*Logan Anderson*

# Contents

# Tasks

## 3.1 Conda Environment

I set up my environment as instructed in the assignment with no issues. The *environment.yml* file will be submitted with the code.

## 3.2 Comparing CIFAR10 with Custom Dataset

### 3.2.1 Steps of Implementing

**1. Prepare CIFAR10 Dataset**: I prepare CIFAR10 as in the code snippet given (but I randomize the classes that are chosen). Additionally, I altered the code so that it actually retrieves 10 elements per class as the base code did not guarantee this.

```python
1. def load_data(data_set=CIFAR10, class_num_filter=5, total_classes=10):
2.     # loads data
3.     transform = transforms.Compose([transforms.ToTensor()])
4.     data = data_set(root='./data', train=True, download=True, transform=transform)
5.
6.     # filters into 5 classes
7.     classes = sorted(sample(range(total_classes), class_num_filter))
8.
9.     # retrieve 10 elements per class
10.    indices = []
11.    class_cnt = {class_: 0 for class_ in classes}
12.    total_required = NUM_SAMPLES // class_num_filter
13.    for i, (_,label) in enumerate(data):
14.      if label not in classes: continue
15.      if class_cnt[label] < total_required:
16.        class_cnt[label] += 1
17.        indices.append(i)
18.
19.    subset_data = Subset(data, indices)
20.    return subset_data
```

**2. Create Custom Dataset:** I took many pictures of a calendar I have in my apartment from several different angles. Below are some examples of images I took for this task:
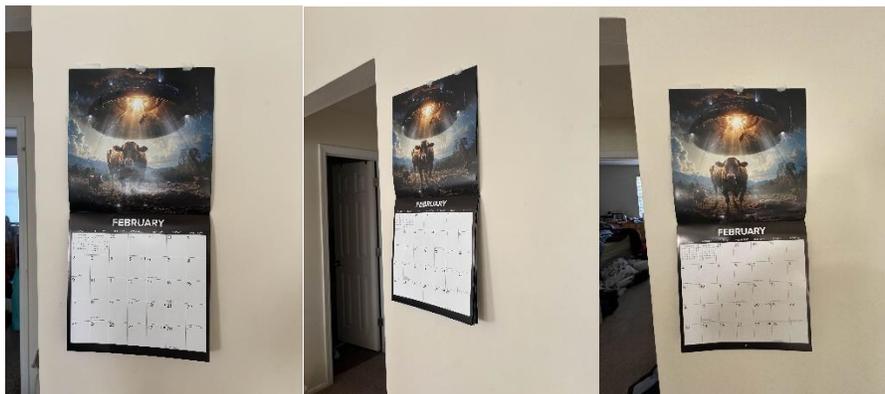
With minor modifications to the skeleton code, I implemented this custom dataset class shown below (I added a method to retrieve the base image given a specific index and modified the *__next__* method to be circular using modulus). Also of note, I return the class as its corresponding base image index for ease of comparison between the base images:

```python
1.  class Custom_Dataset(Dataset):
2.      def __init__(self, root, transform=None):
3.          self.root = root
4.          self.image_paths = [os.path.join(root, img) for img in os.listdir(root)]
5.          self.transform = transform
6.
7.      def get_base_image(self, index):
8.          img_path = self.image_paths[index]
9.          transform = transforms.ToTensor()
10.         return transform(Image.open(img_path).convert("RGB"))
11.
12.     def __len__(self):
13.         return len(self.image_paths)
14.
15.     def __getitem__(self, index):
16.         img_path = self.image_paths[index % len(self.image_paths)]
17.         img = Image.open(img_path).convert("RGB")
18.         if self.transform:
19.             img = self.transform(img)
20.         return img, index % len(self.image_paths)
```

**3. Apply Data Augmentation**: I next applied several different transformations to the input data which I separate between the pre-processing and then dataloader transformations:

*Preprocessing*

- Transformations that I perform to the base images themselves for use in the dataloader
- Purpose: Reduce redundant computation in the dataloader
- Transformations
    - *Resize*: resize the image to a fixed smaller size, ensuring all images are the same size and reducing the overall quantity of pixels to handle

```python
transforms.Resize(new_size, antialias=True) # new_size=256
```

