

CSCI-GA.3033-012
**Graphics Processing Units (GPUs):
Architecture and Programming**
Spring 2012 – May 10th, 2012

NAME:

ID:

- This exam contains 5 questions with a total of 40 points.
- The exam is open book/notes but no electronic devices.
- If you have to make assumptions to continue solving a problem, state your assumptions clearly.
- You answer on the question sheet. You can use extra white papers if you want.

1. [5 points] We have a 2D array with X rows and Y columns. For best performance, the array needs to be memory aligned on 64 byte boundary (as we have seen in pitch). Derive a formula to calculate percentage of wasted memory (the padding) in terms of X and Y , assuming that a single array element is one byte. Assume that the array is NOT already aligned (i.e. ignore the special case where no padding is necessary).

Each row must be aligned at 64. Currently, each row takes Y bytes. So, it needs $(64 - (Y \% 64))$ extra bytes per row to be multiple of 64.

This makes total waste = $X(64 - (Y \% 64))$

The percentage waste is thus: $(64 - (Y \% 64)) / Y$

2. [5 points] Convert **126.0** to IEEE 754 Standard floating point. Show all steps.

We have 3 fields: sign (S), exponent (E), and mantissa (M)

The number is positive $\rightarrow S = 0$

126 in decimal = 1111110 $\rightarrow 1.11111 * 2^6$

$E = 6 + 127 = 133 \rightarrow 10000101$ (8 bits)

$M = 11111000\dots0$ (23 bits and we remove the 1. because it is hidden 1)

So the final answer is: 0 10000100 111110000...000

3. [10 points] Write a kernel that finds the location of all occurrences of zero in array of integers INPUT, and writes the locations of the zero elements to integer array z_locations. For example if INPUT is {0, 0, 4, 0, 8, 9} then z_locations will contain {0,1,3}.

Assumptions: both INPUT and z_locations are array of integers. Each array contains NUM elements. All elements in z_locations are already initialized to -1, so do not include initialization phase for z_locations. INPUT is already pre-loaded with the elements. You just need to write the kernel function and not a full program. Also assume that sizes of blocks and grid have already been set. Both INPUT and z_locations are already in the GPU global memory.

After finishing your code, indicate all the optimizations you did in your code, if any.

A basic, not fairly optimized version is:

```
__global__ zfind(int * INPUT, int * z_location){

    int tid = threadIdx.x + blockIdx.x * blockDim.x;
    int i, j;

    if (tid < NUM )
        if(INPUT[tid] == 0)
            z_locations[tid] = tid;

    __syncthreads();

    if(threadIdx.x == 0 && blockIdx.x == 0)
        for(i = 0; i < NUM; i ++){
            if(z_locations[i] == -1)
                for( j = i+1; j < NUM; j++){
                    if(z_locations[j] != -1){
                        z_locations[i] = z_locations[j];
                        z_locations[j] = -1;
                        break;
                    }
                }
        }
}
```

Example of optimizations that can be done here:

* Each block loads NUM/blockDim.x elements of INPUT and z_locations in shared memory and do the search there, then write it back to global memory.

* The compaction (i.e. the part after syncthreads) can be done in shared memory per block, then written back to global memory.

* Despite the existence of many “if”s, the effect of branch divergence is minimal due to lack of “else”.

4. [4 points] In FERMI memory hierarchy we have 64KB that can be configured as 48K shared memory and 16KB L1 cache or 48KB L1 cache and 16KB shared memory. Indicate when will you use the first configuration and when will you use the second configuration.

If there is a lot of locality in your code regarding data access, then a cache is needed more than shared memory because a cache is very successful to benefit from spatial and temporal locality. But if you want full control of what needs to be close to the SM and do not want it to be evicted, then shared memory is better.

5. [6 points] Assume we have M total threads, each of which need some data from the global memory of the GPU. Those threads can be grouped into X blocks with Y threads each (i.e. $X*Y=M$). Keeping the total number of threads fixed, discuss the effect of increasing X (and decreasing Y to keep M fixed) or increasing Y on bandwidth requirement. Assume the GPU can accommodate M total threads per SM and only 1 block per SM, and the total number of SMs is M (i.e. a maximum of $M*M$ threads can exist in the whole GPU at the same time). Justify your answer.

In general, more threads per block means more warps. This means more opportunities for coalescing and hence potential reduction in bandwidth.