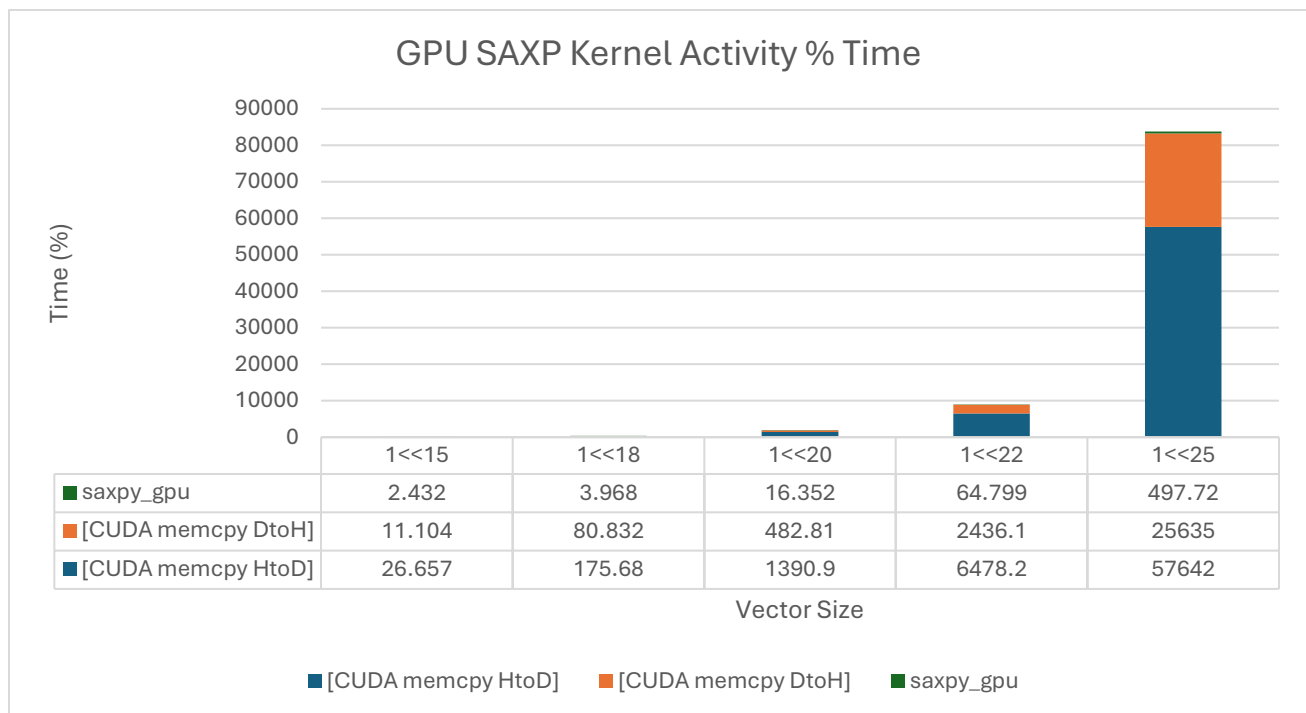
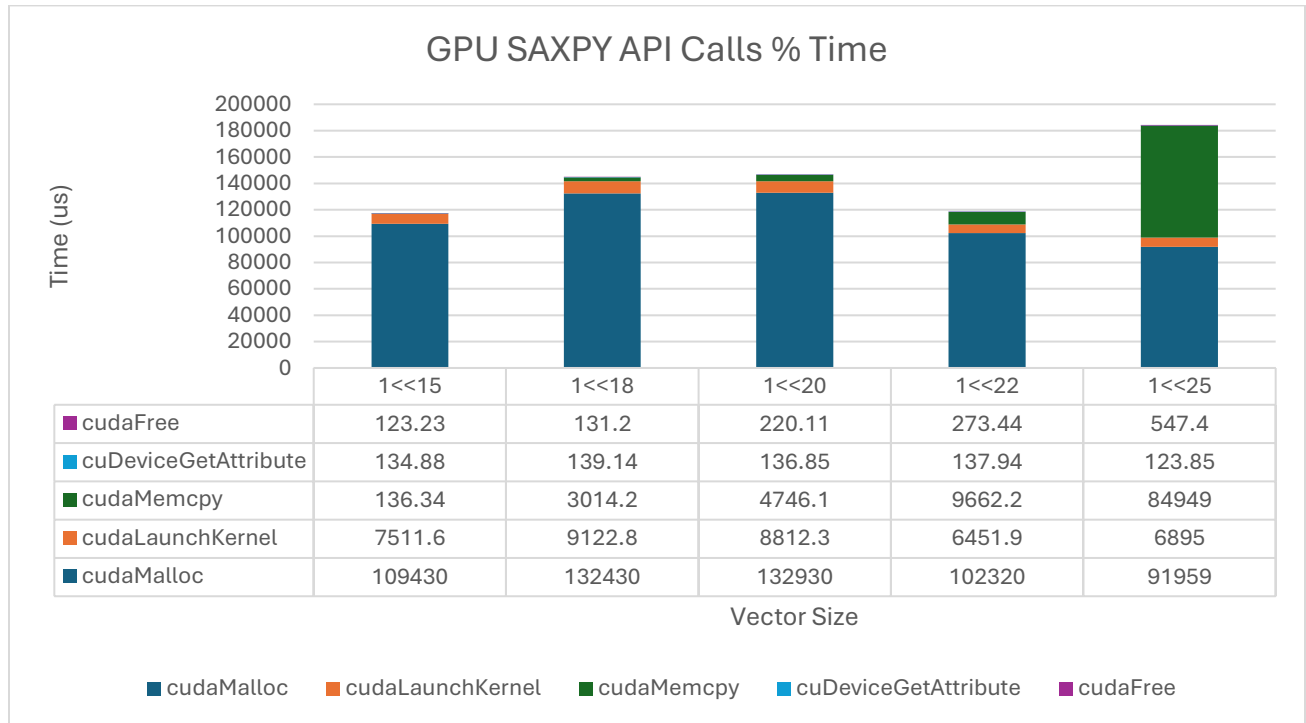


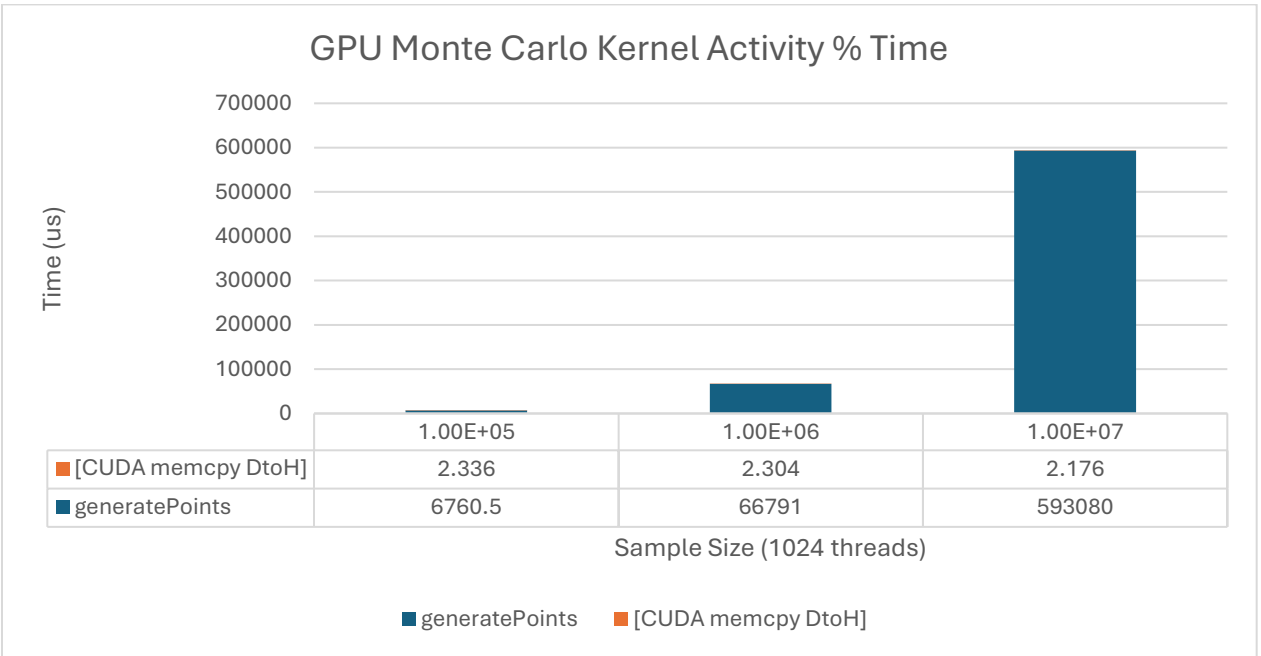
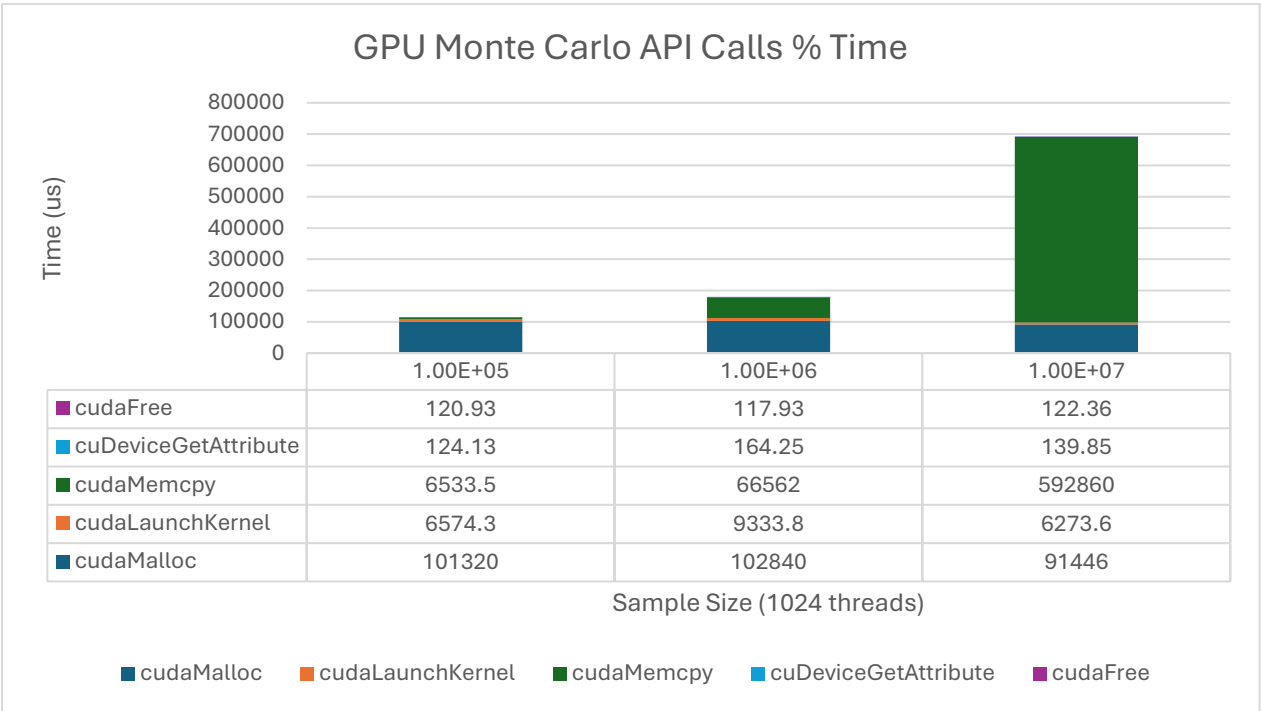
Caroline Jin
 CUDA Programming Assignment 1 (Part A and B)

SAXPY:

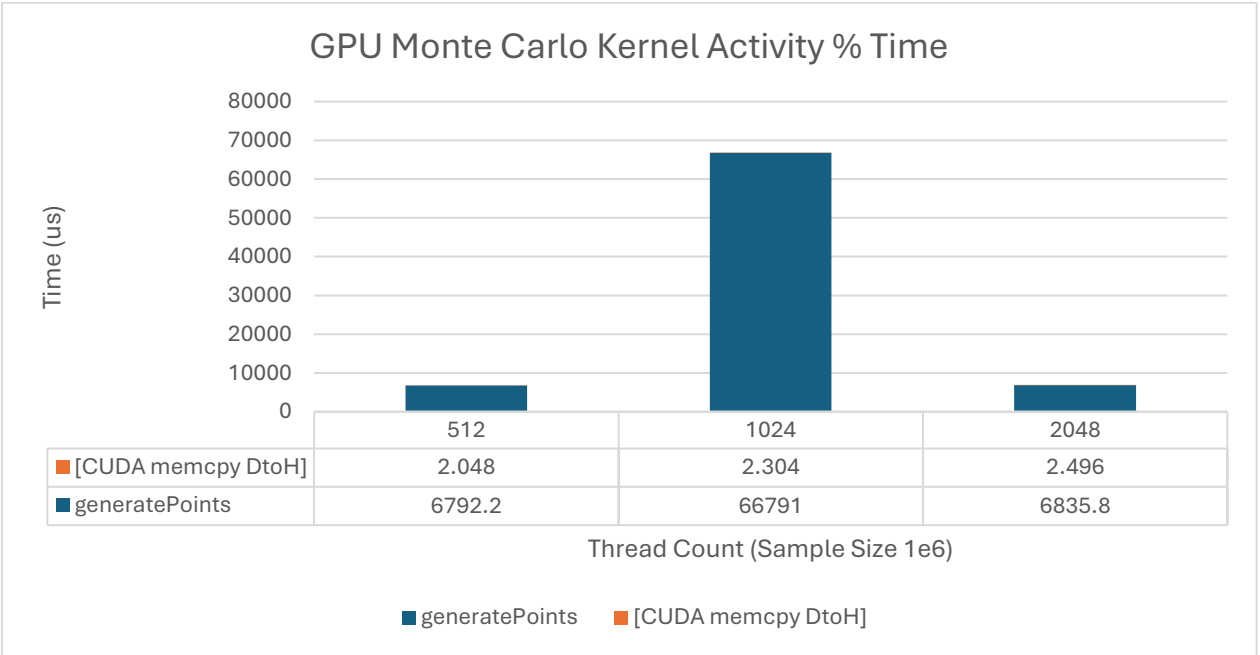
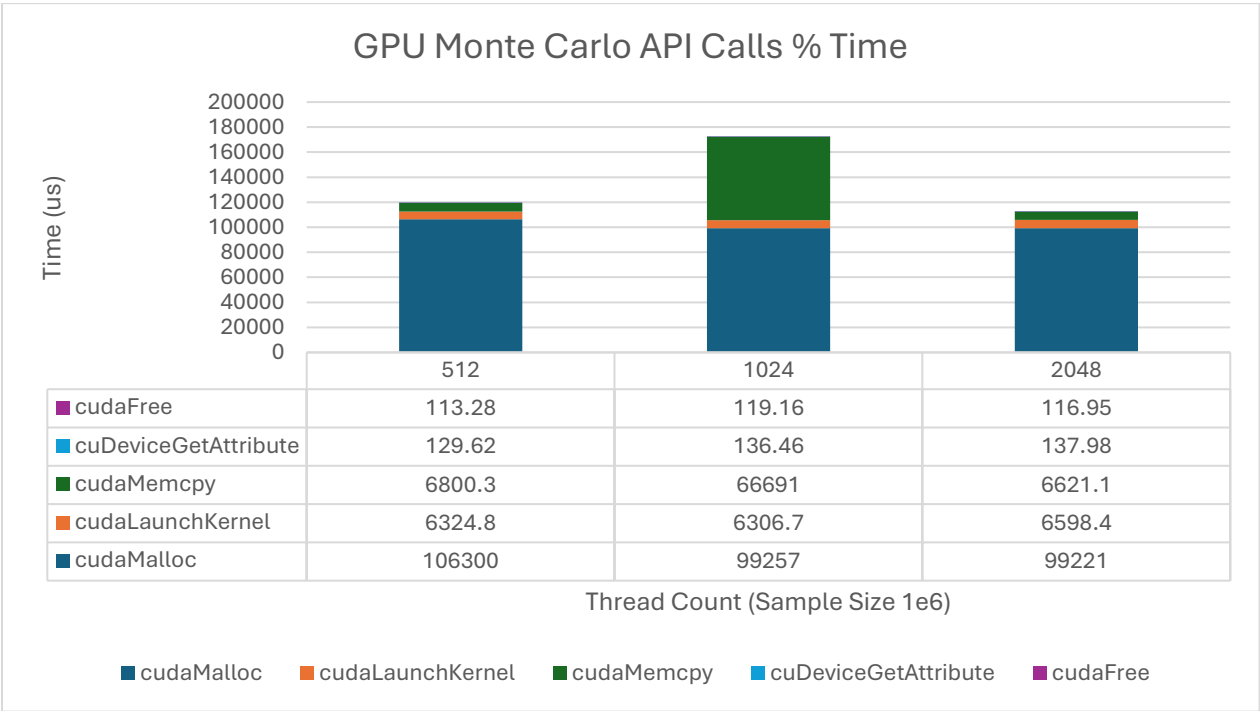


By varying the vector size for the GPU SAXPY code, there is a clear observation that an increased vector size greatly multiplies the time spent in the kernel. Additionally, a larger vector means the host code spends a greater amount of time copying data from the host to device (and vice versa) via cudaMemcpy.

Monte Carlo Pi Approximation:



By varying sample size (the amount of points each thread will generate), the pattern shows exponentially greater time spent in the kernel. In the host code, there was a larger amount of time spent in cudaMemcpy for larger sample size. In the code, data is only transferred one way from device to host. In this transaction, with a greater sample size, the number of hits would be a lot larger per thread, leading to a longer delay.



When varying the thread count, there is a spike in time spent in the kernel and `cudaMemcpy` at 1024 threads compared to the lower and higher number of threads. Based on the number of blocks and how the memory/scheduling is handled in the GPU possibly the higher number of thread count per block is not actually as efficient all the time.