

ELEC 374 – Machine Problem 1

"I do hereby verify that this machine problem submission is my own work and contains my own original ideas, concepts, and designs. No portion of this report or code has been copied in whole or in part from another source, with the possible exception of properly referenced material"

Part 1

```
≡ Outputs.txt
1  Part 1:
2  Device 0: NVIDIA T600
3  Compute Capability: 7.5
4  Clock Rate: 1335.00 MHz
5  Number of Streaming Multiprocessors (SMs): 10
6  Total CUDA Cores (approx.): 640
7  Warp Size: 32
8  Global Memory: 4.00 GB
9  Constant Memory: 64.00 KB
10 Shared Memory per Block: 48.00 KB
11 Registers per Block: 65536
12 Max Threads per Block: 1024
13 Max Block Dimensions: (1024, 1024, 64)
14 Max Grid Dimensions: (2147483647, 65535, 65535)
```

Figure 1: Output data of Part 1 script showing All GPUs and their specifications on the machine

The values above were obtained by running the MP1_Part1.cu code. Some of values above were verified using the data sheet at the following [link \(leads to a pdf\)](#). To obtain the total number of CUDA cores the major and minor number were obtained using the attributes in the deviceProp. That had to be cross-referenced with information obtained online to determine the number of CUDA cores for any given GPU. The program cannot determine the CUDA cores for the new Blackwell architecture as the reference I was using was not up to date.

Part 2

2.1

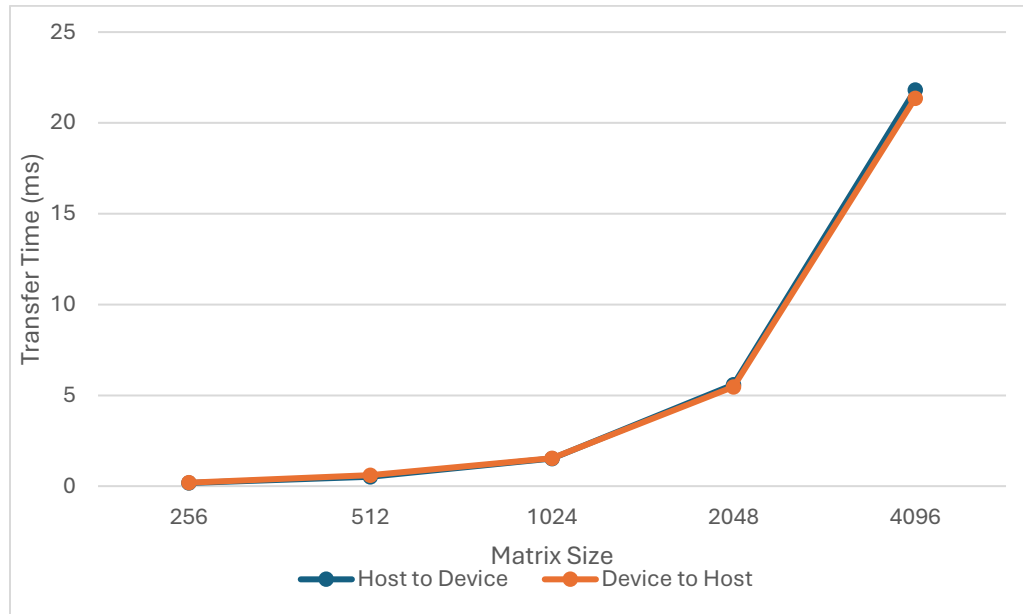


Figure 2: The Effects of Matrix Size on Transfer time (ms) from Host to Device and back

As the graph above illustrates, the Host to Device (H2D) transfer time is relatively the same as Device to Host (D2H). However, upon close inspection H2D seems to be faster for the smaller matrix sizes, and the D2H seems to be faster for larger matrix sizes, though this difference may be insignificant. Seeing no significant difference between H2D and D2H makes sense since the PCIe has equal speeds in either direction.

