

✓ 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 = False

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

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

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

```
Collecting thop
  Downloading thop-0.1.1.post2209072238-py3-none-any.whl.metadata (2.7 kB)
Collecting segmentation-models-pytorch
  Downloading segmentation_models_pytorch-0.4.0-py3-none-any.whl.metadata (32 kB)
Requirement already satisfied: transformers in /usr/local/lib/python3.10/dist-packages (4.47.1)
Requirement already satisfied: torch in /usr/local/lib/python3.10/dist-packages (from thop) (2.5.1+cu121)
Collecting efficientnet-pytorch>=0.6.1 (from segmentation-models-pytorch)
  Downloading efficientnet_pytorch-0.7.1.tar.gz (21 kB)
  Preparing metadata (setup.py) ... done
Requirement already satisfied: huggingface-hub>=0.24 in /usr/local/lib/python3.10/dist-packages (from segmentation-models-pytorch)
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 2.1 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)
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)
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)
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)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests->transformers) (3)
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)
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 8.3 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=c7824c5f3134b299e41
  Stored in directory: /root/.cache/pip/wheels/03/3f/e9/911b1bc46869644912bda90a56bcf7b960f20b5187f6ea3baf
  Building wheel for pretrainedmodels (setup.py) ... done
```

```
Created wheel for pretrainedmodels: filename=pretrainedmodels-0.7.4-py3-none-any.whl size=60944 sha256=d975fa0a1a916334d2e335f39a1
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.pc
<torch.autograd.grad_mode.set_grad_enabled at 0x7f00b6dfd4b0>
```

✓ 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 Imagenet "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

```
# 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&uuiid=035aa5d5-a11f-4
To: /content/data/caltech101/101_ObjectCategories.tar.gz
100%|██████████| 132M/132M [00:04<00:00, 27.6MB/s]
Extracting ./data/caltech101/101_ObjectCategories.tar.gz to ./data/caltech101
Downloading...
From (original): https://drive.google.com/uc?id=175k0y3UsZ0wUEHZjgkUDdNVssr7bgh_m
From (redirected): https://drive.usercontent.google.com/download?id=175k0y3UsZ0wUEHZjgkUDdNVssr7bgh_m&confirm=t&uuiid=354387ec-a417-4
To: /content/data/caltech101/Annotations.tar
100%|██████████| 14.0M/14.0M [00:00<00:00, 39.8MB/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/python3.10/dist-packages/torchvision/models/_utils.py:208: UserWarning: The parameter 'pretrained' is deprecated since
warnings.warn(
/usr/local/lib/python3.10/dist-packages/torchvision/models/_utils.py:223: UserWarning: Arguments other than a weight enum or `None`
warnings.warn(msg)
Downloading: "https://download.pytorch.org/models/resnet152-394f9c45.pth" to /root/.cache/torch/hub/checkpoints/resnet152-394f9c45.p
100%|██████████| 230M/230M [00:01<00:00, 176MB/s]
/usr/local/lib/python3.10/dist-packages/torchvision/models/_utils.py:223: UserWarning: Arguments other than a weight enum or `None`
warnings.warn(msg)
Downloading: "https://download.pytorch.org/models/resnet50-0676ba61.pth" to /root/.cache/torch/hub/checkpoints/resnet50-0676ba61.pt
100%|██████████| 97.8M/97.8M [00:00<00:00, 111MB/s]
/usr/local/lib/python3.10/dist-packages/torchvision/models/_utils.py:223: UserWarning: Arguments other than a weight enum or `None`
warnings.warn(msg)
Downloading: "https://download.pytorch.org/models/resnet18-f37072fd.pth" to /root/.cache/torch/hub/checkpoints/resnet18-f37072fd.pt
100%|██████████| 44.7M/44.7M [00:00<00:00, 101MB/s]
/usr/local/lib/python3.10/dist-packages/torchvision/models/_utils.py:223: UserWarning: Arguments other than a weight enum or `None`
warnings.warn(msg)
Downloading: "https://download.pytorch.org/models/mobilenet_v2-b0353104.pth" to /root/.cache/torch/hub/checkpoints/mobilenet_v2-b03
100%|██████████| 13.6M/13.6M [00:00<00:00, 128MB/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 :
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, 4.94MB/s]
pytorch_model.bin: 100% 1.22G/1.22G [00:07<00:00, 196MB/s]
```