*Dataloader Transformations*

- Transformations I perform upon each successive access of the given iterable
- Purpose: Transform original data in random manner to add uncertainty to inputs and effectively generate extra inputs for the data loader
- *Transformations*

- *RandomHorizontalFlip*: Flips the image horizontally with a given a probability p, allowing for the image to be recognized even if the given camera image is flipped
- *GaussianBlur*: Blurs an image according to a Gaussian distribution to add some distortion to the image, so regardless of clarity, the image can be recognized
- *RandomAffine*: Applies a random affine homography to an image, so the image can be identified even at different angles
- *ToTensor*: A simple transformation that transforms the PIL image representation into a torch tensor representation
- *Normalize*: Translates the pixels so that they are normalized around a given point for ease of use in neural network architectures

```
1. custom_transform = transforms.Compose([
2.                     transforms.RandomHorizontalFlip(p=.2),
3.                     transforms.GaussianBlur(3, (.3,1.5)),
4.                     transforms.RandomAffine(degrees=20, translate=(.1,.1),
scale=(.9, 1.1)),
5.                     transforms.ToTensor(),
6.                     transforms.Normalize((0.5,),(0.5,))])
```

In order to generate the additional samples, similar to the skeleton code, I simply ran:

```
augmented_images = [custom_dataset[i][0] for i in range(50)]
```

## 4. Compare Datasets

The sizes are compared and produce the output shown to the right:

```
1. # compare size of datasets
2. print(f"Size of CIFAR10 subset: {NUM_SAMPLES} images")         # 50
3. print(f"Size of Custom dataset: {len(custom_dataset)} images")  # 20
```

CIFAR10 Subset Samples

FIGURE 2: SELECT CIFAR10 IMAGES
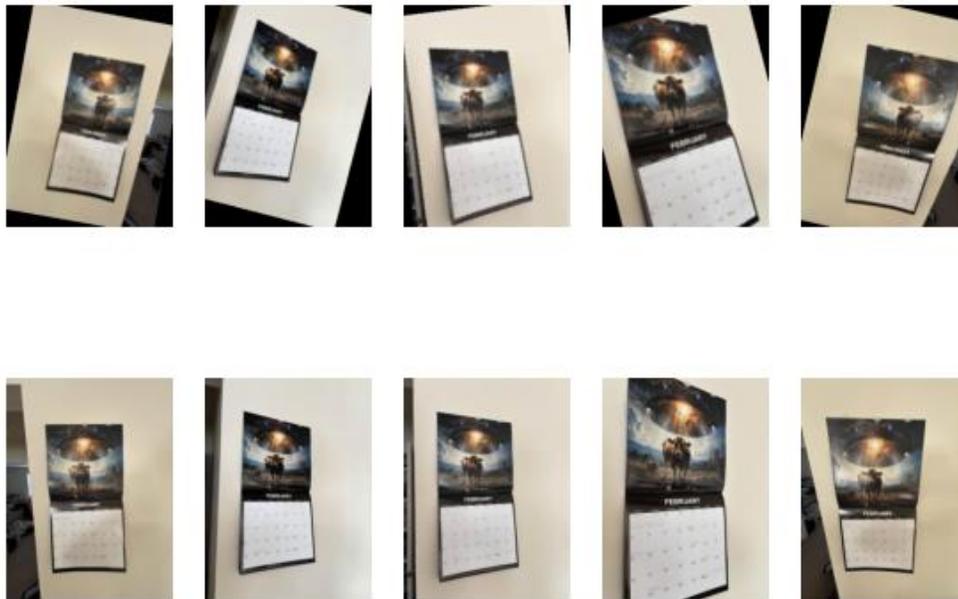


Custom Transformed Images

FIGURE 3: SELECT TRANSFORMED CUSTOM IMAGES (TOP) AND BASE IMAGES (BOTTOM)

## 3.3 Using DataLoader for Parallel Processing

I first wrap the custom dataset I created in a dataloader just as the CIFAR10 dataset was. I set $batch\_size = 4$ and $num\_workers = 2$:

```
custom_loader = DataLoader(custom_dataset, batch_size=4, num_workers=2)
```

