



Diffusion Models

Generative and Deep Learning (GDL)

Daide Bacciu (davide.bacciu@unipi.it)



UNIVERSITÀ DI PISA



Lecture Outline

- ◇ Introduction
 - ◇ Motivations
 - ◇ Learning to generate by denoising
- ◇ Denoising diffusion models
 - ◇ Forward & Reverse process
 - ◇ Training diffusion models
 - ◇ Implementation
- ◇ Advanced & Applications
 - ◇ Conditional generation
 - ◇ Multimodal

Why Diffusion Models?



“Diffusion Models Beat GANs on Image Synthesis” Dhariwal & Nichol, OpenAI, 2021



“a teddy bear on a skateboard in times square”

Ramesh et al.,
2022

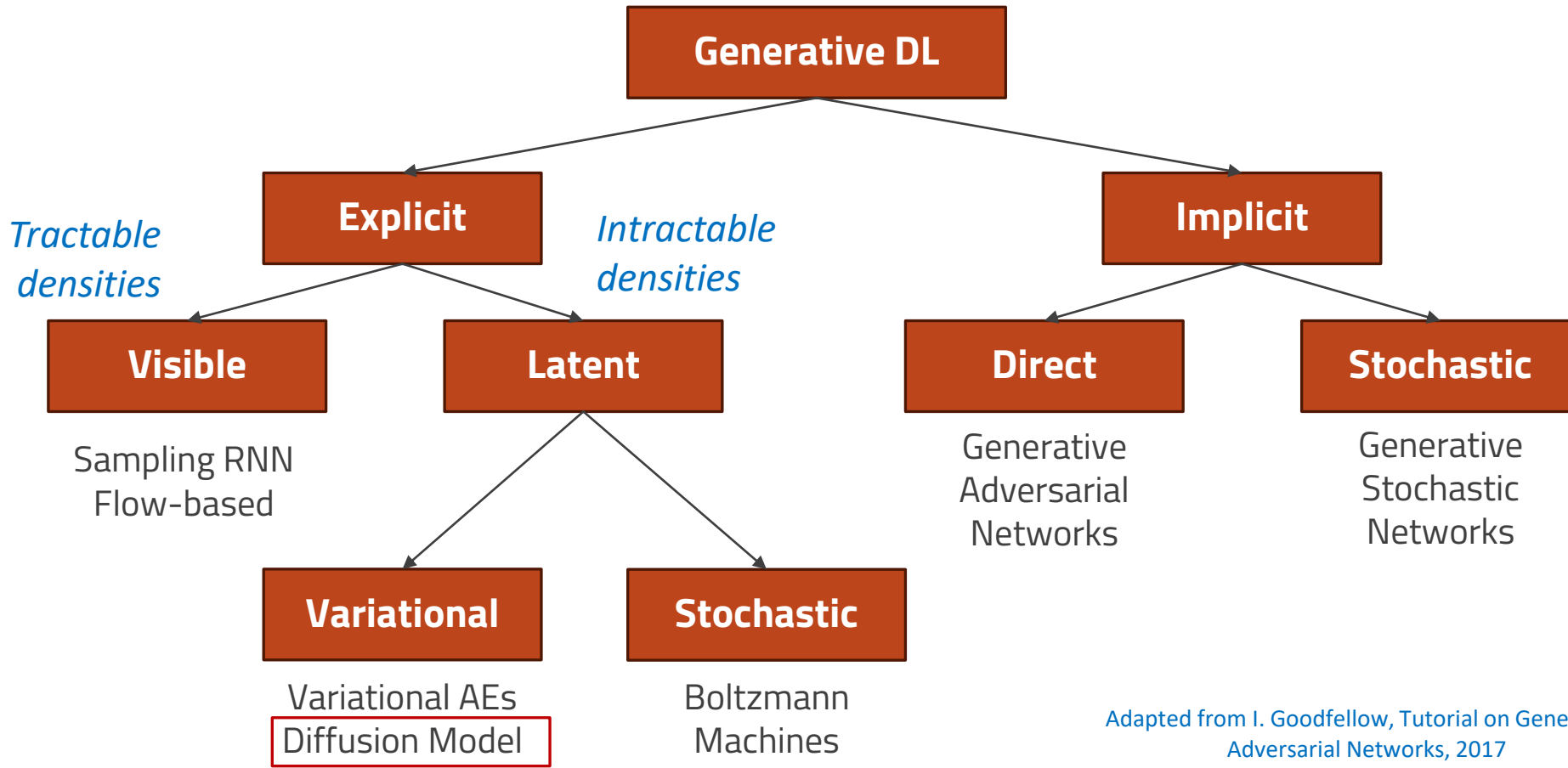
Why Diffusion Models?



