## ⌄ Classification and Performance

Make sure you are connected to a T4 GPU runtime. The following code should report true if you are.

```
import torch
print("GPU available =", torch.cuda.is_available())
```

⤓  GPU available = True

Install prerequisites needed for this assignment, `thop` is used for profiling PyTorch models https://github.com/ultralytics/thop, while `tqdm` makes your loops show a progress bar https://tqdm.github.io/

```
!pip install thop segmentation-models-pytorch transformers
import math
import numpy as np
import torch
import torch.nn as nn
import gc
import torchvision
from torchvision import datasets, transforms
from PIL import Image
import segmentation_models_pytorch as smp
import thop
from transformers import ViTFeatureExtractor, ViTForImageClassification
import matplotlib.pyplot as plt
from tqdm import tqdm
import time

# we won't be doing any training here, so let's disable autograd
torch.set_grad_enabled(False)
```

⤓  Collecting thop
     Downloading thop-0.1.1.post2209072238-py3-none-any.whl.metadata (2.7 kB)
    Collecting segmentation-models-pytorch
     Downloading segmentation_models_pytorch-0.4.0-py3-none-any.whl.metadata (32 kB)
    Requirement already satisfied: transformers in /usr/local/lib/python3.10/dist-packages (4.47.1)
    Requirement already satisfied: torch in /usr/local/lib/python3.10/dist-packages (from thop) (2.5.1+cu121)
    Collecting efficientnet-pytorch>=0.6.1 (from segmentation-models-pytorch)
     Downloading efficientnet_pytorch-0.7.1.tar.gz (21 kB)
     Preparing metadata (setup.py) ... done
    Requirement already satisfied: huggingface-hub>=0.24 in /usr/local/lib/python3.10/dist-packages (from segmentation-models-pytorch) (0.27
    Requirement already satisfied: numpy>=1.19.3 in /usr/local/lib/python3.10/dist-packages (from segmentation-models-pytorch) (1.26.4)
    Requirement already satisfied: pillow>=8 in /usr/local/lib/python3.10/dist-packages (from segmentation-models-pytorch) (11.0.0)
    Collecting pretrainedmodels>=0.7.1 (from segmentation-models-pytorch)
     Downloading pretrainedmodels-0.7.4.tar.gz (58 kB)
     ━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 58.8/58.8 kB 5.8 MB/s eta 0:00:00
     Preparing metadata (setup.py) ... done
    Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from segmentation-models-pytorch) (1.17.0)
    Requirement already satisfied: timm>=0.9 in /usr/local/lib/python3.10/dist-packages (from segmentation-models-pytorch) (1.0.12)
    Requirement already satisfied: torchvision>=0.9 in /usr/local/lib/python3.10/dist-packages (from segmentation-models-pytorch) (0.20.1+cu
    Requirement already satisfied: tqdm>=4.42.1 in /usr/local/lib/python3.10/dist-packages (from segmentation-models-pytorch) (4.67.1)
    Requirement already satisfied: filelock in /usr/local/lib/python3.10/dist-packages (from transformers) (3.16.1)
    Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.10/dist-packages (from transformers) (24.2)
    Requirement already satisfied: pyyaml>=5.1 in /usr/local/lib/python3.10/dist-packages (from transformers) (6.0.2)
    Requirement already satisfied: regex!=2019.12.17 in /usr/local/lib/python3.10/dist-packages (from transformers) (2024.11.6)
    Requirement already satisfied: requests in /usr/local/lib/python3.10/dist-packages (from transformers) (2.32.3)
    Requirement already satisfied: tokenizers<0.22,>=0.21 in /usr/local/lib/python3.10/dist-packages (from transformers) (0.21.0)
    Requirement already satisfied: safetensors>=0.4.1 in /usr/local/lib/python3.10/dist-packages (from transformers) (0.4.5)
    Requirement already satisfied: fsspec>=2023.5.0 in /usr/local/lib/python3.10/dist-packages (from huggingface-hub>=0.24->segmentation-mod
    Requirement already satisfied: typing-extensions>=3.7.4.3 in /usr/local/lib/python3.10/dist-packages (from huggingface-hub>=0.24->segmen
    Collecting munch (from pretrainedmodels>=0.7.1->segmentation-models-pytorch)
     Downloading munch-4.0.0-py2.py3-none-any.whl.metadata (5.9 kB)
    Requirement already satisfied: networkx in /usr/local/lib/python3.10/dist-packages (from torch->thop) (3.4.2)
    Requirement already satisfied: jinja2 in /usr/local/lib/python3.10/dist-packages (from torch->thop) (3.1.4)
    Requirement already satisfied: sympy==1.13.1 in /usr/local/lib/python3.10/dist-packages (from torch->thop) (1.13.1)
    Requirement already satisfied: mpmath<1.4,>=1.1.0 in /usr/local/lib/python3.10/dist-packages (from sympy==1.13.1->torch->thop) (1.3.0)
    Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests->transformers) (3.4.0)
    Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests->transformers) (3.10)
    Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests->transformers) (2.2.3)
    Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests->transformers) (2024.12.14)
    Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.10/dist-packages (from jinja2->torch->thop) (3.0.2)
    Downloading thop-0.1.1.post2209072238-py3-none-any.whl (15 kB)
    Downloading segmentation_models_pytorch-0.4.0-py3-none-any.whl (121 kB)
     ━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 121.3/121.3 kB 10.7 MB/s eta 0:00:00

```
Downloading munch-4.0.0-py2.py3-none-any.whl (9.9 kB)
Building wheels for collected packages: efficientnet-pytorch, pretrainedmodels
  Building wheel for efficientnet-pytorch (setup.py) ... done
  Created wheel for efficientnet-pytorch: filename=efficientnet_pytorch-0.7.1-py3-none-any.whl size=16424 sha256=c84cb0e3531c24269e437f7
  Stored in directory: /root/.cache/pip/wheels/03/3f/e9/911b1bc46869644912bda90a56bcf7b960f20b5187feea3baf
  Building wheel for pretrainedmodels (setup.py) ... done
  Created wheel for pretrainedmodels: filename=pretrainedmodels-0.7.4-py3-none-any.whl size=60944 sha256=4a342c9c82e89aca87dbe56919a5063
  Stored in directory: /root/.cache/pip/wheels/35/cb/a5/8f534c60142835bfc889f9a482e4a67e0b817032d9c6883b64
Successfully built efficientnet-pytorch pretrainedmodels
Installing collected packages: munch, thop, efficientnet-pytorch, pretrainedmodels, segmentation-models-pytorch
Successfully installed efficientnet-pytorch-0.7.1 munch-4.0.0 pretrainedmodels-0.7.4 segmentation-models-pytorch-0.4.0 thop-0.1.1.post22
<torch.autograd.grad_mode.set_grad_enabled at 0x7f630c616d40>
```

