

assignment3-kanishksharma-22b1049

August 5, 2024

```
[1]: # @title Install requirements
!pip install datasets &>> install.log
```

```
[2]: #@title imports and utility functions
from datasets import load_dataset
from PIL import Image
import torch.nn.functional as F
import os
from tqdm.notebook import tqdm
import torch
import numpy as np

def img_to_tensor(im):
    return torch.tensor(np.array(im.convert('RGB'))/255).permute(2, 0, 1).
    ↪unsqueeze(0) * 2 - 1

def tensor_to_image(t):
    return Image.fromarray(np.array(((t.squeeze()).permute(1, 2, 0)+1)/2).clip(0, 1).
    ↪*255).astype(np.uint8))

def gather(consts: torch.Tensor, t: torch.Tensor):
    """Gather consts for $t$ and reshape to feature map shape"""
    c = consts.gather(-1, t)
    return c.reshape(-1, 1, 1, 1)
```

0.1 2.1 Dataset

We'll start with a classic small dataset, with 32px square images from 10 classes. For convenience we just pull a version that is available on the huggingface hub.

```
[3]: #@title cifar10 - 32px images in 10 classes

# Download and load the dataset
cifar10 = load_dataset('cifar10')

# View some examples:
```

```
image = Image.new('RGB', size=(32*5, 32*2))
for i in range(10):
    im = cifar10['train'][i]['img']
    image.paste(im, ( (i%5)*32, (i//5)*32 ))
image.resize((32*5*4, 32*2*4), Image.NEAREST)
```

/usr/local/lib/python3.10/dist-packages/huggingface_hub/utils/_token.py:89:

UserWarning:

The secret `HF_TOKEN` does not exist in your Colab secrets.

To authenticate with the Hugging Face Hub, create a token in your settings tab (<https://huggingface.co/settings/tokens>), set it as secret in your Google Colab and restart your session.

You will be able to reuse this secret in all of your notebooks.

Please note that authentication is recommended but still optional to access public models or datasets.

```
warnings.warn(
```

```
Downloading readme: 0%|          | 0.00/5.16k [00:00<?, ?B/s]
```

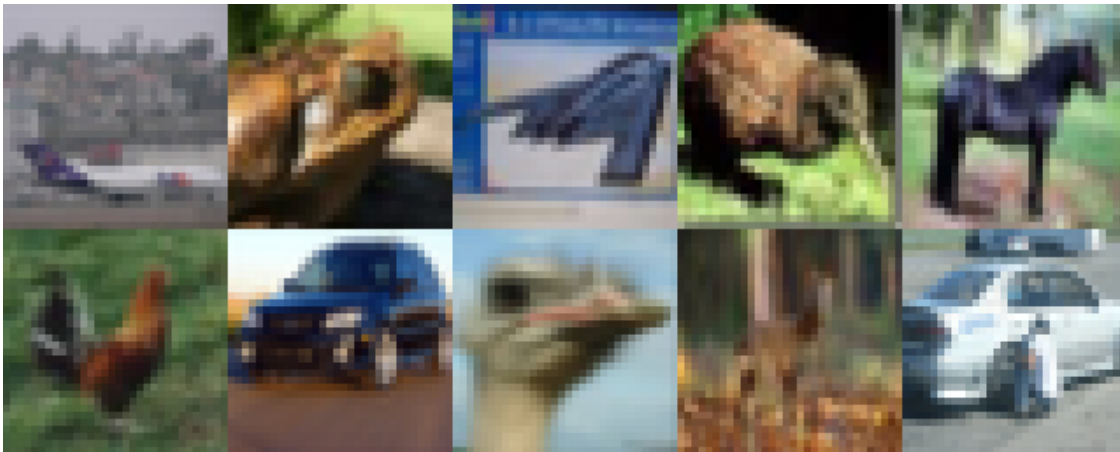
```
Downloading data: 0%|          | 0.00/120M [00:00<?, ?B/s]
```

```
Downloading data: 0%|          | 0.00/23.9M [00:00<?, ?B/s]
```

```
Generating train split: 0%|        | 0/50000 [00:00<?, ? examples/s]
```

```
Generating test split: 0%|         | 0/10000 [00:00<?, ? examples/s]
```

