

Deep Learning: Project 4

Florence Hugh & August Jonasson

January 17, 2025

Introduction and dataset

In this project we explore and fit a diffusion model, namely the denoising diffusion probabilistic model (DDPM) presented in C. F. Higham, D. J. Higham, and Grindrod 2023, or equivalently, in Ho, Jain, and Abbeel 2020, for the purpose of generating new (novel) images. The dataset of our choice contains ≈ 200 images of the fictional Pokéémon creature Bulbasaur¹, and it is thus images of Bulbasaur that we seek to generate using the DDPM.

We chose this dataset as it is relatively small (easing the computational load), and because its individual images are of relatively low complexity (a drawing using only a few shapes and colors as opposed to, say, a photo taken of something in the real world), which we hope will make fitting a model easier. We also like Bulbasaurs, of which some examples from the dataset can be seen in Figure 1.

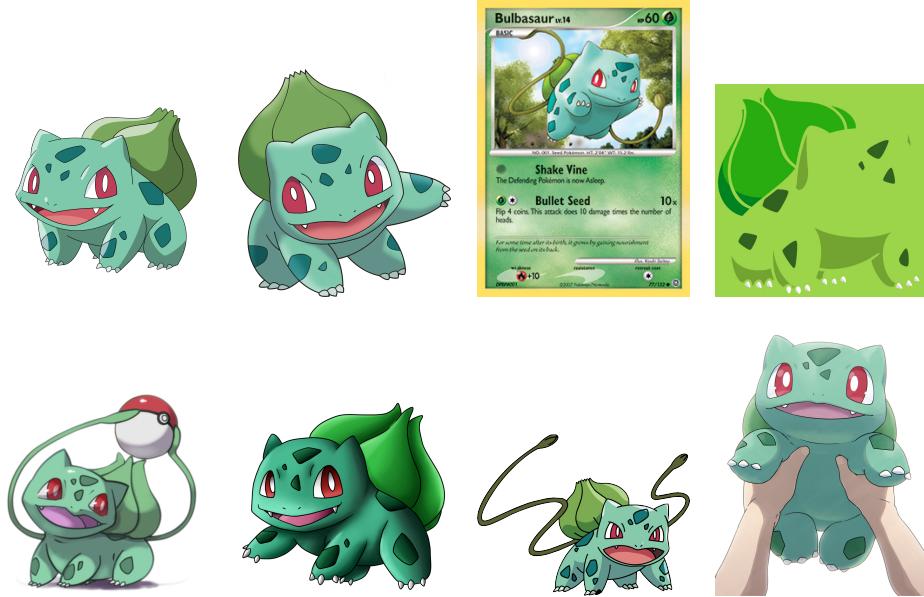


Figure 1: Example samples of images in the dataset.

¹<https://en.wikipedia.org/wiki/Bulbasaur>

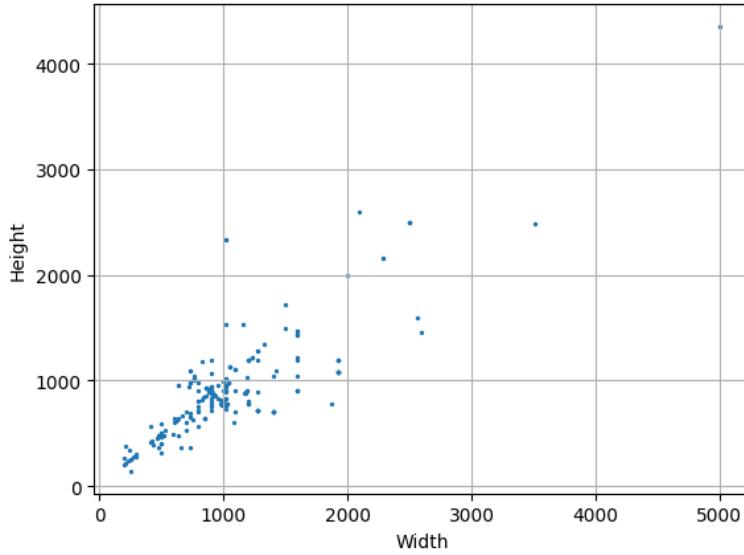


Figure 2: The dimensions of each Bulbasaur image.

Exploratory analysis and data preprocessing

As seen from Figure 2, the resolutions of the individual images within the dataset are not fixed and they also have a varying amount of channels. We thus require some pre-processing of the data before starting to model.

First, we identify along which axis (x or y) that the image is the smallest. We then, as centered as possible, crop out a square of that particular side length using the `image.crop_to_bounding_box`² function, leaving us with a square image. The next step is to resize the image to the desired resolution 64x64 (chosen w.r.t. compute constraints) using *bilinear interpolation*³. The `image.resize` function was used for this. Finally we rescale the pixel values to lie in the $[-1, 1]$ interval using the `image.clip_by_value` function. This is done in order to have the pixel values centered on the same value as the standard normal noise that we are to add in the diffusion process. Finally we have also used the `image.random_flip_left_right` function as a simple form of data augmentation. This is meant to both expand the dataset, but also prevent overfitting in the long run.

²all functions in this section come from tensorflow.image module https://www.tensorflow.org/api_docs/python/tf/image

³read more in this wikipedia article https://en.wikipedia.org/wiki/Bilinear_interpolation

Underlying theory of the DDPM

In order to implement a DDPM, it is important to be familiar with some of the underlying theory, as the hyperparameters aim to fulfill certain conditions, which we will lay bare in the following section. The notation follows C. F. Higham, D. J. Higham, and Grindrod 2023 precisely and the reader can refer to this article for all of the theoretical results that follow.

One crucial result from the probability theory underlying the DDPM is that

$$q(x_t|x_0) \rightarrow \mathcal{N}(x_t; 0, I), \quad \text{as } t \rightarrow \infty, \quad (1)$$

where q is the transition density of the forward process, t is the time-step and \mathcal{N} is a normal density - i.e. that given a starting image x_0 , the diffused image at some late time-step T , x_T , should be indistinguishable from standard normal noise. Without this property, it wouldn't be reasonable to start from pure noise when generating new images from the trained model. Using the definition of the forward process together with the Markovian property, it follows that the transition density is defined as

$$q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, \sqrt{1-\bar{\alpha}_t} I), \quad (2)$$

where

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

and $\alpha_t = 1 - \beta_t$ for monotonically increasing sequence $\{\beta_t \in (0, 1)\}_{t=1}^T$. It is now evident from equations (1) and (2) that controlling the decay of $\sqrt{\bar{\alpha}_t}$ such that it tends to zero for some sufficiently large $t = T$ is of utmost importance when defining the diffusion process.

Setting the distribution scheme of the $\{\beta_t\}$ to a uniform spread on the $[0.0001, 0.02]$ interval (mimicking the authors Ho, Jain, and Abbeel 2020, of the original DDPM article), we can now focus entirely on varying T such that the decay becomes sufficient.

