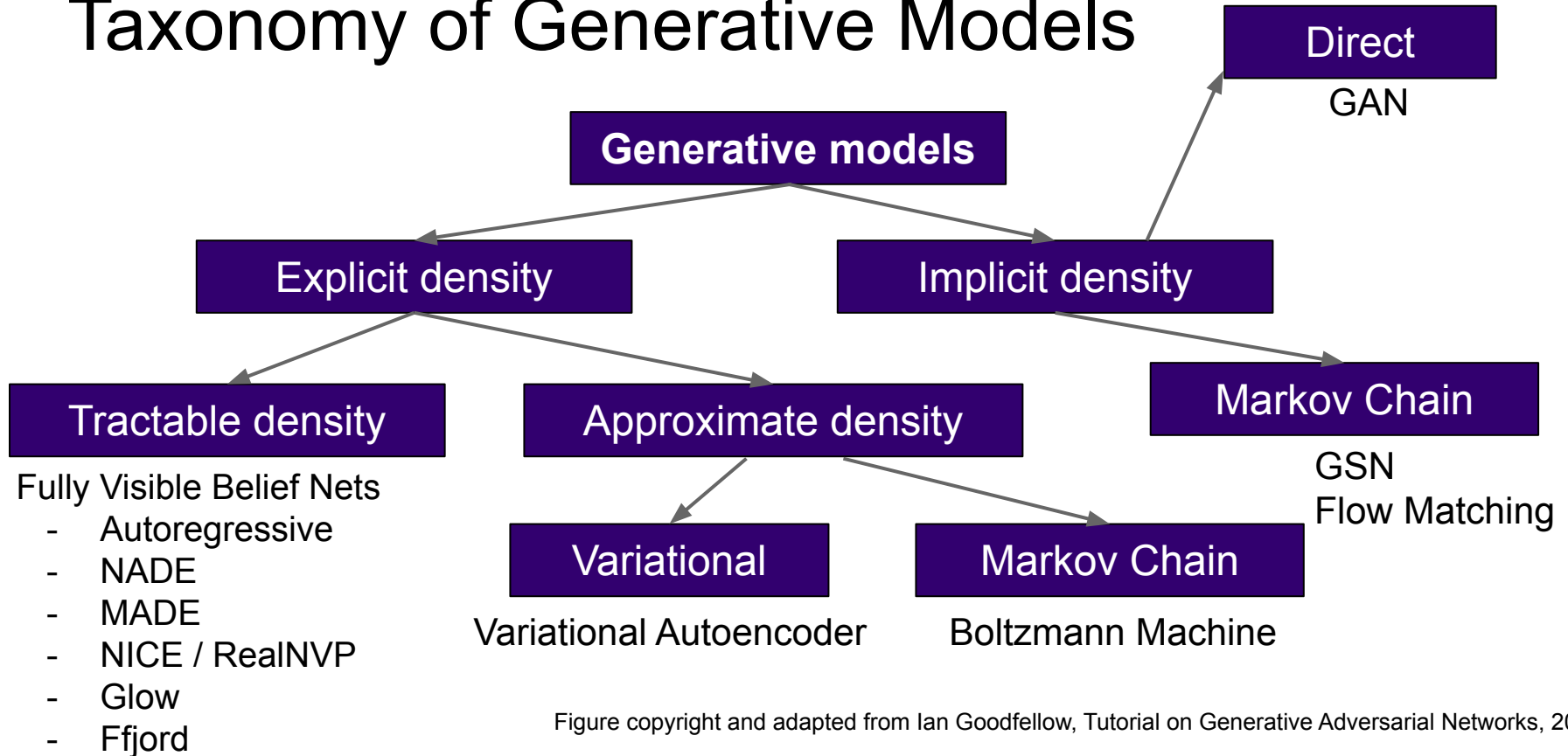


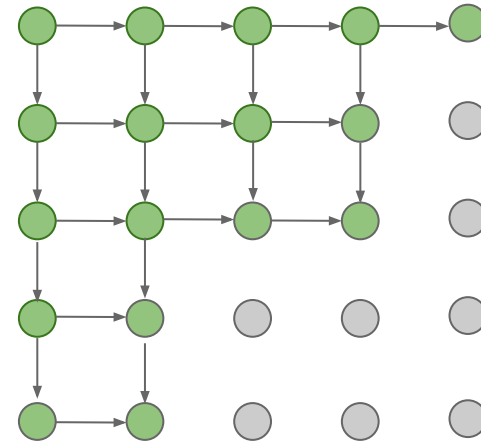
Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

Generative AI so far: Autoregressive models

Generate image pixels starting from corner

Dependency on previous pixels modeled using an RNN (LSTM)

Very slow during both training and testing; $N \times N$ image requires $2N-1$ sequential steps!



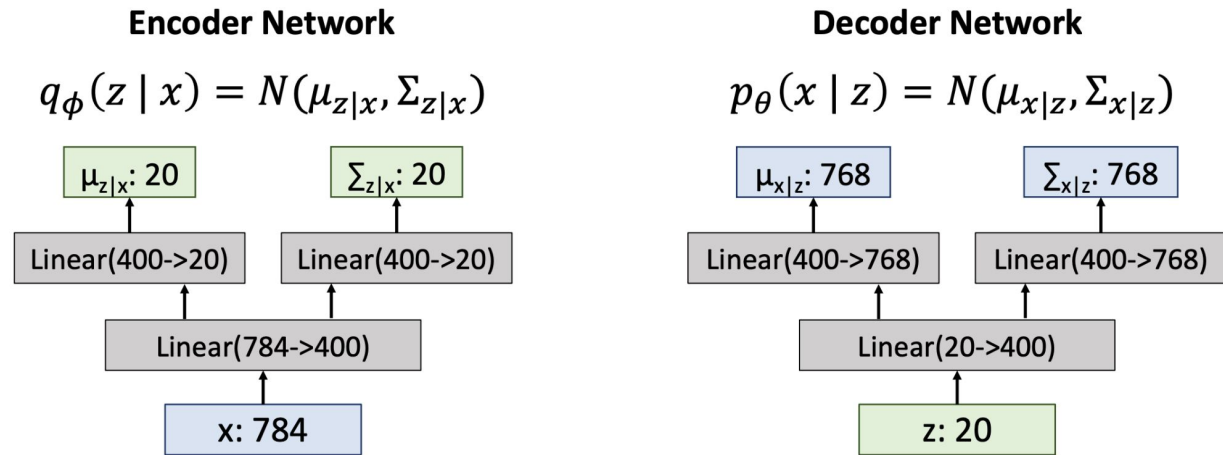
[van der Oord et al. 2016]

Variational Autoencoders: Intractability

Data likelihood: $p_{\theta}(x) = \int p_{\theta}(z)p_{\theta}(x|z)dz$

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

x: 28x28 image = 784-dim vector
z: 20-dim vector



Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Today: implicit density models

Generative Adversarial Networks (GANs)

All the models together

Autoregressive models define tractable density function, optimize likelihood of training data:

$$p_{\theta}(x) = \prod_{i=1}^n p_{\theta}(x_i | x_1, \dots, x_{i-1})$$

VAEs define intractable density function with latent \mathbf{z} :

$$p_{\theta}(x) = \int p_{\theta}(z) p_{\theta}(x|z) dz$$

Cannot optimize directly, derive and optimize lower bound on likelihood instead

So far...

Autoregressive define tractable density function, optimize likelihood of training data:

$$p_{\theta}(x) = \prod_{i=1}^n p_{\theta}(x_i | x_1, \dots, x_{i-1})$$

VAEs define intractable density function with latent \mathbf{z} :

$$p_{\theta}(x) = \int p_{\theta}(z) p_{\theta}(x|z) dz$$

Cannot optimize directly, derive and optimize lower bound on likelihood instead

What if we give up on explicitly modeling density, and just want ability to sample?

So far...

Autoregressive define tractable density function, optimize likelihood of training data:

$$p_{\theta}(x) = \prod_{i=1}^n p_{\theta}(x_i | x_1, \dots, x_{i-1})$$

VAEs define intractable density function with latent \mathbf{z} :

$$p_{\theta}(x) = \int p_{\theta}(z) p_{\theta}(x|z) dz$$

Cannot optimize directly, derive and optimize lower bound on likelihood instead

What if we give up on explicitly modeling density, and just want ability to sample?
GANs: not modeling any explicit density function!

Taxonomy of Generative Models

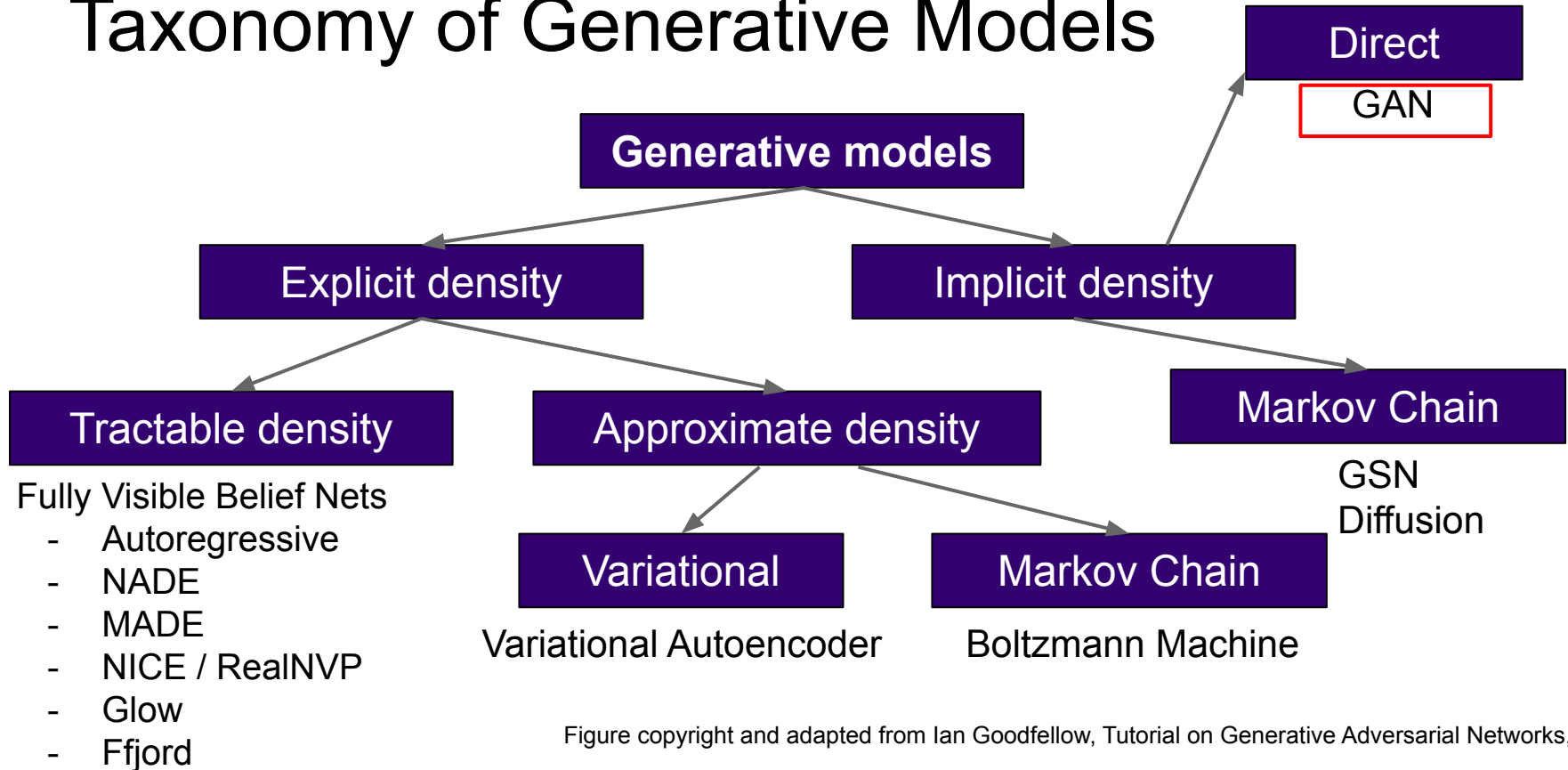


Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

Generative Adversarial Networks

Setup: Assume we have data x_i drawn from distribution $p_{\text{data}}(x)$. Want to sample from p_{data} .

Generative Adversarial Networks

Setup: Assume we have data x_i drawn from distribution $p_{\text{data}}(x)$. Want to sample from p_{data} .

Idea: Introduce a latent variable z with simple prior $p(z)$ (e.g. assume z is a multivariate gaussian).
Sample $z \sim p(z)$ and pass to a Generator Network $x = G(z)$

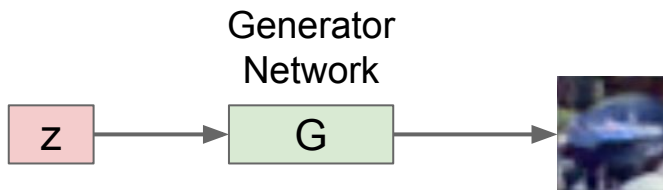
Then x is a sample from the Generator distribution p_G . **We just need to make sure $p_G = p_{\text{data}}$!**

Generative Adversarial Networks

Setup: Assume we have data x_i drawn from distribution $p_{\text{data}}(x)$. Want to sample from p_{data} .

Idea: Introduce a latent variable z with simple prior $p(z)$ (e.g. assume z is a multivariate gaussian). Sample $z \sim p(z)$ and pass to a Generator Network $x = G(z)$

Then x is a sample from the Generator distribution p_G . **We just need to make sure $p_G = p_{\text{data}}$!**



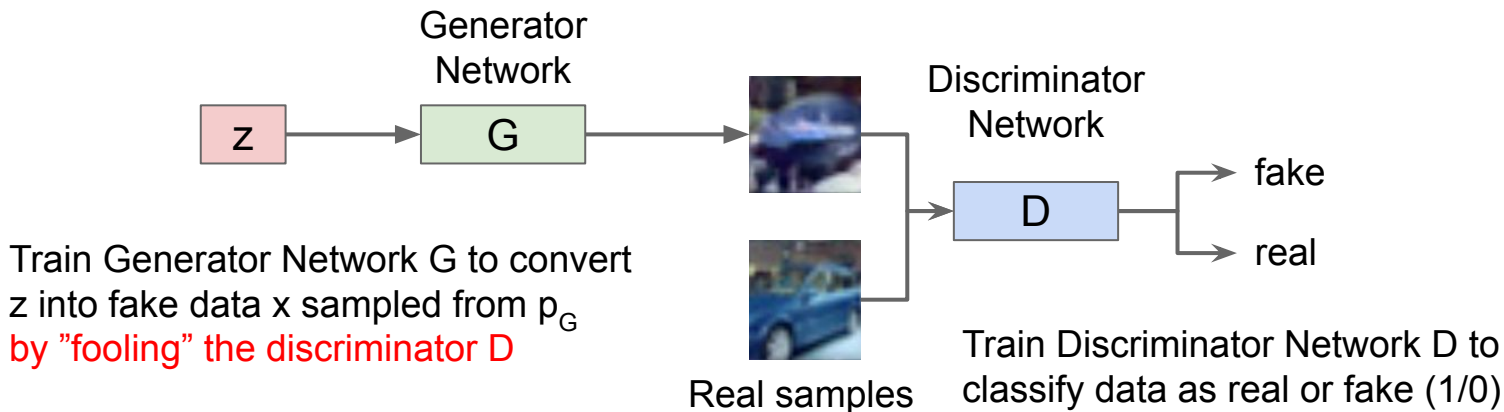
Train Generator Network G to convert z into fake data x sampled from p_G

Generative Adversarial Networks

Setup: Assume we have data x_i drawn from distribution $p_{\text{data}}(x)$. Want to sample from p_{data} .

Idea: Introduce a latent variable z with simple prior $p(z)$ (e.g. assume z is a multivariate gaussian). Sample $z \sim p(z)$ and pass to a Generator Network $x = G(z)$

Then x is a sample from the Generator distribution p_G . **We just need to make sure $p_G = p_{\text{data}}$!**

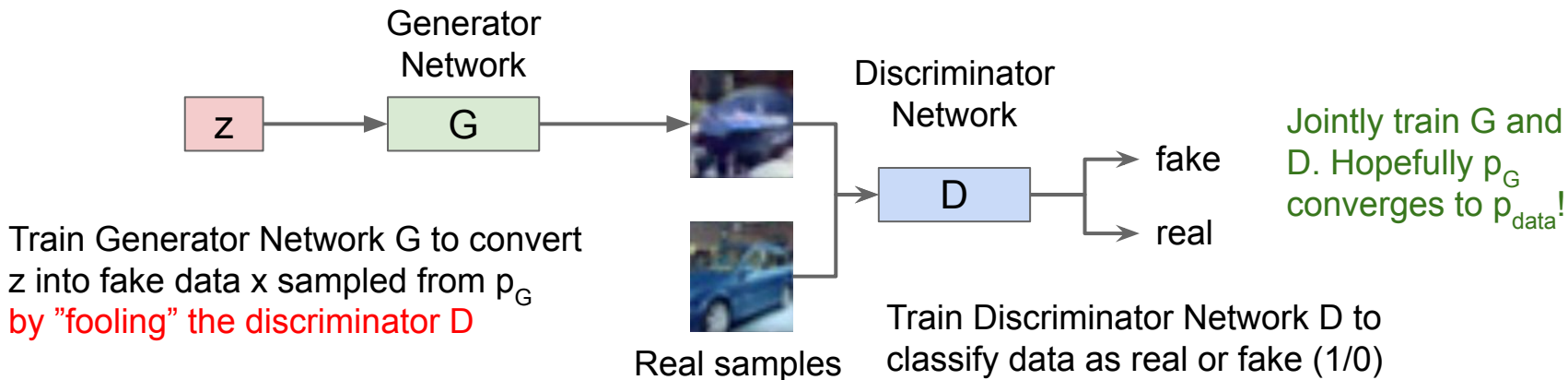


Generative Adversarial Networks

Setup: Assume we have data x_i drawn from distribution $p_{\text{data}}(x)$. Want to sample from p_{data} .

Idea: Introduce a latent variable z with simple prior $p(z)$ (e.g. assume z is a multivariate gaussian). Sample $z \sim p(z)$ and pass to a Generator Network $x = G(z)$

Then x is a sample from the Generator distribution p_G . **We just need to make sure $p_G = p_{\text{data}}$!**



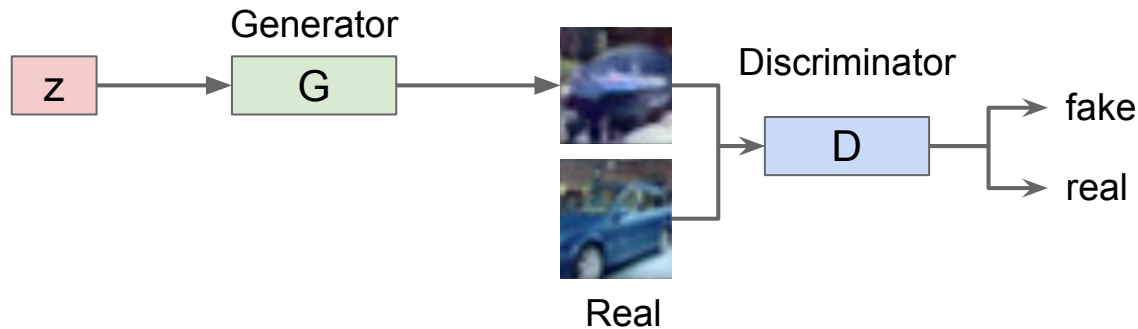
Train Generator Network G to convert z into fake data x sampled from p_G by "fooling" the discriminator D

Train Discriminator Network D to classify data as real or fake (1/0)

Generative Adversarial Networks

Jointly train generator G and discriminator D with a minimax game

$$\min_G \max_D \left(E_{x \sim p_{data}} [\log D(x)] + E_{z \sim p(z)} [\log (1 - D(G(z)))] \right)$$

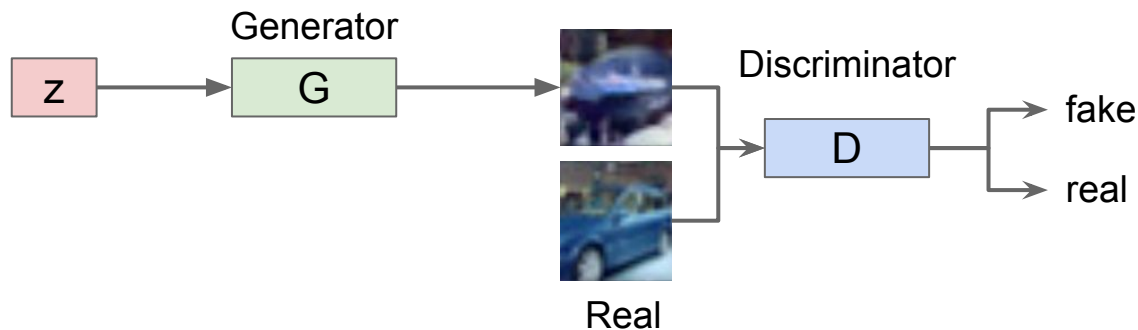


Generative Adversarial Networks

Jointly train generator G and discriminator D with a minimax game

Discriminator wants
 $D(x) = 1$ for real data

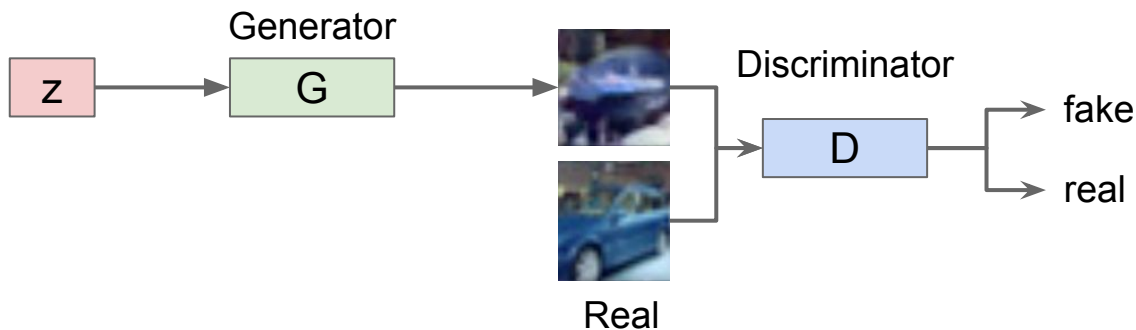
$$\min_G \max_D \left(E_{x \sim p_{data}} [\log D(x)] + E_{z \sim p(z)} [\log (1 - D(G(z)))] \right)$$



Generative Adversarial Networks

Jointly train generator G and discriminator D with a minimax game

$$\min_G \max_D \left(\overbrace{E_{x \sim p_{data}} [\log D(x)]}^{\text{Discriminator wants } D(x) = 1 \text{ for real data}} + \overbrace{E_{z \sim p(z)} [\log (1 - D(G(z)))]}^{\text{Discriminator wants } D(x) = 0 \text{ for fake data}} \right)$$

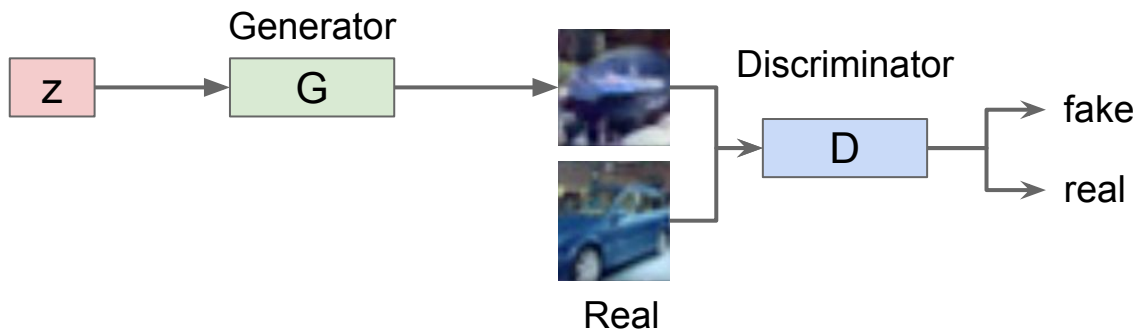


Generative Adversarial Networks

Jointly train generator G and discriminator D with a minimax game

$$\min_G \max_D \left(E_{x \sim p_{data}} [\log D(x)] + E_{z \sim p(z)} [\log (1 - D(G(z)))] \right)$$

Generator wants $D(x) = 1$ for fake data



Generative Adversarial Networks

Jointly train generator G and discriminator D with a minimax game

Train G and D using alternating gradient updates

$$\min_G \max_D \left(E_{x \sim p_{data}} [\log D(x)] + E_{z \sim p(z)} \left[\log \left(1 - D(G(z)) \right) \right] \right)$$

Generative Adversarial Networks

Jointly train generator G and discriminator D with a minimax game

Train G and D using alternating gradient updates

$$\min_G \max_D \left(E_{x \sim p_{data}} [\log D(x)] + E_{z \sim p(z)} \left[\log \left(1 - D(G(z)) \right) \right] \right)$$
$$= \min_G \max_D V(G, D)$$

Generative Adversarial Networks

Jointly train generator G and discriminator D with a minimax game

Train G and D using alternating gradient updates

$$\min_G \max_D \left(E_{x \sim p_{data}} [\log D(x)] + E_{z \sim p(z)} \left[\log \left(1 - D(G(z)) \right) \right] \right)$$

$$= \min_G \max_D V(G, D)$$

For t in $1, \dots, T$:

1. (Update D) $D = D + \alpha_D \frac{\partial V}{\partial D}$
2. (Update G) $G = G - \alpha_G \frac{\partial V}{\partial G}$

Generative Adversarial Networks

Jointly train generator G and discriminator D with a minimax game

Train G and D using alternating gradient updates

$$\min_G \max_D \left(E_{x \sim p_{data}} [\log D(x)] + E_{z \sim p(z)} \left[\log \left(1 - D(G(z)) \right) \right] \right)$$

$$= \min_G \max_D V(G, D)$$

We are not minimizing any overall loss! No training curves to look at!

