

CMU 15-418/618: Exam 1 Practice Exercises

1 Miscellaneous Short Problems

- A. (2 pts) Modern processors have shifted increasingly toward higher core counts rather than simply relying upon faster, more powerful single cores. Which of the following are contributing factors to this industry-wide pivot? (Check all that apply)
- (a) The “thermal wall” resulting from increasingly miniaturized transistors leaking too much current and generating too much heat to avoid damaging themselves
 - (b) Programmers demanding a shift in paradigm and developing new parallel programming techniques compelling hardware designers to keep up
 - (c) Instruction-Level Parallelism (ILP) nearing the practical limit of its potential benefit, due to dependencies in typical instruction streams
 - (d) Improved cross-core interplaning contention management, decreasing the cost of coordination between cores within the same processor.
 - (e) None of the above
- B. Single Instruction, Multiple Data (SIMD) programming includes (Check all that apply):
- (a) SSE instructions
 - (b) Cuda
 - (c) ISPC programming
 - (d) MPI programming
 - (e) Parallelism by Posix threads
 - (f) None of the above
- C. (2 pts) Define latency (Check all that apply):
- (a) The time between the dispatch of an operation and when it starts executing
 - (b) The number of operations that can be completed per unit time
 - (c) The number of operations that can be executed concurrently.
 - (d) The amount of time needed for an operation
 - (e) None of the above
- D. (2 pts) Arithmetic intensity is commonly used as a measure of (Check all that apply):
- (a) The amount of I/O used per computation achieved
 - (b) The ratio of arithmetic instructions to control instructions
 - (c) The amount of computation required to achieve I/O
 - (d) The ratio of I/O instructions to computation instructions
 - (e) None of the above

- E. (2 pts) When developing a parallel solution, one should first develop and benchmark a serial solution. Why? (Check all that apply)
- (a) To provide the simplest path to a correct implementation by enabling easier debugging
 - (b) To provide the fastest development of a baseline against which correctness of more complex solutions can be measured.
 - (c) To provide a baseline, against which to measure performance gains
 - (d) To enable early profiling, so parallelization and optimization can focus on the most impactful parts of the program
 - (e) To prove the exact algorithm that will be used for final parallelization
 - (f) None of the above
- F. (2 pts) Define “cache coherency” (Check only one)
- (a) The discipline that defines what values are read by the processor cores and when, for example, in light of concurrency in a multi-core system with per-core caches.
 - (b) The property such that multiple caches contain the same set of items
 - (c) The property such that multiple caches have the same hit/miss pattern
 - (d) The property such that main memory and any caches contain the same value
 - (e) None of the above.
- G. (2 pts) What is meant by the “inclusion property” of caches? (Check all that apply)
- (a) All values in caches are the same as values in main memory
 - (b) All caches at the same level, e.g. all level-1 caches, contain the same set of values
 - (c) All lines in a higher level cache (closer to processor) are also in lower level caches (farther from processor) caches within the same hierarchy
 - (d) None of the above
- H. (2 pts) What concern is/are raised by “false sharing”? (Check all that apply)
- (a) An item in cache can become stale due to an aliased update
 - (b) The hit rate is lower than it would be absent the “false sharing”, since useful information is removed from a cache
 - (c) Memory bandwidth is wasted as compared to the same situation without “false sharing”, since values that would otherwise be in the cache need to be read into the cache
 - (d) None of the above
- I. (2 pts) “Relaxed consistency” models address which concern(s)? (Check all that apply)
- (a) Performance being limited by strict ordering of memory operations
 - (b) Software complexity being increased by the need for unnecessary synchronization
 - (c) Software architecture being complicated by the need for synchronization libraries
 - (d) None of the above

2 Buying a New Computer

You write a bit of ISPC code that modifies an grayscale image with width 32 and height `height`. The modification is controlled by the contents of a black and white “mask” image of the same size. The code brightens input image pixels by a factor of 1000 if the corresponding pixel of the mask image is white (the mask has value 1.0) and by a factor of 10 otherwise.

The code partitions the image processing work into 64 ISPC tasks, which you can assume balance perfectly onto all available CPU processors.

```
void brighten_image(uniform int height, uniform float image[], uniform float mask_image[])
{
    uniform int NUM_TASKS = 64;
    uniform int rows_per_task = height / NUM_TASKS;
    launch[NUM_TASKS] brighten_chunk(rows_per_task, image, mask_image);
}

void brighten_chunk(uniform int rows_per_task, uniform float image[], uniform float mask_image[])
{
    // 'programCount' is the ISPC gang size.
    // 'programIndex' is a per-instance identifier between 0 and programCount-1.
    // 'taskIndex' is a per-task identifier between 0 and NUM_TASKS-1

    // compute starting image row for this task
    uniform int start_row = rows_per_task * taskIndex;

    // process all pixels in a chunk of rows
    for (uniform int j=start_row; j<start_row+rows_per_task; j++) {
        for (uniform int i=0; i<32; i+=programCount) {

            int idx = j*32 + i + programIndex;
            int iters = (mask_image[idx] == 1.f) ? 1000 : 10;

            float tmp = 0.f;
            for (int j=0; j<iters; j++)
                tmp += image[idx];

            image[idx] = tmp;
        }
    }
}
```

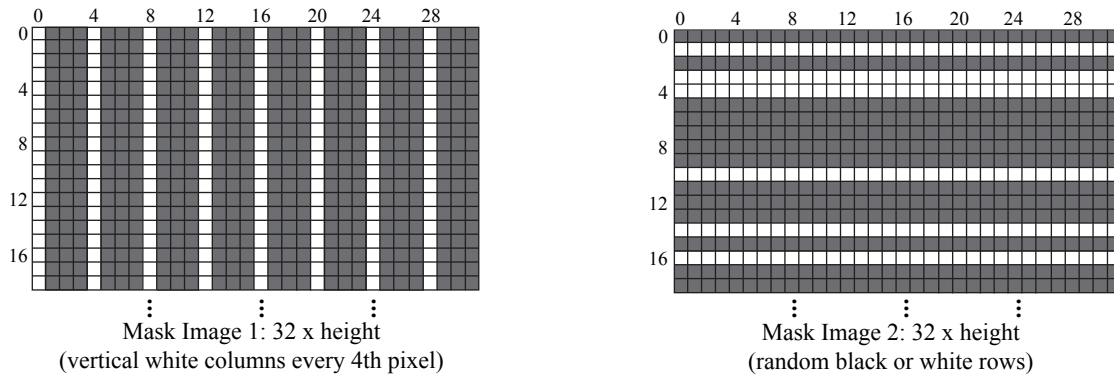


Figure 1: Image masks used to govern image manipulation by `brighten_image`

You go to the store to buy a new CPU that runs this computation as fast as possible. On the shelf you see the following three CPUs on sale for the same price:

- (A) 4 GHz *single core* CPU capable of performing one floating point addition per clock (no parallelism)
- (B) 1 GHz *single core* CPU capable of performing one 32-wide SIMD floating point addition per clock
- (C) 1 GHz *dual core* CPU capable of performing one 4-wide SIMD floating point addition per clock

A. If your only use of the CPU will be to run the above code as fast as possible, and assuming the code will execute using mask image 1 above, rank all three machines in order of performance (from best to worst). Please explain how you determined your ranking by comparing execution times on the various processors. When considering execution time, you may assume that (1) the only operations you need to account for are the floating-point additions in the innermost loop. (2) The ISPC gang size will be set to the SIMD width of the CPU. (3) There are no stalls during execution due to data access.

(Hint: it may be easiest to consider the execution time of each row of the image.)

B. Rank all three machines in order of performance for mask image 2? Please justify your answer, but you are not required to perform detailed calculations like in part A.

3 Buying a New Computer, Again

You plan to port the following sequential C++ code to ISPC so you can leverage the performance benefits of modern parallel processors.

```
float input[LARGE_NUMBER];
float output[LARGE_NUMBER];
// initialize input and output here ...

for (int i=0; i<LARGE_NUMBER; i++) {
    int iters;
    if (i % 16 == 0)
        iters = 256;
    else
        iters = 8;
    for (int j=0; j<iters; j++)
        output[i] += input[i];
}
```

Before sitting down to hack, you go the store, and see the following CPUs all for the same price:

- 4 GHz single core CPU capable of performing one floating point addition per clock (no parallelism)
- 1 GHz quad-core CPU capable of performing one 4-wide SIMD floating point addition per clock
- 1 Ghz dual-core CPU capable of performing one 16-wide SIMD floating point addition per clock

If your only use of the CPU will be to run your future ISPC port of the above code as fast as possible, which machine will provide the best performance for your money? Which machine will provide the least? Please explain why by comparing expected execution times on the various processors. When considering execution time, you may assume that (1) the only operations you need to account for are the floating-point additions in the innermost loop. (2) the ISPC gang size will be set to the SIMD width of the CPU.

(Hint: consider the execution time of groups of 16 elements of the input and output arrays).

4 Angry Students

Your friend is developing a game that features a horde of angry students chasing after professors for making long exams. Simulating students is expensive, so your friend decides to parallelize the computation using one thread to compute and update the student's positions, and another thread to simulate the student's angriness. The state of the game's N students is stored in the global array `students` in the code below).

```
struct Student {
    float position; // assume 1D position for simplicity
    float angriness;
};

Student students[N];

////////////////////////////////////

void update_positions() {
    for (int i=0; i<N; i++) {
        students[i].position = compute_new_position(i);
    }
}

void update_angriness() {
    for (int i=0; i<N; i++) {
        students[i].angriness = compute_new_angriness(i);
    }
}

////////////////////////////////////

// ... initialize students here

pthread_t t0, t1;
pthread_create(&t0, NULL, updatePositions, NULL);
pthread_create(&t1, NULL, updateAngriness, NULL);
pthread_join(t0, NULL);
pthread_join(t1, NULL);
```

- A. Since there is no synchronization between thread 0 and thread 1, your friend expects near a perfect $2\times$ speedup when running on two-core processor that implements invalidation-based cache coherence. She is shocked when she doesn't obtain it. Why is this the case? (For this problem assume that there is sufficient bandwidth to keep two cores busy – “the code is bandwidth bound” is not an answer we are looking for.)

- B. Modify the program to correct the performance problem. You are allowed to modify the code and data structures as you wish, **but you are not allowed to change what computations are performed by each thread and your solution should not substantially increase the amount of memory used by the program.** You only need to describe your solution in pseudocode (compilable code is not required).

5 Oh, the Students Remain Angry

Due to the great success of the hit iPhone app “Angry Students”, your professors decide to release “Angry Students 2: They are Still Angry”, which uses ISPC to take advantage of the SIMD instructions on the iPhone’s ARM processor. The code is written like this:

```
struct Student {
    float position;
    float angriness;
};

Student students[N];

// ispc function
void updateStudents(int N, Student* students) {
    foreach (i = 0 ... N) {
        students[i].position = compute_new_position(i);
        students[i].angriness = compute_new_angriness(i);
    }
}
```

Performance is lower than expected, so the professors change the code to this:

```
float positions[N];
float angriness[N];

// ispc function
void updateStudents(int N, float* positions, float* angriness) {
    foreach (i = 0 ... N) {
        position[i] = compute_new_position(i);
        angriness[i] = compute_new_angriness(i);
    }
}
```

The resulting code runs significantly faster. Why?

6 Parallel Histogram Generation

Your friend implements the following parallel code for generating a histogram from the values in a large input array `input`. For each element of the input array, the code uses the function `bin_func` to compute a “bin” the element belongs to (`bin_func` always returns an integer between 0 and `NUM_BINS-1`), and increments a count of elements in that bin. His port targets a small parallel machine with only two processors. *This machine features 64-byte cache lines and uses an invalidation-based cache coherence protocol.* Your friend’s implementation is given below.

```
float input[N]; // assume input is initialized and N is a very large
int histogram_bins[NUM_BINS]; // output bins
int partial_bins[2][NUM_BINS]; // assume bins are initialized to 0
// assume partial_bins is 64-byte aligned

////////// Code executed by thread 0 //////////
for (int i=0; i<N/2; i++)
    partial_bins[0][bin_func(input[i])]++;

barrier(); // wait for both threads to reach this point

for (int i=0; i<NUM_BINS; i++)
    histogram_bins[i] = partial_bins[0][i] + partial_bins[1][i];

////////// Code executed by thread 1 //////////
for (int i=N/2; i<N; i++)
    partial_bins[1][bin_func(input[i])]++;

barrier(); // wait for both threads to reach this point
```

- A. Your friend runs this code on an input of 1 million elements ($N=1,000,000$) to create a histogram with eight bins ($NUM_BINS=8$). He is shocked when his program obtains far less than a linear speedup, and glumly asserts believe he needs to completely restructure the code to eliminate load imbalance. You take a look and recommend that he not do any coding at all, and just create a histogram with 16 bins instead. Explain why.
- B. Inspired by his new-found great performance, your friend concludes that more bins is better. He tries to use the provided code from part A to compute a histogram of 10,000 elements with 2,000 bins. He is shocked when the speedup obtained by the code drops. Improve the existing code to scale near linearly with the larger number of bins. (Please provide pseudocode as part of your answer – it need not be compilable C code.)
- C. Your friend changes `bin_func` to a function with *extremely high arithmetic intensity*. (The new function requires 100000’s of instructions to compute the output bin for each input element). If the histogram code **provided in part A** is used with this new `bin_func` do you expect scaling to be better, worse, or the same as the scaling you observed using the old `bin_func` in part A? Why? (Please ignore any changes you made to the code in part B for this question.)

7 Reduction with CUDA

You want to write code to sum all of the elements in a vector of length N . You consider two options: *binary* reduction and *square* reduction.

The kernel for binary reduction reduces the size of the vector by a factor of 2 on each step:

```
// Parameter N2 = ceil(N/2.0)
__global__ void
binaryReduceKernel(int N, int N2, float *src, float *dst) {
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    float val;
    if (i < N2) {
        val = src[i];           // 1 cycle
        if (i+N2 < N)
            val += src[i+N2]; // 1 cycle
        dst[i] = val;         // 1 cycle
    }
}
```

The kernel for square reduction treats the vector as an $M \times M$ matrix where $M = \lceil \sqrt{N} \rceil$, and reduces the vector to M elements by summing each row in the array:

```
// Generate position of matrix element i,j
#define RM(i,j,W) ((i)*(W)+(j))

// Parameter M = ceil(sqrt(N))
__global__ void
squareReduceKernel(int N, int M, float *src, float *dst) {
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    float val = 0.0;
    for (int j = 0; j < M; j++) {
        int idx = RM(i,j,M);
        if (idx < N)
            val += src[idx]; // 1 cycle
    }
    if (i < M)
        dst[i] = val;       // 1 cycle
}
```

The general scheme for both kernels is to apply them repeatedly until it becomes better to do a serial summation:

```
for (; N >= NMIN; N = M) {
    M = next(N);           // Either ceil(N/2.0) or ceil(sqrt(N))
    ... reduce vec from N to M ...
}
return serialReduce(vec, N); // N cycles
```

- A. What is the minimum value that can be used for NMIN for performing binary reduction? (and still get a correct result)

- B. What is the minimum value that can be used for NMIN for performing square reduction? (and still get a correct result)

C. In experimenting with the square reduction code, you replace the computation $\text{idx} = \text{RM}(i, j, M)$ with $\text{idx} = \text{RM}(j, i, M)$, so that it computes the sum of each column. The new code runs faster. Explain how this could be.

D. Imagine a hypothetical GPU with the following properties:

- It has one warp of 32 functional units operating in SIMD mode.
- A thread launch by the host requires 20 cycles.
- A memory read or write by the entire warp requires 1 cycle.
- Copying an array of size K from the device to the host and computing the sum of its elements serially on the host requires K cycles.
- All other operations require no cycles.

(a) What would be the cost incurred by reducing a vector of length $N = 1024$ using square reduction, for an optimal setting of parameter N_{MIN} ? (**By optimal we mean the stopping criterion that yields the highest performance.**) Show your work by listing a sequence of steps, describing what would happen with each step and how many cycles it would require.

(b) What would be the cost incurred by reducing a vector of length $N = 1024$ using binary reduction, for an optimal setting of parameter N_{MIN} ? Show your work by listing a sequence of steps, describing what would happen with each step and how many cycles it would require.