



# Diffusion Models

Ganapathy Krishnamurthi

# Overview → Generate images

- Consists of an [encoder] and a [decoder]
  - Encoder takes data sample  $x$  and maps it through a series of intermediate latent variables,  $z_1.. z_T$ .
  - The decoder reverses this process.
    - It starts with  $z_T$  and maps it back through  $z_{T-1}$  to  $z_1$  until it recreates  $x$
  - In both encoder and decoder, the mappings are stochastic rather than deterministic
- $f(x) \rightarrow z_1, z_2, \dots, z_T$
- $x \rightarrow z_1 \rightarrow z_2 \rightarrow \dots \rightarrow z_T$  ... time index
- $x \rightarrow f(x) \rightarrow z_1$   
 $f(z_1) \rightarrow z_2$ ,  $f(z_2) \rightarrow z_3$ , ...
- $z_T \rightarrow z_{T-1} \rightarrow z_{T-2} \dots \rightarrow z_1 \rightarrow x$
- $f^{-1}(z_T) \rightarrow z_{T-1} = f^{-1}(z_{T-1}) \rightarrow z_{T-2}$

# Overview- Encoder

$$\mathbf{x} \rightarrow \mathbf{z}_1 \sim N(0, \sigma) \rightarrow \mathbf{z}_2 \sim N(0, \sigma)$$

- Encoder is pre-specified- it adds noise to the input gradually
- With enough noise addition steps the information in the image is wiped out and we are left only with noise
- All learned parameters are in the decoder

$$\mathbf{z}_1 \rightarrow \mathbf{x}$$

# Overview-Decoder

- A series of networks are trained to map backward between each of the adjacent pair of variables  $z_t$  and  $z_{t-1}$

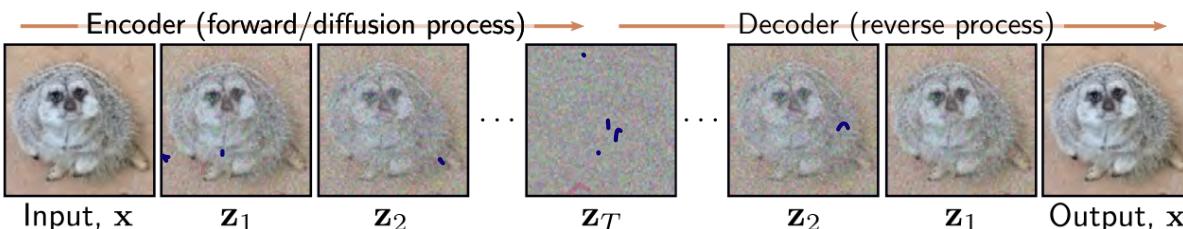
$z_T \xrightarrow{f_1} z_{T-1} \xrightarrow{f_2} z_{T-2} \xrightarrow{f_3} \dots \rightarrow x \sim p_{\theta}(x|z_t)$

- The loss function encourages the network to invert the corresponding encoder step

- Noise is gradually removed until a realistic looking data is obtained 

- To generate a new data example  $x$ , we draw from the noise distribution and pass it through the decoder

# Encoder-Decoder



- Diffusion models. The encoder (forward, or diffusion process) maps the input  $x$  through a series of latent variables  $z_1 \dots z_T$ .
- This process is pre-specified and gradually mixes the data with noise until only noise remains.
- The decoder (reverse process) is learned and passes the data back through the latent variables, removing noise at each stage.
- After training, new examples are generated by sampling noise vectors  $z_T$  and passing them through the decoder.