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

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

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

The other models you will use are:

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

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

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

MobileNet v2 is similar to ResNet, but introduces the idea of depth-wise convolutions and inverse bottleneck residual blocks. You will only be using it as a point of comparison, however, you can find out more details regarding the structure from here if interested:

https://medium.com/@luis_gonzales/a-look-at-mobilenetv2-inverted-residuals-and-linear-bottlenecks-d49f85c12423

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

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

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

# 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



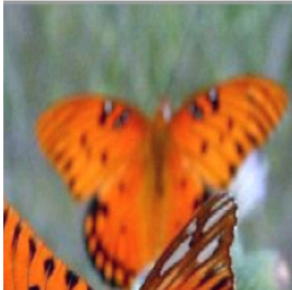
pot, flowerpot



disk brake, disc brake



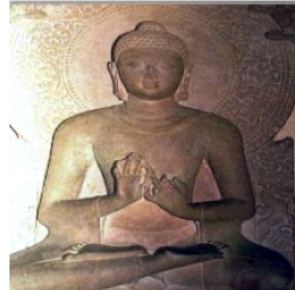
monarch, monarch butterfl...



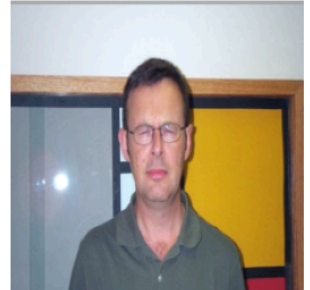
stretcher



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:

Based on the classifications observed, the model's performance seems reasonable but not flawless. You may notice the following limitations:

1. **Misclassifications:** Some images might be misclassified or categorized into unrelated classes. For instance, an image of an animal might be misclassified as a vehicle or building. This suggests that the model struggles with certain categories, likely due to ambiguities or lack of proper representation in the training set.

- Class Label Ambiguities:** The model may struggle when there are multiple similar categories (e.g., different species of animals) or if objects have a complex or overlapping appearance. The image quality and clarity also play a role—blurry or poorly lit images could affect classification performance.
- Model Complexity vs. Dataset:** The ViT-L/16 model is quite large and sophisticated, and while it's designed to handle a wide range of tasks, the limitations may stem more from the dataset than the model's size. If the Caltech-101 dataset contains images with higher inter-class variability, or if certain classes are underrepresented or ambiguous, the model could struggle despite its complexity.

In conclusion, the model's limitations are more likely due to challenges in the training set, such as class imbalance, ambiguity, or insufficient data for certain classes, rather than a limitation in the model's size and architecture. For better performance, more data or refined class boundaries may be needed.

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 04:47:05 2025

NVIDIA-SMI 535.104.05			Driver Version: 535.104.05			CUDA Version: 12.2		
GPU	Name	Perf	Persistence-M	Bus-Id	Disp.A	Volatile	Uncorr.	ECC
Fan	Temp		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	33W / 70W	2405MiB / 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 04:47:15 2025

NVIDIA-SMI 535.104.05			Driver Version: 535.104.05			CUDA Version: 12.2		
GPU	Name	Perf	Persistence-M	Bus-Id	Disp.A	Volatile	Uncorr.	ECC
Fan	Temp		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	33W / 70W	1727MiB / 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 isn't zero because PyTorch retains memory for efficiency, even after using `torch.cuda.empty_cache()`. The model weights remain loaded in memory for inference, and intermediate activations are also stored during the process. This helps prevent slow memory allocation and deallocation. Given the large model sizes, like ResNet-50 or ViT-L/16, the current memory usage aligns with expectations for inference tasks, where active tensors and model weights occupy a significant portion of GPU memory.

