

PARALLEL ALGORITHMS HOMEWORK № 3

Aditya Bhamidipati, Georgia State University

04/18/2020

* Source code snippets displayed are for understanding purpose only and hence are logic centered. For full source code refer the attached .cu files

Problem 1

(50 pts) Write a CUDA program for calculating the prefix sums of an arbitrary-length array as described in Section 8.6 of the Kirk and Hwu book. You should follow their code snippets and input/output conventions.

Listing 1: Source Code – Parallel Prefix Sum method.

```
1
2 void psum(const unsigned int * const h_in, unsigned int * const h_out, ←
   const size_t len)
3 {
4     int nblocks = n/block_size;
5
6     // allocate memories on gpu
7     unsigned int *d_scan, *d_sums, *d_incr;
8
9
10    // scan array by blocks (block_size = 2 * num threads)
11    block_psum<<<nblocks, nthreads, smem>>>(d_in, d_scan, d_sums, block_size←
12      );
13
14    // scan block sums
15    // TODO case when nblocks is bigger than block_size (see TODO 1)
16    block_psum<<<1, nthreads, smem>>>(d_sums, d_incr, NULL, block_size);
17
18    // scatter block sums back to scanned blocks
19    scatter_incr<<<nblocks, nthreads>>>(d_scan, d_incr);
20 }
```

Following Listing 1... Here we initialize our GPU memory to hold our array elements. Scan all the array elements by blocks. Then we perform sum among the blocks separately. The idea is to scatter the block of arrays among GPUs and calculate the prefix-sums in a parallel fashion. Then we scatter the blocks back to their original places so as to give us the original array structure and every index position holds sum till that position.

Listing 2: Source Code – Parallel Block Prefix Sum method.

```
1
2 __global__
3 void block_psum(const unsigned int * const g_in, unsigned int * const ←
4     g_out, unsigned int * const g_sums, const size_t n)
5 {
6     // up sweep
7
8     for (int d = n >> 1; d > 0; d >>= 1)
9     {
10         __syncthreads();
11
12         if (tx < d)
13         {
14             int ai = offset * (2*tx+1) - 1;
15             int bi = offset * (2*tx+2) - 1;
16
17             smem[bi] += smem[ai];
18         }
19         offset <= 1;
20     }
21
22     // save block sum and clear last element
23     if (tx == 0) {
24         if (g_sums != NULL)
25             g_sums[blockIdx.x] = smem[n-1];
26         smem[n-1] = 0;
27     }
28
29     // down sweep
30
31     for (int d = 1; d < n; d <= 1)
32     {
33         offset >= 1;
34         __syncthreads();
35
36         if (tx < d)
37         {
38             int ai = offset * (2*tx+1) - 1;
39             int bi = offset * (2*tx+2) - 1;
40
41             // swap
42             unsigned int t = smem[ai];
43             smem[ai] = smem[bi];
44             smem[bi] += t;
45         }
46     }
47 }
```

```

46    }
47    __syncthreads();
48
49    // save scan result
50    g_out[2*px]    = smem[2*tx];
51    g_out[2*px+1] = smem[2*tx+1];
52 }
```

Following Listing 2... **Upsweep:** This is also known as a reduce operation. The reduction happens in a hierarchical fashion. During the first iteration, the neighbors with one hop away are added. In the subsequent iterations neighbors with 2, 4, 8 etc., hops are added respectively until we are left with one number that is the sum of all elements of the array. **Downsweep:** We start by resetting the right most value by the identity operator. The downsweep creates a mirror image of hierarchy formed by upsweep. If a value is missing for a downsweep operation, it is simply copied from the upsweep operation.

Listing 3: Source Code – Combining the scattered blocks.

```

1
2 __global__
3 void scatter_incr(      unsigned int * const d_array,
4                      const unsigned int * const d_incr)
5 {
6     const size_t bx = 2 * blockDim.x * blockIdx.x;
7     const size_t tx = threadIdx.x;
8     const unsigned int u = d_incr[blockIdx.x];
9     d_array[bx + 2*tx]    += u;
10    d_array[bx + 2*tx+1] += u;
11 }
```