For t in $1, \dots, T$:

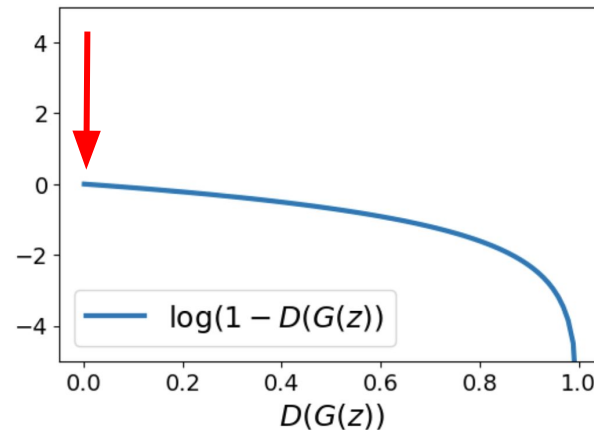
1. (Update D) $D = D + \alpha_D \frac{\partial V}{\partial D}$
2. (Update G) $G = G - \alpha_G \frac{\partial V}{\partial G}$

Generative Adversarial Networks

Jointly train generator G and discriminator D with a minimax game

$$\min_G \max_D \left(E_{x \sim p_{data}} [\log D(x)] + E_{z \sim p(z)} \left[\log \left(1 - D(G(z)) \right) \right] \right)$$

At start of training, generator is very bad and discriminator can easily tell apart real/fake, so $D(G(z))$ close to 0



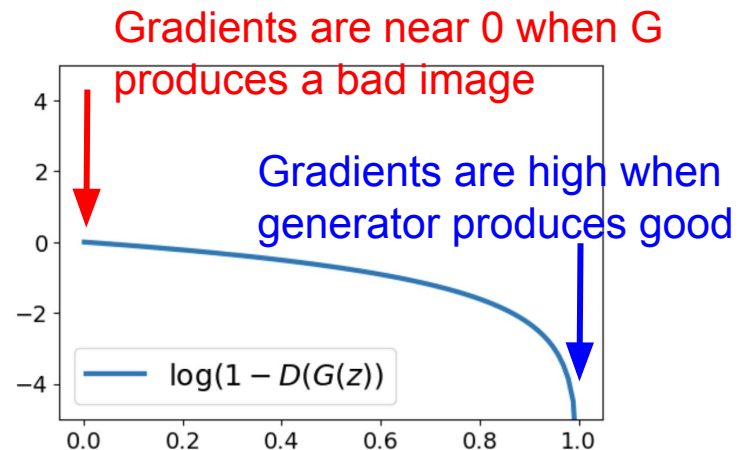
Generative Adversarial Networks

Jointly train generator G and discriminator D with a minimax game

$$\min_G \max_D \left(E_{x \sim p_{data}} [\log D(x)] + E_{z \sim p(z)} \left[\log \left(1 - D(G(z)) \right) \right] \right)$$

At start of training, generator is very bad and discriminator can easily tell apart real/fake, so $D(G(z))$ close to 0

Problem: Why is this a problem?



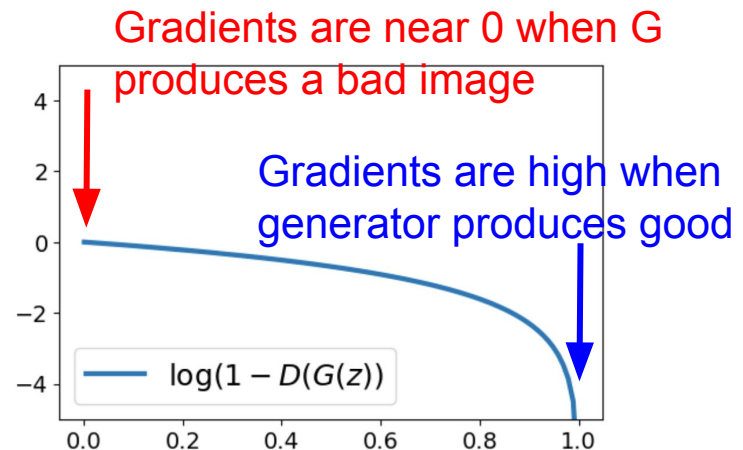
Generative Adversarial Networks

Jointly train generator G and discriminator D with a minimax game

$$\min_G \max_D \left(E_{x \sim p_{data}} [\log D(x)] + E_{z \sim p(z)} \left[\log \left(1 - D(G(z)) \right) \right] \right)$$

At start of training, generator is very bad and discriminator can easily tell apart real/fake, so $D(G(z))$ close to 0

Problem: Vanishing gradients for G
How do we fix this?



Generative Adversarial Networks

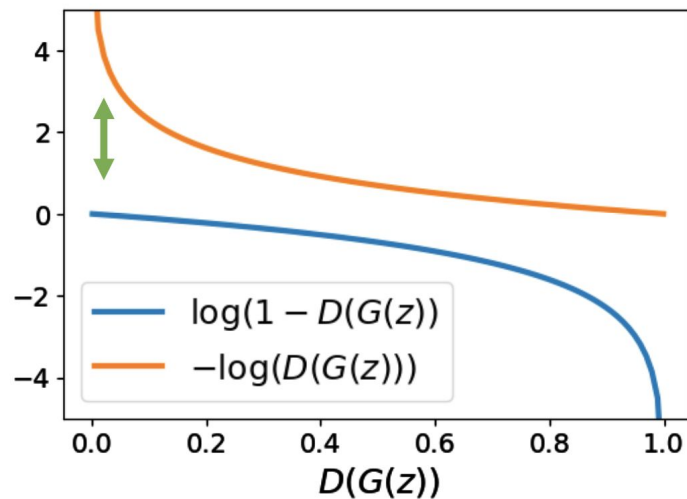
Jointly train generator G and discriminator D with a minimax game

$$\min_G \max_D \left(E_{x \sim p_{data}} [\log D(x)] + E_{z \sim p(z)} \left[\log \left(1 - D(G(z)) \right) \right] \right)$$

At start of training, generator is very bad and discriminator can easily tell apart real/fake, so $D(G(z))$ close to 0

Problem: Vanishing gradients for G

Solution: Train G to minimize $-\log(D(G(z)))$, instead of $\log(1-D(G(z)))$. Then G gets strong gradients at start of training!



Generative Adversarial Networks

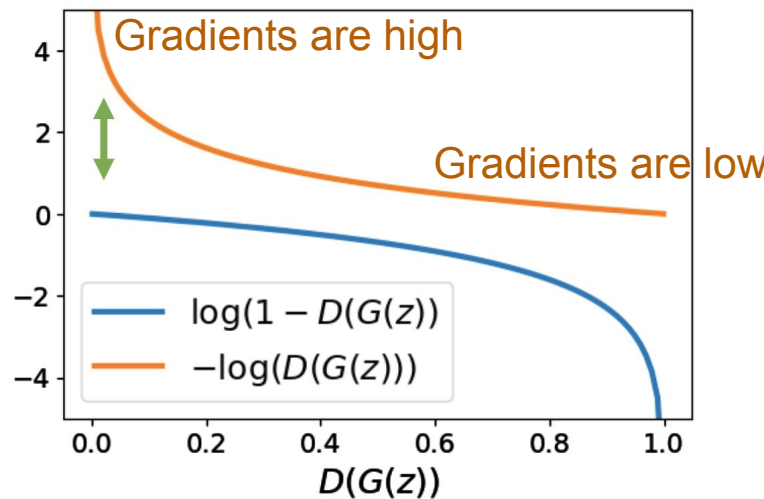
Jointly train generator G and discriminator D with a minimax game

$$\min_G \max_D \left(E_{x \sim p_{data}} [\log D(x)] + E_{z \sim p(z)} \left[\log \left(1 - D(G(z)) \right) \right] \right)$$

At start of training, generator is very bad and discriminator can easily tell apart real/fake, so $D(G(z))$ close to 0

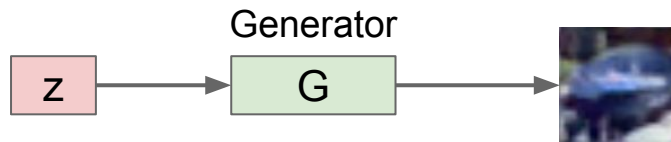
Problem: Vanishing gradients for G

Solution: Train G to minimize $-\log(D(G(z)))$, instead of $\log(1-D(G(z)))$. Then G gets strong gradients at start of training!



Generative adversarial networks

Once trained, throw away the discriminator and use G to generate new images



Generative Adversarial Nets

Generated samples



Nearest neighbor from training set



Figures copyright Ian Goodfellow et al., 2014. Reproduced with permission.

Generative Adversarial Nets

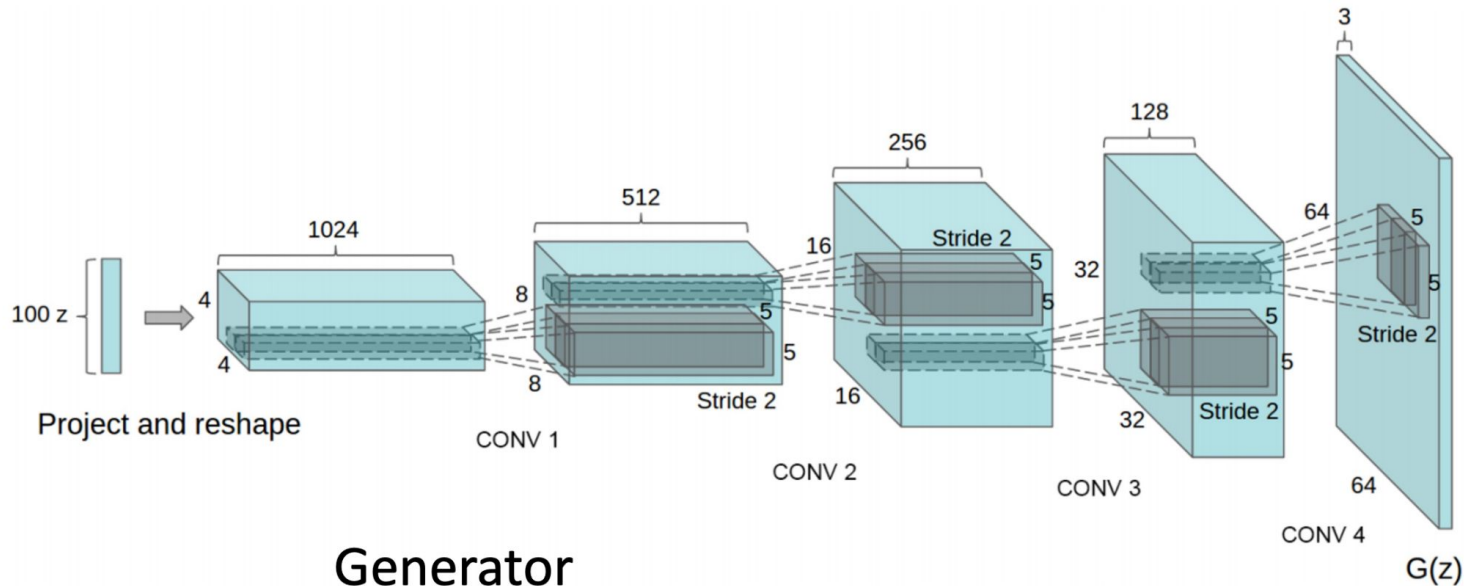
Generated samples (CIFAR-10)



Nearest neighbor from training set

Figures copyright Ian Goodfellow et al., 2014. Reproduced with permission.

Generative Adversarial Nets: Convolutional Architectures



Radford et al, "Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks", ICLR 2016

Generative Adversarial Nets: Convolutional Architectures

Generator is an upsampling network with fractionally-strided convolutions
Discriminator is a convolutional network

Architecture guidelines for stable Deep Convolutional GANs

- Replace any pooling layers with strided convolutions (discriminator) and fractional-strided convolutions (generator).
- Use batchnorm in both the generator and the discriminator.
- Remove fully connected hidden layers for deeper architectures.
- Use ReLU activation in generator for all layers except for the output, which uses Tanh.
- Use LeakyReLU activation in the discriminator for all layers.

Radford et al, "Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks", ICLR 2016

Generative Adversarial Nets: Convolutional Architectures

Samples from the model look much better!



Radford et al,
ICLR 2016

Generative Adversarial Nets: Convolutional Architectures

Interpolating
between
random
points in latent
space



Radford et al,
ICLR 2016

Generative Adversarial Nets: Interpretable Vector Math

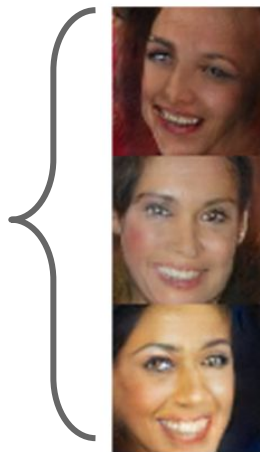
Radford et al, ICLR 2016

Smiling woman

Neutral woman

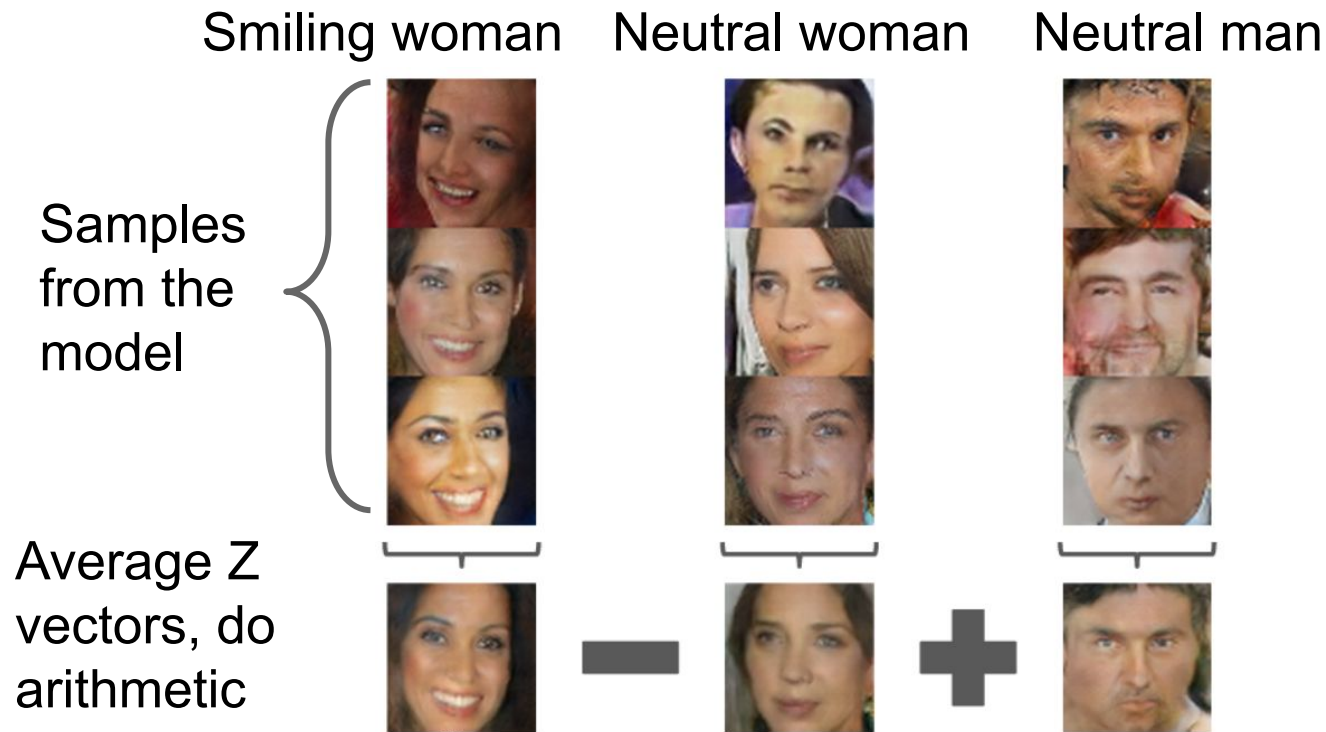
Neutral man

Samples
from the
model



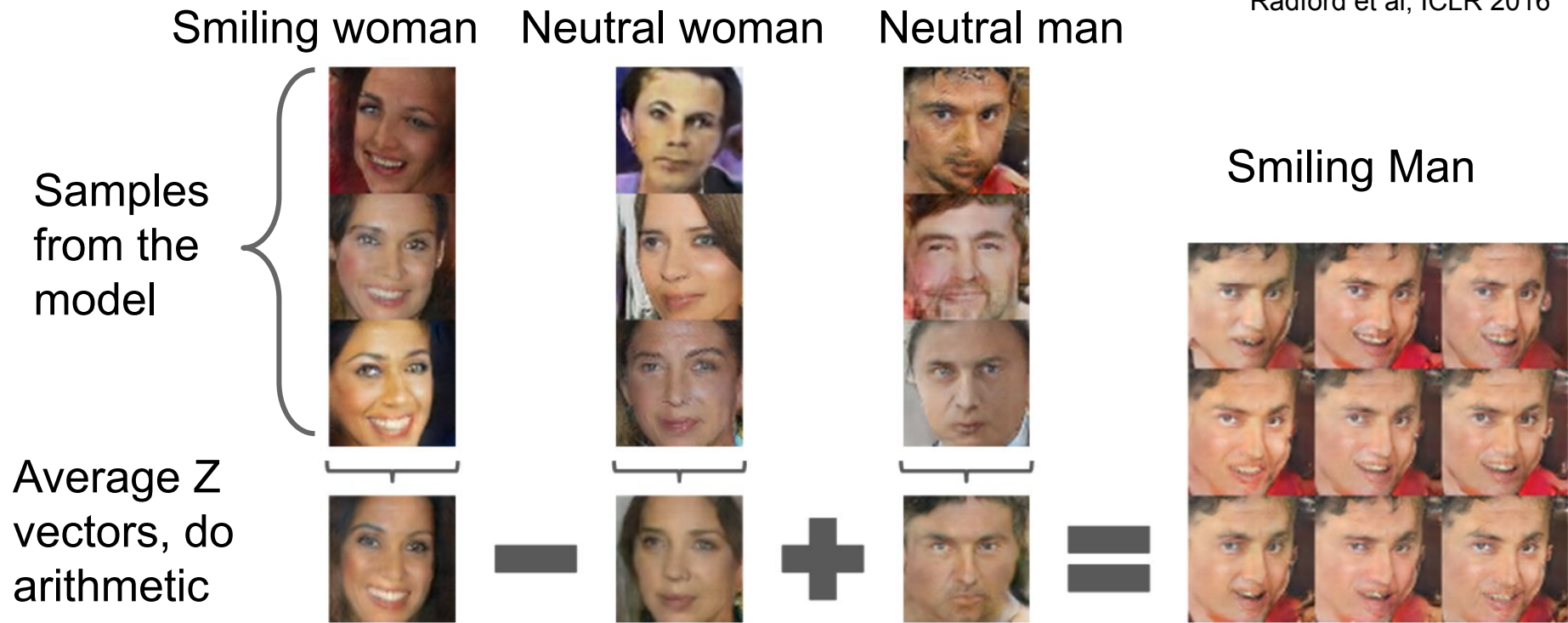
Generative Adversarial Nets: Interpretable Vector Math

Radford et al, ICLR 2016



Generative Adversarial Nets: Interpretable Vector Math

Radford et al, ICLR 2016



Generative Adversarial Nets: Interpretable Vector Math

Glasses man



No glasses man



No glasses woman



-

+

=

Radford et al,
ICLR 2016

Woman with glasses



Since then: Explosion of GANs

“The GAN Zoo”

See also: <https://github.com/soumith/ganhacks> for tips and tricks for trainings GANs

- GAN - Generative Adversarial Networks
- 3D-GAN - Learning a Probabilistic Latent Space of Object Shapes via 3D Generative-Adversarial Modeling
- acGAN - Face Aging With Conditional Generative Adversarial Networks
- AC-GAN - Conditional Image Synthesis With Auxiliary Classifier GANs
- AdaGAN - AdaGAN: Boosting Generative Models
- AEGAN - Learning Inverse Mapping by Autoencoder based Generative Adversarial Nets
- AffGAN - Amortised MAP Inference for Image Super-resolution
- AL-CGAN - Learning to Generate Images of Outdoor Scenes from Attributes and Semantic Layouts
- ALI - Adversarially Learned Inference
- AM-GAN - Generative Adversarial Nets with Labeled Data by Activation Maximization
- AnoGAN - Unsupervised Anomaly Detection with Generative Adversarial Networks to Guide Marker Discovery
- ArtGAN - ArtGAN: Artwork Synthesis with Conditional Categorical GANs
- b-GAN - b-GAN: Unified Framework of Generative Adversarial Networks
- Bayesian GAN - Deep and Hierarchical Implicit Models
- BEGAN - BEGAN: Boundary Equilibrium Generative Adversarial Networks
- BiGAN - Adversarial Feature Learning
- BS-GAN - Boundary-Seeking Generative Adversarial Networks
- CGAN - Conditional Generative Adversarial Nets
- CaloGAN - CaloGAN: Simulating 3D High Energy Particle Showers in Multi-Layer Electromagnetic Calorimeters with Generative Adversarial Networks
- CCGAN - Semi-Supervised Learning with Context-Conditional Generative Adversarial Networks
- CatGAN - Unsupervised and Semi-supervised Learning with Categorical Generative Adversarial Networks
- CoGAN - Coupled Generative Adversarial Networks
- Context-RNN-GAN - Contextual RNN-GANs for Abstract Reasoning Diagram Generation
- C-RNN-GAN - C-RNN-GAN: Continuous recurrent neural networks with adversarial training
- CS-GAN - Improving Neural Machine Translation with Conditional Sequence Generative Adversarial Nets
- CVAE-GAN - CVAE-GAN: Fine-Grained Image Generation through Asymmetric Training
- CycleGAN - Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks
- DTN - Unsupervised Cross-Domain Image Generation
- DCGAN - Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks
- DiscoGAN - Learning to Discover Cross-Domain Relations with Generative Adversarial Networks
- DR-GAN - Disentangled Representation Learning GAN for Pose-Invariant Face Recognition
- DualGAN - DualGAN: Unsupervised Dual Learning for Image-to-Image Translation
- EBGAN - Energy-based Generative Adversarial Network
- f-GAN - f-GAN: Training Generative Neural Samplers using Variational Divergence Minimization
- FF-GAN - Towards Large-Pose Face Frontalization in the Wild
- GAWWN - Learning What and Where to Draw
- GeneGAN - GeneGAN: Learning Object Transfiguration and Attribute Subspace from Unpaired Data
- Geometric GAN - Geometric GAN
- GoGAN - Gang of GANs: Generative Adversarial Networks with Maximum Margin Ranking
- GP-GAN - GP-GAN: Towards Realistic High-Resolution Image Blending
- IAN - Neural Photo Editing with Introspective Adversarial Networks
- iGAN - Generative Visual Manipulation on the Natural Image Manifold
- IcGAN - Invertible Conditional GANs for image editing
- ID-CGAN - Image De-raining Using a Conditional Generative Adversarial Network
- Improved GAN - Improved Techniques for Training GANs
- InfoGAN - InfoGAN: Interpretable Representation Learning by Information Maximizing Generative Adversarial Nets
- LAGAN - Learning Particle Physics by Example: Location-Aware Generative Adversarial Networks for Physics Synthesis
- LAPGAN - Deep Generative Image Models using a Laplacian Pyramid of Adversarial Networks

<https://github.com/hindupuravinash/the-gan-zoo>

GAN improvements: better loss functions



LSGAN, Zhu 2017.

