

2202-COL380 Major

Chinmay Mittal

TOTAL POINTS

40.5 / 77

QUESTION 1

1 Q1 1 / 2

+ 2 pts Correct

✓ + 1 pts *pragma omp critical is correctly explained*

+ 1 pts Significance of Anonymous is explained correctly.

+ 0 pts Incorrect / Unattempted

4 0 / 3

+ 3 pts Correct

+ 2 pts memory flush with barrier

+ 1 pts partially correct

✓ + 0 pts *Incorrect / Unattempted*

QUESTION 2

2 4 / 4

✓ + 1 pts *Correct Explanation for Processor consistent*

✓ + 1 pts *Correct Explanation for Sequentially consistent*

✓ + 2 pts *Correct example*

+ 0 pts wrong

QUESTION 5

5 2 / 2

✓ + 1 pts *Correct Definition of Lock-free*

Lock free protocols ensure that there is system wide progress even if individual threads might be suspended

✓ + 1 pts *Correct Explanation*

Protocols with a critical section may not always be lock free because the entity which has acquired the lock and is in the critical section might be suspended or stuck. In that case no other thread would be able to make progress.

+ 0 pts Incorrect

QUESTION 3

3 1 / 4

+ 1 pts Correct definition of Atomic

+ 1 pts Correct definition of Instantaneous

+ 1 pts explanation for why the instruction is not atomic

✓ + 1 pts *explanation for why the instruction is not instantaneous*

+ 0 pts wrong

QUESTION 6

6 2 / 3

+ 3 pts Fully correct

+ 0 pts Not Attempted/Incorrect

✓ + 1 pts ***Not** lock free since it is a mutual exclusion algorithm. Lock-free by definition must ensure system-wide progress — at least one ready thread, no matter the state of the other threads.*

QUESTION 4

✓ + 1 pts **Not** wait free, since by definition wait-freeness ensures that every ready thread progresses, no matter the state of other threads.

+ 1 pts Starvation free, since algorithm ensures every thread to get an **opportunity** to proceed as long as all threads are scheduled.

QUESTION 7

7 1.5 / 2

+ 2 pts Parallelism at **every level** i.e network, disk, and server level. There shouldn't be a single point of throughput.

✓ + 1 pts Partial correct

+ 0 pts Not attempted/incorrect

+ 0.5 Point adjustment

- 1 It should be parallel at every level, not just for read and writes (see rubrics)

QUESTION 8

8 3 / 3

✓ + 3 pts The map function generates (key, value) pairs. There can be multiple pairs for a single key. The map-reduce framework must sort these pairs into bins by key.

+ 1.5 pts The student mentions the separation of keys, however, the explanation is unclear.

+ 0 pts Incorrect or not attempted.

QUESTION 9

9 3.5 / 4

+ 4 pts Fully correct

✓ + 2 pts Knows linearizability and sequential consistency, and the difference in terms of globally visible time.

+ 1 pts General idea of ordering is there, but the distinction on global time is not clear

✓ + 1 pts Usefulness of linearizability

+ 1 pts Usefulness of sequential consistency

+ 0 pts Incorrect or Didn't attempted

+ 0.5 Point adjustment

- 2 use case of sequential consistency is missing

QUESTION 10

10 4 / 4

✓ + 4 pts Fully Correct answer

+ 1 pts Definition of common CRCW

+ 1 pts Definition of Priority CRCW

+ 1 pts Why common CRCW is also a priority CRCW algorithm

+ 1 pts Why priority CRCW may not be common CRCW algorithm

+ 0 pts No/Incorrect answer

QUESTION 11

11 5 / 5

+ 0 pts Model Solution:

```
void MPI_BARRIER(MPI_Comm comm) {  
    int barrier_value = 1;  
    int result;  
    MPI_Allreduce(&barrier_value, &result, 1,  
        MPI_INT, MPI_SUM, comm);  
}
```

Collectives like MPI_AllReduce, MPI_AllGather and MPI_Alltoall can be used as barriers, as these are blocking calls involving all entities and none may return until the receive buffer has their data —

and data needs to be received from all, which can only happen after all have called that collective.