Following Listing 3... Here we combine all the scattered blocks into single array, where the value at any index gives the prefix sum till that position.

```

File Edit Selection View Go Run Terminal Help
psum.cu - vbhamidipati1 [SSH: hpc_servers] - Visual Studio Code
EXPLORER OPEN EDITORS
stb_image.h HW3
psum.cu HW3
psum_hierarchical.cu HW3
nearest_neighbor_interpolation.cu
VBHAMIDIPATI1 [SSH: HPC SERVERS]
HW3
colordefs.h
decimal.h
image.jpg
knn.h
nearest_neighbor_interpolation.cu
nearest_neighbor_interpolation.exe
output.png
psum_hierarchical.cu
psum_hierarchical.exe
psum.cu
psum.exe
sample.png
stb_image_write.h
stb_image.h
text.h
try.cu
try.exe
utils.h
wtime.h
ondemand / data / sys / dashboard / b...
.bash_history
.bash_logout
.bash_profile
OUTLINE
TIMELINE
SSH: hpc_servers 0 0 0
Type here to search
1: ssh
PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL
[vbhamidipati1@cder01 ~]$ scp "/home/users/vbhamidipati1/HW3/psum.cu" vbhamidipati1@cder01:/home/users/vbhamidipati1/HW3

```

(88).png

```

File Edit Selection View Go Run Terminal Help
psum.cu - vbhamidipati1 [SSH: hpc_servers] - Visual Studio Code
EXPLORER OPEN EDITORS
stb_image.h HW3
psum.cu HW3
psum_hierarchical.cu HW3
nearest_neighbor_interpolation.cu
VBHAMIDIPATI1 [SSH: HPC SERVERS]
HW3
colordefs.h
decimal.h
image.jpg
knn.h
nearest_neighbor_interpolation.cu
nearest_neighbor_interpolation.exe
output.png
psum_hierarchical.cu
psum_hierarchical.exe
psum.cu
psum.exe
sample.png
stb_image_write.h
stb_image.h
text.h
try.cu
try.exe
utils.h
wtime.h
ondemand / data / sys / dashboard / b...
.bash_history
.bash_logout
.bash_profile
OUTLINE
TIMELINE
SSH: hpc_servers 0 0 0
Type here to search
1: ssh
PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL
[vbhamidipati1@cder01 ~]$ nvcc ./HW3/psum.cu -o ./HW3/psum.exe
nvcc warning : The 'compute_20', 'sm_20', and 'sm_21' architectures are deprecated, and may be removed in a future release (Use -Wno-deprecated-gpu-targets to suppress warning).
[vbhamidipati1@cder01 ~]$ 

```

(89).png

The screenshot shows a Visual Studio Code interface with the following details:

- File Explorer:** Shows project files under "OPEN EDITORS" (stb_image.h, psum.cu, psum_hierarchical.cu, nearest_neighbor_interpolation.cu) and "VBHAMIDIPATI1 [SSH: hpc_servers]" (HW3, colordefs.h, decimal.h, image.jpg, knn.h, nearest_neighbor_interpolation.cu, nearest_neighbor_interpolation.exe, output.png, psum.hierarchical.cu, psum.hierarchical.exe, psum.cu, psum.exe, sample.png, stb_image_write.h, stb_image.h, text.h, try.cu, try.exe, utils.h, wtime.h, ondemand/data/sys/dashboard/b..., bash_history, bash_logout, bash_profile). The "psum.cu" file is currently open.
- Code Editor:** Displays the CUDA code for the "block_psum" function. The code includes header imports for , , , , and "utils.h". It defines a constant `MAX_THREADS_PER_BLOCK` and uses `#pragma unroll`. The `block_psum` function takes pointers to input and output arrays and their sizes, and an array of sums. It then iterates over the input array, calculating the sum of elements at indices `g_in[tx], g_in[tx+1], ..., g_in[tx+MAX_THREADS_PER_BLOCK-1]` and storing it in `g_out[tx]`.
- Terminal:** Shows the command-line session on the "hpc_servers" SSH connection. The user has entered several numbers (1, 2, 3, 4, 5, 10) which are being processed by the CUDA code.

