

Denoising Diffusion Probabilistic Model

Shan-Hung Wu & DataLab
Fall 2025

In this lab, we introduce the [Denoising Diffusion Probabilistic Models](#) (DDPMs), which is one of the implementations to demonstrate the usage of the diffusion model generating images, and as an example, dataset [Oxford Flowers 102](#) are used for training the model to generate the images of flower.

Introduction to the diffusion model

Generative modeling experienced tremendous growth in the last five years. Models like VAEs, GANs, and flow-based models proved to be a great success in generating high-quality content, especially images. Diffusion models are a new type of generative model that has proven to be better than previous approaches.

Diffusion models are inspired by non-equilibrium thermodynamics, and they learn to generate by denoising. Learning by denoising consists of two processes, each of which is a Markov Chain. These are:

1. The forward process: In the forward process, we slowly add random noise to the data

in a series of time steps (t_1, t_2, \dots, t_n) . Samples at the current time step are drawn from a Gaussian distribution where the mean of the distribution is conditioned on the sample at the previous time step, and the variance of the distribution follows a fixed schedule. At the end of the forward process, the samples end up with a pure noise distribution.

2. The reverse process: During the reverse process, we try to undo the added noise at

every time step. We start with the pure noise distribution (the last step of the forward process) and try to denoise the samples in the backward direction $(t_n, t_{n-1}, \dots, t_1)$.

We implement the [Denoising Diffusion Probabilistic Models](#) paper or DDPMs for short in this code example. It was the first paper demonstrating the use of diffusion models for generating high-quality images. The authors proved that a certain parameterization of diffusion models reveals an equivalence with denoising score matching over multiple noise levels during training and with annealed Langevin dynamics during sampling that generates the best quality results.

This paper replicates both the Markov chains (forward process and reverse process) involved in the diffusion process but for images. The forward process is fixed and gradually adds Gaussian noise to the images according to a fixed variance schedule denoted by beta in the paper. This is what the diffusion process looks like in case of images: (image -> noise::noise -> image)



The paper describes two algorithms, one for training the model, and the other for sampling from the trained model. Training is performed by optimizing the usual variational bound on negative log-likelihood. The objective function is further simplified, and the network is treated as a noise prediction network. Once optimized, we can sample from the network to generate new images from noise samples. Here is an overview of both algorithms as presented in the paper:

Algorithm 1 Training	Algorithm 2 Sampling
<pre> 1: repeat 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 3: $t \sim \text{Uniform}(\{1, \dots, T\})$ 4: $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 5: Take gradient descent step on $\nabla_{\theta} \ \epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\ ^2$ 6: until converged </pre>	<pre> 1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 2: for $t = T, \dots, 1$ do 3: $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ if $t > 1$, else $\mathbf{z} = \mathbf{0}$ 4: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\bar{\alpha}_t}} (\mathbf{x}_t - \frac{1 - \bar{\alpha}_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t)) + \sigma_t \mathbf{z}$ 5: end for 6: return \mathbf{x}_0 </pre>

Note: DDPM is just one way of implementing a diffusion model. Also, the sampling algorithm in the DDPM replicates the complete Markov chain. Hence, it's slow in generating new samples compared to other generative models like GANs. Lots of research efforts have been made to address this issue. One such example is Denoising Diffusion Implicit Models, or DDIM for short, where the authors replaced the Markov chain with a non-Markovian process to sample faster. You can find the code example for DDIM [here](#)

Implementing a DDPM model is simple. We define a model that takes two inputs: Images and the randomly sampled time steps. At each training step, we perform the following operations to train our model:

1. Sample random noise to be added to the inputs.
2. Apply the forward process to diffuse the inputs with the sampled noise.
3. Your model takes these noisy samples as inputs and outputs the noise prediction for each time step.
4. Given true noise and predicted noise, we calculate the loss values
5. We then calculate the gradients and update the model weights.

Given that our model knows how to denoise a noisy sample at a given time step, we can leverage this idea to generate new samples, starting from a pure noise distribution.

Setup

👉 Please check on the [TensorFlow Addons \(TFA\) Installation](#) for compatibility with your version before installing the tensorflow_addons.

```
In [ ]: import math
import numpy as np
import matplotlib.pyplot as plt

import tensorflow as tf
from tensorflow import keras
from tensorflow.keras import layers
import tensorflow_addons as tfa
import tensorflow_datasets as tfds # If an issue occurs during import, you can try running pip install t;
```

```
c:\Users\USER\anaconda3\envs\LAB13-3-ENV\lib\site-packages\scipy\_init__.py:146: UserWarning: A NumPy version >=1.16.5 and <1.23.0 is required for this version of SciPy (detected version 1.23.0)
  warnings.warn(f"A NumPy version >={np_minversion} and <{np_maxversion}""
c:\Users\USER\anaconda3\envs\LAB13-3-ENV\lib\site-packages\tensorflow_addons\utils\tfa_eol_msg.py:23: UserWarning:
  TensorFlow Addons (TFA) has ended development and introduction of new features.
  TFA has entered a minimal maintenance and release mode until a planned end of life in May 2024.
  Please modify downstream libraries to take dependencies from other repositories in our TensorFlow community (e.g. Keras, Keras-CV, and Keras-NLP).

For more information see: https://github.com/tensorflow/addons/issues/2807
```

```
  warnings.warn(
c:\Users\USER\anaconda3\envs\LAB13-3-ENV\lib\site-packages\tqdm\auto.py:21: TqdmWarning: IPProgress not found. Please update jupyter and ipywidgets. See https://ipywidgets.readthedocs.io/en/stable/user\_install.html
  from .autonotebook import tqdm as notebook_tqdm
```

```
In [2]: gpus = tf.config.experimental.list_physical_devices('GPU')
if gpus:
    try:
        # Restrict TensorFlow to only use the first GPU
        tf.config.experimental.set_visible_devices(gpus[0], 'GPU')

        # Currently, memory growth needs to be the same across GPUs
        for gpu in gpus:
            tf.config.experimental.set_memory_growth(gpu, True)
        logical_gpus = tf.config.experimental.list_logical_devices('GPU')
        print(len(gpus), "Physical GPUs,", len(logical_gpus), "Logical GPUs")
    except RuntimeError as e:
        # Memory growth must be set before GPUs have been initialized
        print(e)
```

1 Physical GPUs, 1 Logical GPUs

Hyperparameters

```
In [3]: batch_size = 32
num_epochs = 1 # Just for the sake of demonstration
total_timesteps = 1000
norm_groups = 8 # Number of groups used in GroupNormalization layer
learning_rate = 2e-4

img_size = 64
img_channels = 3
clip_min = -1.0
clip_max = 1.0

first_conv_channels = 64
channel_multiplier = [1, 2, 4, 8]
widths = [first_conv_channels * mult for mult in channel_multiplier]
has_attention = [False, False, True, True]
num_res_blocks = 2 # Number of residual blocks

dataset_name = "oxford_flowers102"
splits = ["train"]
```

