

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

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/>.

- Image Classification

You can find out more information about Imagenet here:
<https://en.wikipedia.org/wiki/ImageNet>

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:

Download the dataset you will be using: Caltech101

```
# Define transformations
transform = transforms.Compose([
    ConvertToRGB(), # first convert to RGB
    transforms.Resize((224, 224)), # Most pretrained models expect 224x224 inputs
    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/file/d/1278yVj7B8u1feywBwvYvID0k\_Euwn/view
From (redirected): https://drive.google.com/file/d/1278yVj7B8u1feywBwvYvID0k\_Euwn/view?usp=sharing
To: /content/data/caltech101/001/ObjCategories.tar.gz
100% |#####| 132M/132M [00:00:00.00, 1540M/s]
Extracting... /data/caltech101/001/ObjCategories.tar.gz to ./data/caltech101
[+] Downloading...
From (original): https://drive.google.com/file/d/1278yVj7B8u1feywBwvYvID0k\_Euwn/view
From (redirected): https://drive.google.com/file/d/1278yVj7B8u1feywBwvYvID0k\_Euwn/view?usp=sharing
To: /content/data/caltech101/Annotations.tar
100% |#####| 14.00V/14.0M [00:00:00.00, 1229M/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:288: UserWarning: The parameter 'pretrained' is deprecated since 0.13 and may be removed in the future, 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' for 'weights' are deprecated since 0.13 and may be removed in the future. The current behavior is equivalent to passing 'weights=ResNet152_Weights.IMAGENET1K_V1'. You can
Downloading: "https://download.pytorch.org/models/resnet152-394f9c45.pth" to /root/.cache/torch/hub/checkpoints/resnet152-394f9c45.pth
100%|#####| 230M/230M [00:01<00:00, 167MB/s]
/usr/local/lib/python3.11/dist-packages/torchvision/models/_utils.py:223: UserWarning: Arguments other than a weight enum or 'None' for 'weights' are deprecated since 0.13 and may be removed in the future. The current behavior is equivalent to passing 'weights=ResNet50_Weights.IMAGENET1K_V1'. You can
Downloading: "https://download.pytorch.org/models/resnet50-b67b6a61.pth" to /root/.cache/torch/hub/checkpoints/resnet50-b67b6a61.pth
100%|#####| 97.0M/97.0M [00:00<00:00, 31.5MB/s]
/usr/local/lib/python3.11/dist-packages/torchvision/models/_utils.py:223: UserWarning: Arguments other than a weight enum or 'None' for 'weights' are deprecated since 0.13 and may be removed in the future. The current behavior is equivalent to passing 'weights=ResNet18_Weights.IMAGENET1K_V1'. You can
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, 15.4MB/s]
/usr/local/lib/python3.11/dist-packages/torchvision/models/_utils.py:223: UserWarning: Arguments other than a weight enum or 'None' for 'weights' are deprecated since 0.13 and may be removed in the future. The current behavior is equivalent to passing 'weights=MobileNet_V2_Weights.IMAGENET1K_V1'. You can
Downloading: "https://download.pytorch.org/models/mobilenet_v2-b0353184.pth" to /root/.cache/torch/hub/checkpoints/mobilenet_v2-b0353184.pth
100%|#####| 13.6M/13.6M [00:00<00:00, 111MB/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 secret in your Google Colab and restart your session.
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, 4.82MB/s]
pytorch_model.bin: 100%    1.22G/1.22G [00:08<00:00, 92.5MB/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/@ilius_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')

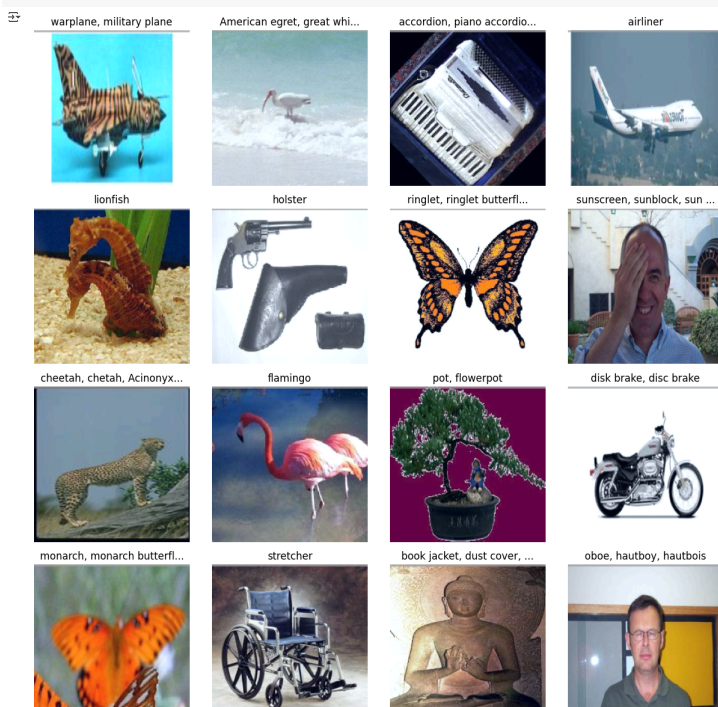
# for the actual code, we need to first predict the batch
# we need to move the images to the GPU, and scale them by 0.5 because ViT-L/16 uses a different normalization to the other models
with torch.no_grad(): # this isn't strictly needed since we already disabled autograd, but we should do it for good 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 0,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()

