

# Problem Collection

## Programming and Architecture of Computing Systems

October 14, 2024

### Contents

<b>Preliminary Notes</b>	<b>1</b>
<b>Small Questions</b>	<b>1</b>
<b>Test Questions</b>	<b>1</b>
<b>Exercises</b>	<b>2</b>

### Preliminary Notes

This brief collection of problems is divided in three parts. The first part covers small questions, the second part includes test questions, and the third part refers to some longer exercises.

To report erratas, typos... please mail either [alvabre@unizar.es](mailto:alvabre@unizar.es), [rgran@unizar.es](mailto:rgran@unizar.es) or [dario@unizar.es](mailto:dario@unizar.es).

### Small Questions

1. Please briefly respond to the following questions: ¿Is a concurrent application always parallel? ¿Is a parallel application always concurrent?
2. According to Amdahl's Law, for a program where the sequential part represents the 15% of the total, what would be the potential speed-up for a 16-core machine.
3. Can a processor execute instructions from two different instruction sets?
4. Enumerate what are the key design features of GPUs to allow a very fast context switching of wavefronts.
5. Explain what the conditional branching problem is on GPUs and how it is solved.
6. Make comparative analysis between a GPU and an ASIC.
7. The iron law of computer performance states that the execution time can be defined as:  $Ex. Time = N \times CPI \times T_{cycle}$ . To improve performance and save energy consumption, a new vector extension has been proposed. The extension reduces both the number of instructions and the frequency by half and 10%, respectively. Since the extra hardware complexity increases the cycles per instruction by 33%, could you please identify which alternative provides the lowest execution time.

### Test Questions

1. The Local Data Share (LDS) cache on a GPU is used to:
  - a) Amplify the regular cache bandwidth
  - b) Execute atomic instructions
  - c) Synchronization of wavefronts
  - d) All of the above

- e) None of the above
- 2. In OpenCL, Local Memory is shared between:
  - a) All the workitems of a global work domain
  - b) Workitems in the same kernel launch
  - c) Local Memory is an abstraction not present in OpenCL
  - d) Workitems in the same workgroup
  - e) None of the previous ones
- 3. In OpenCL, workitems that access global shared variable must explicitly assure memory order in order to avoid race conditions
  - a) Always
  - b) Just in case they do not belong to the same workgroup
  - c) Just in case they do belong to the same workgroup
  - d) Never

## Exercises

1. The dot product algorithm takes two vectors of the same length and returns a single number. The number is the sum of the products of the corresponding entries in the input vectors.

In C++, the algorithm can be coded as follows:

```
template<typename T>
T dot_product(const std::vector<T> &a, const std::vector<T> &b)
{
    if(a.length() != b.length()) {
        error( ... );
    }

    // initialize to 0 regardless the type
    T dot_p{}; // Also T dot_p = T();

    for(size_t i = 0; i < a.length(); ++i) {
        dot_p+=(a[i]*b[i]);
    }
}
```

Please answer the following questions:

- a. Implement the dot product using threads and static partitioning.
- b. Implement the dot product assuming you have the thread pool and the thread-safe queue from Laboratory 4.
- c. For the thread-pool version, would all tasks perform the same amount of work?
2. Given an `std::vector<int>` array, could you please write a parallel algorithm that finds the minimum and maximum values of the array.
3. See Exercise 2 from the collection of exercises referring to metrics.
4. Write an OpenCL program that calculates the dot product of two integer arrays. Additionally to the kernel code, in the host side of the program, just focus on the buffer management, command-queue management and kernel launch.
  - a. Please, analytically model the execution time of this work assuming the computational device has the following characteristics: 8 compute units, each compute unit has 128 parallel cores, each core has two floating-point arithmetic units and, frequency of the computational device is 1.5GHz. Assumption 1: just floating point instructions contribute to the execution time. Assumption 2: each FPU can process a floating point instruction per cycle.

5. Please write a parallel program that given an array of integer values, it finds those values that are prime and larger than a given element. The solution should follow a fork-join parallelism model in C++. To know whether an integer value is prime, you can assume that the function `bool is_prime(int n)` is available:

```
bool is_prime(int n) {
    if (n ≤ 3) {
        return n > 1;
    }

    if (((n % 2) == 0) || ((n % 3) == 0)) {
        return false;
    }

    for(int i = 5; i*i ≤ n; i+=6) {
        return false;
    }
    return true;
}
```

6. Sorting is one of the most important problem in computing. Its computational intensity makes sorting an ideal candidate for parallelization. One of the most common implementation is bucket sort where the input array is split between N buckets that are independently sorted and then concatenated.

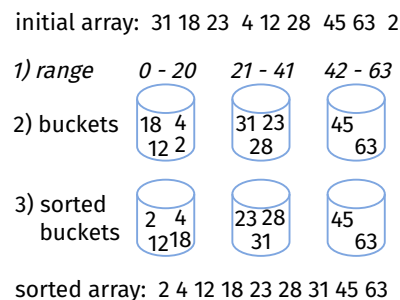


Figure 1: Bucket sort example

The upper figure shows bucket sort main steps: 1) Computation of the ranges. With the maximum value stored in the input array and the number of buckets, you compute the range,  $63/3=21$  in the example, and create the buckets. 2) Bucket insertion. Each value of the array goes to its bucket. 3) Sorting buckets. Each bucket is sorted independently, and 4) the sorted array is assembled by concatenating the arrays of each bucket.