Dataset

We use the [Oxford Flowers 102](#) dataset for generating images of flowers. In terms of preprocessing, we use center cropping for resizing the images to the desired image size, and we rescale the pixel values in the range [-1.0, 1.0]. This is in line with the range of the pixel values that was applied by the authors of the [DDPMs paper](#). For augmenting training data, we randomly flip the images left/right.

```
In [4]: # Load the dataset
(ds,) = tfds.load(dataset_name, split=splits, with_info=False, shuffle_files=True)

def augment(img):
    """Flips an image left/right randomly."""
    return tf.image.random_flip_left_right(img)
```

```

def resize_and_rescale(img, size):
    """Resize the image to the desired size first and then
    rescale the pixel values in the range [-1.0, 1.0].

    Args:
        img: Image tensor
        size: Desired image size for resizing
    Returns:
        Resized and rescaled image tensor
    """

    height = tf.shape(img)[0]
    width = tf.shape(img)[1]
    crop_size = tf.minimum(height, width)

    img = tf.image.crop_to_bounding_box(
        img,
        (height - crop_size) // 2,
        (width - crop_size) // 2,
        crop_size,
        crop_size,
    )

    # Resize
    img = tf.cast(img, dtype=tf.float32)
    img = tf.image.resize(img, size=size, antialias=True)

    # Rescale the pixel values
    img = img / 127.5 - 1.0
    img = tf.clip_by_value(img, clip_min, clip_max)
    return img

def train_preprocessing(x):
    img = x["image"]
    img = resize_and_rescale(img, size=(img_size, img_size))
    img = augment(img)
    return img

train_ds = (
    ds.map(train_preprocessing, num_parallel_calls=tf.data.AUTOTUNE)
    .batch(batch_size, drop_remainder=True)
    .shuffle(batch_size * 2)
    .prefetch(tf.data.AUTOTUNE)
)

```

Gaussian diffusion utilities

We define the forward process and the reverse process as a separate utility. Most of the code in this utility has been borrowed from the original implementation with some slight modifications.

```

In [5]: class GaussianDiffusion:
    """Gaussian diffusion utility.

    Args:
        beta_start: Start value of the scheduled variance
        beta_end: End value of the scheduled variance
        timesteps: Number of time steps in the forward process
    """

    def __init__(
        self,
        beta_start=1e-4,
        beta_end=0.02,
        timesteps=1000,
        clip_min=-1.0,
        clip_max=1.0,
    ):
        self.beta_start = beta_start
        self.beta_end = beta_end

```

```

        self.timesteps = timesteps
        self.clip_min = clip_min
        self.clip_max = clip_max

        # Define the linear variance schedule
        self.betas = betas = np.linspace(
            beta_start,
            beta_end,
            timesteps,
            dtype=np.float64, # Using float64 for better precision
        )
        self.num_timesteps = int(timesteps)

        alphas = 1.0 - betas
        alphas_cumprod = np.cumprod(alphas, axis=0)
        alphas_cumprod_prev = np.append(1.0, alphas_cumprod[:-1])

        self.betas = tf.constant(betas, dtype=tf.float32)
        self.alphas_cumprod = tf.constant(alphas_cumprod, dtype=tf.float32)
        self.alphas_cumprod_prev = tf.constant(alphas_cumprod_prev, dtype=tf.float32)

        # Calculations for diffusion q(x_t | x_{t-1}) and others
        self.sqrt_alphas_cumprod = tf.constant(
            np.sqrt(alphas_cumprod), dtype=tf.float32
        )

        self.sqrt_one_minus_alphas_cumprod = tf.constant(
            np.sqrt(1.0 - alphas_cumprod), dtype=tf.float32
        )

        self.log_one_minus_alphas_cumprod = tf.constant(
            np.log(1.0 - alphas_cumprod), dtype=tf.float32
        )

        self.sqrt_recip_alphas_cumprod = tf.constant(
            np.sqrt(1.0 / alphas_cumprod), dtype=tf.float32
        )
        self.sqrt_recipm1_alphas_cumprod = tf.constant(
            np.sqrt(1.0 / alphas_cumprod - 1), dtype=tf.float32
        )

        # Calculations for posterior q(x_{t-1} | x_t, x_0)
        posterior_variance = (
            betas * (1.0 - alphas_cumprod_prev) / (1.0 - alphas_cumprod)
        )
        self.posterior_variance = tf.constant(posterior_variance, dtype=tf.float32)

        # Log calculation clipped because the posterior variance is 0 at the beginning
        # of the diffusion chain
        self.posterior_log_variance_clipped = tf.constant(
            np.log(np.maximum(posterior_variance, 1e-20)), dtype=tf.float32
        )

        self.posterior_mean_coef1 = tf.constant(
            betas * np.sqrt(alphas_cumprod_prev) / (1.0 - alphas_cumprod),
            dtype=tf.float32,
        )

        self.posterior_mean_coef2 = tf.constant(
            (1.0 - alphas_cumprod_prev) * np.sqrt(alphas) / (1.0 - alphas_cumprod),
            dtype=tf.float32,
        )

    def _extract(self, a, t, x_shape):
        """Extract some coefficients at specified timesteps,
        then reshape to [batch_size, 1, 1, 1, 1, ...] for broadcasting purposes.

        Args:
            a: Tensor to extract from
            t: Timestep for which the coefficients are to be extracted
            x_shape: Shape of the current batched samples
        """
        batch_size = x_shape[0]
        out = tf.gather(a, t)
        return tf.reshape(out, [batch_size, 1, 1, 1])

