## Start by importing necessary packages

You will begin by importing necessary libraries for this notebook. Run the cell below to do so.

# PyTorch and Intro to Training

```
!pip install thop
import math
import numpy as np
import torch
import torch.nn as nn
import torch.optim as optim
from torchvision import datasets, transforms
import thop
import matplotlib.pyplot as plt
from tqdm import tqdm
import time
```

```
Collecting thop
  Downloading thop-0.1.1.post2209072238-py3-none-any.whl.metadata (2.7
kB)
Requirement already satisfied: torch in
/usr/local/lib/python3.11/dist-packages (from thop) (2.5.1+cu121)
Requirement already satisfied: filelock in
/usr/local/lib/python3.11/dist-packages (from torch->thop) (3.16.1)
Requirement already satisfied: typing-extensions>=4.8.0 in
/usr/local/lib/python3.11/dist-packages (from torch->thop) (4.12.2)
Requirement already satisfied: networkx in
/usr/local/lib/python3.11/dist-packages (from torch->thop) (3.4.2)
Requirement already satisfied: jinja2 in
/usr/local/lib/python3.11/dist-packages (from torch->thop) (3.1.5)
Requirement already satisfied: fsspec in
/usr/local/lib/python3.11/dist-packages (from torch->thop) (2024.10.0)
Requirement already satisfied: nvidia-cuda-nvrtc-cu12==12.1.105 in
/usr/local/lib/python3.11/dist-packages (from torch->thop) (12.1.105)
Requirement already satisfied: nvidia-cuda-runtime-cu12==12.1.105
in /usr/local/lib/python3.11/dist-packages (from torch->thop)
(12.1.105)
Requirement already satisfied: nvidia-cuda-cupti-cu12==12.1.105 in
/usr/local/lib/python3.11/dist-packages (from torch->thop) (12.1.105)
Requirement already satisfied: nvidia-cudnn-cu12==9.1.0.70 in
/usr/local/lib/python3.11/dist-packages (from torch->thop) (9.1.0.70)
Requirement already satisfied: nvidia-cublas-cu12==12.1.3.1 in
/usr/local/lib/python3.11/dist-packages (from torch->thop) (12.1.3.1)
Requirement already satisfied: nvidia-cufft-cu12==11.0.2.54 in
/usr/local/lib/python3.11/dist-packages (from torch->thop) (11.0.2.54)
Requirement already satisfied: nvidia-curand-cu12==10.3.2.106 in
```

```
/usr/local/lib/python3.11/dist-packages (from torch->thop)
(10.3.2.106)
Requirement already satisfied: nvidia-cusolver-cu12==11.4.5.107 in
/usr/local/lib/python3.11/dist-packages (from torch->thop)
(11.4.5.107)
Requirement already satisfied: nvidia-cusparse-cu12==12.1.0.106 in
/usr/local/lib/python3.11/dist-packages (from torch->thop)
(12.1.0.106)
Requirement already satisfied: nvidia-nccl-cu12==2.21.5 in
/usr/local/lib/python3.11/dist-packages (from torch->thop) (2.21.5)
Requirement already satisfied: nvidia-nvtx-cu12==12.1.105 in
/usr/local/lib/python3.11/dist-packages (from torch->thop) (12.1.105)
Requirement already satisfied: triton==3.1.0 in
/usr/local/lib/python3.11/dist-packages (from torch->thop) (3.1.0)
Requirement already satisfied: sympy==1.13.1 in
/usr/local/lib/python3.11/dist-packages (from torch->thop) (1.13.1)
Requirement already satisfied: nvidia-nvjitlink-cu12 in
/usr/local/lib/python3.11/dist-packages (from nvidia-cusolver-
cu12==11.4.5.107->torch->thop) (12.6.85)
Requirement already satisfied: mpmath<1.4,>=1.1.0 in
/usr/local/lib/python3.11/dist-packages (from sympy==1.13.1->torch-
>thop) (1.3.0)
Requirement already satisfied: MarkupSafe>=2.0 in
/usr/local/lib/python3.11/dist-packages (from jinja2->torch->thop)
(3.0.2)
Downloading thop-0.1.1.post2209072238-py3-none-any.whl (15 kB)
Installing collected packages: thop
Successfully installed thop-0.1.1.post2209072238
```

# Checking the torch version and CUDA access

Let's start off by checking the current torch version, and whether you have CUDA availablity.

```
print("torch is using version:", torch.__version__, "with CUDA=",
torch.cuda.is_available())
```

```
torch is using version: 2.5.1+cu121 with CUDA= True
```

By default, you will see CUDA = False, meaning that the Colab session does not have access to a GPU. To remedy this, click the Runtime menu on top and select "Change Runtime Type", then select "T4 GPU".

Re-run the import cell above, and the CUDA version / check. It should show now CUDA = True

Sometimes in Colab you get a message that your Session has crashed, if that happens you need to go to the Runtime menu on top and select "Restart session".

You won't be using the GPU just yet, but this prepares the instance for when you will.

**Please note that the GPU is a scarce resource which may not be available at all time. Additionally, there are also usage limits that you may run into (although not likely for this assignment). When that happens you need to try again later/next day/different time of the day. Another reason to start the assignment early!**

# A Brief Introduction to PyTorch

PyTorch, or torch, is a machine learning framework developed my Facebook AI Research, which competes with TensorFlow, JAX, Caffe and others.

Roughly speaking, these frameworks can be split into dynamic and static defintion frameworks.

**Static Network Definition:** The architecture and computation flow are defined simultaneously. The order and manner in which data flows through the layers are fixed upon definition. These frameworks also tend to declare parameter shapes implicitly via the compute graph. This is typical of TensorFlow and JAX.

**Dynamic Network Definition:** The architecture (layers/modules) is defined independently of the computation flow, often during the object's initialization. This allows for dynamic computation graphs where the flow of data can change during runtime based on conditions. Since the network exists independent of the compute graph, the parameter shapes must be declared explitly. PyTorch follows this approach.

All ML frameworks support automatic differentiation, which is necessary to train a model (i.e. perform back propagation).

Let's consider a typical pytorch module. Such modules will inherit from the torch.nn.Module class, which provides many built in functions such as a wrapper for `__call__`, operations to move the module between devices (e.g. `cuda()`, `cpu()`), data-type conversion (e.g. `half()`, `float()`), and parameter and child management (e.g. `state_dict()`, `parameters()`).

```python
# inherit from torch.nn.Module
class MyModule(nn.Module):
  # constructor called upon creation
  def __init__(self):
    # the module has to initialize the parent first, which is what
sets up the wrapper behavior
    super().__init__()

    # we can add sub-modules and parameters by assigning them to self
    self.my_param = nn.Parameter(torch.zeros(4,8)) # this is how you
define a raw parameter of shape 4x5
    self.my_sub_module = nn.Linear(8,12)      # this is how you
define a linear layer (tensorflow calls them Dense) of shape 8x12

    # we can also add lists of modules, for example, the sequential
layer
    self.net = nn.Sequential(  # this layer type takes in a collection
of modules rather than a list
        nn.Linear(4,4),
```

```python
        nn.Linear(4,8),
        nn.Linear(8,12)
    )

    # the above when calling self.net(x), will execute each module in
the order they appear in a list
    # it would be equivelent to x = self.net[2](self.net[1]
(self.net[0](x)))

    # you can also create a list that doesn't execute
    self.net_list = nn.ModuleList([
        nn.Linear(7,7),
        nn.Linear(7,9),
        nn.Linear(9,14)
    ])

    # sometimes you will also see constant variables added to the
module post init
    foo = torch.Tensor([4])
    self.register_buffer('foo', foo) # buffers allow .to(device, type)
to apply

  # let's define a forward function, which gets executed when calling
the module, and defines the forward compute graph
  def forward(self, x):

    # if x is of shape Bx4
    h1 =  x @ self.my_param # tensor-tensor multiplication uses the @
symbol
    # then h1 is now shape Bx8, because my_param is 4x8... 2x4 * 4x8 =
2x8

    h1 = self.my_sub_module(h1) # you execute a sub-module by calling
it
    # now, h1 is of shape Bx12, because my_sub_module was a 8x12
matrix

    h2 = self.net(x)
    # similarly, h2 is of shape Bx12, because that's the output of the
sequence
    # Bx4 -(4x4)-> Bx4 -(4x8)-> Bx8 -(8x12)-> Bx12

    # since h1 and h2 are the same shape, they can be added together
element-wise
    return h1 + h2
```