Problem 2

(20 pts) Write a CUDA program to implement the single-pass scan version of the hierarchical prefix sum algorithm

Listing 4: Parallel Upsweep

```
1 __global__ void prefix_upsweep_kernel (int *b_d, int *a_d, int n, int ←
2   depth, int *blocksum_device) {
3
4   while (tid < n) {
5     smem[threadIdx.x] = a_d[tid];           // each thread copy data to ←
6     shared memory
7     __syncthreads();                      // wait for all threads
8
9     offset = 1;                           //1->2->4->8
10    for (d = 1; d <= depth ; d++) {
11      offset *= 2;
12      if (threadIdx.x % offset == offset-1 ){
13        smem[threadIdx.x]+= smem[threadIdx.x- offset/2];
14        __syncthreads();
15      }
16    } // end for loop
```

```
17
18     b_d[tid] = smem[threadIdx.x];           // *write the result to array←
          b_d[tid] location
19
20     if (threadIdx.x == blockDim.x -1){
21         blocksum_device[blockIdx.x] = smem[threadIdx.x];
22     }
23
24     __syncthreads();                         // wait for all threads to ←
          write results
25
26     tid += blockDim.x * gridDim.x;
27
28 } // end while (tid < n)
29 } // end kernel function
```

Following Listing 4...

Name: prefix_upsweep_kernel

Input: int *b_d, int *a_d, int n, int depth, int *blocksum_device

Output: int *blocksum_device

Operation: Performs the upsweep sum on each block. b_d is updated but not copied to CPU yet.

Note: Size of blocksum_device is same as number of blocks

Listing 5: Parallel downsweep

```
1
2 __global__ void prefix_downsweepsweep_kernel (int *b_d, int *a_d, int n, ←
3   int depth, int *blocksum_device) {
4
5     while (tid < n) {
6       smem[threadIdx.x] = b_d[tid];           // each thread copy data to ←
7         shared memory from previous results b_d
8       //b_d[tid] = smem[threadIdx.x];
9       if (threadIdx.x == blockDim.x - 1 && blockIdx.x != 0){
10         smem[threadIdx.x] = blocksum_device[blockIdx.x-1];
11         //b_d[tid] = blocksum_device[blockIdx.x-1];
12       }
13       if (tid == blockDim.x - 1 ){ // clear last entry in only first ←
14         block - special case
15         smem[threadIdx.x] = 0;
16       }
17
18       __syncthreads();                      // wait for all threads
19
20       offset = blockDim.x;                  //8 -> 4 -> 2
21       for (d = depth; d > 0 ; d--) {        // depth 3 -> 2->1
22         a_d[threadIdx.x] = a_d[threadIdx.x] +
```

```

19
20     if (threadIdx.x % offset == offset-1 ){
21         int tmp3 = smem[threadIdx.x];
22         int tmp1 = smem[threadIdx.x- offset/2];
23         smem[threadIdx.x- offset/2] = tmp3;
24         __syncthreads();
25         smem[threadIdx.x]+= tmp1;
26         __syncthreads();
27     }
28     offset /= 2;
29
30 } // end for loop
31
32 b_d[tid] = smem[threadIdx.x];           // *write the result to array←
33                                b_d[tid] location
34 __syncthreads();                         // wait for all threads to ←
35                                write results
36
37 tid += blockDim.x * gridDim.x;
38
39 } // end while (tid < n)
} // end kernel function

```

Following Listing 5...

Name: prefix_downsweepsweep_kernel

Input: int *b_d, int *a_d, int n, int depth, int *blocksum_device

Output: int *b_d

Operation: Clears the last element in first block. Copies last element in all remaining blocks (nth) with content from corresponding location in blocksum_device[n-1]. Last element in blocksum_device is not needed in exclusive scan.

Note: blocksum_device is updated by cpu with cummulative sums, then updated blocksum_device is given as input to this kernel

The screenshot shows the Visual Studio Code interface with the following details:

- File Explorer:** Shows files in the current workspace, including `stb_image.h`, `psum.cu`, `psum_hierarchical.cu`, and `nearest_neighbor_interpolation.cu`. A folder named `HW3` contains several header files like `colordefs.h`, `decimal.h`, `image.jpg`, `knn.h`, and source files like `psum_hierarchical.cu`.
- Terminal:** The terminal window shows the command being run: `nvcc ./HW3/psum_hierarchical.cu -o ./HW3/psum_hierarchical.exe`. It also displays a warning from nvcc: "warning: The 'compute_20', 'sm_20', and 'sm_21' architectures are deprecated, and may be removed in a future release (Use -Wno-deprecated-gpu-targets to suppress warning)."
- Code Editor:** The main editor pane displays the CUDA code for `psum_hierarchical.cu`. The code implements a hierarchical prefix sum algorithm using shared memory and multiple levels of reduction. It includes comments explaining the thread indexing and shared memory usage.
- Bottom Status Bar:** Shows the current file is `psum_hierarchical.cu`, line 60, column 54. It also shows the number of spaces (4), the current file type (UTF-8), and the CUDA C++ language mode.