Wasserstein GAN, Arjovsky 2017.



Improved Wasserstein GAN, Gulrajani 2017.

GAN improvements: higher resolution

256 x 256 bedrooms



1024 x 1024 faces



Progressive GAN, Karras 2018.

GAN transformations

Source->Target domain transfer



CycleGAN. Zhu et al. 2017.



Pix2pix. Isola 2017. Many examples at <https://phillipi.github.io/pix2pix/>

BigGAN: 512x512 images



Brock et al., 2019

GANs with self-attention mechanism

Goldfish



Indigo bunting



Redshank



Saint Bernard



Zhang et al, "Self-Attention Generative Adversarial Networks", ICML 2019

Controlled generation with GANs

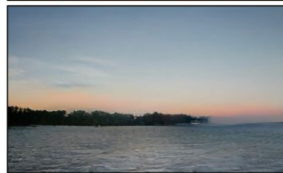
cloud	sky
tree	mountain
sea	grass



Semantic Manipulation Using Segmentation Map →



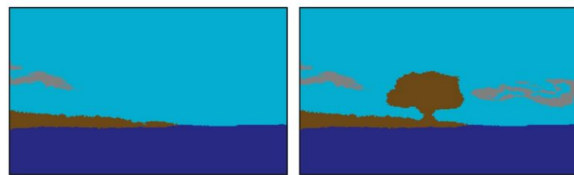
Stylization using Guide Images ↓



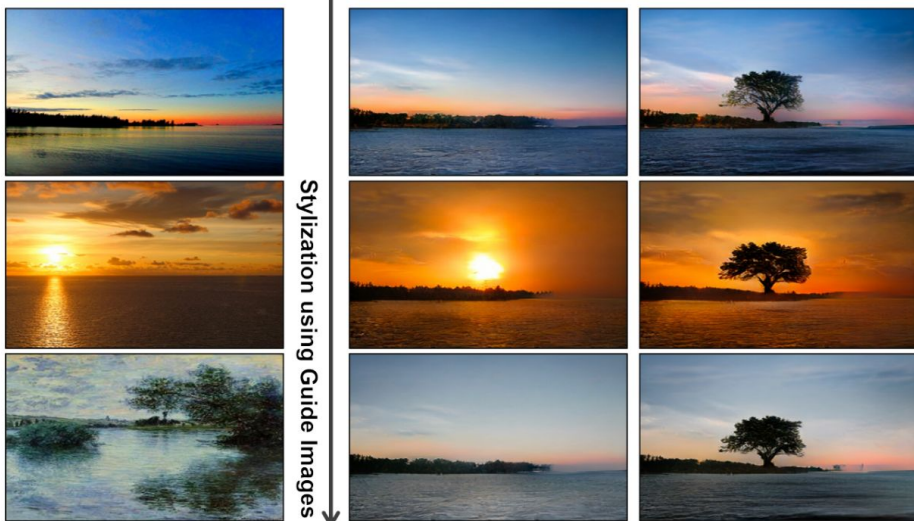
Park et al, "Semantic Image Synthesis with Spatially-Adaptive Normalization", CVPR 2019

Controlled generation with GANs

cloud	sky
tree	mountain
sea	grass



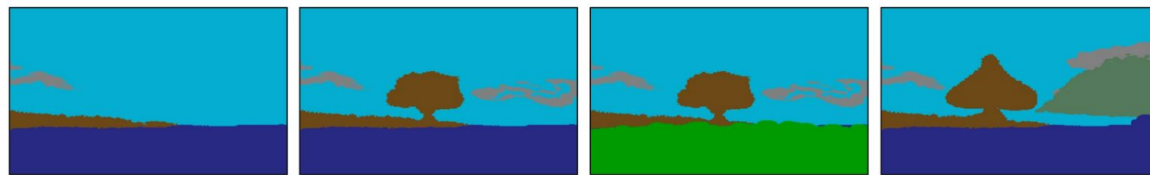
Semantic Manipulation Using Segmentation Map →



Park et al, "Semantic Image Synthesis with Spatially-Adaptive Normalization", CVPR 2019

Controlled generation with GANs

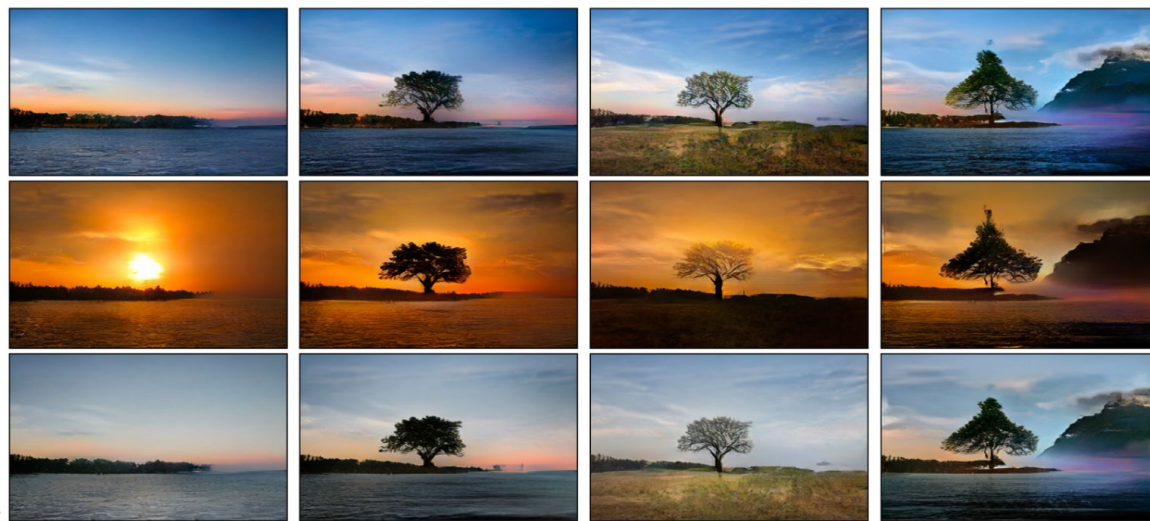
cloud	sky
tree	mountain
sea	grass



Semantic Manipulation Using Segmentation Map →



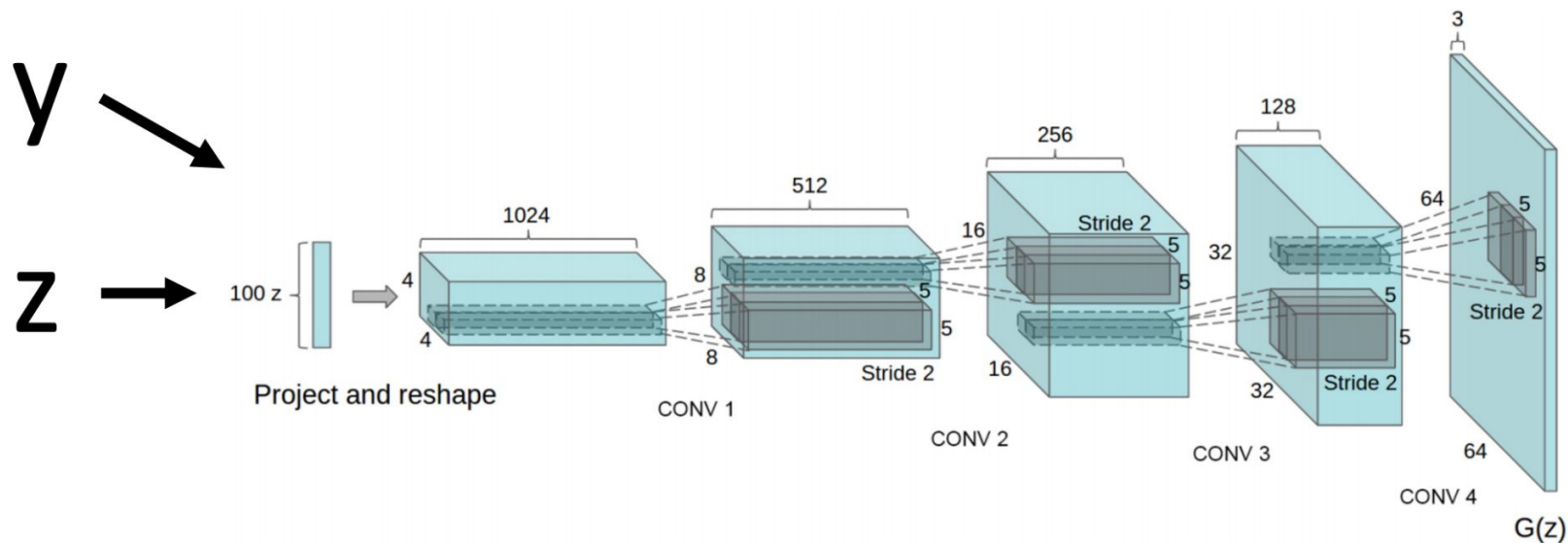
Stylization using Guide Images ↓



Park et al, "Semantic Image Synthesis with Spatially-Adaptive Normalization", CVPR 2019

Conditional GANs: StyleGAN

Y is text that describes the image you want to generate



Karras et al, "Analyzing and Improving the Image Quality of StyleGAN", CVPR 2020

Conditional GANs: StyleGAN

Batch Normalization

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$
$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$
$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$
$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$



Learn a separate
scale and shift
for each
different label y

Conditional Batch Normalization

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$
$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$
$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$
$$y_{i,j} = \gamma_j^y \hat{x}_{i,j} + \beta_j^y$$

Karras et al, "Analyzing and Improving the Image Quality of StyleGAN", CVPR 2020

Conditional GANs: StyleGAN



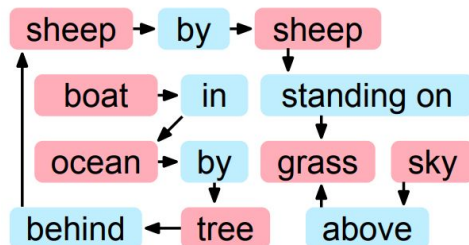
Karras et al, "Analyzing and Improving the Image Quality of StyleGAN", CVPR 2020

Scene graphs to GANs

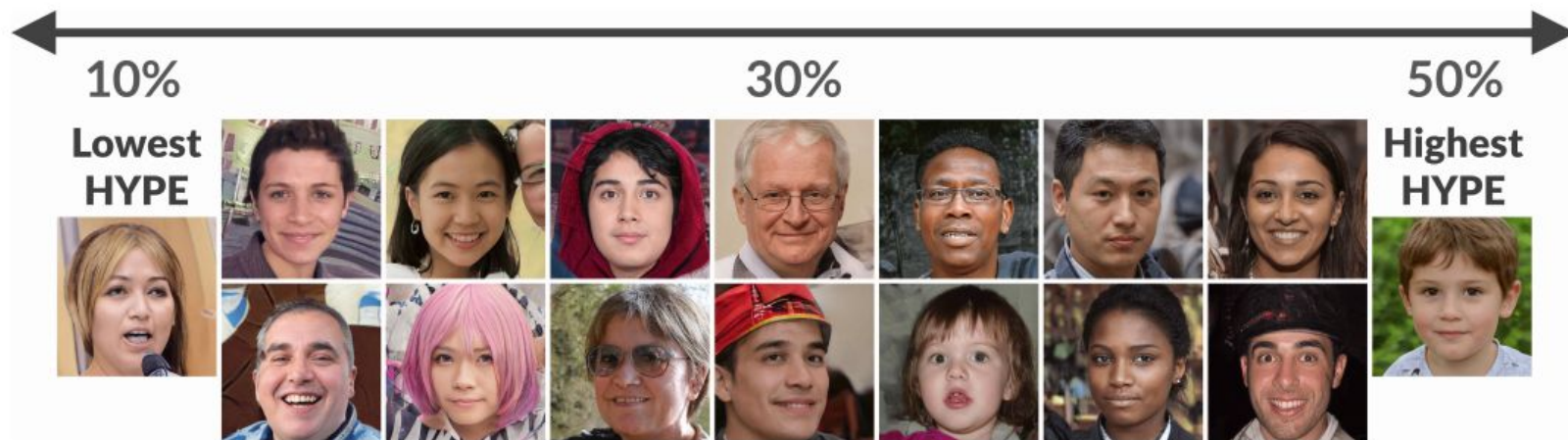
Specifying exactly what kind of image you want to generate.

The explicit structure in scene graphs provides better image generation for complex scenes.

Scene Graph



HYPE: Human eYe Perceptual Evaluations



Zhou, Gordon, Krishna et al. HYPE: Human eYe Perceptual Evaluations, NeurIPS 2019

Figures copyright 2019. Reproduced with permission.

Summary: GANs

Pros:

- Beautiful samples, was state-of-the-art until diffusion models!

Cons:

- Trickier / more unstable to train
- Can't solve inference queries such as $p(x)$, $p(z|x)$

Active areas of research:

- Better loss functions, more stable training (Wasserstein GAN, LSGAN, many others)
- Conditional GANs, GANs for all kinds of applications

Diffusion models

Diffusion Models are outperforming GANs



Dhariwal & Nichol. "Diffusion Models Beat GANs on Image Synthesis", OpenAI 2021



Ho et al. "Cascaded Diffusion Models for High Fidelity Image Generation", Google 2021

Text-to-Image (T2I) Generation

Dall-E2

“a teddy bear on a skateboard in times square”



Ramesh et al. “Hierarchical Text-Conditional Image Generation with CLIP Latents” 2022

Imagen

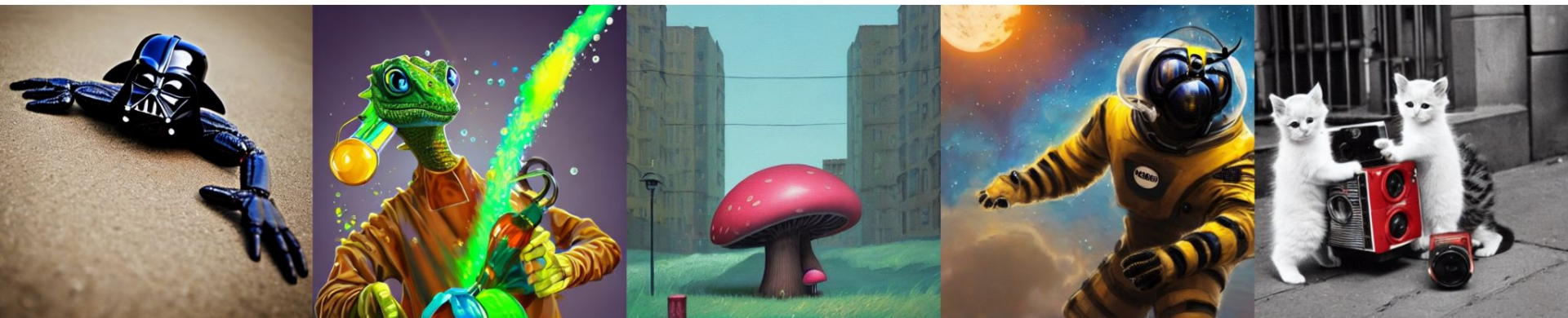
“A group of teddy bears in suit in a corporate office celebrating the birthday of their friend. There is a pizza cake on the desk.”



Saharia et al. “Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding” 2022

Text-to-Image (T2I) Generation

Stable Diffusion



[Mega thread on Twitter/X about Stable Diffusion](#)

Rombach et al. "High-Resolution Image Synthesis with Latent Diffusion Models" 2022

Application of diffusion: Image Super-resolution

Irish Setter

Saharia et al., Image Super-Resolution via Iterative Refinement, ICCV 2021

Gif on this slide is not
displayed in pdf

But what is a diffusion model?

So far...

Autoregressive define tractable density function, optimize likelihood of training data:

$$p_{\theta}(x) = \prod_{i=1}^n p_{\theta}(x_i | x_1, \dots, x_{i-1})$$

VAEs define intractable density function with latent \mathbf{z} :

$$p_{\theta}(x) = \int p_{\theta}(z) p_{\theta}(x|z) dz$$

Cannot optimize directly, derive and optimize lower bound on likelihood instead

GANs give up on explicitly modeling density and just learns to sample “real” data

So far...

Autoregressive define tractable density function, optimize likelihood of training data:

$$p_{\theta}(x) = \prod_{i=1}^n p_{\theta}(x_i | x_1, \dots, x_{i-1})$$

VAEs define intractable density function with latent \mathbf{z} :

$$p_{\theta}(x) = \int p_{\theta}(z) p_{\theta}(x|z) dz$$

Cannot optimize directly, derive and optimize lower bound on likelihood instead

GANs give up on explicitly modeling density and just learns to sample “real” data

All these methods generate data in one forward step! Why this is hard?

Taxonomy of Generative Models

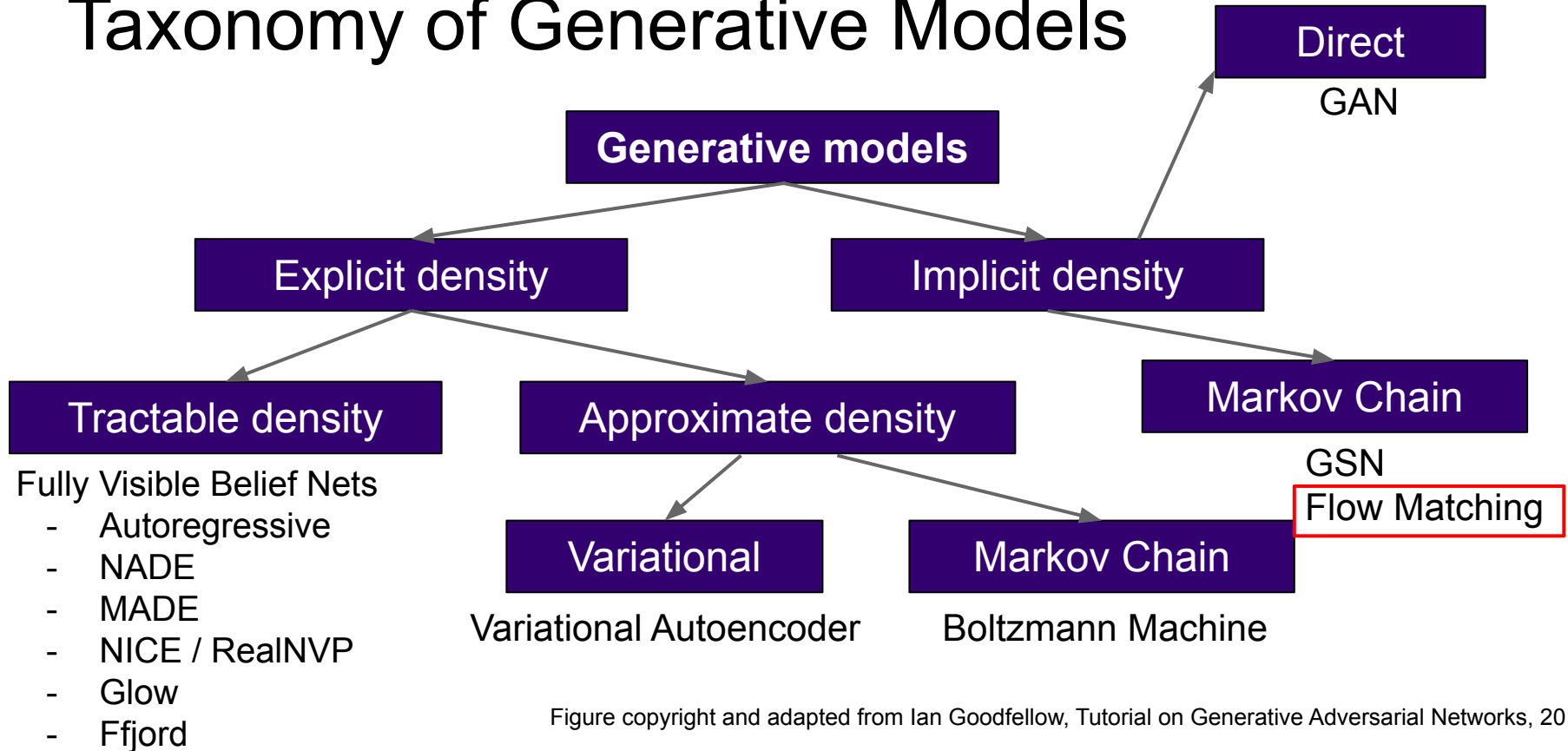
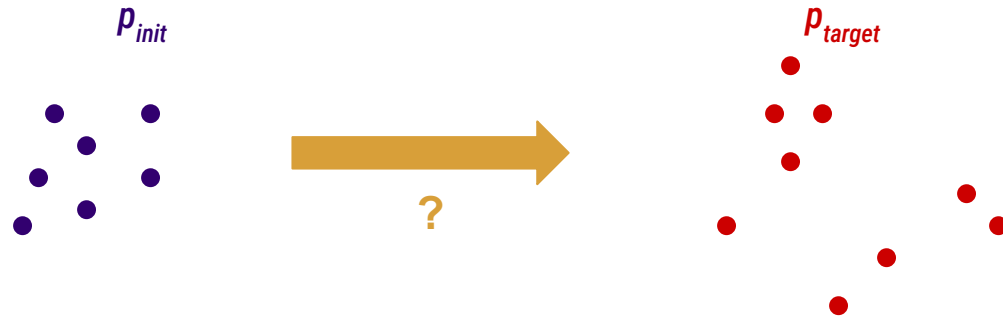


Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

Laying the groundwork - Intro

We aim to transform samples from an *initial distribution* to an *unknown target one*



Flow models map from p_{init} to p_{target} by simulating *ordinary differential equations (ODEs)*¹

and are leveraged in many state-of-the-art models across diverse applications

1. They also simulate stochastic differential equations (SDEs), but this is outside the scope of this lecture

Laying the groundwork - ODEs

What are ordinary differential equations (ODEs)?