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

A Taxonomy

**Diffusion models
latent space has same
size of data!**

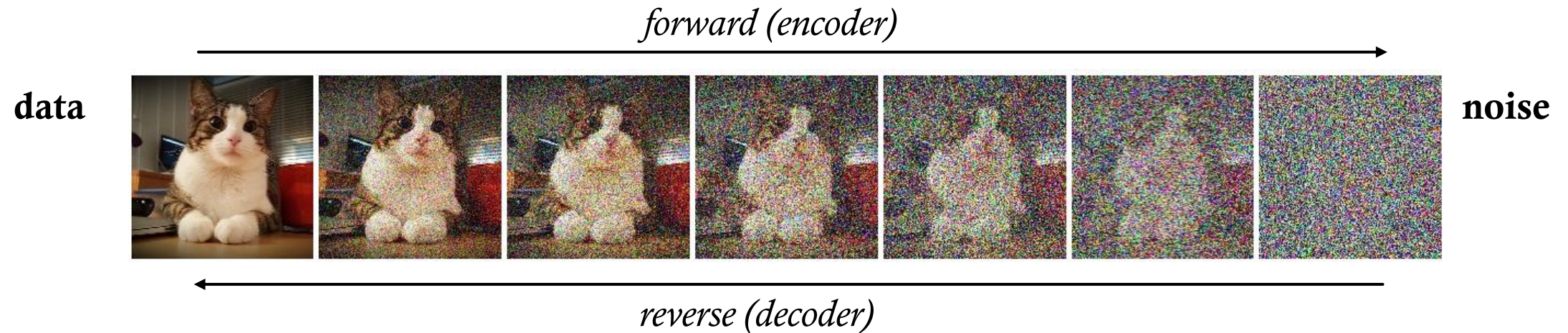


Adapted from I. Goodfellow, Tutorial on Generative Adversarial Networks, 2017

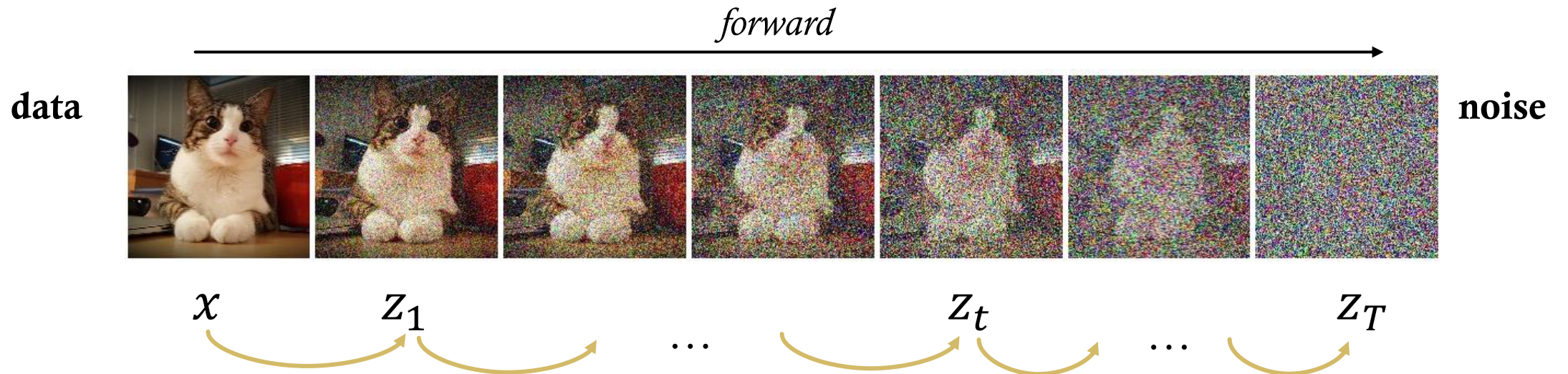
Generation as a diffusion process

Learning to generate by denoising

- ◇ Two processes
 - ◇ **Forward diffusion** gradually adding noise to input
 - ◇ **Reverse process** reconstructs data from noise (**generation**)
- ◇ The key is how to do this **efficiently**

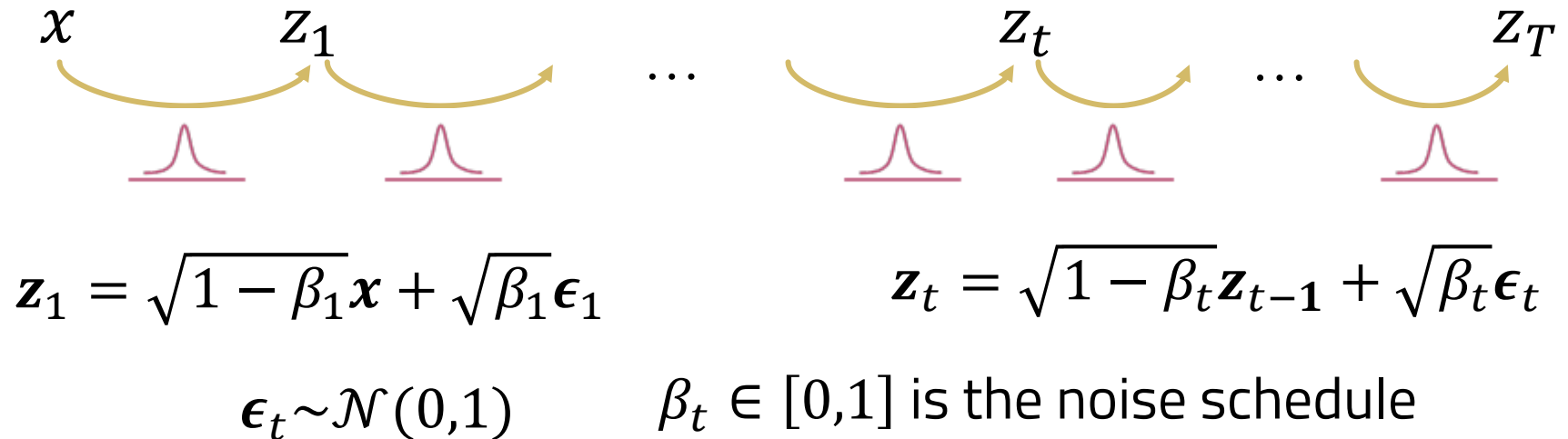
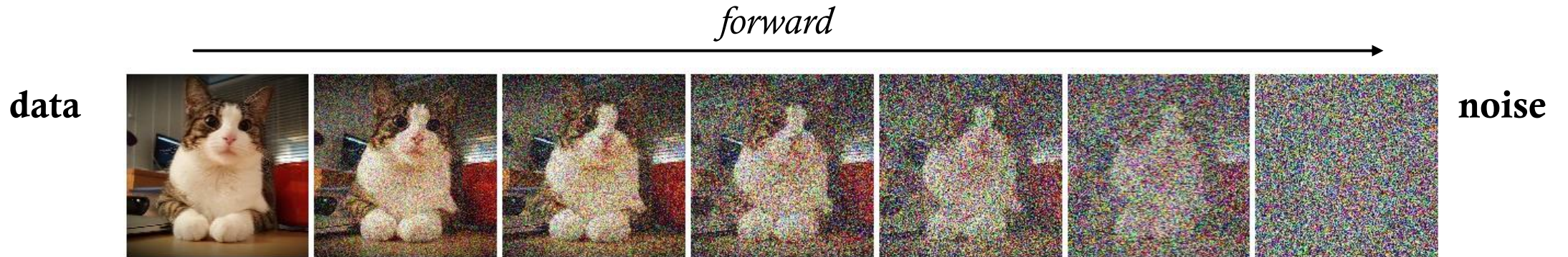