2.2

Is it always beneficial to offload your matrix multiplication to the device?

No, it is not always better to offload matrix multiplication to the GPU. If the transfer times from H2D and D2H plus the execution time on the GPU is longer than the CPU execution time, then it's best to simply run the calculation on the CPU. However, this is only true for smaller matrices, for even medium sized matrices it is better to run the calculation on the GPU which can be seen from the output. The 256x256 matrix for the output of Part 2.3 takes only 1.75 ms to complete and the transfer time is only 0.17 ms, whereas it takes the CPU 35 ms, as can be seen in the output of Part 2.2. And this holds true as matrix sizes increase, since the GPU is made for massively parallel tasks

2.3

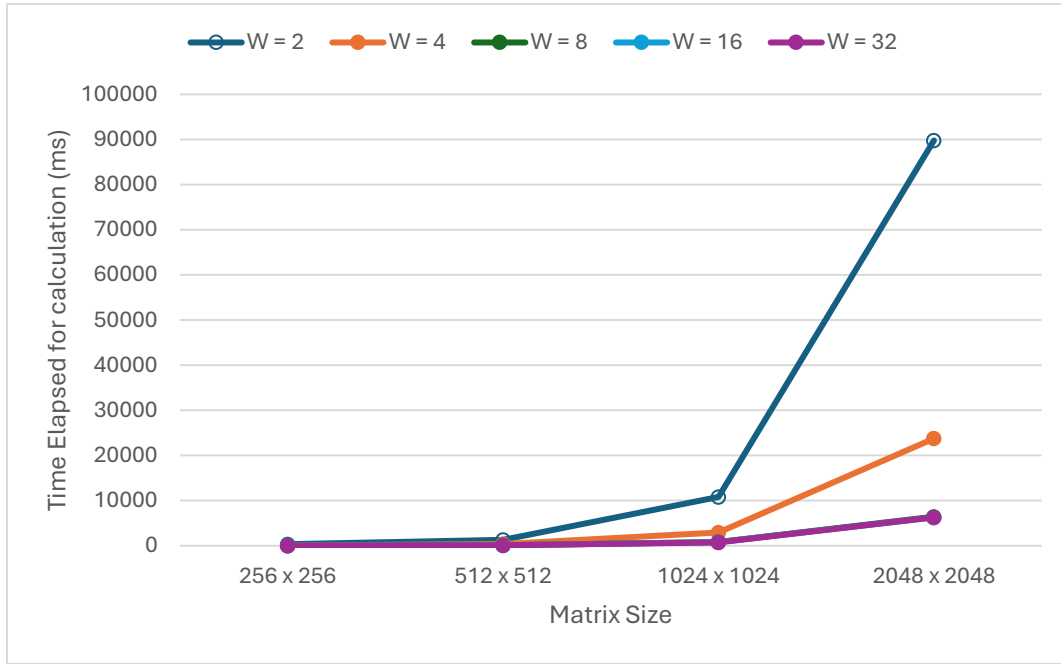


Figure 3: The effect of varying number of blocks/block-width on execution time of Matrices of varying sizes

a) How many times is each element of each input matrix loaded during the execution of the kernel?

In normal matrix multiplication each element of each input matrix should only be loaded once during the execution of the kernel. Since each element of matrix N is loaded once per column in M and each element of matrix M is loaded once per row in N, to complete the matrix multiplication using the naïve approach.

b) What is the floating-point computation to global memory access (CGMA) ratio in each thread? Consider multiply and addition as separate operations and ignore the global memory store at the end. Only count global memory loads towards your off-chip bandwidth.