(91).png

(92).png

The screenshot shows a Visual Studio Code interface with the following details:

- File Explorer:** Shows files in the current workspace, including `stb_image.h`, `psum.cu`, `psum_hierarchical.cu`, and `nearest_neighbor_interpolation.cu`.
- Editor:** Displays the `psum_hierarchical.cu` file, which contains CUDA code for a hierarchical summation kernel.
- Terminal:** Shows the output of the compilation and execution of the code on an NVidia GPU. The terminal output includes:
 - Compilation command: `nvcc ./HW3/psum_hierarchical.cu -o ./HW3/psum_hierarchical.exe`
 - Warning message: "nvcc warning : The 'compute_20', 'sm_20', and 'sm_21' architectures are deprecated, and may be removed in a future release. (Use -Wno-deprecated-gpu-targets to suppress warning.)"
 - Execution command: `./HW3/psum_hierarchical.exe`
 - Input array is: A sequence of 1s and 0s.
 - CPU Reference Result is: A sequence of integers from 0 to 41.
 - Blocksum_CPU Result is: A sequence of integers from 8 to 63.
 - GPU Upsweep Result is: A sequence of integers from 8 to 64.
 - GPU Downsweep (final) Result is: A sequence of integers from 0 to 63.
 - Total entries: 64
- Bottom Status Bar:** Shows the current file is `(93).png`, the search bar, and various system icons.

Problem 3

(30 pts) Write a CUDA program for doing nearest-neighbor interpolation. This program should receive as input an $n \times m$ image and the upsampling factor k (an integer). Note that the resulting image will be $kn \times km$ and should look blocky.

Listing 6: Nearest Neighbor Resize

```
1 uint8_t *cudaNearestNeighborResize(uint8_t *origImageData, int origImgWidth←
    ,int origImgHeight,int newImgWidth,int newImgHeight,decimal_t ←
    resizeRatio)
2 {
3
4
5     int blocksPerLine=ceil(floatDiv(newImgWidth,CUDA_THREADS_PER_BLOCK));
6     int totalNumBlocks=blocksPerLine*newImgHeight; // Each block is ←
            responsible for a single line
7
8     if(totalNumBlocks>65535)
9         return 0;
10
11    size_t origImageSize=origImgWidth*origImgHeight*sizeof(uint8_t);
12
13    size_t newImageSize=newImgWidth*newImgHeight*sizeof(uint8_t);
14
```

```

15     cudaMalloc(&device_origImageData_in,origImageSize);
16     cudaMalloc(&device_newImageData_out,newImageSize);
17
18     cudaMemcpy(device_origImageData_in,origImageData,origImageSize,←
19             cudaMemcpyHostToDevice);
20
21     decimal_t origRatio=decimalDiv(1.0,decimal_t(resizeRatio));
22
23     device_nearestNeighborResize<<<totalNumBlocks,CUDA_THREADS_PER_BLOCK←
24             >>>(origImgWidth,newImgWidth,origRatio,blocksPerLine,←
25             device_origImageData_in,device_newImageData_out);
26
27     uint8_t *newImageData=(uint8_t*)malloc(newImageSize);
28     cudaMemcpy(newImageData,device_newImageData_out,newImageSize,←
29             cudaMemcpyDeviceToHost);
30
31     return newImageData;
32 }
```

Following Listing 6...