Then you can instantiate the module and perform a forward pass by calling it.

```python
# create the module
module = MyModule()
```

```python
# you can print the module to get a high-level summary of it
print("=== printing the module ===")
print(module)
print()
# notice that the sub-module name is in parenthesis, and so are the
list indicies

# let's view the shape of one of the weight tensors
print("my_sub_module weight tensor shape=",
module.my_sub_module.weight.shape)
# the above works because nn.Linear has a member called .weight
and .bias
# to view the shape of my_param, you would use module.my_param
# and to view the shape of the 2nd elment in net_list, you would use
module.net_list[1].weight

# you can iterate through all of the parameters via the state dict
print()
print("=== Listing parameters from the state_dict ===")
for key,value in module.state_dict().items():
  print(f"{key}: {value.shape}")
```

```
=== printing the module ===
MyModule(
  (my_sub_module): Linear(in_features=8, out_features=12, bias=True)
  (net): Sequential(
    (0): Linear(in_features=4, out_features=4, bias=True)
    (1): Linear(in_features=4, out_features=8, bias=True)
    (2): Linear(in_features=8, out_features=12, bias=True)
  )
  (net_list): ModuleList(
    (0): Linear(in_features=7, out_features=7, bias=True)
    (1): Linear(in_features=7, out_features=9, bias=True)
    (2): Linear(in_features=9, out_features=14, bias=True)
  )
)

my_sub_module weight tensor shape= torch.Size([12, 8])

=== Listing parameters from the state_dict ===
my_param: torch.Size([4, 8])
foo: torch.Size([1])
my_sub_module.weight: torch.Size([12, 8])
my_sub_module.bias: torch.Size([12])
net.0.weight: torch.Size([4, 4])
net.0.bias: torch.Size([4])
net.1.weight: torch.Size([8, 4])
net.1.bias: torch.Size([8])
net.2.weight: torch.Size([12, 8])
```

```
net.2.bias: torch.Size([12])
net_list.0.weight: torch.Size([7, 7])
net_list.0.bias: torch.Size([7])
net_list.1.weight: torch.Size([9, 7])
net_list.1.bias: torch.Size([9])
net_list.2.weight: torch.Size([14, 9])
net_list.2.bias: torch.Size([14])

# you can perform a forward pass by first creating a tensor to send
through
x = torch.zeros(2,4)
# then you call the module (this invokes MyModule.forward() )
y = module(x)

# then you can print the result and shape
print(y, y.shape)

tensor([[-0.3016, -0.3592, -0.3739, -0.1964,  0.6734,  0.0330, -
0.2751,  0.0842,
         -0.0989, -0.3719,  0.6114, -0.0605],
        [-0.3016, -0.3592, -0.3739, -0.1964,  0.6734,  0.0330, -
0.2751,  0.0842,
         -0.0989, -0.3719,  0.6114, -0.0605]], grad_fn=<AddBackward0>)
torch.Size([2, 12])
```

Please check the cell below to notice the following:

1. x above was created with the shape 2x4, and in the forward pass, it gets manipulated into a 2x12 tensor. This last dimension is explicit, while the first is called the batch dimmension, and only exists on data (a.k.a. activations). The output shape can be seen in the print statement from y.shape
2. You can view the shape of a tensor by using `.shape`, this is a very helpful trick for debugging tensor shape errors
3. In the output, there's a `grad_fn` component, this is the hook created by the forward trace to be used in back-propagation via automatic differentiation. The function name is `AddBackward`, because the last operation performed was `h1+h2`.

We might not always want to trace the compute graph though, such as during inference. In such cases, you can use the `torch.no_grad()` context manager.

```
# you can perform a forward pass by first creating a tensor to send
through
x = torch.zeros(2,4)
# then you call the module (this invokes MyModule.forward() )
with torch.no_grad():
  y = module(x)

# then you can print the result and shape
print(y, y.shape)
```

```
# notice how the grad_fn is no longer part of the output tensor,
that's because not_grad() disables the graph generation

tensor([[-0.3016, -0.3592, -0.3739, -0.1964,  0.6734,  0.0330, -
0.2751,  0.0842,
         -0.0989, -0.3719,  0.6114, -0.0605],
        [-0.3016, -0.3592, -0.3739, -0.1964,  0.6734,  0.0330, -
0.2751,  0.0842,
         -0.0989, -0.3719,  0.6114, -0.0605]]) torch.Size([2, 12])
```

Aside from passing a tensor through a model with the `no_grad()` context, you can also detach a tensor from the compute graph by calling `.detach()`. This will effectively make a copy of the original tensor, which allows it to be converted to numpy and visualized with matplotlib.

**Note:** Tensors with a `grad_fn` property cannot be plotted and must first be detached.

## Multi-Layer-Perceptron (MLP) Prediction of MNIST

With some basics out of the way, let's create a MLP for training MNIST. You can start by defining a simple torch model.

```
# Define the MLP model
class MLP(nn.Module):
    # define the constructor for the network
    def __init__(self):
        super().__init__()
        # the input projection layer - projects into d=128
        self.fc1 = nn.Linear(28*28, 128)
        # the first hidden layer - compresses into d=64
        self.fc2 = nn.Linear(128, 64)
        # the final output layer - splits into 10 classes (digits 0-9)
        self.fc3 = nn.Linear(64, 10)

    # define the forward pass compute graph
    def forward(self, x):
        # x is of shape BxHxW

        # we first need to unroll the 2D image using view
        # we set the first dim to be -1 meanining "everything else",
the reason being that x is of shape BxHxW, where B is the batch dim
        # we want to maintain different tensors for each training
sample in the batch, which means the output should be of shape BxF
where F is the feature dim
        x = x.view(-1, 28*28)
        # x is of shape Bx784

        # project-in and apply a non-linearity (ReLU activation
function)
        x = torch.relu(self.fc1(x))
```

```
        # x is of shape Bx128

        # middle-layer and apply a non-linearity (ReLU activation
function)
        x = torch.relu(self.fc2(x))
        # x is of shape Bx64

        # project out into the 10 classes
        x = self.fc3(x)
        # x is of shape Bx10
        return x
```