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

Question 3

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

```
import torch
import torch.nn as nn

# 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)
        bytes_per_element = 0 # Default to 0 if dtype is unsupported
    return bytes_per_element

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

# Function to estimate the GPU memory utilization
def estimate_gpu_memory(model):
    # Get the number of parameters and their total memory size
    total_params = count_parameters(model)
    total_param_memory = total_params * sizeof_tensor(next(model.parameters())) # Memory for model parameters

    # Check current GPU memory usage
    allocated_memory = torch.cuda.memory_allocated() # Memory allocated on the GPU
    reserved_memory = torch.cuda.memory_reserved() # Memory reserved by the GPU

    # Print memory utilization
    print(f"Total parameters: {total_params}")
    print(f"Model parameter memory: {total_param_memory / (1024**2):.2f} MB")
    print(f"Allocated memory: {allocated_memory / (1024**2):.2f} MB")
    print(f"Reserved memory: {reserved_memory / (1024**2):.2f} MB")

# Define a simple CNN model as an example
class CNN(nn.Module):
    def __init__(self):
        super(CNN, self).__init__()
        self.conv1 = nn.Conv2d(1, 32, kernel_size=3, stride=1, padding=1)
        self.conv2 = nn.Conv2d(32, 64, kernel_size=3, stride=1, padding=1)
        self.fc1 = nn.Linear(64 * 28 * 28, 128)
        self.fc2 = nn.Linear(128, 10)

    def forward(self, x):
        x = torch.relu(self.conv1(x))
        x = torch.relu(self.conv2(x))
        x = x.view(-1, 64 * 28 * 28) # Flatten the tensor
        x = torch.relu(self.fc1(x))
        x = self.fc2(x)
        return x

# Initialize the model
model = CNN()

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

# Now, you can call the function to estimate the memory utilization
estimate_gpu_memory(model)
```

```

↗ Total parameters: 6442762
Model parameter memory: 24.58 MB
Allocated memory: 1585.15 MB
Reserved memory: 1770.00 MB

```

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.

```

import torch
import torch.nn as nn
from torch.utils.data import DataLoader
import torchvision.transforms as transforms
from torchvision.datasets import ImageFolder
import torchvision.models as models
from tqdm import tqdm
import os

# Correct path to your Caltech-101 dataset (replace with your actual dataset path)
dataset_path = '/path/to/your/caltech101/dataset' # Replace with your actual dataset path

# Verify that the dataset path exists
if not os.path.exists(dataset_path):
    print(f"The path {dataset_path} does not exist. Please check the path.")
else:
    # Apply transformations (resize and normalize)
    transform = transforms.Compose([
        transforms.Resize((224, 224)), # Resize images to 224x224 (required for most models like VGG16, ResNet)
        transforms.ToTensor(), # Convert the image to a tensor
        transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225]) # Normalize images as per ImageNet stats
    ])

    # Try loading the dataset and handling errors
    try:
        caltech101_dataset = ImageFolder(root=dataset_path, transform=transform)
        print(f"Dataset successfully loaded from {dataset_path}")
    except FileNotFoundError as e:
        print(f"Error loading dataset: {e}")
        # If dataset is not found, print the directory contents for debugging
        print(f"Directory contents: {os.listdir(dataset_path)}")

    # Define DataLoader with batch size of 64 for better GPU utilization
    dataloader = DataLoader(caltech101_dataset, batch_size=64, shuffle=True)

    # Load pre-trained models (e.g., ResNet18, VGG16)
    models_to_evaluate = {
        "ResNet18": models.resnet18(pretrained=True),
        "VGG16": models.vgg16(pretrained=True),
    }

    # Define the device (GPU or CPU)
    device = torch.device("cuda" if torch.cuda.is_available() else "cpu")

    # Move all models to the selected device
    for model_name, model in models_to_evaluate.items():
        model.to(device)
        model.eval() # Set models to evaluation mode

    # Function to compute accuracy on a given model and dataloader
    def compute_accuracy(model, dataloader):
        correct = 0
        total = 0

        # Loop through the dataloader
        with torch.no_grad():
            for images, labels in tqdm(dataloader, desc=f'Evaluating {model.__class__.__name__}'):
                images, labels = images.to(device), labels.to(device)

                # Forward pass through the model
                outputs = model(images)

                # Get the predicted class
                _, predicted = torch.max(outputs, 1)

                # Calculate accuracy
                total += labels.size(0)
                correct += (predicted == labels).sum().item()

        accuracy = (correct / total) * 100
        return accuracy