ODEs are defined by a **vector field** u , where $u: \mathbb{R}^d \times [0, 1] \rightarrow \mathbb{R}^d, (x,t) \rightarrow u_t(x)$

Basically, for every time t (bounded between $[0, 1]$) and a location $x \in \mathbb{R}^d$ we get a vector $u_t(x)$ that specifies a **velocity** in the \mathbb{R}^d space

A solution to an ODE is defined by a **trajectory**: $X: [0,1] \rightarrow \mathbb{R}^d, t \rightarrow X_t$

This effectively maps from some time t to some location in the \mathbb{R}^d space

Trajectories are solutions to the equation,

$\frac{d}{dt} X_t = u_t(X_t)$	▶ ODE	The derivative at position X_t should follow the vector field $u_t(x)$
-------------------------------	-------	--

$X_0 = x_0$	▶ initial conditions
-------------	----------------------

At $t=0$, we should start at X_0

Laying the groundwork - ODEs

Okay, but where are we at at time t (what is X_t)?

We define the **flow**, also a solution to the ODE: $\psi : \mathbb{R}^d \times [0, 1] \mapsto \mathbb{R}^d$, $(x_0, t) \mapsto \psi_t(x_0)$

The *flow* models trajectories as function of an input x_0 and whose output is the trajectory position at t .

$$\begin{array}{ll} \frac{d}{dt} X_t = u_t(X_t) & \blacktriangleright \text{ODE} \\ X_0 = x_0 & \blacktriangleright \text{initial conditions} \end{array} \quad \longrightarrow \quad \begin{array}{ll} \frac{d}{dt} \psi_t(x_0) = u_t(\psi_t(x_0)) & \blacktriangleright \text{flow ODE} \\ \psi_0(x_0) = x_0 & \blacktriangleright \text{flow initial conditions} \end{array}$$

Vector fields define ODEs whose solutions are flows

Laying the groundwork - ODEs

If the vector field is *continuously differentiable* and *has bounded derivative*...
... solutions to the ODE are both *guaranteed to exist* and *are unique*

!! Problem: It's not possible to compute the flow explicitly when the vector field is non-linear (usually the case).

 **Solution:** Use numerical methods to simulate the ODE (Euler method)

$$X_0 = x_0 \longrightarrow X_{t+h} = X_t + hu_t(X_t) \quad (t = 0, h, 2h, 3h, \dots, 1 - h)$$

Laying the groundwork - The flow model

Constructing a generative model via an ODE: the **flow** model

$$\begin{array}{ll} X_0 \sim p_{\text{init}} & \blacktriangleright \text{random initialization} \\ \frac{d}{dt} X_t = u_t^\theta(X_t) & \blacktriangleright \text{ODE} \end{array} \quad \longrightarrow \quad \begin{array}{ll} \frac{d}{dt} X_t = u_t(X_t) & \blacktriangleright \text{ODE} \\ X_0 = x_0 & \blacktriangleright \text{initial conditions} \end{array}$$

1. u_t^θ is a neural network with parameters $\theta: \mathbb{R}^d \times [0, 1] \rightarrow \mathbb{R}^d$
2. Objective is to make the endpoint X_1 have the same distribution as p_{target} (p_{data}).

$$X_1 \sim p_{\text{data}} \Leftrightarrow \psi_1^\theta(X_0) \sim p_{\text{data}}$$

ψ_t^θ is the flow induced by u_t^θ

Algorithm 1 Sampling from a Flow Model with Euler method

Require: Neural network vector field u_t^θ , number of steps n

1: Set $t = 0$

2: Set step size $h = \frac{1}{n}$

3: Draw a sample $X_0 \sim p_{\text{init}}$

4: **for** $i = 1, \dots, n - 1$ **do**

5: $X_{t+h} = X_t + hu_t^\theta(X_t)$

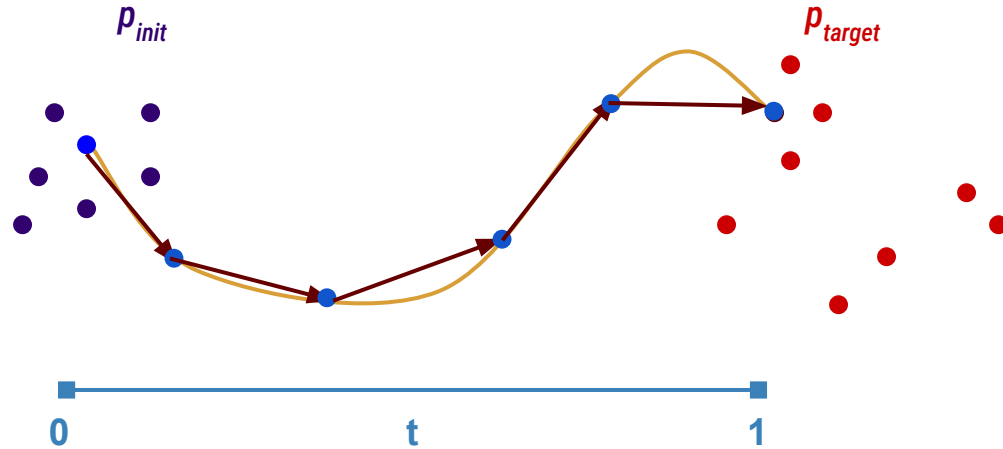
6: Update $t \leftarrow t + h$

7: **end for**

8: **return** X_1

Laying the groundwork - The flow model

Illustrating the flow model sampling



How do we train flow models? The objective

There are many ways to train flow models, yet the **simplest** is sufficient for SoTA

$$\mathcal{L}(\theta) = \|u_t^\theta(x) - \underbrace{u_t^{\text{target}}(x)}_{\text{training target}}\|^2.$$

u_t^{target} should: (1) Be a vector field & (2) Convert from p_{init} to p_{data}

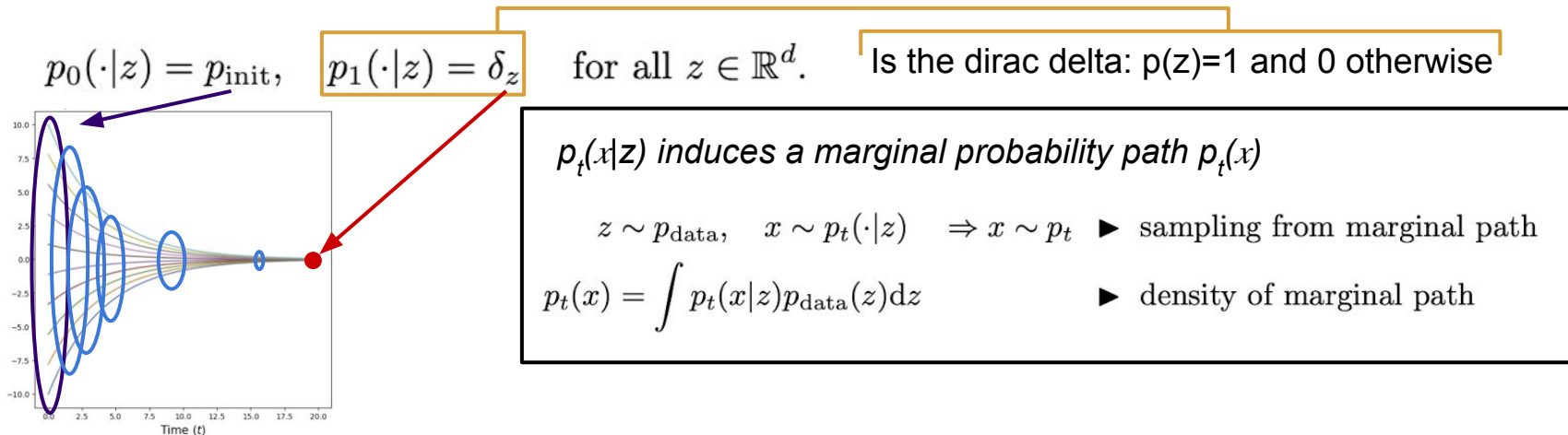
But how do we go about **constructing** u_t^{target} ?

How do we train flow models? The objective

We construct u_t^{target} by specifying a **probability path**

Probability paths describe how to navigate from p_{init} to p_{data} in **probability space**

For a data point $z \in \mathbb{R}^d$, a probability path is denoted by $p_t(x|z)$ over \mathbb{R}^d s.t.,



How do we train flow models? The objective

Probability paths are predominantly specified by *Gaussian distributions*

Let α_t and β_t be two continuously differentiable monotonic functions s.t.,

$$\alpha_0 = \beta_1 = 0 \text{ and } \alpha_1 = \beta_0 = 1$$

Then, the Gaussian probability path is defined as,

$$p_t(\cdot|z) = \mathcal{N}(\alpha_t z, \beta_t^2 I_d) \quad \blacktriangleright \text{ Gaussian conditional path}$$

So its endpoints are represented as,

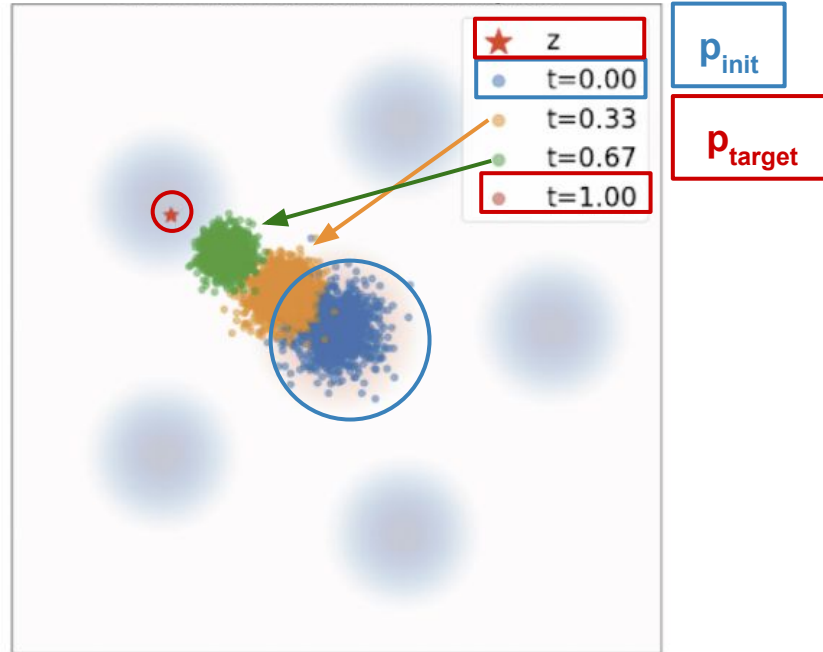
$$p_0(\cdot|z) = \mathcal{N}(\alpha_0 z, \beta_0^2 I_d) = \mathcal{N}(0, I_d), \quad \text{and} \quad p_1(\cdot|z) = \mathcal{N}(\alpha_1 z, \beta_1^2 I_d) = \delta_z,$$

And sampling from p_t is just,

$$z \sim p_{\text{data}}, \epsilon \sim p_{\text{init}} = \mathcal{N}(0, I_d) \Rightarrow x = \alpha_t z + \beta_t \epsilon \sim p_t \quad \blacktriangleright \text{ sampling from marginal Gaussian path}$$

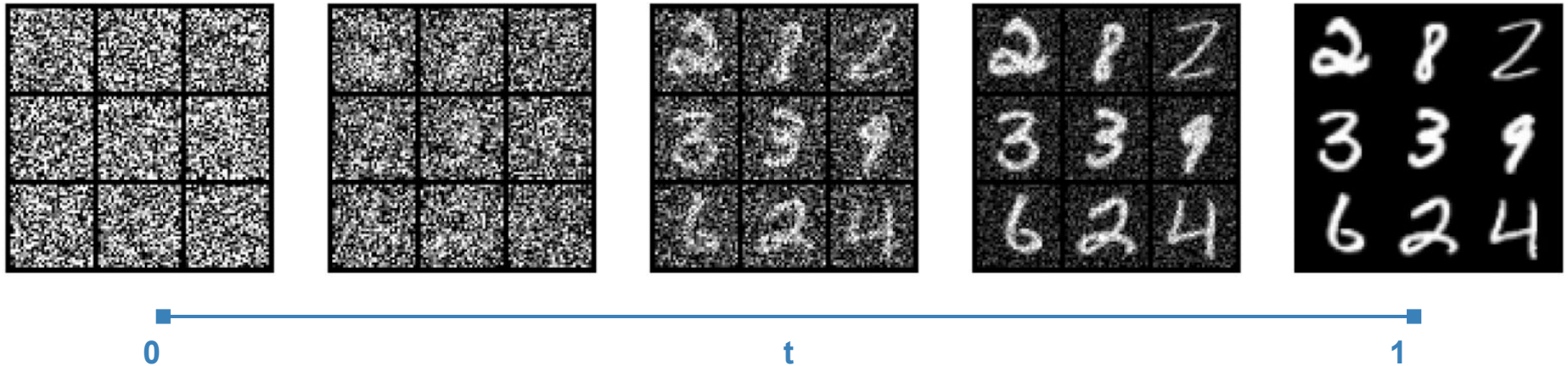
How do we train flow models? The objective

What do *Gaussian Probability paths* look like?



How do we train flow models? The objective

What do *Gaussian Probability paths* look like?



How do we train flow models? The objective

Putting the pieces together: Constructing u_t^{target}

How do we train flow models? The objective

Putting the pieces together: Constructing u_t^{target}

For every data point $z \in \mathbb{R}^d$, let $u_t^{\text{target}}(\cdot|z)$ denote a **conditional vector field**, defined so that the corresponding ODE yields the conditional probability path $p_t(\cdot|z)$, viz.,

$$X_0 \sim p_{\text{init}}, \quad \frac{d}{dt}X_t = u_t^{\text{target}}(X_t|z) \quad \Rightarrow \quad X_t \sim p_t(\cdot|z) \quad (0 \leq t \leq 1). \quad (18)$$

How do we train flow models? The objective

Putting the pieces together: Constructing u_t^{target}

For every data point $z \in \mathbb{R}^d$, let $u_t^{\text{target}}(\cdot|z)$ denote a **conditional vector field**, defined so that the corresponding ODE yields the conditional probability path $p_t(\cdot|z)$, viz.,

$$X_0 \sim p_{\text{init}}, \quad \frac{d}{dt}X_t = u_t^{\text{target}}(X_t|z) \quad \Rightarrow \quad X_t \sim p_t(\cdot|z) \quad (0 \leq t \leq 1). \quad (18)$$

Then the **marginal vector field** $u_t^{\text{target}}(x)$, defined by

$$u_t^{\text{target}}(x) = \int u_t^{\text{target}}(x|z) \frac{p_t(x|z)p_{\text{data}}(z)}{p_t(x)} dz, \quad (19)$$

How do we train flow models? The objective

Putting the pieces together: Constructing u_t^{target}

For every data point $z \in \mathbb{R}^d$, let $u_t^{\text{target}}(\cdot|z)$ denote a **conditional vector field**, defined so that the corresponding ODE yields the conditional probability path $p_t(\cdot|z)$, viz.,

$$X_0 \sim p_{\text{init}}, \quad \frac{d}{dt}X_t = u_t^{\text{target}}(X_t|z) \quad \Rightarrow \quad X_t \sim p_t(\cdot|z) \quad (0 \leq t \leq 1). \quad (18)$$

Then the **marginal vector field** $u_t^{\text{target}}(x)$, defined by

$$u_t^{\text{target}}(x) = \int u_t^{\text{target}}(x|z) \frac{p_t(x|z)p_{\text{data}}(z)}{p_t(x)} dz, \quad (19)$$

follows the marginal probability path, i.e.

$$X_0 \sim p_{\text{init}}, \quad \frac{d}{dt}X_t = u_t^{\text{target}}(X_t) \quad \Rightarrow \quad X_t \sim p_t \quad (0 \leq t \leq 1). \quad (20)$$

In particular, $X_1 \sim p_{\text{data}}$ for this ODE, so that we might say " u_t^{target} converts noise p_{init} into data p_{data} ".

How do we train flow models? The objective

Recall that the overall training objective is, $\mathcal{L}(\theta) = \|u_t^\theta(x) - \underbrace{u_t^{\text{target}}(x)}_{\text{training target}}\|^2$.

An intuitive way to obtain $u_t^\theta \sim u_t^{\text{target}}$ is by minimizing the expectation over x :

$$\mathcal{L}_{\text{FM}}(\theta) = \mathbb{E}_{t \sim \text{Unif}, x \sim p_t} [\|u_t^\theta(x) - u_t^{\text{target}}(x)\|^2]$$

Sample t uniformly on $[0, 1]$

Recall, $z \sim p_{\text{data}}, x \sim p_t(\cdot|z) \Rightarrow x \sim p_t$

However, computing u_t^{target} isn't tractable :(... What do we do?

Leverage the fact that $u_t^{\text{target}}(x|z)$ is! So instead, minimize:

$$\mathcal{L}_{\text{CFM}}(\theta) = \mathbb{E}_{\substack{z \sim p_{\text{data}}, x \sim p_t(\cdot|z) \\ t \sim \text{Unif}}} [\|u_t^\theta(x) - u_t^{\text{target}}(x|z)\|^2]$$

Guarantee: $\mathcal{L}_{\text{FM}}(\theta) = \mathcal{L}_{\text{CFM}}(\theta) + C$, where C is an independent constant!

How do we train flow models? The objective

We aim to minimize, $\mathcal{L}_{\text{CFM}}(\theta) = \mathbb{E}_{\substack{z \sim p_{\text{data}}, x \sim p_t(\cdot|z) \\ t \sim \text{Unif}}} [\|u_t^\theta(x) - u_t^{\text{target}}(x|z)\|^2]$

What exactly is $u_t^{\text{target}}(x|z)$? It's just

Recall VAEs

VAEs define intractable density function with latent \mathbf{z} :

$$p_{\theta}(x) = \int p_{\theta}(z)p_{\theta}(x|z)dz$$

Cannot optimize directly, derive and optimize lower bound on likelihood instead:

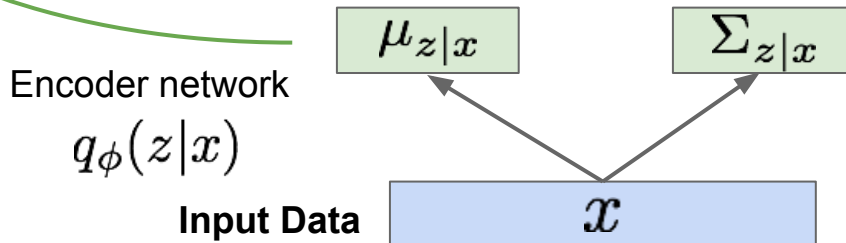
First loss for the encoder

$$\underbrace{\mathbf{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

$$D_{KL}(\mathcal{N}(\mu_{z|x}, \Sigma_{z|x}) || \mathcal{N}(0, I))$$

This equation has an analytical solution

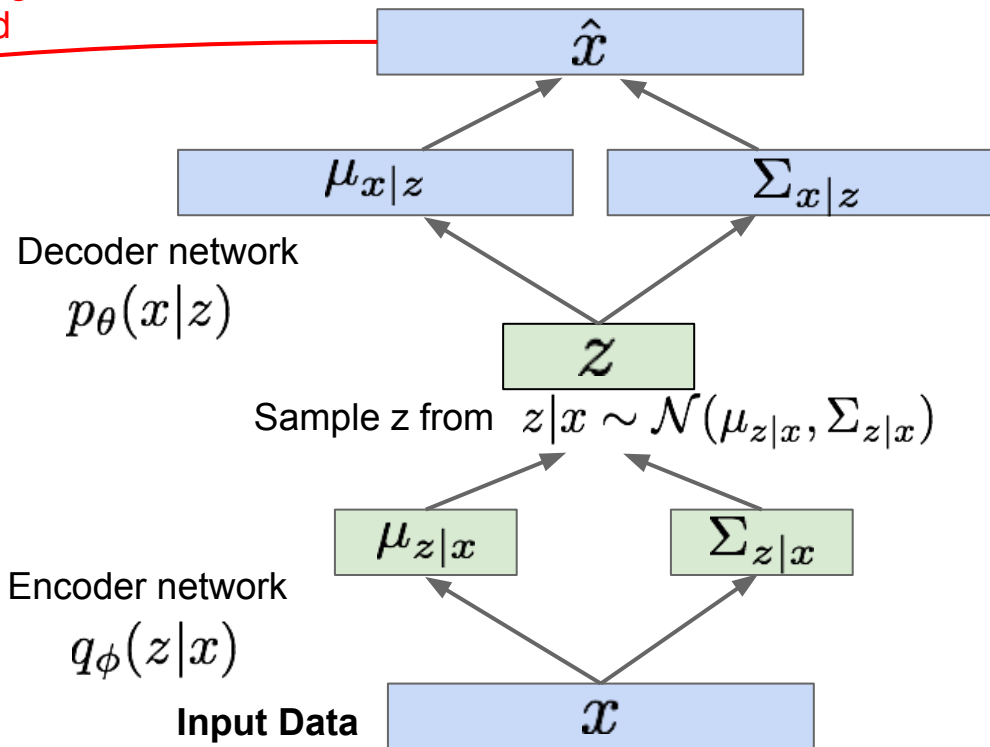
Make the latent variable distribution as similar to a unit normal distribution



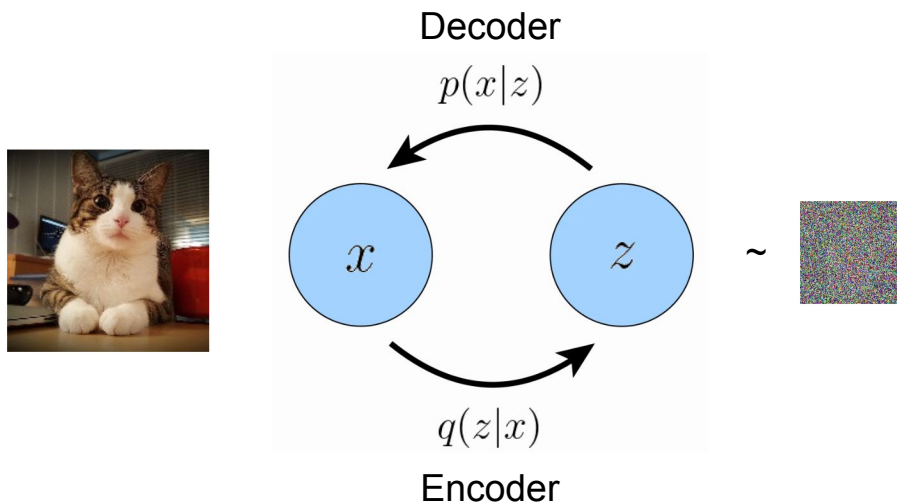
Second loss for both decoder and encoder

