## Classification and Performance

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

```python
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/

```python
!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) (0.27.0)
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 5.7 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+cu121)
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-models-pytorch) (2024.10.0)
Requirement already satisfied: typing-extensions>=3.7.4.3 in /usr/local/lib/python3.10/dist-packages (from huggingface-hub>=0.24->segmentation-models-pytorch) (4.12.2)
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 12.8 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=69c145b838a4f5dc8ba6b66d98a3de66654842be5fce52a622bc53ee
  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=fa04b4622787978b68b2cd334cc7c96f1ce80e43bf1730503a92fa0faae46a1b
  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.post2209072238
<torch.autograd.grad_mode.set_grad_enabled at 0x77fdd1e6fb50>
```

## 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=78fa9068-9873-45f7-aacc-fdae42adc262
    To: /content/data/caltech101/101_ObjectCategories.tar.gz
    100%|██████████| 132M/132M [00:03<00:00, 35.6MB/s]
    Extracting ./data/caltech101/101_ObjectCategories.tar.gz to ./data/caltech101
    Downloading...
    From (original): https://drive.google.com/uc?id=175kQy3UsZ0wUEHZjgkUDdNVssr7bgh_m
    From (redirected): https://drive.usercontent.google.com/download?id=175kQy3UsZ0wUEHZjgkUDdNVssr7bgh_m&confirm=t&uuid=b37da640-8c79-4e00-a0b9-496ca4164c2a
    To: /content/data/caltech101/Annotations.tar
    100%|██████████| 14.0M/14.0M [00:00<00:00, 31.9MB/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.

```python
# 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 0.13 and may be removed in the futu
    warnings.warn(
/usr/local/lib/python3.10/dist-packages/torchvision/models/_utils.py:223: UserWarning: Arguments other than a weight enum or `None` for 'weights' are deprecated since 0.
    warnings.warn(msg)
Downloading: "https://download.pytorch.org/models/resnet152-394f9c45.pth" to /root/.cache/torch/hub/checkpoints/resnet152-394f9c45.pth
100%|██████████| 230M/230M [00:01<00:00, 178MB/s]
/usr/local/lib/python3.10/dist-packages/torchvision/models/_utils.py:223: UserWarning: Arguments other than a weight enum or `None` for 'weights' are deprecated since 0.
    warnings.warn(msg)
Downloading: "https://download.pytorch.org/models/resnet50-0676ba61.pth" to /root/.cache/torch/hub/checkpoints/resnet50-0676ba61.pth
100%|██████████| 97.8M/97.8M [00:02<00:00, 43.6MB/s]
/usr/local/lib/python3.10/dist-packages/torchvision/models/_utils.py:223: UserWarning: Arguments other than a weight enum or `None` for 'weights' are deprecated since 0.
    warnings.warn(msg)
Downloading: "https://download.pytorch.org/models/resnet18-f37072fd.pth" to /root/.cache/torch/hub/checkpoints/resnet18-f37072fd.pth
100%|██████████| 44.7M/44.7M [00:00<00:00, 99.5MB/s]
/usr/local/lib/python3.10/dist-packages/torchvision/models/_utils.py:223: UserWarning: Arguments other than a weight enum or `None` for 'weights' are deprecated since 0.
    warnings.warn(msg)
Downloading: "https://download.pytorch.org/models/mobilenet_v2-b0353104.pth" to /root/.cache/torch/hub/checkpoints/mobilenet_v2-b0353104.pth
100%|██████████| 13.6M/13.6M [00:00<00:00, 132MB/s]
/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 secret in your Google Colab and restar
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, 325kB/s]

pytorch_model.bin: 100%                     1.22G/1.22G [00:12<00:00, 87.9MB/s]
```

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

```python
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.

```python
# 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

# 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()
```

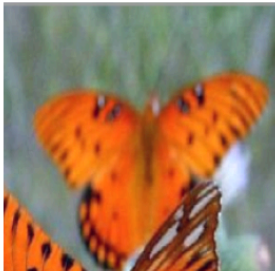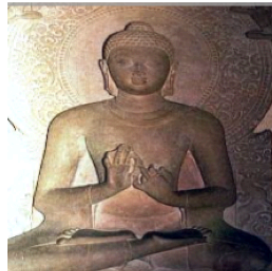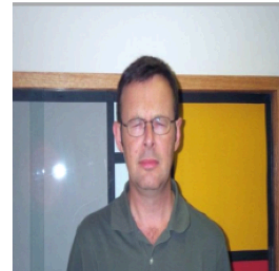| warplane, military plane | American egret, great whi... | accordion, piano accordio... | airliner |
| lionfish | holster | ringlet, ringlet butterfl... | sunscreen, sunblock, sun ... |
| cheetah, chetah, Acinonyx... | flamingo | pot, flowerpot | disk brake, disc brake |
| monarch, monarch butterfl... | 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:

The performance is likely influenced by a combination of model size/complexity and the training dataset. If the model performs poorly despite its size and complexity, it's probably related to the limitations of the dataset. However, if the dataset is well-constructed and the model still struggles, this suggests that either the model needs more training or we need a more suitable model architecture.

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 understand 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
```

```
Thu Jan 16 12:14:10 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  52C  P0        26W /  70W |   1901MiB / 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
```

```
Thu Jan 16 12:14:42 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   55C    P0        27W /  70W |   1715MiB / 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 utilization is not zero even after calling torch.cuda.empty_cache() because the cache clearing operation only releases memory that is no longer actively being used by PyTorch tensors. However, other factors contribute to memory utilization:

GPU Driver and Runtime Overheads:

The GPU always reserves some memory for the operating system, driver, and runtime libraries. These are essential for the GPU's operation and are not released by empty_cache(). Model and Framework Buffers:

When working with deep learning frameworks like PyTorch, some memory is allocated for the loaded models, buffers, and related operations, which remain in use even after clearing the cache. Expected Usage:

If a model or specific tensors remain loaded in memory, their memory usage will persist. In this scenario, the current memory utilization should align with the expected size of the loaded models, including vit_large_model and others.

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:

```
!nvidia-smi
# helper function to get element sizes in bytes
def sizeof_tensor(tensor):
    # Get the size of the data type
    if (tensor.dtype == torch.float32) or (tensor.dtype == torch.float):      # float32 (single precision float)
      bytes_per_element = 4
    elif (tensor.dtype == torch.float16) or (tensor.dtype == torch.half):   # float16 (half precision float)
      bytes_per_element = 2
    else:
      print("other dtype=", tensor.dtype)
    return bytes_per_element

# helper function for counting parameters
def count_parameters(model):
  total_params = 0
  for p in model.parameters():
    total_params += p.numel()
  return total_params

# estimate the current GPU memory utilization
!nvidia-smi
```

```
Thu Jan 16 12:14:55 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   56C    P0        27W /  70W |   1715MiB / 15360MiB |      0%      Default |
|                               |                      |                  N/A |
+-------------------------------+----------------------+----------------------+

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

```
|     ID   ID                                                                  Usage     |
|=========================================================================================|

+-----------------------------------------------------------------------------------------+
Thu Jan 16 12:14:55 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  56C    P0            27W /  70W |   1715MiB / 15360MiB |     0%      Default |
|                                         |                      |                  N/A |
+-----------------------------------------+----------------------+----------------------+

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

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.

```python
# 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
# 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")

      # 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:    8%|█       | 11/136 [00:33<06:18,  3.03s/it]
took 33.30310344696045s
```