Nearest neighbor resize function resizes the image based on the upscaling factor. Copies the size of old image to the memory and creates space for new image based on the image size in the GPU. device_nearestNeighborResize takes the memory values into the CUDA device and performs resize operation

Listing 7: device nearest Neighbor Resize

```

1 __global__ void device_nearestNeighborResize(int origImgWidth,int ←
2             newImgWidth,decimal_t origRatio,int blocksPerLine,uint8_t *origImgData←
3             ,uint8_t *newImgData)
4 {
5     // Each block is responsible for 1024 px of a single row. Do not use ←
6             shared memory, as we cannot predict the distances between two ←
7             lines we need in the source image.
8
9     int newY=(blockId-(blockId%blocksPerLine))/blocksPerLine;
10    int blockIdInLine=blockId-newY*blocksPerLine;
11    int newX=blockIdInLine*CUDA_THREADS_PER_BLOCK+threadId;
12
13    if(newX>=newImgWidth)
14        return;
15
16    int oldX=(int)round(origRatio*((decimal_t)newX));
17    int oldY=(int)round(origRatio*((decimal_t)newY));
```

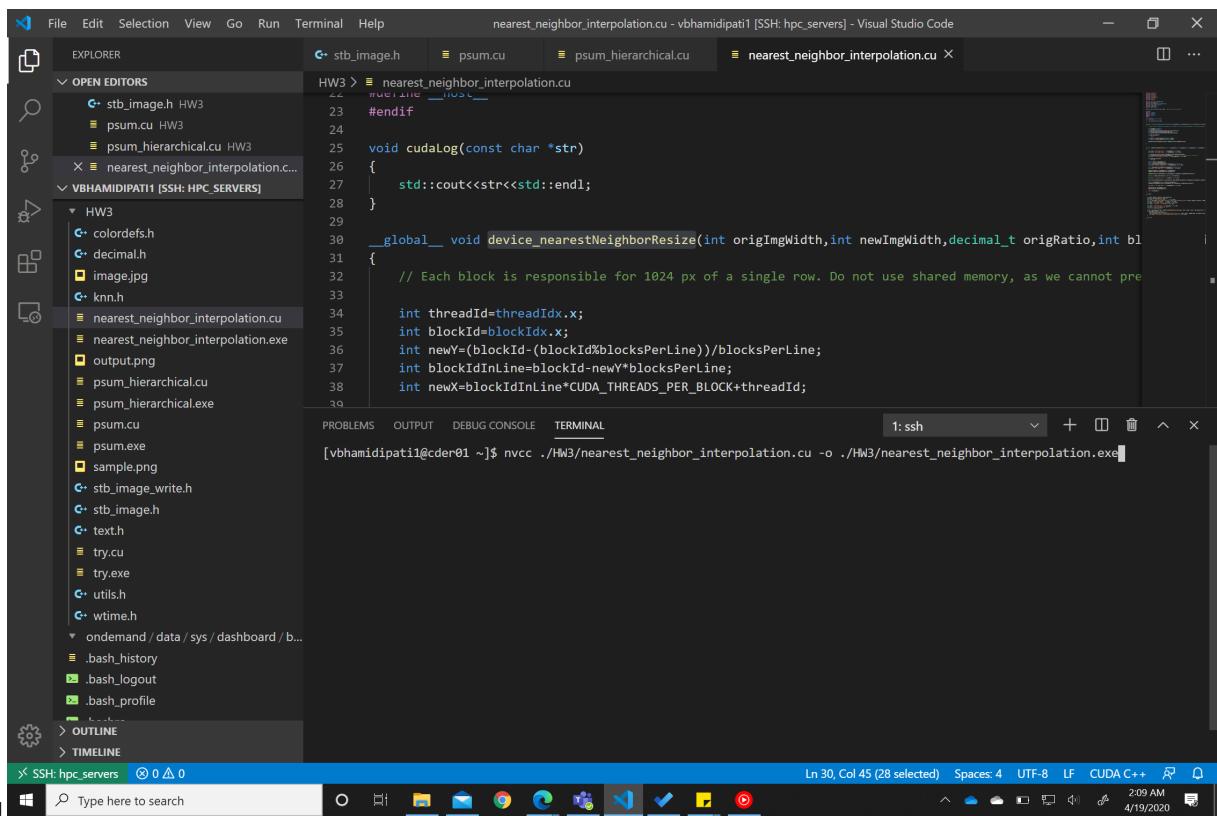
```

15     newImgData[newY*newImgWidth+newX]=origImgData[oldY*origImgWidth+oldX];
16
17 }