The equation below describes the formula for the CGMA ratio:

$$CGMA\ ratio = \frac{FLOP}{Total\ \#\ of\ Memory\ Accesses}$$

For our purposes, each element is only read once, thus there are 2 memory accesses, one for each element in any given atomic calculation. In other words, the element from N is selected and an element from M is selected. $Total\ \#\ of\ Memory\ Accesses = 2S$, where S is the size of the matrix. Similarly, there are two floating point operations done, a multiplication and an addition, and thus this leads to $FLOP = 2S$. Plugging in these values leads to a CGMA ratio of 1.

$$CGMA\ ratio = \frac{2S}{2S} = 1$$

Appendix

Images below show Part 1 code, Part 2 code, and full output, in that order (Hint: Check label at top)

```
1 #include <cuda_runtime.h>
2 #include <stdio.h>
3 // https://docs.nvidia.com/cuda/cuda-runtime-api/structcudaDeviceProp.html
4
5 int getUDACores(cudaDeviceProp p) {
6     int coresPerSM = 0;
7
8     // Determine the number of cores per SM based on compute capability information was found at the link below
9     // https://developer.nvidia.com/cuda-gpus
10    if (p.major == 2) {
11        if (p.minor == 1) coresPerSM = 48;
12        else coresPerSM = 32;
13    }
14    else if (p.major == 3) {
15        coresPerSM = 192;
16    }
17    else if (p.major == 5) {
18        coresPerSM = 128;
19    }
20    else if (p.major == 6) {
21        if (p.minor == 1 || p.minor == 2) coresPerSM = 128;
22        else coresPerSM = 64;
23    }
24    else if (p.major == 7) {
25        coresPerSM = 64;
26    }
27    else if (p.major == 8) {
28        if (p.minor == 0) coresPerSM = 64;
29        else coresPerSM = 128;
30    }
31    else if (p.major == 9) {
32        coresPerSM = 128;
33    }
34    // not including blackwell GPU cause they are too
35    else {
36        printf("Unknown CUDA architecture: Compute Capability %d.%d\n", p.major, p.minor);
37        return 0;
38    }
39
40    return p.multiProcessorCount * coresPerSM;
41 }
42
43 int main() {
44     int nd;
45     cudaGetDeviceCount(&nd);
46     printf("Number of CUDA Devices: %d\n", nd);
47
48     if (nd == 0) {
49         printf("No CUDA devices found. Exiting...\n");
50         return 1;
51     }
52
53     for (int i = 0; i < nd; i++) {
54         cudaDeviceProp p;
55         cudaGetDeviceProperties(&p, i);
56
57         printf("\nDevice %d: %s\n", i, p.name);
58         printf("Compute Capability: %d.%d\n", p.major, p.minor);
59         printf("Clock Rate: %.2f MHz\n", p.clockRate / 1000.0f);
60         printf("Number of Streaming Multiprocessors (SMs): %d\n", p.multiProcessorCount);
61
62         int cudaCores = getUDACores(p);
63         if (cudaCores > 0)
64             printf("Total CUDA Cores (approx.): %d\n", cudaCores);
65         else
66             printf("Could not determine CUDA cores for this device.\n");
67
68         printf("Warp Size: %d\n", p.warpSize);
69         printf("Global Memory: %.2f GB\n", p.totalGlobalMem / (1024.0 * 1024.0 * 1024.0));
70         printf("Constant Memory: %.2f KB\n", p.totalConstMem / 1024.0);
71         printf("Shared Memory per Block: %.2f KB\n", p.sharedMemPerBlock / 1024.0);
72         printf("Registers per Block: %d\n", p.regsPerBlock);
73         printf("Max Threads per Block: %d\n", p.maxThreadsPerBlock);
74         printf("Max Block Dimensions: (%d, %d, %d)\n",
75                p.maxThreadsDim[0], p.maxThreadsDim[1], p.maxThreadsDim[2]);
76         printf("Max Grid Dimensions: (%d, %d, %d)\n",