[3]:



```
[43]: n_steps = 100
beta = torch.linspace(0.0001, 0.04, n_steps)

def q_xt_xtminus1(xtm1, t):
    mean = gather(1. - beta, t) ** 0.5 * xtm1 #  $\sqrt{(1-t)*xtm1}$ 
    var = gather(beta, t) #  $t I$ 
    eps = torch.randn_like(xtm1) # Noise shaped like xtm1
```

```

    return mean + (var ** 0.5) * eps

# Show im at different stages
ims = []
start_im = cifar10['train'][10]['img']
x = img_to_tensor(start_im).squeeze()
for t in range(n_steps):

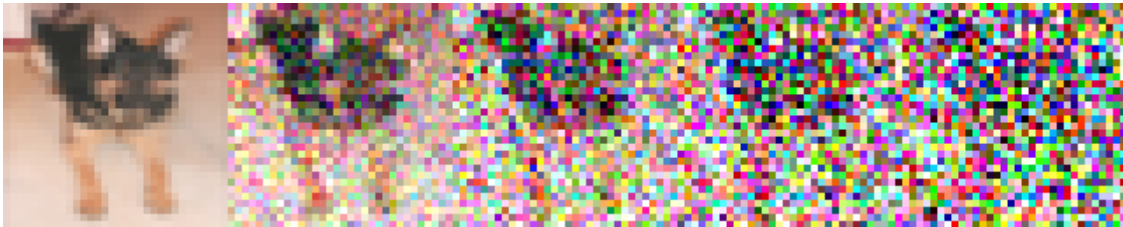
    # Store images every 20 steps to show progression
    if t%20 == 0:
        ims.append(tensor_to_image(x))

    # Calculate Xt given Xt-1 (i.e. x from the previous iteration)
    t = torch.tensor(t, dtype=torch.long) # t as a tensor
    x = q_xt_xtminus1(x, t) # Modify x using our function above

# Display the images
image = Image.new('RGB', size=(32*5, 32))
for i, im in enumerate(ims):
    image.paste(im, ((i%5)*32, 0))
image.resize((32*4*5, 32*4), Image.NEAREST)

```

[43]:



[24]:

```

n_steps = 1000
beta = torch.linspace(0.0001, 0.05, n_steps)
alpha = 1. - beta
alpha_bar = torch.cumprod(alpha, dim=0)

def q_xt_x0(x0, t):
    mean = gather(alpha_bar, t) ** 0.5 * x0 # now alpha_bar
    var = 1-gather(alpha_bar, t) # (1-alpha_bar)
    eps = torch.randn_like(x0)
    return mean + (var ** 0.5) * eps

# Show im at different stages
ims = []
start_im = cifar10['train'][4]['img']
x0 = img_to_tensor(start_im).squeeze()
for t in [0, 200, 400, 600, 800]:

```

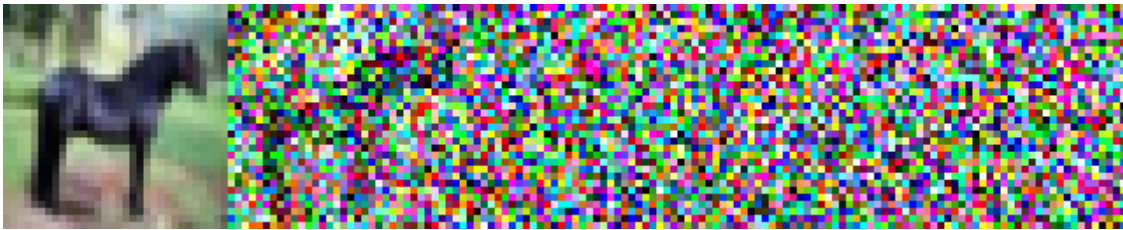
```

x = q_xt_x0(x0, torch.tensor(t, dtype=torch.long)) # TODO move type to gather
ims.append(tensor_to_image(x))

image = Image.new('RGB', size=(32*5, 32))
for i, im in enumerate(ims):
    image.paste(im, ((i%5)*32, 0))
image.resize((32*4*5, 32*4), Image.NEAREST)

```

[24]:



[6]: *#@title Unet Definition*

```

import math
from typing import Optional, Tuple, Union, List

import torch
from torch import nn

# A fancy activation function
class Swish(nn.Module):
    """
    ### Swish actiavation function
    $$x \cdot \sigma(x)$$
    """

    def forward(self, x):
        return x * torch.sigmoid(x)

# The time embedding
class TimeEmbedding(nn.Module):
    """
    ### Embeddings for $t$
    """

    def __init__(self, n_channels: int):
        """
        * `n_channels` is the number of dimensions in the embedding
        """
        super().__init__()
        self.n_channels = n_channels

```

```

    # First linear layer
    self.lin1 = nn.Linear(self.n_channels // 4, self.n_channels)
    # Activation
    self.act = Swish()
    # Second linear layer
    self.lin2 = nn.Linear(self.n_channels, self.n_channels)

def forward(self, t: torch.Tensor):
    # Create sinusoidal position embeddings
    # [same as those from the transformer](../transformers/
    ↪positional_encoding.html)
    #
    # \begin{align}
    # PE^{(1)}_{t,i} &= \sin\Bigg(\frac{t}{10000^{\frac{i}{d-1}}}\Bigg) \backslash
    # PE^{(2)}_{t,i} &= \cos\Bigg(\frac{t}{10000^{\frac{i}{d-1}}}\Bigg)
    # \end{align}
    #
    # where  $d$  is `half_dim`
    half_dim = self.n_channels // 8
    emb = math.log(10_000) / (half_dim - 1)
    emb = torch.exp(torch.arange(half_dim, device=t.device) * -emb)
    emb = t[:, None] * emb[None, :]
    emb = torch.cat((emb.sin(), emb.cos()), dim=1)

    # Transform with the MLP
    emb = self.act(self.lin1(emb))
    emb = self.lin2(emb)

    #
    return emb

# Residual blocks include 'skip' connections
class ResidualBlock(nn.Module):
    """
    ### Residual block
    A residual block has two convolution layers with group normalization.
    Each resolution is processed with two residual blocks.
    """

    def __init__(self, in_channels: int, out_channels: int, time_channels: int, ↵
    ↪n_groups: int = 32):
        """
        * `in_channels` is the number of input channels
        * `out_channels` is the number of input channels
        * `time_channels` is the number channels in the time step ( $t$ ) ↵
        ↪embeddings

```

```

        * `n_groups` is the number of groups for [group normalization](.././
↳normalization/group_norm/index.html)
        """
        super().__init__()
        # Group normalization and the first convolution layer
        self.norm1 = nn.GroupNorm(n_groups, in_channels)
        self.act1 = Swish()
        self.conv1 = nn.Conv2d(in_channels, out_channels, kernel_size=(3, 3),
↳padding=(1, 1))

        # Group normalization and the second convolution layer
        self.norm2 = nn.GroupNorm(n_groups, out_channels)
        self.act2 = Swish()
        self.conv2 = nn.Conv2d(out_channels, out_channels, kernel_size=(3, 3),
↳padding=(1, 1))

        # If the number of input channels is not equal to the number of output
↳channels we have to
        # project the shortcut connection
        if in_channels != out_channels:
            self.shortcut = nn.Conv2d(in_channels, out_channels,
↳kernel_size=(1, 1))
        else:
            self.shortcut = nn.Identity()

        # Linear layer for time embeddings
        self.time_emb = nn.Linear(time_channels, out_channels)

    def forward(self, x: torch.Tensor, t: torch.Tensor):
        """
        * `x` has shape `[batch_size, in_channels, height, width]`
        * `t` has shape `[batch_size, time_channels]`
        """
        # First convolution layer
        h = self.conv1(self.act1(self.norm1(x)))
        # Add time embeddings
        h += self.time_emb(t)[:, :, None, None]
        # Second convolution layer
        h = self.conv2(self.act2(self.norm2(h)))

        # Add the shortcut connection and return
        return h + self.shortcut(x)

# Ahh yes, magical attention...
class AttentionBlock(nn.Module):
    """
    ### Attention block