```

```

def q_mean_variance(self, x_start, t):
    """Extracts the mean, and the variance at current timestep.

    Args:
        x_start: Initial sample (before the first diffusion step)
        t: Current timestep
    """
    x_start_shape = tf.shape(x_start)
    mean = self._extract(self.sqrt_alphas_cumprod, t, x_start_shape) * x_start
    variance = self._extract(1.0 - self.alphas_cumprod, t, x_start_shape)
    log_variance = self._extract(
        self.log_one_minus_alphas_cumprod, t, x_start_shape
    )
    return mean, variance, log_variance

def q_sample(self, x_start, t, noise):
    """Diffuse the data.

    Args:
        x_start: Initial sample (before the first diffusion step)
        t: Current timestep
        noise: Gaussian noise to be added at the current timestep
    Returns:
        Diffused samples at timestep `t`
    """
    x_start_shape = tf.shape(x_start)
    return (
        self._extract(self.sqrt_alphas_cumprod, t, tf.shape(x_start)) * x_start
        + self._extract(self.sqrt_one_minus_alphas_cumprod, t, x_start_shape)
        * noise
    )

def predict_start_from_noise(self, x_t, t, noise):
    x_t_shape = tf.shape(x_t)
    return (
        self._extract(self.sqrt_recip_alphas_cumprod, t, x_t_shape) * x_t
        - self._extract(self.sqrt_recipm1_alphas_cumprod, t, x_t_shape) * noise
    )

def q_posterior(self, x_start, x_t, t):
    """Compute the mean and variance of the diffusion posterior q(x_{t-1} | x_t, x_0).

    Args:
        x_start: Starting point(sample) for the posterior computation
        x_t: Sample at timestep `t`
        t: Current timestep
    Returns:
        Posterior mean and variance at current timestep
    """

    x_t_shape = tf.shape(x_t)
    posterior_mean = (
        self._extract(self.posterior_mean_coef1, t, x_t_shape) * x_start
        + self._extract(self.posterior_mean_coef2, t, x_t_shape) * x_t
    )
    posterior_variance = self._extract(self.posterior_variance, t, x_t_shape)
    posterior_log_variance_clipped = self._extract(
        self.posterior_log_variance_clipped, t, x_t_shape
    )
    return posterior_mean, posterior_variance, posterior_log_variance_clipped

def p_mean_variance(self, pred_noise, x, t, clip_denoised=True):
    x_recon = self.predict_start_from_noise(x, t=t, noise=pred_noise)
    if clip_denoised:
        x_recon = tf.clip_by_value(x_recon, self.clip_min, self.clip_max)

    model_mean, posterior_variance, posterior_log_variance = self.q_posterior(
        x_start=x_recon, x_t=x, t=t
    )
    return model_mean, posterior_variance, posterior_log_variance

def p_sample(self, pred_noise, x, t, clip_denoised=True):
    """Sample from the diffusion model.
    """

```

```

Args:
    pred_noise: Noise predicted by the diffusion model
    x: Samples at a given timestep for which the noise was predicted
    t: Current timestep
    clip_denoised (bool): Whether to clip the predicted noise
        within the specified range or not.
"""
model_mean, _, model_log_variance = self.p_mean_variance(
    pred_noise, x=x, t=t, clip_denoised=clip_denoised
)
noise = tf.random.normal(shape=x.shape, dtype=x.dtype)
# No noise when t == 0
nonzero_mask = tf.reshape(
    1 - tf.cast(tf.equal(t, 0), tf.float32), [tf.shape(x)[0], 1, 1, 1]
)
return model_mean + nonzero_mask * tf.exp(0.5 * model_log_variance) * noise

```

Network architecture

U-Net, originally developed for semantic segmentation, is an architecture that is widely used for implementing diffusion models but with some slight modifications:

1. The network accepts two inputs: Image and time step
2. Self-attention between the convolution blocks once we reach a specific resolution

(16x16 in the paper) 3. Group Normalization instead of weight normalization

We implement most of the things as used in the original paper. We use the `swish` activation function throughout the network. We use the variance scaling kernel initializer.

The only difference here is the number of groups used for the `GroupNormalization` layer. For the flowers dataset, we found that a value of `groups=8` produces better results compared to the default value of `groups=32`. Dropout is optional and should be used where chances of over fitting is high. In the paper, the authors used dropout only when training on CIFAR10.

```

In [6]: # Kernel initializer to use
def kernel_init(scale):
    scale = max(scale, 1e-10)
    return keras.initializers.VarianceScaling(
        scale, mode="fan_avg", distribution="uniform"
    )

class AttentionBlock(layers.Layer):
    """Applies self-attention.

    Args:
        units: Number of units in the dense layers
        groups: Number of groups to be used for GroupNormalization layer
    """

    def __init__(self, units, groups=8, **kwargs):
        self.units = units
        self.groups = groups
        super().__init__(**kwargs)

        self.norm = tfa.layers.GroupNormalization(groups=groups)
        self.query = layers.Dense(units, kernel_initializer=kernel_init(1.0))
        self.key = layers.Dense(units, kernel_initializer=kernel_init(1.0))
        self.value = layers.Dense(units, kernel_initializer=kernel_init(1.0))
        self.proj = layers.Dense(units, kernel_initializer=kernel_init(0.0))

    def call(self, inputs):
        batch_size = tf.shape(inputs)[0]
        height = tf.shape(inputs)[1]
        width = tf.shape(inputs)[2]
        scale = tf.cast(self.units, tf.float32) ** (-0.5)

        inputs = self.norm(inputs)
        q = self.query(inputs)

```

```

k = self.key(inputs)
v = self.value(inputs)

attn_score = tf.einsum("bhwC, bHWc->bhwHW", q, k) * scale
attn_score = tf.reshape(attn_score, [batch_size, height, width, height * width])

attn_score = tf.nn.softmax(attn_score, -1)
attn_score = tf.reshape(attn_score, [batch_size, height, width, height, width])

proj = tf.einsum("bhwHW,bHWc->bhwC", attn_score, v)
proj = self.proj(proj)
return inputs + proj

class TimeEmbedding(layers.Layer):
    def __init__(self, dim, **kwargs):
        super().__init__(**kwargs)
        self.dim = dim
        self.half_dim = dim // 2
        self.emb = math.log(10000) / (self.half_dim - 1)
        self.emb = tf.exp(tf.range(self.half_dim, dtype=tf.float32) * -self.emb)

    def call(self, inputs):
        inputs = tf.cast(inputs, dtype=tf.float32)
        emb = inputs[:, None] * self.emb[None, :]
        emb = tf.concat([tf.sin(emb), tf.cos(emb)], axis=-1)
        return emb

def ResidualBlock(width, groups=8, activation_fn=keras.activations.swish):
    def apply(inputs):
        x, t = inputs
        input_width = x.shape[3]

        if input_width == width:
            residual = x
        else:
            residual = layers.Conv2D(
                width, kernel_size=1, kernel_initializer=kernel_init(1.0)
            )(x)

        temb = activation_fn(t)
        temb = layers.Dense(width, kernel_initializer=kernel_init(1.0))(temb)[
            :, None, None, :
        ]

        x = tfa.layers.GroupNormalization(groups=groups)(x)
        x = activation_fn(x)
        x = layers.Conv2D(
            width, kernel_size=3, padding="same", kernel_initializer=kernel_init(1.0)
        )(x)

        x = layers.Add()([x, temb])
        x = tfa.layers.GroupNormalization(groups=groups)(x)
        x = activation_fn(x)

        x = layers.Conv2D(
            width, kernel_size=3, padding="same", kernel_initializer=kernel_init(0.0)
        )(x)
        x = layers.Add()([x, residual])
        return x

    return apply

def DownSample(width):
    def apply(x):
        x = layers.Conv2D(
            width,
            kernel_size=3,
            strides=2,
            padding="same",
            kernel_initializer=kernel_init(1.0),
        )(x)
        return x