The Communicator and send and receive buffers are the important parameters. (Root is also important for non-barrier based collectives)

Correct Implementation

✓ + 2 pts Implemented the barrier using the above method

+ 2 pts Used other mechanisms which also work (for eg each thread can call broadcast and wait for the root, the root can call gather and call broadcast again when it has received from everyone, to signal they can proceed)

✓ + 2 pts Argued for correctness

✓ + 1 pts Listed important parameters

Partially correct with collectives

+ 2 pts Used collectives but minor mistakes result in an almost correct barrier

+ 1 pts Tried to justify

Without collectives

+ 1.5 pts Used send/receive operations to construct a barrier

+ 0.5 pts Tried to justify

+ 0 pts Incorrect/Unattempted

QUESTION 12

12 0 / 5

+ 5 pts Correct

- 1 pts Early scheduling not ensured (task 4/5/7 go early. 3/6/9 go last)

- 2 pts incorrect use of depend clauses

✓ + 0 pts Incorrect / Unattempted

QUESTION 13

13 5 / 5

✓ + 1 pts Identifying the critical paths: '(4, 5, 8, 9), (4, 7, 8, 9)` and hence should be scheduled accordingly.

✓ + 4 pts **Parallel streams**

Synchronization across streams is guaranteed by busy-waiting on a flag or synchronizing individual streams.

+ 4 pts **Two streams: Serial**

- `(1, 2, 3)` in stream 1. Rest in stream 2.

- Serialize stream 2 as follows:

4 before 5 & 7.

5 & 7 before 8.

8 before 6 & 9.

+ 2 pts **Without streams**

Order guaranteed using global synchronization.

Without using Streams, it becomes all sequential. The sync calls are not useful. While maintaining correct dependency order is okay, it becomes a trivial problem. Scheduling earlier and allowing concurrency is important.

+ 0 pts Incorrect or not attempted.

+ 5 pts Cuda graph

QUESTION 14

14 2 / 6

+ 0 pts Non-answer

✓ + 0.5 pts Identified any race

+ 2 pts Eliminate Race on testing of parent and update.

- 1 pts Elimination of race on parent is not lock-free (Remember critical section/barriers do not

make it lock-free.)

+ 1 pts Shows lock-freeness on parent test correctly

✓ + 2 pts Eliminate race on front

✓ - 1 pts Elimination of Race on front is incorrect/not lock-free

+ 0.5 pts Show lock-free elimination of race on front

+ 0.5 Point adjustment

3 atomic?

efficiency

✓ - 0.5 pts Demonstrates $\log n$ knowledge, but incorrect computation

✓ + 1 pts $\log \log n$ version and its Iso-efficiency

✓ - 0.5 pts Demonstrates $\log \log n$ knowledge, but incorrect computation

✓ + 1 pts Accelerated cascading version and Iso-efficiency

✓ - 0.5 pts Demonstrates accelerated cascading knowledge, but incorrect computation

✓ + 1 pts comparisons

✓ - 0.5 pts Incomplete comparisons

4 Is-eff cannot be a fn of n.

QUESTION 15

15 2.5 / 6

+ 0 pts Non-answer

✓ + 2 pts Problem with Kernel: separate flags per block

+ 1 pts Problem in Kernel: initialized flag in each thread: could be after another thread has entered critical section.

+ 1 pts Problem with Blocks: Initialized flag per thread

+ 2 pts Problem with Warp: Usually no divergence in flag=0, but lock loop can have divergence and could deadlock on CAS if the failed threads continue to schedule. unlock should be updating flag globally -- fence required

- 1 pts Inaccurate/Incomplete Warp explanation

+ 0.5 Point adjustment

QUESTION 17

17 0 / 6

+ 6 pts Completely Correct answer

+ 2 pts O(1) time Correct answer but not work-optimal.

+ 2 pts Correct work optimal but not O(1) time

+ 3 pts Correct answer with partitioning approach but not work optimal.

✓ + 0 pts No/Incorrect answer

QUESTION 18

18 0 / 7

✓ + 0 pts Non-answer

+ 4 pts Data collection and multiplication

- 2 pts Missing collectives

- 2 pts Partially correct data collection with collectives

- 3 pts Partially correct collection without proper collectives.

+ 2 pts Full Cost analysis

QUESTION 16

16 4 / 6

+ 0 pts Non-answer

✓ + 2 pts $t(n,p) = ?$

✓ + 1 pts $\log n$ version and its time on P Iso-

- 1 pts Erroneous analysis
- + 1 pts Result in the right place
- 0.5 pts Redistributing result

-2 marks for missing or illegible entry no./name on any sheet.

