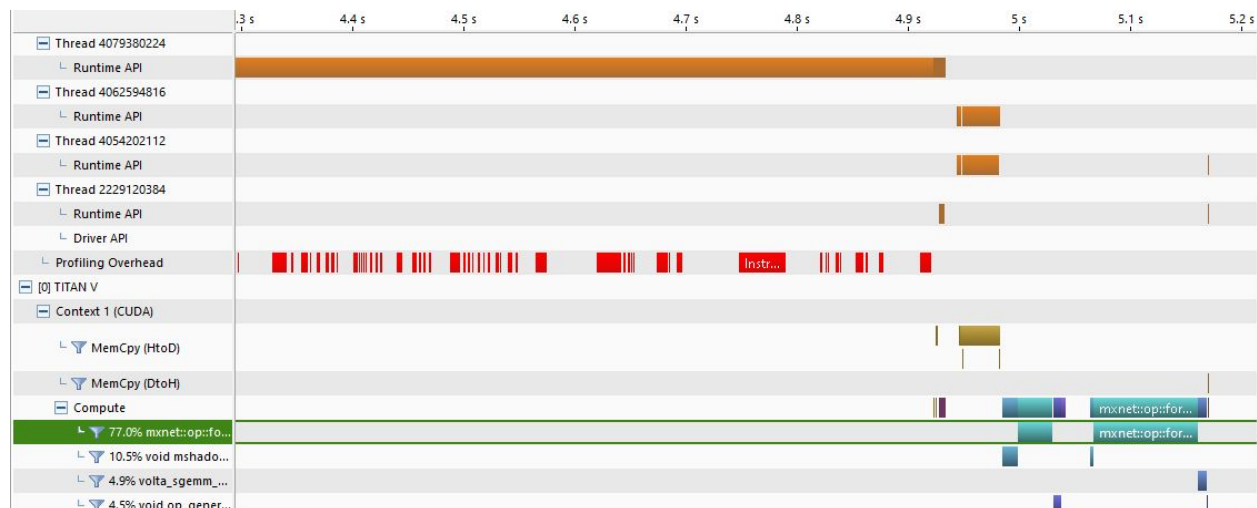


Team: Team
Bryan Lu (bryanlu2)
Omar Mbarki (ombarki2)
UIUC On-campus

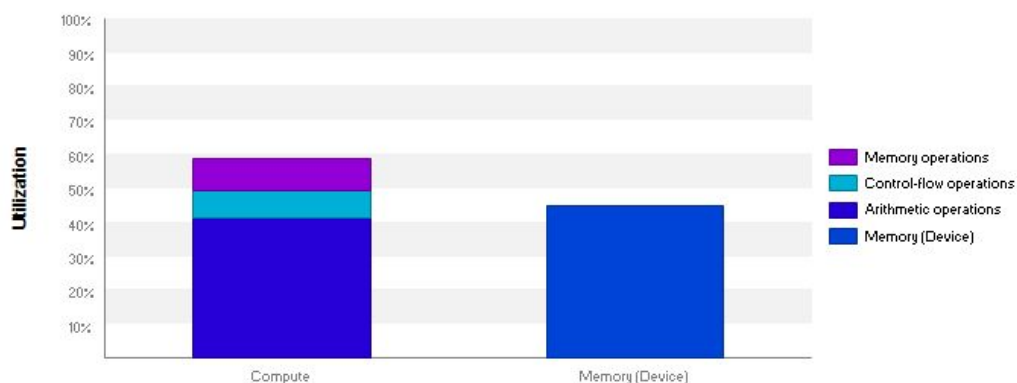
Milestone 3

Timeline.nvprof



i Kernel Performance Is Bound By Instruction And Memory Latency

This kernel exhibits low compute throughput and memory bandwidth utilization relative to the peak performance of "TITAN V". These utilization levels indicate that the performance of the kernel is most likely limited by the latency of arithmetic or memory operations. Achieved compute throughput and/or memory bandwidth below 60% of peak typically indicates latency issues.



It's no surprise that our very basic implementation of the kernel is bounded by memory latency. Using shared memory would solve the global memory bandwidth issue and would be much more efficient in terms of throughput.

Furthermore, we suffer 16.5% control divergence because of the if statement in line 33. There are ways to mitigate this issue, but we will focus our efforts into eliminating the memory bandwidth issue by using shared/constant memory and later down the line employing the unrolling technique to perform matrix multiplications instead of convolutions.

Milestone 2

Kernels that take up more than 90% of program time (ml.2.py):

- CUDA memcpy HtoD
- volta_scudnn_128x64_relu_interior_nn_v1
- volta_gcgemm_64x32_nt
- fft2d_c2r_32x32
- volta_sgemm_128x128_tn
- op_generic_tensor_kernel
- fft2d_r2c_32x32

API calls that take up more than 90% of program time (ml.2.py):

- cudaStreamCreateWithFlags
- cudaMemGetInfo
- cudaFree

CUDA kernels are functions that are run by each individual thread that CUDA creates. The API calls are functions that run on the program host that manage the information needed by the CUDA threads to run (like memory allocations).

Output of rai running MXNet on CPU with program times:

```
% Running /usr/bin/time python ml.1.py
loading fashion-mnist data... done
loading model... done
New Inference
EvalMetric: {'accuracy': 0.8154}
17.78user 4.40system 0:09.57elapsed 231%CPU (0avgtext+0avgdata 6044868maxresident)k
0inputs+2824outputs (0major+1604136minor)pagefault
s 0swaps
% The build folder has been uploaded to http://s3.amazonaws.com/files.rai-project.com/userdata/build-5d8fee98c43a32004f6
9d65.tar.gz. The data will be present for only a short duration of time.
% Server has ended your request.
```

CPU program run times:

- User: 17.78
- System: 4.40
- Elapsed: 9.57

Output of rai running MXNet on GPU with program times:

```

$ Running /usr/bin/time python m1.2.py
loading fashion-mnist data... done
loading model... done
New Inference
EvalMetric: {'accuracy': 0.8154}
4.96user 3.22system 0:04.61elapsed 177%CPU (0avgtext+0avgdata 2974788maxresiden
t)k
0inputs+4536outputs (0major+729927minor)pagefaults 0swaps
^ The build folder has been uploaded to http://s3.amazonaws.com/files.nai-project.com/userdata/build_9d9fef30cf3a32005ee
b26c9.tar.gz. The data will be present for only a short duration of time.
^ Server has ended your request

```

GPU program run times:

- User: 4.96
- System: 3.22
- Elapsed: 4.61

Output of CPU implementation:

```

* Running /usr/bin/time python m2.1.py
Loading fashion-mnist data... done
Loading model... done
New Inference
Op Time: 11.664715
Op Time: 60.728774
Correctness: 0.7653 Model: ece408
87.21user 10.91system 1:16.48elapsed 128%CPU

```

CPU implementation program times:

- User: 87.21
- System: 10.91
- Elapsed: 1:16.48

CPU implementation op times:

- 1: 11.664715
- 2: 60.728774