

Stanford CS149: Parallel Computing
Written Assignment 2
Due October 7th

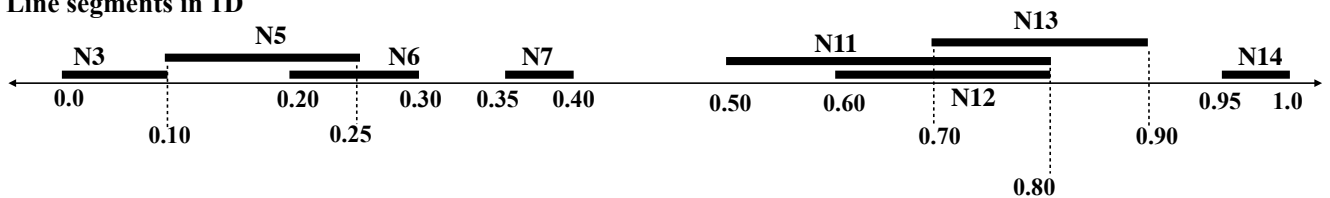
Problem 1: SPMD Tree Search

NOTE: This question is tricky. But don't panic, it's practice and graded on credit! If you can answer this question you really understand SIMD execution!

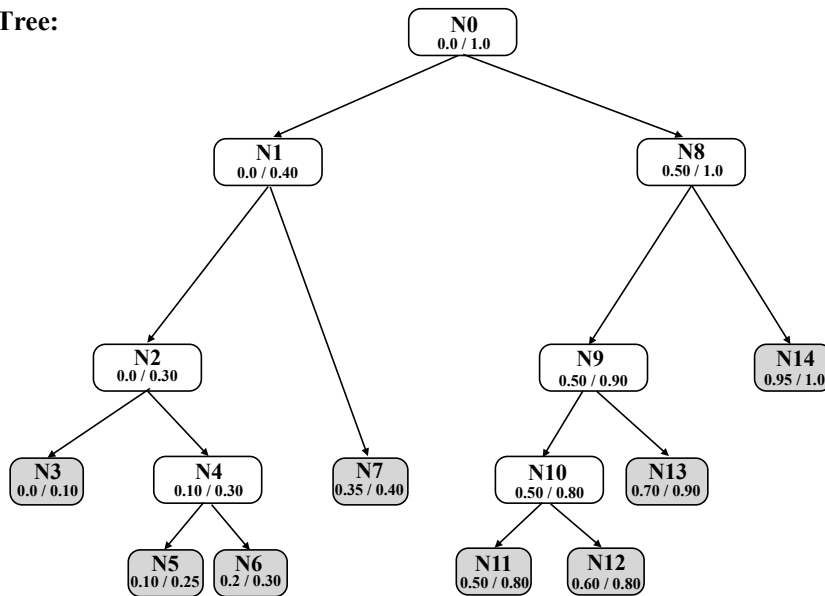
The figure below shows a collection of line segments in 1D. It also shows a binary tree data structure organizing the segments into a hierarchy. Leaves of the tree correspond to the line segments. Each interior tree node represents a spatial extent that bounds all its child segments. Notice that sibling leaves can (and do) overlap. Using this data structure, it is possible to answer the question "what is the largest segment that contains a specified point" without testing the point against all segments in the scene.

For example, the answer for point $p = 0.15$ is segment 5 (in node N5). The answer for the point $p = 0.75$ is segment 11 in node N11.

Line segments in 1D



Binary Search Tree:



```
struct Node {
    float min, max;    // if leaf: start/end of segment, else: bounds on all child segments.
    bool leaf;        // true if nodes is a leaf node
    int segment_id;   // segment id if this is a leaf
    Node* left, *right; // child tree nodes
};
```

On the following two pages, we provide you two ISPC functions, `find_segment_1` and `find_segment_2` that both compute the same thing: they use the tree structure above to find the id of the largest line segment that contains a given query point.

```

struct Node {
    float min, max;    // if leaf: start/end of segment, else: bounds on all child segments.
    bool leaf;        // true if nodes is a leaf node
    int segment_id;   // segment id if this is a leaf
    Node* left, *right; // child tree nodes
};

// -- computes segment id of the largest segment containing points[programIndex]
// -- root_node is the root of the search tree
// -- each program instance processes one query point
export void find_segment_1(uniform float* points, uniform int* results, uniform Node* root_node) {

    Stack<Node*> stack;
    Node* node;
    float max_extent = 0.0;

    // p is point this program instance is searching for
    float p = points[programIndex];
    results[programIndex] = NO_SEGMENT;

    stack.push(root_node);

    while(!stack.size() == 0) {
        node = stack.pop();

        while (!node->leaf) {
            // [I-test]: test to see if point is contained within this interior node
            if (p >= node->min && p <= node->max) {
                // [I-hit]: p is within interior node... continue to child nodes
                push(node->right);
                node = node->left;
            } else {
                // [I-miss]: point not contained within node, pop the stack
                if (stack.size() == 0)
                    return;
                else
                    node = stack.pop();
            }
        }

        // [S-test]: test if point is within segment, and segment is largest seen so far
        if (p >= node->min && p <= node->max && (node->max - node->min) > max_extent) {
            // [S-hit]: mark this segment as "best-so-far"
            results[programIndex] = node->segment_id;
            max_extent = node->max - node->min;
        }
    }
}

```

```

export void find_segment_2(uniform float* points, uniform int* results, uniform Node* root_node) {

    Stack<Node*> stack;
    Node* node;
    float max_extent = 0.0;

    // p is point this program instance is search for
    float p = points[programIndex];

    results[programIndex] = NO_SEGMENT;

    stack.push(root_node);

    while(!stack.size() == 0) {
        node = stack.pop();

        if (!node->leaf) {
            // [I-test]: test to see if point is contained within interior node
            if (p >= node->min && p <= node->max) {
                // [I-hit]: p is within interior node... continue to child nodes
                push(node->right);
                push(node->left);
            }
        } else {
            // [S-test]: test if point is within segment, and segment is largest seen so far
            if (p >= node->min && p <= node->max && (node->max - node->min) > max_extent) {
                // [S-hit]: mark this segment as "best-so-far"
                results[programIndex] = node->segment_id;
                max_extent = node->max - node->min;
            }
        }
    }
}

```

Begin by studying `find_segment_1`.

Given the input $p = 0.1$, the a single program instance will execute the following sequence of steps: (I-test,N0), (I-hit,N0), (I-test, N1), (I-hit, N1), (I-test, N2), (I-hit, N2) (S-test,N3), (S-hit, N3), (I-test, N4), (I-hit, N4), (S-test, N5), (S-hit, N5), (S-test, N6), (S-test,N7), (I-test, N8), (I-miss, N8). Where each of the above "steps" represents reaching a basic block in the code (see comments):

- (I-test, Nx) represents a point-interior node test against node x.
- (I-hit, Nx) represents logic of traversing to the child nodes of node x when p is determined to be contained in x.
- (I-miss, Nx) represents logic of traversing to sibling/ancestor nodes when the point is not contained within node x.
- (S-test, Nx) represents a point-segment (left node) test against the segment represented by node x.
- (S-hit, Nx) represents the basic block where a new largest node is found x.

The question is on the next page...

- A. Confirm you understand the above, then consider the behavior of a **gang of 4 program instances** executing the above two ISPC functions `find_segment_1` and `find_segment_2`. For example, you may wish to consider execution on the following array:

```
points = {0.15, 0.35, 0.75, 0.95}
```

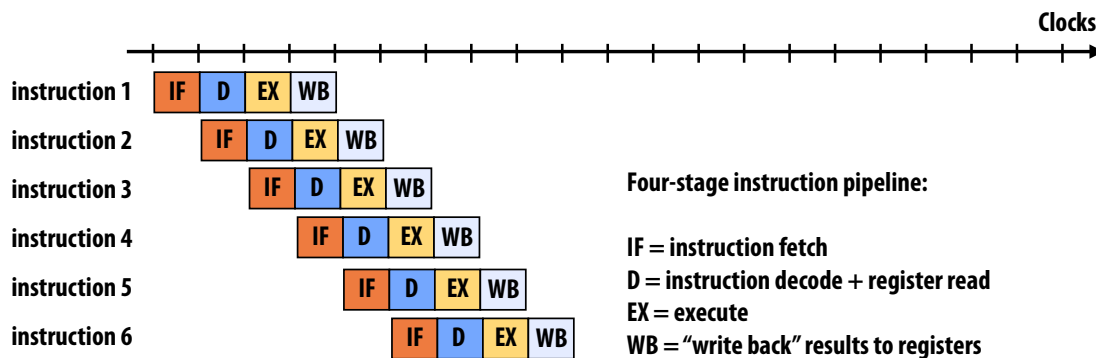
Describe the difference between the traversal approach used in `find_segment_1` and `find_segment_2` in the context of SIMD execution. Your description might want to specifically point out conditions when `find_segment_1` suffers from divergence. (Hint 1: you may want to make a table of four columns, each row is a step by the entire gang and each column shows each program instance's execution. Hint 2: It may help to consider which solution is better in the case of large, heavily unbalanced trees.)

- B. Consider a slight change to the code where as soon as a best-so-far line segment is found (inside [S-hit]) the code makes a call to a **very, very expensive function**. Which solution might be preferred in this case? Why?

Problem 2: A Cardinal Processor Pipeline

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.