- a. Please implement a parallel version of the bucket sort algorithm. For sorting the buckets, you can use any standard sequential sorting algorithm as

```
template<typename T>
void insertion_sort(std::vector<T>& array)
{
    for(size_t i = 1; i < array.size(); ++i) {
        for(size_t j = i; (j > 0) && (array[j-1] > array[j]); --j) {
            std::swap(array[j], array[j-1]);
        }
    }
}
```

- b. ¿Is there any pathological case where the parallel version could not be faster than the sequential version?

Notes: You can concatenate two `std::vector` arrays with the `insert` method; e.g., `dst.insert(dst.end(), src.begin(), src.end())`. The function `float std::floor(float arg)` computes the largest integer value not greater than `arg`.

7. Many Computer Vision applications require the computation of histograms, which help to understand the distributions of a set of numbers. To compute an histogram, you need to visit all elements of the set, compute their bucket, and then increase the corresponding counter of that bucket. For example, if the

input array contains this set of numbers {0, 1, 1, 1, 2, 3}, the output histogram with 4 buckets will be {1, 3, 1, 1}.

- a. Please write a parallel version of an histogram for integer values following a fork-join approximation in C++. The histogram result will be stored in an array of atomic variables and you have to minimize contention on this array.

You can assume the following initial skeleton:

```
int main() {

    const size_t N = 1024*8; // array size
    const size_t m_buckets = 32; // buckets
    const size_t n_threads = 8;

    std::vector<int> array; // please assume this array has been already initialized
    std::vector<std::atomic<int>> histogram(m_buckets);

    // ...
}
```

- b. What could be the maximum speed-up of this implementation?

8. Matrix multiplication belongs to the most used algorithm list in robotics, graphics, and computer vision applications. Therefore, almost every application requires a fast parallel matrix multiplication algorithm.

Assuming a basic 3 nested loop serial implementation in C++:

```
using fmatrix = matrix<float>;

fmatrix matrix_multiply(const fmatrix &a, const fmatrix &b)
{
    fmatrix c{a.rows(), b.cols()};

    for(size_t i = 0; i < a.rows(); ++i) {
        for(size_t j = 0; j < b.cols(); ++j) {
            float val{0.0f};
            for(size_t k = 0; k < a.cols(); ++k) {
                val += a(i, k) * b(k, j);
            }
            c(i, j) = val;
        }
    }
    return c;
}
```

- a. Using `std::async`, write a parallel version `async_col_matrix_multiply` that when the number of concurrent threads supported by the system is 1, only uses 1 thread. Otherwise, the maximum number of concurrent `std::async` will be the number of columns of a matrix, `a.cols()`.
- b. Imagine you have access to a 1024 multi-core machine, and `a.cols()` is always smaller than 128. What would be the maximum speed-up of the parallel version implemented in step a? Could you please write another version, `asyc_matrix_multiply` able to extract all the possible parallelism for this 1024 multi-core machine.
- c. Assuming that every `std::async` creates a new thread on every invocation, could an implementation based on a thread-pool be faster than the version based on `std::async`? Why?

*Note: You can assume that the `matrix<float>` class provides all requirements for storing matrices. If you need extra trivial methods of the class besides `rows` and `cols`, please feel free to use them without writing their implementation.*

9. Alpha compositing is a computer graphics method that combines a foreground and a background images to simulate transparency. With 2D images, alpha compositing extends each pixel with an additional value representing transparency. This new alpha value ranges between 0, fully transparency, and 1 (fully opaque). For example, assuming two images named  $f$  and  $b$ , so that  $f$  is over  $b$ , in another words,  $f$  is the foreground, the over operator can be computed following these equations:

$$\alpha_o = \alpha_f + \alpha_b(1 - \alpha_f)$$

$$p_o = \frac{p_f \alpha_f + p_b \alpha_b (1 - \alpha_f)}{\alpha_o}$$

Where  $p_x$  represents the three color channels (red, green, blue) of each pixel and  $\alpha_x$  represents the alpha value of the output ( $o$ ), foreground ( $f$ ), and background images ( $b$ ).

Assuming a pixel and image classes as follows:

```
struct pixel {
    public:
        uint8_t red, green, blue;
        uint8_t alpha; // 255 corresponds to opaque
};

template <typename T, size_t N, size_t M>
class image {
    using storage_type = std::array<std::array<T, M>, N>;
    public:
        image(){};
        T& operator()(size_t i, size_t j) {return _array[i][j];};
        T operator()(size_t i, size_t j) const {return _array[i][j];};
    private:
        storage_type _array;
};

const size_t height = 128, width = 128;
using alpha_image = image<pixel, height, width>;
```

- a. (1 point) Please write a sequential version of a `alpha_image` `alpha_over_operator(const alpha_image& f, const alpha_image& b)` free function that returns the result of performing an `alpha_over_operator` on two input images. You can use `std::clamp(uint8_t v, uint8_t lo, uint8_t hi)` to clamp the resulting operations if required.
- b. (2.5 points) Please write a parallel version of `alpha_over_operator` that extracts parallelism and pick between data and task level parallelism depending on the regularity of the problem.