# Encoder

Noise Schedule

$\beta \rightarrow$  noise schedule  $0 \cdot 1$

$\beta \rightarrow$  can be time dependent

$$\epsilon_t \sim N(0, \sigma)$$

$$\Rightarrow z_1 = \sqrt{1 - \beta_1} \cdot x + \sqrt{\beta_1} \cdot \epsilon_1$$

$$z_t = \sqrt{1 - \beta_t} \cdot z_{t-1} + \sqrt{\beta_t} \cdot \epsilon_t, \quad \forall t \in \{2, \dots, T\},$$

$$z_T = \sqrt{1 - \beta_T} z_0 + \sqrt{\beta_T} \cdot \epsilon_T$$

$$q(z_1|x) = \text{Norm}_{z_1} \left[ \sqrt{1 - \beta_1} x, \beta_1 I \right]$$

$$q(z_t|z_{t-1}) = \text{Norm}_{z_t} \left[ \sqrt{1 - \beta_t} z_{t-1}, \beta_t I \right] \quad \forall t \in \{2, \dots, T\}.$$

# Decoder

$$\begin{aligned}
 Pr(\mathbf{z}_T) &= \text{Norm}_{\mathbf{z}_T} [\mathbf{0}, \mathbf{I}] \\
 Pr(\mathbf{z}_{t-1} | \mathbf{z}_t, \phi_t) &= \text{Norm}_{\mathbf{z}_{t-1}} \left[ \mathbf{f}_t[\mathbf{z}_t, \phi_t], \sigma_t^2 \mathbf{I} \right] \\
 Pr(\mathbf{x} | \mathbf{z}_1, \phi_1) &= \text{Norm}_{\mathbf{x}} \left[ \mathbf{f}_1[\mathbf{z}_1, \phi_1], \sigma_1^2 \mathbf{I} \right],
 \end{aligned}$$

a.k.

$$\begin{aligned}
 q(\mathbf{z}_1 | \mathbf{x}) &= \text{Norm}_{\mathbf{z}_1} \left[ \sqrt{1 - \beta_1} \mathbf{x}, \beta_1 \mathbf{I} \right] \\
 q(\mathbf{z}_t | \mathbf{z}_{t-1}) &= \text{Norm}_{\mathbf{z}_t} \left[ \sqrt{1 - \beta_t} \mathbf{z}_{t-1}, \beta_t \mathbf{I} \right]
 \end{aligned}
 \quad \forall t \in \{2, \dots, T\}.$$

d

known

Now we replace the model  $\hat{\mathbf{z}}_{t-1} = \mathbf{f}_t[\mathbf{z}_t, \phi_t]$  with a new model  $\hat{\epsilon} = \mathbf{g}_t[\mathbf{z}_t, \phi_t]$  which predicts the noise  $\epsilon$  that was mixed with  $\mathbf{x}$  to create  $\mathbf{z}_t$ :

$$\mathbf{f}_t[\mathbf{z}_t, \phi_t] = \left[ \frac{1}{\sqrt{1 - \beta_t}} \mathbf{z}_t - \frac{\beta_t}{\sqrt{1 - \alpha_t} \sqrt{1 - \beta_t}} \mathbf{g}_t[\mathbf{z}_t, \phi_t] \right].$$

# Training Algorithm

$$x + \epsilon \rightsquigarrow z \rightarrow \square \rightarrow \epsilon$$

$\bar{x} - \epsilon$

**Algorithm 18.1:** Diffusion model training

---

**Input:** Training data  $x$   
**Output:** Model parameters  $\phi_t$

```

repeat
    for  $i \in \mathcal{B}$  do
         $t \sim \text{Uniform}[1, \dots, T]$  // For every training example index in batch
         $\epsilon \sim \text{Norm}[0, I]$  // Sample random timestep
         $\ell_i = \left\| g_t[\sqrt{\alpha_t}x_i + \sqrt{1-\alpha_t}\epsilon, \phi_t] - \epsilon \right\|^2$  // Compute individual loss
    Accumulate losses for batch and take gradient step
until converged
     $\rightarrow$  number of steps