77                p.maxGridSize[0], p.maxGridSize[1], p.maxGridSize[2]);
78     }
79
80     return 0;
81 }
```

```

1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <cuda_runtime.h>
4 #include <device_launch_parameters.h> // forgot to add earlier
5 #include <math.h>
6 #include <time.h>
7
8 // Kernel for matrix multiplication on the GPU
9 __global__ void matMulKernel(const float* M, const float* N, float* P, int width) {
10     // Compute row and column indices for this thread
11     int row = blockIdx.y * blockDim.y + threadIdx.y;
12     int col = blockIdx.x * blockDim.x + threadIdx.x;
13
14     // Out of Bounds check
15     if (row < width && col < width) {
16         float sum = 0.0f;
17         // Perform multiplication
18         for (int k = 0; k < width; k++) {
19             sum += M[row * width + k] * N[k * width + col];
20         }
21         // Store the result
22         P[row * width + col] = sum;
23     }
24 }
25
26 // Single-thread kernel
27 __global__ void matMulKernelSingleThread(const float* M, const float* N, float* P, int width) {
28     if (threadIdx.x == 0 && blockIdx.x == 0 && threadIdx.y == 0 && blockIdx.y == 0) {
29         for (int row = 0; row < width; row++) {
30             for (int col = 0; col < width; col++) {
31                 float sum = 0.0f;
32                 for (int k = 0; k < width; k++) {
33                     sum += M[row * width + k] * N[k * width + col];
34                 }
35                 P[row * width + col] = sum;
36             }
37         }
38     }
39 }
40
41 // Function for matrix multiplication on the CPU
42 void matMulCPU(const float* M, const float* N, float* P, int width) {
43     for (int i = 0; i < width; i++) {
44         for (int j = 0; j < width; j++) {
45             float sum = 0.0f;
46             // Multiply row i of M by column j of N
47             for (int k = 0; k < width; k++) {
48                 sum += M[i * width + k] * N[k * width + j];
49             }
50             // Store the result in P
51             P[i * width + j] = sum;
52         }
53     }
54 }
55
56 // Comparison function to check correct answer
57 bool compareArrays(const float* A, const float* B, int size, float tolerance) {
58     for (int i = 0; i < size; i++) {
59         // Return false if difference > tolerance
60         if (fabs(A[i] - B[i]) > tolerance) {
61             return false;
62         }
63     }
64     return true;
65 }
66
67 // Kernel execution time (data transfer not included)
68 void measureKernelTime(const float* dM, const float* dN, float* dP,
69     int width, int blockDim, float &kernelTimeMs) {
70     // Create the dim3 block and grid
71     dim3 block(blockDim, blockDim);
72     dim3 grid((width + blockDim - 1) / blockDim, (width + blockDim - 1) / blockDim);
73
74     // Create CUDA events to record time taken
75     cudaEvent_t start, stop;
76     cudaEventCreate(&start);
77     cudaEventCreate(&stop);
78
79     // Record the start event, launch kernel
80     cudaEventRecord(start);
81     matMulKernel<<<grid, block>>>>(dM, dN, dP, width);
82     cudaEventRecord(stop);
83
84     // Synchronize to make sure kernel finished, measure time
85     // Was causing error (may not be needed)
86     cudaEventSynchronize(stop);
87
88     // Calculate final value
89     kernelTimeMs = cudaEventElapsedTime(&kernelTimeMs, start, stop);
90
91     // Free memory
92     cudaEventDestroy(start);
93     cudaEventDestroy(stop);
94 }