Maximize likelihood of original input being reconstructed

$$\underbrace{\mathbb{E}_z \left[\log p_\theta(x^{(i)} | z) \right]}_{\mathcal{L}(x^{(i)}, \theta, \phi)} = D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))$$

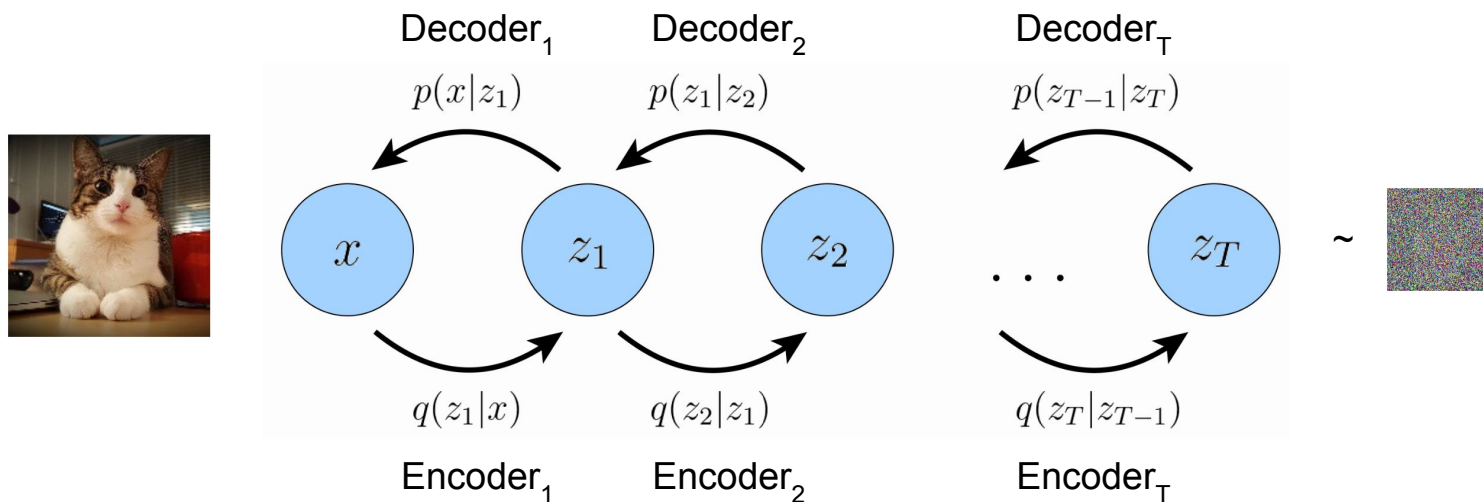


VAEs for images look like this



- We learn **2 networks**, one to encode and one to decode
- We ensure that **z** is similar to a unit normal noise
- To sample new images, we can sample from the unit normal and **decode in 1 step**

Markovian Hierarchical VAEs



- We learn **2T networks**, one to encode and one to decode
- We ensure that z_T is similar to a unit normal noise
- To sample new images, we can sample from the unit normal and **decode in T step**

Markovian Hierarchical VAEs - same derivation

$$\begin{aligned}\log p(x) &= \mathbb{E}_{z_{1:T} \sim q_\phi(z_{1:T}|x)} [\log p_\theta(x^{(i)})] \\ &= \mathbb{E}_{z_{1:T}} \left[\log \frac{p_\theta(x|z_{1:T})p_\theta(z_{1:T})}{p_\theta(z_{1:T}|x)} \right]\end{aligned}$$

$p_\theta(x)$ is independent of $z_{1:T}$

Bayes rule

Markovian Hierarchical VAEs - same derivation

$$\begin{aligned}\log p(x) &= \mathbb{E}_{z_{1:T} \sim q_\phi(z_{1:T}|x)} [\log p_\theta(x^{(i)})] && p_\theta(x) \text{ is independent of } z_{1:T} \\ &= \mathbb{E}_{z_{1:T}} \left[\log \frac{p_\theta(x|z_{1:T}) p_\theta(z_{1:T})}{p_\theta(z_{1:T}|x)} \right] && \text{Bayes rule} \\ &= \mathbb{E}_{z_{1:T}} \left[\log \frac{p_\theta(x|z_{1:T}) p_\theta(z_{1:T})}{p_\theta(z_{1:T}|x)} \frac{q_\phi(z_{1:T}|x)}{q_\phi(z_{1:T}|x)} \right] && \text{Multiplying by a constant} \\ &= \mathbb{E}_{z_{1:T}} [\log p_\theta(x|z_{1:T})] - \mathbb{E}_{z_{1:T}} \left[\log \frac{q_\phi(z_{1:T}|x)}{p_\theta(z_{1:T})} \right] + \mathbb{E}_{z_{1:T}} \left[\log \frac{q_\phi(z_{1:T})}{p(z_{1:T}|x)} \right]\end{aligned}$$

Markovian Hierarchical VAEs - same derivation

$$\begin{aligned}\log p(x) &= \mathbb{E}_{z_{1:T} \sim q_\phi(z_{1:T}|x)} [\log p_\theta(x^{(i)})] && p_\theta(x) \text{ is independent of } z_{1:T} \\ &= \mathbb{E}_{z_{1:T}} \left[\log \frac{p_\theta(x|z_{1:T}) p_\theta(z_{1:T})}{p_\theta(z_{1:T}|x)} \right] && \text{Bayes rule} \\ &= \mathbb{E}_{z_{1:T}} \left[\log \frac{p_\theta(x|z_{1:T}) p_\theta(z_{1:T})}{p_\theta(z_{1:T}|x)} \frac{q_\phi(z_{1:T}|x)}{q_\phi(z_{1:T}|x)} \right] && \text{Multiplying by a constant} \\ &= \mathbb{E}_{z_{1:T}} [\log p_\theta(x|z_{1:T})] - \mathbb{E}_{z_{1:T}} \left[\log \frac{q_\phi(z_{1:T}|x)}{p_\theta(z_{1:T})} \right] + \mathbb{E}_{z_{1:T}} \left[\log \frac{q_\phi(z_{1:T})}{p_\theta(z_{1:T}|x)} \right]\end{aligned}$$

↑
Reconstruction objective maximizes the likelihood of data $p_\theta(x|z)$

↑
This KL term (between Gaussians for encoder and z prior)

↑
 ~~$p_\theta(z|x)$~~ intractable but we know KL divergence always ≥ 0 .

Markovian Hierarchical VAEs

Keeping just the first two terms:

$$\log p(\mathbf{x}) \geq \mathbb{E}_{z_{1:T}}[\log p_{\theta}(\mathbf{x}|z_{1:T})] - \mathbb{E}_{z_{1:T}}\left[\log \frac{q_{\phi}(z_{1:T}|\mathbf{x})}{p_{\theta}(z_{1:T})}\right]$$

Markovian Hierarchical VAEs

Keeping just the first two terms:

$$\begin{aligned}\log p(x) &\geq \mathbb{E}_{z_{1:T}}[\log p_{\theta}(x|z_{1:T})] - \mathbb{E}_{z_{1:T}}\left[\log \frac{q_{\phi}(z_{1:T}|x)}{p_{\theta}(z_{1:T})}\right] \\ &= \mathbb{E}_{z_{1:T}}\left[\log \frac{p_{\theta}(x|z_{1:T})p_{\theta}(z_{1:T})}{q_{\phi}(z_{1:T}|x)}\right] \\ &= \mathbb{E}_{z_{1:T}}\left[\log \frac{p_{\theta}(x, z_{1:T})}{q_{\phi}(z_{1:T}|x)}\right]\end{aligned}$$

Markovian Hierarchical VAEs

Keeping just the first two terms:

$$\begin{aligned}\log p(\mathbf{x}) &\geq \mathbb{E}_{z_{1:T}}[\log p_{\theta}(\mathbf{x}|z_{1:T})] - \mathbb{E}_{z_{1:T}}\left[\log \frac{q_{\phi}(z_{1:T}|\mathbf{x})}{p_{\theta}(z_{1:T})}\right] \\ &= \mathbb{E}_{z_{1:T}}\left[\log \frac{p_{\theta}(\mathbf{x}|z_{1:T})p_{\theta}(z_{1:T})}{q_{\phi}(z_{1:T}|\mathbf{x})}\right] \\ &= \mathbb{E}_{z_{1:T}}\left[\log \frac{p_{\theta}(\mathbf{x}, z_{1:T})}{q_{\phi}(z_{1:T}|\mathbf{x})}\right]\end{aligned}$$

where the joint probability distribution is: $p(\mathbf{x}, z_{1:T}) = p(z_T)p_{\theta}(\mathbf{x} | z_1) \prod_{t=2}^T p_{\theta}(z_{t-1} | z_t)$

This is very similar to the autoregressive model formula

Markovian Hierarchical VAEs

Keeping just the first two terms:

$$\begin{aligned}\log p(\mathbf{x}) &\geq \mathbb{E}_{z_{1:T}}[\log p_{\theta}(\mathbf{x}|z_{1:T})] - \mathbb{E}_{z_{1:T}}\left[\log \frac{q_{\phi}(z_{1:T}|\mathbf{x})}{p_{\theta}(z_{1:T})}\right] \\ &= \mathbb{E}_{z_{1:T}}\left[\log \frac{p_{\theta}(\mathbf{x}|z_{1:T})p_{\theta}(z_{1:T})}{q_{\phi}(z_{1:T}|\mathbf{x})}\right] \\ &= \mathbb{E}_{z_{1:T}}\left[\log \frac{p_{\theta}(\mathbf{x}, z_{1:T})}{q_{\phi}(z_{1:T}|\mathbf{x})}\right]\end{aligned}$$

where the joint probability distribution is: $p(\mathbf{x}, z_{1:T}) = p(z_T)p_{\theta}(\mathbf{x} | z_1) \prod_{t=2}^T p_{\theta}(z_{t-1} | z_t)$

And the encoder posterior is: $q_{\phi}(z_{1:T} | \mathbf{x}) = q_{\phi}(z_1 | \mathbf{x}) \prod_{t=2}^T q_{\phi}(z_t | z_{t-1})$

Markovian Hierarchical VAEs

Keeping just the first two terms:

$$\begin{aligned}\log p(\mathbf{x}) &\geq \mathbb{E}_{z_{1:T}}[\log p_{\theta}(\mathbf{x}|z_{1:T})] - \mathbb{E}_{z_{1:T}}\left[\log \frac{q_{\phi}(z_{1:T}|\mathbf{x})}{p_{\theta}(z_{1:T})}\right] \\ &= \mathbb{E}_{z_{1:T}}\left[\log \frac{p_{\theta}(\mathbf{x}|z_{1:T})p_{\theta}(z_{1:T})}{q_{\phi}(z_{1:T}|\mathbf{x})}\right] \\ &= \mathbb{E}_{z_{1:T}}\left[\log \frac{p_{\theta}(\mathbf{x}, z_{1:T})}{q_{\phi}(z_{1:T}|\mathbf{x})}\right]\end{aligned}$$

Why is this a hard objective to train?

where the joint probability distribution is: $p(\mathbf{x}, \mathbf{z}_{1:T}) = p(\mathbf{z}_T)p_{\theta}(\mathbf{x} | \mathbf{z}_1) \prod_{t=2}^T p_{\theta}(\mathbf{z}_{t-1} | \mathbf{z}_t)$

And the encoder posterior is: $q_{\phi}(\mathbf{z}_{1:T} | \mathbf{x}) = q_{\phi}(\mathbf{z}_1 | \mathbf{x}) \prod_{t=2}^T q_{\phi}(\mathbf{z}_t | \mathbf{z}_{t-1})$

Markovian Hierarchical VAEs

Keeping just the first two terms:

$$\begin{aligned}\log p(\mathbf{x}) &\geq \mathbb{E}_{z_{1:T}}[\log p_{\theta}(\mathbf{x}|z_{1:T})] - \mathbb{E}_{z_{1:T}}\left[\log \frac{q_{\phi}(z_{1:T}|\mathbf{x})}{p_{\theta}(z_{1:T})}\right] \\ &= \mathbb{E}_{z_{1:T}}\left[\log \frac{p_{\theta}(\mathbf{x}|z_{1:T})p_{\theta}(z_{1:T})}{q_{\phi}(z_{1:T}|\mathbf{x})}\right] \\ &= \mathbb{E}_{z_{1:T}}\left[\log \frac{p_{\theta}(\mathbf{x}, z_{1:T})}{q_{\phi}(z_{1:T}|\mathbf{x})}\right]\end{aligned}$$

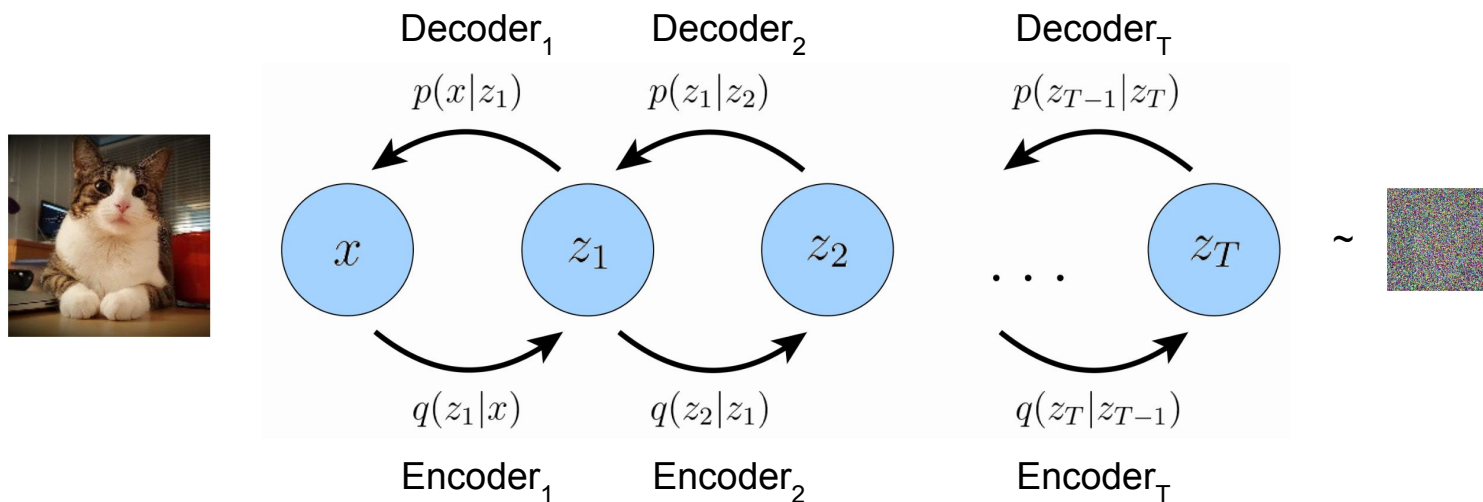
Why is this a hard objective to train?

1. There are too many networks to learn
2. The objective function is expensive!
3. It collapses easily!

where the joint probability distribution is: $p(\mathbf{x}, z_{1:T}) = p(z_T)p_{\theta}(\mathbf{x} | z_1) \prod_{t=2}^T p_{\theta}(z_{t-1} | z_t)$

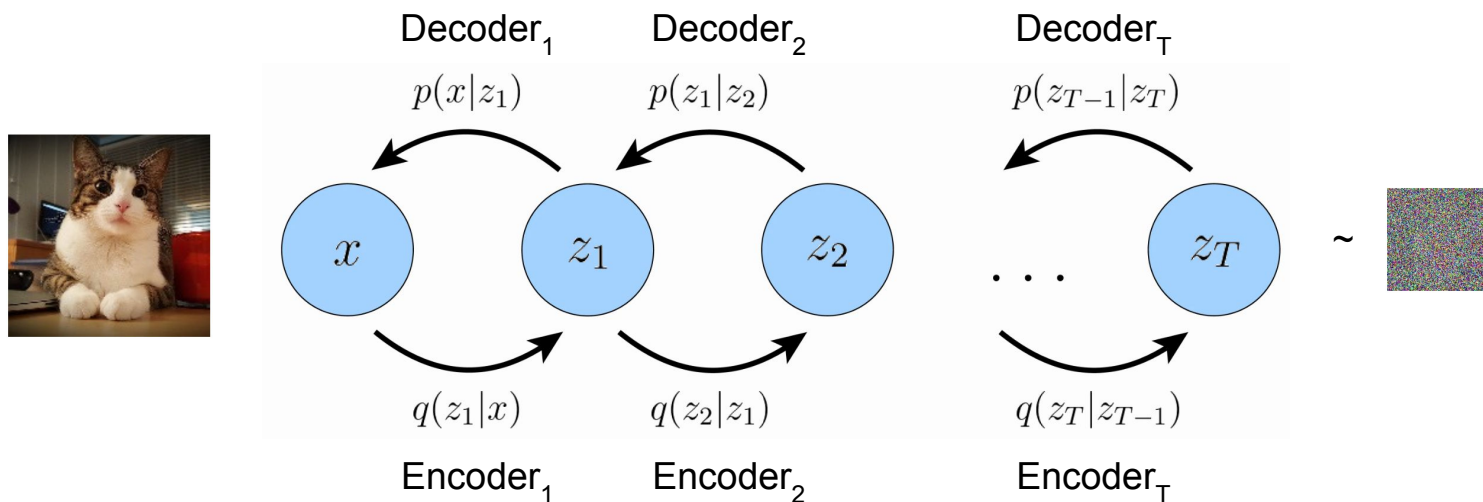
And the encoder posterior is: $q_{\phi}(z_{1:T} | \mathbf{x}) = q_{\phi}(z_1 | \mathbf{x}) \prod_{t=2}^T q_{\phi}(z_t | z_{t-1})$

Markovian Hierarchical VAEs



- Why is this a hard objective to train?
- 1. There are too many networks to learn.
- 2. The objective function is expensive!
- 3. It collapses easily!

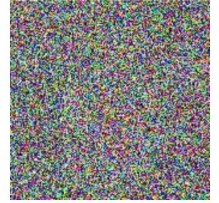
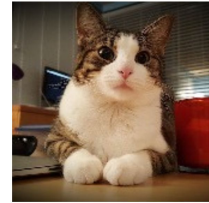
Markovian Hierarchical VAEs



Diffusion models are a special case
With a more interpretable, simpler objective.

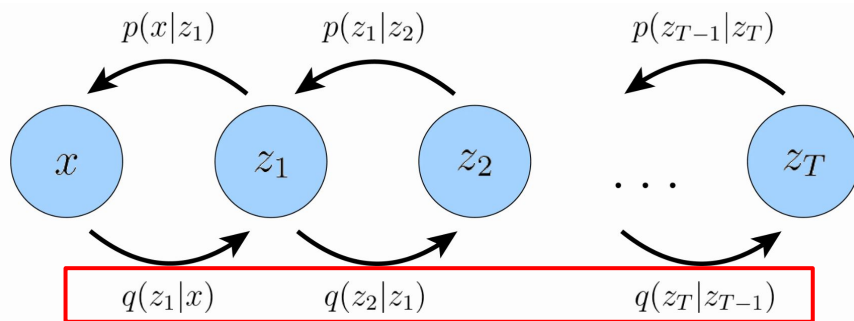
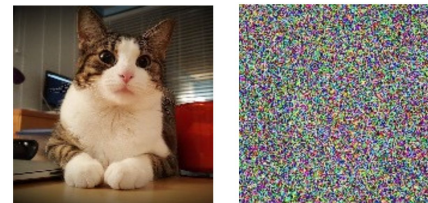
How are diffusion models different?

1. The latent dimension size is exactly equal to the data dimension



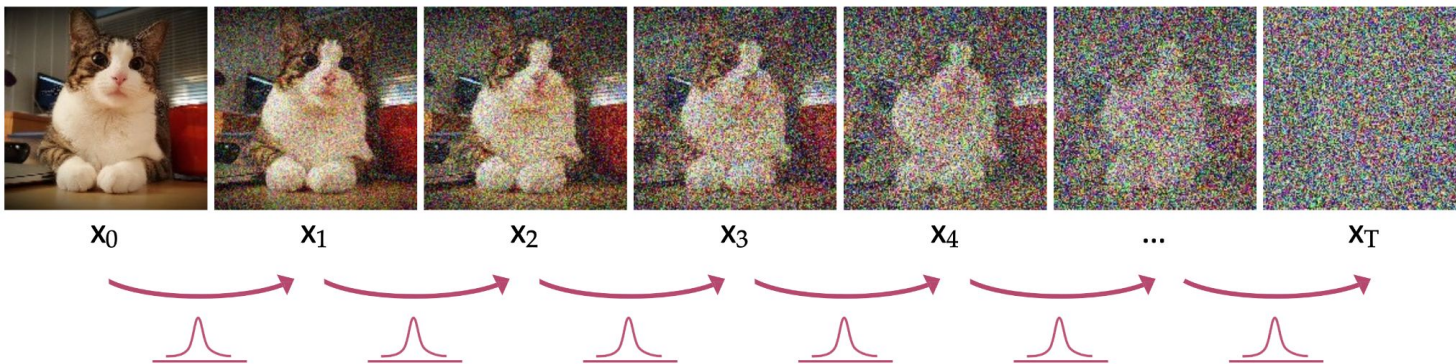
How are diffusion models different?

1. The latent dimension size is exactly equal to the data dimension
2. The encoders are **pre-defined** and **not learned**.



How are diffusion models different?

3. Encoders are designed as a **linear Gaussian model** conditioned on the time step: Add noise at every time step



$$\mathbf{x}_t = \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^*$$

How the forward step was designed:

So, given x_{t-1} we can sample x_t using:

$$x_t \sim \sqrt{\alpha_t}x_{t-1} + (1 - \alpha_t)\epsilon$$

where $\epsilon \sim \mathcal{N}(x; 0, I)$

Why was it designed like this?

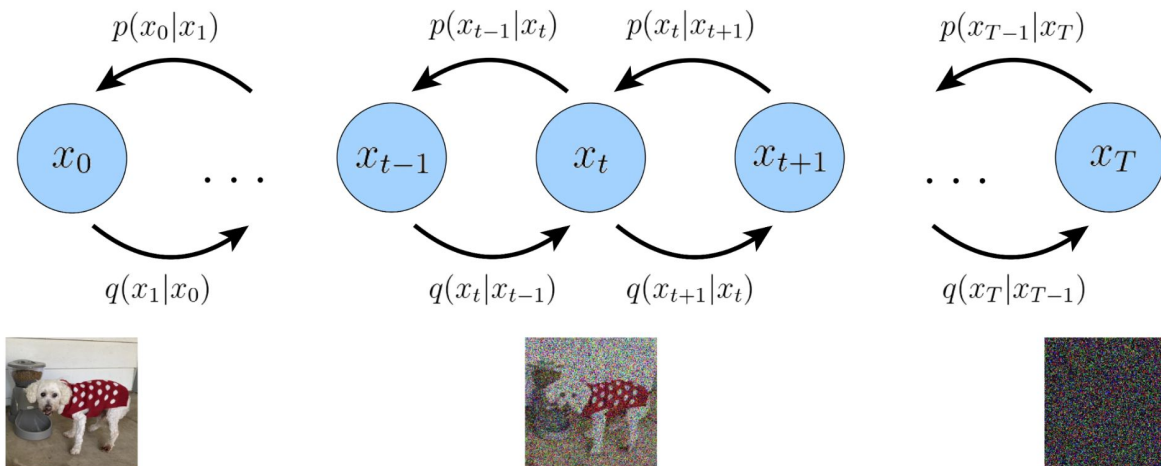
$$\begin{aligned}\mathbf{x}_t &= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_0 \\ &\sim \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})\end{aligned}$$

\mathbf{x}_t is now a Gaussian characterized by \mathbf{x}_0

where
$$\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$$

How are diffusion models different?