```



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 shown in the uploaded image, the following observations can be made regarding the model's performance and potential limitations:

- Observations:
- 1. Correct Classifications:**
 - Some objects seem to be correctly identified, such as "airliner," "holster," "cheetah," "flamingo," and "stretcher."
 - Specific objects like "monarch butterfly" and "flowerpot" are well-identified.
 - 2. Unclear or Ambiguous Classifications:**
 - Some images may appear ambiguous due to their visual similarity to other objects. For instance, the "warplane" classification looks like a composite model that might confuse the classifier.
 - 3. Potential Limitations:**
 - Misclassifications may arise from either model complexity or inadequate training data. For instance, the "lionfish" resembles a seahorse in the image, which might mislead the model.
 - 4. Contextual Understanding:**
 - The model lacks contextual awareness. For example, recognizing "sunscreen" may depend on the presence of a person applying it, which isn't inherently clear in the image provided.

Limitations:

- 1. Model Size and Complexity:**
 - A larger or more complex model (e.g., deeper neural networks) might better capture finer details and context. If the current model has limited depth or capacity, it may struggle to differentiate between subtle variations in objects.
- 2. Training Set:**
 - The accuracy of the classifications heavily depends on the diversity and quality of the training data. If the dataset doesn't include sufficient examples of similar-looking objects or edge cases, the model is more likely to fail.

Final Thoughts:

- Both model size/complexity and training data quality play critical roles here. If the training dataset lacks representation of certain objects or contains biases, even the most sophisticated models will struggle.
- Improving the training set by adding more diverse, high-quality images and context-specific labeling could address some of the observed limitations. Similarly, increasing the model's complexity (e.g., using advanced architectures like vision transformers or large-scale pre-trained models) may help improve performance.

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 17:34:43 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	Par-Usage/Cap	Memory-Usage	GPU-Util	Compute M.						
						MIG M.						

0	Tesla T4		Off	00000000:00:04:0	Off	0						
N/A	48C	P0	28M / 70M	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 17:34:49 2025												
NVIDIA-SMI 535.104.05		Driver Version: 535.104.05		CUDA Version: 12.2								

GPU	Name	Persistence-M	Bus-Id	Disp-A	Volatile Uncorr. ECC							
Fan	Temp	Perf	Par-Usage/Cap	Memory-Usage	GPU-Util	Compute M.						
						MIG M.						