Before you can begin training, you have to do a little boiler-plate to load the dataset. From the previous assignment, you saw how a hosted dataset can be loaded with TensorFlow. With pytorch it's a little more complicated, as you need to manually condition the input data.

```
# define a transformation for the input images. This uses
torchvision.transforms, and .Compose will act similarly to
nn.Sequential
transform = transforms.Compose([
    transforms.ToTensor(), # first convert to a torch tensor
    transforms.Normalize((0.1307,), (0.3081,)) # then normalize the
input
])

# let's download the train and test datasets, applying the above
transform - this will get saved locally into ./data, which is in the
Colab instance
train_dataset = datasets.MNIST('./data', train=True, download=True,
transform=transform)
test_dataset = datasets.MNIST('./data', train=False,
transform=transform)

# we need to set the mini-batch (commonly referred to as "batch"), for
now we can use 64
batch_size = 64

# then we need to create a dataloader for the train dataset, and we
will also create one for the test dataset to evaluate performance
# additionally, we will set the batch size in the dataloader
train_loader = torch.utils.data.DataLoader(train_dataset,
batch_size=batch_size, shuffle=True)
test_loader = torch.utils.data.DataLoader(test_dataset,
batch_size=batch_size, shuffle=False)

# the torch dataloaders allow us to access the __getitem__ method,
which returns a tuple of (data, label)
# additionally, the dataloader will pre-colate the training samples
into the given batch_size
```

```
Downloading http://yann.lecun.com/exdb/mnist/train-images-idx3-
ubyte.gz
Failed to download (trying next):
<urlopen error [Errno 110] Connection timed out>

Downloading https://ossci-datasets.s3.amazonaws.com/mnist/train-
images-idx3-ubyte.gz
Downloading https://ossci-datasets.s3.amazonaws.com/mnist/train-
images-idx3-ubyte.gz to ./data/MNIST/raw/train-images-idx3-ubyte.gz

100%|████████████| 9.91M/9.91M [00:00<00:00, 15.0MB/s]

Extracting ./data/MNIST/raw/train-images-idx3-ubyte.gz to
./data/MNIST/raw

Downloading http://yann.lecun.com/exdb/mnist/train-labels-idx1-
ubyte.gz
Failed to download (trying next):
<urlopen error [Errno 110] Connection timed out>

Downloading https://ossci-datasets.s3.amazonaws.com/mnist/train-
labels-idx1-ubyte.gz
Downloading https://ossci-datasets.s3.amazonaws.com/mnist/train-
labels-idx1-ubyte.gz to ./data/MNIST/raw/train-labels-idx1-ubyte.gz

100%|████████████| 28.9k/28.9k [00:00<00:00, 459kB/s]

Extracting ./data/MNIST/raw/train-labels-idx1-ubyte.gz to
./data/MNIST/raw

Downloading http://yann.lecun.com/exdb/mnist/t10k-images-idx3-ubyte.gz
Failed to download (trying next):
<urlopen error [Errno 110] Connection timed out>

Downloading https://ossci-datasets.s3.amazonaws.com/mnist/t10k-images-
idx3-ubyte.gz
Downloading https://ossci-datasets.s3.amazonaws.com/mnist/t10k-images-
idx3-ubyte.gz to ./data/MNIST/raw/t10k-images-idx3-ubyte.gz

100%|████████████| 1.65M/1.65M [00:00<00:00, 4.19MB/s]

Extracting ./data/MNIST/raw/t10k-images-idx3-ubyte.gz to
./data/MNIST/raw

Downloading http://yann.lecun.com/exdb/mnist/t10k-labels-idx1-ubyte.gz
Failed to download (trying next):
<urlopen error [Errno 110] Connection timed out>

Downloading https://ossci-datasets.s3.amazonaws.com/mnist/t10k-labels-
idx1-ubyte.gz
```

Inspect the first element of the test_loader, and verify both the tensor shapes and data types. You can check the data-type with `.dtype`

**Question 1**

Edit the cell below to print out the first element shapes, dtype, and identify which is the training sample and which is the training label.

```python
# Get the first item
first_item = next(iter(test_loader))

# Extracting the sample and label from the first item
inputs, labels = first_item

# Print out the shapes, dtype, and identify the sample and label
print(f"Training sample shape: {inputs.shape}, dtype: {inputs.dtype}")
print(f"Training label shape: {labels.shape}, dtype: {labels.dtype}")

Training sample shape: torch.Size([64, 1, 28, 28]), dtype:
torch.float32
Training label shape: torch.Size([64]), dtype: torch.int64
```

Now that we have the dataset loaded, we can instantiate the MLP model, the loss (or criterion function), and the optimizer for training.

```python
# create the model
model = MLP()

# you can print the model as well, but notice how the activation
functions are missing. This is because they were called in the forward
pass
# and not declared in the constructor
print(model)

# you can also count the model parameters
param_count = sum([p.numel() for p in model.parameters()])
print(f"Model has {param_count:,} trainable parameters")

# for a critereon (loss) function, you will use Cross-Entropy Loss.
This is the most common criterion used for multi-class prediction,
```

```
# and is also used by tokenized transformer models it takes in an un-
normalized probability distribution (i.e. without softmax) over
# N classes (in our case, 10 classes with MNIST). This distribution is
then compared to an integer label which is < N.
# For MNIST, the prediction might be [-0.0056, -0.2044,  1.1726,
0.0859,  1.8443, -0.9627,  0.9785, -1.0752, 1.1376,  1.8220], with the
label 3.
# Cross-entropy can be thought of as finding the difference between
the predicted distribution and the one-hot distribution

criterion = nn.CrossEntropyLoss()

# then you can instantiate the optimizer. You will use Stochastic
Gradient Descent (SGD), and can set the learning rate to 0.1 with a
momentum
# factor of 0.5. the first input to the optimizer is the list of model
parameters, which is obtained by calling .parameters() on the model
object
optimizer = optim.SGD(model.parameters(), lr=0.01, momentum=0.5)

MLP(
  (fc1): Linear(in_features=784, out_features=128, bias=True)
  (fc2): Linear(in_features=128, out_features=64, bias=True)
  (fc3): Linear(in_features=64, out_features=10, bias=True)
)
Model has 109,386 trainable parameters
```

Finally, you can define a training, and test loop

```
# create an array to log the loss and accuracy
train_losses = []
train_steps = []
test_steps = []
test_losses = []
test_accuracy = []
current_step = 0  # Start with global step 0
current_epoch = 0 # Start with epoch 0

# declare the train function
def cpu_train(epoch, train_losses, steps, current_step):

    # set the model in training mode - this doesn't do anything for us
right now, but it is good practiced and needed with other layers such
as
    # batch norm and dropout
    model.train()

    # Create tqdm progress bar to help keep track of the training
progress
    pbar = tqdm(enumerate(train_loader), total=len(train_loader))
```

```python
    # loop over the dataset. Recall what comes out of the data loader,
and then by wrapping that with enumerate() we get an index into the
    # iterator list which we will call batch_idx
    for batch_idx, (data, target) in pbar:

        # during training, the first step is to zero all of the
gradients through the optimizer
        # this resets the state so that we can begin back propogation
with the updated parameters
        optimizer.zero_grad()

        # then you can apply a forward pass, which includes evaluating
the loss (criterion)
        output = model(data)
        loss = criterion(output, target)

        # given that you want to minimize the loss, you need to
call .backward() on the result, which invokes the grad_fn property
        loss.backward()

        # the backward step will automatically differentiate the model
and apply a gradient property to each of the parameters in the network
        # so then all you have to do is call optimizer.step() to apply
the gradients to the current parameters
        optimizer.step()

        # increment the step count
        current_step += 1

        # you should add some output to the progress bar so that you
know which epoch you are training, and what the current loss is
        if batch_idx % 100 == 0:

            # append the last loss value
            train_losses.append(loss.item())
            steps.append(current_step)

            desc = (f'Train Epoch: {epoch} [{batch_idx *
len(data)}/{len(train_loader.dataset)}'
                    f' ({100. * batch_idx / len(train_loader):.0f}%)]\
tLoss: {loss.item():.6f}')
            pbar.set_description(desc)

    return current_step

# declare a test function, this will help you evaluate the model
progress on a dataset which is different from the training dataset
# doing so prevents cross-contamination and misleading results due to
overfitting
```

```python
def cpu_test(test_losses, test_accuracy, steps, current_step):

    # put the model into eval mode, this again does not currently do
anything for you, but it is needed with other layers like batch_norm
    # and dropout
    model.eval()
    test_loss = 0
    correct = 0

    # Create tqdm progress bar
    pbar = tqdm(test_loader, total=len(test_loader),
desc="Testing...")

    # since you are not training the model, and do not need back-
propagation, you can use a no_grad() context
    with torch.no_grad():
        # iterate over the test set
        for data, target in pbar:
            # like with training, run a forward pass through the model
and evaluate the criterion
            output = model(data)
            test_loss += criterion(output, target).item() # you are
using .item() to get the loss value rather than the tensor itself

            # you can also check the accuracy by sampling the output -
you can use greedy sampling which is argmax (maximum probability)
            # in general, you would want to normalize the logits first
(the un-normalized output of the model), which is done via .softmax()
            # however, argmax is taking the maximum value, which will
be the same index for the normalized and un-normalized distributions
            # so we can skip a step and take argmax directly
            pred = output.argmax(dim=1, keepdim=True)
            correct += pred.eq(target.view_as(pred)).sum().item()

    test_loss /= len(test_loader)

    # append the final test loss
    test_losses.append(test_loss)
    test_accuracy.append(correct/len(test_loader.dataset))
    steps.append(current_step)

    print(f'\nTest set: Average loss: {test_loss:.4f}, Accuracy:
{correct}/{len(test_loader.dataset)}'
          f' ({100. * correct / len(test_loader.dataset):.0f}%)\n')

# train for 10 epochs
for epoch in range(0, 10):
    current_step = cpu_train(current_epoch, train_losses, train_steps,
current_step)
    cpu_test(test_losses, test_accuracy, test_steps, current_step)
```

```
        current_epoch += 1
```