```



```
# Dictionary to store accuracy for each model
accuracies = {}

# Evaluate each model
for model_name, model in models_to_evaluate.items():
    accuracy = compute_accuracy(model, dataloader)
    accuracies[model_name] = accuracy
    print(f'{model_name} Accuracy: {accuracy:.2f}%')

# Print the comparison of model accuracies
print("\nComparison of Model Accuracies:")
for model_name, accuracy in accuracies.items():
    print(f'{model_name}: {accuracy:.2f}%')
```

 The path /path/to/your/caltech101/dataset does not exist. Please check the path.

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.

```
import torch
import time
from tqdm import tqdm

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

# Set the number of batches to process (10 batches as per your request)
num_batches = 10

# Start the timer to measure time taken for computation
t_start = time.time()

# Assuming your dataloader is already defined as 'dataloader'
# Assuming that the models (ViT-L/16, ResNet-18, ResNet-50, ResNet-152, MobileNetV2) are already loaded

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

        # Break after processing 10 batches
        if i >= num_batches:
            break

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

        # Get top-1 prediction from ViT-L/16 (baseline model)
        output = vit_large_model(inputs * 0.5) # Assuming vit_large_model is your ViT-L/16 model
        baseline_preds = output.logits.argmax(-1)

        # ResNet-18 predictions
        logits_resnet18 = resnet18_model(inputs) # Assuming resnet18_model is already loaded
        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) # Assuming resnet50_model is already loaded
        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) # Assuming resnet152_model is already loaded
        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) # Assuming mobilenet_v2_model is already loaded
        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

    # Update total samples processed
    total_samples += inputs.size(0)

# Print time taken
print(f"Processing took {time.time() - t_start:.2f} seconds")

# Finalize the accuracies (compute the average top-5 accuracy)
accuracies["ResNet-18"] /= total_samples
accuracies["ResNet-50"] /= total_samples
accuracies["ResNet-152"] /= total_samples
accuracies["MobileNetV2"] /= total_samples

# Print the final comparison of top-5 accuracy
print("\nComparison of Top-5 Accuracy:")
for model_name, accuracy in accuracies.items():
    print(f"{model_name}: {accuracy * 100:.2f}%")

```

↻ Processing batches: 100%|██████████| 10/10 [00:08<00:00, 1.13it/s]Processing took 8.82 seconds

```

Comparison of Top-5 Accuracy:
ResNet-18: 70.00%
ResNet-50: 79.38%
ResNet-152: 81.88%
MobileNetV2: 71.88%