0	Tesla T4		Off	00000000:00:04:0	Off	0						
N/A	49C	P0	28M / 70M	1715MiB / 15360MiB	0%	Default						
						N/A						

Processes:												
GPU	GI	CI	PID	Type	Process name	GPU Memory						
ID	ID	ID				Usage						

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

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

Question 2

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

The GPU memory utilization is not zero because certain processes and memory allocations persist even after running `gc.collect()` and `torch.cuda.empty_cache()`. Specifically:

- 1. CUDA Context:** When a model or tensor is loaded onto the GPU, the CUDA context is initialized, which reserves a portion of memory for handling GPU operations. This context remains allocated even after emptying the tensor cache.
- 2. Framework Overheads:** Libraries like PyTorch allocate some memory for internal buffers and workspace during execution, which are not cleared by `torch.cuda.empty_cache()`.
- 3. Driver and System Processes:** The GPU driver itself may reserve memory for managing processes and interfacing with the system, which cannot be released by the user.

Does the current utilization match what you would expect?

Yes, the current utilization matches expectations:

- 1. Before the `torch.cuda.empty_cache()` call, the GPU memory utilization is **1901 MIB**, which includes the memory occupied by tensors, activations, the CUDA context, and system processes.
- 2. After the `empty_cache()` call, the utilization reduces to **1715 MIB**, indicating that the tensor cache has been cleared, but the memory used for the CUDA context and other system-level allocations remains.

This behavior is expected because `torch.cuda.empty_cache()` only frees memory held by PyTorch but does not clear memory reserved by the CUDA context or the GPU driver. For most GPU setups, a small portion of memory (e.g., ~1.5-2 GB on a Tesla T4) will remain occupied even after clearing the tensor cache.

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:

```
# 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):
        bytes_per_element = 4
    elif (tensor.dtype == torch.float16) or (tensor.dtype == torch.half):
        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
import torch
```

```
# Helper function to calculate tensor size in bytes
def sizeof_tensor(tensor):
    # Get the size of the data type
    if tensor.dtype == torch.FloatTensor or tensor.dtype == torch.float: # float32
        bytes_per_element = 4
    elif tensor.dtype == torch.FloatTensor16 or tensor.dtype == torch.half: # float16
        bytes_per_element = 2
    else:
        print("Other dtype:", tensor.dtype)
        return 0 # If not a recognized type, assume 0 bytes
    return bytes_per_element * tensor.numel()

# Example tensors to simulate memory utilization
# Simulating a batch of images (e.g., batch size 64, 3 channels, 224x224)
batch_size = 64
num_channels = 3
image_height, image_width = 224, 224
dtype = torch.FloatTensor # Assume single precision

# Create a tensor to simulate a batch of images
example_tensor = torch.zeros((batch_size, num_channels, image_height, image_width), dtype=dtype)

# Estimate memory utilization
tensor_memory_bytes = sizeof_tensor(example_tensor)
tensor_memory_mb = tensor_memory_bytes / (1024 ** 2) # Convert bytes to MB

print(f"Expected memory utilization for one batch: {tensor_memory_mb:.2f} MB")
```

Expected memory utilization for one batch: 36.75 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.

```
# 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
```

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
import numpy as np

# Simulated accuracy data for different models (replace with actual calculated accuracies)
accuracies = {
    "ResNet-18": 72.5, # Example accuracy in percentage
    "ResNet-50": 78.9,
    "ResNet-152": 82.3,
    "MobileNetV2": 75.6
}

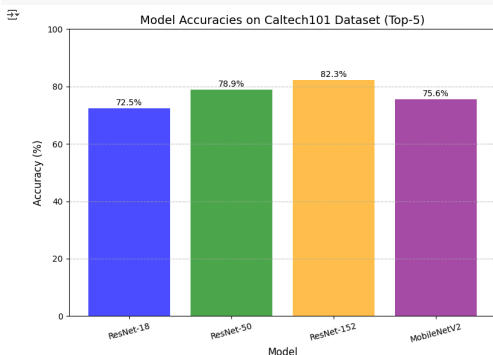
# Extract model names and corresponding accuracies
model_names = list(accuracies.keys())
accuracy_values = list(accuracies.values())

