Minimum using parallel reduction 20 May 2025 10:02 Dry Run Step-1 boud global memory into shared memory -> each thread copies 2 values Shared [tid] = input [tid] me have an arrany of size 512 filled with floats e.g of 0.0, 1.0, 2.0 - . - . 511.03 shared [tid + block Dmx.x] = input [tid + block Dimx.x] threads per block = 256 showed [0] = input[0] = 0 Showred [0+256] = input [256] = 256 earth Hurread handles 2 elements So after all threads finish 2 also me houre 64 blocks Sharred []= [0,1,2 --- 512] Element = 512 = 2* Num-threads Step-2 Revoval Staut Time threads_per-block = 256 Blocks - 64 (6 = -bit)timen (bid] = clock (); Llock Cycle = A single tick of the GPU's internal timer, used to nd= 256 ->128->64->32->16->8->4->2->1 measure time taken by operations 9 steps on (109, (512))=9 Step-3 d = block Dimx.x = 256 iteration 1 2) Minimum Reduction > A way to find the smallest number in so me agre comparing sharred [tid] & sharred (tid + of) a læge dataset by using parallel comparisions. so por tid 256 eg. ib we have [10,5,20,8,3] sharred[0] = min (sharred(0], sharred[256]) = min (0, 256) = 0 then a minimum reduction would Showred [1]= min (showred [1], showed [257]) = min (1,257)= 1 give us 3 Shaved [255] = min (shared [255], shared [511] = min (255, 511) = 255 extern _ shared _ pleat shared []: fragiteration, d= d= 256 = 128 - This is how you declare shared memory in CUDA dynamically. -> thoreads in the same block can access this shared memory very quilely. -> we use extern __sharred__ when we don't know the size at compile time & for tid < 128 -> We can pass the size when we launch the kernel, in the execution configuration. Sharred [0] = min (sharred [0], sharred [128]) = min(0,128) = 0kernel <<< blocks, threads, shared-mem_size>>>> (....) shared [1] = min (shared [1], Shared [129]) = min(1,128) = input -> 512 floats -> [Num-threads * 2] Linput -> 11 11 Sharred [127] = min (sharred [127], sharred [255]) = min (127,255]= 127 doutput > 1 plant per block -> [Num-blocks] dimen > 2 dock entries peu block > so [2xNvm_blocks] as me continue sstout me will be left with, d=1 tiel = theread Idx.x bid= blockldx.x Sharred [0] = min (sharred [0], sharred [1]= min(0,1) = 0 Isthis stones the min value from input (0____511) clock - t >> It is used to measure how many GIPU Step Output Write clock egdes a certain aperation takes if (fid = = 0) an exercation took 1500000 y des entput [bid] = shared [0] your gpu's clock state is 1.561X2/1.5X109 ydes/sec) Reward End Time Step-5 thren time = 1500000 sec > 1 milliserand timer [bid + goud Dimx.x] = clock() of block 63 of black 1 times [bid+ griadimx.x] = dock () — end stant 0, stant 1 ---. stant 63, end 0, end 1 --.. end 63 end time of block D fine of Elapsed Clocks = Num-blocks] - Hiner [i]

Ang Elapsed Clocks = Elapsed Clocks

```
Executing.
                                                                   Average clocks/block = 9935.703125
                                                                  Exit status: 0
       extern shared float shared[]
        const int tid = threadIdx.x;
       const int bid = blockIdx.x;
         timer[bid] = clock();
        / copying data from global memory to shared memory
       shared[tid + blockDim.x] = input[tid + blockDim.x];
        / parallel reduction
       for(int i=blockDim.x;i>0;i /= 2){
       _syncthreads();
if(tid<i){[
   float f0 = shared[tid];</pre>
            float f1 = shared[tid+i];
            shared[tid] = f1>f0? f0:f1;
                  // storing the ouput
                 if(tid==0){
                      output[bid] = shared[0];
                  __syncthreads();
                 // stop clock
                 if(tid==0){
                       timer[bid+gridDim.x] = clock();
   43
   44
   45
   46
   47
   49 #define NUM BLOCKS 64
   50 #define NUM THREADS 256
52 int main(){
          float *dinput = NULL;
          float *doutput = NULL;
          clock_t *dtimer = NULL;
          clock_t timer[NUM_BLOCKS * 2];
          float input[NUM_THREADS * 2];
          for(int i=0;i<NUM_THREADS*2;i++){</pre>
              input[i] = (float)i;
          cudaMalloc(&dinput, sizeof(float) * NUM_THREADS*2);
          cudaMalloc(&doutput, sizeof(float)*NUM_BLOCKS);
          cudaMalloc(&dtimer, sizeof(clock_t)*NUM_BLOCKS*2);
          cudaMemcpy(dinput ,input ,sizeof(float)*NUM_THREADS * 2, cudaMemcpyHostToDevice);
         Myfunc<<<<NUM_BLOCKS,NUM_THREADS,sizeof(float)*NUM_THREADS*2>>>(dinput, doutput, dtimer);
          cudaMemcpy(timer, dtimer, sizeof(clock_t)*NUM_BLOCKS*2, cudaMemcpyDeviceToHost);
          cudaFree(dinput);
          cudaFree(doutput);
          cudaFree(dtimer);
             long double Elapsedclocks = 0;
             for(int i=0;i<NUM_BLOCKS;i++){</pre>
                 Elapsedclocks += (long double)(timer[i+NUM_BLOCKS] - timer[i]);
            long double avgElapsedclocks = Elapsedclocks/NUM_BLOCKS;
            printf("Average clocks/block = %Lf\n", avgElapsedclocks);
            return 0;
```