

The other models you will use:

- 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-visually-inspecting-resnets-4472248311ed>

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, so you can learn more details regarding the structure from here if interested: https://medium.com/@julian_gonzalez/mobilenet-v2-in-depth-a-visual-guide-4e9998c1242a

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 \rightarrow string mapping, which will turn the classes into a human readable form.

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

# define a denorm helper function - this undos the dataloader normalization so we can see the images better
def denorm(tensor):
    tensor = tensor.permute(0, 2, 1) # Change from C,H,W to H,W,C
    tensor -= tensor.mean(dim=0, keepdim=True) # Denormalize if the tensor was normalized
    tensor = tensor*0.2 + 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.show()

# Similarly, let's create an iidx2label helper function
def iidx2label(iidx):
    """Display a tensor of images."""
    tensor = tensor.permute(0, 2, 1) # Change from C,H,W to H,W,C
    tensor -= tensor.mean(dim=0, keepdim=True) # Denormalize if the tensor was normalized
    tensor = tensor*0.2 + 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.show()
```

```
# For our initial look, 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_largemodel(images.cuda()*0.5)
```

We can sample the output using argmax (Find the class with the highest probability)

Here we are calling output.logits.argmax() because huggingface returns a struct rather than a tuple

Also, we are using .cpu() to move the data from the GPU to the CPU. This corresponds to the classes - the shape is B,C

And we also need to re-ew the IDs 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 when calling id2label(), to get the raw contents of the IDs tensor

labels = [i.item() for i in ids]

for id in labels:
 labels.append(vit_largemodel.config.id2label[id.item()])

A quick look - plot the first 4 images

axs=plt.subplots(4, 4, figsize=(12, 12))

for i in range(4):
 for j in range(4):
 id = i*4+j
 axs[i,j].imshow(images[id])
 axs[i,j].set_title(labels[id])

We need to trim the labels because they sometimes are too long

if len(labels[id]) > max_label_len:
 trimmed_label = labels[id][:-max_label_len] + '...

else:
 trimmed_label = labels[id]

axs[i,j].set_title(trimmed_label)

plt.tight_layout()
plt.show()

warplane, military plane



American egret, great whi...



accordion, piano accordio...



airliner



ringlet, ringlet butterfl...

sunscreen, sunblock, sun ...



lionfish



holster



flamingo



pot, flowerpot



cheetah, cheetah, Acinonyx...



disk brake, disc brake



stretcher



book jacket, dust cover, ...



monarch, monarch butterfly...



oboe, hautboy, hautbois



stretcher



book jacket, dust cover, ...



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/weka/datasets/IMAGENET>

Please answer below:

Model Performance: Given the above analysis, the larger models (ViT/L/16, ResNet152) might not significantly outperform smaller models (ResNet50, ResNet18, MobileNetV2) on the Caltech101 dataset due to the limited size and diversity of the training set. In fact, smaller models are likely to perform better or at least more efficiently, as they are less prone to overfitting.

Limitations: These limitations are more likely related to the training set size and model complexity. Models like ViT and ResNet152 are highly complex and require large, diverse datasets to shine. For a smaller dataset like Caltech101, the models may not fully leverage their potential, and their performance may be constrained by the limited data variety.

Add blockquote

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 09:48:26 2025
| NVIDIA-SMI 535.184.05                 Driver Version: 535.184.05      CUDA Version: 12.2 |
| GPU  Name Persistence-M  Bus-Id Disp-A Disp-C Volatile Uncorr. ECC
| Fan  Temp  Perf  Pwr/Usage/Cap | Memory-Usage  GPU-Util  Compute M. 
|=====================================================================
| 0  GeForce RTX 3080 Ti L 00000000-00-00-00  16930MB / 13848MB   0% Default N/A |
+-----+
Processes:
| GPU  GI  CI   PID  Type  Process name                               GPU Memory
|=====================================================================
+-----+
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 09:48:35 2025
| NVIDIA-SMI 535.184.05                 Driver Version: 535.184.05      CUDA Version: 12.2 |
| GPU  Name Persistence-M  Bus-Id Disp-A Disp-C Volatile Uncorr. ECC
| Fan  Temp  Perf  Pwr/Usage/Cap | Memory-Usage  GPU-Util  Compute M. 
|=====================================================================
| 0  GeForce RTX 3080 Ti L 00000000-00-00-00  16930MB / 13848MB   0% Default N/A |
+-----+
```

0 Tesla T4	0/1	00000000:00:04:0 Off	0
N/A GBC P0	MW / 70W	1715MHz / 153MHz	0%
Processes:	ID	Process name	GPU Memory Usage
GPU ID	ID		

