Final Take Home Test

Optimization of Matrix-Matrix Multiplication Using Vector Instructions

Zi Xuan Li

Professor Izidor Gertner

CSC 34200 Spring 2024

May 12th, 2024

Table of Contents

**Objective:**

        The aim of this concluding take-home exam is to enhance compiler-generated code for computing the product of two matrices using vector instructions and DPPS vector instruction. It builds upon the earlier test on matrix multiplication. The objective is to explore the performance contrast between vectorization and non-vectorization, assessed with Chrono, a high-resolution timer for execution time measurement.

**Task#1:** Use CPUID instructions to determine your processor vector processing capabilities.
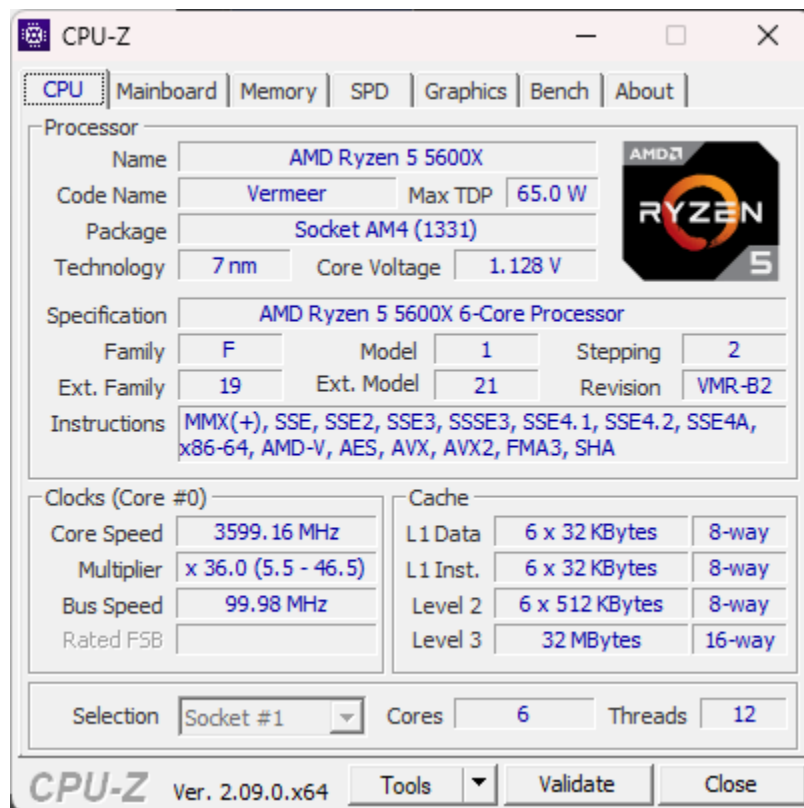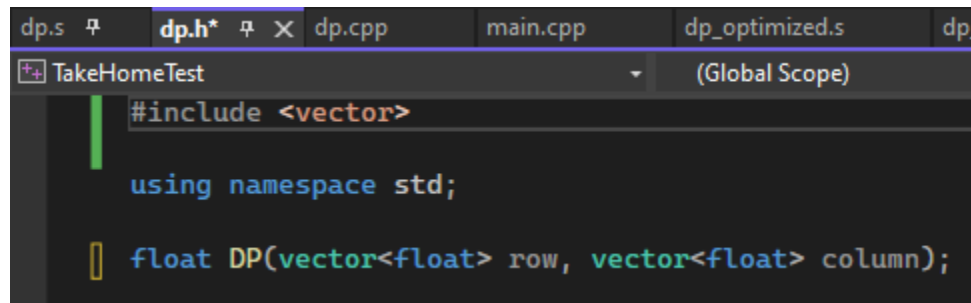


Figure 1: Information of my CPU using CPU-Z from the CPUID instructions

        The processing capabilities of my CPU, the AMD Ryzen 5 5600x, can be found under the instructions tab in the window above. It supports MMX (+), SSE, SSE2, SSE3, SSSE3, SSE4.1, SSE4.2, SSE4A, x86-64, AMD-V, AES, AVX, AVX2, FMA3, and SHA.
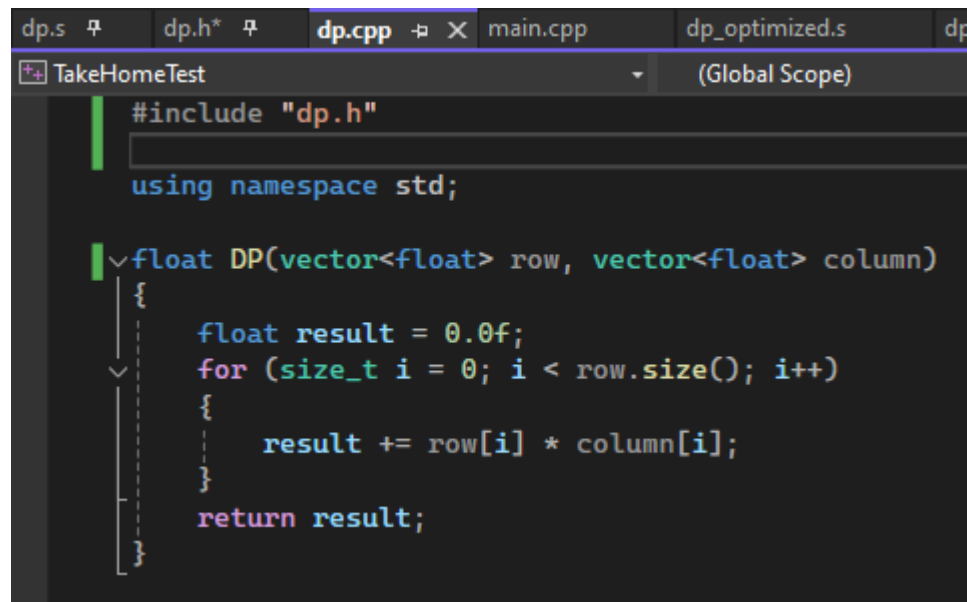
**Task#2:** Write C/C++ main () to compute Matrix-Matrix multiplication, can be taken from previous take-home test. The focus in this take home test is the most inner loop that computes dot product of row and column. You can use a function DP (row, column) from a previous take home test.

Place the function DP (row, column) in a separate file from main () that calls this function. Vector sizes should be powers of 2 (e.g. 16, 32, 64, ....512, ...216 etc.)



Figure 2: Image of the dp.h file



Figure 3: Image of the dp.cpp file

Figure 4: Image of the main.cpp file.

**Task#3:** Compile code in §2 and create assembly code for function DP (row, column) only. Make sure that compiler generated vectorized code.

```asm
	.file	"dp.cpp"
	.text
	.p2align 4
	.globl	_Z2DPSt6vectorIfSaIfEES1_
	.def	_Z2DPSt6vectorIfSaIfEES1_;	.scl	2;	.type	32; .endef
	.seh_proc	_Z2DPSt6vectorIfSaIfEES1_
_Z2DPSt6vectorIfSaIfEES1_:
.LFB1021:
	.seh_endprologue
	movq	8(%rcx), %r8
	movq	(%rcx), %rax
	movq	%r8, %r9
	subq	%rax, %r9
	movq	%r9, %rcx
	sarq	$2, %rcx
	cmpq	%rax, %r8
	je	.L10
	movq	(%rdx), %rdx
	cmpq	$28, %r9
	jbe	.L11
	movq	%rcx, %r9
	xorl	%r8d, %r8d
	vxorps	%xmm0, %xmm0, %xmm0
	shrq	$3, %r9
	salq	$5, %r9
	.p2align 4
	.p2align 3
.L4:
	vmovups	(%rax,%r8), %ymm4
	vmulps	(%rdx,%r8), %ymm4, %ymm1
	addq	$32, %r8
	vshufps	$85, %xmm1, %xmm1, %xmm3
	vshufps	$255, %xmm1, %xmm1, %xmm2
	vaddss	%xmm1, %xmm0, %xmm0
	vaddss	%xmm3, %xmm0, %xmm0
	vunpckhps	%xmm1, %xmm1, %xmm3
	vextractf128	$0x1, %ymm1, %xmm1
	vaddss	%xmm3, %xmm0, %xmm0
	vaddss	%xmm2, %xmm0, %xmm0
	vshufps	$85, %xmm1, %xmm1, %xmm2
	vaddss	%xmm1, %xmm0, %xmm0
	vaddss	%xmm2, %xmm0, %xmm0
	vunpckhps	%xmm1, %xmm1, %xmm2
	vshufps	$255, %xmm1, %xmm1, %xmm1
	vaddss	%xmm2, %xmm0, %xmm0
	vaddss	%xmm1, %xmm0, %xmm0
	cmpq	%r8, %r9
	jne	.L4
	movq	%rcx, %r8
	andq	$-8, %r8
	testb	$7, %cl
```

dp.h  |  dp.cpp  |  main.cpp*  |  **dp.s**  X  |  dp_optimized.s  |  dp_dpps.cpp

```asm
 51         testb   $7, %cl
 52         je    .L22
 53         vzeroupper
 54   .L3:
 55         movq    %rcx, %r9
 56         subq    %r8, %r9
 57         leaq    -1(%r9), %r10
 58         cmpq    $2, %r10
 59         jbe .L9
 60         vmovups (%rax,%r8,4), %xmm5
 61         movq    %r9, %r10
 62         vmulps  (%rdx,%r8,4), %xmm5, %xmm1
 63         andq    $-4, %r10
 64         addq    %r10, %r8
 65         andl    $3, %r9d
 66         vaddss  %xmm1, %xmm0, %xmm0
 67         vshufps $85, %xmm1, %xmm1, %xmm2
 68         vaddss  %xmm2, %xmm0, %xmm0
 69         vunpckhps   %xmm1, %xmm1, %xmm2
 70         vshufps $255, %xmm1, %xmm1, %xmm1
 71         vaddss  %xmm2, %xmm0, %xmm0
 72         vaddss  %xmm1, %xmm0, %xmm0
 73         je    .L1
 74         .p2align 4
 75         .p2align 3
 76   .L9:
 77         vmovss  (%rax,%r8,4), %xmm1
 78         vmulss  (%rdx,%r8,4), %xmm1, %xmm1
 79         incq    %r8
 80         vaddss  %xmm1, %xmm0, %xmm0
 81         cmpq    %rcx, %r8
 82         jb    .L9
 83   .L1:
 84         ret
 85         .p2align 4
 86         .p2align 3
 87   .L10:
 88         vxorps  %xmm0, %xmm0, %xmm0
 89         ret
 90         .p2align 4
 91         .p2align 3
 92   .L22:
 93         vzeroupper
 94         ret
 95   .L11:
 96         xorl    %r8d, %r8d
 97         vxorps  %xmm0, %xmm0, %xmm0
 98         jmp .L3
 99         .seh_endproc
100         .ident   "GCC: (Rev2, Built by MSYS2 project) 12.1.0"
101
```

Figure 5: Image of the compiler generated assembly code from dp.h

**Task#4:** Use high resolution timer to measure execution time (as in previous take-home test). Plot graph: time versus vector size.



Figure 6: Execution times of various 2^n integer matrix sizes from n = 16 to n = 1024
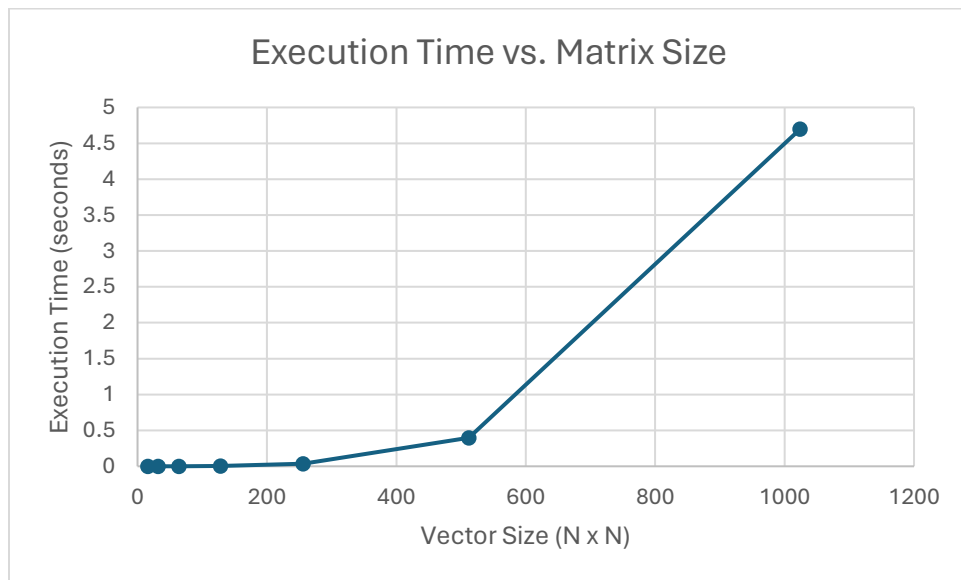


Figure 7: Graph measuring execution time vs. vector size for compiler generated code

**Task#5:** Create assembly code for function DP (row, column). The assembly code should contain vector instructions.

```
dp.h    дp.cpp    main.cpp    dp.s    dp_optimized.s  ×  dp_dpps.cpp

 1          .file    "dp.cpp"
 2          .text
 3          .p2align 4
 4          .globl  _Z2DPSt6vectorIfSaIfEES1_
 5          .def     _Z2DPSt6vectorIfSaIfEES1_;  .scl    2;  .type    32; .endef
 6          .seh_proc   _Z2DPSt6vectorIfSaIfEES1_
 7    _Z2DPSt6vectorIfSaIfEES1_:
 8  ∨.LFB1021:
 9          .seh_endprologue
10          movq    8(%rcx), %rax
11          movq    (%rcx), %rcx
12          movq    %rax, %r9
13          subq    %rcx, %r9
14          movq    %r9, %r8
15          sarq    $2, %r8
16          cmpq    %rcx, %rax
17          je   .L9
18          movq    (%rdx), %rdx
19          cmpq    $28, %r9
20          jbe .L10
21          movq    %r8, %r10
22          xorl    %eax, %eax
23          vxorps  %xmm0, %xmm0, %xmm0
24          shrq    $3, %r10
25          movq    %r10, %r11
26          salq    $5, %r11
27          andl    $1, %r10d
28          je   .L4
29          vmovups (%rcx), %ymm4
30          movl    $32, %eax
31          vmulps  (%rdx), %ymm4, %ymm1
32          vshufps $85, %xmm1, %xmm1, %xmm5
33          vunpckhps   %xmm1, %xmm1, %xmm3
34          vshufps $255, %xmm1, %xmm1, %xmm2
35          vaddss  %xmm1, %xmm0, %xmm0
36          vextractf128    $0x1, %ymm1, %xmm1
37          vaddss  %xmm5, %xmm0, %xmm4
38          vaddss  %xmm3, %xmm4, %xmm0
39          vaddss  %xmm2, %xmm0, %xmm5
40          vshufps $85, %xmm1, %xmm1, %xmm0
41          vunpckhps   %xmm1, %xmm1, %xmm2
42          vaddss  %xmm1, %xmm5, %xmm3
43          vshufps $255, %xmm1, %xmm1, %xmm1
44          vaddss  %xmm0, %xmm3, %xmm5
45          vaddss  %xmm2, %xmm5, %xmm4
46          vaddss  %xmm1, %xmm4, %xmm0
47          cmpq    $32, %r11
48          je   .L23
49          .p2align 4
50          .p2align 3
```

```
dp.h  ⌀    dp.cpp      main.cpp      dp.s      dp_optimized.s  ⌀ ✕  dp_dpps.cpp

50              .p2align 3
51      ∨ .L4:
52              vmovups (%rcx,%rax), %ymm3
53              vmulps  (%rdx,%rax), %ymm3, %ymm5
54              vshufps $85, %xmm5, %xmm5, %xmm1
55              vunpckhps   %xmm5, %xmm5, %xmm3
56              vshufps $255, %xmm5, %xmm5, %xmm2
57              vaddss  %xmm5, %xmm0, %xmm0
58              vextractf128    $0x1, %ymm5, %xmm5
59              vaddss  %xmm1, %xmm0, %xmm4
60              vmovups 32(%rcx,%rax), %ymm1
61              vaddss  %xmm3, %xmm4, %xmm0
62              vaddss  %xmm2, %xmm0, %xmm4
63              vshufps $85, %xmm5, %xmm5, %xmm0
64              vunpckhps   %xmm5, %xmm5, %xmm2
65              vaddss  %xmm5, %xmm4, %xmm3
66              vshufps $255, %xmm5, %xmm5, %xmm5
67              vaddss  %xmm0, %xmm3, %xmm4
68              vaddss  %xmm2, %xmm4, %xmm3
69              vmulps  32(%rdx,%rax), %ymm1, %ymm4
70              addq    $64, %rax
71              vaddss  %xmm5, %xmm3, %xmm0
72              vaddss  %xmm4, %xmm0, %xmm5
73              vshufps $85, %xmm4, %xmm4, %xmm0
74              vunpckhps   %xmm4, %xmm4, %xmm3
75              vshufps $255, %xmm4, %xmm4, %xmm2
76              vextractf128    $0x1, %ymm4, %xmm4
77              vaddss  %xmm0, %xmm5, %xmm1
78              vaddss  %xmm3, %xmm1, %xmm5
79              vaddss  %xmm2, %xmm5, %xmm0
80              vshufps $85, %xmm4, %xmm4, %xmm5
81              vunpckhps   %xmm4, %xmm4, %xmm2
82              vaddss  %xmm4, %xmm0, %xmm3
83              vshufps $255, %xmm4, %xmm4, %xmm4
84              vaddss  %xmm5, %xmm3, %xmm0
85              vaddss  %xmm2, %xmm0, %xmm3
86              vaddss  %xmm4, %xmm3, %xmm0
87              cmpq    %rax, %r11
88              jne .L4
89      ∨ .L23:
90              movq    %r8, %r11
91              andq    $-8, %r11
92              testb   $7, %r8b
93              je   .L26
94              vzeroupper
95      ∨ .L3:
96              movq    %r8, %r9
97              subq    %r11, %r9
98              leaq    -1(%r9), %r10
99              cmpq    $2, %r10
```

dp.h  |  dp.cpp  |  main.cpp  |  dp.s  |  **dp_optimized.s**  ✕  dp_dpps.cpp

```asm
100          jbe .L7
101          vmovups (%rcx,%r11,4), %xmm1
102          movq    %r9, %rax
103          vmulps  (%rdx,%r11,4), %xmm1, %xmm5
104          andq    $-4, %rax
105          addq    %rax, %r11
106          andl    $3, %r9d
107          vaddss  %xmm5, %xmm0, %xmm0
108          vshufps $85, %xmm5, %xmm5, %xmm3
109          vunpckhps   %xmm5, %xmm5, %xmm1
110          vshufps $255, %xmm5, %xmm5, %xmm5
111          vaddss  %xmm3, %xmm0, %xmm4
112          vaddss  %xmm1, %xmm4, %xmm2
113          vaddss  %xmm5, %xmm2, %xmm0
114          je  .L1
115      .L7:
116          vmovss  (%rcx,%r11,4), %xmm3
117          leaq    1(%r11), %r10
118          leaq    0(,%r11,4), %r9
119          vfmadd231ss (%rdx,%r11,4), %xmm3, %xmm0
120          cmpq    %r8, %r10
121          jnb .L1
122          vmovss  4(%rcx,%r9), %xmm4
123          addq    $2, %r11
124          vfmadd231ss 4(%rdx,%r9), %xmm4, %xmm0
125          cmpq    %r8, %r11
126          jnb .L1
127          vmovss  8(%rcx,%r9), %xmm1
128          vfmadd231ss 8(%rdx,%r9), %xmm1, %xmm0
129      .L1:
130          ret
131          .p2align 4
132          .p2align 3
133      .L9:
134          vxorps  %xmm0, %xmm0, %xmm0
135          ret
136          .p2align 4
137          .p2align 3
138      .L26:
139          vzeroupper
140          ret
141      .L10:
142          xorl    %r11d, %r11d
143          vxorps  %xmm0, %xmm0, %xmm0
144          jmp .L3
145          .seh_endproc
146          .ident  "GCC: (Rev2, Built by MSYS2 project) 12.1.0"
147
```

Figure 8: Image of the optimized assembly code of dp.h



```
Time to multiply two 16x16 float matrices: 9.12e-05 seconds
Time to multiply two 32x32 float matrices: 0.0002748 seconds
Time to multiply two 64x64 float matrices: 0.0009819 seconds
Time to multiply two 128x128 float matrices: 0.0056039 seconds
Time to multiply two 256x256 float matrices: 0.0373567 seconds
Time to multiply two 512x512 float matrices: 0.392527 seconds
Time to multiply two 1024x1024 float matrices: 4.4821 seconds
```
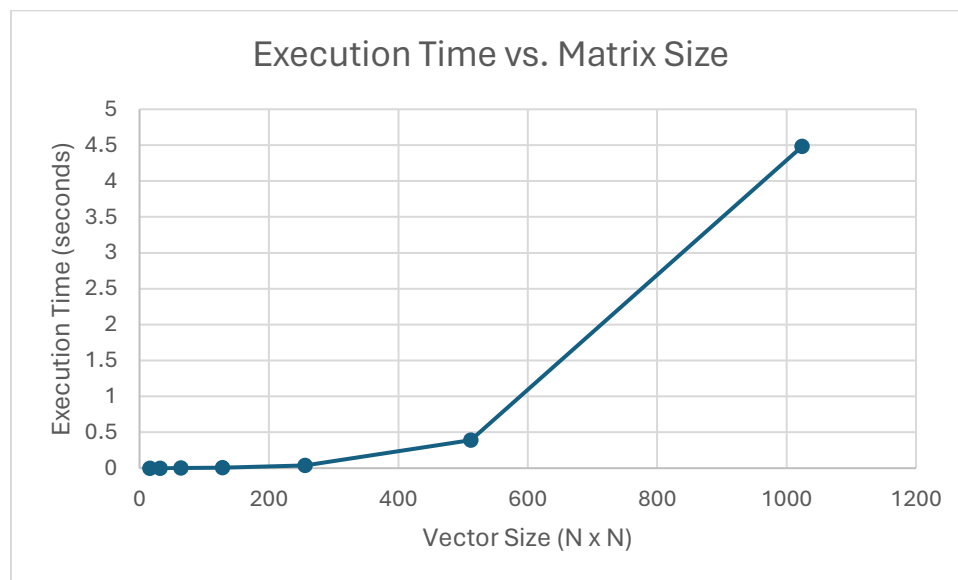
Figure 9: Execution times of various 2^n integer matrix sizes from n = 16 to n = 1024



Figure 10: Graph measuring execution time vs. vector size for optimized code

**Task#6:** To optimize the code further, please try to use machine vector instruction DPPS to compute dot product.

```cpp
#include "dp.h"
#include <immintrin.h>
#include <stdexcept>

using namespace std;

float DP(vector<float> row, vector<float> column)
{
    size_t size = row.size();
    if (size != column.size())
        throw std::invalid_argument("Vectors are not of the same size.");

    float result = 0.0f;
    size_t i = 0;

    __m256 sum_vec = _mm256_setzero_ps();
    for (; i <= size - 8; i += 8)
    {
        __m256 row_vec = _mm256_loadu_ps(&row[i]);
        __m256 col_vec = _mm256_loadu_ps(&column[i]);
        __m256 dp = _mm256_dp_ps(row_vec, col_vec, 0xF1);
        sum_vec = _mm256_add_ps(sum_vec, dp);
    }

    float temp[8];
    _mm256_storeu_ps(temp, sum_vec);
    for (int j = 0; j < 8; ++j)
    {
        result += temp[j];
    }

    for (; i < size; ++i)
    {
        result += row[i] * column[i];
    }
    return result;
}
```

Figure 11: Image of the dp_dpps.cpp file using machine vector instructions

```
Time to multiply two 16x16 float matrices: 9.08e-05 seconds
Time to multiply two 32x32 float matrices: 0.0002639 seconds
Time to multiply two 64x64 float matrices: 0.0007915 seconds
Time to multiply two 128x128 float matrices: 0.0045332 seconds
Time to multiply two 256x256 float matrices: 0.0294053 seconds
Time to multiply two 512x512 float matrices: 0.300061 seconds
Time to multiply two 1024x1024 float matrices: 3.83468 seconds
```

