

./dot_product.cu

Tue Oct 22 22:06:50 2024

1

```

1: #include <stdio.h>
2: #include <time.h>
3:
4:
5: #define BLOCK_SIZE 32
6:
7: const int DSIZE = 256;
8: const int a = 1;
9: const int b = 1;
10:
11: // error checking macro
12: #define cudaCheckErrors() \
13:     do { \
14:         cudaError_t __err = cudaGetLastError(); \
15:         if (__err != cudaSuccess) { \
16:             fprintf(stderr, "Error: %s at %s:%d \n", \
17:                 cudaGetErrorString(__err), __FILE__, __LINE__); \
18:             fprintf(stderr, "*** FAILED - ABORTING***\n"); \
19:             exit(1); \
20:         } \
21:     } while (0)
22:
23:
24: // CUDA kernel that runs on the GPU
25: __global__ void dot_product(const int *A, const int *B, int *C, int N) {
26:
27:     // FIXME
28:     // Use atomicAdd
29:     int idx = threadIdx.x + blockIdx.x * blockDim.x;
30:     if (idx < N) {
31:         atomicAdd(C, A[idx]*B[idx]);
32:     }
33: }
34:
35: int main() {
36:
37:     // Create the device and host pointers
38:     int *h_A, *h_B, *h_C, *d_A, *d_B, *d_C;
39:
40:     // Fill in the host pointers
41:     h_A = new int[DSIZE];
42:     h_B = new int[DSIZE];
43:     h_C = new int;
44:     for (int i = 0; i < DSIZE; i++){
45:         h_A[i] = a;
46:         h_B[i] = b;
47:     }
48:
49:     *h_C = 0;
50:
51:
52:     // Allocate device memory
53:     cudaMalloc(&d_A, DSIZE*sizeof(int));
54:     cudaMalloc(&d_B, DSIZE*sizeof(int));
55:     cudaMalloc(&d_C, sizeof(int));
56:     // Check memory allocation for errors
57:     cudaCheckErrors();
58:     // Copy the matrices on GPU
59:     cudaMemcpy(d_A, h_A, DSIZE*sizeof(int), cudaMemcpyHostToDevice);
60:     cudaMemcpy(d_B, h_B, DSIZE*sizeof(int), cudaMemcpyHostToDevice);
61:     cudaMemcpy(d_C, h_C, sizeof(int), cudaMemcpyHostToDevice);
62:     // Check memory copy for errors
63:     cudaCheckErrors();
64:     // Define block/grid dimentions and launch kernel
65:     dot_product<<<DSIZE/BLOCK_SIZE, BLOCK_SIZE>>>(d_A, d_B, d_C, DSIZE);
66:     // Copy results back to host
67:     cudaMemcpy(h_C, d_C, sizeof(int), cudaMemcpyDeviceToHost);
68:     // Check copy for errors
69:     cudaCheckErrors();
70:     // Verify result
71:     printf("A.B = %d\n", *h_C);
72:     // Free allocated memory
73:     cudaFree(d_A);
74:     cudaFree(d_B);
75:     cudaFree(d_C);
76:
77:     free(h_A);
78:     free(h_B);
79:     free(h_C);
80:
81:     return 0;
82:
83: }

```

```

1: #include <stdio.h>
2: #include <algorithm>
3:
4:
5: using namespace std;
6:
7: #define N 64
8: #define RADIUS 2
9: #define BLOCK_SIZE 32
10:
11:
12: __global__ void stencil_2d(int *in, int *out) {
13:
14:     __shared__ int temp[BLOCK_SIZE + 2 * RADIUS][BLOCK_SIZE + 2 * RADIUS];
15:     int gindex_x = threadIdx.x + blockDim.x * blockIdx.x;
16:     int lindex_x = threadIdx.x + RADIUS;
17:     int gindex_y = threadIdx.y + blockDim.y * blockIdx.y;
18:     int lindex_y = threadIdx.y + RADIUS;
19:
20:     // Read input elements into shared memory
21:     int size = N + 2 * RADIUS;
22:     temp[lindex_x][lindex_y] = in[gindex_x*size + gindex_y];
23:     if (threadIdx.x < RADIUS) {
24:         temp[lindex_x - RADIUS][lindex_y] = in[(gindex_x - RADIUS)*size + gindex_y];
25:         temp[lindex_x + BLOCK_SIZE][lindex_y] = in[(gindex_x + BLOCK_SIZE)*size + gindex_y];
26:     }
27:
28:     if (threadIdx.y < RADIUS) {
29:         temp[lindex_x][lindex_y - RADIUS] = in[(gindex_x)*size + gindex_y - RADIUS];
30:         temp[lindex_x][lindex_y + BLOCK_SIZE] = in[gindex_x*size + gindex_y + BLOCK_SIZE];
31:     }
32:     __syncthreads();
33:     // Apply the stencil
34:     int result = 0;
35:     for (int offset = -RADIUS; offset <= RADIUS; offset++){
36:         result += temp[lindex_x + offset][lindex_y];
37:         result += temp[lindex_x][lindex_y + offset];
38:     }
39:     result -= temp[lindex_x][lindex_y]; // remove center overlap counter twice
40:
41:     // FIXME
42:     // Store the result
43:     out[gindex_y*size*gindex_x] = result;
44: }
45:
46:
47: void fill_ints(int *x, int n) {
48:     // Store the result
49:     // https://en.cppreference.com/w/cpp/algorithm/fill_n
50:     fill_n(x, n, 1);
51: }
52:
53:
54: int main(void) {
55:
56:     int *in, *out; // host copies of a, b, c
57:     int *d_in, *d_out; // device copies of a, b, c
58:
59:     // Alloc space for host copies and setup values
60:     int size = (N + 2*RADIUS)*(N + 2*RADIUS) * sizeof(int);
61:     in = (int *)malloc(size); fill_ints(in, (N + 2*RADIUS)*(N + 2*RADIUS));
62:     out = (int *)malloc(size); fill_ints(out, (N + 2*RADIUS)*(N + 2*RADIUS));
63:
64:     // Alloc space for device copies
65:     cudaMalloc(&d_in, size);
66:     // FIXME
67:     cudaMalloc(&d_out, size);
68:
69:     // Copy to device
70:     cudaMemcpy(d_in, in, size, cudaMemcpyHostToDevice);
71:     // FIXME
72:     cudaMemcpy(d_out, out, size, cudaMemcpyHostToDevice);
73:
74:     // Launch stencil_2d() kernel on GPU
75:     int gridSize = (N + BLOCK_SIZE-1)/BLOCK_SIZE;
76:     dim3 grid(gridSize, gridSize);
77:     dim3 block(BLOCK_SIZE, BLOCK_SIZE);
78:     // Launch the kernel
79:     // Properly set memory address for first element on which the stencil will be applied
80:     stencil_2d<<<grid,block>>>>(d_in + RADIUS*(N + 2*RADIUS) + RADIUS, d_out + RADIUS*(N + 2*RADIUS) +
RADIUS);
81:
82:     // Copy result back to host

```

```
83:         // FIXME
84:         cudaMemcpy(out, d_out, size, cudaMemcpyDeviceToHost);
85:
86:         // Error Checking
87:         for (int i = 0; i < N + 2 * RADIUS; ++i) {
88:             for (int j = 0; j < N + 2 * RADIUS; ++j) {
89:
90:                 if (i < RADIUS || i >= N + RADIUS) {
91:                     if (out[j+i*(N + 2 * RADIUS)] != 1) {
92:                         printf("Mismatch at index [%d,%d], was: %d, should be: %d\n", i, j,
out[j+i*(N + 2 * RADIUS)], 1);
93:                         return -1;
94:                     }
95:                 }
96:                 else if (j < RADIUS || j >= N + RADIUS) {
97:                     if (out[j+i*(N + 2 * RADIUS)] != 1) {
98:                         printf("Mismatch at index [%d,%d], was: %d, should be: %d\n", i, j,
out[j+i*(N + 2 * RADIUS)], 1);
99:                         return -1;
100:                     }
101:                 }
102:                 else {
103:                     if (out[j+i*(N + 2 * RADIUS)] != 1 + 4 * RADIUS) {
104:                         printf("Mismatch at index [%d,%d], was: %d, should be: %d\n", i, j,
out[j+i*(N + 2 * RADIUS)], 1 + 4*RADIUS);
105:                         return -1;
106:                     }
107:                 }
108:             }
109:         }
110:
111:         // Cleanup
112:         free(in);
113:         free(out);
114:         cudaFree(d_in);
115:         cudaFree(d_out);
116:         printf("Success!\n");
117:
118:         return 0;
119:     }
120: }
121:
```