

PAR – Final Exam – Course 2022/23-Q1

January 18th, 2023

Problem 1 (2.5 points)

Given the following code with *tareador* task annotations:

```
#define p ...    // Number of processors
#define NR ...   // Number of rows
#define NC ...   // Number of columns

int M[NR][NC];
int BS = NR/p;  // Assume p divides NR exactly

for (int ii=0; ii<NR; ii+=BS) {                                // Matrix initialization
    tareador_start_task ("init");
    for (int i=ii; i<ii+BS; i++)
        for (int j=0; j<NC; j++)
            M[i][j] = init(i, j);                             // <-- Cost ti
    tareador_end_task ("init");
}

tareador_start_task ("comp1");                                  // Begin computations
for (int i=0; i<BS; i++)
    for (int j=0; j<NC; j++)
        M[i][j] = comp1(M[i][j]);                             // <-- Cost tb
tareador_end_task ("comp1");

for (int ii=BS; ii<NR; ii+=BS) {                                // Final computations
    tareador_start_task ("comp2");
    for (int i=ii, int i0=0;
        i<ii+BS;
        i++, i0++)
        for (int j=0; j<NC; j++)
            M[i][j] = comp2(M[i0][j], M[i][j]);               // <-- Cost tf
    tareador_end_task ("comp2");
}
```

Let us assume the data sharing model explained in class based on a distributed memory architecture in which we consider that local memory accesses have no cost, but an access to data in different processors introduces a data-sharing overhead: the access time to remote data is determined by $t_{comm} = t_s + m \times t_w$, being t_s and t_w the "start-up" and sending time of an element, respectively, and being m the size of the message. Also, according to the data sharing model: at any time, each processor can simultaneously make one remote access to a different processor and serve one remote access from another processor.

Assume that the number of processors p divides the number of rows NR exactly; routines `init`, `comp1` and `comp2` do not modify any memory position; the execution time of one iteration of the body of the most internal loops is t_i , t_b and t_f respectively; tasks are scheduled to processes following the *owner-computes rule* so that a task will be executed by the processor who owns the memory that holds the output of that task; the matrix is already distributed in the memory of each processor when the computation starts; the resulting matrix remains distributed and there is no final communication to a single processor; the strategy used for decomposing matrix M follows a *block row distribution*: the matrix M is distributed so that each processor has $\frac{NR}{p}$ consecutive rows. **We ask** you to:

1. Draw a Task Dependence Graph (TDG) for the case where $p = 4$;
2. Draw a timeline for the execution with $p = 4$;
3. Provide a general expression that determines the parallel execution time on p processors (T_p): express T_p as a function of p , NR , NC , t_i , t_b , t_f , t_s and t_w .

Solution:

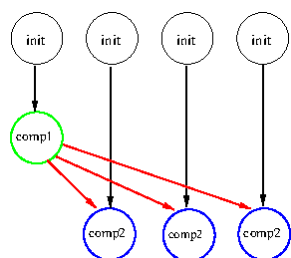


Figure 1: TDG

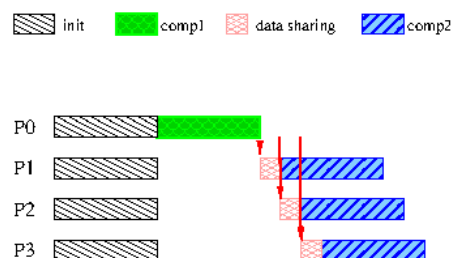


Figure 2: Timeline

$$T_p = t_{data_sharing} + t_{calc}$$

$$t_{data_sharing} = (t_s + \frac{NR}{P} \times NC \times t_w) \times (P - 1)$$

$$t_{calc} = \frac{NR}{P} \times NC \times t_i + \frac{NR}{P} \times NC \times t_b + \frac{NR}{P} \times NC \times t_f = \frac{NR}{P} \times NC \times (t_i + t_b + t_f)$$

Problem 2 (2.5 points)

Assume a NUMA (non-uniform memory architecture) multiprocessor system composed of 4 NUMA-nodes ($node_{0:3}$), each node with 8 GBytes of main memory and 4 cores ($core_{0:3}$). Each core has only one level of cache memory of 2 MBytes (with cache lines of 32Bytes). The system includes all the necessary mechanisms (seen in class) to keep the memory coherence inside a NUMA-node and between NUMA-nodes.

- (0.5 points) In order to support a write-invalidate MSI cache-coherence protocol within each NUMA-node, how many additional bits should be used in each cache line, and what are their role and possible values? How many bits in total per NUMA-node (only within each NUMA-node) are used for that purpose?