COL 380 MAJOR EXAM

SEMESTER II 2022-2023

2 hours, 18 questions, 75 marks (+2 extra credits)

Maximum marks are listed in [] for each question. Justify your answers: most marks are designated for correct justification. Write legibly to receive credit. Follow IITD/COL380 Academic Integrity Code.

1. [2] What does the following stand for in OpenMP? What is the significance of 'Anonymous'?

#pragma omp critical (Anonymous)

These are named critical sections, in open MP. if we use ~~no~~ names for critical sections then at one time only one thread can be ~~execute~~ executing a critical section with a particular name. two threads can execute critical sections with different names simultaneously.

Anonymous means no name. Out of all sections with no names, only one of them can be active at a given time.
(at most)

2. [4] Explain Processor Consistent and Sequentially consistent execution. Show an example that is Processor consistent but sequentially inconsistent.

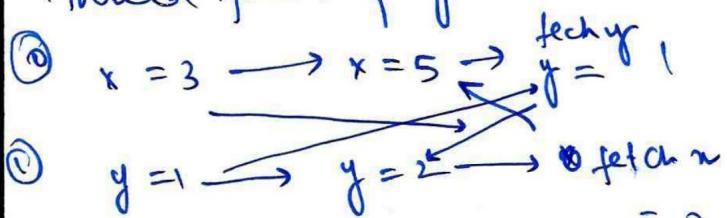
Processor consistent \Rightarrow All writes made by a thread should be seen in that order by all other threads. Also there should be a globally consistent sequential view of writes made to a particular variable by all threads.

Sequentially consistent \Rightarrow There should exist a consistent sequential ordering of instructions of all threads such that the instructions of every thread follow program order in this global sequence.

Sequentially inconsistent because of the cycle.

processor won't

↳ thread 0's view



thread 1's view

$x = 3 \rightarrow x = 5 \rightarrow y = 1 \rightarrow$ fetch $y = 2$

$y = 1 \rightarrow y = 2 \rightarrow n = 3 \rightarrow$ fetch $n = 3 \rightarrow x = 5$

3. [4] Is an architecture (ISA) instruction like `mov (%rbp), %eax` instantaneous? Atomic? Explain.

These instructions are ~~set~~ instantaneous. They not. They take a finite time to complete. Which involves the time the hardware takes to transfer data from one register to the other using the ~~wire~~ wire interconnect between these registers. But the hardware ensures that these instructions are atomic w.r.t. the registers involved. No other thread can access these registers when this instruction is being executed, which makes it atomic w.r.t. these registers and this atomicity is ensured by the hardware.

4. [3] I add a missing `#pragma omp barrier` in a parallel region in my code, but notice a significant increase in the number of cache transaction shown by my profiler. What could the reason(s) be?

Threads executing on different cores can share the same cache and can also have shared memory. If there is no barrier then there is no synchronization and all threads can be executing at their own different pace, and accesses to shared variables by different threads can occur at very different global times not allowing them to benefit from caching. On the other hand barrier imposes synchronization which can potentially lead to shared variables being accessed at similar times (global) and thus can benefit from caching.

say \rightarrow read₀ \rightarrow read₁
When TID₀ executes this instruction first
 x comes into cache, and TID₁ executes it close enough (in global time) it will find it in cache which might not be the case without synchronization.

5. [2] Explain why protocols with critical section are NOT lock-free. (Define lock-free and explain.)

Lock free algorithms \Rightarrow some synchronizer must progress independent of the scheduler.

Consider a program with one thread in the critical section and others waiting to execute the critical section.
 if the scheduler chooses to repeatedly schedule one of the waiting threads, then nobody makes progress.
 Since the waiting threads keep waiting for the thread in the critical section to exit and this thread is never scheduled.
 So depending on who the scheduler ~~schedules~~ schedules, nobody might make progress.

6. [3] Which of the following describe the Bakery algorithm: wait-free/lock-free/starvation-free. Explain

Bakery algorithm relies on the scheduler and is not wait-free or lock-free, $\xrightarrow{\text{it is possible that}}$ nobody can progress depending on the scheduler.
 if the thread executing the critical section is not scheduled repeatedly, then nobody makes progress. All the waiting threads keep waiting for the thread in the critical section which never exits because it is never scheduled by the scheduler.
 But it is starvation free, with the proper scheduling algorithm, the thread in the critical section will exit and all threads will make progress one by one.
 \Rightarrow everyone progresses with the help of the scheduler \rightarrow starvation free.