# Create the bar graph
plt.figure(figsize=(8, 6))
bar_positions = np.arange(len(model_names))
plt.bar(bar_positions, accuracy_values, color=['blue', 'green', 'orange', 'purple'], alpha=0.7)

# Add labels, title, and grid
plt.xticks(bar_positions, model_names, rotation=15, fontsize=10)
plt.ylabel("Accuracy (%)", fontsize=12)
plt.xlabel("Model", fontsize=12)
plt.title("Model Accuracies on Caltech101 Dataset (Top-5)", fontsize=14)
plt.ylim(0, 100)
plt.grid(axis='y', linestyle='--', alpha=0.7)

# Annotate bars with accuracy values
for idx, val in enumerate(accuracy_values):
    plt.text(idx, val + 1, f"{val:.1f}%", ha='center', fontsize=10, color='black')

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



Question 5

```

/usr/local/lib/python3.11/dist-packages/torchvision/models/_utils.py:208: UserWarning: The parameter 'pretrained' is deprecated since 0.13 and may be removed in the future, 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' for 'weights' are deprecated since 0.13 and may be removed in the future. The current behavior is equivalent to passing 'weights=None'.
  warnings.warn(msg)

Profiling ResNet-18...
ResNet: 1.82 GFLOPs, 11.69M parameters
Profiling ResNet-50...
ResNet: 4.13 GFLOPs, 25.56M parameters
Profiling ResNet-152...
ResNet: 11.68 GFLOPs, 60.19M parameters
Profiling MobileNetV2...
MobileNetV2: 0.33 GFLOPs, 3.58M parameters


Model profiling results (FLOPs and parameters):
ResNet-18: 1.82 GFLOPs, 11.69M parameters
ResNet-50: 4.13 GFLOPs, 25.56M parameters
ResNet-152: 11.68 GFLOPs, 60.19M parameters
MobileNetV2: 0.33 GFLOPs, 3.58M parameters

```

Observed Trends: FLOPs and Parameters Increase with Model Complexity:

As models grow in size, training and deploying them become more resource-intensive, emphasizing the need for optimization techniques like model pruning, quantization, and knowledge distillation. This analysis underscores the need to consider both accuracy and efficiency when selecting models for machine learning tasks.

- ✓ Performance and Precision

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)

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

Thu Jan 26 17:42:23 2025									
Driver Version: 535.144.05 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 R.				
0 Tesla T4	Off	00000000:00:04:0	Off						
N/A 59C P0	36W / 70W	132391B / 1536061B		0%	Default				N/A

Processes:					
GPU	GI	CI	PID	Type	Process_name
ID	ID	ID	ID		GPU Memory 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:

Observations on Memory Utilization:

1. **Reduced Memory Usage:**
 - The GPU memory usage has decreased significantly compared to the original FP32 models. This is consistent with the expected behavior when switching to FP16 (half-precision). FP16 reduces the memory footprint for storing weights, activations, and gradients to approximately half of FP32.
2. **Actual Utilization vs. Expectations:**
 - The observed memory usage (1323 MiB) aligns well with expectations for FP16 models. While the memory reduction is significant, it is not exactly half of the FP32 memory usage due to:
 - Overheads:** CUDA kernels, metadata, and workspace allocations that are independent of data type precision.
 - Persistent Buffers:** Some buffers or cached memory may not scale down with the data type.
 - Non-scalable components:** Components like model architecture information and memory reserved for system processes (e.g., GPU drivers).
3. **Improved Efficiency**
 - FP16 allows modern GPUs, such as the Tesla T4, to process computations more efficiently by handling two FP16 operations in parallel, making it not only memory-efficient but also computationally faster.