I then plot 4 images from a single batch of the dataloader:

Custom Transformed Batch of Images



FIGURE 4: SINGLE BATCH OF CUSTOM IMAGES USING DATALOADER

## Performance Comparison

First I check how long it takes to manually augment 1000 images and determine this to be around **3.579 seconds** when running the following code with $num\_images = 1000$

```
[dataset[i][0] for i in range(num_images)]
```

I then try creating 1000 images using the dataloader with varying values for *batch_size* and *num_workers*. The results of this are shown below, and I compute the average time with 5 runs each:

| batch_size\num_workers | 0 | 2 | 4 | 6 |
|---|---|---|---|---|
| 4 | 3.611 s | 6.542 s | 7.665 s | 8.413 s |
| 16 | 2.229 s | 4.694 s | 5.014 s | 6.073 s |
| 64 | 1.054 s | 2.408 s | 2.669 s | 3.187 s |

TABLE 1: TIMES TO GENERATE 1000 SAMPLES WITH VARYING BATCH_SIZE AND NUM_WORKERS

Next, I compare different maximum values of the actual RGB channels using the following code:

```
 1. # get images before and after transformation
 2. for batch_idx, (batch, labels) in enumerate(custom_loader):
 3.    if batch_idx != 0: break
 4.    for trans_img,label in zip(batch, labels):
 5.      base_img = custom_dataset.get_base_image(label)
 6.      max_vals_base = base_img.max(dim=1)[0].max(dim=1)[0]
 7.      max_vals_trans = trans_img.max(dim=1)[0].max(dim=1)[0]
 8.
 9.      print(f"Max values of base: {max_vals_base}")
10.      print(f"Max values of transformed: {max_vals_trans}")
11.      print()
```

The output is as follows for a batch size of 4:

```
Max values of base: tensor([1.0000, 0.9647, 0.9647])
Max values of transformed: tensor([0.9294, 0.8039, 0.8902])

Max values of base: tensor([0.9216, 0.9412, 0.9255])
Max values of transformed: tensor([0.7804, 0.8275, 0.7882])

Max values of base: tensor([0.9686, 0.9451, 0.9020])
Max values of transformed: tensor([0.9137, 0.8667, 0.7804])

Max values of base: tensor([1.0000, 0.9765, 0.9686])
Max values of transformed: tensor([0.9922, 0.9059, 0.8902])
```

## 3.4 Exploring Random Seed and Reproducability

The test is conducted as follows where I first run 2 times with $SEED\_SET = True$ and 2 times with $SEED\_SET = False$:
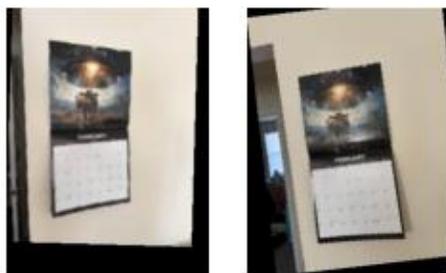
```
1. custom_transform = DEFINED_CUSTOM_TRANSFORMATIONS
2. custom_dataset = Custom_Dataset(root=CUSTOM_OUTPUT_PATH, transform=custom_transform)
3.
4. # get custom loader
5. custom_loader = DataLoader(custom_dataset, batch_size=2, shuffle=True)
6.
7. # plot image
8. if SEED_SET:
9.    plot_data(custom_loader, path_name="custom_seed_set", img_title="Custom Seed Set Images",
num_side_by_side=2, figsize=(3,3))
10. else:
11.    plot_data(custom_loader, path_name=f"custom_seed_not_set{random.randint(0,5)}",
img_title="Custom Seed Not Set Images", num_side_by_side=2, figsize=(3,3))
```

When running the program with a seed set for 2 successive iterations, the result is as follows:



FIGURE 5: CUSTOM SEED SET IMAGES (ITERATION 1 ON THE RIGHT, ITERATION 2 ON THE LEFT)

When a seed is not set, the following is the result:



FIGURE 6: CUSTOM SEED NOT SET IMAGES (ITERATION 1 ON THE RIGHT, ITERATION 2 ON THE LEFT)

Thus, it is clear that when a seed is set, the images produced are identical, i.e. any training done will have randomness but no variation between runs. This is incredible useful when trying to test whether a given model works for predicting outputs as well as any sort of debugging in general with deep learning. When a seed isn't set, the result is different each time and thus is not as good at reproducing the same results.