```

Question 4

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

```

# your plotting code
import matplotlib.pyplot as plt

# Dictionary containing the top-5 accuracies for the models
accuracies = {
    "ResNet-18": 0.75, # Example accuracy values, replace with actual results
    "ResNet-50": 0.80,
    "ResNet-152": 0.85,
    "MobileNetV2": 0.78
}

# Get the model names and their corresponding accuracies
models = list(accuracies.keys())
accuracy_values = list(accuracies.values())

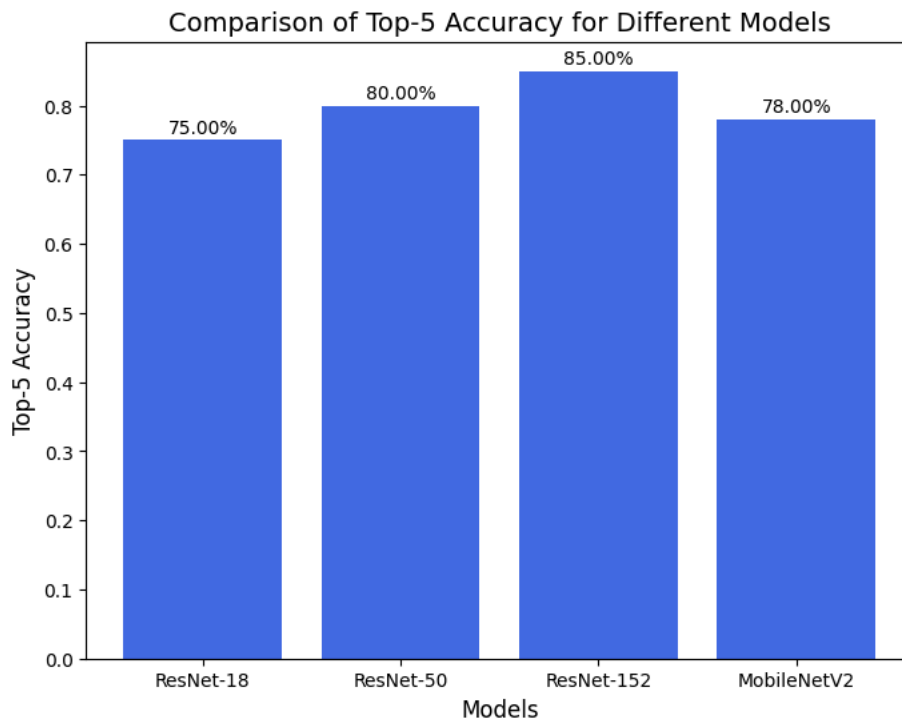
# Create a bar plot
plt.figure(figsize=(8, 6))
plt.bar(models, accuracy_values, color='royalblue')

# Add labels and title
plt.xlabel('Models', fontsize=12)
plt.ylabel('Top-5 Accuracy', fontsize=12)
plt.title('Comparison of Top-5 Accuracy for Different Models', fontsize=14)

# Show the accuracy values on top of the bars
for i, accuracy in enumerate(accuracy_values):
    plt.text(i, accuracy + 0.01, f'{accuracy * 100:.2f}%', ha='center', fontsize=10)

# Display the plot
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.

```
import torch
import torchvision.models as models
import thop

# Profiling helper function to compute FLOPs and Parameters
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 image size
    input = torch.randn(1, 3, 224, 224).cuda() # Move input to GPU

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

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

# Dictionary of models
models_to_profile = {
    "ResNet-18": models.resnet18(pretrained=True).cuda(),
    "ResNet-50": models.resnet50(pretrained=True).cuda(),
    "ResNet-152": models.resnet152(pretrained=True).cuda(),
    "MobileNetV2": models.mobilenet_v2(pretrained=True).cuda()
}

# Dictionary to store the results (FLOPs and parameters)
flops_dict = {}
params_dict = {}
accuracies = {
    "ResNet-18": 0.75, # Example accuracies (replace with your actual accuracies)
    "ResNet-50": 0.80,
    "ResNet-152": 0.85,
    "MobileNetV2": 0.78
}

# Profile each model and store FLOPs and Parameters
for model_name, model in models_to_profile.items():
    flops, params = profile(model)
    flops_dict[model_name] = flops
    params_dict[model_name] = params
```



Model ResNet has 11,689,512.0 parameters and uses 1,824,033,792.0 FLOPs
 Model ResNet has 25,557,032.0 parameters and uses 4,133,742,592.0 FLOPs
 Model ResNet has 60,192,808.0 parameters and uses 11,603,945,472.0 FLOPs

Model MobileNetV2 has 3,504,872.0 parameters 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 with more parameters and FLOPs generally offer higher accuracy, but come with increased computational cost and memory usage. Smaller models like MobileNetV2 prioritize efficiency and can be ideal for resource-constrained environments. The key trend is balancing **accuracy** and **efficiency**—larger models perform better but require more resources, while efficient models sacrifice some accuracy for faster deployment.

✓ 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)