**Question 2**

Using the skills you acquired in the previous assignment edit the cell below to use matplotlib to visualize the loss for training and validation for the first 10 epochs. They should be plotted on the same graph, labeled, and use a log-scale on the y-axis.

```python
# visualize the losses for the first 10 epochs
import matplotlib.pyplot as plt

# Declare the train and test functions (as you've already done)
# ...

# Initialize lists to store losses for plotting
train_losses_list = []
test_losses_list = []

# Train and test for 10 epochs
for epoch in range(0, 10):
    current_step = cpu_train(epoch, train_losses_list, train_steps,
current_step)
    cpu_test(test_losses_list, test_accuracy, test_steps,
current_step)
    current_epoch += 1

# Plot the training and testing losses
plt.figure(figsize=(10, 6))
plt.plot(train_steps, train_losses_list, label='Training Loss',
color='blue')
plt.plot(test_steps, test_losses_list, label='Validation Loss',
color='red')

# Set log scale for the y-axis
plt.yscale('log')

# Labeling the axes
plt.xlabel('Steps')
plt.ylabel('Loss (log scale)')
plt.title('Training and Validation Loss Over Epochs')

# Add a legend
plt.legend()

# Show the plot
plt.show()
```

```
Train Epoch: 0 [57600/60000 (96%)]    Loss: 0.300860: 100%|██████████|
938/938 [00:12<00:00, 73.12it/s]
Testing...: 100%|██████████| 157/157 [00:01<00:00, 82.99it/s]


Test set: Average loss: 0.2794, Accuracy: 9168/10000 (92%)


Train Epoch: 1 [57600/60000 (96%)]    Loss: 0.231814: 100%|██████████|
938/938 [00:12<00:00, 72.29it/s]
Testing...: 100%|██████████| 157/157 [00:01<00:00, 86.64it/s]


Test set: Average loss: 0.1966, Accuracy: 9416/10000 (94%)


Train Epoch: 2 [57600/60000 (96%)]    Loss: 0.149264: 100%|██████████|
938/938 [00:12<00:00, 73.40it/s]
Testing...: 100%|██████████| 157/157 [00:02<00:00, 72.22it/s]


Test set: Average loss: 0.1533, Accuracy: 9536/10000 (95%)


Train Epoch: 3 [57600/60000 (96%)]    Loss: 0.069854: 100%|██████████|
938/938 [00:12<00:00, 73.93it/s]
Testing...: 100%|██████████| 157/157 [00:02<00:00, 77.56it/s]


Test set: Average loss: 0.1302, Accuracy: 9608/10000 (96%)


Train Epoch: 4 [57600/60000 (96%)]    Loss: 0.131344: 100%|██████████|
938/938 [00:12<00:00, 75.95it/s]
Testing...: 100%|██████████| 157/157 [00:01<00:00, 83.79it/s]


Test set: Average loss: 0.1194, Accuracy: 9638/10000 (96%)


Train Epoch: 5 [57600/60000 (96%)]    Loss: 0.080049: 100%|██████████|
938/938 [00:12<00:00, 74.55it/s]
Testing...: 100%|██████████| 157/157 [00:01<00:00, 85.76it/s]


Test set: Average loss: 0.1010, Accuracy: 9698/10000 (97%)


Train Epoch: 6 [57600/60000 (96%)]    Loss: 0.136469: 100%|██████████|
938/938 [00:12<00:00, 75.87it/s]
Testing...: 100%|██████████| 157/157 [00:01<00:00, 88.01it/s]
```

```
Test set: Average loss: 0.0948, Accuracy: 9721/10000 (97%)


Train Epoch: 7 [57600/60000 (96%)]     Loss: 0.046305: 100%|████████████|
938/938 [00:12<00:00, 74.95it/s]
Testing...: 100%|████████████| 157/157 [00:01<00:00, 86.64it/s]


Test set: Average loss: 0.0888, Accuracy: 9731/10000 (97%)


Train Epoch: 8 [57600/60000 (96%)]     Loss: 0.086520: 100%|████████████|
938/938 [00:12<00:00, 75.35it/s]
Testing...: 100%|████████████| 157/157 [00:02<00:00, 73.08it/s]


Test set: Average loss: 0.0849, Accuracy: 9746/10000 (97%)


Train Epoch: 9 [57600/60000 (96%)]     Loss: 0.175178: 100%|████████████|
938/938 [00:12<00:00, 75.22it/s]
Testing...: 100%|████████████| 157/157 [00:02<00:00, 73.91it/s]


Test set: Average loss: 0.0801, Accuracy: 9751/10000 (98%)
```
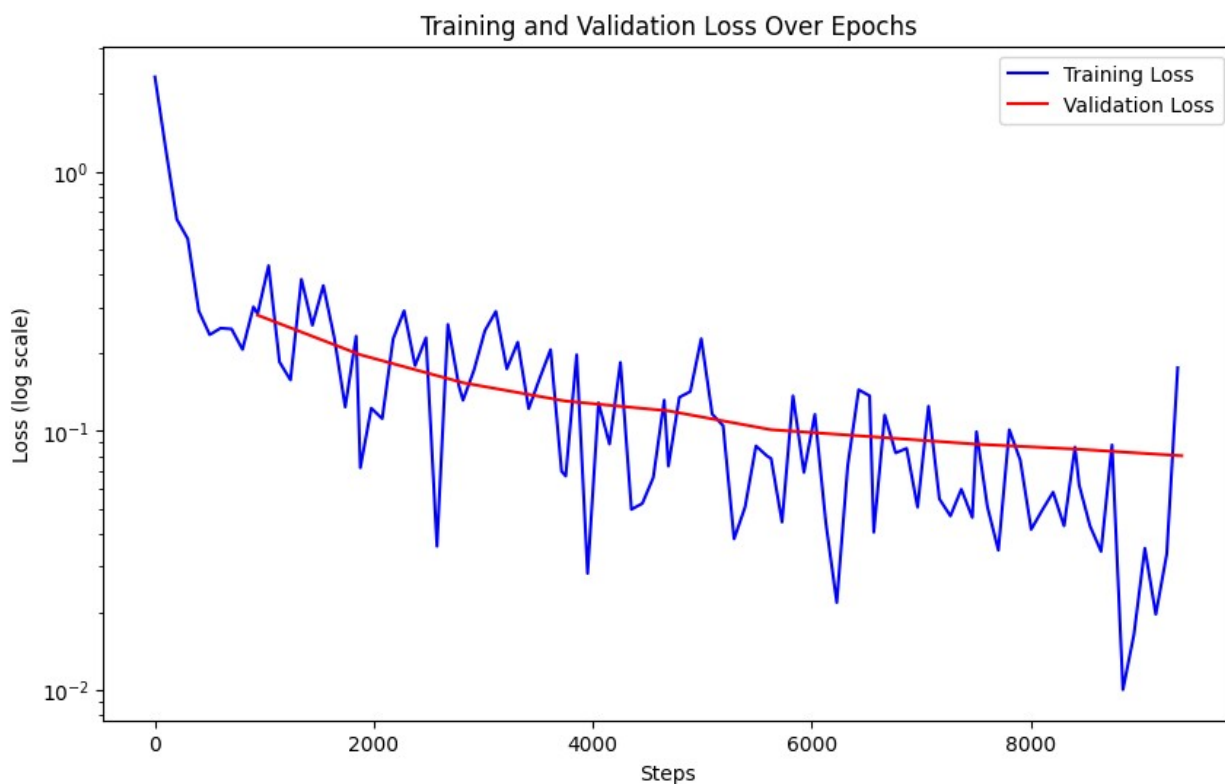