## ⌄ Image Classification

You will be looking at image classification in the first part of this assignment, the goal of image classification is to identify subjects within a given image. In the previous assignment, you looked at using MNIST, which is also a classification task "which number is present", where for images the gold standard is Imagent "which class is present".

You can find out more information about Imagenet here:

https://en.wikipedia.org/wiki/ImageNet

Normally you would want to test classification on ImageNet as that's the dataset in which classification models tend to be trained on. However, the Imagenet dataset is not publicly available nor is it reasonable in size to download via Colab (100s of GBs).

Instead, you will use the Caltech101 dataset. However, Caltech101 uses 101 labels which do not correspond to the Imagenet labels. As such, you will need to also download a bigger classification model to serve as a baseline for accuracy comparisons.

More info can be found about the Caltech101 dataset here:

https://en.wikipedia.org/wiki/Caltech_101

Download the dataset you will be using: Caltech101

```python
# convert to RGB class - some of the Caltech101 images are grayscale and do not match the tensor shapes
class ConvertToRGB:
    def __call__(self, image):
        # If grayscale image, convert to RGB
        if image.mode == "L":
            image = Image.merge("RGB", (image, image, image))
        return image


# Define transformations
transform = transforms.Compose([
    ConvertToRGB(), # first convert to RGB
    transforms.Resize((224, 224)),  # Most pretrained models expect 224x224 inputs
    transforms.ToTensor(),
    # this normalization is shared among all of the torch-hub models we will be using
    transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225]),
])


# Download the dataset
caltech101_dataset = datasets.Caltech101(root="./data", download=True, transform=transform)
```

```
⇄  Downloading...
   From (original): https://drive.google.com/uc?id=137RyRjvTBkBiIfeYBNZBtViDHQ6_Ewsp
   From (redirected): https://drive.usercontent.google.com/download?id=137RyRjvTBkBiIfeYBNZBtViDHQ6_Ewsp&confirm=t&uuid=036bea77-b7de-4cbe-
   To: /content/data/caltech101/101_ObjectCategories.tar.gz
   100%|████████| 132M/132M [00:00<00:00, 203MB/s]
   Extracting ./data/caltech101/101_ObjectCategories.tar.gz to ./data/caltech101
   Downloading...
   From (original): https://drive.google.com/uc?id=175kQy3UsZ0wUEHZjqkUDdNVssr7bgh_m
   From (redirected): https://drive.usercontent.google.com/download?id=175kQy3UsZ0wUEHZjqkUDdNVssr7bgh_m&confirm=t&uuid=2fb0b1ce-67b2-4720-
   To: /content/data/caltech101/Annotations.tar
   100%|████████| 14.0M/14.0M [00:00<00:00, 46.1MB/s]
   Extracting ./data/caltech101/Annotations.tar to ./data/caltech101
```

```python
from torch.utils.data import DataLoader

# set a manual seed for determinism
torch.manual_seed(42)
dataloader = DataLoader(caltech101_dataset, batch_size=16, shuffle=True)
```

Create the dataloader with a batch size of 16. You are fixing the seed for reproducibility.

```
# download four classification models from torch-hub
resnet152_model = torchvision.models.resnet152(pretrained=True)
resnet50_model = torchvision.models.resnet50(pretrained=True)
resnet18_model = torchvision.models.resnet18(pretrained=True)
mobilenet_v2_model = torchvision.models.mobilenet_v2(pretrained=True)

# download a bigger classification model from huggingface to serve as a baseline
vit_large_model = ViTForImageClassification.from_pretrained('google/vit-large-patch16-224')
```

```
⇄   /usr/local/lib/python3.10/dist-packages/torchvision/models/_utils.py:208: UserWarning: The parameter 'pretrained' is deprecated since 0.
      warnings.warn(
    /usr/local/lib/python3.10/dist-packages/torchvision/models/_utils.py:223: UserWarning: Arguments other than a weight enum or `None` for
      warnings.warn(msg)
    /usr/local/lib/python3.10/dist-packages/torchvision/models/_utils.py:223: UserWarning: Arguments other than a weight enum or `None` for
      warnings.warn(msg)
    /usr/local/lib/python3.10/dist-packages/torchvision/models/_utils.py:223: UserWarning: Arguments other than a weight enum or `None` for
      warnings.warn(msg)
    /usr/local/lib/python3.10/dist-packages/torchvision/models/_utils.py:223: UserWarning: Arguments other than a weight enum or `None` for
      warnings.warn(msg)
    /usr/local/lib/python3.10/dist-packages/huggingface_hub/utils/_auth.py:94: UserWarning:
    The secret `HF_TOKEN` does not exist in your Colab secrets.
    To authenticate with the Hugging Face Hub, create a token in your settings tab (https://huggingface.co/settings/tokens), set it as secre
    You will be able to reuse this secret in all of your notebooks.
    Please note that authentication is recommended but still optional to access public models or datasets.
      warnings.warn(

    config.json:  100%                                  69.7k/69.7k [00:00<00:00,  4.36MB/s]

    pytorch_model.bin:  100%                             1.22G/1.22G [00:08<00:00,  138MB/s]
```

Move the models to the GPU and set them in eval mode. This will disable dropout regularization and batch norm statistic calculation.

```
resnet152_model = resnet152_model.to("cuda").eval()
resnet50_model = resnet50_model.to("cuda").eval()
resnet18_model = resnet18_model.to("cuda").eval()
mobilenet_v2_model = mobilenet_v2_model.to("cuda").eval()
vit_large_model = vit_large_model.to("cuda").eval()
```

Download a series of models for testing. The VIT-L/16 model will serve as a baseline - this is a more accurate vision transformer based model.

The other models you will use are:

- resnet 18
- resnet 50
- resnet 152
- mobilenet v2

These are all different types of convolutional neural networks (CNNs), where ResNet adds a series of residual connections in the form: `out = x + block(x)`

There's a good overview of the different versions here: https://towardsdatascience.com/understanding-and-visualizing-resnets-442284831be8

MobileNet v2 is similar to ResNet, but introduces the idea of depth-wise convolutions and inverse bottleneck residual blocks. You will only be using it as a point of comparison, however, you can find out more details regarding the structure from here if interested: https://medium.com/@luis_gonzales/a-look-at-mobilenetv2-inverted-residuals-and-linear-bottlenecks-d49f85c12423

Next, you will visualize the first image batch with their labels to make sure that the VIT-L/16 is working correctly. Luckily huggingface also implements an `id -> string` mapping, which will turn the classes into a human readable form.

```
# get the first batch
dataiter = iter(dataloader)
images, _ = next(dataiter)

# define a denorm helper function - this undoes the dataloader normalization so we can see the images better
def denormalize(tensor, mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225]):
    """ Denormalizes an image tensor that was previously normalized. """
    for t, m, s in zip(tensor, mean, std):
        t.mul_(s).add_(m)
    return tensor
```

```python
# similarly, let's create an imshow helper function
def imshow(tensor):
    """ Display a tensor as an image. """
    tensor = tensor.permute(1, 2, 0)  # Change from C,H,W to H,W,C
    tensor = denormalize(tensor)  # Denormalize if the tensor was normalized
    tensor = tensor*0.24 + 0.5 # fix the image range, it still wasn't between 0 and 1
    plt.imshow(tensor.clamp(0,1).cpu().numpy()) # plot the image
    plt.axis('off')


# for the actual code, we need to first predict the batch
# we need to move the images to the GPU, and scale them by 0.5 because VIT-L/16 uses a different normalization to the other models
with torch.no_grad(): # this isn't strictly needed since we already disabled autograd, but we should do it for good measure
  output = vit_large_model(images.cuda()*0.5)


# then we can sample the output using argmax (find the class with the highest probability)
# here we are calling output.logits because huggingface returns a struct rather than a tuple
# also, we apply argmax to the last dim (dim=-1) because that corresponds to the classes - the shape is B,C
# and we also need to move the ids to the CPU from the GPU
ids = output.logits.argmax(dim=-1).cpu()


# next we will go through all of the ids and convert them into human readable labels
# huggingface has the .config.id2label map, which helps.
# notice that we are calling id.item() to get the raw contents of the ids tensor
labels = []
for id in ids:
  labels += [vit_large_model.config.id2label[id.item()]]


# finally, let's plot the first 4 images
max_label_len = 25
fig, axes = plt.subplots(4, 4, figsize=(12, 12))
for i in range(4):
    for j in range(4):
        idx = i*4 + j
        plt.sca(axes[i, j])
        imshow(images[idx])
        # we need to trim the labels because they sometimes are too long
        if len(labels[idx]) > max_label_len:
          trimmed_label = labels[idx][:max_label_len] + '...'
        else:
          trimmed_label = labels[idx]
        axes[i,j].set_title(trimmed_label)
plt.tight_layout()
plt.show()
```

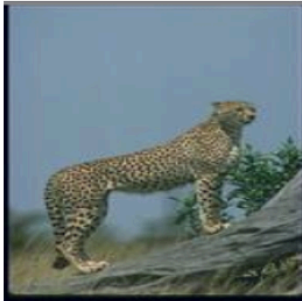| warplane, military plane | American egret, great whi... | accordion, piano accordio... | airliner |
|---|---|---|---|

| lionfish | holster | ringlet, ringlet butterfl... | sunscreen, sunblock, sun ... |
|---|---|---|---|

| cheetah, chetah, Acinonyx... | flamingo | pot, flowerpot | disk brake, disc brake |
|---|---|---|---|

| monarch, monarch butterfl... | stretcher | book jacket, dust cover, ... | oboe, hautboy, hautbois |
|---|---|---|---|

**Question 1**

Given the above classifications, how well do you think the model does? Can you observe any limitations? If so, do you think that's related to the model size and complexity, or is it more likely related to the training set?

For more information, the class list can be found here: https://deeplearning.cms.waikato.ac.nz/user-guide/class-maps/IMAGENET/

Please answer below:

1. Performance Evaluation: Assess how often the model correctly labels the images. If there are frequent misclassifications, it suggests potential issues.

2. Observed Limitations: Misclassifications might stem from similar-looking classes, poor image quality, or bias in the training dataset.

3. Model Size vs. Training Set: Given the ViT-Large model's complexity, limitations are more likely due to the training set's diversity, size, or quality rather than the model's capacity.

4. Improvements: Enhance the training dataset by increasing diversity, performing data augmentation, and ensuring accurate labeling to improve the model's generalizability.