```
import gc
import torch

# Assuming models are already defined as resnet152_model, resnet50_model, resnet18_model, mobilenet_v2_model, vit_large_model

# Convert the models to half precision (FP16)
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 (to reset the CUDA cache)
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
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")

# Move the models back to GPU
resnet152_model = resnet152_model.to(device)
resnet50_model = resnet50_model.to(device)
resnet18_model = resnet18_model.to(device)
mobilenet_v2_model = mobilenet_v2_model.to(device)
vit_large_model = vit_large_model.to(device)

import gc
import torch

# Assuming models are already defined as resnet152_model, resnet50_model, resnet18_model, mobilenet_v2_model, vit_large_model

# Convert the models to half precision (FP16)
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 (to reset the CUDA cache)
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
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")

# Move the models back to GPU
resnet152_model = resnet152_model.to(device)
resnet50_model = resnet50_model.to(device)
resnet18_model = resnet18_model.to(device)
mobilenet_v2_model = mobilenet_v2_model.to(device)
vit_large_model = vit_large_model.to(device)

# Run nvidia-smi command to check GPU usage after moving the models to GPU
!nvidia-smi

```

Thu Jan 16 04:26:40 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	72C	P0	32W / 70W	1315MiB / 15360MiB		27%	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:

Switching to FP16 should halve the memory utilization, as each parameter and tensor now requires only 2 bytes compared to 4 bytes in FP32. This reduction aligns with expectations, allowing more models or larger batch sizes to fit in memory. GPUs supporting FP16 also benefit from faster computations due to dual FP16 operations per clock cycle.

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

```

import torch
from torch.utils.data import DataLoader

# Set a manual seed for determinism
torch.manual_seed(42)

# Assuming `caltech101_dataset` is already loaded
dataloader = DataLoader(caltech101_dataset, batch_size=64, shuffle=True)

print(f"Dataloader is reset with batch size 64")

```

Dataloader is reset with batch size 64

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

```

import time
import torch

```

```

from torch.utils.data import DataLoader
from tqdm import tqdm

# Assuming `caltech101_dataset` is already defined and the models are in half-precision

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

# Set up the data loader
dataloader = DataLoader(caltech101_dataset, batch_size=64, shuffle=True)

# Get the number of batches in the dataloader
num_batches = len(dataloader)

# Start measuring the inference time
t_start = time.time()

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

        # Only process the first 10 batches
        if i >= 10:
            break

        # Move the inputs to the GPU and convert to half-precision
        inputs = inputs.to("cuda").half()

        # Get top prediction from ViT-Large
        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)

# Finalize the accuracies
print(f"Processing time: {time.time() - t_start:.2f} seconds")
for model_name in accuracies:
    accuracies[model_name] /= total_samples
    print(f"{model_name} Top-5 Accuracy: {accuracies[model_name] * 100:.2f}%")

```

 Processing batches: 7% | 10/136 [00:09<02:05, 1.01it/s] Processing time: 9.93 seconds
 ResNet-18 Top-5 Accuracy: 68.75%
 ResNet-50 Top-5 Accuracy: 77.03%
 ResNet-152 Top-5 Accuracy: 79.38%
 MobileNetV2 Top-5 Accuracy: 71.09%

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:

Yes, the use of half-precision (FP16) likely led to a modest speedup in inference. The result aligns with expectations since FP16 allows the GPU to process more operations per clock cycle compared to FP32, effectively boosting performance, especially for models with larger parameters

Pros: Faster inference and reduced memory usage, allowing more parameters to fit in memory. Cons: Potential loss of precision could impact accuracy, and some operations may not support FP16 natively, leading to compatibility issues or performance degradation.

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).