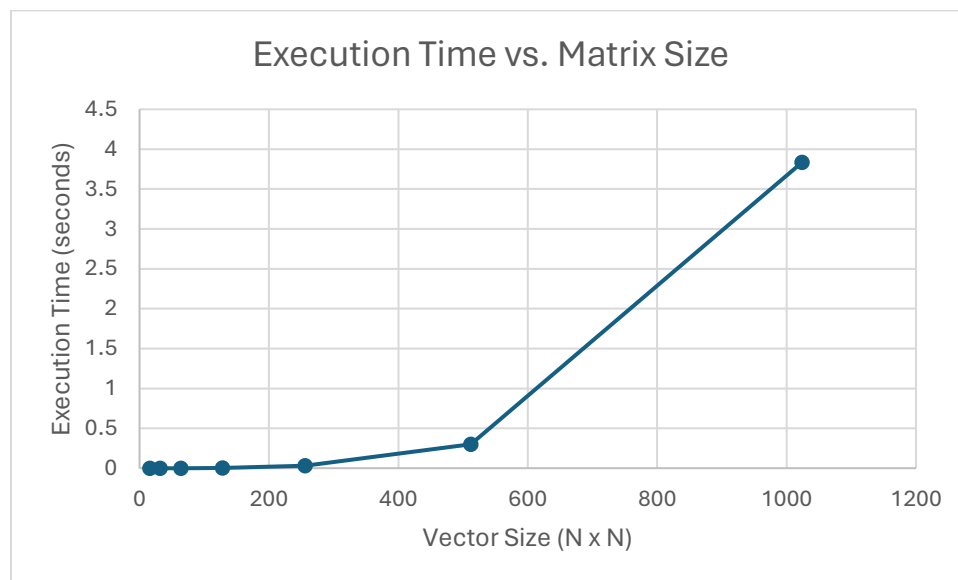Figure 12: Execution times of various 2^n integer matrix sizes from n = 16 to n = 1024.



Figure 13: Graph measuring execution time vs. vector size using machine vector instruction DPPS code

**Task#7:** Compare all plots in one figure. Compare also to the performance plots from the previous take home test.
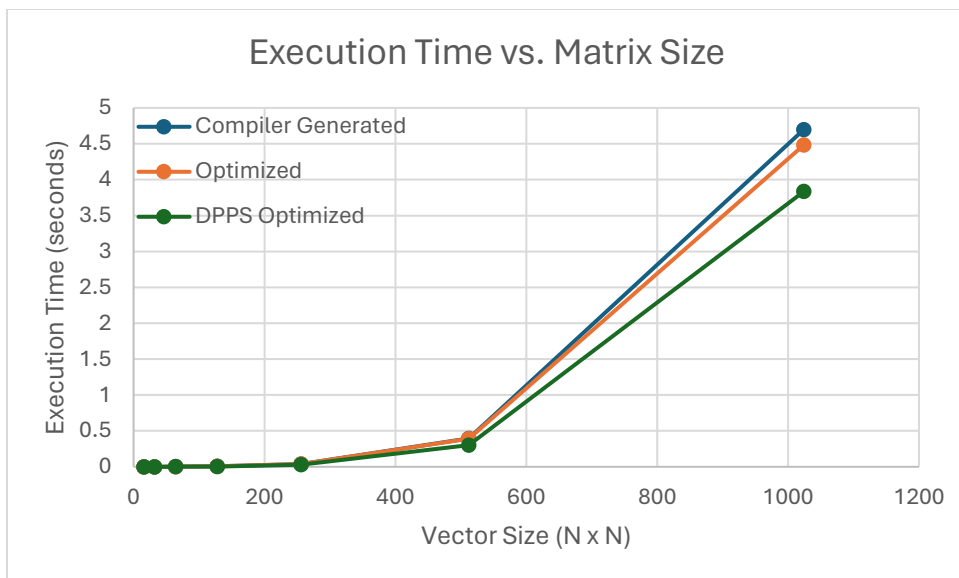
Figure 14: Graph comparing execution time vs. matrix size of compiler generated, optimized, and DPPS optimized assembly code
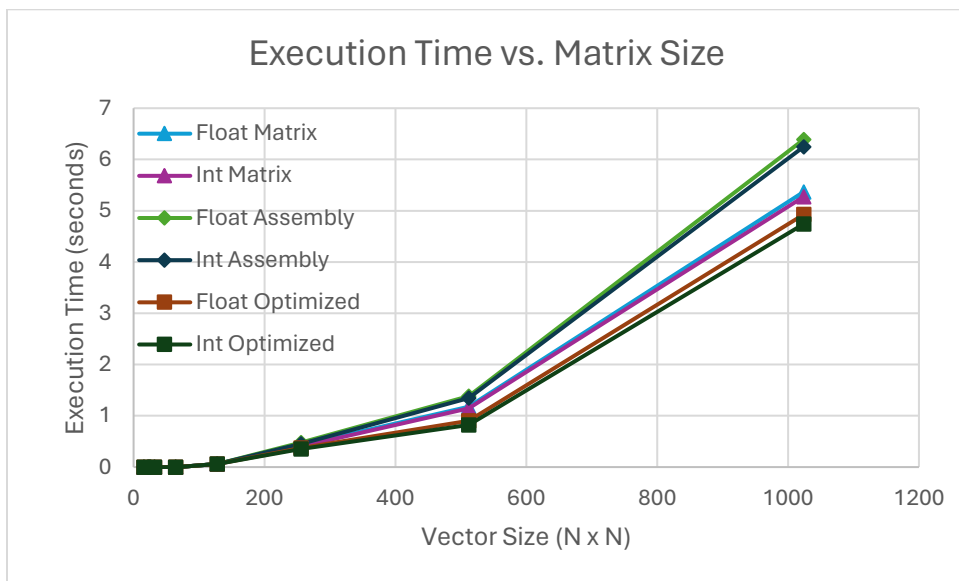