```

This is similar to [transformer multi-head attention](../transformers/mha.html).

```
"""

def __init__(self, n_channels: int, n_heads: int = 1, d_k: int = None,
n_groups: int = 32):
    """
    * `n_channels` is the number of channels in the input
    * `n_heads` is the number of heads in multi-head attention
    * `d_k` is the number of dimensions in each head
    * `n_groups` is the number of groups for [group normalization](../normalization/group_norm/index.html)
    """
    super().__init__()

    # Default `d_k`
    if d_k is None:
        d_k = n_channels
    # Normalization layer
    self.norm = nn.GroupNorm(n_groups, n_channels)
    # Projections for query, key and values
    self.projection = nn.Linear(n_channels, n_heads * d_k * 3)
    # Linear layer for final transformation
    self.output = nn.Linear(n_heads * d_k, n_channels)
    # Scale for dot-product attention
    self.scale = d_k ** -0.5
    #
    self.n_heads = n_heads
    self.d_k = d_k

def forward(self, x: torch.Tensor, t: Optional[torch.Tensor] = None):
    """
    * `x` has shape `[batch_size, in_channels, height, width]`
    * `t` has shape `[batch_size, time_channels]`
    """
    # `t` is not used, but it's kept in the arguments because for the
    ↪ attention layer function signature
    # to match with `ResidualBlock`.
    _ = t
    # Get shape
    batch_size, n_channels, height, width = x.shape
    # Change `x` to shape `[batch_size, seq, n_channels]`
    x = x.view(batch_size, n_channels, -1).permute(0, 2, 1)
    # Get query, key, and values (concatenated) and shape it to
    ↪ `[batch_size, seq, n_heads, 3 * d_k]`
    qkv = self.projection(x).view(batch_size, -1, self.n_heads, 3 * self.
    ↪ d_k)
```

```

        # Split query, key, and values. Each of them will have shape
        ↪ `[batch_size, seq, n_heads, d_k]`
        q, k, v = torch.chunk(qkv, 3, dim=-1)
        # Calculate scaled dot-product  $\frac{Q K^{\top}}{\sqrt{d_k}}$ 
        attn = torch.einsum('bihd,bjhd->bijh', q, k) * self.scale
        # Softmax along the sequence dimension
        ↪  $\underset{\text{seq}}{\text{softmax}}\left(\frac{Q K^{\top}}{\sqrt{d_k}}\right)$ 
        attn = attn.softmax(dim=1)
        # Multiply by values
        res = torch.einsum('bijh,bjhd->bihd', attn, v)
        # Reshape to `[batch_size, seq, n_heads * d_k]`
        res = res.view(batch_size, -1, self.n_heads * self.d_k)
        # Transform to `[batch_size, seq, n_channels]`
        res = self.output(res)

        # Add skip connection
        res += x

        # Change to shape `[batch_size, in_channels, height, width]`
        res = res.permute(0, 2, 1).view(batch_size, n_channels, height, width)

        #
        return res

```

```

class DownBlock(nn.Module):
    """
    ### Down block
    This combines `ResidualBlock` and `AttentionBlock`. These are used in the
    ↪ first half of U-Net at each resolution.
    """

    def __init__(self, in_channels: int, out_channels: int, time_channels: int,
        ↪ has_attn: bool):
        super().__init__()
        self.res = ResidualBlock(in_channels, out_channels, time_channels)
        if has_attn:
            self.attn = AttentionBlock(out_channels)
        else:
            self.attn = nn.Identity()

    def forward(self, x: torch.Tensor, t: torch.Tensor):
        x = self.res(x, t)
        x = self.attn(x)
        return x

```



```

class UpBlock(nn.Module):
    """
    ### Up block
    This combines `ResidualBlock` and `AttentionBlock`. These are used in the
    ↪second half of U-Net at each resolution.
    """

    def __init__(self, in_channels: int, out_channels: int, time_channels: int,
    ↪has_attn: bool):
        super().__init__()
        # The input has `in_channels + out_channels` because we concatenate the
    ↪output of the same resolution
        # from the first half of the U-Net
        self.res = ResidualBlock(in_channels + out_channels, out_channels,
    ↪time_channels)
        if has_attn:
            self.attn = AttentionBlock(out_channels)
        else:
            self.attn = nn.Identity()

    def forward(self, x: torch.Tensor, t: torch.Tensor):
        x = self.res(x, t)
        x = self.attn(x)
        return x


class MiddleBlock(nn.Module):
    """
    ### Middle block
    It combines a `ResidualBlock`, `AttentionBlock`, followed by another
    ↪`ResidualBlock`.
    This block is applied at the lowest resolution of the U-Net.
    """

    def __init__(self, n_channels: int, time_channels: int):
        super().__init__()
        self.res1 = ResidualBlock(n_channels, n_channels, time_channels)
        self.attn = AttentionBlock(n_channels)
        self.res2 = ResidualBlock(n_channels, n_channels, time_channels)

    def forward(self, x: torch.Tensor, t: torch.Tensor):
        x = self.res1(x, t)
        x = self.attn(x)
        x = self.res2(x, t)
        return x