```

```

96 int main()
97
98 // Part 2.1: Transfer Times
99 const int sizes1[] = {256, 512, 1024, 2048, 4096};
100 const int numSizes1 = sizeof(sizes1) / sizeof(int);
101
102 printf("---- Part 2.1: H to D and D to H Transfer Times ----\n");
103 printf("Matrix Sizes: 256, 512, 1024, 2048, 4096\n\n");
104
105 // Arrays to store transfer times
106 float hToDtimes[numSizes1];
107 float dToHtimes[numSizes1];
108
109 // Iterate through varying sizes
110 for (int idx = 0; idx < numSizes1; idx++) {
111     int width = sizes1[idx];
112     size_t bytes = width * (size_t)width * sizeof(float);
113
114     // Allocate host memory
115     float* hM = (float*)malloc(bytes);
116     float* hN = (float*)malloc(bytes);
117     float* hP = (float*)malloc(bytes);
118
119     // Initialize data
120     srand(0);
121     for (int i = 0; i < width * width; i++) {
122         hM[i] = (float)(rand() % 10);
123         hN[i] = (float)(rand() % 10);
124     }
125
126     // Allocate device memory
127     float* dM; float* dN; float* dP;
128     cudaMalloc((void**)&dM, bytes);
129     cudaMalloc((void**)&dN, bytes);
130     cudaMalloc((void**)&dP, bytes);
131
132     // Measure Host to Device transfer time
133     cudaEvent_t startHtoD, stopHtoD;
134     cudaEventCreate(&startHtoD);
135     cudaEventCreate(&stopHtoD);
136     cudaEventRecord(startHtoD);
137     cudaMemcpy(dM, hM, bytes, cudaMemcpyHostToDevice);
138     cudaMemcpy(dN, hN, bytes, cudaMemcpyHostToDevice);
139     cudaEventRecord(stopHtoD);
140     cudaEventSynchronize(stopHtoD); // for error
141
142     float timeHtoD = 0.0f;
143     cudaEventElapsedTime(&timeHtoD, startHtoD, stopHtoD);
144     hToDtimes[idx] = timeHtoD;
145
146     // Measure Device to Host transfer time
147     // Typically you'd measure for P but using N and M for simplicity
148     cudaEvent_t startDtoH, stopDtoH;
149     cudaEventCreate(&startDtoH);
150     cudaEventCreate(&stopDtoH);
151     cudaEventRecord(startDtoH);
152     cudaMemcpy(hP, dM, bytes, cudaMemcpyDeviceToHost);
153     cudaMemcpy(hP, dN, bytes, cudaMemcpyDeviceToHost);
154     cudaEventRecord(stopDtoH);
155     cudaEventSynchronize(stopDtoH);
156
157     float timeDtoH = 0.0f;
158     cudaEventElapsedTime(&timeDtoH, startDtoH, stopDtoH);
159     dToHtimes[idx] = timeDtoH;
160
161     // Free memory
162     cudaFree(dM);
163     cudaFree(dN);
164     cudaFree(dP);
165     free(hM);
166     free(hN);
167     free(hP);
168     cudaEventDestroy(startHtoD);
169     cudaEventDestroy(stopHtoD);
170     cudaEventDestroy(startDtoH);
171     cudaEventDestroy(stopDtoH);
172 }
173
174 // Print results for Part 2.1
175 printf("Host to Device Transfer Times (ms) by Matrix Size:\n");
176 for (int i = 0; i < numSizes1; i++) {
177     printf(" Size %d x %d : %f ms\n", sizes1[i], sizes1[i], hToDtimes[i]);
178 }
179 printf("\nDevice to Host Transfer Times (ms) by Matrix Size:\n");
180 for (int i = 0; i < numSizes1; i++) {
181     printf(" Size %d x %d : %f ms\n", sizes1[i], sizes1[i], dToHtimes[i]);
182 }
183 printf("\n");
184
185 // Part 2.2
186 const int sizes2[] = {256, 512, 1024};
187 const int numSizes2 = sizeof(sizes2) / sizeof(int);
188
189 printf("---- Part 2.2 : CPU vs. GPU (Single Thread) ----\n");
190 printf("Matrix Sizes: 256, 512, 1024\n\n");
191
192 for (int idx = 0; idx < numSizes2; idx++) {
193     int width = sizes2[idx];
194     size_t bytes = width * (size_t)width * sizeof(float);
195
196     float* hM = (float*)malloc(bytes);
197     float* hN = (float*)malloc(bytes);
198     float* hP = (float*)malloc(bytes);
199     float* hRef = (float*)malloc(bytes);
200

```

```

MP1_Part2.cu > ...
96  int main()
192  for (int idx = 0; idx < numSizes2; idx++) {
200
201      srand(0);
202      for (int i = 0; i < width * width; i++) {
203          hM[i] = (float)(rand() % 10);
204          hN[i] = (float)(rand() % 10);
205      }
206
207      float *dM, *dN, *dP;
208      cudaMalloc((void**)&dM, bytes);
209      cudaMalloc((void**)&dN, bytes);
210      cudaMalloc((void**)&dP, bytes);
211
212      cudaEvent_t startHtoD, stopHtoD;
213      cudaEventCreate(&startHtoD);
214      cudaEventCreate(&stopHtoD);
215      cudaEventRecord(startHtoD);
216      cudaMemcpy(dM, hM, bytes, cudaMemcpyHostToDevice);
217      cudaMemcpy(dN, hN, bytes, cudaMemcpyHostToDevice);
218      cudaEventRecord(stopHtoD);
219      cudaEventSynchronize(stopHtoD);
220
221      float hToD_ms = 0.0f;
222      cudaEventElapsedTime(&hToD_ms, startHtoD, stopHtoD);
223
224      // Assign only 1 block with 1 thread as instructed
225      dim3 block(1, 1);
226      dim3 grid(1, 1);
227
228      cudaEvent_t startKernel, stopKernel;
229      cudaEventCreate(&startKernel);
230      cudaEventCreate(&stopKernel);
231
232      // GPU
233      cudaEventRecord(startKernel);
234      matMulKernelSingleThread<<<grid, block>>>>(dM, dN, dP, width);
235      cudaEventRecord(stopKernel);
236      cudaEventSynchronize(stopKernel);
237
238      float gpuKernel_ms = 0.0f;
239      cudaEventElapsedTime(&gpuKernel_ms, startKernel, stopKernel);
240
241      cudaEvent_t startDtoH, stopDtoH;
242      cudaEventCreate(&startDtoH);
243      cudaEventCreate(&stopDtoH);
244
245      cudaEventRecord(startDtoH);
246      cudaMemcpy(hP, dP, bytes, cudaMemcpyDeviceToHost);
247      cudaEventRecord(stopDtoH);
248      cudaEventSynchronize(stopDtoH);
249
250      float dToH_ms = 0.0f;
251      cudaEventElapsedTime(&dToH_ms, startDtoH, stopDtoH);
252
253      // CPU time
254      clock_t cpuStart = clock();
255      matMulCPU(hM, hN, hRef, width);
256      clock_t cpuEnd = clock();
257      float cpu_ms = 1000.0f * (float)(cpuEnd - cpuStart) / CLOCKS_PER_SEC;
258
259      // Getting values for Ignoring transfer vs. including transfer
260      float gpuTotalNoTransfer = gpuKernel_ms;           // GPU only kernel
261      float gpuTotalWithTransfer = gpuKernel_ms + hToD_ms + dToH_ms;
262
263      // Compare correctness tolerance chosen manually
264      bool pass = compareArrays(hRef, hP, width * width, 1e-3f);
265
266      printf("Matrix Size %d x %d\n", width, width);
267      printf("  CPU Time (ms)                : %f\n", cpu_ms);
268      printf("  GPU Time (1 block, 1 thread) (ms) : %f (NO Transfer), %f (WITH Transfer)\n",
269             gpuTotalNoTransfer, gpuTotalWithTransfer);
270      printf("  Transfer Times: H to D = %f ms, D to H = %f ms\n", hToD_ms, dToH_ms);
271      printf("  %s\n", pass ? "Test PASSED" : "Test FAILED");
272
273      // Clean up
274      cudaFree(dM);
275      cudaFree(dN);
276      cudaFree(dP);
277      free(hM);
278      free(hN);
279      free(hP);
280      free(hRef);
281      cudaEventDestroy(startHtoD);
282      cudaEventDestroy(stopHtoD);
283      cudaEventDestroy(startKernel);
284      cudaEventDestroy(stopKernel);
285      cudaEventDestroy(startDtoH);
286      cudaEventDestroy(stopDtoH);
287  }
288
289  // Part 3.3: Vary block width and only measure kernel times
290  printf("--- Part 3.3: Kernel Times with varying Matrix Size & Block Width---\n");
291  printf("MatrixSizes: 256, 512, 1024, 2048, 4096 | BlockWidth: 2,4,8,16,32\n\n");
292
293  int blockWidths[5] = {2, 4, 8, 16, 32};
294  int sizes3[] = {256, 512, 1024, 2048, 4096};
295  int numSizes3 = 5;
296

```

```

MP1_Part2.cu > main()
96  int main()
295      int numSizes3 = 5;
296
297      // Print a table of times for each (matrixSize, blockWidth) using for loop for ease
298      for (int bwIdx = 0; bwIdx < 5; bwIdx++) {
299          int bWidth = blockWidths[bwIdx];
300          printf("BlockWidth = %d\n", bWidth);
301          for (int sIdx = 0; sIdx < numSizes3; sIdx++) {
302              int width = sizes3[sIdx];
303              size_t bytes = width * (size_t)width * sizeof(float);
304
305              float* hM = (float*)malloc(bytes);
306              float* hN = (float*)malloc(bytes);
307              float* hP = (float*)malloc(bytes);
308
309              srand(0);
310              for (int i = 0; i < width * width; i++) {
311                  hM[i] = (float)(rand() % 10);
312                  hN[i] = (float)(rand() % 10);
313              }
314
315              float *dM, *dN, *dP;
316              cudaMalloc((void**)&dM, bytes);
317              cudaMalloc((void**)&dN, bytes);
318              cudaMalloc((void**)&dP, bytes);
319
320              // Transfer time is ignored
321              cudaMemcpy(dM, hM, bytes, cudaMemcpyHostToDevice);
322              cudaMemcpy(dN, hN, bytes, cudaMemcpyHostToDevice);
323
324              // Measure kernel time
325              float kernelMs = 0.0f;
326              measureKernelTime(dM, dN, dP, width, bWidth, kernelMs);
327
328              printf("  Size %d x %d -> Kernel Time = %f ms\n", width, width, kernelMs);
329
330              // Free memory
331              cudaFree(dM);
332              cudaFree(dN);
333              cudaFree(dP);
334              free(hM);
335              free(hN);
336              free(hP);
337          }
338          printf("\n");
339      }
340
341      return 0;
342