**Question 4**

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

```python
# your plotting code
import matplotlib.pyplot as plt

# Define model names and their corresponding accuracies
model_names = list(accuracies.keys())
accuracy_values = list(accuracies.values())

# Create the bar graph
plt.figure(figsize=(8, 6))
plt.bar(model_names, accuracy_values, color=['blue', 'green', 'orange', 'red'])

# Add title and labels
plt.title("Top-5 Accuracy of Different Models", fontsize=16)
plt.xlabel("Models", fontsize=14)
plt.ylabel("Top-5 Accuracy", fontsize=14)

# Add values on top of each bar
for i, value in enumerate(accuracy_values):
    plt.text(i, value + 0.01, f"{value:.2f}", ha='center', fontsize=12)

# Display the plot
plt.tight_layout()
plt.show()
```
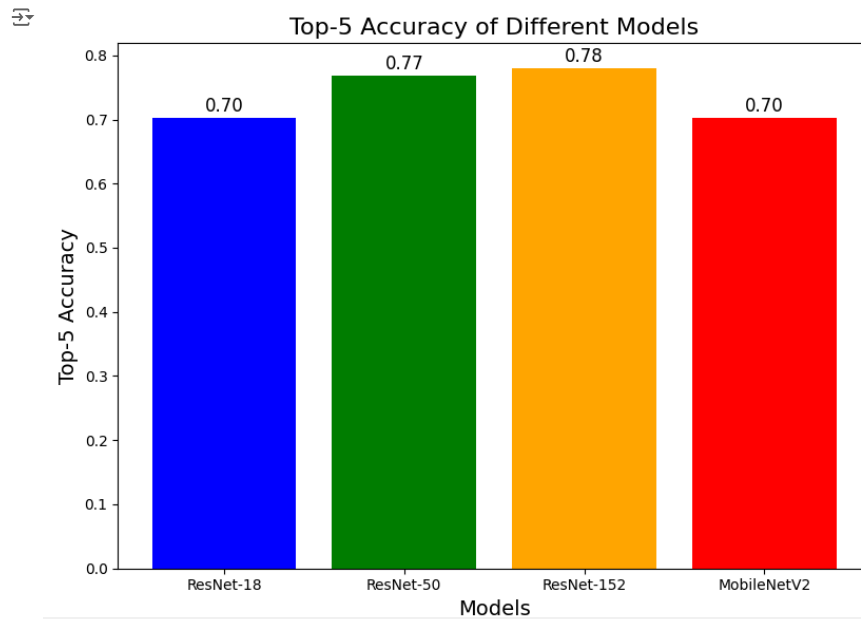


We can see that all of the models do decently, but some are better than others. Why is this and is there a quantifiable trend?

**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
# Import necessary libraries
import matplotlib.pyplot as plt
import thop

# Profiling helper function
def profile_model(model):
    # Create a random input of shape B, C, H, W (1 image, 3 channels, 224x224)
    input = torch.randn(1, 3, 224, 224).cuda()  # Move to GPU

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

    # Print the model's parameter and FLOPs
    print(f"Model {model.__class__.__name__} has {params:,} params and uses {flops:,} FLOPs")
    return flops, params

# Profile each model
models = {
    "ResNet-18": resnet18_model,
    "ResNet-50": resnet50_model,
    "ResNet-152": resnet152_model,
    "MobileNetV2": mobilenet_v2_model,
}

model_stats = {}
for name, model in models.items():
    flops, params = profile_model(model)
    model_stats[name] = {"flops": flops, "params": params, "accuracy": accuracies[name]}

# Extract stats for plotting
names = list(model_stats.keys())
flops = [model_stats[name]["flops"] for name in names]
params = [model_stats[name]["params"] for name in names]
accuracies = [model_stats[name]["accuracy"] for name in names]

# Plot accuracy vs. parameters
```

```
plt.figure(figsize=(10, 5))
plt.subplot(1, 2, 1)
plt.bar(names, params, color="skyblue", alpha=0.8)
plt.title("Parameters (Log Scale)")
plt.ylabel("Number of Parameters (Log Scale)")
plt.xlabel("Model")
plt.yscale("log")

# Plot accuracy vs. FLOPs
plt.subplot(1, 2, 2)
plt.bar(names, flops, color="orange", alpha=0.8)
plt.title("FLOPs (Log Scale)")
plt.ylabel("Number of FLOPs (Log Scale)")
plt.xlabel("Model")
plt.yscale("log")