Solution:

We need 2 additional bits per cache line in order to store the state of the cache line. Those 2 bits can code up to four possible states:

- I invalid
- S shared (two or more nodes may have clean copies)
- M modified (dirty)
- Note that a fourth state can be coded with 2 bits. However, it is not necessary for the MSI protocol.

The number of cache lines is calculated dividing the memory in a cache (2MB) by the cache line size (32B). Therefore: Number of cache lines = $\frac{2MB}{32B} = 2^{16} = 64K$ entries

We need 2 bits per entry. And there are a total of 4 cache memories per NUMAnode.

Thus, the total number of bits amounts to:

$$\frac{2MB}{32B} \text{ entries/cache} \times 2 \text{ bits/entry} \times 4 \text{ caches/NUMAnode} = 2^{19} \text{ bits} = 512K \text{ bits/NUMAnode}$$

- (0.5 points) In order to support a directory-based write-invalidate MSU cache-coherence protocol between NUMA-nodes, how many bits should be used in the directory for each line in main memory, and what are their role and possible values? How many bits in total per NUMA-node (directory) are used for that purpose?

Solution:

The home NUMA-node is in charge of the coherence of its physical memory lines by means of the directory entries. A directory entry stores the line state and the identities of other NUMA-nodes sharing this memory line.

Bits per directory entry:

- Presence bits (nodes currently having the line), 1 bit per NUMA-node: i.e. 4 bits
- State bits (to track the state of memory lines): 2 bits

Those 2 bits can code up to four possible states:

- U uncached, not valid in any cache
- S shared (two or more nodes may have clean copies)
- M modified (dirty)
- Note that a fourth state can be coded with 2 bits. However, it is not necessary for the MSI protocol.

Total: 6 bits per directory entry.

The number of entries in the directory structure per NUMA-node is calculated dividing the memory in a NUMA-node (8GB) by the memory line size (32B). Therefore:

$$\text{Number of Directory Entries} = \frac{8GB}{32B} = 2^{28} = 256 \text{ Mega entries}$$

Thus, with 6 bits per directory entry, the total number of bits per NUMA node will be:

$$\frac{8GB}{32B} \text{ entries/NUMA node} \times 6 \text{ bits/entry} = 3 \times 2^{29} = 1536 \text{ Mega bits/NUMA node}$$

3. (1.5 points) Consider the following parallel program executed on only two cores (processors) inside $node_0$ of the previous multiprocessor system:

```
...  
double A[M][N];  
...
```

Core 0: Update even-numbered rows

```
for ( j = 0 ; j < M ; j += 2 )  
    for ( k = 0 ; k < N ; k++ )  
        A[j][k] = f(j,k);
```

Core 1: Update odd-numbered rows

```
for ( j = 1 ; j < M ; j += 2 )  
    for ( k = 0 ; k < N ; k++ )  
        A[j][k] = g(j,k);
```

Assume that matrix A is stored in main memory of $node_0$ (without copies of any of its elements in the caches of the NUMA nodes), that each element of matrix A is 8 Bytes and that the first element of A is aligned on a cache line boundary, $N = 2$ and M is a value multiple of 2.

- (a) What would be the maximum number of invalidations that would be sent through the bus expressed as a function of M ? What is causing such a memory coherence problem?

Solution:

For $N = 2$, every two rows fall in the same cache line, causing up to 3 invalidations due to false sharing for every pair of lines. Thus, the number of invalidations is $3 \times M/2$.

- (b) Modify the declaration of matrix A so that you do not have to change the code and avoid the previous coherence problem?

Solution:

We add enough padding to the second dimension of matrix A to make the size (in bytes) of a row of A equal or multiple of the number of bytes of a cache line.

```
...  
#define PADDING ((CACHE_LINE_SIZE - (N*sizeof(double) % CACHE_LINE_SIZE))/sizeof(double))  
double A[M][N+PADDING];  
...
```

Problem 3 (2.5 points)

Given the following code:

```
#define SIZE_INDEX 256
#define N 1024*1024*1024

void histogram(unsigned int *S, int n, unsigned int *index) {
    unsigned int i, tmp;
    for (i=0; i<n; i++) {
        tmp = S[i]%SIZE_INDEX;
        index[tmp]++;
    }
}

void main() {
    unsigned int S[N];
    unsigned int index[SIZE_INDEX];
    ... // Here we have initialized S to random numbers and index to 0's
    histogram(S, N, index);
    ...
}
```

Write two different OpenMP parallel implementations of function `histogram` using the strategies presented below. You can modify the sequential code, add local variables and use any OpenMP pragma (except `omp for` worksharing-loop construct) and function you may need. Both implementations should minimize the use of synchronizations and load unbalance between threads during the processing of vector S .