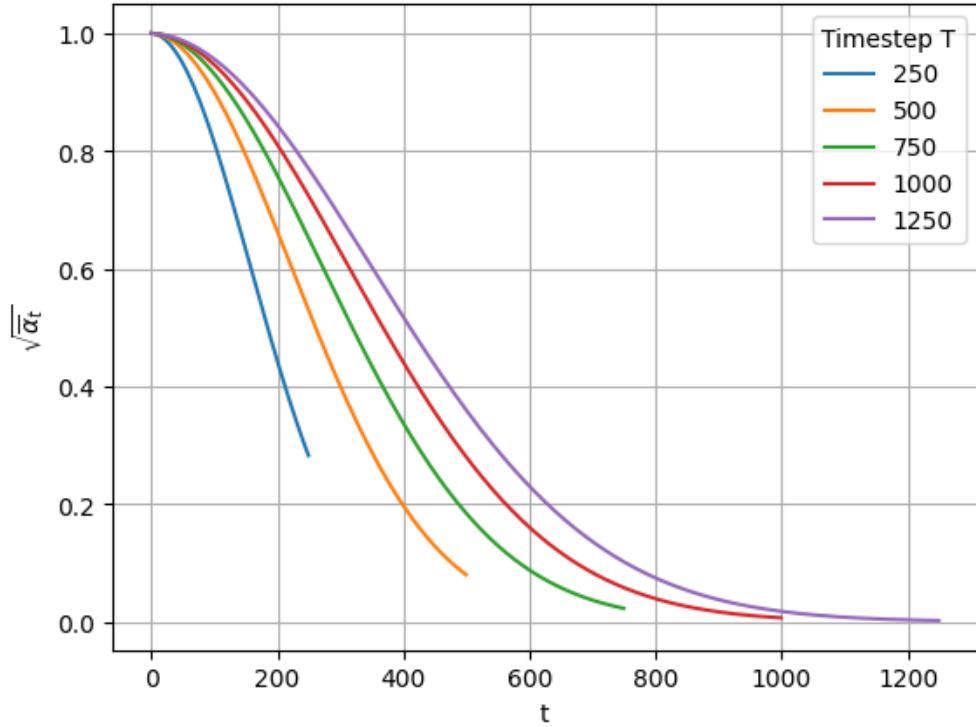


Figure 3: Decays of $\sqrt{\alpha_t}$ for varying number of noise-adding steps T .

Based on the results of Figure 3, we decide to go with $T = 1000$, as this yields $\sqrt{\alpha_T} \approx 10^{-5}$ which we deem to be sufficiently close to zero (this is also what the authors Ho, Jain, and Abbeel 2020 went with, which we have now confirmed the reasons for).

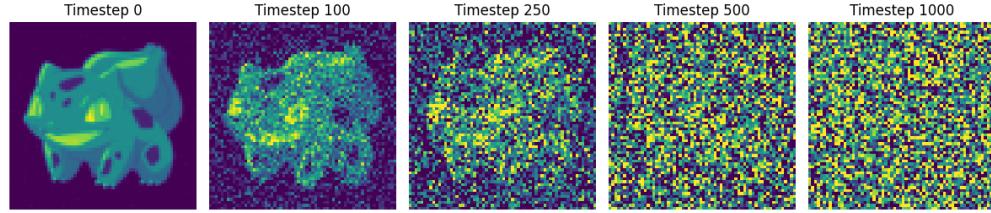


Figure 4: Result of forward process over time on an example of a cropped, rescaled and resized image of Bulbasaur using $T = 1000$.

In Figure 4 we see an example of the forward process applied to a pre-processed image of our dataset, at a select few timesteps. Looking carefully, one might still be able to glimpse some semblance of a structure (not purely noise) at $t = 500$. However faint the structure looks to our eyes, the computer will see it more clearly.

Architecture

Having motivated our decision for the hyperparameters of the diffusion process, we will in this section briefly present the architecture of the neural network that we use to fit the weights and biases of the DDPM.

The U-net architecture, widely used for implementing diffusion models, is effective due to its capability to handle image segmentation tasks (Williams et al. 2023). As our purpose is to generate images of Bulbasaur, U-Nets would be capable of capturing important features, like the red eyes and its leaf-like shell. Therefore we have used a U-Net based architecture with residual and attention blocks following the code from Keras official website⁴.

The encoder part of a U-Net processes the noisy image through a series of convolutional operations, similarly to how a standard convolution neural network works. The difference comes in the decoder part of the network which reconstructs the image by transitioning back to the original resolution. This process takes help from the skip connections which incorporate features from corresponding levels in the encoder to preserve important spatial information (Williams et al. 2023).

The encoder part of the network we used for training consists of two 3x3 convolution layers with a sigmoid linear (SiLu) activation function with added padding to match output pixels with input pixels. During this part, the pixels will be reduced while the channels will increase. The channels are increased to store information from the decreased pixels, this ensures no information gets lost during training.

The network consists of four levels in which the last two levels have an attention block added. These are added to pick up salient inputs. At the bottleneck, convolution layers and attention blocks are applied to generate a feature map of size (8, 8, 256). Finally, the decoder part uses the feature map from the bottleneck to convert the image back to its original input size. This phase works similar as the encoder phase with the difference that it first up samples the images using interpolative resizing of type “nearest” to match the sizing to its opposite level and then follow with matching convolution layers and attention blocks. The final output layer will be of dimensions (64, 64, 3). This architecture landed on around 30 million trainable parameters.

Training and results

Finally, we can use the architecture and the chosen $\{\sqrt{\bar{\alpha}_t}\}$ with $T = 1000$ to train our model. Fitting the model is equivalent to a least-squares regression problem where the randomly sampled noise ϵ acts as response variable, and the timestep $t \in \{1, \dots, T\}$ and image x_t , derived from applying the diffusion process up to time t to an image x_0 from the training data, act as a pair of predictors. It can be shown through the maximization of the ELBO(x) (along with some assumptions of a dominating term) that this really turns into a problem of least squares⁵, hence regression and the minimization of

$$\|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t)\|_2^2, \quad \forall t \in \{1, \dots, T\}.$$

⁴<https://keras.io/examples/generative/ddpm/>

⁵this was shown during a lecture of the course MT7042 at Stockholm University and we have no official source for this

It is this ϵ_θ that our neural network aims to model. In the training process one uniformly samples a t once per image and epoch. Therefore we believe the number of epochs needs to be larger than the number of time-steps T in order to cover most of the t 's for an image. For more information regarding the training of the model and the actual generating of the images from the trained model, see **Algorithm 1 & 2** in C. F. Higham, D. J. Higham, and Grindrod 2023. We now move on to showing some of the results.

Using a slightly lower than recommended learning rate ($= 2 \cdot 10^{-4}$), we land on running 2000 epochs of training (which took around 12 hours to run on our best CPU) using the Adam optimizer with default⁶ settings, and the MSE loss function. The sampling process is then applied to generate 16 new images from this trained model, starting from standard normal noise.



Figure 5: 16 images sampled from a 2000 epoch model using a linear decay on the $\{\alpha_t\}$, $T = 1000$, learning rate= $2 \cdot 10^{-4}$, the mean-squared-error loss function and a default Adam optimizer.

From Figure 5 we can see that the model did a fairly good job on generating Bulbasaurs. With minor deforming shapes, it did capture that it has two red eyes, some spots on its body and a leaf-like shell on its back. Since some of these images look a bit too much like the original images our training might have overfitted. Usually this would be the place for validation of the results, but, due to DDPMs being so new, agreed upon methods of validation are yet to be established. More on this in the next section.

However, there are ways to monitor the training of the model. If we split the data into a training set and a test set, we can still use a training/validation error plot in the usual sense. We do, after all, have a regression problem (even if diffusion models are generally considered as part of unsupervised learning), albeit a bit special, as we are using randomly sampled noise as our response variables. A validation of the training would then amount to, apart from the noise that we have fit the model to, randomly sampling an additional noise, and seeing how well that noise is predicted from the given test data.