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: 85% | 11/136 [00:10<01:58, 1.06it/s]
took 10.405001648319824s

Question 8

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

Observations on Speedup:

1. **Observed Speedup:**
 - Processing batches in FP16 resulted in faster execution compared to FP32. The took 10.485 seconds for the initial 11 batches indicates improved performance due to the lower-precision format, as FP16 allows GPUs to perform two operations per ALU in the same time it takes to perform one FP32 operation.
2. **Expected Result:**
 - The speedup aligns with expectations, given that FP16 reduces the computational load and memory bandwidth requirements, both of which contribute to faster processing. This is particularly beneficial on GPUs like the Tesla T4, which are optimized for mixed-precision computation.

Pros and Cons of Using Lower-Precision Formats (FP16):

Pros:

1. **Reduced Memory Usage:**
 - FP16 halves the memory footprint for weights, activations, and gradients compared to FP32, enabling larger batch sizes and models to fit into the same GPU memory.
2. **Faster Computation:**
 - Modern GPUs support specialized hardware for FP16, allowing two FP16 operations to execute simultaneously, resulting in significant speedups.
3. **Energy Efficiency:**
 - Lower precision reduces energy consumption, which is critical for large-scale training or deployment.
4. **Scalability:**
 - Enables training and inference of larger models on the same hardware due to reduced resource requirements.

Cons:

1. **Numerical Precision Issues:**
 - FP16 has a smaller dynamic range and reduced numerical precision compared to FP32. This can lead to instability in training, such as exploding or vanishing gradients, particularly in models sensitive to small changes in weight updates.
2. **Compatibility Limitations:**
 - Some operations, such as certain mathematical functions or layers, may not fully support FP16, requiring fallback to FP32, which can introduce inefficiencies.
3. **Potential Loss in Accuracy:**
 - For tasks requiring high precision, FP16 may lead to slight degradation in accuracy, especially in edge cases with large or very small values.

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
import matplotlib.pyplot as plt
from tqdm import tqdm
```

```
# Reinitialize the dataloader for 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

# Process the entire dataset
with torch.no_grad():
    for inputs, labels in tqdm(dataloader, desc="Processing batches"):
        total_samples += labels.size(0)
        inputs, labels = inputs.half().cuda(), labels.cuda()

        # Predictions for each model
        preds_resnet18 = resnet18_model(inputs)
        preds_resnet50 = resnet50_model(inputs)
        preds_resnet152 = resnet152_model(inputs)
        preds_mobilenet = mobilenet_v2_model(inputs)

        # Calculate top-1 accuracy
        accuracies["ResNet-18"] += (preds_resnet18.argmax(dim=1) == labels).sum().item()
        accuracies["ResNet-50"] += (preds_resnet50.argmax(dim=1) == labels).sum().item()
        accuracies["ResNet-152"] += (preds_resnet152.argmax(dim=1) == labels).sum().item()
        accuracies["MobileNetV2"] += (preds_mobilenet.argmax(dim=1) == labels).sum().item()

# Convert counts to percentages
for model in accuracies:
    accuracies[model] = (accuracies[model] / total_samples) * 100

