

## CRUZ, AIRON JOHN R. CEPARCO-S11

### Project Specification

Question 1	100 pts
<p>Write the kernel in (1) C program (2) an x86-64 assembly language (3) x86 SIMD using YMM register. The kernel is to perform dot product between vector <i>A</i> and vector <i>B</i> and place the result in <i>sdot</i>.</p> <p>As bonus, (4) CUDA kernel as well. Use the best CUDA implementation.</p> <p><b>Input:</b> Scalar variable <i>n</i> (integer) contains the length of the vector; Vectors <i>A</i> and <i>B</i> are both <b>64-bit integer</b>.</p> <p><b>Input:</b> Memory location <i>n</i> (integer) contains the length of the vector; memory location <i>A</i> and memory location <i>B</i> (both <b>64-bit integer</b>) are vectors with size <i>n</i>.</p> <p><b>Process:</b> <math>sdot = \sum_{i=1}^n a_i b_i = a_1 b_1 + a_2 b_2 + \dots + a_n b_n</math></p> <p><b>Output:</b> store the result in memory location <i>sdot</i>. Display the result for all versions of kernel (i.e., C, x86-64, x86 SIMD using YMM register and CUDA if ever).</p>	
<p>Note:</p> <ol style="list-style-type: none"><li>1a.) Write a C main program to call the kernels of C version, non-SIMD x86-64 assembly language and SIMD YMM x86 assembly language version.</li><li>1b.) For the bonus CUDA program, C main program to call CUDA kernel.</li><li>2.) Time the kernel portion only. Include any overhead time as well (i.e., CUDA data transfer time, page faults etc.)</li><li>3.) For each kernel version, time the process for vector size <math>n = \{2^{20}, 2^{24}, \text{ and } 2^{30}\}</math>. If <math>2^{30}</math> is not possible, you may reduce it to <math>2^{29}</math> or at the least up to <math>2^{28}</math>.</li><li>4.) For each version, you need to run at least 30 times and get the average execution time.</li><li>5.) For the data, you may initialize each vector with the same or different random value.</li><li>6.) You will need to check the correctness of your output. Thus, if C version is your "sanity check answer key", then the output of the x86-64 version and SIMD version have to be checked with the C version and output correspondingly (i.e., x86-64 kernel output is correct, etc.).</li><li>7.) output in Github (make sure that I can access your Github):<ol style="list-style-type: none"><li>a.) Github readme containing the following:<ol style="list-style-type: none"><li>i.) comparative execution time as well as analysis of the performance of different kernels (how many times faster, why is it faster, overheads in calling the kernel, etc.)</li><li>ii.) screenshot of the program output with correctness check (C)</li><li>iii.) screenshot of the program output including correctness check (x86-64)</li><li>iv.) screenshot of the program output including correctness check (SIMD YMM register)</li><li>v.) screenshot of the program output including correctness check (CUDA, optional)</li></ol></li><li>b.) Visual studio project file folder containing <b>complete</b> files (source code: x86-64 non-SIMD, x86-64 SIMD version and all other required files) for others to load and execute your program.</li><li>c.) Link to your colab CUDA program. Make sure I can access it.</li></ol></li></ol>	