⁶for the default settings, see https://www.tensorflow.org/api_docs/python/tf/keras/optimizers/Adam

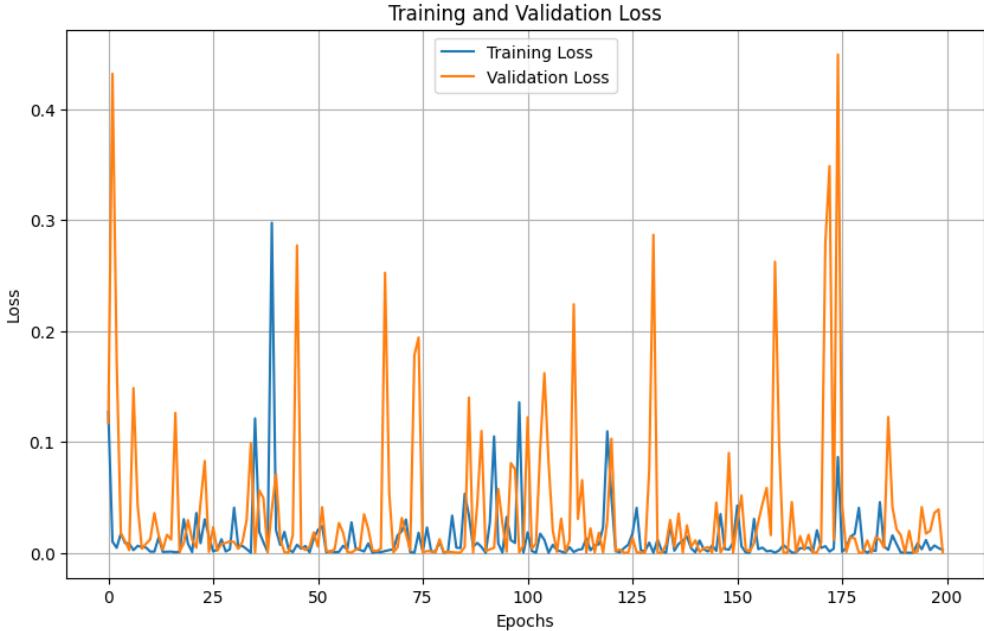


Figure 6: Training and validation error plot of the DDPM using a 80/20 split of the data. Epochs should be multiplied by ten.

Figure 6 represents a model that is different from the one that gave us the results in Figure 5 (which used the whole dataset for training only) in the sense that we have now split the data into a training set and a test set. As such, this should not be interpreted as a complete diagnosis of the training of the previous model, but instead serve as an illustratory example, or at best point out some potential issues with the previous model. The reason for not using this split data in the previous model is that our goal was to generate images of Bulbasaur. Of course we want these images to be as good as possible and therefore we should of course dedicate the entire dataset to training.

Interpreting Figure 6, we see that for the most part, both the training and the validation error is relatively low, the training more so than the validation, which is expected. Were it not for the excessive “spikiness” of the curves, we might have said that the training looks healthy. So where do these spikes come from, and why do they appear to have such a random distribution along the curves? One explanation, which as far as we can tell is also the best, is that as we are trying to fit a model to noise we should expect to see this noise in the model. Usually when performing regression, the predictors and the responses portray some linked pattern. One type of change in the predictor has the tendency for a particular type change in the response. For the DDPM this is not the case. Irregardless of the input, the response is still just standard normal noise. However, just because the response is noise, it doesn’t mean that we don’t have any trends. The standard normal distribution will of course produce more values close to zero, as this is where its mean lies. As such, it would be more beneficial for the model to predict more values close to zero and this is what we expect to see.

When studying these training/validation curves for the DDPM, it is thus not the shape of the curve, but the distribution of the errors that matter. For a healthy model, we expect the errors to behave like the square of a standard normal distribution (since we are using the MSE), i.e. like a chi-squared distribution. This could very well be the case for both the training and validation errors in Figure 6 and we should thus not write it off as a poor fit immediately (unfortunately we could not find any references for this - it might not be true at all even).

Possible improvements

Validation on result

Validating that our result is reasonable is a bit tricky as DDPM is a relatively new field in deep learning and therefore there are no well-established methods. What we want to do is decide the trade-off between plausibility and novelty. But as there is no real method for this we may simply compare the images with our own eyes. We saw that a few of the generated images looked exactly the same as the original images but the computer might detect some small difference that cannot be seen with our eyes only. Thus, this validation technique is not at all reliable.

Another validation could be to compute the Fréchet inception distance to decide the quality of the generated images. This would be done by comparing the distribution of two sets of images, one of the generated images and the other of the real images, to see if they match. This method uses an inception network⁷ to compare the distance of the mean and standard deviation between the two image sets with the defined formula:

$$d_F(\mathcal{N}(\mu, \Sigma), \mathcal{N}(\mu', \Sigma'))^2 = \|\mu - \mu'\|_2^2 + \text{tr}\left(\Sigma + \Sigma' - 2(\Sigma\Sigma')^{\frac{1}{2}}\right),$$

where higher values indicate poor generative images and lower values indicate good generated images.

Monitoring the training

There could be some improvements with our training/validation error. Seeing as our dataset was very small, we could have made a smaller split, in the ratio of 90/10 instead of 80/20 since this made our training set even smaller. However, a better alternative might have been to use a leave-one-out cross-validation, but because this would take far too long seeing as it already takes quite a lot of time to train the dataset already we opted to not apply this.

How to handle overfitting

DDPMs are designed to train noise and predict noise, hence overfitting these models would mean that our generated images are reproducing the exact copy of the original data. Our generated data seem to overfit, if we look at the top-right image in Figure 1 this image can be spotted in our generated images in Figure 5. Since our dataset was fairly small, one way to prevent overfitting would be to incorporate additional data augmentation.

Furthermore, since the problem is regression one might be inclined to regularize using methods such as the L_2 penalty (ridge regression). Apparently, this is not effective when using adaptive learning algorithms (Loshchilov and Hutter 2019) such as Adam, and one would have switch to something like SGD with momentum instead (which apparently performs very well when paired with ridge). One could also try using the L_1 penalization (lasso). In general, most regularization strategies that are applicable to regression problems will be applicable here as well. Another option would be dropout. This seems especially feasible as there are a lot of trainable parameters in our particular U-net.

⁷[https://en.wikipedia.org/wiki/Inception_\(deep_learning_architecture\)](https://en.wikipedia.org/wiki/Inception_(deep_learning_architecture))

Appendix A: training the model

Here, we show some of the generated images that were sampled as the training of the model went on.

