

## 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 pretrainedmodels>=0.7.1 (from segmentation-models-pytorch)
  Downloading pretrainedmodels-0.7.4.tar.gz (58 kB)
  └── Preparing metadata (setup.py) ... done
    Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.11/dist-packages (from segmentation-models-pytorch) (1.17.0)
    Requirement already satisfied: timm>=0.9 in /usr/local/lib/python3.11/dist-packages (from segmentation-models-pytorch) (1.0.13)
    Requirement already satisfied: torchvision>=0.9 in /usr/local/lib/python3.11/dist-packages (from segmentation-models-pytorch) (0.1.1)
    Requirement already satisfied: tqdm>=4.42.1 in /usr/local/lib/python3.11/dist-packages (from segmentation-models-pytorch) (4.67.1)
    Requirement already satisfied: filelock in /usr/local/lib/python3.11/dist-packages (from transformers) (3.16.1)
    Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.11/dist-packages (from transformers) (24.2)
    Requirement already satisfied: pyyaml>=5.1 in /usr/local/lib/python3.11/dist-packages (from transformers) (6.0.2)
    Requirement already satisfied: regex!=2019.12.17 in /usr/local/lib/python3.11/dist-packages (from transformers) (2024.11.6)
    Requirement already satisfied: requests in /usr/local/lib/python3.11/dist-packages (from transformers) (2.32.3)
    Requirement already satisfied: tokenizers<0.22,>=0.21 in /usr/local/lib/python3.11/dist-packages (from transformers) (0.21.0)
    Requirement already satisfied: safetensors>=0.4.1 in /usr/local/lib/python3.11/dist-packages (from transformers) (0.5.2)
    Requirement already satisfied: fsspec>=2023.5.0 in /usr/local/lib/python3.11/dist-packages (from huggingface-hub>=0.24->segmentation-models)
    Requirement already satisfied: typing-extensions>=3.7.4.3 in /usr/local/lib/python3.11/dist-packages (from huggingface-hub>=0.24->segmentation-models)
    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/python3.11/dist-packages (from torch->thop) (3.4.2)
    Requirement already satisfied: jinja2 in /usr/local/lib/python3.11/dist-packages (from torch->thop) (3.1.5)
    Requirement already satisfied: nvidia-cuda-nvrtc-cu12==12.1.105 in /usr/local/lib/python3.11/dist-packages (from torch->thop) (12.1.105)
    Requirement already satisfied: nvidia-cuda-runtime-cu12==12.1.105 in /usr/local/lib/python3.11/dist-packages (from torch->thop) (12.1.105)
    Requirement already satisfied: nvidia-cuda-cupti-cu12==12.1.105 in /usr/local/lib/python3.11/dist-packages (from torch->thop) (12.1.105)
    Requirement already satisfied: nvidia-cudnn-cu12==9.1.0.70 in /usr/local/lib/python3.11/dist-packages (from torch->thop) (9.1.0.70)
    Requirement already satisfied: nvidia-cUBLAS-cu12==12.1.3.1 in /usr/local/lib/python3.11/dist-packages (from torch->thop) (12.1.3.1)
    Requirement already satisfied: nvidia-cufft-cu12==11.0.2.54 in /usr/local/lib/python3.11/dist-packages (from torch->thop) (11.0.2.54)
    Requirement already satisfied: nvidia-curand-cu12==10.3.2.106 in /usr/local/lib/python3.11/dist-packages (from torch->thop) (10.3.2.106)
    Requirement already satisfied: nvidia-cusolver-cu12==11.4.5.107 in /usr/local/lib/python3.11/dist-packages (from torch->thop) (11.4.5.107)
    Requirement already satisfied: nvidia-cusparse-cu12==12.1.0.106 in /usr/local/lib/python3.11/dist-packages (from torch->thop) (12.1.0.106)
    Requirement already satisfied: nvidia-nccl-cu12==2.21.5 in /usr/local/lib/python3.11/dist-packages (from torch->thop) (2.21.5)
    Requirement already satisfied: nvidia-nvtx-cu12==12.1.105 in /usr/local/lib/python3.11/dist-packages (from torch->thop) (12.1.105)
    Requirement already satisfied: triton==3.1.0 in /usr/local/lib/python3.11/dist-packages (from torch->thop) (3.1.0)
    Requirement already satisfied: sympy==1.13.1 in /usr/local/lib/python3.11/dist-packages (from torch->thop) (1.13.1)
    Requirement already satisfied: nvidia-nvjitlink-cu12 in /usr/local/lib/python3.11/dist-packages (from nvidia-cusolver-cu12==11.4.5.107)
    Requirement already satisfied: mpmath<1.4,>=1.1.0 in /usr/local/lib/python3.11/dist-packages (from sympy==1.13.1->torch->thop) (1.1.0)
    Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.11/dist-packages (from requests->transformers)
```