```

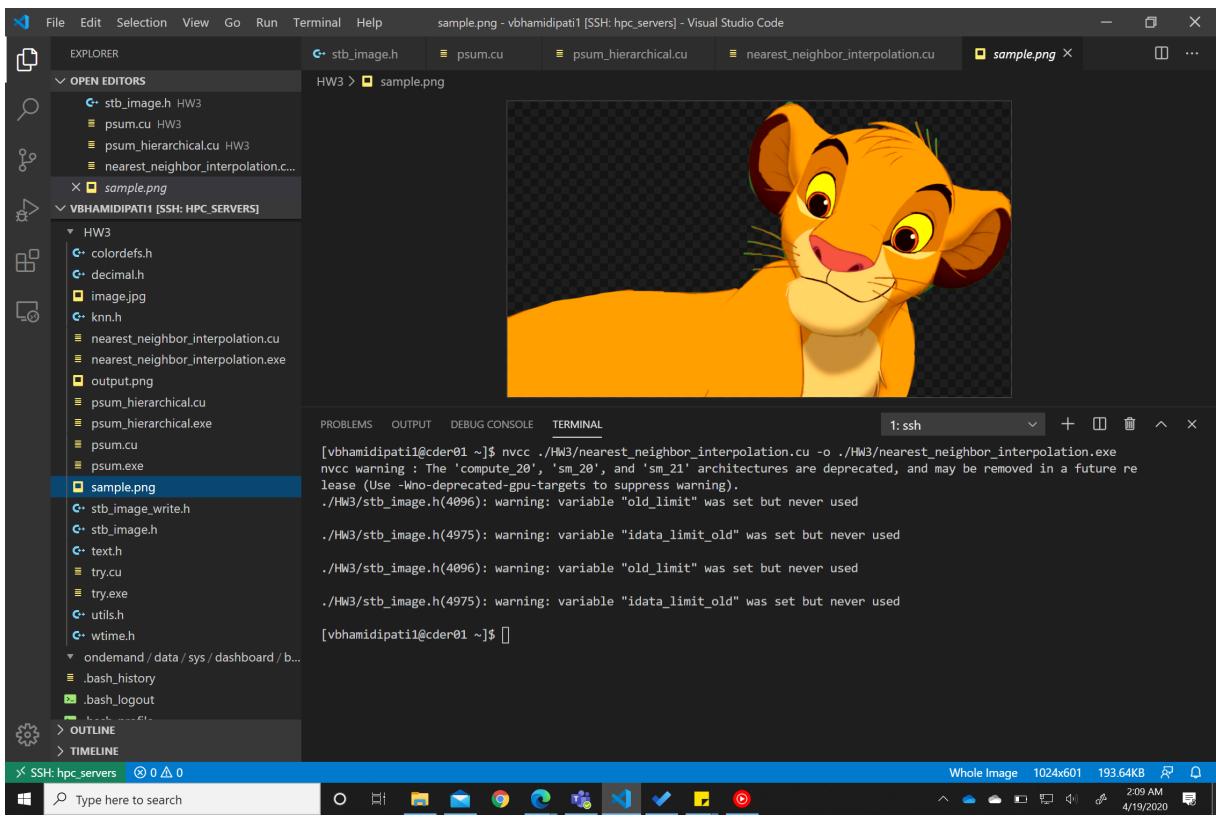
Following Listing 7...

Based on the old image ratio the newer image is adjusted and new data containing the image data is retuned. Here, the image data is stored in the `uint8_t` which specifies that it is stored in unsigned int format of 8 bits.