Forward Diffusion - Intuition

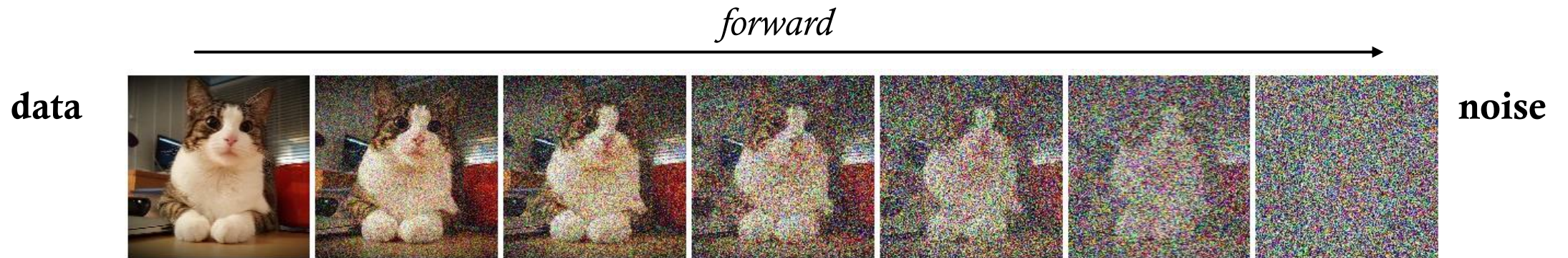


A fixed (i.e. **non-adaptive**) noise process in **T steps** mapping original data x into a **same-sized latent variables** z_t using simple **additive noise**

Forward Diffusion – Noise addition



Forward Diffusion – Distributions



$$q(\mathbf{z}_t | \mathbf{z}_{t-1}) = \mathcal{N}(\sqrt{1 - \beta_t} \mathbf{z}_{t-1}, \beta_t \mathbf{I}) \text{ where } \mathbf{z}_0 = \mathbf{x}$$

$$q(\mathbf{z}_1, \dots, \mathbf{z}_T | \mathbf{x}) = q(\mathbf{z}_1 | \mathbf{x}) \prod_{t=2}^T q(\mathbf{z}_t | \mathbf{z}_{t-1})$$

Diffusion Kernel

Generating \mathbf{z}_t sequentially is time-consuming so we use a **closed-form solution** for $q(\mathbf{z}_t | \cdot)$

$$q(\mathbf{z}_t | \mathbf{x}) = \mathcal{N}(\sqrt{\alpha_t} \mathbf{x}, (1 - \alpha_t) \mathbf{I}) \quad (\text{diffusion kernel})$$

$$\alpha_t = \prod_{s=1}^t (1 - \beta_s)$$

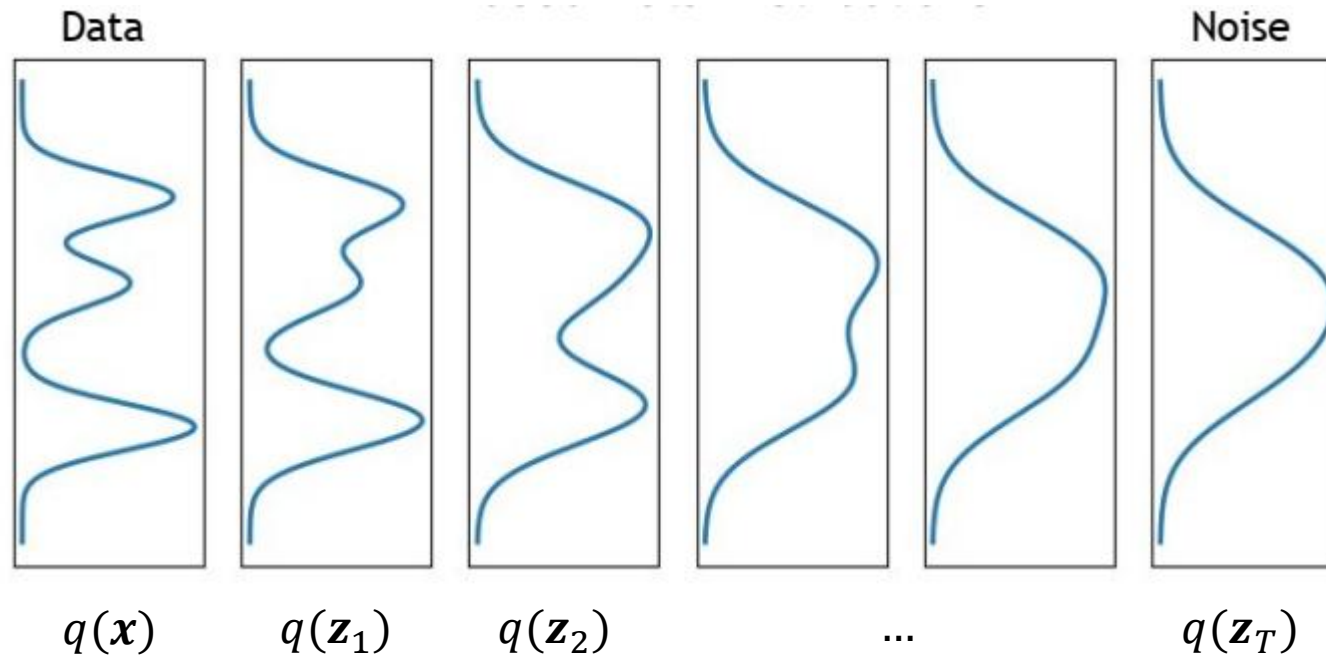
Which allows writing the marginal as

$$q(\mathbf{z}_t) = \int q(\mathbf{z}_t, \mathbf{x}) d\mathbf{x} = \int q(\mathbf{x}) q(\mathbf{z}_t | \mathbf{x}) d\mathbf{x}$$

data distribution

Evolution of diffused data distributions

$$q(\mathbf{z}_t) = \int q(\mathbf{z}_t, \mathbf{x}) d\mathbf{x} = \int q(\mathbf{x}) q(\mathbf{z}_t | \mathbf{x}) d\mathbf{x}$$