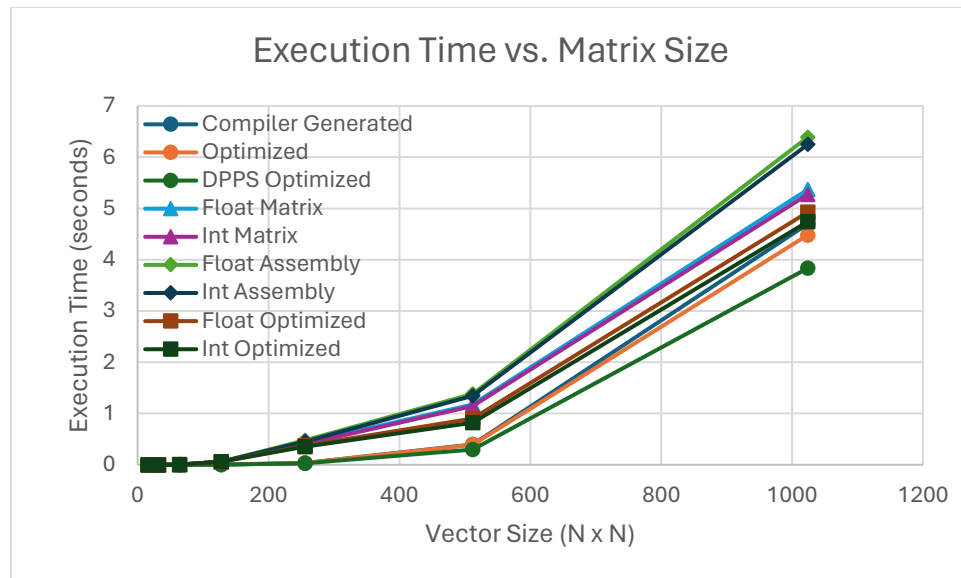Figure 15: Graph comparing execution time vs. matrix size from previous take home test

Figure 16: Graph comparing execution time vs. matrix size from the previous take home test and current take home test

**Conclusion:**

The primary goal of this take-home test was to employ vector instructions to enhance dot product computations in the dp.cpp file. Utilizing vectorization commands in g++, along with DPPS instructions, I compiled and executed the files and observed that the runtime improved with increasing optimization levels. The slowest runtime was without any optimization, followed by optimization with assembly, and finally optimization with DPPS.