4. The Gaussian parameters vary over time in such a way that the distribution of the **latent at final step T is a standard Gaussian**



Terminology: **Forward** and **backward** process

Forward diffusion process (fixed)

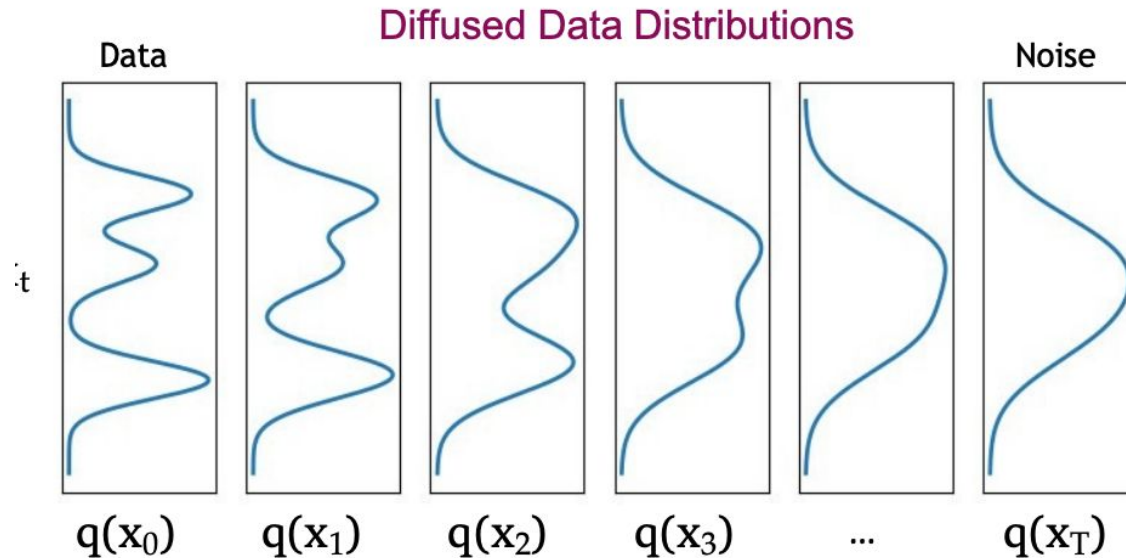


Reverse denoising process (generative)

Note: reverse or backward here doesn't mean the same thing as backpropagation

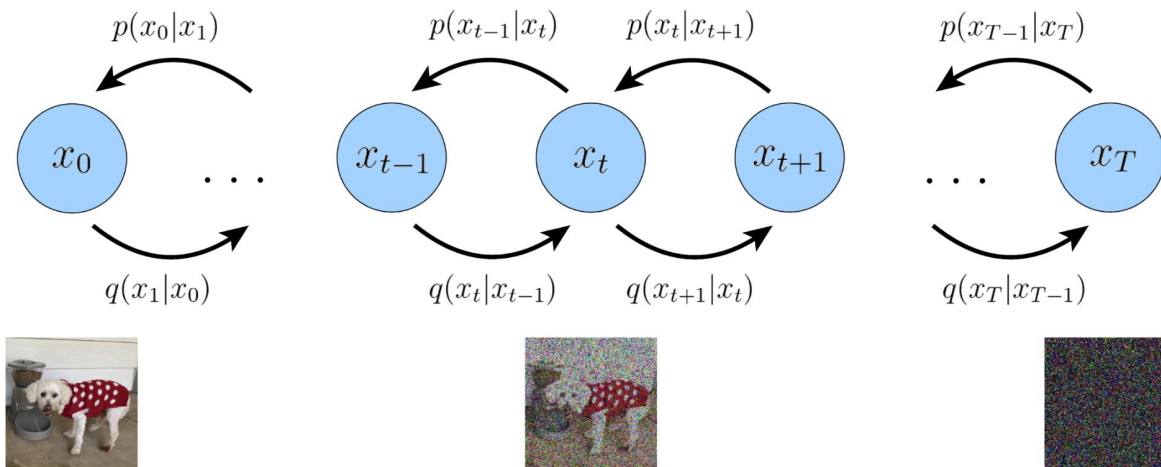
The distribution perspective

Over time, as we add more noise sampled from a Gaussian distribution, it begins to look more like a unit normal



How do we define a loss objective?

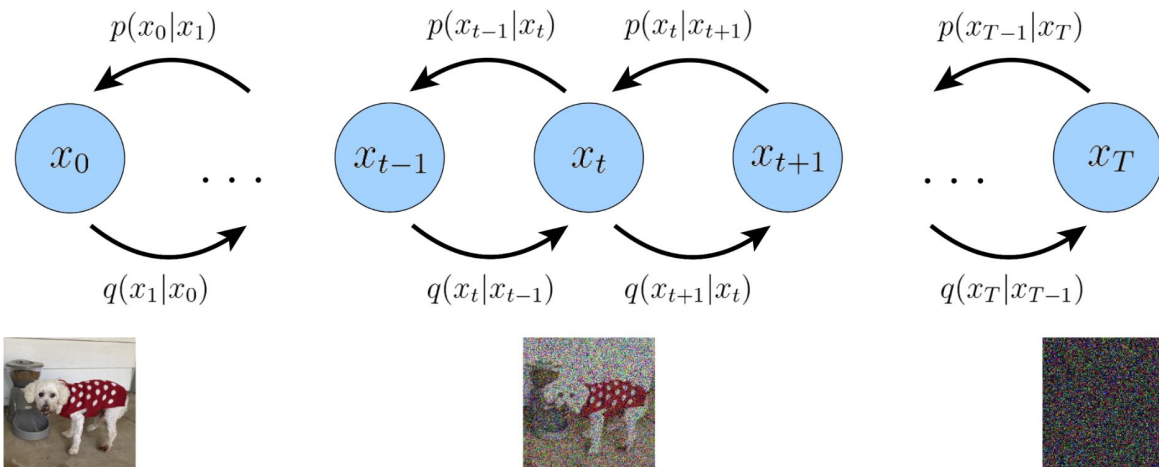
Q. What do we have to learn to generate new samples from noise?



How do we define a loss objective?

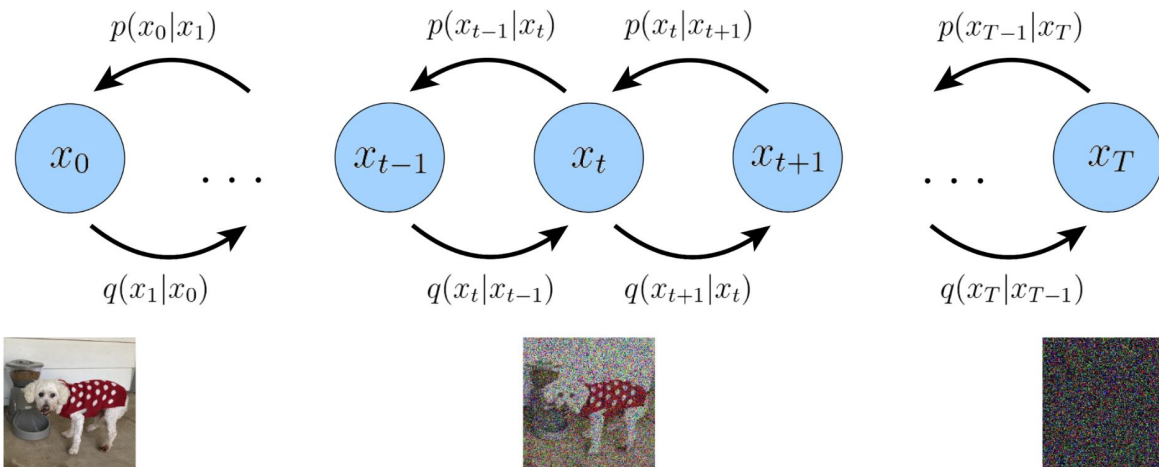
Q. What do we have to learn to generate new samples from noise?

A. We want to define a neural network to predict $p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t)$



How do we define a loss objective?

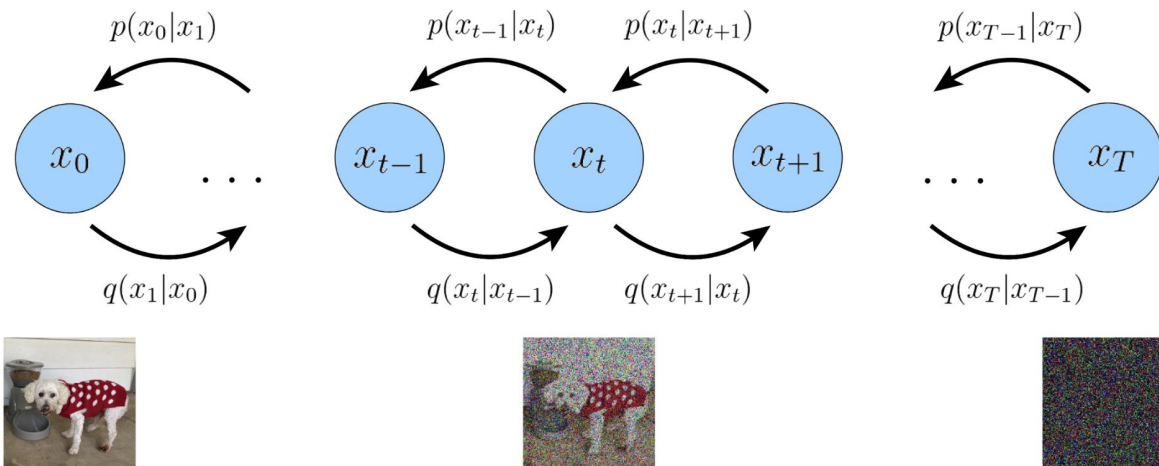
Q. How should we train $p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)$?



How do we define a loss objective?

Q. How should we train $p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)$?

A. We can get it to match $q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$!



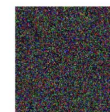
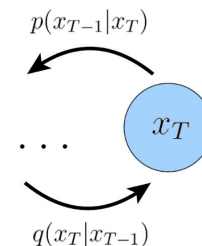
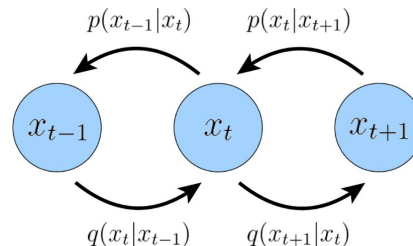
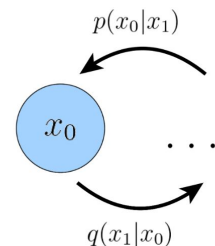
How do we define a loss objective?

Q. How should we train $p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)$?

A. We can get it to match $q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$!

Q. But why $q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$ and not:

$$\begin{aligned} & q(\mathbf{x}_{t-1}) \\ & q(\mathbf{x}_{t-1} | \mathbf{x}_0) \\ & q(\mathbf{x}_{t-1} | \mathbf{x}_t) \end{aligned}$$



Ok so our loss function is:

Q. How should we train $p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t)$?

A. We can get it to match $q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0)$!

Minimize the distance between the two distributions:

$$\arg \min_{\theta} \mathcal{D}_{\text{KL}}(q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t))$$

Ok so our loss function is:

Q. How should we train $p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t)$?

A. We can get it to match $q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0)$!

Minimize the distance between the two distributions:

$$\arg \min_{\theta} \mathcal{D}_{\text{KL}}(q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t))$$

Problem: How do we estimate $q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0)$?

The forward diffusion step

The distribution at step t is a Gaussian

$$q(\mathbf{x}_t \mid \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}_{t-1}, (1 - \alpha_t)\mathbf{I})$$

The forward diffusion step

The distribution at step t is a Gaussian

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}_{t-1}, (1 - \alpha_t)\mathbf{I})$$

The **mean** defined by \mathbf{x}_{t-1} : $\boldsymbol{\mu}_t(\mathbf{x}_t) = \sqrt{\alpha_t}\mathbf{x}_{t-1}$

α_t is a **predefined** value for each step t

The forward diffusion step

The distribution at step t is a Gaussian

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}_{t-1}, (1 - \alpha_t)\mathbf{I})$$

The **mean** defined by \mathbf{x}_{t-1} : $\boldsymbol{\mu}_t(\mathbf{x}_t) = \sqrt{\alpha_t}\mathbf{x}_{t-1}$

α_t is a **predefined** value for each step t

The **covariance** is **independent** of \mathbf{x}_{t-1} (an assumption)

$$\boldsymbol{\Sigma}_t(\mathbf{x}_t) = (1 - \alpha_t)\mathbf{I}$$

How the forward step was designed:

The distribution at step t is a Gaussian

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}_{t-1}, (1 - \alpha_t)\mathbf{I})$$

$$\boldsymbol{\mu}_t(\mathbf{x}_t) = \sqrt{\alpha_t}\mathbf{x}_{t-1}$$

$$\boldsymbol{\Sigma}_t(\mathbf{x}_t) = (1 - \alpha_t)\mathbf{I}$$

So, given \mathbf{x}_{t-1} we can sample \mathbf{x}_t using:

$$\mathbf{x}_t \sim \sqrt{\alpha_t}\mathbf{x}_{t-1} + (1 - \alpha_t)\boldsymbol{\epsilon}$$

where $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{x}; 0, \mathbf{I})$

Why was it designed like this?

$$\mathbf{x}_t = \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^*$$

Why was it designed like this?

$$\begin{aligned}\mathbf{x}_t &= \sqrt{\alpha_t}\mathbf{x}_{t-1} + \sqrt{1 - \alpha_t}\boldsymbol{\epsilon}_{t-1}^* \\ &= \sqrt{\alpha_t} \left(\sqrt{\alpha_{t-1}}\mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}}\boldsymbol{\epsilon}_{t-2}^* \right) + \sqrt{1 - \alpha_t}\boldsymbol{\epsilon}_{t-1}^* \longleftarrow \text{Substituting } \mathbf{x}_{t-1}\end{aligned}$$

Why was it designed like this?

$$\begin{aligned}\mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \\ &= \sqrt{\alpha_t} \left(\sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* \right) + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \quad \longleftarrow \text{Substituting } \mathbf{x}_{t-1} \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \quad \longleftarrow \text{Opening the parentheses}\end{aligned}$$

Why was it designed like this?

$$\begin{aligned}\mathbf{x}_t &= \sqrt{\alpha_t}\mathbf{x}_{t-1} + \sqrt{1 - \alpha_t}\boldsymbol{\epsilon}_{t-1}^* \\ &= \sqrt{\alpha_t} \left(\sqrt{\alpha_{t-1}}\mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}}\boldsymbol{\epsilon}_{t-2}^* \right) + \sqrt{1 - \alpha_t}\boldsymbol{\epsilon}_{t-1}^* \quad \leftarrow \text{Substituting } \mathbf{x}_{t-1} \\ &= \sqrt{\alpha_t\alpha_{t-1}}\mathbf{x}_{t-2} + \boxed{\sqrt{\alpha_t - \alpha_t\alpha_{t-1}}\boldsymbol{\epsilon}_{t-2}^* + \sqrt{1 - \alpha_t}\boldsymbol{\epsilon}_{t-1}^*} \quad \leftarrow \text{Opening the parentheses}\end{aligned}$$

We can interpret this $\sqrt{1 - \alpha_t}\boldsymbol{\epsilon}_{t-1}^*$ as a sample from $\mathcal{N}(\mathbf{0}, (1 - \alpha_t)\mathbf{I})$

We can interpret this $\sqrt{\alpha_t - \alpha_t\alpha_{t-1}}\boldsymbol{\epsilon}_{t-2}^*$ as a sample from $\mathcal{N}(\mathbf{0}, (\alpha_t - \alpha_t\alpha_{t-1})\mathbf{I})$

Why was it designed like this?

$$\begin{aligned} \mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \\ &= \sqrt{\alpha_t} \left(\sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* \right) + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \quad \leftarrow \text{Substituting } \mathbf{x}_{t-1} \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \boxed{\sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^*} \quad \leftarrow \text{Opening the parentheses} \end{aligned}$$

We can interpret this $\sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^*$ as a sample from $\mathcal{N}(\mathbf{0}, (1 - \alpha_t)\mathbf{I})$

We can interpret this $\sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^*$ as a sample from $\mathcal{N}(\mathbf{0}, (\alpha_t - \alpha_t \alpha_{t-1})\mathbf{I})$

Notice that $\sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^*$ is the sum of two Gaussian samples

Why was it designed like this?

$$\begin{aligned} \mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \\ &= \sqrt{\alpha_t} \left(\sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* \right) + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \quad \leftarrow \text{Substituting } \mathbf{x}_{t-1} \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \boxed{\sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^*} \quad \leftarrow \text{Opening the parentheses} \end{aligned}$$

We can interpret this $\sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^*$ as a sample from $\mathcal{N}(\mathbf{0}, (1 - \alpha_t)\mathbf{I})$

We can interpret this $\sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^*$ as a sample from $\mathcal{N}(\mathbf{0}, (\alpha_t - \alpha_t \alpha_{t-1})\mathbf{I})$

Notice that $\sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^*$ is the sum of two Gaussian samples

Using the property: $\mathcal{N}(x; 0, \sigma_1 I) + \mathcal{N}(x; 0, \sigma_2 I) = \mathcal{N}(x; 0, \sqrt{\sigma_1^2 + \sigma_2^2} I)$

Why was it designed like this?

$$\begin{aligned} \mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \\ &= \sqrt{\alpha_t} \left(\sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* \right) + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \quad \leftarrow \text{Substituting } \mathbf{x}_{t-1} \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \boxed{\sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^*} \quad \leftarrow \text{Opening the parentheses} \end{aligned}$$

We can interpret this $\sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^*$ as a sample from $\mathcal{N}(\mathbf{0}, (1 - \alpha_t)\mathbf{I})$

We can interpret this $\sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^*$ as a sample from $\mathcal{N}(\mathbf{0}, (\alpha_t - \alpha_t \alpha_{t-1})\mathbf{I})$

Notice that $\sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^*$ is the sum of two Gaussian samples

Using the property: $\mathcal{N}(x; 0, \sigma_1 I) + \mathcal{N}(x; 0, \sigma_2 I) = \mathcal{N}(x; 0, \sqrt{\sigma_1^2 + \sigma_2^2} I)$

We can rewrite $\sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^*$ as $\sqrt{\sqrt{\alpha_t - \alpha_t \alpha_{t-1}}^2 + \sqrt{1 - \alpha_t}^2} \boldsymbol{\epsilon}_{t-2}$

Why was it designed like this?

$$\begin{aligned} \mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \\ &= \sqrt{\alpha_t} \left(\sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* \right) + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \quad \leftarrow \text{Substituting } \mathbf{x}_{t-1} \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \quad \leftarrow \text{Opening the parentheses} \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\sqrt{\alpha_t - \alpha_t \alpha_{t-1}}^2 + \sqrt{1 - \alpha_t}^2} \boldsymbol{\epsilon}_{t-2}^* \quad \leftarrow \text{Sum of two Gaussians} \end{aligned}$$

Why was it designed like this?

$$\begin{aligned} \mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \\ &= \sqrt{\alpha_t} \left(\sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* \right) + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \quad \leftarrow \text{Substituting } \mathbf{x}_{t-1} \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \quad \leftarrow \text{Opening the parentheses} \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\sqrt{\alpha_t - \alpha_t \alpha_{t-1}}^2 + \sqrt{1 - \alpha_t}^2} \boldsymbol{\epsilon}_{t-2}^* \quad \leftarrow \text{Sum of two Gaussians} \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\alpha_t - \alpha_t \alpha_{t-1} + 1 - \alpha_t} \boldsymbol{\epsilon}_{t-2}^* \quad \leftarrow \text{Squaring the terms} \end{aligned}$$

Why was it designed like this?

$$\begin{aligned} \mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \\ &= \sqrt{\alpha_t} \left(\sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* \right) + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \quad \leftarrow \text{Substituting } \mathbf{x}_{t-1} \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \quad \leftarrow \text{Opening the parentheses} \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\sqrt{\alpha_t - \alpha_t \alpha_{t-1}}^2 + \sqrt{1 - \alpha_t}^2} \boldsymbol{\epsilon}_{t-2}^* \quad \leftarrow \text{Sum of two Gaussians} \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\alpha_t - \alpha_t \alpha_{t-1} + 1 - \alpha_t} \boldsymbol{\epsilon}_{t-2}^* \quad \leftarrow \text{Squaring the terms} \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* \quad \leftarrow \text{Simplifying} \end{aligned}$$

Why was it designed like this?

$$\begin{aligned} \mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \\ &= \sqrt{\alpha_t} \left(\sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* \right) + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \quad \leftarrow \text{Substituting } \mathbf{x}_{t-1} \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \quad \leftarrow \text{Opening the parentheses} \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\sqrt{\alpha_t - \alpha_t \alpha_{t-1}}^2 + \sqrt{1 - \alpha_t}^2} \boldsymbol{\epsilon}_{t-2} \quad \leftarrow \text{Sum of two Gaussians} \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\alpha_t - \alpha_t \alpha_{t-1} + 1 - \alpha_t} \boldsymbol{\epsilon}_{t-2} \quad \leftarrow \text{Squaring the terms} \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2} \quad \leftarrow \text{Simplifying} \\ &= \dots \\ &= \sqrt{\prod_{i=1}^t \alpha_i} \mathbf{x}_0 + \sqrt{1 - \prod_{i=1}^t \alpha_i} \boldsymbol{\epsilon}_0 \quad \leftarrow \text{Substituting till } \mathbf{x}_0 \end{aligned}$$

Why was it designed like this?

$$\begin{aligned}
 \mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \\
 &= \sqrt{\alpha_t} \left(\sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* \right) + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \quad \leftarrow \text{Substituting } \mathbf{x}_{t-1} \\
 &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \quad \leftarrow \text{Opening the parentheses} \\
 &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\sqrt{\alpha_t - \alpha_t \alpha_{t-1}}^2 + \sqrt{1 - \alpha_t}^2} \boldsymbol{\epsilon}_{t-2} \quad \leftarrow \text{Sum of two Gaussians} \\
 &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\alpha_t - \alpha_t \alpha_{t-1} + 1 - \alpha_t} \boldsymbol{\epsilon}_{t-2} \quad \leftarrow \text{Squaring the terms} \\
 &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2} \quad \leftarrow \text{Simplifying} \\
 &= \dots \\
 &= \sqrt{\prod_{i=1}^t \alpha_i} \mathbf{x}_0 + \sqrt{1 - \prod_{i=1}^t \alpha_i} \boldsymbol{\epsilon}_0 \quad \leftarrow \text{Substituting till } \mathbf{x}_0 \\
 &= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_0 \quad \leftarrow \text{Let } \bar{\alpha}_t = \prod_{i=1}^t \alpha_i
 \end{aligned}$$

Why was it designed like this?

$$\begin{aligned}
 \mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \\
 &= \sqrt{\alpha_t} \left(\sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* \right) + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \quad \leftarrow \text{Substituting } \mathbf{x}_{t-1} \\
 &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \quad \leftarrow \text{Opening the parentheses} \\
 &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\sqrt{\alpha_t - \alpha_t \alpha_{t-1}}^2 + \sqrt{1 - \alpha_t}^2} \boldsymbol{\epsilon}_{t-2} \quad \leftarrow \text{Sum of two Gaussians} \\
 &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\alpha_t - \alpha_t \alpha_{t-1} + 1 - \alpha_t} \boldsymbol{\epsilon}_{t-2} \quad \leftarrow \text{Squaring the terms} \\
 &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2} \quad \leftarrow \text{Simplifying} \\
 &= \dots \\
 &= \sqrt{\prod_{i=1}^t \alpha_i} \mathbf{x}_0 + \sqrt{1 - \prod_{i=1}^t \alpha_i} \boldsymbol{\epsilon}_0 \quad \leftarrow \text{Substituting till } \mathbf{x}_0 \\
 &= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_0 \quad \leftarrow \text{Let } \bar{\alpha}_t = \prod_{i=1}^t \alpha_i \\
 &\sim \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}) \quad \leftarrow \mathbf{x}_t \text{ is now a Gaussian characterized by } \mathbf{x}_0
 \end{aligned}$$

Takeaway from the previous slides:

$$\mathbf{x}_t \sim \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$$

We can instantly sample \mathbf{x}_t given any input data \mathbf{x}_0

What about the reverse?

$$q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0)$$

What about the reverse?