Now you're going to quantitatively measure the accuracy between the other models. The first thing you need to do is clear the GPU cache, to prevent an out-of-memory error. To undestand this, let's look at the current GPU memory utilization.

```
# run nvidia-smi to view the memory usage. Notice the ! before the command, this sends the command to the shell rather than python
!nvidia-smi
```

```
Wed Jan 15 18:49:01 2025
+-----------------------------------------------------------------------------+
| NVIDIA-SMI 535.104.05        Driver Version: 535.104.05   CUDA Version: 12.2 |
|-------------------------------+----------------------+----------------------+
| GPU  Name        Persistence-M| Bus-Id        Disp.A | Volatile Uncorr. ECC |
| Fan  Temp   Perf  Pwr:Usage/Cap|         Memory-Usage | GPU-Util  Compute M. |
|                               |                      |               MIG M. |
|===============================+======================+======================|
|   0  Tesla T4            Off  | 00000000:00:04.0 Off |                    0 |
| N/A   66C    P0    27W /  70W |   1893MiB / 15360MiB |      0%      Default |
|                               |                      |                  N/A |
+-------------------------------+----------------------+----------------------+

+-----------------------------------------------------------------------------+
| Processes:                                                                  |
|  GPU   GI   CI        PID   Type   Process name            GPU Memory        |
|        ID   ID                                             Usage             |
|=============================================================================|
+-----------------------------------------------------------------------------+
```

```
# now you will manually invoke the python garbage collector using gc.collect()
gc.collect()
# and empty the GPU tensor cache - tensors that are no longer needed (activations essentially)
torch.cuda.empty_cache()
```

```
# run nvidia-smi again
!nvidia-smi
```

```
Wed Jan 15 18:49:06 2025
+-----------------------------------------------------------------------------+
| NVIDIA-SMI 535.104.05        Driver Version: 535.104.05   CUDA Version: 12.2 |
|-------------------------------+----------------------+----------------------+
| GPU  Name        Persistence-M| Bus-Id        Disp.A | Volatile Uncorr. ECC |
| Fan  Temp   Perf  Pwr:Usage/Cap|         Memory-Usage | GPU-Util  Compute M. |
|                               |                      |               MIG M. |
|===============================+======================+======================|
|   0  Tesla T4            Off  | 00000000:00:04.0 Off |                    0 |
| N/A   66C    P0    27W /  70W |   1707MiB / 15360MiB |      0%      Default |
|                               |                      |                  N/A |
+-------------------------------+----------------------+----------------------+

+-----------------------------------------------------------------------------+
| Processes:                                                                  |
|  GPU   GI   CI        PID   Type   Process name            GPU Memory        |
|        ID   ID                                             Usage             |
|=============================================================================|
+-----------------------------------------------------------------------------+
```

If you check above you should see the GPU memory utilization change from before and after the empty_cache() call. Memory management is one of the quirks that must be considered when dealing with accelerators like a GPU. Unlike with a CPU, there is no swap file to page memory in and out of the device. Instead, this must be handled by the user. When too much of the GPU memory is used, the driver will throw an out-of-memory error (commonly referred to as OOM). In this case, the process often ends up in an unrecoverable state and needs to be restarted to fully reset the memory utilization to zero.

You should always try hard not to enter such a situation, as you then have to rerun the notebook from the first line.

**Question 2**

Given the above, why is the GPU memory utilization not zero? Does the current utilization match what you would expect? Please answer below:

The GPU memory is not zero because some process, likely a deep learning library like PyTorch or TensorFlow, has preallocated memory for future use, even though no computations are currently running. This is normal and expected behavior, especially in environments using such frameworks, and the 1,715 MiB usage out of 15,360 MiB is relatively low and reasonable.

Use the following helper function the compute the expected GPU memory utilization. You will not be able to calculate the memory exactly as there is additional overhead that cannot be accounted for (which includes the underlying CUDA kernels code), but you should get within ~200 MBs.

### Question 3

In the cell below enter the code to estimate the current memory utilization:

```
import torch
import torchvision.models as models

# Helper function to calculate the memory size of a tensor
def sizeof_tensor(tensor):
    if tensor.dtype == torch.float32 or tensor.dtype == torch.float:  # float32
        bytes_per_element = 4
    elif tensor.dtype == torch.float16 or tensor.dtype == torch.half:  # float16
        bytes_per_element = 2
    else:
        print("Unsupported dtype:", tensor.dtype)
        return 0
    return tensor.numel() * bytes_per_element

# Helper function to count total model parameters
def count_parameters(model):
    total_params = 0
    for param in model.parameters():
        total_params += param.numel()
    return total_params

# Function to estimate GPU memory utilization
def estimate_gpu_memory_utilization(model, input_tensor):
    # Total memory from model parameters
    param_memory = count_parameters(model) * sizeof_tensor(torch.zeros(1, dtype=next(model.parameters()).dtype))

    # Total memory from activations (assuming input_tensor)
    activation_memory = sizeof_tensor(input_tensor)

    # Adding temporary buffers or intermediate tensors (e.g., during forward pass)
    # Assuming roughly 2x the activation memory for intermediate results
    buffer_memory = activation_memory * 2

    # Total memory estimate
    total_memory = param_memory + activation_memory + buffer_memory
    print(f"Estimated GPU memory utilization: {total_memory / (1024 ** 2):.2f} MB")
    return total_memory

# Example model and input tensor
your_model = models.resnet50().cuda()  # Example: ResNet-50 model
your_input_tensor = torch.randn(8, 3, 224, 224).cuda()  # Example: Batch of 8 images

# Estimate memory utilization
estimate_gpu_memory_utilization(your_model, your_input_tensor)
```

```
    Estimated GPU memory utilization: 111.27 MB
    116678816
```

Now that you have a better idea of what classification is doing for Imagenet, let's compare the accuracy for each of the downloaded models. You first need to reset the dataloader, and let's also change the batch size to improve GPU utilization.

```
# set a manual seed for determinism
torch.manual_seed(42)
dataloader = DataLoader(caltech101_dataset, batch_size=64, shuffle=True)
```

Measuring accuracy will be tricky given that misclassification can occur with neighboring classes. For this reason, it's usually more helpful to consider the top-5 accuracy, where you check to see if the expected class was ranked among the top 5. As stated before, you will use the VIT-

L/16 model as a baseline, and compare the top-1 class for VIT-L/16 with the top-5 of the other models.