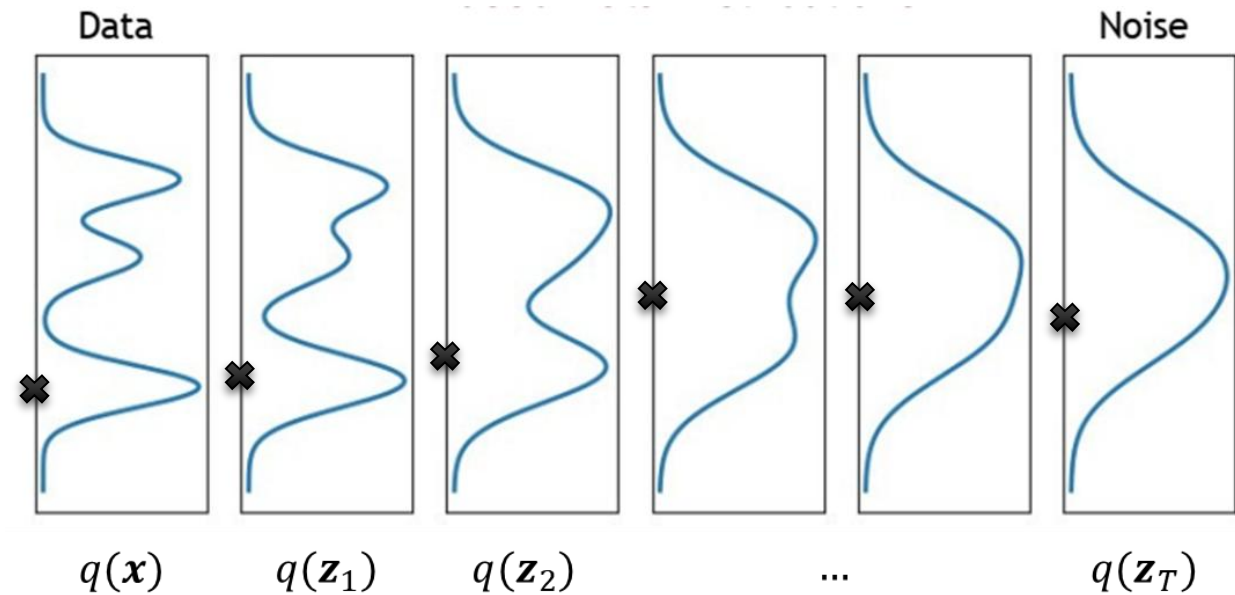
Denoising – Inverting the process

Sample $\mathbf{z}_T \sim \mathcal{N}(0,1)$

Iterate $\mathbf{z}_{t-1} \sim q(\mathbf{z}_{t-1}|\mathbf{z}_t)$

True denoising distribution is intractable

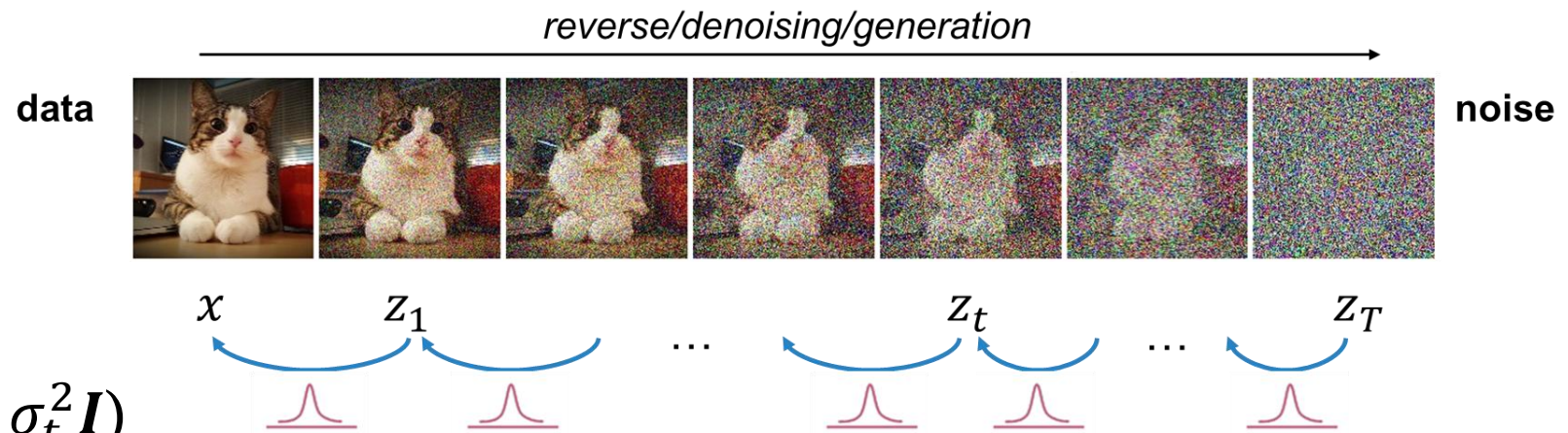
$$q(\mathbf{z}_{t-1}|\mathbf{z}_t) = \frac{q(\mathbf{z}_{t-1})q(\mathbf{z}_t|\mathbf{z}_{t-1})}{q(\mathbf{z}_t)}$$



We cannot de-mix noise if we don't know the starting point \mathbf{x} . If we do, then we can show that $q(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{x})$ is Normal

Denoising – Reverse Process

- ◇ Reverse process learns an approximated denoising distribution (decoder)
- ◇ Assuming reverse distributions are approximately Normal (reasonable if β_t are small and T large).



$$P(\mathbf{z}_T) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$P_{\theta}(\mathbf{z}_{t-1} | \mathbf{z}_t) = \mathcal{N}(\mu_{\theta}(\mathbf{z}_t, t), \sigma_t^2 \mathbf{I})$$

Mean of the denoised image \mathbf{z}_{t-1} predicted by the θ -parameterized model given \mathbf{z}_t (and time encoding)

Training diffusion models

Training

- ◇ Training follows the classical log-likelihood maximization view

$$\begin{aligned}\log P_{\theta}(\mathbf{x}) &= \log \int P_{\theta}(\mathbf{z}_1, \dots, \mathbf{z}_T, \mathbf{x}) d\mathbf{z}_{1\dots T} \\ &= \log \int P_{\theta}(\mathbf{x}|\mathbf{z}_1) \prod_{t=2}^T P_{\theta}(\mathbf{z}_{t-1}|\mathbf{z}_t) P_{\theta}(\mathbf{z}_T) d\mathbf{z}_{1\dots T}\end{aligned}$$

- ◇ ...which is, of course, intractable

Training – ELBO at the rescue

Introduce the encoder distribution q (with $\bar{\mathbf{z}} = \mathbf{z}_1, \dots, \mathbf{z}_T$)

$$\text{(loglik)} \log \int P_\theta(\bar{\mathbf{z}}, \mathbf{x}) d\bar{\mathbf{z}} \geq \int q(\bar{\mathbf{z}}|\mathbf{x}) \log \left[\frac{P_\theta(\bar{\mathbf{z}}, \mathbf{x})}{q(\bar{\mathbf{z}}|\mathbf{x})} \right] d\bar{\mathbf{z}} \quad (\text{ELBO})$$