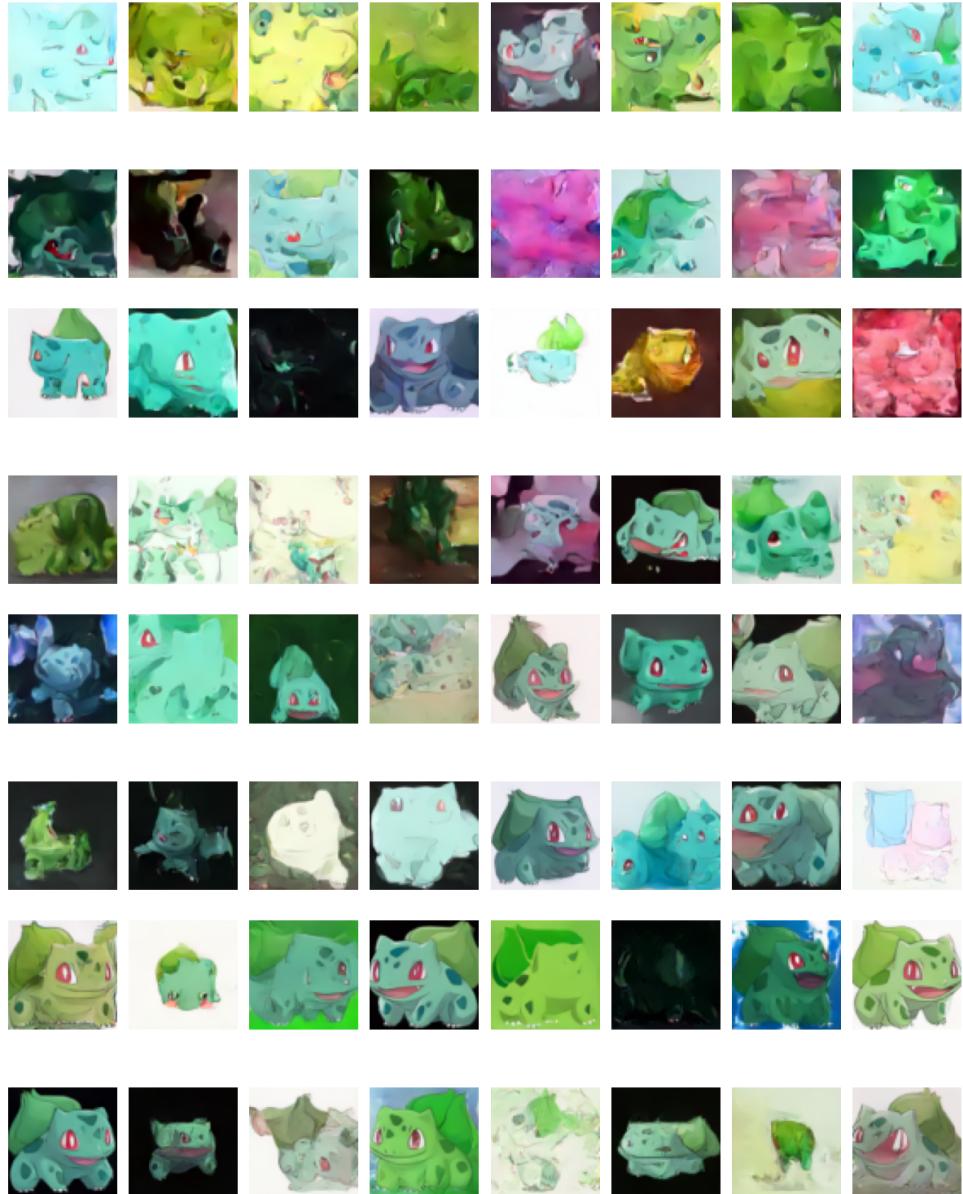


Figure 7: Generated images in sets of 16 from the training of the model. First set corresponds to the model after training for 200 epochs; second: 800 epochs; third: 1400 epochs; fourth: 2000 epochs.

Appendix B: Python Code

Most of the code was inspired by <https://keras.io/examples/generative/ddpm/>, and several functions were used without change.

```
1 #%%
2 import kagglehub
3
4 # Download latest version
5 path = kagglehub.dataset_download("thedagger/pokemon-generation-one")
6
7 print("Path to dataset files:", path)
8
9
10 #%%
11 import math
12 import matplotlib.pyplot as plt
13 import random
14
15 # Requires TensorFlow >=2.11 for the GroupNormalization layer.
16 import tensorflow as tf
17 from tensorflow import keras
18 from tensorflow.keras import layers
19
20 #%%
21 from sklearn.model_selection import train_test_split
22 import numpy as np
23 from PIL import Image
24 import os
25
26
27 bulbasaur = os.path.join(os.getcwd(), "Bulbasaur")
28 file_paths = []
29 labels = []
30
31 for file_name in os.listdir(bulbasaur):
32     file_path = os.path.join(bulbasaur, file_name)
33     if os.path.isfile(file_path):
34         file_paths.append(file_path)
35         labels.append("bulbasaur")
36
37 train_images_ds = tf.data.Dataset.from_tensor_slices(file_paths)
38
39 X_train, X_test, y_train, y_test = train_test_split(file_paths, labels, test_size=0.1,
40                                                     random_state=42)
41
42 train_images_ds = tf.data.Dataset.from_tensor_slices(X_train)
43 test_images_ds = tf.data.Dataset.from_tensor_slices(X_test)
44
45 # %%
46 batch_size = 1
47 num_epochs = 2000 # Just for the sake of demonstration
48 total_timesteps = 1000
49 norm_groups = 8 # Number of groups used in GroupNormalization layer
50 learning_rate = 0.001
51
52 img_size = 64
53 img_channels = 3
54 clip_min = -1.0
55 clip_max = 1.0
56
57 first_conv_channels = 32
58 channel_multiplier = [1, 2, 4, 8]
```

```

59 widths = [first_conv_channels * mult for mult in channel_multiplier]
60 has_attention = [False, False, True, True]
61 num_res_blocks = 2 # Number of residual blocks
62
63 # %%
64 def augment(img):
65     """Flips an image left/right randomly."""
66     return tf.image.random_flip_left_right(img)
67 # tf.image.random_contrast(img)
68
69 def resize_and_rescale(img, size):
70     """Resize the image to the desired size first and then
71     rescale the pixel values in the range [-1.0, 1.0].
72
73     Args:
74         img: Image tensor
75         size: Desired image size for resizing
76     Returns:
77         Resized and rescaled image tensor
78     """
79
80     height = tf.shape(img)[0]
81     width = tf.shape(img)[1]
82     crop_size = tf.minimum(height, width)
83
84     img = tf.image.crop_to_bounding_box(
85         img,
86         (height - crop_size) // 2,
87         (width - crop_size) // 2,
88         crop_size,
89         crop_size,
90     )
91
92     # Resize
93     img = tf.cast(img, dtype=tf.float32)
94     img = tf.image.resize(img, size=size, antialias=True)
95
96     # Rescale the pixel values
97     img = img / 127.5 - 1.0
98     img = tf.clip_by_value(img, clip_min, clip_max)
99     return img
100
101
102 def train_preprocessing(file_path):
103     img = tf.io.read_file(file_path) # Read the image file
104     img = tf.image.decode_image(img, channels=img_channels) # Decode as RGB
105     img = resize_and_rescale(img, size=(img_size, img_size))
106     img = augment(img)
107     return img
108
109 train_ds = (
110     train_images_ds.map(train_preprocessing, num_parallel_calls=tf.data.AUTOTUNE)
111     .batch(batch_size, drop_remainder=True)
112     .shuffle(batch_size * 2)
113     .prefetch(tf.data.AUTOTUNE)
114 )
115
116 test_ds = (
117     test_images_ds.map(train_preprocessing, num_parallel_calls=tf.data.AUTOTUNE)
118     .batch(batch_size, drop_remainder=True)
119     .shuffle(batch_size * 2)
120     .prefetch(tf.data.AUTOTUNE)
121 )
122