# Source Code

**environment.yml**

```
 1. name: ece60146
 2. channels:
 3.   - pytorch
 4.   - defaults
 5. dependencies:
 6.   - _libgcc_mutex=0.1=main
 7.   - _openmp_mutex=5.1=1_gnu
 8.   - blas=1.0=mkl
 9.   - brotli-python=1.0.9=py310h6a678d5_9
10.   - bzip2=1.0.8=h5eee18b_6
11.   - ca-certificates=2024.12.31=h06a4308_0
12.   - certifi=2024.12.14=py310h06a4308_0
13.   - charset-normalizer=3.3.2=pyhd3eb1b0_0
14.   - cuda-cudart=12.4.127=h99ab3db_0
15.   - cuda-cudart_linux-64=12.4.127=hd681fbe_0
16.   - cuda-cupti=12.4.127=h6a678d5_1
17.   - cuda-libraries=12.4.1=h06a4308_1
18.   - cuda-nvrtc=12.4.127=h99ab3db_1
19.   - cuda-nvtx=12.4.127=h6a678d5_1
```

```
20.    - cuda-opencl=12.4.127=h6a678d5_0
21.    - cuda-runtime=12.4.1=hb982923_0
22.    - cuda-version=12.4=hbda6634_3
23.    - ffmpeg=4.3=hf484d3e_0
24.    - filelock=3.13.1=py310h06a4308_0
25.    - freetype=2.12.1=h4a9f257_0
26.    - giflib=5.2.2=h5eee18b_0
27.    - gmp=6.2.1=h295c915_3
28.    - gmpy2=2.1.2=py310heeb90bb_0
29.    - gnutls=3.6.15=he1e5248_0
30.    - idna=3.7=py310h06a4308_0
31.    - intel-openmp=2023.1.0=hdb19cb5_46306
32.    - jinja2=3.1.4=py310h06a4308_1
33.    - jpeg=9e=h5eee18b_3
34.    - lame=3.100=h7b6447c_0
35.    - lcms2=2.16=hb9589c4_0
36.    - ld_impl_linux-64=2.40=h12ee557_0
37.    - lerc=4.0.0=h6a678d5_0
38.    - libcublas=12.4.5.8=h99ab3db_1
39.    - libcufft=11.2.1.3=h99ab3db_1
40.    - libcufile=1.9.1.3=h99ab3db_1
41.    - libcurand=10.3.5.147=h99ab3db_1
42.    - libcusolver=11.6.1.9=h99ab3db_1
43.    - libcusparse=12.3.1.170=h99ab3db_1
44.    - libdeflate=1.22=h5eee18b_0
45.    - libffi=3.4.4=h6a678d5_1
46.    - libgcc-ng=11.2.0=h1234567_1
47.    - libgomp=11.2.0=h1234567_1
48.    - libiconv=1.16=h5eee18b_3
49.    - libidn2=2.3.4=h5eee18b_0
50.    - libjpeg-turbo=2.0.0=h9bf148f_0
51.    - libnpp=12.2.5.30=h99ab3db_1
52.    - libnvfatbin=12.4.127=h7934f7d_2
53.    - libnvjitlink=12.4.127=h99ab3db_1
54.    - libnvjpeg=12.3.1.117=h6a678d5_1
55.    - libpng=1.6.39=h5eee18b_0
56.    - libstdcxx-ng=11.2.0=h1234567_1
57.    - libtasn1=4.19.0=h5eee18b_0
58.    - libtiff=4.5.1=hffd6297_1
59.    - libunistring=0.9.10=h27cfd23_0
60.    - libuuid=1.41.5=h5eee18b_0
61.    - libwebp=1.3.2=h11a3e52_0
62.    - libwebp-base=1.3.2=h5eee18b_1
63.    - llvm-openmp=14.0.6=h9e868ea_0
64.    - lz4-c=1.9.4=h6a678d5_1
65.    - markupsafe=2.1.3=py310h5eee18b_1
66.    - mkl=2023.1.0=h213fc3f_46344
67.    - mkl-service=2.4.0=py310h5eee18b_2
68.    - mkl_fft=1.3.11=py310h5eee18b_0
69.    - mkl_random=1.2.8=py310h1128e8f_0
70.    - mpc=1.1.0=h10f8cd9_1
71.    - mpfr=4.0.2=hb69a4c5_1
72.    - mpmath=1.3.0=py310h06a4308_0
73.    - ncurses=6.4=h6a678d5_0
74.    - nettle=3.7.3=hbbd107a_1
75.    - networkx=3.4.2=py310h06a4308_0
76.    - numpy=1.26.4=py310h5f9d8c6_0
77.    - numpy-base=1.26.4=py310hb5e798b_0
78.    - ocl-icd=2.3.2=h5eee18b_1
79.    - openh264=2.1.1=h4ff587b_0
80.    - openjpeg=2.5.2=he7f1fd0_0
81.    - openssl=3.0.15=h5eee18b_0
82.    - pillow=11.0.0=py310hcea889d_1
83.    - pip=24.2=py310h06a4308_0
```

```
 84.    - pysocks=1.7.1=py310h06a4308_0
 85.    - python=3.10.16=he870216_1
 86.    - pytorch=2.5.1=py3.10_cuda12.4_cudnn9.1.0_0
 87.    - pytorch-cuda=12.4=hc786d27_7
 88.    - pytorch-mutex=1.0=cuda
 89.    - pyyaml=6.0.2=py310h5eee18b_0
 90.    - readline=8.2=h5eee18b_0
 91.    - requests=2.32.3=py310h06a4308_1
 92.    - setuptools=75.1.0=py310h06a4308_0
 93.    - sqlite=3.45.3=h5eee18b_0
 94.    - sympy=1.13.3=py310h06a4308_0
 95.    - tbb=2021.8.0=hdb19cb5_0
 96.    - tk=8.6.14=h39e8969_0
 97.    - torchtriton=3.1.0=py310
 98.    - torchvision=0.20.1=py310_cu124
 99.    - typing_extensions=4.12.2=py310h06a4308_0
100.    - tzdata=2025a=h04d1e81_0
101.    - urllib3=2.3.0=py310h06a4308_0
102.    - wheel=0.44.0=py310h06a4308_0
103.    - xz=5.4.6=h5eee18b_1
104.    - yaml=0.2.5=h7b6447c_0
105.    - zlib=1.2.13=h5eee18b_1
106.    - zstd=1.5.6=hc292b87_0
107. prefix: /home/the_linux_ruler/anaconda3/envs/ece60146
108.
```