```

```

    return apply

def UpSample(width, interpolation="nearest"):
    def apply(x):
        x = layers.UpSampling2D(size=2, interpolation=interpolation)(x)
        x = layers.Conv2D(
            width, kernel_size=3, padding="same", kernel_initializer=kernel_init(1.0)
        )(x)
        return x

    return apply

def TimeMLP(units, activation_fn=keras.activations.swish):
    def apply(inputs):
        temb = layers.Dense(
            units, activation=activation_fn, kernel_initializer=kernel_init(1.0)
        )(inputs)
        temb = layers.Dense(units, kernel_initializer=kernel_init(1.0))(temb)
        return temb

    return apply

def build_model(
    img_size,
    img_channels,
    widths,
    has_attention,
    num_res_blocks=2,
    norm_groups=8,
    interpolation="nearest",
    activation_fn=keras.activations.swish,
):
    image_input = layers.Input(
        shape=(img_size, img_size, img_channels), name="image_input"
    )
    time_input = keras.Input(shape=(), dtype=tf.int64, name="time_input")

    x = layers.Conv2D(
        first_conv_channels,
        kernel_size=(3, 3),
        padding="same",
        kernel_initializer=kernel_init(1.0),
    )(image_input)

    temb = TimeEmbedding(dim=first_conv_channels * 4)(time_input)
    temb = TimeMLP(units=first_conv_channels * 4, activation_fn=activation_fn)(temb)

    skips = [x]

    # DownBlock
    for i in range(len(widths)):
        for _ in range(num_res_blocks):
            x = ResidualBlock(
                widths[i], groups=norm_groups, activation_fn=activation_fn
            )([x, temb])
            if has_attention[i]:
                x = AttentionBlock(widths[i], groups=norm_groups)(x)
            skips.append(x)

        if widths[i] != widths[-1]:
            x = DownSample(widths[i])(x)
            skips.append(x)

    # MiddleBlock
    x = ResidualBlock(widths[-1], groups=norm_groups, activation_fn=activation_fn)(
        [x, temb]
    )
    x = AttentionBlock(widths[-1], groups=norm_groups)(x)
    x = ResidualBlock(widths[-1], groups=norm_groups, activation_fn=activation_fn)(
        [x, temb]
    )

```

```
# UpBlock
for i in reversed(range(len(widths))):
    for _ in range(num_res_blocks + 1):
        x = layers.concatenate(axis=-1)([x, skips.pop()])
        x = ResidualBlock(
            widths[i], groups=norm_groups, activation_fn=activation_fn
        )(x, temb)
        if has_attention[i]:
            x = AttentionBlock(widths[i], groups=norm_groups)(x)

    if i != 0:
        x = UpSample(widths[i], interpolation=interpolation)(x)

# End block
x = tfa.layers.GroupNormalization(groups=norm_groups)(x)
x = activation_fn(x)
x = layers.Conv2D(3, (3, 3), padding="same", kernel_initializer=kernel_init(0.0))(x)
return keras.Model([image_input, time_input], x, name="unet")
```

Training

We follow the same setup for training the diffusion model as described in the paper. We use Adam optimizer with a learning rate of `2e-4`. We use EMA on model parameters with a decay factor of 0.999. We treat our model as noise prediction network i.e. at every training step, we input a batch of images and corresponding time steps to our UNet, and the network outputs the noise as predictions.

The only difference is that we aren't using the Kernel Inception Distance (KID) or Frechet Inception Distance (FID) for evaluating the quality of generated samples during training. This is because both these metrics are compute heavy and are skipped for the brevity of implementation.

Note: We are using mean squared error as the loss function which is aligned with the paper, and theoretically makes sense. In practice, though, it is also common to use mean absolute error or Huber loss as the loss function.

```
In [7]: class DiffusionModel(keras.Model):
    def __init__(self, network, ema_network, timesteps, gdf_util, ema=0.999):
        super().__init__()
        self.network = network
        self.ema_network = ema_network
        self.timesteps = timesteps
        self.gdf_util = gdf_util
        self.ema = ema

    def train_step(self, images):
        # 1. Get the batch size
        batch_size = tf.shape(images)[0]

        # 2. Sample timesteps uniformly
        t = tf.random.uniform(
            minval=0, maxval=self.timesteps, shape=(batch_size,), dtype=tf.int64
        )

        with tf.GradientTape() as tape:
            # 3. Sample random noise to be added to the images in the batch
            noise = tf.random.normal(shape=tf.shape(images), dtype=images.dtype)

            # 4. Diffuse the images with noise
            images_t = self.gdf_util.q_sample(images, t, noise)

            # 5. Pass the diffused images and time steps to the network
            pred_noise = self.network([images_t, t], training=True)

            # 6. Calculate the loss
            loss = self.loss(noise, pred_noise)

            # 7. Get the gradients
            gradients = tape.gradient(loss, self.network.trainable_weights)

            # 8. Update the weights of the network
            self.optimizer.apply_gradients(zip(gradients, self.network.trainable_weights))
```

```

# 9. Updates the weight values for the network with EMA weights
for weight, ema_weight in zip(self.network.weights, self.ema_network.weights):
    ema_weight.assign(self.ema * ema_weight + (1 - self.ema) * weight)

# 10. Return loss values
return {"loss": loss}

def generate_images(self, num_images=16):
    # 1. Randomly sample noise (starting point for reverse process)
    samples = tf.random.normal(
        shape=(num_images, img_size, img_size, img_channels), dtype=tf.float32
    )
    # 2. Sample from the model iteratively
    for t in reversed(range(0, self.timesteps)):
        tt = tf.cast(tf.fill(num_images, t), dtype=tf.int64)
        pred_noise = self.ema_network.predict(
            [samples, tt], verbose=0, batch_size=num_images
        )
        samples = self.gdf_util.p_sample(
            pred_noise, samples, tt, clip_denoised=True
        )
    # 3. Return generated samples
    return samples

def plot_images(
    self, epoch=None, logs=None, num_rows=2, num_cols=8, figsize=(12, 5)
):
    """Utility to plot images using the diffusion model during training."""
    generated_samples = self.generate_images(num_images=num_rows * num_cols)
    generated_samples = (
        tf.clip_by_value(generated_samples * 127.5 + 127.5, 0.0, 255.0)
        .numpy()
        .astype(np.uint8)
    )

    _, ax = plt.subplots(num_rows, num_cols, figsize=figsize)
    for i, image in enumerate(generated_samples):
        if num_rows == 1:
            ax[i].imshow(image)
            ax[i].axis("off")
        else:
            ax[i // num_cols, i % num_cols].imshow(image)
            ax[i // num_cols, i % num_cols].axis("off")

    plt.tight_layout()
    plt.show()

# Build the unet model
network = build_model(
    img_size=img_size,
    img_channels=img_channels,
    widths=widths,
    has_attention=has_attention,
    num_res_blocks=num_res_blocks,
    norm_groups=norm_groups,
    activation_fn=keras.activations.swish,
)
ema_network = build_model(
    img_size=img_size,
    img_channels=img_channels,
    widths=widths,
    has_attention=has_attention,
    num_res_blocks=num_res_blocks,
    norm_groups=norm_groups,
    activation_fn=keras.activations.swish,
)
ema_network.set_weights(network.get_weights()) # Initially the weights are the same

# Get an instance of the Gaussian Diffusion utilities
gdf_util = GaussianDiffusion(timesteps=total_timesteps)

# Get the model
model = DiffusionModel(