Sparing some derivation and simplifications, we approximate ELBO as:

$$\underbrace{\mathbb{E}_{q(\mathbf{z}_1|\mathbf{x})} [\log P_\theta(\mathbf{x}|\mathbf{z}_1)]}_{\text{Reconstruction term}} - \underbrace{\sum_{t=2}^T KL(q(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{x}) || P_\theta(\mathbf{z}_{t-1}|\mathbf{z}_t))}_{\text{Aligns predicted and original (input-conditional) denoising densities}}$$

ELBO Loss Function

Distributions in KL are all Gaussians so can write the full form of the loss

$$\sum_{\mathbf{x}} \left(\underbrace{(-\log \mathcal{N}(\mu_{\theta}(\mathbf{z}_1, t), \sigma_1^2 \mathbf{I}))}_{\text{Reconstruction}} + \sum_{t=2}^T \frac{1}{2\sigma_t^2} \left\| \underbrace{\left(\frac{(1 - \alpha_{t-1})}{(1 - \alpha_t)} \sqrt{1 - \beta_t} \mathbf{z}_t + \frac{\sqrt{\alpha_{t-1} \beta_t}}{(1 - \alpha_t)} \mathbf{x} \right)}_{\text{Target mean of } q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{x})} - \underbrace{\mu_{\theta}(\mathbf{z}_t, t)}_{\text{predicted } \mathbf{z}_{t-1}} \right\|^2 \right)$$

Minimize difference between estimate of \mathbf{z}_{t-1} and the most likely value from ground truth-denoised data

Training – Practical view

Loss can be heavily simplified by reparameterizing so that the model predicts the noise $\epsilon_{\theta}(\cdot)$ that was mixed with the original data, rewriting \mathbf{x} as

$$\mathbf{x} = \frac{1}{\sqrt{\alpha_t}} \mathbf{z}_t - \frac{\sqrt{1 - \alpha_t}}{\sqrt{\alpha_t}} \epsilon_t$$

A network which predicts the unit noise given current noised input \mathbf{z}_t

Inserting \mathbf{x} above into the ELBO yields (after a while)

$$Loss(\theta) = \sum_{\mathbf{x}} \sum_{t=1}^T \left\| \epsilon_{\theta} \left(\underbrace{\sqrt{\alpha_t} \mathbf{x} + \sqrt{1 - \alpha_t} \epsilon_t}_{\mathbf{z}_t}, t \right) - \epsilon_t \right\|^2$$

J. Ho et al, NeurIPS 2020

Implementation - Training

Algorithm 18.1: Diffusion model training

Input: Training data \mathbf{x}

Output: Model parameters θ

repeat

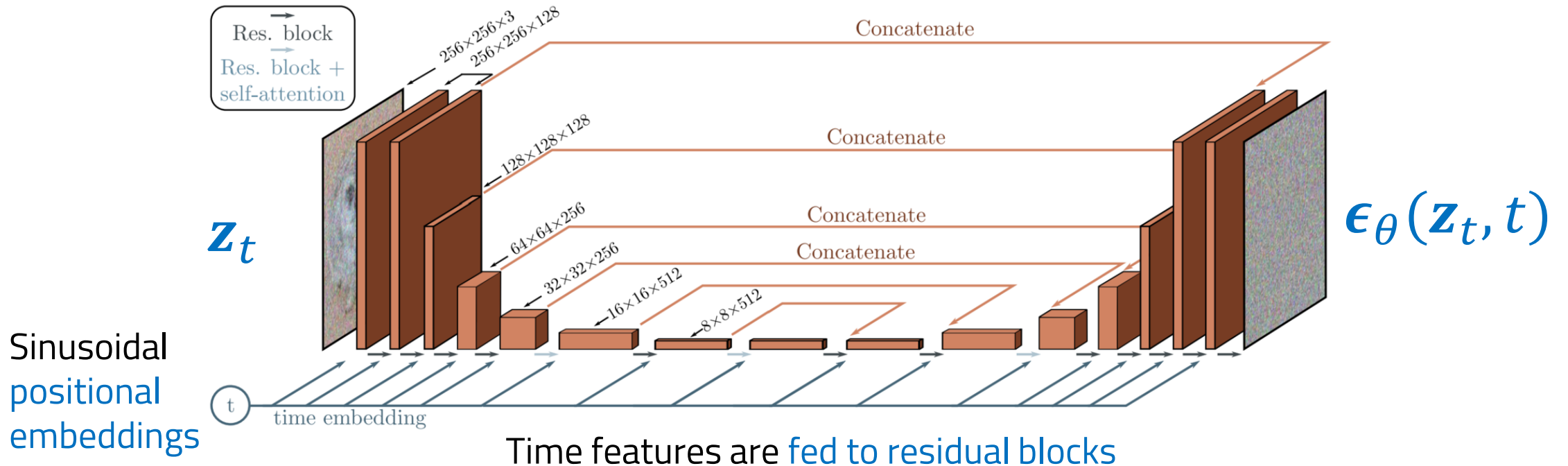
```
  for  $i \in \mathcal{B}$  do                                     // For every training example index in batch
     $t \sim \text{Uniform}[1, \dots, T]$                        // Sample random timestep
     $\epsilon \sim \text{Norm}[\mathbf{0}, \mathbf{I}]$                  // Sample noise
     $l_i = \left\| \epsilon_{\theta}(\sqrt{\alpha_t} \mathbf{x}_i + \sqrt{1 - \alpha_t} \epsilon, t) - \epsilon \right\|^2$  // Compute individual loss
  Accumulate losses for batch and take gradient step
```

until converged

During forward we add noise to image. During reverse we predict that noise with a DNN and then subtract it from the image to denoise it.