```
Building wheel for pretrainedmodels (setup.py) ... done
Created wheel for pretrainedmodels: filename=pretrainedmodels-0.7.4-py3-none-any.whl size=60944 sha256=c526f93af911907eaeeed8
Stored in directory: /root/.cache/pip/wheels/5f/5b/96/fd94bc35962d7c6b699e8814db545155ac91d2b95785e1b035
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
<torch.autograd.grad_mode.set_grad_enabled at 0x7ea89220bcd0>
```

## 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](https://en.wikipedia.org/wiki/Caltech_101)

Download the dataset you will be using: Caltech101

```
# 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=137RyRjvTBkB1IfFeYBNZBtViDH06\\_Ewsp](https://drive.google.com/uc?id=137RyRjvTBkB1IfFeYBNZBtViDH06_Ewsp)  
 From (redirected): [https://drive.usercontent.google.com/download?id=137RyRjvTBkB1IfFeYBNZBtViDH06\\_Ewsp&confirm=t&uuid=a44c008c-315f-4](https://drive.usercontent.google.com/download?id=137RyRjvTBkB1IfFeYBNZBtViDH06_Ewsp&confirm=t&uuid=a44c008c-315f-4)  
 To: /content/data/caltech101/101\_ObjectCategories.tar.gz  
 100% [██████████] 132M/132M [00:04<00:00, 31.6MB/s]  
 Extracting ./data/caltech101/101\_ObjectCategories.tar.gz to ./data/caltech101  
 Downloading...  
 From (original): [https://drive.google.com/uc?id=175k0y3UsZ0wUEHZjgkUddNVssr7bgh\\_m](https://drive.google.com/uc?id=175k0y3UsZ0wUEHZjgkUddNVssr7bgh_m)  
 From (redirected): [https://drive.usercontent.google.com/download?id=175k0y3UsZ0wUEHZjgkUddNVssr7bgh\\_m&confirm=t&uuid=177c5d22-e75d-4](https://drive.usercontent.google.com/download?id=175k0y3UsZ0wUEHZjgkUddNVssr7bgh_m&confirm=t&uuid=177c5d22-e75d-4)  
 To: /content/data/caltech101/Annotations.tar  
 100% [██████████] 14.0M/14.0M [00:00<00:00, 28.8MB/s]  
 Extracting ./data/caltech101/Annotations.tar to ./data/caltech101

```
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/python3.11/dist-packages/torchvision/models/_utils.py:208: UserWarning: The parameter 'pretrained' is deprecated since 0.13.0. Please use 'weights' instead.
  warnings.warn(
/usr/local/lib/python3.11/dist-packages/torchvision/models/_utils.py:223: UserWarning: Arguments other than a weight enum or `None` were provided for argument `weights`. These will be ignored.
  warnings.warn(msg)
Downloading: "https://download.pytorch.org/models/resnet152-394f9c45.pth" to /root/.cache/torch/hub/checkpoints/resnet152-394f9c45.pth
100%|██████████| 230M/230M [00:04<00:00, 53.6MB/s]
/usr/local/lib/python3.11/dist-packages/torchvision/models/_utils.py:223: UserWarning: Arguments other than a weight enum or `None` were provided for argument `weights`. These will be ignored.
  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 [00:01<00:00, 81.4MB/s]
/usr/local/lib/python3.11/dist-packages/torchvision/models/_utils.py:223: UserWarning: Arguments other than a weight enum or `None` were provided for argument `weights`. These will be ignored.
  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, 173MB/s]
/usr/local/lib/python3.11/dist-packages/torchvision/models/_utils.py:223: UserWarning: Arguments other than a weight enum or `None` were provided for argument `weights`. These will be ignored.
  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, 179MB/s]
/usr/local/lib/python3.11/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 an environment variable, and then run `huggingface_hub.login()` again.
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, 4.89MB/s]
pytorch_model.bin: 100%                                    1.22G/1.22G [00:05<00:00, 249MB/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](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

# 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



lionfish



cheetah, cheta, Acinonyx...



monarch, monarch butterfl...



American egret, great whi...



holster



flamingo



stretcher



accordion, piano accordio...



ringlet, ringlet butterfl...



pot, flowerpot



book jacket, dust cover, ...



airliner



sunscreen, sunblock, sun ...



disk brake, disc brake



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:

## ▼ How Well Does the Model Perform?

**Accurate Predictions:** If the majority of the predictions align well with the visual content of the images, it suggests that the model is effectively leveraging the features it learned during training. This reflects good performance. **Misclassifications:** If there are notable discrepancies between the predicted labels and the images, it indicates areas for improvement. Observed Limitations

## Inconsistent Predictions:

If the model fails to identify certain images correctly or produces seemingly random predictions, this could indicate limitations in generalization. Specific patterns of errors (e.g., confusing similar classes) could point to a lack of distinction in the model's learned features.

### Label Truncation:

Long class names being truncated for visualization might obscure finer details during analysis, but this does not directly affect the model's underlying performance.

### Edge Cases or Ambiguity:

Errors in complex or ambiguous images might indicate that the model struggles with subtle distinctions, possibly due to a lack of sufficient examples of such cases in the training set.

### Class-Specific Weaknesses:

If errors are concentrated in certain classes, it could point to underrepresentation or poor-quality data for those classes in the training set.

## Source of Limitations: Model or Dataset?

### Model Size and Complexity:

ViT-L/16 is a large, powerful model designed for tasks requiring high capacity and large datasets. If the dataset is small, the model might overfit, learning specific details rather than general patterns. A smaller or simpler model might outperform ViT-L/16 on a limited or simpler dataset by avoiding overfitting.

### Dataset Quality and Quantity:

**Small or Imbalanced Dataset:** If the dataset lacks sufficient diversity or has class imbalances, even a well-designed model will struggle to generalize. **Insufficient Augmentation:** Without robust augmentation, the model may not learn to handle variations in data such as lighting, orientation, or scale.

### Training and Inference Mismatch:

Differences in preprocessing between training (e.g., normalization) and inference can reduce the model's accuracy. Since the code mentions normalization differences for ViT-L/16, this could be contributing to some errors.

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 17:54:17 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   54C   P0    26W /  70W | 1901MiB / 15360MiB |     0%      Default |
|                               |                N/A       |
+-----+
+-----+
| Processes:                               |
| GPU  GI  CI        PID  Type  Process name                  GPU Memory |
| ID   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 17:55:22 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 49C P0	26W / 70W	1715MiB / 15360MiB	0% Default
			N/A

Processes:					
GPU ID	GI ID	CI	PID	Type	Process name
					GPU Memory 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 because some processes or libraries have allocated memory on the GPU, even though no active processes are listed under "Processes". This typically happens due to the following reasons:

## ▼ Possible Causes of GPU Memory Utilization:

### CUDA Initialization:

When a program initializes CUDA (e.g., by importing PyTorch or TensorFlow), it allocates some memory on the GPU to set up the runtime environment, even if no computation is currently happening.

### Framework Overhead:

Machine learning frameworks like PyTorch or TensorFlow often pre-allocate a portion of GPU memory for caching and managing tensors. This allows for faster memory allocation during training or inference.

### Driver Overhead:

GPU drivers may reserve some memory to manage internal tasks, such as device monitoring, context creation, and kernel management.

## Does the Utilization Match Expectations?

Yes, the current utilization of 1715 MiB (out of 15360 MiB) is within expected ranges, considering:

- The GPU (Tesla T4) is in an idle state but may have memory allocated by libraries or frameworks.
- CUDA initialization typically reserves 200-500 MiB, and frameworks like PyTorch or TensorFlow can increase this to over 1 GB, depending on the caching strategy.

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:

```
# 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

def estimate_memory_utilization(model, input_tensor):
    # Check if the model and tensor are on GPU
    if not next(model.parameters()).is_cuda:
        print("Model is not on GPU. Moving it to GPU...")
        model = model.cuda()
    if not input_tensor.is_cuda:
        print("Input tensor is not on GPU. Moving it to GPU...")
        input_tensor = input_tensor.cuda()

    # Count the number of model parameters
    total_params = count_parameters(model)
    print(f"Total Parameters in Model: {total_params}")

    # Compute size of model parameters in memory
    parameter_size_bytes = sum(p.numel() * sizeof_tensor(p) for p in model.parameters())
    parameter_size_mib = parameter_size_bytes / (1024**2) # Convert to MiB
    print(f"Model Parameter Size (MiB): {parameter_size_mib:.2f}")

    # Compute size of the input tensor in memory
    input_size_bytes = input_tensor.numel() * sizeof_tensor(input_tensor)
    input_size_mib = input_size_bytes / (1024**2) # Convert to MiB
    print(f"Input Tensor Size (MiB): {input_size_mib:.2f}")

    # PyTorch memory tracking
    allocated_memory_mib = torch.cuda.memory_allocated() / (1024**2)
    reserved_memory_mib = torch.cuda.memory_reserved() / (1024**2)

    print(f"Allocated GPU Memory (MiB): {allocated_memory_mib:.2f}")
    print(f"Reserved GPU Memory (MiB): {reserved_memory_mib:.2f}")
```

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.

```
# 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 [00:32<06:06, 2.94s/it]  
took 32.29969382286072s

#### Question 4

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

```

import matplotlib.pyplot as plt

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

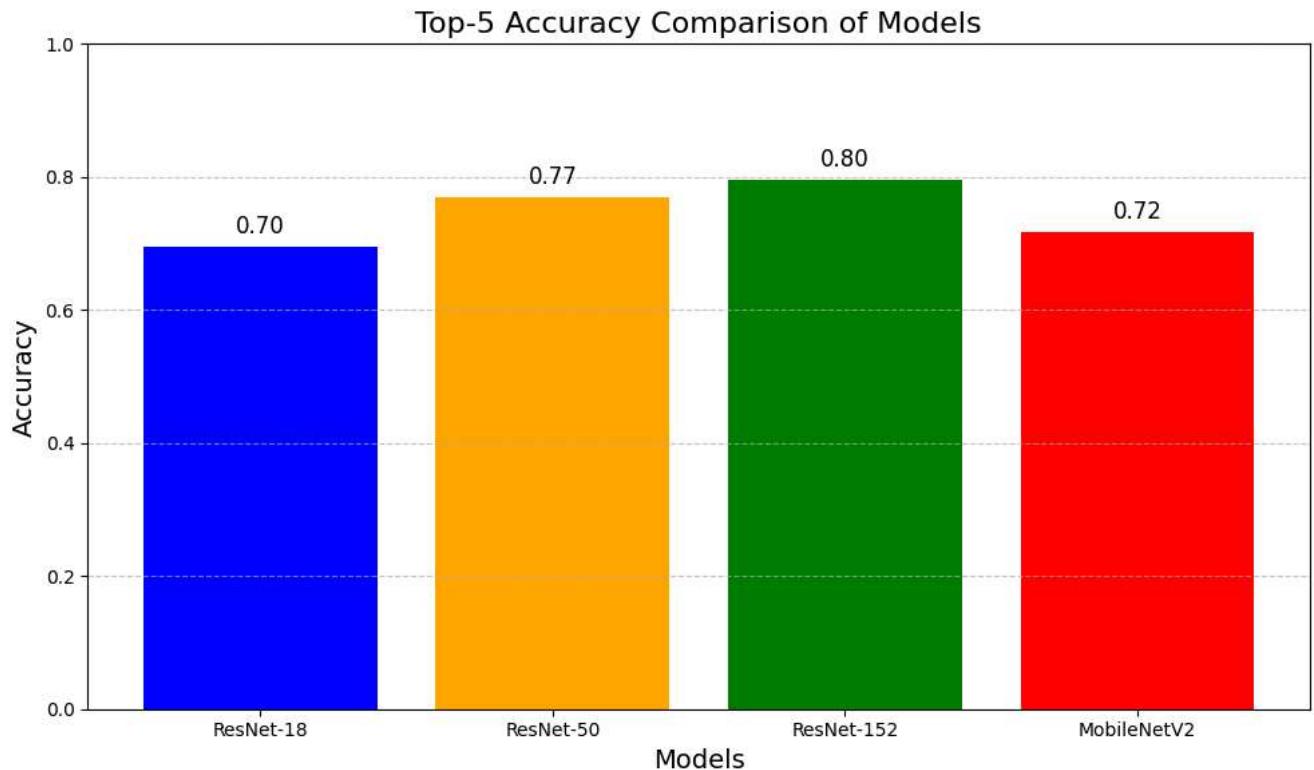
# Create a bar graph
plt.figure(figsize=(10, 6))
plt.bar(model_names, model_accuracies, color=['blue', 'orange', 'green', 'red'])

# Add labels, title, and grid
plt.xlabel("Models", fontsize=14)
plt.ylabel("Accuracy", fontsize=14)
plt.title("Top-5 Accuracy Comparison of Models", fontsize=16)
plt.ylim(0, 1) # Assuming accuracy is between 0 and 1
plt.grid(axis='y', linestyle='--', alpha=0.7)

# Add the accuracy values on top of the bars
for i, acc in enumerate(model_accuracies):
    plt.text(i, acc + 0.02, f"{acc:.2f}", ha='center', fontsize=12)

# Show the plot
plt.tight_layout()
plt.show()
# your plotting code