```



```

E Outputs.txt
1  Part 1:
2  Device 0: NVIDIA T600
3  Compute Capability: 7.5
4  Clock Rate: 1335.00 MHz
5  Number of Streaming Multiprocessors (SMs): 10
6  Total CUDA Cores (approx.): 640
7  Warp Size: 32
8  Global Memory: 4.00 GB
9  Constant Memory: 64.00 KB
10 Shared Memory per Block: 48.00 KB
11 Registers per Block: 65536
12 Max Threads per Block: 1024
13 Max Block Dimensions: (1024, 1024, 64)
14 Max Grid Dimensions: (2147483647, 65535, 65535)
15
16  Part 2:
17  --- Part 2.1: H to D and D to H Transfer Times ---
18  MatrixSizes: 256, 512, 1024, 2048, 4096
19
20  ✓ Host to Device Transfer Times (ms) by Matrix Size:
21  Size 256 x 256 : 0.173728 ms
22  Size 512 x 512 : 0.514912 ms
23  Size 1024 x 1024 : 1.522336 ms
24  Size 2048 x 2048 : 5.598752 ms
25  Size 4096 x 4096 : 21.819136 ms
26
27  ✓ Device to Host Transfer Times (ms) by Matrix Size:
28  Size 256 x 256 : 0.191232 ms
29  Size 512 x 512 : 0.601984 ms
30  Size 1024 x 1024 : 1.534912 ms
31  Size 2048 x 2048 : 5.465600 ms
32  Size 4096 x 4096 : 21.367456 ms
33
34  --- Part 2.2 : CPU vs. GPU (Single Thread) ---
35  MatrixSizes: 256, 512, 1024
36
37  ✓ Matrix Size 256 x 256
38  CPU Time (ms) : 35.000000
39  ✓ GPU Time (1 block,1 thread) (ms) : 8532.343750 (NO Transfer), 8532.598633 (WITH Transfer)
40  Transfer Times: H to D = 0.156704 ms, D to H = 0.098944 ms
41  Test PASSED
42
43  ✓ Matrix Size 512 x 512
44  CPU Time (ms) : 475.000000
45  ✓ GPU Time (1 block,1 thread) (ms) : 78453.343750 (NO Transfer), 78454.273438 (WITH Transfer)
46  Transfer Times: H to D = 0.550944 ms, D to H = 0.372160 ms
47  Test PASSED
48
49  ✓ Matrix Size 1024 x 1024
50  CPU Time (ms) : 3718.000000
51  ✓ GPU Time (1 block,1 thread) (ms) : 656019.062500 (NO Transfer), 656022.812500 (WITH Transfer)
52  Transfer Times: H to D = 1.570624 ms, D to H = 2.213408 ms
53  Test PASSED
54
55  --- Part 3.3: Kernel Times with varying Matrix Size & Block Width---
56  MatrixSizes: 256, 512, 1024, 2048, 4096 | BlockWidth: 2,4,8,16,32
57
58  ✓ BlockWidth = 2
59  Size 256 x 256 -> Kernel Time = 164.616379 ms
60  Size 512 x 512 -> Kernel Time = 353.023529 ms
61  Size 1024 x 1024 -> Kernel Time = 1286.574219 ms
62  Size 2048 x 2048 -> Kernel Time = 10793.482422 ms
63  Size 4096 x 4096 -> Kernel Time = 89778.617188 ms
64
65  ✓ BlockWidth = 4
66  Size 256 x 256 -> Kernel Time = 4.584320 ms
67  Size 512 x 512 -> Kernel Time = 38.927265 ms
68  Size 1024 x 1024 -> Kernel Time = 366.748840 ms
69  Size 2048 x 2048 -> Kernel Time = 2919.217529 ms
70  Size 4096 x 4096 -> Kernel Time = 23755.808594 ms
71
72  ✓ BlockWidth = 8
73  Size 256 x 256 -> Kernel Time = 1.670912 ms
74  Size 512 x 512 -> Kernel Time = 12.938912 ms
75  Size 1024 x 1024 -> Kernel Time = 87.001022 ms
76  Size 2048 x 2048 -> Kernel Time = 791.867615 ms
77  Size 4096 x 4096 -> Kernel Time = 6442.567871 ms
78
79  ✓ BlockWidth = 16
80  Size 256 x 256 -> Kernel Time = 1.344320 ms
81  Size 512 x 512 -> Kernel Time = 11.579136 ms
82  Size 1024 x 1024 -> Kernel Time = 104.591614 ms
83  Size 2048 x 2048 -> Kernel Time = 788.934998 ms
84  Size 4096 x 4096 -> Kernel Time = 6290.804688 ms
85
86  ✓ BlockWidth = 32
87  Size 256 x 256 -> Kernel Time = 1.751808 ms
88  Size 512 x 512 -> Kernel Time = 13.951296 ms
89  Size 1024 x 1024 -> Kernel Time = 89.537376 ms
90  Size 2048 x 2048 -> Kernel Time = 760.267029 ms
91  Size 4096 x 4096 -> Kernel Time = 6298.536133 ms

```