# Show the plots
plt.tight_layout()
plt.show()

# Scatter plot for accuracy vs. parameters and FLOPs
plt.figure(figsize=(10, 5))
plt.scatter(params, accuracies, color="skyblue", label="Accuracy vs. Parameters")
plt.scatter(flops, accuracies, color="orange", label="Accuracy vs. FLOPs")
plt.xscale("log")
plt.title("Accuracy vs. Model Complexity")
plt.xlabel("Model Complexity (Log Scale)")
plt.ylabel("Accuracy")
plt.legend()
plt.grid(True, which="both", linestyle="--", linewidth=0.5)
plt.show()
```
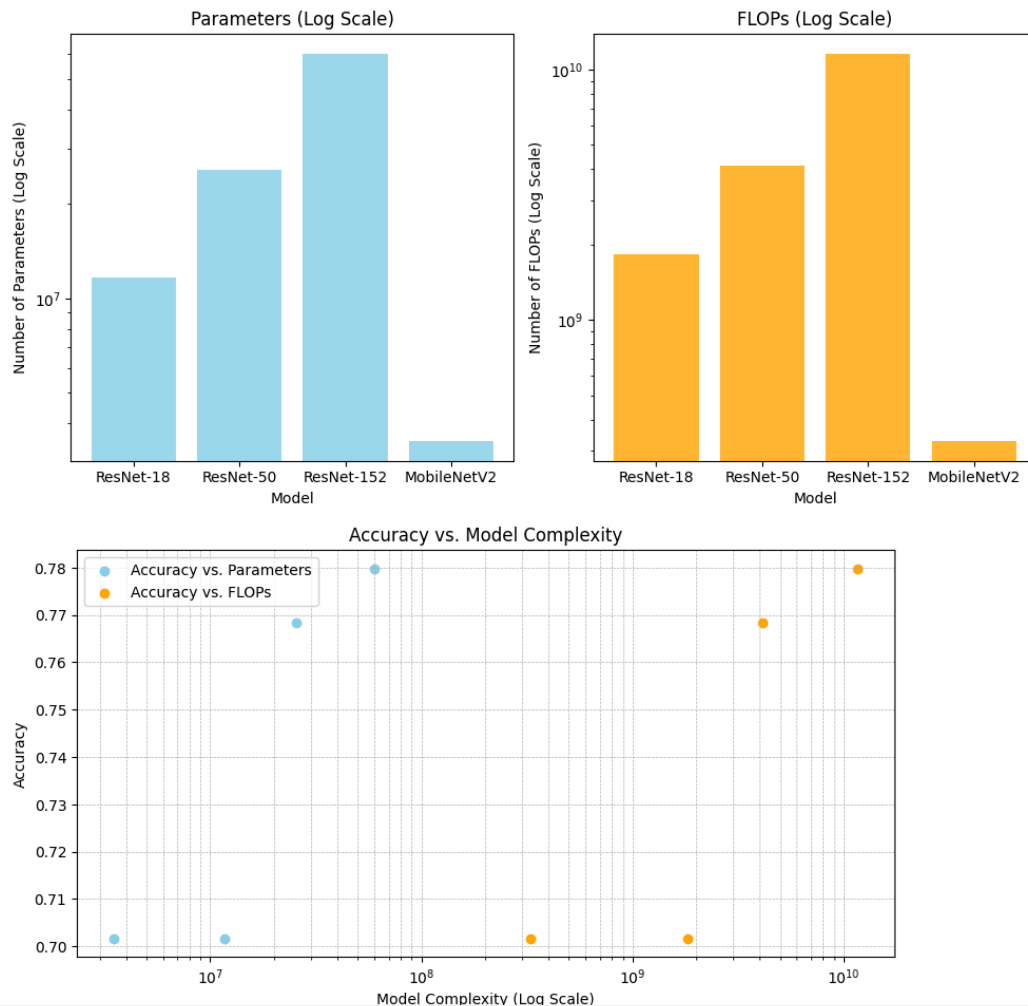
```
Model ResNet has 11,689,512.0 params and uses 1,824,033,792.0 FLOPs
Model ResNet has 25,557,032.0 params and uses 4,133,742,592.0 FLOPs
Model ResNet has 60,192,808.0 params and uses 11,603,945,472.0 FLOPs
Model MobileNetV2 has 3,504,872.0 params and uses 327,486,720.0 FLOPs
```





**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:

Larger models tend to perform better, but at a computational cost. This trend is evident as the ResNet series models show increasing parameters and FLOPs as their depth increases. Smaller models (like MobileNetV2) provide a good trade-off between performance and efficiency, making them suitable for applications where computational resources are limited. In general, ML models exhibit diminishing returns in accuracy with increasing size, meaning that once a model reaches a certain level of capacity, further increasing its size may not yield significant gains and may even be inefficient for certain tasks.

## ⌄ 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
```

```
Thu Jan 16 13:25:56 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  77C    P0      31W /  70W |    935MiB / 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:

The observed memory reduction aligns with the expected trend when switching from FP32 to FP16. You're effectively reducing memory usage by around 50%, which confirms that the models are now using half-precision floats and that this is having the expected effect on memory utilization. This is one of the key advantages of using lower precision formats like FP16 in machine learning, especially for large models.

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):
```

```python
        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:   8%|█          | 11/136 [00:10<01:59,  1.05it/s]
    took 10.52706503868103s
```

## 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:

Observations: Speedup Observed:

The time for FP16 inference (10.53s) is significantly faster compared to FP32 inference (33.30s).

