# EE 5545 Assignment 1

# Ethan Glaser (eg492)

Question 1:

**Name, Year, Type, peak FLOPs, peak memory bandwidth, source**

NVIDIA GeForce 256 SDR, 1999, GPU, 50 GFLOPs, 2.656 GB/s, source source

NVIDIA TITAN Xp, 2017, GPU, 11366.4 GFLOPs, 547.7 GB/s, source

NVIDIA TESLA K80, 2014, GPU, 4.113 TFLOPs, 240.6 GB/s, source

Intel i7 920 2008, CPU, 63 GFLOPs, 25.6 GB/s, source

AMD Ryzen 9 3900X, 2019, CPU, 2649.6 GFLOPs, 47.68 GB/s, source

AMD Radeon RX 6800, 2019, GPU, 18688 GFLOPS, 512 GB/s, source

Apple A14, 2020, SoC, 1000 GFLOPs, 42.7 GB/s, source

Apple M1 Max, 2020, SoC, 2.7 TFLOPs, 330 GB/s, source

TPU v3, 2018, ASIC, 420 TFLOPs, 1 TB/s, source

Most of the results are as expected – the GPU from 1999 and CPU on my laptop are both significantly lower performing (in terms of FLOPS/s and memory bandwidth) compared to the more traditional GPUs. Certain GPUs were higher performance than others, as mid-range and high-end GPUs were selected from both NVIDIA and AMD. The TPU had by far the highest FLOPS/s and memory bandwidth than the GPUs. These results are illustrated in the roofline plot below. No FLOPS/s values were calculated as the data was available on the various datasheets and pages associated with each device.
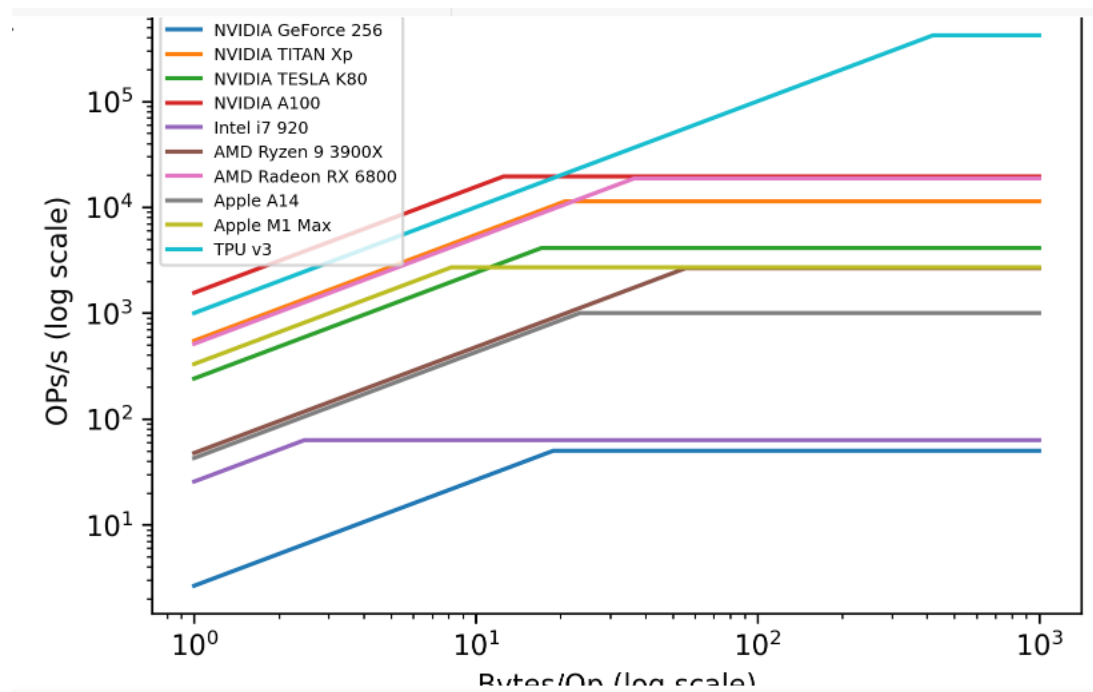