7. [2] What makes a parallel file system parallel?

The ability to read multiple files concurrently, at different locations by different processors.

1

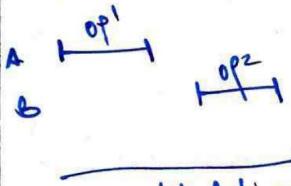
8. [3] In the Map-Reduce programming model, what must the underlying framework do after the end user's Map() calls and before the Reduce() calls (in a message-passing environment).

In a message passing environment each processor will be responsible for reducing a particular key, which processor processes which key can be figured out by the id of that processor. It is possible that the ~~same key~~ ^{key} ~~is~~ ^{is} same id of that processor. Pairs of ~~same key~~ ^{same key} are speed across the processors $\langle \text{key}, \text{value} \rangle$ pairs. Each process needs to transmit a message to the (containing) correct processor responsible for that key, or transmit that it has no such pair. Once each processor gets all pairs of the key for which it is responsible it can perform the reduce function. Each processor needs to ensure that it gets all pairs corresponding to a key from all other processors.

9. [4] What is the difference between Linearizability and Sequential consistency? Explain how they are useful.

Sequential consistency ~~is~~ \Rightarrow there exists a consistent global ordering of all instructions executed by all threads, such that the instructions of every thread appear in program order in this global execution. The write to a variable must appear before the read of that variable of which it is a read effect. Without any writes to the same variable in between. Sequential consistency is about relative ordering without any notion of global time. Linearizability is a stronger condition where for every instruction, there is a start and end time (global clock) and the instruction i takes effect at some time b_i and ends at time e_i .

item.



op_2 is sequentially consistent, but is not linearizable since op_1 must end before op_2 starts. \Rightarrow thus

linearizability leads to additional constraints based on this time. And it is harder to achieve in hardware with real time. And it is harder to achieve in hardware with real time. But for certain applications these type of guarantees are important. e.g. A is check amount ~~no~~ balance and B is credit card payment

10. [4] Explain how a common-CRCW PRAM algorithm is also a Priority-CRCW algorithm, but not vice-versa.

Consider a common-CRCW PRAM algorithm. It is also a valid priority CRCW algorithm. The read steps and computation steps do not pose any problem if they are valid in common CRCW they are valid in priority CRCW. Writes to different locations also do not pose a problem. Consider n threads writing to the same location. All threads try to write the same value v to the memory location. After common CRCW write the memory location will be updated to v . These writes remain valid in priority CRCW, the top priority processor will write v to the memory location and hence all memory updates are correctly preserved.

A priority CRCW algorithm ~~might~~ not be a valid common-CRCW algorithm. At some concurrent write step many processors might be trying to write different values to the same location and the correctness of the algorithm depends on the fact that the highest priority processor writes the correct value. Since they might be writing different values, this is not a valid common CRCW algorithm.

11. [5] Implement *MPI_barrier* using MPI collective message passing calls (list important parameters.) Explain why it is correct.

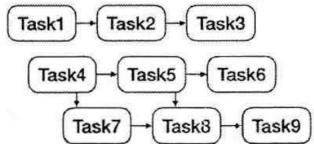
```
func Barrier (MPI_Comm) {  
    int size; int rank; MPI_Comm_rank (MPI_Comm, &rank);  
    MPI_Comm_size (MPI_Comm, &size);  
MPI MPI_Request arr[arr];  
    int recv_arr[arr] = new recv[arr];  
    for (int i = 0; i < arr; i++)  
    {  
        send MPI_Ibsert (&recv[arr], &MPI_Ibsert (&send, 1, MPI_INT, i, MPI_Comm,  
        if (i == rank) {  
            int send = 1;  
            MPI_Ibsert (&send, 1, MPI_INT, i, MPI_Comm,  
        } else {  
            MPI_Ibsert (&recv[i], 1, MPI_INT, i, MPI_Comm,  
        }  
        MPI_Waitall (request_arr, size);  
    }  
    return ;  
}
```

Each process will make n non-blocking broadcast calls each with root $1 \dots n$. Hence n groupwise matching occurs. ~~MPI~~ process will return only when it has received confirmation from all n processes. Using Ibsert & Waitall ~~processes~~ no implicit synchronization used by these calls. Otherwise we would use synchronization to implement synchronization.

