

Lab 3, Context

```
In [1]: from typing import Dict, Tuple
        from tqdm import tqdm
        import torch
        import torch.nn as nn
        import torch.nn.functional as F
        from torch.utils.data import DataLoader
        from torchvision import models, transforms
        from torchvision.utils import save_image, make_grid
        import matplotlib.pyplot as plt
        from matplotlib.animation import FuncAnimation, PillowWriter
        import numpy as np
        from IPython.display import HTML
        from diffusion utilities import *
```

Setting Things Up

```
In [2]: class ContextUnet(nn.Module):
        def __init__(self, in_channels, n_feat=256, n_cfeat=10, height=128):
            super(ContextUnet, self).__init__()

            # number of input channels, number of intermediate features
            self.in_channels = in_channels
            self.n_feat = n_feat
            self.n_cfeat = n_cfeat
            self.h = height #assume h == w. must be divisible by 4,

            # Initialize the initial convolutional layer
            self.init_conv = ResidualConvBlock(in_channels, n_feat, 1)

            # Initialize the down-sampling path of the U-Net with two layers
            self.down1 = UnetDown(n_feat, n_feat) # down1 #[16, 16, 16]
            self.down2 = UnetDown(n_feat, 2 * n_feat) # down2 #[32, 32, 16]

            # original: self.to_vec = nn.Sequential(nn.AvgPool2d(7),
            self.to_vec = nn.Sequential(nn.AvgPool2d((4)), nn.GELU())

            # Embed the timestep and context labels with a one-layer FC
            self.timeembed1 = EmbedFC(1, 2*n_feat)
            self.timeembed2 = EmbedFC(1, 1*n_feat)
            self.contextembed1 = EmbedFC(n_cfeat, 2*n_feat)
            self.contextembed2 = EmbedFC(n_cfeat, 1*n_feat)

            # Initialize the up-sampling path of the U-Net with three layers
            self.up0 = nn.Sequential(
                nn.ConvTranspose2d(2 * n_feat, 2 * n_feat, self.h//4,
                nn.GroupNorm(8, 2 * n_feat), # normalize
                nn.ReLU(),
            )
            self.up1 = UnetUp(4 * n_feat, n_feat)
            self.up2 = UnetUp(2 * n_feat, n_feat)

            # Initialize the final convolutional layers to map to the output
            self.out = nn.Sequential(
```

```

        nn.Conv2d(2 * n_feat, n_feat, 3, 1, 1), # reduce numl
        nn.GroupNorm(8, n_feat), # normalize
        nn.ReLU(),
        nn.Conv2d(n_feat, self.in_channels, 3, 1, 1), # map i
    )

def forward(self, x, t, c=None):
    """
    x : (batch, n_feat, h, w) : input image
    t : (batch, n_cfeat)      : time step
    c : (batch, n_classes)    : context label
    """
    # x is the input image, c is the context label, t is the

    # pass the input image through the initial convolutional
    x = self.init_conv(x)
    # pass the result through the down-sampling path
    down1 = self.down1(x)      #[10, 256, 8, 8]
    down2 = self.down2(down1)  #[10, 256, 4, 4]

    # convert the feature maps to a vector and apply an activ
    hiddenvec = self.to_vec(down2)

    # mask out context if context_mask == 1
    if c is None:
        c = torch.zeros(x.shape[0], self.n_cfeat).to(x)

    # embed context and timestep
    cemb1 = self.contextembed1(c).view(-1, self.n_feat * 2, 1)
    temb1 = self.timeembed1(t).view(-1, self.n_feat * 2, 1, 1)
    cemb2 = self.contextembed2(c).view(-1, self.n_feat, 1, 1)
    temb2 = self.timeembed2(t).view(-1, self.n_feat, 1, 1)
    #print(f"uunet forward: cemb1 {cemb1.shape}. temb1 {temb1

    up1 = self.up0(hiddenvec)
    up2 = self.up1(cemb1*up1 + temb1, down2) # add and mult:
    up3 = self.up2(cemb2*up2 + temb2, down1)
    out = self.out(torch.cat((up3, x), 1))
    return out

```

```
In [3]: # hyperparameters

# diffusion hyperparameters
timesteps = 500
beta1 = 1e-4
beta2 = 0.02

# network hyperparameters
device = torch.device("cuda:0" if torch.cuda.is_available() else
n_feat = 64 # 64 hidden dimension feature
n_cfeat = 5 # context vector is of size 5
height = 16 # 16x16 image
save_dir = './weights/'

# training hyperparameters
batch_size = 100
n_epoch = 32
lr=1e-3

In [4]: # construct DDPM noise schedule
b_t = (beta2 - beta1) * torch.linspace(0, 1, timesteps + 1, device=
a_t = 1 - b_t
ab_t = torch.cumsum(a_t.log(), dim=0).exp()
ab_t[0] = 1

In [5]: # construct model
nn_model = ContextUnet(in_channels=3, n_feat=n_feat, n_cfeat=n_cf
```