Figure 1: Roofline Plot for Q1

```python
1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 def roofline(input_dictionary, input2_dictionary={}, performance_dictionary={}, legend_size=6):
5     plt.figure(dpi=270)
6     for device in input_dictionary.keys():
7         x = np.linspace(1, 1000)
8         current_gf = input_dictionary[device]['GF']
9         current_mb = input_dictionary[device]['MB']
10        meeting = current_gf / current_mb
11        y = np.array([current_mb * current_x if current_x < meeting else current_gf for current_x in x ])
12        #plt.plot(x, y, label=device)
13        plt.plot([1, current_gf / current_mb, 1000], [current_mb, current_gf, current_gf], label=device)
14        if device == 'NVIDIA TESLA K80' and len(input2_dictionary.keys()):
15            for model in input2_dictionary.keys():
16                oi = input2_dictionary[model]['oi']
17                if len(performance_dictionary.keys()):
18                    y_vals = [current_mb * oi if oi < meeting else current_gf]
19                    for batch in [1, 64, 256]:
20                        if 'batch_' + str(batch) + '_time' in performance_dictionary[model]['cpu'].keys():
21                            y_vals.append(performance_dictionary[model]['GF'] / performance_dictionary[model]['gpu']['batch_' + str(batch) + '_time'])
22                    plt.scatter([oi] * len(y_vals), y_vals, marker="*", label=model)
23                else:
24                    plt.scatter([oi], [current_mb * oi if oi < meeting else current_gf], marker="*", label=model)
25        if device == 'Intel i7 920' and len(input2_dictionary.keys()):
26            for model in input2_dictionary.keys():
27                oi = input2_dictionary[model]['oi']
28                if len(performance_dictionary.keys()):
29                    y_vals = [current_mb * oi if oi < meeting else current_gf]
30                    for batch in [1, 64, 256]:
31                        if 'batch_' + str(batch) + '_time' in performance_dictionary[model]['cpu'].keys():
32                            y_vals.append(performance_dictionary[model]['GF'] / performance_dictionary[model]['cpu']['batch_' + str(batch) + '_time'])
33                    plt.scatter([oi] * len(y_vals), y_vals, marker="P", label=model)
34                else:
35                    plt.scatter([oi], [current_mb * oi if oi < meeting else current_gf], marker="P", label=model)
36    plt.xscale('log')
37    plt.yscale('log')
38    plt.xlabel('Bytes/Op (log scale)')
39    plt.ylabel('OPs/s (log scale)')
40    plt.legend(loc=0, prop={'size': legend_size})
41    plt.show()
42
43
44 dict1 = {'NVIDIA GeForce 256': {'GF': 50, 'MB': 2.656},
45          'NVIDIA TITAN Xp': {'MB': 547.7, 'GF': 11366.4},
46          'NVIDIA TESLA K80': {'GF': 4113, 'MB': 240.6},
47          'NVIDIA A100': {'GF': 19500, 'MB': 1555},
48          'Intel i7 920': {'GF': 63, 'MB': 25.6},
49          'AMD Ryzen 9 3900X': {'GF': 2649.6, 'MB': 47.68},
50          'AMD Radeon RX 6800': {'GF': 18688, 'MB': 512},
51          'Apple A14': {'GF': 1000, 'MB': 42.7},
52          'Apple M1 Max': {'GF': 2700, 'MB': 330},
53          'TPU v3': {'GF': 420000, 'MB': 1000}
54          }
55 roofline(dict1)
```

Figure 2: Code to plot roofline, used in Q1 and expanded on for use in Q2 and Q4

Question 2:

Model, FLOPs, Memory footprint

ResNet50, 8.22 GFLOPs, 384.62 MB

AlexNet, 1.43 GFLOPs, 242.03 MB

MobileNet v3 large, 0.45 GFLOPs, 126.9 MB

VGG19, 39.29 GFLOPs, 787.31 MB

SqueezeNet1_0, 1.65 GFLOPs, 97.14 MB

EfficientNetb7, 10.53 GFLOPs, 1.33 GB

GoogleNet, 3.01 GFLOPs, 119.95 MB

ShuffleNet, 0.09 GFLOPs, 30.85 MB

WideResNet101_2, 45.59 GFLOPs, 1.03 GB

Inception v3, 5.69 GFLOPs, 209.48 MB

Every model was compute bound on the CPU, whereas about half (shufflenet, alexnet, efficientnet, mobilenet v3) with relatively lower operational intensities were memory bound on the GPU provided by colab (Nvidia TESLA K80) – as seen in Figure 4.
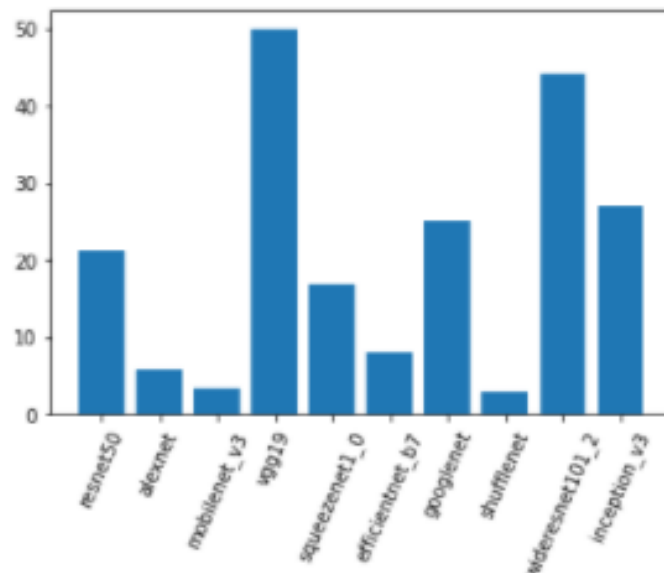


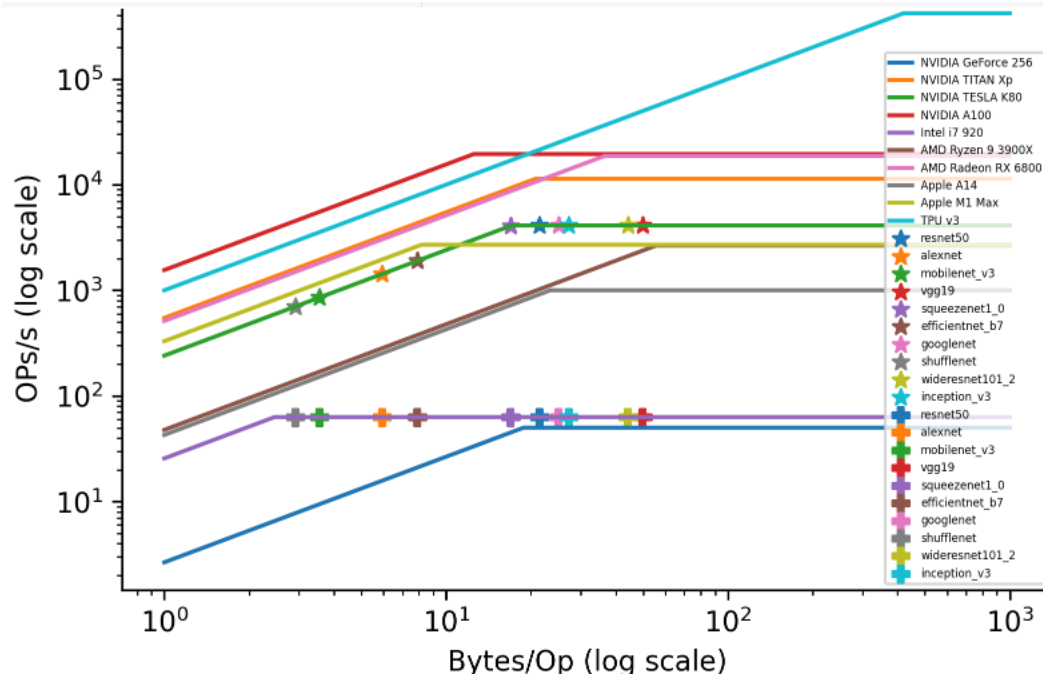Figure 3: Operational Intensity (Bytes/Op) for each of 10 DNN models

Figure 4: Updated Roofline Plot for Q2 with expected model performance for GPU and CPU indicated

```
1
2 def add_to_roofline(input_dict):
3     for model in input_dict.keys():
4         op_intensity = input_dict[model]['GF'] / input_dict[model]['MF'] * 1000
5         input_dict[model]['oi'] = op_intensity
6     plt.bar(input_dict.keys(), [input_dict[i]['oi'] for i in input_dict.keys()])
7     plt.xticks(rotation = 66)
8     plt.show()
9     return input_dict
10
11
12
13 dict2 = {
14     'resnet50': {'GF': 8.22, 'MF': 384.62},
15     'alexnet': {'GF': 1.43, 'MF': 242.03},
16     'mobilenet_v3': {'GF': 0.45, 'MF': 126.9},
17     'vgg19': {'GF': 39.29, 'MF': 787.31},
18     'squeezenet1_0': {'GF': 1.65, 'MF': 97.14},
19     'efficientnet_b7': {'GF': 10.53, 'MF': 1330},
20     'googlenet': {'GF': 3.01, 'MF': 119.950},
21     'shufflenet': {'GF': 0.09, 'MF': 30.85},
22     'wideresnet101_2': {'GF': 45.59, 'MF': 1030},
23     'inception_v3': {'GF': 5.69, 'MF': 209.48}
24 }
25 dict2 = add_to_roofline(dict2)
26 roofline(dict1, input2_dictionary = dict2, legend_size=4)
```