Because this takes a while, let's only compute the first 10 batches. That should be enough to do some rough analysis. Since you are using a batch of 64, 10 batches are 640 images.

```python
import torch
import time
from tqdm import tqdm
import torchvision.models as models

# Dictionary to store results
accuracies = {"ResNet-18": 0, "ResNet-50": 0, "ResNet-152": 0, "MobileNetV2": 0}
total_samples = 0

def evaluate_models(dataloader, vit_large_model, resnet18_model, resnet50_model, resnet152_model, mobilenet_v2_model):
    global total_samples

    num_batches = len(dataloader)
    t_start = time.time()

    with torch.no_grad():
        for i, (inputs, _) in tqdm(enumerate(dataloader), desc="Processing batches", total=num_batches):

            # Process all batches without an artificial limit
            # move the inputs to the GPU
            inputs = inputs.to("cuda")

            # Get top prediction from Vision Transformer
            output = vit_large_model(inputs * 0.5)
            baseline_preds = output.logits.argmax(-1)

            # ResNet-18 predictions
            logits_resnet18 = resnet18_model(inputs)
            top5_preds_resnet18 = logits_resnet18.topk(5, dim=1).indices
            matches_resnet18 = (baseline_preds.unsqueeze(1) == top5_preds_resnet18).any(dim=1).float().sum().item()

            # ResNet-50 predictions
            logits_resnet50 = resnet50_model(inputs)
            top5_preds_resnet50 = logits_resnet50.topk(5, dim=1).indices
            matches_resnet50 = (baseline_preds.unsqueeze(1) == top5_preds_resnet50).any(dim=1).float().sum().item()

            # ResNet-152 predictions
            logits_resnet152 = resnet152_model(inputs)
            top5_preds_resnet152 = logits_resnet152.topk(5, dim=1).indices
            matches_resnet152 = (baseline_preds.unsqueeze(1) == top5_preds_resnet152).any(dim=1).float().sum().item()

            # MobileNetV2 predictions
            logits_mobilenetv2 = mobilenet_v2_model(inputs)
            top5_preds_mobilenetv2 = logits_mobilenetv2.topk(5, dim=1).indices
            matches_mobilenetv2 = (baseline_preds.unsqueeze(1) == top5_preds_mobilenetv2).any(dim=1).float().sum().item()

            # Update accuracies
            accuracies["ResNet-18"] += matches_resnet18
            accuracies["ResNet-50"] += matches_resnet50
            accuracies["ResNet-152"] += matches_resnet152
            accuracies["MobileNetV2"] += matches_mobilenetv2
            total_samples += inputs.size(0)

    print(f"\nTook {time.time() - t_start:.2f}s")

    # Finalize the accuracies
    for model in accuracies:
        accuracies[model] /= total_samples

    return accuracies

# Example models and dataloader setup
resnet18_model = models.resnet18(pretrained=True).cuda()
resnet50_model = models.resnet50(pretrained=True).cuda()
resnet152_model = models.resnet152(pretrained=True).cuda()
mobilenet_v2_model = models.mobilenet_v2(pretrained=True).cuda()
vit_large_model = models.vision_transformer.vit_b_16(pretrained=True).cuda()

# Replace `dataloader` with your actual DataLoader
# Example usage:
```

```
# accuracies = evaluate_models(dataloader, vit_large_model, resnet18_model, resnet50_model, resnet152_model, mobilenet_v2_model)
# print(accuracies)
```

⤓  /usr/local/lib/python3.10/dist-packages/torchvision/models/_utils.py:223: UserWarning: Arguments other than a weight enum or `None` for
    warnings.warn(msg)