Training and Validation Loss Over Epochs

## Question 3

The model may be able to train for a bit longer. Edit the cell below to modify the previous training code to also report the time per epoch and the time for 10 epochs with testing. You can use `time.time()` to get the current time in seconds. Then run the model for another 10 epochs, printing out the execution time at the end, and replot the loss functions with the extra 10 epochs below.

```python
# visualize the losses for 20 epochs
# Get the first item
first_item = next(iter(test_loader))

# Extracting the sample and label from the first item
inputs, labels = first_item

# Print out the shapes, dtype, and identify the sample and label
print(f"Training sample shape: {inputs.shape}, dtype: {inputs.dtype}")
print(f"Training label shape: {labels.shape}, dtype: {labels.dtype}")

# create the model
model = MLP()

# you can print the model as well, but notice how the activation
# functions are missing. This is because they were called in the forward
# pass
# and not declared in the constructor
print(model)

# you can also count the model parameters
param_count = sum([p.numel() for p in model.parameters()])
print(f"Model has {param_count:,} trainable parameters")

# for a critereon (loss) function, you will use Cross-Entropy Loss.
This is the most common criterion used for multi-class prediction,
# and is also used by tokenized transformer models it takes in an un-
normalized probability distribution (i.e. without softmax) over
# N classes (in our case, 10 classes with MNIST). This distribution is
then compared to an integer label which is < N.
# For MNIST, the prediction might be [-0.0056, -0.2044,  1.1726,
0.0859,  1.8443, -0.9627,  0.9785, -1.0752, 1.1376,  1.8220], with the
label 3.
# Cross-entropy can be thought of as finding the difference between
the predicted distribution and the one-hot distribution

criterion = nn.CrossEntropyLoss()

# then you can instantiate the optimizer. You will use Stochastic
Gradient Descent (SGD), and can set the learning rate to 0.1 with a
momentum
# factor of 0.5. the first input to the optimizer is the list of model
```

```python
# parameters, which is obtained by calling .parameters() on the model
# object
optimizer = optim.SGD(model.parameters(), lr=0.01, momentum=0.5)

# create an array to log the loss and accuracy
train_losses = []
train_steps = []
test_steps = []
test_losses = []
test_accuracy = []
current_step = 0  # Start with global step 0
current_epoch = 0 # Start with epoch 0

# declare the train function
def cpu_train(epoch, train_losses, steps, current_step):

    # set the model in training mode - this doesn't do anything for us
    # right now, but it is good practiced and needed with other layers such
    # as
    # batch norm and dropout
    model.train()

    # Create tqdm progress bar to help keep track of the training
    # progress
    pbar = tqdm(enumerate(train_loader), total=len(train_loader))

    # loop over the dataset. Recall what comes out of the data loader,
    # and then by wrapping that with enumerate() we get an index into the
    # iterator list which we will call batch_idx
    for batch_idx, (data, target) in pbar:

        # during training, the first step is to zero all of the
        # gradients through the optimizer
        # this resets the state so that we can begin back propogation
        # with the updated parameters
        optimizer.zero_grad()

        # then you can apply a forward pass, which includes evaluating
        # the loss (criterion)
        output = model(data)
        loss = criterion(output, target)

        # given that you want to minimize the loss, you need to
        # call .backward() on the result, which invokes the grad_fn property
        loss.backward()

        # the backward step will automatically differentiate the model
        # and apply a gradient property to each of the parameters in the network
        # so then all you have to do is call optimizer.step() to apply
        # the gradients to the current parameters
```

```python
        optimizer.step()

        # increment the step count
        current_step += 1

        # you should add some output to the progress bar so that you
know which epoch you are training, and what the current loss is
        if batch_idx % 100 == 0:

            # append the last loss value
            train_losses.append(loss.item())
            steps.append(current_step)

            desc = (f'Train Epoch: {epoch} [{batch_idx *
len(data)}/{len(train_loader.dataset)}'
                    f' ({100. * batch_idx / len(train_loader):.0f}%)]\
tLoss: {loss.item():.6f}')
            pbar.set_description(desc)

    return current_step

# declare a test function, this will help you evaluate the model
progress on a dataset which is different from the training dataset
# doing so prevents cross-contamination and misleading results due to
overfitting
def cpu_test(test_losses, test_accuracy, steps, current_step):

    # put the model into eval mode, this again does not currently do
anything for you, but it is needed with other layers like batch_norm
    # and dropout
    model.eval()
    test_loss = 0
    correct = 0

    # Create tqdm progress bar
    pbar = tqdm(test_loader, total=len(test_loader),
desc="Testing...")

    # since you are not training the model, and do not need back-
propagation, you can use a no_grad() context
    with torch.no_grad():
        # iterate over the test set
        for data, target in pbar:
            # like with training, run a forward pass through the model
and evaluate the criterion
            output = model(data)
            test_loss += criterion(output, target).item() # you are
using .item() to get the loss value rather than the tensor itself

            # you can also check the accuracy by sampling the output -
```

```python
you can use greedy sampling which is argmax (maximum probability)
            # in general, you would want to normalize the logits first
(the un-normalized output of the model), which is done via .softmax()
            # however, argmax is taking the maximum value, which will
be the same index for the normalized and un-normalized distributions
            # so we can skip a step and take argmax directly
            pred = output.argmax(dim=1, keepdim=True)