Figure 5: Code used to generate Q2 results (note function call to main Q1 function with additional argument)

Question 3:

The Spearman Coefficients indicate that there is a strong correlation between FLOPs and CPU latency, suggesting that FLOPs are a good indicator of the CPU's runtime with a coefficient of 0.83. Number of parameters is also somewhat correlated to CPU latency but with a less strong relationship (0.39). These indicators are much less telling for GPU latency, with Spearman coefficients of -0.14 and -0.03 respectively. Therefore there is not a strong relationship between the FLOPs or parameters and the GPU latency, which may be a result of the parallelization of operations running on the GPU, since batches may be handled uniformly instead of in a more sequential manner like with the CPU. The plots of flops and parameters vs. latency aren't particularly revealing aside from showing the general trend that the latency is generally much higher on the CPU and the opposite for throughput.
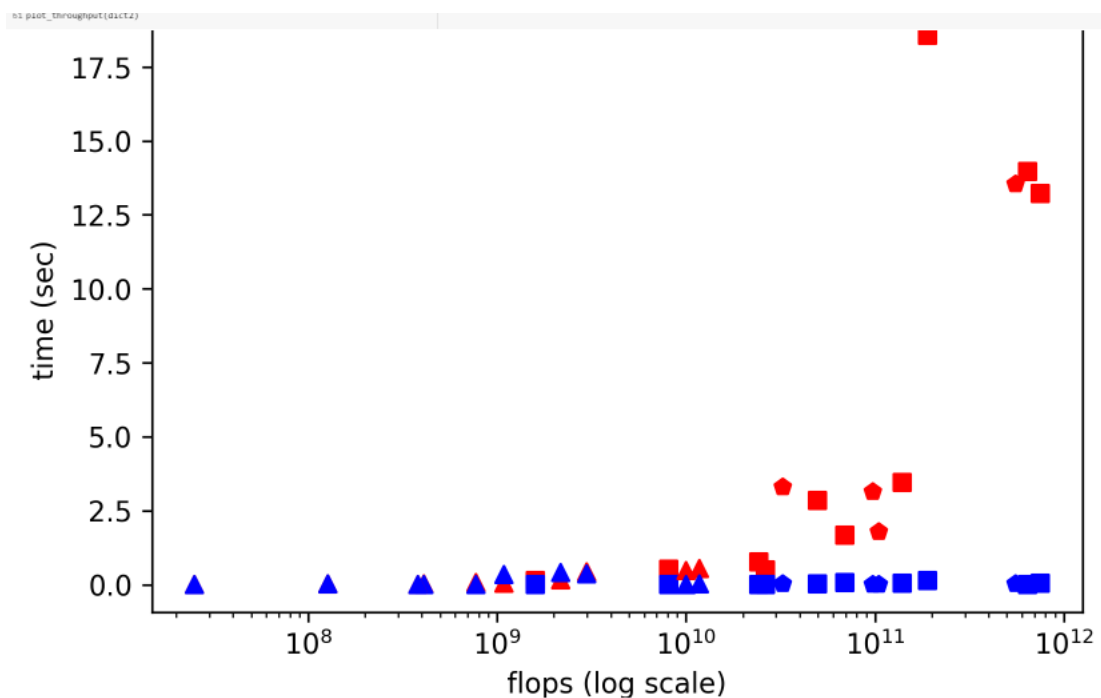