Diffusion model for images

U-Net architectures with ResNet blocks and self-attention layers



J. Ho et al, NeurIPS 2020
Dharivwal and Nichol NeurIPS 2021

Noise Schedules

- ◇ Terms β_t and σ_t control variance of forward diffusion and reverse denoising, respectively
 - ◇ β_t linear schedule
 - ◇ $\sigma_t^2 = \beta_t$
- ◇ **Slowly increase** the amount of added noise (as **high-resolution** information is **corrupted first**)
- ◇ Alternatives
 - ◇ σ_t can be learned by minimizing the bound
 - ◇ β_t can be learned by minimizing the variance of the training objective

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



$$P_\theta(\mathbf{z}_{t-1} | \mathbf{z}_t) = \mathcal{N}(\mu_\theta(\mathbf{z}_t, t), \sigma_t^2 \mathbf{I})$$

Flavours of diffusion

Guided Generation – Classifier Guidance

Guide diffusion process using **auxiliary data c** , using the gradient of a trained classifier as guidance

1. Train the diffusion model unconditionally
2. Train a classifier $P(c|\mathbf{z}_t)$ where \mathcal{C} are conditioning labels
3. Add an extra term when sampling the diffusion model, i.e. when reconstructing \mathbf{z}_{t-1} from \mathbf{z}_t , that modifies the reconstruction in the direction given by the gradient of a classifier

$$\mathbf{z}_{t-1} = \underbrace{\hat{\mathbf{z}}_{t-1} + \sigma_t \boldsymbol{\epsilon}}_{\text{Reversed diffusion}} + \underbrace{\sigma_t^2 \frac{\partial P(c|\mathbf{z}_t)}{\partial \mathbf{z}_t}}_{\text{Classifier guidance}}$$

Classifier Guidance - Issues

- ◇ Classifier guidance comes from **mixing the predicted score function** of the unconditional diffusion model with the **classifier gradients**

$$\frac{\partial \log P_\gamma(\mathbf{z}_t|c)}{\partial \mathbf{z}_t} = \frac{\partial \log P(\mathbf{z}_t)}{\partial \mathbf{z}_t} + \gamma \frac{\partial \log P(c|\mathbf{z}_t)}{\partial \mathbf{z}_t}$$

- ◇ γ guidance scale
- ◇ **Classifier receives a noisy input \mathbf{z}_t** at each step (can't use pretrained ones)
- ◇ Most of \mathbf{z}_t is of no use for predicting $c \Rightarrow$ **arbitrary classifier gradients**

Classifier-free Guidance

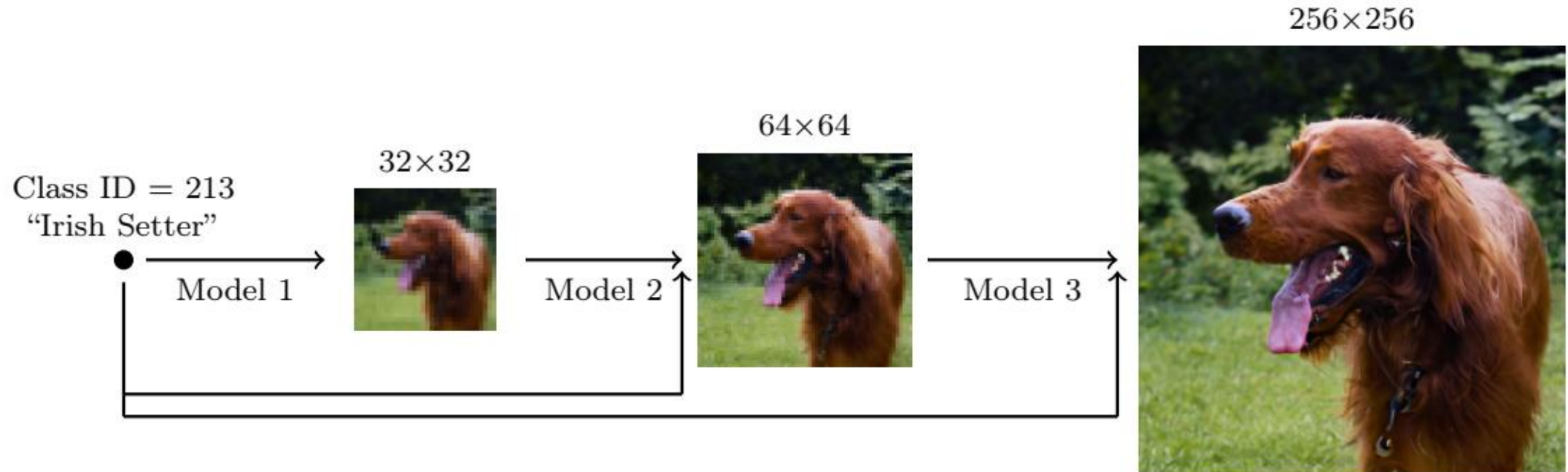
J. Ho, T. Salimans, NeurIPS 2021
Workshop DGMs Applications

Derive guidance from **Bayes rule**

$$\frac{\partial \log P_\gamma(\mathbf{z}_t | c)}{\partial \mathbf{z}_t} = (1 - \gamma) \underbrace{\frac{\partial \log P(\mathbf{z}_t)}{\partial \mathbf{z}_t}}_{\text{Unconditional diffusion score}} + \gamma \underbrace{\frac{\partial \log P(\mathbf{z}_t | c)}{\partial \mathbf{z}_t}}_{\text{Conditional diffusion score}}$$

- ◆ Training **conditional diffusion with dropout** (randomly removing conditioning)
- ◆ Conditioning replaced by flag input (presence/absence of conditioning) => **single model for conditional/unconditional** diffusion

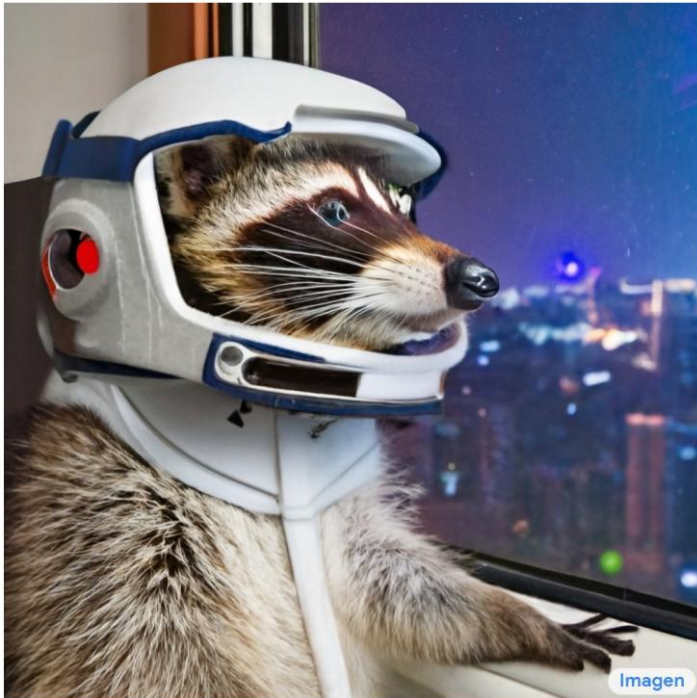
High-Resolution Image Generation



J. Ho et al, JMLR 2022

Conditional Generation

"A photo of a raccoon wearing an astronaut helmet, looking out of the window at night" (IMAGEN)