```

```

class Upsample(nn.Module):
    """
    ### Scale up the feature map by  $2 \times$ 
    """

    def __init__(self, n_channels):
        super().__init__()
        self.conv = nn.ConvTranspose2d(n_channels, n_channels, (4, 4), (2, 2),
↪(1, 1))

    def forward(self, x: torch.Tensor, t: torch.Tensor):
        # `t` is not used, but it's kept in the arguments because for the
↪attention layer function signature
        # to match with `ResidualBlock`.
        _ = t
        return self.conv(x)

class Downsample(nn.Module):
    """
    ### Scale down the feature map by  $\frac{1}{2} \times$ 
    """

    def __init__(self, n_channels):
        super().__init__()
        self.conv = nn.Conv2d(n_channels, n_channels, (3, 3), (2, 2), (1, 1))

    def forward(self, x: torch.Tensor, t: torch.Tensor):
        # `t` is not used, but it's kept in the arguments because for the
↪attention layer function signature
        # to match with `ResidualBlock`.
        _ = t
        return self.conv(x)

# The core class definition (aka the important bit)
class UNet(nn.Module):
    """
    ## U-Net
    """

    def __init__(self, image_channels: int = 3, n_channels: int = 64,
                 ch_mults: Union[Tuple[int, ...], List[int]] = (1, 2, 2, 4),
                 is_attn: Union[Tuple[bool, ...], List[int]] = (False, False,
↪True, True),
                 n_blocks: int = 2):
        """
        * `image_channels` is the number of channels in the image.  $3 \times$  for RGB.

```

```

    * `n_channels` is number of channels in the initial feature map that we
    ↪transform the image into
    * `ch_mults` is the list of channel numbers at each resolution. The
    ↪number of channels is `ch_mults[i] * n_channels`
    * `is_attn` is a list of booleans that indicate whether to use
    ↪attention at each resolution
    * `n_blocks` is the number of `UpDownBlocks` at each resolution
    """
    super().__init__()

    # Number of resolutions
    n_resolutions = len(ch_mults)

    # Project image into feature map
    self.image_proj = nn.Conv2d(image_channels, n_channels, kernel_size=(3,
    ↪3), padding=(1, 1))

    # Time embedding layer. Time embedding has `n_channels * 4` channels
    self.time_emb = TimeEmbedding(n_channels * 4)

    # ##### First half of U-Net - decreasing resolution
    down = []
    # Number of channels
    out_channels = in_channels = n_channels
    # For each resolution
    for i in range(n_resolutions):
        # Number of output channels at this resolution
        out_channels = in_channels * ch_mults[i]
        # Add `n_blocks`
        for _ in range(n_blocks):
            down.append(DownBlock(in_channels, out_channels, n_channels *
    ↪4, is_attn[i]))
            in_channels = out_channels
        # Down sample at all resolutions except the last
        if i < n_resolutions - 1:
            down.append(Downsample(in_channels))

    # Combine the set of modules
    self.down = nn.ModuleList(down)

    # Middle block
    self.middle = MiddleBlock(out_channels, n_channels * 4, )

    # ##### Second half of U-Net - increasing resolution
    up = []
    # Number of channels
    in_channels = out_channels