Figure 6: Latency vs FLOPs (where CPU in red, GPU in blue, batch size 1 as triangles, batch size 64 as squares, and batch size 256 as pentagons)
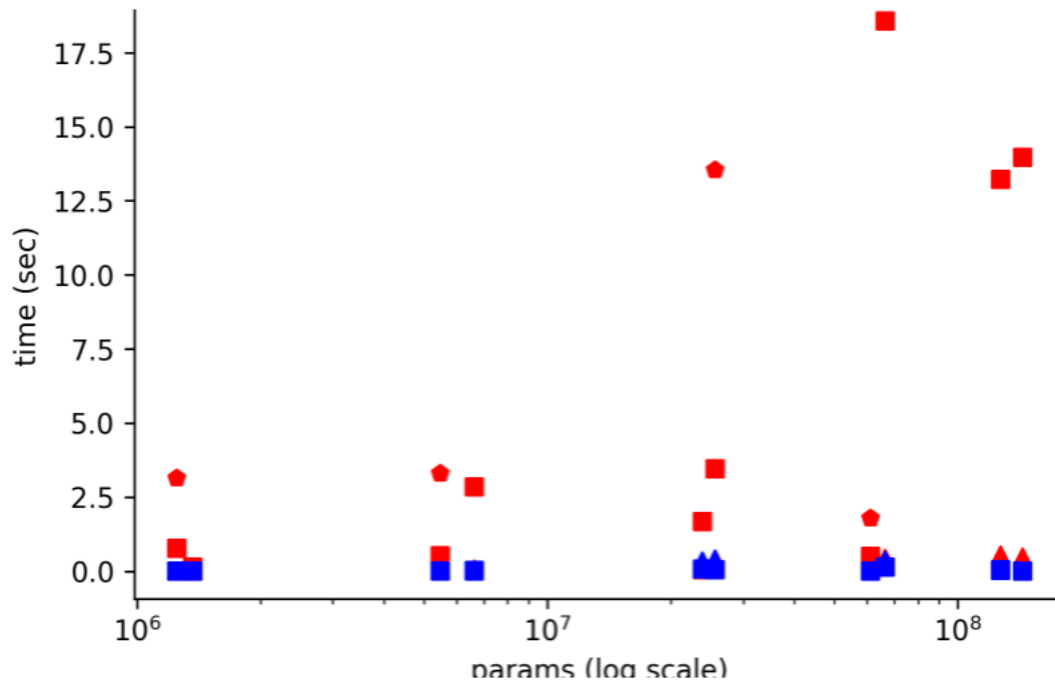
Figure 7: Latency vs parameters (where CPU in red, GPU in blue, batch size 1 as triangles, batch size 64 as squares, and batch size 256 as pentagons)



Figure 8: Spearman Correlation between CPU latency, GPU latency, FLOPs, and number of parameters
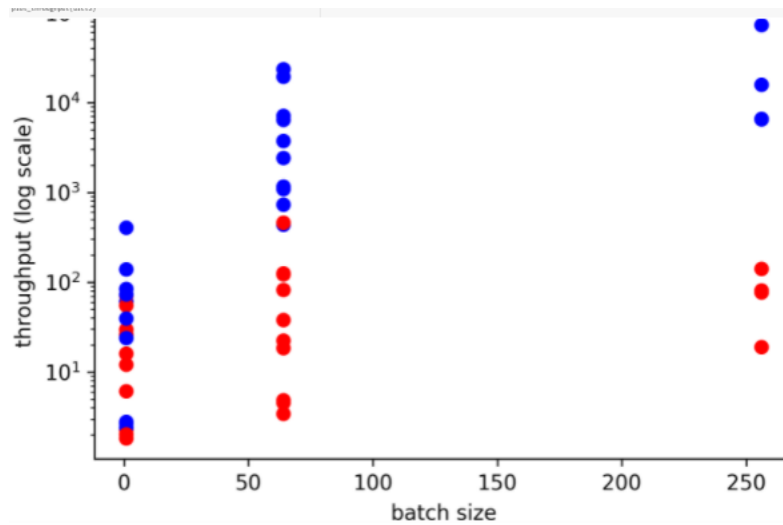
Figure 9: Throughput vs. Batch Size (CPU in red, GPU in blue)

```python
1  import pandas as pd
2  import seaborn as sns
3
4  legend_dict = {'color': {'cpu': 'r', 'gpu': 'b'}, 'shape': {1: '^', 64: 's', 256: 'p'}}
5  def performance_benchmarking_plots(input_dictionary, x_axis):
6      plt.figure(dpi=300)
7      for model in input_dictionary.keys():
8          for device in ['cpu', 'gpu']:
9              for key in input_dictionary[model][device].keys():
10                 if '_time' in key:
11                     current_batch_size = int(key.split('_')[1])
12                     y = input_dictionary[model][device][key]
13                     x = input_dictionary[model][device]['batch_' + str(current_batch_size) + '_' + x_axis]
14                     plt.scatter([x], [y], marker=legend_dict['shape'][current_batch_size], c=legend_dict['color'][device])
15     #plt.yscale('log')
16     plt.xscale('log')
17     plt.xlabel(x_axis + ' (log scale)')
18     plt.ylabel('time (sec)')
19     plt.show()
20
21 def get_spearman(input_dictionary):
22     df_dict = {'CPU Latency': [], 'GPU Latency': [], 'FLOPs': [], 'Parameters': []}
23     for model in input_dictionary.keys():
24         for key in input_dictionary[model]['cpu'].keys():
25             if '_time' in key:
26                 df_dict['CPU Latency'].append(input_dictionary[model]['cpu'][key])
27                 df_dict['FLOPs'].append(input_dictionary[model]['cpu'][key[:-4] + 'flops'])
28                 df_dict['GPU Latency'].append(input_dictionary[model]['gpu'][key])
29                 df_dict['Parameters'].append(input_dictionary[model]['cpu'][key[:-4] + 'params'])
30     df = pd.DataFrame.from_dict(df_dict)
31     spearman = df.corr(method="spearman")
32     plt.figure(figsize=(10,6), dpi=200)
33     heatmap = sns.heatmap(df.corr(), vmin=-1,
34                           vmax=1, annot=True)
35     plt.title("Spearman Correlation")
36     plt.show()
37
38 def plot_throughput(input_dictionary):
39     plt.figure(dpi=300)
40     for model in input_dictionary.keys():
41         for device in ['cpu', 'gpu']:
42             for key in input_dictionary[model][device].keys():
43                 if '_time' in key:
44                     current_batch_size = int(key.split('_')[1])
45                     x = current_batch_size
46                     y = current_batch_size / input_dictionary[model][device][key]
47                     plt.scatter([x], [y], c=legend_dict['color'][device])
48     plt.yscale('log')
49     #plt.xscale('log')
50     plt.xlabel('batch size')
51     plt.ylabel('throughput (log scale)')
52     plt.show()
53
54
55
56 performance_benchmarking_plots(dict2, 'flops')
57 performance_benchmarking_plots(dict2, 'params')
58
59 get_spearman(dict2)
60
61 plot_throughput(dict2)
```