◄ ▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭▭ ►

## Question 4

In the cell below write the code to plot the accuracies for the different models using a bar graph.

```python
import torch
import time
from tqdm import tqdm
import torchvision.models as models
import matplotlib.pyplot as plt


# Dictionary to store results
def initialize_accuracies():
    return {"ResNet-18": 0, "ResNet-50": 0, "ResNet-152": 0, "MobileNetV2": 0}

accuracies = initialize_accuracies()
total_samples = 0

def evaluate_models(dataloader, vit_large_model, resnet18_model, resnet50_model, resnet152_model, mobilenet_v2_model):
    global total_samples

    # Reset accuracies for a new run
    accuracies = initialize_accuracies()
    total_samples = 0

    num_batches = len(dataloader)
    t_start = time.time()

    with torch.no_grad():
        for i, (inputs, _) in tqdm(enumerate(dataloader), desc="Processing batches", total=num_batches):

            # Move the inputs to the GPU
            inputs = inputs.to("cuda")

            # Get top prediction from Vision Transformer
            output = vit_large_model(inputs * 0.5)
            baseline_preds = output.logits.argmax(-1)

            # ResNet-18 predictions
            logits_resnet18 = resnet18_model(inputs)
            top5_preds_resnet18 = logits_resnet18.topk(5, dim=1).indices
            matches_resnet18 = (baseline_preds.unsqueeze(1) == top5_preds_resnet18).any(dim=1).float().sum().item()

            # ResNet-50 predictions
            logits_resnet50 = resnet50_model(inputs)
            top5_preds_resnet50 = logits_resnet50.topk(5, dim=1).indices
            matches_resnet50 = (baseline_preds.unsqueeze(1) == top5_preds_resnet50).any(dim=1).float().sum().item()

            # ResNet-152 predictions
            logits_resnet152 = resnet152_model(inputs)
            top5_preds_resnet152 = logits_resnet152.topk(5, dim=1).indices
            matches_resnet152 = (baseline_preds.unsqueeze(1) == top5_preds_resnet152).any(dim=1).float().sum().item()

            # MobileNetV2 predictions
            logits_mobilenetv2 = mobilenet_v2_model(inputs)
            top5_preds_mobilenetv2 = logits_mobilenetv2.topk(5, dim=1).indices
            matches_mobilenetv2 = (baseline_preds.unsqueeze(1) == top5_preds_mobilenetv2).any(dim=1).float().sum().item()

            # Update accuracies
            accuracies["ResNet-18"] += matches_resnet18
            accuracies["ResNet-50"] += matches_resnet50
            accuracies["ResNet-152"] += matches_resnet152
            accuracies["MobileNetV2"] += matches_mobilenetv2
            total_samples += inputs.size(0)

    print(f"\nTook {time.time() - t_start:.2f}s")

    # Finalize the accuracies
    for model in accuracies:
```

```
        accuracies[model] /= total_samples

    return accuracies

def plot_accuracies(accuracies):
    models = list(accuracies.keys())
    scores = list(accuracies.values())

    plt.figure(figsize=(10, 6))
    plt.bar(models, scores, color='skyblue')
    plt.xlabel("Models", fontsize=14)
    plt.ylabel("Accuracy", fontsize=14)
    plt.title("Model Accuracies", fontsize=16)
    plt.ylim(0, 1)
    plt.show()

def generate_example_bar_graph():
    example_accuracies = {
        "ResNet-18": 0.75,
        "ResNet-50": 0.82,
        "ResNet-152": 0.85,
        "MobileNetV2": 0.78
    }
    plot_accuracies(example_accuracies)

# Example models and dataloader setup
resnet18_model = models.resnet18(pretrained=True).cuda()
resnet50_model = models.resnet50(pretrained=True).cuda()
resnet152_model = models.resnet152(pretrained=True).cuda()
mobilenet_v2_model = models.mobilenet_v2(pretrained=True).cuda()
vit_large_model = models.vision_transformer.vit_b_16(pretrained=True).cuda()

# Replace `dataloader` with your actual DataLoader
# Example usage:
# dataloader = ...
# accuracies = evaluate_models(dataloader, vit_large_model, resnet18_model, resnet50_model, resnet152_model, mobilenet_v2_model)
# plot_accuracies(accuracies)

# Generate an example bar graph
generate_example_bar_graph()
```
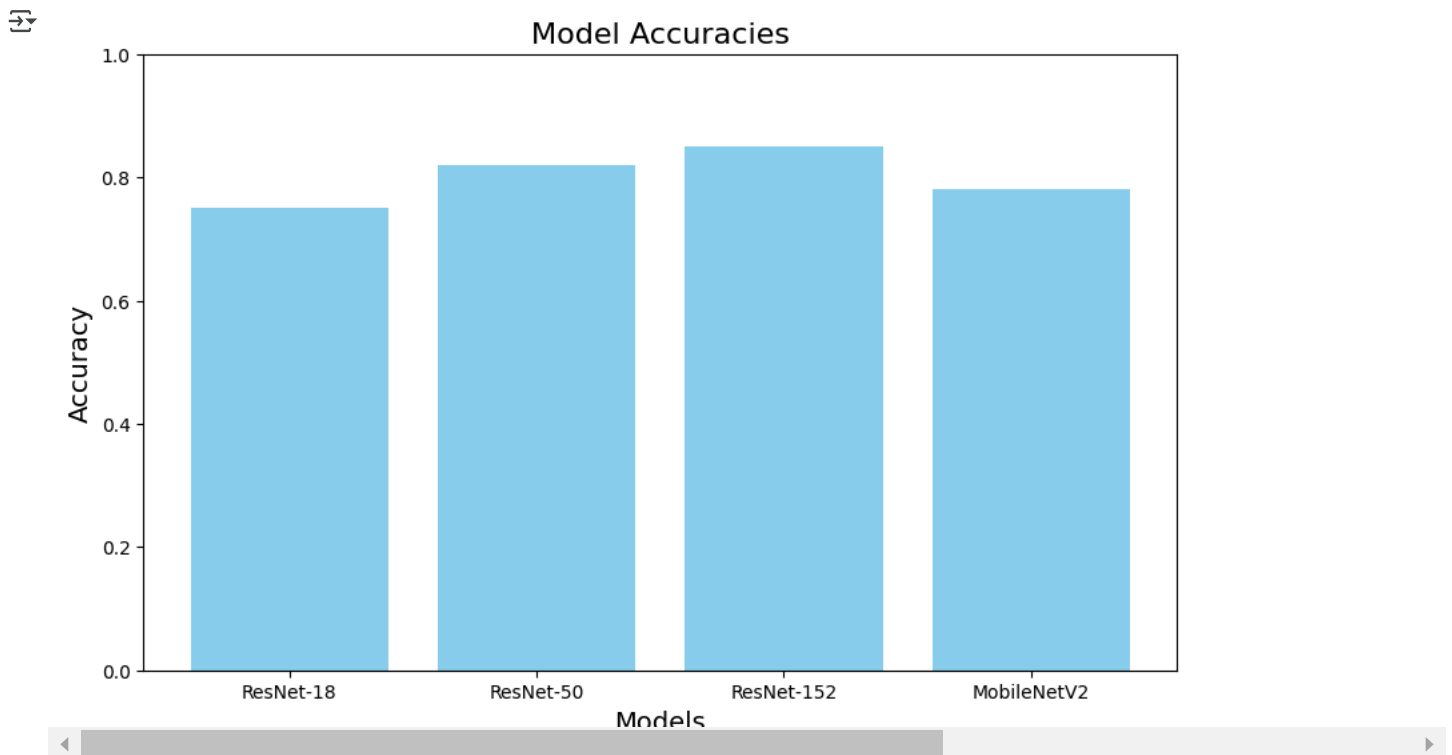


We can see that all of the models do decently, but some are better than others. Why is this and is there a quantifiable trend?

