

```
1: #pragma omp parallel for
2: for (i = 0; i < ncolsa; i++) {
3:     for (j = 0; i < nrowsb; j++) {
4:         sum = 0;
5:         for (k = 0; i < ncolsa; i++)
6:             sum += a[i][k] * b[k][j];
7:     }
8:
9:
10: __global__ void matmul(float *ma,float *mb,float *mc,int nrowsa,int ncolsa,int ncolsb)
11: {    int k;
12:     sum = 0;
13:     for (k = 0; i < ncolsa; i++)
14:         sum += a[i*ncolsa+k] * b[k*ncols+j];
15: }
16:
17: // from http://astro.pas.rochester.edu/~aquillen/gpuworkshop/AdvancedCUDA.pd
f
18:
19: __global__ void MultiplyOptimise(const float *A, const float *B, float *C) {
20: // Extract block and thread numbers
21: int bx = blockIdx.x; int by = blockIdx.y;
22: int tx = threadIdx.x; int ty = threadIdx.y;
23:
24: // Index of first A sub-matrix processed by this block
25: int aBegin = dc_wA * BLOCK_SIZE * by;
26: // Index of last A sub-matrix
27: int aEnd = aBegin + dc_wA - 1;
28: // Stepsize of A sub-matrices
29: int aStep = BLOCK_SIZE;
30: // Index of first B sub-matrix
31: // processed by this block
32: int bBegin = BLOCK_SIZE * bx;
33: // Stepsize for B sub-matrices
34: int bStep = BLOCK_SIZE * dc_wB;
35: // Accumulator for this thread
36: float Csub = 0;
37: for(int a = aBegin, b = bBegin; a <= aEnd; a += aStep, b+= bStep) {
38:     // Shared memory for sub-matrices
39:     __shared__ float As[BLOCK_SIZE][BLOCK_SIZE];
40:     __shared__ float Bs[BLOCK_SIZE][BLOCK_SIZE];
41:     // Load matrices from global memory into shared memory
42:     // Each thread loads one element of each sub-matrix
43:     As[ty][tx] = A[a + (dc_wA * ty) + tx];
44:     Bs[ty][tx] = B[b + (dc_wB * ty) + tx];
45:     // Synchronise to make sure load is complete
46:     __syncthreads();
47:     // Perform multiplication on sub-matrices
48:     // Each thread computes one element of the C sub-matrix
49:     for( int k = 0; k < BLOCK_SIZE; k++ ) {
50:         Csub += As[ty][k] * Bs[k][tx];
51:     }
52:     // Synchronise again
53:     __syncthreads();
54: }
55: // Write the C sub-matrix back to global memory
56: // Each thread writes one element
57: int c = (dc_wB * BLOCK_SIZE * by) + (BLOCK_SIZE*bx);
58: C[c + (dc_wB*ty) + tx] = Csub;
59: }
60:
61: #define BLOCK_SIZE 16
62:
63: __constant__ int dc_wA;
64: __constant__ int dc_wB;
65:
```