Figure 10: Code used to complete Q3 with 3 main functions for each of the plots (first 2 bullet points combined into the first function)

Question 4:

There are many interesting trends that can be observed in Figure 11. One noteworthy observation that doesn't make a ton of sense is that some implementations of vgg19, alexnet, and a few other models outperformed the roofline plot limits, which may be due to incorrect values determined for the vgg19 earlier on, or differences in GPUs that were used in runtime compared to when the Nvidia Tesla K80 was first realized as the Colab GPU. For the most part, achieved performance was distributed in varying degrees of smallness compared to the expected peak performance.
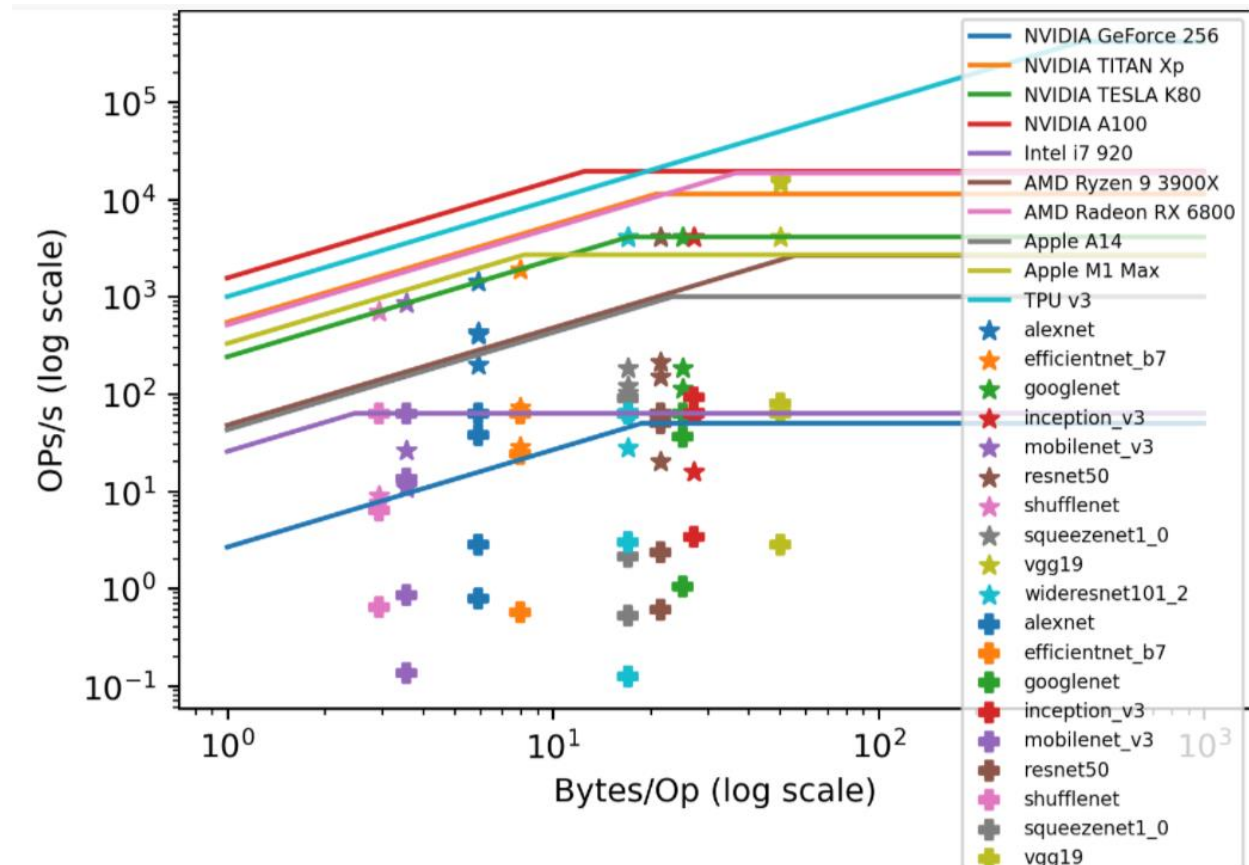