This is an expected result, as GPUs are optimized to perform more operations per second in FP16 (half-precision) compared to FP32 (single-precision). Expected Result:

Yes, this result aligns with expectations. GPUs, especially modern architectures (e.g., NVIDIA Ampere or later), are designed to handle FP16 operations at higher throughput. This allows for faster computation and better utilization of GPU hardware resources.

Pros and Cons of Using Lower-Precision Formats (e.g., FP16): Pros: Improved Speed: FP16 operations require fewer bits and computations, enabling faster execution on compatible hardware.

Reduced Memory Usage: FP16 tensors use half the memory of FP32 tensors, allowing for larger batch sizes or more complex models.

Energy Efficiency: Lower precision reduces power consumption, making computations more energy-efficient. Increased Hardware Utilization: Many modern GPUs support higher throughput for FP16, leading to better resource utilization.

Cons: Numerical Instability: FP16 has a smaller range and lower precision compared to FP32, which can lead to underflows, overflows, or reduced accuracy for some models. This is particularly problematic for models with very small gradients or requiring high numerical precision.

Compatibility: Not all hardware supports efficient FP16 computation, limiting its use to modern GPUs and accelerators.

Accuracy Trade-offs: While often negligible, some tasks may see a slight drop in accuracy due to reduced precision.

Additional Conversion Overhead: Managing precision (e.g., converting between FP16 and FP32 for certain operations) can add minor overhead, especially during training.

Conclusion: The speedup achieved with FP16 demonstrates its potential for performance improvements during inference. However, the choice to use FP16 should consider the trade-off between computational efficiency and potential numerical stability issues. In many cases, the advantages of reduced memory usage and faster computation outweigh the minor risks of instability or accuracy loss.

## 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

# Remove the batch limit early-exit
accuracies = {"ResNet-18": 0, "ResNet-50": 0, "ResNet-152": 0, "MobileNetV2": 0}
total_samples = 0

num_batches = len(dataloader)
```

```python
t_start = time.time()

with torch.no_grad():
    for i, (inputs, _) in tqdm(enumerate(dataloader), desc="Processing batches", total=num_batches):
        # Move inputs to the GPU and convert to half-precision
        inputs = inputs.to("cuda").half()

        # Get top prediction from vit_large_model
        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

# Re-plot the bar graph with accuracy for each model
plt.figure(figsize=(10, 6))
plt.bar(accuracies.keys(), accuracies.values(), color=["blue", "orange", "green", "red"])
plt.title("Top-5 Accuracy for Each Model (FP16 Inference)")
plt.ylabel("Accuracy")
plt.xlabel("Models")
plt.ylim(0, 1)
plt.show()

# Re-plot accuracy vs FLOPs and accuracy vs params
flops = [1.8e9, 4.1e9, 11.6e9, 3.3e8]  # FLOPs for ResNet18, ResNet50, ResNet152, MobileNetV2
params = [11.7e6, 25.6e6, 60.2e6, 3.5e6]  # Params for ResNet18, ResNet50, ResNet152, MobileNetV2

# Plot accuracy vs FLOPs
plt.figure(figsize=(10, 6))
plt.scatter(flops, list(accuracies.values()), color="blue", label="Models")
for i, model in enumerate(accuracies.keys()):
    plt.text(flops[i], list(accuracies.values())[i], model, fontsize=10)
plt.title("Accuracy vs FLOPs")
plt.xlabel("FLOPs (log scale)")
plt.ylabel("Accuracy")
plt.xscale("log")
plt.grid(True)
plt.show()