1. (1 point) Cyclic Data Decomposition of the Output vector `index`.

Solution:

```
#define SIZE_INDEX 256
#define N 1024*1024*1024
...

void histogram(unsigned int *S, int n, unsigned int *index)
{
    unsigned int i, tmp;

    #pragma omp parallel private(i, tmp)
    {
        int id = omp_get_thread_num();
        int num_threads = omp_get_num_threads();

        for (i=0; i<n; i++) {
            tmp = S[i]%SIZE_INDEX;
            if ((tmp%num_threads)==myid)
                index[tmp]++;
        }
    }

    ...
}
```

2. (1.5 points) Block Data Decomposition of the Input vector S .

Solution:

```
#define SIZE_INDEX 256
```

```

#define N 1024*1024*1024
...

void histogram(unsigned int *S, int n, unsigned int *index)
{
    unsigned int i, tmp;

    #pragma omp parallel private(i, tmp)
    {
        int myid          = omp_get_thread_num();
        int num_threads   = omp_get_num_threads();
        int i_start       = myid * (N/num_threads);
        int i_end         = (myid+1) * (N/num_threads);
        int res           = N%num_threads;
        if (res)
        {
            i_start = i_start + (myid<res)?myid:res;
            i_end   = i_end   + (myid<res)?myid+1:res;
        }

        unsigned int local_index[SIZE_INDEX];
        for (i=0;i<SIZE_INDEX;i++)
            local_index[i]=0;

        for (i=i_start;i<i_end;i++) {
            tmp = S[i%SIZE_INDEX];
            local_index[tmp]++;
        }

        for (i=0;i<SIZE_INDEX;i++)
        {
            #pragma omp atomic
            index[i] += local_index[i];
        }
    }
}

```

Problem 4 (2.5 points)

We ask you to write two additional parallel OpenMP implementations of the code that computes the histogram, this time using *task decomposition* strategies:

1. (1 point) Write an efficient OpenMP parallel version of the histogram program in Problem 3 using *explicit tasks*.

Solution: A possible implementation:

```
...
void histogram(unsigned int *S, int n, unsigned int *index)
{
    unsigned int i, tmp;

    #pragma omp parallel
    #pragma omp single
    #pragma omp taskloop private(tmp)
    for (i=0; i<n; i++) {
        tmp = S[i]%SIZE_INDEX;
        #pragma omp atomic
        index[tmp]++;
    }
}
```

2. (1.5 points) Write a *recursive tree task decomposition* with *cutoff* based on the depth of the recursivity for the following code:

```
#define SIZE_INDEX 256
#define N 1024*1024*1024
#define BASE_SIZE 512
#define CUTOFF 3

void histogram(unsigned int *S, int n, unsigned int *index) {
    unsigned int i, tmp;
    for (i=0; i<n; i++) {
        tmp = S[i]%SIZE_INDEX;
        index[tmp]++;
    } }

void histogram_rec(unsigned int *S, int n, unsigned int *index) {
    unsigned int n2=n/2;

    if (n > BASE_SIZE) {
        histogram_rec(S, n2, index);
        histogram_rec(&S[n2], n-n2, index);
    } else {
        histogram(S, n, index);
    }
}

void main() {
    unsigned int S[N];
    unsigned int index[SIZE_INDEX];
    ...
    // Here we have initialized S to random numbers and index to 0's
    histogram_rec(S, N, index);
    ...
}
```

Solution: A possible implementation:

```
...
void histogram(unsigned int *S, int n, unsigned int *index) {
    unsigned int i, tmp;
    for (i=0; i<n; i++) {
        tmp = S[i]%SIZE_INDEX;
        #pragma omp atomic
        index[tmp]++;
    }
}

void histogram_rec(unsigned int *S, int n, unsigned int *index, int depth) {
    unsigned int i, tmp, n2;

    if (n > BASE_SIZE) {
        n2 = n/2;
        if ( !omp_in_final() ) {
            #pragma omp task final(depth>=CUTOFF)
            histogram_rec(S, n2, index, depth+1);
            #pragma omp task final(depth>=CUTOFF)
            histogram_rec(&S[n2], n-n2, index, depth+1);
        }
        else {
            histogram_rec(S, n2, index, depth+1);
            histogram_rec(&S[n2], n-n2, index, depth+1);
        }
    }
    else {
        histogram(S, n, index);
    }
}

void main() {
    unsigned int S[N];
    unsigned int index[SIZE_INDEX];
    ...
    // Here we have initialized S to random numbers and index to 0's
    #pragma omp parallel
    #pragma omp single
    histogram_rec(S, N, index, 0);
    ...
}
```