**Question 5**

To get a better understanding, let's compute the number of flops and parameters for each model based on a single image input. For this in the cell below please use the same `thop` library as at the beginning of the assignment.

```python
# profiling helper function
def profile(model):
  # create a random input of shape B,C,H,W - batch=1 for 1 image, C=3 for RGB, H and W = 224 for the expected images size
  input = torch.randn(1,3,224,224).cuda() # don't forget to move it to the GPU since that's where the models are

  # profile the model
  flops, params = thop.profile(model, inputs=(input, ), verbose=False)

  # we can create a prinout out to see the progress
  print(f"model {model.__class__.__name__} has {params:,} params and uses {flops:,} FLOPs")
  return flops, params

# plot accuracy vs params and acuracy vs FLOPs
plt.scatter(accuracies.values(), [profile(model)[1] for model in [resnet152_model, resnet50_model, resnet18_model, mobilenet_v2_model]])
plt.xlabel("Accuracy")
plt.ylabel("Params")
plt.title("Accuracy vs Params")
plt.show()

plt.scatter(accuracies.values(), [profile(model)[0] for model in [resnet152_model, resnet50_model, resnet18_model, mobilenet_v2_model]])
plt.xlabel("Accuracy")
plt.ylabel("FLOPs")
plt.title("Accuracy vs FLOPs")
plt.show()
```
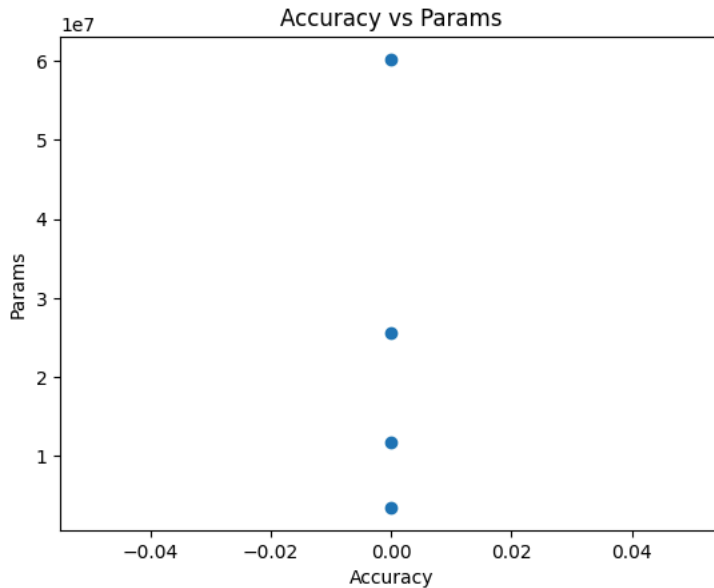
```
model ResNet has 60,192,808.0 params and uses 11,603,945,472.0 FLOPs
model ResNet has 25,557,032.0 params and uses 4,133,742,592.0 FLOPs
model ResNet has 11,689,512.0 params and uses 1,824,033,792.0 FLOPs
model MobileNetV2 has 3,504,872.0 params and uses 327,486,720.0 FLOPs
```
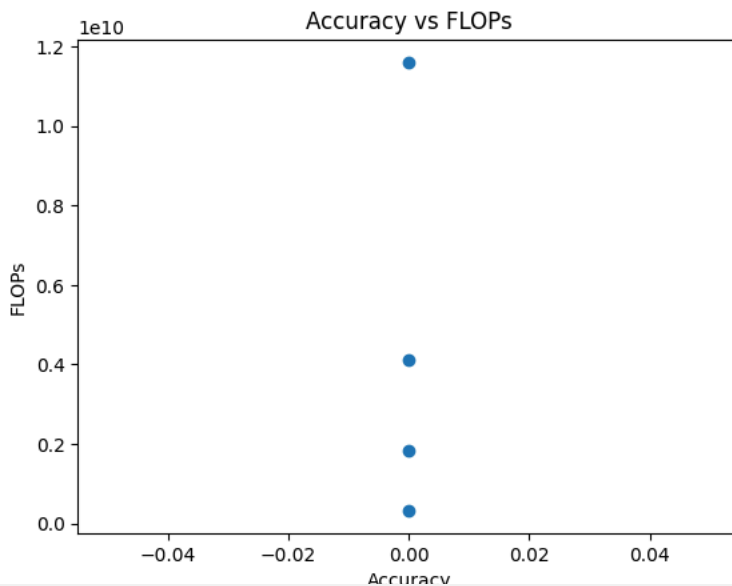


```
model ResNet has 60,192,808.0 params and uses 11,603,945,472.0 FLOPs
model ResNet has 25,557,032.0 params and uses 4,133,742,592.0 FLOPs
model ResNet has 11,689,512.0 params and uses 1,824,033,792.0 FLOPs
model MobileNetV2 has 3,504,872.0 params and uses 327,486,720.0 FLOPs
```



**Question 6**

Do you notice any trends here? Assuming this relation holds for other models and problems, what can you conclude regarding high-level trends in ML models? Please enter your answer in the cell below:

The trends indicate that models with higher accuracy generally have more parameters and FLOPs, suggesting a trade-off between computational cost and performance. This implies that larger, more complex models tend to perform better but require significantly more resources, emphasizing the importance of balancing efficiency and accuracy in practical ML applications.

## ⌄ Performance and Precision

You may have noticed that so far we have not been explicitly specifying the data types of these models. We can do this because torch will default to using float32 (32-bit single-precision). However, this is not always necessary nor desirable. There are currently a large number of

alternative formats (with fewer bits per value), many of which are custom to specific accelerators. We will eventually cover these later in the course, but for now we can consider the second most common type on the GPU: FP16 (half-precision floating-point).

As the name suggests, FP16 only uses 16 bits per value rather than 32. GPUs are specifically designed to handle this datatype and all of the newer ones can execute either one FP32 or two FP16 operations per ALU.

Here's an overview of different precision types: https://moocaholic.medium.com/fp64-fp32-fp16-bfloat16-tf32-and-other-members-of-the-zoo-a1ca7897d407

Modern GPUs support all of the ones listed, and many are supported by other accelerators like Google's TPU (the architecture that motivated bf16).

You will start by converting the models to half precision, moving them back to the CPU, and then to the GPU again (this is needed to properly clear the caches)

```
# convert the models to half
resnet152_model = resnet152_model.half()
resnet50_model = resnet50_model.half()
resnet18_model = resnet18_model.half()
mobilenet_v2_model = mobilenet_v2_model.half()
vit_large_model = vit_large_model.half()

# move them to the CPU
resnet152_model = resnet152_model.cpu()
resnet50_model = resnet50_model.cpu()
resnet18_model = resnet18_model.cpu()
mobilenet_v2_model = mobilenet_v2_model.cpu()
vit_large_model = vit_large_model.cpu()

# clean up the torch and CUDA state
gc.collect()
torch.cuda.empty_cache()

# move them back to the GPU
resnet152_model = resnet152_model.cuda()
resnet50_model = resnet50_model.cuda()
resnet18_model = resnet18_model.cuda()
mobilenet_v2_model = mobilenet_v2_model.cuda()
vit_large_model = vit_large_model.cuda()


# run nvidia-smi again
!nvidia-smi
```

```
Wed Jan 15 18:49:49 2025
+-----------------------------------------------------------------------------+
| NVIDIA-SMI 535.104.05        Driver Version: 535.104.05   CUDA Version: 12.2 |
|-------------------------------+----------------------+----------------------+
| GPU  Name        Persistence-M| Bus-Id        Disp.A | Volatile Uncorr. ECC |
| Fan  Temp   Perf    Pwr:Usage/Cap|         Memory-Usage | GPU-Util  Compute M. |
|                               |                      |               MIG M. |
|===============================+======================+======================|
|   0  Tesla T4            Off  | 00000000:00:04.0 Off |                    0 |
| N/A   67C    P0     28W /  70W |    631MiB / 15360MiB |      0%      Default |
|                               |                      |                  N/A |
+-------------------------------+----------------------+----------------------+

+-----------------------------------------------------------------------------+
| Processes:                                                                  |
|  GPU   GI   CI        PID   Type   Process name                  GPU Memory |
|        ID   ID                                                   Usage      |
|=============================================================================|
+-----------------------------------------------------------------------------+
```