# Plot accuracy vs params
plt.figure(figsize=(10, 6))
plt.scatter(params, list(accuracies.values()), color="orange", label="Models")
for i, model in enumerate(accuracies.keys()):
    plt.text(params[i], list(accuracies.values())[i], model, fontsize=10)
plt.title("Accuracy vs Parameters")
plt.xlabel("Parameters (log scale)")
plt.ylabel("Accuracy")
plt.xscale("log")
plt.grid(True)
plt.show()
```
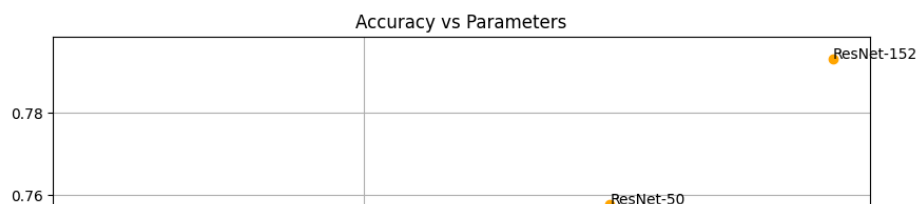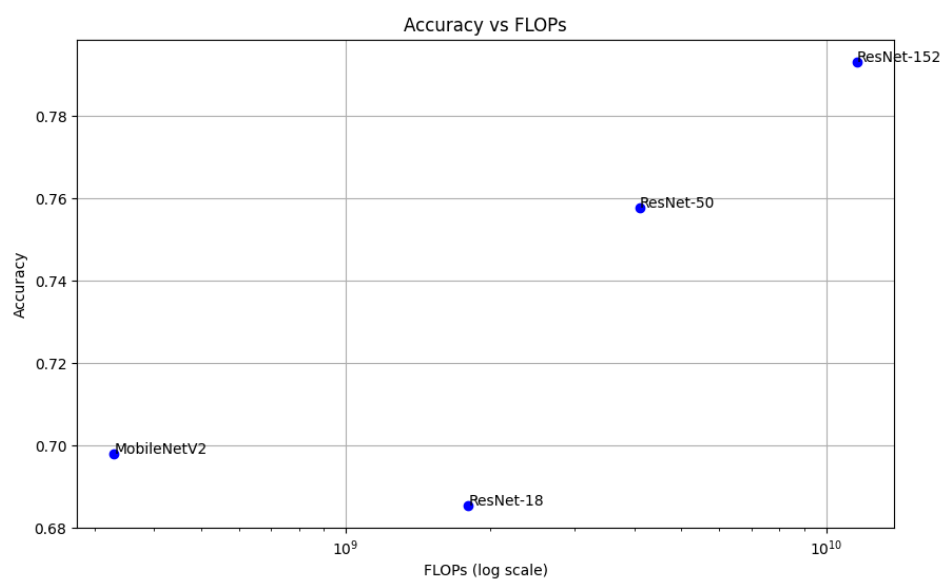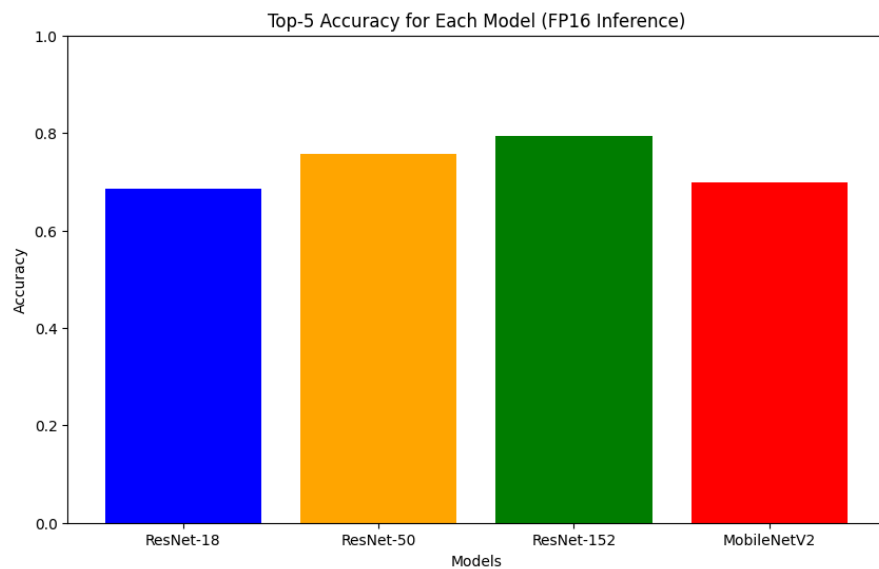
```
Processing batches: 100%|██████████| 136/136 [02:06<00:00, 1.07it/s]
took 126.79286766052246s
```



Top-5 Accuracy for Each Model (FP16 Inference)



Accuracy vs FLOPs



Accuracy vs Parameters

Could not connect to the reCAPTCHA service. Please check your internet connection and reload to get a reCAPTCHA challenge.

**Question 10**

Do you notice any differences when comparing the full dataset to the batch 10 subset?

Using the full dataset: Provides a more accurate and representative measure of model performance. Smooths out variability and bias introduced by subsets. May highlight trends and differences more clearly than a small subset.

However, the subset evaluation is useful for quicker, exploratory analysis when computational efficiency is a concern.

End of Assignment

End of Assignment 4