```

```

        # For each resolution
        for i in reversed(range(n_resolutions)):
            # `n_blocks` at the same resolution
            out_channels = in_channels
            for _ in range(n_blocks):
                up.append(UpBlock(in_channels, out_channels, n_channels * 4,
↪is_attn[i]))
            # Final block to reduce the number of channels
            out_channels = in_channels // ch_mults[i]
            up.append(UpBlock(in_channels, out_channels, n_channels * 4,
↪is_attn[i]))
            in_channels = out_channels
            # Up sample at all resolutions except last
            if i > 0:
                up.append(Upsample(in_channels))

        # Combine the set of modules
        self.up = nn.ModuleList(up)

        # Final normalization and convolution layer
        self.norm = nn.GroupNorm(8, n_channels)
        self.act = Swish()
        self.final = nn.Conv2d(in_channels, image_channels, kernel_size=(3, 3),
↪padding=(1, 1))

    def forward(self, x: torch.Tensor, t: torch.Tensor):
        """
        * `x` has shape `[batch_size, in_channels, height, width]`
        * `t` has shape `[batch_size]`
        """

        # Get time-step embeddings
        t = self.time_emb(t)

        # Get image projection
        x = self.image_proj(x)

        # `h` will store outputs at each resolution for skip connection
        h = [x]
        # First half of U-Net
        for m in self.down:
            x = m(x, t)
            h.append(x)

        # Middle (bottom)
        x = self.middle(x, t)

```

```

        # Second half of U-Net
        for m in self.up:
            if isinstance(m, Upsample):
                x = m(x, t)
            else:
                # Get the skip connection from first half of U-Net and
                ↪ concatenate
                s = h.pop()
                x = torch.cat((x, s), dim=1)
                #
                x = m(x, t)

        # Final normalization and convolution
        return self.final(self.act(self.norm(x)))

```

```

[7]: # Let's see it in action on dummy data:

# A dummy batch of 10 3-channel 32px images
x = torch.randn(10, 3, 32, 32)

# 't' - what timestep are we on
t = torch.tensor([50.], dtype=torch.long)

# Define the unet model
UNET = UNet()

# The forward pass (takes both x and t)
model_output = UNET(x, t)

# The output shape matches the input.
model_output.shape

```

```

[7]: torch.Size([10, 3, 32, 32])

```

```

[47]: # Create the model
UNET = UNet(n_channels=32).cuda()

# Set up some parameters
n_steps = 500
beta = torch.linspace(0.0001, 0.03, n_steps).cuda()
alpha = 1. - beta
alpha_bar = torch.cumprod(alpha, dim=0)

# Modified to return the noise itself as well
def q_xt_x0(x0, t):
    mean = gather(alpha_bar, t) ** 0.5 * x0
    var = 1-gather(alpha_bar, t)