If you check above you should see the GPU memory utilization change from before and after the `empty_cached()` 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 PyTorch retains memory for efficiency, even after using `torch.cuda.empty_cached()`. 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:

```
# helper function to get element sizes in bytes
def size_per_element(tensor):
    # Get 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("Unknown 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
```

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.

```
import torch
import time
from torch import tpu

# Dictionary to store results
accuracies = {"ResNet-50": 0, "ResNet-18": 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, MobileNetV2) are already loaded
with torch.no_grad():
    for i, (inputs, _) in 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 prediction
        logits_resnet18 = resnet18_model(inputs) # Assuming resnet18_model is already loaded
        top1_resnet18 = logits_resnet18.topk(1, dim=1).indices
        matches_resnet18 = (baseline_preds.unsqueezed(1) == top1_resnet18).any(dim=1).float().sum().item()

        # ResNet-50 prediction
        logits_resnet50 = resnet50_model(inputs) # Assuming resnet50_model is already loaded
        top1_resnet50 = logits_resnet50.topk(1, dim=1).indices
        matches_resnet50 = (baseline_preds.unsqueezed(1) == top1_resnet50).any(dim=1).float().sum().item()

        # ResNet-101 prediction
        logits_resnet101 = resnet101_model(inputs) # Assuming resnet101_model is already loaded
        top1_resnet101 = logits_resnet101.topk(1, dim=1).indices
        matches_resnet101 = (baseline_preds.unsqueezed(1) == top1_resnet101).any(dim=1).float().sum().item()

        # MobileNetV2 prediction
        logits_mobilenetv2 = mobilenetv2_model(inputs) # Assuming mobilenet_v2_model is already loaded
        top1_mobilenetv2 = logits_mobilenetv2.topk(1, dim=1).indices
        matches_mobilenetv2 = (baseline_preds.unsqueezed(1) == top1_mobilenetv2).any(dim=1).float().sum().item()

        # Update accuracy
        accuracies["ResNet-18"] += matches_resnet18
        accuracies["ResNet-50"] += matches_resnet50
        accuracies["ResNet-101"] += matches_resnet101
        accuracies["MobileNetV2"] += matches_mobilenetv2

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

    # Print time taken for each batch (optional)
    print(f"Processing batch {i+1}/{num_batches} took {(time.time() - t_start):.2f} seconds")

# Finalize the accuracy by computing the average top-5 accuracy
for model_name, accuracy in accuracies.items():
    print(f"{model_name}: {accuracy / total_samples}")

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

```
Processing batches: 100% | 1/10 [00:03:00:27, 1.46s/it] Processing batch 1/10 took 3.97 seconds
Processing batches: 200% | 2/10 [00:03:00:24, 3.49s/it] Processing batch 2/10 took 6.14 seconds
Processing batches: 300% | 3/10 [00:03:00:21, 3.49s/it] Processing batch 3/10 took 9.16 seconds
Processing batches: 400% | 4/10 [00:03:00:18, 3.49s/it] Processing batch 4/10 took 12.23 seconds
Processing batches: 500% | 5/10 [00:03:00:15, 3.49s/it] Processing batch 5/10 took 15.23 seconds
Processing batches: 600% | 6/10 [00:03:00:12, 3.49s/it] Processing batch 6/10 took 18.23 seconds
Processing batches: 700% | 7/10 [00:03:00:09, 3.49s/it] Processing batch 7/10 took 21.48 seconds
Processing batches: 800% | 8/10 [00:03:00:06, 3.49s/it] Processing batch 8/10 took 24.59 seconds
Processing batches: 900% | 9/10 [00:03:00:03, 3.46s/it] Processing batch 9/10 took 27.14 seconds
Processing batches: 1000% | 10/10 [00:03:00:00, 3.46s/it] Processing batch 10/10 took 30.00 seconds
Overall Top-5 Accuracy:
ResNet-18: 65.00%
ResNet-50: 77.00%
ResNet-101: 76.50%
MobileNetV2: 67.81%
```

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
# Dictionary containing the top-5 accuracies for the models
accuracies = {
    "ResNet-18": 0.75, # Example accuracy values, replace with actual results
    "ResNet-50": 0.88,
    "ResNet-101": 0.75,
    "MobileNetV2": 0.78
}

# Get the model names and their corresponding accuracies
models = [list(accuracies.keys())]
accuracies = [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 + 0.05, f'{accuracy * 100:.2f}%', ha='center', fontsize=10)

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