```

$$g_t(z_t, \phi_t) \rightarrow \epsilon$$

**Algorithm 18.2:** Sampling

---

**Input:** Model,  $g_t[\bullet, \phi_t]$   
**Output:** Sample,  $x$

```

 $z_T \sim \text{Norm}_z[0, I]$  // Sample last latent variable
for  $t = T \dots 2$  do
     $\hat{z}_{t-1} = \frac{1}{\sqrt{1-\beta_t}} z_t - \frac{\beta_t}{\sqrt{1-\alpha_t}\sqrt{1-\beta_t}} g_t[z_t, \phi_t]$  // Predict previous latent variable
     $\epsilon \sim \text{Norm}_\epsilon[0, I]$  // Draw new noise vector
     $z_{t-1} = \hat{z}_{t-1} + \sigma_t \epsilon$  // Add noise to previous latent variable
     $x = \frac{1}{\sqrt{1-\beta_1}} z_1 - \frac{\beta_1}{\sqrt{1-\alpha_1}\sqrt{1-\beta_1}} g_1[z_1, \phi_1]$  // Generate sample from  $z_1$  without noise

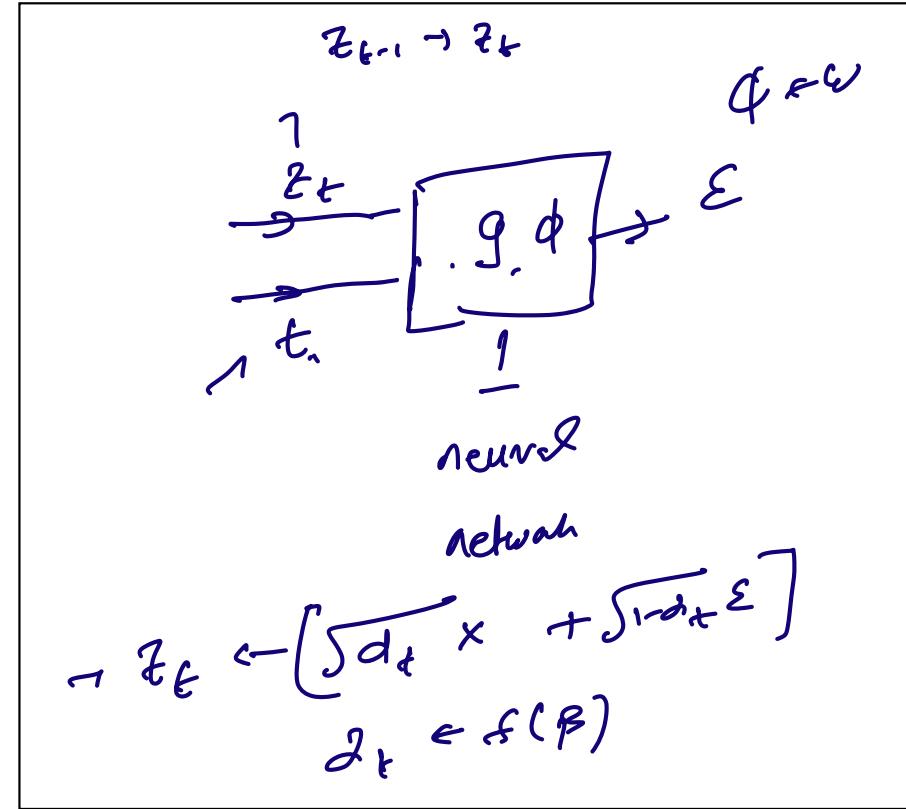
```