```

```

    eps = torch.randn_like(x0).to(x0.device)
    return mean + (var ** 0.5) * eps, eps # also returns noise

# Training params
batch_size = 128 # Lower this if hitting memory issues
lr = 3e-4 # Explore this - might want it lower when training on the full dataset

losses = [] # Store losses for later plotting

dataset = cifar10['train'].select(range(10000)) # to use a 10k subset for demo

optim = torch.optim.AdamW(unet.parameters(), lr=lr) # Optimizer

for i in tqdm(range(0, len(dataset)-batch_size, batch_size)): # Run through the
    ↪dataset
    ims = [dataset[idx]['img'] for idx in range(i,i+batch_size)] # Fetch some
    ↪images
    tims = [img_to_tensor(im).cuda() for im in ims] # Convert to tensors
    x0 = torch.cat(tims) # Combine into a batch
    t = torch.randint(0, n_steps, (batch_size,), dtype=torch.long).cuda() #
    ↪Random 't's
    xt, noise = q_xt_x0(x0, t) # Get the noised images (xt) and the noise (our
    ↪target)
    pred_noise = unet(xt.float(), t) # Run xt through the network to get its
    ↪predictions
    loss = F.mse_loss(noise.float(), pred_noise) # Compare the predictions with
    ↪the targets
    losses.append(loss.item()) # Store the loss for later viewing
    optim.zero_grad() # Zero the gradients
    loss.backward() # Backpropagate the loss (computes and store gradients)
    optim.step() # Update the network parameters (using those gradients)

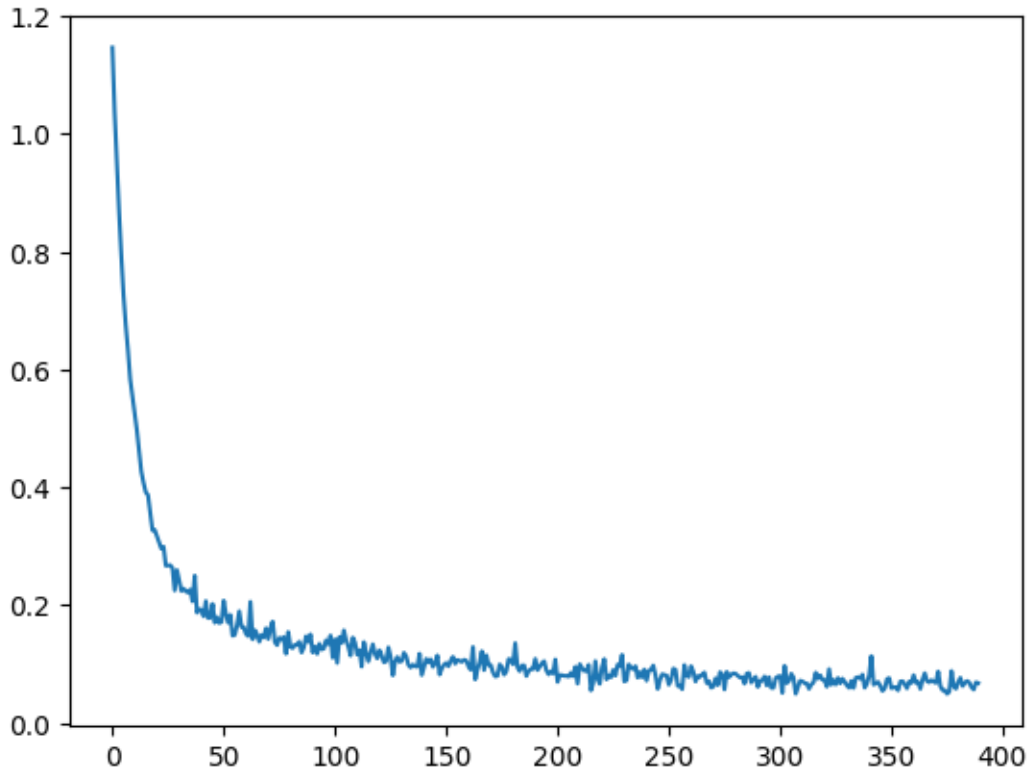
```

```
0%|          | 0/390 [00:00<?, ?it/s]
```

NB: You'd typically create a dataloader to run through a dataset in batches like we did above, but I couldn't remember the syntax and forgot to fix it! Hopefully the way I did it here is clear enough :)

```
[48]: from matplotlib import pyplot as plt
      plt.plot(losses)
```

```
[48]: [<matplotlib.lines.Line2D at 0x7a36606c7700>]
```



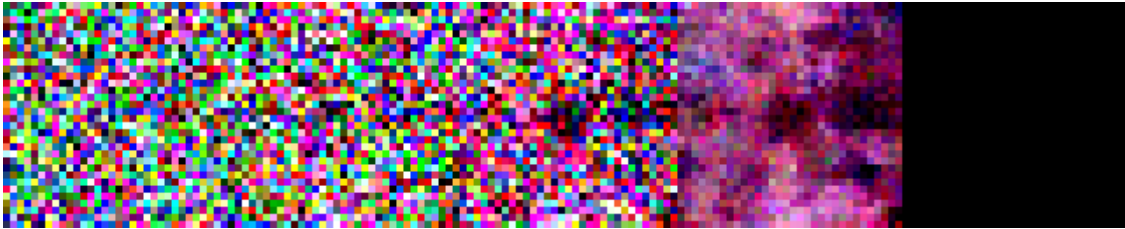
```
[49]: def p_xt(xt, noise, t):
    alpha_t = gather(alpha, t)
    alpha_bar_t = gather(alpha_bar, t)
    eps_coef = (1 - alpha_t) / (1 - alpha_bar_t) ** .5
    mean = 1 / (alpha_t ** 0.5) * (xt - eps_coef * noise) # Note minus sign
    var = gather(beta, t)
    eps = torch.randn(xt.shape, device=xt.device)
    return mean + (var ** 0.5) * eps

x = torch.randn(1, 3, 32, 32).cuda() # Start with random noise
ims = []
for i in range(n_steps):
    t = torch.tensor(n_steps-i-1, dtype=torch.long).cuda()
    with torch.no_grad():
        pred_noise = unet(x.float(), t.unsqueeze(0))
        x = p_xt(x, pred_noise, t.unsqueeze(0))
        if i%160 == 0:
            ims.append(tensor_to_image(x.cpu()))

image = Image.new('RGB', size=(32*5, 32))
for i, im in enumerate(ims[:5]):
    image.paste(im, ((i%5)*32, 0))
```

```
image.resize((32*4*5, 32*4), Image.NEAREST)
```