```

```

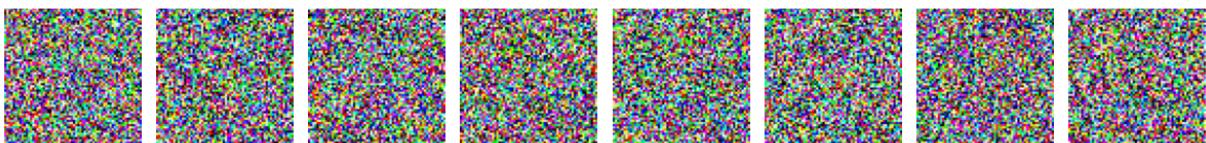
        network=network,
        ema_network=ema_network,
        gdf_util=gdf_util,
        timesteps=total_timesteps,
    )

    # Compile the model
    model.compile(
        loss=keras.losses.MeanSquaredError(),
        optimizer=keras.optimizers.Adam(learning_rate=learning_rate),
    )

    # Train the model
    model.fit(
        train_ds,
        epochs=num_epochs,
        batch_size=batch_size,
        callbacks=[keras.callbacks.LambdaCallback(on_epoch_end=model.plot_images)],
    )

```

31/31 [=====] - ETA: 0s - loss: 0.7792



31/31 [=====] - 204s 6s/step - loss: 0.7711

Out[7]: <keras.callbacks.History at 0x233b7724250>

Results

We trained this model for 800 epochs on a V100 GPU, and each epoch took almost 8 seconds to finish. We load those weights here, and we generate a few samples starting from pure noise.

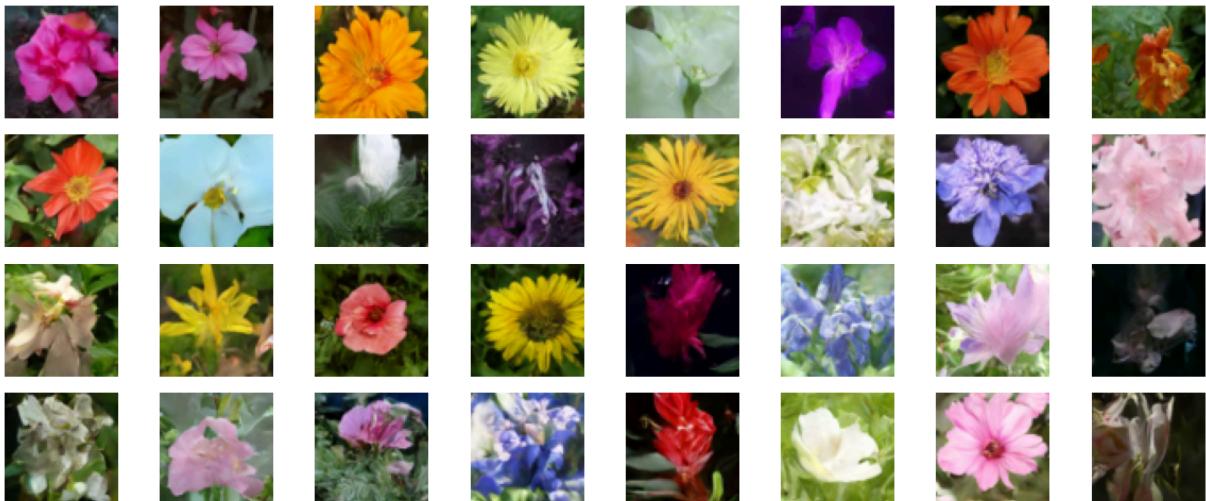
Please extract [checkpoint.zip](#) in the current directory of the code.

```

In [8]: # Load the model weights
model.ema_network.load_weights("./checkpoints/diffusion_model_checkpoint")

# Generate and plot some samples
model.plot_images(num_rows=4, num_cols=8)

```



It takes approximately 3 to 5 minutes to generate flower images.

References

You can find the original introduction by [A_K_Nain here.](#)

1. Denoising Diffusion Probabilistic Models
2. Author's implementation
3. A deep dive into DDPMs
4. Denoising Diffusion Implicit Models
5. Annotated Diffusion Model
6. AIAIART

Assignment

In this assignment, you have to implement the [DDIMs](#).

Denoising Diffusion Implicit Models

The primary differences between DDPM (Denoising Diffusion Probabilistic Models) and DDIM (Denoising Diffusion Implicit Models) lie in their approach, speed, and determinism in the image generation process. DDPM generates images through a gradual, stochastic denoising process, introducing randomness at each step. This approach, while slower, allows for greater sample diversity. In contrast, DDIM employs a deterministic generation process, where the denoising steps are re-designed to achieve faster generation and produce consistent images under the same conditions.

The distinctions between these two models can be summarized as follows:

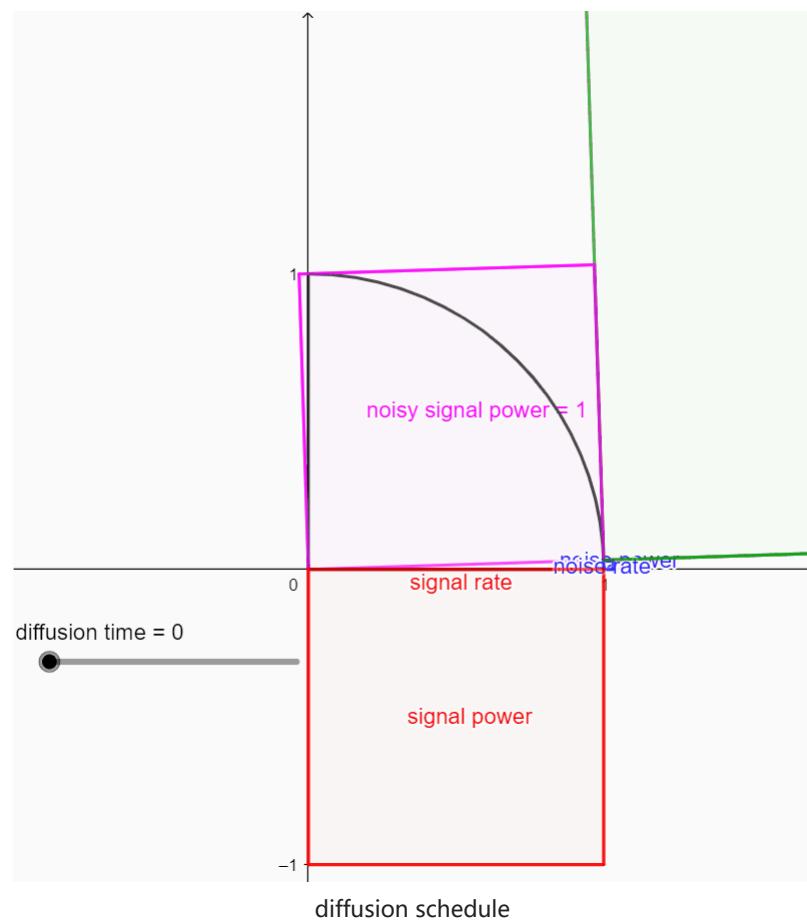
1. **Generation Process:** DDPM generates images through stochastic denoising, gradually removing random noise to yield a final image. Because randomness is introduced at each step, the generated image may vary with each run. DDIM, on the other hand, is deterministic; by adjusting the denoising formula, it can generate images in fewer steps, and identical images can be generated under the same conditions.
2. **Sampling Speed:** Due to its deterministic denoising process, DDIM requires fewer steps, resulting in faster generation. This efficient sampling is especially beneficial in applications where computational resources need to be conserved.
3. **Determinism of the Denoising Process:** The denoising process in DDIM is deterministic, providing better control over the generated images. In contrast, DDPM's stochastic denoising introduces slight variations in the images produced under the same conditions.

In the DDIM diffusion process, we need a function that provides the noise and signal levels at each time step, known as the `diffusion schedule`.

The diffusion schedule provides two key values: `noise_rate` and `signal_rate`, which correspond to $\sqrt{1 - \alpha}$ and $\sqrt{\alpha}$ in the DDIM paper, respectively. To generate a noisy image, we combine random noise and the original image according to these rates.

Since both the random noise and normalized images have a mean of zero and a variance of one, `noise_rate` and `signal_rate` represent the standard deviations of these components in the noisy image, with their squares representing variance (or power in signal processing terms). We set these rates so that their squared sum is always 1, ensuring the noisy image maintains unit variance, consistent with its unscaled components.

In this implementation, we use a continuous version of the cosine schedule, which is widely used in research. This schedule is symmetric, slowing down at the start and end of the diffusion process, and has a clear geometric interpretation based on the trigonometric properties of the unit circle.



Requirements

- Complete the TODO sections in the following code, including `diffusion_schedule`, `denoise`, and `reverse_diffusion`.
- Briefly summarize what you did and explain the performance results.
- This assignment does not specify a clear model loss threshold. However, please train your model to the extent that it can generate images with clearly recognizable flowers. For example:

No description has been provided for this image

Setup

```
In [14]: import math
import numpy as np
import matplotlib.pyplot as plt
import tensorflow as tf

from tensorflow import keras
from tensorflow.keras import layers

import tensorflow_addons as tfa
import tensorflow_datasets as tfds
```

Hyperparameters

```
In [15]: # Feel free to change these parameters according to your system's configuration
# data
dataset_name = "oxford_flowers102"
dataset_repetitions = 1
num_epochs = 10
image_size = 64
shuffle_times = 10

# KID = Kernel Inception Distance, see related section
kid_image_size = 75
```

```

kid_diffusion_steps = 5
plot_diffusion_steps = 20

# sampling
min_signal_rate = 0.02
max_signal_rate = 0.95

# architecture
embedding_dims = 32
embedding_max_frequency = 1000.0
widths = [32, 64, 96, 128]
block_depth = 2

# optimization
batch_size = 32
ema = 0.999
learning_rate = 1e-3
weight_decay = 1e-4

```

Data pipeline

```

In [16]: def preprocess_image(data):
    # center crop image
    height = tf.shape(data["image"])[0]
    width = tf.shape(data["image"])[1]
    crop_size = tf.minimum(height, width)
    image = tf.image.crop_to_bounding_box(
        data["image"],
        (height - crop_size) // 2,
        (width - crop_size) // 2,
        crop_size,
        crop_size,
    )

    # resize and clip
    # for image downsampling it is important to turn on antialiasing
    image = tf.image.resize(image, size=[image_size, image_size], antialias=True)
    return tf.clip_by_value(image / 255.0, 0.0, 1.0)

def prepare_dataset(split):
    # the validation dataset is shuffled as well, because data order matters
    # for the KID estimation
    return (
        tfds.load(dataset_name, split=split, shuffle_files=True)
        .map(preprocess_image, num_parallel_calls=tf.data.AUTOTUNE)
        .cache()
        .repeat(dataset_repetitions)
        .shuffle(shuffle_times * batch_size)
        .batch(batch_size, drop_remainder=True)
        .prefetch(buffer_size=tf.data.AUTOTUNE)
    )

# Load individual datasets
train_part = prepare_dataset("train")
validation_part = prepare_dataset("validation")
test_part = prepare_dataset("test")

# Concatenate all parts into one training dataset
train_dataset = train_part.concatenate(test_part)
val_dataset = validation_part

```

Kernel inception distance

Kernel Inception Distance (KID) is an image quality metric developed as a simpler and computationally lighter alternative to the popular Frechet Inception Distance (FID).

In this example, the images are evaluated at a lower resolution compatible with the Inception network (75x75 instead of the usual 299x299), and the metric is computed only on the validation set to save on

computational cost. Additionally, we restrict the number of sampling steps during evaluation to 5 for efficiency.