## hw2.py

```python
 1. import torch
 2. import torch.backends
 3. import torch.backends.cudnn
 4. from torchvision import transforms
 5. from torchvision.datasets import CIFAR10
 6. from torch.utils.data import DataLoader, Subset, Dataset
 7. from PIL import Image
 8. from random import sample
 9. import random
10. import os
11. import matplotlib.pyplot as plt
12. from math import ceil
13. from time import time
14. import inspect
15. import csv
16. import numpy as np
17.
18. NUM_SAMPLES = 50
19. NUM_SAMPLES_PARALLEL = 1000
20. NUM_ITERATIONS_PER_TEST = 5
21. NUM_WORKERS = [0,2,4,6]
22. BATCH_SIZES = [4,16,64]
23. CUSTOM_IMAGE_SIZE = 256
24. DEFINED_CUSTOM_TRANSFORMATIONS = transforms.Compose([
25.                                                 # transforms.Resize(CUSTOM_IMAGE_SIZE,
antialias=True),
26.                                                 transforms.RandomHorizontalFlip(p=.2),
27.                                                 transforms.GaussianBlur(3, (.3,1.5)),
28.                                                 transforms.RandomAffine(degrees=20,
translate=(.1,.1), scale=(.9, 1.1)),
29.                                                 transforms.ToTensor(),
30.                                                 transforms.Normalize((0.5,),(0.5,))])
31. RESULTS_PATH = '/home/the_linux_ruler/ECE60146/HW2/results/'
32. CUSTOM_INPUT_PATH = "/home/the_linux_ruler/ECE60146/HW2/custom_images"
```