Draw noise  $\epsilon \sim \mathcal{N}$

$$z_T \rightarrow \epsilon \rightarrow z_T - \epsilon = \hat{z}_{T-1}$$

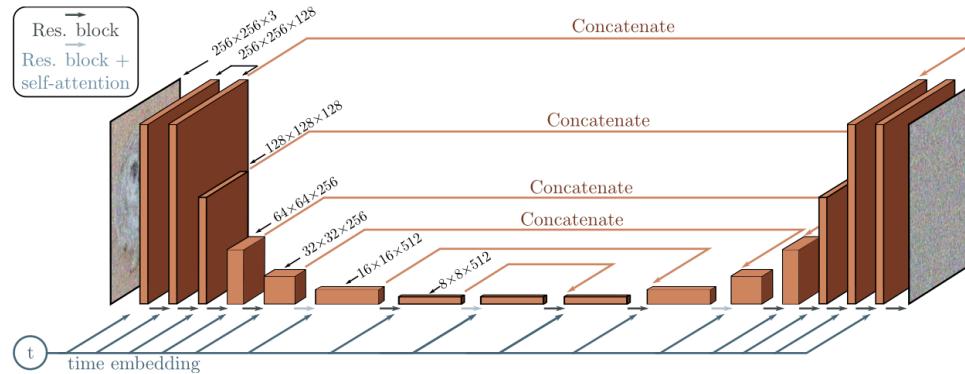
$$z_{T-1} \rightarrow \square + \epsilon$$

$$z_{T-1} - \epsilon = z_{T-2}$$



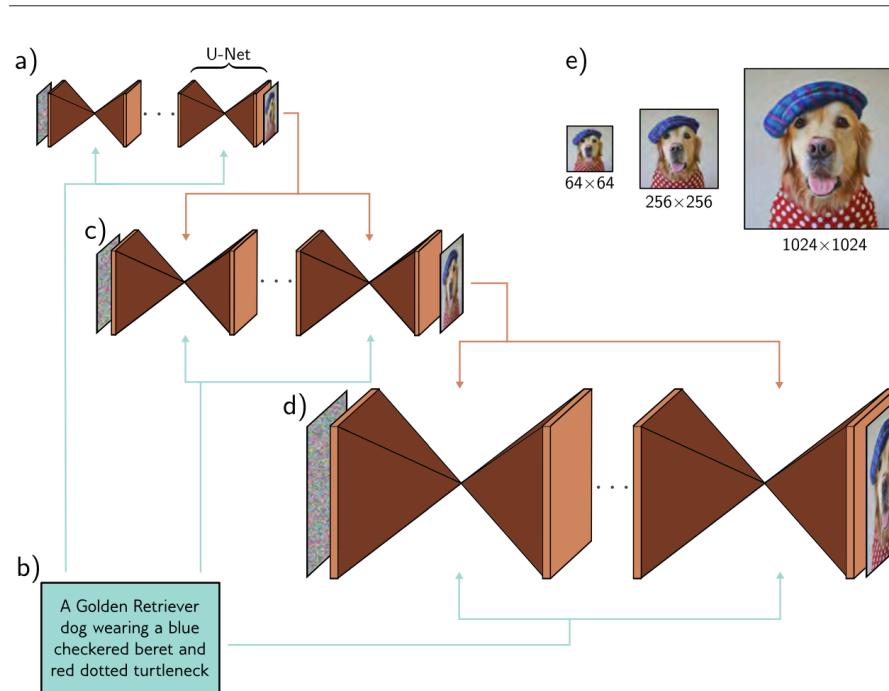


# Network Architecture



- U-Net as used in diffusion models for images. The network aims to predict the noise that was added to the image.
- It consists of an encoder which reduces the scale and increases the number of channels and a decoder which increases the scale and reduces the number of channels.
- The encoder representations are concatenated to their partner in the decoder. Connections between adjacent representations consist of residual blocks, and periodic global self-attention in which every spatial position interacts with every other spatial position.
- A single network is used for all time steps, by passing a sinusoidal time embedding through a shallow neural network and adding the result to the channels at every spatial position at every stage of the U-Net.

# Network Architecture-Conditional Generation



Cascaded conditional generation based on a text prompt.

- A diffusion model consisting of a series of U-Nets is used to generate a 64×64 image.
- This generation is conditioned on a sentence embedding computed by a language model.
- A higher resolution 256×256 image is generated and conditioned on the smaller image and the text encoding.
- This is repeated to create a 1024×1024 image. e) Final image sequence. Adapted from Saharia et al. (2022b).

# Generated outputs - Conditional Generation