Since the dataset is relatively small, we repeat the training and validation splits multiple times within each epoch. This approach helps to mitigate the noise in KID estimation, which is compute-intensive, by ensuring the metric is evaluated after accumulating a substantial number of iterations, providing a more reliable measurement.

```
In [17]: class KID(keras.metrics.Metric):
    def __init__(self, name, **kwargs):
        super().__init__(name=name, **kwargs)

        # KID is estimated per batch and is averaged across batches
        self.kid_tracker = keras.metrics.Mean(name="kid_tracker")

        # a pretrained InceptionV3 is used without its classification layer
        # transform the pixel values to the 0-255 range, then use the same
        # preprocessing as during pretraining
        self.encoder = keras.Sequential(
            [
                keras.Input(shape=(image_size, image_size, 3)),
                layers.Rescaling(255.0),
                layers.Resizing(height=kid_image_size, width=kid_image_size),
                layers.Lambda(keras.applications.inception_v3.preprocess_input),
                keras.applications.InceptionV3(
                    include_top=False,
                    input_shape=(kid_image_size, kid_image_size, 3),
                    weights="imagenet",
                ),
                layers.GlobalAveragePooling2D(),
            ],
            name="inception_encoder",
        )

    def polynomial_kernel(self, features_1, features_2):
        # Use TensorFlow functions instead of ops
        feature_dimensions = tf.cast(tf.shape(features_1)[1], dtype="float32")
        return (tf.matmul(features_1, tf.transpose(features_2)) / feature_dimensions + 1.0) ** 3.0

    def update_state(self, real_images, generated_images, sample_weight=None):
        real_features = self.encoder(real_images, training=False)
        generated_features = self.encoder(generated_images, training=False)

        # compute polynomial kernels using the two sets of features
        kernel_real = self.polynomial_kernel(real_features, real_features)
        kernel_generated = self.polynomial_kernel(generated_features, generated_features)
        kernel_cross = self.polynomial_kernel(real_features, generated_features)

        # estimate the squared maximum mean discrepancy using the average kernel values
        batch_size = tf.shape(real_features)[0]
        batch_size_f = tf.cast(batch_size, dtype="float32")
        mean_kernel_real = tf.reduce_sum(kernel_real * (1.0 - tf.eye(batch_size))) / (
            batch_size_f * (batch_size_f - 1.0)
        )
        mean_kernel_generated = tf.reduce_sum(kernel_generated * (1.0 - tf.eye(batch_size))) / (
            batch_size_f * (batch_size_f - 1.0)
        )
        mean_kernel_cross = tf.reduce_mean(kernel_cross)
        kid = mean_kernel_real + mean_kernel_generated - 2.0 * mean_kernel_cross

        # update the average KID estimate
        self.kid_tracker.update_state(kid)

    def result(self):
        return self.kid_tracker.result()

    def reset_state(self):
        self.kid_tracker.reset_state()
```

Network architecture

The DDIM network architecture is based on a modified U-Net, tailored specifically for efficient and deterministic image generation through diffusion. Here are some key design choices that differentiate it from other diffusion models like DDPM:

1. **Continuous Time Embedding:** The network includes a sinusoidal time embedding that enables continuous representation of the diffusion time. This continuous approach allows DDIM to modify sampling steps flexibly during inference, a significant departure from the discrete time steps in DDPM.
2. **Deterministic Reverse Diffusion:** Unlike DDPM, which uses random noise in each step, DDIM is designed to follow a deterministic reverse process. This change ensures faster generation speeds and consistent outputs under the same conditions, providing better control over the generated images.
3. **Group Normalization and Swish Activation:** Similar to other diffusion models, DDIM utilizes group normalization for stability and the Swish activation function across layers. These are chosen to enhance the model's capacity to handle complex data while maintaining efficient training.
4. **Reduced Sampling Steps:** By adjusting the denoising steps, DDIM achieves faster image generation, as it can effectively use fewer steps compared to DDPM, making it suitable for applications where speed is essential.

These adjustments allow DDIM to strike a balance between performance and efficiency, making it more suitable for high-speed and deterministic generation applications.

```
In [18]: def sinusoidal_embedding(x):
    embedding_min_frequency = 1.0
    frequencies = tf.math.exp(
        tf.linspace(
            tf.math.log(embedding_min_frequency),
            tf.math.log(embedding_max_frequency),
            embedding_dims // 2,
        )
    )
    angular_speeds = tf.cast(2.0 * math.pi * frequencies, "float32")
    embeddings = tf.concat([
        tf.sin(angular_speeds * x), tf.cos(angular_speeds * x)], axis=3)
    return embeddings

def ResidualBlock(width):
    def apply(x):
        input_width = x.shape[3]
        if input_width == width:
            residual = x
        else:
            residual = layers.Conv2D(width, kernel_size=1)(x)
        x = layers.BatchNormalization(center=False, scale=False)(x)
        x = layers.Conv2D(width, kernel_size=3, padding="same", activation="swish")(x)
        x = layers.Conv2D(width, kernel_size=3, padding="same")(x)
        x = layers.Add()([x, residual])
        return x

    return apply

def DownBlock(width, block_depth):
    def apply(x):
        x, skips = x
        for _ in range(block_depth):
            x = ResidualBlock(width)(x)
            skips.append(x)
        x = layers.AveragePooling2D(pool_size=2)(x)
        return x

    return apply

def UpBlock(width, block_depth):
    def apply(x):
        x, skips = x
```

```

x = layers.UpSampling2D(size=2, interpolation="bilinear")(x)
for _ in range(block_depth):
    x = layers.Concatenate()([x, skips.pop()])
    x = ResidualBlock(width)(x)
return x

return apply

def get_network(image_size, widths, block_depth):
    noisy_images = keras.Input(shape=(image_size, image_size, 3))
    noise_variances = keras.Input(shape=(1, 1, 1))

    e = layers.Lambda(sinusoidal_embedding, output_shape=(1, 1, 32))(noise_variances)
    e = layers.UpSampling2D(size=image_size, interpolation="nearest")(e)

    x = layers.Conv2D(widths[0], kernel_size=1)(noisy_images)
    x = layers.Concatenate()([x, e])

    skips = []
    for width in widths[:-1]:
        x = DownBlock(width, block_depth)([x, skips])

    for _ in range(block_depth):
        x = ResidualBlock(widths[-1])(x)

    for width in reversed(widths[:-1]):
        x = UpBlock(width, block_depth)([x, skips])

    x = layers.Conv2D(3, kernel_size=1, kernel_initializer="zeros")(x)

    return keras.Model([noisy_images, noise_variances], x, name="residual_unet")