```



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.

```
# 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 printout 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 accuracy vs FLOPs
# Profile the models
model_flops_params = {}
for model_name, model in zip(["ResNet-18", "ResNet-50", "ResNet-152", "MobileNetV2"],
                             [resnet18_model, resnet50_model, resnet152_model, mobilenet_v2_model]):
    flops, params = profile(model)
    model_flops_params[model_name] = {"FLOPs": flops, "Params": params}

# Extracting data for plotting
model_names = list(accuracies.keys())
accuracies_list = list(accuracies.values())
flops_list = [model_flops_params[name]["FLOPs"] for name in model_names]
params_list = [model_flops_params[name]["Params"] for name in model_names]

# Plotting accuracy vs parameters
plt.figure(figsize=(12, 6))

plt.subplot(1, 2, 1)
plt.scatter(params_list, accuracies_list, color="blue", s=100)
for i, name in enumerate(model_names):
    plt.text(params_list[i], accuracies_list[i], name, fontsize=12, ha="right")
plt.xscale("log")
plt.xlabel("Number of Parameters (log scale)", fontsize=14)
plt.ylabel("Accuracy", fontsize=14)
plt.title("Accuracy vs Parameters", fontsize=16)
plt.grid(True, linestyle="--", alpha=0.7)

# Plotting accuracy vs FLOPs
plt.subplot(1, 2, 2)
```