```

```

123
124 # %%
125 class GaussianDiffusion:
126     """Gaussian diffusion utility.
127
128     Args:
129         beta_start: Start value of the scheduled variance
130         beta_end: End value of the scheduled variance
131         timesteps: Number of time steps in the forward process
132     """
133
134     def __init__(self,
135                  beta_start=1e-4,
136                  beta_end=0.02,
137                  timesteps=500,
138                  clip_min=-1.0,
139                  clip_max=1.0,
140                  ):
141         self.beta_start = beta_start
142         self.beta_end = beta_end
143         self.timesteps = timesteps
144         self.clip_min = clip_min
145         self.clip_max = clip_max
146
147         # Define the linear variance schedule
148         self.betas = betas = np.linspace(
149             beta_start,
150             beta_end,
151             timesteps,
152             dtype=np.float64, # Using float64 for better precision
153         )
154         self.num_timesteps = int(timesteps)
155
156         alphas = 1.0 - betas
157         alphas_cumprod = np.cumprod(alphas, axis=0)
158         alphas_cumprod_prev = np.append(1.0, alphas_cumprod[:-1])
159
160         self.betas = tf.constant(betas, dtype=tf.float32)
161         self.alphas_cumprod = tf.constant(alphas_cumprod, dtype=tf.float32)
162         self.alphas_cumprod_prev = tf.constant(alphas_cumprod_prev, dtype=tf.float32)
163
164         # Calculations for diffusion q(x_t | x_{t-1}) and others
165         self.sqrt_alphas_cumprod = tf.constant(
166             np.sqrt(alphas_cumprod), dtype=tf.float32
167         )
168
169         self.sqrt_one_minus_alphas_cumprod = tf.constant(
170             np.sqrt(1.0 - alphas_cumprod), dtype=tf.float32
171         )
172
173         self.log_one_minus_alphas_cumprod = tf.constant(
174             np.log(1.0 - alphas_cumprod), dtype=tf.float32
175         )
176
177         self.sqrt_recip_alphas_cumprod = tf.constant(
178             np.sqrt(1.0 / alphas_cumprod), dtype=tf.float32
179         )
180         self.sqrt_recipm1_alphas_cumprod = tf.constant(
181             np.sqrt(1.0 / alphas_cumprod - 1), dtype=tf.float32
182         )
183
184         # Calculations for posterior q(x_{t-1} | x_t, x_0)
185         posterior_variance = (

```

```

187         betas * (1.0 - alphas_cumprod_prev) / (1.0 - alphas_cumprod)
188     )
189     self.posterior_variance = tf.constant(posterior_variance, dtype=tf.float32)
190
191     # Log calculation clipped because the posterior variance is 0 at the beginning
192     # of the diffusion chain
193     self.posterior_log_variance_clipped = tf.constant(
194         np.log(np.maximum(posterior_variance, 1e-20)), dtype=tf.float32
195     )
196
197     self.posterior_mean_coef1 = tf.constant(
198         betas * np.sqrt(alphas_cumprod_prev) / (1.0 - alphas_cumprod),
199         dtype=tf.float32,
200     )
201
202     self.posterior_mean_coef2 = tf.constant(
203         (1.0 - alphas_cumprod_prev) * np.sqrt(alphas) / (1.0 - alphas_cumprod),
204         dtype=tf.float32,
205     )
206
207 def _extract(self, a, t, x_shape):
208     """Extract some coefficients at specified timesteps,
209     then reshape to [batch_size, 1, 1, 1, 1, ...] for broadcasting purposes.
210
211     Args:
212         a: Tensor to extract from
213         t: Timestep for which the coefficients are to be extracted
214         x_shape: Shape of the current batched samples
215     """
216     batch_size = x_shape[0]
217     out = tf.gather(a, t)
218     return tf.reshape(out, [batch_size, 1, 1, 1])
219
220 def q_mean_variance(self, x_start, t):
221     """Extracts the mean, and the variance at current timestep.
222
223     Args:
224         x_start: Initial sample (before the first diffusion step)
225         t: Current timestep
226     """
227     x_start_shape = tf.shape(x_start)
228     mean = self._extract(self.sqrt_alphas_cumprod, t, x_start_shape) * x_start
229     variance = self._extract(1.0 - self.alphas_cumprod, t, x_start_shape)
230     log_variance = self._extract(
231         self.log_one_minus_alphas_cumprod, t, x_start_shape
232     )
233     return mean, variance, log_variance
234
235 def q_sample(self, x_start, t, noise):
236     """Diffuse the data.
237
238     Args:
239         x_start: Initial sample (before the first diffusion step)
240         t: Current timestep
241         noise: Gaussian noise to be added at the current timestep
242     Returns:
243         Diffused samples at timestep 't'
244     """
245     x_start_shape = tf.shape(x_start)
246     return (
247         self._extract(self.sqrt_alphas_cumprod, t, tf.shape(x_start)) * x_start
248         + self._extract(self.sqrt_one_minus_alphas_cumprod, t, x_start_shape)
249         * noise
250     )

```

```

251
252     def predict_start_from_noise(self, x_t, t, noise):
253         x_t_shape = tf.shape(x_t)
254         return (
255             self._extract(self.sqrt_recip_alphas_cumprod, t, x_t_shape) * x_t
256             - self._extract(self.sqrt_recipm1_alphas_cumprod, t, x_t_shape) * noise
257         )
258
259     def q_posterior(self, x_start, x_t, t):
260         """Compute the mean and variance of the diffusion
261         posterior  $q(x_{t-1} | x_t, x_0)$ .
262
263         Args:
264             x_start: Starting point(sample) for the posterior computation
265             x_t: Sample at timestep 't'
266             t: Current timestep
267         Returns:
268             Posterior mean and variance at current timestep
269         """
270
271         x_t_shape = tf.shape(x_t)
272         posterior_mean = (
273             self._extract(self.posterior_mean_coef1, t, x_t_shape) * x_start
274             + self._extract(self.posterior_mean_coef2, t, x_t_shape) * x_t
275         )
276         posterior_variance = self._extract(self.posterior_variance, t, x_t_shape)
277         posterior_log_variance_clipped = self._extract(
278             self.posterior_log_variance_clipped, t, x_t_shape
279         )
280         return posterior_mean, posterior_variance, posterior_log_variance_clipped
281
282     def p_mean_variance(self, pred_noise, x, t, clip_denoised=True):
283         x_recon = self.predict_start_from_noise(x, t=t, noise=pred_noise)
284         if clip_denoised:
285             x_recon = tf.clip_by_value(x_recon, self.clip_min, self.clip_max)
286
287         model_mean, posterior_variance, posterior_log_variance = self.q_posterior(
288             x_start=x_recon, x_t=x, t=t
289         )
290         return model_mean, posterior_variance, posterior_log_variance
291
292     def p_sample(self, pred_noise, x, t, clip_denoised=True):
293         """Sample from the diffusion model.
294
295         Args:
296             pred_noise: Noise predicted by the diffusion model
297             x: Samples at a given timestep for which the noise was predicted
298             t: Current timestep
299             clip_denoised (bool): Whether to clip the predicted noise
300                 within the specified range or not.
301         """
302         model_mean, _, model_log_variance = self.p_mean_variance(
303             pred_noise, x=x, t=t, clip_denoised=clip_denoised
304         )
305         noise = tf.random.normal(shape=x.shape, dtype=x.dtype)
306         # No noise when t == 0
307         nonzero_mask = tf.reshape(
308             1 - tf.cast(tf.equal(t, 0), tf.float32), [tf.shape(x)[0], 1, 1, 1]
309         )
310         return model_mean + nonzero_mask * tf.exp(0.5 * model_log_variance) * noise
311
312
313 # %%
314 # Kernel initializer to use