```

Diffusion model

```

In [39]: class DiffusionModel(keras.Model):
    def __init__(self, image_size, widths, block_depth):
        super().__init__()

        self.normalizer = layers.Normalization()
        self.network = get_network(image_size, widths, block_depth)
        self.ema_network = keras.models.clone_model(self.network)

    def compile(self, **kwargs):
        super().compile(**kwargs)

        self.noise_loss_tracker = keras.metrics.Mean(name="n_loss")
        self.image_loss_tracker = keras.metrics.Mean(name="i_loss")
        self.kid = KID(name="kid")

    @property
    def metrics(self):
        return [self.noise_loss_tracker, self.image_loss_tracker, self.kid]

    def denormalize(self, images):
        # convert the pixel values back to 0-1 range
        images = self.normalizer.mean + images * self.normalizer.variance**0.5
        return tf.clip_by_value(images, 0.0, 1.0)

    def diffusion_schedule(self, diffusion_times):
        # TODO: check if the diffusion_times are in the range [0, 1]
        # diffusion_times should be between 0 and 1, corresponding to the start and end of the diffusion

        # angles -> signal and noise rates
        # Use cosine and sine functions to calculate the signal and noise rates based on angles

        noise_rates, signal_rates = None, None #calculate the noise and signal rates
        # note that their squared sum is always: sin^2(x) + cos^2(x) = 1

        return noise_rates, signal_rates

    def denoise(self, noisy_images, noise_rates, signal_rates, training):
        # TODO: implement the denoising network

```

```

# Use the model (self.network or self.ema_network) to predict the noise component in the images

# Predict the noise component and calculate the image component using it
# Here, use the signal_rate and noise_rate to derive the output components

pred_noises, pred_images = None, None #calculate the predicted noises and images
return pred_noises, pred_images

def reverse_diffusion(self, initial_noise, diffusion_steps):
    # reverse diffusion = sampling
    num_images = initial_noise.shape[0]
    step_size = 1.0 / diffusion_steps

    # important line:
    # at the first sampling step, the "noisy image" is pure noise
    # but its signal rate is assumed to be nonzero (min_signal_rate)
    next_noisy_images = initial_noise
    for step in range(diffusion_steps):
        # TODO: implement the reverse diffusion process
        # This process gradually reduces the noise to generate a clearer image

        # remix the predicted components
        # Use the signal_rate and noise_rate from the next step to recombine image and noise components

        pred_noises, pred_images = None, None # network used in eval mode
        next_noisy_images = None # remix the predicted components

    return pred_images

def generate(self, num_images, diffusion_steps):
    # noise -> images -> denormalized images
    initial_noise = tf.random.normal(
        shape=(num_images, image_size, image_size, 3)
    )
    generated_images = self.reverse_diffusion(initial_noise, diffusion_steps)
    generated_images = self.denormalize(generated_images)
    return generated_images

def train_step(self, images):
    # normalize images to have standard deviation of 1, like the noises
    images = self.normalizer(images, training=True)
    noises = tf.random.normal(shape=(batch_size, image_size, image_size, 3))

    # sample uniform random diffusion times
    diffusion_times = tf.random.uniform(
        shape=(batch_size, 1, 1, 1), minval=0.0, maxval=1.0
    )
    noise_rates, signal_rates = self.diffusion_schedule(diffusion_times)
    # mix the images with noises accordingly
    noisy_images = signal_rates * images + noise_rates * noises

    with tf.GradientTape() as tape:
        # train the network to separate noisy images to their components
        pred_noises, pred_images = self.denoise(
            noisy_images, noise_rates, signal_rates, training=True
        )

        noise_loss = self.loss(noises, pred_noises) # used for training
        image_loss = self.loss(images, pred_images) # only used as metric

    gradients = tape.gradient(noise_loss, self.network.trainable_weights)
    self.optimizer.apply_gradients(zip(gradients, self.network.trainable_weights))

    self.noise_loss_tracker.update_state(noise_loss)
    self.image_loss_tracker.update_state(image_loss)

    # track the exponential moving averages of weights
    for weight, ema_weight in zip(self.network.weights, self.ema_network.weights):
        ema_weight.assign(ema * ema_weight + (1 - ema) * weight)

    # KID is not measured during the training phase for computational efficiency
    return {m.name: m.result() for m in self.metrics[:-1]}

def test_step(self, images):
    # normalize images to have standard deviation of 1, like the noises

```

```

        images = self.normalizer(images, training=False)
        noises = tf.random.normal(shape=(batch_size, image_size, image_size, 3))

        # sample uniform random diffusion times
        diffusion_times = tf.random.uniform(
            shape=(batch_size, 1, 1, 1), minval=0.0, maxval=1.0
        )
        noise_rates, signal_rates = self.diffusion_schedule(diffusion_times)
        # mix the images with noises accordingly
        noisy_images = signal_rates * images + noise_rates * noises

        # use the network to separate noisy images to their components
        pred_noises, pred_images = self.denoise(
            noisy_images, noise_rates, signal_rates, training=False
        )

        noise_loss = self.loss(noises, pred_noises)
        image_loss = self.loss(images, pred_images)

        self.image_loss_tracker.update_state(image_loss)
        self.noise_loss_tracker.update_state(noise_loss)

        # measure KID between real and generated images
        # this is computationally demanding, kid_diffusion_steps has to be small
        images = self.denormalize(images)
        generated_images = self.generate(
            num_images=batch_size, diffusion_steps=kid_diffusion_steps
        )
        self.kid.update_state(images, generated_images)

    return {m.name: m.result() for m in self.metrics}

def plot_images(self, epoch=None, logs=None, num_rows=3, num_cols=6):
    # plot random generated images for visual evaluation of generation quality
    generated_images = self.generate(
        num_images=num_rows * num_cols,
        diffusion_steps=plot_diffusion_steps,
    )

    plt.figure(figsize=(num_cols * 2.0, num_rows * 2.0))
    for row in range(num_rows):
        for col in range(num_cols):
            index = row * num_cols + col
            plt.subplot(num_rows, num_cols, index + 1)
            plt.imshow(generated_images[index])
            plt.axis("off")
    plt.tight_layout()
    plt.show()
    plt.close()

```

Training

```
In [34]: # create and compile the model
model = DiffusionModel(image_size, widths, block_depth)
```

```
In [35]: model.compile(
    optimizer=tfa.optimizers.AdamW(
        learning_rate=learning_rate, weight_decay=weight_decay
    ),
    loss=keras.losses.mean_absolute_error,
)
# pixelwise mean absolute error is used as loss
```

```
In [ ]: # save the best model based on the validation KID metric
checkpoint_path = f"checkpoints/DDIM/tf_checkpoint"
checkpoint_callback = keras.callbacks.ModelCheckpoint(
    filepath=checkpoint_path,
    save_weights_only=True,
    monitor="val_kid",
    mode="min",
    save_best_only=True,
)
```

```
In [37]: # calculate mean and variance of training dataset for normalization  
model.normalizer.adapt(train_dataset)
```

```
In [ ]: # run training and plot generated images periodically  
model.fit(  
    train_dataset,  
    validation_data=val_dataset,  
    epochs=num_epochs,  
    callbacks=[  
        keras.callbacks.LambdaCallback(on_epoch_end = model.plot_images ),  
        checkpoint_callback,  
    ],  
)
```

Inference

```
In [38]: # Load the best model and generate images  
model.load_weights(checkpoint_path)  
model.plot_images()
```

