**Part 1: SAXPY**
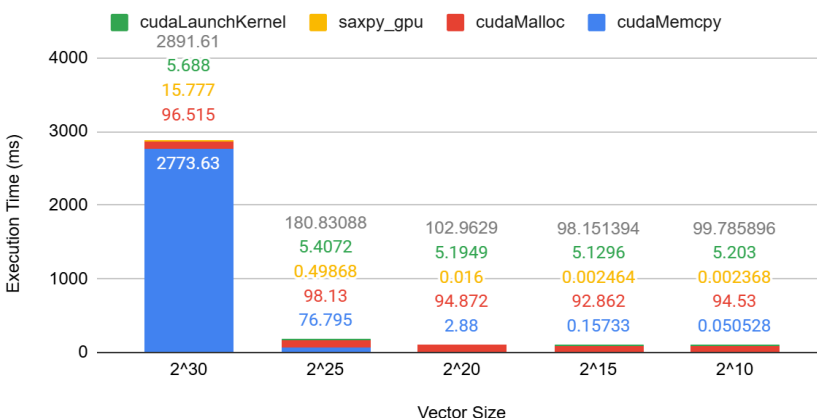
The first chart varies the vector size of the SAXPY operation and measures the change in execution time due to various aspects of the CUDA program. The categories can be split into three distinct types: constant (cudaMalloc, cudaLaunchKernel), compute (saxpy_gpu), and memory (cudaMemcpy). Until the vector size reaches roughly 2^25, the majority of the execution time is spent

Execution Times of API Calls and Kernels on Varying SAXPY Vector Sizes



on constant behaviors, mainly cudaMalloc. However, after the vector size passes 2^25, the execution time is dominated by cudaMemcpy. This shows that execution time is memory-bound, since copying data between device and host uses significantly more execution time than the kernel itself. Both the kernel and cudaMemcpy scale as vector size increases.
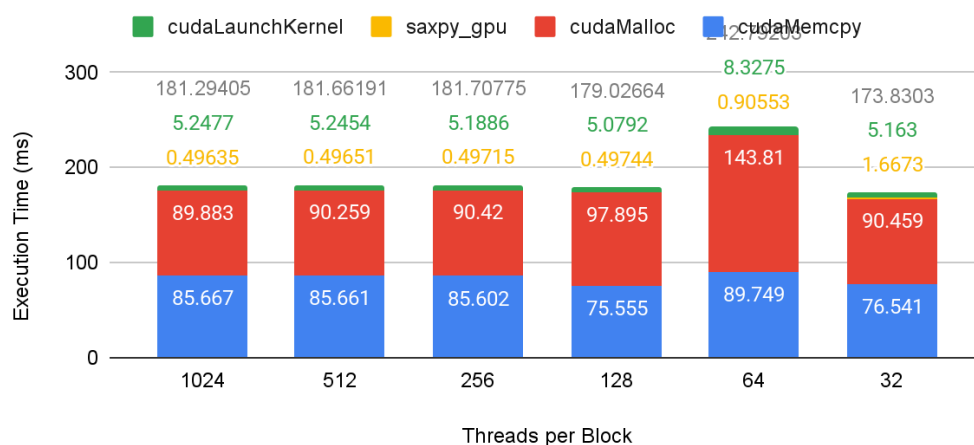
The second chart shows how execution time varies as the threads per block varies, i.e. scheduling more blocks rather than more threads per block. 2^25 vector size was chosen since this was when compute started becoming a more significant factor in execution time. No significant change was observed here, although when threads per block was very low (64 and 32), the kernel time doubled and

Execution Times of API Calls and Kernels on Varying SAXPY Threads per Block

Vector Size = 2^25



quadrupled, respectively. This is likely because more thread blocks needed to be scheduled since there was wasted resource usage within each thread block.
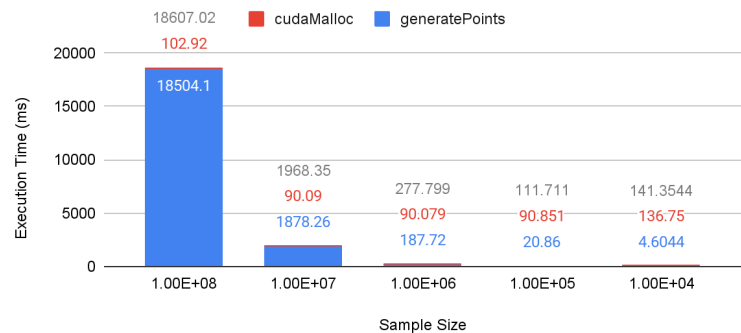
## Part 2: Monte Carlo PI Estimation

The first chart changes the sample size, the number of points calculated by each thread, and measures execution time. Unlike the SAXPY kernel, this program is compute-bound, as execution is heavily dominated by the kernel. The time spent on the generatePoints kernel seems to linearly scale with the sample size. This makes sense since each thread sequentially does the calculation for *sample size* times.