            correct += pred.eq(target.view_as(pred)).sum().item()

    test_loss /= len(test_loader)

    # append the final test loss
    test_losses.append(test_loss)
    test_accuracy.append(correct/len(test_loader.dataset))
    steps.append(current_step)

    print(f'\nTest set: Average loss: {test_loss:.4f}, Accuracy:
{correct}/{len(test_loader.dataset)}'
            f' ({100. * correct / len(test_loader.dataset):.0f}%)\n')

# visualize the losses for the first 10 epochs
import matplotlib.pyplot as plt

# Declare the train and test functions (as you've already done)
# ...

# Initialize lists to store losses for plotting
train_losses_list = []
test_losses_list = []

# Train and test for 20 epochs
for epoch in range(0, 20):
    current_step = cpu_train(epoch, train_losses_list, train_steps,
current_step)
    cpu_test(test_losses_list, test_accuracy, test_steps,
current_step)
    current_epoch += 1

# Plot the training and testing losses
plt.figure(figsize=(10, 6))
plt.plot(train_steps, train_losses_list, label='Training Loss',
color='blue')
plt.plot(test_steps, test_losses_list, label='Validation Loss',
color='red')

# Set log scale for the y-axis
plt.yscale('log')

# Labeling the axes
plt.xlabel('Steps')
```

```python
plt.ylabel('Loss (log scale)')
plt.title('Training and Validation Loss Over Epochs')

# Add a legend
plt.legend()

# Show the plot
plt.show()
```

```
Training sample shape: torch.Size([64, 1, 28, 28]), dtype:
torch.float32
Training label shape: torch.Size([64]), dtype: torch.int64
MLP(
  (fc1): Linear(in_features=784, out_features=128, bias=True)
  (fc2): Linear(in_features=128, out_features=64, bias=True)
  (fc3): Linear(in_features=64, out_features=10, bias=True)
)
Model has 109,386 trainable parameters

Train Epoch: 0 [57600/60000 (96%)]    Loss: 0.313978: 100%|███████████|
938/938 [00:12<00:00, 74.83it/s]
Testing...: 100%|███████████| 157/157 [00:01<00:00, 83.82it/s]


Test set: Average loss: 0.2749, Accuracy: 9179/10000 (92%)


Train Epoch: 1 [57600/60000 (96%)]    Loss: 0.207073: 100%|███████████|
938/938 [00:12<00:00, 73.68it/s]
Testing...: 100%|███████████| 157/157 [00:02<00:00, 64.11it/s]


Test set: Average loss: 0.2054, Accuracy: 9400/10000 (94%)


Train Epoch: 2 [57600/60000 (96%)]    Loss: 0.314983: 100%|███████████|
938/938 [00:13<00:00, 70.04it/s]
Testing...: 100%|███████████| 157/157 [00:01<00:00, 83.77it/s]


Test set: Average loss: 0.1660, Accuracy: 9521/10000 (95%)


Train Epoch: 3 [57600/60000 (96%)]    Loss: 0.150395: 100%|███████████|
938/938 [00:12<00:00, 72.71it/s]
Testing...: 100%|███████████| 157/157 [00:01<00:00, 81.33it/s]


Test set: Average loss: 0.1441, Accuracy: 9571/10000 (96%)
```

```
Train Epoch: 4 [57600/60000 (96%)]     Loss: 0.036602: 100%|████████|
938/938 [00:13<00:00, 71.74it/s]
Testing...: 100%|████████| 157/157 [00:01<00:00, 82.97it/s]


Test set: Average loss: 0.1224, Accuracy: 9644/10000 (96%)


Train Epoch: 5 [57600/60000 (96%)]     Loss: 0.114938: 100%|████████|
938/938 [00:13<00:00, 70.64it/s]
Testing...: 100%|████████| 157/157 [00:02<00:00, 73.71it/s]


Test set: Average loss: 0.1116, Accuracy: 9664/10000 (97%)


Train Epoch: 6 [57600/60000 (96%)]     Loss: 0.180776: 100%|████████|
938/938 [00:13<00:00, 70.93it/s]
Testing...: 100%|████████| 157/157 [00:02<00:00, 70.90it/s]


Test set: Average loss: 0.1011, Accuracy: 9683/10000 (97%)


Train Epoch: 7 [57600/60000 (96%)]     Loss: 0.090131: 100%|████████|
938/938 [00:13<00:00, 72.01it/s]
Testing...: 100%|████████| 157/157 [00:01<00:00, 82.60it/s]


Test set: Average loss: 0.0942, Accuracy: 9702/10000 (97%)


Train Epoch: 8 [57600/60000 (96%)]     Loss: 0.058335: 100%|████████|
938/938 [00:12<00:00, 72.17it/s]
Testing...: 100%|████████| 157/157 [00:01<00:00, 82.96it/s]


Test set: Average loss: 0.0949, Accuracy: 9710/10000 (97%)


Train Epoch: 9 [57600/60000 (96%)]     Loss: 0.044371: 100%|████████|
938/938 [00:12<00:00, 72.28it/s]
Testing...: 100%|████████| 157/157 [00:01<00:00, 83.17it/s]


Test set: Average loss: 0.0858, Accuracy: 9736/10000 (97%)


Train Epoch: 10 [57600/60000 (96%)]     Loss: 0.120272: 100%|████████|
938/938 [00:13<00:00, 71.36it/s]
Testing...: 100%|████████| 157/157 [00:02<00:00, 70.24it/s]
```

```
Test set: Average loss: 0.0843, Accuracy: 9736/10000 (97%)


Train Epoch: 11 [57600/60000 (96%)]    Loss: 0.072290: 100%|██████████|
938/938 [00:14<00:00, 65.14it/s]
Testing...: 100%|██████████| 157/157 [00:01<00:00, 82.88it/s]


Test set: Average loss: 0.0898, Accuracy: 9734/10000 (97%)


Train Epoch: 12 [57600/60000 (96%)]    Loss: 0.036989: 100%|██████████|
938/938 [00:13<00:00, 71.31it/s]
Testing...: 100%|██████████| 157/157 [00:01<00:00, 83.80it/s]


Test set: Average loss: 0.0784, Accuracy: 9765/10000 (98%)


Train Epoch: 13 [57600/60000 (96%)]    Loss: 0.021668: 100%|██████████|
938/938 [00:12<00:00, 72.27it/s]
Testing...: 100%|██████████| 157/157 [00:01<00:00, 84.11it/s]


Test set: Average loss: 0.0772, Accuracy: 9755/10000 (98%)


Train Epoch: 14 [57600/60000 (96%)]    Loss: 0.011514: 100%|██████████|
938/938 [00:12<00:00, 72.93it/s]
Testing...: 100%|██████████| 157/157 [00:01<00:00, 85.02it/s]


Test set: Average loss: 0.0751, Accuracy: 9761/10000 (98%)


Train Epoch: 15 [57600/60000 (96%)]    Loss: 0.048451: 100%|██████████|
938/938 [00:12<00:00, 72.55it/s]
Testing...: 100%|██████████| 157/157 [00:02<00:00, 58.67it/s]


Test set: Average loss: 0.0743, Accuracy: 9771/10000 (98%)


Train Epoch: 16 [57600/60000 (96%)]    Loss: 0.026354: 100%|██████████|
938/938 [00:12<00:00, 72.90it/s]
Testing...: 100%|██████████| 157/157 [00:01<00:00, 84.26it/s]


Test set: Average loss: 0.0739, Accuracy: 9773/10000 (98%)
```

```
Train Epoch: 17 [57600/60000 (96%)]    Loss: 0.098780: 100%|████████|
938/938 [00:12<00:00, 72.56it/s]
Testing...: 100%|████████| 157/157 [00:01<00:00, 85.39it/s]


Test set: Average loss: 0.0735, Accuracy: 9767/10000 (98%)


Train Epoch: 18 [57600/60000 (96%)]    Loss: 0.011532: 100%|████████|
938/938 [00:12<00:00, 73.63it/s]
Testing...: 100%|████████| 157/157 [00:01<00:00, 83.98it/s]


Test set: Average loss: 0.0719, Accuracy: 9781/10000 (98%)


Train Epoch: 19 [57600/60000 (96%)]    Loss: 0.003208: 100%|████████|
938/938 [00:12<00:00, 72.60it/s]
Testing...: 100%|████████| 157/157 [00:01<00:00, 84.21it/s]


Test set: Average loss: 0.0762, Accuracy: 9770/10000 (98%)
```
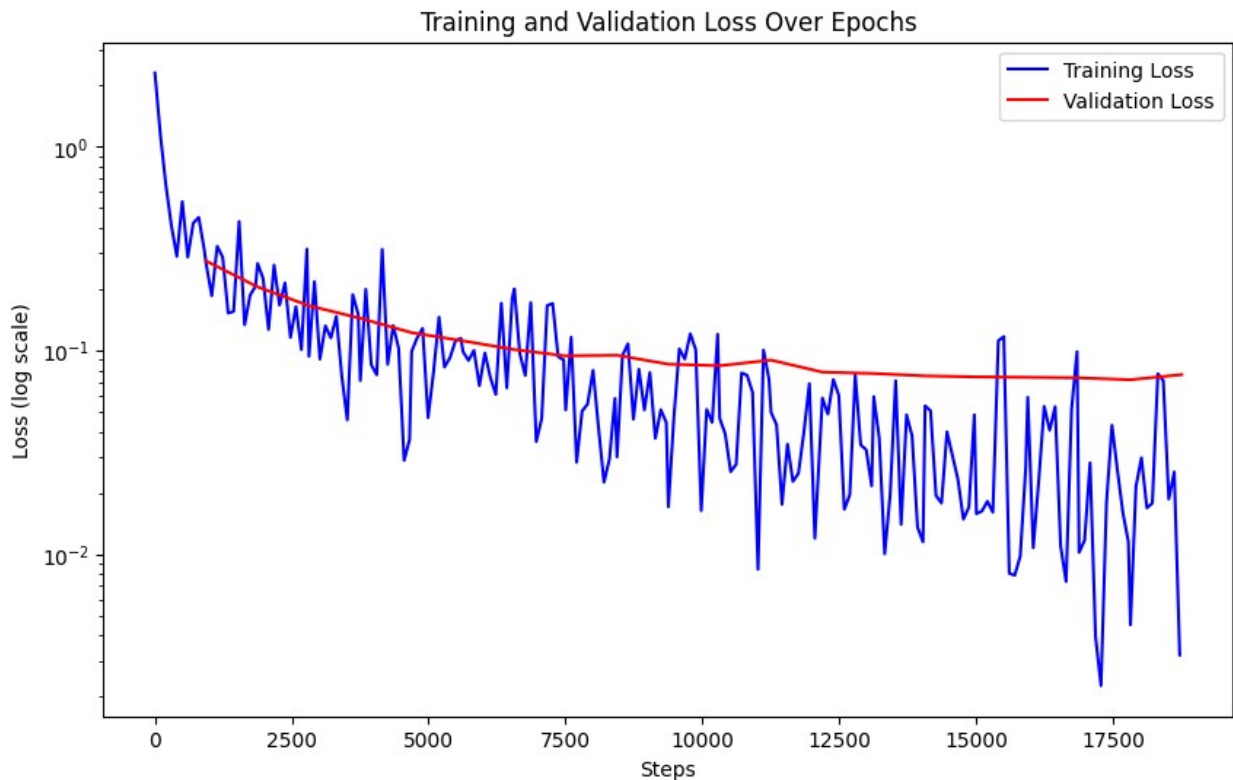