Text-2-image

Panorama completion

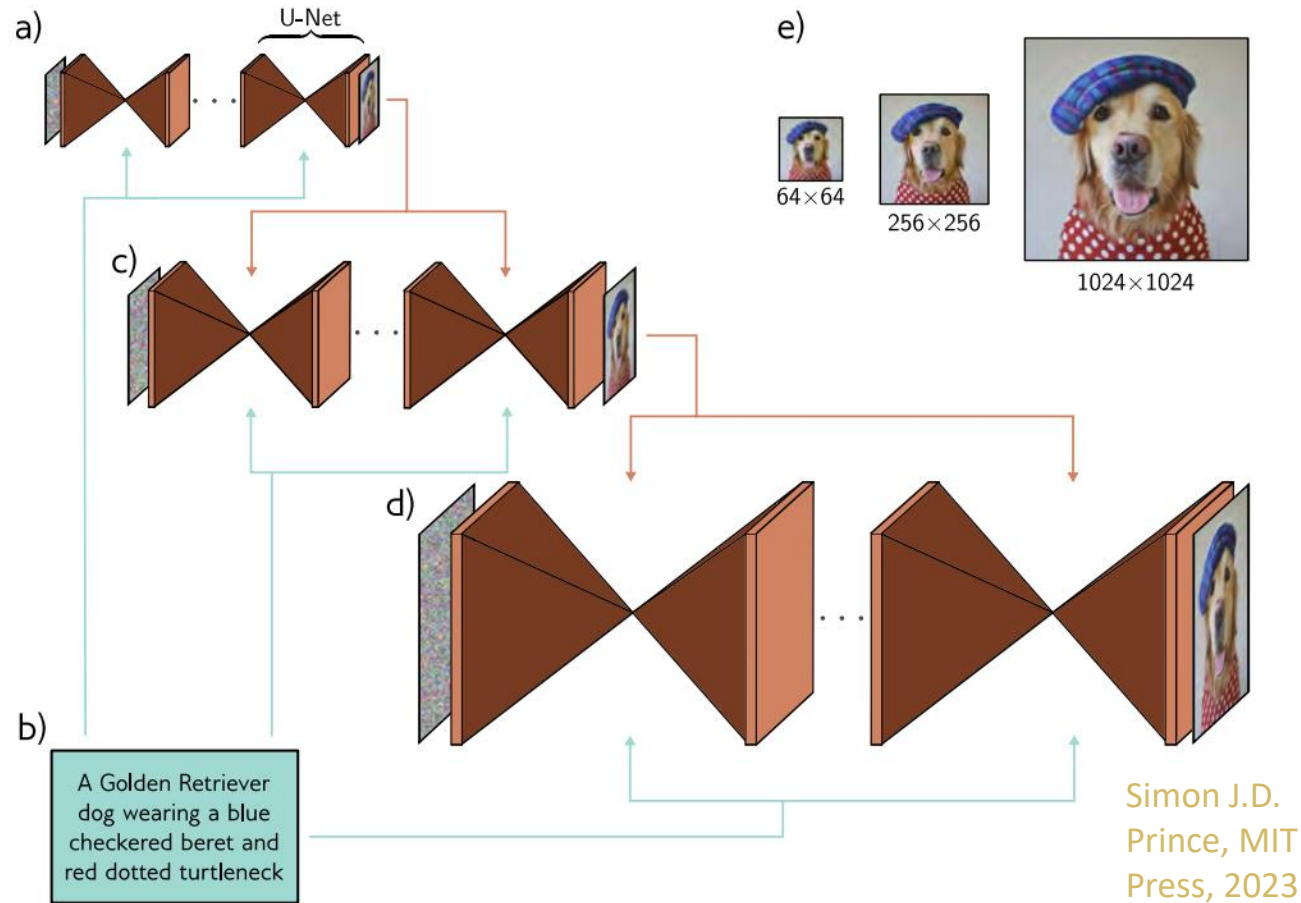


Generated

Input

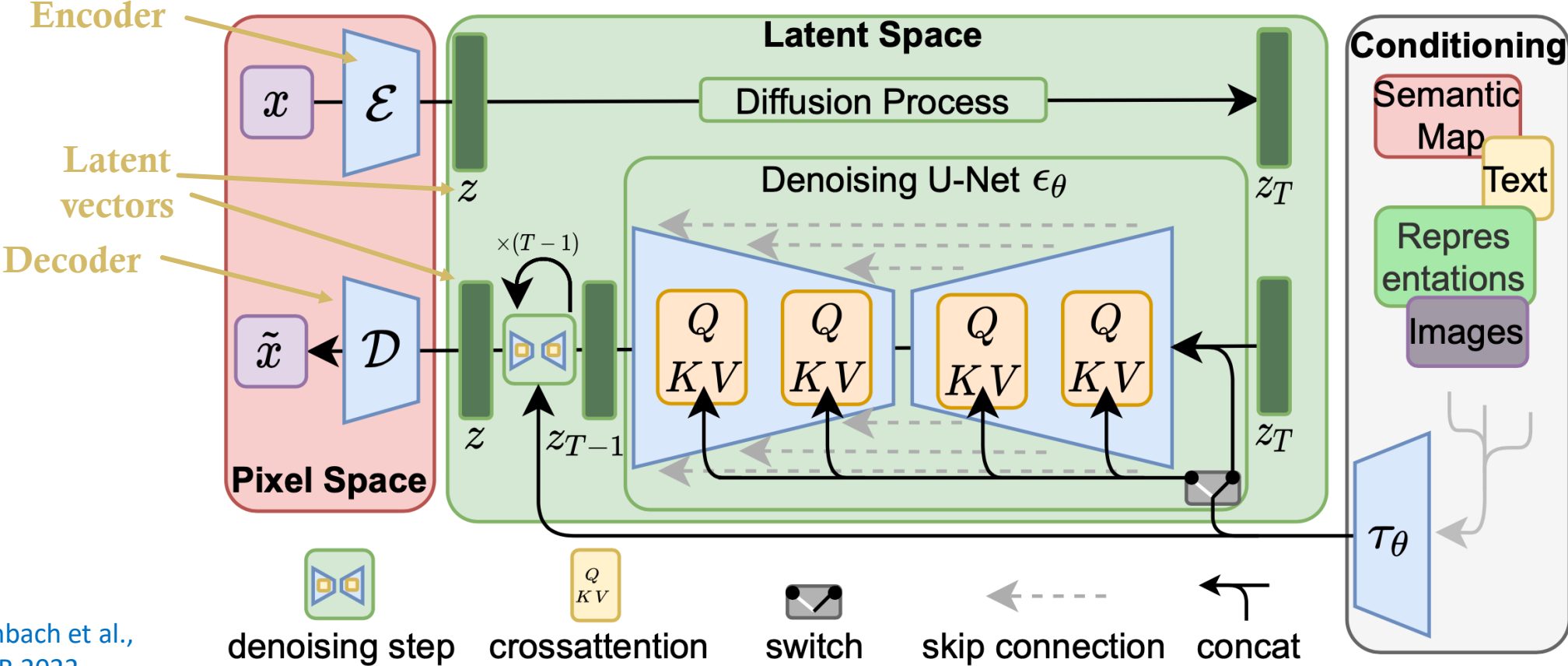
Generated

Cascaded Conditional Generation



- ◇ **Scalar** – vector embedding + spatial addition (or adaptive group normalization)
- ◇ **Image** - channel-wise concatenation of the conditional image
- ◇ **Text** - vector embedding + spatial addition or cross-attention

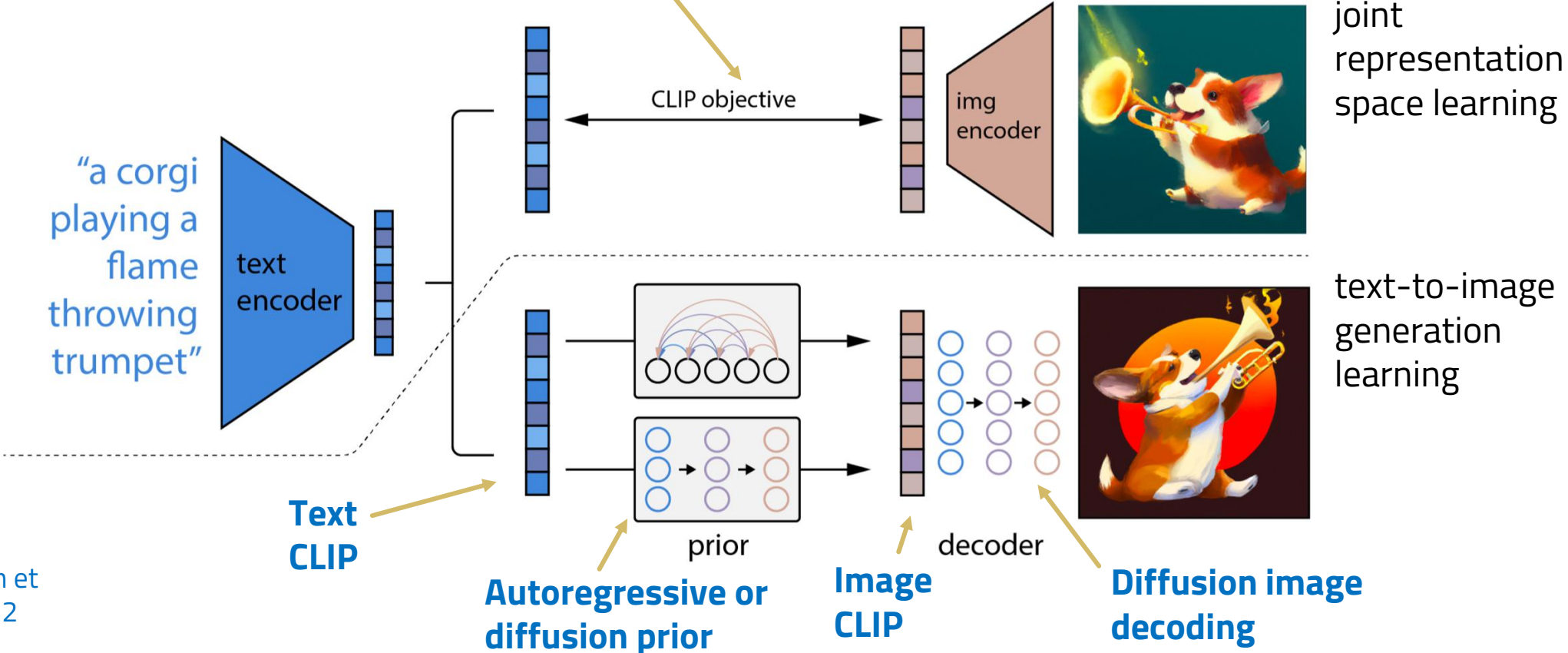
Latent Space Diffusion



Run diffusion in the latent space instead of pixel space for cost saving

DALL-E 2 – Diffusion Model

CLIP => scalable joint txt-img representations



Aditya Ramesh et al, OpenAI 2022



Wrap-up

Take Home Messages

- ◇ Generate data from noise through a **learned incremental denoising with fixed steps**
 - ◇ Diffusion process can be reversed if the variance of the gaussian noise added at each step is small enough
 - ◇ Training goal is to make sure that the predicted noise map at each step is unit gaussian
 - ◇ During generation, subtract the predicted noise from the noisy image at time t to generate the image at time $t-1$
- ◇ Diffusion can be **computationally involved**
 - ◇ Need to take many small steps
 - ◇ Vanilla diffusion on a latent space same size as the original data
- ◇ **Guided generation** can improve sample quality (and reduce diversity)
- ◇ **Latent space diffusion**
 - ◇ Improves efficiency of generation
 - ◇ Generalizes which data that can be used (including discrete objects)
 - ◇ Allows introducing semantic structuring in latent space

Next Lecture

- ◆ Causal Representation Learning
 - ◆ Lecture by Riccardo Massidda
- ◆ Next week
 - ◆ Matching-based approaches: the missing link between VAE, diffusion and normalizing flows
 - ◆ Generative deep learning module wrap-up
 - ◆ Deep Learning for graphs