```

```

315 def kernel_init(scale):
316     scale = max(scale, 1e-10)
317     return keras.initializers.VarianceScaling(
318         scale, mode="fan_avg", distribution="uniform"
319     )
320
321
322 class AttentionBlock(layers.Layer):
323     """Applies self-attention.
324
325     Args:
326         units: Number of units in the dense layers
327         groups: Number of groups to be used for GroupNormalization layer
328     """
329
330     def __init__(self, units, groups=8, **kwargs):
331         self.units = units
332         self.groups = groups
333         super().__init__(**kwargs)
334
335         self.norm = layers.GroupNormalization(groups=groups)
336         self.query = layers.Dense(units, kernel_initializer=kernel_init(1.0))
337         self.key = layers.Dense(units, kernel_initializer=kernel_init(1.0))
338         self.value = layers.Dense(units, kernel_initializer=kernel_init(1.0))
339         self.proj = layers.Dense(units, kernel_initializer=kernel_init(0.0))
340
341     def call(self, inputs):
342         batch_size = tf.shape(inputs)[0]
343         height = tf.shape(inputs)[1]
344         width = tf.shape(inputs)[2]
345         scale = tf.cast(self.units, tf.float32) ** (-0.5)
346
347         inputs = self.norm(inputs)
348         q = self.query(inputs)
349         k = self.key(inputs)
350         v = self.value(inputs)
351
352         attn_score = tf.einsum("bhwc, bHWc->bhwHW", q, k) * scale
353         attn_score = tf.reshape(attn_score, [batch_size, height, width, height * width])
354
355         attn_score = tf.nn.softmax(attn_score, -1)
356         attn_score = tf.reshape(attn_score, [batch_size, height, width, height, width])
357
358         proj = tf.einsum("bhwHW, bHWc->bhwc", attn_score, v)
359         proj = self.proj(proj)
360         return inputs + proj
361
362
363 class TimeEmbedding(layers.Layer):
364     def __init__(self, dim, **kwargs):
365         super().__init__(**kwargs)
366         self.dim = dim
367         self.half_dim = dim // 2
368         self.emb = math.log(10000) / (self.half_dim - 1)
369         self.emb = tf.exp(tf.range(self.half_dim, dtype=tf.float32) * -self.emb)
370
371     def call(self, inputs):
372         inputs = tf.cast(inputs, dtype=tf.float32)
373         emb = inputs[:, None] * self.emb[None, :]
374         emb = tf.concat([tf.sin(emb), tf.cos(emb)], axis=-1)
375         return emb
376

```

```

377
378 def ResidualBlock(width, groups=8, activation_fn=keras.activations.swish):
379     def apply(inputs):
380         x, t = inputs
381         input_width = x.shape[3]
382
383         if input_width == width:
384             residual = x
385         else:
386             residual = layers.Conv2D(
387                 width, kernel_size=1, kernel_initializer=kernel_init(1.0)
388             )(x)
389
390         temb = activation_fn(t)
391         temb = layers.Dense(width, kernel_initializer=kernel_init(1.0))(temb)[
392             :, :
393         ]
394
395         x = layers.GroupNormalization(groups=groups)(x)
396         x = activation_fn(x)
397         x = layers.Conv2D(
398             width, kernel_size=3, padding="same", kernel_initializer=kernel_init(1.0)
399             )(x)
400
401         x = layers.Add()([x, temb])
402         x = layers.GroupNormalization(groups=groups)(x)
403         x = activation_fn(x)
404
405         x = layers.Conv2D(
406             width, kernel_size=3, padding="same", kernel_initializer=kernel_init(0.0)
407             )(x)
408         x = layers.Add()([x, residual])
409         return x
410
411     return apply
412
413
414 def DownSample(width):
415     def apply(x):
416         x = layers.Conv2D(
417             width,
418             kernel_size=3,
419             strides=2,
420             padding="same",
421             kernel_initializer=kernel_init(1.0),
422             )(x)
423         return x
424
425     return apply
426
427
428 def UpSample(width, interpolation="nearest"):
429     def apply(x):
430         x = layers.UpSampling2D(size=2, interpolation=interpolation)(x)
431         x = layers.Conv2D(
432             width, kernel_size=3, padding="same", kernel_initializer=kernel_init(1.0)
433             )(x)
434         return x
435
436     return apply
437
438
439 def TimeMLP(units, activation_fn=keras.activations.swish):
440     def apply(inputs):

```

```

441     temb = layers.Dense(
442         units, activation=activation_fn, kernel_initializer=kernel_init(1.0)
443     )(inputs)
444     temb = layers.Dense(units, kernel_initializer=kernel_init(1.0))(temb)
445     return temb
446
447     return apply
448
449 # %%
450 def build_model(
451     img_size,
452     img_channels,
453     widths,
454     has_attention,
455     num_res_blocks=2,
456     norm_groups=8,
457     interpolation="nearest",
458     activation_fn=keras.activations.swish,
459 ):
460     image_input = layers.Input(
461         shape=(img_size, img_size, img_channels), name="image_input"
462     )
463     time_input = keras.Input(shape=(), dtype=tf.int64, name="time_input")
464
465     x = layers.Conv2D(
466         first_conv_channels,
467         kernel_size=(3, 3),
468         padding="same",
469         kernel_initializer=kernel_init(1.0),
470     )(image_input)
471
472     temb = TimeEmbedding(dim=first_conv_channels * 4)(time_input)
473     temb = TimeMLP(units=first_conv_channels * 4, activation_fn=activation_fn)(temb)
474
475     skips = [x]
476
477     # DownBlock
478     for i in range(len(widths)):
479         for _ in range(num_res_blocks):
480             x = ResidualBlock(
481                 widths[i], groups=norm_groups, activation_fn=activation_fn
482             )([x, temb])
483             if has_attention[i]:
484                 x = AttentionBlock(widths[i], groups=norm_groups)(x)
485             skips.append(x)
486
487             if widths[i] != widths[-1]:
488                 x = DownSample(widths[i])(x)
489             skips.append(x)
490
491     # MiddleBlock
492     x = ResidualBlock(widths[-1], groups=norm_groups, activation_fn=activation_fn)(
493         [x, temb]
494     )
495     x = AttentionBlock(widths[-1], groups=norm_groups)(x)
496     x = ResidualBlock(widths[-1], groups=norm_groups, activation_fn=activation_fn)(
497         [x, temb]
498     )
499
500     # UpBlock
501     for i in reversed(range(len(widths))):
502         for _ in range(num_res_blocks + 1):
503             x = layers.concatenate(axis=-1)([x, skips.pop()])
504             x = ResidualBlock(