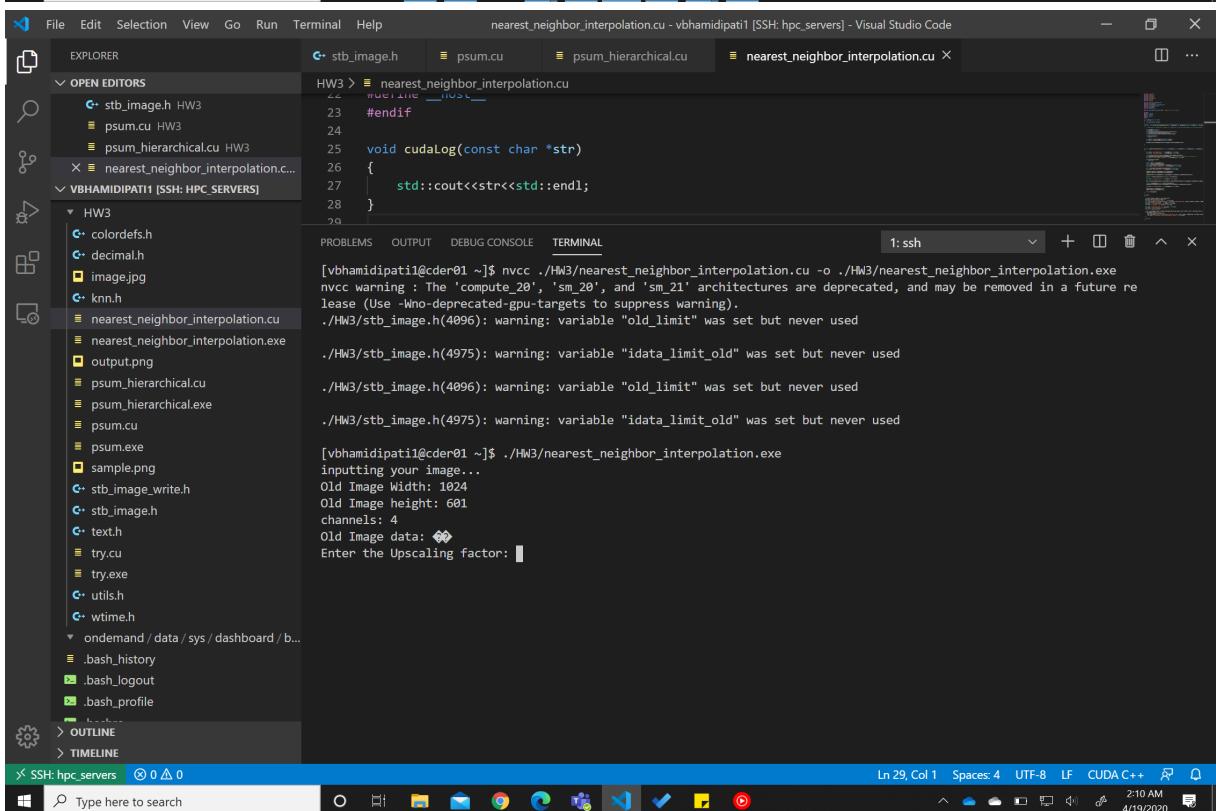


The screenshot shows the Visual Studio Code interface with the following details:

- File Explorer:** Shows a file tree for a project named "nearest_neighbor_interpolation". It includes files like `stb_image.h`, `psum.cu`, `psum_hierarchical.cu`, and `nearest_neighbor_interpolation.cu`. A connection to "VBHAMIDIPATI1 [SSH: HPC_SERVERS]" is visible.
- Editor:** Displays the CUDA C++ code for `nearest_neighbor_interpolation.cu`. The code includes a `cudaLog` function and a `device_nearestNeighborResize` function. The `device_nearestNeighborResize` function calculates new coordinates based on thread and block indices.
- Terminal:** Shows the command: `[vbhamidipati1@cdcer01 ~]$ nvcc ./HW3/nearest_neighbor_interpolation.cu -o ./HW3/nearest_neighbor_interpolation.exe`.
- Status Bar:** Shows "Ln 30, Col 45 (28 selected)" and various settings like "Spaces: 4", "UTF-8", "LF", "CUDA C++".



(95).png



(96).png

```

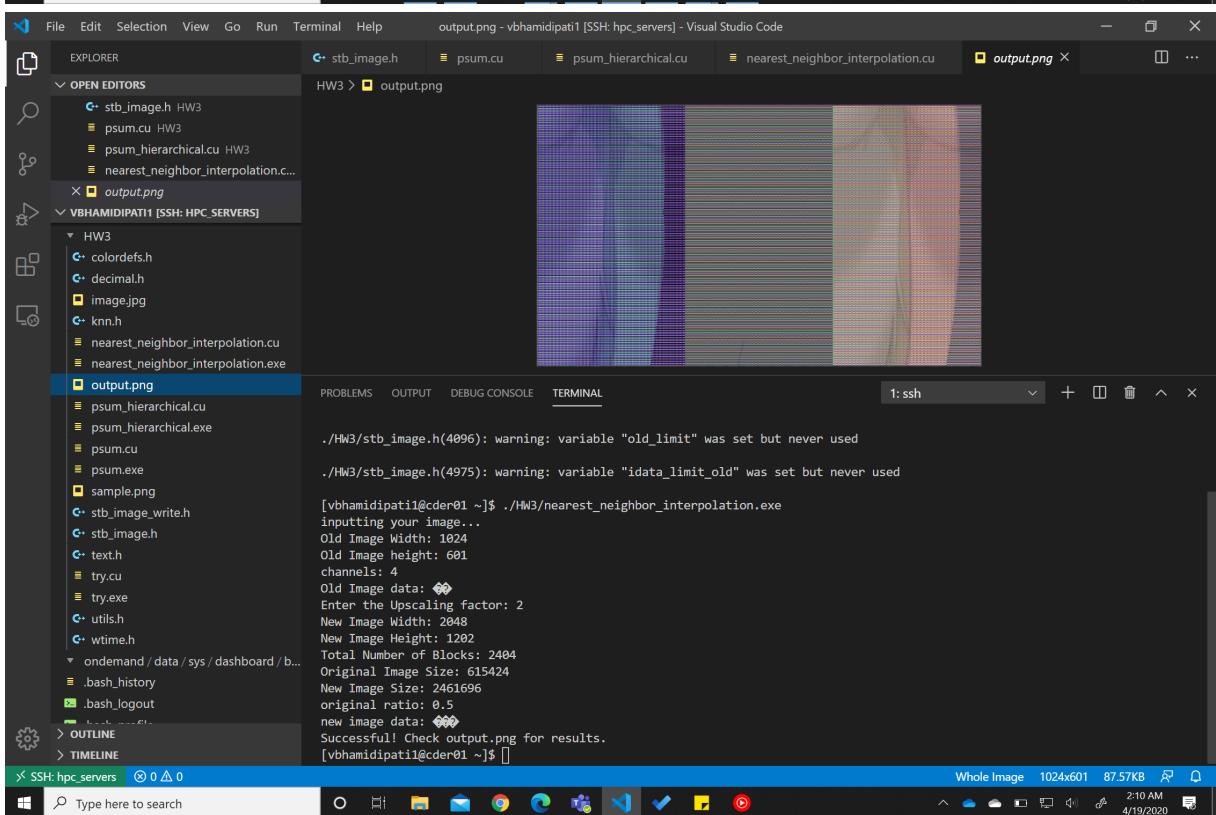
File Edit Selection View Go Run Terminal Help nearest_neighbor_interpolation.cu - vbhamidipati1 [SSH: hpc_servers] - Visual Studio Code
EXPLORER OPEN EDITORS
stb_image.h psum.cu psum_hierarchical.cu nearest_neighbor_interpolation.cu
HW3 > nearest_neighbor_interpolation.cu
23 #endif
24
25 void cudaLog(const char *str)
26 {
27     std::cout<<str<<std::endl;
28 }
29

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL
1:ssh + - x
[vbhamidipati1@cdcer01 ~]$ nvcc ./HW3/nearest_neighbor_interpolation.cu -o ./HW3/nearest_neighbor_interpolation.exe
nvcc warning : The 'compute_20', 'sm_20', and 'sm_21' architectures are deprecated, and may be removed in a future release (Use -Wno-deprecated-gpu-targets to suppress warning).
./HW3/stb_image.h(4096): warning: variable "old_limit" was set but never used
./HW3/stb_image.h(4975): warning: variable "idata_limit_old" was set but never used
./HW3/stb_image.h(4096): warning: variable "old_limit" was set but never used
./HW3/stb_image.h(4975): warning: variable "idata_limit_old" was set but never used
./HW3/stb_image.h(4975): warning: variable "idata_limit_old" was set but never used

[vbhamidipati1@cdcer01 ~]$ ./HW3/nearest_neighbor_interpolation.exe
inputting your image...
Old Image Width: 1024
Old Image height: 601
channels: 4
Old Image data: 
Enter the Upscaling factor: 2
New Image Width: 2048
New Image Height: 1202
Total Number of Blocks: 2404
Original Image Size: 615424
New Image Size: 2461696
original ratio: 0.5
new image data: 
Successful! Check output.png for results.
[vbhamidipati1@cdcer01 ~]$ 

```

(97).png



(98).png

Program 1 -> psum.cu

Program 2 -> psum_hierarchical.cu

Program 3 -> nearest_neighbor_interpolation.cu