Training and Validation Loss Over Epochs

**Question 4**

Make an observation from the above plot. Do the test and train loss curves indicate that the model should train longer to improve accuracy? Or does it indicate that 20 epochs is too long? Edit the cell below to answer these questions.

YES,

# Moving to the GPU

Now that you have a model trained on the CPU, let's finally utilize the T4 GPU that we requested for this instance.

Using a GPU with torch is relatively simple, but has a few gotchas. Torch abstracts away most of the CUDA runtime API, but has a few hold-over concepts such as moving data between devices. Additionally, since the GPU is treated as a device separate from the CPU, you cannot combine CPU and GPU based tensors in the same operation. Doing so will result in a device mismatch error. If this occurs, check where the tensors are located (you can always print `.device` on a tensor), and make sure they have been properly moved to the correct device.

You will start by creating a new model, optimizer, and criterion (not really necessary in this case since you already did this above but it's better for clarity and completeness). However, one change that you'll make is moving the model to the GPU first. This can be done by calling `.cuda()` in general, or `.to("cuda")` to be more explicit. In general specific GPU devices can be targetted such as `.to("cuda:0")` for the first GPU (index 0), etc., but since there is only one GPU in Colab this is not necessary in this case.

```
# create the model
model = MLP()

# move the model to the GPU
model.cuda()

# for a critereon (loss) funciton, we will use Cross-Entropy Loss.
This is the most common critereon used for multi-class prediction, and
is also used by tokenized transformer models
# it takes in an un-normalized probability distribution (i.e. without
softmax) over N classes (in our case, 10 classes with MNIST). This
distribution is then compared to an integer label
# which is < N. For MNIST, the prediction might be [-0.0056, -0.2044,
1.1726,  0.0859,  1.8443, -0.9627,  0.9785, -1.0752, 1.1376,  1.8220],
with the label 3.
# Cross-entropy can be thought of as finding the difference between
what the predicted distribution and the one-hot distribution

criterion = nn.CrossEntropyLoss()

# then you can instantiate the optimizer. You will use Stochastic
Gradient Descent (SGD), and can set the learning rate to 0.1 with a
momentum factor of 0.5
# the first input to the optimizer is the list of model parameters,
```

```
which is obtained by calling .parameters() on the model object
optimizer = optim.SGD(model.parameters(), lr=0.01, momentum=0.5)

# create a new array to log the loss and accuracy
train_losses = []
train_steps = []
test_steps = []
test_losses = []
test_accuracy = []
current_step = 0  # Start with global step 0
current_epoch = 0 # Start with epoch 0
```

Now, copy your previous training code with the timing parameters below. It needs to be slightly modified to move everything to the GPU.

Before the line `output = model(data)`, add:

```
data = data.cuda()
target = target.cuda()
```

Note that this is needed in both the train and test functions.

**Question 5**

Please edit the cell below to show the new GPU train and test fucntions.

```
# the new GPU training functions

# new GPU training for 10 epochs
```

**Question 6**

Is training faster now that it is on a GPU? Is the speedup what you would expect? Why or why not? Edit the cell below to answer.

## Another Model Type: CNN

Until now you have trained a simple MLP for MNIST classification, however, MLPs are not a particularly good for images.

Firstly, using a MLP will require that all images have the same size and shape, since they are unrolled in the input.

Secondly, in general images can make use of translation invariance (a type of data symmetry), but this cannot but leveraged with a MLP.

For these reasons, a convolutional network is more appropriate, as it will pass kernels over the 2D image, removing the requirement for a fixed image size and leveraging the translation invariance of the 2D images.

Let's define a simple CNN below.

```python
# Define the CNN model
class CNN(nn.Module):
    # define the constructor for the network
    def __init__(self):
        super().__init__()
        # instead of declaring the layers independently, let's use the
nn.Sequential feature
        # these blocks will be executed in list order

        # you will break up the model into two parts:
        # 1) the convolutional network
        # 2) the prediction head (a small MLP)

        # the convolutional network
        self.net = nn.Sequential(
          nn.Conv2d(1, 32, kernel_size=3, stride=1, padding=1),  # the
input projection layer - note that a stride of 1 means you are not
down-sampling
          nn.ReLU(),                                             #
activation
          nn.Conv2d(32, 64, kernel_size=3, stride=2, padding=1), # an
inner layer - note that a stride of 2 means you are down sampling. The
output is 28x28 -> 14x14
          nn.ReLU(),                                             #
activation
          nn.Conv2d(64, 128, kernel_size=3, stride=2, padding=1),# an
inner layer - note that a stride of 2 means you are down sampling. The
output is 14x14 -> 7x7
          nn.ReLU(),                                             #
activation
          nn.AdaptiveMaxPool2d(1),                               # a
pooling layer which will output a 1x1 vector for the prediciton head
        )

        # the prediction head
        self.head = nn.Sequential(
          nn.Linear(128, 64),      # input projection, the output from
the pool layer is a 128 element vector
          nn.ReLU(),               # activation
          nn.Linear(64, 10)        # class projection to one of the 10
classes (digits 0-9)
        )


    # define the forward pass compute graph
    def forward(self, x):

        # pass the input through the convolution network
        x = self.net(x)
```

```
        # reshape the output from Bx128x1x1 to Bx128
        x = x.view(x.size(0), -1)

        # pass the pooled vector into the prediction head
        x = self.head(x)

        # the output here is Bx10
        return x

# create the model
model = CNN()

# print the model and the parameter count
print(model)
param_count = sum([p.numel() for p in model.parameters()])
print(f"Model has {param_count:,} trainable parameters")

# the loss function
criterion = nn.CrossEntropyLoss()

# then you can intantiate the optimizer. You will use Stochastic
Gradient Descent (SGD), and can set the learning rate to 0.1 with a
# momentum factor of 0.5
# the first input to the optimizer is the list of model parameters,
which is obtained by calling .parameters() on the model object
optimizer = optim.SGD(model.parameters(), lr=0.01, momentum=0.5)

CNN(
  (net): Sequential(
    (0): Conv2d(1, 32, kernel_size=(3, 3), stride=(1, 1), padding=(1,
1))
    (1): ReLU()
    (2): Conv2d(32, 64, kernel_size=(3, 3), stride=(2, 2), padding=(1,
1))
    (3): ReLU()
    (4): Conv2d(64, 128, kernel_size=(3, 3), stride=(2, 2),
padding=(1, 1))
    (5): ReLU()
    (6): AdaptiveMaxPool2d(output_size=1)
  )
  (head): Sequential(
    (0): Linear(in_features=128, out_features=64, bias=True)
    (1): ReLU()
    (2): Linear(in_features=64, out_features=10, bias=True)
  )
)
Model has 101,578 trainable parameters
```

**Question 7**

Notice that this model now has fewer parameters than the MLP. Let's see how it trains.