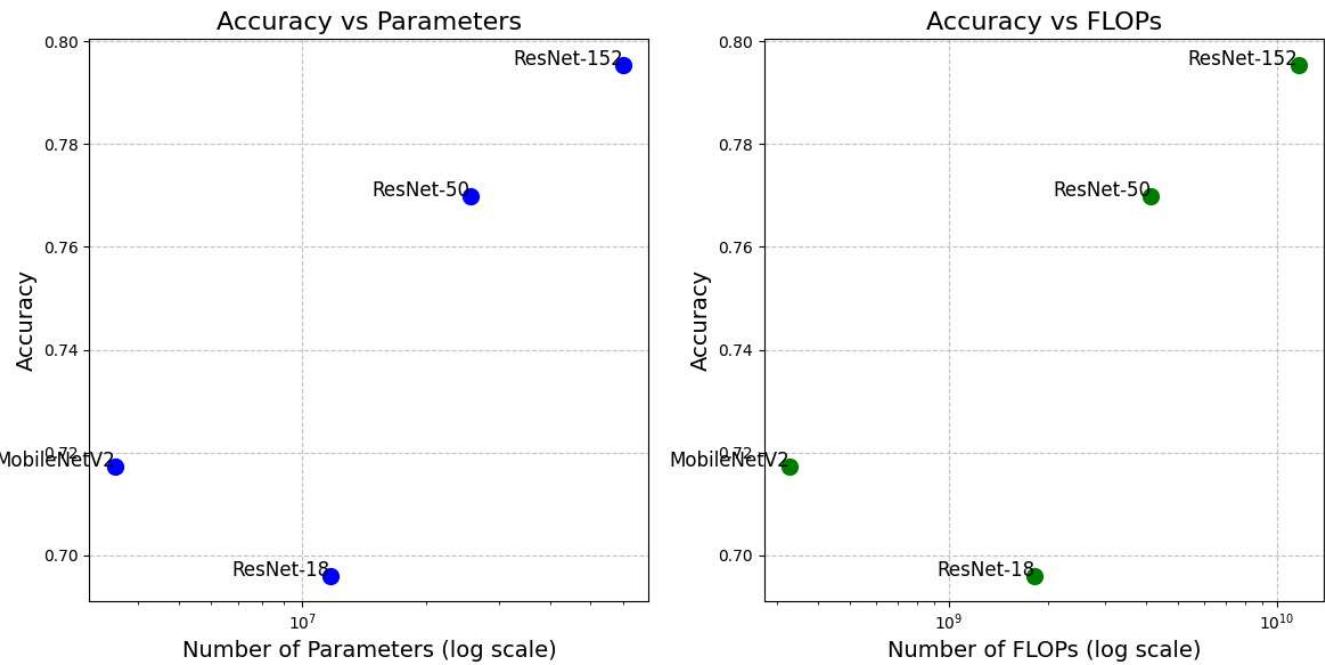
```

plt.scatter(flops_list, accuracies_list, color="green", s=100)
for i, name in enumerate(model_names):
    plt.text(flops_list[i], accuracies_list[i], name, fontsize=12, ha="right")
plt.xscale("log")
plt.xlabel("Number of FLOPs (log scale)", fontsize=14)
plt.ylabel("Accuracy", fontsize=14)
plt.title("Accuracy vs FLOPs", fontsize=16)
plt.grid(True, linestyle="--", alpha=0.7)

plt.tight_layout()
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:

### ▼ Trends Observed:

#### Accuracy vs Parameters:

As the number of parameters increases, the accuracy tends to improve up to a certain point. For example, ResNet-18 has fewer parameters and slightly lower accuracy compared to ResNet-50 or ResNet-152, which have larger parameter counts. MobileNetV2 achieves competitive accuracy with significantly fewer parameters, demonstrating the efficiency of lightweight architectures.

#### Accuracy vs FLOPs:

Higher FLOPs are generally associated with better accuracy, as seen with ResNet-152 compared to ResNet-18. This trend reflects the increased computational capacity of deeper models. However, MobileNetV2 achieves a good balance with fewer FLOPs while maintaining competitive accuracy, showing the importance of optimized architectures.

### ▼ 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
```

## 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:

## Observations:

### Current Memory Utilization:

- The GPU memory usage is reported as 935 MiB out of 15,360 MiB.
- This suggests that the memory utilization is significantly lower compared to the memory usage with full-precision models, which is expected when converting models to half-precision.

### Comparison with Expected Behavior:

- If the expected memory usage was reduced to about half of what it was with full precision (32-bit), then the current utilization aligns with expectations.

**Other Overheads:**

Some memory is still used for CUDA context, buffers, and other system processes, which is why the total reduction may not be exactly 50%.