```

```

505             widths[i], groups=norm_groups, activation_fn=activation_fn
506         )(([x, temb])
507         if has_attention[i]:
508             x = AttentionBlock(widths[i], groups=norm_groups)(x)
509
510         if i != 0:
511             x = UpSample(widths[i], interpolation=interpolation)(x)
512
513     # End block
514     x = layers.GroupNormalization(groups=norm_groups)(x)
515     x = activation_fn(x)
516     x = layers.Conv2D(3, (3, 3), padding="same", kernel_initializer=kernel_init(0.0))(x)
517     return keras.Model([image_input, time_input], x, name="unet")
518
519 # %%
520 class DiffusionModel(keras.Model):
521     def __init__(self, network, ema_network, timesteps, gdf_util, ema=0.999):
522         super().__init__()
523         self.network = network
524         self.ema_network = ema_network
525         self.timesteps = timesteps
526         self.gdf_util = gdf_util
527         self.ema = ema
528
529     def train_step(self, images):
530         # 1. Get the batch size
531         batch_size = tf.shape(images)[0]
532
533         # 2. Sample timesteps uniformly
534         t = tf.random.uniform(
535             minval=0, maxval=self.timesteps, shape=(batch_size,), dtype=tf.int64
536         )
537
538         with tf.GradientTape() as tape:
539             # 3. Sample random noise to be added to the images in the batch
540             noise_train = tf.random.normal(shape=tf.shape(images), dtype=images.dtype)
541
542             # 4. Diffuse the images with noise
543             images_t = self.gdf_util.q_sample(images, t, noise_train)
544
545             # 5. Pass the diffused images and time steps to the network
546             pred_noise = self.network([images_t, t], training=True)
547
548             # 6. Calculate the loss
549             loss = self.loss(noise_train, pred_noise)
550
551             # 7. Get the gradients
552             gradients = tape.gradient(loss, self.network.trainable_weights)
553
554             # 8. Update the weights of the network
555             self.optimizer.apply_gradients(zip(gradients, self.network.trainable_weights))
556
557             # 9. Updates the weight values for the network with EMA weights
558             for weight, ema_weight in zip(self.network.weights, self.ema_network.weights):
559                 ema_weight.assign(self.ema * ema_weight + (1 - self.ema) * weight)
560
561             # 10. Return loss values
562             return {"training_loss": loss}
563
564     def test_step(self, images):
565
566         batch_size = tf.shape(images)[0]
567         t = tf.random.uniform(

```

```

568     minval=0, maxval=self.timesteps, shape=(batch_size,), dtype=tf.int64
569 )
570
571 noise_train = tf.random.normal(shape=tf.shape(images), dtype=images.dtype)
572 images_t = self.gdf_util.q_sample(images, t, noise_train)
573 pred_noise = self.network([images_t, t], training=False)
574 loss = self.loss(noise_train, pred_noise)
575
576     return {"validation_loss": loss}
577
578 def generate_images(self, num_images=16, initial_noise=None):
579
580     if initial_noise is None:
581         # 1. Randomly sample noise (starting point for reverse process)
582         samples = tf.random.normal(
583             shape=(num_images, img_size, img_size, img_channels), dtype=tf.float32
584         )
585         initial_noise = samples
586     else:
587         samples = initial_noise
588
589     # 2. Sample from the model iteratively
590     for t in reversed(range(0, self.timesteps)):
591         tt = tf.cast(tf.fill(num_images, t), dtype=tf.int64)
592         pred_noise = self.ema_network.predict(
593             [samples, tt], verbose=0, batch_size=num_images
594         )
595         samples = self.gdf_util.p_sample(
596             pred_noise, samples, tt, clip_denoised=True
597         )
598     # 3. Return generated samples
599     return samples, initial_noise
600
601 def plot_images(
602     self, epoch=None, logs=None, num_rows=2, num_cols=8, figsize=(12, 5), div_int
603 =5, initial_noise=None
604 ):
605     """Utility to plot images using the diffusion model during training."""
606
607     if (epoch+1) % div_int == 0:
608
609         generated_samples, _ = self.generate_images(num_images=num_rows * num_cols
610 , initial_noise=initial_noise)
611         generated_samples = (
612             tf.clip_by_value(generated_samples * 127.5 + 127.5, 0.0, 255.0)
613             .numpy()
614             .astype(np.uint8)
615         )
616
617         _, ax = plt.subplots(num_rows, num_cols, figsize=figsize)
618         if num_rows == 1 and num_cols == 1:
619             ax.imshow(generated_samples[0])
620             ax.axis("off")
621         else:
622             ax = np.array(ax).reshape(-1)
623             for i, image in enumerate(generated_samples):
624                 ax[i].imshow(image)
625                 ax[i].axis("off")
626
627         plt.tight_layout()
628         plt.show()
629
630 # %%
631 # Build the unet model

```

```

630 network = build_model(
631     img_size=img_size,
632     img_channels=img_channels,
633     widths=widths,
634     has_attention=has_attention,
635     num_res_blocks=num_res_blocks,
636     norm_groups=norm_groups,
637     activation_fn=keras.activations.swish,
638 )
639
640 ema_network = build_model(
641     img_size=img_size,
642     img_channels=img_channels,
643     widths=widths,
644     has_attention=has_attention,
645     num_res_blocks=num_res_blocks,
646     norm_groups=norm_groups,
647     activation_fn=keras.activations.swish,
648 )
649 ema_network.set_weights(network.get_weights()) # Initially the weights are the same
650
651 # Get an instance of the Gaussian Diffusion utilities
652 gdf_util = GaussianDiffusion(timesteps=total_timesteps)
653
654
655 # %%
656 # Get the model
657 model_3 = DiffusionModel(
658     network=network,
659     ema_network=ema_network,
660     gdf_util=gdf_util,
661     timesteps=total_timesteps,
662 )
663
664 # Compile the model
665 model_3.compile(
666     loss=keras.losses.MeanSquaredError(),
667     optimizer=keras.optimizers.Adam(learning_rate=learning_rate),
668 )
669
670 # Train the model
671 history = model_3.fit(
672     train_ds,
673     validation_data = test_ds,
674     epochs=num_epochs,
675     batch_size=batch_size,
676     callbacks=[keras.callbacks.LambdaCallback(on_epoch_end=model_1.plot_images)],
677 )
678
679 # %%
680 training_loss = history.history['training_loss']
681 validation_loss = history.history['val_validation_loss']
682
683 hej = [ training_loss[i] for i in range(len(training_loss)) if i % 10 == 0 ]
684 hi = [ validation_loss[i] for i in range(len(training_loss)) if i % 10 == 0 ]
685
686 plt.figure(figsize=(10, 6))
687 plt.plot(hej, label="Training Loss")
688 plt.plot(hi, label="Validation Loss")
689 plt.xlabel("Epochs")
690 plt.ylabel("Loss")
691 plt.title("Training and Validation Loss")
692 plt.legend()
693 plt.grid(True)

```

```

694 plt.show()
695
696 # %%
697 history2 = model_1.fit(
698     train_ds,
699     validation_data = test_ds,
700     epochs=1000,
701     batch_size=batch_size,
702     callbacks=[keras.callbacks.LambdaCallback(on_epoch_end=model_1.plot_images)],
703 )
704
705 # %%
706 training_loss_2 = history2.history['training_loss']
707 validation_loss_2 = history2.history['val_validation_loss']
708
709 hej_2 = [ training_loss_2[i] for i in range(len(training_loss_2)) if i % 10 == 0 ]
710 hi_2 = [ validation_loss_2[i] for i in range(len(validation_loss_2)) if i % 10 == 0 ]
711
712 training_loss_tot = hej + hej_2
713 validation_loss_tot = hi + hi_2
714
715 # %%
716 plt.figure(figsize=(10, 6))
717 plt.plot(training_loss_tot, label="Training Loss")
718 plt.plot(validation_loss_tot, label="Validation Loss")
719 plt.xlabel("Epochs")
720 plt.ylabel("Loss")
721 plt.title("Training and Validation Loss")
722 plt.legend()
723 plt.grid(True)
724 plt.show()
725
726 # %%
727 # the noise to generate image from
728 num_images = 1
729 initial_noise = tf.random.normal(
730     shape=(num_images, img_size, img_size, img_channels), dtype=tf.float32
731 )
732
733 # Generate and plot some samples from same noise
734 for i in range(10):
735     model_1.plot_images(num_rows=1, num_cols=1, div_int=1, epoch=0, initial_noise=
    initial_noise)
736
737 # %%
738 ## generate model with more pokemons
739
740 bulbasaur = os.path.join(os.getcwd(), "Bulbasaur")
741 pikachu = os.path.join(os.getcwd(), "Pikachu")
742 squirtle = os.path.join(os.getcwd(), "Squirtle")
743
744 file_paths = []
745 labels = []
746
747 for file_name in os.listdir(bulbasaur):
748     file_path = os.path.join(bulbasaur, file_name)
749     if os.path.isfile(file_path):
750         file_paths.append(file_path)
751         labels.append("bulbasaur")
752
753 train_images_ds_new = tf.data.Dataset.from_tensor_slices(file_paths)
754
755 train_ds_new = (
756     train_images_ds_new.map(train_preprocessing, num_parallel_calls=tf.data.AUTOTUNE)