Entry Number: 2010CS10336

Name: CHINMAY MITTAL

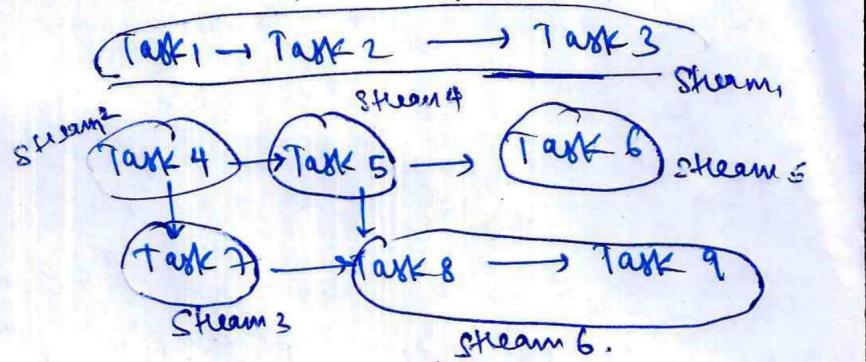
12. [5] Implement the given task graph in OpenMP (approximate code). Assume tasks are functions (with the same name). Do not define variables. Identify and ensure early scheduling of the critical path.



Q12
13. [5] Repeat Q15 for Cuda. Assume kernels with the same names as tasks. (Variables may be defined.)

We can use Cuda Streams to implement this task graph

cudaStream_t stream_1, stream_2,
stream_3, stream_4,
stream_5, stream_6;



cudaStreamCreate(&stream_1); cudaStreamCreate(&stream_2);
cudaStreamCreate(&stream_3); cudaStreamCreate(&stream_4);
cudaStreamCreate(&stream_5); cudaStreamCreate(&stream_6);
cudaStreamCreate(&stream_7); cudaStreamCreate(&stream_8);
cudaStreamCreate(&stream_9); cudaStreamCreate(&stream_10);

task_4 << params, 0, stream_2 >> f();

task_1 << params, 0, stream_1 >> f();

task_2 << params, 0, stream_1 >> f();

task_3 << params, 0, stream_1 >> f();

cudaStreamSynchronize(stream_2);

task_7 << params, 0, stream_3 >> f();

task_5 << params, 0, stream_4 >> f();

cudaStreamSynchronize(stream_3); cudaStreamSynchronize(stream_4);

task_8 << params, 0, stream_6 >> f();

task_9 << params, 0, stream_6 >> f();

cudaStreamSynchronize(stream_6); task_6 << params, 0, stream_5 >> f();

cudaStreamSynchronize(stream_5);

cudaStreamSynchronize(stream_6);

cudaStreamSynchronize(stream_1);

appropriate
block & grid
dim
parens are the
can be set to 1, 1 if only one thread
is required;

14. [6] Consider pseudo-code shown below for breadth-first traversal of a directed graph $G = (V, E)$, where $e \in E = (v_1, v_2)$, where $v_1, v_2 \in V$. Assume *forall* continues to schedule threads in parallel until members remain in the set. Identify and correct the race conditions in the pseudo-code below in a lock-free manner. (Prove lock-freeness.)

BFS(G):

```

forall v ∈ G.V:
    parent[v] = -1
Front = {random vertex r ∈ G.V}
parent[r] = r
Update(Front)
  
```

v ←

Update(front):

```

forall v1 ∈ front:
    forall v2 ∈ G.out_neighbor(v1):
        if (parent(v2) == -1):
            parent[v2] = v1; visit(v2)
            front += v2 // Add to front
  
```

v ←

Since front is a shared variable it can be potentially updated by many threads in parallel and lead to incorrect updates.
We can use CAS to update front in an atomic manner.

~~while~~ do { local ← front 3

new ← front + v2;

(~~while~~^{CAS} (front, local, new));

CAS depends on hardware and is updated independent of the scheduler. Hence remains lock free

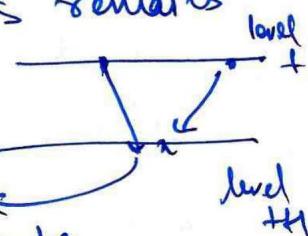
Updating parent or visit is not a race condition. The BFS remains correct irrespective of the order of execution of instructions

CAS ensures that the algorithm remains lock free

The last operation to succeed is correct both can be

valid answers, hence we don't need mutual exclusion.