Figure 11: Roofline Plot for Q4 with actual performance added onto the plot

```
1 roofline(dict1, input2_dictionary=dict2, performance_dictionary=dict2)
```
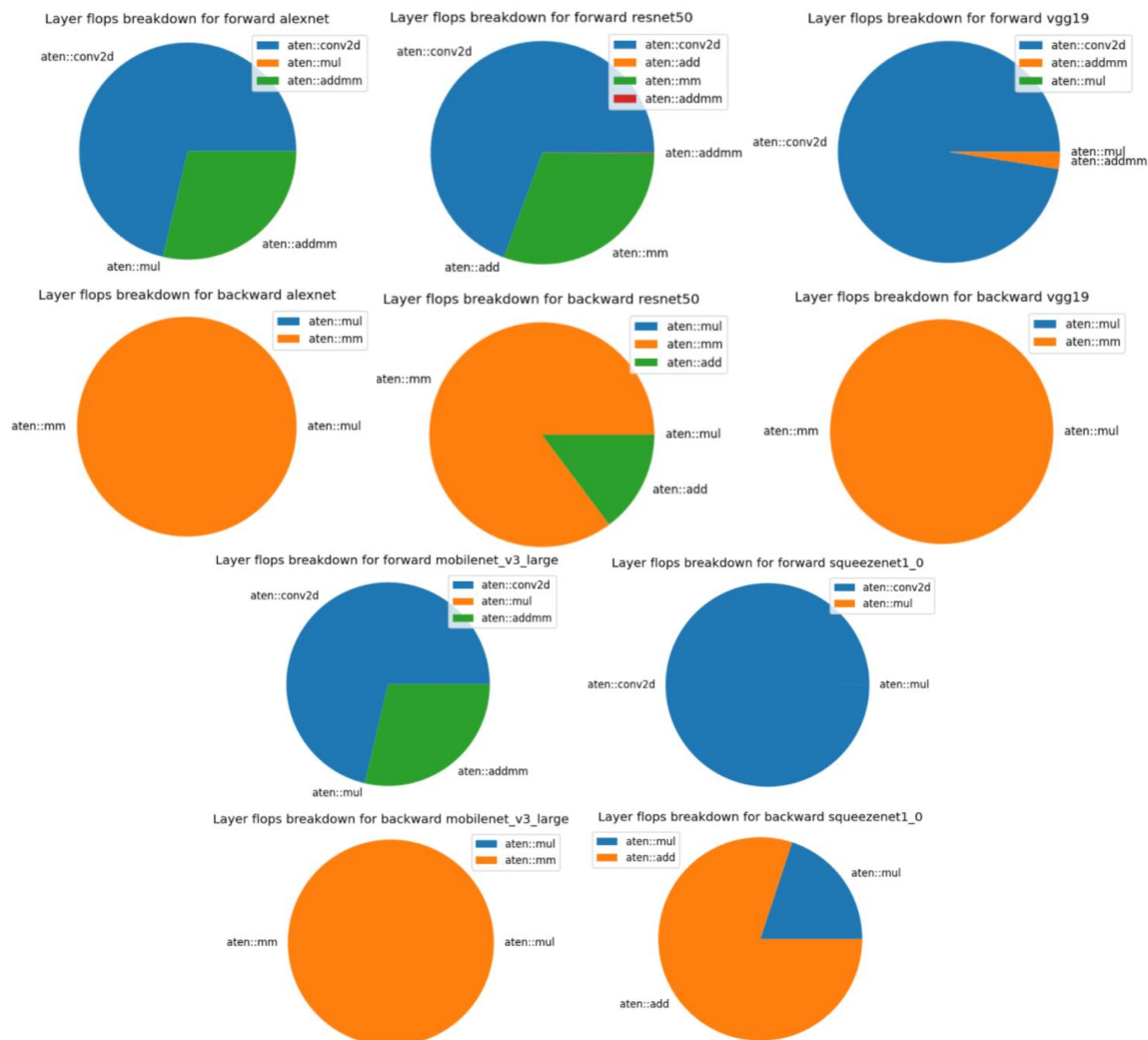
Figure 12: Code used to produce Q4 – again note that most changes were additions to the original function and can be observed in Figure 2

