

Juichi Lee

Email: leejuic@oregonstate.edu

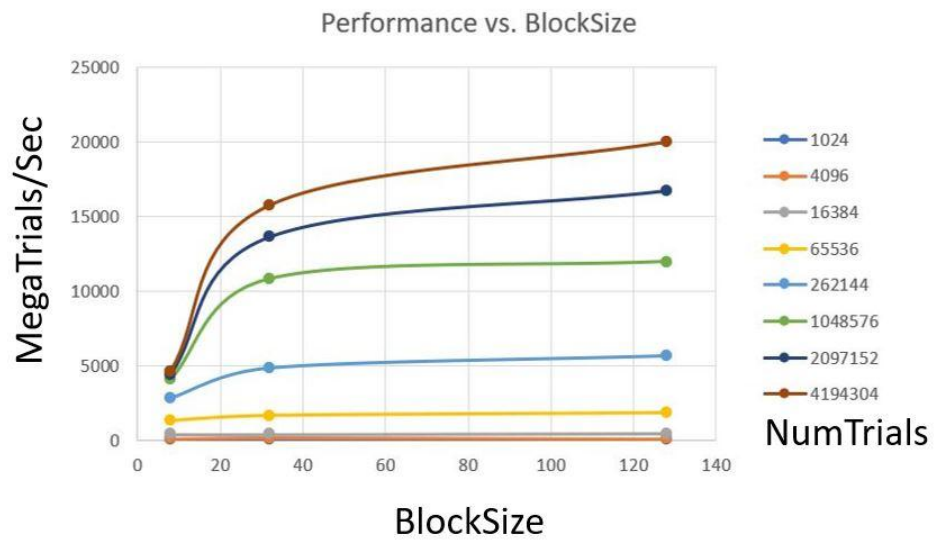
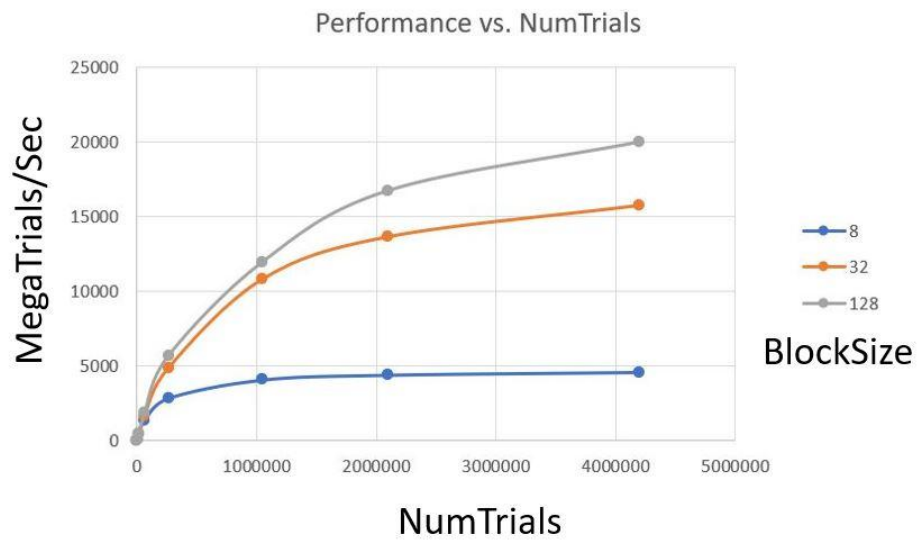
CS 475 Parallel Programming

5/18/2022

### Project #5: CUDA Monte Carlo Simulation

1. What machine you ran this on
  - I ran this on the DGX System.
2. Show the table and the two graphs

NumTrials	BlockSize	MegaTrialsPerSecond	Probability
1024	8	30.303	29.39%
1024	32	29.4118	26.86%
1024	128	30.303	30.47%
4096	8	117.6471	28.54%
4096	32	121.2121	28.76%
4096	128	117.6471	27.42%
16384	8	444.4444	29.08%
16384	32	444.4444	29.51%
16384	128	457.1428	29.09%
65536	8	1338.5621	29.26%
65536	32	1663.6881	28.95%
65536	128	1840.0719	28.93%
262144	8	2840.4992	29.02%
262144	32	4850.2071	29.21%
262144	128	5657.4587	29.08%
1048576	8	4085.2762	29.17%
1048576	32	10807.3877	29.08%
1048576	128	11959.1245	29.11%
2097152	8	4410.2287	29.04%
2097152	32	13641.9656	29.08%
2097152	128	16731.1727	29.09%
4194304	8	4597.4043	29.14%
4194304	32	15740.6032	29.09%
4194304	128	19986.5809	29.07%



3. What patterns are you seeing in the performance curves?
  - Both performance curves appear logarithmic, quickly increasing and then tapering off. The curves of Blocksize and NumTrials also appear to follow a pattern in which the larger Blocksize or NumTrials curve have larger MegaTrials/Sec values than smaller ones.
4. Why do you think the patterns look this way?
  - I think the patterns look this way because BlockSize has an increasing impact on the MegaTrials/Sec. The same could be said for NumTrials, as increasing NumTrials also somewhat increases MegaTrials/Sec.
5. Why is a BLOCKSIZE of 8 so much worse than the others?
  - I think one of the main reasons why the BlockSize of 8 performs much worse than the other sizes is because it is not a multiple of the warp size 32. Since a BlockSize of 8 only has 8 threads per block, there are potentially 24 other threads in a warp that are not being utilized. This results in lower efficiency and performance overall.
6. How do these performance results compare with what you got in Project #1?  
Why?
  - These performance results completely outperform the results I got in Project #1. The main reason is probably because I used the DGX System for this project and that the number of threads used in total in Project #5 was vastly greater than that in Project #1.
7. What does this mean for the proper use of GPU parallel computing?
  - I think this means that GPU parallel computing is best used for computations that require performing repetitive tasks on multiple streams of data that can be split up and worked upon by separate threads.