```

```

757     .batch(batch_size, drop_remainder=True)
758     .shuffle(batch_size * 2)
759     .prefetch(tf.data.AUTOTUNE)
760   )
761
762 # %%
763 betas_250 = np.linspace(
764     0.0001,
765     0.02,
766     250,
767     dtype=np.float64, # Using float64 for better precision
768   )
769 alphas_250 = 1.0 - betas_250
770 alphas_cumprod_250 = np.cumprod(alphas_250, axis=0)
771
772 betas_500 = np.linspace(
773     0.0001,
774     0.02,
775     500,
776     dtype=np.float64, # Using float64 for better precision
777   )
778 alphas_500 = 1.0 - betas_500
779 alphas_cumprod_500 = np.cumprod(alphas_500, axis=0)
780
781
782 betas_750 = np.linspace(
783     0.0001,
784     0.02,
785     750,
786     dtype=np.float64, # Using float64 for better precision
787   )
788 alphas_750 = 1.0 - betas_750
789 alphas_cumprod_750 = np.cumprod(alphas_750, axis=0)
790
791 betas_1000 = np.linspace(
792     0.0001,
793     0.02,
794     1000,
795     dtype=np.float64, # Using float64 for better precision
796   )
797 alphas_1000 = 1.0 - betas_1000
798 alphas_cumprod_1000 = np.cumprod(alphas_1000, axis=0)
799
800 betas_1250 = np.linspace(
801     0.0001,
802     0.02,
803     1250,
804     dtype=np.float64, # Using float64 for better precision
805   )
806 alphas_1250 = 1.0 - betas_1250
807 alphas_cumprod_1250 = np.cumprod(alphas_1250, axis=0)
808
809
810 # %%
811 plt.plot(np.sqrt(alphas_cumprod_250), label="250")
812 plt.plot(np.sqrt(alphas_cumprod_500), label="500")
813 plt.plot(np.sqrt(alphas_cumprod_750), label="750")
814 plt.plot(np.sqrt(alphas_cumprod_1000), label="1000")
815 plt.plot(np.sqrt(alphas_cumprod_1250), label="1250")
816 plt.legend(title="Timestep T")
817 plt.grid(True)
818 plt.xlabel("t")
819 plt.ylabel("$\sqrt{\overline{\alpha_t}}$")
820 # plt.xlim(200, 1300)

```

```

821 # plt.ylim(-0.01, 0.1)
822 plt.show()
823
824
825 # %%
826 # for plotting the forward process
827 # Instantiate the GaussianDiffusion class
828 diffusion = GaussianDiffusion(beta_start=1e-4, beta_end=0.02, timesteps=1001)
829
830 # Load or create an initial image (x_start)
831 # Example: Generate a simple synthetic image (grayscale gradient)
832 x_start = train_preprocessing("bulbasaur_example.png")
833
834 # %%
835 # Normalize and expand dimensions to match expected shape [batch_size, height, width,
# channels]
836 x_start = np.expand_dims(np.expand_dims(x_start, axis=0), axis=-1) # Shape: [1, 64,
# 64, 1]
837
838 # Initialize a random number generator for noise
839 rng = np.random.default_rng(seed=42)
840
841 # Timesteps to visualize (e.g., start, middle, and final timesteps)
842 timesteps_to_visualize = [0, 100, 250, 500, 1000]
843
844 # Plot the diffusion process
845 plt.figure(figsize=(12, 6))
846 for idx, t in enumerate(timesteps_to_visualize):
# Add noise progressively
847     t_tensor = tf.constant([t], dtype=tf.int32) # Current timestep
848     noise = rng.normal(size=x_start.shape, scale=1.0).astype(np.float32) # Gaussian
noise
849     x_t = diffusion.q_sample(x_start, t_tensor, noise) # Diffused image
850
851     # Clip values to [-1, 1] range and rescale to [0, 1] for visualization
852     x_t_clipped = np.clip(x_t.numpy(), -1, 1)
853     x_t_rescaled = (x_t_clipped + 1) / 2.0 # Rescale to [0, 1]
854
855     # Plot the noisy image
856     plt.subplot(1, len(timesteps_to_visualize), idx + 1)
857     plt.imshow(x_t_rescaled[0, :, :, 0])
858     plt.title(f"Timestep {t}")
859     plt.axis("off")
860
861 plt.tight_layout()
862 plt.show()
863
864
865
866 # Exploratory analysis plot
867 # %%
868 def get_dims(filepath):
869     img = Image.open(filepath)
870     arr = np.array(img)
871
872     if arr.ndim == 2: # Grayscale image
873         height, width = arr.shape
874     elif arr.ndim == 3: # Color image
875         height, width, _ = arr.shape
876     else:
877         raise ValueError(f"Unexpected image dimensions: {arr.shape}")
878
879     return height, width
880
881 bulbasaur = os.path.join(os.getcwd(), "Bulbasaur")

```

```
882 filepaths = []
883 for file_name in os.listdir(bulbasaur):
884     filepath = os.path.join(bulbasaur, file_name)
885     if os.path.isfile(filepath):
886         filepaths.append(filepath)
887
888 image_sizes = {}
889 for filepath in filepaths:
890     height, width = get_dims(filepath)
891     image_sizes[filepath] = (height, width)
892
893 height = []
894 width = []
895 for filepath, size in image_sizes.items():
896     height.append(size[0])
897     width.append(size[1])
898
899
900 # %%
901 plt.scatter(width, height, s=2)
902 plt.xlabel("Width")
903 plt.ylabel("Height")
904 plt.grid(True)
905 plt.show()
```

References

- Loshchilov, Ilya and Frank Hutter (2019). *Decoupled Weight Decay Regularization*. arXiv: 1711.05101 [cs.LG]. URL: <https://arxiv.org/abs/1711.05101>.
- Ho, Jonathan, Ajay Jain, and Pieter Abbeel (2020). *Denoising Diffusion Probabilistic Models*. arXiv: 2006.11239 [cs.LG]. URL: <https://arxiv.org/abs/2006.11239>.
- Higham, Catherine F, Desmond J Higham, and Peter Grindrod (2023). “Diffusion Models for Generative Artificial Intelligence: An Introduction for Applied Mathematicians”. In: *arXiv preprint arXiv:2312.14977*.
- Williams, Christopher et al. (2023). “A unified framework for U-Net design and analysis”. In: *Advances in Neural Information Processing Systems* 36, pp. 27745–27782.