$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \frac{q(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{x}_0)q(\mathbf{x}_{t-1} | \mathbf{x}_0)}{q(\mathbf{x}_t | \mathbf{x}_0)}$$

Applying Bayes rule

What about the reverse?

$$\begin{aligned} q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) &= \frac{q(\mathbf{x}_t \mid \mathbf{x}_{t-1}, \mathbf{x}_0)q(\mathbf{x}_{t-1} \mid \mathbf{x}_0)}{q(\mathbf{x}_t \mid \mathbf{x}_0)} \\ &= \frac{\mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}_{t-1}, (1 - \alpha_t)\mathbf{I})}{q(\mathbf{x}_t \mid \mathbf{x}_0)} \end{aligned}$$

The first term is just a single forward diffusion process:

$$q(\mathbf{x}_t \mid \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}_{t-1}, (1 - \alpha_t)\mathbf{I})$$

What about the reverse?

$$\begin{aligned} q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) &= \frac{q(\mathbf{x}_t \mid \mathbf{x}_{t-1}, \mathbf{x}_0) q(\mathbf{x}_{t-1} \mid \mathbf{x}_0)}{q(\mathbf{x}_t \mid \mathbf{x}_0)} \\ &= \frac{\mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t} \mathbf{x}_{t-1}, (1 - \alpha_t) \mathbf{I}) \mathcal{N}(\mathbf{x}_{t-1}; \sqrt{\bar{\alpha}_{t-1}} \mathbf{x}_0, (1 - \bar{\alpha}_{t-1}) \mathbf{I})}{\mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t} \mathbf{x}_0, (1 - \alpha_t) \mathbf{I})} \end{aligned}$$

The second term is also a Gaussian using the formula we just derived:

$$\mathbf{x}_t \sim \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$$

What about the reverse?

$$\begin{aligned} q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) &= \frac{q(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{x}_0)q(\mathbf{x}_{t-1} | \mathbf{x}_0)}{q(\mathbf{x}_t | \mathbf{x}_0)} \\ &= \frac{\mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}_{t-1}, (1 - \alpha_t)\mathbf{I})\mathcal{N}(\mathbf{x}_{t-1}; \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0, (1 - \bar{\alpha}_{t-1})\mathbf{I})}{\mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})} \end{aligned}$$

The third term is also a Gaussian using the same formula:

$$\mathbf{x}_t \sim \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})$$

What about the reverse?

$$\begin{aligned} q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) &= \frac{q(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{x}_0)q(\mathbf{x}_{t-1} | \mathbf{x}_0)}{q(\mathbf{x}_t | \mathbf{x}_0)} \\ &= \frac{\mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}_{t-1}, (1 - \alpha_t)\mathbf{I})\mathcal{N}(\mathbf{x}_{t-1}; \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0, (1 - \bar{\alpha}_{t-1})\mathbf{I})}{\mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})} \\ &\propto \mathcal{N}(\mathbf{x}_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t}}_{\mu_q(\mathbf{x}_t, \mathbf{x}_0)}, \underbrace{\frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}\mathbf{I}}_{\Sigma_q(t)}) \end{aligned}$$

The product of these 3 Gaussian distributions simplify to a Gaussian as well!

Let's call its mean $\mu_q(\mathbf{x}_t, \mathbf{x}_0)$ and variance $\Sigma_q(t)$

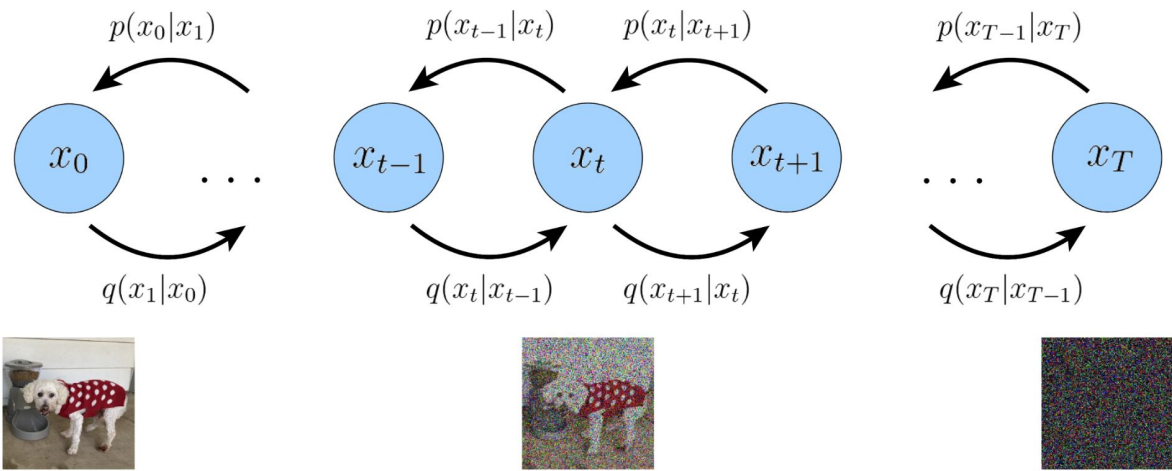
Proof (out of scope for the class)

$$\begin{aligned}
 q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) &= \frac{q(\mathbf{x}_t \mid \mathbf{x}_{t-1}, \mathbf{x}_0)q(\mathbf{x}_{t-1} \mid \mathbf{x}_0)}{q(\mathbf{x}_t \mid \mathbf{x}_0)} \\
 &= \frac{\mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}_{t-1}, (1 - \alpha_t)\mathbf{I})\mathcal{N}(\mathbf{x}_{t-1}; \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0, (1 - \bar{\alpha}_{t-1})\mathbf{I})}{\mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})} \\
 &\propto \exp \left\{ - \left[\frac{(\mathbf{x}_t - \sqrt{\alpha_t}\mathbf{x}_{t-1})^2}{2(1 - \alpha_t)} + \frac{(\mathbf{x}_{t-1} - \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0)^2}{2(1 - \bar{\alpha}_{t-1})} - \frac{(\mathbf{x}_t - \sqrt{\bar{\alpha}_t}\mathbf{x}_0)^2}{2(1 - \bar{\alpha}_t)} \right] \right\} \\
 &= \exp \left\{ - \frac{1}{2} \left[\frac{(\mathbf{x}_t - \sqrt{\alpha_t}\mathbf{x}_{t-1})^2}{1 - \alpha_t} + \frac{(\mathbf{x}_{t-1} - \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0)^2}{1 - \bar{\alpha}_{t-1}} - \frac{(\mathbf{x}_t - \sqrt{\bar{\alpha}_t}\mathbf{x}_0)^2}{1 - \bar{\alpha}_t} \right] \right\} \\
 &= \exp \left\{ - \frac{1}{2} \left[\frac{(-2\sqrt{\alpha_t}\mathbf{x}_t\mathbf{x}_{t-1} + \alpha_t\mathbf{x}_{t-1}^2)}{1 - \alpha_t} + \frac{(\mathbf{x}_{t-1}^2 - 2\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_{t-1}\mathbf{x}_0)}{1 - \bar{\alpha}_{t-1}} + C(\mathbf{x}_t, \mathbf{x}_0) \right] \right\} \\
 &\propto \exp \left\{ - \frac{1}{2} \left[- \frac{2\sqrt{\alpha_t}\mathbf{x}_t\mathbf{x}_{t-1}}{1 - \alpha_t} + \frac{\alpha_t\mathbf{x}_{t-1}^2}{1 - \alpha_t} + \frac{\mathbf{x}_{t-1}^2}{1 - \bar{\alpha}_{t-1}} - \frac{2\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_{t-1}\mathbf{x}_0}{1 - \bar{\alpha}_{t-1}} \right] \right\} \\
 &= \exp \left\{ - \frac{1}{2} \left[\left(\frac{\alpha_t}{1 - \alpha_t} + \frac{1}{1 - \bar{\alpha}_{t-1}} \right) \mathbf{x}_{t-1}^2 - 2 \left(\frac{\sqrt{\alpha_t}\mathbf{x}_t}{1 - \alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0}{1 - \bar{\alpha}_{t-1}} \right) \mathbf{x}_{t-1} \right] \right\} \\
 &= \exp \left\{ - \frac{1}{2} \left[\frac{\alpha_t(1 - \bar{\alpha}_{t-1}) + 1 - \alpha_t}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})} \mathbf{x}_{t-1}^2 - 2 \left(\frac{\sqrt{\alpha_t}\mathbf{x}_t}{1 - \alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0}{1 - \bar{\alpha}_{t-1}} \right) \mathbf{x}_{t-1} \right] \right\}
 \end{aligned}$$

Proof (out of scope for the class)

$$\begin{aligned}
 &= \exp \left\{ -\frac{1}{2} \left[\frac{\alpha_t - \bar{\alpha}_t + 1 - \alpha_t}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})} \mathbf{x}_{t-1}^2 - 2 \left(\frac{\sqrt{\alpha_t} \mathbf{x}_t}{1 - \alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}} \mathbf{x}_0}{1 - \bar{\alpha}_{t-1}} \right) \mathbf{x}_{t-1} \right] \right\} \\
 &= \exp \left\{ -\frac{1}{2} \left[\frac{1 - \bar{\alpha}_t}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})} \mathbf{x}_{t-1}^2 - 2 \left(\frac{\sqrt{\alpha_t} \mathbf{x}_t}{1 - \alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}} \mathbf{x}_0}{1 - \bar{\alpha}_{t-1}} \right) \mathbf{x}_{t-1} \right] \right\} \\
 &= \exp \left\{ -\frac{1}{2} \left(\frac{1 - \bar{\alpha}_t}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})} \right) \left[\mathbf{x}_{t-1}^2 - 2 \frac{\left(\frac{\sqrt{\alpha_t} \mathbf{x}_t}{1 - \alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}} \mathbf{x}_0}{1 - \bar{\alpha}_{t-1}} \right)}{\frac{1 - \bar{\alpha}_t}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}} \mathbf{x}_{t-1} \right] \right\} \\
 &= \exp \left\{ -\frac{1}{2} \left(\frac{1 - \bar{\alpha}_t}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})} \right) \left[\mathbf{x}_{t-1}^2 - 2 \frac{\left(\frac{\sqrt{\alpha_t} \mathbf{x}_t}{1 - \alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}} \mathbf{x}_0}{1 - \bar{\alpha}_{t-1}} \right) (1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{x}_{t-1} \right] \right\} \\
 &= \exp \left\{ -\frac{1}{2} \left(\frac{1}{\frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}} \right) \left[\mathbf{x}_{t-1}^2 - 2 \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1}) \mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t) \mathbf{x}_0}{1 - \bar{\alpha}_t} \mathbf{x}_{t-1} \right] \right\} \\
 &\propto \mathcal{N}(\mathbf{x}_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1}) \mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t) \mathbf{x}_0}{1 - \bar{\alpha}_t}}_{\mu_q(\mathbf{x}_t, \mathbf{x}_0)}, \underbrace{\frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{I}}_{\Sigma_q(t)})
 \end{aligned}$$

Let's go back to the Markovian VAE



We are ready to set up a simple intuitive loss function to train the decoder!

Given an image x_0 :

We want to generate $p_{\theta}(x_{t-1} | x_t)$ to match the Gaussian we just derived: $q(x_{t-1} | x_t, x_0)$

You can show mathematically that:

$q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0)$ Is also a Gaussian:

$$\mathcal{N}(\mathbf{x}_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t}}_{\mu_q(\mathbf{x}_t, \mathbf{x}_0)}, \underbrace{\frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}\mathbf{I}}_{\Sigma_q(t)})$$

You can show mathematically that:

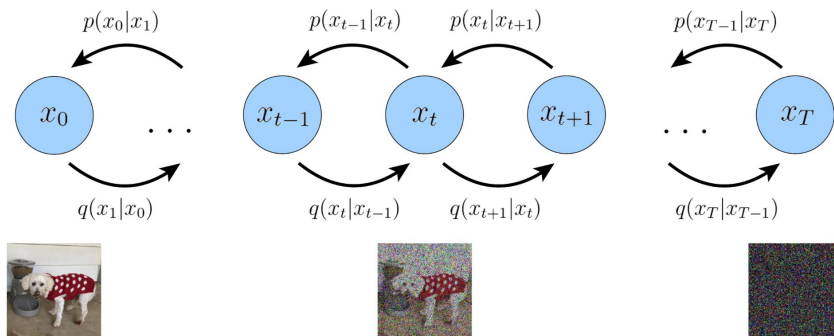
$q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0)$ Is also a Gaussian:

$$\mathcal{N}(\mathbf{x}_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t}}_{\mu_q(\mathbf{x}_t, \mathbf{x}_0)}, \underbrace{\frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}\mathbf{I}}_{\Sigma_q(t)})$$

And that:

$$\begin{aligned}\mu_q(\mathbf{x}_t, \mathbf{x}_0) &= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t} \\ &= \frac{1}{\sqrt{\alpha_t}}\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}\sqrt{\alpha_t}}\boldsymbol{\epsilon}_0\end{aligned}$$

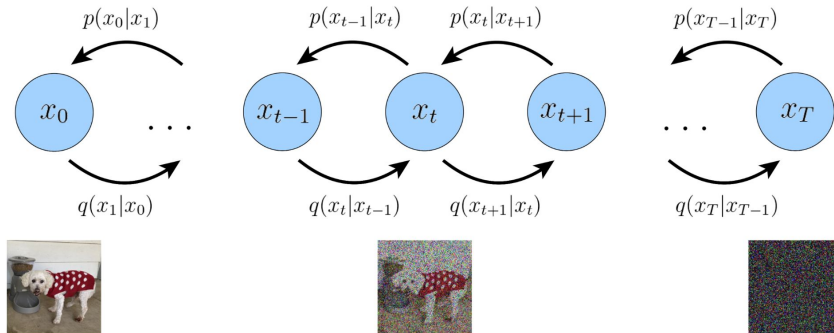
The loss function tries to match distributions



The loss function

$$\arg \min_{\theta} \mathcal{D}_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) || p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t))$$

We can model $p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t)$ as a Gaussian

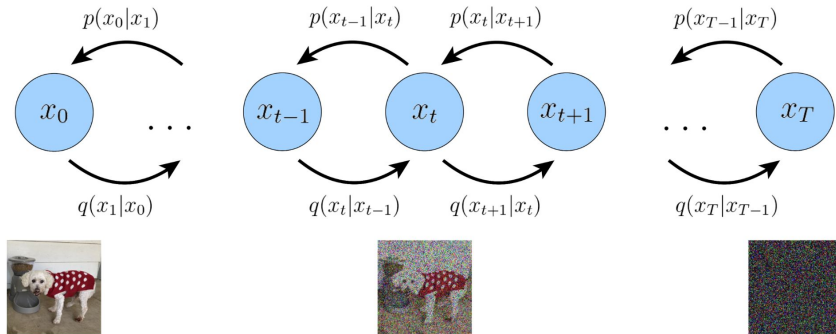


The loss function

$$\arg \min_{\theta} \mathcal{D}_{\text{KL}}(q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t))$$

$$= \arg \min_{\theta} \mathcal{D}_{\text{KL}}(\mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q(t)) \parallel \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\theta}, \boldsymbol{\Sigma}_q(t)))$$

We can model $p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t)$ as a Gaussian



The loss function

$$\begin{aligned} & \arg \min_{\theta} \mathcal{D}_{\text{KL}}(q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t)) \\ &= \arg \min_{\theta} \mathcal{D}_{\text{KL}}(\mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q(t)) \parallel \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\theta}, \boldsymbol{\Sigma}_q(t))) \\ &= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\|\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q\|_2^2 \right] \end{aligned}$$

Proof (out of scope for class)

$$\begin{aligned} & \arg \min_{\theta} \mathcal{D}_{\text{KL}}(q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t)) \\ &= \arg \min_{\theta} \mathcal{D}_{\text{KL}}(\mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q(t)) \parallel \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\theta}, \boldsymbol{\Sigma}_q(t))) \\ &= \arg \min_{\theta} \frac{1}{2} \left[\log \frac{|\boldsymbol{\Sigma}_q(t)|}{|\boldsymbol{\Sigma}_q(t)|} - d + \text{tr}(\boldsymbol{\Sigma}_q(t)^{-1} \boldsymbol{\Sigma}_q(t)) + (\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q)^T \boldsymbol{\Sigma}_q(t)^{-1} (\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q) \right] \\ &= \arg \min_{\theta} \frac{1}{2} \left[\log 1 - d + d + (\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q)^T \boldsymbol{\Sigma}_q(t)^{-1} (\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q) \right] \\ &= \arg \min_{\theta} \frac{1}{2} \left[(\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q)^T \boldsymbol{\Sigma}_q(t)^{-1} (\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q) \right] \\ &= \arg \min_{\theta} \frac{1}{2} \left[(\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q)^T (\sigma_q^2(t) \mathbf{I})^{-1} (\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q) \right] \\ &= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\|\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q\|_2^2 \right] \end{aligned}$$

Ok we are close to the objective:

The loss we want to minimize is $\arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\|\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q\|_2^2 \right]$

Ok we are close to the objective:

The loss we want to minimize is $\arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\|\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q\|_2^2 \right]$

From the previous slide, we got the mean from this:

$$\mathcal{N}(\mathbf{x}_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t}}_{\boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0)}, \underbrace{\frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{I}}_{\boldsymbol{\Sigma}_q(t)})$$

Ok we are close to the objective:

The loss we want to minimize is $\arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\|\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q\|_2^2 \right]$

From the previous slide, we got the mean from this:

$$\mathcal{N}(\mathbf{x}_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t}}_{\boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0)}, \underbrace{\frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{I}}_{\boldsymbol{\Sigma}_q(t)})$$

So, we can write the mean to be: $\boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t}$

Ok we are close to the objective:

The loss we want to minimize is $\arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\|\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q\|_2^2 \right]$

From the previous slide, we got the mean from this:

$$\mathcal{N}(\mathbf{x}_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t}}_{\boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0)}, \underbrace{\frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{I}}_{\boldsymbol{\Sigma}_q(t)})$$

So, we can write the mean to be:

$$\begin{aligned} \boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0) &= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t} \\ &= \frac{1}{\sqrt{\alpha_t}}\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}\sqrt{\alpha_t}}\boldsymbol{\epsilon}_0 \end{aligned}$$

Our neural network can predict **noise** instead!

$$\mu_q(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \epsilon_0$$

We can also set our predicted mean to be:

$$\mu_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \hat{\epsilon}_\theta(\mathbf{x}_t, t)$$

Why is this helpful?

Our neural network can predict **noise** instead!

$$\mu_q(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \epsilon_0$$

We can also set our predicted mean to be:

$$\mu_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \hat{\epsilon}_\theta(\mathbf{x}_t, t)$$

Why is this helpful? Because now our model needs to predict the noise that was injected, which turns out to be empirically more stable of an objective than predicting the image mean.

The two loss objectives are equivalent

The loss function

$$\arg \min_{\theta} \mathcal{D}_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) || p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t))$$

$$= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\|\mu_{\theta} - \mu_q\|_2^2 \right] \quad \text{Instead of predicting the mean image values}$$

The two loss objectives are equivalent

The loss function

$$\arg \min_{\theta} \mathcal{D}_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) || p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t))$$

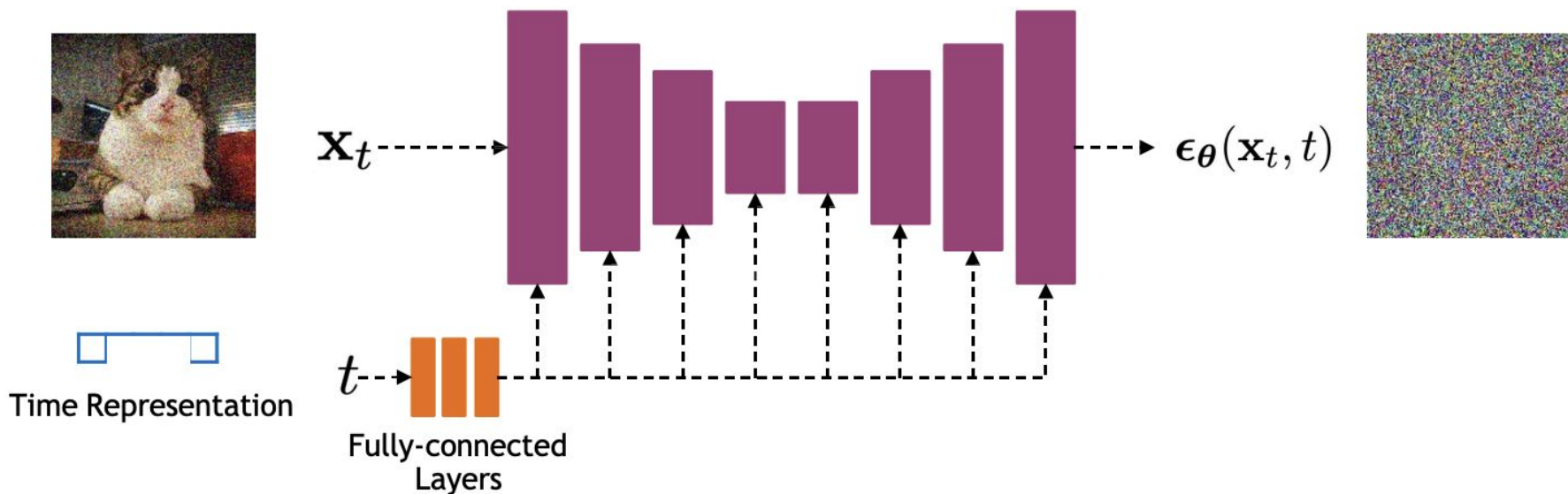
$$= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\|\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q\|_2^2 \right] \quad \text{Instead of predicting the mean image values}$$

$$= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \frac{(1 - \alpha_t)^2}{(1 - \bar{\alpha}_t)\alpha_t} \left[\|\boldsymbol{\epsilon}_0 - \hat{\boldsymbol{\epsilon}}_{\theta}(\mathbf{x}_t, t)\|_2^2 \right] \quad \text{The neural network can predict the added noise}$$

Proof: (out of scope)

$$\begin{aligned} & \arg \min_{\theta} \mathcal{D}_{\text{KL}}(q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t)) \\ &= \arg \min_{\theta} \mathcal{D}_{\text{KL}}(\mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q(t)) \parallel \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\theta}, \boldsymbol{\Sigma}_q(t))) \\ &= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\left\| \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \hat{\boldsymbol{\epsilon}}_{\theta}(\mathbf{x}_t, t) - \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t + \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \boldsymbol{\epsilon}_0 \right\|_2^2 \right] \\ &= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\left\| \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \boldsymbol{\epsilon}_0 - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \hat{\boldsymbol{\epsilon}}_{\theta}(\mathbf{x}_t, t) \right\|_2^2 \right] \\ &= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\left\| \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} (\boldsymbol{\epsilon}_0 - \hat{\boldsymbol{\epsilon}}_{\theta}(\mathbf{x}_t, t)) \right\|_2^2 \right] \\ &= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \frac{(1 - \alpha_t)^2}{(1 - \bar{\alpha}_t)\alpha_t} \left[\|\boldsymbol{\epsilon}_0 - \hat{\boldsymbol{\epsilon}}_{\theta}(\mathbf{x}_t, t)\|_2^2 \right] \end{aligned}$$

The denoising architecture



Time representation: sinusoidal positional embeddings.