*Logan Anderson*

```python
33. CUSTOM_OUTPUT_PATH = "/home/the_linux_ruler/ECE60146/HW2/custom_images2"
34. RUN_SEED_TEST = True
35. SEED_SET = False
36. DEFINED_SEED = 60146
37.
38.
39. # ----------------------------------------------------------------------------------------
40. def current_time():
41.     return round(time() - start_time, 3)
42.
43. def current_time_str():
44.     curr_time = current_time()
45.     if curr_time // 3600:
46.         return f"{int(curr_time) // 3600} hr\t{(int(curr_time) % 3600) // 60} min\t
{round((curr_time - int(curr_time)) + (int(curr_time) % 60),3)} sec"
47.     elif curr_time // 60:
48.         return f"{int(curr_time) // 60} min\t {round((curr_time - int(curr_time)) + (int(curr_time)
% 60), 3)} sec"
49.     else:
50.         return f"{curr_time} sec"
51.
52. def time_diff(time_start):
53.     return round(time() - time_start, 3)
54.
55. def time_func(func, *args, print_params=[], print_extra="", get_time=False,
suppress_output=False):
56.     time_start = time()
57.     params_str = ",".join([f"{names_in_caller(param, depth=3)[0]} = {param}" for param in
print_params])
58.     extra_string = f" with parameter(s) {params_str}" if len(print_params) != 0 else ""
59.     extra_string += " " + print_extra
60.     if not suppress_output: print(f"{func.__name__}{extra_string} starts at time
{current_time_str()}")
61.     ret_val = func(*args)
62.     time_taken = time_diff(time_start)
63.     if not suppress_output: print(f"{func.__name__}{extra_string} ends at time
{current_time_str()} in {time_taken} seconds")
64.     if get_time: return ret_val, time_taken
65.     return ret_val
66.
67. def namestr(obj, namespace):
68.     return [name for name in namespace if namespace[name] is obj]
69.
70. def names_in_caller(obj, depth=2) -> list[str]:
71.     f = inspect.currentframe()
72.     for _ in range(depth): f = f.f_back
73.     return namestr(obj, f.f_locals)
74.
75. def resize_image(input_path, output_path, new_size):
76.     # Open the image file
77.     img = Image.open(input_path)
78.
79.     resize_transform = transforms.Compose([
80.         transforms.Resize(new_size, antialias=True),  # Resize image
81.         transforms.ToTensor()  # Convert to tensor
82.     ])
83.
84.     resized_img = resize_transform(img)
85.     resized_img = transforms.ToPILImage()(resized_img)
86.
87.     # Save resized image
88.     resized_img.save(output_path)
89.
90. def set_seed(seed_num):
```