Using the previous code to train on the CPU with timing, edit the cell below to execute 2 epochs of training.

```python
# create a new array to log the loss and accuracy
train_losses = []
train_steps = []
test_steps = []
test_losses = []
test_accuracy = []
current_step = 0  # Start with global step 0
current_epoch = 0 # Start with epoch 0

# train for 2 epochs on the CPU
num_epochs = 2
for epoch in range(num_epochs):
    model.train()  # Set the model to training mode
    running_loss = 0.0
    correct = 0
    total = 0
    for inputs, labels in train_loader:
        # Zero the gradients
        optimizer.zero_grad()

        # Forward pass
        outputs = model(inputs)

        # Calculate the loss
        loss = criterion(outputs, labels)

        # Backward pass and optimize
        loss.backward()
        optimizer.step()

        # Log the loss
        running_loss += loss.item()

        # Calculate accuracy
        _, predicted = torch.max(outputs, 1)
        total += labels.size(0)
        correct += (predicted == labels).sum().item()

    # Log training loss and accuracy
    avg_loss = running_loss / len(train_loader)
    accuracy = 100 * correct / total
    train_losses.append(avg_loss)
    train_steps.append(current_step)
    print(f"Epoch [{epoch+1}/{num_epochs}], Loss: {avg_loss:.4f}, Accuracy: {accuracy:.2f}%")

    # Test the model
```

```python
    model.eval()  # Set the model to evaluation mode
    test_loss = 0.0
    correct = 0
    total = 0
    with torch.no_grad():  # No need to compute gradients during
testing
        for inputs, labels in test_loader:
            # Forward pass
            outputs = model(inputs)

            # Calculate the loss
            loss = criterion(outputs, labels)
            test_loss += loss.item()

            # Calculate accuracy
            _, predicted = torch.max(outputs, 1)
            total += labels.size(0)
            correct += (predicted == labels).sum().item()

    # Log test loss and accuracy
    avg_test_loss = test_loss / len(test_loader)
    test_accuracy = 100 * correct / total
    test_losses.append(avg_test_loss)
    test_steps.append(current_step)
    print(f"Test Loss: {avg_test_loss:.4f}, Test Accuracy:
{test_accuracy:.2f}%")

    current_epoch += 1
    current_step += 1

# Print the final model parameters count
param_count = sum(p.numel() for p in model.parameters())
print(f"Final model has {param_count:,} trainable parameters")

Epoch [1/2], Loss: 0.0225, Accuracy: 99.47%
Test Loss: 0.0722, Test Accuracy: 97.91%
Epoch [2/2], Loss: 0.0202, Accuracy: 99.57%
Test Loss: 0.0747, Test Accuracy: 97.91%
Final model has 109,386 trainable parameters
```

**Question 8**

Now, let's move the model to the GPU and try training for 2 epochs there.

```python
# create the model
model = CNN()

model.cuda()

# print the model and the parameter count
```

```python
print(model)
param_count = sum([p.numel() for p in model.parameters()])
print(f"Model has {param_count:,} trainable parameters")

# the loss function
criterion = nn.CrossEntropyLoss()

# then you can instantiate the optimizer. You will use Stochastic
Gradient Descent (SGD), and can set the learning rate to 0.1 with a
momentum factor of 0.5
# the first input to the optimizer is the list of model parameters,
which is obtained by calling .parameters() on the model object
optimizer = optim.SGD(model.parameters(), lr=0.01, momentum=0.5)

CNN(
  (net): Sequential(
    (0): Conv2d(1, 32, kernel_size=(3, 3), stride=(1, 1), padding=(1,
1))
    (1): ReLU()
    (2): Conv2d(32, 64, kernel_size=(3, 3), stride=(2, 2), padding=(1,
1))
    (3): ReLU()
    (4): Conv2d(64, 128, kernel_size=(3, 3), stride=(2, 2),
padding=(1, 1))
    (5): ReLU()
    (6): AdaptiveMaxPool2d(output_size=1)
  )
  (head): Sequential(
    (0): Linear(in_features=128, out_features=64, bias=True)
    (1): ReLU()
    (2): Linear(in_features=64, out_features=10, bias=True)
  )
)
Model has 101,578 trainable parameters

# create a new array to log the loss and accuracy
train_losses = []
train_steps = []
test_steps = []
test_losses = []
test_accuracy = []
current_step = 0  # Start with global step 0
current_epoch = 0 # Start with epoch 0

# train for 2 epochs on the GPU

# Train for 2 epochs on the GPU
num_epochs = 2
for epoch in range(num_epochs):
    model.train()  # Set the model to training mode
```

```python
    running_loss = 0.0
    correct = 0
    total = 0
    for inputs, labels in train_loader:
        # Move inputs and labels to the GPU (device)
        inputs, labels = inputs.to(device), labels.to(device)

        # Zero the gradients
        optimizer.zero_grad()

        # Forward pass
        outputs = model(inputs)

        # Calculate the loss
        loss = criterion(outputs, labels)

        # Backward pass and optimize
        loss.backward()
        optimizer.step()

        # Log the loss
        running_loss += loss.item()

        # Calculate accuracy
        _, predicted = torch.max(outputs, 1)
        total += labels.size(0)
        correct += (predicted == labels).sum().item()

    # Log training loss and accuracy
    avg_loss = running_loss / len(train_loader)
    accuracy = 100 * correct / total
    train_losses.append(avg_loss)
    train_steps.append(current_step)
    print(f"Epoch [{epoch+1}/{num_epochs}], Loss: {avg_loss:.4f},
Accuracy: {accuracy:.2f}%")

    # Test the model
    model.eval()  # Set the model to evaluation mode
    test_loss = 0.0
    correct = 0
    total = 0
    with torch.no_grad():  # No need to compute gradients during
testing
        for inputs, labels in test_loader:
            # Move inputs and labels to the GPU (device)
            inputs, labels = inputs.to(device), labels.to(device)

            # Forward pass
            outputs = model(inputs)
```

```
            # Calculate the loss
            loss = criterion(outputs, labels)
            test_loss += loss.item()

            # Calculate accuracy
            _, predicted = torch.max(outputs, 1)
            total += labels.size(0)
            correct += (predicted == labels).sum().item()

    # Log test loss and accuracy
    avg_test_loss = test_loss / len(test_loader)
    test_accuracy = 100 * correct / total
    test_losses.append(avg_test_loss)
    test_steps.append(current_step)
    print(f"Test Loss: {avg_test_loss:.4f}, Test Accuracy:
{test_accuracy:.2f}%")

    current_epoch += 1
    current_step += 1
```

**Question 9**

How do the CPU and GPU versions compare for the CNN? Is one faster than the other? Why do you think this is, and how does it differ from the MLP? Edit the cell below to answer.

The GPU is significantly faster than the CPU for training CNNs due to its ability to handle parallel computation. GPUs excel at tasks like convolutions, where many operations can be performed simultaneously, making them ideal for large-scale deep learning. CPUs, on the other hand, are optimized for sequential tasks and have fewer cores for parallelism, resulting in slower performance for CNNs.

For MLPs, the difference between CPU and GPU is less pronounced. While GPUs still offer speed advantages, especially for large models, the performance boost is smaller compared to CNNs since MLPs involve simpler operations like matrix multiplications, which don't benefit as much from parallelism.

As a final comparison, you can profile the FLOPs (floating-point operations) executed by each model. You will use the thop.profile function for this and consider an MNIST batch size of 1.

```
# the input shape of a MNIST sample with batch_size = 1
input = torch.randn(1, 1, 28, 28)

# create a copy of the models on the CPU
mlp_model = MLP()
cnn_model = CNN()

# profile the MLP
flops, params = thop.profile(mlp_model, inputs=(input, ),
verbose=False)
print(f"MLP has {params:,} params and uses {flops:,} FLOPs")
```

```
# profile the CNN
flops, params = thop.profile(cnn_model, inputs=(input, ),
verbose=False)
print(f"CNN has {params:,} params and uses {flops:,} FLOPs")

MLP has 109,386.0 params and uses 109,184.0 FLOPs
CNN has 101,578.0 params and uses 7,459,968.0 FLOPs
```

**Question 10**

Are these results what you would have expected? Do they explain the performance difference between running on the CPU and GPU? Why or why not? Edit the cell below to answer.

Yes, these results are consistent with expectations and help explain the performance difference between running on the CPU and GPU.

## Results Analysis:

- **MLP**:
  - **Params**: 109,386
  - **FLOPs**: 109,184
  - MLPs typically have a large number of parameters, but their operations mainly consist of dense matrix multiplications. These operations are less computationally intensive than convolutions and don't benefit as much from parallelization. Hence, the FLOPs for an MLP are relatively low.
- **CNN**:
  - **Params**: 101,578
  - **FLOPs**: 7,459,968
  - CNNs have fewer parameters than MLPs but require far more FLOPs due to the convolution operations. These operations involve large input/output tensors and can be highly parallelized, resulting in much higher computational complexity (FLOPs). This explains why CNNs have such a large number of FLOPs compared to MLPs.

## Performance Difference:

- **CPU vs. GPU**: The significant difference in FLOPs between MLP and CNNs explains why CNNs benefit more from the GPU's parallel processing capabilities. CNNs' large number of FLOPs due to convolutions make them more computationally expensive, which is why the GPU provides a much larger speedup compared to the MLP, where the computation is simpler and less parallelizable.