How do we **sample a new image**?

Sample $x_T \sim \mathcal{N}(0, I)$

For $t = T \dots 1$ do

Predict $\hat{\epsilon}_t = p_\theta(x_t)$

$$\mu_{t-1} = \frac{1}{\sqrt{\alpha_t}} x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \sqrt{1 - \bar{\alpha}_{t-1}} \hat{\epsilon}_t$$

Sample $x_{t-1} \sim \mathcal{N}(\mu_{t-1}, \Sigma_{t-1})$

Return x_0



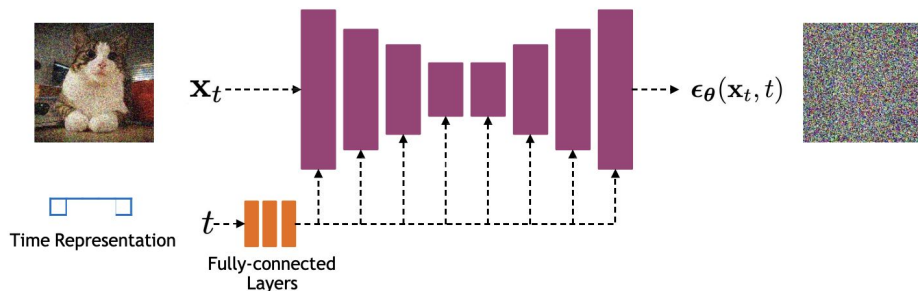
Reverse denoising process (generative)

How is the time step inputted:

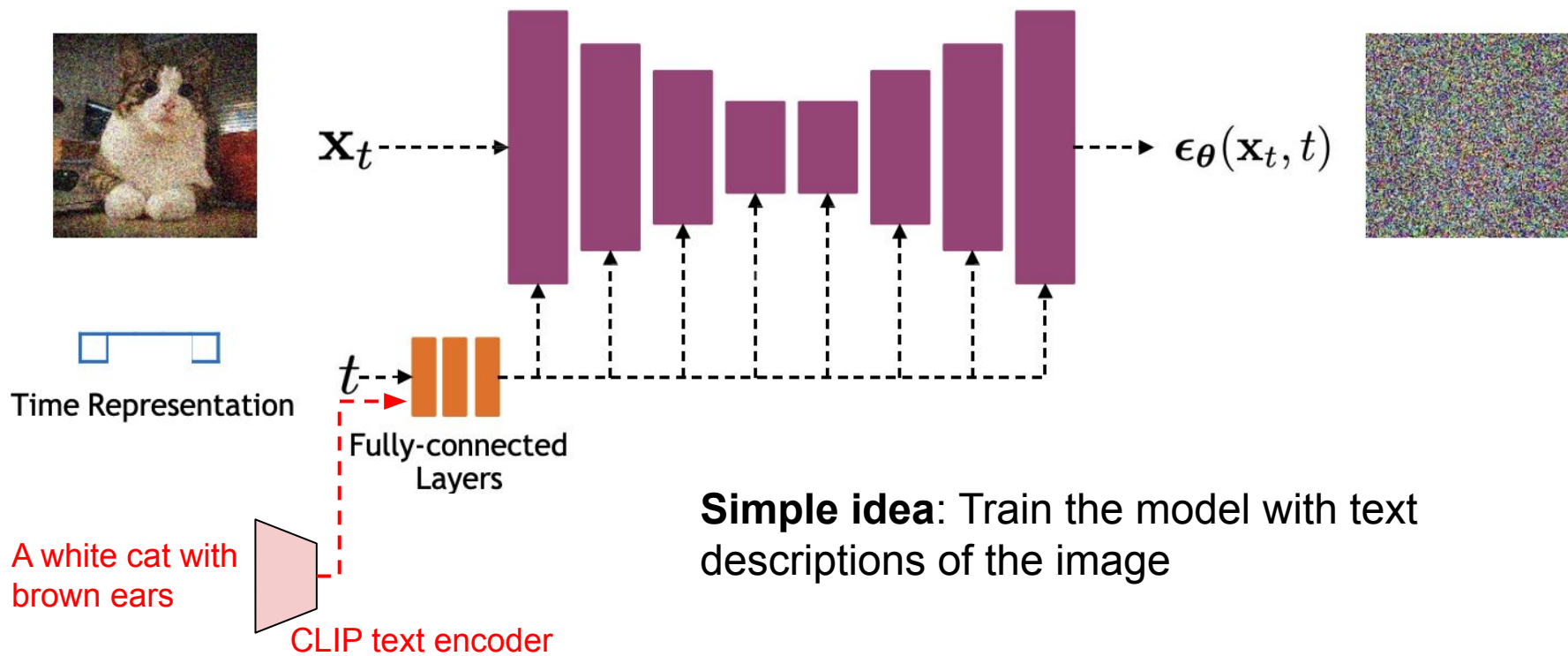
Time representation: sinusoidal positional embeddings.

Added in using: $\text{AdaGN}(h, y) = y_s \text{GroupNorm}(h) + y_b$

- h is the intermediate activations of the residual block following the first convolution in each layer,
- $y = [y_s, y_b]$ is obtained from a linear projection of the timestep

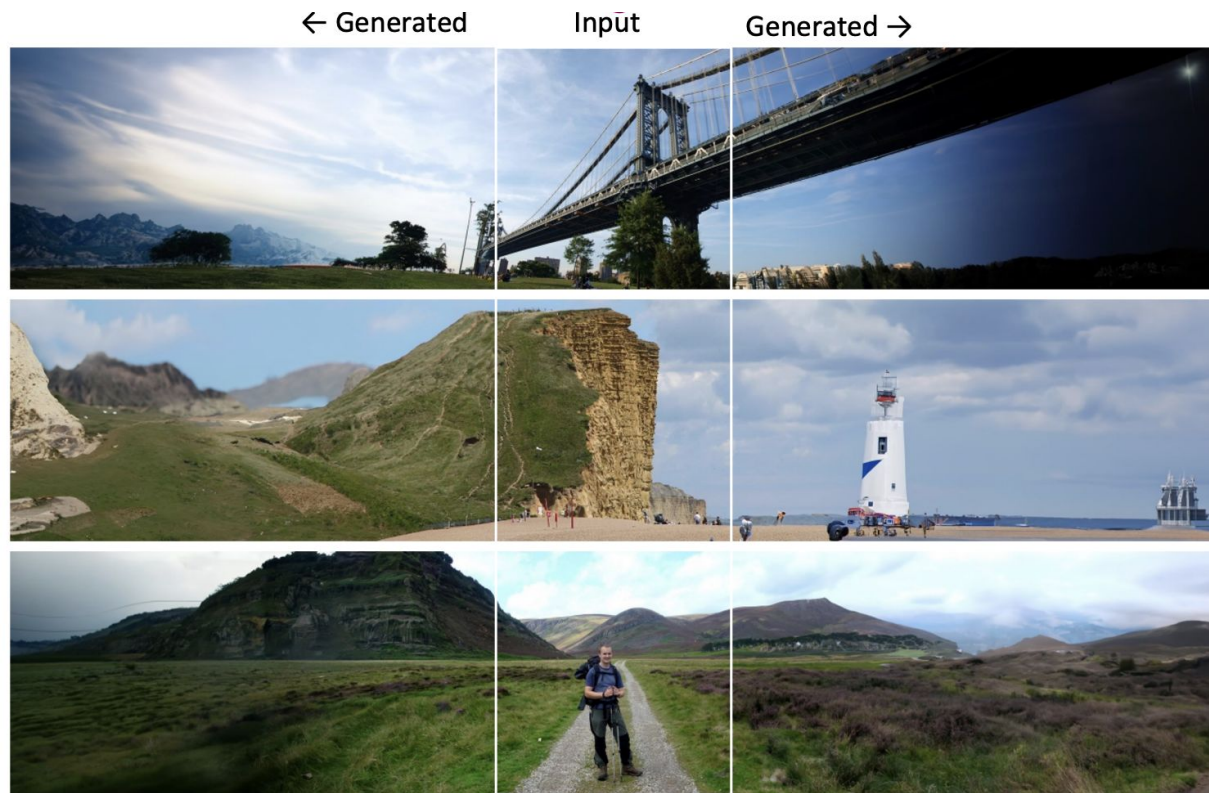


Text-conditioned generation



Simple idea: Train the model with text descriptions of the image

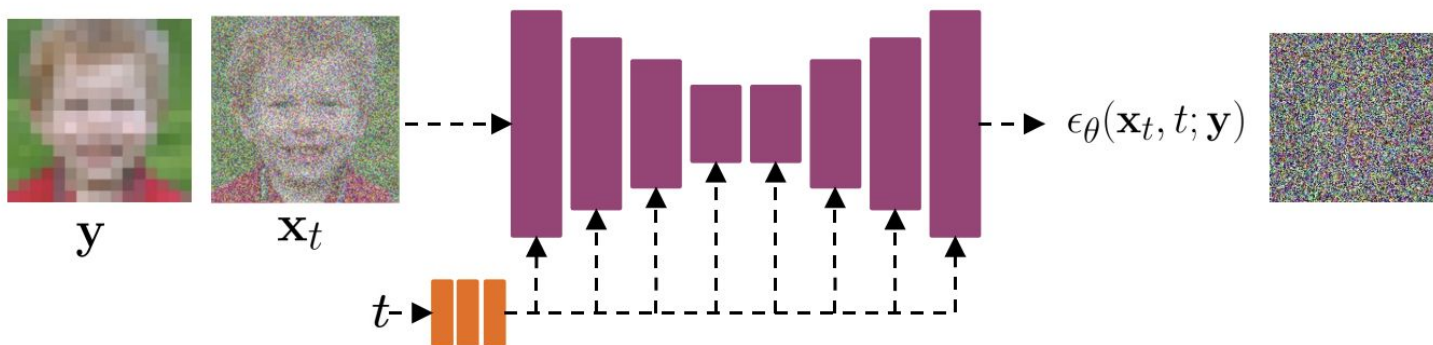
Application: panorama generation



Application: super-resolution

Learn a superresolution diffusion model conditioned on a low resolution image. y is a low resolution input image, x is a high resolution output image

$$\mathbb{E}_{\mathbf{x}, \mathbf{y}} \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \mathbb{E}_t \|\epsilon_{\theta}(\mathbf{x}_t, t; \mathbf{y}) - \epsilon\|_p^p$$



Saharia et al., Image Super-Resolution via Iterative Refinement, 2021

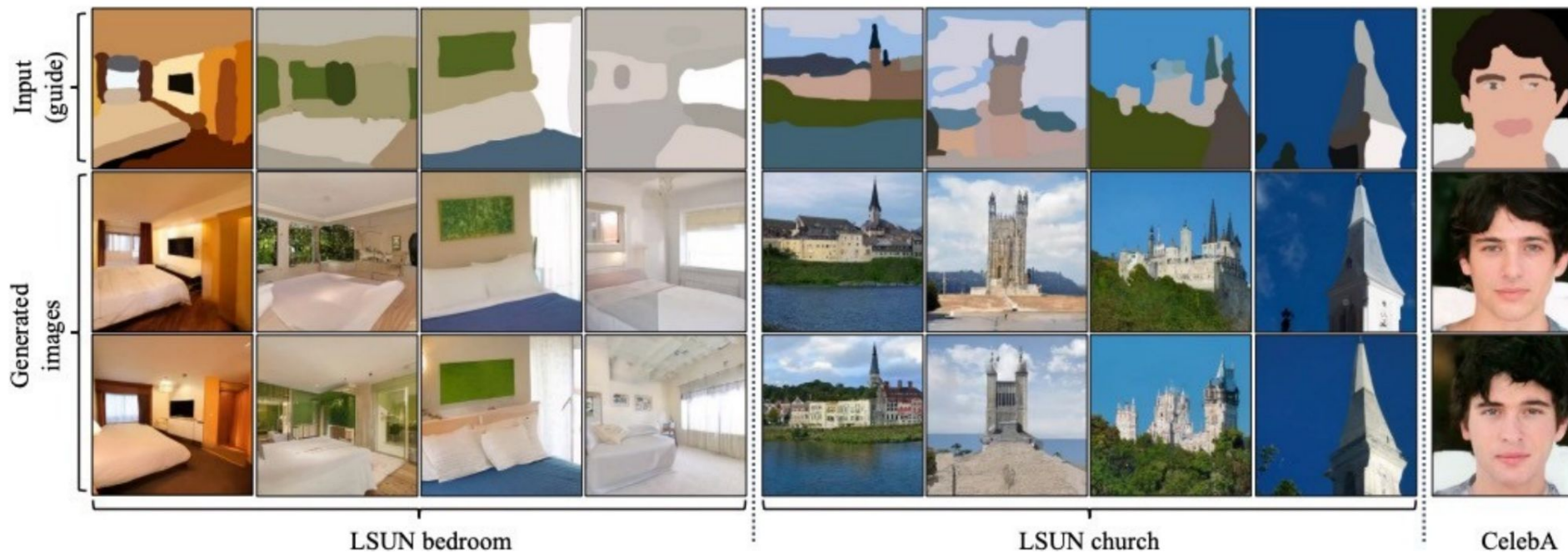
Application: super resolution

Natural Image Super-Resolution $64 \times 64 \rightarrow 256 \times 256$



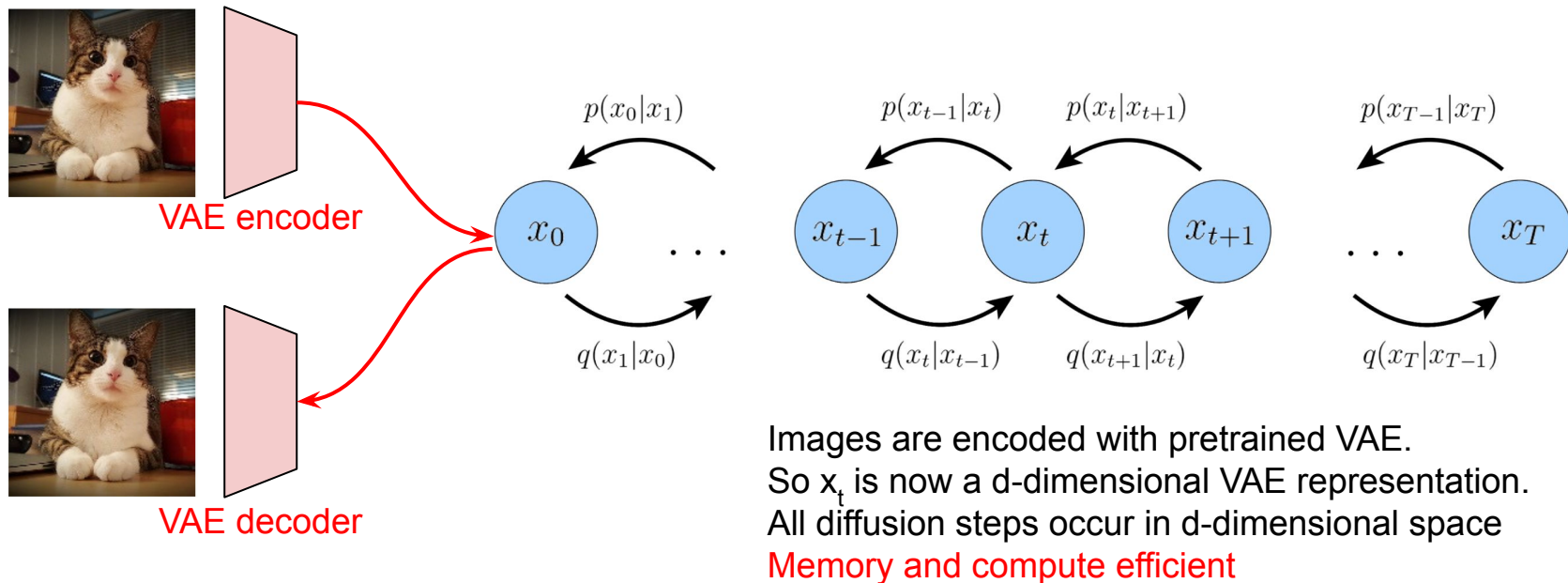
Saharia et al., Image Super-Resolution via Iterative Refinement, 2021

Application: image editing



Meng et al., SDEdit: Guided Image Synthesis and Editing with Stochastic Differential Equations, ICLR 2022

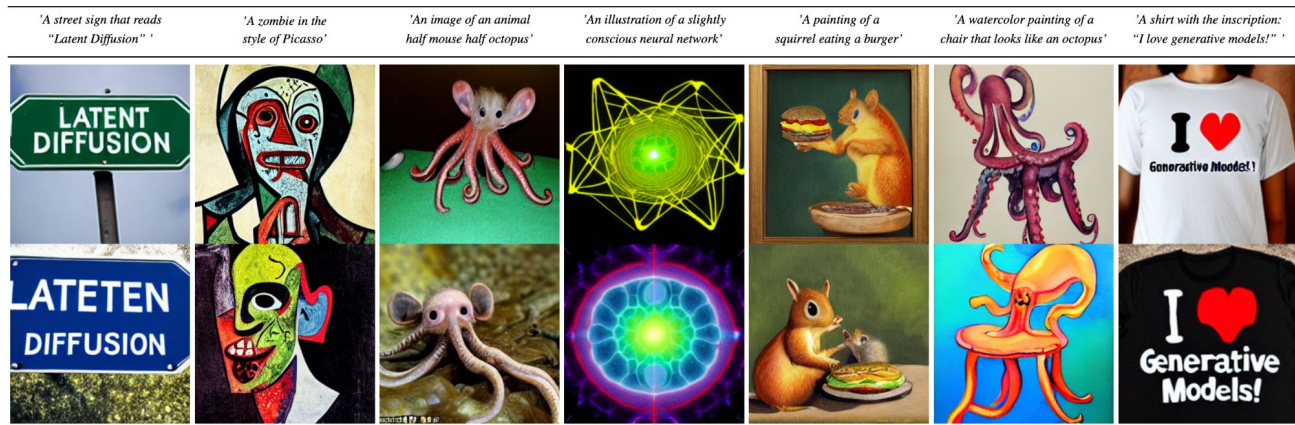
Latent diffusion models: perform diffusion over latent VAE encodings



Rombach et al. High-Resolution Image Synthesis with Latent Diffusion Models ArXiv 2022

Stable diffusion - from Stability AI

- Open sourced diffusion model - main model used for research
- Produces 512x512 images
- UNet with 860M params
- ViT-L text encoder with 123M params
- Fits in 10GB VRAM - fits on most GPUs



Rombach et al. High-Resolution Image Synthesis with Latent Diffusion Models ArXiv 2022

Imagen - Google

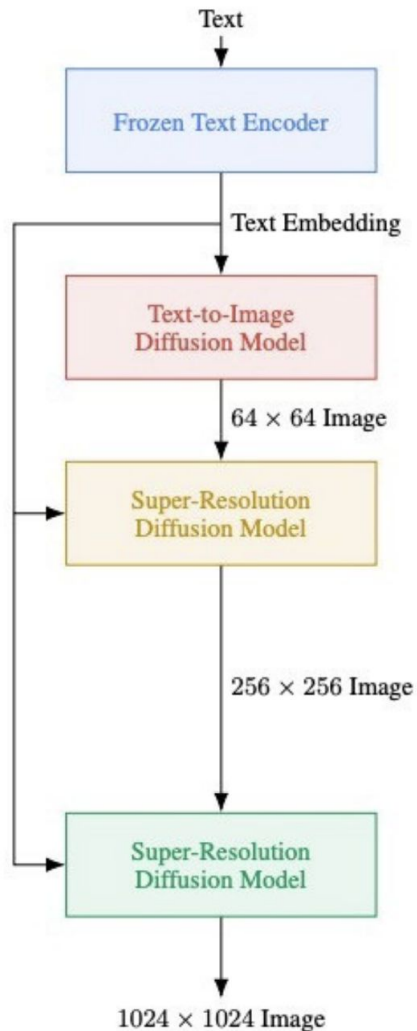
Combines:

- Latent diffusion model
- text conditioning
- 2 super-resolution models

To produce high quality 1024x1024 images

Saharia et al., "Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding", arXiv 2022.

Ranjay Krishna, George Stoica



"A Golden Retriever dog wearing a blue checkered beret and red dotted turtleneck."

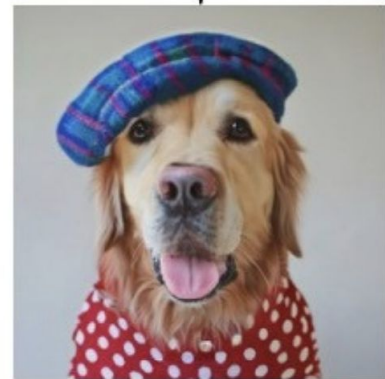


Imagen examples



A dragon fruit wearing karate belt in the snow.



A relaxed garlic with a blindfold reading a newspaper while floating in a pool of tomato soup.



A photo of a Shiba Inu dog with a backpack riding a bike. It is wearing sunglasses and a beach hat.

Last year: Sora video diffusion model

<https://openai.com/sora>

How did they do it?

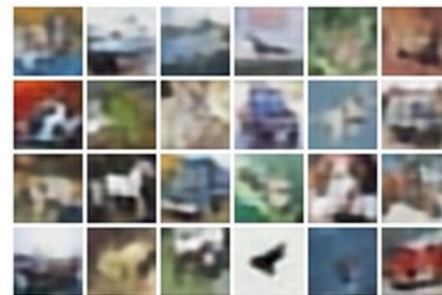
- More data (unknown data source)
- Replaced U-Net architecture with transformers

Comparing the different generative models

Q. Which ones are VAEs good at?

	Autoregressive (VAEs)	GANs	Diffusion
Mode coverage / diversity of generations			
Fast sampling			
High quality samples			

Comparing the different generative



VAEs are bad at generating high quality samples

	Autoregressive (VAEs)	GANs	Diffusion
Mode coverage / diversity of generations	✓		
Fast sampling	✓		
High quality samples	✗		

Comparing the different generative models

Q. Which ones are GANs good at?

	Autoregressive (VAEs)	GANs	Diffusion
Mode coverage / diversity of generations	✓		
Fast sampling	✓		
High quality samples	✗		

Comparing the different generative models

GANs suffer from mode collapse

	Autoregressive (VAEs)	GANs	Diffusion
Mode coverage / diversity of generations	✓	✗	
Fast sampling	✓	✓	
High quality samples	✗	✓	

Comparing the different generative models

Q. Which ones are Diffusion models good at?

	Autoregressive (VAEs)	GANs	Diffusion
Mode coverage / diversity of generations	✓	✗	
Fast sampling	✓	✓	
High quality samples	✗	✓	

Comparing the different generative models

Diffusion models are bad at sampling fast.

	Autoregressive (VAEs)	GANs	Diffusion
Mode coverage / diversity of generations	✓	✗	✓
Fast sampling	✓	✓	✗
High quality samples	✗	✓	✓

Next time

**Reinforcement Learning
+ Quiz!!!**