Question 5:

Based on the results of the 5 indicated DNNs that were tested, the ratio of latency for forward to backward pass can range anywhere from an approximately 2:3 relationship with squeezenet to an approximately 1:8 ratio as seen in alexnet. The flops vary drastically from 1.74:1 all the way to 3866:1. None of the flop ratios experimentally verify the 1:2 ratio that was expected based on the description from lecture.

```
Forward:Backward ratios for latency and flops for alexnet 0.12100476028561714 1.743724687862097
Forward:Backward ratios for latency and flops for resnet50 0.49713446455409416 322.54277922392424
Forward:Backward ratios for latency and flops for vgg19 0.25742234028428507 20.22366663835236
Forward:Backward ratios for latency and flops for mobilenet_v3_large 0.15994074763103963 1.743724687862097
Forward:Backward ratios for latency and flops for squeezenet1_0 0.6965994458567797 3866.176504392086
```

Figure 13: Examples of latency and flop forward pass to backward pass ratios for 5 DNNs



Figures 14-18: Pie Charts showing the proportion of flops from various layers for the forward and backward passes of 5 DNNs

```python
1  from traitlets.traitlets import default
2  from torch.profiler import profile, record_function, ProfilerActivity
3  import time
4  from torchvision import models
5  from collections import defaultdict
6
7  def ratios(model_name, model):
8      device = torch.device('cpu')
9      model = model.to(device)
10     batch_t = torch.rand(size=(1, 3, 112, 112)).to(device)
11     sample = torch.rand(size=(1,1000)).to(device)
12     output = model(batch_t).to(device)
13     loss = (sample - output).sum()
14     with profile(activities=[ProfilerActivity.CPU], with_flops=True, record_shapes=True) as fprof:
15         #with record_function("model_inference"):
16         start = time.time()
17         model(batch_t)
18         end = time.time()
19
20     flops = 0
21     fwd = defaultdict(int)
22     for event in fprof.events():
23         if event.flops is not None and event.flops > 0:
24             fwd[event.name] += event.flops
25             flops += event.flops
26     # Data to plot
27     labels = []
28     sizes = []
29
30     for x, y in fwd.items():
31         labels.append(x)
32         sizes.append(y)
33
34     plt.figure(dpi=120)
35     plt.title('Layer flops breakdown for forward ' + model_name)
36     plt.pie(sizes, labels=labels)
37     plt.axis('equal')
38     plt.legend()
39     plt.show()
40
41     with profile(activities=[ProfilerActivity.CPU], with_flops=True, record_shapes=True) as bprof:
42         start2 = time.time()
43         loss.backward()
44         end2 = time.time()
45     bflops = 0
46     bck = defaultdict(int)
47     for event in bprof.events():
48         if event.flops is not None and event.flops > 0:
49             bck[event.name] += event.flops
50             bflops += event.flops
51
52     labels = []
53     sizes = []
54
55     for x, y in bck.items():
56         labels.append(x)
57         sizes.append(y)
58
59     # Plot
60     plt.figure(dpi=120)
61     plt.title('Layer flops breakdown for backward ' + model_name)
62     plt.pie(sizes, labels=labels)
63
64     plt.axis('equal')
65     plt.legend()
66     plt.show()
67
68     print("Forward:Backward ratios for latency and flops for", model_name, (end - start) / (end2 - start2), flops / bflops)
69  for pair in [('alexnet', models.alexnet(pretrained=True)), ('resnet50', models.resnet50(pretrained=True)), ('vgg19', models.vgg19(pretrained=True)), ('mobilenet_v3_large', models.alexnet(pretrained=True)), ('squeezenet1_0', models.squeezenet1_0(pretrained=True))]:
70     ratios(pair[0], pair[1])
```

Figure 19: Code used to generate ratios and plots for Q5