*Logan Anderson*

```
 91.    random.seed(seed_num)
 92.    torch.manual_seed(seed_num)
 93.    torch.cuda.manual_seed(seed_num)
 94.    np.random.seed(seed_num)
 95.    torch.backends.cudnn.deterministic = True
 96.    torch.backends.cudnn.benchmark = False
 97.    os.environ['PYTHONHASHSEED'] = str(seed_num)
 98. # ------------------------------------------------------------------------------------
 99.
100. class Custom_Dataset(Dataset):
101.    def __init__(self, root, transform=None):
102.      self.root = root
103.      self.image_paths = [os.path.join(root, img) for img in os.listdir(root)]
104.      self.transform = transform
105.
106.    def get_base_image(self, index):
107.      img_path = self.image_paths[index]
108.      transform = transforms.ToTensor()
109.      return transform(Image.open(img_path).convert("RGB"))
110.
111.    def __len__(self):
112.      return len(self.image_paths)
113.
114.    def __getitem__(self, index):
115.      img_path = self.image_paths[index % len(self.image_paths)]
116.      img = Image.open(img_path).convert("RGB")
117.      if self.transform:
118.        img = self.transform(img)
119.      return img, index % len(self.image_paths)
120.
121. def load_data(data_set=CIFAR10, class_num_filter=5, total_classes=10):
122.    # loads data
123.    transform = transforms.Compose([transforms.ToTensor()])
124.    data = data_set(root='./data', train=True, download=True, transform=transform)
125.
126.    # filters into 5 classes
127.    classes = sorted(sample(range(total_classes), class_num_filter))
128.
129.    # retrieve 10 elements per class
130.    indices = []
131.    class_cnt = {class_: 0 for class_ in classes}
132.    total_required = NUM_SAMPLES // class_num_filter
133.    for i, (_,label) in enumerate(data):
134.      if label not in classes: continue
135.      if class_cnt[label] < total_required:
136.        class_cnt[label] += 1
137.        indices.append(i)
138.
139.    subset_data = Subset(data, indices)
140.    return subset_data
141.
142. def transform_back_to_rgb(arr):
143.    arr = np.array(arr)
144.    min_old = np.min(arr)
145.    max_old = np.max(arr)
146.
147.    arr = ((arr - min_old) / (max_old - min_old)) * (255 - 0) + 0
148.    return arr.astype(int)
149.
150.
151. def plot_data(data_loader:DataLoader, path_name:str, img_title:str, num_side_by_side=5,
total_batches=1, figsize=None, plot_base_image=False, dataset=None):
152.    num_up_down = ceil(data_loader.batch_size * total_batches / num_side_by_side)
153.
```

```python
154.    if plot_base_image: num_up_down *= 2
155.
156.    # configures figure size if specified
157.    if figsize == None: fig, axs = plt.subplots(num_up_down, num_side_by_side)
158.    else: fig, axs = plt.subplots(num_up_down, num_side_by_side, figsize=figsize)
159.
160.    # adjusts axs to work for 1-D array
161.    if num_up_down == 1: axs = [axs]
162.
163.    targets = []
164.    for i,(batch, target) in enumerate(data_loader):
165.      offset = i * data_loader.batch_size
166.      if i >= total_batches: break
167.      targets.append(target)
168.      for j,img in enumerate(batch):
169.        y,x = (offset + j) // num_side_by_side, (offset + j) % num_side_by_side
170.        # if num_up_down == 1: x,y = y,x
171.        axs[y][x].imshow(transform_back_to_rgb(img.permute(1,2,0).numpy()))
172.        axs[y][x].axis("off")
173.
174.    if plot_base_image:
175.      for i,target in enumerate(targets):
176.        for j,target_scalar in enumerate(target):
177.          base_img = dataset.get_base_image(target_scalar)
178.          loc = offset + i + j
179.          y,x = loc // num_side_by_side, loc % num_side_by_side
180.          axs[y][x].imshow(transform_back_to_rgb(base_img.permute(1,2,0).numpy()))
181.          axs[y][x].axis("off")
182.
183.    plt.suptitle(img_title)
184.
185.    fig.savefig(f"{RESULTS_PATH}{path_name}")
186.
187. def augment_images(dataset, num_images):
188.    return [dataset[i][0] for i in range(num_images)]
189.
190. def get_images_parallel(data_loader:DataLoader, num_images:int):
191.    imgs = []
192.    iter_data = iter(data_loader)
193.    for _ in range(num_images // data_loader.batch_size):
194.      try:
195.        imgs.append(next(iter_data))
196.      except StopIteration:
197.        iter_data = iter(data_loader)
198.        imgs.append(next(iter_data))
199.
200.    return imgs
201.
202. def create_csv(filename, headers, data):
203.    with open(f"{RESULTS_PATH}{filename}.csv", mode='w', newline='', encoding='utf-8') as file:
204.      writer = csv.writer(file)
205.      writer.writerow(headers)
206.      for data_point in data:
207.        writer.writerow(data_point)
208.
209. if __name__ == "__main__":
210.    start_time = time()
211.
212.    # set seed for same result (for testing purposes)
213.    if SEED_SET: set_seed(DEFINED_SEED)
214.    if RUN_SEED_TEST:
215.      custom_transform = DEFINED_CUSTOM_TRANSFORMATIONS
216.      custom_dataset = Custom_Dataset(root=CUSTOM_OUTPUT_PATH, transform=custom_transform)
217.
```

```
218.     # get custom loader
219.     custom_loader = DataLoader(custom_dataset, batch_size=2, shuffle=True)
220.
221.     # plot image
222.     if SEED_SET:
223.         plot_data(custom_loader, path_name="custom_seed_set", img_title="Custom Seed Set Images",
num_side_by_side=2, figsize=(3,3))
224.     else:
225.         plot_data(custom_loader, path_name=f"custom_seed_not_set{random.randint(0,5)}",
img_title="Custom Seed Not Set Images", num_side_by_side=2, figsize=(3,3))
226.
227.     exit()
228.
229.     data = load_data()
230.
231.     # resize custom images
232.     print("Resizing Images")
233.     for filename in os.listdir(CUSTOM_INPUT_PATH):
234.         resize_image(os.path.join(CUSTOM_INPUT_PATH, filename), os.path.join(CUSTOM_OUTPUT_PATH,
filename), CUSTOM_IMAGE_SIZE)
235.
236.     # create custom transformation and apply with custom dataset
237.     custom_transform = DEFINED_CUSTOM_TRANSFORMATIONS
238.     custom_dataset = Custom_Dataset(root=CUSTOM_OUTPUT_PATH, transform=custom_transform)
239.     augmented_images = augment_images(custom_dataset, NUM_SAMPLES)
240.
241.     # compare size of datasets
242.     print(f"Size of CIFAR10 subset: {NUM_SAMPLES} images")
243.     print(f"Size of Custom dataset: {len(custom_dataset)} images")
244.
245.     # load subset of CIFAR10
246.     cifar_loader = DataLoader(data, batch_size=5, shuffle=False)
247.
248.     # load custom dataset into dataloader
249.     custom_loader = DataLoader(custom_dataset, batch_size=4, num_workers=2)
250.
251.     # get images before and after transformation
252.     for batch_idx, (batch, labels) in enumerate(custom_loader):
253.         if batch_idx != 0: break
254.         for trans_img,label in zip(batch, labels):
255.             base_img = custom_dataset.get_base_image(label)
256.             max_vals_base = base_img.max(dim=1)[0].max(dim=1)[0]
257.             max_vals_trans = trans_img.max(dim=1)[0].max(dim=1)[0]
258.
259.             print(f"Max values of base: {max_vals_base}")
260.             print(f"Max values of transformed: {max_vals_trans}")
261.             print()
262.
263.     # create custom loader with a batch size of 5
264.     custom_loader2 = DataLoader(custom_dataset, batch_size=5)
265.
266.     # plot various examples
267.     plot_data(cifar_loader, "CIFAR10_Examples", "CIFAR10 Subset Samples", figsize=(7,3))
268.     plot_data(custom_loader2, path_name="custom_examples", img_title="Custom Transformed Images",
num_side_by_side=5, figsize=(7,5), plot_base_image=True, dataset=custom_dataset)
269.     plot_data(custom_loader, path_name="custom_example_batch", img_title="Custom Transformed
Batch of Images", num_side_by_side=4, figsize=(7,3))
270.
271.     # augment 1000 images and calculate how long this takes
272.     _, time_taken = time_func(augment_images, custom_dataset, NUM_SAMPLES_PARALLEL,
get_time=True)
273.
274.     # compute times for different batch sizes and num_workers (1000)
275.     all_points = []
```

*Logan Anderson*

```
276.    for batch_size in BATCH_SIZES:
277.        row = []
278.        row.append(batch_size)
279.        for num_workers in NUM_WORKERS:
280.            print(f"Running calculation for batch_size={batch_size}, num_workers={num_workers} at
time {current_time_str()}")
281.            times_taken = []
282.            for _ in range(NUM_ITERATIONS_PER_TEST):
283.                custom_loader = DataLoader(custom_dataset, batch_size=batch_size,
num_workers=num_workers)
284.                _, time_taken = time_func(get_images_parallel, custom_loader, NUM_SAMPLES_PARALLEL,
print_extra=f"with num_workers={num_workers}, batch_size={batch_size}", get_time=True,
suppress_output=True)
285.                times_taken.append(time_taken)
286.            row.append(sum(times_taken) / len(times_taken))
287.        all_points.append(row)
288.
289.    # generate csv file to store this information
290.    create_csv("time_comparisons2", ["batch_size\\num_workers", *NUM_WORKERS], all_points)
```