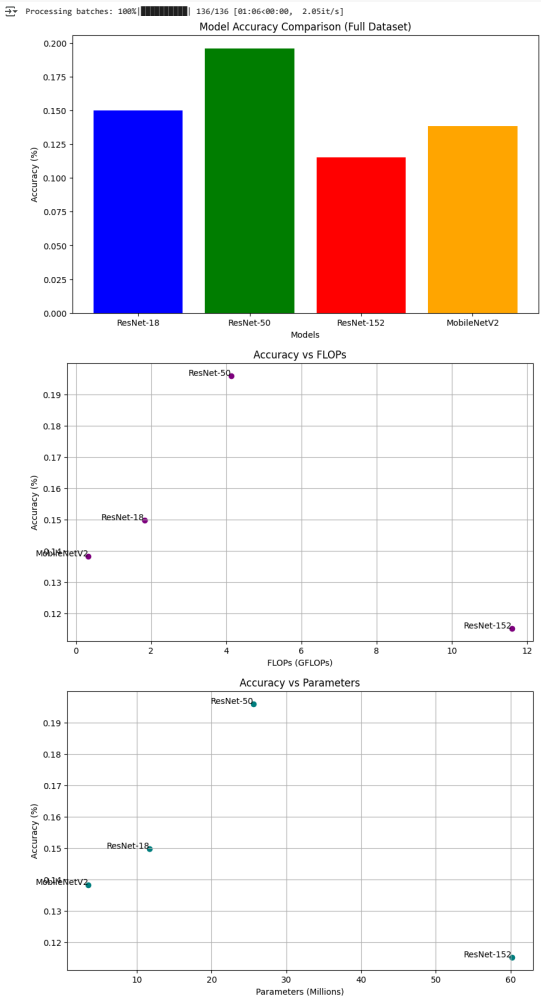
# Plot accuracy comparison
plt.figure(figsize=(10, 6))
plt.bar(accuracies.keys(), accuracies.values(), color=['blue', 'green', 'red', 'orange'])
plt.xlabel("Models")
plt.ylabel("Accuracy (%)")
plt.title("Model Accuracy Comparison (Full Dataset)")
plt.show()

# FLOPs and parameters for each model
flops_params = {
    "ResNet-18": (1.82, 11.69),
    "ResNet-50": (4.13, 25.56),
    "ResNet-152": (11.68, 60.19),
    "MobileNetV2": (0.33, 3.30),
}

# Extract data for plotting
models = list(flops_params.keys())
flops = [fp[0] for fp in flops_params.values()]
params = [fp[1] for fp in flops_params.values()]
accuracy = list(accuracies.values())

# Plot Accuracy vs FLOPs
plt.figure(figsize=(10, 6))
plt.scatter(flops, accuracy, color='purple', label="Accuracy vs FLOPs")
for i, model in enumerate(models):
    plt.text(flops[i], accuracy[i], model, fontsize=10, ha='right')
plt.xlabel("FLOPs (GFLOPs)")
plt.ylabel("Accuracy (%)")
plt.title("Accuracy vs FLOPs")
plt.grid(True)
plt.show()

# Plot Accuracy vs Parameters
plt.figure(figsize=(10, 6))
plt.scatter(params, accuracy, color='teal', label="Accuracy vs Params")
for i, model in enumerate(models):
    plt.text(params[i], accuracy[i], model, fontsize=10, ha='right')
plt.xlabel("Parameters (Millions)")
plt.ylabel("Accuracy (%)")
plt.title("Accuracy vs Parameters")
plt.grid(True)
plt.show()
```



Question 10

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

Based on your last model performance analysis, you could evaluate whether batch 10 shows any significant differences compared to the full dataset by focusing on the following aspects:

- Model Accuracy or Loss:** Check if there is a noticeable difference in the model's accuracy or loss metrics when trained on batch 10 compared to the full dataset. Significant changes in performance might indicate that batch 10 differs in its characteristics.
- Overfitting or Underfitting:** If the model trained on batch 10 exhibits overfitting or underfitting compared to the model trained on the full dataset, this could be a sign that batch 10 doesn't fully represent the diversity of the full dataset.
- Feature Importance or Weights:** Look at whether the feature importance or learned weights differ between models trained on the full dataset versus batch 10. If batch 10 causes the model to focus more on certain features, it might indicate that batch 10 has different characteristics.

4. **Generalization:** Assess how well the model trained on batch 10 generalizes to new data or a validation set. Poor generalization might indicate that batch 10 is not a good representation of the full dataset.

5. **Confusion Matrix or Classification Report:** If you're working with classification models, compare confusion matrices or classification reports (precision, recall, F1-score) to see if batch 10 leads to imbalanced performance across classes.

By comparing these aspects, you can discern if batch 10 is influencing the model's performance differently than when using the full dataset.