We must ensure that depth + starts only when depth + is completely computed. Hence there is a need for synchronization after the inner for loop



15. [6] I wrote the following for mutual exclusion using CAS on Cuda for a kernel, but failed. I tried limiting it to threads within a block and still failed. I finally tried it only within threads of a warp, and failed again. Explain.

```
__device__
void unlock(int &flag) {
    flag = 0;
}
```

```
__device__
void lock(int &flag) {
    int old;
    do old = atomicCAS(&flag, 0, 1);
    while (old != 0);
}
```

```
__global__
void akernel(void) {
    __shared__ int flag;
    flag = 0;
    lock(flag);
    critalsection();
    unlock(flag);
}
```

The flag variable is a shared variable and is shared b/w threads of a block. It is possible that threads of different blocks both acquire the lock since the flag variable copies they maintain are different variables and both can be updated to one.

There is ~~no memory synchronization performed~~ by CAS. It is possible that each thread maintains the shared flag variable in its register and updates it their without other threads in the block noticing that out. We need to use volatile to ensure flag is not kept in registers.

16. [6] Compare the scaling behavior of $O(\log n)$ -time binary-tree based minima-finding, $O(\log \log n)$ non-optimal minima-finding, and the accelerated cascading based minima-finding in terms of their iso-efficiency functions.

Consider the $O(\log(n))$ time algorithm

$$t(n) = O(\log(n))$$

$$w(n) = O(n) \Rightarrow t(n,p) = \left\lfloor \frac{w(n)}{p} \right\rfloor + t(n)$$

$$= O\left(\frac{n}{p} + \log n\right)$$

$$t_1(n) = O(n)$$

for iso efficiency we have

$$\frac{n}{p + \left(\frac{n}{p} + \log n\right)} \approx K$$

~~Since~~ For $O(\log \log n)$ non optimal solution

$$t(n) = O(\log \log n)$$

$$w(n) = O(n \log \log n)$$

$$t(n,p) = \frac{w(n)}{p} + t(n) \approx O\left(\log \log n + \frac{n}{p} \log \log n\right)$$

$$\frac{o(n)}{O\left(\frac{n}{p} \log \log n + \log \log n\right)} \approx \frac{o(n)}{n \log \log n + p \log \log n}$$

$$\frac{t_1(n,p)}{p + t(n,p)} \approx \Theta(p) \approx K$$

$\log n$ is $\omega(n)$

p is $\Omega\left(\frac{n}{\log n}\right)$

since $t_1(n,p)$ is $O(n)$
iso efficiency $\approx \omega\left(\frac{n}{\log n}\right)$

Weakly
scalable

not scalable
irrespective of p will keep O dening.

Accelerated Cascading

$$\frac{t(n,p)}{p(t(n,p))} \approx t$$

$$t(n) \approx O(\log \log n)$$

$$w(n) \approx O(n)$$

$$\frac{n}{n + p \log \log n} \approx t$$

$$t(n,p) = \frac{w(n)}{p} + t(n)$$

$$\approx \Theta\left(\frac{n}{p} + \log \log n\right)$$

$\log \log n$ is $\omega(n)$
 p is $\Omega\left(\frac{n}{\log \log n}\right)$

weakly scalable
less scalable

than $O(\log n)$
algorithm

17. [6] Provide a work-optimal $O(1)$ -time common-CRCW PRAM algorithm to find the first occurrence (i.e., the smallest index) of value v in array A. (v may appear 0 or more times in A.)

A sequential algorithm can simply iterate through the array and keep track of the first occurrence of $v \Rightarrow$ best sequential time $O(n) \Rightarrow$ optimal work $O(n)$.

Assuming priority CRCW per index
each processor tries to write its index (if it has the element)
on the answer location
↓
read $A[i]$

Entry Number:

2020CS10336

Name:

18. [7] Consider two $n \times n$ matrices A and B, initially distributed *block-wise* among $p \times p$ ranks of a communicator. Each block is $b \times b$ (and $n = bp$). Provide MPI-like code to implement $C = Ax B$. C also must be distributed block-wise like A and B. Use collectives where possible (list important parameters). Measure message passing cost.

ROUGH PAGE. WILL NOT BE GRADED