Context

```
In [6]: # reset neural network
nn_model = ContextUnet(in_channels=3, n_feat=n_feat, n_cfeat=n_cf

# re setup optimizer
optim = torch.optim.Adam(nn_model.parameters(), lr=lr)
```

```
# training with context code
# set into train mode
nn_model.train()

for ep in range(n_epoch):
    print(f'epoch {ep}')

    # linearly decay learning rate
    optim.param_groups[0]['lr'] = lr*(1-ep/n_epoch)

    pbar = tqdm(dataloader, mininterval=2 )
    for x, c in pbar: # x: images c: context
        optim.zero_grad()
        x = x.to(device)
        c = c.to(x)

        # randomly mask out c
        context_mask = torch.bernoulli(torch.zeros(c.shape[0]) + 0.9).to(device)
        c = c * context_mask.unsqueeze(-1)

        # perturb data
        noise = torch.randn_like(x)
```

```

t = torch.randint(1, timesteps + 1, (x.shape[0],)).to(device)
x_pert = perturb_input(x, t, noise)

# use network to recover noise
pred_noise = nn_model(x_pert, t / timesteps, c=c)

# loss is mean squared error between the predicted and true noise
loss = F.mse_loss(pred_noise, noise)
loss.backward()

optim.step()

# save model periodically
if ep%4==0 or ep == int(n_epoch-1):
    if not os.path.exists(save_dir):
        os.mkdir(save_dir)
    torch.save(nn_model.state_dict(), save_dir + f"context_model_{ep}.pth")
    print('saved model at ' + save_dir + f"context model {ep}.pth")
In [7]: # load in pretrain model weights and set to eval mode
nn_model.load_state_dict(torch.load(f"{save_dir}/context_model_{t}"))
nn_model.eval()
print("Loaded in Context Model")

```

Loaded in Context Model

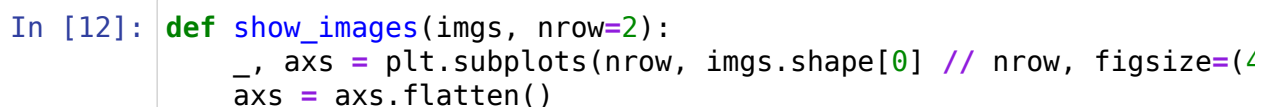
Sampling with context

```

In [8]: # helper function; removes the predicted noise (but adds some noise)
def denoise_add_noise(x, t, pred_noise, z=None):
    if z is None:
        z = torch.randn_like(x)
    noise = b_t.sqrt()[t] * z
    mean = (x - pred_noise * ((1 - a_t[t]) / (1 - ab_t[t])).sqrt())
    return mean + noise

```

```
In [*]: # visualize samples with randomly selected context
plt.clf()
ctx = F.one_hot(torch.randint(0, 5, (32,)), 5).to(device=device).
samples, intermediate = sample_ddpm_context(32, ctx)
animation_ddpm_context = plot_sample(intermediate, 32, 4, save_dir,
HTML(animation_ddpm_context.to_ishtml()))
```



```

for img, ax in zip(imgs, axs):
    img = (img.permute(1, 2, 0).clip(-1, 1).detach().cpu().numpy()
    ax.set_xticks([])
    ax.set_yticks([])
    ax.imshow(img)

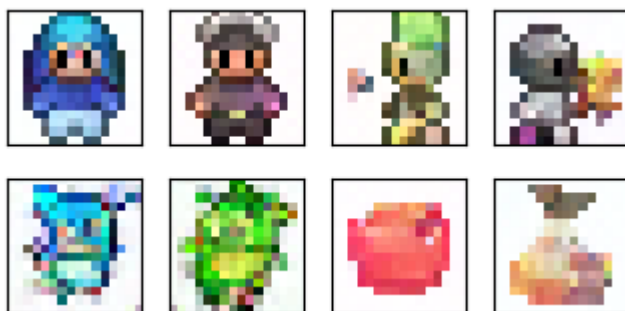
```

```

In [13]: # user defined context
ctx = torch.tensor([
    # hero, non-hero, food, spell, side-facing
    [1,0,0,0,0],
    [1,0,0,0,0],
    [0,0,0,0,1],
    [0,0,0,0,1],
    [0,1,0,0,0],
    [0,1,0,0,0],
    [0,0,1,0,0],
    [0,0,1,0,0],
]).float().to(device)
samples, _ = sample_ddpm_context(ctx.shape[0], ctx)
show_images(samples)

```

sampling timestep 1



```

In [14]: # mix of defined context
ctx = torch.tensor([
    # hero, non-hero, food, spell, side-facing
    [1,0,0,0,0], #human
    [1,0,0.6,0,0],
    [0,0,0.6,0.4,0],
    [1,0,0,0,1],
    [1,1,0,0,0],
    [1,0,0,1,0]
]).float().to(device)
samples, _ = sample_ddpm_context(ctx.shape[0], ctx)
show_images(samples)

```

sampling timestep 1



Acknowledgments

Sprites by ElvGames, [FrootsnVeggies](https://zrghr.itch.io/foots-and-veggies-culinary-pixels) (<https://zrghr.itch.io/foots-and-veggies-culinary-pixels>) and [kyrise](https://kyrise.itch.io/) (<https://kyrise.itch.io/>)

This code is modified from, <https://github.com/cloneofsimo/minDiffusion> (<https://github.com/cloneofsimo/minDiffusion>)

Diffusion model is based on [Denoising Diffusion Probabilistic Models](https://arxiv.org/abs/2006.11239) (<https://arxiv.org/abs/2006.11239>) and [Denoising Diffusion Implicit Models](https://arxiv.org/abs/2010.02502) (<https://arxiv.org/abs/2010.02502>)

In []: