

Execution Time and Vector Size:

The execution time does not scale linearly with the size of the vectors. For example, increasing the vector size from 5 to 50,000 and then to 100,000 elements doesn't result in a proportional increase in execution time. This demonstrates the parallel processing power of the GPU, which can handle larger data sets efficiently.

Threads Per Block:

The number of threads per block varies from 2 to 1024 in these runs. There is not a consistent trend in execution time with the number of threads per block, which indicates that there may be other limiting factors in the computation, such as memory bandwidth or kernel launch overhead, which are not directly related to the number of threads.

Blocks Per Grid:

Similar to threads per block, the blocks per grid don't show a consistent trend with execution time. It's important to note that optimal values for threads per block and blocks per grid depend on the specific GPU architecture and the nature of the problem.

Kernel Overhead:

Even for a very small vector size (5 elements), there is an execution time of 0.01715 ms, which suggests that there is a certain amount of overhead associated with launching the kernel and the CUDA API calls, regardless of the workload size.

Data Transfer Overhead:

The execution times reported do not consider the time taken to transfer data between the host and the device. In practical applications, this time can be significant, especially for large vectors.

Efficiency Plateau:

There's a plateau in efficiency gains; after a certain point, merely increasing threads per block or blocks per grid does not guarantee better performance. For example, the execution time for 50,000 elements with different threads per block and blocks per grid configurations (256 threads/196 blocks vs. 1024 threads/49 blocks vs. 64 threads/780 blocks) is around the same (0.029723 ms, 0.030046 ms, and 0.030677 ms, respectively), suggesting that the GPU resources are being similarly utilized across these configurations.

Correctness of Results:

The correctness of the output is verified as the results are consistent across different runs.