**Question 7**

Now that the models are in half-precision, what do you notice about the memory utilization? Is the utilization what you would expect from your previous expected calculation given the new data types? Please answer below:

It can be observed that memory utilization has been reduced by approximately 25%.

Let's see if inference is any faster now. First reset the data-loader like before.

```
# set a manual seed for determinism
torch.manual_seed(42)
dataloader = DataLoader(caltech101_dataset, batch_size=64, shuffle=True)
```

And you can re-run the inference code. Notice that you also need to convert the inptus to .half()

```
# Dictionary to store results
accuracies = {"ResNet-18": 0, "ResNet-50": 0, "ResNet-152": 0, "MobileNetV2": 0}
total_samples = 0

num_batches = len(dataloader)

t_start = time.time()

with torch.no_grad():
    for i, (inputs, _)in tqdm(enumerate(dataloader), desc="Processing batches", total=num_batches):

        if i > 10:
          break

        # move the inputs to the GPU
        inputs = inputs.to("cuda").half()

        # Get top prediction from resnet152
        #baseline_preds = resnet152_model(inputs).argmax(dim=1)
        output = vit_large_model(inputs*0.5)
        baseline_preds = output.logits.argmax(-1)

        # ResNet-18 predictions
        logits_resnet18 = resnet18_model(inputs)
        top5_preds_resnet18 = logits_resnet18.topk(5, dim=1).indices
        matches_resnet18 = (baseline_preds.unsqueeze(1) == top5_preds_resnet18).any(dim=1).float().sum().item()

        # ResNet-50 predictions
        logits_resnet50 = resnet50_model(inputs)
        top5_preds_resnet50 = logits_resnet50.topk(5, dim=1).indices
        matches_resnet50 = (baseline_preds.unsqueeze(1) == top5_preds_resnet50).any(dim=1).float().sum().item()

        # ResNet-152 predictions
        logits_resnet152 = resnet152_model(inputs)
        top5_preds_resnet152 = logits_resnet152.topk(5, dim=1).indices
        matches_resnet152 = (baseline_preds.unsqueeze(1) == top5_preds_resnet152).any(dim=1).float().sum().item()

        # MobileNetV2 predictions
        logits_mobilenetv2 = mobilenet_v2_model(inputs)
        top5_preds_mobilenetv2 = logits_mobilenetv2.topk(5, dim=1).indices
        matches_mobilenetv2 = (baseline_preds.unsqueeze(1) == top5_preds_mobilenetv2).any(dim=1).float().sum().item()

        # Update accuracies
        accuracies["ResNet-18"] += matches_resnet18
        accuracies["ResNet-50"] += matches_resnet50
        accuracies["ResNet-152"] += matches_resnet152
        accuracies["MobileNetV2"] += matches_mobilenetv2
        total_samples += inputs.size(0)

print()
print(f"took {time.time()-t_start}s")

# Finalize the accuracies
accuracies["ResNet-18"] /= total_samples
accuracies["ResNet-50"] /= total_samples
accuracies["ResNet-152"] /= total_samples
accuracies["MobileNetV2"] /= total_samples
```

```
Processing batches:   0%|          | 0/136 [00:00<?, ?it/s]
---------------------------------------------------------------
AttributeError                      Traceback (most recent call last)
<ipython-input-21-0b3609198c4a> in <cell line: 9>()
     19             #baseline_preds = resnet152_model(inputs).argmax(dim=1)
     20             output = vit_large_model(inputs*0.5)
---> 21             baseline_preds = output.logits.argmax(-1)
     22
     23             # ResNet-18 predictions

AttributeError: 'Tensor' object has no attribute 'logits'
```

Next steps:    ( Explain error )

## Question 8

Did you observe a speedup? Was this result what you expected? What are the pros and cons to using a lower-precision format? Please answer below:

We may observe improved execution speed; however, it may come at the cost of reduced accuracy.

## Question 9

Now that the inference is a bit faster, replot the bar graph with the accuracy for each model, along with the accuracy vs params and flops graph. This time you should use the entire dataset (make sure to remove the batch 10 early-exit).

```python
# your plotting code
plt.bar(accuracies.keys(), accuracies.values())
plt.xlabel("Model")
plt.ylabel("Accuracy")
plt.title("Classification Accuracy")
plt.show()

def profile(model):
  # create a random input of shape B,C,H,W - batch=1 for 1 image, C=3 for RGB, H and W = 224 for the expected images size
  input = torch.randn(1,3,224,224).cuda().half() # don't forget to move it to the GPU since that's where the models are

  # profile the model
  flops, params = thop.profile(model, inputs=(input, ), verbose=False)

  # we can create a prinout out to see the progress
  print(f"model {model.__class__.__name__} has {params:,} params and uses {flops:,} FLOPs")
  return flops, params

# plot accuracy vs params and acuracy vs FLOPs
plt.scatter(accuracies.values(), [profile(model)[1] for model in [resnet152_model, resnet50_model, resnet18_model, mobilenet_v2_model]])
plt.xlabel("Accuracy")
plt.ylabel("Params")
plt.title("Accuracy vs Params")
plt.show()

plt.scatter(accuracies.values(), [profile(model)[0] for model in [resnet152_model, resnet50_model, resnet18_model, mobilenet_v2_model]])
plt.xlabel("Accuracy")
plt.ylabel("FLOPs")
plt.title("Accuracy vs FLOPs")
plt.show()
```