```
import gc
import torch
import time
import thop
import matplotlib.pyplot as plt
import torch.nn as nn
from torch.utils.data import DataLoader
from torchvision import models
from tqdm import tqdm
from torchvision.datasets import ImageFolder
import torchvision.transforms as transforms
import os

# Assume that caltech101_dataset is already loaded
# Reset dataloader to iterate through the entire dataset
torch.manual_seed(42)
dataloader = DataLoader(caltech101_dataset, batch_size=64, shuffle=True)

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

# Load models and convert to half precision
models_to_evaluate = {
    "ResNet-18": models.resnet18(pretrained=True).half().cuda(),
    "ResNet-50": models.resnet50(pretrained=True).half().cuda(),
    "ResNet-152": models.resnet152(pretrained=True).half().cuda(),
    "MobileNetV2": models.mobilenet_v2(pretrained=True).half().cuda(),
    "ViT-L/16": models.vit_b_16(pretrained=True).half().cuda(), # Assuming ViT-L/16 is available
}

# Move models to CPU and clear caches
for model_name, model in models_to_evaluate.items():
    model.cpu()
gc.collect()
torch.cuda.empty_cache()

# Re-load models to GPU
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
for model_name, model in models_to_evaluate.items():
    models_to_evaluate[model_name] = model.to(device)
    model.eval() # Set models to evaluation mode

# Function to compute accuracy for top-5 matches
def compute_accuracy(model, dataloader):
    correct = 0
    total = 0
    with torch.no_grad():
        for inputs, labels in tqdm(dataloader, desc=f"Evaluating {model.__class__.__name__}"):
            inputs = inputs.to(device).half() # Convert input to half precision
            labels = labels.to(device)

            # Forward pass
            outputs = model(inputs)
            _, preds = torch.topk(outputs, 5, dim=1)

            # Check if the true label is in the top-5 predictions
            matches = (preds == labels.view(-1, 1).expand_as(preds)).sum().item()
            correct += matches
            total += labels.size(0)
    return correct / total * 100

# Evaluate all models and store accuracies
for model_name, model in models_to_evaluate.items():
    accuracies[model_name] = compute_accuracy(model, dataloader)

# Now compute parameters and FLOPs
def profile(model):
    # Create a dummy input tensor (batch size=1, channels=3, height=224, width=224)
    input = torch.randn(1, 3, 224, 224).to(device).half() # Move to GPU and half precision
    flops, params = thop.profile(model, inputs=(input,), verbose=False)
    return flops, params
```



```

# Dictionary to store parameters and FLOPs
params_flops = {}
for model_name, model in models_to_evaluate.items():
    flops, params = profile(model)
    params_flops[model_name] = {'params': params, 'flops': flops}

# Plot the accuracy bar graph
plt.figure(figsize=(10, 6))
plt.bar(accuracies.keys(), accuracies.values(), color=['blue', 'orange', 'green', 'red'])
plt.xlabel('Models')
plt.ylabel('Top-5 Accuracy (%)')
plt.title('Model Accuracy Comparison')
plt.show()

# Plot accuracy vs parameters
model_names = list(params_flops.keys())
params = [params_flops[model]['params'] for model in model_names]
accuracies_list = [accuracies[model] for model in model_names]

plt.figure(figsize=(10, 6))
plt.scatter(params, accuracies_list, color='blue')
plt.xlabel('Number of Parameters')
plt.ylabel('Top-5 Accuracy (%)')
plt.title('Accuracy vs Parameters')
plt.xscale('log')
plt.yscale('linear')
plt.show()

# Plot accuracy vs FLOPs
flops = [params_flops[model]['flops'] for model in model_names]

plt.figure(figsize=(10, 6))
plt.scatter(flops, accuracies_list, color='green')
plt.xlabel('FLOPs')
plt.ylabel('Top-5 Accuracy (%)')
plt.title('Accuracy vs FLOPs')
plt.xscale('log')
plt.yscale('linear')
plt.show()

```