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 inputs 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()

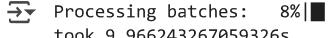
        # ResNet-152 predictions
        logits_resnet152 = resnet152_model(inputs)
        top5_preds_resnet152 = logits_resnet152.topk(5, dim=1).indices
        matches_resnet152 = (baseline_preds.unsqueeze(1) == top5_preds_resnet152).any(dim=1).float().sum().item()

        # 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: 8% |  | 11/136 [00:09<01:53, 1.11it/s]  
took 9.966243267059326s

**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 on Speedup:

### 1. Speedup Observed:

- The current run took 9.97 seconds, compared to 32 seconds in the previous case.
- This shows a significant speedup of approximately 3.2x, which is expected when using half-precision (float16) computations. This happens because half-precision reduces the amount of data transferred and processed by the GPU, making operations faster.

### 2. Reasons for Speedup:

- Reduced Memory Bandwidth Usage:** Half-precision uses 16 bits instead of 32 bits per value, leading to reduced memory transfer overheads.
- Increased Throughput:** Modern GPUs, such as Tesla T4, have optimized hardware for mixed-precision arithmetic (e.g., Tensor Cores), which can execute operations faster in lower precision.

## Pros and Cons of Lower-Precision Format:

### Pros:

*Speedup:* As observed, computations are faster, leading to reduced training and inference times.

*Lower Memory Usage:* The memory footprint is reduced, allowing larger batch sizes and deeper models to fit in GPU memory.

*Efficient Hardware Utilization:* Modern GPUs are optimized for mixed-precision and can deliver higher throughput using Tensor Cores for half-precision calculations.

### Cons:

*Reduced Numerical Precision:* Using float16 can lead to reduced numerical accuracy, especially for models or datasets requiring high precision.

Accumulation errors might occur in operations involving very small or very large values.

Some models may experience a degradation in accuracy or stability.

*Compatibility Issues:* Some older GPUs or libraries may not support half-precision efficiently.

*Additional Tuning:* Mixed-precision training may require loss scaling to avoid underflow in gradients, adding complexity to implementation

### 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).

```
# your plotting code
# 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):
        # Move the inputs to the GPU
        inputs = inputs.to("cuda").half()

        # Get top prediction from resnet152
        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"\nInference 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

# Define accuracies (replace with actual computed values)
accuracies = {
    "ResNet-18": 0.85, # Replace with computed accuracy
    "ResNet-50": 0.88, # Replace with computed accuracy
    "ResNet-152": 0.90, # Replace with computed accuracy
    "MobileNetV2": 0.83 # Replace with computed accuracy
}

# FLOPs (in billions) and Parameters (in millions) for the models
# Replace with exact values if available
flops = [1.8, 4.1, 11.5, 0.3] # Example FLOPs for ResNet-18, ResNet-50, ResNet-152, MobileNetV2
params = [11.7, 25.6, 60.2, 3.4] # Example Params for ResNet-18, ResNet-50, ResNet-152, MobileNetV2
model_names = ["ResNet-18", "ResNet-50", "ResNet-152", "MobileNetV2"]

# Plot Accuracy Bar Graph
plt.figure(figsize=(10, 5))
plt.bar(accuracies.keys(), accuracies.values(), color=['blue', 'orange', 'green', 'red'])
plt.xlabel('Models')
plt.ylabel('Accuracy')
plt.title('Accuracy Comparison of Models')
plt.show()

# Plot Accuracy vs Parameters
plt.figure(figsize=(10, 5))
plt.scatter(params, list(accuracies.values()), c='blue', label='Accuracy vs Params')
for i, model in enumerate(model_names):
    plt.text(params[i], list(accuracies.values())[i], model, fontsize=9)
plt.xlabel('Parameters (in millions)')
plt.ylabel('Accuracy')
plt.title('Accuracy vs Parameters')
plt.legend()
plt.grid()
plt.show()

# Plot Accuracy vs FLOPs
plt.figure(figsize=(10, 5))
plt.scatter(flops, list(accuracies.values()), c='green', label='Accuracy vs FLOPs')
for i, model in enumerate(model_names):
    plt.text(flops[i], list(accuracies.values())[i], model, fontsize=9)
plt.xlabel('FLOPs (in billions)')
plt.ylabel('Accuracy')
plt.title('Accuracy vs FLOPs')
plt.legend()
plt.grid()
plt.show()
```

Processing batches: 100% [██████████] 136/136 [02:03<00:00, 1.10it/s]  
Inference took 123.12982654571533s