[49]:



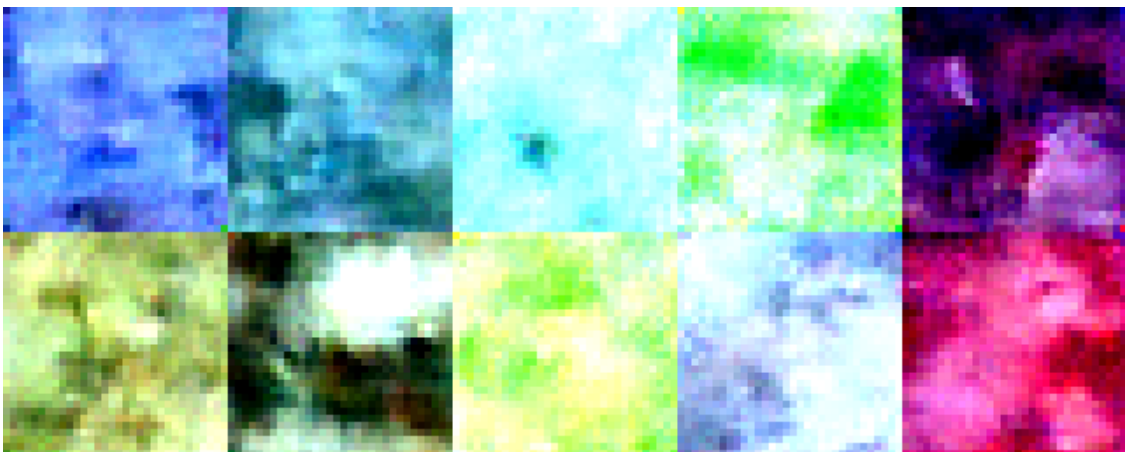
Perfect? No... Oh well, let's try a few more times:

```
[50]: #@title Make and show 10 examples:
x = torch.randn(10, 3, 32, 32).cuda() # Start with random noise
ims = []
for i in range(n_steps):
    t = torch.tensor(n_steps-i-1, dtype=torch.long).cuda()
    with torch.no_grad():
        pred_noise = unet(x.float(), t.unsqueeze(0))
        x = p_xt(x, pred_noise, t.unsqueeze(0))

for i in range(10):
    ims.append(tensor_to_image(x[i].unsqueeze(0).cpu()))

image = Image.new('RGB', size=(32*5, 32*2))
for i, im in enumerate(ims):
    image.paste(im, ((i%5)*32, 32*(i//5)))
image.resize((32*4*5, 32*4*2), Image.NEAREST)
```

[50]:




```
[51]: #@title Start with a heavily noised horse (t=50, top left = starting point):
horse = cifar10['train'][4]['img']
x0 = img_to_tensor(horse)
x = torch.cat([q_xt_x0(x0.cuda(), torch.tensor(50, dtype=torch.long).cuda())[0],
↳for _ in range(10)] )
example_start = q_xt_x0(x0.cuda(), torch.tensor(50, dtype=torch.long).cuda())[0]
print(x.shape)
ims = []
for i in range(50, n_steps):
    t = torch.tensor(n_steps-i-1, dtype=torch.long).cuda()
    with torch.no_grad():
        pred_noise = unet(x.float(), t.unsqueeze(0))
        x = p_xt(x, pred_noise, t.unsqueeze(0))

for i in range(10):
    ims.append(tensor_to_image(x[i].unsqueeze(0).cpu()))

image = Image.new('RGB', size=(32*5, 32*2))
for i, im in enumerate(ims):
    image.paste(im, ((i%5)*32, 32*(i//5)))
    if i==0: image.paste(tensor_to_image(example_start.unsqueeze(0).cpu()),
↳((i%5)*32, 32*(i//5))) # Show the heavily noised starting point top left
image.resize((32*4*5, 32*4*2), Image.NEAREST)
```

torch.Size([10, 3, 32, 32])

[51]:

