

Stanford CS149: Parallel Computing Written Assignment 2

Miscellaneous Short Problems

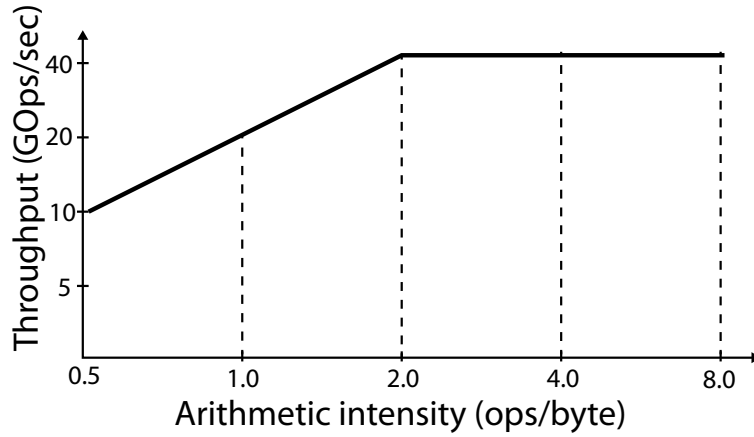
Problem 1:

- A. Your friend is designing a soccer playing robot for the next RoboCup competition. The software on the robot must send a torque request to the robot's motors every 1 ms in order for the robot to successfully locomote. Unfortunately in your friend's implementation, computing torques takes 10 ms when running serially on a single core of the robot's 1 GHz single-core CPU. As a result your friend's robot constantly falls over without being touched. You laugh at your friend, and call their robot "Neymar"!

You look at your friend's code and notice that 20% of the code is inherently serial. The rest is perfectly parallelizable. You dig into your desk drawer and find two processors: Processor A is a 32-core processor that runs at 1 GHz. Processor B is a 8-core processor that runs at 4 GHz. Can you solve your friend's performance problem? Justify your answer by computing the maximum performance they can achieve.

- B. When we discussed Cilk, we emphasized how `cilk_spawn foo()` differs from a normal C function call `foo()` in that the Cilk call can run asynchronously with the caller. Notice that Cilk doesn't explicitly state that the callee function runs **in parallel with the caller**. Give one reason why the designers of Cilk intentionally designed a language that does not specify when the call will run relative to the caller?

- C. The figure below shows a “roofline” plot for a specific processor. Each point on the graph corresponds to the performance (the Y axis is ops/sec) of a different program with the given arithmetic intensity (the X-axis gives the arithmetic ops per byte read by the program. **What is the bandwidth available to the processor? Please give a number, and a 1-2 sentence explanation of your answer.**
- Hint: consider why the curve slants diagonally upward in the left part of the graph. Why does the curve flatten out on the right side?**



- D. Your friend suspects that their program is suffering from high communication overhead, so to overlap the sending of multiple messages, they try to change their code to use asynchronous, non-blocking sends instead of synchronous, blocking sends. The result is this code (assume it is run by thread 1 in two-thread program).

```
float mydata[ARRAY_SIZE];
int dst_thread = 2;

update_data_1(mydata); // updates contents of mydata
async_send(dst_thread, mydata, sizeof(float) * ARRAY_SIZE);

update_data_2(mydata); // updates contents of mydata
async_send(dst_thread, mydata, sizeof(float) * ARRAY_SIZE);
```

Your friend runs to you to say “my program no longer gives the correct results.” What is their bug?

Another exercise on understanding the behavior of caches!

Problem 2:

Consider the following C code run by one thread on a CPU featuring a cache with 64 byte (16 float) cache lines and a total capacity of 128 KB. Assume the cache employs a least recently used (LRU) eviction policy and it is fully associative (any cache line can be placed in any available location in the cache). **The processor does not perform any data pre-fetching.**

```
const int SIZE = 1024 * 1024 * 1024;
float input[SIZE];          // sizeof(float) * SIZE bytes
float output[SIZE];        // sizeof(float) * SIZE bytes

// assume input is appropriately initialized here

for (int i=0; i<SIZE; i++) {
    output[i] = input[i] * 2.0;
}
```

- A. In the for loop, what fraction of loads from input are cache misses? Briefly explain why? (**Hint: Please keep in mind the cache line size.**)

- B. Imagine the code is changed in the following way to have a nested set of for loops:

```
const int SIZE = 8 * 1024;
float input[SIZE];        // sizeof(float) * SIZE bytes
float output[SIZE];      // sizeof(float) * SIZE bytes

// assume input is appropriately initialized here

for (int j=0; j<10000; j++) {
    for (int i=0; i<SIZE; i++) {
        output[i] = output[i] + input[i] * 2.0;
    }
}
```

Assume the processor's cache is the same as in part A. In the "common case" where $j > 0$, what fraction of accesses to input and output (both loads and stores) are cache misses? Briefly explain why?

C. Now consider the following code which works on 2D arrays:

```
const int SIZE_1 = 1024 * 1024 * 1024;
const int SIZE_2 = 1024 * 8;

float input[SIZE_1][SIZE_2];
float output[SIZE_1-1][SIZE_2];

// assume input is appropriately initialized here

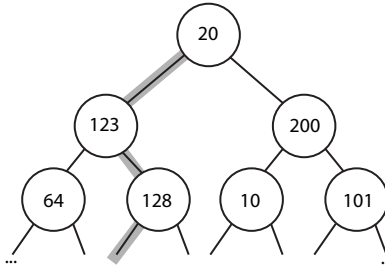
for (int j=0; j<SIZE_1-1; j++) {
    for (int i=0; i<SIZE_2; i++) {
        output[j][i] = input[j+1][i] + input[j][i];
    }
}
```

Now assume that the cache still has the 64 byte line size as before, but **is now only 32 KB in size**. Please rewrite the code so that you minimize the number of cache misses. You can only modify the loop structure (you can add/remove loops and adjust loop bounds) and modify array indexing, but you cannot change the number of mathematical operations performed. **Pseudocode is fine, it need not compile. But to help the grader understand your code, please briefly (in a sentence or two) describe the rough approach used in your solution.**

Running CUDA Code on a GPU

Problem 3:

Consider a complete binary tree of depth 16 that holds one floating point number at each node as shown below. (The figure shows up to depth two.)



Now imagine a CUDA program where each CUDA thread computes the sum of results obtained by applying the functions $f()$ or $g()$ to all numbers on a path from the root to a leaf. The path taken through the tree is determined by the thread's id, as given in the code below. (In the figure above we highlight the path for thread id 2 which in binary is ...00000010 or left-right-left-left-left...)

```
struct Node {
    Node *left, *right;
    float value;
};

void traverse(Node* root, float* output) {
    int threadId = blockDim.x * blockIdx.x + threadIdx.x; // compute 1D thread id
    float sum = 0.0;
    int pathBits = threadId;
    int depth = 0;
    Node* curNode = root;
    while (curNode != NULL) {

        // Consider this a single load of 12 bytes
        // Assume processor doesn't use arithmetic cycles to issue loads
        float val = curNode->value; // 4 bytes
        Node* left = curNode->left; // 4 bytes
        Node* right = curNode->right; // 4 bytes

        if (depth % 2 == 0)
            sum += f(val); // 7 arithmetic instructions
        else
            sum += g(val); // 7 arithmetic instructions

        // ***** count the lines below as 3 arithmetic instructions
        curNode = (pathBits & 1) ? right : left; // *** line "A" ***
        pathBits >> 1; // shift right by 1 bit
        depth++;
    }

    output[threadId] = sum; // each thread writes its result
}
```

Questions are on the next page.

Assume that the program runs on a GPU with a SIMD width (aka CUDA warp size) of 32. Does the program suffer from low utilization due to SIMD divergence, why or why not? (For simplicity, please assume that the conditional “?” operator in line A is a single statement with no divergence.)

An Exercise in Data-Parallel Thinking

Problem 4:

Assume you are given a library that can execute a bulk launch of N independent invocations of an application-provided function using the following CUDA-like syntax:

```
my_function<<<N>>>(arg1, arg2, arg3...);
```

For example the following code would output: (id is a built-in id for the current function invocation)

```
void foo(int* x) {
    printf("Instance %d : %d\n", id, x[id]);
}
int A[] = {10,20,30}
foo<<<3>>>(A);
```

```
"Instance 0 : 10"
"Instance 1 : 20"
"Instance 2 : 30"
```

The library also provides the data-parallel function `exclusive_scan` (using the + operator) that works as discussed in class.

```
exclusive_scan(N, in, out);
```

Example usage:

```
N      = 6
in     = {1, 2, 3, 4, 5, 6}
=====
out    = {0, 1, 3, 6, 10, 15}
```

In this problem, we'd like you to design a data-parallel implementation of `largest_segment_size()`, which, given an array of flags that denotes a partitioning of an array into segments, computes the size of the longest segment in the array.

```
int largest_segment_size(int N, int* flags);
```

The function takes as input an array of N flags (flags) (with 1's denoting the start of segments), and returns the size of the largest segment. The first element of flags will always be 1. For example, the following flags array describes five segments of lengths 4, 2, 2, 1, and 1.

```
N      = 10
flags  = {1, 0, 0, 0, 1, 0, 1, 0, 1, 1}
=====
result: = 4
```

Questions on next page...

- A. The first step in your implementation should be to compute the size of each segment. Please use the provided library functions (bulk launch of a function of your choice + `exclusive_scan` to implement the function `segment_sizes()` below. *Hint: We recommend that you get a basic solution done first, then consider the edge cases like how to compute the size of the last segment.*

```
// Example output of segment_sizes(N, flags, num_segs, sizes):  
// N = 8  
// flags = {1, 0, 1, 0, 0, 0, 1, 0}  
// =====  
// num_segs = 3  
// sizes = {2, 4, 2}
```

```
// you may wish to define functions used in bulk launches here
```

```
// You can allocate any required intermediate arrays in this function  
// You may assume that 'seg_sizes' is pre-allocated to hold N elements,  
// which is enough storage for the worse case where the flags array  
// is all 1's.  
void segment_sizes(int N, int* flags, int* num_segs, int* seg_sizes) {
```

```
}
```

B. Now implement `largest_segment_size()` using `segment_sizes()` as a subroutine. **NOTE: this problem can be answered even without a valid answer to Part A.** Your implementation may assume that the number of segments described by `flags` is always a power of two. A full credit implementation will maximize parallelism and minimize work when computing the maximum segment size from an array of segment sizes. *Hint: we are looking for solutions with $\lg_2(\text{num_segs})$ span.*

// you may want to implement helper functions here that are called via bulk launch

```
int largest_segment_size(int N, int* flags) {  
  
    int num_segs;  
    int seg_sizes[N];  
    segment_sizes(N, flags, &num_segs, seg_sizes);
```

```
}
```


- C. The professors start getting anxious because the CAs haven't completed grading, so they send an angry email to the staff. "Let's speed it up already!" they write. The CA's look at Pranil and Drew and, say "Quit surfing the internet reading articles about ML accelerators and please come help!" Assuming that Pranil and Drew grade questions at exactly the same speed as the other CAs, which question should they help with grading? (please choose one) Describe why? What is the new **steady state throughput** of the staff in terms of exams per hour? (Regardless of the question chosen, Assume that Pranil/Drew's help is going to come in the form of grading a different exam in parallel with the other CAs working on the same question. Their help doesn't reduce the amount of time it takes to grade one question on one exam.)

Misc Problems

PRACTICE PROBLEM 2:

- A. A key idea in this course is the difference between *abstraction* and *implementation*. Consider two abstractions we've studied: ISPC's `foreach` and Cilk's `spawn` construct. **Briefly describe how these two abstractions have similar semantics.** (Hint: what do the constructs declare about the associated loop iterations? Be precise!). **Then briefly describe how their implementations are quite different** (Hint: consider their mapping to modern CPUs). As a reminder, we give you two syntax examples below:

ISPC `foreach`:

=====

```
foreach (i = 0 ... 100) {  
    x[i] = y[i];  
}
```

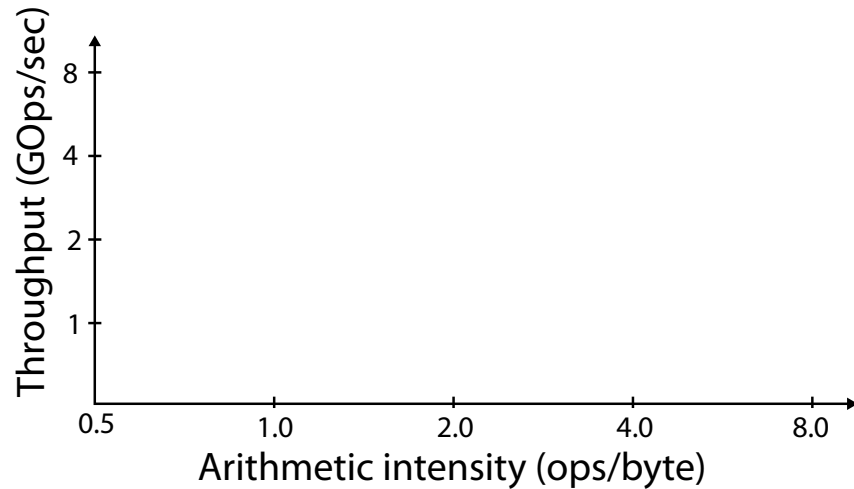
Cilk:

=====

```
void f(int i, float* x, float* y) {  
    x[i] = y[i]  
}
```

```
for (int i=0; i<100; i++) {  
    cilk_spawn f(i, x, y);  
}
```

- B. In class we described the usefulness of making roofline graphs, which plots the instruction throughput of a machine (gigaops/sec) as a function of a program's arithmetic intensity (ops performed per byte transferred from memory). Note moving along the X axis is changing the properties of the code being run. The Y axis plots the performance of the machine when running a specified program. Consider the roofline plot below. Please plot the roofline curve for a machine featuring a **1 GHz dual-core processor. Each core can execute one 4-wide SIMD instruction per clock.** This processor is connected to a memory system providing 4 GB/sec of bandwidth. *Hint: what is the peak throughput of this processor? What are its bandwidth requirements when running a piece of code with a specified arithmetic intensity? Recall ops/second \times bytes/op is bytes/sec. Arithmetic intensity is 1/(bytes/op).* Plot the expected throughput of the processor when running code at each arithmetic intensity on the X axis, and draw a line between the points.



- C. Consider a cache that contains 32 KB of data, has a cache line size of 4 bytes, is fully associative (meaning any cache line can go anywhere in the cache), and uses an LRU (least recently used—the line evicted is the line that was last accessed the longest time ago) replacement policy. Please describe why the following code will take a cache miss on every data access to the array A.

```
const int SIZE = 1024 * 64;
float A[SIZE];
float sum = 0.0;
for (int reps=0; reps<32; reps++)
    for (int i=0, i<SIZE; i++)
        sum += A[i];
```

D. Consider the following piece of C code.

```
float A[VERY_LARGE];
float B[VERY_LARGE];
float C[VERY_LARGE];
float D[VERY_LARGE];
float E[VERY_LARGE];

for (int i=0; i<VERY_LARGE; i++)
    C[i] = A[i] * B[i];
for (int i=0; i<VERY_LARGE; i++)
    D[i] = C[i] + B[i];
for (int i=0; i<VERY_LARGE; i++)
    E[i] = D[i] - A[i];
```

Assume that VERY_LARGE is so large that the arrays are hundreds of MBs in size, and that the code is run on a single-core processor with a 8 MB cache. Please modify the program to maximize its arithmetic intensity. You only need to write to the output array E, you don't need to fill in C and D if it is not necessary. However, please **DO NOT CHANGE** the number of math operations performed. **If we assume the program before and after the modification is bandwidth bound, how much does your modification improve its performance?**

A Barrier is Worth a 1000 Locks

PRACTICE PROBLEM 3:

Consider the following code written in an SPMD style. Note this is not ISPC code. It's C-like code, but assume that N threads are running the code. The threads cooperate to compute a histogram for the values in an input array `data`. You can assume that `data` contains random numbers between 0 and 100, the histogram has 10 bins, and that `bins[i]` is supposed to contain the count of the number of elements in `data` that fall between $10 \times i$ and $10 \times (i + 1)$

```
// These variables are global variables accessible to all threads.

const int N = VERY_LARGE_NUMBER; // assume N is a very large number
const int NUM_THREADS = 4;
int data[N];
int bins[10]; // assume initialized to 0
Lock myLock;

// This function is run in SPMD fashion by all threads

void run(int threadId) {
    int elsPerThread = N / NUM_THREADS;
    int start = threadId * elsPerThread;
    int end = start + elsPerThread;

    for (int i=start; i<end; i++) {
        int binId = data[i] / 10;
        myLock.lock();
        bins[binId]++;
        myLock.unlock();
    }
}
```

- A. You run the program on a four core processor, and observe that it gets the correct answer, and that work is well distributed among the threads. However you don't observe a great speedup compared to a single threaded version of the code. What is a potential significant performance problem?

- B. Imagine that instead of locks, you are allowed to use a single `barrier()` in the code. Please give a solution that yields good work distribution onto all 4 threads, uses no locks, and uses only a single call to `barrier()`. Your solution is allowed to allocate new global or per-thread variables. **Hint: Keep in mind that N is assumed to be much, much larger than the number of bins in the histogram.**

```
const int N = VERY_LARGE_NUMBER; // assume N is a very large number
const int NUM_THREADS = 4;
int data[N];
int bins[10]; // assume initialized to 0
```

```
void run(int threadId) {

    int elsPerThread = N / NUM_THREADS;
    int start = threadId * elsPerThread;
    int end = start + elsPerThread;

    for (int i=start; i<end; i++) {

        int binId = data[i] / 10;

    }

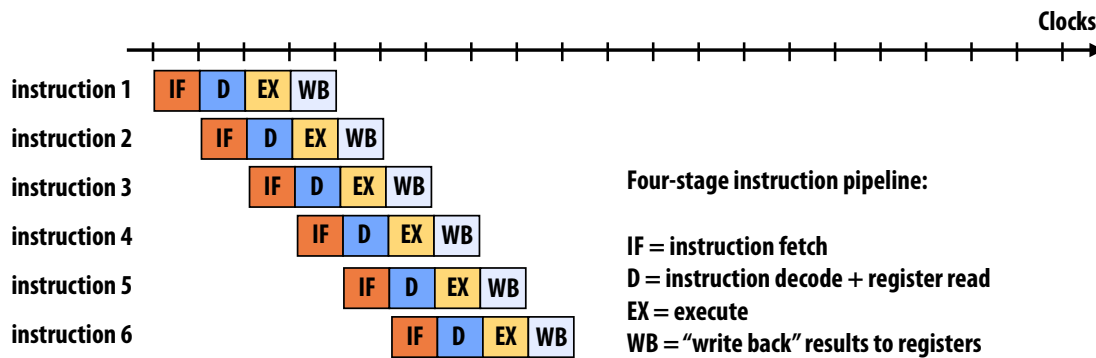
}
```

A Cardinal Processor Pipeline

PRACTICE PROBLEM 4:

The fast-growing startup Cardinal Processors, Inc. builds a single core, single threaded processor that executes instructions using a simple four-stage pipeline. As shown in the figure below, each unit performs its work for an instruction **in one clock**. To keep things simple, assume this is the case for all instructions in the program, including loads and stores (memory is infinitely fast).

The figure shows the execution of a program with six **independent instructions** on this processor. *However, if instruction B depends on the results of instruction A, instruction B will not begin the IF phase of execution until the clock after WB completes for A.*



- A. Assuming all instructions in a program are **independent** (yes, a bit unrealistic) what is the instruction throughput of the processor?

- B. Assuming all instructions in a program are **dependent** on the previous instruction, what is the instruction throughput of the processor?

- C. What is the latency of completing an instruction?

- D. Imagine the IF stage is modified to improve its throughput to fetch TWO instructions per clock, but no other part of the processor is changed. What is the new overall maximum instruction throughput of the processor?

E. Consider the following C program:

```
float A[500000];
float B[500000];
// assume A is initialized here

for (int i=0; i<500000; i++) {
    float x1 = A[i];
    float x2 = 6 * x1;
    float x3 = 4 + x2;
    B[i] = x3;
}
```

Assuming that we consider only the four instructions in the loop body (for simplicity, disregard instructions for managing the loop or calculating load/store addresses), what is the average instruction throughput of this program? (Hint: You should probably consider instruction dependencies, and at least two loop iterations worth of work).

F. Modify the program to achieve peak instruction throughput on the processor. Please give your answer in C-pseudocode.

- G. Now assume the program is reverted to the original code from part E, but the for loop is parallelized using OpenMP. (Recall from written assignment 1 is that openMP is a set of C++ compiler extensions that enable thread-parallel execution. Iterations of the for loop will be carried out in parallel by a pool of worker threads.)

```
// assume iterations of this FOR LOOP are parallelized across multiple
// worker threads in a thread pool.
#pragma omp parallel for
for (int i=0; i<1000000; i++) {
    float x1 = A[i];
    float x2 = 2*x1;
    float x3 = 3 + x2;
    B[i] = x3;
}
```

Given this program, imagine you wanted to add multi-threading to the **single-core processor** to obtain **peak instruction throughput** (100% utilization of execution resources). What is the smallest number of threads your processor could support and still achieve this goal? You may not change the program.

Particle Simulation

PRACTICE PROBLEM 5:

Consider the following code that uses a simple $O(N^2)$ algorithm to compute forces due to gravitational interactions between all N particles in a particle simulation. One important detail of this algorithm is that force computation is symmetric ($\text{gravity}(i, j) = \text{gravity}(j, i)$). Therefore, iteration i only needs to compute interactions with particles with index j , where $i < j$. As a result, there are $N^2/2$ calls to gravity rather than N^2 .

In this problem, **assume the code is run on a dual-core processor, with infinite memory bandwidth. The processor implements invalidation-based cache coherence across the cores. The cache line size is 64 bytes.**

```
struct Particle {
    float force; // for simplicity, assume force is represented as a single float
    Lock l;
};

Particle particles[N];

void compute_forces(int threadId) {

    // thread 0 takes first half, thread 1 takes second half
    int start = threadId * N/2;
    int end = start + N/2;

    for (int i=start; i<end; i++) {

        // only compute forces for each pair (i,j) once, then accumulate force
        // into *both* particle i and j

        for (int j=i+1; j<N; j++) {
            float force = gravity(i, j);

            lock(particles[i].l);
            particles[i] += force;
            unlock(particles[i].l);

            lock(particles[j].l);
            particles[j] += force;
            unlock(particles[j].l);
        }
    }
}
```

Question is on the next page...

A. Although the code makes $N^2/2$ calls to `gravity()` it takes N^2 locks. Modify the code so that the number of lock/unlock operations is reduced by $2\times$. You may not allocate additional variables or change how lock iterations are mapped to the threads.

B. **(This questions can be answered independently from part A)** Looking at the original code, there another major performance problem that does not have to do with the number of lock/unlock operations. Please describe the problem and then describe a solution. Clearly describing an implementable solution strategy is fine, you do not need to write precise pseudocode.

Because The Professor with the Most ALUs (Sometimes) Wins

PRACTICE PROBLEM 6:

Consider the following ISPC code that computes $ax^2 + bx + c$ for elements x of an entire input array.

```
void polynomial(float a, float b, float c,
               uniform float x[], uniform float output[], int elementsPerTask) {
    uniform int start = taskIndex * elementsPerTask;
    uniform int end = start + elementsPerTask;

    foreach (i = start ... end) {
        output[i] = (a * x[i] * x[i]) + (b * x[i]) + c;    // 5 arithmetic ops
    }
}

// assume N is very, very large, and is a multiple of 1024
void run(int N, float a, float b, float c, float* input, float* output) {
    uniform int elementsPerTask = 1024;
    launch[N/elementsPerTask] polynomial(a, b, c, input, output, elementsPerTask);
}
```

Professor Kayvon, seeking to capture the highly lucrative polynomial evaluation market, builds a multi-core CPU packed with ALUs. "The professor with the most ALUs wins, he yells!" The processor has:

- 4 cores clocked at 1 GHz, capable of one 32-wide SIMD floating-point instruction per clock (1 addition, 1 multiply, etc.)
- Two hardware execution contexts per core
- A 1 MB cache per core with 128-byte cache lines (In this problem assume allocations are cache-line aligned so that each SIMD vector load or store instruction will load one cache line). Assume cache hits are 0 cycles.
- The processor is connected to a memory system providing a whopping 512 GB/sec of BW
- The latency of memory loads is 95 cycles. (There is no prefetching.) For simplicity, assume the latency of stores is 0.

A. What is the peak arithmetic throughput of Prof. Kayvon's processor?

B. What should Prof. Kayvon set the ISPC gang size to when running this ISPC program on this processor?

C. Prof. Kayvon runs the ISPC code on his new processor, the performance of the code is not good. What fraction of peak performance is observed when running this code? Why is peak performance not obtained?

D. Prof. Olukuton sees Kayvon's struggles, and sees an opportunity to start his own polynomial computation processor company that achieves double the performance of Prof. Kayvon's chip. "Oh shucks, now I'll have to double the number of cores in my chip, that will cost a fortune." Kayvon says.

TA Mario writes Kayvon an email that reads "There's another way to achieve peak performance with your original design, and it doesn't require adding cores." Describe a change to Prof. Kayvon's processor that causes it to obtain peak performance on the original workload. Be specific about how you'd realize peak performance (give numbers).

The following year Prof. Kayvon makes a new version of his processor. The new version is the **exact same quad-core processor** as the one described at the beginning of this question, except now the chip **supports 64 hardware execution contexts per core**. Also, the ISPC code is changed to compute a more complex polynomial. In the code below assume that `coeffs` is an array of a few hundred polynomial coefficients and that `expensive_polynomial` involves 100's of arithmetic operations.

```
void polynomial(uniform float coeffs[], uniform float input[],
               uniform float output[], int elementsPerTask) {
    uniform int start = taskIndex * elementsPerTask;
    uniform int end = start + elementsPerTask;
    foreach (i = start ... end) {
        output[i] = expensive_poly(coeffs, input[i]); // 100's of arithmetic ops
    }
}

void run(int N, float* coeffs, float* input, float* output) {
    uniform int elementsPerTask = 1024;
    launch[N/elementsPerTask] polynomial(coeffs, input, output, elementsPerTask);
}
```

- E. What is the peak arithmetic throughput of Prof. Kayvon's new processor?
- F. Imagine running the program with $N=8 \times 1024$ and $N = 64 \times 1024$. Assuming that the system schedules worker threads onto available execution contents in an efficient manner, do either of the two values of N result in the program achieving near peak utilization of the machine? Why or why not? (For simplicity, assume task launch overhead is negligible.)
- G. Now consider the case where $N=9 \times 1024$. Now what is the performance problem? Describe a simple code change that results in the program obtaining close to peak utilization of the machine. (Assume task launch overhead is negligible.)

Sending Messages

PRACTICE PROBLEM 7:

In class we talked about the `barrier()` synchronization primitive. No thread proceeds past a barrier until all threads in the system have reached the barrier. (In other words, the call to `barrier()` will not return to the caller until its known that all threads have called `barrier()`). Consider implementing a barrier in the context of a message passing program that is only allowed to communicate via **blocking sends and receives**. Using only the helper functions defined below, implement a barrier. Your solution should make no assumptions about the number of threads in the system. **Keep in mind that all threads in a message passing program execute in their own address space—there are no shared variables.**

```
// send msg with id msgId and contents msgValue to thread dstThread
void blockingSend(int dstThread, int msgId, int value);

// recv message from srcThread. Upon return, msgId and msgValue are populated
void blockingRecv(int srcThread, int* msgId, int* msgValue);

// returns the id of the calling thread
int getThreadId();

// returns the number of threads in the program
int getNumThreads();
```