This chart shows execution time as the number of generation threads increases. An interesting finding is that execution time stays the same from 2^10 to 2^15 threads. This likely because as threads increase, the program is increasing compute in a parallel fashion, rather than a sequential fashion. Thus, it makes good use of the GPU's resources. However, once threads increase further, it saturates the compute of the GPU and execution time explodes. It is also worth noting that accuracy of the pi estimation increases as (generation threads * sample size) increases. Thus, it seems more feasible to increase the generation threads, since this will increase the accuracy more without as big of a penalty on execution time.

The last comparison done was changing the *reduce size* parameter. However, this led to no significant changes in execution time.
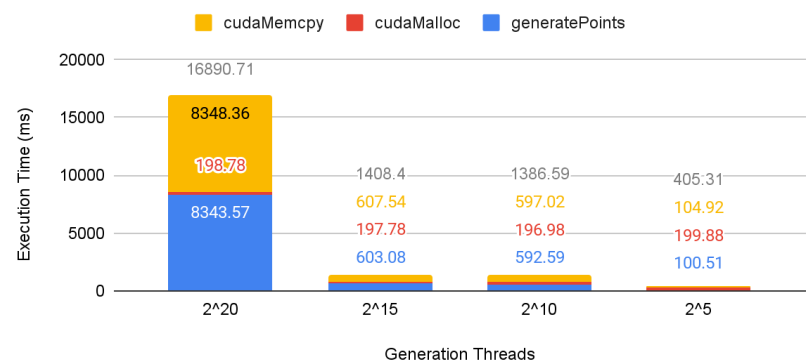
**Execution of Monte Carlo PI Estimation with Varying Sample Size**

Generation Threads = 1024, Reduce Size = 32

| | cudaMalloc | generatePoints |

| Sample Size | 1.00E+08 | 1.00E+07 | 1.00E+06 | 1.00E+05 | 1.00E+04 |
|---|---|---|---|---|---|
| total (gray) | 18607.02 | 1968.35 | 277.799 | 111.711 | 141.3544 |
| cudaMalloc (red) | 102.92 | 90.09 | 90.079 | 90.851 | 136.75 |
| generatePoints (blue) | 18504.1 | 1878.26 | 187.72 | 20.86 | 4.6044 |

**Execution Time of Monte Carlo PI Estimation with Varying Number of Generation Threads**

Sample Size = 1e6, Reduce Size = 32

| | cudaMemcpy | cudaMalloc | generatePoints |

| Generation Threads | 2^20 | 2^15 | 2^10 | 2^5 |
|---|---|---|---|---|
| total (gray) | 16890.71 | 1408.4 | 1386.59 | 405.31 |
| cudaMemcpy (orange) | 8348.36 | 607.54 | 597.02 | 104.92 |
| cudaMalloc (red) | 198.78 | 197.78 | 196.98 | 199.88 |
| generatePoints (blue) | 8343.57 | 603.08 | 592.59 | 100.51 |

**Execution Time of Monte Carlo PI Estimation as Reduce Size Varies**

Sample Size = 1e6, Generation Threads = 1024

| | cudaMemcpy | cudaMalloc | generatePoints |

| Reduce Size | 1024 | 256 | 32 | 8 | 1 |
|---|---|---|---|---|---|
| total (gray) | 925.81 | 925.9 | 926.13 | 927.06 | 925.7 |
| cudaMemcpy (orange) | 397.57 | 397.58 | 397.45 | 397.56 | 397.56 |
| cudaMalloc (red) | 132.99 | 132.99 | 133.46 | 134.18 | 132.82 |
| generatePoints (blue) | 395.25 | 395.33 | 395.22 | 395.32 | 395.32 |